

UNIVERSIDAD DE GRANADA



APORTACIONES A LOS MÉTODOS DE
ESTIMACIÓN DE PARÁMETROS
LINEALES Y NO LINEALES CON
INFORMACIÓN AUXILIAR

TESIS DOCTORAL

Directora:

Profa. Dra. D^a. María del Mar Rueda García

Doctorando:

Juan Francisco Muñoz Rosas

DEPARTAMENTO DE ESTADÍSTICA E INVESTIGACIÓN OPERATIVA

Granada 2006

APORTACIONES A LOS MÉTODOS DE ESTIMACIÓN DE PARÁMETROS LINEALES Y NO LINEALES CON INFORMACIÓN AUXILIAR

Memoria presentada por Juan Francisco Muñoz Rosas para aspirar al Doctorado Europeo en Ciencias Estadísticas.

Vº.Bº. de la Directora de Tesis

Fdo.: Profa. Dra. D^a. María del Mar Rueda García
Departamento de Estadística e Investigación Operativa

UNIVERSIDAD DE GRANADA

Granada 2006

AGRADECIMIENTOS

Estas líneas están dedicadas a todas aquellas personas que me han apoyado y ayudado a realizar la presente tesis doctoral.

Quisiera manifestar mi más profundo agradecimiento a la profesora María del Mar Rueda García por su constante e inestimable aportación a este trabajo. El asesoramiento científico y el constante ánimo que me ha transmitido han sido fundamentales para llevar a cabo esta tesis doctoral, puesto que gracias a su ayuda he adquirido todos mis conocimientos sobre muestreo en poblaciones finitas.

Agradezco también enormemente la confianza que ha depositado en mi, el tiempo dedicado y la paciencia mostrada durante este periodo. Asimismo, me siento agradecido por haberme enseñado a trabajar e investigar y espero seguir mejorando estas aptitudes junto a ella. Realmente es para mi un honor haber podido contar con María del Mar Rueda García como directora de tesis. Sin su ayuda y orientación este trabajo nunca hubiera sido posible.

Gracias a mis compañeros del Departamento de Estadística e Investigación Operativa y del Departamento de Métodos Cuantitativos para la Economía y la Empresa por su apoyo.

Finalmente agradezco a mi familia el apoyo que me han dado. Una mención especial se merece mi mujer por su continuo ánimo, apoyo y paciencia que ha mostrado conmigo. Sobretudo ella sabe el esfuerzo y trabajo que ha requerido la elaboración de esta tesis doctoral.

Índice general

1. Introducción	1
1.1. Problemas planteados	1
1.2. Antecedentes	2
1.3. Objetivos científicos y aportes a la teoría del muestreo	5
1.4. Notación y conceptos básicos	10
2. El método de verosimilitud empírica	13
2.1. Introducción	14
2.2. Estimación de la media poblacional	18
2.2.1. Estimadores basados en el diseño muestral	18
2.2.2. Propiedades teóricas	37
2.2.3. Estimadores modelo-calibrados	42
2.2.4. Propiedades teóricas	46
2.3. Tratamiento de datos faltantes	48
2.3.1. Introducción	48
2.3.2. Estimador propuesto	52
2.3.3. Propiedades teóricas	54
2.3.4. Propiedades empíricas	57
2.4. Estimación de la función de distribución	65
2.4.1. Introducción	65
2.4.2. Algunos estimadores de la función de distribución	67
2.4.3. Estimador propuesto modelo-asistido	78
2.4.4. Propiedades teóricas	81
2.4.5. Propiedades empíricas	88
3. Aportaciones a la estimación de cuantiles	99
3.1. Introducción	99
3.2. Estimadores bajo muestreo bifásico	102
3.2.1. Introducción	103
3.2.2. Estimadores propuestos	104
3.2.3. Propiedades teóricas	107

3.2.4.	Propiedades empíricas	112
3.2.5.	Aplicación al muestreo estratificado	120
3.2.6.	Propiedades teóricas	125
3.2.7.	Propiedades empíricas	127
3.3.	Estimadores bajo muestreo en dos ocasiones sucesivas	133
3.3.1.	Introducción	133
3.3.2.	Muestreo con probabilidades desiguales	135
3.3.3.	Propiedades teóricas	138
3.3.4.	Propiedades empíricas	143
3.3.5.	Generalización a múltiples variables auxiliares	150
3.3.6.	Propiedades teóricas	152
3.3.7.	Propiedades empíricas	155
3.4.	Estimadores bajo el método de verosimilitud empírica	164
3.4.1.	Antecedentes	164
3.4.2.	Aplicación a la estimación de líneas de pobreza	166
3.4.3.	Estimadores propuestos modelo-asistidos	170
3.4.4.	Propiedades. Estimación de la varianza	173
3.4.5.	Propiedades empíricas	174
4.	Discusión	180
4.1.	Conclusiones y valoración de resultados	180
4.2.	Perspectivas y futuras líneas de investigación	184
5.	Redacción para aspirar a la mención europea en el título de Doctor	186
5.1.	Abstract	186
5.2.	Pseudo empirical likelihood method in the presence of missing data	188
5.2.1.	Introduction	188
5.2.2.	The proposed class of estimators	190
5.2.3.	Asymptotic properties	194
5.2.4.	The optimal estimators of the proposed class	196
5.2.5.	An empirical study	197
5.3.	Quantile estimation under two phase sampling	204
5.3.1.	Introduction	204
5.3.2.	Quantile direct estimation	205
5.3.3.	Estimation using auxiliary information	207
5.3.4.	Quantile estimation under two-phase sampling for stratification	212
5.3.5.	An empirical study	214
5.4.	Conclusions	225

Bibliografía	229
A. Descripción de poblaciones finitas	246
A.1. Poblaciones naturales	246
A.1.1. Fam1500	246
A.1.2. Counties	247
A.1.3. Hospitals	250
A.1.4. Murthy	251
A.1.5. Turismos	252
A.1.6. ECPF1997	254
A.2. Poblaciones simuladas	256
A.2.1. Pop06, Pop07, Pop08 y Pop09	256
A.2.2. Pob098 y Pob080	258
B. Procedimientos de estimación y selección de unidades	261
B.1. Métodos de muestreo probabilísticos	261
B.1.1. Muestreo aleatorio simple	261
B.1.2. Muestreo de Midzuno	262
B.1.3. Muestreo de Lahiri	262
B.1.4. Muestreo de Poisson	263
B.2. Diseños muestrales	264
B.2.1. Muestreo estratificado	264
B.2.2. Muestreo bifásico	265
B.2.3. Muestreo bifásico aplicado a la estratificación	266
B.2.4. Muestreo en ocasiones sucesivas	266
B.3. Tipos de inferencia en poblaciones finitas	267
B.3.1. Aproximación basada en el diseño	268
B.3.2. Aproximación basada en modelos	268
B.3.3. Aproximación modelo-asistida	269
B.3.4. Aproximación modelo-calibrada	270
C. Programación de estimadores mediante el software R	272
C.1. Introducción	273
C.1.1. Funciones complementarias	273
C.1.2. Métodos de muestreo	275
C.1.3. Funciones de distribución básicas	277
C.1.4. Cuantiles básicos	281
C.2. El método de verosimilitud empírica	286
C.2.1. Tratamiento de datos faltantes	286
C.2.2. Estimación modelo-asistida de la función de distribución	290
C.3. Aportaciones a la estimación de cuantiles	300

C.3.1. Muestreo en dos ocasiones sucesivas con probabilidades desiguales	300
C.3.2. Muestreo en dos ocasiones sucesivas con múltiples vari- ables auxiliares	313
C.3.3. Muestreo bifásico	317
C.3.4. Muestreo bifásico aplicado a la estratificación	322
C.3.5. Estimación modelo-asistida usando verosimilitud empírica	336

Índice de figuras

2.1.	Eficiencia Relativa para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se toman muestras de tamaño $n = 200$	59
2.2.	Eficiencia Relativa para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se considera la población Fam1500 y muestras de tamaño $n = 150$	60
2.3.	Eficiencia Relativa para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se considera la población Hospitals y muestras de tamaño $n = 100$	60
2.4.	Sesgo Relativo para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_w^{AC} (estándar), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se toman muestras de tamaño $n = 200$	63
2.5.	Sesgo Relativo para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_w^{AC} (estándar), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se considera la población Fam1500 y muestras de tamaño $n = 150$	64
2.6.	Sesgo Relativo para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_w^{AC} (estándar), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se considera la población Hospitals y muestras de tamaño $n = 100$	64
2.7.	Eficiencia Relativa de distintos estimadores en las poblaciones Pob098 y Pob080.	91
2.8.	Eficiencia Relativa de distintos estimadores en la población Murthy.	92
2.9.	Sesgo Relativo Medio de distintos estimadores en las poblaciones Pob098, Pob080 y Murthy.	94
2.10.	Eficiencia Relativa Media de distintos estimadores en las poblaciones Pob098, Pob080 y Murthy.	95

2.11. Diagramas de cajas con bigotes de las Desviaciones Absolutas Medias de distintos estimadores en las poblaciones Pob098 (con $n = 100$), Pob080 (con $n = 100$) y Murthy (con $n = 50$).	97
3.1. Eficiencia Relativa para la población Fam1500 y bajo el diseño muestral <i>Mas.Midzuno</i> . $n' = 150$	114
3.2. Eficiencia Relativa para la población Fam1500 y bajo el diseño muestral <i>Mas.Poisson</i> . $n' = 150$	115
3.3. Eficiencia Relativa para la población Counties y bajo el diseño muestral <i>Mas.Midzuno</i> . $n' = 150$	116
3.4. Eficiencia Relativa para la población Counties y bajo el diseño muestral <i>Mas.Poisson</i> . $n' = 150$	117
3.5. Sesgo Relativo en porcentaje para la población Fam1500 cuando x_1 se usa como variable auxiliar y x_2 para asignar probabilidades. $n' = 150$	118
3.6. Sesgo Relativo en porcentaje para la población Counties cuando x_1 se usa como variable auxiliar y x_2 para asignar probabilidades. Los valores SR para el estimador directo en (**) son mayores de 97.6 %, 74.6 % y 21.5 % para $\beta = 0,25, 0,5$ y 0.75 , respectivamente, y están omitidos. $n' = 150$	119
3.7. Eficiencia Relativa para el diseño muestral <i>SMS</i>	145
3.8. Eficiencia Relativa para el diseño muestral <i>MSS</i>	145
3.9. Eficiencia Relativa para el diseño muestral <i>MMM</i>	146
3.10. Sesgo Relativo para el diseño muestral <i>SMS</i>	147
3.11. Sesgo Relativo para el diseño muestral <i>MSS</i>	147
3.12. Sesgo Relativo para el diseño muestral <i>MMM</i>	148
3.13. Diagrama de caja con bigotes para los valores de los distintos estimadores. Se asume el diseño muestral <i>SMS</i> y tamaños muestrales $n' = 75$ y $n = 50$	149
3.14. Ratios Teóricos entre la varianza del estimador óptimo propuesto y la varianza del estimador estándar bajo la población Counties y el cuantil de orden $\beta = 0,5$	156
3.15. Ratios Teóricos entre la varianza del estimador óptimo propuesto y la varianza del estimador estándar bajo la población Turismos y el cuantil de orden $\beta = 0,5$	157
3.16. Eficiencia Relativa para los estimadores óptimo propuesto y estándar en la población Counties y para el cuantil de orden $\beta = 0,5$	159
3.17. Eficiencia Relativa para los estimadores óptimo propuesto y estándar en la población Turismos y para el cuantil de orden $\beta = 0,5$	160

3.18. Evolución de los valores W_{opt} usados por el estimador óptimo propuesto en la población Counties y para el cuantil de orden $\beta = 0,5$	161
3.19. Evolución de los valores W_{opt} usados por el estimador óptimo propuesto en la población Turismos y para el cuantil de orden $\beta = 0,5$	162
5.1. Relative Efficiency for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Pop06, Pop07, Pop08 and Pop09 populations and <i>SRSWOR</i> with $n = 200$ are considered.	200
5.2. Relative Efficiency for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Fam1500 population and <i>SRSWOR</i> with $n = 150$ are considered.	201
5.3. Relative Efficiency for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Hospitals population and <i>SRSWOR</i> with $n = 100$ are considered.	201
5.4. Relative Bias for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_w^{AC} (Standard), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Pop06, Pop07, Pop08 and Pop09 populations and <i>SRSWOR</i> with $n = 200$ are considered.	202
5.5. Relative Bias for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_w^{AC} (Standard), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Fam1500 population and <i>SRSWOR</i> with $n = 150$ are considered.	203
5.6. Relative Bias for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_w^{AC} (Standard), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Hospitals population and <i>SRSWOR</i> with $n = 100$ are considered.	203
5.7. Relative Efficiency for Fam1500 population and under <i>SRSWOR.M</i> sampling design. $n' = 150$	215
5.8. Relative Efficiency for Fam1500 population and under <i>SRSWOR.P</i> sampling design. $n' = 150$	216
5.9. Relative Efficiency for Counties population and under <i>SRSWOR.M</i> sampling design. $n' = 150$	217
5.10. Relative Efficiency for Counties population and under <i>SRSWOR.P</i> sampling design. $n' = 150$	218

5.11. Relative Bias in percent for Fam1500 population when x_1 is used as an auxiliary variable and x_2 is used to assign probabilities. $n' = 150$	221
5.12. Relative Bias in percent for Counties population when x_1 is used as an auxiliary variable and x_2 is used to assign probabilities. The RB 's values for the direct estimator in (**) are larger than 97.6 %, 74.6 % and 21.5 % for $\beta = 0,25, 0,5$ and 0.75, respectively, and are omitted. $n' = 150$	222
5.13. Relative Efficiency for Fam1500 and Counties populations and under ST.M sampling design. $n' = 150$	223
5.14. Relative Bias in percent for Fam1500 and Counties populations under ST.M sampling design and when the variable x_1 is used. $n' = 150$	224
A.1. Diagramas de dispersión de la población Fam1500	247
A.2. Diagramas de dispersión de las poblaciones Counties70 y Counties60.	249
A.3. Diagrama de dispersión de la población Hospitals.	250
A.4. Diagrama de dispersión de la población Murthy.	251
A.5. Diagramas de dispersión de la población Turismos.	253
A.6. Diagrama de dispersión de la población ECPF1997.	255
A.7. Diagramas de dispersión de las poblaciones Pop06, Pop07, Pop08 y Pop09	257
A.8. Diagramas de dispersión de la población Pob098	259
A.9. Diagramas de dispersión de la población Pob080	260

Índice de Tablas

3.1.	Esperanza empírica de $\lim_{t \rightarrow \infty} \widehat{F}_{st}^*(t)$ para varios diseños muestrales y considerando la variable x_1	127
3.2.	Esperanza empírica de $\lim_{t \rightarrow \infty} \widehat{F}_{st}^*(t)$ para varios diseños muestrales y considerando la variable x_2	127
3.3.	Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_1 . $\beta = 0,5$ y $n' = 150$	128
3.4.	Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_1 . $\beta = 0,5$ y $n' = 300$	129
3.5.	Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_2 . $\beta = 0,5$ y $n' = 150$	129
3.6.	Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_2 . $\beta = 0,5$ y $n' = 300$	129
3.7.	Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_1 . $\beta = 0,5$ y $n' = 150$	131
3.8.	Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_1 . $\beta = 0,5$ y $n' = 300$	131
3.9.	Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_2 . $\beta = 0,5$ y $n' = 150$	132
3.10.	Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_1 . $\beta = 0,5$ y $n' = 300$	132
3.11.	Combinaciones de diseños muestrales usados en muestreo con dos ocasiones sucesivas y probabilidades desiguales.	144

3.12. Medidas globales medias de precisión y eficiencia basadas en cuantiles de órdenes $\beta = 0,1, 0,3, 0,5, 0,7, 0,9$, y muestras de tamaño $n = 500$	176
3.13. Medidas de precisión y eficiencia para la línea de bajos ingresos cuando $\alpha = 0,6$, $\beta = 0,5$ y se toman muestras de tamaño $n = 500$.	177
3.14. Medidas de precisión y eficiencia para la razón de cuantiles cuando $\beta_1 = 0,5$, $\beta_2 = 0,25$, y se toman muestras de tamaño $n = 500$.	178
3.15. Medidas de precisión y eficiencia para la razón de cuantiles cuando $\beta_1 = 0,95$, $\beta_2 = 0,2$, y se toman muestras de tamaño $n = 500$.	178
5.1. Description and references of populations.	214
A.1. Análisis descriptivo para las variables de la población Fam1500 .	247
A.2. Análisis descriptivo para las variables de la población Counties60	248
A.3. Análisis descriptivo para las variables de la población Counties70	248
A.4. Análisis descriptivo para las variables de la población Hospitals	250
A.5. Análisis descriptivo para las variables de la población Murthy .	251
A.6. Análisis descriptivo para las variables de la población Turismos .	252
A.7. Análisis descriptivo para las variables de la población ECPF1997	254
A.8. Análisis descriptivo para las variables de la población Pop06 . .	256
A.9. Análisis descriptivo para las variables de la población Pop07 . .	256
A.10. Análisis descriptivo para las variables de la población Pop08 . .	256
A.11. Análisis descriptivo para las variables de la población Pop09 . .	257
A.12. Análisis descriptivo para las variables de la población Pob098 . .	258
A.13. Análisis descriptivo para las variables de la población Pob080 . .	258

Capítulo 1

Introducción

1.1. Problemas planteados

En el campo del muestreo en poblaciones finitas son numerosas las aportaciones que pueden hacerse a los métodos de estimación con información auxiliar de parámetros lineales y no lineales. Por ejemplo, en los últimos años han surgido nuevas metodologías para obtener estimadores más precisos usando información auxiliar. Estas nuevas metodologías son los estimadores de calibración (Deville y Särndal, 1992) y el método de verosimilitud empírica (Chen y Sitter, 1999). De estas metodologías, el método de verosimilitud empírica tiene un buen comportamiento asintótico y empírico, pero a causa de su reciente aparición, existen bastantes situaciones donde no ha sido analizado. En este trabajo se plantean diversos escenarios (presencia de datos faltantes, estimación de la función de distribución bajo un enfoque basado en el diseño muestral, etc) donde este método no había sido examinado, se estudian sus propiedades más importantes y se comprueba su eficiencia desde el punto de vista teórico y empírico.

Por otro lado, los métodos clásicos estudiados en muestreo de poblaciones finitas se han centrado en la estimación de parámetros lineales como la media o el total. En las últimas décadas se ha estado tratando el problema de la estimación de la función de distribución por diversos autores, pero este no es el caso de la estimación de los cuantiles, los cuales no han sido definidos ni analizados en algunas situaciones, como por ejemplo en los diseños muestrales más complejos, etc. De este modo, en este trabajo se pretende plantear y estudiar la estimación de los cuantiles en aquellas situaciones que aunque son más

complejas no son las menos utilizadas, puesto que son los diseños muestrales empleados por la mayoría de los organismos y agencias estadísticas, investigaciones sociales y económicas, etc. Además, los cuantiles son muy utilizados en estos organismos por la información que recogen y para obtener medidas de gran importancia para el interés de una nación, como por ejemplo la estimación de las líneas de pobreza, proporción de bajos ingresos, etc.

Existen determinados problemas para algunos de los estimadores de cuantiles que han sido propuestos en la literatura del muestreo. En primer lugar, varios de los estimadores de la función de distribución no cumplen las propiedades de una verdadera función de distribución, mientras que existen otros estimadores que dependen estrictamente de un modelo de superpoblación. En algunas ocasiones, puede ocurrir que no exista ningún modelo que se ajuste suficientemente bien a la población en estudio, por lo que una perspectiva basada en el diseño muestral resultaría más apropiada.

En resumen, los objetivos que se persiguen en este trabajo son: (i) analizar el método de verosimilitud empírica en campos no tratados (estimación de la función de distribución desde una perspectiva basada en el diseño muestral y usando una aproximación modelo-asistida, en presencia de datos faltantes, estimación de cuantiles, etc), (ii) estudiar el comportamiento de los cuantiles en diseños más complejos (muestreo en dos ocasiones con probabilidades desiguales o con múltiples variables auxiliares, muestreo bifásico, etc).

1.2. Antecedentes

El método de verosimilitud empírica para la estimación de parámetros se propone en Chen y Qin (1993) bajo muestreo aleatorio simple, mientras que Zhong y Rao (1996) lo aplican bajo muestreo estratificado aleatorio, pero es en 1999 cuando Chen y Sitter lo definen para cualquier diseño muestral y obtienen las propiedades más importantes. A raíz de este trabajo el método se generaliza y son numerosos los autores que investigan en esta línea de investigación. De este modo, Wu y Sitter (2001a) usan esta metodología asumiendo modelos de superpoblación y usando una nueva aproximación llamada modelo-calibrada. Chen y Wu (2002) obtienen estimadores para la función de distribución y cuantiles, pero sólo bajo la perspectiva de los modelos de superpoblación y usando la aproximación modelo-calibrada, Wu (2003, 2004b, 2005) trata otros aspectos de esta metodología, como propiedades asintóticas, optimalidad de estimadores, aspectos computacionales para el cálculo de los estimadores, etc.

La pérdida o falta de información es una propiedad común en las investigaciones por muestreo. Esta pérdida de información puede ocurrir por varias razones. Los individuos muestreados pueden negarse a participar en el estudio, los entrevistadores no puedan contactar con los individuos del estudio, pérdida accidental de información, etc. Tratar con datos faltantes en una investigación por muestreo no es un asunto relativamente sencillo. Existen una gran variedad de métodos en el caso de existir valores perdidos en los datos muestrales. La solución más simple es eliminar las unidades con falta de respuesta y aplicar un determinado método de muestreo a las unidades restantes. Este método, el cual Rubin (1987) llamó análisis de casos completos, puede producir sesgo en las estimaciones y varianzas muestrales más grandes (ver Rubin, 1987 o Little y Rubin, 1987).

La imputación es otra técnica que puede usarse en los individuos con falta de respuesta (Little y Rubin, 1987, Rao y Toutenburg, 1995, Särndal, 1992), aunque en algunas ocasiones conduce a inferencias no válidas en la etapa de estimación. Por ejemplo, la varianza puede resultar seriamente subestimada cuando la proporción de valores perdidos no es pequeña (Rao y Shao, 1992, Särndal, 1990, 1992).

Algunos autores han definido estimadores de tipo razón en presencia de datos faltantes. Estos estimadores solamente han sido definidos para una clase limitada de diseños muestrales. Por ejemplo, Tracy y Osahan (1994), Toutenburg y Srivastava (1998, 1999, 2000) desarrollaron estimadores de tipo razón para muestreo aleatorio simple sin reemplazamiento. Rueda y González (2004) plantean estimadores de tipo regresión en la presencia de datos faltantes.

En resumen, el tratamiento de datos faltantes en las encuestas por muestreo mediante el método de verosimilitud empírica es un campo que no ha sido investigado.

Como se ha comentado, el problema de la estimación de la función de distribución se ha discutido mediante el método de verosimilitud empírica, pero tan sólo desde la perspectiva modelo-calibrada (véase Wu y Sitter, 2001a). No se conocen antecedentes de la resolución de este problema desde un enfoque modelo-asistido.

La técnica usual en muestreo de poblaciones finitas para estimar cuantiles es la de obtener un estimador eficiente de la función de distribución e invertir este estimador para obtener una estimación del cuantil correspondiente. Sin ningún tipo de información auxiliar, los estimadores básicos para la estimación de la función de distribución son los de tipo Horvitz y Thompson (Horvitz y

Thompson, 1952) y de tipo Hájek (Rao, 1966, Basu, 1971, Särndal *et al.*, 1992). Nótese que los estimadores de tipo Hájek son verdaderas funciones de distribución y permiten una estimación más coherente de cuantiles.

Asumiendo información auxiliar en la etapa de estimación, son numerosos los estimadores propuestos para la función de distribución y/o cuantiles (Gross, 1980, Sedransk y Meyer, 1978, Chambers y Dunstan, 1986, Sedransk y Smith, 1988, Kuk y Mak, 1989, Francisco y Fuller, 1991, Dorfman y Hall, 1993, Mak y Kuk, 1993 Rao, 1994, etc.), aunque la mayoría de estos estimadores están diseñados para esquemas de muestreo más simples como muestreo aleatorio simple, muestreo estratificado, etc. Rao *et al.* (1990) proponen estimadores de tipo razón y diferencia para la estimación de la mediana poblacional usando una aproximación modelo-asistida basada en el diseño. Más recientemente, Rueda *et al.* (1998) y Rueda y Arcos (2001) proponen intervalos de confianza para cuantiles basados en estimadores de tipo razón y diferencia para la función de distribución. En Rueda *et al.* (2003, 2004) se utilizan estimadores de tipo diferencia para la estimación de cuantiles usando cuantiles poblacionales de la variable auxiliar.

La estimación de la mediana en muestreo de dos fases o bifásico fue desarrollada por Singh *et al.* (2001), Singh (2003) y Allen *et al.* (2002). El inconveniente de estos trabajos es que han sido desarrollados exclusivamente para muestreo aleatorio simple y la extensión para un determinado cuantil tampoco se investiga.

En muestreo de dos ocasiones sucesivas, la teoría desarrollada por Jessen (1942) y Patterson (1950) permite obtener estimadores de la media en la segunda ocasión combinando un estimador de tipo razón basado en las muestras solapadas y un estimador directo basado en la muestra no solapada de la ocasión más reciente. El muestreo en ocasiones sucesivas ha sido también discutido en Narain (1953), Adhvaryu (1978), Eckler (1955), Gordon (1983), Arnab y Okafor (1992), Singh y Srivastava (1973), Singh *et al.* (1992) y Singh (2003), el cual proporciona una extensa bibliografía de este tema. En todos los estudios anteriores, el parámetro de interés es la media poblacional. En ningún caso se trata el problema de la estimación de cuantiles.

Al igual que para el problema de la estimación de la función de distribución, la estimación de cuantiles usando el método de verosimilitud empírica se ha estudiado únicamente mediante la aproximación modelo-calibrada. Este enfoque presenta el inconveniente de proporcionar inferencias no válidas cuando no resulta posible adaptar un modelo de superpoblación a los datos del estudio. Bajo estas circunstancias, la aproximación modelo-asistida ofrecería resultados

más eficientes. Lamentablemente, en la literatura del muestreo en poblaciones finitas no se dispone de tal aproximación bajo el método de verosimilitud empírica.

Por otro lado, las numerosas medidas de pobreza (tal como líneas de pobreza, proporción de bajos ingresos, etc) que se manejan en los organismos estadísticos de numerosos países dependen de cuantiles. Los estimadores que se usan para el cálculos de estas medidas están basados en las técnicas tradicionales, y no se hace un uso efectivo de la información auxiliar. El estudio de la eficiencia de estimadores más complejos en la determinación de estas medidas de pobreza es una labor que tampoco ha sido extensamente analizada.

1.3. Objetivos científicos y aportes a la teoría del muestreo

A continuación se indica como se distribuye el presente texto y se comenta de forma breve los principales objetivos científicos y las aportaciones a la teoría del muestreo en poblaciones finitas.

En la siguiente sección se describe el marco de trabajo general seguido a lo largo del texto y se dan algunos conceptos básicos de la teoría del muestreo en población finitas. El objetivo es familiarizarse con la notación y conceptos que van a ser usados en todo el texto.

En la teoría de muestreo en poblaciones finitas el objetivo principal de cualquier metodología o de cualquier diseño muestral es el de mejorar las estimaciones de los parámetros en estudio en el sentido de construir nuevos estimadores que, para el mismo tamaño muestral, tengan menor error de estimación, lo que implica mayor precisión en las estimaciones de los parámetros, o equivalentemente, tengan el mismo error que los ya conocidos pero con un menor tamaño muestral, lo que produce una disminución en el coste real de la realización de la encuesta. Existen dos procedimientos para intentar mejorar las precisiones de las estimaciones. Por un lado, se pueden emplear nuevas técnicas de estimación y por otro, usar métodos de muestreo más complejos que utilicen más información (muestreo en ocasiones sucesivas, etc), o que la información auxiliar sea más fiable (muestreo bifásico), etc. La primera de estas técnicas se lleva a cabo en el Capítulo 2, en donde se pone a prueba el método de verosimilitud empírica como método de estimación, mientras que la segunda técnica se aplica en el Capítulo 3 para el problema de la estimación

de cuantiles, campo que no ha sido lo suficientemente tratado (véase la sección anterior) en la teoría del muestreo para los diseños muestrales considerados en dicho capítulo.

Como se ha comentado, el método de verosimilitud empírica se desarrolla en el Capítulo 2 para distintos escenarios. Esta reciente metodología obtiene estimadores tan eficientes (ver Chen y Sitter, 1999, Wu, 2002, Rueda *et al.*, 2006b, etc.) como los utilizados clásicamente en muestreo de poblaciones finitas, lo que lo convierte en una alternativa válida a usar en las encuestas por muestreo, puesto que si el escenario es el apropiado puede ayudar a obtener estimaciones más eficientes, reducir costes en las encuestas, etc. En la Sección 2.2 se recopilan los principales aspectos y resultados del método de verosimilitud empírica. Además, bajo esta metodología se plantean varias situaciones de un interés relevante en la teoría del muestreo, de los que destacan el problema de los datos faltantes y la estimación de la función de distribución y cuantiles.

Cuando se realiza un estudio mediante encuestas o cualquier otro procedimiento, es usual encontrarse ante la presencia de datos faltantes que vienen dados por parte del entrevistado o por cualquier otra circunstancia (pérdida casual de información, errores en la etapa de manipulación de datos, etc). Ante tal problema, una técnica frecuentemente utilizada es eliminar del estudio a aquellos individuos que presentan datos faltantes en alguna de sus variables. El inconveniente principal de esta técnica es el incremento del sesgo en las estimaciones. Otra técnica habitualmente utilizada es la imputación, que presenta el inconveniente de obtener en algunas ocasiones inferencias no válidas como consecuencia de considerar los valores imputados como si éstos fueran valores verdaderos.

En la Sección 2.3 se propone un camino alternativo para el tratamiento de los datos faltantes que no necesita eliminar del estudio a ningún individuo, aprovechando toda la información que se tiene en la muestra. Este procedimiento se desarrolla bajo el método de verosimilitud empírica. Se estudian las propiedades teóricas y mediante un estudio de simulación, se contrasta la precisión de los estimadores propuestos con otros estimadores conocidos y también diseñados para el tratamiento de datos faltantes. Véase también Rueda *et al.* (2006b).

El problema de la estimación de la función de distribución es un tema actual y muy importante del muestreo en poblaciones finitas, por tratarse de una función que permite determinar las características más importantes de la población en estudio, proporcionando información relevante acerca del comportamiento global de la población. Sin duda, los estimadores estudiados

clásicamente en la teoría del muestreo, como totales, medias, proporciones y varianzas, no ofrecen tanta información como la función de distribución. El problema de la estimación de cuantiles y de otros parámetros de tipo no funcional queda resuelto con el conocimiento de la función de distribución, puesto que éstos pueden obtenerse mediante inversión directa de la función de distribución. Además, permite obtener medidas importantes como la determinación de las líneas de pobreza, proporción de bajos ingresos, etc. y son muy útiles en investigaciones de tipo social o económico. Debido a la importancia de estos parámetros en algunas investigaciones o estudios, se debe de disponer de buenos métodos y técnicas para obtener las mejores estimaciones posibles.

Asumiendo una aproximación modelo-calibrada (véase Apéndice B), Chen y Wu (2002) propusieron estimadores de la función de distribución usando el método de verosimilitud empírica. Por otro lado, estos estimadores están basados en información auxiliar a través de un único punto del conjunto de valores para los que se define la función de distribución, presentando el problema de obtener estimaciones menos precisas cuando el argumento en el que se evalúa la función de distribución está bastante alejado del punto considerado para la variable auxiliar. Por tanto, estos estimadores presentan dos inconvenientes principalmente: (i) es necesario el conocimiento y el uso de un modelo de superpoblación para los datos muestrales del estudio y (ii) se hace un uso poco eficiente de la información auxiliar.

Asumiendo el método de verosimilitud empírica, en la Sección 2.4 se propone un estimador modelo-asistido para la función de distribución basado en un uso efectivo de la información auxiliar. Este estimador será más eficiente cuanto mayor sea la correlación entre las variables auxiliares y la variable principal. Además, no resulta necesario el conocimiento de un modelo de superpoblación, puesto que el estimador propuesto no es dependiente del modelo. El uso efectivo de la información auxiliar se justifica porque el estimador propuesto está basado en tres puntos perfectamente repartidos en el recorrido de valores en donde se define la función de distribución, de modo que, independientemente del valor donde se evalúe la función de distribución, éste valor estará cercano a alguno de los tres puntos, obteniendo estimaciones más precisas para la función de distribución. Esto permitirá también mejorar la calidad de la estimación de los cuantiles y de aquellos otros parámetros relacionados con éstos y que suelen obtenerse en las grandes instituciones estadísticas. Una propiedad deseable de un estimador de la función de distribución, es que éste sea por si mismo una verdadera función de distribución. Este es otro punto importante a la hora de obtener estimadores eficientes para los cuantiles poblacionales. Notamos que el estimador propuesto también posee esta propiedad.

En el Capítulo 3 se analiza el problema de la estimación de cuantiles en distintos esquemas de muestreo frecuentemente usados en la práctica, varios métodos de estimación y por último, usando el método de verosimilitud empírica.

La Sección 3.2 resuelve el problema de la estimación de cuantiles en muestreo bifásico cuando las muestras en cada una de las fases son seleccionadas mediante cualquier diseño muestral, con probabilidades iguales o desiguales. Se proponen varios estimadores de tipo directo, razón y exponencial que proporcionan estimaciones óptimas para un determinado cuantil. Se analizan importantes propiedades de estos estimadores tales como la insesgadez, estimación de varianzas, etc. Como caso particular, se investiga también el muestreo bifásico aplicado a la estratificación, diseño muestral que ofrece importantes ganancias en eficiencia debido a los beneficios que produce el muestreo estratificado. Todas estas propiedades se ven desde un punto teórico, aunque el análisis de los estimadores se completa con estudio empíricos llevados a cabo para los cuantiles y bajos distintos diseños muestrales con probabilidades desiguales. En términos de sesgo y de eficiencia relativa, estos estudios reflejan que los estimadores propuestos mejoran a otros estimadores diseñados en muestreo bifásico.

La mayoría de las investigaciones llevadas a cabo por los organismos nacionales de estadística son periódicas, es decir, se repiten a intervalos regulares de tiempo. Bajo este escenario, es aplicable la metodología propuesta en la Sección 3.3 para estimar cuantiles en muestreo en dos ocasiones, lo que puede permitir obtener una mejor precisión en la etapa de estimación como se ha comprobado desde el punto de vista teórico y práctico. Esta investigación se ha llevado a cabo, por un lado, bajo el uso de un diseño muestral arbitrario y por otro para el caso de múltiples variables auxiliares, siendo varios los objetivos científicos y aportes a la teoría del muestreo, puesto que los métodos tradicionales de estimación en muestreo de ocasiones sucesivas se han centrado en el problema de la estimación de parámetros lineales. Para el caso de la estimación de cuantiles, la situación es bastante diferente, y sólo recientemente este campo ha sido tratado por los estudios de investigación. En cualquier caso, los estudios existentes están basados únicamente en muestreo aleatorio simple y utilizan sólo la variable de interés en la fase de estimación, o bien sólo están diseñados para una única variable auxiliar.

En la Sección 3.4 se plantea el problema de la estimación de cuantiles mediante estimadores modelo-asistidos basados en el método de verosimilitud empírica. La aplicación de estos estimadores a la estimación de algunas medidas de pobreza también se discute dentro de esta sección. Debido a la

complejidad natural de los cuantiles y principalmente de las medidas de pobreza que se manejan, se propone usar la técnica bootstrap para el problema de la estimación de la varianzas de los estimadores. En los numerosos estudios empíricos llevados a cabo, puede observarse que tanto los estimadores propuestos como las estimaciones de las varianzas presentan un buen cumplimiento en términos de sesgo y eficiencia relativa.

Una valoración global de los resultados obtenidos así como las principales conclusiones de todos los estudios de este texto se resumen en el Capítulo 4. Además, se detallan las perspectivas y futuras líneas de investigación que se derivan a partir de las aportaciones y resultados concluyentes de los estudios llevados a cabo.

Según la legislación vigente para los estudios de Tercer Ciclo y del Título de Doctor, en el Capítulo 5 se hace una redacción de una parte de la tesis doctoral en la lengua inglesa. En concreto, se hace un resumen de todo el trabajo (Sección 5.1), se describen los estimadores propuestos bajo el método de verosimilitud empírica y en presencia de datos faltantes (Sección 5.2), se plantea la estimación de cuantiles en muestreo bifásico (Sección 5.3), y se comentan las principales conclusiones obtenidas a partir de la presente tesis doctoral (Sección 5.4). Notamos que las Secciones 5.2 y 5.3 se corresponden, respectivamente, con las Secciones 2.3 y 3.2. Estos trabajos también pueden consultarse en Rueda *et al.* (2006a, 2006b).

El texto se completa con una serie de apéndices de consulta sobre varios aspectos relacionados con los estudios llevados a cabo. Así, el Apéndice A recoge las principales propiedades y características de las poblaciones finitas que han sido usadas en los estudios de simulación. Además de un breve resumen estadístico de los datos de estas poblaciones, se muestran los diagramas de dispersión de tales poblaciones. El objetivo es que de forma personalizada se puedan interpretar los distintos estudios de simulación y poder obtener, por tanto, conclusiones subjetivas. Para ello, resulta apropiado un mejor conocimiento de los datos que se utilizan en los estudios de simulación.

Son numerosos los métodos de extracción de unidades y diseños muestrales que se utilizan en este texto. En el Apéndice B se recogen todos ellos con el fin de poder obtener una rápida consulta sobre la notación o sobre el propio método de muestreo o diseño muestral.

Por último, notar que todos los estudios de simulación se han llevado a cabo mediante el lenguaje de programación *R*. Todos los procedimientos y funciones para obtener en *R* tanto los estimadores propuestos en este texto como el resto

de estimadores para cada diseño muestral están disponibles en el Apéndice C.

Son numerosas las razones por las que se ha usado este software. En primer lugar, es un lenguaje intuitivo con una gran cantidad de argumentos estadísticos que facilitan la implementación de los estimadores propuestos. Otros programas como *Mathematica*, *Matlab*, *C++*, etc., carecen de tales procedimientos estadísticos. Por otro lado, es un paquete que destaca por su rapidez y que permite obtener el mayor número de simulaciones en menor tiempo. *R* es un lenguaje de programación gratuito y disponible a cualquier usuario, al contrario de otros específicos de estadística como *SAS*, que debido a sus altas licencias está únicamente disponible, en la mayoría de los casos, a las grandes empresas. El dispositivo gráfico que dispone *R* y su compatibilidad con *S-PLUS* son otros argumentos hacen que la mayoría de los investigadores en el campo del muestreo en poblaciones finitas prefieran el uso de este software. Sirva de ejemplo los artículos publicados en este sentido (por ejemplo Wu, 2005) así como las conferencias internacionales sobre el programa *R* que también se están abriendo paso, como la segunda conferencia internacional de usuarios de *R* que se celebrará del 15 al 17 de junio de 2006 en Viena, Austria. De hecho, el gran auge que está teniendo este software hace que se estén introduciendo día a día nuevos procedimientos y paquetes estadísticos.

1.4. Notación y conceptos básicos

En esta sección se describe el marco de trabajo usual en el ámbito del muestreo de poblaciones finitas. Además, se introducen algunos conceptos básicos y la notación común que se sigue a lo largo del texto.

Se denomina *población* a un conjunto de unidades del que se desea obtener cierta información. Esta población se denota como U , es finita y contiene N elementos distintos e identificados, es decir, $U = \{1, \dots, i, \dots, N\}$

En la población U es posible medir o contar en cada unidad una o varias *características* o *variables*, o clasificar sus unidades de acuerdo a ellas. A partir de estos resultados se puede llegar al conocimiento de valores como la media, el total, la proporción, función de distribución, etc., a los que se denomina *parámetros poblacionales*. La media, el total, etc., son parámetros lineales, mientras que la función de distribución, cuantiles, etc., son parámetros no lineales.

Existen dos estrategias posibles para la recopilación de datos: (i) examinar

todas las unidades de la población, es decir, realizar un censo, y (ii) examinar, según unos planes establecidos con anterioridad, unas pocas unidades de la población que son representativas, es decir, obtener una muestra, y suponer que de los resultados obtenidos se infieren las características de toda la población.

En la práctica, determinados parámetros poblacionales son desconocidos y no pueden calcularse mediante un censo. Por esta razón, se recurre a una muestra para estimar estos parámetros poblacionales. Así, una muestra es un subconjunto de unidades, s , de U seleccionados de acuerdo con un diseño de muestreo específico, d , que asigna una probabilidad conocida, $p(s)$, tal que $p(s) > 0$ para todo $s \in S$, donde S es el conjunto de las posibles muestras s y $\sum_{s \in S} p(s) = 1$. El valor de la media, total, proporción o función de distribución obtenido a partir de la muestra se denomina *estimador* del correspondiente parámetro poblacional.

Dentro de esta población interesa estudiar ciertas características de una variable de *estudio*, *interés* o *principal* denominada y . Las variables *auxiliares* son aquellas, que sin ser objeto de estudio, son usadas para varios fines, como por ejemplo, para la selección de unidades en la muestra, mejorar las estimaciones, etc. Asociado al elemento i de la muestra se conoce exactamente y sin error el valor de la característica de interés, esta cantidad se denotará como y_i . Para P variables auxiliares, el vector de variables auxiliares viene dado por $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_p, \dots, \mathbf{x}_P)$, donde $\mathbf{x}_p = (x_{1p}, \dots, x_{ip}, \dots, x_{Np})^t$. Se asume que estas variables auxiliares también son conocidas para aquellos individuos seleccionados en la muestra. En algunas ocasiones, se supone que los totales o medias poblacionales de las variables auxiliares son conocidos, es decir, las cantidades $\mathbf{X} = (X_1, \dots, X_P)$ o $\bar{\mathbf{X}} = (\bar{X}_1, \dots, \bar{X}_P)$ son conocidas, donde $X_p = \sum_{i=1}^N x_{ip}$ y $\bar{X}_p = N^{-1} \sum_{i=1}^N x_{ip}$.

La probabilidad de inclusión de primer orden asociadas al plan de muestreo d para un individuo i , π_i , indica la probabilidad que tiene este individuo de pertenecer a la muestra s . Asimismo, π_{ij} indica la probabilidad de que ambas unidades i y j pertenezcan a la muestra s . A esta cantidad se le llama probabilidad de inclusión de segundo orden. Otras cantidades que serán usadas son los pesos básicos del diseño $d_i = \pi_i^{-1}$, $\Delta_{ij} = \pi_{ij} - \pi_i \pi_j$, etc.

De este modo, los principales parámetros poblacionales desconocidos en la práctica y que habrá que estimar son la media poblacional de la variable de interés,

$$\bar{Y} = \frac{1}{N} \sum_{i=1}^N y_i,$$

el total poblacional,

$$Y = \sum_{i=1}^N y_i,$$

la función de distribución,

$$F_y(t) = \frac{1}{N} \sum_{i=1}^N \delta(t - y_i)$$

y el cuantil para un orden β ($0 < \beta < 1$),

$$Q_y(\beta) = F_y^{-1}(\beta) = \inf\{t \mid F_y(t) \geq \beta\},$$

donde $\delta(\cdot)$ es la función indicadora que toma el valor $\delta(a) = 1$ si $a \geq 0$ y $\delta(a) = 0$ en otro caso y $F_y^{-1}(t)$ denota la función inversa de $F_y(t)$.

Sin ningún tipo de información auxiliar, la media poblacional de la variable de interés, \bar{Y} , suele estimarse mediante el estimador de tipo Hortviz-Thompson

$$\widehat{Y}_{HT} = \frac{1}{N} \sum_{i \in s} d_i y_i. \quad (1.1)$$

Para el caso de la estimación de la función de distribución, este estimador viene dado por

$$\widehat{F}_{HTy}(t) = \frac{1}{N} \sum_{i \in s} d_i \delta(t - y_i), \quad (1.2)$$

aunque suele usarse el estimador de tipo Hájek que es una verdadera función de distribución. Este estimador viene dado por

$$\widehat{F}_{HKy}(t) = \sum_{i \in s} d_i^* \delta(t - y_i), \quad (1.3)$$

donde $d_i^* = d_i / \sum_{j \in s} d_j$. El cuantil de orden β puede estimarse directamente mediante la inversión de este último estimador, esto es,

$$\widehat{Q}_{HKy}(\beta) = \widehat{F}_{HKy}^{-1}(\beta) = \inf\{t \mid \widehat{F}_{HKy}(t) \geq \beta\}. \quad (1.4)$$

Capítulo 2

El método de verosimilitud empírica

El método de verosimilitud empírica para la estimación de parámetros fue propuesto en Chen y Qin (1993), aunque fueron Chen y Sitter (1999) quienes establecieron las bases teóricas principales de este método, y partir de las cuales se han basado todos los estudios posteriores. En este capítulo se investiga esta técnica reciente en diferentes campos del muestreo en poblaciones finitas.

En la Sección 2.2 se recogen los principales aspectos de esta metodología para el caso de la estimación de la media poblacional, pueden verse las propiedades asintóticas más importantes y los diferentes tipos de estimadores basados en cada una de las perspectivas de estimación.

En cualquier estudio es usual encontrarse con el problema de datos faltantes. En la Sección 2.3 se propone usar un estimador basado en el método de verosimilitud empírica como solución al problema de la existencia de datos faltantes (véase también Rueda *et al.*, 2006b).

La estimación de la función de distribución mediante el método de verosimilitud empírica se estudia en la Sección 2.4. Se propone usar la aproximación modelo-asistida para obtener tal estimador, y se hace un uso eficiente de la información auxiliar al estar basado el estimador en varias variables auxiliares y en varios puntos de estimación.

2.1. Introducción

En la teoría del muestreo en poblaciones finitas, el objetivo principal de un método determinado para obtención de estimadores o de cualquier diseño muestral es el de mejorar las estimaciones de los parámetros en estudio en el sentido de construir nuevos estimadores que, para el mismo tamaño muestral, tengan menor error de estimación, lo que implica mayor precisión en las estimaciones de los parámetros, o equivalentemente, tengan el mismo error que los ya conocidos pero con un menor tamaño muestral, lo que produce una disminución en el coste real de la realización de la encuesta.

Por estas razones fundamentalmente, la metodología del muestreo en poblaciones finitas precisa de nuevas aportaciones que abaraten los costes de los estudios o investigaciones estadísticas, se mejoren las estimaciones desde el punto de vista de la eficiencia o sesgidez y se dispongan, en general, de mejores propiedades.

Es conocido que según la información que se utilice en la etapa de estimación de parámetros, se tienen dos caminos para intentar mejorar la precisión de las estimaciones: por un lado, utilizar diseños muestrales más complejos (muestreos estratificados, por conglomerados, polietápicos, adaptativos, etc.) basados únicamente en los datos de la característica de interés, y por otro lado, emplear las metodologías propias de la teoría del muestreo en poblaciones finitas basadas en el uso de información auxiliar. Esta información auxiliar, dada a través de un vector de variables auxiliares, debe estar altamente correlacionada con la característica de interés para poder obtener mayor precisión en la etapa de estimación. Estas dos alternativas se pueden combinar para perseguir el objetivo de obtener mejores estimaciones, es decir, usar diseños muestrales más complejos en métodos de estimación de parámetros que utilicen información auxiliar es una opción muy atractiva en la materia que nos ocupa.

El método de verosimilitud empírica, que se desarrolla a largo de este capítulo, permite combinar las dos ideas anteriores y es bastante eficiente como se ha comprobado tanto desde el punto de vista teórico como empírico (véase por ejemplo, Chen y Qin, 1993, Chen y Sitter, 1999, Chen y Wu, 2002, Wu, 2003, Rueda y Muñoz, 2005, 2006a, 2006d, etc.).

Los primeros métodos que incorporan información auxiliar en la fase de estimación son los llamados métodos indirectos de estimación, entre los que destacan los conocidos métodos de razón, diferencia y regresión. Estos estimadores no siempre garantizan que se produzca una disminución del error de muestreo respecto a los estimadores que no usan información auxiliar. Esta

ganancia en precisión depende en mayor medida de la relación entre las variables auxiliares y la variable objeto de estudio, del buen uso de las hipótesis que se supongan para emplear un procedimiento u otro, y de que dichas hipótesis se ajusten en mayor o menor medida al problema real.

Los estimadores anteriores se basan únicamente en los datos muestrales, es decir, utilizan un enfoque basado en el diseño muestral. Recientemente, en muestreo se está utilizando la perspectiva basada en modelos (ver p.e. Pérez, 2002 y Sánchez-Crespo, 2002) y la nueva aproximación modelo-calibrada (Wu y Sitter, 2001a). Estas aproximaciones se basan en modelos de superpoblación y son dependientes de dichos modelos. El objetivo de estos métodos es obtener estimaciones más precisas, resultados más concluyentes en la comparación de estrategias, producir estrategias óptimas, obtener propiedades asintóticas más atractivas, etc., pero cuando el esquema de trabajo está perfectamente identificado con un modelo de superpoblación. Bajo esta perspectiva cobra especial importancia el uso de variables auxiliares cuyos valores tienen que ser conocidos para todos los individuos de la población. Por tanto, para poder usar este enfoque se debe conocer el adecuado modelo de superpoblación asociado a los datos de la población en estudio. En resumen, estas aproximaciones son más eficientes que el enfoque basado en el diseño muestral cuando el modelo de superpoblación se ajusta bien, y pueden llegar a obtener propiedades no deseables, como inferencias no válidas, cuando se usa un modelo de superpoblación erróneo. En consecuencia, para llegar a cabo estas aproximaciones, sería conveniente obtener más información: el modelo de superpoblación apropiado y todos los valores de las variables auxiliares para todos los individuos de la población. Cuando no se dan estas circunstancias, puede resultar más apropiado un método de estimación basado en el diseño muestral.

Una alternativa intermedia entre los métodos anteriores y la clásica estimación basada en diseños, es la aproximación modelo-asistida. Ésta consiste en usar un modelo de superpoblación para obtener una estimación de un determinado parámetro poblacional, y entonces, usar éste último en la etapa de estimación. Sin pérdida de eficiencia, la ventaja de este estimador es que sus estimaciones no son dependientes del modelo de superpoblación, permitiendo obtener inferencias válidas independientemente de si el modelo resulta ser apropiado o no para los datos de la población de estudio. El conocido estimador de regresión generalizado (Cassel *et al.*, 1976, Särndal, 1980), los estimadores de calibración (Deville y Särndal, 1992) y el propio estimador de verosimilitud empírica (Chen y Qin, 1993, Chen y Sitter, 1999) pueden ser categorizados como aproximaciones modelo-asistidas

Son dos los métodos para obtener estimadores que han aparecido recién-

temente: los estimadores de calibración y los de verosimilitud empírica. Los primeros fueron propuestos por Deville y Särndal (1992), y desde entonces se han comprobado sus propiedades teóricas, se han obtenido numerosas modificaciones, y se ha extendido el método a diversos esquemas de muestreo, siendo todos los resultados obtenidos bastante satisfactorios.

El método de verosimilitud empírica para la estimación de parámetros es más novedoso que el método de calibración. Fue propuesto en Chen y Qin (1993) para muestreo aleatorio simple, aunque el auge y el interés de esta metodología se produce en 1999 cuando Chen y Sitter plantean el método para cualquier diseño muestral. Al igual que el método de calibración, este método permite incorporar información auxiliar de una o varias variables adicionales, y se puede plantear tanto desde una perspectiva modelo-asistida, como desde la reciente aproximación modelo-calibrada (Wu y Sitter, 2001a).

Los estimadores de verosimilitud empírica para la media poblacional basados en el diseño muestral y bajo la aproximación modelo-calibrada, serán vistos en la Sección 2.2. Las principales propiedades asintóticas de estos estimadores podrán también consultarse en esta sección. Nótese que el método de verosimilitud empírica usa la aproximación modelo-asistida para determinar un determinado parámetro o variable, y posteriormente se basa en el diseño muestral para determinar los estimadores. Por simplicidad y sin pérdida de generalidad, en este caso nos referiremos como aproximación modelo-asistida o aproximación basada en el diseño muestral.

Todos los métodos generales de estimación de parámetros asumen que no existen datos faltantes en la muestra. Cuando existen observaciones perdidas en la muestra, la solución más simple es eliminar aquellos individuos con observaciones incompletas y restringir el estudio a los individuos que presentan observaciones completas para todas las variables. De este modo, con este conjunto de observaciones se puede aplicar cualquier técnica de estimación de parámetros. Una consecuencia de este método es la reducción de individuos en la muestra respecto a la muestra planificada, lo que produce mayores sesgos en las estimaciones y mayor varianza muestral. Usando el método de verosimilitud empírica, en la Sección 2.3 se proponen estimadores para el tratamiento de datos faltantes con buenas propiedades asintóticas y empíricas y que aprovechan todas las observaciones muestrales, estén éstas completas o incompletas para todas las variables del estudio.

Otro tema de actualidad en muestreo es el problema de la estimación de la función de distribución. Los estudios se han centrado clásicamente en la estimación de parámetros poblacionales de tipo puntual, como totales, medias,

proporciones y varianzas. La estimación de la función de distribución es un campo muy importante al tratarse de una función que permite determinar las características más importantes de la población en estudio, proporcionando información relevante acerca del comportamiento global de la población. Obtener buenos estimadores para tal función no es tan simple como en el caso de los estimadores puntuales. Para este problema, un buen estimador, $\widehat{F}(t)$, ha de cumplir las propiedades básicas de una verdadera función de distribución:

1. $\lim_{t \rightarrow -\infty} \widehat{F}(t) = 0$; $\lim_{t \rightarrow +\infty} \widehat{F}(t) = 1$.
2. $\widehat{F}(t)$ es no decreciente, es decir, $\forall t_1 < t_2$ se verifica $\widehat{F}(t_1) \leq \widehat{F}(t_2)$.
3. Dado $t > t^*$, $\lim_{t \rightarrow t^*} \widehat{F}(t) = \widehat{F}(t^*)$.

Varios de los estimadores propuestos en la literatura del muestreo en poblaciones finitas no satisfacen todas estas propiedades y no son, por tanto, funciones de distribución. Por ejemplo, la función de distribución estimada mediante el método de calibración no cumple los requisitos necesarios para ser una verdadera función de distribución.

En la Sección 2.4 se propone un estimador modelo-asistido para la función de distribución basado en el diseño muestral que cumple estas propiedades y goza de una excelente ganancia en eficiencia como consecuencia de un uso efectivo de la información auxiliar. Éstas son dos ventajas importantes de este estimador propuesto basado en el método de verosimilitud empírica. En esta sección, también pueden consultarse los principales estimadores de verosimilitud pseudo empírica modelo-calibrados para la función de distribución.

En resumen, este capítulo ofrece una descripción detallada del método de verosimilitud empírica en la estimación de la media o total de la población. El objetivo de este análisis es mostrar de forma sencilla cómo se construye este estimador en distintos diseños muestrales y para los distintos enfoques existentes en muestreo, cuáles son sus propiedades más importantes y la relación que tiene con otros estimadores más conocidos. Usando este esquema teórico, se aportan nuevas soluciones al problema de los datos faltantes y a la estimación de la función de distribución.

2.2. Estimación de la media poblacional

2.2.1. Estimadores basados en el diseño muestral

La metodología de verosimilitud empírica fue usada por Owen (1988, 1990) como un método para la construcción de regiones de confianza con observaciones independientes. Owen afirmó que el estadístico de verosimilitud empírica tiene una distribución asintótica χ^2 , y por tanto se puede usar para la estimación de intervalos de confianza y contraste de hipótesis. Qin y Lawless (1994, 1995) usan el método de verosimilitud empírica para la estimación puntual cuando la información se incorpora a través de la maximización de la función de verosimilitud empírica. A raíz de aquí, este método se popularizó y una gran gama de desarrollos sobre verosimilitud empírica han sido descritos en el reciente libro de Owen (2001) para distintos ámbitos.

Históricamente el uso de verosimilitud empírica fue propuesto por Hartley y Rao (1968), pero la primera aplicación formal en muestreo para poblaciones finitas del método de verosimilitud empírica se debe a Chen y Qin (1993), que lo estudiaron bajo muestreo aleatorio simple.

A continuación se detalla de forma breve la idea principal del método de verosimilitud empírica para el problema de la estimación de la media muestral de y , $\bar{Y} = N^{-1} \sum_{i=1}^N y_i$, y para muestreo aleatorio simple. En este caso, el estimador usual es el estimador de tipo Horvitz-Thompson, dado por

$$\bar{y} = \frac{1}{n} \sum_{i \in s} y_i = \sum_{i \in s} \frac{1}{n} y_i. \quad (2.1)$$

En la expresión (2.1) se observa que el estimador usa n puntos y_i de la muestra con el mismo peso ($1/n$) para estimar el parámetro. Puede ocurrir que ciertas observaciones y_i sean más determinantes que otras para el cálculo del parámetro. Bajo estas circunstancias es conveniente darle a las observaciones más determinantes un mayor peso que aquellas que son menos influyentes para estimar el valor del parámetro. Esta es la idea de los estimadores de verosimilitud empírica, es decir, pretenden cambiar los pesos $1/n$ por otros pesos \hat{p}_i , $i = \{1, \dots, n\}$, con el objetivo de mejorar la estimación del parámetro. Las variables auxiliares juegan un papel importante en este método, puesto que son usadas para obtener los nuevos pesos.

Sea p_i la masa de probabilidad de y_i , con $i \in s$. El estimador máximo

verosímil empírico de \bar{Y} se define como

$$\hat{y}_{PE} = \sum_{i \in s} \hat{p}_i y_i,$$

donde \hat{p}_i , $i = \{1, \dots, n\}$, maximiza la función de verosimilitud empírica, $L(\mathbf{p}) = \prod_{i \in s} p_i$ sujeta a las restricciones

$$\sum_{i \in s} p_i = 1 \quad (p_i > 0), \quad (2.2)$$

$$\sum_{i \in s} p_i \mathbf{x}_i = \bar{\mathbf{X}}. \quad (2.3)$$

La información auxiliar se incorpora en la segunda restricción. Esta expresión se justifica al asumir que los pesos que dan una estimación perfecta para $\bar{\mathbf{X}}$, deberían de dar una buena precisión en la estimación de \bar{Y} . Resulta razonable asumir que las estimaciones serán más eficientes a medida que y y \mathbf{x} presenten una relación lineal más fuerte.

Este problema de maximización con restricciones puede resolverse mediante el método de los multiplicadores de Lagrange.

Los estimadores de verosimilitud empírica se pueden diseñar desde distintas perspectivas, siendo el investigador quien debe decidir el modo de aplicar el método de verosimilitud empírica. Algunos de los distintos enfoques a través de los cuales se puede diseñar esta metodología son los siguientes:

(E1). Sustitución de $L(\mathbf{p})$.

En Chen y Qin (1993) se usa la función $L(\mathbf{p})$ para obtener los estimadores de verosimilitud empírica, mientras que Chen y Sitter (1999) usaron el logaritmo de esta función a nivel poblacional, esto es, propusieron usar

$$l(\mathbf{p}) = \log \prod_{i=1}^N p_i = \sum_{i=1}^N \log(p_i).$$

Notamos que el hecho de utilizar logaritmos no produce ningún cambio en las estimaciones al tratarse la función logaritmo de una función estrictamente creciente que conserva los puntos extremos de la función original. La ventaja es una mayor facilidad para obtener estimaciones. El problema que se plantea es cómo estimar $l(\mathbf{p})$ a través de una función eficiente $\hat{l}(\mathbf{p})$. Tomando $\log(p_i)$ como una variable de la que se pretende

estimar su total, este planteamiento presenta fácil solución. Como se detalla en Chen y Sitter (1999) y para un determinado diseño general, $l(\mathbf{p})$ se puede estimar a través de la denominada log-función de verosimilitud pseudo empírica, dada por:

$$\widehat{l}(\mathbf{p}) = \sum_{i \in s} d_i \log(p_i),$$

donde d_i son pesos básicos que hacen que $\widehat{l}(\mathbf{p})$ sea insesgada bajo el diseño para $l(\mathbf{p})$, es decir

$$E \left[\widehat{l}(\mathbf{p}) \right] = E \left[\sum_{i \in s} d_i \log(p_i) \right] = \sum_{i=1}^N \log(p_i) = l(\mathbf{p}).$$

Este cambio en la función de verosimilitud empírica hace que esta técnica sea aplicable bajo un diseño muestral general, a diferencia del método original propuesto por Chen y Qin (1993) que está diseñado exclusivamente para muestreo aleatorio simple. Bajo este método de muestreo, ambas perspectivas del método de verosimilitud empírica producen las mismas estimaciones.

(E2). Sustitución de la restricción $\sum_{i \in s} p_i \mathbf{x}_i = \overline{\mathbf{X}}$.

Al imponer que $\sum_{i \in s} p_i \mathbf{x}_i = \overline{\mathbf{X}}$, se están considerando valores para p_i que proporcionan estimaciones perfectas para $\overline{\mathbf{X}}$, y podemos plantearnos cómo de efectivo es el uso que se está haciendo de la información adicional a través de la condición anterior. Por este motivo, si la información auxiliar $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_P)$ es conocida, una cuestión a preguntarse sería: ¿Cual es la mejor expresión a usar en la restricción (2.3) para hallar el estimador de verosimilitud empírica? . Para resolver esta pregunta se ha definido la cantidad $\mathbf{u}_i = u(y_i, \mathbf{x}_i)$, con $i = \{1, \dots, N\}$, siendo $u(\cdot)$ una función conocida de y_i y de \mathbf{x}_i y que verifica

$$\frac{1}{N} \sum_{i=1}^N \mathbf{u}_i = \mathbf{0}.$$

De este modo, \mathbf{u}_i es una variable de calibración que reemplaza la expresión (2.3) por

$$\sum_{i \in s} p_i \mathbf{u}_i = \frac{1}{N} \sum_{i=1}^N \mathbf{u}_i = \mathbf{0}, \quad (2.4)$$

donde $\mathbf{u}_i = \mathbf{x}_i - \overline{\mathbf{X}}$. La cuestión que surge ahora es cómo escoger $u(\cdot)$ para obtener estimadores más eficientes. En resumen, este método dispone

de numerosas alternativas o soluciones dependiendo de la función $u(\cdot)$ escogida. Una elección apropiada de esta función supondrá más exactitud en las estimaciones. El uso de la aproximación modelo-calibrada es una solución óptima a este problema cuando no pueda asumirse una relación lineal entre y y \mathbf{x} .

(E3). Utilización de la aproximación modelo-calibrada.

En (E2) se usa una aproximación modelo-asistida, esto es, se asume una relación lineal (aunque pueden establecerse relaciones de otro tipo) para determinar unos valores \mathbf{u}_i apropiados, y posteriormente, se realizan estimaciones basadas en el diseño. Si la relación entre la variable de interés y y el vector de variables auxiliares $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_P)$ puede ser descrita a través de un modelo de superpoblación con una buena bondad de ajuste, puede resultar útil el uso de estimadores modelo-calibrados (Chen y Sitter, 2001) frente a los estimadores basados en el diseño. Esta aproximación consiste en asumir un determinado modelo de superpoblación, obtener los valores estimados para la variable y mediante este modelo, y a continuación usarlos en la etapa de estimación.

En este sentido, se han propuesto varios modelos que dan lugar a los estimadores óptimos modelo-calibrados. Éstos usan el criterio de mínima esperanza bajo el modelo de superpoblación de la varianza basada en el diseño para obtener la solución óptima (véase por ejemplo los trabajos de Godambe, 1955, Godambe y Thompson, 1973 y Cassel *et al.*, 1976). Los estimadores modelo-calibrados se desarrollan con detalle en la Sección 2.2.3.

La perspectiva de Chen y Sitter (véase (E1)) es más apropiada como se ha comprobado en las investigaciones posteriores. Además, puede ser aplicada a cualquier diseño muestral, no estando limitada exclusivamente al muestreo aleatorio simple. De este modo, los primeros pasos antes de aplicar el método de verosimilitud empírica son:

1. Enfocar el problema bajo un modelo de población fija, es decir, basado en el diseño muestral y aplicando la aproximación modelo-asistida, o bien, asumir un modelo de superpoblación para poder aplicar el enfoque modelo-calibrado.
2. Determinar la función $u(\cdot)$ utilizada en la restricción (2.4). Para el enfoque basado en el diseño muestral se suele usar $\mathbf{u}_i = \mathbf{x}_i - \bar{\mathbf{X}}$, mientras que bajo el enfoque modelo-calibrado, la función $u(\cdot)$ es única y fácilmente deducible a partir del modelo de superpoblación.

Estimadores bajo muestreo aleatorio simple

Una vez tenidas en cuenta estas consideraciones previas, empezaremos analizando el método de verosimilitud empírica según Chen y Qin (1993), el cual está diseñado para muestreo aleatorio simple.

Este estimador fue la primera aplicación formal del método de verosimilitud empírica en poblaciones finitas para la estimación de parámetros lineales y usando información auxiliar. Este planteamiento no se puede extender a diseños muestrales más complejos.

Según Chen y Qin (1993), el uso de verosimilitud empírica en el contexto de poblaciones finitas se puede plantear de dos formas diferentes:

1. Si todos los valores de y_i están disponibles para la población en estudio, la función de verosimilitud se define como $L^*(\mathbf{p}) = \prod_{i=1}^N p_i$, donde p_i es la densidad de la observación y_i . En la práctica esta situación no se va a presentar y lo más usual es que y_i sea conocida para los individuos de la muestra s . En tal caso la función de verosimilitud empírica para cualquier muestra s , con $s \subseteq S$, se define como $L(\mathbf{p}) = \prod_{i \in s} p_i$, donde se requiere que $\sum_{i=1}^n p_i \leq 1$. Este planteamiento fue propuesto por Jagers (1986) y es el que se sigue en varios estudios de estimación de parámetros en muestreo de poblaciones finitas mediante verosimilitud empírica (Chen y Qin, 1993, Zhong y Rao, 1996, etc).
2. Según el esquema de muestreo propuesto por Hartley y Rao (1968), los cuales consideraban que la variable de interés sólo puede tomar un número finito de valores, es decir, y_i , con $i = \{1, \dots, I\}$. Bajo esta situación, la población media se define como:

$$\bar{Y} = \sum_{i=1}^I \frac{N_i}{N} y_i,$$

donde N_i es el número de unidades en la población con característica y_i . Bajo muestreo aleatorio simple de tamaño n , la verosimilitud basada en el diseño está dada por una distribución hipergeométrica multidimensional:

$$L(N_1, \dots, N_I) = \prod_{i=1}^I \frac{\binom{N_i}{n_i}}{\binom{N}{n}},$$

donde n_i es el número de unidades en la muestra con la característica y_i . Cuando $N \rightarrow \infty$, $N_i/N \rightarrow p_i$, y $n/N \rightarrow 0$, la verosimilitud se puede aproximar por una función de verosimilitud de una distribución multinomial, a saber:

$$\frac{n!}{\prod_{i=1}^I n_i!} \prod_{i=1}^I p_i^{n_i}.$$

Utilizando el primer planteamiento propuesto por Jagers (1986), al maximizar $L(\mathbf{p})$ sin usar información auxiliar, resulta $\hat{p}_i = 1/n$ para cada $i \in s$, y el estimador de verosimilitud empírica está dado por

$$\hat{y}_{EL} = \sum_{i \in s} \hat{p}_i y_i = \frac{1}{n} \sum_{i \in s} y_i = \bar{y},$$

que coincide con el estimador directo usual para la media poblacional.

Cuando se dispone de alguna información auxiliar, ésta puede usarse en la etapa de maximización de la función de verosimilitud para obtener nuevos pesos p_i que produzcan estimaciones más eficientes para la media. Se asume que la información auxiliar disponible para la población verifica

$$\frac{1}{N} \sum_{i=1}^N \mathbf{u}_i = \mathbf{0},$$

donde $\mathbf{u}_i = u(y_i, \mathbf{x}_i)$ es una función conocida de y_i y de \mathbf{x}_i de vectores valuados. De este modo, el nuevo problema consiste en maximizar $L(\mathbf{p})$ sujeto a las restricciones:

$$\sum_{i \in s} p_i = 1 \quad (p_i \geq 0), \quad (2.5)$$

$$\sum_{i \in s} p_i \mathbf{u}_i = \mathbf{0}. \quad (2.6)$$

Usando el método de los multiplicadores de Lagrange, los valores esperados para p_i , con $i \in s$, están dados por:

$$\hat{p}_i^* = \frac{1}{n(1 + \lambda^t \mathbf{u}_i)}, \quad (2.7)$$

donde λ es la solución de la ecuación

$$\sum_{i \in s} \frac{\mathbf{u}_i}{1 + \lambda^t \mathbf{u}_i} = \mathbf{0}. \quad (2.8)$$

El estimador de verosimilitud empírica para la media poblacional bajo muestreo aleatorio simple y usando la metodología de Chen y Qin (1993) está dado por

$$\widehat{y}_{EL} = \sum_{i \in s} \widehat{p}_i^* y_i. \quad (2.9)$$

Asumiendo que la relación entre y y el vector \mathbf{x} es lineal, la función de calibración usual viene dada por $\mathbf{u}_i = \mathbf{x}_i - \overline{\mathbf{X}}$, en cuyo caso la restricción (2.6) resulta ser

$$\begin{aligned} \sum_{i \in s} p_i \mathbf{u}_i &= \sum_{i \in s} p_i (\mathbf{x}_i - \overline{\mathbf{X}}) = \sum_{i \in s} p_i \mathbf{x}_i - \sum_{i \in s} p_i \overline{\mathbf{X}} = \sum_{i \in s} p_i \mathbf{x}_i - \overline{\mathbf{X}} = \mathbf{0} \Rightarrow \\ &\Rightarrow \sum_{i \in s} p_i \mathbf{x}_i = \overline{\mathbf{X}}, \end{aligned} \quad (2.10)$$

que indica que las cantidades p_i dan estimaciones perfectas para $\overline{\mathbf{X}}$, y por tanto, deberían dar una buena aproximación para la media de variable de interés si la relación entre y y \mathbf{x} es lineal.

Cuando $\mathbf{u}_i = \mathbf{x}_i - \overline{\mathbf{X}}$, las soluciones a las ecuaciones (2.7) y (2.8) también son obtenidas por Hartley y Rao (1968) a través de una aproximación similar. Estos autores demostraron que el estimador de regresión es asintóticamente equivalente al estimador dado en (2.9). Un resultado similar puede hacerse para el estimador de la mediana propuesto por Kuk y Mak (1989) cuando $\mathbf{u}_i = \delta(\mathbf{x} \leq M_{\mathbf{x}}) - 0,5$, siendo $M_{\mathbf{x}}$ la mediana de \mathbf{x} , y $\delta(\cdot)$ la función indicadora que toma el valor $\delta(a) = 1$ si $a \geq 0$ y el valor 0 en otro caso.

Puede ocurrir que las ecuaciones (2.7) y (2.8) no tengan solución. Esta situación surge cuando el conjunto convexo $\{\mathbf{u}_i, i \in s\}$ no contiene al $\mathbf{0}$. Se han planteado dos soluciones para este problema:

1. Usar la verosimilitud euclídea propuesta por Owen (1991):

$$\frac{1}{2} \sum_{i \in s} (1 - np_i)^2,$$

y no requerir que $0 \leq p_i \leq 1$.

2. Reemplazar la restricción (2.6) por

$$\sum_{i \in s} p_i \mathbf{u}_i = \tilde{\mathbf{u}},$$

tal que $\tilde{\mathbf{u}}$ está dentro del conjunto convexo y tiende a $\mathbf{0}$.

En cualquier caso, cuando n es grande, la situación en la cual las ecuaciones (2.7) y (2.8) no tienen solución es poco probable.

Existen situaciones extremas en las cuales el método de verosimilitud empírica es incapaz de usar la información auxiliar, como por ejemplo, cuando \mathbf{x} es dicotómica y todas las observaciones son $\mathbf{x}_i = 1$. Estos casos también son poco probables en la práctica.

Estimadores bajo un diseño muestral general

El estimador del apartado anterior está diseñado sólo para muestreo aleatorio simple, y su metodología no se puede extender a otros diseños muestrales más complejos. Chen y Sitter (1999) proponen una aproximación de verosimilitud pseudo empírica que es aplicable a cualquier diseño muestral y coincide bajo muestreo aleatorio simple con el estimador propuesto en Chen y Qin (1993).

El método de verosimilitud empírica para un diseño muestral general asume que la muestra s es seleccionada usando algún diseño muestral, $p(\cdot)$, es decir, la muestra $s \subseteq S$ es extraída con probabilidad $p(s)$. El objetivo es maximizar la verosimilitud de la población en estudio, es decir, maximizar $L^*(\mathbf{p}) = \prod_{i=1}^N p_i$. Por conveniencia, y teniendo en cuenta la monotonía de la función logaritmo, se considera el objetivo de maximizar $l(\mathbf{p}) = \log L^*(\mathbf{p}) = \sum_{i=1}^N \log p_i$. En la práctica, solo se disponen de los valores y_i para las unidades de la muestra, pudiéndose, por tanto, utilizar únicamente las cantidades p_i para $i \in s$. Esto provoca que se necesite una estimación eficiente para $l(\mathbf{p})$. Esta estimación viene dada por la llamada función de verosimilitud pseudo empírica

$$\widehat{l}(\mathbf{p}) = \sum_{i \in s} d_i \log p_i, \quad (2.11)$$

que tiene la propiedad de ser una estimación insesgada bajo el diseño para $l(\mathbf{p})$, esto es

$$E[\widehat{l}(\mathbf{p})] = E \left[\sum_{i \in s} d_i \log p_i \right] = \sum_{i=1}^N \log p_i = l(\mathbf{p}),$$

donde $E[\cdot]$ denota la esperanza bajo el diseño muestral.

La información auxiliar se incorpora a través de la función de calibración

$\mathbf{u}_i = u(y_i, \mathbf{x}_i)$, donde $u(\cdot)$ es una función de y_i y de \mathbf{x}_i que debe satisfacer:

$$\frac{1}{N} \sum_{i=1}^N \mathbf{u}_i = \mathbf{0}.$$

Las cantidades \hat{p}_i necesarias para obtener el *estimador de verosimilitud pseudo empírica* (*PEMLE*) se obtienen maximizando la función dada en (2.11) sujeta a las restricciones (2.5) y (2.6).

Usando el método de los multiplicadores de Lagrange para resolver este problema, se obtiene, para $i \in s$, las cantidades

$$\hat{p}_i = \frac{d_i^*}{1 + \lambda^t \mathbf{u}_i}, \quad (2.12)$$

donde el vector de multiplicadores de Lagrange, λ , es la solución de la expresión:

$$\sum_{i \in s} \frac{d_i^* \mathbf{u}_i}{1 + \lambda^t \mathbf{u}_i} = \mathbf{0}, \quad (2.13)$$

siendo $d_i^* = d_i / \sum_{j \in s} d_j$. El *PEMLE* para la media poblacional se define entonces como

$$\hat{y}_{PE} = \sum_{i \in s} \hat{p}_i y_i. \quad (2.14)$$

Se recuerda que asumiendo una relación lineal entre y y \mathbf{x} se suele considerar la función de calibración $\mathbf{u}_i = \mathbf{x}_i - \bar{\mathbf{X}}$. En este caso, la restricción (2.6) puede expresarse como:

$$\sum_{i \in s} p_i \mathbf{x}_i = \bar{\mathbf{X}}.$$

En el caso de no disponer de información auxiliar, en cuyo caso se toma $\mathbf{u}_i = \mathbf{0}$, el método de verosimilitud empírica produce $\hat{p}_i = d_i^*$, y el *PEMLE* viene dado por

$$\hat{y}_{PE} = \sum_{i \in s} d_i^* y_i,$$

que coincide con el estimador directo para la media poblacional de tipo Hájek. En general, este estimador no coincide con el estimador directo usual de tipo Horvitz-Thompson, aunque se demuestra que disfruta de buenas propiedades respecto a este último (véase Rao, 1966, Basu, 1977 y Sándal *et al.*, 1992). Respecto al problema de la estimación de la función de distribución, el estimador de tipo Hájek disfruta de mejores propiedades, puesto que el estimador de tipo Horvitz-Thompson no cumple las propiedades para ser una verdadera función

de distribución (en concreto $\lim_{t \rightarrow +\infty} \widehat{F}_{HTy}(t) \neq 1$), propiedades que si posee el estimador de tipo Hájek.

Esta propiedad para la función de distribución también se cumple para cualquier función de calibración, y no tan solo para $\mathbf{u}_i = \mathbf{0}$. Esto es, las cantidades \widehat{p}_i dadas en (2.12) son estrictamente positivas y satisfacen $\sum_{i \in s} \widehat{p}_i = 1$ (como puede comprobarse en (2.5)), condiciones necesarias para estimar una verdadera función de distribución, hecho que no sucede, por ejemplo, con los estimadores de regresión generalizados (*GREG*) definidos en Cassel *et al.* (1976) y Särndal (1980) o los estimadores de calibración propuestos en Deville y Särndal (1992).

A continuación, se dan expresiones del *PEMLE* para algunos diseños muestrales más simples y conocidos. De estos ejemplos se desprende que la aplicabilidad de esta metodología no es tan complicada y que estos estimadores están relacionados con otros estimadores tradicionales.

Ejemplo 2.1 *Muestreo Aleatorio Simple.*

Bajo este diseño $\pi_i = n/N$, $d_i = 1/\pi_i = N/n$ y $\sum_{j \in s} d_j = N$, obteniéndose

$$d_i^* = \frac{d_i}{\sum_{j \in s} d_j} = \frac{1}{n}. \quad (2.15)$$

Si no se dispone de información auxiliar, $\mathbf{u}_i = \mathbf{0}$, $\widehat{p}_i = d_i^*$ y el *PEMLE* para la media poblacional está dado por

$$\widehat{y}_{PE} = \sum_{i \in s} \widehat{p}_i y_i = \frac{1}{n} \sum_{i \in s} y_i, \quad (2.16)$$

que coincide con el estimador usual bajo muestreo aleatorio simple (\bar{y}) y con el estimador \widehat{y}_{EL} propuesto en Chen y Qin (1993).

Usando la información auxiliar, el *PEMLE* viene dado por

$$\widehat{y}_{PE} = \sum_{i \in s} \widehat{p}_i y_i, \quad (2.17)$$

donde

$$\widehat{p}_i = \frac{1}{n(1 + \lambda^t \mathbf{u}_i)}, \quad (2.18)$$

y λ es la solución de la ecuación

$$\sum_{i \in s} \frac{\mathbf{u}_i}{1 + \lambda^t \mathbf{u}_i} = \mathbf{0}. \quad (2.19)$$

Puede observarse que este estimador coincide, de nuevo, con el estimador \widehat{y}_{EL} .

Ejemplo 2.2 Muestreo con probabilidades iguales y con reemplazamiento.

En los métodos de muestreo con reemplazamiento se demuestra (véase Hansen y Hurwitz, 1943) que $d_i = 1/(n\alpha_i)$, donde α_i es la probabilidad de que la unidad i -ésima sea seleccionada. Además, al tratarse de un muestreo con probabilidades iguales se tiene que $\alpha_i = 1/N$ y por tanto $d_i = N/n$, que coincide con los pesos básicos en un muestreo aleatorio simple. En consecuencia, las expresiones (2.15), (2.16), (2.17), (2.18) y (2.19) coinciden en este diseño. La única diferencia está en la muestra, es decir, el método para seleccionarla es distinto y además aquí es posible tener unidades repetidas.

Ejemplo 2.3 Muestreo con probabilidades desiguales y sin reemplazamiento.

Se tiene que $d_i = 1/\pi_i$,

$$\widehat{p}_i = \frac{d_i^*}{1 + \lambda^t u_i}, \text{ donde } d_i^* = \frac{1/\pi_i}{\sum_{j \in s} 1/\pi_j},$$

y λ es solución de la ecuación (2.13). Sabido esto, el PEMLE se construye según (2.14).

Bajo este muestreo existen muchos procedimientos para extraer una muestra (consúltese, por ejemplo, Chaudhuri y Vos, 1988). Todos ellos poseen expresiones que permiten calcular las cantidades π_i , necesarias para obtener el PEMLE. En este texto se usan los métodos de Lahiri, Midzuno y Poisson. En el Apéndice B se describen con detalle la puesta en práctica de estos métodos de muestreo y se facilitan las expresiones para las probabilidades de inclusión. Además, en el Apéndice ?? pueden verse funciones en el lenguaje de programación R, material que permite extraer muestras basadas en estos procedimientos de muestreo con probabilidades desiguales.

Ejemplo 2.4 Muestreo con probabilidades desiguales y con reemplazamiento.

Es sabido que en este caso $d_i = 1/(n\alpha_i)$, donde α_i es la probabilidad de que la unidad i -ésima sea seleccionada en cada extracción y por tanto

$$d_i^* = \frac{d_i}{\sum_{j \in s} d_j} = \frac{1/\alpha_i}{\sum_{j \in s} 1/\alpha_j}. \quad (2.20)$$

Y así, el PEMLE se construye mediante la expresión (2.14). En el caso particular de usar el tamaño de cada unidad como una variable auxiliar para la asignación de probabilidades, se tiene que $\alpha_i = M_i/M$, donde M_i es el tamaño de la unidad i , y $M = \sum_{i=1}^N M_i$. Sustituyendo este valor en la expresión (2.20), se obtiene una expresión más simple para el PEMLE.

Una cuestión sin resolver hasta el momento es el procedimiento a seguir para despejar λ en la expresión (2.13), donde además, se ha de verificar que las cantidades \hat{p}_i sean positivas. La resolución de este problema no es tan simple al tratarse de ecuaciones no lineales, debiéndose emplear métodos específicos para la resolución de ecuaciones no lineales, como el de bisección o el de Newton-Raphson. A continuación se describe una modificación del método de Newton-Raphson, propuesto en Chen *et al.* (2002), para el cálculo del PEMLE en caso de que este problema tenga una única solución y ésta exista.

Sea

$$g(\lambda) = \sum_{i \in s} \frac{d_i^* \mathbf{u}_i}{1 + \lambda^t \mathbf{u}_i}.$$

Para una muestra dada, s , el conjunto de valores factibles de λ tal que $\hat{p}_i > 0$ está dado por el conjunto convexo $A = \{\lambda : 1 + \lambda^t \mathbf{u}_i > 0, i \in s\}$. El problema de maximizar la función $\hat{l}(\mathbf{p})$, definida en (2.11), sujeta a las restricciones (2.5) y (2.6) es similar al problema de maximizar la función cóncava

$$\tilde{l}(\lambda) = \sum_{i \in s} d_i^* \log(1 + \lambda^t \mathbf{u}_i),$$

con respecto a λ sobre el conjunto convexo A , puesto que $\partial \tilde{l}(\lambda) / \partial \lambda = g(\lambda)$. Si la única solución de $g(\lambda) = 0$ existe, ésta puede encontrarse aplicando la siguiente modificación del algoritmo de Newton-Raphson:

Algoritmo 2.1

Paso 0: Sea $\lambda_0 = \mathbf{0}$, $k = 0$, $\gamma_0 = 1$ y $\epsilon = 10^{-8}$.

Paso 1: Calcular $\Delta(\lambda_k)$ donde

$$\Delta(\lambda) = \left\{ \frac{\partial}{\partial \lambda} g^*(\lambda) \right\}^{-1} \quad ; \quad g^*(\lambda) = \left\{ - \sum_{i \in s} \frac{d_i^* \mathbf{u}_i \mathbf{u}_i^t}{(1 + \lambda^t \mathbf{u}_i)^2} \right\}^{-1} \sum_{i \in s} \frac{d_i^* \mathbf{u}_i}{1 + \lambda^t \mathbf{u}_i}.$$

Si $\|\Delta(\lambda_k)\| < \epsilon$, se detiene el algoritmo y la solución es λ_k . En otro caso ir al Paso 2

Paso 2: Calcular $\delta_k = \gamma_k \Delta(\lambda_k)$. Si $1 + (\lambda_k - \delta_k)^t \mathbf{u}_i \leq 0$ para algún i o $\tilde{l}(\lambda_k - \delta_k) < \tilde{l}(\lambda_k)$, entonces tomar $\gamma_k = \gamma_k/2$ y repetir el Paso 2.

Paso 3: Considerar $\lambda_{k+1} = \lambda_k - \delta_k$, $k = k + 1$ y $\gamma_{k+1} = (k + 1)^{-1/2}$. Ir al Paso 1.

La expresión $\|\cdot\|$ denota la norma euclídea.

La demostración de este resultado puede consultarse en Chen *et al.* (2002). Así mismo, puede comprobarse que este algoritmo es similar a la modificación del método de Newton descrito en Polyak (1987). Los cambios del paso 2 aseguran que en cada iteración el valor de λ sigue dentro del rango de A y que la función cóncava $\tilde{l}(\lambda)$ se mueve alrededor del punto máximo. El algoritmo es simple, eficiente y la convergencia está garantizada, lo cual indica que, salvo en casos extraños, el *PEMLE* puede siempre obtenerse.

Estimadores bajo muestreo estratificado

La metodología de verosimilitud empírica para obtener estimadores en muestreo de poblaciones finitas se extiende a diseños muestrales más complejos, como por ejemplo muestreo estratificado. Siguiendo la notación clásica del muestreo estratificado y descrita en el Apéndice B, se define la log-función de verosimilitud en muestreo estratificado como

$$l(\mathbf{p}) = \sum_{h=1}^L \sum_{i=1}^{N_h} \log(p_{hi}), \quad (2.21)$$

que puede verse como un total poblacional, cuya estimación insesgada a partir de la muestra s y bajo un diseño muestral específico está dada por

$$\hat{l}(\mathbf{p}) = \sum_{h=1}^L \sum_{i \in s_h} d_{hi} \log(p_{hi}). \quad (2.22)$$

En este caso, d_{hi} son los pesos diseñados básicos que hacen que $\hat{l}(\mathbf{p})$, denominada log-función de verosimilitud pseudo empírica, sea insesgada bajo el diseño para $l(\mathbf{p})$. Por ejemplo, asumiendo muestreo aleatorio simple en cada estrato, se tiene $d_{hi} = N_h/n_h$.

En muestreo estratificado, el *PEMLE* se obtiene maximizando la función (2.22) sujeta a las restricciones

$$\sum_{i \in s_h} p_{hi} = 1 \quad (p_{hi} > 0), \quad h = \{1, \dots, L\}, \quad (2.23)$$

$$\sum_h W_h \sum_{i \in s_h} p_{hi} \mathbf{x}_{hi} = \bar{\mathbf{X}}. \quad (2.24)$$

En la restricción (2.24) se ha considerado por comodidad una relación lineal entre y y \mathbf{x} , aunque es posible modificar esta restricción en caso de existir o considerar oportuno asumir otro tipo de relación entre y y \mathbf{x} .

Una vez obtenidas todas las soluciones \hat{p}_{hi} de este problema, el *PEMLE* bajo muestreo estratificado está dado por

$$\hat{y}_{PEst} = \sum_{h=1}^L W_h \sum_{i \in s_h} \hat{p}_{hi} y_{hi}. \quad (2.25)$$

Dependiendo de si las cantidades $\bar{\mathbf{X}}_h = N_h^{-1} \sum_{i=1}^{N_h} \mathbf{x}_{hi}$ son conocidas o no, el cálculo de este estimador se puede orientar en dos caminos distintos.

En primer lugar, si las cantidades $\bar{\mathbf{X}}_h$ son conocidas para $h = \{1, \dots, L\}$, y asumiendo una relación lineal, la restricción (2.24) puede sustituirse por la restricción

$$\sum_{i \in s_h} p_{hi} \mathbf{x}_{hi} = \bar{\mathbf{X}}_h, \quad h = \{1, \dots, L\}, \quad (2.26)$$

y el problema que se plantea en este caso es maximizar (2.22) sujeta a las restricciones (2.23) y (2.26). Según este planteamiento, el cálculo del *PEMLE* bajo muestreo estratificado es bastante simple, esto es, se calcula el *PEMLE* para cada estrato, \hat{y}_{PEh} , y el estimador final viene dado por

$$\hat{y}_{PEst} = \sum_{h=1}^L W_h \hat{y}_{PEh}.$$

Por otro lado, cuando $\bar{\mathbf{X}}_h$ son desconocidas para cualquier h , la restricción (2.26) no puede establecerse, y el problema de maximizar (2.22) sujeto a las restricciones (2.23) y (2.24) no es una cuestión tan simple. Incluso resulta imposible aplicar el Algoritmo 2.1 bajo muestreo estratificado debido a que la función (2.22) y la restricción (2.24) están formuladas para el conjunto de los estratos, esto es, contienen dobles sumatorias, mientras que la restricción (2.23) está formulada a nivel del estrato, es decir, contiene una sola sumatoria. Existen dos estrategias a seguir para buscar una solución óptima:

- (G1). Considerar en lugar de la restricción (2.24) otra restricción arbitraria para cada estrato y buscar la solución intermedia bajo esta situación. La solución final se obtiene a través del método de verosimilitud empírica.
- (G2). Reemplazar las restricciones de modo que las nuevas estén todas formuladas a nivel del conjunto de los estratos, y por tanto el Algoritmo 2.1 pueda ser aplicado.

La estrategia (G1) fue seguida por Chen y Sitter (1999). El planteamiento que se propuso es el siguiente. El *PEMLE* bajo muestreo estratificado se calcula considerando los pesos \hat{p}_{hi} obtenidos al maximizar la función (2.22) sujeta a las restricciones

$$\begin{aligned}\sum_h \sum_{i \in s_h} p_{hi} &= 1, \\ \sum_h \sum_{i \in s_h} p_{hi} \mathbf{x}_{hi} &= \bar{\mathbf{X}}.\end{aligned}\tag{2.27}$$

Estas restricciones surgen al incorporar la información auxiliar contenida en el tamaño de cada estrato, es decir, toda la información auxiliar usada para construir el *PEMLE* se puede incluir en los vectores $\mathbf{u}_i = \mathbf{U}_i^* - \bar{\mathbf{U}}^*$, donde $i = \{1, \dots, N\}$, $\mathbf{U}_i^* = (\mathbf{x}_i, \nu_{1i}, \dots, \nu_{Li})^t$, $\bar{\mathbf{U}}^* = (\bar{\mathbf{X}}, W_1, \dots, W_L)^t$ y ν_{hi} vale 1 si $i \in h$ y 0 en otro caso. En este sentido, la información de los tamaños de los estratos se usa de forma efectiva, lo cual no ocurre ni con el estimador de regresión generalizado (*GREG*) ni con el estimador óptimo de regresión (*ORE*) propuesto en Rao (1994), y esto hace que se obtengan mejores estimaciones. A su vez, bajo muestreo estratificado, el *ORE* es más eficiente que el *GREG* porque usa la correlación entre y y \mathbf{x} . Asumiendo muestreo estratificado aleatorio, el *PEMLE* es equivalente al *ORE* (y ambos mejores que el *GREG*) puesto que los pesos muestrales son constantes dentro de cada estrato e incluyen el tamaño del estrato que es equivalente a incluir la correlación. No obstante, asumiendo otro diseño muestral, por ejemplo muestreo estratificado con probabilidades proporcionales al tamaño en cada estrato, el *PEMLE* es más eficiente que el *ORE* debido a que usa los tamaños de los estratos que contienen información importante que no es suministrada ni por los pesos muestrales ni por la correlación. En resumen, bajo muestreo estratificado, el *PEMLE* gana en eficiencia respecto a otros estimadores (véase, por ejemplo, Chen y Sitter, 1999, Zhong y Rao, 1996, Zhong y Rao, 2000).

Según lo descrito, se ha de resolver el problema de maximizar (2.22) sujeta a las restricciones (2.27). Como las restricciones

$$\begin{aligned}\sum_{i \in s_h} p_{hi} &= W_h, \quad \forall h = \{1, \dots, L\}, \\ \sum_{i \in s_h} p_{hi} \mathbf{x}_{hi} &= W_h \tilde{\mathbf{x}}_h, \quad \forall h = \{1, \dots, L\},\end{aligned}\tag{2.28}$$

son equivalentes a las dadas en (2.27), el problema se resuelve buscando las cantidades

$$\tilde{\mathbf{x}}_h, \quad h = \{1, \dots, L\}, \quad (2.29)$$

tal que $\sum_h W_h \tilde{\mathbf{x}}_h = \bar{\mathbf{X}}$ y maximizando (2.22) sujeta a las nuevas restricciones (2.28). Aplicando el método de los multiplicadores de Lagrange, la solución que se obtiene es

$$\hat{p}_{hi} = \frac{W_h d_{hi}}{d_h + \lambda_h^t(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)}, \quad (2.30)$$

donde λ_h para $h = \{1, \dots, L\}$, se obtiene de la ecuación

$$\sum_{i \in s_h} \frac{d_{hi}(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)}{d_h + \lambda_h^t(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)} = \mathbf{0}, \quad (2.31)$$

y $d_h = \sum_{i \in s_h} d_{hi}$. Sabido esto, el valor máximo para la función (2.22) es

$$\begin{aligned} & \sum_h \sum_{i \in s_h} d_{hi} \log(\hat{p}_{hi}) = \\ & = - \sum_h \sum_{i \in s_h} d_{hi} \log [d_h + \lambda_h^t(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)] + \end{aligned} \quad (2.32)$$

$$+ \sum_h \sum_{i \in s_h} d_{hi} [\log(d_{hi}) + \log(W_h)]. \quad (2.33)$$

Como (2.33) es constante, se puede maximizar (2.32) respecto a $\tilde{\mathbf{x}}_h$ y bajo la condición $\sum_h W_h \tilde{\mathbf{x}}_h = \bar{\mathbf{X}}$. Notamos que λ_h es una función que depende de $\tilde{\mathbf{x}}_h$. Usando de nuevo el método de Lagrange, se tiene

$$l(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_L, \mathbf{t}) = - \sum_h \sum_{i \in s_h} d_{hi} \log [d_h + \lambda_h^t(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)] - \mathbf{t}^t \left(\sum_{h=1}^L W_h \tilde{\mathbf{x}}_h - \bar{\mathbf{X}} \right).$$

Tomando derivadas respecto a $\tilde{\mathbf{x}}_h$ e igualando al vector de ceros se obtiene

$$- \sum_{i \in s_h} \frac{d_{hi} \left[\frac{\partial \lambda_h^t}{\partial \tilde{\mathbf{x}}_h}(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h) - \lambda_h^t \right]}{d_h + \lambda_h^t(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)} - \mathbf{t}^t W_h = -\lambda_h^t - \mathbf{t}^t W_h = \mathbf{0},$$

y por tanto $\lambda_h^t = W_h \mathbf{t}^t$. La expresión (2.31) puede expresarse como

$$\sum_{i \in s_h} \frac{d_{hi}(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)}{d_h + W_h \mathbf{t}^t(\mathbf{x}_{hi} - \tilde{\mathbf{x}}_h)} = \mathbf{0}. \quad (2.34)$$

Debido a estos desarrollos, puede emplearse el siguiente algoritmo para la búsqueda de los pesos \hat{p}_{hi} necesarios para obtener el *PEMLE* en muestreo estratificado.

Algoritmo 2.2

Paso 1. Fijar un vector \mathbf{t} y obtener las cantidades $\tilde{\mathbf{x}}_h$, $h = \{1, \dots, L\}$, soluciones de la expresión (2.34).

Paso 2. Si $\sum_h W_h \tilde{\mathbf{x}}_h = \bar{\mathbf{X}}$, se calculan las cantidades \hat{p}_{hi} según (2.30), donde $\lambda_h = W_h \mathbf{t}$. En caso contrario, elegir otro \mathbf{t} y volver al paso anterior.

Una vez que hemos calculadas las cantidades \hat{p}_{hi} , con $i \in s_h$ y $h = \{1, \dots, L\}$, mediante el algoritmo anterior, el *PEMLE* está dado por

$$\hat{y}_{PE} = \sum_{h=1}^L \sum_{i \in s_h} \hat{p}_{hi} y_{hi}.$$

Se deben de tener en cuenta las siguientes observaciones cuando se aplica el Algoritmo 2.2:

- Las cantidades $\tilde{\mathbf{x}}_h$ se pueden ver como funciones que dependen de \mathbf{t} , según la expresión (2.34).
- Se tiene que $\sum_h W_h \tilde{\mathbf{x}}_h$ es monótona respecto \mathbf{t} . Esto es importante para determinar las soluciones $\tilde{\mathbf{x}}_h$, puesto que aumentando o disminuyendo el valor \mathbf{t} , es posible llegar fácilmente a ellas.
- La unicidad de la solución está asegurada como consecuencia de la monotonía de $\sum_h W_h \tilde{\mathbf{x}}_h$ respecto \mathbf{t} .

Este algoritmo, que también ha sido descrito en Zhong y Rao (2000), es más eficiente cuando la variable auxiliar \mathbf{x} es unidimensional, puesto que en este caso puede encontrarse la solución incrementando o disminuyendo el valor de \mathbf{t} , el cual es unidimensional. Cuando se tiene más de una variable auxiliar, buscar la solución es un problema más complejo al tener que aumentar o disminuir un vector \mathbf{t} . Además, el cálculo de \hat{p}_{hi} requiere resolver repetidamente sistemas no-lineales de grandes dimensiones según la expresión (2.34), y esto en la práctica es difícil de calcular. Por estas razones, se han buscado aproximaciones alternativas, que sean eficientes y fáciles de llevar a la práctica tanto si se dispone de una variable auxiliar como si son varias.

En Wu (2004b) se detalla el siguiente planteamiento que resuelve los inconvenientes anteriores y se basa en la estrategia (G2).

El objetivo que se persigue es poder aplicar el Algoritmo 2.1 de Chen *et al.* (2002). Para ello, tanto la log-función de verosimilitud pseudo empírica como las restricciones deben estar formuladas para el conjunto de los estratos, esto es, todas deben tener dobles sumatorias. Para este propósito, se tiene que reemplazar la expresión (2.23) por otra similar formulada a nivel poblacional. Sean las restricciones

$$\sum_{h=1}^L W_h \sum_{i \in s_h} p_{hi} = 1, \quad (2.35)$$

$$\sum_{i \in s_h} p_{hi} = 1, \quad h = \{1, \dots, L-1\}. \quad (2.36)$$

Manteniendo al margen (2.35), se combinan (2.36) y (2.24) añadiendo en el vector de variable auxiliares $L-1$ variables indicadoras para cada estrato. Esto es, si $\mathbf{x}_{hi} = (x_{hi1}, \dots, x_{hiP})$, se define

$$\begin{aligned} \mathbf{z}_{1i} &= (1, 0, \dots, 0, x_{1i1}, \dots, x_{1iP})^t, \\ \mathbf{z}_{2i} &= (0, 1, \dots, 0, x_{2i1}, \dots, x_{2iP})^t, \\ &\vdots \\ \mathbf{z}_{(L-1)i} &= (0, 0, \dots, 1, x_{(L-1)i1}, \dots, x_{(L-1)iP})^t, \\ \mathbf{z}_{Li} &= (0, 0, \dots, 0, x_{Li1}, \dots, x_{LiP})^t, \end{aligned}$$

y $\bar{\mathbf{Z}} = (W_1, \dots, W_{L-1}, \bar{X}_1, \dots, \bar{X}_P)^t$, siendo $(\bar{X}_1, \dots, \bar{X}_P)^t = \bar{\mathbf{X}}$. Así, las restricciones (2.36) y (2.24) se pueden combinar mediante la restricción

$$\sum_{h=1}^L W_h \sum_{i \in s_h} p_{hi} \mathbf{z}_{hi} = \bar{\mathbf{Z}}. \quad (2.37)$$

El problema de maximizar $\hat{l}(p)$ sujeta a (2.23) y (2.24) es equivalente a maximizar $\hat{l}(p)$ sujeta a (2.35) y (2.37). Usando el método de los multiplicadores de Lagrange a éste último planteamiento, se obtiene

$$\hat{p}_{hi} = \frac{d_{hi}^*}{1 + \lambda^t \mathbf{u}_{hi}} = \mathbf{0},$$

donde

$$d_{hi}^* = \frac{d_{hi}}{W_h \sum_{h=1}^L \sum_{i \in s_h} d_{hi}}, \quad \mathbf{u}_{hi} = \mathbf{z}_{hi} - \bar{\mathbf{Z}},$$

y λ es solución de

$$\sum_{h=1}^L \sum_{i \in s_h} \frac{d_{hi} \mathbf{u}_{hi}}{1 + \lambda^t \mathbf{u}_{hi}} = \mathbf{0}. \quad (2.38)$$

En esta situación es posible aplicar el Algoritmo 2.1, estando garantizada la convergencia a la única solución, si tal solución existe.

Ejemplo 2.5 *Estimadores bajo muestreo bifásico.*

Los estimadores comentados hasta el momento en esta sección están basados en un diseño muestral general y utilizan el vector media poblacional de las variables auxiliares para obtener las estimaciones. Cuando este vector es desconocido, ni los estimadores de verosimilitud empírica ni cualquier otro estimador basado en información auxiliar puede ser utilizado, puesto que la mayoría de éstos se construyen con ayuda de $\bar{\mathbf{X}}$ para mejorar la precisión en la estimación de parámetros de la variable de interés. Véase, por ejemplo, Cochran (1977) y Särndal et al. (1992) para consultar los numerosos estimadores en la literatura del muestreo de poblaciones finitas que hacen uso de la información auxiliar.

En la situación anterior, donde tan solo se conocen los datos muestrales de las variables auxiliares, es necesario estimar $\bar{\mathbf{X}}$ o intentar dar una buena aproximación mediante alguna técnica o recurso. El muestreo bifásico (también denominado muestreo doble o en dos fases) permite estimar estas cantidades desconocidas y por tanto, es posible utilizar todos los métodos basados en información auxiliar. En el Apéndice B puede consultarse una completa definición del muestreo bifásico y una lista de los principales conceptos de este diseño muestral.

De este modo, en este ejemplo resuelve el problema de la estimación de parámetros lineales en muestreo bifásico con diseños muestrales arbitrarios en cada una de las dos fases y aplicando el método de verosimilitud empírica.

En muestreo bifásico, el método de verosimilitud empírica puede ser aplicado como sigue. El PEMLE viene dado por

$$\hat{\mathbf{y}}_{PEb} = \sum_{i \in s} \hat{p}_i y_i \quad (2.39)$$

donde los pesos \hat{p}_i maximizan la log-función de verosimilitud pseudo empírica

$$\hat{l}(p) = \sum_{i \in s} d_i \log(p_i) \quad (2.40)$$

sujeta a las restricciones

$$\sum_{i \in s} p_i = 1 \quad (p_i \geq 0) \quad (2.41)$$

$$\sum_{i \in s} p_i \mathbf{u}'_i = \mathbf{0} \quad (2.42)$$

donde para todo $i \in s$, $d_i = d'_i d_{i/s'}$, y \mathbf{u}'_i es una función que depende de y y de los valores de \mathbf{x} obtenidos en la muestra de la primera fase, s' . Además, esta función ha de verificar

$$\frac{1}{n'} \sum_{i \in s'} \mathbf{u}'_i = \mathbf{0}.$$

Asumiendo relación lineal entre y y \mathbf{x} , es usual considerar $\mathbf{u}'_i = \mathbf{x}_i - \bar{\mathbf{x}}'$, y la restricción (2.42) se puede expresar como

$$\sum_{i \in s} p_i \mathbf{x}_i = \frac{1}{n'} \sum_{i \in s'} \mathbf{x}_i = \bar{\mathbf{x}}',$$

que viene a indicar que si los pesos que van a ser estimados se ponderan sobre los datos muestrales del vector de variables auxiliares de la segunda fase, se obtendrá la cantidad $\bar{\mathbf{x}}'$, es decir, la media muestral del vector de las variables auxiliares obtenida a partir de la muestra de la primera fase. De ahí la importancia de realizar un gran esfuerzo para obtener una buena estimación para $\bar{\mathbf{X}}$ con los datos de la muestra de la primera fase.

La solución del problema planteado se resuelve por el método de los multiplicadores de Lagrange, obteniendo como solución para todo $i \in s$ las cantidades

$$\hat{p}_i = \frac{d_i^*}{1 + \lambda^t \mathbf{u}'_i},$$

donde

$$d_i^* = \frac{d_i}{\sum_{j \in s} d_j} = \frac{d'_i d_{i/s'}}{\sum_{j \in s} d'_j d_{j/s'}},$$

y λ es el vector de multiplicadores de Lagrange que se obtiene de la ecuación

$$\sum_{i \in s} \frac{d_i^* \mathbf{u}'_i}{1 + \lambda^t \mathbf{u}'_i} = \mathbf{0}.$$

2.2.2. Propiedades teóricas

En esta sección se describen las propiedades asintóticas más importantes de los estimadores de verosimilitud empírica basados en el diseño muestral. En primer lugar, se describen las propiedades teóricas más importantes del estimador de verosimilitud empírica propuesto en Chen y Qin (1993) bajo muestreo aleatorio simple. A continuación, se demuestra la relación que tiene el *PEMLE* con los conocidos estimadores de regresión. Esta sección de propiedades se completa con el estudio de los estimadores de verosimilitud empírica en muestreo estratificado y su relación con otros estimadores.

Propiedades en muestreo aleatorio simple

A continuación se estudian las propiedades asintóticas del estimador de verosimilitud empírica descrito en Chen y Qin (1993). Asumamos muestreo aleatorio simple, donde el tamaño de la muestra, n , y el tamaño de la población, N , tienden a infinito cuando un cierto índice, ν , tiende a infinito, es decir, existe una secuencia de poblaciones finitas indexadas por ν , donde $\pi_\nu = \{(x_{1\nu}, y_{1\nu}), \dots, (x_{N\nu}, y_{N\nu})\}$ y el tamaño poblacional N_ν tiende a infinito. Por comodidad, se suprime el índice ν siempre que sea posible y se considera sólo una variable auxiliar. Sea

$$\sigma_x^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2, \quad \sigma_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2,$$

$$\sigma_{xy} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})(y_i - \bar{Y}),$$

y $\bar{x}, \bar{y}, s_x^2, s_y^2$ y s_{xy} sus correspondientes versiones muestrales. Se considera que la función de calibración satisface $\sum_{i=1}^N u_i = 0$ y se tiene que

$$\sigma_u^2 = \frac{1}{N-1} \sum_{i=1}^N u_i^2, \quad \sigma_{yu} = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})u_i.$$

La media poblacional de variable de interés se estima a través del estimador $\widehat{y}_{EL} = \sum_{i \in s} \widehat{p}_i y_i$. Los siguientes teoremas pueden ser definidos.

Teorema 2.1 *Suponiendo que cuando $\nu \rightarrow \infty$, el tamaño poblacional N , el tamaño muestral n , y $N - n$ tienden a infinito, y*

$$\frac{1}{N} \left\{ \sum_{i=1}^N |u_i|^3 \right\}, \quad \frac{1}{N} \left\{ \sum_{i=1}^N |y_i|^3 \right\},$$

tienen una cota superior independiente de ν , entonces se verifica

$$\frac{n^{1/2}(\widehat{y}_{EL} - \bar{Y})}{\sigma_\nu} \rightarrow N(0, 1),$$

donde $\sigma_\nu^2 = \left(1 - \frac{n}{N}\right) \left(\sigma_y^2 - \frac{\sigma_{yu}^2}{\sigma_u^2}\right)$.

La demostración de este resultado puede consultarse en Chen y Qin (1993). Una consecuencia importante que puede observarse de este teorema es que a mayor correlación entre u e y , mayor será la ganancia en precisión. Se demuestra que la eficiencia asintótica del método es equivalente a la del método de regresión.

En la práctica, la cantidad σ_v^2 es desconocida, con lo que se tiene que buscar un buen estimador. Una alternativa es la estimación de σ_y^2 , σ_{yu} y σ_u^2 por separado, aunque para tamaños muestrales moderados trabaja mejor el estimador jackknife para la varianza. En el siguiente teorema, debido a Chen y Qin (1993), se demuestra que el estimador jackknife es un buen estimador para σ_v^2 .

Teorema 2.2 *Bajo las mismas condiciones del Teorema 2.1, si $\widehat{y}_{EL}(-j)$ es el estimador cuando la observación j -ésima es eliminada y*

$$\widehat{\sigma}_J^2 = \left(1 - \frac{n}{N}\right) (n-1) \sum_{i \in s} (\widehat{y}_{EL}(-j) - \widehat{y}_{EL})^2,$$

entonces,

$$\widehat{\sigma}_J^2 - \sigma_v^2 = o_p(1).$$

Propiedades para un diseño muestral general

En lo que sigue, se asume una sola variable auxiliar y la función de calibración $u_i = x_i - \bar{X}$. Consideremos también las siguientes condiciones

$$(C2.1). \quad u^* = \max_{i \in s} |u_i| = o_p(n^{1/2}),$$

$$(C2.2). \quad \frac{\sum_{i \in s} d_i u_i}{\sum_{i \in s} d_i u_i^2} = O_p(n^{-1/2}).$$

El siguiente teorema, debido a Chen y Sitter (1999), puede establecerse

Teorema 2.3 *Bajo las condiciones (C2.1) y (C2.2), el PEMLE de \bar{Y} cuando \bar{X} es conocida, es asintóticamente equivalente al estimador de regresión generalizado (GREG). Es decir,*

$$\lambda = \frac{\bar{x}_w - \bar{X}}{\sum_{i \in s} d_i^* (x_i - \bar{x}_w)^2} + o_p(n^{-1/2}),$$

y así $\widehat{y}_{PE} = \widehat{y}_{GREG} + o_p(n^{-1/2})$, donde

$$\widehat{y}_{GREG} = \sum_{i \in s} \widetilde{d}_i y_i, \quad \widetilde{d}_i = d_i \left[1 - \frac{(x_i - \bar{x}_w)(\bar{x}_w - \bar{X})}{\sum_{i \in s} d_i (x_i - \bar{x}_w)^2} \right],$$

$$\bar{y}_w = \sum_{i \in s} d_i y_i, \quad \bar{x}_w = \sum_{i \in s} d_i x_i \quad y \quad d_i^* = \frac{d_i}{\sum_{j \in s} d_j}.$$

Las condiciones (C2.1) y (C2.2) deben satisfacerse para que este teorema pueda establecerse. Sin embargo, estas condiciones no son muy restrictivas y los diseños muestrales más conocidos las satisfacen. En Chen y Sitter (1999) se demuestra cómo estas condiciones se cumplen en tres diseños comunes, como son, el muestreo con probabilidades proporcionales al tamaño con reemplazamiento, el método de Rao-Hartley-Cochran y el muestreo por conglomerados.

Un punto importante es la estimación de la varianza del estimador \widehat{y}_{PE} . Según el Teorema 2.3, resulta evidente que cualquier estimador de la varianza consistente para \widehat{y}_{GREG} será consistente para el *PEMLE*. Aunque esto es asintóticamente válido, no es atractivo usar un estimador de la varianza del *GREG* para estimar la varianza del *PEMLE*. Una óptima alternativa es aplicar estimadores de la varianza remuestreados, tales como jackknife, bootstrap y replicaciones de muestras repetidas balanceadas (ver Shao y Wu (1989, 1992), Chen y Qin (1993) y Shao (1994)) sobre \widehat{y}_{PE} , recalculando \widehat{p}_i en cada remuestra.

Propiedades en muestreo estratificado

La primera propiedad del *PEMLE* en muestreo estratificado se basa en el Teorema 2.3.

Corolario 2.1 *Bajo las condiciones (C2.1) y (C2.2) se tiene*

$$\widehat{y}_{PE} = \bar{y}_w - \frac{\sum_{h=1}^L \sum_{i \in s_h} d_{hi}^* (x_{hi} - \bar{x}_w) y_{hi}}{\sum_{h=1}^L \sum_{i \in s_h} d_{hi}^* (x_{hi} - \bar{x}_w)^2} (\bar{x}_w - \bar{X}) + o_p(n^{-1/2})$$

donde

$$n = \sum_{h=1}^L n_h, \quad \bar{y}_w = \sum_{h=1}^L \sum_{i \in s_h} d_{hi}^* y_{hi},$$

$$\bar{x}_w = \sum_{h=1}^L \sum_{i \in s_h} d_{hi}^* x_{hi} \quad y \quad d_{hi}^* = \frac{d_{hi}}{\sum_{h=1}^L \sum_{j \in s_h} d_{hj}}.$$

Considerando muestreo aleatorio estratificado, es decir, cuando $d_{hi} = N_h/n_h$, la expresión anterior se reduce a

$$\begin{aligned} \widehat{\bar{y}}_{PE} &= \bar{y}_{st} - \frac{\sum_{h=1}^L \sum_{i \in s_h} W_h (x_{hi} - \bar{x}_{st}) y_{hi} / n_h}{\sum_{h=1}^L \sum_{i \in s_h} W_h (x_{hi} - \bar{x}_{st})^2 / n_h} (\bar{x}_{st} - \bar{X}) + o_p(n^{-1/2}) = \\ &= \bar{y}_{GREG} + o_p(n^{-1/2}), \end{aligned}$$

donde

$$\bar{y}_{st} = \sum_{h=1}^L W_h \bar{y}_h \quad y \quad \bar{x}_{st} = \sum_{h=1}^L W_h \bar{x}_h.$$

Esta no es la mejor aproximación posible, puesto que se sabe que el estimador de regresión óptimo (*ORE*), definido en Rao (1994), funciona mejor que el *GREG* en muestreo estratificado. Por este motivo, en Chen y Sitter (1999) se busca una mejor aproximación. En el siguiente corolario se relaciona el *PEMLE* con el *ORE* bajo muestreo aleatorio estratificado. Para ello, se asume que existe una secuencia de poblaciones finitas indexadas por ν , tal que cuando $\nu \rightarrow \infty$ se verifican las condiciones

$$(C2.3). \quad 0 \leq c_1 \leq \sum_{h=1}^L W_h \sigma_h^2 \leq c_2 \leq \infty,$$

$$(C2.4). \quad \max\{n_h^{-1} W_h\} = O(n^{-1}),$$

$$(C2.5). \quad N^{-1} \sum_{h=1}^L \sum_{i=1}^{N_h} |x_{hi}|^3 = O(1),$$

$$(C2.6). \quad N^{-1} \sum_{h=1}^L \sum_{i=1}^{N_h} |y_{hi}|^3 = O(1).$$

Corolario 2.2 *Bajo muestreo aleatorio estratificado y las condiciones (C2.3), (C2.4), (C2.5) y (C2.6), el PEMLE de \bar{Y} , cuando \bar{X} es conocida, es asintóticamente equivalente a \bar{y}_{st}^* , esto es, $\widehat{\bar{y}}_{PE} = \bar{y}_{st}^* + o_p(n^{-1/2})$, donde*

$$\bar{y}_{st}^* = \bar{y}_{st} - \frac{\sum_{h=1}^L W_h \sum_{i \in s_h} (x_{hi} - \tilde{x}_h) y_{hi} / n_h}{\sum_{h=1}^L W_h \sum_{i \in s_h} (x_{hi} - \tilde{x}_h)^2 / n_h} (\bar{x}_{st} - \bar{X}),$$

y las cantidades \tilde{x}_h están definidas en (2.29). Cuando L permanece finito, $\tilde{x}_h - \bar{x}_h = O_p(n^{-1/2})$ y el estimador $\widehat{\bar{y}}_{PE}$ es asintóticamente equivalente al estimador lineal óptimo dado en Rao (1994).

Asumiendo otros diseños muestrales en cada estrato, las comparaciones con respecto otros estimadores son demasiado dificultosas y se ha de recurrir a la simulación para realizar las comparaciones.

En este caso, la estimación de la varianza se obtiene también a través de estimadores de la varianza remuestreados. En Chen y Sitter (1999), se demuestra que bajo muestreo aleatorio estratificado el estimador de la varianza jackknife para el *PEMLE* es consistente.

2.2.3. Estimadores modelo-calibrados

Una de las restricciones considerada en los estimadores de verosimilitud empírica viene dada por

$$\sum_{i \in s} p_i \mathbf{u}_i = \mathbf{0}, \quad (2.43)$$

donde $\mathbf{u}_i = u(y_i, \mathbf{x}_i)$ y $u(\cdot)$ es una función conocida de y y de \mathbf{x} que verifica

$$\frac{1}{N} \sum_{i=1}^N \mathbf{u}_i = \mathbf{0}. \quad (2.44)$$

Asumiendo una relación lineal entre la característica de interés y el vector auxiliar de variables, se utiliza frecuentemente la expresión $\mathbf{u}_i = \mathbf{x}_i - \bar{\mathbf{X}}$, y se plantea la cuestión de cómo efectivo es el uso que se está haciendo de la información auxiliar. Si tal relación no es lineal, los estimadores de verosimilitud empírica obtenidos a partir de la expresión $\mathbf{u}_i = \mathbf{x}_i - \bar{\mathbf{X}}$ pueden resultar ineficaces y surge, por tanto, el problema de encontrar una función de calibración apropiada para los datos del estudio, es decir, que se adapte a cada situación para poder usar la información auxiliar de la mejor manera posible. Una alternativa eficiente para resolver este problema es el uso de los estimadores modelo-calibrados, los cuales están basados en modelos de superpoblación.

Recientemente, en la literatura del muestreo se están utilizando a menudo estimaciones que no están basadas en el diseño muestral, sino que dependen de un determinado modelo de superpoblación que relaciona la variable de interés a través de las variables auxiliares. Tales procedimientos son los estimadores basados en modelos y los estimadores modelo-calibrados. Con la aparición de los modelos de superpoblación la teoría de muestreo tuvo un gran empuje pues se le dotó de un instrumento muy valioso que permitió obtener resultados más concluyentes en la comparación de estrategias y eventualmente producir estrategias óptimas en varias situaciones. Ejemplos e información sobre modelos

de superpoblación pueden consultarse, por ejemplo, en Godambe (1955), Godambe y Thompson (1973), Cassel *et al.* (1976), Pérez (2002) y Sánchez-Crespo (2002).

Los estimadores modelo-calibrados están propuestos en Wu y Sitter (2001a), y se obtienen adaptando un modelo de superpoblación, y a continuación, usar los valores estimados mediante este modelo en la etapa de estimación. Así, se obtiene una función eficiente de calibración, y además es posible encontrar la mejor función $u(\cdot)$ en el sentido de mínima esperanza bajo un modelo de superpoblación de la varianza asintótica basada en el diseño.

Los valores \mathbf{u}_i pueden expresarse como

$$\mathbf{u}_i = \mathbf{w}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i,$$

donde \mathbf{w}_i es una función conocida. Es fácil demostrar que bajo esta situación también se verifica (2.44), y por tanto, se cumplen las condiciones necesarias para aplicar la metodología de verosimilitud empírica. Operando en la restricción (2.43) se llega a la restricción alternativa

$$\sum_{i \in s} p_i \mathbf{w}_i = \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i, \quad (2.45)$$

que es la que suele usarse en los estimadores modelo-calibrados de verosimilitud empírica. Por tanto, el problema de buscar unos valores óptimos \mathbf{u}_i para obtener estimadores más eficientes, es similar al de encontrar la cantidades \mathbf{w}_i , para $i \in s$.

La idea de definir estimadores óptimos bajo un modelo y asumiendo el criterio de mínima esperanza bajo un modelo de superpoblación de la varianza asintótica basada en el diseño ha sido discutida por diversos autores, véase, por ejemplo, Godambe (1955), Godambe y Thompson (1973) y Cassel *et al.* (1976).

Un primer estimador modelo-calibrado surge cuando se asume el siguiente esquema asintótico. Existe una secuencia de poblaciones finitas indexadas por ν . El tamaño poblacional y el tamaño muestral para la población ν -ésima se denotan como N_ν y n_ν . Cuando $\nu \rightarrow \infty$, $N_\nu \rightarrow \infty$ y $n_\nu \rightarrow \infty$. El índice ν se suprimirá para simplificar notación. Por ejemplo, véase Isaki y Fuller (1982) para un mayor detalle de este esquema asintótico. Por último, sea y_1, y_2, \dots, y_N una muestra aleatoria de un modelo de superpoblación ξ tal que

$$E_\xi(y_i) = \mu_i, \quad V_\xi(y_i) = \sigma_i^2, \quad i = \{1, 2, \dots, N\}, \quad (2.46)$$

y y_1, y_2, \dots, y_N son independientes entre ellos. E_ξ y V_ξ denotan la esperanza y la varianza bajo el modelo de superpoblación.

Sea \tilde{y}_{C_w} el estimador de verosimilitud pseudo empírica modelo-calibrado de \bar{Y} cuando se usa $C_w = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\}$ en la restricción (2.45) y L^* un conjunto de secuencias $C_w = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\}$ que verifican

$$\frac{1}{N} \sum_{i=1}^N (\mathbf{w}_i)^6 = O(1) \quad y \quad \frac{1}{N} \sum_{i=1}^N (\mathbf{w}_i)^2 \rightarrow \mathbf{c} \neq \mathbf{0} \text{ cuando } N \rightarrow \infty.$$

Estas condiciones sobre la secuencia $C_w \in L^*$ no son muy restrictivas y se usan para facilitar las demostraciones. Asumiremos que $\{\mu_1, \dots, \mu_N\} \in L^*$.

Se dice que un diseño muestral es regular si el diseño que resulta de un tamaño de muestra indexado tiene probabilidades de inclusión π_i y π_{ij} independientes de la característica y_i dada \mathbf{x}_i , y además satisface las siguientes condiciones:

$$(C2.7). \quad \max_{i \in s} \left(\frac{nd_i}{N} \right) = O(1).$$

$$(C2.8). \quad \frac{1}{N} \sum_{i \in s} d_i \mathbf{w}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i = O_p(n^{-1/2}) \text{ para cualquier secuencia de funciones } (\mathbf{w}_1, \dots, \mathbf{w}_N) \in L^*.$$

En Wu (2003) se demuestra que entre todas las clases de estimadores \tilde{y}_{C_w} con $C_w = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\} \in L^*$, el valor $C_\mu = \{\mu_1, \dots, \mu_N\}$ como variable de calibración en (2.45) minimiza $E_\xi[AV_p(\tilde{y}_{C_w})]$ bajo el modelo (2.46) y para cualquier diseño muestral regular. AV_p denota la varianza asintótica bajo el diseño. Así, el estimador de verosimilitud pseudo empírica modelo-calibrado (*MCPE*) que presenta la propiedad arriba comentada, se construye tomando $\mathbf{w}_i = \mu_i$, o lo que es lo mismo, tomando $\mathbf{u}_i = \mu_i - N^{-1} \sum_{i=1}^N \mu_i$. Sustituyendo estas cantidades \mathbf{u}_i en las expresiones (2.12) y (2.13) se obtiene un primer estimador de verosimilitud empírica basado en la aproximación modelo-calibrada.

Otra alternativa para construir estimadores modelo-calibrados es asumir que y_1, y_2, \dots, y_N es una muestra aleatoria de un modelo de superpoblación semiparamétrico ξ tal que

$$E_\xi(y_i|\mathbf{x}_i) = \mu_i = \mu(\mathbf{x}_i, \theta), \quad V_\xi(y_i|\mathbf{x}_i) = \nu_i^2 \sigma^2, \quad i = \{1, \dots, N\}, \quad (2.47)$$

donde $\theta = (\theta_0, \theta_1, \dots, \theta_P)^t$ y σ^2 son parámetros poblacionales desconocidos, $\mu(\mathbf{x}, \theta)$ es una función conocida de \mathbf{x} y de θ , ν_i es una función conocida de

\mathbf{x}_i o bien de $\mu_i = \mu(\mathbf{x}_i, \theta)$ y E_ξ y V_ξ denotan la esperanza y la varianza con respecto al modelo de superpoblación. Además, se asume que los pares $(y_1, \mathbf{x}_1); (y_2, \mathbf{x}_2); \dots; (y_N, \mathbf{x}_N)$ son mutuamente independientes.

Este modelo es bastante general, e incluye dos casos muy importantes:

1. El modelo de regresión lineal o no lineal

$$y_i = \mu(\mathbf{x}_i, \theta) + \nu_i \epsilon_i \quad i = \{1, \dots, N\},$$

donde ϵ_i son variables aleatorias independientes e idénticamente distribuidas, con $E_\xi(\epsilon_i) = 0$, $V_\xi(\epsilon_i) = \sigma^2$ y $\nu_i = \nu(x_i)$ una función conocida y estrictamente positiva que depende de \mathbf{x}_i .

2. El modelo lineal generalizado

$$g(\mu_i) = \mathbf{x}_i^t \theta \quad V_\xi(y_i | \mathbf{x}_i) = \nu(\mu_i) \quad i = \{1, \dots, N\},$$

donde $\mu_i = E_\xi(y_i | \mathbf{x}_i)$, $g(\cdot)$ es una función de enlace y $\nu(\cdot)$ es la función de varianza.

Los verdaderos parámetros del modelo son desconocidos, aunque pueden estimarse mediante cualquier método basado en el diseño. Asumiendo una aproximación basada en el modelo, la dupla (y_i, \mathbf{x}_i) con $i \in s$ puede verse como una muestra independiente idénticamente distribuida del modelo de superpoblación. Los parámetros θ se pueden estimar usando procedimientos estándares. Bajo el enfoque basado en el diseño, los datos muestrales pueden no seguir la misma estructura del modelo que la población finita completa bajo un esquema muestral complejo, y θ puede carecer de sentido desde el punto de vista del diseño. En este caso, θ se reemplaza por θ_N , una estimación de θ basada en los datos de la población completa. θ_N se reemplaza entonces por $\hat{\theta}$, una estimación basada en el diseño de los datos muestrales (véase Godambe y Thompson, 1986).

Asumiendo el modelo (2.47), el estimador de verosimilitud pseudo empírico modelo-calibrado se construye tomando $\mathbf{w}_i = \mu(\mathbf{x}_i, \hat{\theta})$. Los valores \mathbf{u}_i vienen dados por $\mathbf{u}_i = \hat{\mu}_i - N^{-1} \sum_{i=1}^N \hat{\mu}_i$, donde $\hat{\mu}_i = \mu(\mathbf{x}_i, \hat{\theta})$. Considerando estas cantidades en las expresiones (2.12) y (2.13) se obtiene el *MCPE*.

Al igual que ocurre bajo el primer *MCPE* que se ha definido, en Wu (2003) se demuestra que entre todas las clases de estimadores \tilde{y}_{C_w} , donde $C_w = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\} \in L^*$, el valor $C_\mu = \{\mu(\mathbf{x}_1, \theta), \dots, \mu(\mathbf{x}_N, \theta)\}$ como

variable de calibración en (2.45) minimiza $E_{\xi}[AV_p(\tilde{y}_{C_w})]$ bajo el modelo (2.47) y para cualquier diseño muestral regular.

A continuación se resumen las observaciones más importantes sobre los estimadores de verosimilitud empírica basados en una aproximación modelo-calibrada.

1. En Wu y Sitter (2001a) se demuestra que reemplazar θ por $\hat{\theta}$ en $\mu_i = \mu(\mathbf{x}_i, \theta)$, no cambia asintóticamente el estimador resultante.
2. Con probabilidad tendiendo a uno, el *MCPE* existe y se puede calcular usando el algoritmo 2.1 de Chen *et al.* (2002).
3. El uso efectivo de la información auxiliar depende los los parámetros estimados y de la relación entre la variable respuesta y las covarianzas. Por tanto, usar la calibración sobre las variables auxiliares sin un estudio exhaustivo previo no es usualmente una buena aproximación.
4. Es sabido que para una relación lineal entre y y el vector de variables auxiliares, se toma $\mathbf{u}_i = \mathbf{x}_i - \bar{\mathbf{X}}$ para la construcción del *PEMLE*. En esta situación, el *PEMLE* y el *MCPE* son asintóticamente equivalentes si se considera $\hat{\mu}_i = \mathbf{x}_i^t \hat{\theta}$ como variable de calibración para el cálculo de la aproximación modelo-calibrada. La demostración de este resultado puede consultarse en Wu y Sitter (2001a).
5. Si la relación entre y y \mathbf{x} es lineal, tan sólo el conocimiento de $\bar{\mathbf{X}}$ es suficiente para obtener estimadores eficientes para la media o el total poblacional. Si dicha relación no es lineal o el parámetro de interés no es una función lineal, una información auxiliar completamente disponible y/o más datos sobre el modelo son esenciales para una estimación óptima.
6. Al igual que se ha comentado anteriormente, las cantidades \hat{p}_i son positivas. Esta propiedad no se cumple ni en los estimadores de calibración ni en calculo del *GREG* y juega un papel muy importante en la estimación de otros parámetros de interés en el muestreo, como son la función de distribución, cuantiles, varianza y otras funciones cuadráticas.

2.2.4. Propiedades teóricas

Sea el esquema asintótico siguiente: se asume que existe una secuencia de diseños muestrales y una secuencia de poblaciones finitas indexadas por ν . El

tamaño muestral n_ν y el tamaño poblacional N_ν se aproximan a infinito cuando $\nu \rightarrow \infty$.

Las condiciones siguientes son necesarias para poder aplicar el Teorema 2.4.

$$(C2.9). \quad \widehat{\theta} = \theta_N + O_p(n^{-1/2}) \text{ y } \theta_N \rightarrow \theta.$$

$$(C2.10). \quad \text{Para cada } \mathbf{x}_i, \frac{\partial \mu(\mathbf{x}_i, \mathbf{t})}{\partial \mathbf{t}} \text{ es continua en } \mathbf{t} \text{ y}$$

$$\left| \frac{\partial \mu(\mathbf{x}_i, \mathbf{t})}{\partial \mathbf{t}} \right| \leq h(\mathbf{x}_i, \theta)$$

para \mathbf{t} en un entorno de θ , y $N^{-1} \sum_{i=1}^N h(\mathbf{x}_i, \theta) = O_p(1)$.

(C2.11). Los pesos básicos muestrales, $d_i = \pi_i^{-1}$, hacen que los estimadores de Horvitz-Thompson para ciertas medias muestrales estén asintóticamente normalmente distribuidos.

$$(C2.12). \quad \mathbf{u}^* = \max_{i \in s} |\mathbf{u}_i| = o_p(n^{1/2}), \text{ donde } \mathbf{u}_i = \mu(\mathbf{x}_i, \theta_N) - \frac{1}{N} \sum_{i=1}^N \mu(\mathbf{x}_i, \theta_N).$$

$$(C2.13). \quad \frac{\sum_{i \in s} d_i \mathbf{u}_i}{\sum_{i \in s} d_i \mathbf{u}_i^2} = O_p(n^{1/2}).$$

$$(C2.14). \quad h^* = \max_{i \in s} |h_i| = o_p(n), \text{ siendo } h_i = h(\mathbf{x}_i, \theta_N).$$

El siguiente teorema puede establecerse.

Teorema 2.4 *Bajo el esquema asintótico descrito y las condiciones anteriores (C2.9)~(C2.14), se tiene que*

$$\widehat{y}_{MCPE} = \widehat{y}_{MC} + o_p(n^{-1/2}),$$

donde \widehat{y}_{MC} es el estimador modelo-calibrado para la media obtenido mediante el modelo de calibración y cuya expresión viene dada por

$$\widehat{y}_{MC} = \widehat{y}_{HT} + \left\{ \frac{1}{N} \sum_{i=1}^N \widehat{\mu}_i - \frac{1}{N} \sum_{i \in s} d_i \widehat{\mu}_i \right\} \widehat{B}_N,$$

con

$$\widehat{B}_N = \frac{\sum_{i \in s} d_i q_i (\widehat{\mu}_i - \bar{\mu})(y_i - \bar{y})}{\sum_{i \in s} d_i q_i (\widehat{\mu}_i - \bar{\mu})^2}, \quad \bar{y} = \frac{\sum_{i \in s} d_i q_i y_i}{\sum_{i \in s} d_i q_i} \quad \text{y} \quad \bar{\mu} = \frac{\sum_{i \in s} d_i q_i \widehat{\mu}_i}{\sum_{i \in s} d_i q_i}.$$

Las cantidades q_i son constantes positivas.

Puesto que \widehat{y}_{MCPE} es asintóticamente equivalente al \widehat{y}_{MC} , las mismas expresiones de la varianza y del estimador de la varianza de \widehat{y}_{MC} pueden usarse para \widehat{y}_{MCPE} . Estas varianzas asintóticas basadas en el diseño vienen dadas por

$$V(\widehat{y}_{MCPE}) = \frac{1}{N^2} \sum_{i < j}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{U_i}{\pi_i} - \frac{U_j}{\pi_j} \right)^2,$$

donde π_{ij} son las probabilidades de inclusión de segundo orden, $U_i = y_i - \mu_i B_N$, $\mu_i = \mu(\mathbf{x}_i, \theta_N)$,

$$B_N = \frac{\sum_{i=1}^N q_i (\mu_i - \bar{\mu}_N) (y_i - \bar{Y})}{\sum_{i=1}^N q_i (\mu_i - \bar{\mu}_N)^2} \quad \text{y} \quad \bar{\mu}_N = \frac{1}{N} \sum_{i=1}^N \mu_i.$$

Un estimador para esta varianza viene dado por

$$\widehat{V}(\widehat{y}_{MCPE}) = \frac{1}{N^2} \sum_{i < j}^n \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{u_i}{\pi_i} - \frac{u_j}{\pi_j} \right)^2,$$

donde $u_i = y_i - \widehat{\mu}_i \widehat{B}_N$.

Estas varianzas asintóticas y la demostración del teorema se pueden consultar en Wu y Sitter (2001a).

Aunque estas aproximaciones son asintóticamente válidas, resulta más atractivo usar estimadores de varianzas remuestreados sobre el *MCPE*.

2.3. Tratamiento de datos faltantes

En esta sección se propone un estimador para la media poblacional cuando algunas observaciones de la variable de estudio o de las variables auxiliares están perdidas en la muestra. El nuevo estimador es valido para cualquier diseño muestral con probabilidades desiguales y está basado en el método de verosimilitud empírica. Este estimador se compara con otros conocidos estimadores en un estudio empírico.

2.3.1. Introducción

En la práctica, es común el uso de información auxiliar poblacional en la etapa de estimación. Esta técnica tiene muchas ventajas. Por ejemplo, una

adecuada información auxiliar puede producir una reducción considerable en el sesgo y el error muestral.

Cuando una o más variables auxiliares correlacionadas con la variable de estudio están disponibles, el método de calibración (Huang y Fuller, 1978, Deville y Särndal, 1992) y el método de verosimilitud pseudo empírica (Chen y Qin, 1993, Chen y Sitter, 1999, Wu y Sitter, 2001a, Wu, 2002) pueden usarse para estimar el total poblacional, la media poblacional, funciones de distribución y cuantiles. Ambos métodos usan información auxiliar de una o más variables auxiliares.

Generalmente, estas técnicas proporcionan estimadores que son más eficientes que los estimadores tradicionales, tales como el estimador de Horvitz y Thompson (1952) y el estimador tipo Hájek para la media (Rao, 1966, Basu, 1971, Särndal *et al.*, 1992). Sin embargo, el método de verosimilitud empírica asume respuesta completa sin valores perdidos, esto es, se asume que ninguna unidad muestral falla para proporcionar información en las variables de estudio y auxiliares.

La pérdida de información es una propiedad común en las investigaciones por muestreo. Esta pérdida de información puede ocurrir por varias razones. Los individuos muestreados pueden negarse a participar en el estudio, los entrevistadores no puedan contactar con los individuos del estudio, pérdida accidental de información, etc. En esta sección, se asume que si hay falta de respuesta, ésta es uniforme. Tratar con datos faltantes en una investigación por muestreo no es un asunto relativamente sencillo. Existen una gran variedad de métodos para usar en el caso de existir valores perdidos en los datos muestrales.

Ante la presencia de datos faltantes, la solución más simple es eliminar las unidades con falta de respuesta y aplicar el método de verosimilitud empírica a las unidades restantes. Sin embargo, este método, el cual Rubin (1987) llamó análisis de casos completos, puede producir sesgo en las estimaciones y varianzas muestrales más grandes (ver Rubin, 1987 o Little y Rubin, 1987).

La imputación es otra técnica que puede usarse en los individuos con falta de respuesta (Little y Rubin, 1987, Rao y Toutenburg, 1995, Särndal, 1992). La imputación consiste en sustituir los valores perdidos por un valor adecuado. Tratar los valores imputados como si estos fueran valores verdaderos y posteriormente usar el método de verosimilitud empírica puede dirigir a inferencias no válidas. Por ejemplo, la varianza puede resultar seriamente subestimada cuando la proporción de valores perdidos no es pequeña (Rao y Shao, 1992,

Särndal, 1990, 1992). Además, en algunas encuestas realizadas por organismos oficiales de estadística (como por ejemplo en la Oficina de Estadística de Suecia) está prohibida la imputación como solución al problema de datos faltantes.

Otra opción es intentar mejorar la precisión de las estimaciones incluyendo los valores observados de la variable auxiliar donde la variable de estudio está perdida. Así, aunque se tenga un valor perdido para y , el valor de x es observado y utilizado en el proceso de estimación.

Los estimadores de tipo razón, diferencia o producto también asumen respuesta completa. Algunos autores han definido estimadores de tipo razón en presencia de datos faltantes. Estos estimadores solamente han sido definidos para una clase limitada de diseños muestrales. Por ejemplo, Tracy y Osahan (1994), Toutenburg y Srivastava (1998, 1999, 2000) desarrollaron estimadores de tipo razón para muestreo aleatorio simple sin reemplazamiento.

En esta sección se propone modificar el estimador de verosimilitud pseudo empírica (*PEMLE*), el cual puede obtenerse bajo cualquier diseño muestral con probabilidades iguales o desiguales. El estimador propuesto usa toda la información muestral recogida para la variable de estudio y una variable auxiliar x , esto es, el estimador propuesto es función de los valores de x para las unidades con datos y perdidos, y función de los valores de y para las unidades con valores x perdidos.

Se considera la situación en la cual para algunos individuos existen observaciones perdidas en una de las características, pero no en la otra, es decir, la pérdida de información se produce para ambas características separadamente, pero no simultáneamente. De este modo, sea p ($p \geq 0$) el número de unidades que responden a x pero no a y , es decir, asumimos que tenemos p datos perdidos para la variable y . También se tiene información auxiliar incompleta, esto es, q ($q \geq 0$) unidades muestrales responden a y pero no a x . Notamos que p y q son números enteros. Así, se tiene un conjunto de $n - p - q$ unidades ($p + q \leq n$) que responden a ambas variables y y x . Con este esquema, los datos muestrales presentan la siguiente estructura

y_1	\dots	y_{n-p-q}	perdido	\dots	perdido	y_{n-q+1}	\dots	y_n
x_1	\dots	x_{n-p-q}	$x_{n-p-q+1}$	\dots	x_{n-q}	perdido	\dots	perdido

Sean los tres siguientes conjuntos disjuntos de unidades muestrales

$$\begin{aligned} s_A &= \{i \in s \mid x_i, y_i \text{ no están perdidos}\}, \\ s_B &= \{i \in s \mid x_i \text{ no está perdido, } y_i \text{ está perdido}\}, \\ s_C &= \{i \in s \mid y_i \text{ no está perdido, } x_i \text{ está perdido}\}. \end{aligned}$$

Asumiendo muestreo aleatorio simple sin reemplazamiento, Toutenburg y Srivastava (2000) propusieron cuatro estimadores para la media poblacional de y :

$$\bar{y}_{T1} = \bar{y}^A \left[\frac{(n-p-q)\bar{x}^A + p\bar{x}^B}{(n-q)\bar{x}^A} \right], \quad (2.48)$$

$$\bar{y}_{T2} = \bar{y}^A \left[\frac{(n-q)\bar{x}^A}{(n-p-q)\bar{x}^A + p\bar{x}^B} \right], \quad (2.49)$$

$$\bar{y}_{T3} = \left[\frac{((n-p-q)\bar{x}^A + p\bar{x}^B)((n-p-q)\bar{y}^A + q\bar{y}^C)}{(n-q)(n-p)\bar{x}^A} \right], \quad (2.50)$$

$$\bar{y}_{T4} = \left[\frac{(n-p-q)\bar{y}^A + q\bar{y}^C}{(n-p-q)\bar{x}^A + p\bar{x}^B} \right] \left[\frac{n-q}{n-p}\bar{x}^A \right], \quad (2.51)$$

donde \bar{y}^i y \bar{x}^i son las medias muestrales basadas en s_i , con $i = A, B, C$.

Los estimadores \bar{y}_{T1} y \bar{y}_{T2} dependen de las muestras s_A y s_B , y no dependen de la muestra s_C . Sin embargo, \bar{y}_{T3} y \bar{y}_{T4} dependen de las muestras s_A , s_B y s_C . Toutenburg y Srivastava (2000) demostraron que ninguno de estos estimadores es uniformemente superior a otro. Una elección apropiada del estimador requiere el conocimiento de parámetros poblacionales.

Rueda y González (2004) propusieron varios estimadores que pueden usarse bajo cualquier diseño muestral en la presencia de datos faltantes. Estos estimadores están basados en métodos razón, diferencia y regresión. Por ejemplo, el estimador siguiente es asintóticamente insesgado, bajo muestreo aleatorio simple es asintóticamente normal y es mejor, en el sentido de error cuadrático medio, que el resto de estimadores propuestos.

$$\bar{y}_{Reg} = \hat{\alpha}_{reg}\bar{y}_{HT}^A + (1-\hat{\alpha}_{reg})\bar{y}_{HT}^C + \frac{\widehat{Cov}_{i \in s_A}(x, y)}{\widehat{Var}_{i \in s_A}(x)} \left[\bar{X} - \left(\hat{\beta}_{reg}\bar{x}_{HT}^A + (1-\hat{\beta}_{reg})\bar{x}_{HT}^B \right) \right], \quad (2.52)$$

donde \bar{y}_{HT}^i y \bar{x}_{HT}^i son los estimadores de Horvitz-Thompson (1952) basados en s_i ($i = A, B, C$), $\widehat{Cov}_{i \in s_A}(x, y)$ y $\widehat{Var}_{i \in s_A}(x)$ denotan los estimadores de la covarianza y varianza basados en s_A . Véase Rueda y González (2004) para consultar los valores óptimos $\hat{\alpha}_{reg}$ y $\hat{\beta}_{reg}$.

2.3.2. Estimador propuesto

Sean los siguientes estimadores de tipo Hájek. Propiedades de este estimador están descritos en Rao (1966), Basu (1971) y Särndal *et al.* (1992).

$$\bar{y}_w^A = \sum_{i \in s_A} d_i^{A*} y_i \quad ; \quad \bar{y}_w^C = \sum_{i \in s_C} d_i^{C*} y_i \quad ; \quad \bar{y}_w^{AC} = \sum_{i \in s_A \cup s_C} d_i^{AC*} y_i; \quad (2.53)$$

$$\bar{x}_w^A = \sum_{i \in s_A} d_i^{A*} x_i \quad ; \quad \bar{x}_w^B = \sum_{i \in s_B} d_i^{B*} x_i \quad ; \quad \bar{x}_w^{AB} = \sum_{i \in s_A \cup s_B} d_i^{AB*} x_i; \quad (2.54)$$

con

$$d_i^{A*} = \frac{d_i^A}{\sum_{j \in s_1} d_j^A} \quad , \quad d_i^{B*} = \frac{d_i^B}{\sum_{j \in s_B} d_j^B} \quad , \quad d_i^{C*} = \frac{d_i^C}{\sum_{j \in s_C} d_j^C}, \quad (2.55)$$

$$d_i^{AB*} = \frac{d_i^{AB}}{\sum_{j \in s_A \cup s_B} d_j^{AB}} \quad , \quad d_i^{AC*} = \frac{d_i^{AC}}{\sum_{j \in s_A \cup s_C} d_j^{AC}}, \quad (2.56)$$

$$d_i^A = 1/\pi_i^A, \quad d_i^B = 1/\pi_i^B, \quad d_i^C = 1/\pi_i^C, \quad d_i^{AB} = 1/\pi_i^{AB}, \quad d_i^{AC} = 1/\pi_i^{AC}. \quad (2.57)$$

Las cantidades π_i^A , π_i^B , π_i^C , π_i^{AB} y π_i^{AC} son, respectivamente, las probabilidades de inclusión de primer orden de las muestras s_A , s_B , s_C , $s_A \cup s_B$ y $s_A \cup s_C$.

Cuando $u_i = 0$ (sin usar información auxiliar), se obtiene $\hat{p}_i = d_i^*$ y el estimador de verosimilitud pseudo empírico (*PEMLE*) coincide con el estimador de tipo Hájek dado por $\sum_{i \in s} d_i^* y_i$. Este estimador no usa la variable auxiliar x .

Sea el *PEMLE* de \bar{Y} dado por

$$\bar{y}_{PE}^A = \sum_{i \in s_A} \hat{p}_i^A y_i,$$

donde \hat{p}_i^A maximiza $l(p^A) = \sum_{i \in s_A} d_i^A \log p_i^A$ sujeta a

$$\sum_{i \in s_A} p_i^A = 1 \quad (0 \leq p_i^A \leq 1), \quad (2.58)$$

$$\sum_{i \in s_A} p_i^A u_i = 0. \quad (2.59)$$

Considerando el método de multiplicadores de Lagrange, \hat{p}_i^A está dado por

$$\hat{p}_i^A = \frac{d_i^{A*}}{1 + \lambda^A u_i}, \quad \text{for } i \in s_A, \quad (2.60)$$

donde el vector de multiplicadores de Lagrange, λ^A , se obtiene de la ecuación

$$\sum_{i \in s_A} \frac{d_i^{A*} u_i}{1 + \lambda^A u_i} = 0. \quad (2.61)$$

El estimador \bar{y}_{PE}^A no usa la información de la muestra s_B y s_C . A continuación se define un *PEMLE* que considera la información de s_A y s_B . Como la variable de interés contiene $n - p - q$ valores, el nuevo vector de pesos \hat{p}_i^{AB} debe definirse con dimensión $n - p - q$. Así, el nuevo estimador está dado por

$$\bar{y}_{PE}^{AB} = \sum_{i \in s_A} \hat{p}_i^{AB} y_i,$$

donde \hat{p}_i^{AB} ($i \in s_A$) se obtienen como \hat{p}_i^A (el cual tiene dimensión $n - p - q$), aunque en este caso se usa el vector de multiplicadores de Lagrange λ^{AB} , el cual está basado en las muestras s_A y s_B , en la expresión (2.60). λ^{AB} se obtiene de (2.61) después de sustituir d_i por d_i^{AB} en (2.61).

Pueden usarse otros métodos como el de imputación para obtener el *PEMLE* basado en las muestras s_A y s_B , aunque estos no están relacionados con el método de verosimilitud empírica.

Aunque \bar{y}_{PE}^{AB} parece mejor estimador que \bar{y}_{PE}^A al usar información de las muestras s_A y s_B , este estimador no resulta apropiado porque las condiciones $\sum_{i \in s_A} \hat{p}_i^{AB} = 1$ y $\sum_{i \in s_A} \hat{p}_i^{AB} u_i = 0$ no se cumplen. El estimador no queda bien construido y las ventajosas propiedades del método de verosimilitud empírica no sostienen. En el estudio empírica de la Sección 2.3.4 se confirma este comentario.

Desafortunadamente, el estimador propuesto \bar{y}_{PE}^A no usa información de la variable de estudio y proporcionada por la muestra s_C . Para resolver este problema, se propone una clase de estimadores que usan toda la información de la variable y incluida en las muestras s_A y s_C (véase también Rueda *et al.*, 2006b). Esta clase viene dada por

$$\bar{y}_{PE\alpha} = \alpha \bar{y}_{PE}^A + (1 - \alpha) \bar{y}_w^C, \quad (2.62)$$

donde α es una constante debidamente escogida que verifica $0 < \alpha < 1$. En la Sección 2.3.3, se proponen valores apropiados para α . El estimador \bar{y}_w^C está definido en (2.53).

Se observa que si $\alpha = 1$, el estimador resultante es \bar{y}_{PE}^A , y por tanto, este estimador está incluido en la clase $\bar{y}_{PE\alpha}$.

Cualquier estimador de esta clase usa toda la información disponible de las muestras s_A y s_C sin usar técnicas de imputación. Los valores de x de la muestra s_B no se usan para la estimación. No obstante, los valores de la variable y están perdidos para $i \in s_B$. Incluir esta información en la clase considerando \bar{y}_{PE}^{AB} en lugar de \bar{y}_{PE}^A empeoraría las estimaciones. En la Sección 2.3.4, un estudio de simulación muestra que los estimadores de la clase propuesta son tan eficientes como otros estimadores que usan información de cada muestra s_A , s_B y s_C .

2.3.3. Propiedades teóricas

En esta sección se demuestra que el estimador $\bar{y}_{PE\alpha}$ propuesto en (2.62) es asintóticamente insesgado. La varianza asintótica de $\bar{y}_{PE\alpha}$ también se deriva.

Sean las siguientes condiciones.

$$(C2.15). \quad u^{A*} = \max_{i \in s_A} |u_i| = o_p(n^{1/2}).$$

$$(C2.16). \quad \frac{\sum_{i \in s_A} d_i^A u_i}{\sum_{i \in s_A} d_i^A u_i^2} = O_p(n^{1/2}).$$

Estas condiciones fueron usadas por Chen y Sitter (1999), los cuales demuestran que varios diseños muestrales más comunes las satisfacen. Dadas estas condiciones, el siguiente resultado puede obtenerse.

Corolario 2.3 *Bajo las condiciones (C2.15) y (C2.16), se tiene que*

$$\bar{y}_{PE\alpha} = \alpha \bar{y}_{GREG}^A + (1 - \alpha) \bar{y}_w^C + o_p(n^{-1/2}) \quad (2.63)$$

donde

$$\bar{y}_{GREG}^A = \bar{y}_w^A + (\bar{X} - \bar{x}_w^A) b, \quad (2.64)$$

con

$$b = \frac{\sum_{i \in s_A} d_i^{A*} x_i y_i - \bar{y}_w^A \bar{x}_w^A}{\sum_{i \in s_A} d_i^{A*} (x_i - \bar{x}_w^A)^2}. \quad (2.65)$$

Demostración

Chen y Sitter (1999) demostraron que \bar{y}_{PE}^A es asintóticamente equivalente a \bar{y}_{GREG}^A . Sabido esto, este resultado se sigue fácilmente. \square

Teorema 2.5 *Bajo las condiciones (C2.15) y (C2.16), se tiene que*

$$\bar{y}_{GREG}^A \simeq \bar{y}_{GREG}^{A2},$$

donde

$$\bar{y}_{GREG}^{A2} = \bar{y}_w^A + (\bar{X} - \bar{x}_w^A)B, \quad (2.66)$$

con

$$B = \frac{Cov(x, y)}{Var(x)}. \quad (2.67)$$

Demostración

Para establecer este resultado, se asume que la población finita envuelve una secuencia de poblaciones donde n y N aumentan de modo que $n/N \rightarrow f$ cuando $n \rightarrow \infty$ y donde f es una constante.

Randles (1982) estudió el comportamiento asintótico de algunas familias comunes de estadísticos podía establecerse aunque algunos parámetros vitales en la formulación del estadístico fuesen desconocidos. Este autor demostró que si $T_n(\hat{\lambda})$ es una función de datos que usa el estimador $\hat{\lambda}$, el cual también es una función de los datos que estima consistentemente el parámetro λ , entonces $T_n(\hat{\lambda})$ y $T_n(\lambda)$ tienen la misma distribución límite y se verifica

$$\left. \frac{\partial \mu(\gamma)}{\partial \gamma} \right|_{\gamma=\lambda} = 0,$$

donde $\mu(\gamma) = \lim_{n \rightarrow +\infty} E_\lambda[T_n(\gamma)]$ y la esperanza es considerada cuando el verdadero parámetro es λ .

Sea $T_n(\gamma) = \bar{y}_w^A + (\bar{X} - \bar{x}_w^A)\gamma$. Notamos que $T_n(b) = \bar{y}_{GREG}^A$ ha sido establecido en (2.64). Consideremos $\mu(\gamma) = \lim_{n \rightarrow \infty} E_\gamma[T_n(\gamma)]$. Notamos que cuando $\gamma = B$, el cual está definido en (2.67), se obtiene $\mu(B) = \tilde{Y}$ donde $\tilde{Y} = \lim_{n \rightarrow \infty} \bar{Y}$. Puesto $\mu(\gamma)$ verifica

$$\left. \frac{\partial \mu(\gamma)}{\partial \gamma} \right|_{\gamma=B} = 0,$$

esto implica que $\bar{y}_{GREG}^A \simeq \bar{y}_{GREG}^{A2}$. Esto completa la demostración. \square

Usando el corolario 2.3 y el teorema 2.6 se obtiene

$$\bar{y}_{PE\alpha} \simeq \alpha \bar{y}_{GREG}^{A2} + (1 - \alpha) \bar{y}_w^C, \quad (2.68)$$

el cual implica que $\bar{y}_{PE\alpha}$ es asintóticamente insesgado.

Teorema 2.6 *Bajo las condiciones (C2.15) y (C2.16), la varianza asintótica de $\bar{y}_{PE\alpha}$ está dada por*

$$\begin{aligned} AV(\bar{y}_{PE\alpha}) &= \alpha^2 \left[V(\bar{y}_w^A) + B^2 V(\bar{x}_w^A) - 2BCov(\bar{y}_w^A, \bar{x}_w^A) \right] + \\ &+ (1 - \alpha)^2 V(\bar{y}_w^C) + 2\alpha(1 - \alpha) \left[Cov(\bar{y}_w^A, \bar{y}_w^C) - BCov(\bar{x}_w^A, \bar{y}_w^C) \right]. \end{aligned} \quad (2.69)$$

Demostración

La aproximación (2.68) implica que la varianza asintótica de $\bar{y}_{PE\alpha}$ está dada por

$$V\left(\alpha \bar{y}_{GREG}^{A2} + (1 - \alpha) \bar{y}_w^C\right) = \alpha^2 V(\bar{y}_{GREG}^{A2}) + (1 - \alpha)^2 V(\bar{y}_w^C) + 2\alpha(1 - \alpha) Cov(\bar{y}_{GREG}^{A2}, \bar{y}_w^C). \quad (2.70)$$

Usando (2.66), la varianza de \bar{y}_{GREG}^{A2} es

$$\begin{aligned} V(\bar{y}_{GREG}^{A2}) &= V\left(\bar{y}_w^A + (\bar{X} - \bar{x}_w^A)B\right) \\ &= V\left(\bar{y}_w^A - \bar{x}_w^A B\right) \\ &= V(\bar{y}_w^A) + B^2 V(\bar{x}_w^A) - 2BCov(\bar{y}_w^A, \bar{x}_w^A). \end{aligned} \quad (2.71)$$

El valor $Cov(\bar{y}_{GREG}^{A2}, \bar{y}_w^C)$ es

$$Cov(\bar{y}_{GREG}^{A2}, \bar{y}_w^C) = Cov(\bar{y}_w^A, \bar{y}_w^C) - BCov(\bar{x}_w^A, \bar{y}_w^C). \quad (2.72)$$

Así de (2.70), (2.71) y (3.32), la varianza asintótica de $\bar{y}_{PE\alpha}$ está dada por (2.69). El Teorema 2.6 se sigue fácilmente. \square

El estimador óptimo de la clase propuesta está dado por el estimador definido en (2.62) con un valor α que minimize la varianza asintótica dada por (2.69).

La varianza asintótica (2.69) puede expresarse como

$$AV(\bar{y}_{PE\alpha}) = \alpha^2 M^* + (1 - \alpha)^2 N^* + 2\alpha(1 - \alpha) L^*,$$

donde

$$M^* = V(\bar{y}_w^A) + B^2 V(\bar{x}_w^A) - 2BCov(\bar{y}_w^A, \bar{x}_w^A), \quad (2.73)$$

$$N^* = V(\bar{y}_w^C), \quad (2.74)$$

$$L^* = Cov(\bar{y}_w^A, \bar{y}_w^C) - BCov(\bar{x}_w^A, \bar{y}_w^C). \quad (2.75)$$

El valor α_{opt} que minimiza la varianza asintótica es solución de la ecuación

$$\left. \frac{\partial AV(\bar{y}_{PE\alpha})}{\partial \alpha} \right|_{\alpha=\alpha_{opt}} = 2\alpha_{opt}M^* - 2(1 - \alpha_{opt})N^* + 2(1 - 2\alpha_{opt})L^* = 0,$$

la cual implica

$$\alpha_{opt} = \frac{N^* - L^*}{M^* + N^* - 2L^*}. \quad (2.76)$$

Substituyendo α_{opt} en (2.69), se obtiene la varianza asintótica más pequeña dada por

$$AV(\bar{y}_{PE\alpha_{opt}}) = \alpha_{opt}^2 M^* + (1 - \alpha_{opt})^2 N^* + 2\alpha_{opt}(1 - \alpha_{opt})L^*. \quad (2.77)$$

Desafortunadamente, el valor óptimo α_{opt} depende de parámetros poblacionales desconocidos, los cuales pueden estimarse a partir de los datos muestrales.

Bajo muestreo aleatorio simple y muestreo estratificado, $\sum_{i \in s} d_i = N$, esto es, el estimador de Horvitz-Thompson y el estimador de tipo Hájek son idénticos, y por tanto, los estimadores de las varianzas y covarianzas de las expresiones (2.73), (2.74) y (2.75) pueden obtenerse fácilmente. Una expresión analítica para (2.73), (2.74) y (2.75) bajo muestreo aleatorio simple puede encontrarse en Rueda y González (2004).

Con estas estimaciones, puede obtenerse una aproximación $\tilde{\alpha}_{opt}$ de α_{opt} . Por lo tanto, la expresión del estimador propuesto viene dada por

$$\tilde{y}_{PE\alpha_{opt}} = \tilde{\alpha}_{opt} \bar{y}_{PE}^A + (1 - \tilde{\alpha}_{opt}) \bar{y}_w^C. \quad (2.78)$$

También es posible establecer la insesgadez asintótica de $\tilde{y}_{PE\alpha_{opt}}$.

2.3.4. Propiedades empíricas

En esta sección se comparan los estimadores propuestos con otros estimadores alternativos usando un estudio empírico basado en poblaciones reales y simuladas, usadas previamente en estudios de estimadores de regresión y razón, estimación de la varianza e intervalos de confianza.

Las poblaciones naturales usadas en este estudio son la Fam1500 y Hospitals (véase Apéndice A). Se recuerda que los coeficientes de correlación están dados

por $\rho_{y,x_1} = 0,848$ y $\rho_{y,x_2} = 0,546$ en la población Fam1500 y $\rho_{y,x} = 0,911$ en la población Hospitals.

Paralelamente a Wu y Sitter (2001a), se han generado cuatro poblaciones de $N = 2000$ unidades mediante muestras independientes e idénticamente distribuidas mediante el modelo

$$y = \theta_0 + \theta_1 x + \epsilon, \quad (2.79)$$

donde $x \sim \text{Gamma}(1, 1)$, $\epsilon \sim N(0, \sigma^2)$ y $\theta_0 = \theta_1 = 1$. Los coeficientes de correlación están dados por 0.6, 0.7, 0.8 y 0.9, y las poblaciones se llaman Pop06, Pop07, Pop08 y Pop09, respectivamente. Ver más detalles de estas poblaciones en Apéndice A.

Se analiza la precisión de los estimadores propuestos por medio de un estudio empírico donde para cada población se han representado tres números diferentes de valores perdidos para la variable x , p . Varios valores perdidos de y , q , se han representado en el eje de abscisas. De este modo, el comportamiento de los estimadores puede observarse para relaciones fuertes y débiles entre variables y diferentes situaciones de datos perdidos.

El comportamiento de los estimadores \bar{y}_{PE}^A y $\tilde{\bar{y}}_{PE\alpha_{opt}}$ se compara con los siguientes estimadores: (i) el estimador estándar de tipo Hájek para la media poblacional basado en las muestras s_A y s_C , es decir, \bar{y}_w^{AC} ; (ii) \bar{y}_{T1} , \bar{y}_{T2} , \bar{y}_{T3} y \bar{y}_{T4} , los estimadores propuestos en Toutenburg y Srivastava (2000); (iii) \bar{y}_{PE}^{AB} , el *PEMLE* basado en las muestras s_A y s_B . Aunque se ha señalado que los pesos no quedan bien definidos, se usa en el estudio de simulación para observar su comportamiento; (iv) \bar{y}_{Reg} , el estimador propuesto en Rueda y González (2004) basado en las muestras s_A , s_B y s_C .

Para cada una de las seis poblaciones, se han generado $B = 1000$ muestras independientes bajo muestreo aleatorio simple con tamaño muestral n . A continuación, se eliminan de la muestra p elementos de la variable auxiliar y q elementos de la variable de estudio de forma aleatoria. Bajo este escenario, las submuestras s_A , s_B y s_C pueden definirse fácilmente. El cumplimiento de todos los estimadores se mide en términos de Sesgo Relativo (*SR*) y de Eficiencia Relativa (*ER*), donde

$$SR_j = \frac{1}{B} \sum_{b=1}^B \frac{|\hat{\bar{y}}_j(b) - \bar{y}|}{\bar{Y}} \quad ; \quad ER_j = \frac{ECM(\hat{\bar{y}}_j)}{ECM(\bar{y}_w^{AC})},$$

b indica la b -ésima simulación, el Error Cuadrático Medio empírico está dado por $ECM(\hat{\bar{y}}_j) = B^{-1} \sum_{b=1}^B (\hat{\bar{y}}_j(b) - \bar{Y})^2$, y $j = 1, \dots, 8$ se refiere a los estimadores \bar{y}_{PE}^A , \bar{y}_{PE}^{AB} , $\tilde{\bar{y}}_{PE\alpha_{opt}}$, \bar{y}_{Reg} , \bar{y}_{T1} , \bar{y}_{T2} , \bar{y}_{T3} y \bar{y}_{T4} .

Figura 2.1: Eficiencia Relativa para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se toman muestras de tamaño $n = 200$.

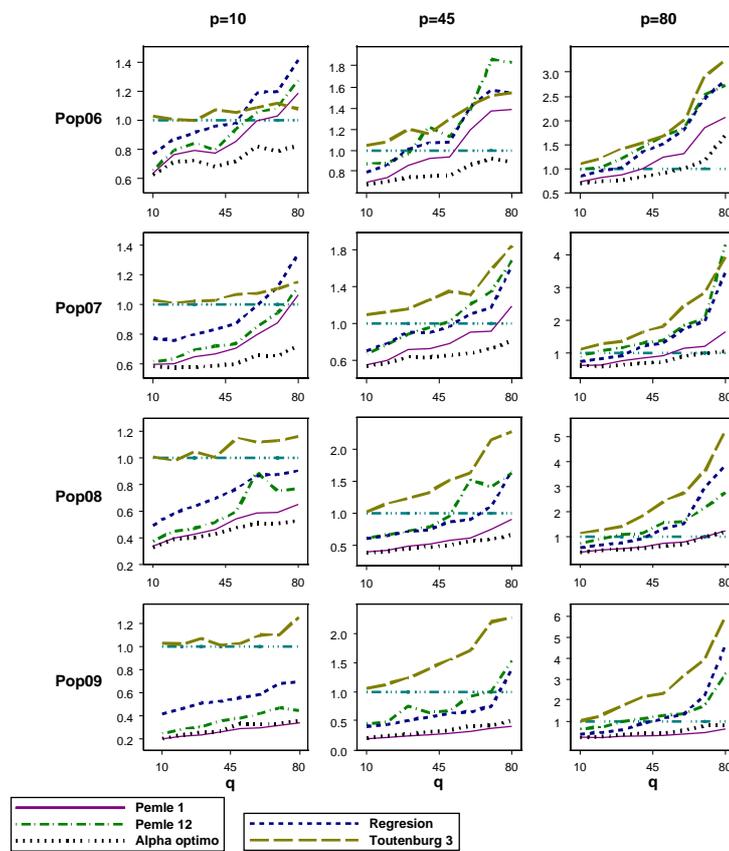


Figura 2.2: Eficiencia Relativa para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_{Reg} (Regresión) y \bar{y}_{T_3} (Toutenburg 3). Se considera la población Fam1500 y muestras de tamaño $n = 150$.

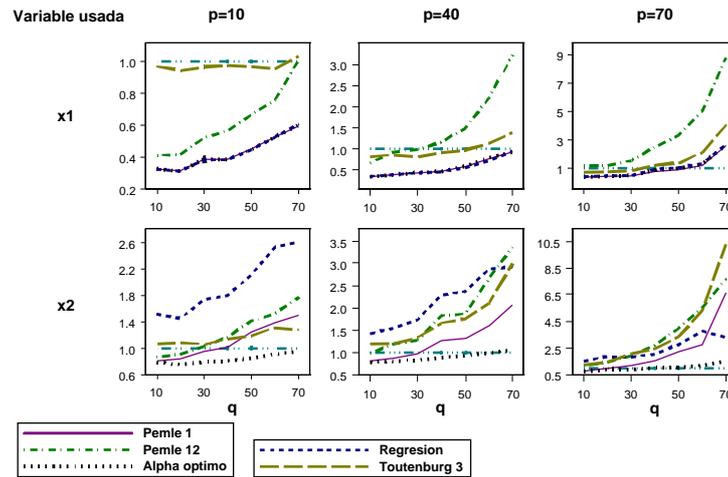
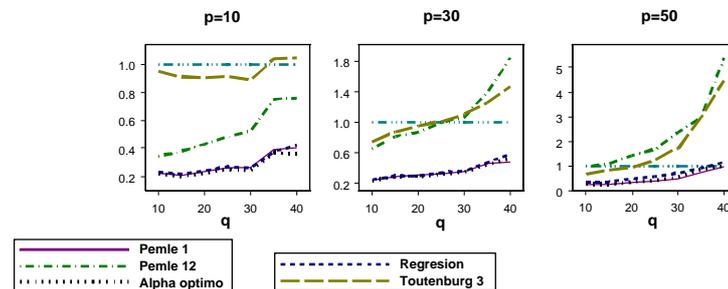


Figura 2.3: Eficiencia Relativa para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_{Reg} (Regresión) y \bar{y}_{T_3} (Toutenburg 3). Se considera la población Hospitals y muestras de tamaño $n = 100$.



Las simulaciones se han llevado a cabo en R y los códigos se encuentran en el Apéndice C.

En primer lugar, se observa que el estimador \bar{y}_{T3} posee una considerable ganancia en precisión respecto a los estimadores \bar{y}_{T1} , \bar{y}_{T2} y \bar{y}_{T4} . Con el fin de obtener más claridad en las figuras, las líneas correspondientes a los estimadores \bar{y}_{T1} , \bar{y}_{T2} y \bar{y}_{T4} no se han incluido.

Las Figuras 2.1, 2.2 y 2.3 representan los valores de la Eficiencia Relativa (eje de ordenadas) para los estimadores \bar{y}_{PE}^A , \bar{y}_{PE}^{AB} , $\tilde{\bar{y}}_{PE\alpha_{opt}}$, \bar{y}_{Reg} y \bar{y}_{T3} bajo muestreo aleatorio simple y diferentes valores de p y q . Las líneas horizontales en el punto 1 representan la ER para \bar{y}_w^{AC} , el estimador estándar.

De estas figuras, se puede llegar a las siguientes conclusiones generales:

1. Si aumenta la relación entre y y x y, además, el número de datos faltantes es escaso, todos los estimadores (excepto \bar{y}_{T3}) obtienen mejores estimaciones con respecto al estimador estándar. Cuando ambos p y q incrementan, las estimaciones son peores con respecto a \bar{y}_w^{AC} , y de ahí, que todas las líneas sean crecientes.
2. Los mejores resultados se consiguen a través del estimador $\tilde{\bar{y}}_{PE\alpha_{opt}}$, esto es, el ECM es siempre menor que el resto de estimadores y siempre mejora considerablemente los resultados proporcionados por el estimador directo \bar{y}_w^{AC} .
3. El peor comportamiento lo muestra el estimador de Toutenburg y Srivastava (2000). Esto puede deberse al hecho de que este estimador no usa \bar{X} como información auxiliar.

Comparando entre los estimadores basados en el método de verosimilitud empírica, se observa

1. Los estimadores \bar{y}_{PE}^A y $\tilde{\bar{y}}_{PE\alpha_{opt}}$ son equivalentes cuando existe una fuerte relación entre y y x y el número de datos perdidos es pequeño. La ganancia en eficiencia de $\tilde{\bar{y}}_{PE\alpha_{opt}}$ con respecto a \bar{y}_{PE}^A es mayor en el caso contrario.
2. \bar{y}_{PE}^{AB} nunca es mejor que los estimadores \bar{y}_{PE}^A o $\tilde{\bar{y}}_{PE\alpha_{opt}}$ en términos de eficiencia. La razón para esto es que sus pesos no están bien definidos.

Un estimador que usa la información de s_A , s_B y s_C es \bar{y}_{Reg} . En las poblaciones Hospitals y Fam1500 (cuando se usa x_1), \bar{y}_{PE}^A , $\tilde{\bar{y}}_{PE\alpha_{opt}}$ y \bar{y}_{Reg} son equivalentes. En los otros casos, \bar{y}_{Reg} nunca mejora en eficiencia a $\tilde{\bar{y}}_{PE\alpha_{opt}}$. Aunque \bar{y}_{Reg} usa información de s_A , s_B y s_C , $\tilde{\bar{y}}_{PE\alpha_{opt}}$ es considerablemente más eficiente cuando la correlación entre y y x es baja y aumentan los valores de p y q .

Finalmente, si comparamos el estimador propuesto con el estimador estándar

1. \bar{y}_w^{AC} es únicamente más eficiente que $\tilde{\bar{y}}_{PE\alpha_{opt}}$ cuando la relación entre variables es débil y el número total de datos perdidos, $p + q$, es alto. En este caso, el resto de estimadores obtienen significativamente peores estimaciones. Esto ocurre, por ejemplo, en Pop06, $p = 80$, $q = 60$, esto es, el 70 % de la muestra son valores perdidos. En la práctica, esta situación es improbable o inaceptable. No obstante, este caso se muestra para poder revelar el comportamiento de los estimadores en situaciones extremas.
2. Como se esperaba, cuando el número de valores de x perdidos, p , incrementa, la ganancia en precisión del estimador propuesto con respecto a \bar{y}_w^{AC} es menor. Equivalentemente, cuando p permanece fijo, la ganancia en precisión decrece cuando el número de valores perdidos q aumenta. Este resultado es lógico porque si p/q es pequeño, se proporciona más información por la muestra s_C en relación con la muestra s_B , y \bar{y}_w^{AC} también usa la información de s_C .

Las Figuras 2.4, 2.5 y 2.6 muestran los valores del Sesgo Relativo (SR) para todos los estimadores. Puede observarse que los valores SR están todos en un rango razonable, teniendo los estimadores \bar{y}_{PE}^A y $\tilde{\bar{y}}_{PE\alpha_{opt}}$ el mejor comportamiento en términos de SR . Estas figuras presentan similares resultados que la ER , y por tanto, se puede llegar a las mismas conclusiones.

En resumen, estas simulaciones muestran como un apropiado uso de las muestras s_A y s_C por el estimador propuesto puede reducir el error de los estimadores directo, regresión, de verosimilitud pseudo empírica, etc. Por tanto, el estimador propuesto $\tilde{\bar{y}}_{PE\alpha_{opt}}$ es una óptima alternativa para la estimación de parámetros lineales en presencia de datos faltantes y con un buen uso de la información auxiliar.

Figura 2.4: Sesgo Relativo para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_w^{AC} (estándar), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se toman muestras de tamaño $n = 200$.

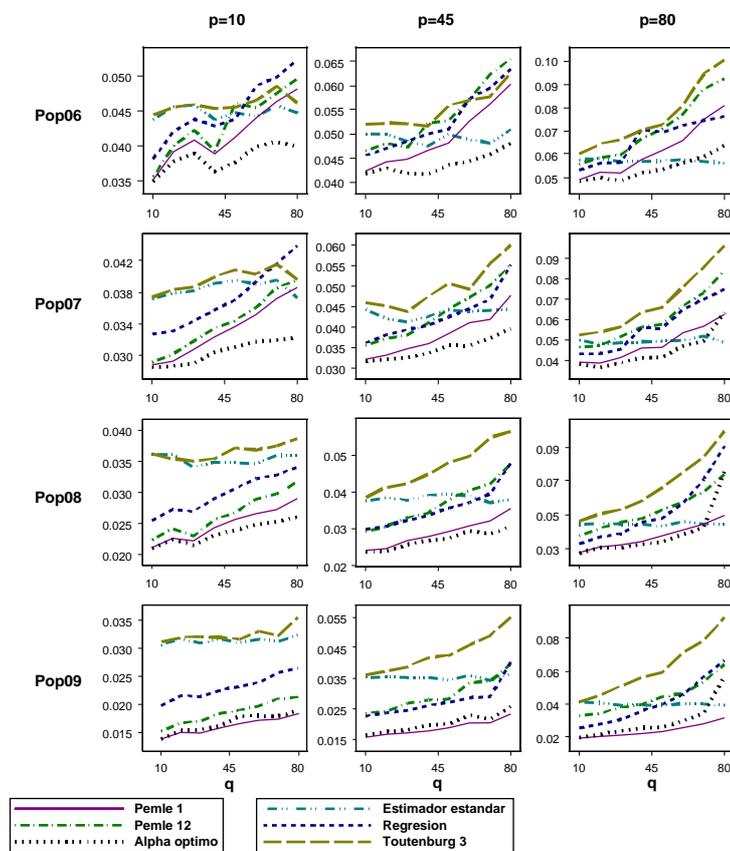


Figura 2.5: Sesgo Relativo para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{y}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_w^{AC} (estándar), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se considera la población Fam1500 y muestras de tamaño $n = 150$.

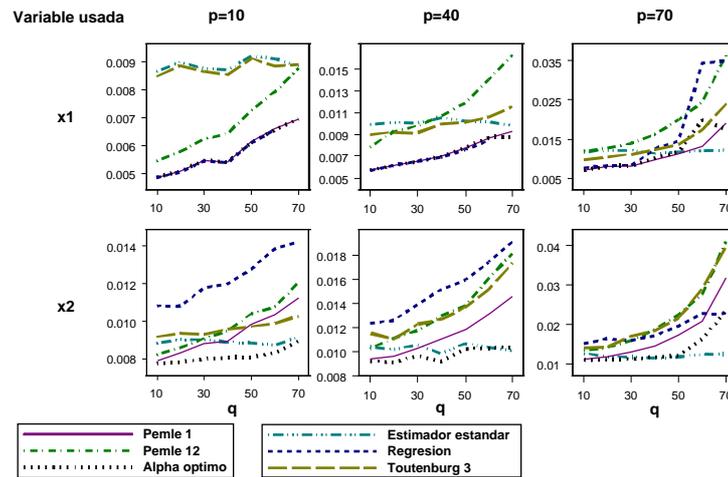
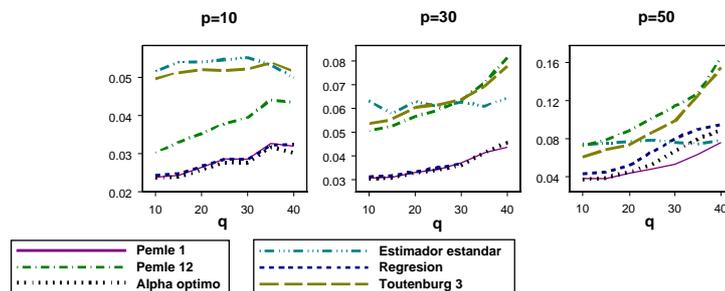


Figura 2.6: Sesgo Relativo para los estimadores \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{y}_{PE\alpha_{opt}}$ (Alpha óptimo), \bar{y}_w^{AC} (estándar), \bar{y}_{Reg} (Regresión) y \bar{y}_{T3} (Toutenburg 3). Se considera la población Hospitals y muestras de tamaño $n = 100$.



2.4. Estimación de la función de distribución

2.4.1. Introducción

El problema de la estimación de la función de distribución es un tema actual y muy importante del muestreo en poblaciones finitas, por tratarse de una función que permite determinar las características más importantes de la población en estudio, proporcionando información relevante acerca del comportamiento global de la población. Sin duda, los estimadores estudiados clásicamente en la teoría del muestreo, como totales, medias, proporciones y varianzas, no ofrecen tanta información como la función de distribución, aunque obtener estimadores eficientes para tal función no es tan simple como en el caso de los estimadores puntuales.

La estimación de cuantiles y de otros parámetros de tipo no funcional como el índice de Gini o la curva de Lorentz también queda resuelto con el conocimiento de la función de distribución. Los cuantiles, por ejemplo, pueden obtenerse mediante inversión directa de la función de distribución. Además, permite obtener medidas importantes como la determinación de las líneas de pobreza, proporción de bajos ingresos, etc. y son muy útiles en investigaciones de tipo social o económico. Debido a la importancia de estos parámetros en algunas investigaciones o estudios, se debe disponer de buenos métodos y técnicas para obtener las mejores estimaciones posibles.

Recordemos que la función de distribución para una variable de interés, y , y una población finita, U , es la proporción de unidades en U para las cuales el valor de y es menor o igual que t . El problema de la estimación de la función de distribución en la presencia de información auxiliar ha recibido recientemente mucha atención debido a las importantes propiedades que posee, el interés considerable que tiene cuando, por ejemplo, y es una medida de gastos o ingresos, etc.

La función de distribución poblacional,

$$F_y(t) = \frac{1}{N} \sum_{i=1}^N \delta(t - y_i), \quad (2.80)$$

satisface las siguientes condiciones:

$$(C2.17). \quad \lim_{t \rightarrow -\infty} F_y(t) = 0 \quad ; \quad \lim_{t \rightarrow +\infty} F_y(t) = 1.$$

$$(C2.18). \quad F_y(t) \text{ es monótona no-decreciente: } \forall t_1 < t_2, F_y(t_1) \leq F_y(t_2).$$

(C2.19). $F_y(t)$ es continua por la derecha: Dado $t > t^*$, $\lim_{t \rightarrow t^*} F_y(t) = F_y(t^*)$.

Varios de los estimadores propuestos en la literatura del muestreo en poblaciones finitas no satisfacen todas estas propiedades y no son, por tanto, funciones de distribución. Por ejemplo, la función de distribución estimada mediante el método de calibración no cumple los requisitos necesarios para ser una verdadera función de distribución.

Asumamos que la variable de estudio, y , está altamente asociada con un vector auxiliar de variables, $\mathbf{x}_i = (x_{i1}, \dots, x_{ij}, \dots, x_{iJ})^t$, donde los valores $\mathbf{x}_1, \dots, \mathbf{x}_N$ son conocidos para toda la población. Como se ha comentado en varias ocasiones, en las investigaciones por muestreo es común el uso de esta información poblacional auxiliar en la etapa de estimación para incrementar la precisión de los estimadores de una media o un total. Bajo este escenario, el uso de la información auxiliar a sido extensamente estudiado, pero bastante menos ha sido el esfuerzo por aplicarlo a la estimación de la función de distribución y cuantiles poblacionales. Notamos que la aplicación de las técnicas usuales para la estimación de medias y totales en el escenario de la estimación de la función de distribución producen resultados no deseables y, en general, con una pérdida significativa en eficiencia.

Por otro lado, el número de variables auxiliares a usar en la etapa de estimación es otro punto de vista interesante en la estimación de la función de distribución. Algunos de los estimadores en la literatura están contruidos para una única variable auxiliar, y el uso de otras variables auxiliares resulta imposible o con un alto coste computacional. Si estas variables presentan una fuerte relación con la variable de estudio, éstas deberían incluirse en el estudio y parece razonable asumir que podrían obtenerse mejores propiedades. Estos estimadores tienen la desventaja de la pérdida de eficiencia provocada por el hecho de no poder usar esta información auxiliar multivariante. Estas consideraciones sugieren que un uso más eficiente de la información auxiliar en la etapa de estimación es posible en el problema de la estimación de la función de distribución.

Sabemos que el método de verosimilitud pseudo empírica es una técnica reciente que puede usarse para la estimación de medias o totales poblacionales (Chen y Qin, 1993, Chen y Sitter, 1999), funciones de distribución (Chen y Wu, 2002, Wu, 2003) y otros parámetros. Asumiendo este método, Chen y Wu (2002) propusieron estimadores modelo-calibrados para estimar la función de distribución. Estos estimadores están contruidos por medio de restricciones que requieren el uso de un valor fijado t_0 . Estos estimadores sufren una considerable pérdida de eficiencia cuando t_0 se encuentra alejado de t , el punto

donde se evalúa la función de distribución. El estimador propuesto en la sección siguiente emplea el método de verosimilitud empírica y permite el uso de información auxiliar multivariante. Este estimador está basado en una aproximación modelo-asistida. Además, se usa un conjunto apropiado de puntos en las restricciones para evitar el problema de la pérdida de eficiencia.

2.4.2. Algunos estimadores de la función de distribución

En este apartado se describen los principales trabajos y enfoques relacionados con la estimación de la función de distribución poblacional. Destacamos las propiedades más importantes de estos estimadores, prestando especial interés en los estimadores modelo-calibrados de verosimilitud empírica. Estos últimos presentan bastantes similitudes con el estimador propuesto en la Sección 2.4.3, por lo que señalaremos las principales diferencias entre unos y otros. Todos los estimadores que se exponen a continuación están basados en distintas aproximaciones. Aprovecharemos la ocasión para describir los tipos de inferencias que existen recientemente en muestreo de poblaciones finitas (véase también Apéndice B).

En la expresión (2.80) se observa que la función de distribución puede verse como una media poblacional de la variable $z_i = \delta(t - y_i)$, y por tanto, sin utilizar ningún tipo de información auxiliar, la estimación de la función de distribución es un caso especial de la estimación de la media poblacional. Haciendo uso de esta perspectiva, los estimadores más conocidos son el de Horvitz y Thompson (1952), dado por

$$\widehat{F}_{HTy}(t) = \frac{1}{N} \sum_{i \in s} d_i \delta(t - y_i),$$

y el estimador de tipo Hájek dado por

$$\widehat{F}_{HKy}(t) = \frac{\sum_{i \in s} d_i \delta(t - y_i)}{\sum_{j \in s} d_j} = \sum_{i \in s} d_i^* \delta(t - y_i),$$

donde $d_i = 1/\pi_i$. Nótese que el estimador de Horvitz y Thompson puede usarse únicamente cuando el tamaño poblacional es conocido, mientras que el de tipo Hájek puede emplearse en ambas situaciones. Bajo cualquier diseño muestral en el cual $\sum_{i \in s} d_i = N$, puede demostrarse que $\widehat{F}_{HTy}(t) = \widehat{F}_{HKy}(t)$.

En presencia de información auxiliar, Rao *et al.* (1990) propusieron dos nuevos estimadores basados en el diseño muestral: el estimador de tipo razón

dado por

$$\widehat{F}_r(t) = \frac{1}{N} \frac{\sum_{i \in s} d_i \delta(t - y_i)}{\sum_{i \in s} d_i \delta(t - \widehat{R}x_i)} \sum_{i \in U} \delta(t - \widehat{R}x_i), \quad (2.81)$$

y el estimador diferencia dado por

$$\widehat{F}_d(t) = \frac{1}{N} \left\{ \sum_{i \in s} d_i \delta(t - y_i) + \sum_{i \in U} \delta(t - \widehat{R}x_i) - \sum_{i \in s} d_i \delta(t - \widehat{R}x_i) \right\}, \quad (2.82)$$

donde $\widehat{R} = \frac{\sum_{i \in s} d_i y_i}{\sum_{i \in s} d_i x_i}$. Se observa que ambos estimadores utilizan como información auxiliar la variable $\delta(t - \widehat{R}x)$.

Al no utilizar ningún tipo de información auxiliar, los estimadores $\widehat{F}_{HTy}(t)$ y $\widehat{F}_{HKy}(t)$ son menos eficientes que $\widehat{F}_r(t)$ y $\widehat{F}_d(t)$, pero sin embargo, éstos últimos tienen el inconveniente de dar valores, por lo general, fuera del rango $[0, 1]$ y no siempre son funciones monótonas respecto a t , con lo que no cumplen las propiedades de la función de distribución. Por este motivo, son numerosos los casos en los que la inversión directa de $\widehat{F}_r(t)$ y $\widehat{F}_d(t)$ no produce buenas estimaciones para los cuantiles.

En Rao *et al.* (1990) y en Francisco y Fuller (1991) se propone transformar $\widehat{F}_d(t)$ en una función monótona antes de obtener estimaciones para los cuantiles. Estos procesos tienen básicamente dos inconvenientes: (i) no son transformaciones triviales y (ii) se desconoce la pérdida de eficiencia al realizar la transformación.

Otro estimador para la función de distribución bastante reciente es el obtenido mediante el método de calibración descrito en Deville y Särndal (1992). Al igual que los anteriores que utilizan información auxiliar tienen la propiedad no deseable de no ser una auténtica función de distribución. Esto se debe a que los pesos que se utilizan para ponderar las unidades muestrales de la variable de interés, $\delta(t - y_i)$, pueden ser negativos, y por tanto, el estimador resultante puede llegar a ser decreciente. Además se demuestra que su límite cuando $t \rightarrow +\infty$ es distinto de 1.

Por tanto, es deseable requerir que un estimador para la función de distribución sea por sí mismo una verdadera función de distribución. Nótese, que una verdadera función de distribución debe satisfacer las condiciones (C2.17), (C2.18) y (C2.19).

El conocido estimador de regresión generalizado (*GREG*) (Cassel *et al.*, 1976, Särndal, 1980) es un estimador modelo-asistido que está basado en un

modelo lineal. Más recientemente, son dos los principales métodos en la literatura que están categorizados como aproximaciones modelo-asistidas. Estos procedimientos son el de calibración (Deville y Särndal, 1992) y el de verosimilitud empírica (Chen y Qin, 1993, Chen y Sitter, 1999). Notamos que estos procedimientos no son dependientes de un modelo, aunque usan uno de ellos para construir el estimador. En otras palabras, los estimadores modelo-asistidos son aproximadamente (asintóticamente) insesgados bajo el diseño, independientemente de si el modelo es correcto o no, y son particularmente eficientes si el modelo en el que se basa es correcto. Así, la aproximación modelo-asistida proporciona inferencias válidas bajo el modelo asumido y al mismo tiempo, está protegido contra una mala especificación del modelo en el sentido de proporcionar inferencias válidas basadas en el diseño, independientemente de la relación de la variable de interés con la variable auxiliar. Un ejemplo de estimadores modelo-asistidos para la función de distribución son los estimadores $\hat{F}_r(t)$ y $\hat{F}_d(t)$.

Otro procedimiento para estimar parámetros lineales o no lineales en poblaciones finitas es la aproximación basada en modelos, la cual asume un modelo de superpoblación y donde los estimadores son dependientes del modelo. Chambers y Dunstan (1986) y Dorfman y Hall (1993) propusieron estimadores basados en modelos para la función de distribución. El estimador de Chambers y Dunstan presenta el inconveniente de ser inconsistente bajo el diseño. Además, se necesita llevar a cabo un cuidadoso contraste sobre el modelo antes de que estos estimadores sean usados. Todos estos métodos presentan un grado de dificultad en la computación y un pobre cumplimiento cuando el modelo especificado es incorrecto. Bajo muestreo aleatorio simple, Wang y Dorfman (1996) combinaron los estimadores de Chambers y Dunstan (1986) con estimadores de tipo diferencia basados en el diseño en un estimador híbrido, que bajo ciertas condiciones, es más eficiente que ambos estimadores. No obstante, este estimador hereda las desventajas de ambos estimadores y tiene una complicada generalización a diseños muestrales más complejos. Silva y Skinner (1995) llevaron a cabo un estudio exhaustivo de las propiedades del estimador, y destacaron algunos problemas importantes, como por ejemplo, la pérdida en eficiencia cuando este estimador se usa en la estimación de cuantiles.

Finalmente, la recientemente desarrollada aproximación modelo-calibrada (Wu y Sitter, 2001a) puede también usarse en las investigaciones por muestreo. Estos estimadores se obtienen, en primer lugar, adaptando un modelo de superpoblación, y a continuación, usando los valores estimados mediante este modelo en la etapa de estimación. Por tanto, si para una población dada se conoce el modelo de superpoblación asociado o un modelo que se ajuste bastante bien a dicha población, entonces puede resultar interesante utilizar la

perspectiva modelo-calibrada para la estimación de la función de distribución poblacional mediante el método de verosimilitud empírica.

Chen y Wu (2002) plantean una aproximación modelo-calibrada para obtener tres estimadores de la función de distribución usando el método de verosimilitud empírica y tres modelos de superpoblación distintos. Estos modelos son bastantes generales, e incluyen los casos más importantes usados en muestreo. Bajo los modelos que se describen, estos estimadores tienen mínima esperanza bajo el modelo de la varianza asintótica basada en el diseño entre una clase de estimadores, es decir, son óptimos dentro de esa clase. Además, estos estimadores son asintóticamente insesgados bajo el diseño si se satisface el modelo y aproximadamente insesgados bajo el modelo. Por último, los estimadores resultantes son verdaderas funciones de distribución y permiten obtener cuantiles eficientemente mediante inversión directa.

Sea un modelo de superpoblación semi-paramétrico, ξ , en el cual se supone que la relación entre y y \mathbf{x} puede describirse de la forma siguiente

$$E_{\xi}(y_i|\mathbf{x}_i) = \mu(\mathbf{x}_i, \theta), \quad V_{\xi}(y_i|\mathbf{x}_i) = \sigma_i^2, \quad \text{con } i = \{1, \dots, N\},$$

donde θ es un vector de parámetros de la superpoblación. Para este vector, se puede obtener un estimador basado en el diseño, $\hat{\theta}$, utilizando métodos generales para la estimación de ecuaciones (véase por ejemplo Godambe y Thompson, 1986 y Wu y Sitter, 2001a).

Dado el modelo ξ , el estimador modelo-calibrado de verosimilitud empírica (*MCPE*) para la función de distribución viene dado por

$$\hat{F}_{MCPE}(t) = \sum_{i \in s} \hat{p}_i \delta(t - y_i) = \sum_{i \in s} \hat{p}_i z_i, \quad (2.83)$$

donde los pesos \hat{p}_i maximizan la función (2.11) sujeta a las restricciones (2.5) y (2.45). La función w_i de la restricción (2.45) viene dada por $w_i = E_{\xi}(z_i|\mathbf{x}_i) = E_{\xi}(\delta(t_0 - y_i)|\mathbf{x}_i) = P(y_i \leq t_0|\mathbf{x}_i)$. El valor t_0 en la segunda restricción se considera fijo para conseguir que el estimador $\hat{F}_{MCPE}(t)$ sea una verdadera función de distribución. Se pueden proponer otras expresiones para w_i , pero se ha considerado $w_i = E_{\xi}(z_i|\mathbf{x}_i)$ porque de entre todos los posibles valores $w_i = w(\mathbf{x}_i)$, el valor $w_i = E_{\xi}(z_i|\mathbf{x}_i)$ minimiza la esperanza bajo el modelo de la varianza asintótica basada en el diseño muestral.

En lo que sigue, se describen tres estimadores de verosimilitud pseudo empírica modelo-calibrados distintos para la función de distribución basados en diferentes modelos de superpoblación (véase Chen y Wu, 2002). Wu (2003) proporciona resultados de optimalidad para estos estimadores.

Estimadores bajo un modelo de regresión

Un modelo de superpoblación comúnmente usado en poblaciones finitas es el modelo de regresión, que viene dado por

$$y_i = \mu(\mathbf{x}_i, \theta) + \nu_i \varepsilon_i, \quad i = \{1, \dots, N\}, \quad (2.84)$$

donde ν_i es una función conocida de \mathbf{x}_i , y ε_i , con $i = \{1, \dots, N\}$, son variables aleatorias independientes e idénticamente distribuidas con media 0 y varianza σ^2 .

Para un modelo de regresión lineal se tiene que $\mu(\mathbf{x}_i, \theta) = \mathbf{x}_i^t \theta$, aunque se puede considerar cualquier otro modelo no lineal. Sea θ_N y σ_N los estimadores de θ y σ , respectivamente, basados en los datos poblacionales. Se sabe que bajo un modelo de regresión lineal con varianzas homogéneas y θ de dimensión J , $\theta_N = (\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t \mathbf{y}$, donde \mathbf{x} es la matriz de orden $N \times J$, $\mathbf{y} = (y_1, \dots, y_N)^t$, y

$$\sigma_N^2 = \frac{(\mathbf{y} - \mathbf{x}\theta_N)^t (\mathbf{y} - \mathbf{x}\theta_N)}{(N - J)}.$$

Bajo el modelo (2.84), las cantidades w_i en (2.45) vienen dadas por

$$\begin{aligned} w_i^* &= E_{\xi}(z_i | \mathbf{x}_i) = P(y_i \leq t_0 | \mathbf{x}_i) = P(\mu(\mathbf{x}_i, \theta_N) + \nu_i \varepsilon_i \leq t_0) = \\ &= G\left(\frac{t_0 - \mu(\mathbf{x}_i, \theta_N)}{\nu_i}\right), \end{aligned} \quad (2.85)$$

donde $G(\cdot)$ es la función de distribución de los términos ε_i , esto es,

$$G(t) = \frac{1}{N} \sum_{i=1}^N \delta(t - \varepsilon_i).$$

Como el vector θ_N es desconocido, es necesario buscar una estimación eficiente para poder obtener las cantidades w_i^* . Para este propósito, también es necesario una estimación de $G(\cdot)$. Una posible estimación viene dada por los residuos estimados, $\hat{\varepsilon}_i$, y la función $G_n(\cdot)$, donde

$$\hat{\varepsilon}_i = \frac{y_i - \mu(\mathbf{x}_i, \hat{\theta})}{\nu_i}, \quad G_n(t) = \sum_{i \in s} d_i^* \delta(t - \hat{\varepsilon}_i) = \frac{\sum_{i \in s} d_i \delta(t - \hat{\varepsilon}_i)}{\sum_{j \in s} d_j},$$

y $\hat{\theta}$ es la estimación basada en el diseño para θ_N . En conclusión, se llega a que las cantidades w_i de la restricción (2.45) vienen dadas por

$$w_i = G_n\left(\frac{t_0 - \mu(\mathbf{x}_i, \hat{\theta})}{\nu_i}\right). \quad (2.86)$$

En algunas situaciones, resulta razonable asumir que los términos de error ε_i en el modelo (2.84) están normalmente distribuidos. En este caso, se llega a que

$$w_i^* = \Phi \left(\frac{t_0 - \mu(\mathbf{x}_i, \theta_N)}{\nu_i \sigma_N} \right), \quad (2.87)$$

donde $\Phi(\cdot)$ es la función de distribución de la ley de probabilidad normal estándar. Se observa que se considera θ_N y no θ en la definición de w_i^* . Esto se hace para que las cantidades w_i^* estén bien definidas sobre la población y puedan tomar todos los argumentos posibles basados en el diseño. En la práctica, se sustituye θ_N y σ_N por $\hat{\theta}$ y $\hat{\sigma}$ respectivamente, donde éstas últimas cantidades son las estimaciones basadas en el diseño muestral de los parámetros desconocidos del modelo. De este modo, se llega a la expresión

$$w_i = \Phi \left(\frac{t_0 - \mu(\mathbf{x}_i, \hat{\theta})}{\nu_i \hat{\sigma}} \right). \quad (2.88)$$

En resumen, el estimador *MCPE* según el modelo (2.84) está dado por $\hat{F}_{MCPE}^{(1)}(t) = \sum_{i \in s} \hat{p}_i \delta(t - y_i)$, donde los pesos \hat{p}_i maximizan la función (2.11) sujeta a las restricciones (2.5) y (2.45). Las cantidades w_i de la segunda restricción vienen dadas por (2.86), o por los valores (2.88) en caso de existir normalidad en los errores del modelo de superpoblación.

Estimadores bajo un modelo lineal generalizado

Resulta atractivo adaptar un modelo lineal generalizado a las cantidades $w_i = E_{\xi}(z_i | \mathbf{x}_i) = P(y_i \leq t_0 | \mathbf{x}_i)$. Para ello se considera el modelo de regresión logístico

$$\log \left(\frac{w_i}{1 - w_i} \right) = \mathbf{x}_i^t \theta, \quad (2.89)$$

con función varianza $V(w) = w(1 - w)$. Bajo este modelo, el parámetro poblacional θ_N puede definirse como una solución de las ecuaciones de estimación óptimas basadas en la población, esto es, $\sum_{i=1}^N \mathbf{x}_i (z_i^* - w_i) = \mathbf{0}$, donde $z_i^* = \delta(t_0 - t)$. Así,

$$w_i^* = \frac{\exp(\mathbf{x}_i^t \theta_N)}{1 + \exp(\mathbf{x}_i^t \theta_N)}. \quad (2.90)$$

Un estimador basado en el diseño, $\hat{\theta}$, para el parámetro poblacional θ_N puede obtenerse resolviendo la correspondiente versión muestral del sistema anterior, esto es, $\sum_{i \in s} d_i \mathbf{x}_i (z_i^* - w_i) = \mathbf{0}$. De este modo, un segundo *MCPE*, esta vez

bajo el modelo (2.89), viene dado por $\widehat{F}_{MCPE}^{(2)}(t) = \sum_{i \in s} \widehat{p}_i \delta(t - y_i)$, donde los pesos \widehat{p}_i se obtienen considerando

$$w_i = \frac{\exp(\mathbf{x}_i^t \widehat{\theta})}{1 + \exp(\mathbf{x}_i^t \widehat{\theta})}. \quad (2.91)$$

El modelo de regresión logístico da una razonable estimación en la mayoría de las estimaciones.

Estimadores bajo valores pseudo estimados de un modelo semi-paramétrico

La variable $z_i = \delta(t - y_i)$ toma solamente valores 0 ó 1, pero los valores estimados w_i construidos bajo los modelos (2.84) y (2.89) están siempre entre 0 y 1. También es posible utilizar los llamados valores pseudo estimados $w_i = \delta(t_0 - \widehat{y}_i)$, los cuales también son variables dicotómicas y donde \widehat{y}_i son valores estimados para y_i .

Bajo un modelo semi-paramétrico, $E_\xi(y_i | \mathbf{x}_i) = \mu_i$ y $V_\xi(y_i | \mathbf{x}_i) = v(\mu_i)$, donde $\mu_i = \mu(\mathbf{x}_i, \theta)$ y $v(\cdot)$ es una función varianza. Los valores estimados \widehat{y}_i están dados por $\mu(\mathbf{x}_i, \widehat{\theta})$. Sea $h(\cdot)$ una conocida función de enlace tal que $h(\mu_i) = \mathbf{x}_i \theta$. $\widehat{\theta}$ es el estimador máximo verosímil que se obtiene del sistema de ecuaciones

$$\sum_{i \in s} \frac{d_i \mathbf{x}_i (y_i - \mu_i)}{v(\mu_i) h'(\mu_i)} = \mathbf{0},$$

donde $h'(u) = \partial h(u) / \partial u$. θ_N es la solución a

$$\sum_{i=1}^N \frac{\mathbf{x}_i (y_i - \mu_i)}{v(\mu_i) h'(\mu_i)} = \mathbf{0}.$$

El estimador $\widehat{F}_{MCPE}^{(3)}(t) = \sum_{i \in s} \widehat{p}_i \delta(t - y_i)$ cuando los pesos \widehat{p}_i se obtienen usando los valores pseudo estimados

$$w_i = \delta(t_0 - \widehat{y}_i). \quad (2.92)$$

En la práctica se usan estas cantidades debido a que los valores

$$w_i^* = \delta(t_0 - \mu(\mathbf{x}_i, \theta_N)), \quad (2.93)$$

son desconocidos.

Bajo un modelo lineal simple con una única variable auxiliar, $\mu(\mathbf{x}, \theta) = \theta_0 + \theta_1 x_i$, y

$$\frac{1}{N} \sum_{i=1}^N w_i = \frac{1}{N} \sum_{i=1}^N \delta(t_0 - (\theta_0 + \theta_1 x_i)) = F_x \left(\frac{t_0 - \theta_0}{\theta_1} \right),$$

donde $F_x(t)$ es la función de distribución de la variable x . La restricción (2.45) se resume a

$$\sum_{i \in s} p_i \delta(t_0 - (\hat{\theta}_0 + \hat{\theta}_1 x_i)) = F_x \left(\frac{t_0 - \hat{\theta}_0}{\hat{\theta}_1} \right),$$

con lo que solamente se debe conocer la distribución de frecuencias de x para obtener $\hat{F}_{MCPE}^{(3)}(t)$.

Notamos que puede usarse cualquier modelo de superpoblación. Si el modelo de superpoblación asociado a la población en estudio es otro distinto a cualquiera de estos tres, el planteamiento para el cálculo del estimador de verosimilitud pseudo empírica modelo-calibrado es similar a lo comentado. Bastaría con obtener las cantidades w_i óptimas bajo el modelo de superpoblación asociado.

La elección del valor t_0 es un aspecto importante, puesto que los estimadores son más precisos para estimar $F_y(t)$ cuando t está en las cercanías del punto t_0 . En consecuencia, ningún w_i con un valor fijo t_0 puede ser uniformemente óptimo para $F_y(t)$ en todos los valores de t . El problema de encontrar un valor óptimo t_0 no se discute en Chen y Wu (2002). De hecho, sus correspondientes estudios empíricos usan cuantiles poblacionales de la variable y para obtener el valor t_0 . Esta elección no puede realizarse en la práctica debido que los cuantiles poblacionales de la variable de estudio son desconocidos. En resumen, estos estimadores presentan dos inconvenientes principalmente: (i) es necesario el conocimiento de un modelo de superpoblación para los datos muestrales del estudio y (ii) se hace un uso poco eficiente de la información auxiliar, puesto que sería posible definir los estimadores usando más de un punto t_0 , utilizando de este modo más información auxiliar, lo que conlleva esperar estimaciones más precisas. Estos problemas puede solventarse en gran medida mediante la metodología propuesta en la Sección 2.4.3, donde se usa un vector \mathbf{t}_0 para obtener estimaciones más eficientes para cualquier t .

El estimador que se propone para la función de distribución usa una aproximación modelo-asistida y el método de verosimilitud empírica. Con el objetivo de que este estimador sea más eficiente para cualquier t , éste usa un vector \mathbf{t}_0 basado en los cuantiles poblacionales de una pseudo-variable que es conocida en la práctica. Además, este estimador es una verdadera función de distribución y

goza de una excelente ganancia en eficiencia como consecuencia de un uso efectivo de la información auxiliar. Éstas son dos de las ventajas más importantes del estimador propuesto.

Como los estimadores de verosimilitud pseudo empírica modelo-calibrados van a ser usados en los estudios de simulación para estudiar la eficiencia de los estimadores propuestos modelos-asistidos, y además, estos dos tipos de estimadores presentan una gran similitud entre ellos, resulta interesante destacar las propiedades asintóticas más importantes que se han discutido para los *MCPE*. Estas son las siguientes.

El primer resultado asintótico demuestra que el estimador $\widehat{F}_{MCPE}^{(1)}(t)$ está relacionado con un estimador de tipo regresión generalizado. Para ello, se debe requerir el siguiente esquema asintótico.

Se asume que hay una secuencia de poblaciones finitas $\{U_\nu, \nu = 1, 2, \dots\}$. $F_\nu(t)$ denota $F_y(t)$ para la población U_ν . El índice ν se suprime y todos los procesos límites se construyen cuando $\nu \rightarrow \infty$. Denotando $\mu'(\mathbf{x}, \theta) = \partial\mu(\mathbf{x}, \theta)/\partial\theta$, las siguientes condiciones de regularidad deben satisfacerse:

$$(C2.20). \quad |\widehat{\theta} - \theta_N| = O_p(n^{-1/2}), \text{ y } \widehat{\sigma} - \sigma_N = O_p(n^{-1/2}).$$

$$(C2.21). \quad \max_{i=1, \dots, N}(nd_i/N) = O(1).$$

(C2.22). Para cada \mathbf{x}_i , $\mu(\mathbf{x}_i, \theta)$ es doble diferenciable, y además

$$\frac{1}{N} \sum_{i=1}^N \left(\frac{\mu'(\mathbf{x}_i, \theta_N)}{\nu_i} \right)^2 = O(1), \quad \frac{1}{N} \sum_{i=1}^N \left(\frac{\mu(\mathbf{x}_i, \theta_N)}{\nu_i} \right)^2 = O(1),$$

$$\frac{1}{N} \sum_{i=1}^N \frac{1}{\nu_i^2} = O(1).$$

Teorema 2.7 *Bajo las condiciones de regularidad (C2.20), (C2.21) y (C2.22), el estimador de verosimilitud pseudo empírica modelo calibrado $\widehat{F}_{MCPE}^{(1)}(t)$ basado en el modelo (2.84) y con*

$$w_i = \Phi \left(\frac{t_0 - \mu(\mathbf{x}_i, \widehat{\theta})}{\nu_i \widehat{\sigma}} \right), \text{ o } w_i = G \left(\frac{t_0 - \mu(\mathbf{x}_i, \theta)}{\nu_i} \right),$$

es asintóticamente equivalente a un estimador tipo regresión generalizado, esto es:

$$\widehat{F}_{MCPE}^{(1)}(t) = \widehat{F}_{HKy}(t) + B_N(\bar{w}_N^* - \bar{w}_d^*) + o_p(n^{-1/2}), \quad (2.94)$$

$$\text{donde } \widehat{F}_{HKy}(t) = \sum_{i \in s} d_i^* \delta(t - y_i), \quad \bar{w}_N^* = \frac{1}{N} \sum_{i=1}^N w_i^*, \quad \bar{w}_d^* = \sum_{i \in s} d_i^* w_i^*,$$

$$w_i^* = \Phi \left(\frac{t_0 - \mu(\mathbf{x}_i, \theta_N)}{\nu_i \sigma_N} \right) \quad \text{y} \quad B_N = \frac{\sum_{i=1}^N (w_i^* - \bar{w}_N^*) \delta(t_0 - y_i)}{\sum_{i=1}^N (w_i^* - \bar{w}_N^*)^2}.$$

Teorema 2.8 *Bajo las condiciones de regularidad (C2.20), (C2.21) y (C2.22), el estimador $\widehat{F}_{MCPE}^{(1)}(t)$ es asintóticamente un estimador consistente bajo el diseño para $F_y(t)$, y por tanto, aproximadamente insesgado bajo el modelo (2.84) cuando $t = t_0$.*

La demostración de estos resultados pueden consultarse en Chen y Wu (2002). Bajo algunas condiciones de regularidad bastantes débiles, este teorema también es valido cuando $G(\cdot)$ se reemplaza por su estimación $G_n(\cdot)$, puesto que reemplazar $\widehat{\theta}$ por θ no cambia asintóticamente el estimador resultante.

El estimador tipo regresión generalizado referido en la expresión (2.94) se utiliza tan sólo para la comparación asintótica, puesto que no es un estimador real. Esta equivalencia hace que se herede la eficiencia asintótica de un estimador tipo regresión generalizado. Por otro lado, se eliminan las desventajas del estimador $\widehat{F}_d(t)$ dado en (2.82). También es importante recordar que $\widehat{F}_{MCPE}^{(1)}(t)$ será más eficiente cuando t está próximo del valor t_0 considerado en la restricción (2.45).

De este teorema, los siguientes corolarios pueden derivarse.

Corolario 2.4 *Bajo las condiciones de regularidad (C2.20), (C2.21) y (C2.22), el MCPE para la función distribución $\widehat{F}_{MCPE}^{(2)}(t)$ es asintóticamente equivalente a un estimador de tipo regresión generalizado, esto es,*

$$\widehat{F}_{MCPE}^{(2)}(t) = \widehat{F}_{HKy}(t) + B_N(\bar{w}_N^* - \bar{w}_d^*) + o_p(n^{-1/2}).$$

Corolario 2.5 *Bajo las condiciones de regularidad (C2.20), (C2.21) y (C2.22), $\widehat{F}_{MCPE}^{(2)}(t)$ es asintóticamente un estimador consistente bajo el diseño para $F_y(t)$. Además, es aproximadamente insesgado bajo el modelo (2.89) y para $t = t_0$.*

Corolario 2.6 *Bajo las condiciones de regularidad (C2.20), (C2.21) y (C2.22), el MCPE para la función distribución $\widehat{F}_{MCPE}^{(3)}(t)$ es asintóticamente equivalente a un estimador de tipo regresión generalizado, esto es,*

$$\widehat{F}_{MCPE}^{(3)}(t) = \widehat{F}_{HKy}(t) + B_N(\overline{w}_N^* - \overline{w}_d^*) + o_p(n^{-1/2}).$$

Notamos que $\widehat{F}_{MCPE}^{(3)}(t)$ es asintóticamente un estimador insesgado bajo el diseño para $F_y(t)$, aunque no es aproximadamente un estimador insesgado bajo el modelo puesto que $E_\xi(\delta(t - y_i)|\mathbf{x}) \neq \delta(t - \mu(\mathbf{x}_i, \theta))$. No obstante, $\widehat{F}_{MCPE}^{(3)}(t)$ posee propiedades no disfrutadas por $\widehat{F}_{MCPE}^{(1)}(t)$ y $\widehat{F}_{MCPE}^{(2)}(t)$. Por ejemplo, un argumento a favor de $\widehat{F}_{MCPE}^{(3)}(t)$ es que si el modelo se ajusta perfectamente a la población, esto es, $y_i = \mu(\mathbf{x}_i, \theta)$ para $i = \{1, \dots, N\}$, entonces, $w_i = \delta(t_0 - y_i)$, y $\widehat{F}_{MCPE}^{(3)}(t)$ se reduce al valor exacto de $F_y(t_0)$. Entonces, es de esperar que en el caso de fuerte información auxiliar, la correlación entre y_i y \widehat{y}_i sea alta, y consecuentemente $\widehat{F}_{MCPE}^{(3)}(t)$ cumplirá bien.

En los tres casos, los estimadores son asintóticamente equivalentes a un estimador de tipo regresión generalizado y son verdaderas funciones de distribución.

El siguiente paso es conocer las varianzas asintóticas de los estimadores $\widehat{F}_{MCPE}^{(1)}(t)$, $\widehat{F}_{MCPE}^{(2)}(t)$ y $\widehat{F}_{MCPE}^{(3)}(t)$. Sea $\widehat{F}_{MCPE}(t)$ uno de estos tres estimadores. La varianza asintótica basada en el diseño de $\widehat{F}_{MCPE}(t)$ es la misma que el estimador tipo razón dado por $\sum_{i \in s} d_i^*(\delta(t - y_i) - w_i^* B_N)$, donde w_i^* están dadas en (2.85), (2.87), (2.90) o (2.93), dependiendo del caso considerado. Denotando por V_p la varianza basada en el diseño, el siguiente teorema puede establecerse (véase Chen y Wu, 2002).

Teorema 2.9 *Bajo las condiciones (C2.20), (C2.21) y (C2.22) y un diseño con tamaño muestral fijado, la varianza asintótica basada en el diseño de $\widehat{F}_{MCPE}(t)$ viene dada por*

$$V_p(\widehat{F}_{MCPE}(t)) = \frac{1}{N^2} \sum_{i < j} \sum_{j=1}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{U_i}{\pi_i} - \frac{U_j}{\pi_j} \right)^2 + o(n^{-1}),$$

donde $U_i = \delta(t - y_i) - F_y(t) - (w_i^* - \overline{w}^*) B_N$ y $\overline{w}^* = N^{-1} \sum_{i=1}^N w_i^*$. $V_p(\widehat{F}_{MCPE}(t))$ puede estimarse consistentemente por

$$\widehat{V}_p(\widehat{F}_{MCPE}(t)) = \frac{1}{N^2} \sum_{i < j} \sum_{j \in s} \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{u_i}{\pi_i} - \frac{u_j}{\pi_j} \right)^2,$$

donde $u_i = \delta(t - y_i) - \widehat{F}_{MCPE}(t) - (w_i - \bar{w})\widehat{B}$, y \bar{w} y \widehat{B} son estimadores basadas en la muestra para \bar{w}^* y B_N , respectivamente. Dependiendo del modelo, los valores w_i están dados en (2.86), (2.88), (2.91) o (2.92).

2.4.3. Estimador propuesto modelo-asistido

En esta sección se propone usar la aproximación modelo-asistida basada en el método de verosimilitud empírica para construir un estimador de la función de distribución poblacional finita. Información auxiliar multivariante puede incorporarse en la etapa de estimación y se hace un uso efectivo de la información auxiliar. Este estimador basado en el diseño muestral es una auténtica función de distribución que disfruta de varias propiedades importantes.

Para construir el nuevo estimador para $F_y(t)$, se modifican los pesos del estimador $\widehat{F}_{HKy}(t)$, es decir d_i^* , por unos nuevos pesos \widehat{p}_i . Este conjunto de pesos se determina por medio de una aproximación modelo-asistida (véase Apéndice B) y usando las técnicas de verosimilitud empírica (Sección 2.2).

Se considera la estimación modelo-asistida porque esta aproximación proporciona un esquema de trabajo conveniente en el cual se pueden desarrollar estimadores muy precisos. A través de un modelo de superpoblación se construyen estimadores basados en la muestra que mejoran la precisión de las estimaciones cuando el modelo es correcto, pero que también mantiene propiedades importantes, tales como consistencia y una varianza estimable, cuando el modelo es incorrecto.

Se considera el usual modelo de regresión dado por

$$y_i = \beta^t \mathbf{x}_i + v_i \varepsilon_i, \quad i = 1, \dots, N \quad (2.95)$$

donde v_i es una función conocida de x_i y los valores ε_i son variables aleatorias independientes e idénticamente distribuidas con media 0 y varianza σ^2 .

En la práctica, los valores del vector β son desconocidos. Mediante la teoría de regresión, puede deducirse que el estimador de mínimos cuadrados de β (Särndal *et al.*, 1992)

$$\mathbf{B} = \left(\sum_{i \in U} \frac{\mathbf{x}_i \mathbf{x}_i^t}{\sigma^2} \right)^{-1} \cdot \sum_{i \in U} \frac{\mathbf{x}_i y_i}{\sigma^2} \quad (2.96)$$

es el mejor estimador insesgado lineal de β bajo el modelo (3.49). B es un parámetro poblacional desconocido, pero puede estimarse usando los datos

muestrales y aplicando el principio de estimación de las probabilidades de inclusión, esto es

$$\hat{\beta} = \left(\sum_{i \in s} \frac{d_i \mathbf{x}_i \mathbf{x}_i^t}{\sigma^2} \right)^{-1} \cdot \sum_{i \in s} \frac{d_i \mathbf{x}_i y_i}{\sigma^2}. \quad (2.97)$$

El estimador propuesto modelo-asistido basado en el método de verosimilitud empírica se obtiene definiendo la pseudo-variable g , donde $g_i = \hat{\beta}^t \mathbf{x}_i$, para $i \in s$. Esta variable puede considerarse como una predicción para y_i bajo el anterior modelo lineal.

Sean $t_{g25} = Q_g(0,25)$, $t_{g50} = Q_g(0,5)$ y $t_{g75} = Q_g(0,75)$ los cuartiles poblacionales de la variable g , donde $Q_g(\alpha) = \inf\{t \mid F_g(t) \geq \alpha\} = F_g^{-1}(\alpha)$. Bajo nuestro marco de trabajo, estas cantidades están disponibles, puesto que asumimos que la información auxiliar poblacional es conocida. El estimador de verosimilitud pseudo empírica modelo-asistido para la función de distribución se define como $\hat{F}_{MA}(t) = \sum_{i \in s} \hat{p}_i \delta(t - y_i)$, donde los nuevos pesos \hat{p}_i se obtienen maximizando $\hat{l}(\mathbf{p})$ sujeta a (??) y a las siguientes condiciones

$$\sum_{i \in s} p_i \delta(t_{g25} - g_i) = \frac{1}{N} \sum_{i=1}^N \delta(t_{g25} - g_i) = F_g(t_{g25}) = 0,25, \quad (2.98)$$

$$\sum_{i \in s} p_i \delta(t_{g50} - g_i) = \frac{1}{N} \sum_{i=1}^N \delta(t_{g50} - g_i) = F_g(t_{g50}) = 0,5, \quad (2.99)$$

$$\sum_{i \in s} p_i \delta(t_{g75} - g_i) = \frac{1}{N} \sum_{i=1}^N \delta(t_{g75} - g_i) = F_g(t_{g75}) = 0,75. \quad (2.100)$$

Nótese que la idea de usar $\delta(t - g_i)$, para algún t , como una variable de calibración para construir restricciones tales como (3.55), (3.56) y (3.57) fue discutida, en primer lugar, por Wu y Sitter (2001a) y posteriormente elaborada en Chen y Wu (2002).

Existen dos aspectos importantes relacionados con este o cualquier otro procedimiento de estimación. Éstos son la eficiencia y la consistencia. La eficiencia se refiere al cumplimiento del estimador en términos de sesgo y error cuadrático medio. En la Sección 2.4.5, se realiza una comparación de la eficiencia de $\hat{F}_{MA}(t)$ con respecto otros estimadores conocidos. Las restricciones (3.55), (3.56) y (3.57) son requerimientos de consistencia altamente usados y son impuestos en la práctica, porque resulta razonable pensar que los pesos que dan estimaciones perfectas para las variables auxiliares, deberían también dar una buena estimación para la variable de estudio.

La elección de t_{g25} , t_{g50} y t_{g75} en (3.55), (3.56) y (3.57) se realiza por varias razones. En primer lugar, esto está altamente relacionado con la existencia de la solución del método de verosimilitud empírica. Si se usaran más de tres valores t_0 , esto es, un mayor número de restricciones, se podría llegar a problemas de existencia de solución (véase la Sección 2.4.4 para un mayor detalle). Estos puntos están también especificados por motivos de eficiencia. Si se usa un único punto t_0 , $\widehat{F}_{MA}(t)$ será más eficiente para t en las proximidades de t_0 . Para varios valores de t_0 , es razonable asumir que si éstos están perfectamente distribuidos dentro del posible rango de valores de t , entonces, $\widehat{F}_{MA}(t)$ será más eficiente. Los valores t_{g25} , t_{g50} y t_{g75} exhiben una buena distribución y por tanto, $\widehat{F}_{MA}(t)$ será más preciso cuando t se encuentre en los alrededores de los cuartiles poblacionales de la variable g . Esto afecta a un alto rango de valores de la variable de estudio.

$\widehat{F}_{MA}(t)$ será, especialmente, más eficiente cuando t es igual a uno de los valores t_{g25} , t_{g50} ó t_{g75} . Esto implica que no hay una elección óptima de valores para todo t . Por otro lado, para t igual a t_{g25} , t_{g50} y t_{g75} y si el modelo (3.49) se ajusta perfectamente a la población de estudio, esto es, $y_i = \beta^t \mathbf{x}_i = g_i$, $i = 1, \dots, N$, entonces $\delta(t - g_i) = \delta(t - y_i)$ y $\widehat{F}_{MA}(t)$ se reduce al valor exacto de $F_y(t)$. Es de esperar, que en el caso de una información auxiliar fuertemente relacionada con la variable de estudio, la correlación entre y_i y g_i será mayor, y consecuentemente, $\widehat{F}_{MA}(t)$ cumplirá mejor en el sentido de obtener estimaciones más precisas para $F_y(t)$.

Denotando por $\mathbf{t}_g = (t_{g25}, t_{g50}, t_{g75})^t$,

$$\delta(\mathbf{t}_g - g_i) = (\delta(t_{g25} - g_i), \delta(t_{g50} - g_i), \delta(t_{g75} - g_i))^t$$

y $\mathbf{K} = (0,25, 0,50, 0,75)^t$, las restricciones anteriores pueden expresarse por

$$\sum_{i \in s} p_i \delta(\mathbf{t}_g - g_i) = \mathbf{K} \quad (2.101)$$

o también como

$$\sum_{i \in s} p_i \mathbf{u}_i = \mathbf{0}, \quad (2.102)$$

donde $\mathbf{u}_i = \delta(\mathbf{t}_g - g_i) - \mathbf{K}$.

Mediante el conocido método de multiplicadores de Lagrange, puede demostrarse que la solución del problema de maximización sujeto a las condiciones (??) y (2.102) está dado por

$$\widehat{p}_i = \frac{d_i^*}{1 + \lambda^t \mathbf{u}_i}, \quad (2.103)$$

donde el multiplicador de Lagrange λ , cuya dimensión es tres, se obtiene de la ecuación

$$h(\lambda) = \sum_{i \in s} \frac{d_i^* \mathbf{u}_i}{1 + \lambda^t \mathbf{u}_i} = \mathbf{0}. \quad (2.104)$$

Puede demostrarse que, con probabilidad tendiendo a uno cuando el tamaño muestral va a infinito, existe una única solución a $h(\lambda) = \mathbf{0}$. Si tal solución existe, ésta puede encontrarse, por ejemplo, con el Algoritmo 2.1, el cual tiene garantizada la convergencia a la solución.

2.4.4. Propiedades teóricas

Un estimador modelo-asistido para la función de distribución se ha definido en la Sección 2.4.3. A continuación estudiamos varias propiedades de este estimador, las cuales pueden ser importantes en la práctica. En concreto, se estudia la existencia del estimador, se demuestra que $\hat{F}_{MA}(t)$ es una verdadera función de distribución, se obtiene otra propiedad relacionada con la eficiencia del estimador propuesto y se establecen algunos resultados asintóticos.

Existencia del estimador

Existen dos aspectos computacionales por los cuales el estimador $\hat{F}_{MA}(t)$ no pueda existir: (i) en la obtención del vector $\hat{\beta}$ y (ii) para encontrar la solución a $h(\lambda) = \mathbf{0}$ en (2.104).

En el punto (i), $\hat{\beta}$ siempre existe cuando se aplica información auxiliar univariante. En otro caso, $(\sum_{i \in s} d_i \mathbf{x}_i \mathbf{x}_i^t)^{-1}$ no puede calcularse si no es de rango completo. Esta situación es poco probable cuando $n \geq J$.

Respecto a la segunda cuestión, se ha comentado que puede emplearse el Algoritmo 2.1.

Para el caso de la estimación de la media poblacional, la variable \mathbf{u}_i que usualmente se toma es $\mathbf{u}_i = \mathbf{x}_i - \bar{\mathbf{X}}$ (Chen y Sitter, 1999), la cual está también justificada por un modelo lineal. Bajo esta situación y usando el Algoritmo 2.1, $h(\lambda) = \mathbf{0}$ falla para proporcionar la solución si: (i) el vector de medias $\bar{\mathbf{X}}$ no es un punto interior del conjunto convexo formado por $\{\mathbf{x}_i, i \in s\}$, ó (ii) la matriz $\sum_{i \in s} d_i \mathbf{u}_i \mathbf{u}_i^t$ no es de rango completo.

En (i), el estimador de verosimilitud pseudo empírica no existe. Para el

caso de estimar la media poblacional, esto ocurre con una probabilidad tendiendo a cero cuando el tamaño muestral tiende a infinito. En el escenario de la estimación de la función de distribución, la situación es bastante diferente. En particular, para el procedimiento propuesto, el vector \mathbf{K} es siempre un punto interior del conjunto formado por $\{\delta(\mathbf{t}_g - g_i), i \in s\}$, puesto que los componentes de este vector son 0 ó 1, mientras que los componentes de \mathbf{K} toman valores dentro de $[0, 1]$. Notamos que los componentes del vector $\delta(\mathbf{t}_g - g_i)$ no pueden ser todos 0 ó 1 para $i \in s$, salvo para situaciones extremas.

Sea $\mathbf{t}_0 = (t_{0(1)}, \dots, t_{0(h)}, \dots, t_{0(H)})^t$ otro vector diferente de \mathbf{t}_g con similar o diferente dimensión y que puede usarse en restricciones como la dada por (2.101). Respecto al punto (ii), decir que resulta necesario una cuidadosa elección del vector \mathbf{t}_0 para evitar o eliminar el problema de multicolinealidad. En lo que sigue, se justifica la elección $\mathbf{t}_g = (t_{g25}, t_{g50}, t_{g75})^t$. En primer lugar, si se toman valores de $t_{0(h)}$ con dos ellos muy cercanos, entonces, resulta más probable que surja el problema de la multicolinealidad. Si se usan valores extremos de t_0 (o muy elevados o demasiados pequeños), la variable indicadora $\delta(\mathbf{t}_0 - g_i)$ podría tener todos sus elementos iguales a cero o a uno para $i \in s$, y por tanto, el método de verosimilitud empírica no tendría solución. Teniendo estas consideraciones en cuenta, la elección $\mathbf{t}_g = (t_{g25}, t_{g50}, t_{g75})^t$ resulta apropiada, puesto que cada punto está alejado del resto y además, estos puntos no se encuentran cercanos a los valores extremos de la variable g , evitando que la variable indicadora $\delta(\mathbf{t}_g - g_i)$ pueda contener valores que sean todos iguales a cero o a uno para $i \in s$. Bajo este planteamiento, el problema de la multicolinealidad es improbable. Notamos que este problema decrece conforme aumenta el tamaño muestral. Por ejemplo, no se ha observado problemas de multicolinealidad para el estimador $\widehat{F}_{MA}(t)$ en los estudios de simulación de la Sección 2.4.5, mientras que cuando se usa un vector \mathbf{t}_0 con dimensión mayor de 5, nos encontramos problemas de multicolinealidad para tamaños muestrales mayores de 50.

Como se comentó en la Sección 2.4.3, la elección $\mathbf{t}_g = (t_{g25}, t_{g50}, t_{g75})^t$ está también especificada por motivos de eficiencia. Además, el estimador $\widehat{F}_{MA}(t)$ es fácilmente computable debido a que el vector \mathbf{t}_g es de dimension igual a 3 y por tanto, el sistema (2.104) presenta un número pequeño de ecuaciones.

$\widehat{F}_{MA}(t)$ es una auténtica función de distribución

La siguiente cuestión es comprobar si el estimador propuesto es una verdadera función de distribución. Para determinar esto, debemos verificar si se satisfacen, para $\widehat{F}_{MA}(t)$, las condiciones (C2.17), (C2.18) y (C2.19) de la Sección 2.4.1.

Resultado 2.1 *El estimador $\widehat{F}_{MA}(t)$ es una verdadera función de distribución.*

Demostración

Resulta fácil demostrar que la condición (C2.17) siempre se satisface si los pesos \widehat{p}_i , para $i = 1, \dots, n$, son independientes de t :

$$\lim_{t \rightarrow -\infty} \widehat{F}_{MA}(t) = \lim_{t \rightarrow -\infty} \sum_{i \in s} \widehat{p}_i \delta(t - y_i) = \sum_{i \in s} \widehat{p}_i \lim_{t \rightarrow -\infty} \delta(t - y_i) = \sum_{i \in s} \widehat{p}_i 0 = 0.$$

$$\lim_{t \rightarrow +\infty} \widehat{F}_{MA}(t) = \lim_{t \rightarrow +\infty} \sum_{i \in s} \widehat{p}_i \delta(t - y_i) = \sum_{i \in s} \widehat{p}_i \lim_{t \rightarrow +\infty} \delta(t - y_i) = \sum_{i \in s} \widehat{p}_i = 1.$$

Por otro lado, $\widehat{F}_{MA}(t)$ es una función continua por la derecha y monótona no decreciente para unos pesos \widehat{p}_i que sean independientes de t :

- Sea $t_1 < t_2$, entonces $\delta(t_1 - y_i) \leq \delta(t_2 - y_i)$ para $i \in s$ y $\widehat{F}_{MA}(t_1) = \sum_{i \in s} \widehat{p}_i \delta(t_1 - y_i) \leq \sum_{i \in s} \widehat{p}_i \delta(t_2 - y_i) = \widehat{F}_{MA}(t_2)$, puesto que \widehat{p}_i son los mismos valores positivos para t_1 y t_2 .
- Sea $t > t^*$, $\lim_{t \rightarrow t^*} \widehat{F}_{MA}(t) = \lim_{t \rightarrow t^*} \sum_{i \in s} \widehat{p}_i \delta(t - y_i) = \sum_{i \in s} \widehat{p}_i \lim_{t \rightarrow t^*} \delta(t - y_i) = \sum_{i \in s} \widehat{p}_i \delta(t^* - y_i) = \widehat{F}_{MA}(t^*)$.

Por tanto, las condiciones (C2.17), (C2.18) y (C2.19) se satisfacen para $\widehat{F}_{MA}(t)$ si el mismo conjunto de valores \widehat{p}_i son usados para cada argumento t . Como $\widehat{F}_{MA}(t)$ asume un vector fijo \mathbf{t}_g , entonces, $\widehat{F}_{MA}(t)$ es una verdadera función de distribución. \square

$\widehat{F}_{MA}(t)$ es igual a $F_y(t)$ cuando $x_i = y_i$

En las investigaciones por muestreo que incorporan muestreo sucesivo, la variable auxiliar es la misma que la variable principal, pero medida en un periodo anterior. En este caso, la información auxiliar incluye valores poblacionales de la variable x , los cuales pueden estar próximos a los valores de y . En tal situación, resulta razonable esperar que un estimador de $F_y(t)$ debería de aproximarse a $F_y(t)$ a medida que x se aproxima a y . Esta propiedad no se satisface para el estimador estándar, puesto que éste no hace uso de la información auxiliar.

Si $y_i = x_i$, puede verse que $\widehat{\beta} = \mathbf{1}$, $g_i = y_i$ y segunda restricción planteada para el estimador $\widehat{F}_{MA}(t)$ está dada por $\sum_{i \in s} p_i \delta(\mathbf{t}_g - y_i) = F_y(\mathbf{t}_g)$. Así, $\widehat{F}_{MA}(t) = \sum_{i \in s} \widehat{p}_i \delta(t - y_i)$ es exactamente igual a $F_y(t)$ si t coincide con uno de los valores de vector \mathbf{t}_g . Si esto no sucede, la igualdad, en general, no se cumple, aunque se esperan que las desviaciones sean pequeñas si el argumento t está próximo a un componente de \mathbf{t}_g .

Comportamiento asintótico

El siguiente paso es establecer el comportamiento asintótico del estimador $\widehat{F}_{MA}(t)$. Lamentablemente, este estimador usa los vectores \mathbf{t}_g y $\widehat{\beta}$, que son dependientes de la muestra, lo que dificulta la obtención del comportamiento asintótico de este estimador. No obstante, es posible obtener algunos resultados para el estimador $\widehat{F}_{MA1}(t)$ que es muy similar al estimador propuesto aunque menos eficiente al utilizar menos información auxiliar. Este estimador se obtiene equivalentemente que el estimador propuesto, con la diferencia de que los pesos \widehat{p}_i están basados en las restricciones ?? y

$$\sum_{i \in s} p_i \delta(t_0 - g_i) = \frac{1}{N} \sum_{i=1}^N \delta(t_0 - g_i) = F_g(t_0), \quad (2.105)$$

para un valor cualquiera t_0 especificado.

Nota 2.1 *En caso de haber establecido propiedades asintóticas como la equivalencia con otros estimadores o la determinación de la varianza del estimador $\widehat{F}_{MA}(t)$, estas expresiones serían solamente válidas para muestras de gran tamaño y por tanto, serían poco útiles en la práctica. Habitualmente, la replicación de algún tipo, como Bootstrap, Jackknife o replicación mediante muestras balanceadas (Shao y Tu, 1995), es una alternativa que se usa en la etapa*

de estimación de la varianza, particularmente para la estimación de varianzas de funciones de distribución que son especialmente dificultosas. Tales procedimientos son fáciles de computar (Dalglish, 1995) y además, han demostrado un buen cumplimiento para el método de verosimilitud empírica (Chen y Sitter, 1999) y para la estimación de la función de distribución (Lombardia et al., 2003, Lombardia et al., 2004).

Teorema 2.10 Cuando el vector $\hat{\beta}$ se reemplaza por el parámetro \mathbf{B} dado en (3.50), el correspondiente estimador de verosimilitud pseudo empírica modelado-asistido, $\hat{F}_{MA1}^B(t)$, cuando se usa el punto $t_0 = t$, es asintóticamente equivalente a un estimador de tipo regresión generalizado:

$$\hat{F}_{MA1}^B(t) = \hat{F}_{HKy}(t) + (F_b(t) - \hat{F}_b(t))\hat{D} + o_p(n^{-1/2}),$$

donde

$$\hat{D} = \frac{\hat{\sigma}_{z,w}}{\hat{\sigma}_w^2} = \frac{\sum_{i \in s} d_i^* [\delta(t - y_i) - \hat{F}_{HKy}(t)] [\delta(t - b_i) - \hat{F}_b(t)]}{\sum_{i \in s} d_i^* [\delta(t - b_i) - \hat{F}_b(t)]^2},$$

$b_i = \mathbf{B}^t \mathbf{x}_i$, $F_b(t)$ es la función de distribución de la variable b y $\hat{F}_b(t)$ denota el estimador de tipo Hájek para la función de distribución de b en el punto t . z y w denotan las variables $\delta(t - y)$ y $\delta(t - b)$, respectivamente. Por tanto, $\hat{F}_{MA1}^B(t)$ es asintóticamente insesgado bajo el diseño y tiene la misma varianza asintótica que el estimador de tipo regresión generalizado.

Demostración

Para demostrar este teorema, asumimos que la población finita está envuelta en una secuencia de poblaciones donde n y N aumentan de tal forma que $(n/N) \rightarrow f$ cuando $n \rightarrow \infty$. Además, se considera la variable de calibración $\delta(t - b_i)$ en (2.105) para construir $\hat{F}_{MA1}^B(t)$. Sea $u_i = \delta(t - b_i) - F_b(t)$. Puesto que $|u_i| \leq 1$, las condiciones (C2.1) y (C2.2) del Teorema 2.3 se satisfacen y por tanto

$$\lambda = \frac{\sum_{i \in s} d_i^* u_i}{\sum_{i \in s} d_i^* u_i^2} + o_p(n^{-1/2}),$$

y $\hat{p}_i = d_i^* (1 - \lambda u_i) + o_p(n^{-1/2})$. Así:

$$\hat{F}_{MA1}^B(t) = \sum_{i \in s} \hat{p}_i \delta(t - y_i) + o_p(n^{-1/2}) =$$

$$\begin{aligned}
&= \sum_{i \in s} d_i^* \left[1 - \frac{(\widehat{F}_b(t) - F_b(t)) u_i}{\sum_{i \in s} d_i^* u_i^2} \right] \delta(t - y_i) + o_p(n^{-1/2}) = \\
&= \sum_{i \in s} d_i^* \delta(t - y_i) - \frac{\widehat{F}_b(t) - F_b(t)}{\sum_{i \in s} d_i^* u_i^2} \sum_{i \in s} d_i^* u_i \delta(t - y_i) + o_p(n^{-1/2}) = \\
&= \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t)) \frac{\sum_{i \in s} d_i^* u_i \delta(t - y_i)}{\sum_{i \in s} d_i^* u_i^2} + o_p(n^{-1/2}) = \\
&= \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t)) \widehat{D} + o_p(n^{-1/2}).
\end{aligned}$$

□

El resultado anterior es válido cuando se usa el parámetro poblacional \mathbf{B} . El siguiente resultado garantiza que el Teorema 2.10 también se cumple cuando usamos el parámetro muestral $\widehat{\beta}$, el usado por el estimador $\widehat{F}_{MA1}(t)$.

Teorema 2.11 *Los estimadores $\widehat{F}_{MA1}(t)$ y $\widehat{F}_{MA1}^B(t)$ tienen la misma distribución límite.*

Demostración

Denotemos los estimadores modelo-asistidos de verosimilitud pseudo empírica por $\widehat{F}_{MA1}(t) = T_n(\widehat{\beta})$ y $\widehat{F}_{MA1}^B(t) = T_n(\mathbf{B})$. La expresión $T_n(\widehat{\beta})$ depende del estimator $\widehat{\beta}$, es cual es función de los datos muestrales y estima consistentemente el vector de parámetros β . Reemplazando el estimator $\widehat{\beta}$ en $T_n(\cdot)$ por γ y denotándolo por $T_n(\gamma)$, es posible encontrar la distribución límite de la media de esta expresión cuando el valor actual del parámetro es β : $\mu(\gamma) = \lim_{n \rightarrow \infty} E_\beta[T_n(\gamma)] = \widetilde{F}_y(t)$, donde $\widetilde{F}_y(t)$ es el valor límite de $F_y(t)$ cuando $N \rightarrow \infty$. Por tanto

$$\frac{\partial \mu(\gamma)}{\partial \gamma} \Big|_{\gamma=\beta} = \left(\frac{\partial \mu(\gamma)}{\partial \gamma_1} \Big|_{\gamma=\beta}, \frac{\partial \mu(\gamma)}{\partial \gamma_2} \Big|_{\gamma=\beta}, \dots, \frac{\partial \mu(\gamma)}{\partial \gamma_J} \Big|_{\gamma=\beta} \right) = (0, 0, \dots, 0).$$

Randles (1982) demostró que bajo esta condición, la distribución límite de $T_n(\widehat{\beta})$ ($= \widehat{F}_{MA1}(t)$) y $T_n(\mathbf{B})$ ($= \widehat{F}_{MA1}^B(t)$) son idénticas. □

Teorema 2.12 *El comportamiento asintótico del estimador $\widehat{F}_y^{D1}(t) = \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t)) \widehat{D}$ es el mismo del estimador $\widehat{F}_y^{D2}(t) = \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t)) D$, con*

$$D = \frac{\sigma_{z,w}}{\sigma_w^2} = \frac{\sum_{i \in U} d_i^* [\delta(t - y_i) - F_y(t)] [\delta(t - b_i) - F_b(t)]}{\sum_{i \in U} d_i^* [\delta(t - b_i) - F_b(t)]^2}.$$

Consecuentemente, $\widehat{F}_{MA1}^B(t)$ es asintóticamente normal y asintóticamente insesgado bajo el diseño. Su correspondiente varianza asintótica está dada por

$$AV(\widehat{F}_{MA1}^B(t)) = \sum_{i \in U} \sum_{l \in U} \Delta_{il} (d_i^* E_i)(d_l^* E_l), \quad (2.106)$$

donde $\Delta_{il} = \pi_{il} - \pi_i \pi_l$ y $E_i = \delta(t - y_i) - \delta(t - b_i)D$.

Demostración

$\widehat{F}_y^{D1}(t)$ puede expresarse como sigue:

$$\begin{aligned} \widehat{F}_y^{D1}(t) &= \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t))\widehat{D} = \\ &= \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t))(\widehat{D} - D + D) = \\ &= \widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t))D + (F_b(t) - \widehat{F}_b(t))(\widehat{D} - D) = \\ &= \widehat{F}_y^{D2}(t) + (F_b(t) - \widehat{F}_b(t))(\widehat{D} - D). \end{aligned}$$

$\widehat{F}_b(t)$ y \widehat{D} son asintóticamente insesgados bajo el diseño para $F_b(t)$ y D , respectivamente, y por tanto el producto $(F_b(t) - \widehat{F}_b(t))(\widehat{D} - D)$ será de menor orden que $\widehat{F}_b(t)$. Consecuentemente, el término $(F_b(t) - \widehat{F}_b(t))(\widehat{D} - D)$ tiene menor orden que $\widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t))D$. Entonces, $\widehat{F}_y^{D1}(t)$ es asintóticamente insesgado y puesto que los estimadores $\widehat{F}_{HKy}(t)$ y $\widehat{F}_b(t)$ son asintóticamente normales, el estimador $\widehat{F}_y^{D1}(t)$ es asintóticamente normal.

La varianza asintótica de $\widehat{F}_y^{D1}(t)$ coincide con la varianza del estadístico $\widehat{F}_y^{D2}(t)$, la cual está dada por

$$\begin{aligned} V\left(\widehat{F}_{HKy}(t) + (F_b(t) - \widehat{F}_b(t))D\right) &= V\left(\widehat{F}_{HKy}(t) + F_b(t)D - \widehat{F}_b(t)D\right) = \\ &= V\left(\widehat{F}_{HKy}(t) - \widehat{F}_b(t)D\right), \end{aligned}$$

puesto que $F_b(t)D$ es un término constante. Ahora

$$\begin{aligned} \widehat{F}_{HKy}(t) - \widehat{F}_b(t)D &= \sum_{i \in s} d_i^* \delta(t - y_i) - \sum_{i \in s} d_i^* \delta(t - b_i)D = \\ &= \sum_{i \in s} d_i^* [\delta(t - y_i) - \delta(t - b_i)D] = \sum_{i \in s} d_i^* E_i, \end{aligned}$$

con $E_i = \delta(t - y_i) - \delta(t - b_i)D$.

Así, la varianza asintótica de $\widehat{F}_{MA1}^B(t)$ está dada por

$$\begin{aligned} AV(\widehat{F}_{MA1}^B(t)) &= V(\widehat{F}_{HKy}(t) - \widehat{F}_b(t)D) = V\left(\sum_{i \in s} d_i^* E_i\right) = \\ &= \sum_{i \in U} \sum_{l \in U} \Delta_{il}(d_i^* E_i)(d_l^* E_l). \end{aligned}$$

□

Considerando el Teorema 2.11, el resultado anterior también sostiene para $\widehat{F}_{MA1}(t)$ en lugar de $\widehat{F}_{MA1}^B(t)$. Por tanto, asumiendo el estimador $\widehat{F}_{MA1}(t)$, la varianza (2.106) puede estimarse por

$$\widehat{V}(\widehat{F}_{MA1}(t)) = \sum_{i \in s} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} (\widehat{p}_i e_i)(\widehat{p}_j e_j),$$

donde $e_i = \delta(t - y_i) - \delta(t - g_i)\widehat{G}$, con

$$\widehat{G} = \frac{\widehat{\sigma}_{z,v}}{\widehat{\sigma}_v^2} = \frac{\sum_{i \in s} d_i^* [\delta(t - y_i) - \widehat{F}_{HKy}(t)] [\delta(t - g_i) - \widehat{F}_g(t)]}{\sum_{i \in s} d_i^* [\delta(t - g_i) - \widehat{F}_g(t)]^2},$$

y donde v denota a la variable $\delta(t - g)$.

Nota 2.2 Algunos autores, tal como Rao et al. (1990), usan la pseudo-variable $g_i = \widehat{\mathbf{R}}^t \mathbf{x}_i$, para $i = 1, \dots, N$, para construir estimadores modelo-asistidos para la función de distribución, donde $\widehat{\mathbf{R}} = (\sum_{i \in s} d_i \mathbf{x}_i)^{-1} (\sum_{i \in s} d_i y_i)$. El problema de esta pseudo-variable es que puede únicamente usarse para una variable auxiliar. Bajo tal situación, $\widehat{\mathbf{R}}$ ó $\widehat{\beta}$ pueden usarse.

Nota 2.3 El estimador $\widehat{F}_{MA}(t)$ es computacionalmente simple y no depende de parámetros desconocidos, puesto que el vector \mathbf{t}_g puede calcularse fácilmente a través de \mathbf{x} , el cual asumimos es conocido. Cuando esta información no está disponible, el muestreo bifásico es una técnica apropiada para poder aplicar el estimador propuesto. Este muestreo consiste en tomar una primera muestra más grande, donde se recogen los datos de la variable auxiliar. Esto servirá como información auxiliar completa en una segunda muestra más pequeña. Véase el Apéndice B para un mayor detalle de este esquema de muestreo.

2.4.5. Propiedades empíricas

Las principales propiedades del estimador $\widehat{F}_{MA}(t)$ han sido establecidas en la Sección 2.4.4. El siguiente paso es analizar la precisión de este estimador

por medio de un estudio empírico. Por tanto, en esta sección se llevan a cabo estudios de simulación para investigar el cumplimiento muestral de varios estimadores de la función de distribución existentes en la literatura del muestreo en poblaciones finitas.

Para realizar estos estudios se han usado dos poblaciones simuladas generadas bajo una relación lineal entre y y \mathbf{x} , y una población natural, en la cual no se sostiene una relación de este tipo.

Las poblaciones simuladas, de tamaño $N = 1000$, se han generado mediante el modelo

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i, \quad (2.107)$$

donde las variables x_{1i} y x_{2i} se han generado de distribuciones Gamma y las cantidades ϵ_i son variables aleatorias independientes e idénticamente distribuidas con distribución Normal de parámetros 0 y σ^2 . El valor de σ^2 se escoge de modo que el coeficiente de correlación entre y_i y $\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$ es 0.98 para la primera población (Pob098) y 0.80 para la segunda población (Pob080). Como población natural se emplea la población Murthy, la cual presenta un comportamiento exponencial en sus datos. En el Apéndice A están disponibles las propiedades más importantes de estas poblaciones así como sus respectivos diagramas de dispersión.

La precisión del estimador $\widehat{F}_{MA}(t)$ es comparada con los siguientes estimadores: $\widehat{F}_{HTy}(t)$ (el estimador convencional), $\widehat{F}_{CD}(t)$ (Chambers y Dunstan, 1986), $\widehat{F}_{RKM}(t)$ (Rao *et al.*, 1990), $\widehat{F}_{MC}^{(1)}(t)$ (Chen y Wu, 2002), $\widehat{F}_r(t)$ y $\widehat{F}_d(t)$.

Notamos que el modelo (A.2) fue también usado por Chen y Wu (2002), teniendo el estimador $\widehat{F}_{MC}^{(1)}(t)$ el mejor cumplimiento en la mayoría de los casos. En este estudio, también se usa el estimador $\widehat{F}_{MA}(t)$ cuando se considera un valor t_0 en las restricciones. Este estimador será nombrado como $\widehat{F}_{MA1}(t)$. Esto nos permitirá comprobar la ganancia de precisión de usar un vector en las restricciones en lugar de usar un único valor. Así, el mismo punto $t_0 = Q_g(0,5)$ es usado por los estimadores $\widehat{F}_{MC}^{(1)}(t)$ y $\widehat{F}_{MA1}(t)$ para cada t , puesto que esto es necesario para obtener una auténtica función de distribución.

Se llevan a cabo dos estudios de simulación. Por un lado, se evalúan los estimadores en los puntos $t = Q_y(0,25)$, $t = Q_y(0,50)$ y $t = Q_y(0,75)$. Con el fin de revelar el comportamiento medio de los distintos estimadores en diferentes valores de t , se realiza otro estudio de simulación para los argumentos $t = Q_y(0,1), Q_y(0,2), \dots, Q_y(0,9)$. Éste último nos permitirá observar el comportamiento del estimador $\widehat{F}_{MA}(t)$ cuando se usan valores de t alejados de

$$\mathbf{t}_g = (t_{g25}, t_{g50}, t_{g75})^t.$$

Primera simulación

Esta primera simulación consiste en tomar una muestra aleatoria simple de las anteriores poblaciones y estimar la función de distribución en los puntos $t = Q_y(0,25)$, $t = Q_y(0,50)$ y $t = Q_y(0,75)$ mediante los distintos estimadores. Este proceso se repite $B = 1000$ veces para diferentes tamaños muestrales. A continuación, el cumplimiento de todos los estimadores se compara en términos de Sesgo Relativo (SR) y de Eficiencia Relativa (ER), con

$$SR(t) = \frac{1}{B} \sum_{b=1}^B \frac{\widehat{F}(t)_b - F_y(t)}{F_y(t)} \quad ; \quad ER(t) = \frac{ECM[\widehat{F}(t)]}{ECM[\widehat{F}_{HTy}(t)]}, \quad (2.108)$$

donde b expresa la b -ésima simulación, $\widehat{F}(t)$ es un estimador cualquiera de la función de distribución, $ECM[\widehat{F}(t)] = B^{-1} \sum_{b=1}^B [\widehat{F}(t)_b - F_y(t)]^2$ es el Error Cuadrático Medio empírico para $\widehat{F}(t)$, y $ECM[\widehat{F}_{HTy}(t)]$ se define de modo similar para el estimador estándar. Notamos que valores de ER menores de 1 indican que el estimador $\widehat{F}(t)$ es mejor que $\widehat{F}_{HTy}(t)$ en términos de error cuadrático medio.

Las funciones que permiten llevar a cabo este estudio pueden consultarse en el Apéndice ???. La función de R usada para encontrar la solución de la ecuación $h(\lambda) = \mathbf{0}$ puede también verse en Wu (2005).

Las Figuras 2.7 y 2.8 muestran la ER para las tres poblaciones cuando se evalúan en los cuartiles poblacionales de la variable de interés. En los casos donde un estimador cumpla peor que el estimador estándar, su correspondiente línea estará omitida. Los valores absolutos de las cantidades SR para $\widehat{F}_{MA}(t)$ están todas dentro de un rango razonable y son todos menores del 1%. Esto sostiene para el resto de estimadores en la mayoría de los casos. De este modo, estos valores no se muestran.

De las Figuras 2.7 y 2.8 se pueden obtener las siguientes conclusiones:

1. $\widehat{F}_{MA}(t)$ es considerablemente más preciso que el resto de estimadores en $t = Q_y(0,25)$ y $t = Q_y(0,75)$, y exhibe la más baja ER en estos casos. Cuando se estima la mediana de la variable de interés, la situación es diferente, es decir, otros estimadores presentan un similar comportamiento que $\widehat{F}_{MA}(t)$. Por ejemplo, uno de estos estimadores es $\widehat{F}_{MC}^{(1)}(t)$

Figura 2.7: Eficiencia Relativa de distintos estimadores en las poblaciones Pob098 y Pob080.

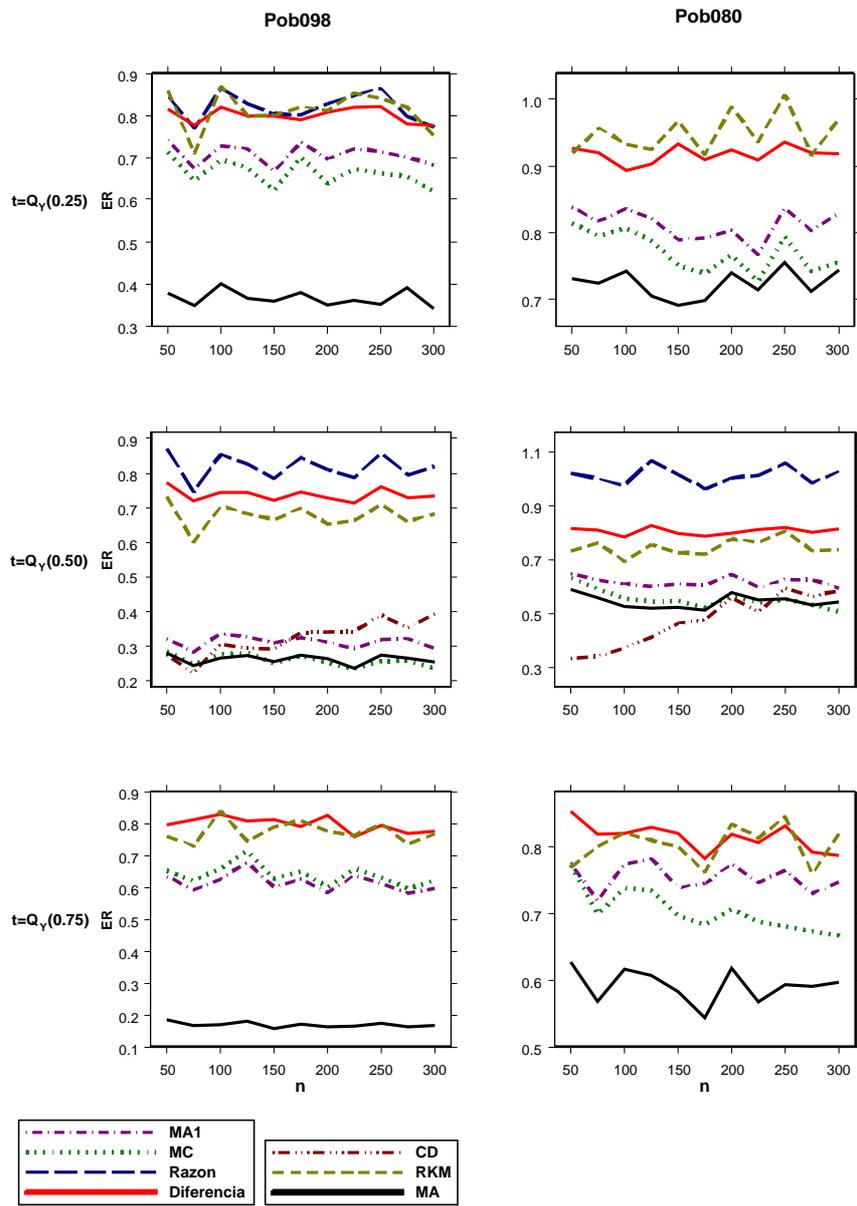
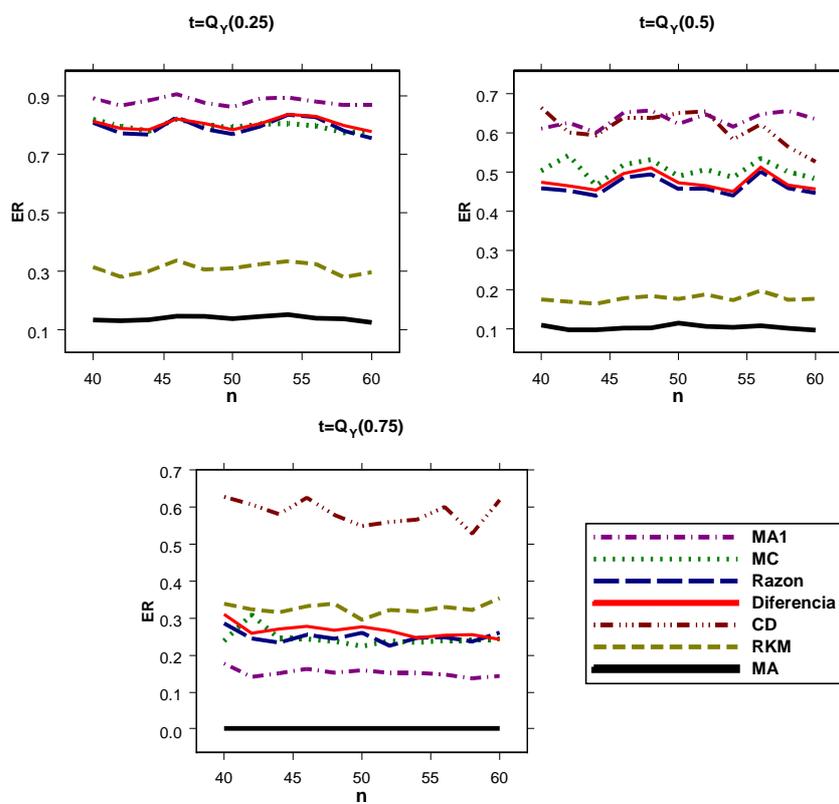


Figura 2.8: Eficiencia Relativa de distintos estimadores en la población Murthy.



en las poblaciones Pop098 y Pop080. Este estimador muestra una mayor ER en los puntos $t = Q_y(0,25)$ y $t = Q_y(0,75)$ debido a que t_0 está alejado de t . El conocimiento del modelo correcto maximiza la eficiencia de $\widehat{F}_{MC}^{(1)}(t)$, pero solamente cuando t está próximo a t_0 .

2. En los casos donde hay una fuerte información auxiliar (Pop098), la ganancia de usar $\widehat{F}_{CD}(t)$, $\widehat{F}_{MC}^{(1)}(t)$, $\widehat{F}_{MA}(t)$ y $\widehat{F}_{MA1}(t)$ puede ser substancial comparada con el estimador estándar.
3. La débil linealidad en la población Murthy afecta especialmente a $\widehat{F}_{MC}^{(1)}(t)$ y $\widehat{F}_{CD}(t)$, los cuales son más eficientes cuando los datos se rigen por un modelo lineal (Pop098 y Pop080).
4. $\widehat{F}_{CD}(t)$ es menos eficiente que el estimador estándar de tipo Horvitz-Thompson cuando la función de distribución se estima en los puntos $t = Q_y(0,25)$ y $t = Q_y(0,75)$. Este estimador es bastante preciso cuando t está próximo a $Q_y(0,5)$, aunque llega a ser considerablemente menos eficiente cuando t está alejado de $Q_y(0,5)$.
5. $\widehat{F}_{MA1}(t)$ es siempre menos preciso que $\widehat{F}_{MA}(t)$. Esto revela que la ganancia de usar el vector \mathbf{t}_g en lugar de un valor t_0 . En cualquier caso, $\widehat{F}_{MA1}(t)$ tiene un buen comportamiento y es siempre más eficiente que el estimador estándar.
6. En términos de ER , el estimador más eficiente para $F_y(t)$ se obtiene por $\widehat{F}_{MA}(Q_y(0,75))$ en la población Murthy. En este caso, los estimadores modelo-calibrados y basados en modelo no tienen un buen comportamiento. Esto puede deberse a que no existe una buena linealidad y a que t está alejado de t_0 .
7. Los estimadores $\widehat{F}_r(t)$ y $\widehat{F}_d(t)$ son siempre considerablemente menos eficientes que $\widehat{F}_{MA}(t)$.

Segunda simulación

La simulación anterior se ha realizado en los puntos $t = Q_y(0,25)$, $t = Q_y(0,50)$ y $t = Q_y(0,75)$. Puede observarse que el orden de estos cuantiles coincide con el orden de los cuantiles del vector \mathbf{t}_g . Es esperable que $\widehat{F}_{MA}(t)$ cumpla bien en esta situación. Por este motivo, usaremos otro estudio de simulación para medir la precisión de los distintos estimadores en los puntos $t = Q_y(0,1), Q_y(0,2), \dots, Q_y(0,9)$.

Figura 2.9: Sesgo Relativo Medio de distintos estimadores en las poblaciones Pob098, Pob080 y Murthy.

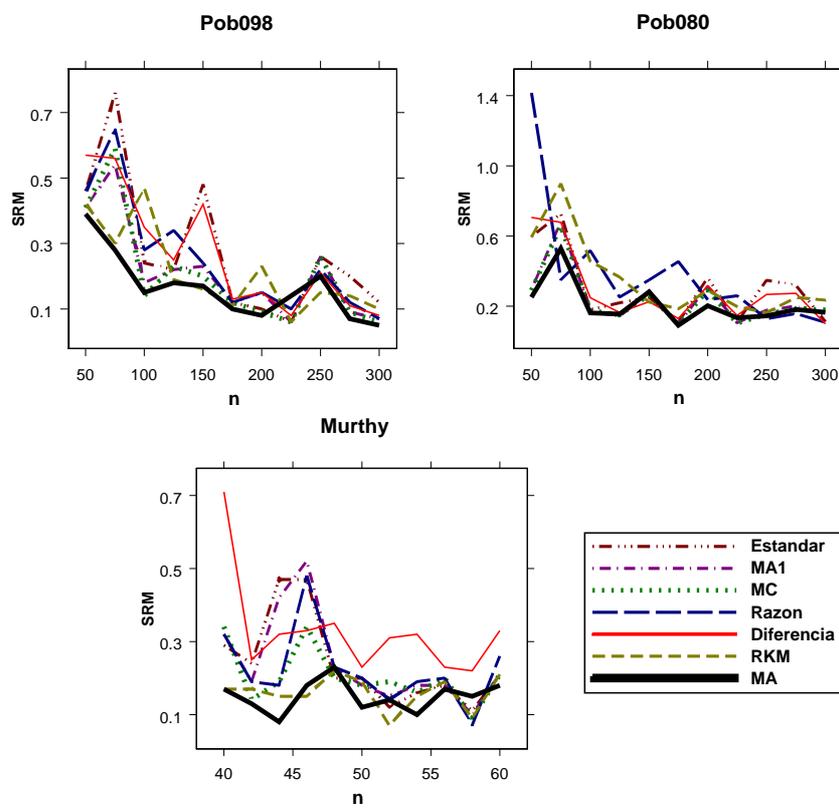
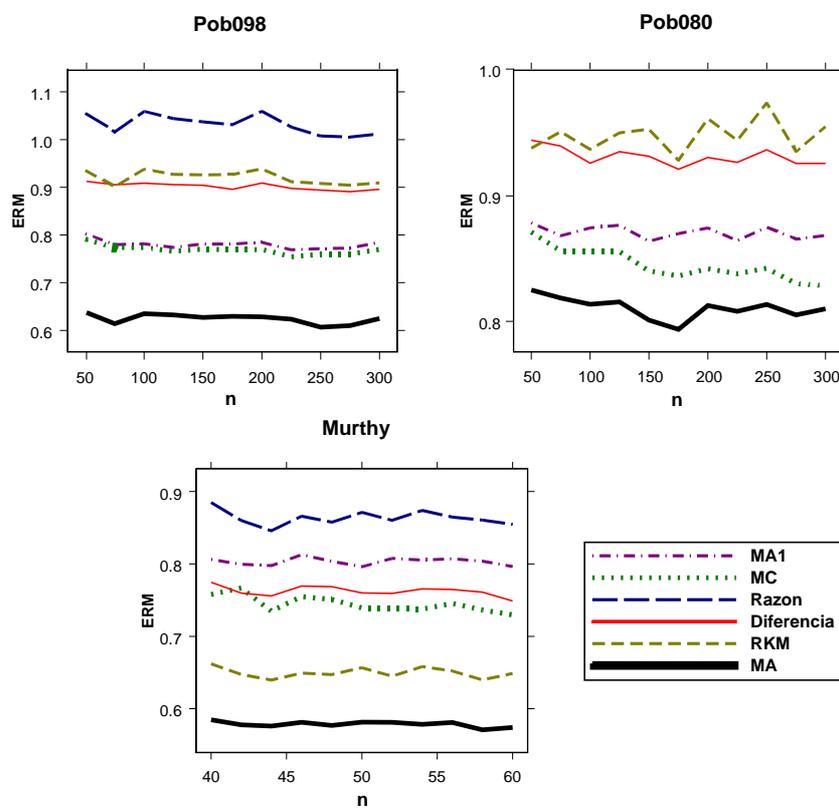


Figura 2.10: Eficiencia Relativa Media de distintos estimadores en las poblaciones Pob098, Pob080 y Murthy.



En este caso, el cumplimiento de los estimadores es medido mediante el Sesgo Relativo Medio (SRM) y la Eficiencia Relativa Media (ERM), dados respectivamente por

$$SRM = \frac{1}{9} \sum_{q=1}^9 |SR(t_q)| \quad ; \quad ERM = \sqrt{\frac{1}{9} \sum_{q=1}^9 ER(t_q)},$$

donde $SR(t)$ y $ER(t)$ están definidos en (2.108) y t_q es el q -ésimo decil para la variable de estudio.

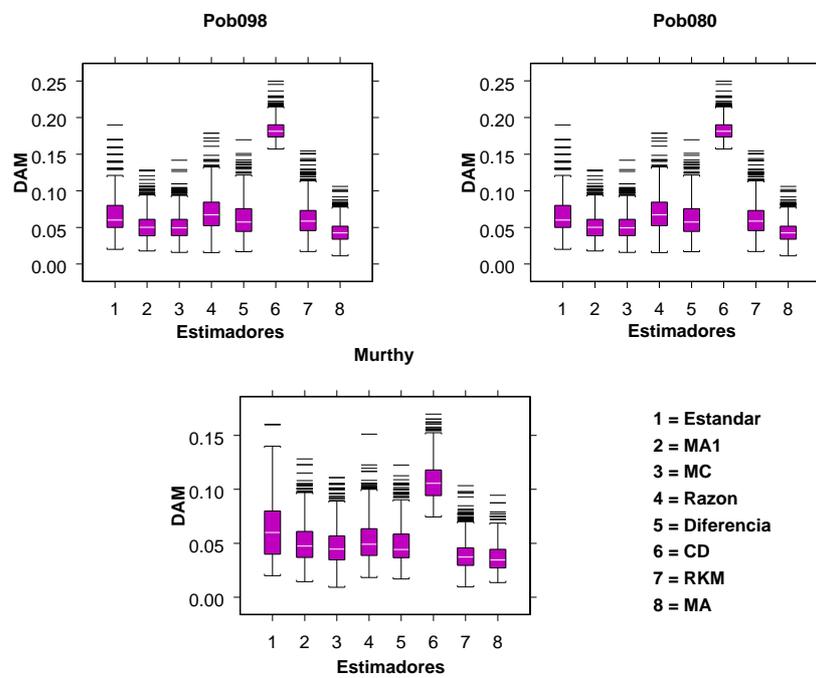
Consideramos también una medida global del cumplimiento de los estimadores a través de los 9 cuantiles para cada muestra obtenida de las $B = 1000$ simulaciones. Esta medida es la Desviación Absoluta Máxima (DAM) que está dada por: $DAM(b) = \max_q |\widehat{F}(t_q)_b - F(t_q)|$, para $b = 1, \dots, B$. Notamos que las medidas SRM , ERM y DAM se han usado también por Silva y Skinner (1995).

La Figura 2.9 muestra los valores SRM (%) para las tres poblaciones. Puede observarse que todos los estimadores exhiben valores SRM menores del 1% para las poblaciones Pob098 y Murthy. Asumiendo una relación más débil (Pob080), el estimador de tipo razón presenta el peor comportamiento (su SRM es sobre el 1.4%). En la mayoría de los casos, puede observarse que los valores SRM son decrecientes según el tamaño muestral y que el estimador $\widehat{F}_{MA}(t)$ presenta el menor sesgo.

Los valores ERM para las tres poblaciones están mostrados en la Figura 2.10. Estos resultados revelan que hay una razonable ganancia de eficiencia al usar $\widehat{F}_{MA}(t)$ con respecto a otros estimadores. $\widehat{F}_{MC}^{(1)}(t)$ muestra el segundo mejor comportamiento en las poblaciones Pob098 y Pob080, las cuales están basadas en un modelo lineal. A pesar de esta relación lineal entre y y \mathbf{x} , la pérdida de eficiencia de $\widehat{F}_{MC}^{(1)}(t)$ comparada con $\widehat{F}_{MA}(t)$ se debe al hecho de que el estimador $\widehat{F}_{MC}^{(1)}(t)$ usa un único valor fijo $t_0 = 0,5$, y éste es menos preciso cuando t está alejado de t_0 . En términos de ERM , $\widehat{F}_{CD}(t)$ muestra el peor comportamiento de todos los estimadores considerados. $\widehat{F}_{CD}(t)$ es más preciso cuando t está cercano a $Q_y(0,5)$, aunque este estimador sufre una considerable pérdida de eficiencia en cuantiles extremos (de bajo o alto orden).

La Figura 2.11 muestra los diagramas de cajas con bigotes de las distribuciones de los valores DAM obtenidos para las tres poblaciones y muestras de tamaño 100 para las poblaciones Pob098 y Pob080 y 50 para la población Murthy. Estos diagramas confirman el análisis anterior: $\widehat{F}_{CD}(t)$ presenta la máxima desviación absoluta mientras que $\widehat{F}_{MA}(t)$ muestra el mejor compor-

Figura 2.11: Diagramas de cajas con bigotes de las Desviaciones Absolutas Medias de distintos estimadores en las poblaciones Pob098 (con $n = 100$), Pob080 (con $n = 100$) y Murthy (con $n = 50$).



tamiento en todos los casos.

En todos los estudios (*ER*, *SR*, *SRM*, *ERM* y *DAM*), el estimador propuesto, $\hat{F}_{MA}(t)$, proporciona una buena mejoría sobre $\hat{F}_{MA1}(t)$, el cual usa un único punto t_0 . Esto confirma la ganancia en eficiencia de usar el vector \mathbf{t}_g , especialmente cuando t está alejado de t_0 .

Capítulo 3

Aportaciones a la estimación de cuantiles

3.1. Introducción

El problema de la estimación de la totales y medias poblacionales en presencia de variables auxiliares ha sido extensamente discutido en la literatura del muestreo de poblaciones finitas. Para el problema de la estimación de la mediana y otros cuantiles, la situación es bastante diferente y tan solo en la actualidad este problema está siendo discutido, debido en parte, al creciente interés de este tipo de medidas. Notamos que los distintos estimadores y métodos propuestos para la media y el total de una variable no tienen una extensión obvia al problema de la estimación de cuantiles.

Un ejemplo del uso de cuantiles y otras medidas relacionadas en muestreo de poblaciones finitas es el siguiente. Frecuentemente, los organismos nacionales de estadística y otras agencias se encuentran con variables, tales como ingresos, gastos, etc., que presentan distribuciones con una alta asimetría. Bajo estas circunstancias, la mediana resulta más apropiada que la media poblacional. De este modo, asumiendo datos de Encuestas Continuas de Presupuestos Familiares, los gobiernos de diferentes países obtienen numerosas medidas de pobreza, tal como la proporción de bajos ingresos, que dependen directamente de determinados cuantiles. Un ejemplo de este tipo de medidas viene dado por Eurostat (2000), en donde se define que un salario es bajo si éste está por debajo del 60 % del salario mediano mensual, es decir, el cuantil de orden $\beta = 0,5$ se emplea en Eurostat. A nivel nacional, el Instituto Nacional de Estadística y

sus correspondientes organismos autónomos, definen una medida similar para determinar el índice de pobreza, aunque en este caso la variable principal es el gasto producido en los hogares españoles. Otros estudios de tipo económico también usan cuantiles para estudiar la relación entre gastos en alimentación de los hogares y los correspondientes ingresos, análisis de salarios y gastos, impacto de varias características demográficas, calidad en la escuela, análisis de demanda, etc. Una extensa bibliografía sobre estas y otras aplicaciones en estudios de tipo económico puede consultarse en Koenker y Hallock (2001).

Al igual que para el caso de la estimación de parámetros lineales como medias o totales, las estimaciones serán más eficientes si se incorpora información auxiliar, altamente correlacionada con la variable de interés, en la etapa de estimación. En la estimación de cuantiles, existen dos grandes métodos que incorporan la información auxiliar de forma eficiente:

M1. Estimación de cuantiles indirectos: consiste en construir estimadores de tipo razón, diferencia o regresión, tal como se construyen para la media o el total. Ejemplos de este tipo de estimación pueden verse en Kuk y Mak (1989), Rueda, *et al.* (1998, 2003, 2004), etc. Notamos que para formular estos estimadores, se requiere conocer los cuantiles poblacionales de las variables auxiliares.

M2. Estimación a través de la función de distribución: La técnica habitual en muestreo de poblaciones finitas es invertir la función de distribución para obtener la estimación de un determinado cuantil. Se requiere, por tanto, usar eficientemente la información auxiliar en la etapa de estimación de la función de distribución. El inconveniente de esta técnica es que el estimador de la función de distribución debe ser una verdadera función de distribución para estimar cuantiles con mayor precisión. Aunque este hecho resulta imprescindible, existen varios estimadores en la literatura que no cumplen tal propiedad. Chambers y Dunstan (1986) fueron de los primeros investigadores en considerar la inversa de la función de distribución para obtener cuantiles. Otras importantes referencias son Rao *et al.* (1990), Wang y Dorfman (1996), Dorfman y Hall (1993), Kuo (1988), Silva y Skinner (1995).

Notamos que durante el desarrollo de este capítulo se tratarán exclusivamente con estimadores derivados del método *M2*, el cual es más usado por su calidad de estimación y eficiencia.

Los primeros trabajos relacionados con el problema de la estimación de parámetros de posición, como la mediana y los cuantiles se deben a Woodruff

(1952) donde se construyen intervalos de confianza bajo muestreo aleatorio simple. Posteriormente, Hill (1968) utiliza un enfoque bayesiano para la construcción de sus estimadores, mientras que Sendrask y Meyer (1978) se basan en un enfoque puramente probabilístico de distribución de estadísticos ordenados para muestreo aleatorio simple y estratificado. Pero los estimadores más eficientes y con mejores propiedades se desarrollan posteriormente bajo aproximaciones modelo-asistidas, basadas en el modelo y modelo-calibradas. También se han propuestos estimadores de cuantiles mediante intervalos de confianza basados en estimadores de razón, regresión y diferencia y usando información auxiliar multivariante (Rueda, Arcos y Artés, 1998, Rueda y Arcos, 2001, Rueda y Arcos, 2002a, Rueda y Arcos, 2002b).

En la literatura, los estimadores de cuantiles más conocidos son los siguientes. En primer lugar, citamos el estimador de Chambers y Dunstan (1986) para la función de distribución, el cual está basado en un modelo de superpoblación. La inversión directa de esta función puede usarse para la obtención de cuantiles. Siguiendo esta técnica, Rao *et al.* (1990) propusieron estimadores de tipo razón y diferencia usando una aproximación basada en el diseño. Kuk y Mak (1989) propusieron dos estimadores para los cuales solamente es necesario conocer a nivel poblacional el valor de la mediana de una variable auxiliar. Más recientemente, Rueda *et al.* (1998) y Rueda y Arcos (2001) propusieron intervalos de confianza para los cuantiles basados en estimadores de tipo razón y diferencia de la función de distribución. En Rueda *et al.* (2003, 2004) se plantea la estimación de cuantiles mediante estimadores de tipo diferencia usando cuantiles poblacionales del mismo orden de la variable auxiliar. La estimación de cuantiles usando técnicas recientes de estimación también ha sido investigada. Por ejemplo, Chen y Wu (2002) proponen la estimación de cuantiles usando la aproximación modelo-calibrada.

Existe otro gran número de estimadores de cuantiles propuestos para distintos diseños muestrales. Los estimadores más importantes se irán citando a lo largo del presente capítulo, en el cual se trata el problema de la estimación de cuantiles desde distintos enfoques. Por un lado, se desarrollan nuevos estimadores en diseños muestrales más complejos, y por otro, se proponen estimadores asumiendo el reciente método de verosimilitud empírica.

Para formular la mayoría de los estimadores de cuantiles, ya sean a través del método $M1$ o del método $M2$, es necesario conocer los valores poblacionales de las variables auxiliares, aunque esto es poco usual en la práctica. La solución a este problema se trata en la Sección 3.2 mediante el uso del muestreo bifásico (véase Apéndice B), en el cual la información auxiliar poblacional puede estimarse usando la muestra de la primera fase. Por tanto, en esta sección se

proponen estimadores de cuantiles en muestreo bifásico y asumiendo que las unidades muestrales se extraen mediante métodos de muestreo con probabilidades desiguales en cada una de las dos fases. La eficiencia de estos estimadores puede mejorarse si se usa un muestreo estratificado en la primera fase. Asumiendo este último diseño muestral, denominado muestreo bifásico aplicado a la estratificación, se comprueba que los estimadores propuestos pueden llegar a ser más precisos con respecto a otros existentes en la literatura.

Por otro lado, en la Sección 3.3 se plantean nuevos estimadores de cuantiles bajo muestreo en ocasiones sucesivas (véase Apéndice B para un mayor detalle sobre este diseño). El método propuesto se basa en el caso de que las muestras son seleccionadas mediante muestreos probabilísticos con probabilidades desiguales (por ejemplo, con unidades proporcionales al tamaño). Notamos que éste es el caso de los organismos nacionales y agencias de estadística que realizan encuestas continuas a lo largo del tiempo. Dentro de esta sección también se definen estimadores de cuantiles basados en múltiples variables auxiliares. La introducción de tal información proporciona un marco de estimación apropiado que permite obtener estimadores más precisos. El comportamiento de todos los estimadores propuestos se analiza desde el punto de vista teórico (mediante aproximaciones asintóticas), y desde una perspectiva empírica (analizando los resultados obtenidos a partir de una serie de poblaciones).

Para cerrar este capítulo, en la Sección 3.4 se proponen estimadores para cuantiles asumiendo el método de verosimilitud empírica, expuesto con detalle en el capítulo anterior. Los estimadores propuestos usan de manera eficiente la información auxiliar, lo que se traduce en una mejoría de la precisión. Esta precisión de los estimadores propuestos se ha evaluado para el cálculo de algunas medidas de pobreza oficiales, las cuales dependen de forma directa de cuantiles. Este estudio se ha llevado a cabo asumiendo distintos estimadores de cuantiles. Los resultados obtenidos reflejan que los estimadores propuestos proporcionan estimaciones más precisas para las medidas de pobreza involucradas en tal estudio.

3.2. Estimadores bajo muestreo bifásico

En esta sección se resuelve el problema de la estimación de cuantiles bajo muestreo en dos fases o muestreo bifásico con diseños muestrales arbitrarios en cada una de las dos fases. Se proponen varios estimadores de tipo directo, razón y exponencial que proporcionan estimaciones óptimas para un determinado cuantil. Se analizan propiedades importantes de estos estimadores

tales como la insesgadez, estimación de varianzas, etc. Como caso particular, se investiga también el muestreo bifásico aplicado a la estratificación, diseño muestral que ofrece importantes ganancias en eficiencia debido a los beneficios que produce el muestreo estratificado. Todas estas propiedades se ven desde un punto teórico, aunque el análisis de los estimadores se completa con un estudio empírico llevado a cabo para los cuartiles y bajos distintos diseños muestrales con probabilidades desiguales. Este estudio refleja que los estimadores propuestos mejoran a otros estimadores conocidos en términos de sesgo y eficiencia relativa.

Notamos que la mayor ventaja al usar muestreo bifásico es una alta ganancia en precisión sin un sustancial incremento en costes. De hecho, este diseño muestral se usa frecuentemente en numerosas encuestas por razones de coste y eficiencia.

3.2.1. Introducción

Para el problema de la estimación de un determinado parámetro en muestreo de poblaciones finitas, la información auxiliar juega un papel muy importante en la precisión de los estimadores. La mayoría de los estimadores basados en información auxiliar se basan en el conocimiento a nivel poblacional de las variables auxiliares. En la práctica, esta cantidad no tiene porque ser conocida. De hecho, son muy poco frecuentes las encuestas que disponen de esta información, por lo que resulta imposible obtener estos estimadores basados en información auxiliar. Una alternativa es estimar los parámetros poblacionales que usan los estimadores, aunque esto conlleva a importantes errores en la etapa de la estimación de la varianza (véase Berger, Muñoz y Rancourt, 2006). Bajo esta situación, el uso de un muestreo bifásico es la técnica más apropiada para resolver este problema.

Por tanto, el muestreo bifásico es una herramienta útil para aquellas investigaciones en las cuales no existe conocimiento previo de las variables auxiliares a nivel poblacional. Otro punto a favor del muestreo bifásico es la creación de un esquema importante de información que permite la selección probabilística de sub-muestras. Para una mayor profundización sobre el muestreo bifásico en la estimación de medias o totales puede consultarse, por ejemplo, Särndal *et al.* (1992), Fernández y Mayor (1994) y Artés y García (2002).

En lo que respecta a la estimación de cuantiles en muestreo bifásico, los primeros autores en realizar investigaciones en este sentido fueron Singh *et al.* (2001), Singh (2003) y Allen *et al.* (2002) para el problema de la estimación

de la mediana poblacional. Estos trabajos fueron desarrollados exclusivamente para muestreo aleatorio simple. Con el fin de completar estos estudios, en esta sección se proponen numerosos estimadores para un determinado cuantil cuando se lleva a cabo un muestreo bifásico con diseños muestrales arbitrarios en cada una de las dos fases.

A continuación se describe brevemente en que consiste un muestreo bifásico. Suponemos que tenemos una población U compuesta por N unidades de la que se extrae en una primera fase una muestra, s' , de tamaño, n' , bastante grande y de bajo costo, según cierto criterio muestral, d_1 , tal que $p_{d1}(s')$ será la probabilidad de que s' sea seleccionada y donde las correspondientes probabilidades de inclusión de primer y segundo orden se denotan, respectivamente, como π'_i y π'_{ij} para i y $j \in U$. En esta muestra, una o varias variables auxiliares pueden ser recogidas fácilmente, es decir, dicha muestra permite obtener la información auxiliar necesaria para todo el proceso. Dada s' , una segunda muestra s de tamaño n es seleccionada en la segunda fase mediante un diseño d_2 , tal que $p(s/s')$ es la probabilidad condicional de escoger s . Las probabilidades de inclusión bajo este diseño se denotan como $\pi_{i/s'}$ y $\pi_{ij/s'}$. Notamos que $\Delta'_{ij} = \pi'_{ij} - \pi'_i \pi'_j$ y $\Delta^{s'}_{ij} = \pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}$. Véase también Apéndice B para un mayor detalle sobre el muestreo bifásico.

3.2.2. Estimadores propuestos

Sin usar ningún tipo de información auxiliar, el candidato natural para estimar el cuantil β es $\widehat{Q}_y(\beta) = \inf\{t \mid \widehat{F}_{HTy}(t) \geq \beta\} = \widehat{F}_{HTy}^{-1}(\beta)$, donde

$$\widehat{F}_{HTy}(t) = \frac{1}{N} \sum_{i \in s} \frac{\delta(t - y_i)}{\pi_i}$$

es el estimador de tipo Horvitz y Thompson (1952) de $F_y(t)$, y las probabilidades de inclusión están dadas por $\pi_i = \sum_{s' \ni i} p_{d1}(s') \pi_{i/s'}$.

Como puede observarse, para determinar π_i se deben conocer las probabilidades $\pi_{i/s'}$ para cada s' , las cuales no se conocen generalmente porque $\pi_{i/s'}$ pueden depender del diseño de la primera fase, por ejemplo si la muestra de la segunda fase es diseñada mediante un muestreo proporcional a una variable auxiliar.

Notamos que el estimador de tipo Horvitz-Thompson para la media poblacional tampoco puede obtenerse en la práctica bajo este muestreo. Por esta razón, Särndal *et al.* (1992) propusieron el uso de estimadores π^* . Usando es-

ta idea, se definen las cantidades $\pi_i^* = \pi'_i \pi_{i/s'}$ y $\pi_{ij}^* = \pi'_{ij} \pi_{ij/s'}$, que permiten definir el π^* -estimador de la función de distribución como

$$\widehat{F}_{HTy}^*(t) = \frac{1}{N} \sum_{i \in s} \frac{\delta(t - y_i)}{\pi_i^*},$$

y así, el estimador directo propuesto para un cuantil β esta dado por

$$\widehat{Q}_y^*(\beta) = \widehat{F}_{HTy}^{*-1}(\beta). \quad (3.1)$$

Notamos que $\widehat{Q}_y^*(\beta)$ no coincide generalmente con el estimador $\widehat{Q}_y(\beta)$ excepto en casos excepcionales, aunque la principal ventaja del estimador directo propuesto sobre el estándar comentado es su aplicabilidad para cualesquiera que sean los diseños muestrales usados en cada fase.

El estimador (3.1) se ha definido sin usar ninguna información auxiliar. Si esta información está disponible, el uso de estimadores indirectos nos puede ayudar a obtener estimaciones más precisas para los cuantiles en muestreo bifásico. De este modo, el siguiente paso es definir una clase de estimadores que usen información auxiliar. En primer lugar mostraremos los principales antecedentes relacionados con el tema que nos ocupa.

Asumiendo muestreo aleatorio simple y que la mediana de la variable x es conocida, Kuk y Mak (1989) propusieron el siguiente estimador de tipo razón para la mediana

$$\widehat{Q}_y^r(0,5) = \widehat{Q}_y(0,5) \frac{Q_x(0,5)}{\widehat{Q}_x(0,5)}.$$

Además, Kuk y Mak (1989) propusieron otros estimadores de cuantiles bajo muestreo aleatorio simple llamados estimadores de posición y de estratificación, pero la extensión de cualquiera de ellos a otros diseños muestrales más complejos no ha sido posible

Rueda *et al.* (2003, 2004) propusieron para cualquier diseño muestral d y para cualquier β , métodos de diferencia y exponenciales para estimar un cuantil β . Singh *et al.* (2001) sugirieron estimadores de tipo razón, regresión, posición y estratificación de la mediana cuando la muestra es seleccionada en dos fases y usando muestreo aleatorio simple en cada una de ellas. Bajo muestreo bifásico y muestreo aleatorio simple en cada fase, Allen *et al.* (2002) propusieron dos clases de estimadores para la mediana poblacional. Estos estimadores usan la información proporcionada por dos variables auxiliares, x y z , donde se asume que la mediana de z es conocida.

A continuación se presenta una clase de estimadores para cuantiles poblacionales finitos cuando las muestras en ambas fases son seleccionadas bajo cualquier esquema de muestreo:

$$\widehat{Q}_y^{\mathcal{H}}(\beta) = H(\widehat{Q}_y^*(\beta), t^*), \quad (3.2)$$

donde $t^* = \widehat{Q}_x^*(\beta)/\widehat{Q}'_x(\beta)$, y $\widehat{Q}'_x(\beta)$ es el estimador de $Q_x(\beta)$ basado en la muestra de la primera fase, esto es, $\widehat{Q}'_x(\beta) = \inf\{t \mid \widehat{F}'_{HTx}(t) \geq \beta\}$, donde

$$\widehat{F}'_{HTx}(t) = \frac{1}{N} \sum_{i \in s'} \frac{\delta(t - x_i)}{\pi'_i}.$$

La función H satisface las siguientes condiciones:

(C3.1). Asume valores en un subconjunto convexo cerrado $\mathcal{C} \subset \mathbb{R}_2$, el cual contiene el punto $(Q_y(\beta), 1)$.

(C3.2). H es una función continua en \mathcal{C} , tal que $H(Q_y(\beta), 1) = Q_y(\beta)$.

(C3.3). Las primeras y segundas derivadas parciales de H existen y son continuas en \mathcal{C} , con

$$H_{10}(Q_y(\beta), 1) = \left. \frac{\partial H(q, t)}{\partial q} \right|_{(q,t)=(Q_y(\beta),1)} = 1.$$

Un caso particular dentro de la clase general de estimadores \mathcal{H} es el estimador tipo razón dado por:

$$\widehat{Q}_{yr}^*(\beta) = \widehat{Q}_y^*(\beta) \frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)},$$

el cual corresponde a la elección $H(q, t) = q/t$.

Otro estimador para el cuantil β , llamado el estimador exponencial, está dado por:

$$\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^\alpha,$$

siendo α una constante fija. Este estimador también se encuentra dentro de la clase \mathcal{H} , puesto que se corresponde con la elección $H(q, t) = q/t^\alpha$. Notamos que estos estimadores se han definido en Rueda et al (2006a)

Nota 3.1 Si $\alpha = 0$, entonces $\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_y^*(\beta)$, esto es, $\widehat{Q}_{ye}^*(\beta)$ coincide con el estimador π^* . Por otro lado, si $\alpha = 1$, entonces $\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_{yr}^*(\beta)$. Por último, puede comprobarse que si $\alpha = -1$, entonces $\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_{yp}^*(\beta)$, el cual puede definirse como un estimador producto.

Nota 3.2 Bajo muestreo aleatorio simple en cada fase y $\beta = 0,5$, los estimadores propuestos $\widehat{Q}_{yr}^*(\beta)$ y $\widehat{Q}_{ye}^*(\beta)$ se corresponden, respectivamente, con los estimadores $\widehat{M}_y^{(a)}$ y $\widehat{M}_y^{(b)}$ propuestos por Singh et al. (2001).

3.2.3. Propiedades teóricas

En este apartado se estudian las principales propiedades del estimador $\widehat{Q}_y^*(\beta)$ y de los estimadores basados en la clase \mathcal{H} . Debido a que estos estimadores no son funciones continuas, serán necesarias aproximaciones lineales.

Teorema 3.1 El estimador $\widehat{Q}_y^*(\beta)$ es asintóticamente insesgado para $Q_y(\beta)$

Demostración

En primer lugar, el estimador $\widehat{Q}_y^*(\beta)$ puede expresarse asintóticamente como una función lineal de la función de distribución estimada y evaluada en el punto $Q_y(\beta)$ mediante la representación Bahadur (véase, por ejemplo, Chambers y Dunstan, 1986, Kuk y Mak, 1989, Chen y Wu, 2002, etc):

$$\widehat{Q}_y^*(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))}(\beta - \widehat{F}_{HTy}^*(Q_y(\beta))) + O(n^{-1/2}), \quad (3.3)$$

donde $f_y(\cdot)$ denota la derivada del valor límite de $F_y(\cdot)$ cuando $N \rightarrow \infty$.

Además, es sabido que el estimador $\widehat{F}_{HTy}^*(t)$ es insesgado de $F(t)$. En consecuencia, se tiene que $E(\beta - \widehat{F}_{HTy}^*(Q_y(\beta))) = 0$ y basándose en la ecuación (3.3), puede verse fácilmente que $E(\widehat{Q}_y^*(\beta)) = Q_y(\beta) + O(n^{-1/2})$, esto es, el estimador $\widehat{Q}_y^*(\beta)$ es asintóticamente insesgado de $Q_y(\beta)$. \square

Corolario 3.1 De la ecuación (3.3) y puesto que $\widehat{Q}_y^*(\beta)$ es asintóticamente insesgado de $Q_y(\beta)$, la varianza asintótica de $\widehat{Q}_y^*(\beta)$, al primer grado de aproximación, está dada por:

$$V(\widehat{Q}_y^*(\beta)) = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in U} (\pi'_{ij} - \pi'_i \pi'_j) \frac{\delta(Q_y(\beta) - y_i)}{\pi'_i} \frac{\delta(Q_y(\beta) - y_j)}{\pi'_j} + \right.$$

$$+ E_{d1} \left[\sum_{i,j \in s'} (\pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}) \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_y(\beta) - y_j)}{\pi_j^*} \right].$$

Corolario 3.2 *Un estimador insesgado de $V(\widehat{Q}_y^*(\beta))$ está dado por*

$$\begin{aligned} \widehat{V}(\widehat{Q}_y^*(\beta)) &= \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in s} \frac{\pi'_{ij} - \pi'_i \pi'_j}{\pi_{ij}^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i'} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j'} + \right. \\ &\quad \left. + \sum_{i,j \in s} \frac{\pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j^*} \right). \end{aligned}$$

En la práctica, la cantidad $f_y(Q_y(\beta))$ es desconocida. Un valor aproximado de $f_y(Q_y(\beta))$ puede obtenerse aplicando métodos estándares tal como el kernel (Silverman, 1986). Este estimador de la varianza no depende de esperanzas relacionadas con el diseño de la primera fase, haciendo posible su cálculo en la práctica.

Teorema 3.2 *Cualquier estimador dentro de la clase \mathcal{H} es asintóticamente insesgado para $Q_y(\beta)$.*

Demostración

Para obtener este resultado nos basaremos en las siguientes aproximaciones lineales:

$$\widehat{Q}_y^*(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))} (\beta - \widehat{F}_{HTy}^*(Q_y(\beta))) + O(n^{-1/2}),$$

$$\widehat{Q}_x^*(\beta) - Q_x(\beta) = \frac{1}{f_x(Q_x(\beta))} (\beta - \widehat{F}_{HTx}^*(Q_x(\beta))) + O(n^{-1/2}),$$

$$\widehat{Q}_x'(\beta) - Q_x(\beta) = \frac{1}{f_x(Q_x(\beta))} (\beta - \widehat{F}'_{HTx}(Q_x(\beta))) + O(n'^{-1/2}),$$

y usando la expansión de la serie de Taylor de primer orden para H sobre el punto $(Q_y(\beta), 1)$:

$$\begin{aligned} \widehat{Q}_y^{\mathcal{H}}(\beta) &= H((Q_y(\beta), 1)) + \left(\widehat{Q}_y^*(\beta) - Q_y(\beta) \right) H_{10}(Q_y(\beta), 1) + \\ &\quad + (t-1)H_{01}(Q_y(\beta), 1) + O(n^{-1}), \end{aligned} \quad (3.4)$$

donde H_{10} y H_{01} denotan las derivadas parciales de primer orden de H con respecto a q y t , respectivamente. Como $\widehat{F}_{HTy}^*(t)$ y $\widehat{F}_{HTx}^*(t)$ son estimadores insesgados de $F_y(t)$ y $F_x(t)$, respectivamente, puede observarse que cualquier estimador en \mathcal{H} será asintóticamente insesgado para $Q_y(\beta)$. \square

Para obtener las expresiones asintóticas de las varianzas, consideraremos la expansión de la serie de Taylor dada en (3.4), que da lugar a la expresión:

$$\widehat{Q}_y^{\mathcal{H}}(\beta) - Q_y(\beta) = \left(\widehat{Q}_y^*(\beta) - Q_y(\beta) \right) + \frac{\widehat{Q}_x^*(\beta)}{\widehat{Q}_x'(\beta)} H_{01}(Q_y(\beta), 1) + O(n^{-1}).$$

Desarrollando se obtiene

$$\widehat{Q}_y^{\mathcal{H}}(\beta) - Q_y(\beta) \simeq Q_y(\beta)e_0 + (e_1 - e_2)H_{01}(Q_y(\beta), 1) - e_2(e_1 - e_2)H_{01}(Q_y(\beta), 1), \quad (3.5)$$

donde:

$$e_0 = \frac{\widehat{Q}_y^*(\beta)}{Q_y(\beta)} - 1, \quad e_1 = \frac{\widehat{Q}_x^*(\beta)}{Q_x(\beta)} - 1 \quad \text{y} \quad e_2 = \frac{\widehat{Q}_x'(\beta)}{Q_x'(\beta)} - 1.$$

Introduciendo varianzas en (3.5) y bajo una aproximación de primer orden, se llega a la expresión:

$$V(\widehat{Q}_y^{\mathcal{H}}(\beta)) = Q_y(\beta)^2 V(e_0) + H_{01}(Q_y(\beta), 1)^2 V(e_1 - e_2) + 2H_{01}(Q_y(\beta), 1) Cov(e_0, e_1 - e_2).$$

Por otro lado, bajo muestreo bifásico:

$$V(\widehat{Q}_y^{\mathcal{H}}(\beta)) = E_{d1} V(\widehat{Q}_y^{\mathcal{H}}(\beta)/s') + V_{d1} E(\widehat{Q}_y^{\mathcal{H}}(\beta)/s')$$

refleja la variación debida a cada una de las dos fases de muestreo. Usando las propiedades conocidas del estimador de Horvitz-Thompson y su varianza, se obtiene

$$V_{d1} E(\widehat{Q}_y^{\mathcal{H}}(\beta)/s') = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi'_i} \frac{\delta(Q_y(\beta) - y_j)}{\pi'_j} \right)$$

y

$$E_{d1} V(\widehat{Q}_y^{\mathcal{H}}(\beta)/s') = E_{d1} \left(\frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \sum_{i,j \in s'} \Delta^{s'}_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_y(\beta) - y_j)}{\pi_j^*} \right) + \\ + \frac{H_{01}^2(Q_y(\beta), 1)}{Q_x^2(\beta)} \frac{1}{N^2} \frac{1}{f_x^2(Q_x(\beta))} \sum_{i,j \in s'} \Delta^{s'}_{ij} \frac{\delta(Q_x(\beta) - x_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*} +$$

$$+2 \frac{H_{01}(Q_y(\beta), 1)}{Q_x(\beta)} \frac{1}{N^2} \frac{1}{f_y(Q_y(\beta)) f_x(Q_x(\beta))} \sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*},$$

donde $\Delta'_{ij} = \pi'_{ij} - \pi'_i \pi'_j$ y $\Delta^{s'}_{ij} = \pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}$. Esta expresión no puede obtenerse en la práctica, así que para ello

$$\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi'_i} \frac{\delta(Q_y(\beta) - y_j)}{\pi'_j}$$

se estima por

$$\sum_{i,j \in s} \frac{\Delta'_{ij}}{\pi_{ij}^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j^*},$$

y

$$E_{d1} \left(\sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_y(\beta) - y_j)}{\pi_j^*} \right)$$

por

$$\sum_{i,j \in s} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j^*}.$$

Las funciones $f_x(Q_x(\beta))$ y $f_y(Q_y(\beta))$ pueden calcularse siguiendo Silverman (1986).

Las varianzas asintóticas de los estimadores de tipo razón, producto y exponencial se derivan a partir de $H(q, t) = q/t$, $H(q, t) = q * t$ y $H(q, t) = q/t^\alpha$, respectivamente.

Una vez que la clase y sus propiedades principales han sido definidas, el siguiente paso es obtener el estimador óptimo en la clase $\widehat{Q}_{ye}^*(\beta)$. La idea de optimalidad se define en el sentido de minimizar la varianza asintótica de estos estimadores.

El valor óptimo de α está dado por

$$\alpha_{opt} = \frac{Q_x(\beta) \text{Cov}(\widehat{Q}_y(\beta), \widehat{Q}_x(\beta)) - \text{Cov}(\widehat{Q}_y(\beta), \widehat{Q}'_x(\beta))}{Q_y(\beta) V(\widehat{Q}_x(\beta)) + \widehat{Q}'_x(\beta) - 2\text{Cov}(\widehat{Q}_x(\beta), \widehat{Q}'_x(\beta))}.$$

Usando las propiedades de muestreo bifásico, se obtiene:

$$\alpha_{opt} = \frac{Q_x(\beta) f_x(Q_x(\beta))}{Q_y(\beta) f_y(Q_y(\beta))} \frac{E_{d1} \left(\sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*} \right)}{E_{d1} \left(\sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_x(\beta) - x_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*} \right)},$$

y el estimador óptimo está dado por

$$\widehat{Q}_y^{\alpha_{opt}}(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\alpha_{opt}}.$$

Puede verse que:

$$\begin{aligned} V(\widehat{Q}_y^{\mathcal{H}}(\beta)) &\geq V(\widehat{Q}_y^{\alpha_{opt}}(\beta)) = V(\widehat{Q}_y(\beta)) - K_1 = \\ &= V(\widehat{Q}_y(\beta)) - \frac{(Cov(\widehat{Q}_y(\beta), \widehat{Q}_x(\beta)) - Cov(\widehat{Q}_y(\beta), \widehat{Q}'_x(\beta)))^2}{V(\widehat{Q}_x(\beta)) + \widehat{Q}'_x(\beta) - 2Cov(\widehat{Q}_x(\beta), \widehat{Q}'_x(\beta))}, \end{aligned} \quad (3.6)$$

esto es, el valor más bajo de la varianza de $\widehat{Q}_y^{\mathcal{H}}(\beta)$ está dado por el estimador exponencial con α_{opt} .

La ecuación (3.6) demuestra que el estimador propuesto $\widehat{Q}_y^{\alpha_{opt}}(\beta)$ es siempre más eficiente que el estimador más simple $\widehat{Q}_y(\beta)$. Puede observarse que K_1 es la cantidad que se reduce de varianza cuando se usa el estimador exponencial con el valor óptimo de α en lugar de usar el estimador $\widehat{Q}_y(\beta)$.

En la práctica, el valor de α es desconocido. Sin embargo, los datos muestrales podrán usarse para obtener un estimador para este parámetro. Un posible estimador para el valor óptimo de α está dado por

$$\widehat{\alpha} = \frac{\widehat{Q}_x^*(\beta) f_x(Q_x(\beta))}{\widehat{Q}_y^*(\beta) f_y(Q_y(\beta))} \frac{\sum_{i,j \in s} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*}}{\sum_{i,j \in s} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(Q_x(\beta) - x_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*}}. \quad (3.7)$$

De este modo, se puede definir un estimador óptimo para el cuantil β como:

$$\widehat{Q}_y^{\widehat{\alpha}}(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\widehat{\alpha}}.$$

Siguiendo el procedimiento discutido en Allen *et al.* (2002), puede demostrarse que $E(\widehat{Q}_y^{\widehat{\alpha}}(\beta)) = Q_y(\beta) + o(n^{-1})$ y al primer grado de aproximación, $V(\widehat{Q}_y^{\widehat{\alpha}}(\beta)) = V(\widehat{Q}_y^{\alpha_{opt}}(\beta))$, esto es, los estimadores $\widehat{Q}_y^{\widehat{\alpha}}(\beta)$ y $\widehat{Q}_y^{\alpha_{opt}}(\beta)$ son asintóticamente equivalentes.

3.2.4. Propiedades empíricas

Se han propuesto varios estimadores para cuantiles en muestreo bifásico cuando las muestras en ambas fases se seleccionan con probabilidades desiguales. A continuación se lleva a cabo un estudio de simulación con el objetivo de observar el comportamiento de estos estimadores y destacar el más eficiente entre ellos. En este estudio se han considerado las poblaciones Fam1500 y Counties (véase Apéndice A).

Se han generado 1000 muestras independientes bajo diferentes métodos de muestreo en cada fase. El tamaño muestral en la primera fase, n' , se ha fijado en 150, mientras que el tamaño de la muestra de la segunda fase, n , varía entre 10 y 100. Los casos considerados son los siguientes:

1. (*Mas.Midzuno*): Las muestras en la primera fase han sido seleccionadas mediante muestreo aleatorio simple de tamaño n' , mientras que las muestras de la segunda fase se han tomado mediante el método de Midzuno (véase Apéndice B). Las probabilidades de inclusión en este caso vienen dadas por:

$$\pi'_i = \frac{n'}{N}, \quad \pi_{i/s'} = \frac{n' - n}{n' - 1} \frac{x_i}{\sum_{j \in s'} x_j} + \frac{n - 1}{n' - 1} \quad \rightarrow \quad \pi_i^* = \pi'_i \pi_{i/s'}.$$

2. (*Mas.Poisson*): En la primera fase se usa muestreo aleatorio simple de tamaño n' , y las muestras de la segunda fase son seleccionadas mediante el método de Poisson (véase Apéndice B), de modo que las probabilidades de inclusión están dadas por:

$$\pi'_i = \frac{n'}{N}, \quad \pi_{i/s'} = n \frac{x_i}{\sum_{j \in s'} x_j} \quad \rightarrow \quad \pi_i^* = \pi'_i \pi_{i/s'}.$$

El cumplimiento de los estimadores propuestos en muestreo bifásico para un determinado cuantil se evalúa para los tres cuartiles, $\beta = 0,25, 0,50, 0,75$, en términos de Sesgo Relativo (%) (SR) y Eficiencia Relativa (ER) mediante aproximaciones Monte Carlo derivadas de $B = 1000$ muestras independientes. Estas medidas vienen dadas por:

$$SR_i = 100 \times \frac{1}{B} \sum_{b=1}^B \frac{\widehat{Q}_y^i(\beta)_b - Q_y(\beta)}{Q_y(\beta)} \quad ; \quad ER_i = \frac{ECM[\widehat{Q}_y^i(\beta)]}{ECM[\widehat{Q}_y^*(\beta)]},$$

donde b indica la b -ésima simulación y $\widehat{Q}_y^i(\beta)$ denota el i -ésimo estimador propuesto, con

- $\widehat{Q}_y^1(\beta) = \widehat{Q}_y^*(\beta) \frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)},$
- $\widehat{Q}_y^2(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\widehat{\alpha}},$ donde $\widehat{\alpha}$ está dado en (3.7),
- $\widehat{Q}_y^3(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\alpha_{opt}}.$

$ECM[\widehat{Q}_y^i(\beta)] = B^{-1} \sum_{b=1}^B [\widehat{Q}_y^i(\beta)_b - Q_y(\beta)]^2$ es el Error Cuadrático Medio empírico y $ECM[\widehat{Q}_y^*(\beta)]$ se define análogamente para $\widehat{Q}_y^*(\beta)$, el estimador directo definido en (3.1). Se recuerda que este estimador no usa información auxiliar.

Las Figuras 3.1, ..., 3.4 representan la eficiencia relativa para los estimadores $\widehat{Q}_y^1(\beta)$, $\widehat{Q}_y^2(\beta)$ y $\widehat{Q}_y^3(\beta)$ en las diferentes poblaciones y bajo los diseños *Mas.Midzuno* y *Mas.Poisson*. Estas figuras muestran el comportamiento de los estimadores cuando aumenta el tamaño muestral en la segunda fase, mientras que el tamaño muestral de la primera fase permanece constante.

Cuando existe alta correlación lineal entre y y la variable auxiliar, todos los estimadores son más eficientes que el estimador $\widehat{Q}_y^*(\beta)$, mostrado con líneas horizontales. La ganancia en eficiencia relativa decrece cuando aumenta el tamaño muestral en la muestra de la segunda fase. Este resultado resulta lógico porque si el tamaño muestral en la segunda fase es pequeño, entonces la muestra tendrá menos información de la variable y , y el estimador $\widehat{Q}_y^*(\beta)$ presentará mayor grado de error, mientras que los estimadores de tipo razón y exponencial son más eficientes porque usan más información. Cuando n incrementa, $\widehat{Q}_y^*(\beta)$ obtiene mejores estimaciones y más cercanas a las estimaciones de los estimadores de tipo razón y exponencial.

$\widehat{Q}_y^3(\beta)$ es el estimador más eficiente en la mayoría de los casos. Este resultado era deseable porque este estimador es asintóticamente óptimo en la clase (3.2). Sin embargo, el estimador $\widehat{Q}_y^2(\beta)$ presenta valores bastantes similares y no depende de valores desconocidos. Se observa que $\widehat{Q}_y^1(\beta)$ es el estimador menos eficiente de entre los estimadores propuestos. Cuando la relación lineal entre las variable es más débil, $\widehat{Q}_y^1(\beta)$ es incluso menos eficiente que el estimador directo, mientras que $\widehat{Q}_y^2(\beta)$ y $\widehat{Q}_y^3(\beta)$ continúan teniendo un buen comportamiento. En resumen, el uso del estimador exponencial mejora las estimaciones, especialmente si la relación lineal entre las variables es débil.

Figura 3.1: Eficiencia Relativa para la población Fam1500 y bajo el diseño muestral *Mas.Midzuno*. $n' = 150$.

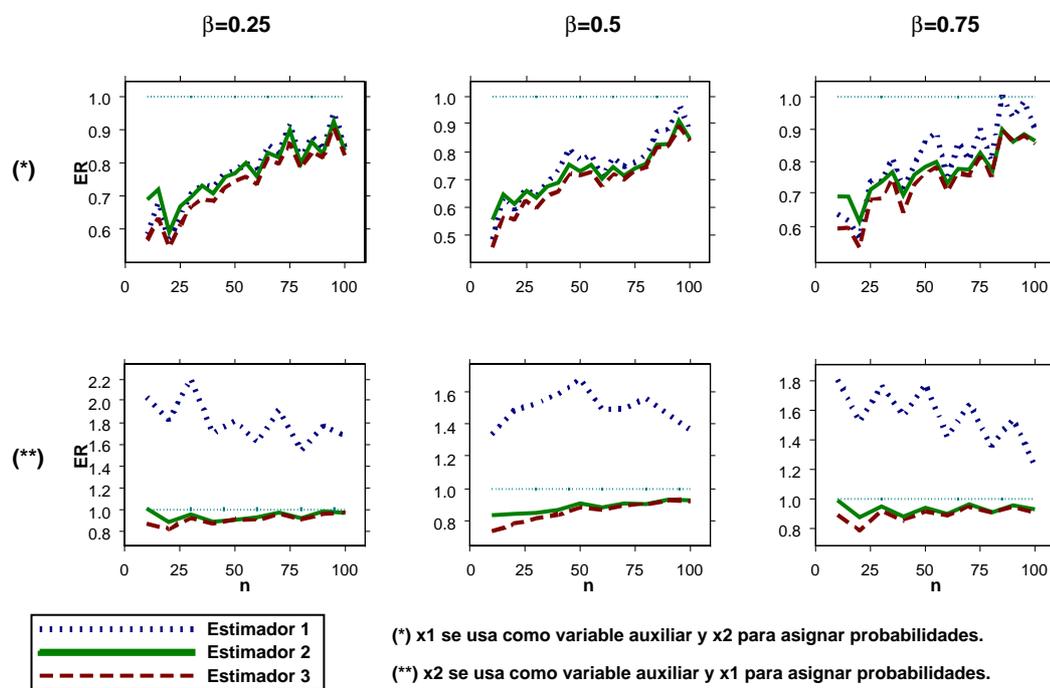


Figura 3.2: Eficiencia Relativa para la población Fam1500 y bajo el diseño muestral *Mas.Poisson*. $n' = 150$.

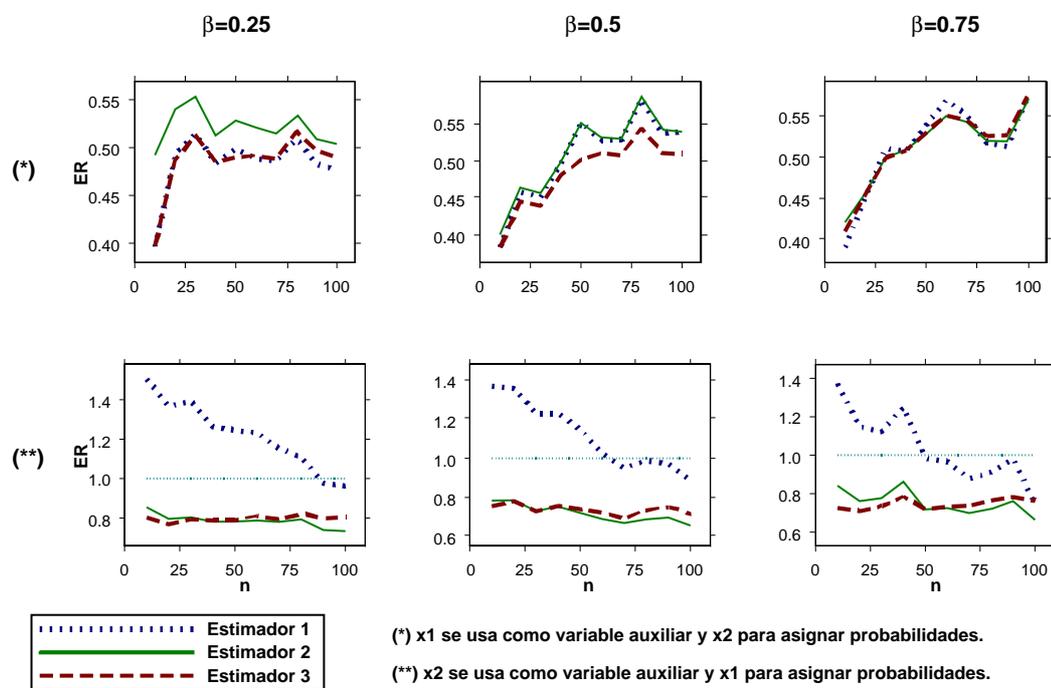


Figura 3.3: Eficiencia Relativa para la población Counties y bajo el diseño muestral *Mas.Midzuno*. $n' = 150$.

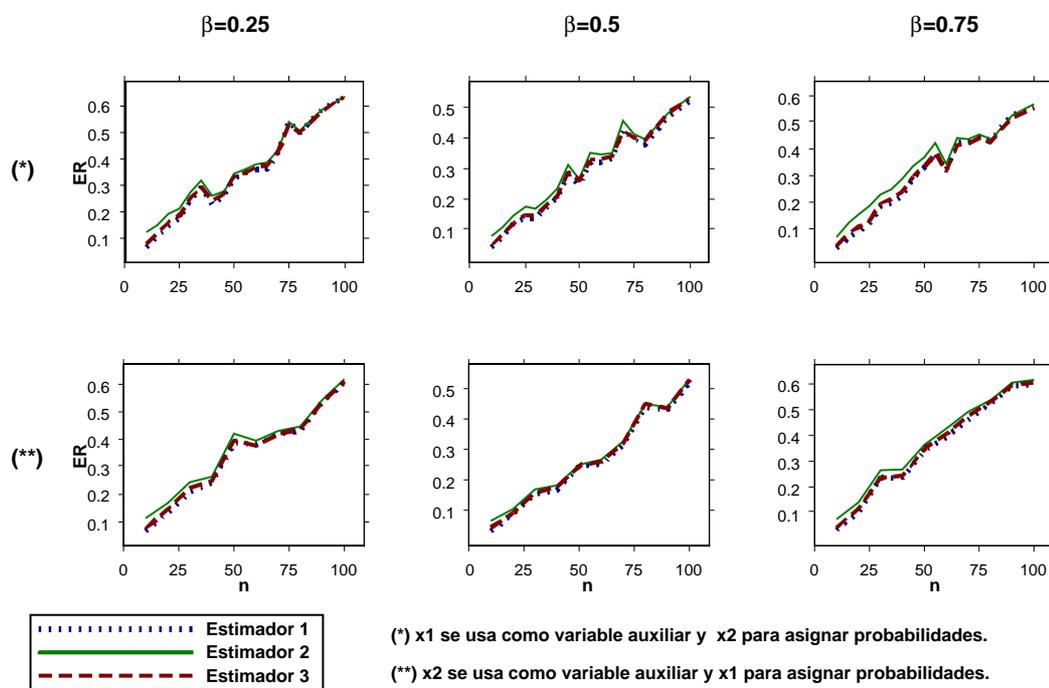


Figura 3.4: Eficiencia Relativa para la población Counties y bajo el diseño muestral *Mas.Poisson*. $n' = 150$.

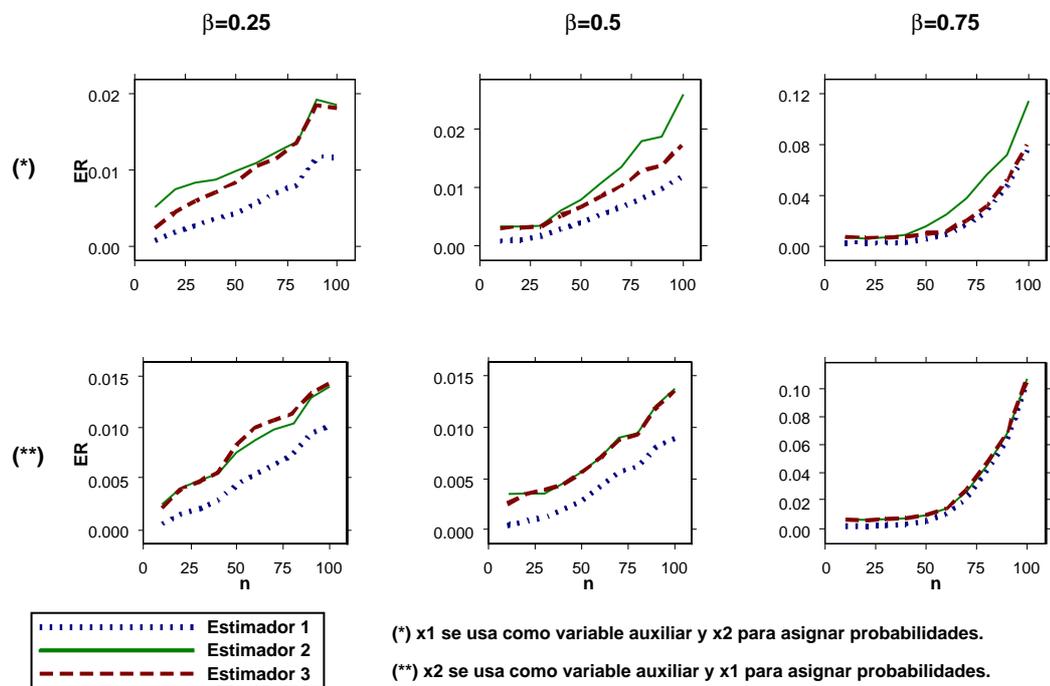


Figura 3.5: Sesgo Relativo en porcentaje para la población Fam1500 cuando x_1 se usa como variable auxiliar y x_2 para asignar probabilidades. $n' = 150$.

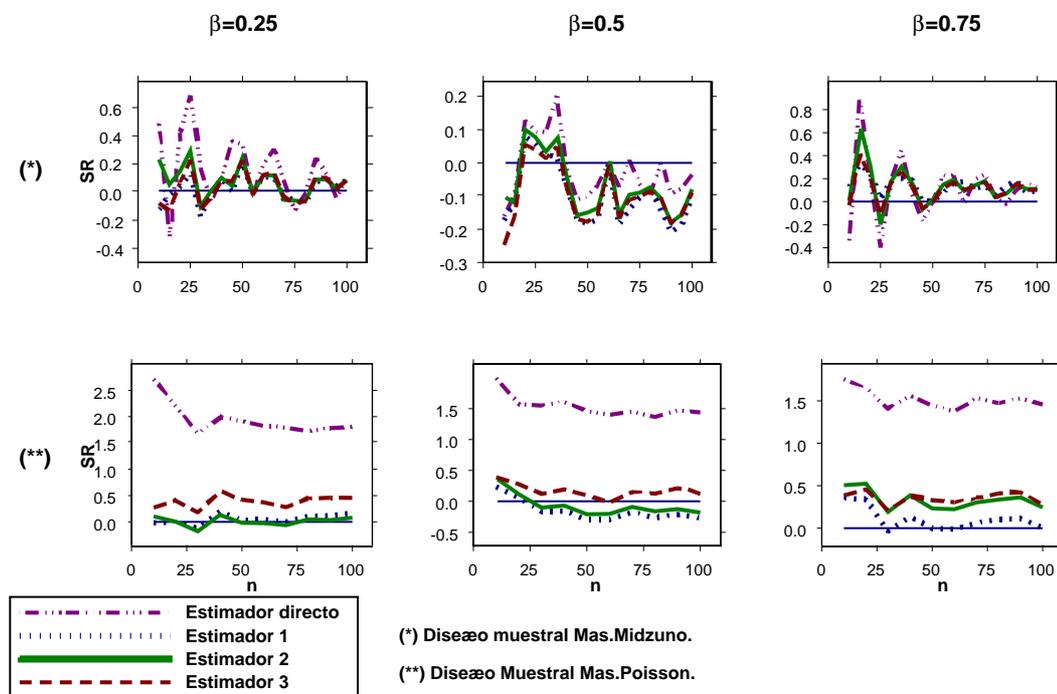
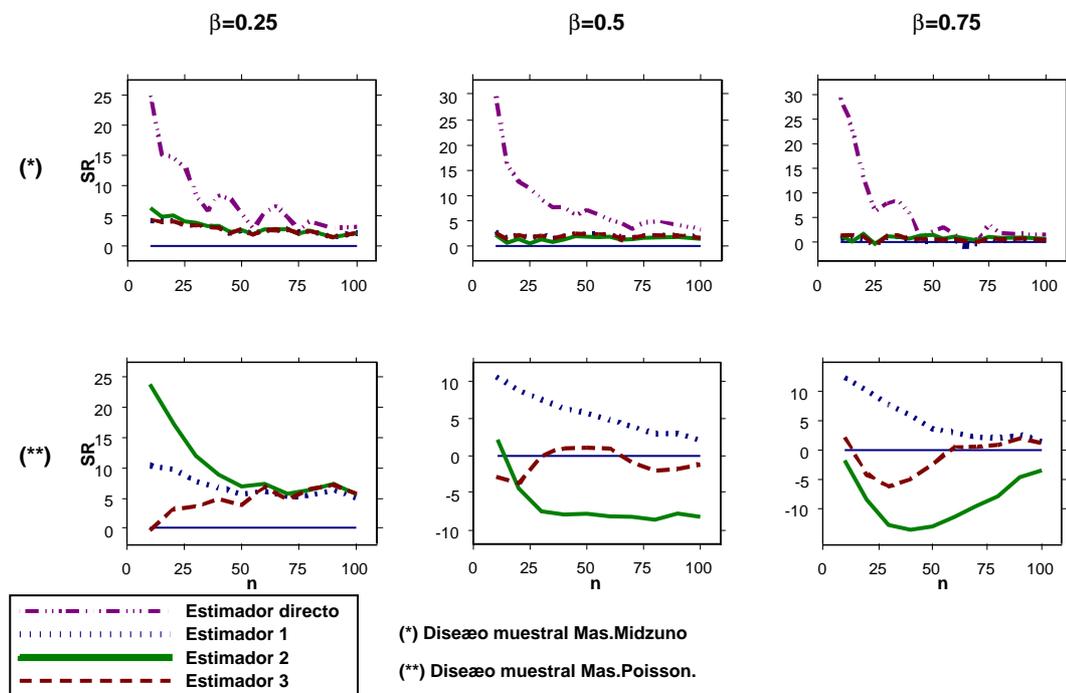


Figura 3.6: Sesgo Relativo en porcentaje para la población Counties cuando x_1 se usa como variable auxiliar y x_2 para asignar probabilidades. Los valores SR para el estimador directo en (**) son mayores de 97.6%, 74.6% y 21.5% para $\beta = 0,25, 0,5$ y 0.75 , respectivamente, y están omitidos. $n' = 150$.



Por otro lado, el método de Poisson produce resultados más eficientes en el sentido de ER que el método de Midzuno y con respecto al estimador $\widehat{Q}_y^*(\beta)$. Esto se debe a que el estimador directo presenta estimaciones muy dispersas bajo el método de Poisson causadas por la heterogeneidad de las probabilidades de inclusión.

Los estimadores propuestos son casi equivalentes en la población Counties porque los coeficientes de correlación lineal están más cercanos a 1. De hecho, la ER de los estimadores propuestos en esta población es mejor que la ER en la población Fam1500.

El estudio del sesgo es otro aspecto importante, particularmente para estimadores de tipo razón, que puede probar la existencia de sub-estimaciones o sobre-estimaciones en los estimadores. Los valores SR en la población Fam1500 están todos dentro de un rango razonable, teniendo el estimador $\widehat{Q}_y^*(\beta)$ el mayor valor en torno al 3%, como puede verse en la Figura 3.5. Los valores de SR para la población Counties cuando x_1 se usa como variable auxiliar y x_2 para asignar probabilidades están mostrados en la Figura 3.6. El estimador $\widehat{Q}_y^*(\beta)$ obtiene, claramente, sobre-estimación, especialmente cuando el tamaño muestral en la segunda fase es pequeño y bajo el diseño muestral Mas.Poisson, mientras que el valor absoluto de los valores SR para los estimadores propuestos son menores de 7% para el diseño Mas.Midzuno y menores de 13% para el diseño Mas.Poisson, excepto en muestras pequeñas para el estimador $\widehat{Q}_y^2(\beta)$, el cual no supera el 25%. En resumen, el estudio de los valores SR revela que los estimadores propuestos presentan un menor sesgo que el estimador directo.

3.2.5. Aplicación al muestreo estratificado

Es sabido que el muestreo estratificado es una potente técnica que proporciona resultados eficientes cuando la población está adecuadamente estratificada y las variables auxiliares y principal presentan una alta correlación. Sin embargo, el muestreo bifásico es la herramienta más apropiada cuando la información auxiliar poblacional no está disponible, que es lo que ocurre en la mayoría de los casos. Estas dos técnicas pueden combinarse en el llamado muestreo bifásico aplicado a la estratificación. Asumiendo este diseño muestral, en esta sección se define un estimador para la función de distribución y se estudian sus principales propiedades. Este estimador se usará para construir nuevos estimadores de cuantiles, y aplicando la relación entre ambos parámetros, será posible también determinar la expresión asintótica de la

varianza del estimador propuesto. La estimación de la varianza es un aspecto muy importante con un alto número de aplicaciones, tal como la construcción de intervalos de confianza, obtención del tamaño muestral óptimo, etc. Por esta razón, tanto el estimador propuesto como su varianza se analizan mediante un estudio de simulación. Los resultados de este estudio reflejan algunas útiles ganancias en eficiencia del estimador propuesto y de su varianza sobre otros estimadores.

La única diferencia de este método de muestreo con respecto al expuesto en la Sección 3.2.2, es el uso adicional del muestreo estratificado. Bajo determinadas condiciones, esta técnica es particularmente eficiente, siendo frecuentemente utilizada en la práctica por diferentes razones: (i) administrativas, cuando el marco de trabajo está dividido en varios distritos geográficos, (ii) importante ganancia en eficiencia sobre diseños muestrales no estratificados, etc.

En resumen, el muestreo bifásico aplicado a la estratificación combina las principales ventajas del muestreo bifásico y muestreo estratificado. Esta técnica consiste en tomar a primera gran muestra de la población en estudio según un diseño muestral determinado. En esta muestra, se observa una variable auxiliar, la cual se usa para estratificar dicha muestra en H estratos. De cada estrato, se selecciona una muestra y se observa la variable de interés. Asumiendo este diseño muestral, se demuestra que, para los datos de la población Fam1500, el estimador propuesto y su varianza pueden proporcionar estimaciones más precisas.

A continuación se describe el muestreo bifásico aplicado a la estratificación y el estimador natural para estimar la función de distribución. Además, se propone un estimador para la función de distribución basado en estimadores π^* .

La notación seguida para el muestreo bifásico aplicado a la estratificación es la siguiente. Una primera muestra s' de tamaño n' es diseñada según el diseño muestral d_1 , de modo que $p_{d_1}(s')$ es la probabilidad de que s' sea seleccionada y donde las correspondientes probabilidades de inclusión de primer y segundo orden se denotan como π'_i y π'_{ij} , para $i, j \in U$. Para los elementos en s' , se recoge la información de una variable auxiliar, x . Esta variable se usa para dividir s' en H pre-especificados estratos denotados como s'_h , ($h = 1, \dots, H$), con n'_h elementos en el estrato h . De este modo, de s'_h se puede seleccionar una muestra s_h de tamaño n_h mediante un diseño $p_h(/s')$. La muestra final será $s = \bigcup_{h=1}^H s_h$. La probabilidades de inclusión para las unidades de la segunda fase se denotan como $\pi_{i/s'}$ y $\pi_{ij/s'}$, para $i, j \in s'$. Notamos que $\Delta'_{ij} = \pi'_{ij} - \pi'_i \pi'_j$ y

$$\Delta_{ij}^{s'} = \pi_{ij/s'} - \pi_{i/s'}\pi_{j/s'}.$$

El primer paso para estimar un determinado cuantil es obtener un buen estimador para la función de distribución con propiedades deseables. El candidato natural (estimador de tipo Horvitz y Thompson) para estimar la función de distribución bajo la técnica de muestreo en estudio es:

$$\widehat{F}_{st}(t) = \frac{1}{N} \sum_{h=1}^H \sum_{i \in s_h} \frac{\delta(t - y_i)}{\pi_i},$$

donde las probabilidades de inclusión están dadas por $\pi_i = \sum_{s' \ni i} p_{d1}(s')\pi_{i/s'}$. Este estimador no puede obtenerse siempre en la práctica debido a que las probabilidades $\pi_{i/s'}$, para cada s' , deben de conocerse para poder determinar π_i . Esto no es siempre posible porque $\pi_{i/s'}$ puede depender del resultado de la primera fase (por ejemplo si la muestra de la segunda fase se selecciona mediante un muestreo proporcional a una variable auxiliar).

En la práctica, el uso del estimador de tipo Horvitz-Thompson no resulta posible ni para el problema de la estimación de la media poblacional. Por esta razón, Särndal *et al.* (1992) propusieron el uso de π^* -estimadores. Usando esta idea, se introducen las cantidades $\pi_i^* = \pi'_i \pi_{i/s'}$ y $\pi_{ij}^* = \pi'_{ij} \pi_{ij/s'}$ para definir el π^* -estimador de la función de distribución como

$$\widehat{F}_{st}^*(t) = \frac{1}{N} \sum_{h=1}^H \sum_{i \in s_h} \frac{\delta(t - y_i)}{\pi_i^*}. \quad (3.8)$$

La calidad de un estimador puede medirse a través de diversas propiedades deseables. A continuación se analizan algunas de las más importantes para el estimador dado por (3.8).

Simplicidad

El cálculo de un estimador de la función de distribución, $\widehat{F}_y(t)$, será particularmente simple si

$$\widehat{F}_y(t) = \frac{1}{N} \sum_{i \in s} w_i \delta(t - y_i),$$

donde los pesos w_i dependen sólo de la etiqueta i . Esto es particularmente deseable para investigaciones con múltiples características. Puede comprobarse fácilmente que el $\widehat{F}_{st}^*(t)$ posee esta propiedad.

Unicidad en la definición

El estimador propuesto es un estimador basado en el diseño muestral, el cual no depende de la elección de un modelo. Además se ha asumido que los estratos están pre-especificados. De este modo, la expresión para $\widehat{F}_{st}^*(t)$ es única.

Sesgo

Una medida importante de la calidad de un estimador es la insesgadez. Särndal *et al.* (1992) establecieron que, para el caso de estimar el total poblacional, el π^* -estimador es insesgado. Este resultado puede extenderse fácilmente al problema de la estimación de la función de distribución, esto es, asumiendo que $z_i = \delta(t - y_i)$ es la variable de interés, el estimador (3.8) puede verse como un problema de estimación de la media poblacional de la variable z_i .

Disponibilidad de la varianza

La varianza de $\widehat{F}_{st}^*(t)$ está dada por

$$V(\widehat{F}_{st}^*(t)) = \frac{1}{N^2} \left(\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(t - y_i)}{\pi'_i} \frac{\delta(t - y_j)}{\pi'_j} + \right. \\ \left. + E_{d1} \left[\sum_{h=1}^H \sum_{i,j \in s'_h} \Delta'_{ij} \frac{\delta(t - y_i)}{\pi_i^*} \frac{\delta(t - y_j)}{\pi_j^*} \right] \right). \quad (3.9)$$

De este modo, un estimador insesgado de esta varianza viene dado por:

$$\widehat{V}(\widehat{F}_{st}^*(t)) = \frac{1}{N^2} \left(\sum_{i,j \in s} \frac{\Delta'_{ij}}{\pi_{ij}^*} \frac{\delta(t - y_i)}{\pi_i'} \frac{\delta(t - y_j)}{\pi_j'} + \right. \\ \left. + \sum_{h=1}^H \sum_{i,j \in s_h} \frac{\Delta'_{ij}}{\pi_{ij}/s'} \frac{\delta(t - y_i)}{\pi_i^*} \frac{\delta(t - y_j)}{\pi_j^*} \right), \quad (3.10)$$

puesto que cada componente de (3.10) es insesgado de su correspondiente componente en la ecuación (3.9).

$\widehat{F}_{st}^*(t)$ es una verdadera función de distribución

En primer lugar, notamos que varios de los estimadores propuestos en la literatura no son verdaderas funciones de distribución. Por ejemplo, ninguno de los conocidos estimadores de tipo razón y diferencia propuestos por Rao *et al.* (1990) es una función de distribución en general (véase Kuk, 1993, Mukhopadhyay, 2000).

Las condiciones (C2.18) y (C2.19) siempre se satisfacen para $\widehat{F}_{st}^*(t)$ y el valor límite de $\widehat{F}_{st}^*(t)$ es también igual a 0. En general, $\lim_{t \rightarrow +\infty} \widehat{F}_{st}^*(t)$ no es igual a 1, aunque esto se verifica para algunos diseños muestrales tal como muestreo aleatorio simple. En la Sección 3.2.7 se analiza $\lim_{t \rightarrow +\infty} \widehat{F}_{st}^*(t)$ para algunos diseños muestrales mediante un estudio de simulación. Los resultados obtenidos para la población Fam1500 sostienen que este valor está bastante próximo a 1. En resumen, el estimador $\widehat{F}_{st}^*(t)$ mantiene todas las condiciones para ser una verdadera función de distribución, excepto en $\lim_{t \rightarrow +\infty} \widehat{F}_{st}^*(t) = 1$, la cual se verifica para algunos diseños muestrales y está bastante próximo a 1 en otros.

La mayoría de los estimadores de cuantiles se obtiene mediante la inversión de la función de distribución. Asumiendo muestreo bifásico, Singh *et al.* (2001) propusieron el siguiente estimador:

$$\widehat{F}_{SJT}(t) = \frac{n'_x \widetilde{F}_{YA}^*(t)}{n'} + \frac{(n' - n'_x) \widetilde{F}_{YB}^*(t)}{n'},$$

donde n'_x es el número de unidades en la primera muestra con $x \leq \widehat{Q}'_x(0,5)$ y $\widetilde{F}_{YA}^*(t)$ y $\widetilde{F}_{YB}^*(t)$ denotando la proporción de unidades en la muestra de la segunda fase para las cuales $x \leq \widehat{Q}'_x(0,5)$ y $x > \widehat{Q}'_x(0,5)$, respectivamente, que tiene valores de y menores o iguales que t . $\widehat{Q}'_x(0,5)$ es el estimador de tipo Horvitz-Thompson para $Q_x(0,5)$ basado en la primera muestra. De este modo, se definió el siguiente estimador para la mediana

$$\widehat{Q}_{SJT}(0,5) = \widehat{F}_{SJT}^{-1}(0,5) = \inf\{t | \widehat{F}_{SJT}(t) \geq 0,5\} \quad (3.11)$$

Siguiendo esta técnica, el cuantil de orden β puede estimarse a partir de $\widehat{F}_{st}^*(t)$ como

$$\widehat{Q}_{st}^*(\beta) = \widehat{F}_{st}^{*-1}(\beta) = \inf\{t | \widehat{F}_{st}^*(t) \geq \beta\}. \quad (3.12)$$

3.2.6. Propiedades teóricas

A continuación se estudian las propiedades del estimador $\widehat{Q}_{st}^*(\beta)$. Para ello, se necesita una aproximación lineal debido a que $\widehat{Q}_{st}^*(\beta)$ no es una función continua.

Teorema 3.3 *El estimador $\widehat{Q}_{st}^*(\beta)$ es asintóticamente insesgado para $Q_y(\beta)$.*

Demostración

El estimador $\widehat{Q}_{st}^*(\beta)$ puede expresarse asintóticamente como una función lineal de la función de distribución estimada evaluada en el cuantil $Q_y(\beta)$ mediante la representación de Bahadur (véase Chambers y Dunstan, 1986):

$$\widehat{Q}_{st}^*(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))}(\beta - \widehat{F}_{st}^*(Q_y(\beta))) + O(n^{-1/2}), \quad (3.13)$$

donde $f_y(\cdot)$ denota la derivada del valor límite de $F_y(\cdot)$ cuando $N \rightarrow \infty$. Como $\widehat{F}_{st}^*(t)$ es un estimador insesgado de $F(t)$, se tiene que $E(\beta - \widehat{F}_{st}^*(Q_y(\beta))) = 0$ y considerando la expresión (3.13), puede comprobarse fácilmente que

$$E(\widehat{Q}_{st}^*(\beta)) = Q_y(\beta) + O(n^{-1/2}).$$

□

Asumiendo la insesgades del estimador $\widehat{Q}_{st}^*(\beta)$ y la expresión (3.13), es posible determinar fácilmente la varianza de dicho estimador al primer grado de aproximación. Esta varianza esta descrita en el siguiente corolario.

Corolario 3.3 *La varianza asintótica del estimador $\widehat{Q}_{st}^*(\beta)$ viene dada por*

$$V(\widehat{Q}_{st}^*(\beta)) = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_i)}{\pi'_i} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_j)}{\pi'_j} + E_{d1} \left[\sum_{h=1}^H \sum_{i,j \in s'_h} \Delta'_{ij} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_j)}{\pi_j^*} \right] \right).$$

Un estimador insesgado para esta varianza viene dado por:

$$\widehat{V}(\widehat{Q}_{st}^*(\beta)) = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in s} \frac{\Delta'_{ij}}{\pi_{ij}^*} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_i)}{\pi'_i} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_j)}{\pi'_j} + \right.$$

$$+ \sum_{h=1}^H \sum_{i,j \in s_h} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_j)}{\pi_j^*} \Big). \quad (3.14)$$

El valor $f_y(Q_y(\beta))$ suele ser desconocido, aunque puede aproximarse aplicando métodos estándares como el método Kernel (Silverman, 1986).

Este estimador para la varianza del estimador propuesto presenta una forma explícita, lo que permite que pueda obtenerse siempre en la práctica, es decir, la expresión (3.14) no depende del valor esperado sobre el diseño de la primera fase, haciendo posible los cálculos directos.

Una vez que la varianza del estimador ha sido determinada, intervalos de confianza y otras importantes aplicaciones derivadas de la varianza podrán también obtenerse.

En el siguiente ejemplo se determina las expresiones del estimador propuesto $\widehat{Q}_{st}^*(\beta)$ y de su correspondiente varianza estimada para el caso de selección de unidades mediante muestreo aleatorio simple.

Ejemplo 3.1 *Asumiendo muestreo aleatorio simple en cada fase, el π^* -estimador viene dado por*

$$\widehat{Q}_{st}^*(\beta) = \inf\{t \mid \sum_{h=1}^H \frac{n'_h}{n'} \sum_{i \in s_h} \frac{\delta(t - y_i)}{n_h} \geq \beta\},$$

y el estimador de su varianza puede obtenerse de (3.14) después de sustituir las probabilidades $\pi_{i/s'}$, π_i^* , $\pi_{ij/s'}$ y π_{ij}^* por

$$\begin{aligned} \pi_{i/s'} &= \frac{n_h}{n'} \quad ; \quad \pi_i^* = \frac{n'_h n_h}{N n'_h}, \quad \text{para } i \in s_h, \\ \pi_{ij/s'} &= \begin{cases} \frac{n_h(n_h - 1)}{n'_h(n'_h - 1)} & \text{si } i, j \in s'_h \\ \frac{n_h n_l}{n'_h n'_l} & \text{si } i \in s'_h \text{ y } j \in s'_l \end{cases} \\ \pi_{ij}^* &= \begin{cases} \frac{n_h(n_h - 1) n'(n' - 1)}{n'_h(n'_h - 1) N(N - 1)} & \text{si } i, j \in s'_h \\ \frac{n_h n_l n'(n' - 1)}{n'_h n'_l N(N - 1)} & \text{si } i \in s'_h \text{ y } j \in s'_l \end{cases} \end{aligned}$$

Tabla 3.1: Esperanza empírica de $\lim_{t \rightarrow \infty} \widehat{F}_{st}^*(t)$ para varios diseños muestrales y considerando la variable x_1 .

n'	n	d_{SS}	d_{SM}	d_{SP}	d_{MS}	d_{MM}	d_{MP}	d_{PS}	d_{PM}	d_{PP}
150	30	1.000	1.010	1.000	1.001	1.011	1.000	1.000	1.000	1.000
	50	1.000	1.005	1.000	1.001	1.006	1.000	1.000	1.000	0.999
	70	1.000	1.003	1.000	1.001	1.004	1.000	1.000	1.000	1.000
	90	1.000	1.002	1.000	1.001	1.002	1.000	0.999	1.000	1.000
300	60	1.000	1.005	1.000	1.000	1.005	1.000	0.999	1.000	1.000
	100	1.000	1.003	1.000	1.000	1.003	1.000	1.000	1.000	1.000
	140	1.000	1.001	1.000	1.000	1.002	1.000	1.000	1.000	1.000
	180	1.000	1.001	1.000	1.000	1.001	1.000	1.000	1.000	1.000

Tabla 3.2: Esperanza empírica de $\lim_{t \rightarrow \infty} \widehat{F}_{st}^*(t)$ para varios diseños muestrales y considerando la variable x_2 .

n'	n	d_{SS}	d_{SM}	d_{SP}	d_{MS}	d_{MM}	d_{MP}	d_{PS}	d_{PM}	d_{PP}
150	30	1.000	1.011	1.002	1.001	1.011	0.998	1.001	1.002	1.002
	50	1.000	1.005	1.002	1.001	1.006	1.001	1.000	1.001	0.999
	70	1.000	1.003	0.999	1.001	1.004	0.999	1.000	1.000	0.999
	90	1.000	1.002	1.000	1.001	1.002	0.999	1.000	1.001	0.999
300	60	1.000	1.005	1.000	1.000	1.005	0.999	1.000	1.000	0.999
	100	1.000	1.003	1.000	1.000	1.003	1.000	1.000	1.000	0.999
	140	1.000	1.001	1.000	1.000	1.002	1.000	0.999	1.000	0.999
	180	1.000	1.001	1.000	1.000	1.001	1.000	1.000	1.000	1.000

3.2.7. Propiedades empíricas

Asumiendo muestreo bifásico aplicado a la estratificación, se ha propuesto un estimador para un determinado cuantil poblacional y su correspondiente varianza asintótica ha sido determinada. La insesgidez del estimador de cuantiles también ha sido discutida. El siguiente paso será analizar, mediante un estudio de simulación, éstas y otras medidas importantes de calidad para los dos estimadores propuestos. Los resultados se compararan sobre otros estimadores conocidos en la literatura del muestreo en poblaciones finitas.

En este estudio se usa la población Fam1500 (véase Apéndice A), donde recordamos que las correlaciones entre la variable principal y las auxiliares vienen dadas por $\rho_{y,x_1} = 0,848$ y $\rho_{y,x_2} = 0,546$.

En primer lugar, analizaremos $\lim_{t \rightarrow \infty} F_{st}^*(t)$ para poder comprobar como

Tabla 3.3: Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_1 . $\beta = 0,5$ y $n' = 150$.

n	ER				SR (%)				$RECMR$ (%)			
	30	50	70	90	30	50	70	90	30	50	70	90
\widehat{Q}_{st}^*	0.59	0.69	0.59	0.68	-0.1	-0.1	-0.1	0.0	2.7	2.2	1.7	1.5
\widehat{Q}_y^*	1.00	1.00	1.00	1.00	0.2	-0.1	0.0	0.0	3.5	2.6	2.2	1.9
\widehat{Q}_{SJT}	0.64	0.66	0.67	0.74	-0.2	-0.1	-0.1	0.0	2.8	2.1	1.8	1.6
$\widehat{V}(\widehat{Q}_{st}^*)$	0.32	0.42	0.42	0.26	-5.2	9.2	13.2	7.4	15.8	12.7	14.9	8.6
$\widehat{V}(\widehat{Q}_y^*)$	1.00	1.00	1.00	1.00	-16.6	-13.5	-13.5	-11.3	16.6	13.5	13.5	11.3
$\widehat{V}(\widehat{Q}_{SJT})$	1.11	2.18	2.37	2.29	27.4	30.1	31.1	23.2	27.4	30.1	31.1	23.2

de cercano se encuentra de 1. Recordamos que $F_{st}^*(t)$ será una verdadera función de distribución si este valor es igual a 1. Se ha considerado muestreo aleatorio simple (S), el método de Midzuno (M) y el método de Poisson (P). Las diferentes combinaciones de diseños muestrales se van a denotar como d_{ij} , para $i, j = \{S, M, P\}$, donde i y j van a expresar los diseños muestrales usados en la primera y segunda fase, respectivamente. Este estudio se ha llevado a cabo usando aproximaciones Monte Carlo derivadas de 1000 muestras independientes, para $\beta = 0,5$, $n' = 150$ y 300 y varios valores de n .

Para cada diseño muestral, las Tablas 3.1 y 3.2 muestran la esperanza empírica de $\lim_{t \rightarrow \infty} \widehat{F}_{st}^*(t)$ basada en 1000 muestras de la población Fam1500. Puede observarse que todos los resultados están cercanos a 1, obteniéndose mejores resultados cuando la muestra de la segunda fase es mayor. Como esperábamos, asumiendo muestreo aleatorio simple en cada una de las fases, siempre se obtiene que $\lim_{t \rightarrow \infty} \widehat{F}_{st}^*(t) = 1$. Esto también ocurre en la mayoría de los casos cuando se considera el método de Poisson en alguna de las dos fases. En general, la variable x_1 (para correlaciones altas) obtiene mejores resultados que la variable x_2 .

El siguiente paso es comparar el comportamiento del estimador propuesto para cuantiles y de su varianza con respecto a otros estimadores. En este estudio, se ha incluido el estimador (3.11) y su correspondiente estimador de la varianza propuesto en Singh *et al.* (2001). La ganancia en eficiencia sobre muestreo no estratificado puede contrastarse si comparamos el estimador propuesto con el estimador basado en la segunda fase, sin considerar estratos en la primera fase. Este estimador será denotado como $\widehat{Q}_y^*(\beta)$ y lo usaremos como el estimador base en las comparaciones.

Tabla 3.4: Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_1 . $\beta = 0,5$ y $n' = 300$.

n	ER				SR (%)				$RECMR$ (%)			
	60	100	140	180	60	100	140	180	60	100	140	180
\widehat{Q}_{st}^*	0.55	0.61	0.73	0.76	-0.1	0.0	-0.1	-0.1	1.8	1.4	1.3	1.1
\widehat{Q}_y^*	1.00	1.00	1.00	1.00	0.1	0.1	0.0	-0.1	2.5	1.8	1.5	1.3
\widehat{Q}_{SJT}	0.58	0.62	0.73	0.80	0.0	0.0	0.0	-0.1	1.9	1.4	1.3	1.1
$\widehat{V}(\widehat{Q}_{st}^*)$	0.10	0.09	0.33	0.13	-4.8	-4.1	-9.9	-4.2	11.7	8.0	10.7	5.0
$\widehat{V}(\widehat{Q}_y^*)$	1.00	1.00	1.00	1.00	-20.2	-16.2	-13.4	-10.4	20.2	16.2	13.4	10.4
$\widehat{V}(\widehat{Q}_{SJT})$	1.18	2.10	1.68	2.38	37.7	37.6	23.7	20.2	37.7	37.6	23.7	20.2

Tabla 3.5: Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_2 . $\beta = 0,5$ y $n' = 150$.

n	ER				SR (%)				$RECMR$ (%)			
	30	50	70	90	30	50	70	90	30	50	70	90
\widehat{Q}_{st}^*	0.59	0.60	0.72	0.77	-0.1	0.0	0.1	-0.1	2.7	2.1	1.8	1.7
\widehat{Q}_y^*	1.00	1.00	1.00	1.00	0.2	0.1	0.0	-0.1	3.5	2.7	2.1	1.9
\widehat{Q}_{SJT}	0.78	0.84	0.90	0.94	-0.1	0.0	0.0	-0.1	3.1	2.5	2.0	1.9
$\widehat{V}(\widehat{Q}_{st}^*)$	0.27	0.12	0.28	0.24	-8.1	-1.8	-2.1	-8.6	17.5	10.4	6.7	9.5
$\widehat{V}(\widehat{Q}_y^*)$	1.00	1.00	1.00	1.00	-19.8	-18.3	-9.0	-14.9	19.8	18.3	9.0	14.9
$\widehat{V}(\widehat{Q}_{SJT})$	0.01	0.01	0.18	0.13	0.9	-1.7	4.2	-5.7	0.9	1.8	4.2	5.7

Tabla 3.6: Medidas de eficiencia y precisión para los estimadores de cuantiles y sus varianzas asumiendo el diseño muestral d_{SM} y la variable x_2 . $\beta = 0,5$ y $n' = 300$.

n	ER				SR (%)				$RECMR$ (%)			
	60	100	140	180	60	100	140	180	60	100	140	180
\widehat{Q}_{st}^*	0.57	0.57	0.66	0.73	-0.1	0.0	-0.1	0.0	1.8	1.4	1.2	1.1
\widehat{Q}_y^*	1.00	1.00	1.00	1.00	0.0	-0.1	-0.1	-0.1	2.4	1.8	1.5	1.3
\widehat{Q}_{SJT}	0.80	0.84	0.89	0.90	-0.1	-0.1	-0.1	0.0	2.1	1.7	1.4	1.2
$\widehat{V}(\widehat{Q}_{st}^*)$	0.29	0.09	0.06	0.08	0.7	3.1	-3.2	-4.8	12.0	8.4	5.8	5.7
$\widehat{V}(\widehat{Q}_y^*)$	1.00	1.00	1.00	1.00	-12.8	-17.0	-15.5	-14.5	12.8	17.0	15.5	14.5
$\widehat{V}(\widehat{Q}_{SJT})$	0.42	0.03	0.01	0.13	10.3	3.3	2.0	5.9	10.3	3.3	2.1	5.9

La precisión de todos los estimadores de cuantiles y sus respectivas varianzas se miden para $\beta = 0,5$ mediante el Sesgo Relativo (SR), la Eficiencia Relativa (ER) y la Raíz cuadrada del Error Cuadrático Medio Relativo ($RECMR$). Para un cuantil, $\widehat{Q}_y(\beta)$, están medidas están dadas por

$$\begin{aligned} SR[\widehat{Q}_y(\beta)] &= \left(E[\widehat{Q}_y(\beta)] - Q_y(\beta) \right) / Q_y(\beta), \\ ER[\widehat{Q}_y(\beta)] &= ECM[\widehat{Q}_y(\beta)] / ECM[\widehat{Q}_y^*(\beta)], \\ RECMR[\widehat{Q}_y(\beta)] &= \left(ECM[\widehat{Q}_y(\beta)] \right)^{1/2} / Q_y(\beta), \end{aligned}$$

y para el estimador de la varianza de un cuantil, $\widehat{V}(\widehat{Q}_y(\beta))$, las medidas son

$$\begin{aligned} SR[\widehat{V}(\widehat{Q}_y(\beta))] &= \left(E[\widehat{V}(\widehat{Q}_y(\beta))] - V[Q_y(\beta)] \right) / V[Q_y(\beta)], \\ ER[\widehat{V}(\widehat{Q}_y(\beta))] &= ECM[\widehat{V}(\widehat{Q}_y(\beta))] / ECM[\widehat{V}(\widehat{Q}_y^*(\beta))], \\ RECMR[\widehat{V}(\widehat{Q}_y(\beta))] &= \left(ECM[\widehat{V}(\widehat{Q}_y(\beta))] \right)^{1/2} / V[Q_y(\beta)], \end{aligned}$$

donde $E[\cdot]$, $ECM[\cdot]$ y $V[\cdot]$ denotan las Esperanzas, Errores Cuadráticos Medios y Varianzas empíricas basadas en 1000 muestras. Notamos que valores de $ER[\widehat{Q}_y(\beta)]$ y $ER[\widehat{V}(\widehat{Q}_y(\beta))]$ menores de 1 indican que $\widehat{Q}_y(\beta)$ y $\widehat{V}(\widehat{Q}_y(\beta))$ son más precisos que $\widehat{Q}_y^*(\beta)$ y $\widehat{V}(\widehat{Q}_y^*(\beta))$, respectivamente. También se ha calculado la Cobertura de los intervalos de confianza al 95 % (asumiendo distribución normal) y la longitud media de los intervalos basados en 1000 muestras.

Asumiendo muestreo aleatorio simple para obtener la muestra de la primera fase y el método de Midzuno para obtener la segunda muestra, en las Tablas 3.3 y 3.4 pueden observarse los resultados de las distintas medidas de precisión para los estimadores y asumiendo la variable x_1 . En este caso (para una alta correlación), tanto el estimador propuesto como su correspondiente varianza son mas precisos, en términos de ER , que sus competidores. Los valores absolutos de las medidas SR , para todos los cuantiles, son siempre menores de 0,2 %. Respecto a las varianzas, se observa que $\widehat{V}(\widehat{Q}_y^*)$ presenta subestimación, mientras que $\widehat{V}(\widehat{Q}_{SJT})$ claramente arrastra una seria sobreestimación. Los estimadores propuestos también presentan la mejor precisión en términos de $RECMR$.

A continuación se analiza la precisión de los estimadores usando una menor correlación entre la variable principal y auxiliar. Para ello, observamos las Tablas 3.5 y 3.6. El estimador propuesto para estimar cuantiles es más preciso que el resto en términos de ER . Respecto a la estimación de varianzas, $\widehat{V}(\widehat{Q}_{SJT})$ parece tener el mejor comportamiento, aunque esto sólo ocurre para una escasa correlación entre las variables (situación no deseada en la práctica)

Tabla 3.7: Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_1 . $\beta = 0,5$ y $n' = 150$.

n	Cobertura (%)				Longitud Media			
	30	50	70	90	30	50	70	90
\widehat{Q}_{st}^*	94.1	93.4	96.6	95.3	828	656	566	512
\widehat{Q}_y^*	92.2	92.5	92.8	93.9	1010	772	646	564
\widehat{Q}_{SJT}	96.9	97.3	97.4	96.8	998	771	650	571

Tabla 3.8: Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_1 . $\beta = 0,5$ y $n' = 300$.

n	Cobertura (%)				Longitud Media			
	60	100	140	180	60	100	140	180
\widehat{Q}_{st}^*	94.4	93.9	93.7	93.2	568	447	385	347
\widehat{Q}_y^*	92.1	93.1	93.0	93.1	701	534	444	385
\widehat{Q}_{SJT}	96.8	98.1	96.9	97.0	703	541	454	398

y para el caso de varianzas. Conclusiones similares pueden obtenerse a partir del sesgo y del error cuadrático medio. Como resulta razonable, éstas últimas medidas mejoran para cada estimador a medida que se aumenta el tamaño de la muestra de cualquiera de las dos fases.

Por último, se analiza la cobertura y la longitud media de los intervalos de confianza de cada estimador. Estas medidas vienen dadas por las Tablas 3.7 y 3.8 para la variable x_1 y las Tablas 3.9 y 3.10 para la variable x_2 . En todos los casos se observa que el estimador propuesto tiene la menor longitud media empírica para el intervalo de confianza. Para altas correlaciones, la cobertura del estimador propuesto es mejor que la del resto de estimadores, puesto que se obtienen valores más próximos al 95%. Para bajas correlaciones, la cobertura del estimador propuesto se ve ligeramente superada por la cobertura de \widehat{Q}_{SJT} , aunque éste último estimador tiene el inconveniente de presentar intervalos de confianza mucho más amplios.

Tabla 3.9: Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_2 . $\beta = 0,5$ y $n' = 150$.

n	Cobertura (%)				Longitud Media			
	30	50	70	90	30	50	70	90
\widehat{Q}_{st}^*	93.7	94.0	94.7	93.8	830	655	567	512
\widehat{Q}_y^*	90.7	93.5	94.1	92.8	1010	772	646	565
\widehat{Q}_{SJT}	93.8	94.7	95.4	94.5	1001	775	654	576

Tabla 3.10: Cobertura y Longitud Media de Intervalos de Confianza de los distintos estimadores bajo el diseño d_{SM} y asumiendo la variable x_1 . $\beta = 0,5$ y $n' = 300$.

n	Cobertura (%)				Longitud Media			
	60	100	140	180	60	100	140	180
\widehat{Q}_{st}^*	94.8	95.7	94.8	92.4	568	447	385	347
\widehat{Q}_y^*	92.7	92.8	92.6	92.4	701	534	444	385
\widehat{Q}_{SJT}	96.3	95.1	94.8	94.7	707	541	461	406

3.3. Estimadores bajo muestreo en dos ocasiones sucesivas

El muestreo en ocasiones sucesivas es una técnica muy conocida que puede emplearse en las investigaciones longitudinales para estimar determinados parámetros poblacionales y medidas de diferencia o cambio de una variable objeto de estudio. En esta sección se discute la estimación de cuantiles en la ocasión más reciente bajo un muestreo en dos ocasiones sucesivas. Este estudio se realiza, por un lado, para el caso de un muestreo con probabilidades de selección de unidades desiguales y por otro, para hacer un uso más efectivo de la información auxiliar, considerando varias variables auxiliares en la etapa de estimación. Se estudian las propiedades más importantes y se deducen las expresiones de las varianzas. Como es habitual, se mide la precisión de los estimadores propuestos en estudios de simulación basados en varias poblaciones.

3.3.1. Introducción

En numerosas investigaciones por muestreo, una misma población puede ser muestreada repetidamente y la misma variable de estudio es medida en cada ocasión, de modo que se sigue el desarrollo de ésta sobre el tiempo. Por ejemplo, las encuestas de presupuestos familiares son llevadas a cabo periódicamente para estimar el número de empleados, las encuestas de opinión se llevan a cabo a intervalos regulares de tiempo para medir las preferencias de los votantes, etc. En estos casos, el uso de la teoría de un esquema de muestreo sucesivo puede ser una alternativa atractiva para mejorar las estimaciones de nivel en un punto en el tiempo, el cambio entre dos puntos, etc. (véase por ejemplo Cochran, 1977).

El muestreo en ocasiones sucesivas ha sido extensamente usado en las ciencias sociales y aplicadas para estimar medidas de nivel, cambios de un parámetro lineal tal como la media o el total, y contrastar la dirección de este cambio. Otros ejemplos del uso de encuestas longitudinales pueden consultarse en Ruspini (1999) para el análisis en el cambio social, Solga (2001) para el estudio de movilidad laboral, etc.

Asumiendo muestreo en dos ocasiones sucesivas, la teoría desarrollada por Jessen (1942) y Patterson (1950) proporciona el estimador óptimo de la media poblacional en la segunda ocasión, combinando dos estimadores distintos de esta media. Por un lado, se usa un estimador de tipo regresión basado en la

muestra solapada de la muestra, considerando que la variable auxiliar es el valor de la variable principal en la primera ocasión. Por último, se considera un estimador simple de la media basado en una muestra aleatoria de la porción no solapada de la segunda ocasión. El muestreo en ocasiones sucesivas también ha sido discutido en Narain (1953), Adhvaryu (1978), Eckler (1955), Gordon (1983), Arnab y Okafor (1992), Singh y Srivastava (1973), Singh *et al.* (1992) y Singh (2003), el cual proporciona una extensa bibliografía sobre este tópico. En todos los estudios anteriores, el parámetro considerado para su estimación es la media poblacional.

La mayoría de los estudios relacionados con la estimación de la mediana han sido desarrollados asumiendo muestreo aleatorio simple o muestreo estratificado (Gross, 1980, Sedransk y Meyer, 1978, Sedransk y Smith, 1988) y considerando sólo la variable de interés, sin hacer un uso explícito de las variables auxiliares para la construcción de los estimadores. Así, la literatura relacionada con la estimación de medianas y otros cuantiles usando una variable auxiliar es considerablemente menos extensa que para el problema de la estimación de medias y totales poblacionales.

Recientemente, Martínez *et al.* (2004) propusieron una metodología de estimación de cuantiles en muestreo en ocasiones sucesivas usando el valor de la variable principal en una ocasión anterior como variable auxiliar. Este estudio fue desarrollado bajo muestreo aleatorio simple y asumiendo que sobre la ocasión más reciente se toma una submuestra a partir de las unidades previamente seleccionadas, y que ciertas de estas unidades son reemplazadas por otras nuevas unidades seleccionadas independientemente de la muestra solapada. Más recientemente, el problema de la estimación de cuantiles en muestreo con dos ocasiones sucesivas puede también consultarse en Rueda y Muñoz (2006b) y Rueda, Muñoz y Arcos (2006).

En esta sección, se propone un estimador para un cuantil de orden β en muestreo de ocasiones sucesivas con diseños muestrales arbitrarios en cada una de las dos fases que consta este esquema de muestreo. Se usará un estimador de tipo razón en la porción de muestra solapada para proporcionar el estimador óptimo de un cuantil. Para ello, se pondera las estimaciones inversamente a sus varianzas. Las propiedades del estimador propuesto se estudian bajo aproximaciones basadas en muestras de gran tamaño. El comportamiento de este estimador también se estudia bajo los datos de una población real.

La notación habitual a seguir en muestreo en ocasiones sucesivas es la siguiente (véase también Apéndice B). Consideramos que estamos haciendo un seguimiento continuo de la población U , de tamaño N , sobre dos, o más,

periodos de tiempo con valores y_i en el periodo u ocasión más reciente. Se asume que una muestra de tamaño n' está diseñada en la ocasión anterior. En la ocasión reciente, una submuestra (llamada muestra solapada) de tamaño m es diseñada de las n' unidades seleccionadas previamente, y $u = n - m$ unidades son reemplazadas por nuevas unidades seleccionadas de la población restante. $\chi = m/n$ será la fracción de solapamiento.

La muestra de la primera fase s' con tamaño n' está diseñada según un diseño muestral d_1 , tal que $p_{d_1}(s')$ es la probabilidad de que s' sea escogida. Las correspondientes probabilidades de inclusión de primer y segundo orden vienen dadas por π'_i, π'_{ij} , para $i, j \in U$. Dada s' , en la segunda ocasión, una muestra solapada s_m con tamaño m , es diseñada según un diseño d_2 , tal que $p_m(s_m/s')$ es la probabilidad condicional de escoger s_m . Las probabilidades de inclusión bajo este diseño se denotan como $\pi_{i/s'}$ y $\pi_{ij/s'}$. La muestra no solapada s_u es por tanto seleccionada de $U - s' = s'^c$ según el diseño d_3 , tal que $p_u(s_u/s'^c)$ es la probabilidad condicional de escoger s_u . Las probabilidades de inclusión bajo este diseño se denotarán como π_{i/s'^c} y π_{ij/s'^c} .

En muestreo con dos ocasiones sucesivas, el estimador habitual para la estimación de cuantiles se construye como sigue. En primer lugar se estima la función de distribución a partir de la muestra s obtenida en la ocasión más reciente. Este estimador viene dado por $\hat{F}_{yn}(t) = n^{-1} \sum_{i \in s} \delta(t - y_i)$, el cual coincide con el estimador de tipo Horvitz-Thompson bajo muestreo aleatorio simple. A continuación se estima el cuantil de orden β a partir de esta función de distribución, es decir:

$$\hat{Q}_{yn}(\beta) = \hat{F}_{yn}^{-1}(\beta) = \inf \left\{ t : \hat{F}_{yn}(t) \geq \beta \right\}. \quad (3.15)$$

3.3.2. Muestreo con probabilidades desiguales

Asumiendo muestreo en dos ocasiones sucesivas, Särndal *et al.* (1992) demostró que el estimador de tipo Horvitz-Thompson de una media no puede siempre usarse en la práctica debido a que el estimador requiere el cálculo de las probabilidades de inclusión π_i , y esto no es posible para las unidades de la muestra s_u o para las unidades de la muestra s_m .

En esta sección y en las dos siguientes asumiremos que se dispone de una única variable auxiliar, x , que serán los valores de la variable principal que toman los individuos en el primer periodo u ocasión. También puede considerarse que x es una variable auxiliar altamente correlacionada con la variable principal, aunque en la práctica esto no es lo habitual.

A continuación se define un estimador compuesto basado en estimadores π^* (véase Särndal *et al.*, 1992, p.347) y que combina un estimador construido en la muestra solapada con otro estimador basado en la muestra no solapada.

Así, usando la muestra no solapada, s_u , nosotros podemos obtener el siguiente estimador para la función de distribución

$$\widehat{F}_{yu}(t) = \frac{1}{N} \sum_{i \in s_u} \frac{\delta(t - y_i)}{\pi_i^{t/c} \pi_{i/s^c}},$$

el cual es un estimador π^* . El correspondiente estimador para el cuantil de orden β viene por tanto dado por

$$\widehat{Q}_{yu}(\beta) = \inf\{t : \widehat{F}_{yu}(t) \geq (\beta)\}. \quad (3.16)$$

A partir de la muestra solapada pueden construirse los siguientes estimadores de la función de distribución

$$\widehat{F}_{ym}(t) = \frac{1}{N} \sum_{i \in s_m} \frac{\delta(t - y_i)}{\pi_i' \pi_{i/s'}}, \quad (3.17)$$

$$\widehat{F}_{xm}(t) = \frac{1}{N} \sum_{i \in s_m} \frac{\delta(t - x_i)}{\pi_i' \pi_{i/s'}}, \quad (3.18)$$

los cuales son estimadores π^* basados en la segunda y primera ocasión respectivamente. Usando también la muestra de la primera fase, es posible construir un estimador de tipo Horvitz-Thompson para la variable auxiliar como sigue

$$\widehat{F}_x(t) = \frac{1}{N} \sum_{i \in s'} \frac{\delta(t - x_i)}{\pi_i'}. \quad (3.19)$$

Usando los estimadores dados en (3.17), (3.18) y (3.19) y basándonos en la muestra solapada y en la muestra de la primera fase, se propone el siguiente estimador de tipo razón

$$\widehat{Q}_{ym}^r(\beta) = \widehat{Q}_{ym}(\beta) \frac{\widehat{Q}_x(\beta)}{\widehat{Q}_{xm}(\beta)}, \quad (3.20)$$

donde

$$\widehat{Q}_{ym}(\beta) = \inf\{t : \widehat{F}_{ym}(t) \geq \beta\}, \quad (3.21)$$

$$\widehat{Q}_{xm}(\beta) = \inf\{t : \widehat{F}_{xm}(t) \geq \beta\}, \quad (3.22)$$

$$\widehat{Q}_x(\beta) = \inf\{t : \widehat{F}_x(t) \geq \beta\}. \quad (3.23)$$

Seguindo a Jessen (1942), se propone el estimador compuesto $\widehat{Q}_y^R(\beta)$ para $Q_y(\beta)$ como combinación lineal del estimador (3.16) y el estimador (3.20). Este estimador viene dado por

$$\widehat{Q}_y^R(\beta) = w\widehat{Q}_{ym}^r(\beta) + (1-w)\widehat{Q}_{yu}(\beta), \quad (3.24)$$

donde w es un peso constante y no negativo. El siguiente paso será determinar w de modo que se minimice la varianza del estimador compuesto $\widehat{Q}_y^R(\beta)$.

Teorema 3.4 *La varianza mínima del estimador $\widehat{Q}_y^R(\beta)$ viene dada por*

$$V_{min}(\widehat{Q}_y^R(\beta)) = \frac{V_1V_2 - C^2}{V_1 + V_2 - 2C}. \quad (3.25)$$

Demostración

La varianza de $\widehat{Q}_y^R(\beta)$ viene dada por

$$\begin{aligned} V(\widehat{Q}_y^R(\beta)) &= w^2V(\widehat{Q}_{ym}^r(\beta)) + (1-w)^2V(\widehat{Q}_{yu}(\beta)) \\ &+ 2w(1-w)Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta)) = w^2V_1 + (1-w)^2V_2 + 2w(1-w)C = \\ &(V_1 + V_2 - 2C)\left\{w - \frac{V_2 - C}{V_1 + V_2 - 2C}\right\}^2 + \frac{V_1V_2 - C^2}{V_1 + V_2 - 2C} \geq \\ &\frac{V_1V_2 - C^2}{V_1 + V_2 - 2C} = V_{min}(\widehat{Q}_y^R(\beta)), \end{aligned}$$

puesto que $V_1 + V_2 - 2C > 0$, y donde

$$V_1 = V(\widehat{Q}_{ym}^r(\beta)),$$

$$V_2 = V(\widehat{Q}_{yu}(\beta)),$$

$$C = Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta)).$$

□

Por tanto el valor de w que hace mínima la varianza de $\widehat{Q}_y^R(\beta)$ viene dado por

$$w = \frac{V_2 - C}{V_1 + V_2 - 2C}. \quad (3.26)$$

Partiendo de este resultado, el estimador propuesto será más eficiente que el estimador habitual $\widehat{Q}_{yu}(\beta)$ y el estimador de tipo razón $\widehat{Q}_{ym}^r(\beta)$.

3.3.3. Propiedades teóricas

En esta sección se estudian las propiedades asintóticas del estimador propuesto en (3.24). Los resultados obtenidos se derivan asumiendo las siguientes condiciones:

(C3.4). Asumimos que s' es una muestra aleatoria simple de U , lo cual implica que la muestra complementaria s'^c es también una muestra aleatoria simple de U . Finalmente, asumiremos que s_m es una muestra aleatoria simple de s' y s_u es otra muestra aleatoria simple de s'^c . Bajo estas condiciones, las probabilidades de inclusión vienen dadas por: $\pi'_i = \frac{n'}{N}$,

$$\pi'_{ij} = \frac{n' n' - 1}{N N - 1}, \quad \pi_{i/s'} = \frac{m}{n'}, \quad \pi_{ij/s'} = \frac{m(m-1)}{n'(n'-1)}, \quad \pi_{i/s'^c} = \frac{u}{N-n'}, \quad \pi_{ij/s'^c} = \frac{u(u-1)}{(N-n')(N-n'-1)}.$$

(C3.5). Suponemos que la población finita está envuelta en una secuencia de poblaciones $\{U_\nu\}$, donde n_ν y N_ν aumentan de modo que $(n_\nu/N_\nu) \rightarrow f$ cuando $n_\nu \rightarrow \infty$.

(C3.6). Se asume que cuando $N_\nu \rightarrow \infty$ la distribución bivariable formada por (x, y) puede aproximarse por una distribución continua con densidades marginales $f_x(\cdot)$ y $f_y(\cdot)$ para x e y respectivamente, siendo $f_x(Q_x(\beta))$ y $f_y(Q_y(\beta))$ positivas.

Teorema 3.5 *El estimador compuesto $\widehat{Q}_y^R(\beta)$ es asintóticamente insesgado para $Q_y(\beta)$.*

Demostración

Para demostrar este resultado usaremos la insesgidez de los dos estimadores en los que se basa el estimador propuesto. En primer lugar, es sabido que el cuantil muestral $\widehat{Q}_{yu}(\beta)$ es asintóticamente insesgado para $Q_y(\beta)$ (véase por ejemplo Särndal *et al.*, 1992), por lo que pasamos a estudiar si dicha propiedad la satisface el estimador de tipo razón $\widehat{Q}_{ym}^r(\beta)$. Para ello, usaremos una aproximación lineal debido a que $\widehat{Q}_{ym}^r(\beta)$ no es una función continua.

El estimador $\widehat{Q}_{ym}^r(\beta)$ puede expresarse asintóticamente como una función lineal de la función de distribución estimada evaluada en el cuantil $Q_y(\beta)$

mediante la representación Bahadur (véase por ejemplo Chambers y Dunstan, 1986):

$$\widehat{Q}_{ym}^r(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))}(\beta - \widehat{F}_{ym}^r(Q_y(\beta))) + o_p(n^{-1/2}), \quad (3.27)$$

donde $f_y(\cdot)$ denota la derivada del valor límite de $F_y(\cdot)$ cuando $N \rightarrow \infty$ y $\widehat{F}_{ym}^r(t)$ denota un estimador de tipo razón para $F_y(t)$, es decir

$$\widehat{F}_{ym}^r(t) = \frac{\widehat{F}_{ym}(t)}{\widehat{F}_{xm}(t)} \widehat{F}_x(t).$$

El estimador $\widehat{Q}_{ym}^r(t)$ es asintóticamente insesgado debido a que $\widehat{F}_{ym}^r(t)$ es un estimador insesgado de $F_y(t)$ (véase Rao *et al.*, 1990). De este modo, $E(\beta - \widehat{F}_{ym}^r(Q_y(\beta))) = 0$, y usando (3.27) puede verse que $E(\widehat{Q}_{ym}^r(\beta)) = Q_y(\beta) + O(n^{-1/2})$.

Puesto que $\widehat{Q}_{ym}^r(\beta)$ y $\widehat{Q}_{yu}(\beta)$ son asintóticamente insesgados para $Q_y(\beta)$, el estimador propuesto $\widehat{Q}_y^R(\beta)$ también lo será. \square

Teorema 3.6 *El estimador compuesto $\widehat{Q}_y^R(\beta)$ es asintóticamente normal.*

Demostración

La normalidad asintótica de la clase propuesta se deriva fácilmente a partir de la expresión (3.24).

En primer lugar, bajo las condiciones (C3.4), (C3.5) y (C3.6), el cuantil muestral $\widehat{Q}_{yu}(\beta)$ es asintóticamente normal. Este resultado puede consultarse en Gross (1980).

Por otro lado, es sabido que el estimador $\widehat{F}_{ym}^r(t)$ es asintóticamente normal. Asumiendo además la aproximación lineal (3.27), puede derivarse fácilmente la normalidad del estimador $\widehat{Q}_{ym}^r(\beta)$.

Por último, usando los dos resultados anteriores, la linealidad de la expresión (3.24) nos permite establecer la normalidad del estimador compuesto propuesto. \square

El siguiente paso en el estudio asintótico del estimador propuesto es la determinación de una expresión para la varianza de dicho estimador. La expresión (3.24) del estimador propuesto nos va a permitir computar su varianza

asintótica a partir de la varianza del estimador basado en la muestra solapada, la varianza del estimador basado en la muestra no solapada y la covarianza entre ambos. Así

$$V(\widehat{Q}_y^R(\beta)) = w^2V_1 + (1-w)^2V_2 + 2w(1-w)C. \quad (3.28)$$

Estas varianzas y covarianzas toman una forma simple cuando la unidades muestrales se seleccionan mediante muestreo aleatorio simple.

Gross (1980) demostró que una expresión asintótica para la varianza del estimador $\widehat{Q}_{yu}(\beta)$ está dada por

$$V(\widehat{Q}_{yu}(\beta)) = \frac{N-u}{N}\beta(1-\beta)(u)^{-1}\{f_y(Q_y(\beta))\}^{-2}. \quad (3.29)$$

Teorema 3.7 *La varianza del estimador de razón propuesto está dada por*

$$\begin{aligned} V(\widehat{Q}_{ym}^r(\beta)) &= \frac{\beta(1-\beta)}{f_y(Q_y(\beta))^2} \left[\left(\frac{1}{m} - \frac{1}{N} \right) + \left(\frac{1}{m} - \frac{1}{n'} \right) R \frac{f_y(Q_y(\beta))}{f_x(Q_x(\beta))} \right] \times \\ &\times \left\{ R \frac{f_y(Q_y(\beta))}{f_x(Q_x(\beta))} + 2 \left(1 - \frac{P_{11}(x,y)}{\beta(1-\beta)} \right) \right\}, \end{aligned} \quad (3.30)$$

donde $P_{11}(x,y)$ denota la proporción de valores en la población para los cuales $x \leq Q_x(\beta)$ e $y \leq Q_y(\beta)$, y $R = Q_y(\beta)/Q_x(\beta)$.

Demostración

Usando propiedades del muestreo bifásico, la expresión asintótica para $V(\widehat{Q}_{ym}^r(\beta))$ puede obtenerse de

$$\begin{aligned} \widehat{Q}_{ym}^r(\beta) - Q_y(\beta) &\cong \widehat{Q}_{ym}(\beta) - Q_y(\beta) + \left(\frac{\widehat{Q}_{xm}(\beta)}{\widehat{Q}_x(\beta)} - 1 \right) (-Q_y(\beta)) \\ &= Q_y(\beta)e_0 + (e_1 - e_2)(-Q_y(\beta)) - e_2(e_1 - e_2)(-Q_y(\beta)), \end{aligned} \quad (3.31)$$

con la notación: $e_0 = \frac{\widehat{Q}_{ym}(\beta)}{Q_y(\beta)} - 1$, $e_1 = \frac{\widehat{Q}_{xm}(\beta)}{Q_x(\beta)} - 1$ y $e_2 = \frac{\widehat{Q}_x(\beta)}{Q_x(\beta)} - 1$.

La expresión asintótica de la varianza del estimador $\widehat{Q}_{ym}^r(\beta)$ se obtiene elevando al cuadrado los dos miembros de (3.31) y posteriormente tomando esperanzas (Notamos que solamente se han considerado términos de orden

uno):

$$\begin{aligned}
V(\widehat{Q}_{ym}^r(\beta)) &= \frac{\beta(1-\beta)}{f_y(Q_y(\beta))^2} \left[\left(\frac{1}{m} - \frac{1}{N} \right) + \left(\frac{1}{m} - \frac{1}{n'} \right) \frac{f_y(Q_y(\beta))}{Q_x(\beta)f_x(Q_x(\beta))} \right] \times \\
&\times (-Q_y(\beta)) \left\{ \frac{f_y(Q_y(\beta))}{Q_x(\beta)f_x(Q_x(\beta))} (-Q_y(\beta)) + 2 \left(\frac{P_{11}(x,y)}{\beta(1-\beta)} - 1 \right) \right\} \\
&= \frac{\beta(1-\beta)}{f_y(Q_y(\beta))^2} \left[\left(\frac{1}{m} - \frac{1}{N} \right) + \left(\frac{1}{m} - \frac{1}{n'} \right) R \frac{f_y(Q_y(\beta))}{f_x(Q_x(\beta))} \right] \times \\
&\times \left\{ R \frac{f_y(Q_y(\beta))}{f_x(Q_x(\beta))} + 2 \left(1 - \frac{P_{11}(x,y)}{\beta(1-\beta)} \right) \right\}.
\end{aligned}$$

Los valores de $E[e_0^2]$, $E[e_1^2]$, $E[e_2^2]$, $E[e_0e_1]$ y $E[e_0e_2]$ pueden verse en Allen *et al.* (2002) y Singh (2003). \square

Teorema 3.8 La covarianza entre los estimadores $\widehat{Q}_{yu}(\beta)$ y $\widehat{Q}_{ym}^r(\beta)$ está dada por

$$Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta)) = \frac{1}{f_y(Q_y(\beta))^2} \frac{-n}{N-n} \left(1 - \frac{n'}{N} \right) \frac{\beta(1-\beta)}{n'}. \quad (3.32)$$

Demostración

Para obtener la covarianza entre los estimadores $\widehat{Q}_{yu}(\beta)$ y $\widehat{Q}_{ym}^r(\beta)$ al primer orden de aproximación, nos basaremos en la propia definición de varianza :

$$\begin{aligned}
&Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta)) = \\
&= Cov(E(\widehat{Q}_{yu}(\beta)/s'), E(\widehat{Q}_{ym}^r(\beta)/s')) + E(Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta)/s')). \quad (3.33)
\end{aligned}$$

Debido a la independencia entre s_u y s_m , el segundo término es cero. En lo que respecta al primer término

$$E(\widehat{Q}_{ym}^r(\beta)/s') = \widehat{Q}_{ys'}(\beta) + o(m^{-1})$$

y

$$E(\widehat{Q}_{yu}(\beta)/s') = \widehat{Q}_{ys'^c}(\beta),$$

donde

$$\widehat{Q}_{ys'}(\beta) = \inf\{t : \widehat{F}_{ys'}(t) \geq \beta\} \quad ; \quad \widehat{Q}_{ys'^c}(\beta) = \inf\{t : \widehat{F}_{ys'^c}(t) \geq \beta\},$$

$$\widehat{F}_{ys'}(t) = \frac{1}{N} \sum_{i \in s'} \frac{\delta(t - y_i)}{\pi_i} \quad ; \quad \widehat{F}_{ys'^c}(t) = \frac{1}{N} \sum_{i \in s'^c} \frac{\delta(t - y_i)}{\pi_i^c}$$

Por otro lado, la representación Bahadur da (véase Kuk y Mak, 1989)

$$\widehat{Q}_{ys^{t'c}}(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))}(\beta - \widehat{F}_{ys^{t'c}}(Q_y(\beta))) + o_p(n^{-1/2}),$$

$$\widehat{Q}_{ys'}(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))}(\beta - \widehat{F}_{ys'}(Q_y(\beta))) + o_p(n^{-1/2}),$$

y de este modo se obtiene

$$Cov(\widehat{Q}_{ys^{t'c}}(\beta), \widehat{Q}_{ys'}(\beta)) \simeq \frac{1}{f_y(Q_y(\beta))^2} Cov(\widehat{F}_{ys'}(Q_y(\beta)), \widehat{F}_{ys^{t'c}}(Q_y(\beta))) = \quad (3.34)$$

$$\frac{1}{f_y(Q_y(\beta))^2} \frac{-n}{N-n} V(\widehat{F}_{ys'}(Q_y(\beta))) = \frac{1}{f_y(Q_y(\beta))^2} \frac{-n}{N-n} \left(1 - \frac{n'}{N}\right) \frac{\beta(1-\beta)}{n'}, \quad (3.35)$$

obteniendo así el resultado (3.32). \square

Sustituyendo los valores (3.29), (3.30) y (3.32) en (3.28), se obtiene la siguiente expresión para la varianza del estimador propuesto

$$V(\widehat{Q}_y^R(\beta)) = C_1 \frac{\left(\frac{n}{1-\chi} - \frac{1}{N}\right) \left[\left(\frac{1}{n\chi} - \frac{1}{N}\right) + C_2 \left(\frac{1}{n\chi} - \frac{1}{n'}\right)\right] - \left[C_1 \frac{-n}{N-n} \left(\frac{1}{n'} - \frac{1}{N}\right)\right]^2}{\left(\frac{n}{1-\chi} - \frac{1}{N}\right) + \left[\left(\frac{1}{n\chi} - \frac{1}{N}\right) + C_2 \left(\frac{1}{n\chi} - \frac{1}{n'}\right)\right] + 2C_1 \frac{-n}{N-n} \left(\frac{1}{n'} - \frac{1}{N}\right)}, \quad (3.36)$$

donde $\chi = m/n$ es la fracción de solapamiento, $C_1 = \frac{\beta(1-\beta)}{f_y(Q_y(\beta))^2}$ y

$$C_2 = R \frac{f_y(Q_y(\beta))}{f_x(Q_x(\beta))} \left\{ R \frac{f_y(Q_y(\beta))}{f_x(Q_x(\beta))} + 2\left(1 - \frac{P_{11}(x, y)}{\beta(1-\beta)}\right) \right\}.$$

El valor $f_y(Q_y(\beta))$ puede estimarse aplicando métodos estándares tal como el kernel (Silverman, 1986). El estimador resultante para $f_y(Q_y(\beta))$ junto con $p_{11}(x, y)$ (la proporción de valores en la muestra para los cuales $x \leq \widehat{Q}_x(\beta)$ y $y \leq \widehat{Q}_y(\beta)$) pueden usarse para proporcionar un estimador consistente de las varianzas asintóticas y los valores óptimos w y $1-w$.

Para completar el estudio asintótico en esta sección, analizaremos la ganancia en precisión del estimador propuesto sobre el estimador $\widehat{Q}_{yn}(\beta)$, el cual está basado exclusivamente en las n unidades muestrales para la segunda ocasión. La varianza de este estimador está dada por

$$V(\widehat{Q}_{yn}(\beta)) = \frac{N-n}{N} \beta(1-\beta)(n)^{-1} \{f_y(Q_y(\beta))\}^{-2}. \quad (3.37)$$

De este modo, la ganancia en precisión, G_1 , de $\widehat{Q}_y^R(\beta)$ sobre $\widehat{Q}_{yn}(\beta)$ está dada por

$$G_1 = \frac{V(\widehat{Q}_{yn}(\beta)) - V(\widehat{Q}_y^R(\beta))}{V(\widehat{Q}_y^R(\beta))}. \quad (3.38)$$

Esta ganancia en precisión dependerá de los tamaños muestrales, del orden del cuantil y de la población objeto de estudio.

El valor óptimo de u que maximiza (3.38) coincide con el valor que minimiza la varianza asintótica (3.36).

Por tanto, el problema es obtener el mínimo en χ de la función $\phi(\chi) = V(\widehat{Q}_y^R(\beta))$ y verificando la condición natural $0 < \chi < 1$. Esta función es monótona en el intervalo $(0, 1)$. El crecimiento o decrecimiento depende del orden del cuantil y de la población en estudio. Por tanto, los valores óptimos para χ estarán próximos a cero (cuando se renueva completamente la muestra al pasar de una ocasión a otra), o bien, estarán próximos a uno (cuando la misma muestra se conserva de una ocasión a otra). Todos estos resultados asintóticos pueden también consultarse en Rueda y Muñoz (2006b).

3.3.4. Propiedades empíricas

El siguiente paso en el análisis de estimador propuesto en muestreo con dos ocasiones sucesivas y usando diseños probabilísticos desiguales consiste en llevar a cabo un estudio de simulación asumiendo distintos tamaños muestrales en todas las muestras y bajo distintos esquemas de muestreo. Para este análisis se usará la población Counties (véase Apéndice A para una descripción completa de esta población).

Como se ha podido comprobar, para la puesta en práctica de un muestreo con dos ocasiones sucesivas es necesario seleccionar tres muestras diferentes, las cuales pueden obtenerse a partir de distintos diseños muestrales. En concreto, estas tres muestras son la muestra de la primera fase, la muestra solapada y la muestra no solapada. En el estudio de simulación de esta sección se usarán las distintas combinaciones de esquemas de muestreo descritas en la Tabla 3.11. El método de Midzuno (véase Apéndice B) se emplea como método de extracción de unidades con probabilidades desiguales, aunque es posible la aplicación del estimador propuesto bajo cualquier otro diseño muestral.

Para cada esquema de muestreo se han generado $B = 1000$ simulaciones con tamaños muestrales $n' = 75$, $n = 25$, $m = 5, \dots, 15$ y $n' = 75$, $n = 50$,

Tabla 3.11: Combinaciones de diseños muestrales usados en muestreo con dos ocasiones sucesivas y probabilidades desiguales.

Acrónimo	Muestra	Tipo de muestreo
<i>SMS</i>	s'	Muestreo aleatorio simple (S)
	s_m	Método de Midzuno (M)
	s_u	Muestreo aleatorio simple (S)
<i>MSS</i>	s'	Método de Midzuno (M)
	s_m	Muestreo aleatorio simple (S)
	s_u	Muestreo aleatorio simple (S)
<i>MMM</i>	s'	Método de Midzuno (M)
	s_m	Método de Midzuno (M)
	s_u	Método de Midzuno (M)

$m = 5, \dots, 30$. El cumplimiento del estimador propuesto se evalúa para los tres cuartiles, $\beta = 0,25, 0,50, 0,75$, en términos de Sesgo Relativo (*SR*) y Eficiencia Relativa (*ER*), donde

$$SR = \frac{1}{B} \sum_{b=1}^B \frac{|\widehat{Q}_y^R(\beta)_b - Q_y(\beta)|}{Q_y(\beta)} \quad ; \quad ER = \frac{ECM[\widehat{Q}_y^R(\beta)]}{ECM[\widehat{Q}_{yn}(\beta)]},$$

siendo b la b -ésima simulación, $\widehat{Q}_y^R(\beta) = w\widehat{Q}_{ym}^r(\beta) + (1-w)\widehat{Q}_{yu}(\beta)$, $ECM[\widehat{Q}_y^R(\beta)] = B^{-1} \sum_{b=1}^B [\widehat{Q}_y^R(\beta)_b - Q_y(\beta)]^2$, y $ECM[\widehat{Q}_{yn}(\beta)]$ se define análogamente para $\widehat{Q}_{yn}(\beta)$, el estimador estándar para el cuantil poblacional basado en la ocasión más reciente.

Notamos que el valor óptimo para la constante w (3.26) depende de varianzas y covarianzas desconocidas, en concreto depende de $V(\widehat{Q}_{ym}^r(\beta))$, $V(\widehat{Q}_{yu}(\beta))$ y $Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta))$. Se usarán técnicas Jackknife (Efron y Tibshirani, 1993) para la estimación de estas expresiones.

Por otro lado, la constante w depende de covarianzas porque la muestra solapada y la no solapada son dependientes, aunque algunos autores ignoran este hecho y consideran tales muestras como independientes, es decir, emplearían la constante

$$w^* = \frac{V(\widehat{Q}_{yu}(\beta))}{V(\widehat{Q}_{ym}^r(\beta)) + V(\widehat{Q}_{yu}(\beta))},$$

donde $Cov(\widehat{Q}_{yu}(\beta), \widehat{Q}_{ym}^r(\beta))$ estaría omitida. Con el fin de analizar este hecho en la práctica, el estimador propuesto basado en la constante w^* (asumiendo

Figura 3.7: Eficiencia Relativa para el diseño muestral *SMS*.

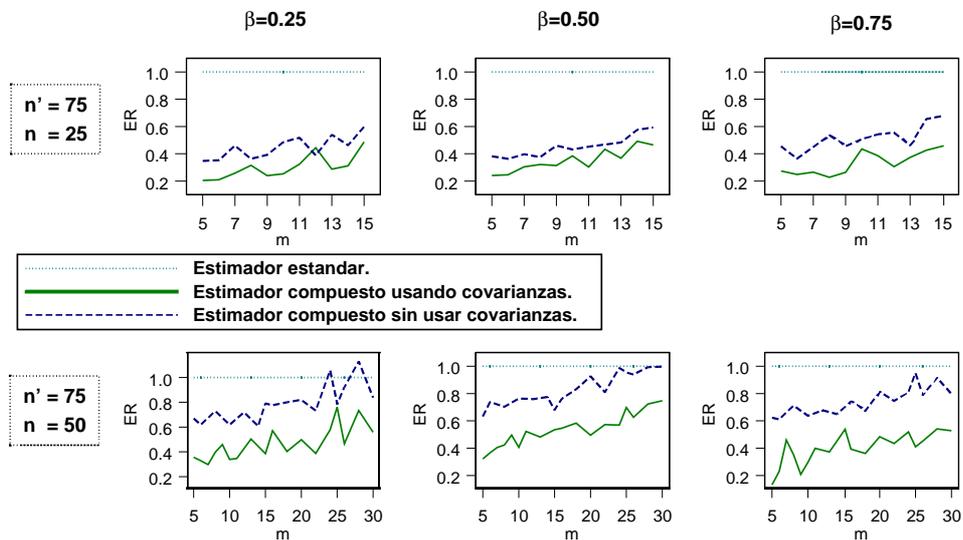


Figura 3.8: Eficiencia Relativa para el diseño muestral *MSS*.

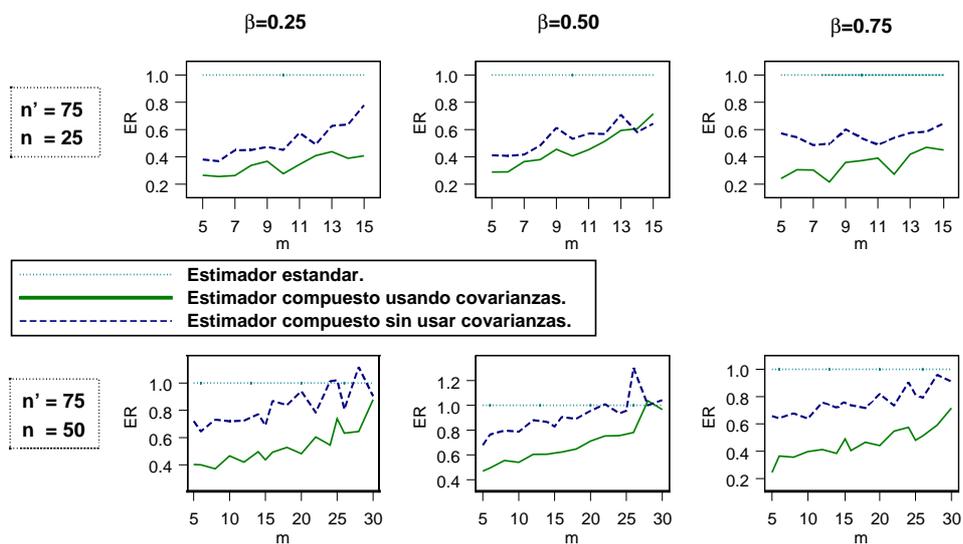
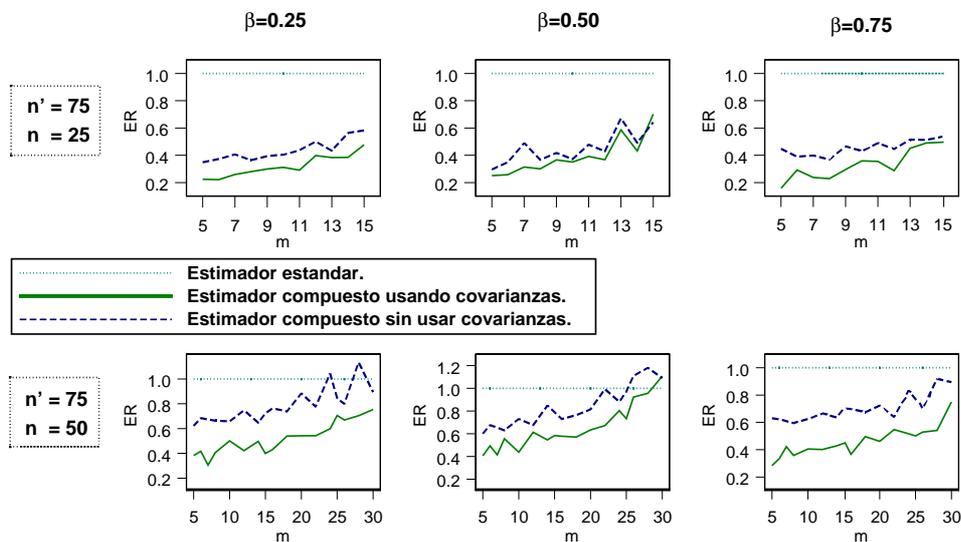


Figura 3.9: Eficiencia Relativa para el diseño muestral *MMM*.



que existe independencia entre las muestras, por lo que se ignoran las covarianzas) ha sido incluido en el estudio de simulación.

En primer lugar analizaremos la eficiencia de los estimadores, la cual puede observarse en las Figuras 3.7, 3.8 y 3.9, en donde se representa la Eficiencia Relativa de los distintos estimadores y combinaciones de diseños y tamaños muestrales. La variación en el cumplimiento de los estimadores desde distintas perspectivas puede por tanto observarse. Notamos que las curvas continuas corresponden al estimador propuesto (usando covarianzas), mientras que las curvas discontinuas corresponden al estimador compuesto que no emplea covarianzas. Las líneas horizontales representan al estimador estándar.

En los tres casos, los resultados obtenidos muestran un buen cumplimiento del estimador propuesto, el cual es siempre más eficiente que el estimador estándar, excepto para el caso de fracciones de solapamiento elevadas. Cuando la fracción de solapamiento aumenta, decrece la Eficiencia Relativa para el estimador propuesto en comparación con el estimador estándar.

En lo que respecta al comportamiento del uso o no de covarianzas en el estimador propuesto, puede comprobarse que se obtiene una ligera mejoría en eficiencia cuando se tiene en cuenta las covarianzas en la construcción del estimador, teniendo por tanto sentido la hipótesis de dependencia entre el es-

Figura 3.10: Sesgo Relativo para el diseño muestral *SMS*.

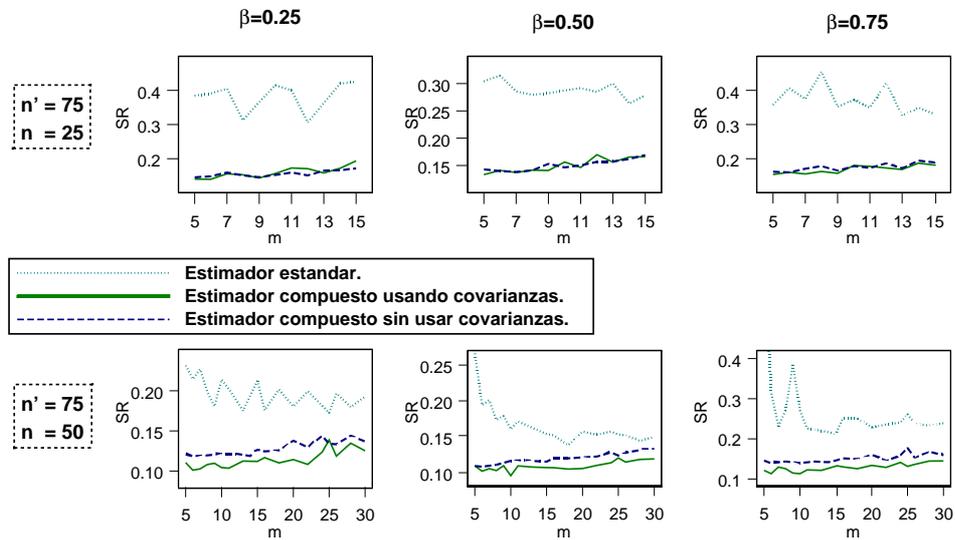


Figura 3.11: Sesgo Relativo para el diseño muestral *MSS*.

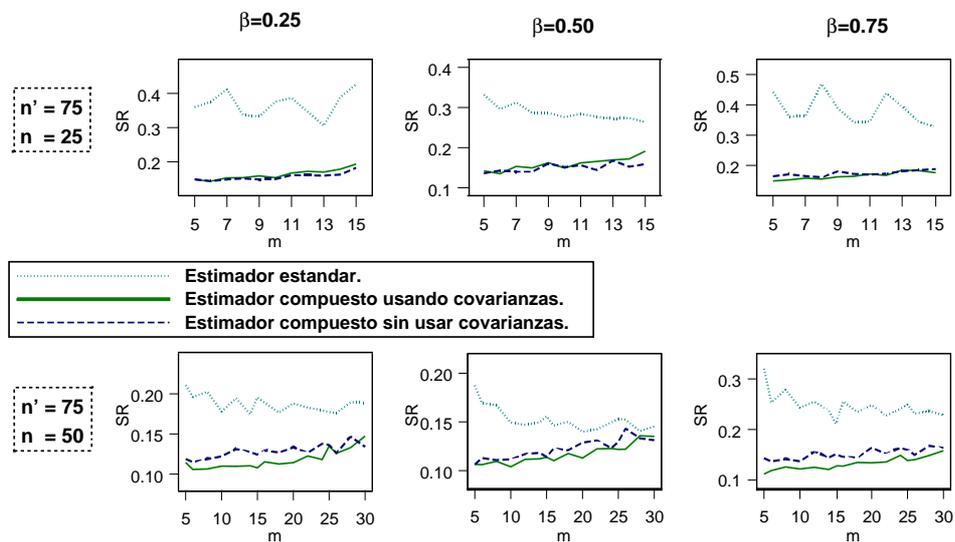
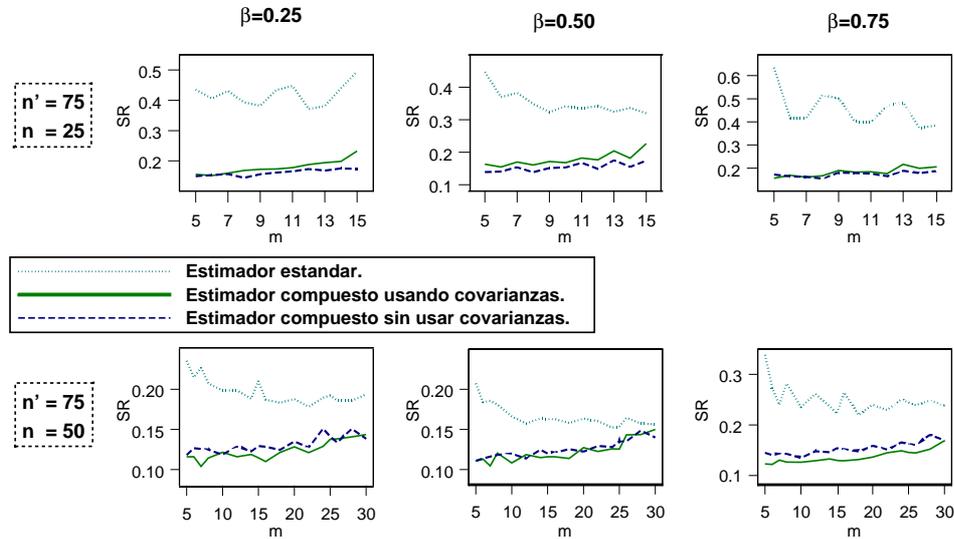


Figura 3.12: Sesgo Relativo para el diseño muestral *MMM*.



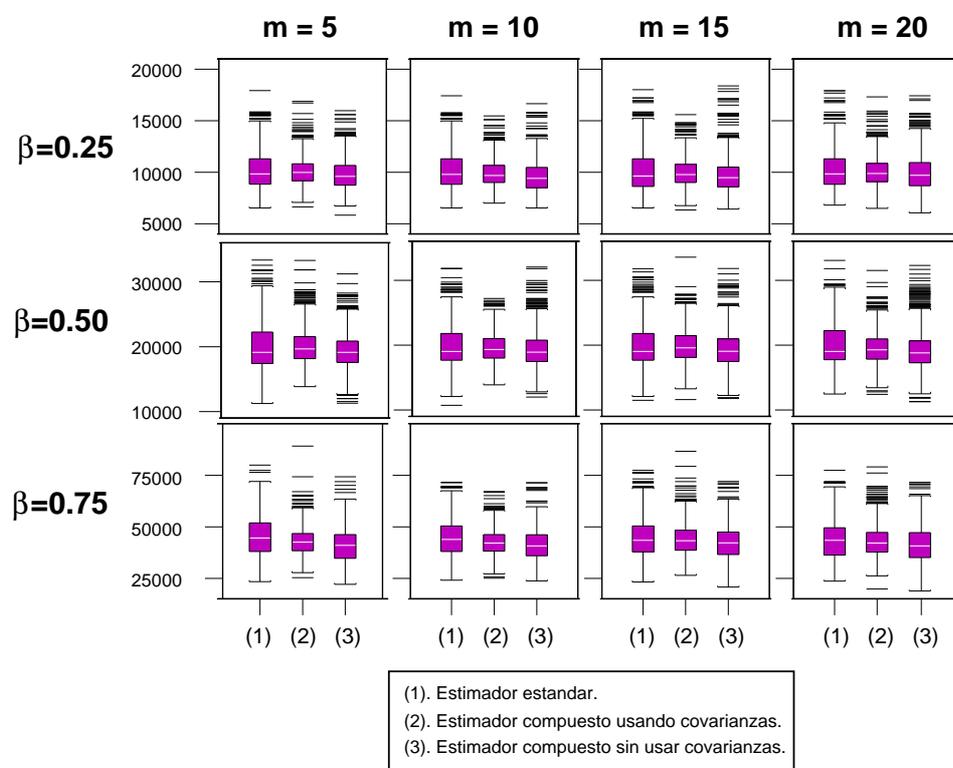
timador de la muestra no solapada y el estimador propuesto para la parte solapada. Puede además observarse que la ganancia en precisión sobre el estimador que omite las covarianzas es mayor a medida que aumenta el tamaño muestral de la ocasión más reciente. En resumen, estos resultados recomiendan el uso de covarianzas en el estimador propuesto para la estimación de cuantiles bajo un muestreo con dos ocasiones sucesivas y probabilidades desiguales.

El análisis del Sesgo Relativo de los distintos estimadores puede seguirse en las Figuras 3.10, 3.11 y 3.12. Esta es otra medida que nos permitirá evaluar el comportamiento del estimador propuesto y estudiar las propiedades de dicho estimador con respecto al estimador estándar.

A partir de estas figuras puede observarse un similar comportamiento de los estimadores al obtenido en el estudio de la Eficiencia Relativa. Los valores del Sesgo Relativo para los estimadores propuestos están siempre por debajo de 0.2, y en algunas ocasiones son inferiores a 0.1, mientras que el Sesgo Relativo para el estimador estándar es bastante mayor llegando incluso a 0.6.

Por último, analizaremos los valores observados de los estimadores mediante diagramas de cajas con bigotes. Por brevedad se ha considerado el diseño muestral *SMS* y los tamaños muestrales $n' = 75$, $n = 50$ y $m = 5, 10, 15, 20$. La Figura 3.13 nos da tal información para los tres cuantiles. También en este

Figura 3.13: Diagrama de caja con bigotes para los valores de los distintos estimadores. Se asume el diseño muestral *SMS* y tamaños muestrales $n' = 75$ y $n = 50$.



estudio se comprueba que el estimador propuesto presenta el mejor comportamiento, al obtenerse estimadores menos dispersos en comparación con el estimador estándar y el estimador que omite las covarianzas.

Notamos que se han realizado otras simulaciones con distintos tamaños muestrales a los usados en los estudios anteriores. En todos los casos los resultados confirman el buen comportamiento del estimador propuesto frente a sus competidores. También se ha observado que la ganancia en precisión del estimador propuesto es mejor a medida que el tamaño muestral en la primera ocasión aumenta con respecto al tamaño de la segunda ocasión. Por otro lado, cuando el tamaño muestral en la primera ocasión es menor que el tamaño en la segunda, se obtiene una menor ganancia en precisión, y esta ganancia disminuye a medida que aumenta la diferencia entre tamaños muestrales. Este resultado es lógico porque si n' es mayor en comparación con n , la primera muestra proporcionará mayor información, y el estimador de tipo razón basado en la muestra solapada presentará un menor grado de error, por lo que es de esperar que el estimador propuesto mejore también en precisión.

3.3.5. Generalización a múltiples variables auxiliares

En las secciones anteriores se han definido y examinado los estimadores de cuantiles en muestreo con dos ocasiones sucesivas y para un diseño muestral arbitrario. Tanto estos estimadores como los ya definidos en la literatura del muestreo en poblaciones finitas, consideran el uso de una única variable auxiliar, aunque la investigación dispongan de otras variables que también estén altamente correlacionadas con la variable principal. Esto sugiere que se podría hacer un uso más efectivo de la información auxiliar. Asumiendo muestreo aleatorio simple, en este apartado se define una clase de estimadores que pueden obtenerse a partir de un vector multivariante de variables auxiliares. En concreto, esta clase está formada por un estimador de tipo razón construido a partir de todas las variables auxiliares disponibles en las muestras que están solapadas y por un estimador de la variable de interés en la muestra no solapada de la ocasión más reciente. El estimador óptimo en el sentido de minimizar la varianza de esta clase también será calculado.

Por tanto, en la presente sección y en 3.3.6 y 3.3.7 asumiremos que en la primera ocasión se dispone de p variables auxiliares (denotadas por x_1, \dots, x_p) en lugar de una única variable auxiliar x . Con esta información y con la obtenida a partir de la muestra s' , es posible calcular los estimadores $\hat{Q}_{xim}(\beta)$, $i = 1, \dots, p$, que son los cuantiles de orden β estimados a partir de la mues-

tra solapada, y los estimadores $\widehat{Q}_{xi}(\beta)$, $i = 1, \dots, p$, que serán las estimaciones basadas en la muestra s' . Estos estimadores se construyen como en (3.22) y (3.23), respectivamente, pero sustituyendo la variable x por x_i , para $i = 1, \dots, p$

Siguiendo a Olkin (1958), se propone el siguiente estimador de tipo razón multivariante de $Q_y(\beta)$ basado en la parte solapada:

$$\widehat{Q}_{ym}^{MR}(\beta) = \sum_{1 \leq i \leq p} w_i \frac{\widehat{Q}_{ym}(\beta)}{\widehat{Q}_{xim}(\beta)} \widehat{Q}_{xi}(\beta) = \sum_{1 \leq i \leq p} w_i \widehat{Q}_{yrim}(\beta). \quad (3.39)$$

Los pesos w_i (verificando $\sum_{1 \leq i \leq p} w_i = 1$) se obtienen de modo que maximizan la precisión del estimador $\widehat{Q}_{ym}^{MR}(\beta)$. Se usa el criterio de mínima varianza para obtener estas cantidades. Sabido esto, la varianza de este estimador viene dada por

$$V(\widehat{Q}_{ym}^{MR}(\beta)) = \sum_{1 \leq i \leq p} w_i^2 V(\widehat{Q}_{yrim}(\beta)) + 2 \sum_{i < j} w_i w_j Cov(\widehat{Q}_{yrim}(\beta), \widehat{Q}_{yrjm}(\beta)).$$

Esta última ecuación puede escribirse como $V(\widehat{Q}_{ym}^{MR}(\beta)) = w' B w$, donde $w = (w_1, \dots, w_p)'$, $B = (b_{ij})$ y $b_{ij} = Cov(\widehat{Q}_{yrim}(\beta), \widehat{Q}_{yrjm}(\beta))$ para $i, j = 1, \dots, p$. Para obtener el valor extremo usaremos la desigualdad de Cauchy-Schwarz, y puesto que B es semidefinida positiva, se obtiene que el valor óptimo w está dado por

$$w_{opt} = \frac{B^{-1}e}{e' B^{-1}e},$$

donde $e = (1, \dots, 1)'$. Por tanto, la mínima varianza obtenida a partir de w_{opt} será

$$V_{min}(\widehat{Q}_{ym}^{MR}(\beta)) = \frac{1}{e' B^{-1}e}.$$

Asumiendo muestreo en dos ocasiones sucesivas, se propone el siguiente estimador compuesto que combina el anterior estimador de tipo razón múltiple basado en la muestra solapada con el estimador de la muestra no solapada:

$$\widehat{Q}_y(\beta) = W \widehat{Q}_{ymopt}^{MR}(\beta) + (1 - W) \widehat{Q}_{yu}(\beta), \quad (3.40)$$

donde $\widehat{Q}_{ymopt}^{MR}(\beta)$ está dado por el estimador $\widehat{Q}_{ym}^{MR}(\beta)$ cuando se considera el valor óptimo de w , esto es w_{opt} , mientras que W es una constante que satisface $0 < W < 1$ y que es escogida de modo que el estimador $\widehat{Q}_y(\beta)$ presente la mínima varianza dentro la clase anterior. Un simple cálculo demuestra que

$$W_{opt} = \frac{V(\widehat{Q}_{yu}(\beta))}{V(\widehat{Q}_{yu}(\beta)) + V(\widehat{Q}_{ymopt}^{MR}(\beta))}. \quad (3.41)$$

En resumen, el estimador propuesto que presenta las propiedades óptimas en términos de mínima varianza está dado por

$$\widehat{Q}_{yopt}(\beta) = W_{opt} \widehat{Q}_{ymopt}^{MR}(\beta) + (1 - W_{opt}) \widehat{Q}_{yu}(\beta), \quad (3.42)$$

y varianza viene dada por

$$V(\widehat{Q}_{yopt}(\beta)) = W_{opt}^2 V(\widehat{Q}_{ymopt}^{MR}(\beta)) + (1 - W_{opt})^2 V(\widehat{Q}_{yu}(\beta)), \quad (3.43)$$

la cual puede también escribirse como

$$V(\widehat{Q}_{yopt}(\beta)) = \frac{V(\widehat{Q}_{yu}(\beta))V(\widehat{Q}_{ymopt}^{MR}(\beta))}{V(\widehat{Q}_{yu}(\beta)) + V(\widehat{Q}_{ymopt}^{MR}(\beta))}. \quad (3.44)$$

3.3.6. Propiedades teóricas

El siguiente paso en el estudio del estimador propuesto $\widehat{Q}_{yopt}(\beta)$ es la determinación de sus propiedades más importantes, además de la propiedad de mínima varianza ya comentada. En concreto se establece la normalidad de dicho estimador y su correspondiente varianza exacta.

Los resultados que se establecen se derivan asumiendo las condiciones (C3.4), (C3.5) y (C3.6).

Teorema 3.9 *El estimador de razón multivariante $\widehat{Q}_{ym}^{MR}(\beta)$ dado por (3.39) y la clase propuesta de estimadores $\widehat{Q}_y(\beta)$ dada por (3.40) son asintóticamente normales.*

Demostración

En primer lugar, los cuantiles muestrales $\widehat{Q}_{yu}(\beta)$, $\widehat{Q}_{ym}(\beta)$, $\widehat{Q}_{xi}(\beta)$ y $\widehat{Q}_{xim}(\beta)$ son asintóticamente normales como se demostró en Gross (1980).

Sean las siguientes funciones de este estimador

$$H_1(\widehat{Q}_{ym}(\beta), \widehat{Q}_{x1}(\beta), \dots, \widehat{Q}_{xp}(\beta), \widehat{Q}_{x1m}(\beta), \dots, \widehat{Q}_{xpm}(\beta)) = \sum_{1 \leq i \leq p} w_i \frac{\widehat{Q}_{ym}(\beta)}{\widehat{Q}_{xim}(\beta)} \widehat{Q}_{xi}(\beta).$$

H_1 es una función continua con derivadas parciales de primer y segundo orden y continuas en las cercanías de: $(Q_y, Q_{x1}, \dots, Q_{xp})$. Bajo esta situación y usando los resultados de Cramer (1946), $\widehat{Q}_{ym}^{MR}(\beta)$ es asintóticamente normal.

La normalidad asintótica de la clase propuesta de estimadores se deriva fácilmente como consecuencia de la expresión lineal de la clase. \square

La normalidad asintótica del estimador $\widehat{Q}_{yopt}(\beta)$ también se deriva al pertenecer este estimador a la clase (3.40).

La linealidad de la clase de estimadores también nos permitirá computar sus varianzas. Para ello, será necesario conocer las varianzas del estimador de razón multivariante basado en la muestra solapada y el estimador que solamente envuelve a la muestra no solapada, $\widehat{Q}_{yu}(\beta)$, como puede verse en (3.43) y (3.44).

Gross (1980) demostró que una expresión asintótica para la varianza del estimador $\widehat{Q}_{yu}(\beta)$ está dada por

$$V(\widehat{Q}_{yu}(\beta)) = \frac{N-u}{N} \beta(1-\beta)(u)^{-1} \{f_y(Q_y(\beta))\}^{-2}. \quad (3.45)$$

Teorema 3.10 *La varianza de $V(\widehat{Q}_{yrim}(\beta))$, con $i = 1, \dots, p$, y la covarianza entre $\widehat{Q}_{yrim}(\beta)$ y $\widehat{Q}_{yrjm}(\beta)$, con $i, j = 1, \dots, p$ vienen dadas por*

$$V(\widehat{Q}_{yrim}(\beta)) = \frac{\beta(1-\beta)}{f_y(Q_y(\beta))^2} \left[\left(\frac{1}{m} - \frac{1}{N} \right) + \left(\frac{1}{m} - \frac{1}{n'} \right) R_i \frac{f_y(Q_y(\beta))}{f_{xi}(Q_{xi}(\beta))} \times \right. \\ \left. \times \left\{ R_i \frac{f_y(Q_y(\beta))}{f_{xi}(Q_{xi}(\beta))} + 2 \left(1 - \frac{P_{11}(y, x_i)}{\beta(1-\beta)} \right) \right\} \right], \quad (3.46)$$

y

$$Cov(\widehat{Q}_{yrim}(\beta), \widehat{Q}_{yrjm}(\beta)) = \frac{\beta(1-\beta)}{f_y(Q_y(\beta))^2} \left[\left(\frac{1}{m} - \frac{1}{N} \right) + \right. \\ \left(\frac{1}{n'} - \frac{1}{m} \right) R_i \frac{f_y(Q_y(\beta))}{f_{xi}(Q_{xi}(\beta))} \left(\frac{P_{11}(y, x_i)}{\beta(1-\beta)} - 1 \right) + \\ \left(\frac{1}{n'} - \frac{1}{m} \right) R_j \frac{f_y(Q_y(\beta))}{f_{xj}(Q_{xj}(\beta))} \left(\frac{P_{11}(y, x_j)}{\beta(1-\beta)} - 1 \right) - \\ \left. \left(\frac{1}{n'} - \frac{1}{m} \right) R_i R_j \frac{f_y^2(Q_y(\beta))}{f_{xi}(Q_{xi}(\beta)) f_{xj}(Q_{xj}(\beta))} \left(\frac{P_{11}(x_i, x_j)}{\beta(1-\beta)} - 1 \right) \right], \quad (3.47)$$

donde $P_{11}(y, x_i)$ denota la proporción de valores en la población para los cuales $y \leq Q_y(\beta)$ y $x_i \leq Q_{xi}(\beta)$, y $R_i = Q_y(\beta)/Q_{xi}(\beta)$.

Demostración

El estimador $\widehat{Q}_{yrim}(\beta)$ puede expresarse como

$$\widehat{Q}_{yrim}(\beta) = Q_y(\beta)(1 + e_0)(1 + e_{2i})(1 - e_{1i} + e_{1i}^2 + \dots), \quad (3.48)$$

donde $e_0 = \frac{\widehat{Q}_{ym}(\beta)}{Q_y(\beta)} - 1$, $e_{1i} = \frac{\widehat{Q}_{xim}(\beta)}{Q_{xi}(\beta)} - 1$ y $e_{2i} = \frac{\widehat{Q}_{xi}(\beta)}{Q_{xi}(\beta)} - 1$, $i = 1, \dots, p$.

Considerando la expansión de series de Taylor se obtiene la expresión

$$\begin{aligned} & (\widehat{Q}_{yrim}(\beta) - Q_y(\beta))(\widehat{Q}_{yrim}(\beta) - Q_y(\beta)) \cong \\ & \cong Q_y(\beta)^2(e_0 + e_{2i} - e_{1i} + e_{1i}^2 - e_{1i}e_{2i} - e_{1i}e_0 + e_0e_{2i} + \dots) \\ & (e_0 + e_{2j} - e_{1j} + e_{1j}^2 - e_{1j}e_{2j} - e_{1j}e_0 + e_0e_{2j} + \dots). \end{aligned}$$

La expresión asintótica de la covarianza de los estimadores $\widehat{Q}_{yrim}(\beta)$ y $\widehat{Q}_{yrim}(\beta)$ se obtiene tomando esperanzas (se han considerado solamente términos de orden uno). Las esperanzas de las variables e_i pueden derivarse de Singh (2003):

$$E[e_0^2] = \frac{N - m}{Nm} \beta(1 - \beta)(Q_y(\beta)f_y(Q_y(\beta)))^{-2},$$

$$E[e_{1i}^2] = \frac{N - m}{Nm} \beta(1 - \beta)(Q_{xi}(\beta)f_{xi}(Q_{xi}(\beta)))^{-2},$$

$$E[e_{2i}^2] = E[e_{1i}e_{2i}] = \frac{N - n'}{Nn'} \beta(1 - \beta)(Q_{xi}(\beta)f_{xi}(Q_{xi}(\beta)))^{-2},$$

$$E[e_0e_{1i}] = \frac{N - m}{Nm} (P_{11}(y, x_i) - \beta(1 - \beta))(Q_{xi}(\beta)Q_y(\beta)f_{xi}(Q_{xi}(\beta))f_y(Q_y(\beta)))^{-1},$$

$$E[e_0e_{2i}] = \frac{N - n'}{Nn'} (P_{11}(y, x_i) - \beta(1 - \beta))(Q_{xi}(\beta)Q_y(\beta)f_{xi}(Q_{xi}(\beta))f_y(Q_y(\beta)))^{-1},$$

$$\begin{aligned} E[e_{1j}e_{2i}] &= E[e_{2j}e_{2i}] = \frac{N - n'}{Nn'} (P_{11}(x_j, x_i) - \beta(1 - \beta)) \times \\ & \times (Q_{xj}(\beta)f_{xj}(Q_{xj}(\beta))Q_{xi}(\beta)f_{xi}(Q_{xi}(\beta)))^{-1}, \end{aligned}$$

$$\begin{aligned} E[e_{1j}e_{1i}] &= \frac{N - m}{Nm} (P_{11}(x_j, x_i) - \beta(1 - \beta)) \times \\ & \times (Q_{xj}(\beta)f_{xj}(Q_{xj}(\beta))Q_{xi}(\beta)f_{xi}(Q_{xi}(\beta)))^{-1}. \end{aligned}$$

Sustituyendo estos valores y operando adecuadamente, se obtiene la expresión dada en (3.47). \square

Por tanto, usando las expresiones (3.45) (3.46) y (3.47), la matriz B , la varianza del estimador propuesto dado en (3.43) o (3.44) y el valor W_{opt} definido en (3.41) quedan determinadas.

3.3.7. Propiedades empíricas

En la Sección 3.3.5 se ha definido un estimador óptimo dentro de la clase (3.40). La normalidad y la varianza asintótica de este estimador se ha demostrado en la Sección 3.3.6. El siguiente paso en este estudio es comprobar la exactitud de este estimador. En este apartado, la eficiencia del estimador propuesto y su varianza serán analizadas. En primer lugar, se analiza la ganancia en eficiencia de la varianza asintótica del estimador $\widehat{Q}_{yopt}(\beta)$ con la varianza de $\widehat{Q}_{yn}(\beta)$, el estimador estándar basado en la ocasión más reciente y el cual está dado en (3.15). A continuación, el comportamiento de estos estimadores serán contrastados en una situación real mediante un estudio empírico.

En ambos estudios se usaran dos poblaciones naturales: la población Counties y la población Turismos (véase Apéndice A). La población turismos resulta interesante en este caso porque dispone de cuatro variables auxiliares. Se pueden comparar los varios estimadores usando un número distinto de variables auxiliares, de modo que pueda observarse la evolución de la ganancia en precisión al aumentar el número de variables auxiliares usadas en la etapa de estimación.

Comparaciones teóricas

El primer estudio consiste en comparar la varianza del estimador óptimo propuesto dado en (3.44) con la varianza del estimador frecuentemente usado, $\widehat{Q}_{yn}(\beta)$. Este estudio nos permitirá conocer el comportamiento de las varianzas teóricas de los estimadores. Gross (1980) comprobó que una expresión asintótica para la varianza del estimador $\widehat{Q}_{yn}(\beta)$ está dada por

$$V(\widehat{Q}_{yn}(\beta)) = \frac{N-n}{N} \beta(1-\beta)(n)^{-1} \{f_y(Q_y(\beta))\}^{-2}.$$

En las Figuras 3.14 y 3.15, las varianzas teóricas de los estimadores $\widehat{Q}_{yopt}(\beta)$ y $\widehat{Q}_{yn}(\beta)$ son comparadas por medio de sus cocientes, esto es, las figuras mues-

Figura 3.14: Ratios Teóricos entre la varianza del estimador óptimo propuesto y la varianza del estimador estándar bajo la población Counties y el cuantil de orden $\beta = 0,5$.

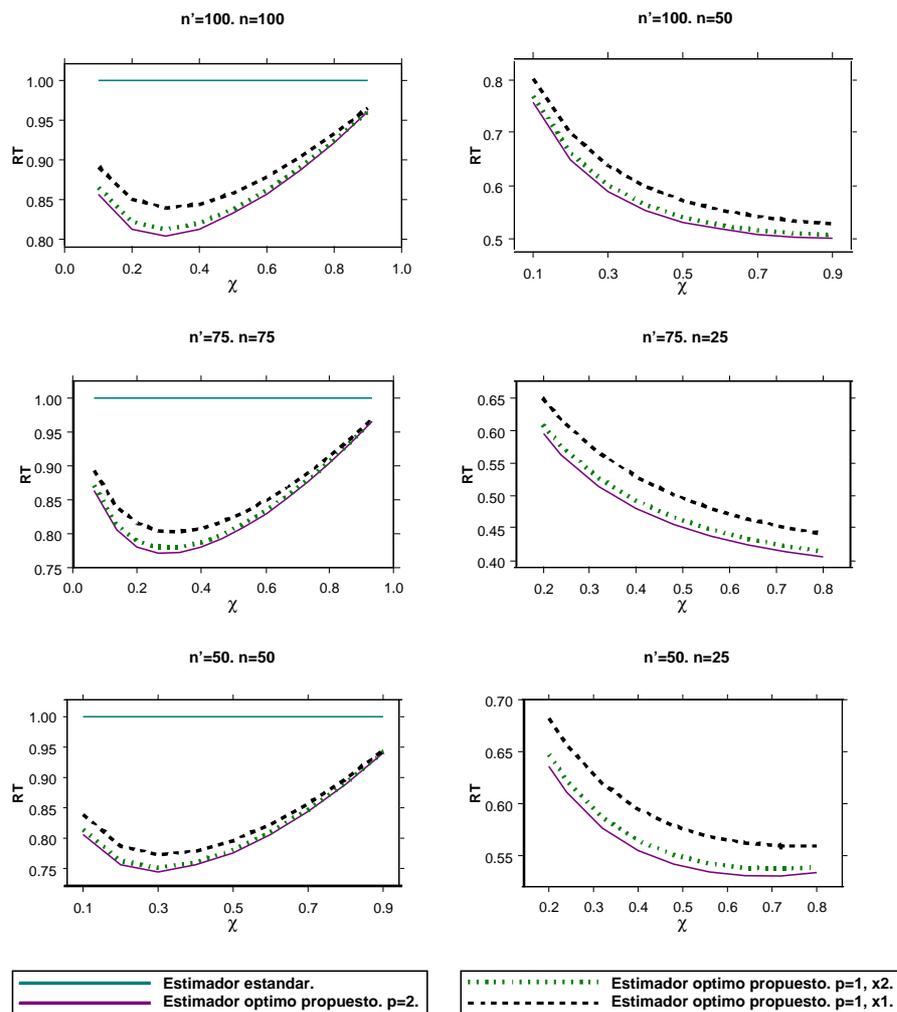
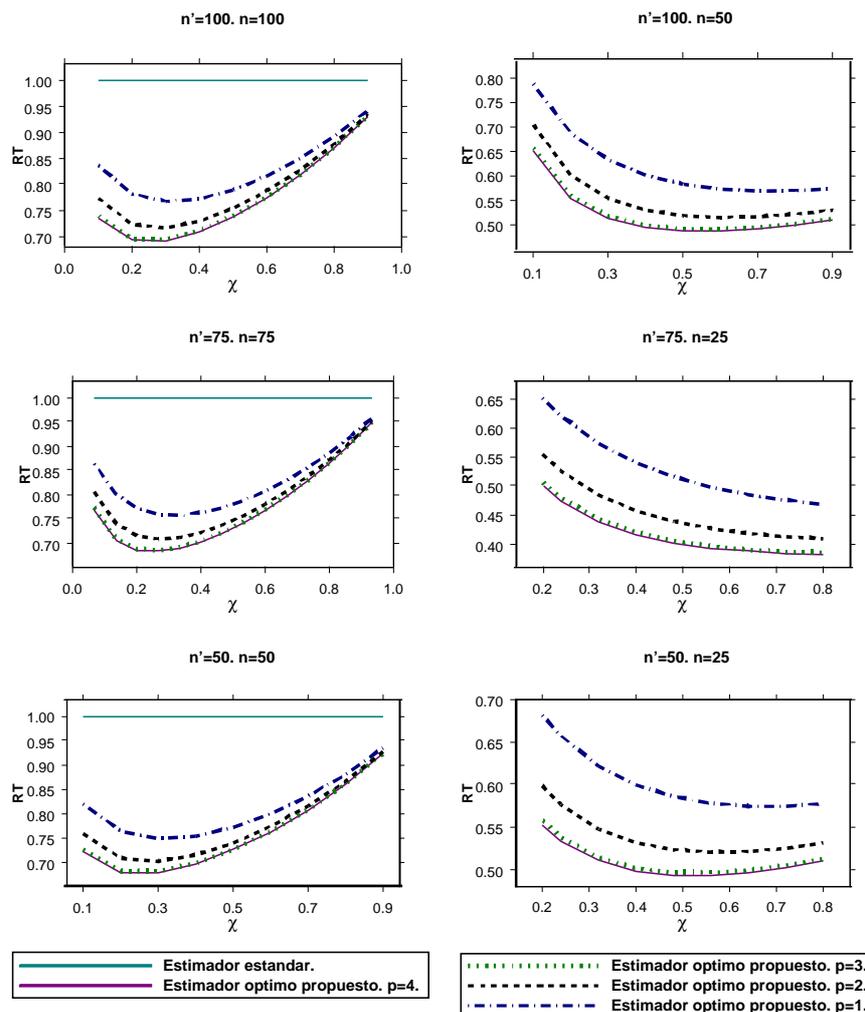


Figura 3.15: Ratios Teóricos entre la varianza del estimador óptimo propuesto y la varianza del estimador estándar bajo la población Turismos y el cuantil de orden $\beta = 0,5$.



tran los Ratio Teóricos $RT = V(\widehat{Q}_{yopt}(\beta))/V(\widehat{Q}_{yn}(\beta))$. En este estudio, se representan diferentes valores de m en el eje de abscisas y el estimador propuesto se ha obtenido para cada valor de p ($p = 1, 2$ en la población Counties y $p = 1, 2, 3, 4$ en la población Turismos). Las líneas horizontales muestran los RT para el estimador $\widehat{Q}_{yn}(\beta)$. Notamos que valores de RT por debajo de 1 indican que $V(\widehat{Q}_{yopt}(\beta))$ es menor que $V(\widehat{Q}_{yn}(\beta))$, y por tanto el estimador propuesto es más eficiente.

De estas comparaciones teóricas, se pueden destacar la siguientes conclusiones:

1. Para ambas poblaciones, el estimador propuesto parece tener uniformemente menor varianza que el estimador estándar, $\widehat{Q}_{yn}(\beta)$, y a su vez menor varianza que el estimador propuesto cuando éste utiliza una única variable auxiliar.
2. Las mejores propiedades se obtienen cuando se usan todas las variables auxiliares.
3. Cuando los tamaños muestrales en ambas ocasiones son iguales, la fracción de solapamiento óptima está entre 0.2 y 0.4. Una fracción de solapamiento más alta resulta apropiada cuando el tamaño muestral en la ocasión reciente es menor que el tamaño muestral de la primera ocasión.
4. En ambas poblaciones, los ratios más bajos se obtienen cuando los tamaños muestrales son $n' = 75$ y $n = 25$, en cuyo caso los RT , para valores grandes de χ , son aproximadamente iguales a 0.4, esto es, la varianza asintótica del estimador propuesto presenta una mejoría del 60% con respecto a la varianza asintótica del estimador estándar.

Estudio empírico

El siguiente paso consiste en llevar a cabo un estudio de simulación con el fin de revelar la ganancia en eficiencia de $\widehat{Q}_{yopt}(\beta)$ con respecto a $\widehat{Q}_{yn}(\beta)$ en una situación real. De nuevo, las poblaciones Counties y Turismos serán usadas. Este estudio también muestra el comportamiento de $\widehat{Q}_{yopt}(\beta)$ cuando este estimador usa un número diferente de variables auxiliares.

Se generan $B = 1000$ muestras independientes bajo muestreo con dos ocasiones sucesivas. Todas las muestras (solapadas y no solapadas) se obtienen bajo muestreo aleatorio simple. El cumplimiento de estos estimadores se evalúa

Figura 3.16: Eficiencia Relativa para los estimadores óptimo propuesto y estándar en la población Counties y para el cuantil de orden $\beta = 0,5$.

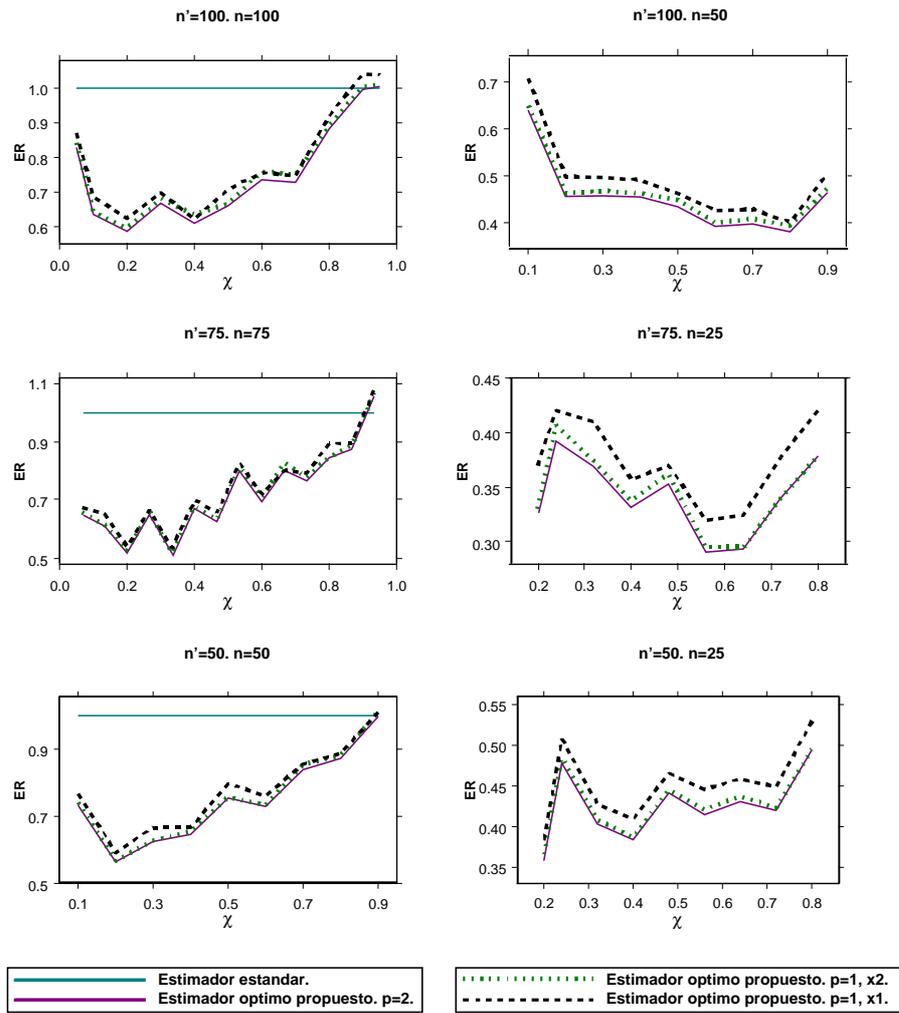


Figura 3.17: Eficiencia Relativa para los estimadores óptimo propuesto y estándar en la población Turismos y para el cuantil de orden $\beta = 0,5$.

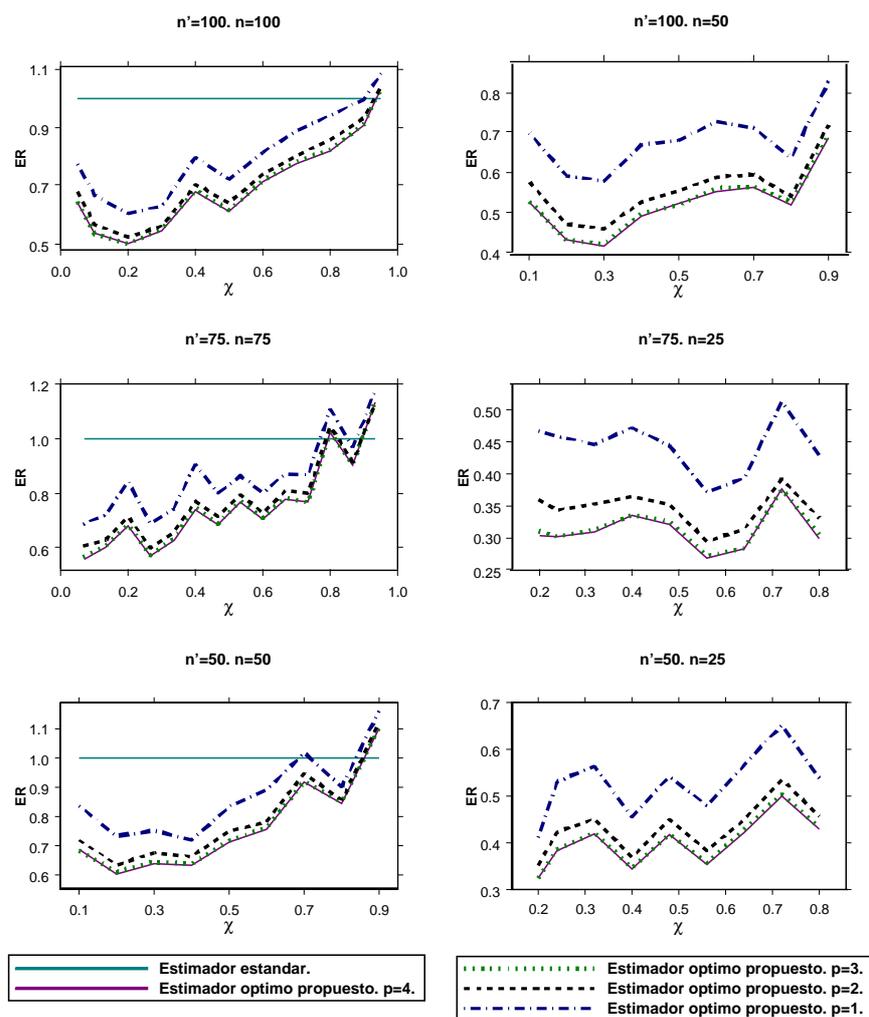


Figura 3.18: Evolución de los valores W_{opt} usados por el estimador óptimo propuesto en la población Counties y para el cuantil de orden $\beta = 0,5$.

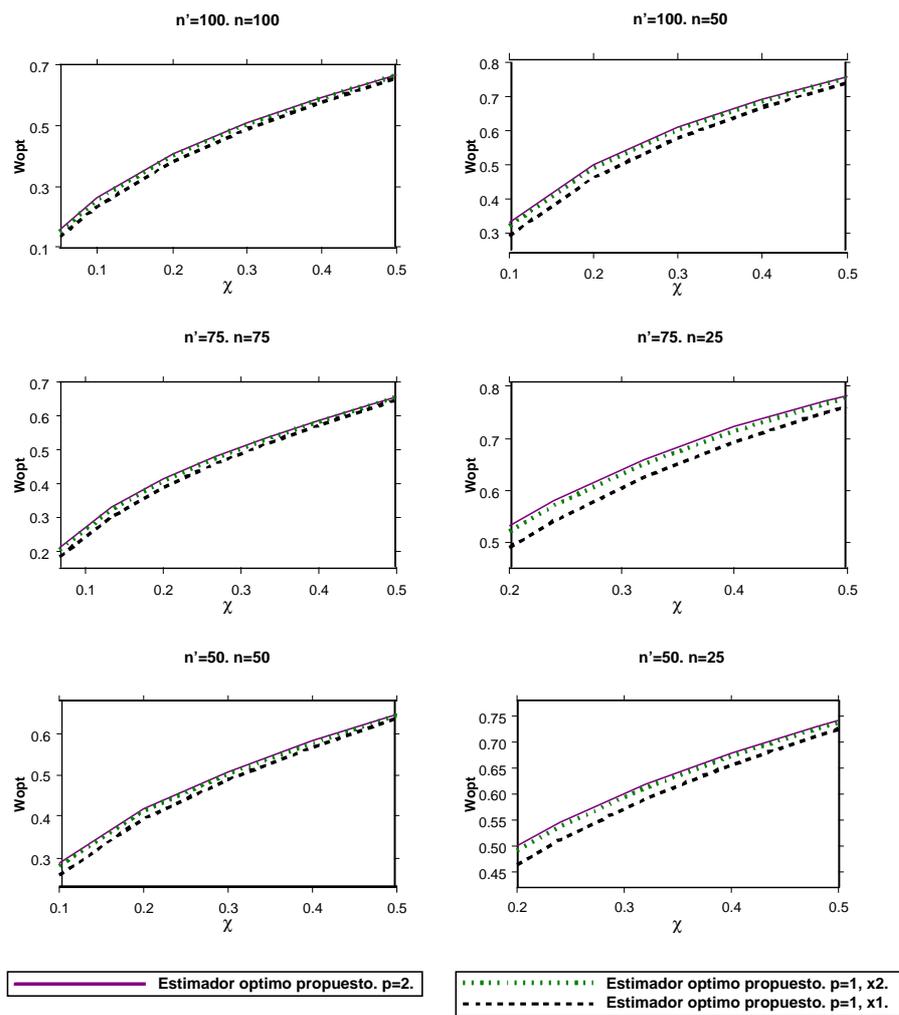
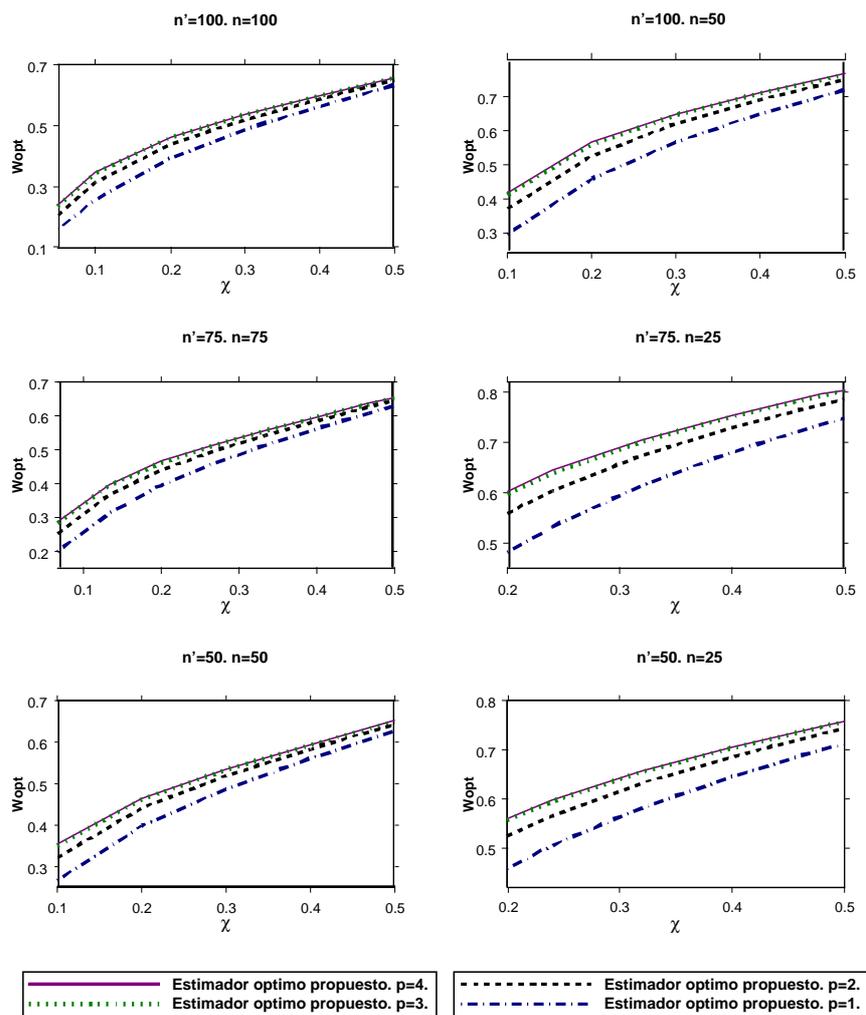


Figura 3.19: Evolución de los valores W_{opt} usados por el estimador óptimo propuesto en la población Turismos y para el cuantil de orden $\beta = 0,5$.



para el cuantil de orden $\beta = 0,5$ en términos de Sesgo Relativo (SR) y Eficiencia Relativa (ER), con

$$SR = \frac{1}{B} \sum_{b=1}^B \frac{\widehat{Q}_{yopt}(\beta)_b - Q_y(\beta)}{Q_y(\beta)} \quad ; \quad ER = \frac{ECM[\widehat{Q}_{yopt}(\beta)]}{ECM[\widehat{Q}_{yn}(\beta)]},$$

donde b indica la b ésima simulación, el Error Cuadrático Medio empírico está dado por $ECM[\widehat{Q}_{yopt}(\beta)] = B^{-1} \sum_{b=1}^B [\widehat{Q}_{yopt}(\beta)_b - Q_y(\beta)]^2$, y donde $ECM[\widehat{Q}_{yn}(\beta)]$ se define de modo similar para $\widehat{Q}_{yn}(\beta)$. Por tanto, el comportamiento empírico del estimador propuesto se compara con el estimador estándar mediante diferentes valores de p .

Las generaciones aleatorias, cálculos y obtención de estimadores se han obtenido mediante el programa R . Los detalles de la programación están disponibles en el Apéndice C.

Las Figuras 3.16 y 3.17 representan la ER obtenida en el estudio de simulación. En la Figuras 3.18 y 3.19 se muestra la evolución de los valores óptimos W_{opt} con respecto a la fracción de solapamiento. Los valores SR están todos dentro de un rango razonable y por tanto se han omitido.

De las Figuras 3.16, 3.17, 3.18 y 3.19 se pueden hacer las siguientes observaciones:

1. Los resultados confirman un buen comportamiento por parte del estimador óptimo propuesto en comparación con el estimador estándar, y a su vez con respecto al estimador óptimo simple, es decir, el estimador propuesto óptimo basado en una única variable auxiliar.
2. Este estudio también nos muestra que se obtienen estimaciones más precisas cuando se usa un mayor número de variables auxiliares.
3. Cuando los tamaños muestrales en ambas ocasiones son iguales, la fracción de solapamiento óptima está entre 0,2 y 0,4. En otro caso, no puede observarse una fracción de solapamiento óptima.
4. Los valores W_{opt} son crecientes con respecto a la fracción de solapamiento. Este resultado era predecible puesto que a medida que aumenta el tamaño muestral de la parte solapada con respecto al tamaño de la muestra no solapada, el estimador de razón multivariante debería tener un mayor peso dentro del estimador propuesto. En todos los casos, los valores más altos de W_{opt} se obtienen cuando se usan todas las variables auxiliares

en la etapa de estimación. Este resultado demuestra que se obtienen estimaciones más precisas cuando se usan todas las variables auxiliares: de la expresión (3.41) puede observarse que W_{opt} es mayor si $V(\hat{Q}_{ymopt}(\beta))$ tiene valores más pequeños, y bajo esta situación, el estimador óptimo propuesto obtiene estimaciones más precisas.

5. Cuando el tamaño muestral en la segunda ocasión es menor que el tamaño en la primera ocasión, se obtiene una mayor ganancia en precisión, y esta ganancia aumenta a medida que crece la diferencia entre los tamaños muestrales. Este resultado es razonable porque si n es pequeño en relación con n' , entonces, la primera muestra proporcionará mayor información, y el estimador de razón múltiple basado en la muestra solapada presentará también un menor grado de error.

3.4. Estimadores bajo el método de verosimilitud empírica

En este apartado se utiliza el método de verosimilitud empírica para la estimación de cuantiles. Para ello, usaremos el estimador de verosimilitud empírica para la función de distribución definido en la Sección 2.4.3. Tomando la inversa de este estimador, podremos obtener estimadores de cuantiles fácilmente. Estos estimadores también se utilizarán para el análisis de algunas medidas de pobreza.

Bajo datos de la Encuesta Continua de Presupuestos Familiares para el primer trimestre del año 1997, mostraremos como tanto el estimador propuesto para los cuantiles como el método bootstrap para la estimación de la varianza, exhiben un buen comportamiento en comparación con otros estimadores alternativos.

3.4.1. Antecedentes

Asumiendo el método de verosimilitud empírica, los únicos estimadores conocidos para cuantiles en la literatura se basan en la aproximación modelo-calibrada, es decir, se usan los estimadores modelo-calibrados para la función de distribución descritos en la Sección 2.4.2. Sea $\hat{F}_{MCPE}(t)$ uno de estos estimadores cuando se usa el punto $t_0 = \hat{Q}_{HKy}(\beta)$. Notamos que $\hat{F}_{MCPE}(t)$ será más eficiente que $\hat{F}_{HKy}(t)$ para t en las cercanías de $Q_y(\beta)$.

El cuantil $Q_y(\beta)$ puede estimarse mediante inversión directa de $\widehat{F}_{MCPE}(t)$, esto es, $\widehat{Q}_{MCPE}(\beta) = \widehat{F}_{MCPE}^{-1}(\beta)$ para $\beta \in (0, 1)$. Puesto que $\widehat{F}_{MCPE}(t)$ es una verdadera función de distribución, esta inversión es computacionalmente simple.

Notamos que tanto este estimador como su correspondiente varianza asintótica serán usadas en la Sección 3.4.5 para su comparación empírica con el estimador propuesto bajo el método de verosimilitud empírica. Por esta razón, a continuación se resumen las principales propiedades asintóticas de este estimador. Para ello, asumimos que hay una secuencia de poblaciones finitas $\{U_\nu, \nu = 1, 2, \dots\}$. $F_\nu(t)$ y $Q_\nu(\beta)$ denotan respectivamente $F_y(t)$ y $Q_y(\beta)$, para la población U_ν . Además, sean los diseños muestrales siguientes

- (i) Muestreo aleatorio simple con o sin reemplazamiento.
- (ii) Muestreo estratificado aleatorio simple con o sin reemplazamiento.
- (iii) Muestreo con probabilidades desiguales de una etapa con reemplazamiento.
- (iv) Muestreo de varias etapas con reemplazamiento en la primera etapa.

Notamos que en el caso de diseños con reemplazamiento se usa el estimador de tipo Hansen-Hurwitz, esto es $\pi_i = nq_i$, donde q_i es la probabilidad de seleccionar la i -ésima unidad.

Una representación Bahadur para el cuantil $\widehat{Q}_{MCPE}(\beta)$ puede establecerse para estos diseños muestrales. Sean también las condiciones (C2.20), (C2.21) y (C2.22) dadas en la Sección 2.4.2 junto a las siguientes:

- (C3.7). Existe una función de distribución $F(t)$ doble diferenciable con función de densidad $f(t)$, tal que $F_\nu(t) - F(t) = o(1)$, y para cualquier $a_\nu = O(n^{-1/2})$

$$\sup_{|\delta| \leq a_\nu} |[F_\nu(t + \delta) - F_\nu(t)] - [F(t + \delta) - F(t)]| = o(n_\nu^{-1/2}),$$

donde el tamaño muestral $n_\nu \rightarrow \infty$ cuando $\nu \rightarrow \infty$.

- (C3.8). Para un fijado $\beta \in (0, 1)$, $Q_\nu(\beta) \rightarrow Q_0(\beta)$, donde $Q_0(\beta)$ es el cuantil β de $F(t)$ y $f(Q_0(\beta)) > 0$.

El siguiente teorema puede establecerse.

Teorema 3.11 *Bajo los diseños muestrales (i)~(iv) y las condiciones (C2.20), (C2.21), (C2.22), (C3.7) y (C3.8), se tiene que*

$$\widehat{Q}_{MCPE}(\beta) - Q_y(\beta) = \frac{1}{f(Q_y(\beta))} \left(\beta - \widehat{F}_{MCPE}(Q_y(\beta)) \right) + o_p(n^{-1/2}),$$

donde $f(\cdot)$ es la función densidad de la función de distribución límite de $F_y(t)$ cuando $N \rightarrow \infty$.

En consecuencia, la varianza asintótica de $\widehat{Q}_{MCPE}(\beta)$ puede aproximarse por

$$\begin{aligned} V(\widehat{Q}_{MCPE}(\beta)) &\simeq \frac{1}{f(Q_y(\beta))^2} V(\widehat{F}_{MCPE}(Q_y(\beta))) = \\ &= \frac{1}{f(Q_y(\beta))^2} \frac{1}{N^2} \sum_{i < j} \sum_{j=1}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{U_i}{\pi_i} - \frac{U_j}{\pi_j} \right)^2 + o(n^{-1}), \end{aligned}$$

donde $U_i = \delta(Q_y(\beta) - y_i) - F_y(Q_y(\beta)) - (w_i^* - \bar{w}^*)B_N$ y $\bar{w}^* = N^{-1} \sum_{i=1}^N w_i^*$. w_i^* viene dada por (2.85), (2.87), (2.90) o (2.93) cuando $t_0 = Q_y(\beta)$.

Esta varianza puede estimarse mediante

$$\begin{aligned} \widehat{V}(\widehat{Q}_{MCPE}(\beta)) &\simeq \frac{1}{f(Q_y(\beta))^2} V(\widehat{F}_{MCPE}(\widehat{Q}_{MCPE}(\beta))) = \\ &= \frac{1}{f(Q_y(\beta))^2} \frac{1}{N^2} \sum_{i < j} \sum_{j=1}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{u_i}{\pi_i} - \frac{u_j}{\pi_j} \right)^2 + o(n^{-1}), \end{aligned}$$

donde $u_i = \delta(\widehat{Q}_{MCPE}(\beta) - y_i) - \beta - (w_i - \bar{w})B_N$ y $\bar{w} = N^{-1} \sum_{i=1}^N w_i$. w_i viene dada por (2.86), (2.88), (2.91) o (2.92) cuando $t_0 = \widehat{Q}_{HKy}(\beta)$. $f(Q_y(\beta))$ puede estimarse mediante procedimientos estándares (Silverman, 1986).

La ganancia en eficiencia al usar $\widehat{Q}_{MCPE}(\beta)$ sobre $\widehat{Q}_{HKy}(\beta)$ es comparable a la ganancia de $\widehat{F}_{MCPE}(t)$ sobre $\widehat{F}_{HKy}(t)$. Con la óptima elección $w_i = E_\xi(z_i | \mathbf{x}_i)$, la ganancia máxima de la eficiencia asintótica está garantizada. Así, este método puede aplicarse en diseños muestrales complejos y para un vector multivariante de variables auxiliares.

3.4.2. Aplicación a la estimación de líneas de pobreza

El análisis de las líneas de pobreza es un tema reciente y de gran interés en la sociedad. La proporción oficial de pobreza y el número de personas en

pobreza son importantes medidas para el bienestar económico de un país.

El análisis de la estructura de los ingresos y la desigualdad de ingresos son los principales objetivos en los estudios de pobreza. Esto se debe a que la desigualdad de los ingresos puede afectar a la eficiencia del mercado laboral, y a que esto conlleva a una serie de problemas relacionados con la igualdad social, tal como la incidencia de la pobreza o la estratificación social.

La aplicación de una medida de pobreza requiere la especificación de una línea de pobreza, la cual separe a la población en pobres y no pobres. En la literatura, existen distintas formas de especificar una línea de pobreza. Por ejemplo, La Organización para la Cooperación Económica y el Desarrollo (OECD, acrónimo de *Organization for Economic Cooperation and Development*) en el año 1997, definió la línea de bajos ingresos como dos tercios del salario mediano, de modo que un empleado se consideraba que tenía ingresos bajos si recibía un salario inferior al anterior umbral señalado. Sin embargo, Eurostat (2000) define que un empleado en la Unión Europea percibe un salario bajo si su salario mensual es inferior al 60 % del salario mediano de su correspondiente país.

Los empleados con bajos ingresos, en particular, ha sido un centro de investigación con alto interés político (Lucifora y Salverda, 1998). Por un lado, a un nivel macroeconómico, los empleados con bajos ingresos es claramente relevante para la igualdad social, como lo demuestran las razones con alta pobreza en los países donde los empleados con bajos ingresos es relativamente alto (OECD, 1997). Por otro lado, desde una perspectiva microeconómica, existe una relación entre salarios bajos y estado de pobreza de los hogares (OECD, 1997, Eurostat, 2000).

En la literatura, existen tres tipos de métodos para determinar las líneas de pobreza: los métodos absolutos, relativos y los subjetivos. Los métodos absolutos obtienen la línea de pobreza como una cantidad mínima de fuentes en un punto del tiempo y ponen al día la línea solamente para cambios de precio sobre el tiempo. La línea de pobreza usada por el estadístico oficial de pobreza de Estados Unidos es un ejemplo de línea de pobreza absoluta. El método relativo especifica la línea de pobreza como un punto en la distribución de ingresos o gastos y, por lo tanto, la línea puede estar sin fecha automáticamente sobre el tiempo para cambios en niveles de vida. En la práctica, los investigadores a menudo especifican la línea de pobreza relativa como un porcentaje del ingreso o gasto medio (Wolfson y Evans, 1989, Johnson y Webb, 1992), como un porcentaje del ingreso o gasto mediano (Smeeding, 1991, Eurostat, 2000) o simplemente como un cuantil (OECD, 1982). El método subjetivo deriva

de la línea de pobreza basada en la opinión pública. Comparada con las dos primeras aproximaciones, el método subjetivo es relativamente menos popular y raramente se usa.

Mientras que las líneas de pobreza absolutas han sido usadas en la mayoría de los estadísticos de pobreza de los gobiernos, las líneas de pobreza relativas han ganado recientemente en popularidad y uso tanto en las comparaciones internacionales de pobreza como en análisis nacionales de pobreza a través del tiempo. Preston (1995) estableció las distribuciones muestrales de los estadísticos de pobreza relativos.

La desigualdad de salario es requerida a menudo en estudios de pobreza o distribución de la riqueza. Tradicionalmente, La oficina censal de Estados Unidos ha empleado un determinado número de percentiles límite y razones para estudiar cambios en la desigualdad de salarios de los hogares. Entre ellos encontramos la razón de ingresos para un determinado hogar entre el percentil 95 y el percentil 20, el percentil 95 con respecto a la mediana, etc. Derivadas de estos percentiles son también bastantes usados en la literatura de ingresos. Algunos investigadores han propuestos otras medidas alternativas como la razón entre los percentiles 90 y 10 o la razón entre los percentiles de orden 50 y 10. Eurostat (2000) también emplea el salario mediano con respecto al primer decil. Estos valores dan una idea de la extensión de las desigualdades entre salarios. Por ejemplo, la razón entre los percentiles de orden 50 y 10 nos permite ver si la incidencia de empleos con bajos ingresos está fuertemente relacionada con la dispersión de salarios en el tallo izquierdo de la distribución.

La atención dada a este tipo de estadísticos en los medios de comunicación y en los círculos de política es considerable, hasta el punto de que importantes decisiones políticas pueden verse influenciadas por estas medidas.

La característica común de estas medidas es su complejidad. Éstas son funciones no lineales de las observaciones y un alto número de éstas dependen de cuantiles. Como se ha comentado, la literatura relacionada a la estimación de medianas y otros cuantiles, los cuales usan una variable auxiliar, es considerablemente menos extenso que en el caso de medias y totales, y las técnicas habituales, tal como el método de regresión, no tienen una extensión obvia a la estimación de cuantiles. Por tanto, la mayoría de los estudios relacionados con cuantiles han sido desarrollados asumiendo muestreo aleatorio simple o muestreo estratificado (Gross, 1980, Sedransk y Meyer, 1978, Sedransk y Smith, 1988, Kuk y Mak, 1989, Singh *et al.*, 2001), o bien considerando aproximaciones basadas en el modelo (Chambers y Dunstan, 1986, Dorfman y Hall, 1993, Mak y Kuk, 1993), las cuales asumen un modelo de superpoblación, los

estimadores son dependientes de dicho modelos y puede llegarse a obtener un pobre cumplimiento de los estimadores bajo una inapropiada especificación del modelo. En la práctica, estas situaciones no son usuales, especialmente para el caso de datos relacionados con ingresos o gastos, los cuales se analizan asumiendo diseños muestrales complejos con probabilidades desiguales y cuyos datos, además, exhiben una alta asimetría, lo que hace muy difícil asociar un modelo de superpoblación a los datos en estudio. El uso de estimadores de cuantiles eficientes basados en información auxiliar y aproximaciones independientes del modelo, puede ayudarnos a obtener una mejoría en la estimación de medidas de pobreza. Notamos que la mayoría de los estudios relacionados con medidas de pobreza han sido llevados a cabo usando estimadores clásicos de la literatura del muestreo en poblaciones finitas.

El propósito de esta sección es desarrollar un estimador de cuantiles que pueda aplicarse a diferentes medidas de pobreza. Para ello, usaremos la aproximación modelo-asistida. Esta aproximación se usa habitualmente en las encuestas por muestreo y tiene una buena aceptación. Ejemplos de estimadores modelo-asistidos son el estimador de regresión generalizado (*GREG*) (Cassel *et al.*, 1976, Särndal, 1980) para el caso de estimar medias y totales, y estimadores de tipo razón y diferencia (Rao *et al.*, 1990) para el caso de estimar funciones de distribución y cuantiles. La principal ventaja de la aproximación modelo-asistida es que proporciona inferencias válidas bajo el modelo asumido y al mismo tiempo está protegido contra una inapropiada especificación del modelo en el sentido de proporcionar inferencias basadas en el diseño válidas, independientemente de los valores poblacionales para la variable de estudio.

En la presencia de información auxiliar, existen varios procedimientos de estimación para obtener estimadores más eficientes. Recientemente, se han propuestos los estimadores de calibración (Deville y Särndal, 1992) y los estimadores de verosimilitud empírica (Chen y Qin, 1993, Chen y Sitter, 1999). En las siguientes secciones se usará el método de verosimilitud empírica para construir nuevos estimadores para un determinado cuantil. En lo que respecta a la estimación de cuantiles usando el método de verosimilitud empírica (véase la Sección 3.4.1), Chen y Wu (2002) propusieron estimadores modelo-calibrados (Wu y Sitter, 2001a). Estos estimadores requieren el uso de un modelo de superpoblación apropiado, y son por tanto dependientes de dicho modelo. Además, estos estimadores se construyen por medio de restricciones que requieren el uso de un único valor fijado. Una importante pérdida de eficiencia puede llegar a obtenerse cuando dicho valor fijado se encuentra alejado del cuantil que va a ser estimado.

El estimador propuesto usa de modo efectivo la información auxiliar en la

etapa de estimación porque éste está basado en tres valores fijados contruidos a partir de la información auxiliar. Estos valores se encuentran bien repartidos dentro de la distribución de datos, resolviendo de este modo la perdida de eficiencia provocada por la elección de un valor fijado situado a gran distancia de cuantil que se va a estimar. Este estimador propuesto está basado en el estimador para la función de distribución descrito en la Sección 2.4.3.

Debido a la naturaleza específica de los cuantiles y a la complejidad de algunas medidas de pobreza, las varianzas de esto estadísticos complejos no pueden expresarse por simples formulas. Mostraremos como la técnica bootstrap es una posible alternativa en la estimación de la varianza del estimador propuesto.

3.4.3. Estimadores propuestos modelo-asistidos

En este epígrafe se describe el estimador propuesto usando la metodología de verosimilitud empírica. Como se ha comentado, usaremos una perspectiva modelo-asistida debido a que esta proporciona un enfoque en el cual se pueden desarrollar estimadores eficientemente. Para ello, necesitaremos un modelo de superpoblación que describa la relación entre la variable de interés y las variables auxiliares. Este modelo será posteriormente usado para construir estimadores basados en el diseño.

Como resulta habitual, consideraremos el modelo regresión lineal dado por

$$y_i = \beta^t \mathbf{x}_i + v_i \varepsilon_i, \quad i = 1, \dots, N \quad (3.49)$$

donde v_i es una función conocida de x_i y las cantidades ε_i son variables aleatorias independientes e idénticamente distribuidas con media 0 y varianza σ^2 . Notamos que en la práctica los valores del vector β son desconocidos, aunque es sabido que este parámetro puede estimarse eficientemente por mínimos cuadrados (véase por ejemplo Särndal *et al.*, 1992) como

$$B = \left(\sum_{i \in U} \frac{\mathbf{x}_i \mathbf{x}_i^t}{\sigma^2} \right)^{-1} \cdot \sum_{i \in U} \frac{\mathbf{x}_i y_i}{\sigma^2}. \quad (3.50)$$

Este estimador es óptimo en el sentido de ser el mejor estimador lineal e insesgado para β bajo el modelo (3.49). A su vez, B es una característica poblacional finita, aunque puede estimarse usando los datos muestrales. Esta estimación viene dada por

$$\hat{\beta} = \left(\sum_{i \in s} \frac{d_i \mathbf{x}_i \mathbf{x}_i^t}{\sigma^2} \right)^{-1} \cdot \sum_{i \in s} \frac{d_i \mathbf{x}_i y_i}{\sigma^2}. \quad (3.51)$$

Como ya sabemos, el método de verosimilitud empírica presenta buenas propiedades asintóticas y empíricas para el problema de la estimación de medias o totales (Chen y Qin, 1993, Chen y Sitter, 1999), funciones de distribución (Chen y Wu, 2002), estimación en presencia de datos faltantes (Rueda *et al.*, 2006b, Leung y Qin, 2006), etc. Chen y Wu (2002) propusieron estimadores de verosimilitud empírica modelo-calibrados que requieren el uso de un único valor prefijado. La aplicación de estos estimadores a la estimación de cuantiles resulta posible, aunque este proceso arrastra una importante pérdida de eficiencia cuando dicho valor prefijado está alejado de cuantil que va a ser estimado. Con el propósito de reducir esta pérdida en eficiencia, se proponen estimadores modelo-asistidos para cuantiles usando el método de verosimilitud empírica y tres valores prefijados que ayudarán a reducir tal pérdida de eficiencia.

Asumiendo el método de verosimilitud empírica (Chen y Sitter, 1999), el estimador propuesto para el cuantil β está dado por

$$\widehat{Q}_{MA}(\beta) = \inf\{t : \widehat{F}_{MA}(t) \geq \beta\}, \quad (3.52)$$

donde

$$\widehat{F}_{MA}(t) = \sum_{i \in s} \widehat{p}_i \delta(t - y_i), \quad (3.53)$$

y las cantidades \widehat{p}_i son las soluciones al problema de maximización de la función de verosimilitud pseudo empírica $\widehat{l}(\mathbf{p}) = \sum_{i \in s} d_i \log(p_i)$ sujeta a

$$\sum_{i \in s} p_i = 1, \quad (p_i > 0), \quad (3.54)$$

$$\sum_{i \in s} p_i \delta(t_{g25} - g_i) = \frac{1}{N} \sum_{k=1}^N \delta(t_{g25} - g_k) = F_g(t_{g25}) = 0,25, \quad (3.55)$$

$$\sum_{i \in s} p_i \delta(t_{g50} - g_i) = \frac{1}{N} \sum_{k=1}^N \delta(t_{g50} - g_k) = F_g(t_{g50}) = 0,5, \quad (3.56)$$

$$\sum_{i \in s} p_i \delta(t_{g75} - g_i) = \frac{1}{N} \sum_{k=1}^N \delta(t_{g75} - g_k) = F_g(t_{g75}) = 0,75, \quad (3.57)$$

donde $t_{g25} = Q_g(0,25)$, $t_{g50} = Q_g(0,50)$, $t_{g75} = Q_g(0,75)$, y $Q_g(\alpha)$ es el cuantil α para la variable $g_i = \widehat{\beta}^t \mathbf{x}_i$.

Notamos que la idea de usar $\delta(t - g_i)$ para cualquier t como una variable de calibración para formar restricciones como las dadas en (3.55), (3.56) y (3.57) fue en primer lugar discutida en Wu y Sitter (2001a) y posteriormente elaborada en Chen y Wu (2002).

La elección de los valores t_{g25} , t_{g50} y t_{g75} en (3.55), (3.56) y (3.57) está motivada por varias razones. En primer lugar esta elección está altamente relacionada con la existencia de solución en el método de verosimilitud empírica, puesto que tomando valores más próximos entre ellos podría dirigir a problemas de multicolinealidad. Además, si se toman valores extremos, la variable indicadora podría tomar únicamente ceros o bien valores todos iguales a uno, en cuyo caso el método de verosimilitud empírica no tendría solución. Tomando estas consideraciones en cuenta, resulta razonable asumir que la elección t_{g25} , t_{g50} y t_{g75} es apropiada. En cualquier caso, notamos que este problema de multicolinealidad decrece a medida que aumenta el tamaño muestral.

Por otro lado, la elección de los valores t_{g25} , t_{g50} y t_{g75} se ha determinado por razones de eficiencia. Así, si se usara un único punto t_0 , $\widehat{F}_{MA}(t)$ sería menos eficiente a medida que t se aleje de t_0 . Asumiendo varios valores como en nuestro caso, resulta razonable asumir que si estos puntos están perfectamente distribuidos dentro del posible rango de valores de t , entonces $\widehat{F}_{MA}(t)$ será más eficiente. Los valores t_{g25} , t_{g50} y t_{g75} verifican tal propiedad y $\widehat{F}_{MA}(t)$ será más preciso en t en las cercanías de los cuantiles poblacionales de la variable g . Esto concierne a un rango amplio de valores para la variable de interés.

Una vez que se ha definido el estimador de cuantiles, las medidas de pobreza que dependan de tales parámetros podrán ser estimadas. Por ejemplo, la línea de bajos ingresos puede definirse como la fracción α de un cuantil β (Eurostat, 2000, Blackburn, 1990, 1994, Smeeding, 1991, etc.):

$$L_{\alpha,\beta} = \alpha Q_y(\beta), \quad (3.58)$$

y las medidas para cuantificar la desigualdad de gastos están dadas por la razón entre los cuantiles de órdenes β_1 y β_2 (Eurostat, 2000, U.S. Census Bureau, etc):

$$r_{\beta_1,\beta_2} = Q_y(\beta_1)/Q_y(\beta_2). \quad (3.59)$$

Estas medidas pueden estimarse fácilmente por

$$\widehat{L}_{\alpha,\beta} = \alpha \widehat{Q}_y(\beta), \quad (3.60)$$

para la medida dada en (3.58), y por

$$\widehat{r}_{\beta_1,\beta_2} = \widehat{Q}_y(\beta_1)/\widehat{Q}_y(\beta_2), \quad (3.61)$$

para la medida dada en (3.59).

3.4.4. Propiedades. Estimación de la varianza

El estudio de las propiedades asintóticas del estimador propuesto pasa por analizar tales propiedades para el estimador $\hat{F}_{MA}(t)$, las cuales se han establecido en la Sección 2.4.4. Queda por tanto describir una expresión para la varianza del estimador propuesto para cuantiles. La determinación de tal expresión es posible, aunque tendría únicamente validez asintótica, es decir, para tamaños muestrales bastantes elevados, situación no siempre presente en la práctica. Además, por la estructura no lineal del cuantil, se requiere el uso de una aproximación lineal que emplea parámetros poblacionales, por ejemplo densidades, que también tendrían que ser estimados, lo que conlleva a otra pérdida de eficiencia en la etapa de estimación de la varianza.

Si aplicamos el estimador propuesto a la estimación de medidas de pobreza, la determinación de dicha expresión asintótica resulta aún más difícil, puesto que la característica común de las medidas de pobreza, como por ejemplo (3.58) y (3.59), es su complejidad. Además, los datos de ingresos y gastos provienen usualmente de encuestas complejas (muestreos con probabilidades desiguales de tipo estratificado, con múltiple etapas, por conglomerados, etc), lo que también dificulta la determinación de expresiones asintóticas bajo estas situaciones. La única alternativa en estos casos es el uso de métodos especiales para la estimación de varianzas.

Por estas razones, proponemos el uso de técnicas alternativas para la estimación de la varianza del estimador propuesto. En concreto, se propone la técnica bootstrap que frecuentemente se usa en la estimación de cuantiles, y en particular, para la estimación de las medidas de pobreza. Este hecho queda justificado por los estudios ya llevados a cabo y los cuales resumiremos brevemente a continuación. Puesto que el estudio empírico que llevamos a cabo está basado en algunas medidas de pobreza, centraremos nuestra atención a la estimación de la varianza de medidas de pobreza.

En primer lugar, notamos que en los estudios de pobreza, la variabilidad muestral de las diferentes medidas estimadas presentan un interés particular cuando éstas son comparadas entre países, a través del tiempo o entre subgrupos dentro de un país.

Los métodos tradicionales para aproximar la varianza de un estimador (véase Wolter, 1985), envuelven una de las siguientes estrategias: linealización de Taylor o métodos de replicación tal como bootstrap, jackknife, etc. En los casos donde los estimadores presentan una forma compleja (como en el caso de cuantiles), los métodos de replicación son preferidos por ser más fáciles de

implementar, aunque para el caso de cuantiles, el clásico método jackknife da estimadores inconsistentes para la varianza (Kovar *et al.*, 1988, Shao y Wu, 1989). También pueden usarse para la estimación de la varianza otros métodos alternativos tal como linealización y técnicas residuales (Deville, 1999). Una complicación al aplicar el método de linealización en la estimación de cuantiles es que éste requiere la estimación de funciones de densidad de probabilidad para la variable de interés.

Los métodos bootstraps están ganando en popularidad en las investigaciones empíricas. Por ejemplo, en el Instituto Estadístico de Canadá se llevó a cabo un estudio de simulación para comparar la eficiencia de varios métodos de remuestreo con respecto al método de estimación de ecuaciones (véase Kovacevic, Yung y Pandher, 1995) en el caso de medidas de desigualdad de ingresos. Para algunos cuantiles, el estimador bootstrap exhibía el menor sesgo relativo, mientras que el método de estimación de ecuaciones junto con el método bootstrap eran los óptimos en el sentido de estabilidad. Estos resultados confirman la ventaja al usar el método bootstrap sobre el resto de aproximaciones. Asumiendo también medidas de pobreza, Shao y Chen (1998) también demostraron la consistencia del método bootstrap para la estimación de la varianza.

3.4.5. Propiedades empíricas

En esta sección se evalúa la precisión del estimador propuesto junto con otros estimadores conocidos. Además, se estudia la eficiencia de estos procedimientos cuando se aplica la estimación de cuantiles a diversas medidas de pobreza. El comportamiento del método bootstrap para la estimación de varianzas será también analizado. Para ello, se calculan las estimaciones bootstrap para los distintos estimadores y comparamos estos resultados con los obtenidos a través de las correspondientes expresiones para la varianza de cada estimador, en aquellos casos que se disponga de tales expresiones. Por simplicidad, se asume muestreo aleatorio simple.

En este estudio se usa la población ECPF1997 (véase Apéndice A) que está formada por los datos de ingresos y gastos de 3000 familias extraídas de la Encuesta Continua de Presupuestos Familiares del año 1997. Estos datos se han duplicado tres veces para crear una población artificial de $N = 9000$ individuos, a partir de los cuales nos basaremos para llevar a cabo el presente estudio de simulación. Como variable principal se han tomado los ingresos, mientras que como variable auxiliar se consideran los gastos familiares.

El cumplimiento del estimador de cuantiles propuesto y su correspondiente estimación de la varianza obtenida mediante bootstrap se comparará con los estimadores de cuantiles obtenidos a partir de las siguientes funciones de distribución: el clásico estimador de tipo Horvitz-Thompson, $\widehat{F}_{HTy}(t)$, el cual lo usaremos como estimador de comparación para todos los estimadores, los estimadores de tipo razón y diferencia ($\widehat{F}_r(t)$, $\widehat{F}_d(t)$, $\widehat{F}_{dm}(t)$) propuestos en Rao *et al.* (1990), el estimador de Chambers y Dunstan (1986), $\widehat{F}_{CD}(t)$, y $\widehat{F}_{MC}(t)$, el estimador propuesto en Chen y Wu (2002). Además, calcularemos el estimador modelo-asistido asumiendo un único valor prefijado. Esto nos permitirá conocer la ganancia en precisión al usar más de un valor prefijado.

Dado un cuantil de orden β , el comportamiento de todos los estimadores de cuantiles y sus varianzas están medidos por medio del Sesgo Relativo, (SR) y Eficiencia Relativa (ER). Así, para un determinado cuantil, $\widehat{Q}_y(\beta)$, calcularemos

$$\begin{aligned} ER[\widehat{Q}_y(\beta)] &= ECM[\widehat{Q}_y(\beta)]/ECM[\widehat{Q}_{HTy}(\beta)], \\ SR[\widehat{Q}_y(\beta)] &= 100 \times \left(E[\widehat{Q}_y(\beta)] - Q_y(\beta) \right) / Q_y(\beta), \end{aligned} \quad (3.62)$$

y para un estimador de la varianza, $\widehat{V}(\widehat{Q}_y(\beta))$, se obtendrá las medidas dadas por (3.62) después de sustituir $\widehat{Q}_y(\beta)$ y $Q_y(\beta)$ por $\widehat{V}(\widehat{Q}_y(\beta))$ y $V[Q_y(\beta)]$ respectivamente. $E[\cdot]$, $ECM[\cdot]$ y $V[\cdot]$ son las Esperanzas Empíricas, Error Cuadrático Medio y Varianzas basadas en 500 muestras. Notamos que valores de $ER[\widehat{Q}_y(\beta)]$ y $ER[\widehat{V}(\widehat{Q}_y(\beta))]$ menores de 1 indican que $\widehat{Q}_y(\beta)$ y $\widehat{V}(\widehat{Q}_y(\beta))$ son más precisos que $\widehat{Q}_{HTy}(\beta)$ y $\widehat{V}(\widehat{Q}_{HTy}(\beta))$, respectivamente. Asumiendo normalidad, también se ha obtenido la Cobertura de los Intervalos de Confianza (CI) al 95 % y la Longitud Media de cada Intervalo (LI). Todos los estudios se han basado en muestras de tamaño $n = 500$.

Notamos que la precisión de cada estimador depende directamente del cuantil que va a ser estimado. Por ejemplo, el estimador de Chambers y Dunstan es muy eficiente en la estimación de la mediana, aunque generalmente sufre de importantes sesgos en las estimaciones a medida que se estiman cuantiles más alejados de la mediana (véase Rao *et al.*, 1990, Chambers *et al.*, 1993, y Dorfman, 1993). Por este motivo, el primer estudio desarrollado intenta medir la precisión media global de cada estimador a partir de los resultados obtenidos en las estimaciones de los cuantiles de órdenes $\beta = 0,1, 0,3, 0,5, 0,7, 0,9$. Las medidas usadas para realizar tal medición son el Sesgo Absoluto Medio (SAM), dado por

$$SAM = \frac{1}{5} \sum_{i=1}^5 |SR[\widehat{Q}_y(\beta_i)]|,$$

Tabla 3.12: Medidas globales medias de precisión y eficiencia basadas en cuantiles de órdenes $\beta = 0,1, 0,3, 0,5, 0,7, 0,9$, y muestras de tamaño $n = 500$.

Est.	Varianzas bootstrap						Varianzas asintóticas			
	ERM	SAM	ERM	SAM	CIM	LIM	ERM	SAM	CIM	LIM
MA	0.86	0.25	0.82	14.05	92.9	550.96	–	–	–	–
MA1	0.89	0.23	0.83	12.65	93.2	561.62	–	–	–	–
MC	0.92	0.25	0.86	8.72	92.9	563.18	0.78	7.16	93.9	553.87
HK	1.00	0.26	1.00	9.97	92.8	622.32	1.00	9.52	94.0	616.53
r	1.04	0.23	1.08	9.87	93.3	654.58	1.01	3.96	93.2	646.85
d	1.05	0.25	1.06	7.32	92.9	651.83	1.02	3.67	93.3	650.31
dm	0.87	0.21	0.81	12.17	92.7	556.01	0.70	5.27	93.9	548.07
CD	3.58	12.44	0.48	10.24	17.1	436.84	–	–	–	–

la raíz cuadrada del valor medio de las medidas ER , es decir,

$$ERM = \sqrt{\frac{1}{5} \sum_{i=1}^5 ER[\widehat{Q}_y(\beta_i)]},$$

y por último, los valores medios para las medidas CI y LI . Dichas medidas se denotarán como CIM y LIM respectivamente. En la Tabla 3.12 puede observarse las distintas medidas globales para todos los estimadores. A partir de la eficiencia relativa media, podemos comprobar que el estimador propuesto presenta el mejor comportamiento, seguido del estimador de diferencia óptimo (dm). El estimador de Chambers y Dunstan es el menos eficiente, mientras que los estimadores de tipo razón y diferencia también funcionan peor que el estimador estándar. En el estudio de las varianzas observamos que las expresiones asintóticas funcionan ligeramente mejor que la técnica bootstrap, por lo que a tenor de los resultados sería aceptable recurrir a tal procedimiento para la estimación de la varianza. Por último, al estimar todas las varianzas de los estimadores mediante bootstrap, se observa que el estimador propuesto presenta el mejor comportamiento, al estimar los intervalos de confianza con menor longitud y una cobertura similar al resto de estimadores.

El siguiente paso en esta sección es el análisis de la eficiencia del estimador propuesto cuando se aplica a la estimación de medidas de pobreza. En primer lugar analizamos los resultados obtenidos para la estimación de las líneas de bajos ingresos (Tabla 3.13) y a continuación describiremos las conclusiones

Tabla 3.13: Medidas de precisión y eficiencia para la línea de bajos ingresos cuando $\alpha = 0,6$, $\beta = 0,5$ y se toman muestras de tamaño $n = 500$.

Est.	Varianzas bootstrap						Varianzas asintóticas			
	<i>ER</i>	<i>SR</i>	<i>ER</i>	<i>SR</i>	<i>CI</i>	<i>LI</i>	<i>ER</i>	<i>SR</i>	<i>CI</i>	<i>LI</i>
MA	0.70	-0.10	0.57	16.59	93.8	391.54	–	–	–	–
MA1	0.79	-0.08	0.63	13.03	94.2	410.32	–	–	–	–
MC	0.78	-0.11	0.65	14.87	94.0	412.62	0.53	15.81	94.8	423.94
HK	1.00	-0.24	1.00	17.09	93.4	470.88	1.00	18.41	94.2	482.73
r	1.09	-0.00	0.98	7.77	94.6	473.71	0.81	6.97	93.8	481.26
d	1.11	0.01	0.97	6.40	93.8	474.52	0.87	7.45	93.8	486.03
dm	0.74	-0.07	0.49	7.39	93.6	388.18	0.37	8.17	94.8	398.41
CD	1.11	2.23	0.09	0.65	77.2	313.01	–	–	–	–

más importantes en la estimación de razones entre cuantiles para el análisis de la desigualdad entre ingresos (Tablas 3.14 y 3.15).

En primer lugar, notamos que al tratarse de medidas relativas, los resultados obtenidos para las líneas de bajos ingresos en la Tabla 3.13 serán los mismos si se usaran otros valores de α , o bien si se considera la propia mediana. Por tanto, las conclusiones que puedan extraerse de esta tabla se podrían hacer para estos casos comentados.

En la Tabla 3.13 observamos que el estimador propuesto es el más eficiente en términos de eficiencia relativa. Todos los sesgos relativos se encuentran dentro de un rango razonable, excepto el de Chambers y Dunstan con un valor superior al resto, en torno al 2.23%. Un aspecto importante a tener en cuenta en la estimación de la varianza es que las estimaciones bootstrap son, en términos generales, más precisas que las obtenidas mediante las expresiones asintóticas, puesto que se obtienen para cada estimador sesgos más reducidos, e intervalos de confianza menos amplios con idénticas coberturas. Este resultado nos confirma que la técnica bootstrap es un procedimiento óptimo en la estimación de la varianza de la mediana, y en particular, la estimación de la varianza de las líneas de bajos ingresos. Observando las estimaciones bootstrap podemos comprobar que el estimador diferencia óptimo y el estimador propuesto obtiene las mejores estimaciones para la varianza.

Las Tablas 3.14 y 3.15 nos dan las distintas medidas de precisión y eficiencia para medidas de pobreza dadas por razones de cuantiles. De nuevo, el estimador propuesto se muestra más eficiente en términos de eficiencia relativa.

Tabla 3.14: Medidas de precisión y eficiencia para la razón de cuantiles cuando $\beta_1 = 0,5$, $\beta_2 = 0,25$, y se toman muestras de tamaño $n = 500$.

Est.	ER	SR	Varianzas bootstrap			
			ER	SR	CI	LI
MA	0.93	0.05	0.92	18.18	93.6	0.18
MA1	1.04	0.14	1.07	17.75	95.2	0.19
MC	1.00	-0.01	1.01	14.68	93.8	0.19
HK	1.00	0.05	1.00	15.91	95.2	0.19
r	1.62	0.34	2.53	14.78	94.4	0.24
d	1.65	0.29	2.16	11.45	94.2	0.23
dm	0.90	0.06	0.80	15.69	93.8	0.18
CD	21.07	14.10	0.05	23.43	0.0	0.08

Tabla 3.15: Medidas de precisión y eficiencia para la razón de cuantiles cuando $\beta_1 = 0,95$, $\beta_2 = 0,2$, y se toman muestras de tamaño $n = 500$.

Est.	ER	SR	Varianzas bootstrap			
			ER	SR	CI	LI
MA	0.93	0.56	1.01	-0.70	91.4	0.92
MA1	14.66	1.70	–	-82.28	91.4	1.06
MC	1.02	0.61	1.07	-3.21	91.6	0.96
HK	1.00	0.27	1.00	-3.04	91.4	0.95
r	1.40	0.95	2.15	0.30	92.6	1.14
d	1.38	0.72	2.01	-3.69	91.4	1.11
dm	1.03	0.61	1.12	-6.12	90.8	0.95
CD	46.52	43.58	–	–	2.4	1.33

Conclusiones similares pueden derivarse de los resultados obtenidos en la etapa de la estimación de la varianza mediante bootstrap. El estimador de Chambers y Dunstan ofrece el peor comportamiento con importantes sobreestimaciones en la estimación de las razones. Esto se debe a que se están estimando cuantiles alejados de la mediana.

Capítulo 4

Discusión

En este capítulo se hace una discusión conjunta de los resultados obtenidos en todos los capítulos anteriores, resumiendo las principales conclusiones y ofreciendo las perspectivas y futuras líneas de investigación que se derivan de estos resultados.

4.1. Conclusiones y valoración de resultados

El presente trabajo se divide en dos grandes bloques: estimación bajo el método de verosimilitud empírica (Capítulo 2) y la estimación de cuantiles (Capítulo 3). En estos dos capítulos se han planteado nuevos estimadores en situaciones reales del muestreo en poblaciones finitas.

Así, asumiendo el método de verosimilitud empírica se han propuesto estimadores en presencia de datos faltantes, situación muy usual en la práctica y que no se tiene en cuenta en la mayoría de las investigaciones por muestreo. Las aportaciones hechas en este sentido dan una alternativa para la solución de este problema, puesto que se ha comprobado que puede existir una importante ganancia en eficiencia en las estimaciones de los parámetros desconocidos.

En concreto, se ha usado el método de verosimilitud empírica para estimar una media poblacional cuando en la encuesta nos encontramos con información faltante tanto en la variable de estudio como en la variable auxiliar. Se ha asumido que la muestra puede ser seleccionada mediante un diseño muestral arbitrario, con probabilidades iguales o desiguales.

El estimador propuesto se basa en una clase de estimadores formada por un estimador de verosimilitud empírica y por un estimador de tipo Hájek. Se han derivado las propiedades asintóticas de estos estimadores y el estimador óptimo dentro de la clase propuesta en el sentido de minimizar la varianza asintótica.

El estimador propuesto se ha comparado con otros estimadores en un estudio de simulación, donde se ha comprobado que el estimador óptimo presenta el mejor comportamiento con respecto a sus competidores cuando el número de valores perdidos es relativamente elevado y la relación lineal entre la variable principal y la auxiliar es débil.

Asumiendo el método de verosimilitud empírica también se han propuesto estimadores modelo-asistidos para la función de distribución. El estimador propuesto posee un importante número de propiedades deseables. Por ejemplo:

- Puede aplicarse fácilmente a diseños muestrales con probabilidades desiguales.
- No es dependiente de un modelo de superpoblación como le ocurre por ejemplo a los estimadores basados en modelo o a los estimadores modelo-calibrados.
- Se establecen las condiciones para la existencia del estimador.
- Bajo ciertas condiciones, el estimador es una verdadera función de distribución. Notamos que esta propiedad no se satisface para un gran número de estimadores en la literatura.
- Se satisfacen también otras propiedades importantes como la insesgadez asintótica, normalidad asintótica, disponibilidad de un estimador de la varianza, etc.

La precisión del estimador propuesto se ha comparado mediante varias medidas con otros estimadores conocidos. Estos estudios han mostrado un comportamiento óptimo por parte del estimador propuesto modelo-asistido. También se ha visto que el estimador de Chambers y Dunstan puede llegar a ser muy eficiente cuando el modelo en el que se basa es apropiado, aunque como se discutió en Rao *et al.* (1990), Chambers *et al.* (1993) y Dorfman (1993), este estimador cumple pobremente cuando se tiene una mala especificación del modelo. Un comentario similar puede hacerse sobre el estimador de verosimilitud empírica modelo-calibrado. Este estimador también sufre una importante

pérdida de eficiencia cuando se considera un valor fijado alejado del punto donde va a ser estimada la función de distribución.

Otra propiedad importante que caracteriza al estimador propuesto es el uso eficiente que se hace de la información auxiliar: por un lado porque pueden usarse múltiples variables auxiliares en la etapa de estimación, y por otro porque se usan un conjunto de valores prefijados que poseen una buena distribución y ayudan a mejorar la estimación de la función de distribución, especialmente en las proximidades de algunos de estos puntos. Recordamos también que el hecho de considerar \mathbf{t}_g y \mathbf{x} como valores fijados hacen que los pesos \hat{p}_i sean independientes de t y puedan establecerse mejores propiedades para el estimador propuesto.

En conclusión, el método de verosimilitud empírica modelo-asistido es una aproximación práctica y simple que incorpora fácilmente información auxiliar en la estimación de la función de distribución. Este estimador presenta un buen cumplimiento y puede ser una alternativa válida a otros estimadores de la función de distribución.

El estudio de la estimación de cuantiles se ha llevado a cabo en el Capítulo 3. Los aportes a la teoría de la estimación de cuantiles se han centrado en tres aspectos: estimación en muestreo con dos ocasiones sucesivas, estimación en muestreo bifásico y estimación usando el ya comentado método de verosimilitud empírica.

El muestreo en ocasiones sucesivas es una técnica muy conocida que puede usarse en encuestas continuas para estimar parámetros poblacionales y medidas de diferencia o cambio de una variable de interés. Las encuestas de tipo económico o social llevadas a cabo por la agencias nacionales y otros organismos estadísticos usan este diseño muestral, y la estimación de cuantiles es un problema común en la mayoría de estos estudios. Dentro del muestreo en dos ocasiones sucesivas se han planteado estimadores desde dos perspectivas bastantes usadas dentro del muestreo en poblaciones finitas: bajo diseños muestrales probabilísticos con probabilidades desiguales y asumiendo múltiples variables auxiliares.

Asumiendo diseños muestrales con probabilidades desiguales en cada ocasión se ha propuesto un estimador compuesto por un estimador de tipo razón (en la porción solapada por ambas muestras) y otro de tipo Hájek (en la parte no solapada de la muestra más reciente). El estimador propuesto es fácil de computar y se ha mostrado bastante preciso en los estudios de simulación.

Asumiendo muestreo aleatorio simple en cada una de las dos ocasiones,

se ha obtenido la normalidad asintótica del estimador, la cual nos sirve, por ejemplo, para construir intervalos de confianza para los cuantiles.

Por otro lado, asumiendo múltiples variables auxiliares y muestreo aleatorio simple en cada una de las dos ocasiones, se ha propuesto una clase de estimadores para cuantiles basados en un estimador de tipo razón multivariante y construido a partir de la información obtenida en la parte solapada. Bajo la clase propuesta se ha obtenido la expresión del estimador óptimo en el sentido de mínima varianza asintótica. El estimador propuesto posee un buen número de propiedades deseables, tal como normalidad asintótica, disponibilidad de la varianza del estimador, simplicidad de computación, etc. En los estudios empíricos y teóricos que se han llevado a cabo, el estimador se muestra más preciso que otros estimadores conocidos.

La mayoría de los procedimientos de muestreo que usan información auxiliar se basan en estimadores que requieren el uso de variables conocidas a nivel poblacional, siendo este hecho poco frecuente en la práctica. Una solución a este problema se presenta con la aplicación de un muestreo bifásico. Por tanto, el problema de la estimación de cuantiles basados en información auxiliar queda resuelto con los estimadores propuestos en este sentido. Con el fin de obtener unas estimaciones más precisas en poblaciones heterogéneas, con una posible distribución en grupos homogéneos, también se han propuesto estimadores para cuantiles en muestreo bifásico y usando un muestreo estratificado en la muestra de la primera fase.

Asumiendo muestreo bifásico bajo cualquier método de extracción de unidades en cada una de las dos fases, se han propuesto estimadores de tipo razón y exponencial. Se ha demostrado la insesgadez de estos estimadores y se han proporcionado expresiones para sus varianzas. Estos resultados nos han servido para poder obtener un estimador óptimo en el estimador de tipo exponencial. Bajo distintos esquemas de muestreo y varios estudios de simulación, se ha comprobado que los estimadores propuestos pueden obtener estimaciones más precisas que el resto de estimadores existentes en la literatura.

Los estimadores propuestos en muestreo bifásico, cuando se usa un muestreo estratificado en la primera fase, están basados en un estimador eficiente para la función de distribución. Se han establecido varias propiedades para este estimador de la función de distribución, por lo que el estimador propuesto para cuantiles posee mejores propiedades. Los resultados teóricos y empíricos que se han llevado a cabo han demostrado que el estimador propuesto puede proporcionar resultados óptimos en este esquema de muestreo.

Por último, se han propuesto estimadores para cuantiles desde una perspectiva modelo-asistida y considerando el método de verosimilitud empírica. La aplicación de estos estimadores a la estimación de algunas medidas de pobreza también ha sido analizada. Se ha propuesto usar la técnica bootstrap para la estimación de la varianza los estimadores propuestos. La precisión de todos estos procedimientos nuevos ha sido confirmada en estudios de simulación y para el problema de la estimación de cuantiles y medidas de pobreza usadas por numerosos organismos de estadística internacionales y de varios países.

4.2. Perspectivas y futuras líneas de investigación

La reciente implementación del método de verosimilitud empírica en la teoría del muestreo en poblaciones finitas y la compleja definición de un cuantil, hace que se tengan futuras líneas de investigación en los dos campos tratados a lo largo de este texto.

En lo que respecta al método de verosimilitud empírica, es un procedimiento que puede ser investigado en la mayoría de los campos de los que se compone el muestreo. En concreto, basándonos en la metodología propuesta para el tratamiento de datos faltantes en la estimación de la media poblacional, se puede aplicar dicha metodología al problema de la estimación de otros parámetros tal como razones, varianzas y cuantiles.

El método de verosimilitud empírica requiere el uso de técnicas de resolución de ecuaciones no lineales como bisección o Newton-Raphson. Se puede investigar el uso de técnicas nuevas más potentes en el método de verosimilitud empírica. Esto ayudará a una solución fiable y obtenida a partir de un menor número de iteraciones.

La precisión del método de verosimilitud empírica en diseños muestrales complejos (bifásico, en ocasiones sucesivas, por conglomerados, etc), así como en la estimación de parámetros en subpoblaciones o áreas pequeñas son temas de gran interés que pueden proporcionar soluciones aceptables a la teoría del muestreo.

Para la estimación de cuantiles se han propuesto estimadores en muestreo con dos ocasiones sucesivas y extracción de unidades con probabilidades desiguales. La extensión del estimador propuesto al caso de más de dos ocasiones sucesivas tiene una interesante aplicabilidad, puesto que es el caso de la ma-

yoría de las encuestas continuas. En estos estudios también se puede tener en cuenta el uso de una información auxiliar multivariante.

Los estimadores propuestos en muestreo con dos ocasiones sucesivas se han basado en estimadores de tipo razón, diferencia y exponenciales en la parte solapada. El estudio de otro tipo de estimadores, tal como el de diferencia combinado propuesto por Rao *et al.* (1990), puede proporcionar estimadores aún más precisos. Las investigaciones en este sentido serían bastante interesantes.

Existen un gran número de diseños muestrales complejos (por conglomerados, polietápicos, etc) donde la estimación de cuantiles no ha sido investigada. El uso de estos diseños en la práctica es bastante frecuente y el estudio de estimadores más precisos resulta también interesante y con una alta aplicabilidad.

El estudio de parámetros lineales en subpoblaciones ha sido extensamente estudiado en los últimos años. Sin embargo, este no es el caso del problema de la estimación de cuantiles. El uso obligado de los cuantiles en las distintas encuestas de tipo económico y social invita a que este campo sea investigado, y se propongan estimadores de cuantiles eficientes. Un procedimiento para llevar a cabo tal tarea es usar técnicas similares a las usadas en el problema de la estimación de medias o totales. El estudio consistirá por tanto en conocer bajo que condiciones estas técnicas pueden aplicarse a la estimación de cuantiles y en dicho caso, cuales de ellas será la que nos proporcione estimaciones más precisas.

Capítulo 5

Redacción para aspirar a la mención europea en el título de Doctor

5.1. Abstract

The present work deals two topics of survey sampling: the recent pseudo empirical likelihood method (Chapter 2) and the estimation of quantiles (Chapter 3).

The notation followed in this text and other basic definitions in survey sampling are described in Chapter 1. Another information such as the previous bibliography is also commented.

The pseudo empirical likelihood method is described in Chapter 2 under different settings. In Section 2.2 the most important properties and results are summarized.

The problem of missing information is a common feature in practice when a survey is carried out. In this case, a procedure to solve this problem is to eliminate to the individuals that have missing data for any variable. The main disadvantage from this technique is a large bias for the estimators. Another customary technique is the imputation. This presents sometimes the disadvantage of providing invalid inferences due to that the imputed values are considered like real values. In Section 2.3 a method is proposed for the estimation in presence of missing data. This does not need to eliminate to any

individual and all the sample information is used in the estimation stage. This procedure uses the pseudo empirical likelihood method. The asymptotic properties are established and empirical results obtained via simulation studies are also shown. These results give a gain in precision for the proposed estimator compared with other known estimators in presence of missing data.

The problem of estimating the distribution function is a important topic due to that this function allows us to analyze the most important characteristics of the population. Using this function, the problem of estimating quantiles and other parameters is also solved. Assuming the pseudo empirical likelihood method, in Section 2.4 a model-assisted estimator based on an effective use of auxiliary information is proposed for the distribution function. This estimator is not dependent model. Under certain conditions, the proposed estimators are themselves distribution functions. Empirical results confirm the good behaviour for the proposed estimator.

The problem of estimating quantiles under different settings is analyzed in Chapter 3. Thus, the estimation of quantiles under sampling on two occasions is investigated in Section 3.3. On the one hand, the use of an arbitrary sampling design is studied and on the other hand we also obtain more efficient estimators assuming multivariate auxiliary information. Both theoretical and empirical justifications have demonstrated that the proposed estimators can be more accuracy than other known estimators.

In Section 3.2 we deal the estimation of quantiles under two-phase sampling when the samples are selected under arbitrary sampling designs. Direct, ratio and exponential type estimators are proposed and important properties such as unbiasedness or variance estimation are also established. In particular, we also investigated the estimation of quantiles under two-phase sampling for stratification. We have observed that this sampling design provides important gain in efficiency due to the properties of stratified sampling. The proposed estimators are more accuracy than other known estimators in terms of relative bias and relative efficiency.

Finally, assuming the pseudo empirical likelihood method, model-assisted estimators for quantiles have been proposed. The application of these estimators to the estimation of several poverty measures have been also discussed. The bootstrap method for variance estimation has been considered and empirical studies have shown that the proposed estimators have a better precision than other estimators.

A brief overview and some conclusions are summarized in Chapter 4. We

also point out some future researches related to the topics discussed in this text.

The text finishes with a set of appendices. Appendix A describes the most important properties related to the populations used in the text. The notation and definition of the several sampling designs used are described in Appendix B. Finally, note that all the empirical studies have been carried out using the software *R*. The *R* codes are available in Appendix C.

5.2. Pseudo empirical likelihood method in the presence of missing data

In this section an estimator for the population mean when some observations on the study and auxiliary variables are missing from the sample is proposed. The proposed estimator is valid for any unequal probability sampling design, and is based upon the pseudo empirical likelihood method. The proposed estimator is compared with other estimators in a simulation study.

5.2.1. Introduction

It is common practice to use auxiliary population information at the estimation stage. This technique has many advantages. For example, suitable auxiliary information can produce a considerable reduction in bias and sampling error.

When one or more auxiliary variables correlated with the study variable are available, calibration (Huang and Fuller, 1978; Deville and Särndal, 1992) and pseudo empirical likelihood (Chen and Qin, 1993; Chen and Sitter, 1999; Wu and Sitter, 2001a; Wu, 2002) are two methods that can be used to estimate population total, population mean, distribution functions and quantiles. Both methods use auxiliary information on one or several auxiliary variables.

These techniques usually provide estimators that are more efficient than traditional estimators, such as the Horvitz and Thompson (1952) estimator and the Hájek ratio estimator (Rao, 1966; Basu, 1971; Särndal *et al.*, 1992). However, pseudo empirical likelihood assumes complete response with no missing values, that is, it is assumed that no sampled units fails to provide information on the study and auxiliary variables.

Missingness is a common feature in survey data. Item nonresponse occurs when a sampled unit fails to provide information on some variables of interest (study and auxiliary variables). This may occur for various reasons. Sampled individuals may refuse to participate in the study, maybe not contactable by interviewers. Accidental loss of information may also occur. Dealing with missing data in survey research is no simple matter. A variety of methods can be used to adjust for missing values, such as imputation and weighting.

Unit nonresponse occurs when sampled unit fail to provide information on the study variable and on auxiliary variables. In this paper, we assume that unit nonresponse is uniform and part of the sampling design.

With item nonresponse, the simplest solution is eliminate the nonresponding units and apply the pseudo empirical likelihood method to the subset of responding units. However, this method, which Rubin (1987) called complete case analysis, can produce bias in the estimations and lead to greater sampling variance (see Rubin, 1987 or Little and Rubin, 1987).

Imputation is another technique that could be used for item non-response (Little and Rubin, 1987; Rao and Toutenburg, 1995; Särndal, 1992). Imputation consists in substituting missing values by suitable values. Treating the imputed values as if they are true values, and use naively the pseudo empirical likelihood method may lead to invalid inference. For example, the variance can be seriously underestimated when the proportion of missing values is not small (Rao and Shao, 1992; Särndal, 1990, 1992).

Another option is to try to improve the precision of point estimation by including the observed values of the auxiliary variable x for the units where the values of the study variable y is missing. Indeed, we could observe a missing value for y , although the a value for x is observed.

Ratio, difference and product estimators assume complete response. These estimators have only been developed for a limited class of sampling designs. For example, Tracy and Osahan (1994), Toutenburg and Srivastava (1998, 1999, 2000) developed ratio estimators for missing values under Simple Random Sampling Without Replacement (*SRSWOR*) sampling.

We propose modify pseudo empirical likelihood estimator which can be implemented with any unequal probability sampling design. The proposed estimator uses all the sample information collected for the study (y) and the auxiliary (x) variable, as the proposed estimators is function of the values of x for the units with y missing, and function of the values of y for the units with x missing.

5.2.2. The proposed class of estimators

In this section, after introducing some notation, we introduce several estimators that can be used to adjust for missingness. We also introduce the pseudo empirical likelihood method proposed by Chen and Sitter (1999) and the proposed estimator in the presence of missing data.

Consider a population $U = \{1, 2, \dots, N\}$ of N units from which a random sample s of fixed size n is selected according to a specified sampling design with first order inclusion probabilities π_i . Let y_i and x_i be respectively the values of a study variable y and an auxiliary variable, for unit i . The parameter of interest is the population mean $\bar{Y} = \sum_{i=1}^N y_i/N$. It is only possible to know the values of y_i and x_i for $i \in s$. However, some of these values will be missing. Suppose that the population mean $\bar{X} = \sum_{i=1}^N x_i/N$ is known without sampling error from other source or a census.

We consider a situation in which there are some observations missing on one of the characteristics at a time. Thus the missingness occurs for both the characteristics separately but not simultaneously. Suppose that p ($p \geq 0$) units respond to x but not to y ; that is, suppose that we have p missing values for the study variable y . We also have incomplete auxiliary information, that is, q ($q \geq 0$) units respond to y but not to x . Note that p and q are integer. In addition, we have a set of $n - p - q$ units ($p + q \leq n$) that respond to both y and x variables. Thus, the sample data have the following structure:

y_1	\dots	y_{n-p-q}	Missing	\dots	Missing	y_{n-q+1}	\dots	y_n
x_1	\dots	x_{n-p-q}	$x_{n-p-q+1}$	\dots	x_{n-q}	Missing	\dots	Missing

Consider three disjoint sets of sampled units.

$$\begin{aligned}
 s_A &= \{i \in s \mid x_i, y_i \text{ are non-missing}\}, \\
 s_B &= \{i \in s \mid x_i \text{ are non-missing, } y_i \text{ is missing}\}, \\
 s_C &= \{i \in s \mid y_i \text{ are non-missing, } x_i \text{ is missing}\}.
 \end{aligned}$$

Assuming *SRSWOR*, Toutenburg and Srivastava (2000) proposed four es-

timators for the population mean \bar{Y}

$$\bar{y}_{T1} = \bar{y}^A \left[\frac{(n-p-q)\bar{x}^A + p\bar{x}^B}{(n-q)\bar{x}^A} \right], \quad (5.1)$$

$$\bar{y}_{T2} = \bar{y}^A \left[\frac{(n-q)\bar{x}^A}{(n-p-q)\bar{x}^A + p\bar{x}^B} \right], \quad (5.2)$$

$$\bar{y}_{T3} = \left[\frac{((n-p-q)\bar{x}^A + p\bar{x}^B)((n-p-q)\bar{y}^A + q\bar{y}^C)}{(n-q)(n-p)\bar{x}^A} \right], \quad (5.3)$$

$$\bar{y}_{T4} = \left[\frac{(n-p-q)\bar{y}^A + q\bar{y}^C}{(n-p-q)\bar{x}^A + p\bar{x}^B} \right] \left[\frac{n-q}{n-p} \bar{x}^A \right], \quad (5.4)$$

where \bar{y}^i and \bar{x}^i are the sample means based on s_i , with $i = A, B, C$.

The estimators \bar{y}_{T1} and \bar{y}_{T2} are function of the data from s_A and s_B , and do not depend on the data from s_C . However, \bar{y}_{T3} and \bar{y}_{T4} are function of the data from s_A , s_B and s_C . Toutenburg and Srivastava (2000) show that none of the estimators defined above is uniformly superior to other. An appropriate choice of estimator requires the knowledge of population parameters.

Rueda and González (2004) propose several estimators that can be used under any sampling design in the presence of missing values. These estimators are based upon ratio, difference and regression methods. For example, the following estimator is asymptotically unbiased, under *SRSWOR* is asymptotically normal and is better, in the sense of mean squared error, than the other proposed estimators

$$\bar{y}_{Reg} = \hat{\alpha}_{reg} \bar{y}_{HT}^A + (1 - \hat{\alpha}_{reg}) \bar{y}_{HT}^C + \frac{\widehat{Cov}_{i \in s_A}(x, y)}{\widehat{Var}_{i \in s_A}(x)} \left[\bar{X} - \left(\hat{\beta}_{reg} \bar{x}_{HT}^A + (1 - \hat{\beta}_{reg}) \bar{x}_{HT}^B \right) \right] \quad (5.5)$$

where \bar{y}_{HT}^i and \bar{x}_{HT}^i are the Horvitz-Thompson (1952) estimators based on s_i ($i = A, B, C$), $\widehat{Cov}_{i \in s_A}(x, y)$ and $\widehat{Var}_{i \in s_A}(x)$ represent the estimators of the population covariance and variance based on s_A and the optimum values $\hat{\alpha}_{reg}$ and $\hat{\beta}_{reg}$ can be seen in Rueda and González (2004).

In the following, the pseudo empirical likelihood method is described. The pseudo empirical maximum likelihood estimator is given by $\bar{y}_{PE} = \sum_{i \in s} \hat{p}_i y_i$. For non-stratified sampling designs, Chen and Sitter (1999) recommended to find the \hat{p}_i that maximize the so-called pseudo empirical (log) likelihood function $l(p) = \sum_{i \in s} d_i \log p_i$, subject to

$$\sum_{i \in s} p_i = 1 \quad (0 \leq p_i \leq 1), \quad (5.6)$$

$$\sum_{i \in s} p_i u_i = 0, \quad (5.7)$$

where $d_i = 1/\pi_i$ and u_i are known quantities verifying $N^{-1} \sum_{i=1}^N u_i = 0$. Appropriate u_i are necessary to obtain efficient estimators. Wu and Sitter (2001a), Chen and Wu (2002) and Wu (2002) proposed several expressions for u_i , based on superpopulation models. In this paper, we used the most commonly used u_i given by $u_i = x_i - \bar{X}$, which can be justified by a linear model between y and x . The Lagrange multiplier method can be applied to show that

$$\hat{p}_i = \frac{d_i^*}{1 + \lambda u_i}, \text{ for } i \in s, \quad (5.8)$$

where $d_i^* = d_i / \sum_{j \in s} d_j$, and the Lagrange multiplier, λ , is the solution of

$$\sum_{i \in s} \frac{d_i^* u_i}{1 + \lambda u_i} = 0. \quad (5.9)$$

Note that expression (5.8) is valid for any u_i . Chen and Sitter (1999) showed that if $u_i = x_i - \bar{X}$ is used as the calibration variable, the pseudo empirical maximum likelihood estimator is asymptotically equivalent to a regression estimator under mild conditions. If u_i is obtained using superpopulation models, an analogous result can be found for the model-calibration estimator defined in Wu and Sitter (2001a).

Another advantage to the pseudo empirical maximum likelihood is that the resulting weights are always positive, which may not be true for other methods, such as calibration. The intrinsic positive weights associated with this approach make the technique generally applicable to the estimation of distribution functions and quantiles (Chen and Wu, 2002).

Consider the following Hájek ratio estimators. Properties of this estimator are described in Rao (1966), Basu (1971) and Särndal *et al.* (1992).

$$\bar{y}_w^A = \sum_{i \in s_A} d_i^{A*} y_i \quad ; \quad \bar{y}_w^C = \sum_{i \in s_C} d_i^{C*} y_i \quad ; \quad \bar{y}_w^{AC} = \sum_{i \in s_A \cup s_C} d_i^{AC*} y_i; \quad (5.10)$$

$$\bar{x}_w^A = \sum_{i \in s_A} d_i^{A*} x_i \quad ; \quad \bar{x}_w^B = \sum_{i \in s_B} d_i^{B*} x_i \quad ; \quad \bar{x}_w^{AB} = \sum_{i \in s_A \cup s_B} d_i^{AB*} x_i; \quad (5.11)$$

with

$$d_i^{A*} = \frac{d_i^A}{\sum_{j \in s_1} d_j^A}, \quad d_i^{B*} = \frac{d_i^B}{\sum_{j \in s_B} d_j^B}, \quad d_i^{C*} = \frac{d_i^C}{\sum_{j \in s_C} d_j^C}, \quad (5.12)$$

$$d_i^{AB*} = \frac{d_i^{AB}}{\sum_{j \in s_A \cup s_B} d_j^{AB}}, \quad d_i^{AC*} = \frac{d_i^{AC}}{\sum_{j \in s_A \cup s_C} d_j^{AC}}, \quad (5.13)$$

$$d_i^A = 1/\pi_i^A, \quad d_i^B = 1/\pi_i^B, \quad d_i^C = 1/\pi_i^C, \quad d_i^{AB} = 1/\pi_i^{AB}, \quad d_i^{AC} = 1/\pi_i^{AC}, \quad (5.14)$$

The quantities π_i^A , π_i^B , π_i^C , π_i^{AB} and π_i^{AC} are respectively the first order inclusion probabilities of the samples s_A , s_B , s_C , $s_A \cup s_B$ and $s_A \cup s_C$.

Note that when $u_i = 0$, we have $\hat{p}_i = d_i^*$ and the pseudo empirical maximum likelihood estimator is the Hájek ratio estimator give by $\sum_{i \in s} d_i^* y_i$. This estimator does not use the auxiliary variable x .

Consider the pseudo empirical maximum likelihood estimator of \bar{Y} given by

$$\bar{y}_{PE}^A = \sum_{i \in s_A} \hat{p}_i^A y_i,$$

where \hat{p}_i^A maximizes $l(p^A) = \sum_{i \in s_A} d_i^A \log p_i^A$ subject to (5.6) and (5.7) after substituting p_i by p_i^A into (5.6) and (5.7). Considering the Lagrange multiplier method, \hat{p}_i^A is given by (5.8) and (5.9) after substituting d_i by d_i^A into (5.8) and (5.9). Note that all these new expressions are defined for $i \in s_A$.

The estimator \bar{y}_{PE}^A does not use the information provided by the samples s_B and s_C . We now define a pseudo empirical maximum likelihood estimator that takes the information of s_A and s_B into account. Because the interest variable contains $n - p - q$ values, the new vector of weights \hat{p}_i^{AB} must be defined with dimension $n - p - q$. Thus, the new estimator is given by

$$\bar{y}_{PE}^{AB} = \sum_{i \in s_A} \hat{p}_i^{AB} y_i,$$

where the \hat{p}_i^{AB} ($i \in s_A$) are obtained as \hat{p}_i^A (which has dimension $n - p - q$), although we now use the Lagrange multiplier λ^{AB} based on the samples s_A and s_B into (5.8). The quantity λ^{AB} is obtained from (5.9) after substituting d_i by d_i^{AB} into (5.9).

We can use other methods (such as imputation) to obtain a pseudo empirical maximum likelihood estimator that uses the information given by the samples s_A and s_B . For simplicity, imputation is not considered in this paper.

Although \bar{y}_{PE}^{AB} seems better than \bar{y}_{PE}^A , as it uses information from the samples s_A and s_B . However, \bar{y}_{PE}^{AB} is not recommended, as the constraints $\sum_{i \in s_A} \hat{p}_i^{AB} = 1$ and $\sum_{i \in s_A} \hat{p}_i^{AB} u_i = 0$ are not fulfilled. The advantageous properties of this pseudo empirical likelihood method do not hold. In Section 5.2.5,

\bar{y}_{PE}^{AB} and \bar{y}_{PE}^A are compared with the proposed estimator defined in Section 5.2.4.

Unfortunately, the proposed estimator \bar{y}_{PE}^A does not use the information on the study variable y provided by the sample s_C . In order to solve this problem, we propose a class of estimators, which use of all the information on the variable y included in s_A and s_C . An estimator of this class is defined by

$$\bar{y}_{PE\alpha} = \alpha\bar{y}_{PE}^A + (1 - \alpha)\bar{y}_w^C, \quad (5.15)$$

where α is a suitable constant, which is such that $0 < \alpha < 1$. In Section 5.2.4, we propose suitable values for α . The estimator \bar{y}_w^C is defined in (5.10)

An estimator of this class uses all the information available from s_A and s_C . The values of x from the sample s_B are not used for estimation, as y is missing for $i \in s_B$. This would also worsen the estimation. In Section 5.2.5, a simulation study shows that the estimators of the proposed class are as efficient as other estimators that uses information from each s_A , s_B and s_C sample.

5.2.3. Asymptotic properties

We now show that $\bar{y}_{PE\alpha}$ defined by (5.15), is asymptotically unbiased. We also derive the asymptotic variance of $\bar{y}_{PE\alpha}$.

The following assumptions are made.

$$(A1) \quad u^{A*} = \max_{i \in s_A} |u_i| = o_p(n^{1/2}).$$

$$(A2) \quad \frac{\sum_{i \in s_A} d_i^A u_i}{\sum_{i \in s_A} d_i^A u_i^2} = O_p(n^{1/2}).$$

These assumptions are given by Chen and Sitter (1999) who show that many commonly used sampling designs satisfy the above assumptions. Given the above assumptions, the following results are obtained.

Result 1. *Under assumptions (A1) and (A2), we have*

$$\bar{y}_{PE\alpha} = \alpha\bar{y}_{GREG}^A + (1 - \alpha)\bar{y}_w^C + o_p(n^{-1/2}), \quad (5.16)$$

where

$$\bar{y}_{GREG}^A = \bar{y}_w^A + (\bar{X} - \bar{x}_w^A)b, \quad (5.17)$$

with

$$b = \frac{\sum_{i \in s_A} d_i^{A*} x_i y_i - \bar{y}_w^A \bar{x}_w^A}{\sum_{i \in s_A} d_i^{A*} (x_i - \bar{x}_w^A)^2}. \quad (5.18)$$

Proof of Result 1: Chen and Sitter (1999) showed that \bar{y}_{PE}^A is asymptotically equivalent to \bar{y}_{GREG}^A . The results easily follows. \square

Result 2. Under assumptions (A1) and (A2), we have

$$\bar{y}_{GREG}^A \simeq \bar{y}_{GREG}^{A2},$$

where

$$\bar{y}_{GREG}^{A2} = \bar{y}_w^A + (\bar{X} - \bar{x}_w^A)B, \quad (5.19)$$

with

$$B = \frac{Cov(x, y)}{Var(x)}. \quad (5.20)$$

Proof of Result 2: To establish this result, we assume the finite population embeds in a sequence of populations where n and N increase such that $n/N \rightarrow f$ when $n \rightarrow \infty$ and where f is a constant.

We now use a result due to Randles (1982) who shows that if a statistic $T_n(\hat{\lambda})$ and $\hat{\lambda}$ are function of the data and $\hat{\lambda}$ is an estimator of λ , then $T_n(\hat{\lambda})$ and $T_n(\lambda)$ have the same limiting distribution provided

$$\left. \frac{\partial \mu(\gamma)}{\partial \gamma} \right|_{\gamma=\lambda} = 0,$$

where $\mu(\gamma) = \lim_{n \rightarrow +\infty} E_\lambda[T_n(\gamma)]$ and the expectation is taken when the true parameter is λ .

Let $T_n(\gamma) = \bar{y}_w^A + (\bar{X} - \bar{x}_w^A)\gamma$. Note that $T_n(b) = \bar{y}_{GREG}^A$ defined by (5.17). Consider $\mu(\gamma) = \lim_{n \rightarrow \infty} E_\gamma[T_n(\gamma)]$. Note that when $\gamma = B$ defined by (5.20), we have $\mu(B) = \tilde{Y}$ where $\tilde{Y} = \lim_{n \rightarrow \infty} \bar{Y}$. Since $\mu(\gamma)$ verifies

$$\left. \frac{\partial \mu(\gamma)}{\partial \gamma} \right|_{\gamma=B} = 0,$$

this implies that $\bar{y}_{GREG}^A \simeq \bar{y}_{GREG}^{A2}$. This completes the proof. \square

Using results 1 and 2, we have

$$\bar{y}_{PE\alpha} \simeq \alpha \bar{y}_{GREG}^{A2} + (1 - \alpha) \bar{y}_w^C, \quad (5.21)$$

which implies that $\bar{y}_{PE\alpha}$ is asymptotically design unbiased.

Result 3. *Under assumptions (A1) and (A2), the asymptotic variance of $\bar{y}_{PE\alpha}$ is given by*

$$\begin{aligned} AV(\bar{y}_{PE\alpha}) &= \alpha^2 \left[V(\bar{y}_w^A) + B^2 V(\bar{x}_w^A) - 2BCov(\bar{y}_w^A, \bar{x}_w^A) \right] + \\ &+ (1 - \alpha)^2 V(\bar{y}_w^C) + 2\alpha(1 - \alpha) \left[Cov(\bar{y}_w^A, \bar{y}_w^C) - BCov(\bar{x}_w^A, \bar{y}_w^C) \right]. \end{aligned} \quad (5.22)$$

Proof of Result 3: The approximation (5.21) implies that the asymptotic variance of $\bar{y}_{PE\alpha}$ is given by

$$V\left(\alpha \bar{y}_{GREG}^{A2} + (1 - \alpha) \bar{y}_w^C\right) = \alpha^2 V(\bar{y}_{GREG}^{A2}) + (1 - \alpha)^2 V(\bar{y}_w^C) + 2\alpha(1 - \alpha) Cov(\bar{y}_{GREG}^{A2}, \bar{y}_w^C). \quad (5.23)$$

Using (5.19), the variance of \bar{y}_{GREG}^{A2} is

$$\begin{aligned} V(\bar{y}_{GREG}^{A2}) &= V\left(\bar{y}_w^A + (\bar{X} - \bar{x}_w^A)B\right) \\ &= V\left(\bar{y}_w^A - \bar{x}_w^A B\right) \\ &= V(\bar{y}_w^A) + B^2 V(\bar{x}_w^A) - 2BCov(\bar{y}_w^A, \bar{x}_w^A). \end{aligned} \quad (5.24)$$

The value $Cov(\bar{y}_{GREG}^{A2}, \bar{y}_w^C)$ is

$$Cov(\bar{y}_{GREG}^{A2}, \bar{y}_w^C) = Cov(\bar{y}_w^A, \bar{y}_w^C) - BCov(\bar{x}_w^A, \bar{y}_w^C). \quad (5.25)$$

Thus, from (5.23), (5.24) and (5.25), the asymptotic variance of $\bar{y}_{PE\alpha}$ is given by (5.22). Result 3 follows. \square

5.2.4. The optimal estimators of the proposed class

The optimal estimator of the proposed class is the estimator defined by (5.15) with a value α that minimizes the asymptotic variance given by (5.22).

The asymptotic variance (5.22) can be re-written as

$$AV(\bar{y}_{PE\alpha}) = \alpha^2 M^* + (1 - \alpha)^2 N^* + 2\alpha(1 - \alpha) L^*,$$

where

$$M^* = V(\bar{y}_w^A) + B^2V(\bar{x}_w^A) - 2BCov(\bar{y}_w^A, \bar{x}_w^A), \quad (5.26)$$

$$N^* = V(\bar{y}_w^C), \quad (5.27)$$

$$L^* = Cov(\bar{y}_w^A, \bar{y}_w^C) - BCov(\bar{x}_w^A, \bar{y}_w^C). \quad (5.28)$$

The value α_{opt} that minimizes the asymptotic variance is the solution of the equation

$$\left. \frac{\partial AV(\bar{y}_{PE\alpha})}{\partial \alpha} \right|_{\alpha=\alpha_{opt}} = 2\alpha_{opt}M^* - 2(1 - \alpha_{opt})N^* + 2(1 - 2\alpha_{opt})L^* = 0,$$

which implies

$$\alpha_{opt} = \frac{N^* - L^*}{M^* + N^* - 2L^*}. \quad (5.29)$$

By substituting α_{opt} into (5.22), we obtain the smallest asymptotic variance given by

$$AV(\bar{y}_{PE\alpha_{opt}}) = \alpha_{opt}^2 M^* + (1 - \alpha_{opt})^2 N^* + 2\alpha_{opt}(1 - \alpha_{opt})L^*. \quad (5.30)$$

The optimum value α_{opt} depends on unknown population parameters which can be estimated from sample data.

Under *SRSWOR* and stratified sampling, we have $\sum_{i \in s} d_i = N$, that is, the Horvitz-Thompson estimator and the Hájek ratio estimator are equal, and so the estimators of variances and covariances in (26), (27) and (28) can be easily obtained. An analytic expression for (5.26), (5.27) and (5.28) under *SRSWOR* can be found in Rueda and González (2004).

With these estimates, we obtain an approximation $\tilde{\alpha}_{opt}$ of α_{opt} . The proposed estimator is therefore

$$\tilde{y}_{PE\alpha_{opt}} = \tilde{\alpha}_{opt}\bar{y}_{PE}^A + (1 - \tilde{\alpha}_{opt})\bar{y}_w^C. \quad (5.31)$$

It is also possible to establish the asymptotic unbiasedness of $\tilde{y}_{PE\alpha_{opt}}$.

5.2.5. An empirical study

In this section, the proposed estimators are compared with other alternative estimators with a simulation based on simulated and real populations.

The Fam1500 population consists on $N = 1500$ Households in Andalusia (Spain) where the study variable is the feeding expenses (y). We have two

auxiliary variables: the incomes (x_1) and other expenses (x_2). This population was first used by Fernández *et al.* (1994). The correlation coefficients are given by $\rho_{y,x_1} = 0,848$ and $\rho_{y,x_2} = 0,546$.

The second population is a national sample of hospitals in the U.S. (Royall and Cumberland, 1981 and Valliant *et al.*, 2000). The population size is $N = 393$, the study variable is the number of patients discharged (y), whereas the auxiliary variable is the number of beds (x). The correlation coefficient is given by $\rho_{y,x} = 0,911$.

We also generate four finite populations consisting of $N = 2000$ units using the model

$$y = \theta_0 + \theta_1 x + \epsilon, \quad (5.32)$$

where $x \sim \text{Gamma}(1, 1)$, $\epsilon \sim N(0, \sigma^2)$ and $\theta_0 = \theta_1 = 1$. This model has also been used by Wu and Sitter (2001a). These four populations are obtained by choosing different values of σ^2 and different values for the correlation coefficients between y and x (0.6, 0.7, 0.8, and 0.9). These populations are called Pop06, Pop07, Pop08 and Pop09, respectively.

Let us compare the behaviour of \bar{y}_{PE}^A and $\tilde{\bar{y}}_{PE\alpha_{opt}}$ with the following estimators: (i) the standard estimator \bar{y}_w^{AC} , which is the Hájek type estimator for the mean based on the samples s_A and s_C ; (ii) \bar{y}_{T1} , \bar{y}_{T2} , \bar{y}_{T3} and \bar{y}_{T4} , the estimators proposed by Toutenburg and Srivastava (2000); (iii) \bar{y}_{PE}^{AB} , the pseudo empirical maximum likelihood estimator based on the samples s_A and s_B and (iv) \bar{y}_{Reg} , the estimator proposed by Rueda and Gonzalez (2004) based on the samples s_A , s_B and s_C .

For each of the six populations, we generated $B = 1000$ independent samples under *SRSWOR* with a sample size n . For some values of q and p , we generated randomly missing values for x and y . The performance of the estimators is measured in terms of the Relative Bias (*RB*) and the Relative Efficiency (*RE*) defined by

$$RB_j = \frac{1}{B} \sum_{b=1}^B \frac{|\hat{y}_j(b) - \bar{Y}|}{\bar{Y}} \quad ; \quad RE_j = \frac{MSE(\hat{y}_j)}{MSE(\bar{y}_w^{AC})},$$

where $\hat{y}_j(b)$ is the value of an estimate for the b th sample selected, the empirical Mean Square Error is given by $MSE(\hat{y}_j) = B^{-1} \sum_{b=1}^B (\hat{y}_j(b) - \bar{Y})^2$. The index $j = 1, \dots, 8$ refers to the \bar{y}_{PE}^A , \bar{y}_{PE}^{AB} , $\tilde{\bar{y}}_{PE\alpha_{opt}}$, \bar{y}_{Reg} , \bar{y}_{T1} , \bar{y}_{T2} , \bar{y}_{T3} and \bar{y}_{T4} estimators.

The simulation was programmed in *R* and the source codes are available

from the authors upon request.

We have observed that the \bar{y}_{T3} is always better than \bar{y}_{T1} , \bar{y}_{T2} and \bar{y}_{T4} . For clarity in Figures 5.1, . . . , 5.6, the lines corresponding to \bar{y}_{T1} , \bar{y}_{T2} and \bar{y}_{T4} estimators are not included. The comparisons between these estimators are also available from the authors.

Figures 5.1, 5.2 and 5.3 represent the *RE* (vertical axis) for \bar{y}_{PE}^A , \bar{y}_{PE}^{AB} , $\tilde{y}_{PE\alpha_{opt}}$, \bar{y}_{Reg} and \bar{y}_{T3} estimators under *SRSWOR* and different values of p and q . Horizontal dotted lines represent the *RE* for \bar{y}_w^{AC} , the standard estimator.

From Figures 5.1, 5.2 and 5.3, we observe that all the estimators (except \bar{y}_{T3}) are better than the standard estimator if y and x have a stronger correlation and if p and q are small. If p and q are large, all the estimators are worst than \bar{y}_w^{AC} . The estimator $\tilde{y}_{PE\alpha_{opt}}$ has the smallest *RE*. The Toutenburg estimator has the largest *RE*. This might be due to the fact that this latter estimator does not use \bar{X} as auxiliary information.

The estimators \bar{y}_{PE}^A and $\tilde{y}_{PE\alpha_{opt}}$ have the same precision under a strong relationship between y and x and when the number of missing values is small. The gain in efficiency from $\tilde{y}_{PE\alpha_{opt}}$ with regard to \bar{y}_{PE}^A is larger in the contrary case. The estimator \bar{y}_{PE}^{AB} never outperforms \bar{y}_{PE}^A or $\tilde{y}_{PE\alpha_{opt}}$. The reason for this is that its weights are not well defined.

In the Hospitals and Fam1500 (when we use x_1) populations, \bar{y}_{PE}^A , $\tilde{y}_{PE\alpha_{opt}}$ and \bar{y}_{Reg} have the same precision. In the simulated and Fam1500 (when we use x_2) populations, \bar{y}_{Reg} never outperforms $\tilde{y}_{PE\alpha_{opt}}$. Although \bar{y}_{Reg} uses information from s_A , s_B and s_C , $\tilde{y}_{PE\alpha_{opt}}$ is considerably more efficient when the relationship between y and x is weak and when p and q are large.

The estimator \bar{y}_w^{AC} is more efficient than $\tilde{y}_{PE\alpha_{opt}}$ only when the relationship is weak and the total number of missing data, $p+q$, is large. When the number of missing x -values, p , is large, the gain in precision of the proposed estimator with respect to \bar{y}_w^{AC} is smaller. Similarly, when p is fixed, the gain in precision decreases as the number of missing q grows. This result can be explained by the fact that if p/q is small, more information is provided by the sample s_C and less information is provided by s_B , and the proposed estimator only used the information given by s_C like \bar{y}_w^{AC} .

Figures 5.4, 5.5 and 5.6 give the *RB* for the estimators. We observe that the *RB* are all within a reasonable range, with the \bar{y}_{PE}^A and $\tilde{y}_{PE\alpha_{opt}}$ estimators having the smallest *RB*.

Figure 5.1: Relative Efficiency for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Pop06, Pop07, Pop08 and Pop09 populations and *SRSWOR* with $n = 200$ are considered.

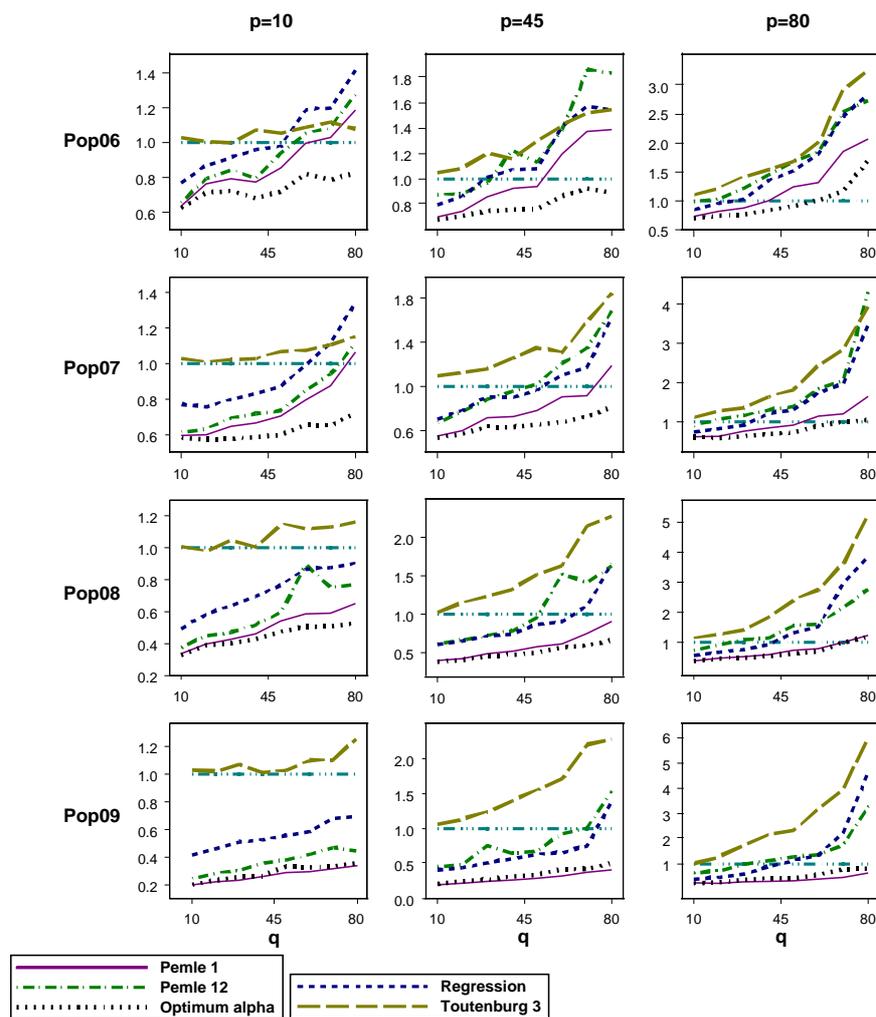


Figura 5.2: Relative Efficiency for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Fam1500 population and *SRSWOR* with $n = 150$ are considered.

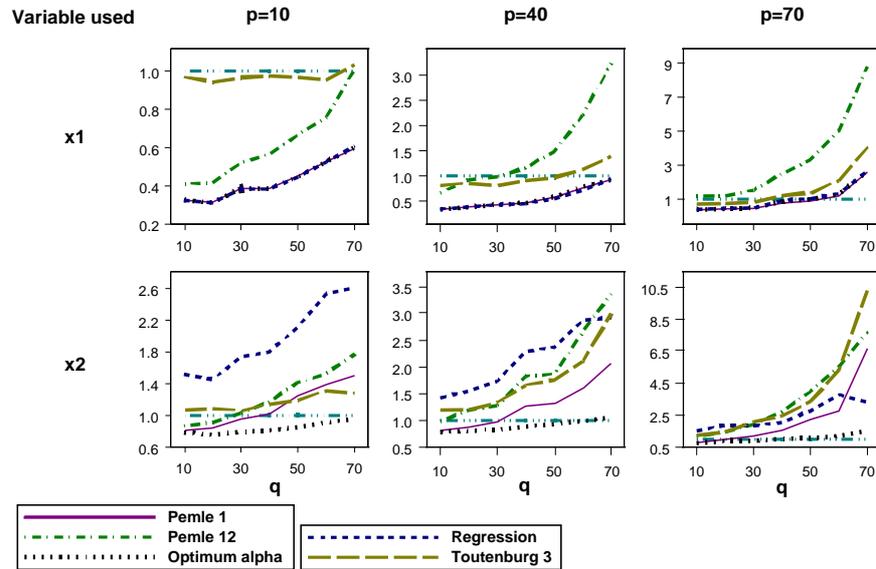


Figura 5.3: Relative Efficiency for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Hospitals population and *SRSWOR* with $n = 100$ are considered.

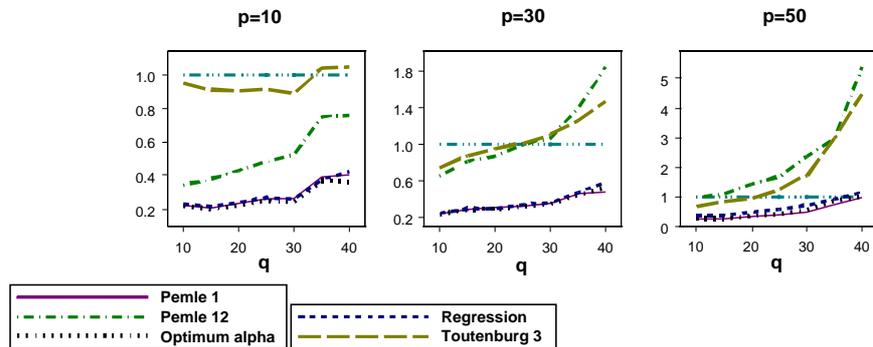


Figura 5.4: Relative Bias for \bar{y}_{PE}^A (Pemle 1), \bar{Y}_{PE}^{AB} (Pemle 12), $\tilde{\bar{y}}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_w^{AC} (Standard), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Pop06, Pop07, Pop08 and Pop09 populations and *SRSWOR* with $n = 200$ are considered.

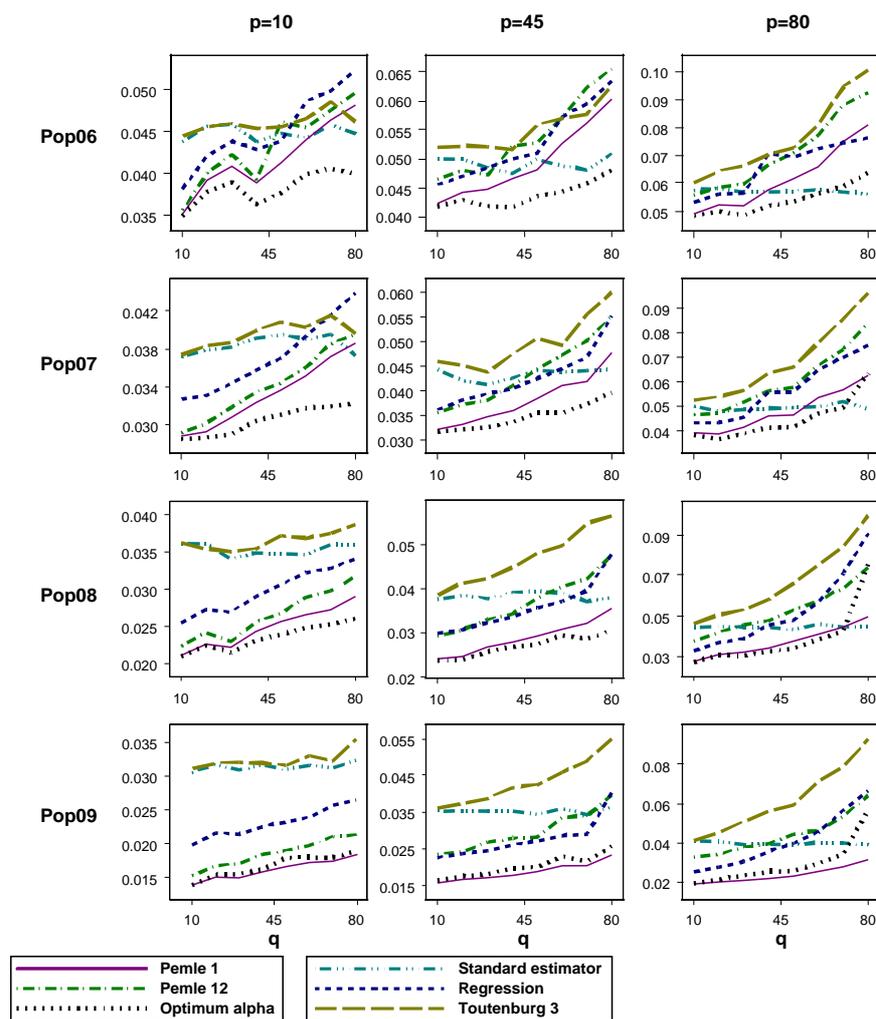


Figura 5.5: Relative Bias for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{y}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_w^{AC} (Standard), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Fam1500 population and *SRSWOR* with $n = 150$ are considered.

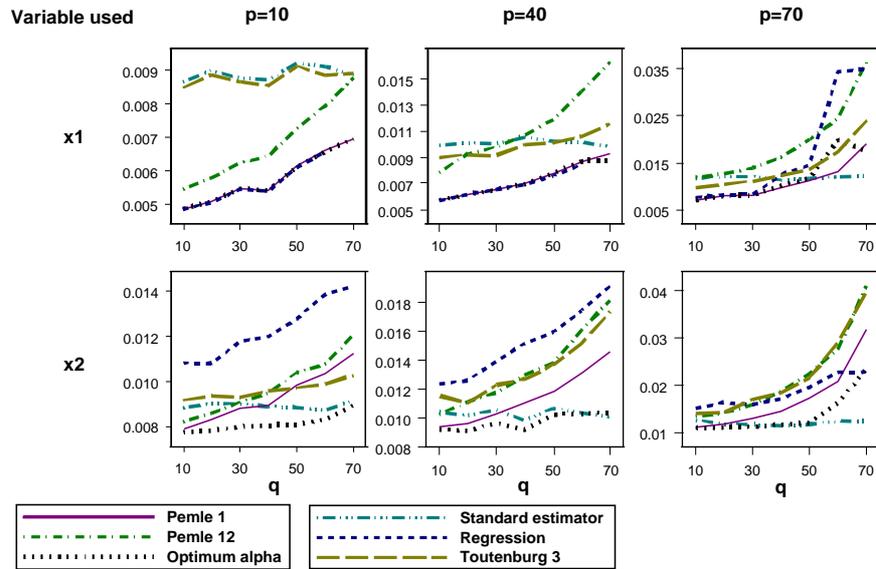
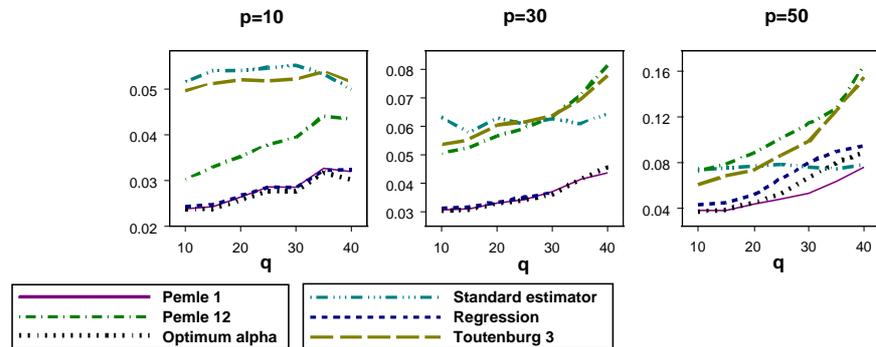


Figura 5.6: Relative Bias for \bar{y}_{PE}^A (Pemle 1), \bar{y}_{PE}^{AB} (Pemle 12), $\tilde{y}_{PE\alpha_{opt}}$ (Optimum alpha), \bar{y}_w^{AC} (Standard), \bar{y}_{Reg} (Regression) and \bar{y}_{T3} (Toutenburg 3) estimators. The Hospitals population and *SRSWOR* with $n = 100$ are considered.



5.3. Quantile estimation under two phase sampling

The estimation of quantiles in two-phase sampling with arbitrary sampling design in each of the two phases is investigated. Several ratio and exponentiation type estimators that provide the optimum estimate of a quantile based on an optimum exponent α are proposed. Properties of these estimators are studied under large sample size approximation and the use of double sampling for stratification to estimate quantiles can also be seen. The real performance of these estimators will be evaluated for the three quartiles on the basis of data from two real populations using different sampling designs. The simulation study shows that proposed estimators can be very satisfactory in terms of relative bias and efficiency.

5.3.1. Introduction

The problem of estimating a population mean in the presence of an auxiliary variable has been widely discussed in the finite population sampling literature. However, for the problem of estimating a population median, the situation is quite different and only recently has this problem been discussed. Rao *et al.* (1990) proposed ratio and difference estimators for the median using a design-based approach. Kuk and Mak (1989) proposed two estimators for which it was only necessary to know the values of the median of the auxiliary variable for the whole population. More recently, Rueda *et al.* (1998) and Rueda and Arcos (2001) proposed confidence intervals for quantiles based on ratio and difference estimators of the distribution function. In Rueda *et al.* (2003, 2004) the population information is used through a quantile of the auxiliary variable with the same or different order as that of the quantile of the main variable considered for estimation using difference type estimators.

The above estimators are based on prior knowledge of the median $Q_x(0,5)$ of the auxiliary characteristic. In many cases $Q_x(0,5)$ may not be known, and it may be seen that taking the sample selection in two phases is an attractive solution.

Two-phase sampling is a good compromise for surveys in which no prior knowledge is available about the population. A key to successful two-phase sampling is the creation of a highly informative frame for the part of the population from which the subsample is drawn. The estimation of the median

in two-phase sampling is developed by Singh *et al.* (2001), Singh (2003) and Allen *et al.* (2002). Swamy *et al.* (2005) have shown that auxiliary information, without knowing its true functional form, can also be used to reduce the bias while estimating the relation among the federal funds and the Federal Reserve's expectations about future values of certain policy variables is considered.

These papers have been developed using simple random sampling. Sampling surveys for economic variables (as income) that possess highly skewed distributions are almost always complex in structure, and methods such as stratification and probability proportional to size are common place.

In this section we propose various estimators of a β -quantile in two-phase sampling with arbitrary sampling designs in each of the two phases.

5.3.2. Quantile direct estimation

This study has been carried out under the fixed population approach. Let U be a finite population with N different elements where y_1, \dots, y_N are the values of the variable of interest y , and $F_y(t) = N^{-1} \sum_{i=1}^N \delta(t - y_i)$, $(-\infty < t < \infty)$, is the population distribution function, where $\delta(a)$ takes the value 1 if $a \geq 0$ and the value 0 otherwise. Let x be an auxiliary variable and x_i ($i = 1, \dots, N$) be the value of its i th population unit.

The first-phase sample s' of size n' is drawn according to a sampling design d_1 , such that $p_{d1}(s')$ is the probability that s' is chosen and where the corresponding first and second order probabilities are π'_i and π'_{ij} for i and $j \in U$. For the elements in s' , information of the auxiliary variable can be recorded. Given s' , the second-phase sample s of size n is drawn according to the design d_2 such that $p(s/s')$ is the conditional probability of choosing s . The inclusion probabilities under this design are denoted by $\pi_{i/s'}$ and $\pi_{ij/s'}$.

A particular case is presented when the variable x is used to stratify s' into L strata denoted by s'_h , ($h = 1, \dots, L$), with n'_h elements in the h th stratum. In this way, a sample s_h of size n_h can be drawn from s'_h according to a design $p_h(/s')$ independently from each stratum. The final sample is $s = \bigcup_{h=1}^L s_h$. This particular design is called *Two-phase sampling for stratification*.

Without using auxiliary information, the natural candidate to estimate the β -quantile $Q_y(\beta)$ is $\hat{Q}_y(\beta) = \inf\{t \mid \hat{F}_{HTy}(t) \geq \beta\} = \hat{F}_{HTy}^{-1}(\beta)$, where $\hat{F}_{HTy}(t) = N^{-1} \sum_{i \in s} \delta(t - y_i)/\pi_i$ is the Horvitz and Thompson (1952) type estimator of $F_y(t)$ and the inclusion probability of the i th element is given by

$$\pi_i = \sum_{s' \ni i} p_{d1}(s') \pi_{i/s'}.$$

Consequently, to determine π_i we must know the probabilities $\pi_{i/s'}$ for every s' , which we ordinarily do not, because $\pi_{i/s'}$ may depend on the outcome of phase one (for example if the second-phase sample is drawn by a sampling proportional to an auxiliary variable).

Because the Horvitz-Thompson estimator of a mean cannot always be used in practice, in two phase sampling, Särndal *et al.* (1992) proposed the use of π^* estimators. Using this idea, we introduce the quantities

$$\pi'_i = \sum_{s' \ni i} p_{d1}(s'), \quad \pi'_{ij} = \sum_{s' \ni i, j} p_{d1}(s'), \quad \pi_i^* = \pi'_i \cdot \pi_{i/s'} \quad \text{and} \quad \pi_{ij}^* = \pi'_{ij} \cdot \pi_{ij/s'},$$

to define the π^* -estimator of the distribution function as

$$\widehat{F}_{HTy}^*(t) = \frac{1}{N} \sum_{i \in s} \frac{\delta(t - y_i)}{\pi_i^*},$$

and thus, we suggest the following direct estimator of the β -quantile:

$$\widehat{Q}_y^*(\beta) = \widehat{F}_{HTy}^{*-1}(\beta). \quad (5.33)$$

Note that $\widehat{Q}_y^*(\beta)$ does not generally agree with the estimator $\widehat{Q}_y(\beta)$ except in rare cases, but it makes direct calculation possible for all sample designs d_1 and d_2 used in each phase.

We now study the properties of the $\widehat{Q}_y^*(\beta)$ estimator. For this, a linear approximation is needed because $\widehat{Q}_y^*(\beta)$ is not a continuous function.

The estimator $\widehat{Q}_y^*(\beta)$ can be expressed asymptotically as a linear function of the estimated distribution function evaluated at the quantile $Q_y(\beta)$ by the Bahadur representation (see Chambers and Dunstan, 1986):

$$\widehat{Q}_y^*(\beta) - Q_y(\beta) = \frac{1}{f_y(Q_y(\beta))} (\beta - \widehat{F}_{HTy}^*(Q_y(\beta))) + O(n^{-1/2}), \quad (5.34)$$

where $f_y(\cdot)$ denotes the derivative of the limiting value of $F_y(\cdot)$ as $N \rightarrow \infty$. This linear approximation previously used by Kuk and Mak (1989) and Chen and Wu (2002) helps to study the asymptotic properties of the estimator.

On the one hand, the estimator $\widehat{Q}_y^*(\beta)$ is asymptotically unbiased because $\widehat{F}_{HTy}^*(t)$ is an unbiased estimator of $F(t)$. In this way, $E(\beta - \widehat{F}_{HTy}^*(Q_y(\beta))) = 0$, and by using (5.34) it can be seen that $E(\widehat{Q}_y^*(\beta)) = Q_y(\beta) + O(n^{-1/2})$.

On the other hand, from (5.34) we obtain the asymptotic variance of $\widehat{Q}_y^*(\beta)$, to the first degree of approximation, as:

$$V(\widehat{Q}_y^*(\beta)) = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in U} (\pi'_{ij} - \pi'_i \pi'_j) \frac{\delta(Q_y(\beta) - y_i)}{\pi'_i} \frac{\delta(Q_y(\beta) - y_j)}{\pi'_j} + E_{d1} \left[\sum_{i,j \in s'} (\pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}) \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_y(\beta) - y_j)}{\pi_j^*} \right] \right),$$

and one can construct an unbiased estimator of the variance as:

$$\widehat{V}(\widehat{Q}_y^*(\beta)) = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in s} \frac{\pi'_{ij} - \pi'_i \pi'_j}{\pi_{ij}^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi'_i} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi'_j} + \sum_{i,j \in s} \frac{\pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j^*} \right).$$

An approximate value of $f_y(Q_y(\beta))$ can be obtained by applying standard methods such as the kernel or the k th nearest neighbour methods (Silverman, 1986). The variance estimator is stated in an explicit form (it does not depend on the expected value over the first phase design), thus making direct calculation possible.

5.3.3. Estimation using auxiliary information

In the previous section an estimator is defined without using auxiliary information. We now define a class of estimators that takes the auxiliary variable into account.

Assuming simple random and without replacement (*SRSWOR*) sampling and the median of the variable x is known, Kuk and Mak (1989) proposed a ratio estimator for the population median as:

$$\widehat{Q}_y^r(0,5) = \widehat{Q}_y(0,5) \frac{Q_x(0,5)}{\widehat{Q}_x(0,5)}.$$

Furthermore, Kuk and Mak (1989) proposed other estimators of quantiles under *SRSWOR* design called position and stratification estimators, but the extension of them to more complex sampling designs is very difficult.

Rueda *et al.* (2003, 2004) proposed, for any sampling design d and for any β , difference and exponentiation methods to estimate a β -quantile. Singh *et al.* (2001) suggested ratio, regression, position and stratification estimators of the median when the sample is drawn in two phases, using *SRSWOR* in both phases. Under this sampling design, Allen *et al.* (2002) proposed two classes of estimators for the population median using information on two auxiliary variables x and z in double sampling when the population median of z is known.

Proposed estimators

Here, we present a class of estimators of finite population quantiles when the sample is drawn using a general two-phase sampling, described earlier, as:

$$\widehat{Q}_y^{\mathcal{H}}(\beta) = H(\widehat{Q}_y^*(\beta), t^*), \quad (5.35)$$

with $t^* = \widehat{Q}_x^*(\beta)/\widehat{Q}_x'(\beta)$, and $\widehat{Q}_x'(\beta)$ being the estimator of $Q_x(\beta)$ from the first stage of sampling, i.e. $\widehat{Q}_x'(\beta) = \inf\{t \mid \widehat{F}_{HTx}'^{-1}(t) \geq \beta\}$, where $\widehat{F}_{HTx}'(t) = N^{-1} \sum_{i \in s'} \delta(t - x_i)/\pi_i'$. The function H satisfies the following conditions:

1. It assumes values in a closed convex subset $\mathcal{C} \subset \mathbb{R}_2$ which contains the point $(Q_y(\beta), 1)$;
2. H is a continuous function in \mathcal{C} such that $H(Q_y(\beta), 1) = Q_y(\beta)$, and
3. The first and second order partial derivatives of H exist and are also continuous in \mathcal{C} , with

$$H_{10}(Q_y(\beta), 1) = \left. \frac{\partial H(q, t^*)}{\partial q} \right|_{(q, t^*) = (Q_y(\beta), 1)} = 1.$$

A particular case within the general class of estimators \mathcal{H} is the ratio type estimator:

$$\widehat{Q}_{yr}^*(\beta) = \widehat{Q}_y^*(\beta) \frac{\widehat{Q}_x'(\beta)}{\widehat{Q}_x^*(\beta)},$$

which corresponds to the choice $H(q, t^*) = q/t^*$.

Another estimator of the β -quantile, called the exponentiation estimator, can be derived from:

$$\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}_x'(\beta)}{\widehat{Q}_x^*(\beta)} \right)^\alpha,$$

with α as a fixed constant, which corresponds to the choice of $H(q, t^*) = q/(t^*)^\alpha$.

Note 1. If $\alpha = 0$ then $\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_y^*(\beta)$, i.e. $\widehat{Q}_{ye}^*(\beta)$ coincides with the π^* -estimator, if $\alpha = 1$ then $\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_{yr}^*(\beta)$, and if $\alpha = -1$ then $\widehat{Q}_{ye}^*(\beta) = \widehat{Q}_{yp}^*(\beta)$. This we can define as a product estimator.

Note 2. If *SRSWOR* sampling is used in each phase and $\beta=0.5$, the proposed estimators $\widehat{Q}_{yr}^*(\beta)$ and $\widehat{Q}_{ye}^*(\beta)$ lead, respectively, to the estimators $\widehat{M}_y^{(a)}$ and $\widehat{M}_y^{(b)}$ proposed by Singh *et al.* (2001).

Properties of the class of estimators

Any estimator in \mathcal{H} is asymptotically unbiased for $Q_y(\beta)$. This result can be obtained from the following expressions:

$$\begin{aligned}\widehat{Q}_y^*(\beta) - Q_y(\beta) &= \frac{1}{f_y(Q_y(\beta))}(\beta - \widehat{F}_{HTy}^*(Q_y(\beta))) + O(n^{-1/2}), \\ \widehat{Q}_x^*(\beta) - Q_x(\beta) &= \frac{1}{f_x(Q_x(\beta))}(\beta - \widehat{F}_{HTx}^*(Q_x(\beta))) + O(n^{-1/2}), \\ \widehat{Q}'_x(\beta) - Q_x(\beta) &= \frac{1}{f_x(Q_x(\beta))}(\beta - \widehat{F}'_{HTx}(Q_x(\beta))) + O(n^{-1/2}),\end{aligned}$$

and by using the first order Taylor's series expansion for H about the point $(Q_y(\beta), 1)$:

$$\begin{aligned}\widehat{Q}_y^{\mathcal{H}}(\beta) &= H((Q_y(\beta), 1)) + \left(\widehat{Q}_y^*(\beta) - Q_y(\beta)\right) H_{10}(Q_y(\beta), 1) + \\ &+ (t^* - 1)H_{01}(Q_y(\beta), 1) + O(n^{-1}),\end{aligned}\tag{5.36}$$

where H_{10} and H_{01} denote the first order partial derivatives of H with respect to q and t^* , respectively.

When $\widehat{F}_{HTy}^*(t)$ and $\widehat{F}_{HTx}^*(t)$ are unbiased estimators of $F_y(t)$ and $F_x(t)$, respectively, any estimator in \mathcal{H} is asymptotically unbiased for $Q_y(\beta)$.

Asymptotic expression of variances

Consider the Taylor's series expansion (5.36) and consequently the expression

$$\widehat{Q}_y^{\mathcal{H}}(\beta) - Q_y(\beta) = \left(\widehat{Q}_y^*(\beta) - Q_y(\beta)\right) + \left(\frac{\widehat{Q}_x^*(\beta)}{\widehat{Q}'_x(\beta)} - 1\right) H_{01}(Q_y(\beta), 1) + O(n^{-1}).$$

Then, we have

$$\begin{aligned}\widehat{Q}_y^{\mathcal{H}}(\beta) - Q_y(\beta) &= Q_y(\beta)e_0 + \frac{e_1 - e_2}{1 + e_2}H_{01}(Q_y(\beta), 1) \\ &\simeq Q_y(\beta)e_0 + (e_1 - e_2)(1 - e_2)H_{01}(Q_y(\beta), 1) \\ &= Q_y(\beta)e_0 + (e_1 - e_2)H_{01}(Q_y(\beta), 1) - e_2(e_1 - e_2)H_{01}(Q_y(\beta), 1),\end{aligned}$$

where: $e_0 = \frac{\widehat{Q}_y^*(\beta)}{Q_y(\beta)} - 1$, $e_1 = \frac{\widehat{Q}_x^*(\beta)}{Q_x(\beta)} - 1$ and $e_2 = \frac{\widehat{Q}_x'(\beta)}{Q_x(\beta)} - 1$, and we obtain, to the first order of approximation, the variance:

$$\begin{aligned}V(\widehat{Q}_y^{\mathcal{H}}(\beta)) &= Q_y(\beta)^2V(e_0) + H_{01}(Q_y(\beta), 1)^2V(e_1 - e_2) + \\ &\quad + 2H_{01}(Q_y(\beta), 1)Q_y(\beta)Cov(e_0, e_1 - e_2).\end{aligned}$$

On the other hand, in two phase sampling:

$$V(\widehat{Q}_y^{\mathcal{H}}(\beta)) = E_{d1}V(\widehat{Q}_y^{\mathcal{H}}(\beta)/s') + V_{d1}E(\widehat{Q}_y^{\mathcal{H}}(\beta)/s')$$

reflects the variation due to each of the two phases of sampling. Using the known properties of the Horvitz-Thompson estimator and its variance and by denoting $\Delta'_{ij} = \pi'_{ij} - \pi'_i\pi'_j$ and $\Delta^{s'}_{ij} = \pi_{ij/s'} - \pi_{i/s'}\pi_{j/s'}$, we obtain:

$$V_{d1}E(\widehat{Q}_y^{\mathcal{H}}(\beta)/s') = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi'_i} \frac{\delta(Q_y(\beta) - y_j)}{\pi'_j} \right)$$

and

$$\begin{aligned}E_{d1}V(\widehat{Q}_y^{\mathcal{H}}(\beta)/s') &= E_{d1} \left(\frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \sum_{i,j \in s'} \Delta^{s'}_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_y(\beta) - y_j)}{\pi_j^*} \right) + \\ &\quad + \frac{H_{01}^2(Q_y(\beta), 1)}{Q_x^2(\beta)} \frac{1}{N^2} \frac{1}{f_x^2(Q_x(\beta))} \sum_{i,j \in s'} \Delta^{s'}_{ij} \frac{\delta(Q_x(\beta) - x_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*} + \\ &\quad + 2 \frac{H_{01}(Q_y(\beta), 1)}{Q_x(\beta)} \frac{1}{N^2} \frac{1}{f_y(Q_y(\beta))f_x(Q_x(\beta))} \sum_{i,j \in s'} \Delta^{s'}_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*}.\end{aligned}$$

The last variance is not stated explicitly, but as an expected value over the first phase design. This causes no problem for the variance estimation,

$$\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(Q_y(\beta) - y_i)}{\pi'_i} \frac{\delta(Q_y(\beta) - y_j)}{\pi'_j}$$

which can be estimated by

$$\sum_{i,j \in s} \frac{\Delta'_{ij}}{\pi_{ij}^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i'} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j'},$$

and

$$E_{d1} \left(\sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_y(\beta) - y_j)}{\pi_j^*} \right)$$

by

$$\sum_{i,j \in s} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_y^*(\beta) - y_j)}{\pi_j^*}$$

and $f_x(Q_x(\beta))$ and $f_y(Q_y(\beta))$ by following Silverman (1986).

The asymptotic variances of ratio, product and exponentiation estimators corresponding to $H(q, t^*) = q/t^*$, $H(q, t^*) = qt^*$ and $H(q, t^*) = q/(t^*)^\alpha$, respectively, can be derived.

Optimal estimators

In this section we derive the expression of the optimal estimator in the class $\widehat{Q}_{ye}^*(\beta)$. Again the optimality is defined in the sense of minimizing the (asymptotic) variance of these estimators.

This leads to the optimal value of α given by

$$\alpha_{opt} = \frac{Q_x(\beta) \text{Cov}(\widehat{Q}_y(\beta), \widehat{Q}_x(\beta)) - \text{Cov}(\widehat{Q}_y(\beta), \widehat{Q}'_x(\beta))}{Q_y(\beta) V(\widehat{Q}_x(\beta)) + \widehat{Q}'_x(\beta) - 2\text{Cov}(\widehat{Q}_x(\beta), \widehat{Q}'_x(\beta))}.$$

By using the properties of two-phase sampling the next expression can be obtained:

$$\alpha_{opt} = \frac{Q_x(\beta) f_x(Q_x(\beta)) E_{d1} \left(\sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_y(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*} \right)}{Q_y(\beta) f_y(Q_y(\beta)) E_{d1} \left(\sum_{i,j \in s'} \Delta_{ij}^{s'} \frac{\delta(Q_x(\beta) - x_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*} \right)},$$

and then

$$\widehat{Q}_y^{\alpha_{opt}}(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\alpha_{opt}}.$$

It can be easily seen:

$$\begin{aligned} V(\widehat{Q}_y^{\mathcal{H}}(\beta)) &\geq V(\widehat{Q}_y^{\alpha_{opt}}(\beta)) = V(\widehat{Q}_y(\beta)) - K_1 = \\ &= V(\widehat{Q}_y(\beta)) - \frac{(Cov(\widehat{Q}_y(\beta), \widehat{Q}_x(\beta)) - Cov(\widehat{Q}_y(\beta), \widehat{Q}'_x(\beta)))^2}{V(\widehat{Q}_x(\beta)) + \widehat{Q}'_x(\beta) - 2Cov(\widehat{Q}_x(\beta), \widehat{Q}'_x(\beta))}, \end{aligned} \quad (5.37)$$

that is, the lower bound of the variance of $\widehat{Q}_y^{\mathcal{H}}(\beta)$ is the variance of the exponentiation estimator with α_{opt} .

Equation (5.37) shows that the proposed estimator $\widehat{Q}_y^{\alpha_{opt}}(\beta)$ always remains more efficient than the simple estimator $\widehat{Q}_y(\beta)$. Specifically, K_1 is the amount by which the variance is reduced when we use the exponentiation estimator with an optimal α instead of the $\widehat{Q}_y(\beta)$ estimator.

In practice the optimal value of α is unknown. Nevertheless, the sample data can be used to calculate its estimator. Thus, an estimator of the optimal value of α is given by

$$\widehat{\alpha} = \frac{\widehat{Q}_x^*(\beta) f_x(Q_x(\beta))}{\widehat{Q}_y^*(\beta) f_y(Q_y(\beta))} \frac{\sum_{i,j \in s} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_y^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*}}{\sum_{i,j \in s} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(Q_x(\beta) - x_i)}{\pi_i^*} \frac{\delta(Q_x(\beta) - x_j)}{\pi_j^*}}. \quad (5.38)$$

We can define an optimal estimator of the β -quantile as:

$$\widehat{Q}_y^{\widehat{\alpha}}(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\widehat{\alpha}}.$$

Following the procedure discussed in Allen *et al.* (2002) it can be shown that $E(\widehat{Q}_y^{\widehat{\alpha}}(\beta)) = Q_y(\beta) + o(n^{-1})$ and to the first degree of approximation, $V(\widehat{Q}_y^{\widehat{\alpha}}(\beta)) = V(\widehat{Q}_y^{\alpha_{opt}}(\beta))$, i.e., the estimators $\widehat{Q}_y^{\widehat{\alpha}}(\beta)$ and $\widehat{Q}_y^{\alpha_{opt}}(\beta)$ are asymptotically equivalent.

5.3.4. Quantile estimation under two-phase sampling for stratification

In Section 5.3.2 we show a particular case of two-phase sampling where the first phase sample is stratified using the auxiliary variable. This sampling

design was called two-phase sampling for stratification. We now define an estimator for the quantile $Q_y(\beta)$ under this sampling design and analyze several of its properties.

In the first place, we define the following estimator for the distribution function:

$$\widehat{F}_{st}^*(t) = \frac{1}{N} \sum_{h=1}^L \sum_{i \in s_h} \frac{\delta(t - y_i)}{\pi_i^*},$$

and we suggest estimating the quantile $Q_y(\beta)$ by $\widehat{Q}_{st}^*(\beta) = \widehat{F}_{st}^{*-1}(\beta)$, where the inverse \widehat{F}_{st}^{*-1} exists in the same way as \widehat{F}_{HTy}^{-1} above.

To study the properties of the $\widehat{Q}_{st}^*(\beta)$ estimator, we will first analyze the properties of the $\widehat{F}_{st}^*(t)$ estimator.

Note that $\widehat{F}_{st}^*(t)$ is unbiased and its variance is given by:

$$\begin{aligned} V(\widehat{F}_{st}^*(t)) &= \frac{1}{N^2} \left(\sum_{i,j \in U} \Delta'_{ij} \frac{\delta(t - y_i)}{\pi'_i} \frac{\delta(t - y_j)}{\pi'_j} + \right. \\ &\quad \left. + E_{d1} \left[\sum_{h=1}^L \sum_{i,j \in s'_h} \Delta'_{ij} \frac{\delta(t - y_i)}{\pi_i^*} \frac{\delta(t - y_j)}{\pi_j^*} \right] \right). \end{aligned} \quad (5.39)$$

Thus, an unbiased estimator of variance is given by:

$$\begin{aligned} \widehat{V}(\widehat{F}_{st}^*(t)) &= \frac{1}{N^2} \left(\sum_{i,j \in s} \frac{\Delta'_{ij}}{\pi_{ij}^*} \frac{\delta(t - y_i)}{\pi'_i} \frac{\delta(t - y_j)}{\pi'_j} + \right. \\ &\quad \left. + \sum_{h=1}^L \sum_{i,j \in s_h} \frac{\Delta'_{ij}}{\pi_{ij}/s'} \frac{\delta(t - y_i)}{\pi_i^*} \frac{\delta(t - y_j)}{\pi_j^*} \right), \end{aligned} \quad (5.40)$$

because each component of (5.40) is unbiased for its counterpart in equation (5.39).

Similar to Section 5.3.2, the $\widehat{Q}_{st}^*(\beta)$ estimator can be expressed as a linear function of $\widehat{F}_{st}^*(Q_y(\beta))$. In addition, because $\widehat{F}_{st}^*(t)$ is unbiased of $F_y(t)$, we deduce that $\widehat{Q}_{st}^*(\beta)$ is asymptotically unbiased. An approximate unbiased estimator of the variance is given by:

$$\widehat{V}(\widehat{Q}_{st}^*(\beta)) = \frac{1}{N^2} \frac{1}{f_y^2(Q_y(\beta))} \left(\sum_{i,j \in s} \frac{\Delta'_{ij}}{\pi_{ij}^*} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_i)}{\pi'_i} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_j)}{\pi'_j} + \right.$$

Tabla 5.1: Description and references of populations.

Population	Variables	ρ_{yx}	References
Fam1500 (N=1500)	y :Feeding expenses		Fernández <i>et al.</i> (1994)
	x_1 :Family incomes	0.848	
	x_2 :Other expenses	0.546	
Counties (N=304)	y :Population in 1970		Royall <i>et al.</i> (1981)
	x_1 :Population in 1960	0.982	Valliant <i>et al.</i> (2000)
	x_2 :Households in 1960	0.982	

$$+ \sum_{h=1}^L \sum_{i,j \in s_h} \frac{\Delta_{ij}^{s'}}{\pi_{ij/s'}} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_i)}{\pi_i^*} \frac{\delta(\widehat{Q}_{st}^*(\beta) - y_j)}{\pi_j^*} \Bigg).$$

5.3.5. An empirical study

The present investigation proposes several estimators for quantiles in sampling in two phases with unequal probabilities. The use of two-phase sampling for stratification has also been considered for estimating quantiles. In this section we carried out a simulation study to reveal the behaviour of these estimators and to point out the most efficient estimator. For this purpose, we examined two natural populations, used previously for finite population sampling. The populations in question are Fam1500 and Counties. A brief description and the references of these populations can be seen in Table 5.1.

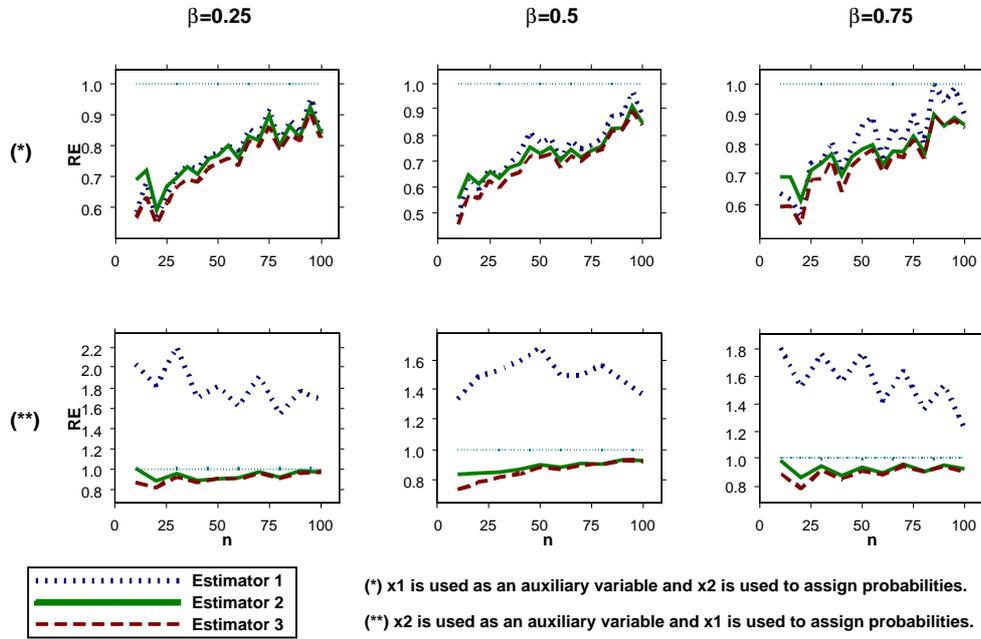
In these populations there are several auxiliary variables having different linear correlation coefficients with the variable of interest y . In this study the behaviour of estimators can be observed when strong and weak relationships between variables are considered.

We have generated 1000 independent samples under different methods in each phase. The first phase sample size, n' , is fixed at 150 and the second phase sample size, n , is allowed to change from 10 to 100. The methods used are:

1. (*SRSWOR.M*) The first phase is *SRSWOR* of size n' . The second phase is carried out using the Midzuno-Sen method (Singh, 2003, pg. 390) to extract samples with unequal probabilities:

$$\pi'_i = \frac{n'}{N}, \quad \pi_{i/s'} = \frac{n' - n}{n' - 1} \frac{x_i}{\sum_{j \in s'} x_j} + \frac{n - 1}{n' - 1} \quad \rightarrow \quad \pi_i^* = \pi'_i \pi_{i/s'}.$$

Figura 5.7: Relative Efficiency for Fam1500 population and under *SRSWOR.M* sampling design. $n' = 150$.



2. (*SRSWOR.P*) The first phase is *SRSWOR* of size n' . The second phase is carried out by Poisson sampling (Singh, 2003, pg. 499) such that the conditional inclusion probability is proportional to x :

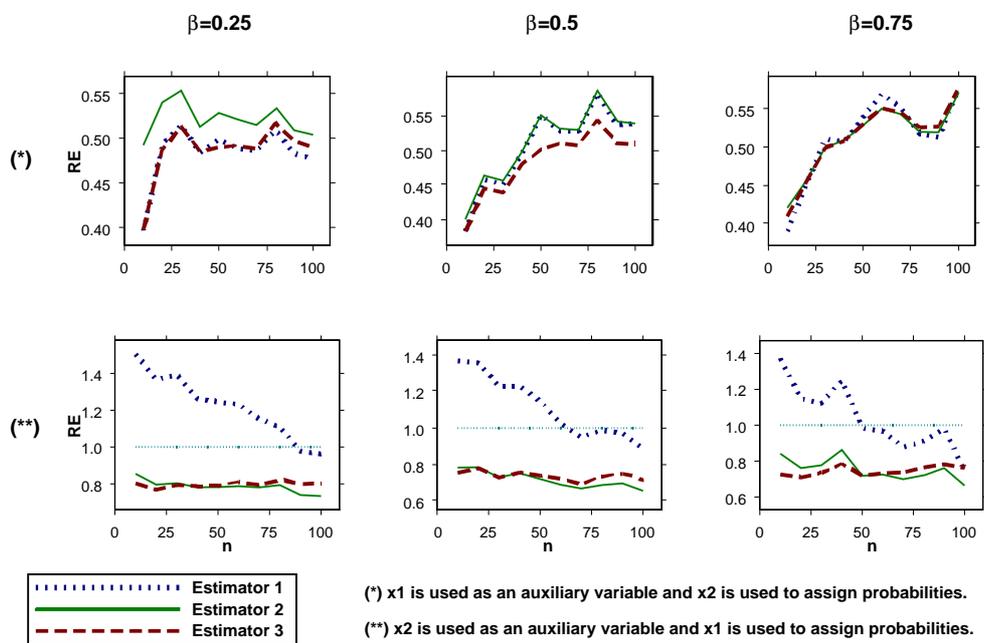
$$\pi'_i = \frac{n'}{N}, \quad \pi_{i/s'} = n \frac{x_i}{\sum_{j \in s'} x_j} \quad \rightarrow \quad \pi_i^* = \pi'_i \pi_{i/s'}.$$

3. (*ST.M*) Two-phase sampling for stratification: In the first phase, a sample is drawn according *SRSWOR*. For the elements in s' information is recorded that will permit a stratification. From stratum h , a sample s_h of size n_h is drawn with unequal probabilities using the Midzuno-Sen Method:

$$\pi'_i = \frac{n'}{N}, \quad \pi_{i/s'_h} = \frac{n'_h - n_h}{n'_h - 1} \frac{x_i}{\sum_{j \in s'_h} x_j} + \frac{n_h - 1}{n'_h - 1} \quad \rightarrow$$

$$\rightarrow \pi_i^* = \pi'_i \pi_{i/s'_h} \text{ for } i \in s'_h.$$

Figura 5.8: Relative Efficiency for Fam1500 population and under *SRSWOR.P* sampling design. $n' = 150$.



(*) x_1 is used as an auxiliary variable and x_2 is used to assign probabilities.

(**) x_2 is used as an auxiliary variable and x_1 is used to assign probabilities.

Figura 5.9: Relative Efficiency for Counties population and under *SRSWOR.M* sampling design. $n' = 150$.

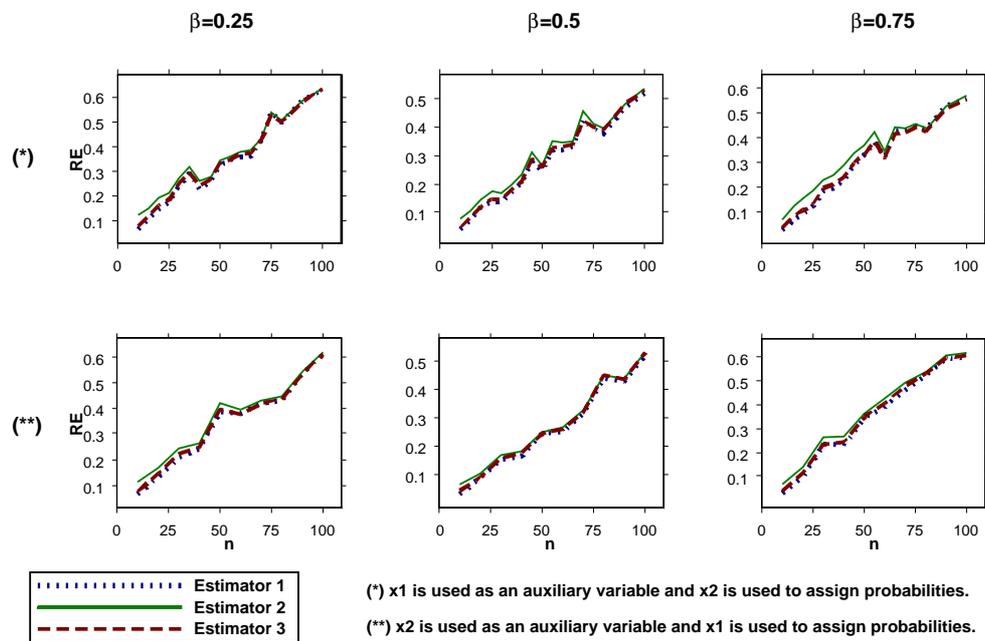
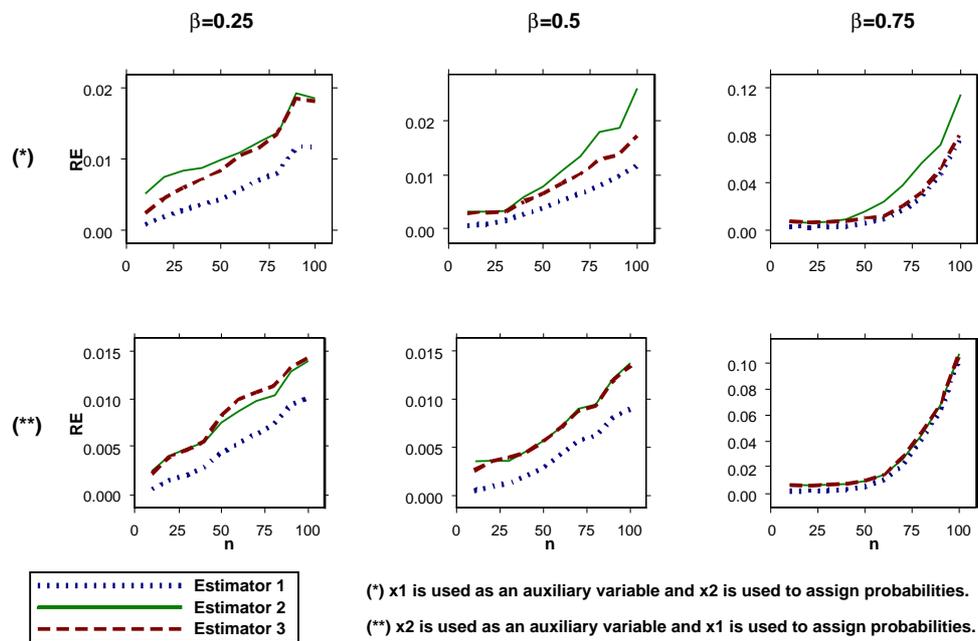


Figura 5.10: Relative Efficiency for Counties population and under *SRSWOR.P* sampling design. $n' = 150$.



The performance of the proposed estimators is evaluated for the three quartiles, $\beta = 0,25, 0,50, 0,75$, in terms of Relative Bias (%) (RB) and Relative Efficiency (RE) with Monte Carlo approximations derived from the $B = 1000$ independent samples:

$$RB_i = 100 \times \frac{1}{B} \sum_{b=1}^B \frac{\widehat{Q}_y^i(\beta)_b - Q_y(\beta)}{Q_y(\beta)} \quad ; \quad RE_i = \frac{MSE[\widehat{Q}_y^i(\beta)]}{MSE[\widehat{Q}_y^*(\beta)]},$$

where b indexes the b th simulation run and $\widehat{Q}_y^i(\beta)$ denotes the i th proposed estimator with

- $\widehat{Q}_y^1(\beta) = \widehat{Q}_y^*(\beta) \frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)},$
- $\widehat{Q}_y^2(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\widehat{\alpha}},$ where $\widehat{\alpha}$ can be seen in (5.38),
- $\widehat{Q}_y^3(\beta) = \widehat{Q}_y^*(\beta) \left(\frac{\widehat{Q}'_x(\beta)}{\widehat{Q}_x^*(\beta)} \right)^{\alpha_{opt}},$
- $\widehat{Q}_y^4(\beta) = \widehat{Q}_{st}^*(\beta).$

$MSE[\widehat{Q}_y^i(\beta)] = B^{-1} \sum_{b=1}^B [\widehat{Q}_y^i(\beta)_b - Q_y(\beta)]^2$ and $MSE[\widehat{Q}_y^*(\beta)]$ is similarly defined for $\widehat{Q}_y^*(\beta)$, the direct estimator defined in (5.33). This does not use the auxiliary information.

The random generations, calculations and all the estimators were obtained using the R program. Programming details are available in Appendix C.

Figures 5.7, ..., 5.10 represent the RE for $\widehat{Q}_y^1(\beta)$, $\widehat{Q}_y^2(\beta)$ and $\widehat{Q}_y^3(\beta)$ estimators in different populations and the $SRSWOR.M$ and $SRSWOR.P$ designs. These figures show the behaviour of the estimators when the sample size in the second phase increases, while the first phase sample size remains fixed.

If there is a high linear correlation coefficient between y and the auxiliary variable, then all estimators are more efficient than the $\widehat{Q}_y^*(\beta)$ estimator (shown with horizontal dotted lines). The gain in relative efficiency decreases if the sample size in the second phase increases. This result is logical because if the sample size in the second phase is small, the sample will have less information of the y variable, and the $\widehat{Q}_y^*(\beta)$ estimator will present a larger degree of error,

while the ratio estimators are more efficient because more information is used. As n increases, $\widehat{Q}_y^*(\beta)$ obtains better estimator which is closer to the ratio estimator. Note that for the Fam1500 population and under *SRSWOR.P* sampling design with the first phase sample size $n'=150$, as the second phase sample size n increases from 10 to 100, the RE shows two pecks: if $n=25$ and $n=80$ for $\beta=0.25$; if $n=55$ and $n=80$ for $\beta=0.5$; and if $n=60$ and $n=100$ for $\beta=0.75$. It looks that if we are estimating higher quartile then a large second phase sample size may be required so long as the efficiency of the proposed estimators is concerned.

$\widehat{Q}_y^3(\beta)$ is the most efficient estimator in many cases. This is expected because this estimator is asymptotically optimum in the class (5.35). Nevertheless, the estimator $\widehat{Q}_y^2(\beta)$ has very similar values and does not depend on unknown values. $\widehat{Q}_y^1(\beta)$ is usually less efficient than other proposed estimators. When the linear relation between the variables is weaker, $\widehat{Q}_y^1(\beta)$ is even less efficient than the direct estimator, while $\widehat{Q}_y^2(\beta)$ and $\widehat{Q}_y^3(\beta)$ continue to perform better. In short, the use of the exponentiation estimator improves the estimates, especially if there is a weak relationship between the study and auxiliary variables.

On the other hand, the Poisson method of sampling produces more efficient results in the sense of *RE* than the Midzuno-Sen method and with regard to $\widehat{Q}_y^*(\beta)$ because the direct estimator present disperses estimates under the Poisson method caused by the heterogeneity of the inclusion probabilities.

Proposed estimators are almost equivalents in the Counties population because the linear correlation coefficients are larger. In fact, the *RE* of the proposed estimators in this population is better then those in the Fam1500 population.

Bias is another important aspect, particularly for ratio estimator that can show the underestimation or overestimation. The *RB*'s values in the Fam1500 population are all within a reasonable range, with the $\widehat{Q}_y^*(\beta)$ having the largest at 3% as seen in Figure 5.11. The *RB*'s values for the Counties population when x_1 is used as an auxiliary variable and x_2 is used to assign probabilities are shown in Figure 5.12. The $\widehat{Q}_y^*(\beta)$ estimator clearly leads to serious overestimation, especially when the sample size is small and under the *SRSWOR.P* sampling design, whereas the absolute *RB*'s of the proposed estimators are less than 7% for the *SRSWOR.M* sampling design and less than 13% for the *SRSWOR.P* sampling design, except on small sample sizes for the $\widehat{Q}_y^2(\beta)$ estimator, which has the largest at 25%. In short, the study of the *RB* reveals

Figure 5.11: Relative Bias in percent for Fam1500 population when x_1 is used as an auxiliary variable and x_2 is used to assign probabilities. $n' = 150$.

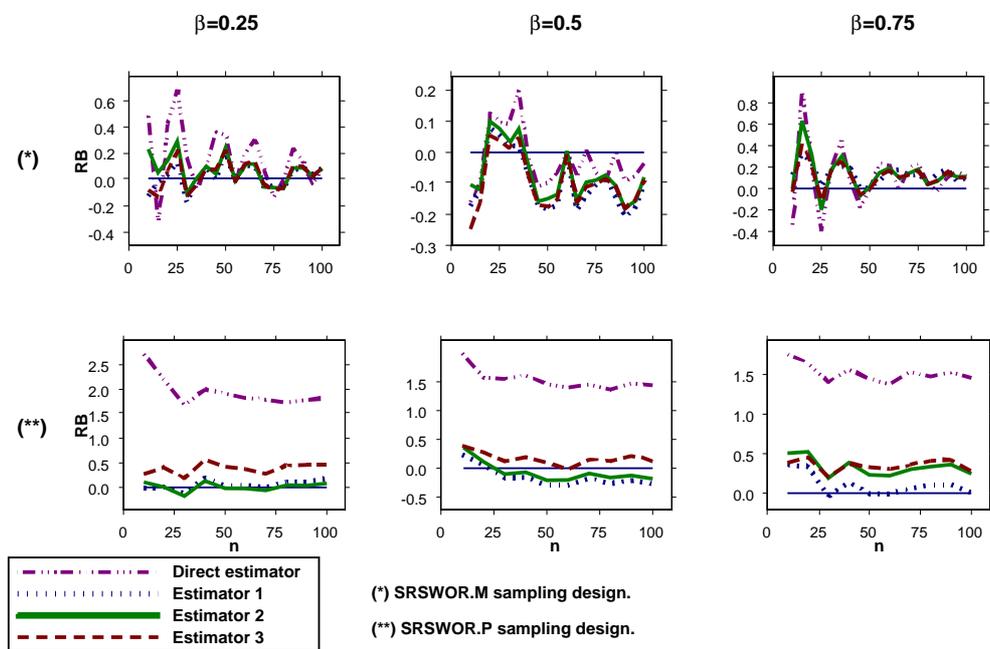


Figure 5.12: Relative Bias in percent for Counties population when x_1 is used as an auxiliary variable and x_2 is used to assign probabilities. The RB 's values for the direct estimator in (**) are larger than 97.6%, 74.6% and 21.5% for $\beta = 0,25, 0,5$ and $0,75$, respectively, and are omitted. $n' = 150$.

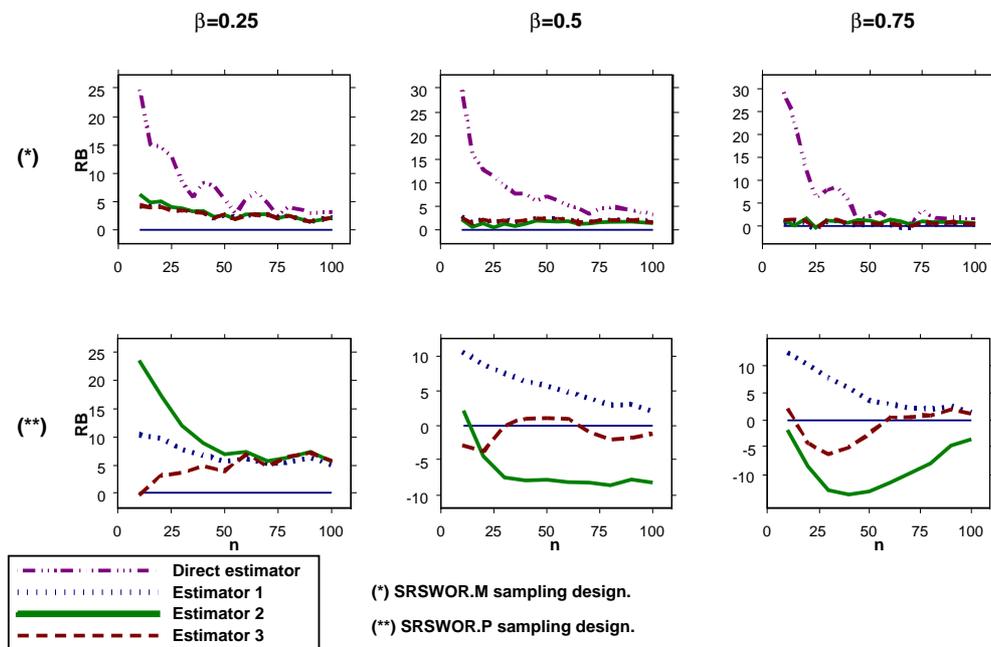
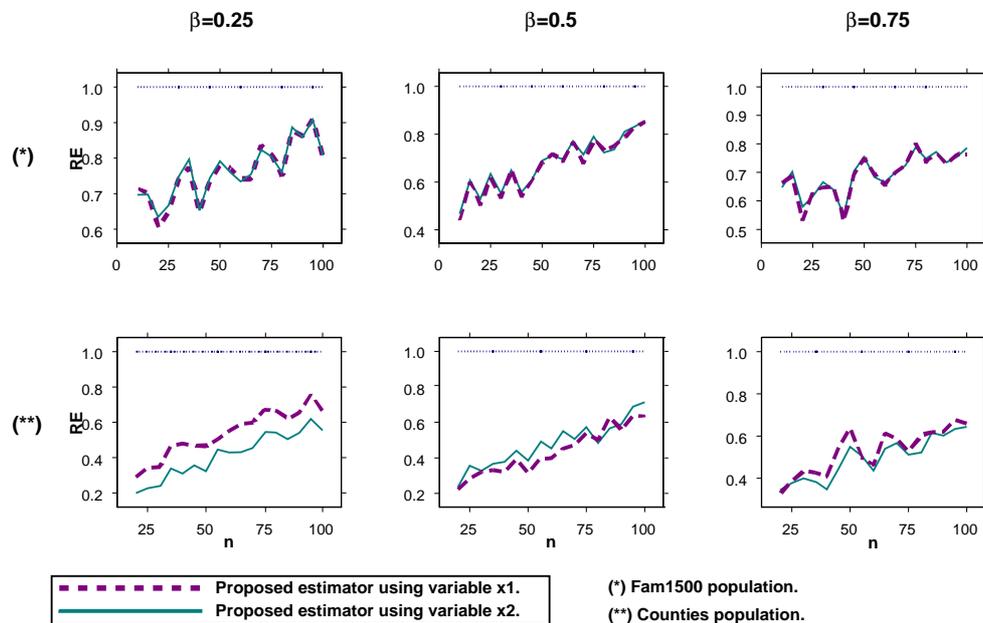


Figure 5.13: Relative Efficiency for Fam1500 and Counties populations and under ST.M sampling design. $n' = 150$.

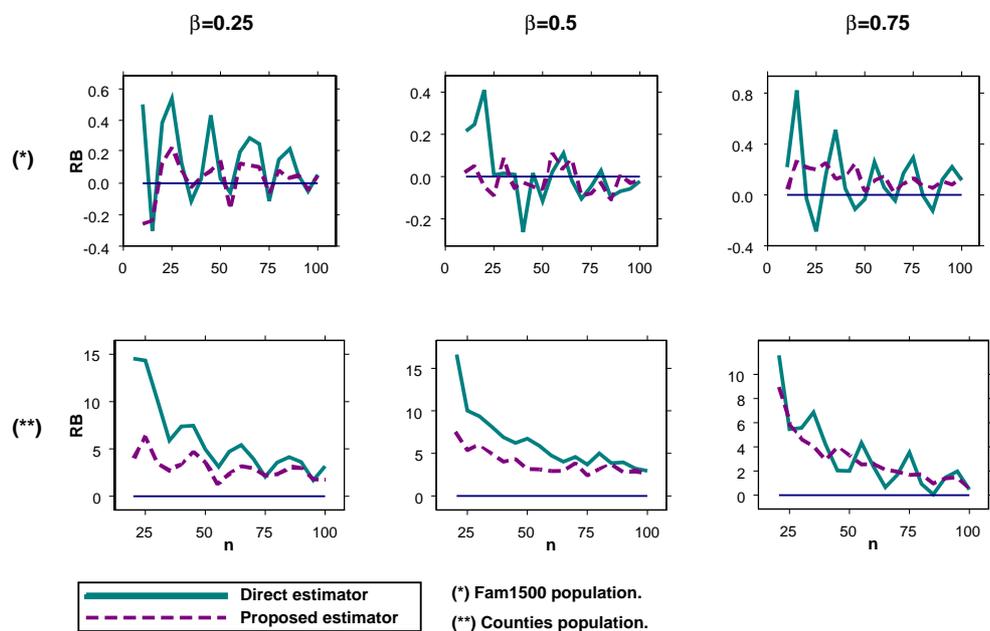


that the proposed estimators are better than the direct estimator.

Figure 5.13 is an example of two-phase sampling for stratification. The proposed estimator is compared with the direct estimator if the strata are not considered. It can also be observed that the use of stratification is recommended because the estimates are more efficient, especially if the sample size in the second phase of the sample decreases. In all cases the proposed estimators show improvement over the direct estimator irrespective of the linear relationship between variables, although the gain in RE is better if this coefficient is larger. In reality, the gain in efficiency is guaranteed because the strata are well designed, i.e., the strata are homogeneous inside and heterogeneous among them.

As far as the RB is concerned, the proposed estimator, $\widehat{Q}_{st}^*(\beta)$ is better than $\widehat{Q}_y^*(\beta)$ as can be observed in Figure 5.14. The RB 's values of $\widehat{Q}_{st}^*(\beta)$ are less than 10%, whereas the $\widehat{Q}_y^*(\beta)$ estimator leads to a weak overestimation for the Counties population. In fact, the Fam1500 population produces better estimates than the Counties population in terms of RB . The estimators are

Figura 5.14: Relative Bias in percent for Fam1500 and Counties populations under ST.M sampling design and when the variable x_1 is used. $n' = 150$.



showing similar behaviour when the variable x_2 is used and consequently these figures are not shown.

5.4. Conclusions

In Section 2.3 we use the pseudo empirical likelihood method to estimate a population mean when observations are missing either for the study variable or for the auxiliary variable. We assume that the sample is selected according to an unequal probability sampling design.

We propose a class of pseudo empirical maximum likelihood estimators, using the data from the units with non-missing (or missing) values for the study variable and missing (or non-missing) values for the auxiliary variables. We derive asymptotic properties of these estimators. We derive the optimal estimator of the proposed class as being the one with the smallest asymptotic variance.

The proposed estimator is compared with other estimators in a simulation study. The proposed optimum estimator has the smallest empirical variance when the number of missing values is high and the relationship between the auxiliary and principal variables is weak.

The practical advantage of the method proposed is its breadth of applicability. The proposed estimator may be extended in a number of ways. For example, we concentrate on the estimation of population means. However, the method proposed can be extended to other parameters such as ratios, variances and quantiles.

Next, in Section 2.4 the empirical likelihood method is used to provide a model-assisted estimator for the distribution function. The proposed estimator possesses a number of desirable features, namely that it can easily be used under unequal probability sampling designs, the estimator is not model dependent, the conditions for the existence of the estimator are established, the estimator is a genuine distribution function, the facts of asymptotic unbiasedness, asymptotic normality, availability of a variance estimator, computation simplicity, etc.

The accuracy of the proposed estimator has been compared in terms of various measures with other known estimators. These studies have shown the behaviour of the model-assisted pseudo empirical likelihood estimator to be optimal. While Chambers and Dunstan's estimator can be very efficient when

the model upon which it is based is appropriate, nevertheless, as noted by Rao *et al.* (1990), Chambers *et al.* (1993) and Dorfman (1993), this estimator can perform poorly under model misspecification. A similar remark can be made about the model-calibrated pseudo empirical likelihood estimator. This estimator also suffers an important loss of efficiency when the value t_0 used in the benchmark constraints is far away from t .

Another important property is the efficient use of the auxiliary information: on the one hand because several auxiliary variables can be used at the estimation stage and on the other because the good distribution of the set of points \mathbf{t}_g allow that more accuracy estimates will be possible for any t .

In conclusion, we suggest that in many standard survey settings, the model-assisted empirical likelihood method provides a simple and practical approach by incorporating multivariate auxiliary information into the estimation of the distribution function. The model-assisted empirical likelihood method exhibits good performance and this can be a valid alternative to other estimators of the distribution function.

The quantile estimation is discussed in Chapter 3. In the first place, new estimators are defined under successive sampling (Section 3.3). This is a known technique that can be used in longitudinal surveys to estimate population parameters and measurements of difference or change of a study variable. Economics and social surveys carried out by nacional agencies and other statistical offices use this sampling scheme. The quantile estimation is a common problem in most of these studies.

Assuming sampling on two occasions, an estimator for population quantiles is proposed when the samples are obtained by a general sampling design in each occasion. The most important properties of this estimator are also discussed. The ratio estimator proposed is very easy to compute and highly efficient.

When simple random sampling is used on both occasions, we obtain the asymptotic normality of the estimator, which is used to construct approximate confidence intervals for this parameter. Under unequal probability sampling, this expression becomes more complicated than in the preceding case, although it can be obtained by means of the variances and covariances of the Horvitz-Thompson estimators.

Let us assume that on the second occasion, two samples are drawn: a matched sample is selected from the first sample and an unmatched one is selected from the complemented first sample, independently of the matched portion (the matched and unmatched samples are conditionally independent).

The proposed estimator can also be defined if we use a sampling scheme in which the unmatched sample is selected from the population U , i.e. if the matched and unmatched samples are unconditionally independent. The properties of this new estimator can be studied in a similar manner, and this is simpler because the conditioning disappears.

The extension of the estimator to the case where we have three or more occasions is also straightforward.

Assuming successive sampling on two occasions, we also propose a class of estimators for quantiles using a multiple ratio estimator based on the matched samples on the first and the current occasions, selected using simple random sampling. On the proposed class of estimators, we derive the expression of the optimal estimator in the sense of minimizing the asymptotic variance. The proposed estimator possesses a number of desirable properties, such as asymptotic normality, availability of variance estimator, simplicity of computation, etc. In our theoretical and empirical studies, the new estimator appears to be more accuracy than the standard and univariate ratio estimators.

The proposed estimator can be generalized to a sampling survey performed on more than two occasions and to other sampling designs different to sample random sampling. These are areas for further researches.

Most of the procedures in survey sampling that use auxiliary information are based on estimators that require the use of known population variables, although this situation is unlikely in practice. A solution to this problem is the use of the two-phase sampling.

Assuming two-phase sampling, with arbitrary sampling design for the selection of units in each phase, ratio and exponential type estimators have been proposed. The unbiasedness of these estimators have been discussed and expressions for the variance estimation have been also established. These results have allowed us to find an optimum estimator into the exponential type estimator. Assuming several sampling design and simulation studies, we have observed that the proposed estimators can provide more accuracy estimates than other known estimators.

We have also proposed estimators for a β -quantile and its variance when a two-phase sampling for stratification with arbitrary sampling designs in each of the two phases is used. First, a distribution function estimator has been defined and we have seen that this holds many desirable properties. In addition, efficient estimators for a β -quantile and its variance have been defined. These estimators have been analyzed via a simulation study and some useful

gains in efficiency about other estimators have been revealed. Thus, theoretical and empirical justifications have demonstrated that the proposed quantile estimator and its variance estimation handle well the difficult aspects of this sampling design.

In conclusion, the two-phase sampling is a simple and practical technique that is widely recommended when population auxiliary information is not available. Assuming that the first sample is stratified, we have proposed estimators that exhibit a better performance with regard to estimators based on a non-stratified sampling. These results suggest that the proposed method is an attractive alternative for estimating quantiles when population auxiliary information is not available.

Finally, quantile estimators have been proposed by using both model-assisted approach and the pseudo empirical likelihood method. The application of these estimators to the estimation of several poverty measures is also discussed. We propose using the bootstrap technique for the proposed estimators in the variance estimation stage. The accuracy of these new procedures have been analyzed in simulation studies for the problem of the quantile estimation and for some poverty measures used by several statistical agencies. Results show that the proposed estimators can provide more efficient estimates than other known procedures.

Bibliografía

- [1] **Adhvaryu, D.** (1978) Successive sampling using multi-auxiliary information. *Sankhya* **40**, 167-173.
- [2] **Aitchison, J. y Silvey, S.D.** (1958) Maximum-likelihood estimation of parameter subject to restraints. *Ann. Math. Statist.* **29**, 813-888.
- [3] **Allen, J., Singh, H.P., Singh, S. y Smarandache, F.** (2002) A general class of estimators of population median using two auxiliary variables in double sampling. INTERSTAT.
- [4] ¹ **Arcos, A., Rueda, M. y Muñoz, J.F.**(2006) An improved class of estimators of a finite population quantile in sample surveys. *Applied Mathematics Letters*. En prensa.
- [5] **Arnab, R. y Okafor, F.C.** (1992) A note on double sampling over two occasions. *Pakistan Journal of Statistics* **8**, 9-18.
- [6] **Artés Rodríguez, E.M. y García Luengo A.V.** (2002) *Diseños muestrales en el tiempo*. Monografías, Universidad de Almería.
- [7] **Bahadur, R.R** (1966) A note on quantiles in large samples. *Annals of Mathematical Statistics* **37**, 577-580.
- [8] **Basu, D.** (1971) *Foundations of statistical inference. A Symposium*, eds. V.P. Godambe and D. A. Sprott, Toronto: Holt Rinehart and Winston.
- [9] **Berger, Y.G.** (2004) Variance estimation for measures of change in probability sampling. *The Canadian Journal of Statistics* **32**, 451-467.
- [10] ¹ **Berger, Y.G., Muñoz, J.F. y Rancourt, E.** (2006) Variance estimation of regression estimators when control total are estimated: an application to the composite estimator. *Survey Methodology*. En revisión.

¹Bibliografía correspondiente al doctorando.

- [11] **Berger, Y.G. y Rao, J.N.K.** (2006) Adjusted jackknife for imputation under unequal probability sampling without replacement. *Journal of the Royal Statistical Society Series B* **68**, 531-547.
- [12] **Berger, Y.G. y Skinner, C.J.** (2003) Variance estimation for a low income proportion. *Journal of the Royal Statistical Society, Series C* **52**, 457-468.
- [13] **Berger, Y.G. y Skinner, C.J.** (2005) A jackknife variance estimator for unequal probability sampling. *Journal of the Royal Statistical Society Series B* **67**, 79-89.
- [14] **Bethlehem, J.G. y Keller, W.J.** (1987) Linear weighting of sample survey data. *Journal of Official Statistics* **3**, 141-153.
- [15] **Bickel, P.J. y Freedman, D.A.** (1984) Asymptotic normality and the bootstrap in stratified sampling. *The Annals of Statistics* **12**, 470-482.
- [16] **Binder, D.A.** (1982) Non-parametric bayesian models for samples from finite population. *Journal of the Royal Statistical Society, Series B* **44** (3), 388-393.
- [17] **Binder, D.A.** (1983) On the variances of asymptotically normal estimator from complex surveys. *International Statistical Review* **51**, 279-292.
- [18] **Binder, D.A. y Kovačević** (1995) Estimating some measures of income inequality from survey data: an application of the estimating equation approach. *Survey Methodology* **21**, 137-145.
- [19] **Blackburn, M.** (1990) Trends in poverty in the United States, 1967-84. *Review of Income and Wealth* **36**, 53-66.
- [20] **Blackburn, M.** (1994) International comparisons of poverty. *American Economic Review* **84**, 371-374.
- [21] **Brewer, K.R.W.** (1999) Cosmetic calibration with unequal probability sampling. *Survey Methodology* **25**, 205-212.
- [22] **Brewer, K.R.W., Early, L.J. y Joyce, S.F.** (1972) Selecting several samples from a single population. *Australian Journal of Statistics* **14**, 231-239.
- [23] **Casell, C.M., Särndal, C.E. y Wretman, J.H.** (1976) Some results on generalized difference estimation and generalized regression estimation for finite populations. *Biometrika* **63**, 615-620.

- [24] **Casell, C.M., Särndal, C.E. y Wretman, J.H.** (1977) *Foundations of Inference in Survey Sampling*. New York: Wiley.
- [25] **Chambers, R.L., Dorfman, A.H. y Hall, P.** (1992) Properties of estimator of the finite population distribution function. *Biometrika* **79**, 577-582.
- [26] **Chambers, R.L., Dorfman, A.H. y Wehrly, T.E.** (1993) Bias robust estimation in finite population using nonparametric calibration. *Journal of the American Statistical Association* **88**, 268-277.
- [27] **Chambers, R.L. y Dunstan, R.** (1986) Estimating distribution functions from survey data. *Biometrika* **73**, 597-604.
- [28] **Chaudhuri, A. y Vos, J.W.E.** (1988) *Unified theory and strategies of survey sampling*. North-Holland, Amsterdam.
- [29] **Chen, H. y Chen, J.** (2000) Bahadur representations of the empirical likelihood quantile processes. *Nonparametric Statist.* **12**, 645-660.
- [30] **Chen, J. y Qin, J.** (1993) Empirical likelihood estimation for finite populations and the effective usage of auxiliary information. *Biometrika* **80**, 107-116.
- [31] **Chen, J., Rao, J.N.K. y Sitter, R.R.** (2000) Efficient random imputation for missing data in complex surveys. *Statistica Sinica* **10**, 1153-1169.
- [32] **Chen, J. y Sitter, R.R.** (1999) A pseudo empirical likelihood approach to the effective use of auxiliary information in complex surveys. *Statistica Sinica* **9**, 385-406.
- [33] **Chen, J., Sitter, R.R. y Wu, C.** (2002) Using empirical likelihood methods to obtain range restricted weights in regression estimators for surveys. *Biometrika* **89**, 230-237.
- [34] **Chen, J., Thompson, M.E. y Wu, C.** (2004) Estimation of fish abundance indices based on scientific research trawl surveys. *Biometrics* **60**, 116-123.
- [35] **Chen, J. y Wu, C.** (2002) Estimation of distribution function and quantiles using the model-calibrated pseudo empirical likelihood method. *Statistica Sinica* **12**, 1223-1239.
- [36] **Cochran, W.G.** (1977) *Sampling Techniques*. 3rd ed. New York: Wiley

- [37] **Cramer, H.** (1946) *Mathematical methods of statistics*. Princeton University Press. Princeton.
- [38] **Dalgleish, L. I.** (1995) Software review: Bootstrapping and jackknifing with BOJA. *Statistics and Computing* **5**, 165-174.
- [39] **Das, A.K. y Tripathi, T.P.** (1978) Use of auxiliary information in estimating the finite population variance. *Sankhya, Series C* **40**, 139-148.
- [40] **Datta, G.S. y Ghosh, M.** (1993) Bayesian estimation of finite population variances with auxiliary information. *Sankhya, Series B* **55**, 156-170.
- [41] **Deng, L.Y. y Wu, C.F.J.** (1987) Estimation of variance of the regression estimator. *Journal of the American Statistical Association* **82**, 568-576.
- [42] **Deville, J.C.** (1999) Variance estimation for complex statistics and estimators: linearization and residual techniques. *Survey Methodology* **25**, 193-203.
- [43] **Deville, J.C. y Särndal, C.E.** (1992) Calibration estimators in survey sampling. *Journal of the American Statistical Association* **87**, 376-382.
- [44] **Dickens, R. y Manning, A.** (2004) Has the national minimum wage reduced UK wage inequality?. *Journal of the Royal Statistical Society, Series A* **167**, 613-626.
- [45] **Dorfman, A.H.** (1993). A comparison of design-based and model-based estimators of the finite population distribution function. *The Australian Journal of Statistics* **35**, 29-41.
- [46] **Dorfman, A.H. y Hall, P.** (1993) Estimators of the finite population distribution function using nonparametric regression. *The Annals of Statistics* **21** (3), 1452-1475.
- [47] **Eckler, A.R.** (1955) Rotation Sampling. *The Annals of Mathematical Statistics* **26** 664-685.
- [48] **Efron, B. y Tibshirani, R.J.** (1993) *An introduction to the Bootstrap*. Chapman & Hall, London.
- [49] **Eurostat.** (2000) Low-wage employees in EU countries. Statistics in Focus: Population and Social Conditions. Theme 3 – 11/2000. *Office for Official Publications of the EC*, Luxemburgo.

- [50] **Fernández García, F.R. y Mayor Gallego, J.A.** (1994) *Muestreo en Poblaciones Finitas: Curso Básico*. P.P.U., Barcelona.
- [51] **Fernández Sánchez, M.P., Hernández Bastida, A. y Sánchez González, C.** (2004) Análisis de los ingresos y gastos trimestrales de los hogares españoles usando verosimilitud empírica. *Estudios de Economía Aplicada* **22**, 139-150.
- [52] **Francisco, C.A. y Fuller, W.A.** (1991) Quantiles estimation with a complex survey design. *The Annals of Statistics* **19**, 454-469.
- [53] **Fuller, W.A., Loughin, M.M. y Baker, H.D.** (1994) Regression weighting in the presence of nonresponse with application to the 1987-1988 nationwide food consumption survey. *Survey Methodology* **20**, 75-85.
- [54] **Godambe, V.P.** (1955) A unified theory of sampling from finite populations. *Journal of the Royal Statistical Society, Series B* **17**, 269-278.
- [55] **Godambe, V.P. y Thompson, M.E.** (1973) Estimation in sampling theory with exchangeable prior distributions. *The Annals of Statistics* **1**, 1212-1221.
- [56] **Godambe, V.P. y Thompson, M.E.** (1986) Parameters of superpopulation and survey population: Their relationships and estimation. *International Statistical Review* **54**, 127-138.
- [57] **Gordon, L.** (1983) Successive sampling in finite populations. *The Annals of Statistics* **11**, 702-706.
- [58] **Gross, S.T.** (1980) Median estimation in sample survey. *Proc. Surv. Res. Meth. Sect. Amer. Statist. Ass.* 181-184.
- [59] **Hájek, J.** (1959) Optimum strategies and other problems in probability sampling. *Casopis Pest. Mat.* **84**, 387-423.
- [60] **Hájek, J.** (1964) Asymptotic theory of rejective sampling with varying probabilities from a finite population. *Annals of Mathematical Statistics* **35**, 1491-1523.
- [61] **Hall, P.** (1990) Pseudo-likelihood theory for empirical likelihood. *The Annals of Statistics* **18**, 121-140.
- [62] **Hall, P. y La Scala, B.** (1990) Methodology and algorithms of empirical likelihood. *International Statistical Review* **58**, 109-127.

- [63] **Hansen, M.H. y Hurwitz, W.N.** (1943) On the theory of sampling from finite populations. *Annals of Mathematical Statistics* **14**, 333-362.
- [64] **Hanurav, T.V.** (1966) Some aspects of unified sampling theory. *Sankhya, Series A* **28**, 175-204.
- [65] **Hartley, H.O. y Rao, J.N.K.** (1968) A new estimation theory for sample surveys. *Biometrika* **55**, 547-557.
- [66] **Hartley, H.O. y Ross, A.** (1954) Unbiased ratio estimators. *Nature* **174**, 270-271.
- [67] **Hedayat, A.S. y Sinha, B.K.** (1991) *Design and Inference in Finite Population Sampling*. John Wiley and Sons.
- [68] **Hill, B.M.** (1968) Posterior distribution of percentiles: Bayes theorem for sampling from a population. *Journal of the American Statistical Association* **63**, 677-691.
- [69] **Horvitz, D.G. y Thompson, D.J.** (1952) A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association* **47**, 663-685.
- [70] **Huang, E.T. y Fuller, W.A.** (1978) Nonnegative regression estimation for sample survey data. *In Proc. Social Statistics Sec., Am. Statist. Assoc., 300-305* Washington, D.C: American Statistical Association.
- [71] **Instituto Nacional de Estadística.** (1992) Encuesta Continua de Presupuestos Familiares. Metodología. *Instituto Nacional de Estadística. Madrid.*
- [72] **Isaki, C.T.** (1983) Variance estimation using auxiliary information. *Journal of the American Statistical Association* **78**, 117-123.
- [73] **Isaki, C.T. y Fuller, W.A.** (1982) Survey design under the regression superpopulation model. *Journal of the American Statistical Association* **77**, 89-96.
- [74] **Jagers, P.** (1986) Post-stratification against bias in sampling. *International Statistical Review* **54**, 159-167.
- [75] **Jessen, R.J.** (1942) Statistical investigation of a sample survey for obtaining farm facts. *Iowa Agricultural Experiment Statistical Research Bulletin*, 304.

- [76] **Jonhson, P. y Webb, S.** (1992) Official statistics on poverty in the United Kingdom. Poverty measurement for economies in transition in eastern european countries. Polish Statistical Association and Polish Central Statistica Office, Warsaw. *Journal of Economics Perspectives* **15**, 143-156.
- [77] **Koenker, R. y Hallock, K.F.** (2001) Quantile regression. *Journal of Economics Perspectives* **15**, 143-156.
- [78] **Kovačevik, M.S. y Binder, D. A.** (1997) Variance estimation for measures of income inequality and polarization - The estimating equations approach. *Journal of Official Statistics* **13**, 41-58.
- [79] **Kovačevik, M.S. y Yung, W.** (1997) Variance estimation for measures of income inequality and polarization - an empirical study. *Survey Methodology* **23**, 41-52.
- [80] **Kovačevik, M.S., Yung, W. y Pandher** (1995) Estimating the sampling variances of measures of income inequality and polarization - an empirical study. *Statistic Canada, Methodology Branch Working Paper*, HSMD-95-007E.
- [81] **Kovar, J.G., Rao, J.N.K. y Wu, C.F.J.** (1988) Bootstrap and other methods to measure errors in survey estimates. *The Canadian Journal of Statistics* **16**, 25-45.
- [82] **Kuk, A.Y.C.** (1993) A kernel method for estimating finite population distribution functions using auxiliary information. *Biometrika* **80**, 385-392.
- [83] **Kuk, A.Y.C. y Mak, T.K.** (1989) Median estimation in the presence of auxiliary information. *Journal of the Royal Statistical Society, Series B* **51**, 261-269.
- [84] **Kuk, A.Y.C. y Mak, T.K.** (1994) A functional approach to estimating finite population distribution functions. *Theory Meth.* **23 (3)**, 883-896.
- [85] **Kuo, L.** (1988) Classical and Prediction Approaches to Estimating Distribution Functions from Survey Data. Proceeding of the Section on Survey Research Methods. *American Statistical Association*, 280-285.
- [86] **Lahiri, D.B.** (1951) A method of sample selection providing unbiased ratio estimates. *Bulletin of the International Statistical Institute* **33**, 133-140.

- [87] **Leung, D.H.Y. y Qin, J.** (2006) Analysing survey data with incomplete responses by using a method based on empirical likelihood. *Journal of the Royal Statistical Society, Series C* **55**, 379-396.
- [88] **Little, R.J.A. y Rubin, D.B.** (1987) *Statistical analysis with missing data*. John Wiley, New York.
- [89] **Liu, T.P. y Thompson, M.E.** (1983) Properties of estimators of quadratic finite population functions: the batch approach. *The Annals of Statistics* **11**, 275-285.
- [90] **Lombardía, M. J., González-Manteiga, W. y Prada-Sánchez, J.M.** (2003) Bootstrapping the Chambers-Dunstan estimate of a finite population distribution function. *J. Statistical Planning and Inference* **116**, 367-388.
- [91] **Lombardía, M. J., González-Manteiga W., y Prada-Sánchez, J.M.** (2004) Bootstrapping the Dorfman-Hall-Chambers-Dunstan estimator of a finite population distribution function. *Journal of Nonparametric Statistics* **16**, 63-90.
- [92] **Lucifora, C. y Salverda, W.** (1998) *Policies for low wage employment and social exclusion*. Ed. FrancoAngeli.
- [93] **Mak, T.K. y Kuk, A.Y.C.** (1993) A new method for estimating finite-population quantiles using auxiliary information. *The Canadian Journal of Statistics* **25**, 29-38.
- [94] **Midzuno, H.** (1952) On the sampling system with probability proportional to sum of sizes. *Annals of Institute of Statistical Mathematics* **3**, 99-107.
- [95] **Molina, C.E.A. y Skinner, C.J.** (1992) Pseudo-likelihood and Quasi-likelihood estimation for complex sampling schemes. *Computational Statistics and Data Analysis* **13**, 395-405.
- [96] **Mukhopadhyay, P.** (2000) *Topics in Survey Sampling* Springer.
- [97] ¹ **Muñoz, J.F., Sánchez, I. y Álvarez, E.** (2006) Revisión de las principales técnicas de estimación de parámetros lineales en subpoblaciones. *Metodología de Encuestas*. En prensa
- [98] **Murthy, M.N.** (1967) *Sampling theory and method*. Calcutta: Statistical Publishing Society.

¹Bibliografía correspondiente al doctorando.

- [99] **Narain, R.D.** (1953) On the recurrence formula in sampling on successive occasions. *Journal of the Indian Society of Agricultural Statistics* **5**, 96-99.
- [100] **OECD** (1982) The OECD list of social indicators, Paris.
- [101] **OECD** (1997) Labour market policies: new challenges policies for low-paid workers and unskilled job seekers. *OECD Working Papers. vol 5*, n° 86 .
- [102] **Ogus, J.K. y Clark, D.F.** (1971) The annual survey of manufacturers: A report on methodology. Technical Report No. 2, U.S. Bureau of Census, Washington D.C.
- [103] **Olkin, I.** (1958) Multivariate ratio estimation for finite population. *Biometrika* **45**, 154-165.
- [104] **Owen, A.B.** (1988) Empirical likelihood ratio confidence intervals for a single functional. *Biometrika* **75**, 237-249.
- [105] **Owen, A.B.** (1990) Empirical likelihood confidence regions. *The Annals of Statistics* **18**, 90-120.
- [106] **Owen, A.B.** (1991) Empirical likelihood for linear models. *The Annals of Statistics* **19**, 1725-1747.
- [107] **Owen, A.B.** (2001) *Empirical likelihood*. Chapman y Hall/CRC.
- [108] **Patterson, H.D.** (1950) Sampling on successive occasions with partial replacement of units. *Journal of the Royal Statistical Society, Series B* **12**, 241-255.
- [109] **Pérez, R.A.** (2002) ¿Qué es un modelo de superpoblación?. *Metodología de Encuestas* **4 (1)**, 79-86.
- [110] **Pfefferman, D. y Krieger, A.M.** (1991) Poststratification using regression estimators when information on strata means and sizes is missing. *Biometrika* **78**, 409-419.
- [111] **Polyak, B.T.** (1987) *Introduction to Optimization*. New York: Optimization Software, Inc. Publications Division.
- [112] **Prasad, N.G.N. y Thach, T.** (2001) Variance estimation under two-phase sampling. *Working paper, Department of Mathematical Sciences, University of Alberta*.

- [113] **Preston, I.** (1995) Sampling distributions of relative poverty statistics. *Journal of the Royal Statistical Society, Series C* **44**, 91-99.
- [114] **Qin, J. y Lawless, J.F.** (1994) Empirical likelihood and general estimating equations. *The Annals of Statistics* **22**, 300-325.
- [115] **Qin, J. y Lawless, J.F.** (1995) Estimating equations, empirical likelihood and constraints on parameters. *The Canadian Journal of Statistics* **23**, 145.
- [116] **Randles, R.H.** (1982) On the asymptotic normality of statistics with estimated parameters. *The Annals of Statistics* **10**, 462-474.
- [117] **Rao, J.N.K.** (1966) Alternative estimators in PPS sampling for multiple characteristics. *Sankhya Series A* **28**, 47-60.
- [118] **Rao, J.N.K.** (1994) Estimating totals and distribution functions using auxiliary information at the estimation stage. *Journal of Official Statistics* **10**, 153-165.
- [119] **Rao, J.N.K., Hartley, H.O. y Cochran, W.G.** (1962) A simple procedure of unequal probability sampling without replacement. *Journal of the Royal Statistical Society, Series B* **24**, 482-491.
- [120] **Rao, J.N.K., Kovar, J.G. y Mantel, H.J.** (1990) On estimating distribution functions and quantiles from survey data using auxiliary information. *Biometrika* **77**, 365-375.
- [121] **Rao, J.N.K. y Liu, J.** (1992) On estimating distribution functions from sample survey data using supplementary information at the estimation stage. *In Nonparametric Statistics and Related Topics, Ed A.K. Md. E. Saleh*, 399-407 New York:Elsevier.
- [122] **Rao, J.N.K. y Shao, A.J.** (1966) Jackknife variance estimation with survey data under hotdeck imputation. *Biometrika* **53**, 811-822.
- [123] **Rao, J.N.K. y Singh, A.C.** (1997) A ridge-shrinkage method for range-restricted weight calibration in survey sampling. *In Proc. Survey Res. Meth. Sect., Am. Statist. Assoc.*, 57-65 Washington, D.C: American Statistical Association.
- [124] **Rao, C.R. y Toutenburg, H.** (1995) *Linear Models: Least Squares and Alternatives*. Springer, New York.

- [125] **Rao, J.N.K. y Wu, C.F.J.** (1985) Inference from stratified samples: Second-order analysis of three methods for nonlinear statistics. *Journal of the American Statistical Association* **80**, 620-630.
- [126] **Robinson, J.** (1987) Conditioning ratio estimates under simple random sampling. *Journal of the American Statistical Association* **82**, 826-831.
- [127] **Royall, R.M. y Cumberland, W.G.** (1981) An empirical study of the ratio estimator and estimator of its variance. *Journal of the American Statistical Association* **76**, 66-88.
- [128] **Rubin, D.B.** (1987) *Multiple imputation for nonresponse in sample surveys*. Wiley, New York.
- [129] **Rueda, M. y Arcos, A.** (2001) On estimating the median from survey data using multiple auxiliary information. *Metrika* **4**, 161-173.
- [130] **Rueda, M. y Arcos, A.** (2002a) The use of quantiles of auxiliary variables to estimate medians. *Biometrical Journal* **44** (5), 619-632.
- [131] **Rueda, M. y Arcos, A.** (2002b). Estimación por intervalos de la mediana con estimadores de razón y diferencia. *Estudios de Economía Aplicada* **20**, 241-260.
- [132] **Rueda, M., Arcos, A. y Artés, E.** (1997) Improvement on Estimating Quantiles in Finite Population Using Indirect Methods of Estimation. *Lecture Notes in Computer Science* **1280**, 491-500.
- [133] **Rueda, M., Arcos, A. y Artés, E.** (1998) Quantile Interval Estimation in Finite Population using a Multivariate Ratio Estimator. *Metrika* **47**, 203-213.
- [134] **Rueda, M., Arcos, A. y Martínez-Miranda, M.D.** (2003) Difference estimators of quantiles in finite populations. *Test* **12**, 481-496.
- [135] **Rueda, M., Arcos, A., Martínez-Miranda, M.D. y Román, Y.** (2004) Some improved estimators of finite population quantile using auxiliary information in sample surveys. *Computational Statistics and Data Analysis* **45**, 825-848.
- [136] ¹ **Rueda, M., Arcos, A., Muñoz, J.F. y Singh, S.** (2006a) Quantile estimation in two-phase sampling. *Computational Statistics and Data Analysis*. En prensa.

¹Bibliografía correspondiente al doctorando.

- [137] **Rueda, M. y González, S.** (2004) Missing data and auxiliary information in surveys. *Computational Statistic* **19**, 551-567.
- [138] **Rueda, M. y Martínez, S.** (2002) Estimadores de calibración: Una nueva metodología para el uso de la información auxiliar. *Metodología de Encuestas*, **4**, 161-173.
- [139] ¹ **Rueda, M., Martínez, S., Arcos, A., Martínez, H. y Muñoz, J.F.**(2006c) Mean estimation under successive sampling with calibration estimators. *Australian and New Zealand Journal of Statistics*. En revisión.
- [140] ¹ **Rueda, M. y Muñoz, J.F.** (2005) Una revisión del método de verosimilitud empírica en las encuestas por muestreo. *Investigación Operacional* **26**, 225-237.
- [141] ¹ **Rueda, M. y Muñoz, J.F.**(2006a) A model-assisted estimator for the distribution function using the pseudo empirical likelihood method. *Statistics and Computing*. En revisión
- [142] ¹ **Rueda, M. y Muñoz, J.F.** (2006b) Estimating quantiles under sampling in two occasions with unequal probabilities. *Computational Statistics and Data Analysis*. Aceptado bajo revisión.
- [143] ¹ **Rueda, M. y Muñoz, J.F.** (2006c) Estimating quantiles under two-phase sampling for stratification. *Statistics and Probability Letters*. En revisión.
- [144] ¹ **Rueda, M. y Muñoz, J.F.** (2006d) Model-assisted estimation of quantiles using empirical likelihood. Applications to different poverty measures. *Journal of the Royal Statistical Society, Series C*. En revisión.
- [145] ¹ **Rueda, M., Muñoz, J.F., y Arcos, A.** (2006) Estimating quantiles under sampling on two occasions with P auxiliary variables. *Quality and Quantity*. En prensa.
- [146] ¹ **Rueda, M., Muñoz, J.F., Berger, Y.G., Arcos, A. y Martínez, S.**(2006b) Pseudo empirical likelihood method in the presence of missing data. *Metrika*. En prensa.
- [147] **Ruspini, E.** (1999) Longitudinal research and the analysis of social change. *Quality and Quantity* **33**, 219-227.

¹Bibliografía correspondiente al doctorando.

- [148] **Sánchez-Crespo, G.** (2002) Introducción a los modelos de superpoblación en las técnicas de muestreo con probabilidades desiguales. *Metodología de Encuestas* **4 (1)**, 87-104.
- [149] **Särndal, C.E.** (1980) On π -inverse weighting versus best linear weighting in probability sampling. *Biometrika* **67**, 639-650.
- [150] **Särndal, C.E.** (1990) Methods for estimating the precision of survey estimates when imputation has been used. Proceedings of Symposium 1990: Measurement and improvement of data quality, Ottawa, 337-347.
- [151] **Särndal, C.E.** (1992) Methods for estimating the precision of survey estimates when imputation has been used. *Survey Methodology* **18**, 241-252.
- [152] **Särndal, C.E., Swensson, B. y Wretman, J.H.** (1989) The weighted technique for estimating the variance of the general regression estimator of the finite population total. *Biometrika* **76**, 527-537.
- [153] **Särndal, C.E., Swensson, B. y Wretman, J.H.** (1992) *Model Assisted Survey Sampling*. Springer-Verlag, New York
- [154] **Sedransk, J. y Meyer, J.** (1978) Confidence Intervals for the quantiles of a finite populations: simple random and stratified simple random sampling. *Journal of the Royal Statistical Society, Series B* **40**, No2, 239-252.
- [155] **Sedransk, J. y Smith, P.J.** (1988) Inference for finite population quantiles. In: Krishnaiah, P.R. and Rao, C. R. (eds.) *Handbook of Statistics* **6**, Cap11, 267-289. North-Holland.
- [156] **Sen, A.R.** (1972) Successive sampling with p ($p \geq 1$) auxiliary variables. *The Annals of Mathematical Statistics* **43 (6)**, 2031-2034.
- [157] **Sen, A.R.** (1973) Some theory of sampling on successive occasions. *The Australian Journal of Statistics* **15 (2)**, 105-110.
- [158] **Serfling, R.J.** (1980) *Approximation Theorems of Mathematical Statistics* John Wiley, New York.
- [159] **Shah, D.N. y Patel, P.A.** (1996) Asymptotic properties of a generalized regression-type predictor of a finite population variance in probability sampling. *The Canadian Journal of Statistics* **24**, 373-384.

- [160] **Shao, J.** (1994) L-statistics in complex survey problems. *The Annals of Statistics* **22**, 946-967.
- [161] **Shao, J. y Chen, Y.** (1998) Bootstrapping sample quantiles based on complex survey data under hot deck imputation. *Statistica Sinica* **8**, 1071-1085.
- [162] **Shao, J. y Rao, J.N.K.** (1993) Standard errors for low income proportions estimated from stratified multi-stage samples. *Sankhya Series B* **55**, 393-414.
- [163] **Shao, J. y Tu, D.** (1995) *The Jackknife and Bootstrap*. Springer-Verlag, New York.
- [164] **Shao, J. y Wu, C.F.J.** (1989) A general theory for jackknife variance estimation. *The Annals of Statistics* **17**, 1176-1197.
- [165] **Shao, J. y Wu, C.F.J.** (1992) Asymptotic properties of the balanced repeated replication method for sample quantiles. *The Annals of Statistics* **20**, 1571-1593.
- [166] **Shi, X., Wu, C.F.J. y Chen, J.** (1990) Weak and strong representations for quantile processes from finite populations with application to simulation size in resampling inference. *The Canadian Journal of Statistics* **18**, 141-148.
- [167] **Shorack, G.R. y Wellner, J.A.** (1986) *Empirical Processes with Applications to Statistics*. Wiley, New York.
- [168] **Silva, P.L.D. y Skinner, C.J.** (1995) Estimating distribution functions with auxiliary information using poststratification. *Journal of Official Statistics* **11 (3)**, 277-294.
- [169] **Silverman, B.W.** (1986) *Density estimation for statistics and data analysis*. Chapman and Hall.
- [170] **Singh, S.** (2003) *Advanced sampling theory with applications: How Michael Selected Amy*, Kluwer Academic Publishers, The Netherlands.
- [171] **Singh, S., Horn, S. y Yu, F.** (1998) Estimation of variance of general regression estimator: high level calibration approach. *Survey Methodology* **24**, 41-50.
- [172] **Singh, S., Joarder, A.H. y Tracy, D.S.** (2001) Median estimation using double sampling. *Aust. N. Z. J. Stat.* **43**, 33-46.

- [173] **Singh, A.C. y Mohl, C.A.** (1996) Understanding calibration estimators in survey sampling. *Survey Methodology* **22**, 107-115.
- [174] **Singh, H.P., Singh, H.P. y Singh, V.P.** (1992) A generalized efficient class of estimators of population mean in two phase and successive sampling. *Inter. J. Mgmt. Syst.* **8 (2)**, 173-183.
- [175] **Singh, S. y Srivastava, A.K.** (1973) Use of auxiliary information in two stage successive sampling. *Journal of Indian Society of Agricultural Statistic* **25**, 101-104.
- [176] **Sitter, R.R y Wu, C.F.J.** (2001) A note on Woodruff confidence intervals for quantiles. *Statistics and Probability Letters* **52**, 353-358.
- [177] **Sitter, R.R y Wu, C.F.J.** (2002) Efficient estimation of quadratic finite population functions in the presence of auxiliary information. *Journal of the American Statistical Association* **97**, 535-543.
- [178] **Smeeding, T.M.** (1991) Cross-national comparisons of inequality and poverty position. In: Osberg, L. (Ed.), *Economic Inequality and Poverty: International Perspectives*, M.E. Sharpe, Inc., Armonk.
- [179] **Solga, H.** (2001) Longitudinal surveys and the study of occupational mobility: Panel and retrospective design in comparison. *Quality and Quantity* **35**, 291-309.
- [180] **Sukhatme, P.V., Sukhatme, B.V., Sukhatme, S. y Asok, C.** (1984) *Sampling Theory of Surveys with Applications*. Iowa State University Press. Iowa.
- [181] **Swamy, P.A.V.B., Tavlas, G.S. y Chang, I.L.** (2005) How stable are monetary policy rules: estimating the time-varying coefficient in monetary policy reaction function for the U.S. *Computational Statistics and Data Analysis* **49**, 575-590.
- [182] **Théberge, A.** (1999) Extensions of calibration estimators in survey sampling. *Journal of the American Statistical Association* **94**, 635-644.
- [183] **Toutenburg, H. y Srivastava, V.K.** (1998) Estimation of ratio of population means in survey sampling when some observations are missing. *Metrika* **48**, 177-187.
- [184] **Toutenburg, H. y Srivastava, V.K.** (1999) Amputation versus imputation of missing values through ratio method in sample surveys. Unpublished document.

- [185] **Toutenburg, H. y Srivastava, V.K.** (2000) Efficient estimation of population mean using incomplete survey data on study and auxiliary characteristic. Unpublished document.
- [186] **Tracy, D.S. y Osahan, S.S.** (1994) Random nonresponse on study variable versus on study as well as auxiliary variables. *Statistica* **54**, 163-168.
- [187] **Valliant, R., Dorfman, A.H. y Royall, R.M.** (2000) *Finite population sampling and inference: A prediction approach*. Wiley Series in Probability and Statistics, Survey Methodology Section. New York. John Wiley and Sons, Inc.
- [188] **Wang, S. y Dorfman, A.H.** (1996) A new estimator for the finite population distribution function. *Biometrika* **83**, 639-652.
- [189] **Wolfson, M. y Evans, J.M.** (1989) Statistics Canada's low income cut-offs: metodological concerns and possibilities - a discussion paper. Research Paper Series, Statistical Canada, Ottawa. distribution function. *Biometrika* **83**, 639-652.
- [190] **Wolter, K.M.** (1985) *Introduction to Variance Estimation*. Springer-Verlag.
- [191] **Woodruff, R.S.** (1952) Confidence intervals for medians and other position measures. *Journal of the American Statistical Association* **47**, 635-646.
- [192] **Wu, C.F.J.** (1982) Estimation of variance of the ratio estimator. *Biometrika* **69**, 183-189.
- [193] **Wu, C.** (2002) Empirical likelihood method for finite populations. *Recent Advances in Statistical Methods*, Y.P. Chaubey, Ed., Imperial College Press, London, 339-351.
- [194] **Wu, C.** (2003) Optimal calibration estimators in survey sampling. *Biometrika* **90**, 937-951.
- [195] **Wu, C.** (2004a) Weighted empirical likelihood inference. *Statistics and Probability Letters* **66/1**, 67-79.
- [196] **Wu, C.** (2004b) Some algorithmic aspects of the empirical likelihood method in survey sampling. *Statistica Sinica* **14**, 1057-1067.

- [197] **Wu, C.** (2004c) Combining information from multiple surveys through empirical likelihood method. *The Canadian Journal of Statistics* **32**, 15-26.
- [198] **Wu, C.** (2005) Algorithms and R Codes for the Pseudo Empirical Likelihood Method in Survey Sampling. *Survey Methodology*, **31**, 239-243.
- [199] **Wu, C. y Luan, Y .** (2003) Optimal calibration estimators under two-phase sampling. *Journal of Official Statistics* **19**, 119-131.
- [200] **Wu, C. y Sitter, R.R.** (2001a) A model-calibration approach to using complete auxiliary information from survey data. *Journal of the American Statistical Association* **96**, 185-193.
- [201] **Wu, C. y Sitter, R.R.** (2001b) Variance estimation for the finite population distribution function with complete auxiliary information. *The Canadian Journal of Statistics* **29**, 289-307.
- [202] **Yates, F. y Grundy, P.M.** (1953) Selection without replacement form within strata with probability proportional to size. *Journal of the Royal Statistical Society, Series B* **15**, 253-261.
- [203] **Zhang, B.** (1996) Estimating a population variance with known mean. *International Statistical Review* **64**, 215-229.
- [204] **Zheng, B.** (2001) Statistical inference for poverty measures with relative poverty lines. *Journal of Econometrics* **101**, 337-356.
- [205] **Zhong, C.X.B., Chen, J. y Rao, J.N.K.** (2000) Empirical likelihood inference in the presence of measurement error. *The Canadian Journal of Statistics* **28**, 841.
- [206] **Zhong, C.X.B. y Rao, J.N.K.** (1996) Empirical likelihood inference for finite populations with auxiliary information using stratified random sampling. *Proceeding of the Section on Survey Research Methods, Am. Statist. Assoc.*, 793-803. Washington, DC: American Statistical Association.
- [207] **Zhong, C.X.B. y Rao, J.N.K.** (2000) Empirical likelihood inference under stratified random sampling using auxiliary information. *Biometrika* **87**, 929-938.

Apéndice A

Descripción de poblaciones finitas

En este apéndice se detallan las distintas poblaciones que han sido usadas en este trabajo con objeto de estudiar el comportamiento de los estimadores propuestos y su precisión con respecto a otros estimadores existentes en la literatura. Notamos que las poblaciones basadas en datos reales han sido utilizadas por otros autores en diferentes estudios de simulación, siendo estas poblaciones fiables y apropiadas para el estudio del comportamiento de estimadores en muestreo de poblaciones finitas. Las poblaciones que han sido simuladas siguen los modelos propuestos por otros autores, o bien, se han simulado de manera que pueda ser posible la extracción de muestras en los diseños muestrales más complejos que han sido tratados en este trabajo. De esta forma, se dispone de una estructura de datos apropiada para la obtención de tanto los estimadores propuestos como del resto de estimadores existentes en la literatura.

A.1. Poblaciones naturales

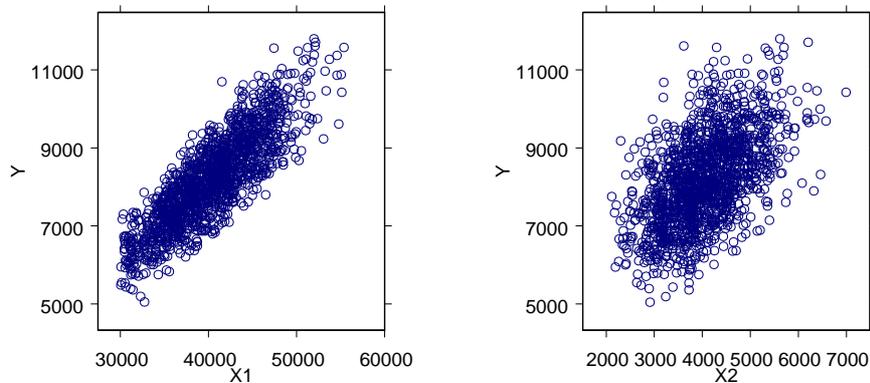
A.1.1. Fam1500

Esta población consta de $N = 1500$ familias de Andalucía y fue usada por primera vez por Fernández y Mayor (1994). Numerosos estudios posteriores (por ejemplo, Rueda *et al.*, 2006a, 2006b, Rueda y González, 2004, etc.) han usado esta población en sus estudios de simulación. La característica de interés,

Tabla A.1: Análisis descriptivo para las variables de la población Fam1500

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	5045	7358	8136	8181.94	8941	11795	0.14	
x_1	30052	36660	40200	40283.96	43700	55379	0.12	0.848
x_2	2116	3515	4001	4044.40	4538	6990	0.19	0.546

Figura A.1: Diagramas de dispersión de la población Fam1500



y , son los gastos de alimentación, mientras que las variables auxiliares x_1 y x_2 son, respectivamente, los ingresos familiares y otros gastos. En la Tabla A.1 puede consultarse información adicional sobre las variables de la población Fam1500, mientras que la Figura A.1 muestra los diagramas de dispersión correspondientes a dichas variables.

A.1.2. Counties

Las poblaciones Counties60 y Counties70 son poblaciones habitualmente usadas en muestreo de poblaciones finitas. Fueron usadas por primera vez en Royall y Cumberland (1981). Posteriormente, se ha usado en numerosos trabajos, como por ejemplo en Valliant *et al.* (2000). La población Counties60 consta de $N = 304$ ciudades de Carolina del Norte, Carolina del Sur y Georgia con menos de 100000 hogares en el año 1960. La variable y es la población

Tabla A.2: Análisis descriptivo para las variables de la población Counties60

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	1876	9787	18330	32916	38690	266623	1.24	
x	482	2502	4886	8931	10410	76887	1.30	0.998

Tabla A.3: Análisis descriptivo para las variables de la población Counties70

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	1924	9613	19080	36984	42560	409644	1.38	
x_1	482	2502	4886	8931	10410	76887	1.30	0.982
x_2	1876	9787	18330	32916	38690	266623	1.24	0.982

de cada ciudad, excluyendo los barrios de grupos de residentes. Como variable auxiliar, x , se tiene el número de hogares en 1960.

Por otro lado, la población Counties70 está formada por la variable de interés y que denota la población de 304 ciudades de Carolina del Norte, Carolina del Sur y Georgia con menos de 100000 hogares en el año 1970, excluyendo los barrios de grupos de residentes y por las variables auxiliares x_1 y x_2 , que coinciden con las variables x e y , respectivamente, de la población anterior.

Los datos de esta población pueden descargarse de:

ftp://ftp.wiley.com/public/sci_tech_med/finite_populations

Además, un breve resumen descriptivo de estas poblaciones puede consultarse en las Tablas A.2 y A.3. La Figura A.2 nos da los diagramas de dispersión entre las distintas variables de estas poblaciones. Puede observarse que estas poblaciones exhiben una mejor relación lineal entre las variables que la población Fam1500, lo que nos ha permitido comprobar en los distintos estudios el grado de ganancia en precisión en función de una mayor o menor relación lineal entre la variable principal y las auxiliares.

Figura A.2: Diagramas de dispersión de las poblaciones Counties70 y Counties60.

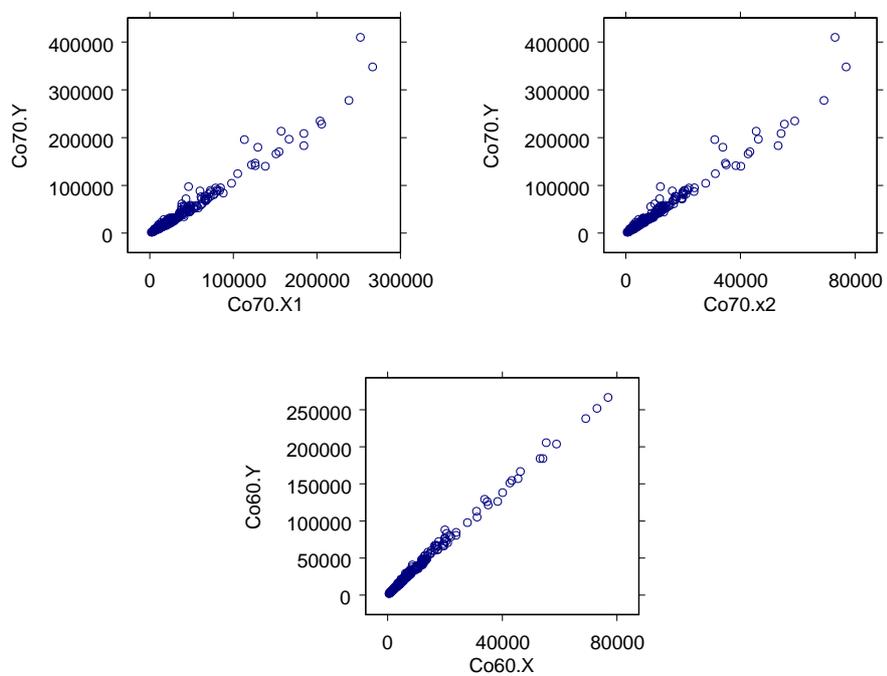
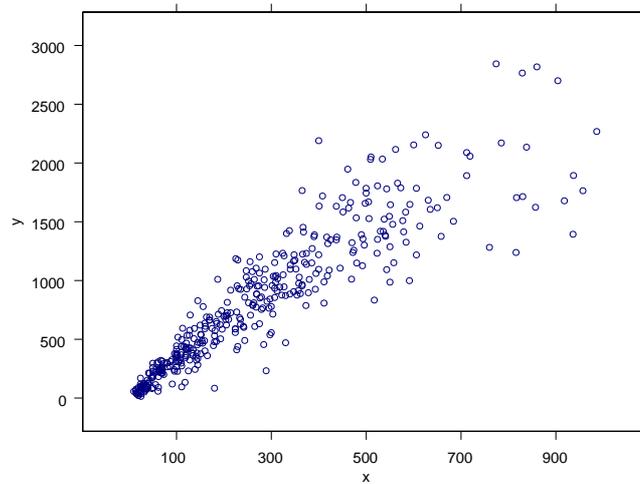


Tabla A.4: Análisis descriptivo para las variables de la población Hospitals

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	14	311	713	814.65	1186	2844	0.72	
x	1	102	233	274.70	393	986	0.78	0.911

Figura A.3: Diagrama de dispersión de la población Hospitals.



A.1.3. Hospitals

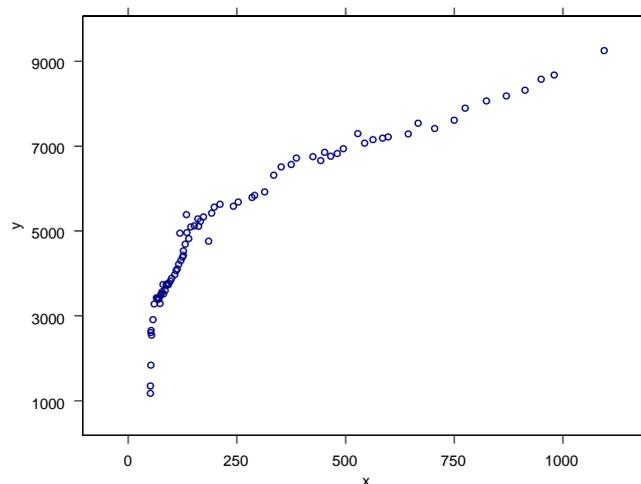
Esta población es una muestra nacional de hospitales en Estados Unidos. Esta muestra también fue considerada como una población en los estudios llevados a cabo por Royall y Cumberland (1981) y Valliant *et al.* (2000). El tamaño poblacional es de $N = 393$ hospitales de corta estancia con menos de 1000 camas, la variable de interés, y , es el número de pacientes dados de alta, mientras que la variable auxiliar es el número de camas que dispone el hospital.

El resumen descriptivo de las variables de esta población puede consultarse en la Tabla A.4. El diagrama de dispersión dado por la Figura A.3 nos permite profundizar en la estructura que presentan los datos de las variables de la población Hospitals.

Tabla A.5: Análisis descriptivo para las variables de la población Murthy

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	1176	3727.0	5105	5183.0	6754.0	9250	0.35	
x_1	51	86.5	148	285.1	445.3	1095	0.94	0.915

Figura A.4: Diagrama de dispersión de la población Murthy.



A.1.4. Murthy

La población Murthy es apropiada para observar el efecto de una mala especificación de un modelo de superpoblación en los estimadores, y poder proporcionar, por tanto, una indicación de la robustez de tales estimadores. Esta población consta de 80 fábricas donde la variable de interés, y , es la producción, y como variable auxiliar, x , se ha considerado el número de trabajadores. Esta población se usó previamente en Murthy (1967), Kuk y Mak (1989) y Kuk y Mak (1994).

En la Figura A.4 puede comprobarse que una hipótesis de linealidad no sería válida para las variables de esta población. Un estudio más exhaustivo sobre las características de las variables de la población Murthy puede obtenerse a partir de la Tabla A.5.

Tabla A.6: Análisis descriptivo para las variables de la población Turismos

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	11	343.3	894.0	3967.8	2483.5	308738	4.23	
x_1	5	73.0	176.5	810.2	464.0	61176	4.41	0.994
x_2	4	101.0	263.0	1313.7	749.3	111977	4.55	0.998
x_3	1	123.0	338.0	1373.1	957.5	102710	4.04	0.998
x_4	0	22.0	61.0	295.9	174.8	24023	4.26	0.961

A.1.5. Turismos

Esta población se ha obtenido a partir del número de turismos recogidos en los años 2002 y 2003 por el Instituto de Estadística de Andalucía en los distintos municipios de Andalucía. Estos datos pueden descargarse en la página web del Instituto de Estadística de Andalucía:

<http://www.juntadeandalucia.es/institutodeestadistica>

Por tanto, La población Turismos está formada por el número de turismos en $N = 770$ municipios de Andalucía. La variable principal, y , es el número de turismos por municipio en el año 2003. Se dispone de cuatro variables auxiliares: x_1 , x_2 , x_3 y x_4 que corresponden al número de turismos en el año 2002 con capacidad cilíndrica de clase 1, 2, 3 y 4, respectivamente.

El objetivo que tiene el uso de esta población es comprobar la ganancia en eficiencia de las estimaciones cuando se aumenta de manera paulatina el número de variables auxiliares.

En el análisis descriptivo de la Tabla A.6 se muestran las características más importantes de las variables de la población Turismos. En estas variables destaca la presencia de una alta asimetría y una importante variabilidad en los datos, como reflejan los correspondientes coeficientes de variación. Los diagramas de dispersión asociados a estas variables están disponibles en la Figura A.5.

Figura A.5: Diagramas de dispersión de la población Turismos.

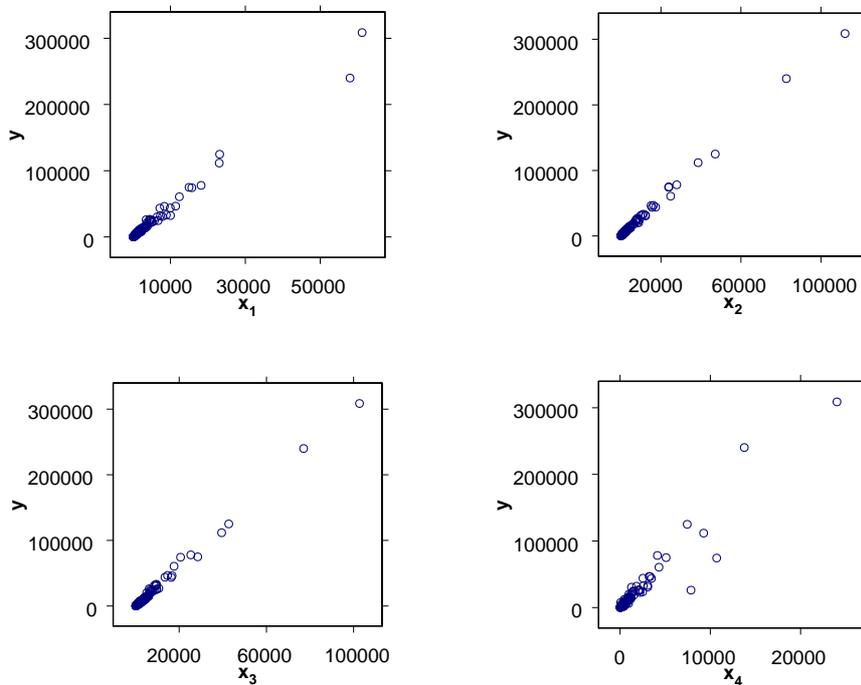


Tabla A.7: Análisis descriptivo para las variables de la población ECPF1997

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	240.4	2745	4037	4660	5842	61320	0.67	
x	107.6	2609	3845	4527	5654	27730	0.66	0.594

A.1.6. ECPF1997

La última población natural que se ha considerado en este trabajo se corresponde con los datos muestrales procedentes del primer trimestre del año 1997 de la Encuesta Continua de Presupuestos Familiares (*ECPF*). Véase Instituto Nacional de Estadística (1992) para una consulta detallada de la metodología. Esta población ha sido también analizada en Fernández *et al.* (2004).

Notamos que el objetivo de esta encuesta es proporcionar estimaciones acerca de los gastos de consumo y de los ingresos para el conjunto nacional, según varias variables de clasificación. La población consta de $N = 3000$ hogares españoles, donde se ha considerado que la variable de interés, y , son los ingresos totales trimestrales por hogar (en euros), mientras que los gastos trimestrales por hogar (en euros) será la variable auxiliar.

El correspondiente análisis descriptivo de las variables de esta población está dado por la Tabla A.7. Observamos que en este caso no existe una fuerte relación lineal entre la variable principal y la auxiliar. Este hecho es frecuente entre datos correspondientes a variables tales como ingresos o gastos, donde la alta presencia de valores extremos habitualmente dificulta la interpretación de algunas medidas como la media.

En cualquier caso, el objetivo al usar esta población es comprobar el comportamiento real de distintos estimadores en situaciones donde no pueda aceptarse una fuerte relación lineal entre las variables. En la Figura A.6 se muestra el correspondiente diagrama de dispersión.

Figura A.6: Diagrama de dispersión de la población ECPF1997.

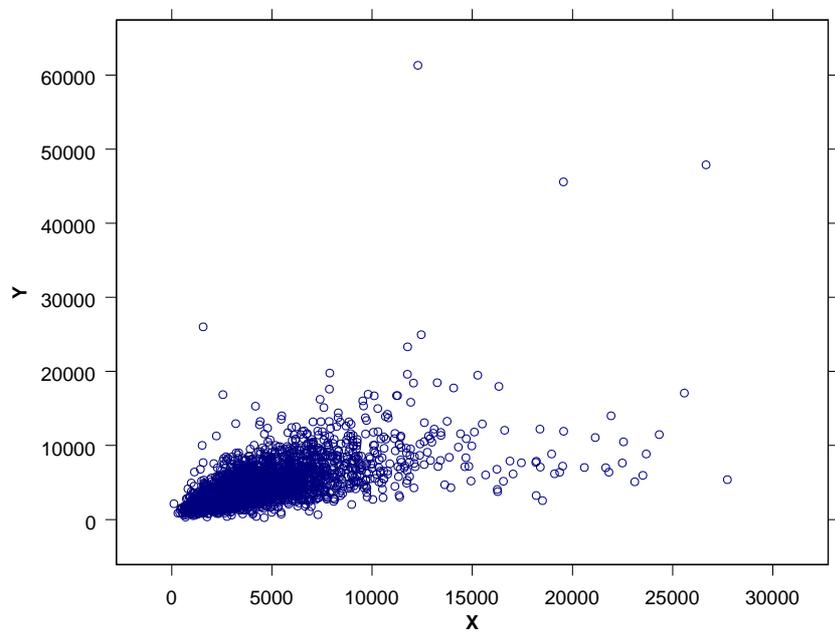


Tabla A.8: Análisis descriptivo para las variables de la población Pop06

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	-2.4588	0.87	1.93	1.98	2.96	9.33	0.81	
x	0.0008	0.27	0.66	0.96	1.32	8.10	1.03	0.6

Tabla A.9: Análisis descriptivo para las variables de la población Pop07

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	-2.349	1.02	1.88	2.00	2.86	10.03	0.71	
x	0.001	0.30	0.70	0.99	1.36	8.22	0.98	0.7

A.2. Poblaciones simuladas

A.2.1. Pop06, Pop07, Pop08 y Pop09

Paralelamente a Wu y Sitter (2001a), se han generado cuatro poblaciones de $N = 2000$ unidades mediante muestras independientes e idénticamente distribuidas mediante el modelo

$$y = \theta_0 + \theta_1 x + \epsilon, \quad (\text{A.1})$$

donde $x \sim \text{Gamma}(1, 1)$, $\epsilon \sim N(0, \sigma^2)$ y $\theta_0 = \theta_1 = 1$. Estas poblaciones se han generado escogiendo diferentes valores de σ^2 , de modo que los coeficientes de correlación entre y y x están dados por 0.6, 0.7, 0.8 y 0.9. Las poblaciones se han llamado Pop06, Pop07, Pop08 y Pop09, respectivamente. La Figura A.7 muestra los diagramas de dispersión de estas poblaciones, mientras que los distintos estudios descriptivos están dados por las Tablas A.8, A.9, A.10 y A.11.

Tabla A.10: Análisis descriptivo para las variables de la población Pop08

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	-2.243	1.15	1.81	1.99	2.63	8.54	0.64	
x	0.001	0.25	0.67	0.98	1.34	7.36	1.04	0.8

Tabla A.11: Análisis descriptivo para las variables de la población Pop09

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	-0.374	1.23	1.73	1.96	2.43	11.80	0.57	
x	0.002	0.29	0.67	0.98	1.33	10.51	1.02	0.9

Figura A.7: Diagramas de dispersión de las poblaciones Pop06, Pop07, Pop08 y Pop09

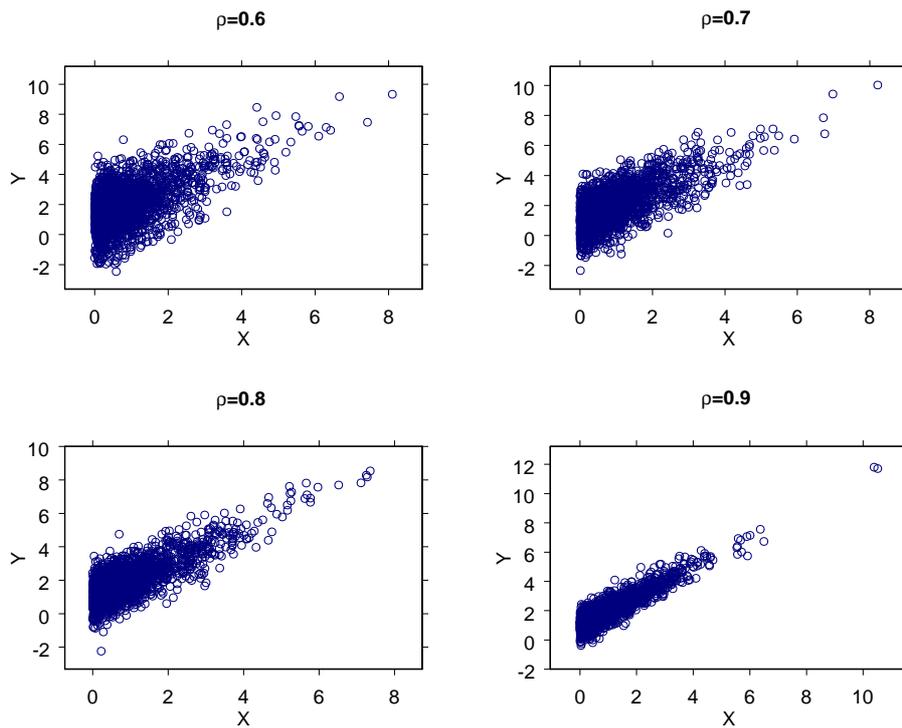


Tabla A.12: Análisis descriptivo para las variables de la población Pob098

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	-0.207	5.07	7.33	7.99	9.97	25.65	0.52	
x_1	0.003	0.90	2.26	3.08	4.37	22.32	0.96	0.71
x_2	0.081	1.80	3.17	3.85	5.34	17.55	0.72	0.67
\hat{y}	1.615	4.97	7.23	7.93	10.03	25.08	0.51	0.98

Tabla A.13: Análisis descriptivo para las variables de la población Pob080

V.	Min	Q_1	Me	Media	Q_3	Max	Cv	ρ_{yx}
y	-0.097	6.61	8.69	8.89	11.00	19.98	0.37	
x_1	0.480	2.46	3.67	3.98	5.15	11.86	0.50	0.60
x_2	0.417	2.54	3.59	3.89	5.00	12.20	0.48	0.53
\hat{y}	3.316	6.88	8.65	8.87	10.47	20.84	0.30	0.80

A.2.2. Pob098 y Pob080

Por último, se han generado dos poblaciones (Pob098 y Pob080) de tamaño $N = 1000$ mediante el modelo

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i, \quad (\text{A.2})$$

donde $\beta_0 = \beta_1 = \beta_2 = 1$ y las variables x_{1i} y x_{2i} se han generado de distribuciones Gamma con parámetros de forma y escala dados por 4 y 1, respectivamente. Las cantidades ϵ_i son variables aleatorias independientes e idénticamente distribuidas con distribución Normal de parámetros 0 y σ^2 . El valor de σ^2 se ha seleccionado de modo que el coeficiente de correlación entre y_i e $\hat{y}_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$ es 0.98 para la primera población (Pob098) y 0.80 para la segunda población (Pob080). Los análisis descriptivos de estas poblaciones están dados por las Tablas A.12 y A.13, mientras que los diagramas de dispersión los encontramos en las Figuras A.8 y A.9.

Figura A.8: Diagramas de dispersión de la población Pob098

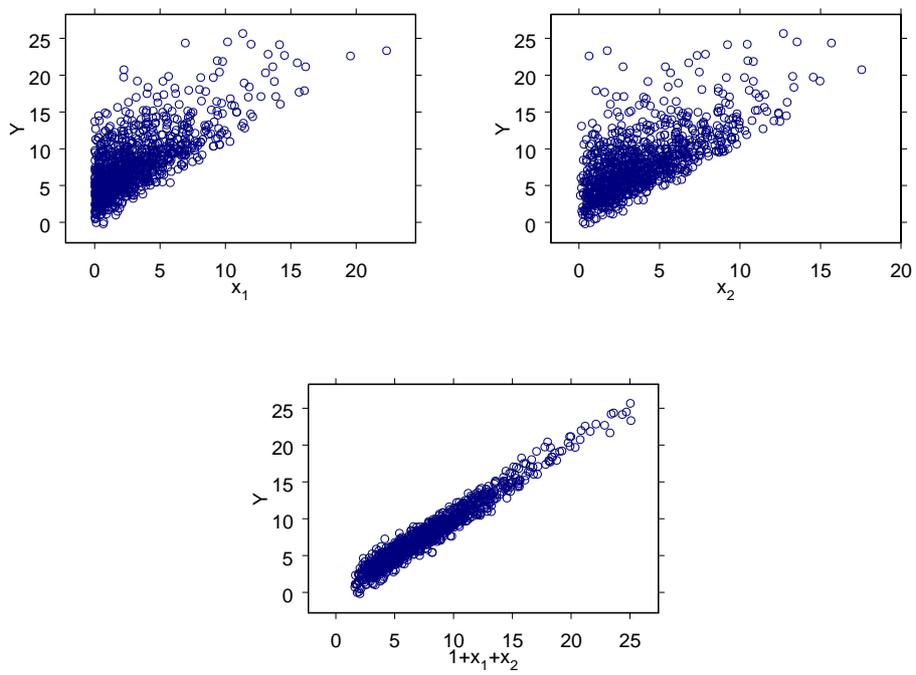
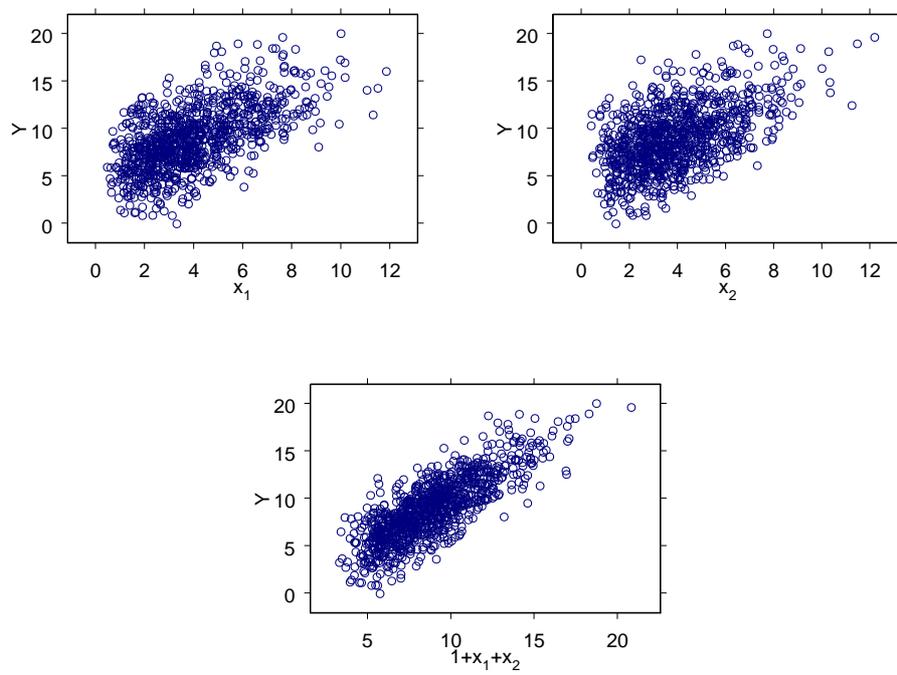


Figura A.9: Diagramas de dispersión de la población Pob080



Apéndice B

Procedimientos de estimación y selección de unidades

El objetivo de este apéndice es describir los diseños muestrales que han sido usados a lo largo de este texto. Las distintas perspectivas de estimación que existen en muestreo de poblaciones finitas están también definidas.

B.1. Métodos de muestreo probabilísticos

B.1.1. Muestreo aleatorio simple

Sea U una población finita con N unidades distintas e identificables. El muestreo aleatorio simple consiste en extraer sucesiva e independientemente unidades de esta población, sin reposición, con probabilidades iguales en cada extracción, hasta completar una muestra s con tamaño n . Las probabilidades de inclusión de primer y segundo orden están dadas por

$$\pi_i = \frac{n}{N}$$
$$\pi_{ij} = \frac{n}{N} \frac{n-1}{N-1}$$

B.1.2. Muestreo de Midzuno

Para obtener una muestra s con n elementos, el método de Midzuno (véase Singh, 2003, pg.390, para un mayor detalle) consiste en extraer un elemento de la población U con una determinada probabilidad, y los $n - 1$ restantes elementos mediante muestreo aleatorio simple de entre los $N - 1$ elementos que quedan en la población. La primera unidad suele seleccionarse mediante probabilidades proporcionales a una variable auxiliar, x , es decir, las probabilidades de selección del primer elemento vienen dadas por $\alpha_i = x_i/X$, con $i \in U$ y $X = \sum_{i=1}^N x_i$. En este caso, las probabilidades de inclusión de primer y segundo orden están dadas por:

$$\pi_i = \frac{N-n}{N-1}\alpha_i + \frac{n-1}{N-1}$$
$$\pi_{ij} = \frac{n-1}{N-1} \left[\frac{N-n}{N-2}(\alpha_i + \alpha_j) + \frac{n-2}{N-2} \right], \quad i \neq j$$

B.1.3. Muestreo de Lahiri

Supongamos que para la población U , con N elementos, se desea seleccionar un elemento de manera que las probabilidades de selección de cada uno de ellos estén dadas por p_1, \dots, p_N , con $\sum_{i=1}^N p_i = 1$. El método de Lahiri (Lahiri, 1951) obtiene una muestra con probabilidades proporcionales a una variable auxiliar x . El algoritmo a seguir es el siguiente:

Algoritmo B.1 *Extracción de unidades mediante el muestreo de Lahiri.*

Paso 1. *Obtener el mínimo valor C que verifique $\sum_{i \in s} x_i \leq C$, para cualquier muestra s . Es decir, C será la suma de los mayores valores de la variable auxiliar.*

Paso 2. *Generar un número aleatorio e mediante una distribución Uniforme entre 0 y C .*

Paso 3. *Seleccionar una muestra s_j mediante muestreo aleatorio simple.*

Paso 4. *Si $\sum_{i \in s_j} x_i \geq e$, el método de Lahiri proporcina la muestra s_j . En caso contrario ir al Paso 2.*

La probabilidades de inclusión de primer orden para las unidades seleccionadas mediante el método de Lahiri viene dadas por

$$\pi_i = \frac{nx_i}{X},$$

donde X es el total poblacional de la variable auxiliar.

B.1.4. Muestreo de Poisson

El muestreo de Poisson (llamado así por Hájek, 1964) consiste en seleccionar independientemente unidades de muestreo, con probabilidad $\pi_i = n\alpha_i$, de una población de tamaño N . Debido a la independencia de las extracciones, las probabilidades de inclusión satisfacen $\pi_{ij} = \pi_i\pi_j$. El tamaño muestral ν es una variable aleatoria tal que $E(\nu) = \sum_{i=1}^N \pi_i = n$ y $Var(\nu) = n - n^2 \sum_{i=1}^N \alpha_i^2 = \sum_{i=1}^N \pi_i(1 - \pi_i)$. Su varianza para $\pi_i = n/N$ vale $n(N - n)/N$.

Notamos que mediante el muestreo de Poisson es posible obtener una muestra vacía, cuya probabilidad viene dada por:

$$P_0 = \prod_{i=1}^N (1 - n\alpha_i),$$

que es máxima para $\alpha_i = 1/N$, en cuyo caso P_0 tiende a e^{-n} cuando N tiende a infinito.

El muestreo de Poisson modificado fue introducido por Ogus y Clark (1971) para excluir el caso de muestras vacías. En Chaudhuri y Vos (1988) se puede ver una descripción de esta técnica.

Brewer *et al.* (1972) dieron el siguiente algoritmo práctico para llevar a cabo un muestreo de Poisson realizado mediante la rotación muestral:

Algoritmo B.2 *Extracción de unidades mediante el muestreo de Poisson.*

Paso 1. *Se elige un número arbitrario c , $0 \leq c < 1$.*

Paso 2. *Se extraen N números aleatorios R_i independientes según una $U[0, 1)$.*

Paso 3. *Para $i = 1, \dots, N$ se selecciona la unidad i si:*

- (a) $n\alpha_i \leq 1 - c$ si $c \leq R_i < c + n\alpha_i$.

(b) $n\alpha_i > 1 - c$ si $R_i < c + n\alpha_i - 1$ o $R_i \geq c$.

En otro caso se rechaza la unidad.

En Singh (2003), pg. 499, puede obtenerse más información acerca de este método de muestreo.

B.2. Diseños muestrales

B.2.1. Muestreo estratificado

La población U consta de N unidades y es dividida en L subpoblaciones de tamaños N_1, N_2, \dots, N_L . Estas subpoblaciones, que reciben el nombre de estratos, no se superponen y juntas forman la totalidad de la población, es decir, $\sum_{h=1}^L N_h = N$. Llamaremos S_h al conjunto de todas las posibles muestras de tamaño n_h para el estrato h . Una vez que han sido determinados los estratos se extrae una muestra $s_h \in S_h$ de tamaño n_h de cada uno mediante un diseño de muestreo específico que asigna una probabilidad conocida $p(s_h)$ a cada $s_h \in S_h$ tal que $p(s_h) > 0$ y $\sum_{s_h \in S_h} p(s_h) = 1$. La muestra final está compuesta por el conjunto de estas submuestras y su tamaño será $n = \sum_{h=1}^L n_h$. El proceso de muestreo se realiza de modo independiente en cada estrato, lo que permite la aplicación simultánea de métodos de muestreo diferentes de acuerdo con la información disponible, el coste y las razones que motivaron la estratificación. Asumiendo que la variable principal se denota como y , y se dispone de un vector de variables auxiliares, \mathbf{x} , a continuación se detallan algunos conceptos utilizados en muestreo estratificado:

- y_{hi} , valor obtenido en la unidad i del estrato h para la característica y .
- \mathbf{x}_{hi} , valor obtenido en la unidad i del estrato h para el vector \mathbf{x} .
- $W_h = \frac{N_h}{N}$, ponderación del estrato h .
- $f_h = \frac{n_h}{N_h}$, fracción de muestreo en el estrato h .
- $\bar{Y}_h, \bar{\mathbf{X}}_h$, medias poblacionales del estrato h de las variables y y \mathbf{x} .
- $\bar{y}_h, \bar{\mathbf{x}}_h$, medias muestrales del estrato h de las variables y y del vector \mathbf{x} .

- π_{hi} , probabilidad de inclusión de primer orden de la unidad i en el estrato h .
- p_{hi} , densidad de la observación i en el estrato h

Si tomamos una muestra aleatoria simple de cada estrato, el procedimiento se conoce como muestreo estratificado aleatorio. En tal caso las probabilidades de inclusión de primer orden vienen dadas por

$$\pi_{hi} = \frac{n_h}{N_h},$$

con $h = \{1, \dots, L\}$.

B.2.2. Muestreo bifásico

El muestreo bifásico es apropiado cuando quiere hacerse un uso eficiente de la información auxiliar, pero no se dispone de los datos poblacionales necesarios para construir los estimadores basados en información auxiliar.

El muestreo bifásico consiste en extraer en una primera fase una muestra grande y poco costosa mediante cualquier diseño muestral, con el objetivo de observar o medir tan sólo el vector de variables auxiliares \mathbf{x} . Con la información obtenida se pretende dar una buena aproximación de los totales y medias poblacionales, o de cualquier otra información poblacional necesaria. En la segunda fase se extrae otra muestra mas pequeña y menos costosa a partir de la muestra anterior que ahora juega el papel de población. En esta segunda muestra la variable de interés, y , es observada, y el método de muestreo para extraerla puede ser el mismo o distinto del método de extracción de unidades usado en la primera fase.

Para una mayor profundización sobre el muestreo bifásico puede consultarse Särndal *et al.* (1992), Fernández y Mayor (1994) y Artés y García (2002).

La notación que se sigue bajo este muestreo es la siguiente. De la población U compuesta por N unidades se extrae en la primera fase una muestra, s' , de tamaño, n' , bastante grande y de bajo costo, según cierto criterio muestral, d_1 , tal que $p_{d_1}(s')$ será la probabilidad de que s' sea seleccionada y donde las correspondientes probabilidades de inclusión de primer y segundo orden se denotan, respectivamente, como π'_i y π'_{ij} para i y $j \in U$. En esta muestra, una o varias variables auxiliares pueden ser recogidas fácilmente, es decir, dicha muestra permite obtener la información auxiliar necesaria para todo el proceso.

Dada s' , una segunda muestra s de tamaño n es seleccionada en la segunda fase mediante un diseño d_2 , tal que $p(s/s')$ es la probabilidad condicional de escoger s . Las probabilidades de inclusión bajo este diseño se denotan como $\pi_{i/s'}$ y $\pi_{ij/s'}$. Notamos que $\Delta'_{ij} = \pi'_{ij} - \pi'_i \pi'_j$ y $\Delta^{s'}_{ij} = \pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}$.

B.2.3. Muestreo bifásico aplicado a la estratificación

Los dos diseños muestrales anteriores pueden combinarse en el llamado muestreo bifásico aplicado a la estratificación, el cual tiene las principales ventajas de cada uno de ellos.

Asumiendo la notación del apartado anterior, el muestreo bifásico aplicado a la estratificación se presenta cuando la variable auxiliar se usa para estratificar s' en H estratos denotados como s'_h , ($h = 1, \dots, H$) y con n'_h elementos en el estrato h . De este modo, de s'_h se puede seleccionar una muestra s_h de tamaño n_h mediante un diseño arbitrario $p_h(/s')$ independiente para cada estrato. La muestra final es $s = \bigcup_{h=1}^H s_h$. La probabilidades de inclusión para las unidades de la segunda fase se denotan como $\pi_{i/s'}$ y $\pi_{ij/s'}$, para $i, j \in s'$. Notamos que $\Delta'_{ij} = \pi'_{ij} - \pi'_i \pi'_j$ y $\Delta^{s'}_{ij} = \pi_{ij/s'} - \pi_{i/s'} \pi_{j/s'}$. Este diseño muestral permite combinar las buenas propiedades del muestreo estratificado en términos de eficiencia con el muestreo bifásico.

Notamos que cualquier estimador basado en el diseño muestral o en la aproximación modelo-asistida (véase Sección B.3) que quiera construirse bajo el muestreo bifásico o el muestreo bifásico aplicado a la estratificación deberá usar las probabilidades de inclusión de primer orden dadas por $\pi_i = \sum_{s' \ni i} p_{d1}(s') \pi_{i/s'}$. Estas probabilidades de inclusión no podrán calcularse siempre en la práctica debido a que las probabilidades $\pi_{i/s'}$, para cada s' , deben de conocerse para determinar π_i . Esto no será siempre posible porque $\pi_{i/s'}$ puede depender del resultado de la primera fase, como por ejemplo cuando la muestra de la segunda fase está diseñada mediante un muestreo proporcional a una variable auxiliar. Por esta razón, Särndal *et al.* (1992) propusieron el uso de los estimadores π^* . Las probabilidades de inclusión que usan estos estimadores están dadas por $\pi_i^* = \pi'_i \pi_{i/s'}$ y $\pi_{ij}^* = \pi'_{ij} \pi_{ij/s'}$.

B.2.4. Muestreo en ocasiones sucesivas

Consideramos que estamos haciendo un seguimiento continuo de la población U , de tamaño N , sobre dos, o más, periodos de tiempo con valores y_i en el

periodo u ocasión más reciente y x_i en una ocasión anterior (o conjunto de valores basados en ocasiones anteriores). Se asume que una muestra de tamaño n' está diseñada en la ocasión anterior. En la ocasión reciente, una submuestra (llamada muestra solapada) de tamaño m es diseñada de las n' unidades seleccionadas previamente, y $u = n - m$ unidades son reemplazadas por nuevas unidades seleccionadas de la población restante. $\chi = m/n$ será la fracción de solapamiento.

La muestra de la primera fase s' con tamaño n' está diseñada según un diseño muestral d_1 , tal que $p_{d_1}(s')$ es la probabilidad de que s' sea escogida. Las correspondientes probabilidades de inclusión de primer y segundo orden vienen dadas por π'_i, π'_{ij} , para $i, j \in U$. Dada s' , en la segunda ocasión, una muestra solapada s_m con tamaño m , es diseñada según un diseño d_2 , tal que $p_m(s_m/s')$ es la probabilidad condicional de escoger s_m . Las probabilidades de inclusión bajo este diseño se denotan como $\pi_{i/s'}$ y $\pi_{ij/s'}$. La muestra no solapada s_u es por tanto seleccionada de $U - s' = s'^c$ según el diseño d_3 , tal que $p_u(s_u/s'^c)$ es la probabilidad condicional de escoger s_u . Las probabilidades de inclusión bajo este diseño se denotarán como π_{i/s'^c} y π_{ij/s'^c} .

Notamos que todas las muestras y submuestras pueden seleccionarse bajo cualquier diseño muestral. En este trabajo se han asumido solamente dos ocasiones, aunque la generalización de los estimadores propuestos a un número más elevado de ocasiones también puede plantearse fácilmente. Los tamaños muestrales pueden ser diferentes en cada ocasión.

Para el caso de múltiples variables auxiliares, denotaremos como y a la variable principal, con valores y_1, \dots, y_N para los N elementos poblacionales. En la ocasión anterior se dispone de p variables auxiliares, x_1, x_2, \dots, x_p , las cuales serán usadas para construir estimadores multivariantes basados en información auxiliar.

Para un mayor detalle del muestreo en ocasiones sucesivas puede consultarse Jessen (1942), Patterson (1950), Narain (1953), Adhvaryu (1978), Eckler (1955), Gordon (1983), Arnab y Okafor (1992), Singh (2003), etc.

B.3. Tipos de inferencia en poblaciones finitas

En muestreo de poblaciones finitas existen diferentes perspectivas de estimación. Este epígrafe resulta apropiado no solo porque se describen cada uno de estos procedimientos de estimación, sino que también se describen las

principales ventajas e inconvenientes de cada uno de ellos.

Notamos que si bien la elección del diseño muestral y del método de selección de unidades eran dos aspectos importantes en muestreo de poblaciones finitas, no menos resulta la elección de la apropiada perspectiva de estimación, puesto que una mala elección podría dirigir, entre otros aspectos, a inferencias no validas, estimaciones inconsistentes e importantes sesgos.

B.3.1. Aproximación basada en el diseño

La aproximación basada en el diseño es el procedimiento de estimación más conocido y usado en muestreo de poblaciones finitas, aunque no necesariamente el más eficiente en el sentido de obtener estimaciones más precisas. El resto de métodos de estimación están basados, de forma directa o indirecta, en modelos de superpoblación, por lo que destacan por una mayor precisión en las estimaciones cuando el modelo es el correcto para los datos en estudio. En cualquier caso, la aproximación basada en el diseño es preferida por la mayoría de los principales institutos y organismos estadísticos por su simplicidad y menor coste computacional frente al resto de aproximaciones, aunque recientemente éstas están cobrando especial interés por la demostrada calidad en sus estimaciones.

Ejemplos de estimadores basados en el diseño son el clásico estimador de tipo Horvitz-Thompson para la media poblacional, o el estimador de tipo Hájek para la función de distribución.

B.3.2. Aproximación basada en modelos

Otro procedimiento para estimar parámetros lineales o no lineales en poblaciones finitas es la aproximación basada en modelos, la cual asume un modelo de superpoblación y donde los estimadores son dependientes del modelo, es decir, para cada estudio, el paso previo antes de poner en práctica esta aproximación es el análisis de los datos para ajustar el mejor modelo a éstos. Una vez seleccionado el modelo, se estiman los parámetros desconocidos, y las estimaciones de los parámetros poblacionales dependerán de forma directa del modelo de superpoblación estimado. Un cuidadoso contraste sobre el modelo estimado debería también llevarse a cabo con el fin de evitar inferencias no válidas.

Notamos que con la aparición de los modelos de superpoblación, la teoría de muestreo tuvo un gran empuje pues se le dotó de un instrumento muy valioso que permitió obtener resultados más concluyentes en la comparación de estrategias y, eventualmente, producir estrategias óptimas en varias situaciones. Véase Pérez (2002) y Sánchez-Crespo (2002) para un seguimiento más exhaustivo de la metodología de los modelos de superpoblación.

En los modelos de superpoblación, cobra especial importancia el uso de variables auxiliares, cuyos valores deberían ser conocidos a nivel poblacional para obtener estimaciones más precisas.

Algunos ejemplos de estimadores basados en modelos son los propuestos por Chambers y Dunstan (1986) o Dorfman y Hall (1993) para el caso de la estimación de la función de distribución.

Un inconveniente que pueden presentar los estimadores obtenidos a partir de esta aproximación es la inconsistencia bajo el diseño, como le ocurre por ejemplo al estimador de Chambers y Dunstan (1986).

Sin duda, los principales inconvenientes de estos métodos es el grado de dificultad en la computación y un pobre cumplimiento cuando el modelo especificado es incorrecto. La generalización de algunos estimadores basados en modelos a diseños muestrales más complejos también puede resultar en ocasiones bastante costosa.

B.3.3. Aproximación modelo-asistida

La aproximación modelo-asistida es otro procedimiento basado en modelos de superpoblación, aunque en este caso las estimaciones no son dependientes del modelo ajustado. Esta aproximación se usa habitualmente en las encuestas por muestreo y tiene una buena aceptación.

Ejemplos de estimadores modelo-asistidos son el estimador de regresión generalizado (*GREG*) (Cassel *et al.*, 1976, Särndal, 1980, Hedayat y Sinha, 1991) para el caso de estimación de medias y totales poblacionales, y estimadores de tipo razón y diferencia (Rao *et al.*, 1990) para el caso de estimar funciones de distribución y cuantiles. Más recientemente, son dos los principales métodos en la literatura que están categorizados como aproximaciones modelo-asistidas. Estos procedimientos son el de calibración (Deville y Särndal, 1992) y el de verosimilitud empírica (Chen y Qin, 1993, Chen y Sitter, 1999).

La aproximación modelo-asistida consiste en asumir un modelo de superpoblación. A partir de este modelo, se estiman determinados parámetros necesarios para construir el estimador. Este estimador conserva la misma estructura, independientemente de la relación que exista entre las variables involucradas en el estudio. Por este motivo, la aproximación modelo-asistida no es dependiente del modelo. El ejemplo más claro de aproximación modelo-asistida es el estimador *GREG*. En este caso se asume un modelo de regresión que estima de manera óptima el parámetro b de regresión usado por este estimador. Independientemente del tipo de relación existente entre las variables del estudio, el estimador *GREG* presenta la misma definición.

Por tanto, estos procedimientos no son dependientes de un modelo, aunque requieren el uso de uno para construir el estimador. En otras palabras, los estimadores modelo-asistidos son aproximadamente (asintóticamente) insesgados bajo el diseño, independientemente de si el modelo es correcto o no, y son particularmente eficientes si el modelo en el que se basa es correcto. Así, la aproximación modelo-asistida proporciona inferencias válidas bajo el modelo asumido y al mismo tiempo, está protegido contra una mala especificación del modelo en el sentido de proporcionar inferencias basadas en el diseño válidas para cualesquiera que sean los valores de la variable principal.

En resumen, la principal ventaja de la aproximación modelo-asistida es que proporciona inferencias válidas bajo el modelo ajustado y al mismo tiempo está protegido contra una inapropiada especificación del modelo en el sentido de proporcionar inferencias válidas basadas en el diseño, independientemente de los valores poblacionales para la variable de estudio.

B.3.4. Aproximación modelo-calibrada

La perspectiva de estimación más reciente en muestreo de poblaciones finitas es la llamada aproximación modelo-calibrada, la cual fue desarrollada por Wu y Sitter (2001a).

Estos estimadores se obtienen, en primer lugar, adaptando un modelo de superpoblación, y a continuación, usar los valores estimados mediante este modelo en la etapa de estimación. Se trata por tanto, al igual que los estimadores basados en modelos, de estimadores dependientes del modelo.

Otro ejemplo de estimadores modelo-calibrados son los propuestos por Chen y Wu (2002). Asumiendo esta aproximación, estos autores proponen estimadores para la función de distribución bajo el método de verosimilitud

empírica.

Notamos que al ser estimadores dependientes del modelo, presentarán similares propiedades a las destacadas en los estimadores basados en modelos. Por ejemplo, tienen un destacado grado de dificultad en la computación y un pobre cumplimiento cuando el modelo ajustado es incorrecto.

Apéndice C

Programación de estimadores mediante el software R

Todos los estudios de simulación en este texto se han llevado a cabo mediante el lenguaje de programación *R*. Todos los procedimientos y funciones para obtener en *R* tanto los estimadores propuestos en este texto como el resto de estimadores para cada diseño muestral están disponibles en el presente apéndice.

Son numerosas las razones por las que se ha usado este software. En primer lugar, es un lenguaje intuitivo con una gran cantidad de argumentos estadísticos que facilitan la implementación de los estimadores propuestos. Otros programas como *Mathematica*, *Matlab*, *C++*, etc., carecen de tales procedimientos estadísticos. Por otro lado, es un paquete que destaca por su rapidez y que permite obtener el mayor número de simulaciones en menor tiempo. *R* es un lenguaje de programación gratuito y disponible a cualquier usuario, al contrario de otros específicos de estadística como *SAS*, que debido a sus altas licencias está disponible, en la mayoría de los casos, a grandes empresas. El dispositivo gráfico que dispone *R* y su compatibilidad con *S-PLUS* son otros argumentos hacen que la mayoría de los investigadores en el campo del muestreo en poblaciones finitas se decanten por este software. Sirva de ejemplo los artículos publicados en este sentido (por ejemplo Wu, 2005) así como las conferencias internacionales sobre el programa *R* que también se están abriendo paso, como la segunda conferencia internacional de usuarios de *R* que se celebró recientemente del 15 al 17 de junio de 2006 en Viena, Austria. De hecho, el gran auge que está teniendo este software hace que se estén introduciendo día a día nuevos procedimientos y paquetes estadísticos.

C.1. Introducción

C.1.1. Funciones complementarias

```
Delta<-function(t,vector)
# Obtiene un vector "resul" que indica los valores del vector "vector" que son
# menores o mayores que "t".
{resul<-t>=vector
resul
}
```

```
Prob.inclusion.Mas<-function(N,n)
# Obtiene las probabilidades de inclusión de primer y segundo orden de una
# muestra obtenida bajo MAS cuando los tamaños poblacionales y muestrales son
# "N" y "n", respectivamente.
{ Pi<-vector(len=n)
  for (i in 1:n) Pi[i]<-n/N
  Pij<-matrix(nro=n, nco=n)
  for (i in 1:n)
    for (j in 1:n)
      if (i != j)
        Pij[i,j]<-(n/N)*((n-1)/(N-1))
      else Pij[i,i]<- Pi[i]
list(Pi=Pi, Pij=Pij)
}
```

```
simulaNR<-function(f,df)
# Calcula la solución de "f=0" mediante Newton-Raphson. "df" es la derivada de
# "f". Se hacen un máximo de 100 iteraciones
{tol<-10^(-10)
iter<-100
x0<-0
i<-1
while (i<iter)
  { x1<-x0-f(x0)/df(x0)
    if (f(x1)==0) i<-iter
    if (abs(x1-x0)<tol) i<-iter
    x0<-x1
    i<-i+1
  }
x1
}
```

```
simulaNRv2<-function(f,df)
# Idem, pero evita que se interrumpa el programa en caso de no existir solución
{tol<-10^(-10)
iter<-100
x0<-0
i<-1
while (i<iter)
  { x1<-x0-f(x0)/df(x0)
    if (is.na(f(x1))==TRUE)
      {
        Estado<-0
        i<-iter
        x1<-0
      }
    else

```

```

    {if (f(x1)==0) i<-iter
    if (abs(x1-x0)<tol) i<-iter
    x0<-x1
    i<-i+1
    Estado<-1
    }
  }
list(x1=x1, Estado=Estado)
}

```

```

Lag2<-function(u, ds, mu)
# Idem del anterior, pero para un vector de p-variables auxiliares. "u" es una
# matriz de orden n x p, "ds" contiene las cantidades di, y "mu" contiene los
# totales en las restricciones (véase la metodología de verosimilitud empírica)
{length(ds)->n
  if (length(u[,1]) != n) u<-t(u)
  ds1<-matrix(nr=n, nc=1)
  ds1[,1]<-ds
  length(mu)->p
  mu1<-matrix(nr=p, nc=1)
  mu1[,1]<-mu
  u<-u-rep(1,n)%*%t(mu1)
  M<-0*mu1
  dif<-1
  tol<-1e-08
  while(dif>tol)
  { D1<-0*mu1
    DD<-D1%*%t(D1)
    for(i in 1:n)
    {aa<-as.numeric(1+t(M)%*%u[i,])
      D1<-D1+ds1[i]*u[i,]/aa
      DD<-DD-ds1[i]*(u[i,]%*%t(u[i,]))/aa^2
    }
    D2<-solve(DD,D1,tol=1e-100)
    dif<-max(abs(D2))
    rule<-1
    while(rule>0)
    {rule<-0
      if(min(1+t(M-D2)%*%t(u))<=0 ) rule<-rule+1
      if(rule>0) D2<-D2/2
    }
    M<-M-D2
  }
}
M
}

```

```

densidad<-function(vector, punto)
# Obtiene la función densidad del vector "vector" en el punto "punto".
{
  N <- length(vector)
  L <- floor(N^(0.5))
  k <- 2 * L + 1 # ahora he de calcular el d_k(punto)
  x <- sort(vector)
  dis <- double(N)
  for(i in 1:N) dis[i] <- abs(abs(punto) - abs(x[i]))
  dis <- sort(dis)
  dk <- dis[k]
  final <- (k - 1)/(2 * N * dk)
  final
}

```

```

divideestratos<-function(muestral,Nh)
# Divide "muestral" (que son las etiquetas de las unidades muestrales) en
# estratos con ayuda de "Nh" (tamaños poblacionales de los estratos).
{L<-length(Nh)
muestral<-sort(muestral)
tamlh<-c()
Trues<-muestral<=Nh[1]
tamlh<-c(tamlh, sum(Trues))
muestral<-setdiff(muestral,muestral[Trues])
for (h in 2:L)
{Trues<-muestral<=sum(Nh[1:h])
tamlh<-c(tamlh, sum(Trues))
muestral<-setdiff(muestral,muestral[Trues])
}
tamlh
}

afijacion.proporcional<-function(n, Nh)
# Obtiene los tamaños muestrales de cada estrato mediante afijación
# proporcional, donde "n" es el tamaño de la muestra y "Nh" son los
# tamaños poblacionales de los estratos.
{
  (n/sum(Nh))*Nh->mivector
  redondeo<-round(mivector)
  suma<-sum(redondeo)
  while (suma!=n)
  { if (suma<n)
    {
      mivector-redondeo->decimales
      order(decimales)->orden
      length(decimales)->le
      orden[le]->posicion
      redondeo[posicion]<-redondeo[posicion]+1
      mivector[posicion]<-redondeo[posicion]
      suma<-sum(redondeo)
    }
    else if (suma>n)
    {
      redondeo-mivector->decimales
      order(decimales)->orden
      length(decimales)->le
      orden[le]->posicion
      redondeo[posicion]<-redondeo[posicion]-1
      mivector[posicion]<-redondeo[posicion]
      suma<-sum(redondeo)
    }
  } #end while
redondeo
#### EJEMPLO
# afijacion.proporcional(100, c(517,633,350))
}

```

C.1.2. Métodos de muestreo

```

EtiquetasHH<-function(n,x)
# Extrae una muestra de tamaño "n" con reemplazo y con probabilidades
# proporcionales a "x".
{N<-length(x)
runif(n)->Nalea
# Nalea es un vector con n-números aleatorios entre 0 y 1.
p<-vector(len=N)

```

```

Tx<-sum(x)
for (i in 1:N) p[i]<-(x[i]/Tx)
B<-vector(len=N)
B[1]<-p[1]
for (i in 2:N) B[i]<-p[i]+B[i-1]
eti<-vector(len=n)
for (i in 1:n)
  {j<-1;
   eti[i]<-1
   while (Nalea[i]>B[j])
     # Cuando B[j-1]<Nalea<=B[j], se considera la unidad j
     { j<-j+1
       eti[i]<-j
     }
  }
# eti contiene las etiquetas de las unidades a seleccionar.
eti
}

lahiri<-function(n,x)
# Extrae una muestra de tamaño "n" por el metodo de Lahiri.
# "x" es la variable auxiliar para extraer la muestra.
{
  sort(x)->X
  subx<-vector(len=n)
  N<-length(x)
  m1<-N-n
  for(i in 1:n) subx[i]<-X[i+m1]
  C<-sum(subx) # C es la suma de los n-mayores valores de x
  muestrax<-vector(len=n)
  sumax<-0
  e<-1
  while(sumax<e)
    { e<-runif(1,0,C)
      Emas<-sample(N,n)
      muestrax<-x[Emas]
      sumax<-sum(muestrax)
    }
  Emas # Emas contiene las etiquetas para extraer la muestra.
}

Etiquetas.midzuno<-function(n,z)
# Extrae muestra de tamaño "n" por el método de Midzuno. "z" es la variable
# auxiliar para extraer la primera unidad.
{ j<-EtiquetasHH(1,z) # Función que extrae una unidad proporcional a z.
  N<-length(z)
  E.todas<-1:N
  E.resto<-setdiff(E.todas,j)
  Emidzuno<-1:n-1
  Emidzuno<-sample(E.resto,n-1)
  Emidzuno<-c(j,Emidzuno)
  Emidzuno
}

Etiquetas.Poisson<-function(n,x)
# Extrae 1 muestra de tamaño "n" mediante un muestreo de Poisson.
# x: Es la variable auxiliar para extraer la muestra.
{ length(x)->N
  valor.c<-runif(1,0,1)
  valores.R<-runif(N,0,1)
  muestra<-c()

```

```

sum(x)->Total.X
z<-x/Total.X
for (i in 1:N)
  { if ( valor.c<=valores.R[i] && valores.R[i]<valor.c+n*z[i] && n*z[i]<=1-
valor.c) muestra<-c(muestra,i)
    else if (valores.R[i]<valor.c+n*z[i]-1 | valores.R[i]>=valor.c ) if
(n*z[i]>1-valor.c) muestra<-c(muestra,i)
    }
  length(muestra)->n.obtenido
muestra
#### EJEMPLO
#Etiquetas.Poisson(50,Fam1500[,2])
}

```

C.1.3. Funciones de distribución básicas

```

F.distribucion.HT<-function(N,t,Pi,vector)
# Estimación Horvitz_Thompson de la F.distribución. Los parámetros son:
# "N" : Tamaño de la Población
# "t" : Para que calcule la F.d en el punto "t".
# "Pi": Prob. de inclusión. Puede ser un número o vector de longitud n.
# "vector": Datos muestrales de una variable. Vector de longitud "n".
# Nota: Si Pi=1 y longitud de vector=N, obtiene la verdadera F.d.
{
length(vector)->n
Delta(t,vector)->vectorD
Fd<-(1/N)*sum(vectorD/Pi)
Fd
}

F.distribucion.Hk<-function(N,t,Pi,vector)
# Estimación tipo Hájek para la F.d. Los parámetros de entrada son los mismos
# que los definidos en la función anterior.
{
length(vector)->n
Delta(t,vector)->vectorD
Fd<-sum(vectorD/Pi)/sum(1/Pi)
Fd
}

F.distribucion.r<-function(N,t,Pi,muestray,muestrax,datosx)
# Obtiene la estimación tipo razón para la función de distribución.
# "datosx" son los datos poblacionales de la variable auxiliar x.
{
length(muestray)->zn
Delta(t, muestray)->vectorDy
if (sum(vectorDy)==zn) Fd<-1
else if (sum(vectorDy)==0) Fd<-0
else
{ Rgorro<-sum(muestray/Pi)/sum(muestrax/Pi)
Delta(t, Rgorro*muestrax)->vectorDRx
if (sum(vectorDRx)==0) Fd<-0
else
{Delta(t, Rgorro*datosx)->vectorDRX
Fd<- (1/N) * ( sum(vectorDy/Pi)/sum(vectorDRx/Pi) ) * sum(vectorDRX)
}
}
}
Fd
}

```

```

F.distribucion.d1<-function(N,t,Pi,muestray,muestrax,datosx)
# Estimación tipo diferencia para la F.d.El factor de corrección es d=1.
{
  length(muestray)->n
  Rgorro<-sum(muestray/Pi)/sum(muestrax/Pi)
  Delta(t, muestray)->vectorDy
  Delta(t, Rgorro*muestrax)->vectorDRx
  Delta(t, Rgorro*datosx)->vectorDRX
  Fd<-(1/N)*(sum(vectorDy/Pi)+(sum(vectorDRX)-sum(vectorDRx/Pi)))
Fd
}

F.distribucion.dopt.est<-function(N,t,Pi,muestray,muestrax,datosx)
# F.d. de tipo diferencia con factor de corrección óptimo.
{
  length(muestray)->n
  Rgorro<-sum(muestray/Pi)/sum(muestrax/Pi)
  Delta(t, muestray)->vectorDy
  Delta(t, Rgorro*muestrax)->vectorDRx
  Delta(t, Rgorro*datosx)->vectorDRX
  Fd.Hk<-sum(vectorDy/Pi)/sum(1/Pi)
  Fd.RX<-(1/N)*sum(vectorDRX)
  SY<-(1/(n-1))*sum( (vectorDy-Fd.Hk)^2 )
  SX<-(1/(n-1))*sum( (vectorDRx-Fd.RX)^2 )
  Ro<-cov(vectorDy, vectorDRx)/var(vectorDRx)
  if (SX==0) d.optimo<-1 else d.optimo<-Ro*SY/SX
  di<-1/Pi
  Fd<-(1/N)*( sum(di*vectorDy)+d.optimo*(sum(vectorDRX)-sum(di*vectorDRx)) )
Fd
}

F.distribucion.CD<-function(N,t,Etiquetas,muestray,muestrax,datosx)
# Función de Distribución de Chambers y Dunstan
{
  b<-sum(muestrax*muestray/sqrt(muestrax))/sum(muestrax^2/sqrt(muestrax))
  U<-(muestray-b*muestrax)/sqrt(muestrax)
  Eti.N<-1:N
  Eti2<-setdiff(Eti.N, Etiquetas)
  muestrax2<-datosx[Eti2]
  n<-length(muestray)
  n2<-N-n
  q<-(t-b*muestrax2)/sqrt(muestrax2)
  k<-vector(len=n2)
  for (j in 1:n2) k[j]<-sum(Delta(q[j],U) )
  vectorD<-Delta(t,muestray)
  Fd<-(1/N)*(sum(vectorD)+(1/n)*sum(k) )
Fd
# EJEMPLO
# F.distribucion.CD(N=304, t=median(counties[,1]), sample(counties[,1],100),
# sample(counties[,2],100), counties[,2] )
}

F.distribucion.RKM<-function(N,t,Pi,Pij,muestray,muestrax,datosx)
# Función de Distribución de Rao, Kovar y Mantel
{
  n<-length(muestray)
  Rgorro<-sum(muestray/Pi)/sum(muestrax/Pi)
  destrella<-(1/Pi)/sum(1/Pi)
  vectorG<-vector(len=N)
  qN<-(t-Rgorro*datosx)/sqrt(datosx)
  qn<-(muestray-Rgorro*muestrax)/sqrt(muestrax)
  for (j in 1:N) vectorG[j]<-sum(destrella*Delta(qN[j],qn))
  vectorG2<-vector(len=n)
  vn1<-(t-Rgorro*muestrax)/sqrt(muestrax)
  vn2<-(muestray-Rgorro*muestrax)/sqrt(muestrax)
}

```

```

for (j in 1:n) vectorG2[j]<-(1/sum(Pi[j]/Pij[j,]))*
      sum( (Pi[j]/Pij[j,])*Delta(vn1[j],vn2) )
vectorDy<-Delta(t,muestray)
Fd<-(1/N)*(sum(vectorDy/Pi) + sum(vectorG) - sum(vectorG2/Pi) )
Fd
}

F.distribucion.v1<-function(Tipo.de.Fd,N,t,Pi,vectory)
# En función del valor asignado a "TipoFd" obtiene:
# TipoFd = 1:Estimador de Horwitz-Thompson      (HT)
# TipoFd = 2:Estimador de Hayek                 (Hk)
{
switch(Tipo.de.Fd,
      Fd<-F.distribucion.HT(N,t,Pi,vectory),
      Fd<-F.distribucion.Hk(N,t,Pi,vectory)
)
Fd
}

F.distribucion<-
function(TipoFd,N,t,Pi, Pij, Etiquetas, muestray, muestrax, datosx)
# En función del valor asignado a "TipoFd" obtiene:
# TipoFd = 1:Estimador de Horwitz-Thompson      (HT)
# TipoFd = 2:Estimador de Hayek                 (Hk)
# TipoFd = 3:Estimador de Razón                 (r)
# TipoFd = 4:Estimador de Diferencia            (d1)
# TipoFd = 5:Estimador de Chambers y Dunstan   (CD)
# TipoFd = 6:Estimador de Rao, Kovar y Mantel   (RKM)
{
switch(TipoFd,
      Fd<-F.distribucion.HT(N,t,Pi,muestray),
      Fd<-F.distribucion.Hk(N,t,Pi,muestray),
      Fd<-F.distribucion.r(N,t,Pi,muestray,muestrax,datosx),
      Fd<-F.distribucion.d1(N,t,Pi,muestray,muestrax,datosx),
      Fd<-F.distribucion.CD(N,t,Etiquetas, muestray,muestrax,datosx),
      Fd<-F.distribucion.RKM(N,t,Pi,Pij,muestray,muestrax,datosx),
)
Fd
}

F.distribucion.PS<-function(Nh, nh, t , Pi , vectory)
# F.d. estimada de Silva y Skinner. Los parámetros de entrada son:
# "Nh" :Tamaños poblacionales de estratos.
# "nh" :Tamaños muestrales de estratos.
# "t"  :Para que calcule la F.d en el punto "t".
# "Pi" :Prob. de inclusión. Puede ser un N° o un vector de longitud "n".
# vectory: Datos muestrales de una variable. Vector de longitud "n".
{
N<-sum(Nh)
L<-length(Nh)
Nh.gorro<-vector(len=L)
Fd.g<-vector(len=L)
suma<-0
for (h in 1:L)
{
if (h !=1 )
a2<-sum(nh[1:(h-1)])+1
else
a2<-1
b2<- sum(nh[1:h])
Nh.gorro[h]<-sum(1/Pi[a2:b2])
}
}

```

```

    Fd.g[h]<- sum(Delta(t,vectorx[a2:b2]) / Pi[a2:b2] ) / Nh.gorro[h]
    suma<-suma+ Nh[h]*Fd.g[h]
  }
Fd<-(1/N)*suma
list(Fd=Fd, Fd.g=Fd.g, Nh.gorro=Nh.gorro)
}

F.distribucion.SJT<-function(Qx.1,t,vectorx.2,vectorx.1,vectorx.2)
# F.distribucion estimada de Singh, Joarder y Tracy.
# Muestras obtenidas bajo M. bifásico. "vectorx.2" y "vectorx.1" son los datos
# de y y x de la muestra de la 2ª fase. "vectorx.1" son los datos de la primera
# fase. "Qx.1" es el estimador del cuantil de x obtenido a partir de la primera
# muestra y "t" es punto para evaluar la F.d.
{
nx.1<-sum(vectorx.1<=Qx.1)
nx.2<-sum(vectorx.2<=Qx.1)
tam1<-length(vectorx.1)
tam2<-length(vectorx.2)
TruesA.x2<- vectorx.2<=Qx.1
Fd.YA<- sum( vectorx.2[TruesA.x2]<=t ) / nx.2
TruesB.x2<- vectorx.2>Qx.1
Fd.YB<- sum( vectorx.2[TruesB.x2]<=t ) / (tam2-nx.2)
Fd<- ( (nx.1*Fd.YA) / tam1) + ( (tam1-nx.1)*Fd.YB/tam1 )
Fd
}

varianza.PS<-function(N,a, Pi, Pij)
# Obtiene la varianza estimada del estimador de la F.d. de Silva y Skinner. "N"
# es el tamaño poblacional, "a" está definido en Silva y Skinner (1995), y "Pi"
# y "Pij" son las probabilidades de inclusión.
{
n<-length(a)
suma<-0
Indices<-1:n
for (j in 2:n)
  {
    Ind.i<- Indices < j
    Ind.i <- Indices[Ind.i]
    for (i in Ind.i)
      suma<-suma+ ( (Pi[i]*Pi[j]-Pij[i,j])/Pij[i,j] )*(a[i]/Pi[i]-a[j]/Pi[j])^2
  }
var<-(1/N^2)*suma
var
}

covarianza.PS<-function(N,a1,a2, Pi, Pij)
# Covarianza estimada del estimador de la F.d. de Silva y Skinner
{
n<-length(a1)
suma<-0
Indices<-1:n
for (j in 2:n)
  {
    Ind.i<- Indices < j
    Ind.i <- Indices[Ind.i]
    for (i in Ind.i)
      suma<-suma+ ( (Pi[i]*Pi[j]-Pij[i,j])/Pij[i,j] )*(a1[i]/Pi[i]-a2[j]/Pi[j])^2
  }
var<-(1/N^2)*suma
var
}

```

C.1.4. Cuantiles básicos

```
cuantil.HT.Rapido<-function(TipoFd,beta,N,Pi,vector)
# Obtiene el cuantil estimado de orden "beta" cuando la F.d es de HT o Hk. Si
# Pi=1 y longitud de "vector"="N", calcula el verdadero cuantil. Se usa un
# algoritmo de bisección para una búsqueda más eficiente.
{
if (beta<0 | beta>1) stop("Introduce un cuantil beta entre [0 , 1]")
length(Pi)->n1
length(vector)->n
if (n1 != 1)
{
if (n1 != n) stop("Longitud distinta. \n", "vector: \n",vector,"\n", "Pi:
\n",Pi,"\n")
}
sort(vector)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<-F.distribucion.v1(TipoFd,N,t,Pi,datos)
if (Fd.n>beta)
{ while(B-A != 1)
{ t<-datos[M]
Fd<-F.distribucion.v1(TipoFd,N,t,Pi,datos)
if (Fd>beta)
{ B<-M
M<-floor((A+B)/2)
if (A==B) B<-A+1
t<-datos[B]
}
else if (Fd<beta)
{ A<-M
M<-floor((A+B)/2)
if (A==B) { B<-A+1; t<-datos[M] }
else t<-datos[B]
}
}
else {B<-2 ; A<-1 ; t<-datos[M]}
}
}
else { t<-datos[n] }
t
# EJEMPLO
# cuantil.HT.Rapido(1,beta=0.47,N=100,Pi=1,vector=1:100)
}

cuantil<-
function(TipoFd,beta,N,Pi,Pij,Etiquetas, muestray, muestrax, datosx)
# Obtiene el cuantil estimado de orden beta cuando se usan distintos tipos de
# F.d. Se hace uso de un metodo de bisección para su cálculo.
{
if (beta<0 | beta>1) stop("Introduce un cuantil beta entre [0 , 1]")
length(Pi)->n1
length(muestray)->n
if (n1 != 1)
{ if (n1 != n) stop("Longitud distinta. \n", "muestray: \n",muestray,"\n",
"Pi: \n",Pi,"\n")
}
sort(muestray)->datos
t<-datos[n]
```

```

A<-1
B<-n
M<-floor(n/2)
Fd.n<-F.distribucion(1,N,t,Pi,Pij, Etiquetas, muestray, muestrax, datosx)
if (Fd.n>beta)
  { while(B-A != 1)
    { t<-datos[M]
    Fd<-F.distribucion(TipoFd,N,t,Pi,Pij,Etiquetas,muestray,muestrax,datosx)
      if (Fd>beta)
        { B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        { A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
else { t<-datos[n] }
t
# EJEMPLO:
# cuantil(TipoFd=1,beta=0.47,N=100,Pi=1,muestray=1:100)
}

cuantil.HK<-function(TipoFd=2,beta,N,Pi,muestray)
# Obtiene el cuantil de orden "beta" estimado de tipo Hájek
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<-F.distribucion(TipoFd=TipoFd,N=N,t=t,Pi=Pi,muestray=muestray)
if (Fd.n>beta)
  { while(B-A != 1)
    { t<-datos[M]
    Fd<-F.distribucion(TipoFd=TipoFd,N=N,t=t,Pi=Pi,muestray=muestray)
      if (Fd>beta)
        { B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        { A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
else { t<-datos[n] }
t
# cuantil.HK(TipoFd=2,beta=0.47,N=100,Pi=1,muestray=1:100)
}

```

```

cuantil.r<-function(beta,N,Pi,muestray, muestrax,x)
# Obtiene el cuantil de orden "beta" estimado de tipo razón
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<- F.distribucion.r(N,t,Pi,muestray,muestrax, x)
if (Fd.n>beta)
  { while(B-A != 1)
    { t<-datos[M]
      Fd<-F.distribucion.r(N,t,Pi,muestray,muestrax, x)
      if (Fd>beta)
        { B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        { A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
else { t<-datos[n] }
t
# eti<-sample(1500,100)
# cuantil.r(beta=0.50,N=1500,Pi=100/1500,muestray=Fam1500[eti,1],
# muestrax=Fam1500[eti,2],x=Fam1500[,2])
}

cuantil.d1<-function(beta,N,Pi,muestray, muestrax,x)
# Obtiene el cuantil de orden "beta" de tipo diferencia (con factor d=1)
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<- F.distribucion.d1(N,t,Pi,muestray,muestrax, x)
if (Fd.n>beta)
  { while(B-A != 1)
    { t<-datos[M]
      Fd<- F.distribucion.d1(N,t,Pi,muestray,muestrax, x)
      if (Fd>beta)
        { B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        { A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
}

```

```

    }
  }
else { t<-datos[n] }
t
}

cuantil.dopt.est2<-function(beta,N,Pi,muestray, muestrax,x)
# Cuantil de orden "beta" de tipo diferencia (con factor d óptimo)
{
length(muestray)->n
sort(muestray)->datos
t<-datos[1]
Fd<- beta-1
i<-0
while (Fd<beta)
  {i<-i+1
  if(i>n)
  {
t<-datos[n]
Fd<-beta+1
}
else
{
Fd<- F.distribucion.dopt.est(N,t,Pi,muestray,muestrax, x)
t<-datos[i]
}
}
}

cuantil.CD<-function(beta,N,etiquetas,muestray, muestrax,x)
# Obtiene el cuantil de orden "beta" de tipo Chambers y Dunstan.
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<- F.distribucion.CD(N,t,etiquetas, muestray,muestrax, x)
if (Fd.n>beta)
  {
while(B-A != 1)
  { t<-datos[M]
Fd<- F.distribucion.CD(N,t,etiquetas, muestray,muestrax, x)
if (Fd>beta)
  { B<-M
M<-floor((A+B)/2)
if (A==B) B<-A+1
t<-datos[B]
}
else if (Fd<beta)
  { A<-M
M<-floor((A+B)/2)
if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
}
else {B<-2 ; A<-1 ; t<-datos[M]}
}
}
else { t<-datos[n] }
t
}

```

```

cuantil.RKM<-function(beta,N,Pi, Pij, muestray, muestrax,x)
# Obtiene el cuantil de orden "beta" de tipo Rao, Kovar, Mantel.
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<- F.distribucion.RKM(N,t,Pi,Pij, muestray,muestrax, x)
if (Fd.n>beta)
  { while(B-A != 1)
    { t<-datos[M]
      Fd<- F.distribucion.RKM(N,t,Pi,Pij, muestray,muestrax, x)
      if (Fd>beta)
        { B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        { A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
else { t<-datos[n] }
t
}

cuantil.PS<-function(beta,Nh,nh,Pi,vector)
# Obtiene el cuantil estimado utilizando la F.d de Silva y Skinner.
{
if (beta<0 | beta>1) stop("Introduce un cuantil beta entre [0 , 1]")
length(Pi)->n1
length(vector)->n
if (n1 != 1)
  if (n1 != n) stop("Longitud distinta. \n", "vector: \n",vector,"\n",
"Pi: \n",Pi,"\n")
Fd<- -1
sort(vector)->datos.sort
i<-0
n<-sum(nh)
tn<-datos.sort[n]
Salida<- F.distribucion.PS(Nh, nh, tn , Pi , vector)
Fdn<-Salida$Fd
if (Fdn<beta) t<-tn
else
  while (Fd<beta)
  {
i<-i+1
t<-datos.sort[i]
Salida<- F.distribucion.PS(Nh, nh, t , Pi , vector)
Fd<-Salida$Fd
}
t
}

```

```

cuantil.SJT<-function(beta,Qx.1,vectory.2,vectorx.1,vectorx.2)
# Estima el cuantil de orden "beta" para la Fd de Singh, Joarder y Tracy. Se
# tienen los mismos parámetros que "F.distribucion.SJT".
{
if (beta<0 | beta>1) stop("Introduce un cuantil beta entre [0 , 1]")
Fd<- -1
sort(vectory.2)->datos.sort
i<-0
tam2<-length(vectory.2)
tn<-datos.sort[tam2]
Fdn<-F.distribucion.SJT( Qx.1, tn , vectory.2,vectorx.1,vectorx.2)
  if (Fdn<beta) t<-tn
  else
    while (Fd<beta)
      {
        i<-i+1
        t<-datos.sort[i]
        Fd<-F.distribucion.SJT( Qx.1, t , vectory.2,vectorx.1,vectorx.2)
      }
t
}

```

C.2. El método de verosimilitud empírica

C.2.1. Tratamiento de datos faltantes

```

Faltantes.MAS<-function(NombreP, R, n, y, x, p, q)
# Simulación para la estimación de la media poblacional bajo MAS en presencia
# de datos faltantes. Obtención de varios estimadores.
# "NombreP": Nombre de la población
# "R": Número de muestras a extraer (Replicas).
# "n": Tamaño de la muestra.
# "y": Datos poblacionales de la característica de interés.
# "x": Datos poblacionales de la variable auxiliar.
# "p", "q": Número de datos faltantes.
{
# Cantidades generales
N<-length(y)
N2<-length(x)
if(N != N2) stop("Los vectores $y$ y $x$ tienen distinta longitud")
if(n > N) stop("El tamaño de muestra es mayor que N. Cambia $n$")
Xmedia<-mean(x)
Ymedia<-mean(y)
n1<-n-p-q
u1<-vector(len=n1)
destrella1<-1/n1
destrella12<-1/(n1+p)
pesos1<-vector(len=n1)
pesos12<-vector(len=n1)
psumal<-0
psumal2<-0
# Archivos para escribir resultados.
salidaEstimador<-paste("C:\\PEMLE\\Faltantes\\est_Mas_",
NombreP, "_R", R, "n", n, "p", p, "q", q, ".txt", sep="")
salidaError<-paste("C:\\PEMLE\\Faltantes\\err_Mas_", NombreP, "_R",
R, "n", n, "p", p, ".txt", sep="")
ecmPemle1<-0 ; sesPemle1<-0 ; ses2Pemle1<-0

```

```

ecmPemle12<-0 ; sesPemle12<-0 ; ses2Pemle12<-0
ecmT1<-0 ; sesT1<-0 ; ses2T1<-0
ecmT2<-0 ; sesT2<-0 ; ses2T2<-0
ecmT3<-0 ; sesT3<-0 ; ses2T3<-0
ecmT4<-0 ; sesT4<-0 ; ses2T4<-0
ecmYalfa<-0 ; sesYalfa<-0 ; ses2Yalfa<-0
ecmYtp<-0 ; sesYtp<-0 ; ses2Ytp<-0
ecmYreg<-0 ; sesYreg<-0 ; ses2Yreg<-0
ecmMas13<-0 ; sesMas13<-0 ; ses2Mas13<-0
for (i in 1:R)
{
suma<-0
repeticiones<-0
while (suma !=1)
{
sample(N,n)->etiquetas
sample(etiquetas,p)->etiquetas2
setdiff(etiquetas,etiquetas2)->etiquetasnp
sample(etiquetasnp, q)->etiquetas3
setdiff(etiquetasnp,etiquetas3)->etiquetas1
# Muestras de $y$ y de $x$:
muestray1<-y[etiquetas1]
muestrax1<-x[etiquetas1]
muestrax2<-x[etiquetas2]
muestray3<-y[etiquetas3]
muestray13<-c(muestray1,muestray3)
muestrax12<-c(muestrax1,muestrax2)
# Cálculo de los pesos del PEMLE
u1<-muestrax1-Xmedia
u2<-muestrax12-Xmedia
f<-function(x) sum(destrella1*u1/(1+x*u1))
df<-function(x) -sum(destrella1*u1^2*(1+x*u1)^(-2))
simulaNR(f,df)->lambda1
f<-function(x) sum(destrella12*u12/(1+x*u12))
df<-function(x) -sum(destrella12*u12^2*(1+x*u12)^(-2))
simulaNR(f,df)->lambda12
for (j in 1:n1)
{
# Valores del vector "p" de cada muestra para calcular PEMLE.
pesos1[j]<-destrella1/(1+lambda1*u1[j])
pesos12[j]<-destrella1/(1+lambda12*u1[j])
}
repeticiones<-repeticiones+1
if (repeticiones>100) stop("No es posible obtener sum(p)=1.")
suma<-sum(pesos1)
} # Fin del While
psuma1<-psuma1+sum(pesos1)
psuma12<-psuma12+sum(pesos12)
# Valores de las varianzas y Covarianzas
# Para alfa.optimo
s2y<-var(muestray13)
s2x<-var(muestrax12)
sxy<-cov(muestrax1, muestray1)
Bgorro<-sum(as.numeric(muestrax1)*as.numeric(muestray1))/
sum(as.numeric(muestrax1)^2)
VarY1<-(s2y/n1)*(1-n1/N)
VarX1<-(s2x/n1)*(1-n1/N)
VarY3<-(s2y/q)*(1-q/N)
CovY1X1<-(1/n1-1/N)*sxy
A1<-VarY1+Bgorro^2*VarX1-2*Bgorro*CovY1X1
C1<-VarY3
ifelse( (n-p)/2>=q , CovY1Y3<-(1/n1-1/N)*s2y, CovY1Y3<-(1/q-1/N)*s2y )

```

```

ifelse( (n-p)/2>=q , CovX1Y3<-(1/n1-1/N)*sxy, CovX1Y3<-(1/q-1/N)*sxy )
D1<-CovY1Y3-Bgorro*CovX1Y3
alfaOPT<-(C1-D1)/(A1+C1-2*D1)
# Para Ymedia.theta.Phi
A3<-VarY3-CovY1Y3
ifelse( (n-q)/2>=p , CovY1X2<-(1/n1-1/N)*sxy, CovY1X2<-(1/p-1/N)*sxy )
ifelse( p>=q , CovY3X2<-(1/p-1/N)*sxy, CovY3X2<-(1/q-1/N)*sxy )
C3<-CovY1X1-CovY1X2-CovX1Y3+CovY3X2
D3<-CovY1X2-CovY3X2
E3<-VarY1+VarY3-2*CovY1Y3
ifelse( (n-q)/2>=p , CovX1X2<-(1/n1-1/N)*s2x, CovX1X2<-(1/p-1/N)*s2x )
VarX2<-(s2x/p)*(1-p/N)
F3<-VarX1+VarX2-2*CovX1X2
G3<-VarX2-CovX1X2
H3<-CovX1Y3-CovY3X2
nume<-C3*(A3-Bgorro*D3)+E3*(Bgorro*G3+H3)
deno<-E3*F3-C3^2
Phi<-(1/Bgorro)*(nume/deno)
Theta<-(A3+Bgorro*(Phi*C3-D3))/E3
# Para estimador Ytotal.regresion2 de Rueda-González (2004)
Areg<-2*VarY1*N^2 + 2*VarY3*N^2 - 4*CovY1Y3*N^2
Breg<-2*Bgorro* N^2*( (-1)*CovY1X1+CovY1X2+CovX1Y3-CovY3X2 )
Creg<-(-2)*VarY3*N^2+2*CovY1Y3*N^2+2*Bgorro*N^2*((-1)*CovY1X2 +CovY3X2)
Dreg<- 2*Bgorro^2*N^2*(VarX1+VarX2-2*CovX1X2)
Ereg<-2*Bgorro*N^2*((-1)*Bgorro*VarX2+Bgorro*CovX1X2-CovX1Y3+CovY3X2 )
betaREG<- ((Breg*Creg-Areg*Ereg)/(Areg*Dreg-Breg^2))
alfaREG<- -(Creg/Areg -(Breg/Areg)*betaREG)
##### ESTIMADORES
# Estimador de verosimilitud empírica. Muestra 1
Pemle1<-sum(pesos1*muestray1)
ecmPemle1<-ecmPemle1+(Pemle1-Ymedia)^2
sesPemle1<-sesPemle1+abs(Pemle1-Ymedia)
ses2Pemle1<-ses2Pemle1+(Pemle1-Ymedia)
# Estimador de verosimilitud empírica. Muestra 12
Pemle12<-sum(pesos12*muestray1)
ecmPemle12<-ecmPemle12+(Pemle12-Ymedia)^2
sesPemle12<-sesPemle12+abs(Pemle12-Ymedia)
ses2Pemle12<-ses2Pemle12+(Pemle12-Ymedia)
# Estimador directo. Muestra 1
Mas1<-mean(muestray1)
# Estimador directo. Muestra 13 (Base)
Mas13<-mean(muestray13)
ecmMas13<-ecmMas13+(Mas13-Ymedia)^2
sesMas13<-sesMas13+abs(Mas13-Ymedia)
ses2Mas13<-ses2Mas13+(Mas13-Ymedia)
# Para obtener el estimador T1
nume<- ( n1*mean(muestrax1) + p*mean(muestrax2) )
deno<-(n-q)*mean(muestrax1)
T1<-Mas1*(nume/deno)
ecmT1<-ecmT1+(T1-Ymedia)^2
sesT1<-sesT1+abs(T1-Ymedia)
ses2T1<-ses2T1+(T1-Ymedia)
# Para obtener el estimador T2
nume2<-deno
deno2<-nume
T2<-Mas1*(nume2/deno2)
ecmT2<-ecmT2+(T2-Ymedia)^2
sesT2<-sesT2+abs(T2-Ymedia)
ses2T2<-ses2T2+(T2-Ymedia)
# Para obtener el estimador T3
nume<- (n1*mean(muestrax1)+p*mean(muestrax2))*(n1*Mas1+q*mean(muestray3))

```

```

deno<-(n-q)*(n-p)*mean(muestrax1)*Mas1
T3<-Mas1*(nume/deno)
ecmT3<-ecmT3+(T3-Ymedia)^2
sesT3<-sesT3+abs(T3-Ymedia)
ses2T3<-ses2T3+(T3-Ymedia)
# Para obtener el estimador T4
nume<-( n1*Mas1+q*mean(muestray3) )*(n-q)*mean(muestrax1)
deno<- ( n1*mean(muestrax1) + p*mean(muestrax2) )*(n-p)
T4<-(nume/deno)
ecmT4<-ecmT4+(T4-Ymedia)^2
sesT4<-sesT4+abs(T4-Ymedia)
ses2T4<-ses2T4+(T4-Ymedia)
# Para Obtener YmediaAlfa y el Ymedia_theta_Phi
Mas3<-mean(muestray3)
Yalfa<-alfaOPT*Pemle1+(1-alfaOPT)*Mas3
ecmYalfa<-ecmYalfa+(Yalfa-Ymedia)^2
sesYalfa<-sesYalfa+abs(Yalfa-Ymedia)
ses2Yalfa<-ses2Yalfa+(Yalfa-Ymedia)
Ytp<-Theta*Mas1+(1-Theta)*Mas3+
Bgorro*(Xmedia-Phi*mean(muestrax1)-(1-Phi)*mean(muestrax2))
ecmYtp<-ecmYtp+(Ytp-Ymedia)^2
sesYtp<-sesYtp+abs(Ytp-Ymedia)
ses2Ytp<-ses2Ytp+(Ytp-Ymedia)
# Para obtener YmediaREG
Yreg<-alfaREG*Mas1+(1-alfaREG)*Mas3+ Bgorro*(Xmedia-
betaREG*mean(muestrax1)-(1-betaREG)*mean(muestrax2))
ecmYreg<-ecmYreg+(Yreg-Ymedia)^2
sesYreg<-sesYreg+abs(Yreg-Ymedia)
ses2Yreg<-ses2Yreg+(Yreg-Ymedia)
## Se escriben estimadores
linea<-paste(i,Ymedia, Mas13, Pemle1,Pemle12, Yalfa, Ytp,Yreg,
T1,T2,T3,T4, sep="\t")
write(linea, file=salidaEstimador, ncolumns=1, append=T)
}
## AHORA SE ESCRIBEN LOS ERRORES
ecmMas13<-ecmMas13/R
RecmMas13<-1
RecmPemle1<-(ecmPemle1/R)/ecmMas13
RecmPemle12<-(ecmPemle12/R)/ecmMas13
RecmT1<-(ecmT1/R)/ecmMas13
RecmT2<-(ecmT2/R)/ecmMas13
RecmT3<-(ecmT3/R)/ecmMas13
RecmT4<-(ecmT4/R)/ecmMas13
RecmYalfa<-(ecmYalfa/R)/ecmMas13
RecmYtp<-(ecmYtp/R)/ecmMas13
RecmYreg<-(ecmYreg/R)/ecmMas13
RsesMas13<-sesMas13/(R*Ymedia)
RsesPemle1<-sesPemle1/(R*Ymedia)
RsesPemle12<-sesPemle12/(R*Ymedia)
RsesT1<-sesT1/(R*Ymedia)
RsesT2<-sesT2/(R*Ymedia)
RsesT3<-sesT3/(R*Ymedia)
RsesT4<-sesT4/(R*Ymedia)
RsesYalfa<-sesYalfa/(R*Ymedia)
RsesYtp<-sesYtp/(R*Ymedia)
RsesYreg<-sesYreg/(R*Ymedia)
# SESGOS RELATIVOS
Rses2Mas13<-ses2Mas13/(R*Ymedia)
Rses2Pemle1<-ses2Pemle1/(R*Ymedia)
Rses2Pemle12<-ses2Pemle12/(R*Ymedia)
Rses2T1<-ses2T1/(R*Ymedia)

```

```

Rses2T2<-ses2T2/(R*Ymedia)
Rses2T3<-ses2T3/(R*Ymedia)
Rses2T4<-ses2T4/(R*Ymedia)
Rses2Yalfa<-ses2Yalfa/(R*Ymedia)
Rses2Ytp<-ses2Ytp/(R*Ymedia)
Rses2Yreg<-ses2Yreg/(R*Ymedia)
linea<-paste(psumal,psumal2, n, p, q,
  ecmMas13,RecmMas13,RecmPemle1,RecmPemle12,
  RecmYalfa,RecmYtp,RecmYreg,RecmT1,RecmT2,RecmT3,RecmT4,
  sesMas13, RsesMas13, RsesPemle1,RsesPemle12, RsesYalfa, RsesYtp,
  RsesYreg,RsesT1, RsesT2, RsesT3, RsesT4,
  ses2Mas13, Rses2Mas13, Rses2Pemle1,Rses2Pemle12, Rses2Yalfa, Rses2Ytp,
  Rses2Yreg,Rses2T1, Rses2T2, Rses2T3, Rses2T4, sep="\t")
write(linea, file=salidaError, ncolumns=1, append=T)
# EJEMPLO
# Faltantes.MAS(NombreP, R=30,n=200,Fam1500[,1],Fam1500[,2],20,15)
}

```

C.2.2. Estimación modelo-asistida de la función de distribución

```
## FUNCIONES COMPLEMENTARIAS PARA LOS ESTIMADORES DE CALIBRACION ##
```

```

estimacionHP<-function(s,y,d)
{
# "s" es la muestra, "d" es el vector con la probabilidades de inclusion de las
# unidades incluidas en "s", e "y" es la variable de estudio.
ys<-y[s]
w<-d*ys
return(sum(w))
}

calibracion<-function(s,y,d,X,q)
{
# "X" es una matriz que contiene las v.v. auxiliares, cada fila de la matriz
# indica una v. auxiliar. "q" es el vector de constantes en la calibración.
ys<-y[s]
Xs<-X[,s]
qs<-q[s]
xs<-t(Xs)*(d*qs)
td<-Xs%*%xs
ty<-estimacionHP(s,y,d)
if (prod(svd(td)$d)==0) return(ty)
else {
tx<-c()
Tx<-c()
for (j in 1:dim(X)[1]){
tx[j]<-estimacionHP(s,X[j,],d)
Tx[j]<-sum(X[j,])}
ts<-solve(td)
r<-Xs%*%(d*qs*ys)
B<-ts%*%(r)
u<-(Tx-tx)%*%B
tyr<-ty+u
tyrr<-tyr[1,]
return(tyrr)}
}

```

```

distribucion<-function(t,y)
{ k<-t>=y
  f<-sum(k)/length(y)
  f
}

distribucionv<-function(t,y)
{F<-c()
  for (i in 1:length(t) ) F[i]<-distribucion(t[i],y)
  F
}

dis.invertv<-function(y, beta)
{t<-sort(y)
  f<-c()
  F<-distribucionv(t,y)
  for( j in 1:length(beta) )
  {
    w<-beta[j]<=F
    i<-length(F)-sum(w)+1
    f[j]<-t[i]
  }
  f
}

calibracionfd1<-function(t,to,s,d,y,X,q)
# Obtiene estimadores de calibración para un vector de variables auxiliares.
# "t": vector de puntos donde queremos estimar la función de distribución.
# "to": vector de puntos de la variable auxiliar empleados en la calibración.
{
  ys<-y[s] ; Xs<-X[,s] ; qs<-q[s] ; xs<-t(Xs)*(d*qs) ; td<-Xs%%xs
  ts<-solve(td); r<-Xs%%(d*qs*ys) ; B<-ts%%(r); u<-t(B)%%X ; z<-u[1,]
  l<-dis.invertv(z,to)
  H<-array(0,dim=c(length(l),length(y)))
  for (j in 1:length(l)) H[j,]<-l[j]>=z
  FW<-c()
  for (j in 1:length(t))
  {
    deltax<-t[j]>=y
    FW[j]<-calibracion(s,deltax,d,H,q)}
  return(FW/length(y))
}

calibracionfd2<-function(t,to,s,d,y,x,q)
# Obtiene estimadores de calibración para una sola variable auxiliar.
# "x": es la variable auxiliar.
{
  ys<-y[s] ; xs<-x[s] ; qs<-q[s]; B<-sum(d*qs*xs*ys)/sum(d*qs*xs*xs)
  z<-B*x ; l<-to ; H<-array(0,dim=c(length(l),length(y)))
  for (j in 1:length(l)) H[j,]<-l[j]>=z
  FW<-c()
  for (j in 1:length(t))
  {
    deltax<-t[j]>=y
    FW[j]<-calibracion(s,deltax,d,H,q)}
  return(FW/length(y))
}

## ESTIMADORES DE VEROSIMILITUD EMPÍRICA ##

```

```

F.distribucion.ve.RX.vec<-function(N, t, t0, Pi, muestray, muestrax, datosx)
# Estimación de la F.d con una v.a. y un pto "t0". Se calcula el coeficiente R
# (ver Rao et al, 1990). "t" : vector de ptos donde se evalua la función.
{
  length(muestray)->zn
  psuma<-1
  Rgorro<-sum(muestray/Pi)/sum(muestrax/Pi)
  Delta(t0, Rgorro*muestrax)->vectorDRx
  Delta(t0, Rgorro*datosx)->vectorDRX
  vector.u<- vectorDRx-(1/N)*sum(vectorDRX)
  destrella<- (1/Pi)/sum(1/Pi)
  f<-function(x) sum( (destrella*vector.u)/(1+x*vector.u))
  df<-function(x) -sum(destrella*vector.u^2*(1+x*vector.u)^(-2))
  resul<-simulaNRv2(f,df)
  Estado.NR<-resul$Estado
  lambda<-resul$x1
  vector.p<-destrella/(1+lambda*vector.u)
  psuma<-sum(vector.p)
  length(t)->T
  Fd<-vector(len=T)
  for (i in 1:T)
  { Delta(t[i], muestray)->vectorDy
    Fd[i]<-sum(vector.p*vectorDy)
  }
}
list(Fd=Fd, psuma=psuma, vector.p=vector.p, vector.g=Rgorro*muestrax
, Estado.NR=Estado.NR )
} # end F.distribucion.ve.RX.vec

```

```

F.distribucion.ve.t0x.vec<-function(N,t,t0.x, Pi,muestray,muestrax, datosx)
# Estimación de la F.d con una v.a. y un punto "t0" basado en "x".
# "t" es un vector.
{
  length(muestray)->zn
  psuma<-1
  vectorDelta.X<-Delta(t0.x,datosx)
  Fd.X<-(1/N)*sum(vectorDelta.X)
  v.delta.n<-Delta(t0.x,muestrax)
  vector.u<- v.delta.n-Fd.X
  destrella<- (1/Pi)/sum(1/Pi)
  f<-function(x) sum( (destrella*vector.u)/(1+x*vector.u))
  df<-function(x) -sum(destrella*vector.u^2*(1+x*vector.u)^(-2))
  resul<-simulaNRv2(f,df)
  Estado.NR<-resul$Estado
  lambda<-resul$x1
  vector.p<-destrella/(1+lambda*vector.u)
  psuma<-sum(vector.p)
  length(t)->T
  Fd<-vector(len=T)
  for (i in 1:T)
  { Delta(t[i], muestray)->vectorDy
    Fd[i]<-sum(vector.p*vectorDy)
  }
}
list(Fd=Fd, psuma=psuma, Estado.NR=Estado.NR)
} # end F.distribucion.ve.t0x.vec

```

```

F.distribucion.ve.BX.vec<-function(N,t,t0, Pi,muestray,muestrax, datosx)
# Estimación de la F.d con una o varias v.a. y un pto "t0". Se calcula el
# coeficiente B de regresión. "t" : vector de ptos donde se evalua la función.
{
  length(muestray)->zn
  psuma<-1
  di<-1/Pi
  if (length(muestrax)==zn)
  {

```

```

muestrax<-as.matrix(muestrax)
muestrax<-t(muestrax)
datosx<-as.matrix(datosx)
datosx<-t(datosx)
}
J<-length(muestrax[,1])
Betagorro<-matrix(nr=J,nc=1)
if ( length(muestrax)==zn )
Betagorro[1,1]<-sum((1/Pi)*muestrax*muestray)/sum((1/Pi)*muestrax^2)
else
{ xs<-t(muestrax)*(1/Pi)
  td<-muestrax%%xs
  ts<-solve(td)
  r<-muestrax%%((1/Pi)*muestray)
  Betagorro<-ts%%(r)
}
g<-t(Betagorro)%%datosx
vectorDBX<-Delta(t0,g)
vectorDBx<-Delta(t0,t(Betagorro)%%muestrax)
vector.u<- vectorDBx-(1/N)*sum(vectorDBX)
destrella<- (1/Pi)/sum(1/Pi)
f<-function(x) sum( (destrella*vector.u)/(1+x*vector.u))
df<-function(x) -sum(destrella*vector.u^2*(1+x*vector.u)^(-2))
  resul<-simulaNRv2(f,df)
  Estado.NR<-resul$Estado
  lambda<-resul$x1
  vector.p<-destrella/(1+lambda*vector.u)
  psuma<-sum(vector.p)
length(t)->T
Fd<-vector(len=T)
for (i in 1:T)
  { Delta(t[i], muestray)->vectorDy
    Fd[i]<-sum(vector.p*vectorDy)
  }
list(Fd=Fd, psuma=psuma, vector.p=vector.p,
vector.g=t(Betagorro)%%muestrax,Estado.NR=Estado.NR )
} # end F.distribucion.ve.BX

```

```

F.distribucion.Mve.2t.vec<-function(N,t,t0, Pi,muestray,muestrax, datosx)
# F.d. de Verosimilitud Empírica basada en la aproximación modelo-calibrada. Se
# basa en un modelo de regresión lineal. Se usa una v.auxiliar y un punto "t".
{
  length(muestray)->n
  di<-1/Pi
  psuma<-1
  muestray<-as.numeric(muestray)
  muestrax<-as.numeric(muestrax)
  betal<-sum(di*muestrax*muestray)/ sum(di*muestrax^2)
  beta0<- (1/N)*(sum(di*muestray)-betal*sum(di*muestrax) )
  vector.mu<-beta0 + betal*muestrax
  errores.n<-muestray-vector.mu
  var.e<-sum(errores.n^2)/(N-2)
  puntos.n<-(t0-vector.mu)/var.e
  vector.g<-pnorm(puntos.n)
  vector.MU<-beta0+betal*datosx
  puntos.N<-(t0-vector.MU)/var.e
  vector.G<-pnorm(puntos.N)
  vector.u<- vector.g-(1/N)*sum(vector.G)
  destrella<- (di)/sum(di)
  f<-function(x) sum( (destrella*vector.u)/(1+x*vector.u))
  df<-function(x) -sum(destrella*vector.u^2*(1+x*vector.u)^(-2))
}

```

```

        resul<-simulaNRv2(f,df)
        Estado.NR<-resul$Estado
        lambda<-resul$x1
        vector.p<-destrella/(1+lambda*vector.u)
        psuma<-sum(vector.p)
        length(t)->T
        Fd<-vector(len=T)
        for (i in 1:T)
            { Delta(t[i], muestray)->vectorDy
              Fd[i]<-sum(vector.p*vectorDy)
            }
    list(Fd=Fd, psuma=psuma, vector.g=vector.g,Estado.NR=Estado.NR)
} # end F.distribucion.Mve.2t.vec

F.distribucion.Mve.Jvar.vec<-function(N,t,t0, Pi,muestray,muestrax, datosx)
# F.d. de Verosimilitud Empírica basada en la aproximación modelo-calibrada. Se
# usan J-variables auxiliares y un modelo de regresión lineal con término
# independiente. Cada fila de "muestrax" o "datosx" es una variable auxiliar.
{
    length(muestray)->n
    di<-1/Pi
    psuma<-1
    muestray<-as.numeric(muestray)
    J<-length(muestrax[,1])
    X<-matrix(nr=n,nc=(J+1))
    X[,1]<-rep(1,n)
    for (j in 1:J) X[, (j+1)]<-as.numeric(muestrax[j,])
    Y<-as.matrix(muestray)
    Coeficientes<-solve(crossprod(X))%*%t(X)%*%Y ## Es una matriz columna.
    vector.mu<-X%*%Coeficientes ### Matriz columna
    vector.mu<-as.vector(vector.mu)
    errores.n<-muestray-vector.mu
    var.e<-sum(errores.n^2)/(N-J-1)
    puntos.n<-(t0-vector.mu)/var.e
    vector.g<-pnorm(puntos.n)
    X2<-matrix(nr=N,nc=(J+1))
    X2[,1]<-rep(1,N)
    for (j in 1:J) X2[, (j+1)]<-as.numeric(datosx[j,])
    vector.MU<-X2%*%Coeficientes
    vector.MU<-as.vector(vector.MU)
    puntos.N<-(t0-vector.MU)/var.e
    vector.G<-pnorm(puntos.N)
    vector.u<- vector.g-(1/N)*sum(vector.G)
    destrella<- (di)/sum(di)
    f<-function(x) sum( destrella*vector.u)/(1+x*vector.u)
    df<-function(x) -sum(destrella*vector.u^2*(1+x*vector.u)^(-2))
        resul<-simulaNRv2(f,df)
        Estado.NR<-resul$Estado
        lambda<-resul$x1
        vector.p<-destrella/(1+lambda*vector.u)
        psuma<-sum(vector.p)
        length(t)->T
        Fd<-vector(len=T)
        for (i in 1:T)
            {
                Delta(t[i], muestray)->vectorDy
                Fd[i]<-sum(vector.p*vectorDy)
            }
    list(Fd=Fd, psuma=psuma, vector.g=vector.g,Estado.NR=Estado.NR)
} # end F.distribucion.Mve.Jvar.vec

```

```

F.distribucion.RX.Tt.vec<-function(N,t,valores.t, Pi,muestray,muestrax, datosx)
# Estimación modelo-asistida de la F.d. de verosimilitud empírica. "valores.t"
#son los puntos usados en las restricciones. Se usa MAS
# "muestrax" : matriz (J x n). Se usa el coeficiente R.
# J: es el número de variables auxiliares.
# P: número de puntos t.
{
  length(muestray)->zn
  if (length(muestrax)==zn)
    {
      muestrax<-as.matrix(muestrax)
      muestrax<-t(muestrax)
      datosx<-as.matrix(datosx)
      datosx<-t(datosx)
    }
  J<-length(muestrax[,1])
  Rgorro<-matrix(nr=J,nc=1)
  for (j in 1:J) Rgorro[j,1]<-sum(muestray/Pi)/sum(muestrax[j,]/Pi)
  g.n<-t(Rgorro)%*%muestrax
  g.N<-t(Rgorro)%*%datosx
  P<-length(valores.t)
  vector.t<-c()
  for (p in 1:P) vector.t<-c(vector.t, cuantil(1, valores.t[p], N,
Pi=1,muestray=g.N) )
  matrizDRx<-matrix(nr=P, nco=zn)
  matrizDRX<-matrix(nr=P, nco=N)
  vectorF.RX<-double(P)
  for (p in 1:P)
    {
      Delta(vector.t[p], g.n)->matrizDRx[p,]
      Delta(vector.t[p], g.N)->matrizDRX[p,]
      vectorF.RX[p]<-(1/N)*sum(matrizDRX[p,])
    }
  lambda<-Lag2(u=matrizDRx, ds=(1/Pi), mu=vectorF.RX)
  lambda<-as.matrix(lambda)
  destrella<-(1/Pi)/sum(1/Pi)
  vector.p<-vector(len=zn)
  matriz.u<- matrizDRx-vectorF.RX
  for (i in 1:zn) vector.p[i]<-destrella[i]/(1+t(lambda)%*%matriz.u[,i])
  psuma<-sum(vector.p)
  length(t)->T
  Fd<-vector(len=T)
  for (i in 1:T)
    {
      Delta(t[i], muestray)->vectorDy
      Fd[i]<-sum(vector.p*vectorDy)
    }
  list(Fd=Fd, psuma=psuma, vector.p=vector.p, vector.g=g.n )
} # end F.distribucion.RX.Tt.vec

F.distribucion.BX.T.vec<-function(N,t,valores.t, Pi,muestray,muestrax, datosx)
# Estimación modelo-asistida de la F.d. de verosimilitud empírica. "valores.t"
#son los puntos usados en las restricciones. Se usa MAS
# "muestrax" : matriz (J x n). Se usa el coeficiente B de regresión.
# J: es el número de variables auxiliares.
# P: número de puntos t.
{
  length(muestray)->zn
  if (length(muestrax)==zn)
    {
      muestrax<-as.matrix(muestrax)
      muestrax<-t(muestrax)
      datosx<-as.matrix(datosx)
      datosx<-t(datosx)
    }
  J<-length(muestrax[,1])

```

```

betagorro<-matrix(nr=J,nc=1)
  if ( length(muestrax)==zn )
    betagorro[1,1]<-sum((1/Pi)*muestrax*muestray)/sum((1/Pi)*muestrax^2)
  else
    { xs<-t(muestrax)*(1/Pi)
      td<-muestrax%%xs
      ts<-solve(td)
      r<-muestrax%%((1/Pi)*muestray)
      betagorro<-ts%%(r)
    }
  g.n<-t(betagorro)%%muestrax
  g.N<-t(betagorro)%%datosx
  P<-length(valores.t)
  vector.t<-c()
  for (p in 1:P) vector.t<-c(vector.t, cuantil(1, valores.t[p], N,
Pi=1,muestray=g.N) )
  matrizDRx<-matrix(nr=P, nco=zn)
  matrizDRX<-matrix(nr=P, nco=N)
  vectorF.RX<-double(P)
  for (p in 1:P)
    { Delta(vector.t[p], g.n)->matrizDRx[p,]
      Delta(vector.t[p], g.N)->matrizDRX[p,]
      vectorF.RX[p]<-(1/N)*sum(matrizDRX[p,])
    }
  lambda<-Lag2(u=matrizDRx, ds=(1/Pi), mu=vectorF.RX)
  lambda<-as.matrix(lambda)
  destrella<- (1/Pi)/sum(1/Pi)
  vector.p<-vector(len=zn)
  matriz.u<- matrizDRx-vectorF.RX
  for (i in 1:zn) vector.p[i]<-destrella[i]/(1+t(lambda)%%matriz.u[,i])
  psuma<-sum(vector.p)
  length(t)->T
  Fd<-vector(len=T)
  for (i in 1:T)
    { Delta(t[i], muestray)->vectorDy
      Fd[i]<-sum(vector.p*vectorDy)
    }
  vector.t<-c(vector.t, max(g.N) )
  list(Fd=Fd, psuma=psuma, vector.t=vector.t, vector.p=vector.p, vector.g=g.n)
} # end F.distribucion.BX.T.vec

```

```

FdPemle.MAS.T.vec<-function(n,y,x, puntos.t, t0,t0.x)
# Asumiendo Muestreo Aleatorio Simple, obtiene diferentes estimadores de la
#función de distribución. Los argumentos son:
# "n": Tamaño de la muestra
# "y": Característica de interés
# "x": Variable auxiliar. Cada fila es una variable.
# "puntos.t": Vector de puntos donde evaluar la Fd.
# "t0": cuantil para utilizar en las restricciones.
# "t0.x": cuantil para utilizar en Fd.X.
{ N<-length(y)
  sample(N,n)->etiquetas
  muestray<-y[etiquetas]
  muestrax2<-x[,etiquetas]
  muestrax<-muestrax2[1,] # muestrax no es una matriz, es un vector.
  x2<-x
  x<-x[1,] # x es un vector.
  Prob.inclusion.Mas(N,n)->prob
  Pi<-prob$Pi
  Pij<-prob$Pij
  di<-1/Pi

```

```

T<-length(puntos.t)
MatrizSumas<-vector(len=15) # Las columnas son los distintos estimadores
MatrizFd<-matrix(nr=T,nc=24)
Fd.Y<-vector(len=T)
for (j in 1:T)
{
t<-puntos.t[j]
DY<-Delta(t,y)
#### F.d CON UN SOLO PUNTO t EN LAS RESTRICCIONES.
F.distribucion.Hk(N,t,Pi,muestray)->MatrizFd[j,1] ## Fd.base
Fd.Y[j]<-(1/N)*sum(DY) ## Fd.Y
F.distribucion.r(N,t,Pi,muestray,muestrax, x)->MatrizFd[j,8]
F.distribucion.dl(N,t,Pi,muestray,muestrax, x)->MatrizFd[j,9]
F.distribucion.dopt.est(N,t,Pi,muestray,muestrax, x)->MatrizFd[j,10]
F.distribucion.CD(N,t,etiquetas, muestray,muestrax, x)->MatrizFd[j,11]
F.distribucion.RKM(N,t,Pi,Pij, muestray,muestrax, x) ->MatrizFd[j,12]
} # end for (j in 1:T)
F.distribucion.ve.t0x.vec(N,puntos.t,t0.x, Pi,muestray,muestrax, x)->zdatos
MatrizFd[,2]<-zdatos$Fd
MatrizSumas[1]<-zdatos$psuma
### RX solo puede usarse para lv.auxiliar, con un solo punto t o varios.
F.distribucion.ve.RX.vec(N,puntos.t,t0, Pi,muestray,muestrax, x)->zdatos
MatrizFd[,3]<-zdatos$Fd
MatrizSumas[2]<-zdatos$psuma
F.distribucion.ve.BX.vec(N,puntos.t,t0, Pi,muestray,muestrax, x)->zdatos
MatrizFd[,4]<-zdatos$Fd
MatrizSumas[3]<-zdatos$psuma
F.distribucion.ve.BX.vec(N,puntos.t,t0, Pi,muestray,muestrax2, x2)->zdatos
MatrizFd[,5]<-zdatos$Fd
MatrizSumas[4]<-zdatos$psuma
F.distribucion.Mve.2t.vec(N,puntos.t,t0, Pi,muestray,muestrax, x)->zdatos
MatrizFd[,6]<-zdatos$Fd
MatrizSumas[5]<-zdatos$psuma
F.distribucion.Mve.Jvar.vec(N,puntos.t,t0, Pi,muestray,muestrax2, x2)->zdatos
MatrizFd[,7]<-zdatos$Fd
MatrizSumas[6]<-zdatos$psuma
## Programas para P-puntos t.
valores.t<-c(0.25,0.5,0.75)
F.distribucion.RX.Tt.vec(N,puntos.t,valores.t,Pi,muestray,muestrax,x)->zdatos
MatrizFd[,13]<-zdatos$Fd
MatrizSumas[7]<-zdatos$psuma
F.distribucion.BX.T.vec(N,puntos.t,valores.t,Pi,muestray,muestrax,x)->zdatos
MatrizFd[,14]<-zdatos$Fd
MatrizSumas[8]<-zdatos$psuma
### Este es el vector para el estimador de calibracion
vector.t<-zdatos$vector.t
MatrizFd[,15]<-
calibracionfd2(puntos.t,to=vector.t,s=etiquetas,d=1/Pi,y,x,q=rep(1,N) )
##### aqui con 2 variables auxiliares
F.distribucion.BX.T.vec(N,puntos.t,valores.t,Pi,muestray,muestrax2,x2)->zdatos
MatrizFd[,16]<-zdatos$Fd
MatrizSumas[9]<-zdatos$psuma
MatrizFd[,17]<-
calibracionfd1(puntos.t,to=c(valores.t,1),s=etiquetas,d=1/Pi,y,X=x2,q=rep(1,N))
MatrizFd[,18]<-regresionfijo3v2(
puntos.t,to=c(valores.t,1),s=etiquetas,d=1/Pi,y,X=x2,q=rep(1,N))
### aqui con otro valor de t.
valores.t<-c(0.25,0.4,0.6,0.75)
F.distribucion.RX.Tt.vec(N,puntos.t,valores.t,Pi,muestray,muestrax, x)->zdatos
MatrizFd[,19]<-zdatos$Fd
MatrizSumas[10]<-zdatos$psuma

```

```

F.distribucion.BX.T.vec(N,puntos.t,valores.t, Pi,muestray,muestrax, x)->zdatos
  MatrizFd[,20]<-zdatos$Fd
  MatrizSumas[11]<-zdatos$psuma
## 2 variables auxiliares
F.distribucion.BX.T.vec(N,puntos.t,valores.t,Pi,muestray,muestrax2, x2)->zdatos
  MatrizFd[,21]<-zdatos$Fd
  MatrizSumas[12]<-zdatos$psuma
### Último vector t0
  valores.t<-c(0.25,0.4,0.5,0.6,0.75)
F.distribucion.RX.Tt.vec(N,puntos.t,valores.t, Pi,muestray,muestrax, x)->zdatos
  MatrizFd[,22]<-zdatos$Fd
  MatrizSumas[13]<-zdatos$psuma
F.distribucion.BX.T.vec(N,puntos.t,valores.t, Pi,muestray,muestrax, x)->zdatos
  MatrizFd[,23]<-zdatos$Fd
  MatrizSumas[14]<-zdatos$psuma
### con 2 v.a
F.distribucion.BX.T.vec(N,puntos.t,valores.t,Pi,muestray,muestrax2,x2)->zdatos
  MatrizFd[,24]<-zdatos$Fd
  MatrizSumas[15]<-zdatos$psuma
list(MatrizSumas=MatrizSumas, MatrizFd=MatrizFd, Fd.Y=Fd.Y)
# EJEMPLO:
#FdPemle.MAS.T.vec(n=100,y=Fam1500[,1],x=t(Fam1500[,2:3]),
#puntos.t=c(6000,8135.5,10000) , t0=8135.5,t0.x=40200.5)
}

```

```

SimulaFdPemle.v4<-funcion(NombreP, B, y, x, vector.n, vector.beta, beta0)
# Simulaciones bajo MAS con diferentes estimadores de la F.d. ARGUMENTOS:
# "NombreP": Nombre de la población. Se usa para darle nombre a ficheros.
# "B": Número de replicas o de experimentos aleatorios.
# "y": Datos poblacionales de la variable de interés.
# "x": Datos poblacionales de las variables auxiliar. En cada fila una v.a.
# "vector.n": Vector de los tamaños muestrales de las muestras.
# "vector.beta": Orden de cuantiles.
# "beta0": Orden del cuantil a usar en restricciones.
{ N<-length(y)
  lT<-length(vector.beta)
  puntos.t<-vector(len=lT)
  for (t in 1:lT) puntos.t[t]<-cuantil(1, vector.beta[t], N, Pi=1,muestray=y)
  t0.x<-cuantil(1, beta0, N, Pi=1,muestray=x[1,])
  t0<-cuantil(1, beta0, N, Pi=1,muestray=y)
  salidaError<-c()
  salidaSumas<-c()
  lineal<-paste("n", "t", "Denominador", "Recm.Base",
    "Recm.t0x", "Recm.Rx", "Recm.Bx", "Recm.Bx2", "Recm.g.est", "Recm.g.Jvar",
    "Recm.Razon", "Recm.dif.d1", "Recm.dif.dopt.est", "Recm.CD", "Recm.RKM",
    "Recm.RXT", "Recm.BXT", "Recm.cal", "Recm.BXTv2", "Recm.cal2", "Recm.cal2v2",
    "Recm.RXT2", "Recm.BXT2", "Recm.BXT2v2",
    "Recm.RXT3", "Recm.BXT3", "Recm.BXT3v2",
    "Rses.Base",
    "Rses.t0x", "Rses.Rx", "Rses.Bx", "Rses.Bx2", "Rses.g.est", "Rses.g.Jvar",
    "Rses.Razon", "Rses.dif.d1", "Rses.dif.dopt.est", "Rses.CD", "Rses.RKM",
    "Rses.RXT", "Rses.BXT", "Rses.cal", "Rses.BXTv2", "Rses.cal2", "Rses.cal2v2",
    "Rses.RXT2", "Rses.BXT2", "Rses.BXT2v2",
    "Rses.RXT3", "Rses.BXT3", "Rses.BXT3v2",
    "RsesRe.Base",
    "RsesRe.t0x", "RsesRe.Rx", "RsesRe.Bx", "RsesRe.Bx2", "RsesRe.g.est",
    "RsesRe.g.Jvar",
    "RsesRe.Razon", "RsesRe.dif.d1", "RsesRe.dif.dopt.est",
    "RsesRe.CD", "RsesRe.RKM",

```

```

        "RsesRe.RXT",          "RsesRe.BXT",          "RsesRe.cal", "RsesRe.BXTv2",
"RsesRe.cal2", "RsesRe.cal2v2",
        "RsesRe.RXT2", "RsesRe.BXT2", "RsesRe.BXT2v2",
        "RsesRe.RXT3", "RsesRe.BXT3", "RsesRe.BXT3v2",
        sep="\t")
linea2<-paste("n",          "SumaT.t0x",          "SumaT.Rx",          "SumaT.Bx",
"SumaT.Bx2", "SumaT.g.est", "SumaT.g.Jvar",
        "SumaT.RXT",          "SumaT.BXT", "SumaT.BXTv2",
        "SumaT.RXT2", "SumaT.BXT2", "SumaT.BXT2v2",
        "SumaT.RXT3", "SumaT.BXT3", "SumaT.BXT3v2", sep="\t")
linea3<-paste("b", "t0", "t0.x", "t", "Fd.Y", "Fd.HK",
        "Fd.t0x", "Fd.Rx", "Fd.Bx", "Fd.Bx2", "Fd.g.est", "Fd.g.Jvar",
        "Fd.Razon", "Fd.dif.d1", "Fd.dif.dopt.est", "Fd.CD", "Fd.RKM",
        "Fd.RXT", "Fd.BXT", "Fd.cal", "Fd.BXTv2", "Fd.cal2", "Fd.cal2v2",
        "Fd.RXT2", "Fd.BXT2", "Fd.BXT2v2",
        "Fd.RXT3", "Fd.BXT3", "Fd.BXT3v2",
        sep="\t")
salidaSumas<- paste("C:\\Pemle\\Fd\\Sumas",NombreP,"FdPemleB",B,".xls", sep="")
write(linea2, file=salidaSumas, ncolumns=1, append=T)
for (t in 1:lT)
{
    salidaError<-c(salidaError,
paste("C:\\Pemle\\Fd\\err",NombreP,"FdPemleB",B,"beta",vector.beta[t], ".xls",
sep="") )
    write(linea1, file=salidaError[t], ncolumns=1, append=T)
}
salidaVarErrores<-
paste("C:\\Pemle\\Fd\\VariablesErrores",NombreP,"FdPemleB",B,".xls", sep="")
write(linea1, file=salidaVarErrores, ncolumns=1, append=T)
salidaVarEstimadores<-
paste("C:\\Pemle\\Fd\\VariablesEstimadores",NombreP,"FdPemleB",B,".xls",
sep="")
write(linea3, file=salidaVarEstimadores, ncolumns=1, append=T)
for (n in vector.n)
{salidaEstimador<-c()
for (t in 1:lT)
{
    salidaEstimador<-c(salidaEstimador,
paste("C:\\Pemle\\Fd\\est",NombreP,"FdPemleB",B,"n",n,"beta",vector.beta[t], ".x
ls", sep="") )
    write(linea3, file=salidaEstimador[t], ncolumns=1, append=T)
}
}
M.Sumas<-double(15)
M.ECM<-matrix(0,nr=lT, nc=24)
M.ses<-matrix(0,nr=lT, nc=24)
M.sesRe<-matrix(0,nr=lT, nc=24)
for (b in 1:B)
{
    FdPemle.MAS.T.vec(n,y,x, puntos.t, t0,t0.x)->Datos
    M.Sumas<-M.Sumas+Datos$MatrizSumas
    Fd.Y<-Datos$Fd.Y
    for(i in 1:24)
    {
        M.ECM[,i]<-M.ECM[,i]+(Datos$MatrizFd[,i]-Fd.Y)^2
        M.ses[,i]<-M.ses[,i]+abs(Datos$MatrizFd[,i]-Fd.Y)
        M.sesRe[,i]<-M.sesRe[,i]+(Datos$MatrizFd[,i]-Fd.Y)
    }
}
### Se escriben los estimadores
for (t in 1:lT)
{ linea<-paste(b, t0, t0.x, puntos.t[t], Fd.Y[t],sep="\t")
for (i in 1:24) linea<-paste(linea, Datos$MatrizFd[t,i],sep="\t")
}

```

```

        write(linea,file=salidaEstimador[t], ncolumns=1, append=T)      }
    } # end del for (b in 1:B)
### Se escriben SUMAS
    linea<-paste(n,sep="\t")
    for (i in 1:15) linea<-paste(linea, M.Sumas[i],sep="\t")
    write(linea,file=salidaSumas, ncolumns=1, append=T)
### Se escribe ECM, ses y sesRe
    ### ECM:
    RE<-matrix(nr=1T, nc=24)
    Denominador<-M.ECM[,1]/B
    for(i in 1:24) RE[,i]<-(M.ECM[,i]/B)/Denominador
    ### Sesgos Absolutos
    RBa<-matrix(nr=1T, nc=24)
    for(i in 1:24) RBa[,i]<-M.ses[,i]/(B*Fd.Y)
    ### Sesgos relativos
    RBr<-matrix(nr=1T, nc=24)
    for(i in 1:24) RBr[,i]<-100*M.sesRe[,i]/(B*Fd.Y) ### En tanto por ciento.
    for (t in 1:1T)
    { linea<-paste(n,puntos.t[t], Denominador[t], sep="\t")
      for (i in 1:24) linea<-paste(linea, RE[t,i],sep="\t")
      for (i in 1:24) linea<-paste(linea, RBa[t,i],sep="\t")
      for (i in 1:24) linea<-paste(linea, RBr[t,i],sep="\t")
      write(linea,file=salidaError[t], ncolumns=1, append=T)
    }
} # end del for (n in vector.n)
## EJEMPLOS ##
#SimulaFdPemle.v4(NombreP=c("Fam1500Yx1"),B=100,y=Fam1500[,1],
#x=t(Fam1500[,2:3]),vector.n=seq(50,200,by=50),vector.beta=c(0.25,0.5,0.75),
#beta0=0.5)
#SimulaFdPemle.v4(NombreP=c("Murthydec"),B=1000,y=Murthy[,1],
#x=t(Murthy[,2:3]),vector.n=seq(42,40,by=-2),
#vector.beta=c(0.25,0.75,seq(0.1,0.9,by=0.1)),beta0=0.5)
}

```

C.3. Aportaciones a la estimación de cuantiles

C.3.1. Muestreo en dos ocasiones sucesivas con probabilidades desiguales

```
### ALGUNOS DISEÑOS MUESTRALES CON PROBABILIDADES DESIGUALES ###
```

```

muestra.2ocas.M.M.M<-function(N, tam1, tam2, m)
# Función que devuelve 2 muestras seleccionadas según muestreo en 2 ocasiones.
# Al principio de cada vector aparece la parte apareada y a #continuación la
# parte que se renueva. También se proporciona las #probabilidades de
# inclusión. Las muestras se extraen mediante M.M.M = Mas.Mas.Mas:
# Muestra 1: tamaño n´= tam1. Mediante MAS
# Muestra 2: tamaño m. Mediante MAS
# Muestra 3: tamaño u = tam2-m. Mediante MAS
# Los parámetros de entrada son:
# "N": es el tamaño de la población o marco desde el que generamos la muestra
# "tam1": es el tamaño muestral en la primera ocasión
# "tam2": es el tamaño muestral en la segunda ocasión
# "m": es el número de unidades que se mantienen
{
  poblacion <- c(1:N)
  muestral <- sample(N,tam1)
  muestra2apareada <- sample(tam1,m)
  muestra2apareada <- muestral[muestra2apareada]
}

```

```

u <- tam2 - m
# calculo la población nueva desde la que se genera la parte no apareada
poblacionrestante <- setdiff(poblacion, muestra1)
muestra2nueva <- sample(N-tam1,u)
muestra2nueva <- poblacionrestante[muestra2nueva]
#MUESTRA2: en la primera parte de muestra2 aparece la parte apareada
#y a continuación la nueva
muestra2 <- c(muestra2apareada, muestra2nueva)
#MUESTRA1: en la primera parte ponemos la parte apareada
#y a continuación la que se renueva
muestralrenovar <- setdiff(muestra1, muestra2apareada)
muestra1 <- c(muestra2apareada, muestralrenovar)
### Calculo de las probabilidades de Inclusión
#### De la Muestra 1.
PItam1<-vector(len=tam1)
for (i in 1:tam1) PItam1[i]<-tam1/N
#### De la Muestra u
PIu<-vector(len=u)
for (i in 1:u) PIu[i]<-u/N
#### De la Muestra m
PIm<-vector(len=m)
for (i in 1:m) PIm[i]<-m/N
resultado <- list(muestra1 = muestra1, muestra2 = muestra2, muestram=
muestra2apareada, muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
resultado
}

muestra.2ocas.M.L.M<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante M.L.M = Mas.Lahiri.Mas:
# Muestra 1: tamaño n´= tam1. Mediante MAS
# Muestra 2: tamaño m. Mediante Lahiri
# Muestra 3: tamaño u = tam2-m. Mediante MAS
{
N1<-length(z)
if (N != N1) stop("La longitud del vector z es dintinta a N")
poblacion <- c(1:N)
muestra1 <- sample(N,tam1)
datos.z.muestra1<-z[muestra1]
muestra2apareada <- lahiri(m,datos.z.muestra1)
muestra2apareada <- muestra1[muestra2apareada]
u <- tam2 - m
# calculo la población nueva desde la que se genera la parte no apareada
poblacionrestante <- setdiff(poblacion, muestra1)
muestra2nueva <- sample(N-tam1,u)
muestra2nueva <- poblacionrestante[muestra2nueva]
muestra2 <- c(muestra2apareada, muestra2nueva)
muestralrenovar <- setdiff(muestra1, muestra2apareada)
muestra1 <- c(muestra2apareada, muestralrenovar)
### Calculo de las probabilidades de Inclusión
datos.z.m<-z[muestra2apareada]
#### De la Muestra 1.
PItam1<-vector(len=tam1)
for (i in 1:tam1) PItam1[i]<-tam1/N
#### De la Muestra u
PIu.Condicionado<- u/(N-tam1)
PIu.Complementario<-1-(tam1/N)
PIu<-vector(len=u)
for (i in 1:u) PIu[i]<-PIu.Complementario*PIu.Condicionado
#### De la Muestra m
PIm.Partel<-PItam1[1:m]
PIcondicionada<-vector(len=m)
sum(datos.z.muestra1)->TotalZtam1

```

```

        for (i in 1:m) PIcondicionada[i]<-m*datos.z.m[i]/TotalZtam1
        PIm<-PIm.Partel*PIcondicionada
resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
muestra2apareada,muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
        resultado
}

```

```

muestra.2ocas.L.L.L<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante L.L.L = Lahiri.Lahiri.Lahiri:
# Muestra 1: tamaño n´= tam1. Mediante Lahiri
# Muestra 2: tamaño m. Mediante Lahiri
# Muestra 3: tamaño u = tam2-m. Mediante Lahiri
{
    N1<-length(z)
    if (N != N1) stop("La longitud del vector z es dintinta a N")
    poblacion <- c(1:N)
    muestral <- lahiri(tam1,z)
    datos.z.muestral<-z[muestral]
    muestra2apareada <- lahiri(m,datos.z.muestral)
    muestra2apareada <- muestral[muestra2apareada]
    u <- tam2 - m
    # calculo la población nueva desde la que se genera la parte no apareada
    poblacionrestante <- setdiff(poblacion, muestral)
    datos.z.restante<-z[poblacionrestante]
    Eti.muestra2nueva <-lahiri(u,datos.z.restante)
    muestra2nueva <- poblacionrestante[Eti.muestra2nueva]
    muestra2 <- c(muestra2apareada, muestra2nueva)
    muestralrenovar <- setdiff(muestral, muestra2apareada)
    muestral <- c(muestra2apareada, muestralrenovar)
    ### Calculo de las probabilidades de Inclusión
    datos.z.muestral<-z[muestral]
    datos.z.u<-z[muestra2nueva]
    datos.z.m<-z[muestra2apareada]
    #####De la muestra 1.
    PItam1<-vector(len=tam1)
    TotalZ<-sum(z)
    for (i in 1:tam1) PItam1[i]<-tam1*datos.z.muestral[i]/TotalZ
    ##### De la muestra u.
    PIu.Condicionado<-vector(len=u)
    PIu.Complementario<-vector(len=u)
    TotalZrestante<-sum(datos.z.restante)
    for (i in 1:u)
        { PIu.Condicionado[i]<-u*datos.z.u[i]/TotalZrestante
          PIu.Complementario[i]<-1-(tam1*datos.z.u[i]/TotalZ)
        }
    PIu<-PIu.Complementario*PIu.Condicionado
    ##### De la Muestra m
    PIm.Condicionada<-vector(len=m)
    sum(datos.z.muestral)->TotalZtam1
    for (i in 1:m) PIm.Condicionada[i]<-m*datos.z.m[i]/TotalZtam1
    PIm.Partel<-PItam1[1:m]
    PIm<-PIm.Partel*PIm.Condicionada
    ### Salida del Programa
resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
muestra2apareada,muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
        resultado
}

```

```

muestra.2ocas.L.M.L<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante L.M.L = Lahiri.MAS.Lahiri:
# Muestra 1: tamaño n´= tam1. Mediante Lahiri
# Muestra 2: tamaño m. Mediante MAS

```

```

# Muestra 3: tamaño u = tam2-m. Mediante Lahiri
{
  N1<-length(z)
  if (N != N1) stop("La longitud del vector z es diferente a N")
  poblacion <- c(1:N)
  muestral <- lahiri(tam1,z)
  muestra2apareada <- sample(tam1, m)
  muestra2apareada <- muestral[muestra2apareada]
  u <- tam2 - m
  # calculo la población nueva desde la que se genera la parte no apareada
  poblacionrestante <- setdiff(poblacion, muestral)
  datos.z.restante<-z[poblacionrestante]
  Eti.muestra2nueva <-lahiri(u,datos.z.restante)
  muestra2nueva <- poblacionrestante[Eti.muestra2nueva]
  muestra2 <- c(muestra2apareada, muestra2nueva)
  muestralrenovar <- setdiff(muestral, muestra2apareada)
  muestral <- c(muestra2apareada, muestralrenovar)
###Calculo de las probabilidades de Inclusión
  datos.z.muestral<-z[muestral]
  datos.z.u<-z[muestra2nueva]
####De la muestra 1.
  PItam1<-vector(len=tam1)
  TotalZ<-sum(z)
  for (i in 1:tam1) PItam1[i]<-tam1*datos.z.muestral[i]/TotalZ
#### De la muestra u.
  PIu.Condicionado<-vector(len=u)
  PIu.Complementario<-vector(len=u)
  TotalZrestante<-sum(datos.z.restante)
  for (i in 1:u)
    { PIu.Condicionado[i]<-u*datos.z.u[i]/TotalZrestante
      PIu.Complementario[i]<-1-(tam1*datos.z.u[i]/TotalZ)
    }
  PIu<-PIu.Complementario*PIu.Condicionado
#### De la Muestra m
  PIm.Condicionada<-m/tam1
  PIm.Partel<-PItam1[1:m]
  PIm<-PIm.Partel*PIm.Condicionada
#### Salida del Programa
resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
muestra2apareada,muestra2nueva=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
  resultado
}

muestra.2ocas.L.M.M<-funcion(N, z, tam1, tam2, m)
# Las muestras se extraen mediante L.M.L = Lahiri.MAS.MAS:
# Muestra 1: tamaño n´= tam1. Mediante Lahiri
# Muestra 2: tamaño m. Mediante MAS
# Muestra 3: tamaño u = tam2-m. Mediante MAS
{
  N1<-length(z)
  if (N != N1) stop("La longitud del vector z es diferente a N")
  poblacion <- c(1:N)
  muestral <- lahiri(tam1,z)
  muestra2apareada <- sample(tam1, m)
  muestra2apareada <- muestral[muestra2apareada]
  u <- tam2 - m
# calculo la población nueva desde la que se genera la parte no apareada
  poblacionrestante <- setdiff(poblacion, muestral)
  Eti.muestra2nueva <- sample(N-tam1,u)
  muestra2nueva <- poblacionrestante[Eti.muestra2nueva]
  muestra2 <- c(muestra2apareada, muestra2nueva)
  muestralrenovar <- setdiff(muestral, muestra2apareada)
  muestral <- c(muestra2apareada, muestralrenovar)
}

```

```

### Calculo de las probabilidades de Inclusión
  datos.z.muestral<-z[muestral]
  datos.z.u<-z[muestra2nueva]
### De la Muestral
  PItam1<-vector(len=tam1)
  TotalZ<-sum(z)
  for (i in 1:tam1) PItam1[i]<-tam1*datos.z.muestral[i]/TotalZ
### De la Muestra u
  PIu.Condicionado<-u/(N-tam1)
  PIu.Complementario<-vector(len=u)
  for (i in 1:u) PIu.Complementario[i]<-1-(tam1*datos.z.u[i]/TotalZ)
  PIu<-PIu.Complementario*PIu.Condicionado
### De la Muestra m
  PIm.Partel<-PItam1[1:m]
  PIm.Condicionada<-m/tam1
  PIm<-PIm.Partel*PIm.Condicionada
resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
muestra2apareada,muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
  resultado
}

muestra.2ocas.M.Z.M<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante M.Z.M = MAS.Midzuno.MAS:
# Muestra 1: tamaño n´= tam1. Mediante MAS
# Muestra 2: tamaño m. Mediante Midzuno
# Muestra 3: tamaño u = tam2-m. Mediante MAS
{
  N1<-length(z)
  if (N != N1) stop("La longitud del vector z es dintinta a N")
  poblacion <- c(1:N)
  muestral <- sample(N,tam1)
  datos.z.muestral<-z[muestral]
  muestra2apareada <- Etiquetas.midzuno(m,datos.z.muestral)
  muestra2apareada <- muestral[muestra2apareada]
  u <- tam2 - m
  # calculo la población nueva desde la que se genera la parte no apareada
  poblacionrestante <- setdiff(poblacion, muestral)
  muestra2nueva <- sample(N-tam1,u)
  muestra2nueva <- poblacionrestante[muestra2nueva]
  muestra2 <- c(muestra2apareada, muestra2nueva)
  muestralrenovar <- setdiff(muestral, muestra2apareada)
  muestral <- c(muestra2apareada, muestralrenovar)
### Calculo de las probabilidades de Inclusión
  datos.z.m<-z[muestra2apareada]
  PItam1<-vector(len=tam1)
  for (i in 1:tam1) PItam1[i]<-tam1/N
#### De la Muestra u
  PIu.Condicionado<- u/(N-tam1)
  PIu.Complementario<-1-(tam1/N)
  PIu<-vector(len=u)
  for (i in 1:u) PIu[i]<-PIu.Complementario*PIu.Condicionado
#### De la Muestra m
  PIm.Partel<-PItam1[1:m]
  PIcondicionada<-vector(len=m)
  sum(datos.z.muestral)->TotalZtam1
for (i in 1:m)
PIcondicionada[i]<-((tam1-m)/(tam1-1))*(datos.z.m[i]/TotalZtam1)+(m-1)/(tam1-1)
  PIm<-PIm.Partel*PIcondicionada
resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
muestra2apareada,muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
  resultado
}

```

```

muestra.2ocas.Z.M.M<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante Z.M.M = Midzuno.MAS.MAS:
# Muestra 1: tamaño n´= tam1. Mediante Midzuno
# Muestra 2: tamaño m. Mediante MAS
# Muestra 3: tamaño u = tam2-m. Mediante MAS
{
  N1<-length(z)
  if (N != N1) stop("La longitud del vector z es dintinta a N")
  poblacion <- c(1:N)
  muestral <- Etiquetas.midzuno(tam1,z)
  muestra2apareada <- sample(tam1, m)
  muestra2apareada <- muestral[muestra2apareada]
  u <- tam2 - m
  # calculo la población nueva desde la que se genera la parte no apareada
  poblacionrestante <- setdiff(poblacion, muestral)
  Eti.muestra2nueva <- sample(N-tam1,u)
  muestra2nueva <- poblacionrestante[Eti.muestra2nueva]
  muestra2 <- c(muestra2apareada, muestra2nueva)
  muestralrenovar <- setdiff(muestral, muestra2apareada)
  muestral <- c(muestra2apareada, muestralrenovar)
  ### Calculo de las probabilidades de Inclusión
  datos.z.muestral<-z[muestral]
  datos.z.u<-z[muestra2nueva]
  ##### De la Muestral
  PItam1<-vector(len=tam1)
  TotalZ<-sum(z)
  for (i in 1:tam1) PItam1[i]<-((N-tam1)/(N-1) )*(datos.z.muestral[i]/TotalZ)+
    (tam1-1)/(N-1)
  ##### De la Muestra u
  PIu.Condicionado<-u/(N-tam1)
  PIu.Complementario<-vector(len=u)
  for (i in 1:u) PIu.Complementario[i]<-
  1-(((N-tam1)/(N-1))*(datos.z.u[i]/TotalZ)+(tam1-1)/(N-1) )
  PIu<-PIu.Complementario*PIu.Condicionado
  ##### De la Muestra m
  PIm.Partel<-PItam1[1:m]
  PIm.Condicionada<-m/tam1
  PIm<-PIm.Partel*PIm.Condicionada
  resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
  muestra2apareada,muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
  resultado
}

```

```

muestra.2ocas.Z.M.Z<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante Z.M.M = Midzuno.MAS.Midzuno:
# Muestra 1: tamaño n´= tam1. Mediante Midzuno
# Muestra 2: tamaño m. Mediante MAS
# Muestra 3: tamaño u = tam2-m. Mediante Midzuno
{
  N1<-length(z)
  if (N != N1) stop("La longitud del vector z es dintinta a N")
  poblacion <- c(1:N)
  muestral <- Etiquetas.midzuno(tam1,z)
  muestra2apareada <- sample(tam1, m)
  muestra2apareada <- muestral[muestra2apareada]
  u <- tam2 - m
  # calculo la población nueva desde la que se genera la parte no apareada
  poblacionrestante <- setdiff(poblacion, muestral)
  datos.z.restante<-z[poblacionrestante]
  Eti.muestra2nueva <-Etiquetas.midzuno(u,datos.z.restante)
  muestra2nueva <- poblacionrestante[Eti.muestra2nueva]
  muestra2 <- c(muestra2apareada, muestra2nueva)
  muestralrenovar <- setdiff(muestral, muestra2apareada)

```

```

    muestral <- c(muestra2apareada, muestralrenovar)
### Calculo de las probabilidades de Inclusión
    datos.z.muestral<-z[muestral]
    datos.z.u<-z[muestra2nueva]
###De la muestra 1.
    PItam1<-vector(len=tam1)
    TotalZ<-sum(z)
for (i in 1:tam1) PItam1[i]<-((N-tam1)/(N-1) )*(datos.z.muestral[i]/TotalZ)+
    (tam1-1)/(N-1)

#### De la muestra u.
    PIu.Condicionado<-vector(len=u)
    PIu.Complementario<-vector(len=u)
    TotalZrestante<-sum(datos.z.restante)
    for (i in 1:u)
    {
PIu.Condicionado[i]<-
    ( (N-tam1-u)/(N-tam1-1) )*(datos.z.u[i]/TotalZrestante)+( (u-1)/(N-tam1-1) )
PIu.Complementario[i]<-
    1-((N-tam1)/(N-1) )*(datos.z.u[i]/TotalZ)+(tam1-1)/(N-1) )
    }
    PIu<-PIu.Complementario*PIu.Condicionado
#### De la Muestra m
    PIm.Condicionada<-m/tam1
    PIm.Partel<-PItam1[1:m]
    PIm<-PIm.Partel*PIm.Condicionada
#### Salida del Programa
resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
muestra2apareada,muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
    resultado
}

muestra.2ocas.Z.Z.Z<-function(N, z, tam1, tam2, m)
# Las muestras se extraen mediante Z.M.M = Midzuno.Midzuno.Midzuno:
# Muestra 1: tamaño n´= tam1. Mediante Midzuno
# Muestra 2: tamaño m. Mediante Midzuno
# Muestra 3: tamaño u = tam2-m. Mediante Midzuno
{
    N1<-length(z)
    if (N != N1) stop("La longitud del vector z es dintinta a N")
    poblacion <- c(1:N)
    muestral <- Etiquetas.midzuno(tam1,z)
    datos.z.muestral<-z[muestral]
    muestra2apareada <- Etiquetas.midzuno(m,datos.z.muestral)
    muestra2apareada <- muestra1[muestra2apareada]
    u <- tam2 - m

# calculo la población nueva desde la que se genera la parte no apareada
    poblacionrestante <- setdiff(poblacion, muestra1)
    datos.z.restante<-z[poblacionrestante]
    Eti.muestra2nueva <-Etiquetas.midzuno(u,datos.z.restante)
    muestra2nueva <- poblacionrestante[Eti.muestra2nueva]
    muestralrenovar <- setdiff(muestral, muestra2apareada)
    muestral <- c(muestra2apareada, muestralrenovar)

#### Calculo de las probabilidades de Inclusión
    datos.z.muestral<-z[muestral]
    datos.z.u<-z[muestra2nueva]
    datos.z.m<-z[muestra2apareada]

### De la muestra 1.
    PItam1<-vector(len=tam1)
    TotalZ<-sum(z)
    for (i in 1:tam1)
PItam1[i]<-((N-tam1)/(N-1) )*(datos.z.muestral[i]/TotalZ)+(tam1-1)/(N-1)
#### De la muestra u.

```

```

    PIu.Condicionado<-vector(len=u)
    PIu.Complementario<-vector(len=u)
    TotalZrestante<-sum(datos.z.restante)
    for (i in 1:u)
    {
        PIu.Condicionado[i]<-
        ((N-tam1-u)/(N-tam1-1))*(datos.z.u[i]/TotalZrestante)+((u-1)/(N-tam1-1))
        PIu.Complementario[i]<-
        1-((N-tam1)/(N-1))*(datos.z.u[i]/TotalZ)+(tam1-1)/(N-1)
    }
    PIu<-PIu.Complementario*PIu.Condicionado
## De la Muestra m
    PIm.Partel<-PItam1[1:m]
    PIcondicionada<-vector(len=m)
    sum(datos.z.muestral)->TotalZtam1
    for (i in 1:m) PIcondicionada[i]<-
    ((tam1-m)/(tam1-1))*(datos.z.m[i]/TotalZtam1)+(m-1)/(tam1-1)
    PIm<-PIm.Partel*PIcondicionada
#### Salida del Programa
    resultado <- list(muestral = muestral, muestra2 = muestra2, muestram=
    muestra2apareada,muestrau=muestra2nueva, PItam1 = PItam1, PIm = PIm, PIu=PIu )
    resultado
}

muestra.2ocas<-funcion(TIPO,N, z, tam1, tam2, m)
# En función del valor asignado a "TIPO", realiza los siguientes muestreos:
# TIPO=1: M.M.M, TIPO=2: M.L.M, TIPO=3: L.M.M, TIPO=4: L.M.L, TIPO=5: L.L.L,
# TIPO=6: M.Z.M, TIPO=7: Z.M.M, TIPO=8: Z.M.Z, TIPO=9: Z.Z.Z
{
switch(TIPO,
    resultado<- muestra.2ocas.M.M.M(N, tam1, tam2, m),
    resultado<- muestra.2ocas.M.L.M(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.L.M.M(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.L.M.L(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.L.L.L(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.M.Z.M(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.Z.M.M(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.Z.M.Z(N, z, tam1, tam2, m),
    resultado<- muestra.2ocas.Z.Z.Z(N, z, tam1, tam2, m)
)
resultado
}

estimadoresJackknife<-funcion(Tipo.Fd,Xm,Ym,PIm,QgorroX, beta, fx, fy, N)
# Obtiene estimaciones Jackknife para las varianzas de los estimadores
{
    tam <- length(Ym)
    xi <- double(tam - 1)
    yi <- double(tam - 1)
    Pii <-double(tam - 1)
    estimi.diferencia<-double(tam)
    estimi.razon<-double(tam)
    for (i in 1:tam)
    {
        if (i==1)
        {
            for (j in 1:(tam-1))
            {
                xi[j]<-Xm[j+1]
                yi[j]<-Ym[j+1]
            }
            Pii[j]<-PIm[j+1]
        }
        else if(i==tam)
        {
            for (j in 1:(tam-1))

```

```

                                {      xi[j]<-Xm[j]
                                {      yi[j]<-Ym[j]
                                {      PIi[j]<-PIIm[j]
                                }
                                }
                                }
else
{
  for (j in 1:(i-1))
  {
    xi[j]<-Xm[j]
    yi[j]<-Ym[j]
    PIi[j]<-PIIm[j]
  }
  for (j in (i+1):tam)
  {
    xi[j-1]<-Xm[j]
    yi[j-1]<-Ym[j]
    PIi[j-1]<-PIIm[j]
  }
}
QgorroXi<-cuantil.HT.Rapido(Tipo.Fd, beta, N, PIi, xi)
QgorroYi<-cuantil.HT.Rapido(Tipo.Fd, beta, N, PIi, yi)
### Para el calculo de QgorroDiferencia
tami <- length(yi)
suma <- 0
for(h in 1:tami)
  {
    if((xi[h] <= QgorroXi) && (yi[h] <= QgorroYi))
    {
      suma <- suma + 1
    }
  }
p1lx <- suma/tami
dxopt <- (fx/fy) * (p1lx/(beta * (1 - beta)) - 1)
estimi.diferencia[i] <- QgorroYi+dxopt*(QgorroX-QgorroXi)
estimi.razon[i]<-QgorroYi*(QgorroX/QgorroXi)
}
### Cálculo de la estimación jackknife del SE según Efron
QgorroXm<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PIIm,Xm)
QgorroYm<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PIIm,Ym)
tam <- length(Ym)
suma <- 0
for(i in 1:tam)
  {
    if((Xm[i] <= QgorroXm) && (Ym[i] <= QgorroYm))
    suma <- suma + 1
  }
p1lx <- suma/tam
dxopt <- (fx/fy) * (p1lx/(beta * (1 - beta)) - 1)
QgorroDifm <- QgorroYm + dxopt * (QgorroX - QgorroXm)
vari.diferencia<-((tam-1)/tam)*sum( (estimi.diferencia-QgorroDifm)^2 )
desvi.diferencia<-sum(estimi.diferencia-QgorroDifm)
QgorroRazon<-QgorroYm*(QgorroX/QgorroXm)
vari.razon<-((tam-1)/tam)*sum( (estimi.razon-QgorroRazon)^2 )
desvi.razon<-sum(estimi.razon-QgorroRazon)
list(vari.Dif=vari.diferencia, vari.Razon=vari.razon, QgorroDifm=QgorroDifm,
QgorroRazon=QgorroRazon, desvi.Razon=desvi.razon,desvi.Dif=desvi.diferencia)
}
estimJackknife.directo.HT.Rapido<-function(Tipo.Fd,Yu, PIu, beta, N)
# Obtiene estimaciones Jackknife para las varianzas del estimador directo
{
  tam <- length(Yu)
  yi <- double(tam - 1)
  PIi <-double(tam - 1)
  estimi <- double(tam)
  for (i in 1:tam)

```

```

{
  if (i==1)
  {
    for (j in 1:(tam-1))
    {
      yi[j]<-Yu[j+1]
      PIi[j]<-PIu[j+1]
    }
  }
  else if(i==tam)
  {
    for (j in 1:(tam-1))
    {
      yi[j]<-Yu[j]
      PIi[j]<-PIu[j]
    }
  }
  else
  {
    for (j in 1:(i-1))
    {
      yi[j]<-Yu[j]
      PIi[j]<-PIu[j]
    }
    for (j in (i+1):tam)
    {
      yi[j-1]<-Yu[j]
      PIi[j-1]<-PIu[j]
    }
  }
  estimi[i] <- cuantil.HT.Rapido(Tipo.Fd,beta,N,PIi,yi)
}
QgorroYu <- cuantil.HT.Rapido(Tipo.Fd,beta,N,PIu,Yu)
vari<-((tam-1)/tam)*sum( (estimi-QgorroYu)^2 )
desvi.directo= sum(estimi-QgorroYu)
list(vari=vari, QgorroYu=QgorroYu, desvi.directo=desvi.directo)
}

```

Datos.X.HT.Rapido<-

```

function(Tipo.Muestreo,Tipo.Fd,datosx,datosy,z, N, tam1, tam2, m, beta, Qy, Qx)
# Obtiene información necesaria para calcular estimadores y sus varianzas.
# "Tipo.Muestreo" es un valor entre 1 y 9, y "Tipo.Fd" es para usar el
# estimador de tipo HT (1) o de tipo Hájek (2). "datosx" son los datos de la
# primera ocasión, "datosy" son los datos de la segunda ocasión y "z" es un
# vector de longitud "N" de la variable auxiliar.
{
  muestras <- muestra.2ocas(Tipo.Muestreo,N, z, tam1, tam2, m)
  muestral <- muestras$muestral
  muestra2 <- muestras$muestra2
  muestrau<-muestras$muestrau
  muestram<-muestras$muestram
  Xm <- double(m)
  Ym <- double(m)
  Xm<-datosx[muestram]
  Ym<-datosy[muestram]
  X <- double(tam1)
  X<-datosx[muestral]
  u <- tam2 - m
  Yu <- double(u)
  Yu<-datosy[muestrau]
  fx <- densidad(datosx, Qx) ### estimación de la densidad: fx, fy
  fy <- densidad(datosy, Qy)
### Probabilidades de Inclusión y QgorroX
  PITam1<-muestras$PITam1
  PIm<-muestras$PIm
  PIu<-muestras$PIu
  QgorroX<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PITam1,X)
resultado <- list(X = X, Xm = Xm, Ym = Ym, Yu = Yu, fx = fx, fy = fy,
  QgorroX=QgorroX, PITam1=PITam1, PIm=PIm, PIu=PIu )
  resultado
}

```

```

EstimadoresDyR.CovRepe<-function(Tipo.Fd,Xm,Ym,PIm,QgorroX,Yu,PIu,beta,fx,fy,N)
# Obtiene los estimadores propuestos de tipo razón y diferencia.
{
  estimJackknife.directo.HT.Rapido(Tipo.Fd,Yu, PIu, beta, N)->A
  estimadoresJackknife(Tipo.Fd,Xm,Ym,PIm,QgorroX, beta, fx, fy, N)->B
  A$vari->vari.directo
  A$QgorroYu->QgorroYu
  A$desvi.directo->desvi.directo
  B$vari.Dif->vari.diferencia
  B$vari.Razon->vari.razon
  B$QgorroDifm->QgorroDifm
  B$QgorroRazon->QgorroRazon
  B$desvi.Dif->desvi.Dif
  B$desvi.Razon->desvi.Razon
  #calculamos el W óptimo usando estimaciones jackknife de las varianzas
  VqYu <- vari.directo
  VqYdif <- vari.diferencia
## FORMA 1
  W.d <- VqYu/(VqYu + VqYdif)
  if (is.na(W.d)) {W.d<-0.5}
  estimador.D <- QgorroDifm * W.d + (1 - W.d) * QgorroYu
## FORMA 2 (Con covarianzas)
  m<-length(Ym)
  u<-length(Yu)
  CovDif.estimada<- ( (m*u-1)/(m*u) )*desvi.directo*desvi.Dif
  W.d.Cov <- (VqYu-CovDif.estimada)/(VqYu + VqYdif-2*CovDif.estimada)
  if (is.na(W.d.Cov)) { W.d.Cov<-99}
  if (W.d.Cov < 0) { W.d.Cov<-99}
  if (W.d.Cov > 1) { W.d.Cov<-99}
  estimador.D.Cov <- QgorroDifm * W.d.Cov + (1 - W.d.Cov) * QgorroYu
  VqYr <- vari.razon
## FORMA 1
  W.r <- VqYu/(VqYu + VqYr)
  if (is.na(W.r)) {W.r<-0.5}
  estimador.R <- QgorroRazon * W.r + (1 - W.r) * QgorroYu
## FORMA 2 (Con covarianzas)
  CovR.estimada<- ( (m*u-1)/(m*u) )*desvi.directo*desvi.Razon
  W.r.Cov <- (VqYu-CovR.estimada)/(VqYu + VqYr-2*CovR.estimada)
  if (is.na(W.r.Cov)) {W.r.Cov<-99}
  if (W.r.Cov < 0) { W.r.Cov<-99}
  if (W.r.Cov >1) { W.r.Cov<-99}
  estimador.R.Cov <- QgorroRazon * W.r.Cov + (1 - W.r.Cov) * QgorroYu
list(estimador.D=estimador.D, estimador.R=estimador.R,
  estimador.D.Cov=estimador.D.Cov, estimador.R.Cov=estimador.R.Cov,
  W.d.Cov=W.d.Cov, W.r.Cov=W.r.Cov, W.r=W.r, W.d=W.d)
}

Simula.2ocas.CovySinRepe<-function(NombreP, Tipo.Muestreo, Tipo.Fd, DatosTodos,
variableY, variableX, variableZ, tam1, tam2, m, replicas, beta)
## Programa que realiza el proceso de simulación
{
  N<-length(DatosTodos[,1])
  variables<-length(DatosTodos[1,])
  z<-DatosTodos[,variableZ]
  datosx<-DatosTodos[,variableX]
  datosy<-DatosTodos[,variableY]
  salidaError <- paste("D:\\OcasionesSucesivas\\err",NombreP,"_tam", tam1, "y",
  tam2, "beta_", beta, ".txt", sep = "")
  salidaEstimador <- paste("D:\\OcasionesSucesivas\\est",NombreP,"_m", m,
  "tam", tam1, "y", tam2, "beta_", beta, ".txt", sep = "")
  Qy<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,vectory=datosy)
  Qx<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,vectory=datosx)
  ecmDirecto <- 0
}

```

```

ecmDirectoSinRepe <- 0
ecmRatio <- 0
ecmDiferenciaJack <- 0
ecmRatio.Cov <- 0
ecmDiferenciaJack.Cov <- 0
ecmRatioSinRepe <- 0
ecmDiferenciaJackSinRepe <- 0
sesDirecto <- 0
sesRatio <- 0
sesDiferenciaJack <- 0
sesRatio.Cov <- 0
sesDiferenciaJack.Cov <- 0
sesRatioSinRepe <- 0
sesDiferenciaJackSinRepe <- 0
##### Sesgos Relativos
ResesDirecto <- 0
ResesRatio <- 0
ResesDiferenciaJack <- 0
ResesRatio.Cov <- 0
ResesDiferenciaJack.Cov <- 0
ResesRatioSinRepe <- 0
ResesDiferenciaJackSinRepe <- 0
for(b in 1:replicas)
  { W<-1
    Repeticiones<-0
    while (W !=0)
      { W<-0
        Repeticiones<-Repeticiones+1
        datos <- Datos.X.HT.Rapido(Tipo.Muestreo,Tipo.Fd, datosx, datosy,
z, N, tam1, tam2, m, beta,Qy,Qx)
        Xm <- datos$Xm
        Ym <- datos$Ym
        X <- datos$X
        fx <- datos$fx
        fy <- datos$fy
        Yu <- datos$Yu
        QgorroX<-datos$QgorroX
        PItam1<-datos$PItam1
        PIm<-datos$PIm
        PIu<-datos$PIu
EstimadoresDyR.CovRepe(Tipo.Fd,Xm,Ym,PIm, QgorroX, Yu, PIu, beta, fx, fy, N)->
Estimadores
        if (Repeticiones == 1)
          {
# el estimador directo Sin Repetir Muestra. Se obtiene solo el ecmDirecto
            Ytotal <- c(Ym, Yu)
#####SE CONSIDERA OBTENIDO MEDIANTE "MAS".
            PIn<-vector(len=tam2)
            for (i in 1:tam2) PIn[i]<-tam2/N
            estimadorDirecto <- cuantil.HT.Rapido(Tipo.Fd,beta, N, PIn, Ytotal)
            ecmDirectoSinRepe <- ecmDirectoSinRepe + (estimadorDirecto - Qy)^2
##### cálculos para el estimador de razón Sin repetir muestra
            estimadorRatioSinRepe <- Estimadores$estimador.R
            ecmRatioSinRepe <- ecmRatioSinRepe + (estimadorRatioSinRepe - Qy)^2
            sesRatioSinRepe <- sesRatioSinRepe + abs(estimadorRatioSinRepe - Qy)
            ResesRatioSinRepe <- ResesRatioSinRepe + (estimadorRatioSinRepe - Qy)
###cálculos para el estimador de diferencia Sin repetir muestra
            estimadorDiferenciaJackSinRepe <- Estimadores$estimador.D
            ecmDiferenciaJackSinRepe <- ecmDiferenciaJackSinRepe +
(estimadorDiferenciaJackSinRepe - Qy)^2

```

```

        sesDiferenciaJackSinRepe <- sesDiferenciaJackSinRepe +
abs(estimadorDiferenciaJackSinRepe - Qy)
        ResesDiferenciaJackSinRepe <- ResesDiferenciaJackSinRepe +
(estimadorDiferenciaJackSinRepe - Qy)
    }
    if ( Estimadores$W.d.Cov == 99) W<-1
    if ( Estimadores$W.r.Cov == 99) W<-1
    } ## end while (W !=0)
#### el estimador directo para comparar:
        Ytotal <- c(Ym, Yu)
        #####SE CONSIDERA OBTENIDO MEDIANTE "MAS".
        PIn<-vector(len=tam2)
        for (i in 1:tam2) PIn[i]<-tam2/N
        estimadorDirecto <- cuantil.HT.Rapido(Tipo.Fd,beta, N, PIn, Ytotal)
        ecmDirecto <- ecmDirecto + (estimadorDirecto - Qy)^2
        sesDirecto <- sesDirecto + abs(estimadorDirecto - Qy)
        ResesDirecto <- ResesDirecto + (estimadorDirecto - Qy)
#####cálculos para el estimador de razón
        estimadorRatio <- Estimadores$estimador.R
        ecmRatio <- ecmRatio + (estimadorRatio - Qy)^2
        sesRatio <- sesRatio + abs(estimadorRatio - Qy)
        ResesRatio <- ResesRatio + (estimadorRatio - Qy)
### cálculos para el estimador de razón (v.Covarianzas)
        estimadorRatio.Cov <- Estimadores$estimador.R.Cov
        ecmRatio.Cov <- ecmRatio.Cov + (estimadorRatio.Cov - Qy)^2
        sesRatio.Cov <- sesRatio.Cov + abs(estimadorRatio.Cov - Qy)
        ResesRatio.Cov <- ResesRatio.Cov + (estimadorRatio.Cov - Qy)
### cálculos para el estimador de diferencia
        estimadorDiferenciaJack <- Estimadores$estimador.D
        ecmDiferenciaJack <- ecmDiferenciaJack + (estimadorDiferenciaJack - Qy)^2
        sesDiferenciaJack <- sesDiferenciaJack +
abs(estimadorDiferenciaJack - Qy)
        ResesDiferenciaJack <- ResesDiferenciaJack +
(estimadorDiferenciaJack - Qy)
### cálculos para el estimador de diferencia (v.Covarianzas)
        estimadorDiferenciaJack.Cov <- Estimadores$estimador.D.Cov
        ecmDiferenciaJack.Cov <- ecmDiferenciaJack.Cov +
(estimadorDiferenciaJack.Cov - Qy)^2
        sesDiferenciaJack.Cov <- sesDiferenciaJack.Cov +
abs(estimadorDiferenciaJack.Cov - Qy)
        ResesDiferenciaJack.Cov <- ResesDiferenciaJack.Cov +
(estimadorDiferenciaJack.Cov - Qy)
#### escribimos el valor de los estimadores
        linea <- paste(b, Qy, estimadorDirecto, estimadorRatio,
estimadorDiferenciaJack,estimadorRatio.Cov, estimadorDiferenciaJack.Cov,
        estimadorRatioSinRepe, estimadorDiferenciaJackSinRepe, sep = "\t")
        write(linea, file = salidaEstimador, ncolumns = 1, append = T)
    } ## end for (b in 1: replicas)
    ecmDirecto <- (ecmDirecto/replicas)
    RecmDirecto<- 1
    RecmRatio <- (ecmRatio/replicas)/ecmDirecto
    RecmDiferenciaJack <- (ecmDiferenciaJack/replicas)/ecmDirecto
    RecmRatio.Cov <- (ecmRatio.Cov/replicas)/ecmDirecto
    RecmDiferenciaJack.Cov <- (ecmDiferenciaJack.Cov/replicas)/ecmDirecto
    RecmRatioSinRepe <- (ecmRatioSinRepe/replicas)/ecmDirecto
    RecmDiferenciaJackSinRepe <- (ecmDiferenciaJackSinRepe/replicas)/ecmDirecto
    sesDirecto <- sesDirecto/(replicas * Qy)
    sesRatio <- sesRatio/(replicas * Qy)
    sesDiferenciaJack <- sesDiferenciaJack/(replicas * Qy)
    sesRatio.Cov <- sesRatio.Cov/(replicas * Qy)
    sesDiferenciaJack.Cov <- sesDiferenciaJack.Cov/(replicas * Qy)

```

```

sesRatioSinRepe <- sesRatioSinRepe/(replicas * Qy)
sesDiferenciaJackSinRepe <- sesDiferenciaJackSinRepe/(replicas * Qy)
##### Sesgos Relativos
ResesDirecto <- ResesDirecto/(replicas * Qy)
ResesRatio <- ResesRatio/(replicas * Qy)
ResesDiferenciaJack <- ResesDiferenciaJack/(replicas * Qy)
ResesRatio.Cov <- ResesRatio.Cov/(replicas * Qy)
ResesDiferenciaJack.Cov <- ResesDiferenciaJack.Cov/(replicas * Qy)
ResesRatioSinRepe <- ResesRatioSinRepe/(replicas * Qy)
ResesDiferenciaJackSinRepe <- ResesDiferenciaJackSinRepe/(replicas * Qy)
### vamos a escribir la matriz de resultados
linea <- paste(m,sesDirecto,ecmDirecto,ecmDirectoSinRepe,RecmDirecto,sesRatio,
RecmRatio,sesDiferenciaJack,RecmDiferenciaJack,sesRatio.Cov,RecmRatio.Cov,
sesDiferenciaJack.Cov,RecmDiferenciaJack.Cov,sesRatioSinRepe,RecmRatioSinRepe,
sesDiferenciaJackSinRepe, RecmDiferenciaJackSinRepe, ResesDirecto, ResesRatio,
ResesDiferenciaJack,ResesRatio.Cov, ResesDiferenciaJack.Cov,
ResesRatioSinRepe, ResesDiferenciaJackSinRepe, sep = "\t")
write(linea, file = salidaError, ncolumns = 1, append = T)
}

```

C.3.2. Muestreo en dos ocasiones sucesivas con múltiples variables auxiliares

Datos.X.Pvar.v3<-

```

function(Tipo.Muestreo,Tipo.Fd,datosx,datosy,z,N,tam1, tam2, m, beta, Qy, Qx)
# Obtiene datos necesarios para calcular los estimadores propuestos. "Qx" es un
# vector de longitud "p". "datosx" es una matriz con p columnas.
{
muestras <- muestra.2ocas(Tipo.Muestreo,N, z, tam1, tam2, m)
muestra1 <- muestras$muestra1
muestra2 <- muestras$muestra2
muestraU<-muestras$muestraU
muestraM<-muestras$muestraM
p<-length(datosx[1,])
##### Datos de la submuestra para las variables X
Xm<-matrix(nr=m, nc=p)
for (k in 1:p) Xm[,k] <- datosx[muestraM, k]
##### Datos de la submuestra para la variable Y
Ym<-datosy[muestraM]
##### Datos de la primera ocasion (para las variables X)
X<-matrix(nro=tam1, nco=p)
for (k in 1:p) X[,k]<-datosx[muestra1,k]
##### Datos de la muestra no solapada en la ocasión más reciente
u <- tam2 - m
Yu <- double(u)
Yu<-datosy[muestraU]
##### Probabilidades de Inclusión y QgorroX
PItam1<-muestras$PItam1
PIm<-muestras$PIm
PIu<-muestras$PIu
QgorroX<-vector(len=p)
for (k in 1:p) QgorroX[k]<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PItam1,X[,k])
Ytotal <- c(Ym, Yu)
### SE CONSIDERA OBTENIDO MEDIANTE "MAS".
PIn<-vector(len=tam2)
for (i in 1:tam2) PIn[i]<-tam2/N
estimadorDirecto <- cuantil.HT.Rapido(Tipo.Fd,beta, N, PIn, Ytotal)
## Estimación de la densidad: fx, fy
fx<-vector(len=p)
for (k in 1:p) fx[k]<-densidad(datosx[,k], Qx[k])
}

```

```

        fy <- densidad(datosy, Qy)
### Resultado del Programa
resultado<-list(X=X,Xm=Xm,Ym = Ym, Yu = Yu, fx = fx, fy = fy, QgorroX=QgorroX,
PItam1=PItam1, PIm=PIm, PIu=PIu , estimadorDirecto=estimadorDirecto )
        resultado
## NOTA: X y Xm son matrices y fx y QgorroX son vectores.
}

est.Pvar.v4<-function(Tipo.Fd,Xm,Ym,PIm,QgorroX,beta,fx,fy,
N,tam1,tam2,Qy,Qx,y,x,X,estimadorDirecto,p1lx, p1lxx, Yu, PIu)
# Este programa obtiene los estimadores propuestos bajo P-variables auxiliares
# vectores: QgorroX, fx. Matrices: X, Xm
{
    valores.p<-length(Xm[,])
    QgorroXm<-vector(len=valores.p)
    for (k in 1:valores.p)
        QgorroXm[k]<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PIm,Xm[,k])
    QgorroYm<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PIm,Ym)
    tam <- length(Ym)
##### Calculo de la matriz B original y V(Qu)
    m<-length(PIm)
    u<-tam2-m
    vector.R<-Qy/Qx
    VqYu.ver<-((N-u)/(N*u) ) * beta*(1-beta)*fy^(-2)
    matriz.B.ver<-matrix(nr=valores.p, nc=valores.p) # La llamada matriz B
    for ( i in 1:valores.p)
        for (j in 1:valores.p)
            { if (i==j) matriz.B.ver[i,j]<- (beta*(1-beta)/fy^2)*
                ( (1/m - 1/N) + (1/m - 1/tam1) * vector.R[i] * (fy/fx[i]) *
                    (vector.R[i]*(fy/fx[i]) + 2*(1-(p1lx[i]/(beta*(1-beta)))))) )
                else matriz.B.ver[i,j]<- (beta*(1-beta)/fy^2)*
                ( (1/m - 1/N) +
                    (1/tam1 - 1/m) * vector.R[i] * (fy/fx[i]) * ( p1lx[i]/(beta*(1-beta)) -1) +
                    (1/tam1 - 1/m) * vector.R[j] * (fy/fx[j]) * ( p1lx[j]/(beta*(1-beta)) -1) -
                    (1/tam1 - 1/m) * vector.R[j] * vector.R[i] * (fy^2/(fx[j]*fx[i])) * (
                    p1lxx[i,j]/(beta*(1-beta)) -1) ) )
            }
    QgorroRazon<-vector(len=valores.p)
    for (k in 1:valores.p) QgorroRazon[k]<-QgorroYm*(QgorroX[k]/QgorroXm[k])
    QgorroYu <- cuantil.HT.Rapido(Tipo.Fd,beta,N,PIu,Yu)
    if (valores.p==2)
    {
        W.r.ver<-vector(len=3)
        estimador.R.ver<-vector(len=3)
        W.r.ver[1] <- VqYu.ver/(VqYu.ver + matriz.B.ver[1,1])
        if (is.na(W.r.ver[1])) {W.r.ver[1]<-0.5}
        estimador.R.ver[1] <- QgorroRazon[1] * W.r.ver[1] + (1 - W.r.ver[1]) * QgorroYu
        W.r.ver[2] <- VqYu.ver/(VqYu.ver + matriz.B.ver[2,2])
        if (is.na(W.r.ver[2])) {W.r.ver[2]<-0.5}
        estimador.R.ver[2] <- QgorroRazon[2] * W.r.ver[2] + (1 - W.r.ver[2]) * QgorroYu
        e<-matrix(1,nr=2,nc=1)
        inversa.B.ver<-solve(matriz.B.ver, tol=1e-1000)
        deno2<-t(e) %*% inversa.B.ver %*% e
        eBe2<-as.vector(deno2)
        w.opt.ver<- ( inversa.B.ver %*% e) /eBe2
        w.opt.ver<-as.vector(w.opt.ver)
        QgorroYmRazon.ver<-sum(w.opt.ver*QgorroRazon) # Con la B.verdadera
#Calculamos W.opt usando estimador de formulas verdaderas
        VqYr.ver <-(1/eBe2)
        W.r.ver[3] <- VqYu.ver/(VqYu.ver + VqYr.ver)
        if (is.na(W.r.ver[3])) {W.r.ver[3]<-0.5}
        estimador.R.ver[3]<-QgorroYmRazon.ver* W.r.ver[3] + (1 - W.r.ver[3]) * QgorroYu
    } # end if(valores.p==2)
    else

```

```

{ W.r.ver<-vector(len=valores.p)
  estimador.R.ver<-vector(len=valores.p)
for (p in 1:valores.p)
  { if (p==1)
    { W.r.ver[p] <- VqYu.ver/(VqYu.ver + matriz.B.ver[p,p])
      if (is.na(W.r.ver[p])) {W.r.ver[p]<-0.5}
    estimador.R.ver[p] <- QgorroRazon[p] * W.r.ver[p] + (1 - W.r.ver[p]) * QgorroYu
    }
    else
    { e<-matrix(1,nr=p,nc=1)
      matriz.B.verP<-matriz.B.ver[1:p,1:p]
      inversa.B.ver<-solve(matriz.B.verP, tol=1e-1000)
      deno2<-t(e) %*% inversa.B.ver %*% e
      eBe2<-as.vector(deno2)
      w.opt.ver<- ( inversa.B.ver %*% e) /eBe2
      w.opt.ver<-as.vector(w.opt.ver)
      QgorroRazonP<-QgorroRazon[1:p]
      QgorroYmRazon.ver<-sum(w.opt.ver*QgorroRazonP) # Con la B.verdadera
#Calculamos W.opt usando estimador de formulas verdaderas
      VqYr.ver <-(1/eBe2)
      W.r.ver[p] <- VqYu.ver/(VqYu.ver + VqYr.ver)
      if (is.na(W.r.ver[p])) {W.r.ver[p]<-0.5}
    estimador.R.ver[p]<-QgorroYmRazon.ver*W.r.ver[p] + (1 - W.r.ver[p]) * QgorroYu
    } # end else if (p==1)
  } # end for (p in valores.p)
} # end else if (valores.p==2)
list(estimador.R.ver=estimador.R.ver, W.r.ver=W.r.ver)
### estas dos salidas son vectores.
} # end function

```

```

Simula.2ocas.Pvar.v4<-function(NombreP, Tipo.Muestreo, Tipo.Fd, DatosTodos,
variableY, variablesX, variableZ, tam1, tam2, m, replicas, beta)
# Realiza la simulación en M.Ocasiones sucesivas con P-variables auxiliares
{
  N<-length(DatosTodos[,1])
  p<-length(variablesX)
  z<-DatosTodos[,variableZ]
  datosx<-DatosTodos[,variablesX]
  datosy<-DatosTodos[,variableY]
  salidaError <- paste("C:\\OcasionesSucesivas\\Pvariables\\err",NombreP,"_tam",
tam1, "y", tam2, "beta_", beta, ".txt", sep = "")
  salidaEstimador<- paste("C:\\OcasionesSucesivas\\Pvariables\\est",NombreP,"_m",
m, "tam", tam1, "y", tam2, "beta_", beta, ".txt", sep = "")
#AHORA EMPIEZA EL PROCESO DE REPLICACIÓN####
  Qy<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,vectory=datosy)
  Qx<-vector(le=p)
  for (k in 1:p) Qx[k]<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,datosx[,k])
  if(p==2) Num.est<-3 else Num.est<-p
  ecmDirecto <- 0
  ecmRatio.ver <- double(Num.est)
  sesDirecto <- 0
  sesRatio.ver <- double(Num.est)
  ResesDirecto <- 0
  ResesRatio.ver <- double(Num.est)
### Calculo de pllx (La verdadera)
  pllx<-vector(len=p)
  for (k in 1:p)
  { suma <- 0
    for(i in 1:N) # length(x) = length(y) = N
      { if((datosx[i,k] <= Qx[k]) && (datosy[i] <= Qy))
        suma <- suma + 1
      }
  }
}

```

```

    }
    p1lx[k] <- suma/N
  }
### Calculo de p1lxx (La original)
  p1lxx<-matrix(nr=p,nc=p)
  for (k in 1:p)
    for (l in 1:p)
      if (k<l)
        { suma <- 0
          for(i in 1:N) ### length(x) = N
            { if((datosx[i,k] <= Qx[k]) && (datosx[i,l] <= Qx[l]))
              suma <- suma + 1
            }
          p1lxx[k,l] <- suma/N
          p1lxx[l,k] <- suma/N # puesto que es simetrica
        }
  }
  for(b in 1:replicas)
    {
      datos <- Datos.X.Pvar.v3(Tipo.Muestreo,Tipo.Fd, datosx, datosy, z, N,
tam1, tam2, m, beta,Qy,Qx)
      Xm <- datos$Xm
      Ym <- datos$Ym
      X <- datos$X
      fx <- datos$fx
      fy <- datos$fy
      Yu <- datos$Yu
      QgorroX<-datos$QgorroX
      PITam1<-datos$PITam1
      PIm<-datos$PIm
      PIu<-datos$PIu
      estimadorDirecto<-datos$estimadorDirecto
est.Pvar.v4(Tipo.Fd,Xm,Ym,PIm,QgorroX, beta, fx, fy, N, tam1, tam2, Qy, Qx,
datosy, datosx,X,estimadorDirecto,p1lx, p1lxx, Yu, PIu)->Estimadores
      W.r.ver<-Estimadores$W.r.ver
      estimador.R.ver<-Estimadores$estimador.R.ver
      ecmDirecto <- ecmDirecto + (estimadorDirecto - Qy)^2
      sesDirecto <- sesDirecto + abs(estimadorDirecto - Qy)
      ResesDirecto <- ResesDirecto + (estimadorDirecto - Qy)
    }
  }
  for (i in 1:Num.est)
    { #### Usando varianzas verdaderas
      estimadorRatio.ver <-estimador.R.ver[i]
      ecmRatio.ver[i] <- ecmRatio.ver[i] + (estimadorRatio.ver - Qy)^2
      sesRatio.ver[i] <- sesRatio.ver[i] + abs(estimadorRatio.ver - Qy)
      ResesRatio.ver[i] <- ResesRatio.ver[i] +(estimadorRatio.ver - Qy)
    }
  }
## Escribimos el valor de los estimadores
  estim<-paste(W.r.ver[1], estimador.R.ver[1], sep="\t")
  for (i in 2:Num.est)
    estim<-paste(estim,W.r.ver[i], estimador.R.ver[i], sep="\t" )
    linea <- paste(b, Qy, estimadorDirecto, estim, sep = "\t")
    write(linea, file = salidaEstimador, ncolumns = 1, append = T)
  } # end for (b in 1:replicas)
### VAMOS A ESCRIBIR LOS RESULTADOS DEL PROCESO DE REPLICACIÓN ###
  ecmDirecto <- (ecmDirecto/replicas)
  RecmDirecto<- 1
  sesDirecto <- sesDirecto/(replicas * Qy)
  ResesDirecto <- ResesDirecto/(replicas * Qy)
  RecmRatio.ver<-vector(len=Num.est)
  for (i in 1:Num.est)
    {
      ### RE

```

```

RecmRatio.ver[i] <- (ecmRatio.ver[i]/replicas)/ecmDirecto
### Sesgos absolutos
sesRatio.ver[i] <- sesRatio.ver[i]/(replicas * Qy)
### Sesgos Relativos
ResesRatio.ver[i] <- ResesRatio.ver[i]/(replicas * Qy)
}
estim<-paste(W.r.ver[1],sesRatio.ver[1],ResesRatio.ver[1],RecmRatio.ver[1],
sep="\t")
for (i in 2:Num.est)
{
estim<-paste(estim,W.r.ver[i],sesRatio.ver[i],ResesRatio.ver[i],
RecmRatio.ver[i], sep="\t" )
}
linea<-paste(m,sesDirecto,ResesDirecto,ecmDirecto,RecmDirecto,estim,sep = "\t")
write(linea, file = salidaError, ncolumns = 1, append = T)
## NOTA
## Si p=2, escribe en el sigte. orden: p1.x1, p1.x2, p2
## Si p>2, escribe en el sigte. orden: p1, p2, p3, p4,.....
## EJEMPLO
#Simula.2ocas.Pvar.v4(NombreP=c("Turismos"),Tipo.Muestreo=1,Tipo.Fd=1,
#DatosTodos=Turismos,variableY=3,variablesX=c(9,10,11,12),variableZ=1,tam1=100,
#tam2=100, m=50, replicas=1000, beta=0.5)
}

```

C.3.3. Muestreo bifásico

```
## Algunos diseños muestrales asumiendo muestreo bifásico
```

```

muestra.bifasico.Inclusion.M.M<-function(N, tam1, tam2)
# Extraccion de muestras en muestreo bifásico. M.M = MAS.MAS = TIPO=1.
# Muestra 1: tamaño n´= tam1. Mediante MAS
# Muestra 2: tamaño n = tam 2. Mediante MAS
{
muestra1 <- sample(N,tam1)
muestra2 <- sample(tam1,tam2)
muestra2 <- muestra1[muestra2]
### Calculo de las probabilidades de Inclusión
# De primer orden
#### De la Muestra 1.
PItam1<-vector(len=tam1)
for (i in 1:tam1) PItam1[i]<-tam1/N
#### De la Muestra 2.
PItam2<-vector(len=tam2)
PIcondi.tam2<-vector(len=tam2)
for (i in 1:tam2)
{ PItam2[i]<-tam2/N
PIcondi.tam2[i]<-tam2/tam1
}
# De segundo orden
M.disenomatrix(nro=tam2, nco=tam2)
for (i in 1:tam2)
for (j in 1:tam2)
if (i != j) M.disenom[i,j]<-(tam2/tam1)*((tam2-1)/(tam1-1))
for (i in 1:tam2) M.disenom[i,i]<-PIcondi.tam2[i]
resultado <- list(muestra1 = muestra1, muestra2 = muestra2, PItam1 = PItam1,
PItam2 = PItam2, M.disenom=M.disenom, PIcondi.tam2=PIcondi.tam2)
resultado
}

```

```

muestra.bifasico.Inclusion.M.Z<-function(N, x, tam1, tam2)
# Extraccion de muestras en muestreo bifásico. M.Z = MAS.Midzuno = TIPO=5.
{
  N1<-length(x)
  if (N != N1) stop("La longitud del vector x es dintinta a N")
  muestral <- sample(N,tam1)
  datos.x.muestral<-x[muestral]
  muestra2 <- Etiquetas.midzuno(tam2,datos.x.muestral)
  muestra2 <- muestral[muestra2]
### Calculo de las probabilidades de Inclusión
# De primer orden
  datos.x.tam2<-x[muestra2]
#### De la Muestra 1.
  PItam1<-vector(len=tam1)
  for (i in 1:tam1) PItam1[i]<-tam1/N
#### De la Muestra 2
  PI.Partel<-PItam1[1:tam2]
  PIcondi.tam2<-vector(len=tam2)
  sum(datos.x.muestral)->TotalXtam1
  for (i in 1:tam2) PIcondi.tam2[i]<-
  ((tam1-tam2)/(tam1-1) )*(datos.x.tam2[i]/TotalXtam1)+(tam2-1)/(tam1-1)
  PItam2<-PI.Partel*PIcondi.tam2
# De segundo orden
M.disen0<-matrix(nro=tam2, nco=tam2)
  for (i in 1:tam2)
  { for (j in 1:tam2)
    { if (i != j)
      M.disen0[i,j]<-(tam2-1)/(tam1-1)*((tam1-tam2)/(tam1-2))*
      (datos.x.tam2[i]/TotalXtam1+datos.x.tam2[j]/TotalXtam1)+
      (tam2-2)/(tam1-2)
    }
  }
  for (i in 1:tam2) M.disen0[i,i]<-PIcondi.tam2[i]
  resultado <- list(muestral=muestral,muestra2=muestra2,PItam1=PItam1,
    PItam2 = PItam2, M.disen0=M.disen0, PIcondi.tam2=PIcondi.tam2)
  resultado
}

muestra.bifasico.Inclusion.M.P<-function(N, x, tam1, tam2)
# Extraccion de muestras en muestreo bifásico. M.P = MAS.Poisson = TIPO=8.
{
  N1<-length(x)
  if (N != N1) stop("La longitud del vector x es dintinta a N")
  muestral <- sample(N,tam1)
  datos.x.muestral<-x[muestral]
  salida <- Etiquetas.Poisson(tam2,datos.x.muestral)
  muestra2<-salida$Etiquetas
  muestra2 <- muestral[muestra2]
## Bajo este muestreo cambia el tamaño muestral:
  tam2<-salida$n.obtenido
### Calculo de las probabilidades de Inclusión
# De primer orden
  datos.x.tam2<-x[muestra2]
#### De la Muestra 1.
  PItam1<-vector(len=tam1)
  for (i in 1:tam1) PItam1[i]<-tam1/N
#### De la Muestra 2
  PI.Partel<-PItam1[1:tam2]
  PIcondi.tam2<-vector(len=tam2)
  sum(datos.x.muestral)->TotalXtam1
  for (i in 1:tam2) PIcondi.tam2[i]<- tam2*(datos.x.tam2[i]/TotalXtam1)
  PItam2<-PI.Partel*PIcondi.tam2
# De segundo orden

```

```

M.disenos<-matrix(nro=tam2, nco=tam2)
  for (i in 1:tam2)
    { for (j in 1:tam2)
      { if (i != j)
        M.disenos[i,j]<-PIcondi.tam2[i]*PIcondi.tam2[j]
      }
    }
  for (i in 1:tam2) M.disenos[i,i]<-PIcondi.tam2[i]
resultado <- list(muestral=muestral,muestra2=muestra2,PItam1=PItam1,
                 PITam2 = PITam2, M.disenos=M.disenos, PIcondi.tam2=PIcondi.tam2)
  resultado
}

muestra.bifasico.Inclusion<-function(TIPO,N, x, tam1, tam2)
# Para seleccionar un tipo de muestreo y llevarlo a la práctica.
{
if (tam2>tam1) stop("tam2 ha de ser mas pequeño que tam1")
switch(TIPO,
  resultado<- muestra.bifasico.Inclusion.M.M(N, tam1, tam2),
  resultado<- stop("Da otro N°"),#muestra.bifasico.M.L(N, x, tam1, tam2),
  resultado<- stop("Da otro N°"),#muestra.bifasico.L.M(N, x, tam1, tam2),
  resultado<- stop("Da otro N°"),#muestra.bifasico.L.L(N, x, tam1, tam2),
  resultado<- muestra.bifasico.Inclusion.M.Z(N, x, tam1, tam2),
  resultado<- stop("Da otro N°"),#muestra.bifasico.Z.M(N, x, tam1, tam2),
  resultado<- stop("Da otro N°"),#muestra.bifasico.Z.Z(N, x, tam1, tam2),
  resultado<- muestra.bifasico.Inclusion.M.P(N, x, tam1, tam2)
)
resultado
}

Datos.X.HT.ICB<-function(Tipo.Muestreo,Tipo.Fd, datosx, datosy,datosz, N, tam1,
tam2, beta, Qy, Qx)
# Obtiene datos de muestras, Probabilidades de inclusión, estimadores de
#cuantiles básicos y funciones de densidad bajo un muestreo bifásico.
{
# Para extraer muestras y obtener P. inclusion se utiliza variable z.
muestras <- muestra.bifasico.Inclusion(Tipo.Muestreo,N,datosz,tam1,tam2)
muestral <- muestras$muestral
muestra2 <- muestras$muestra2
### Se pone el nuevo tam2 en el caso de M. de Poisson
length(muestra2)->tam2
Xtam2 <- double(tam2)
Ytam2 <- double(tam2)
Ztam2<-double(tam2)
Xtam2<-datosx[muestra2]
Ytam2<-datosy[muestra2]
Ztam2<-datosz[muestra2]
X <- double(tam1)
X<-datosx[muestral]
## Estimación de las densidades
fx<-densidad(datosx, Qx)
fy<-densidad(datosy, Qy)
### Probabilidades de Inclusión y QgorroX
PItam1<-muestras$PItam1
PItam2<-muestras$PItam2
QgorroXtam1<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PItam1,X)
M.disenos<-muestras$M.disenos
PIcondi.tam2<-muestras$PIcondi.tam2
### Resultado del Programa
resultado <- list(X = X, Xtam2 = Xtam2, Ytam2 = Ytam2, Ztam2=Ztam2,
                 QgorroXtam1=QgorroXtam1, PITam1=PItam1, PITam2=PItam2,

```

```

M.diseno=M.diseno, fx=fx, fy=fy, PIcondi.tam2=PIcondi.tam2)
resultado
}

est.IC.Bifasico<-
function(Tipo.Muestreo,Tipo.Fd,Xtam2,Ytam2,Ztam2,PItam1,PItam2,QgorroXtam1,
beta,N,M.diseno,fx,fy,PIcondi.tam2)
# Obtiene una lista con diferentes estimadores de cuantiles
{
tam2<-length(Ytam2)
tam1<-length(PItam1)
QgorroXtam2<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PItam2,Xtam2)
QgorroYtam2<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PItam2,Ytam2)
# Estimadores con Probabilidades de Incursión .I
nume<-0
deno<-0
for (i in 1:tam2)
for(j in 1: tam2)
{
nume<-nume+
((M.diseno[i,j]-PIcondi.tam2[i]*PIcondi.tam2[j])/(M.diseno[i,j]))*
(Delta(QgorroYtam2,Ytam2[i])/PItam2[i])*
(Delta(QgorroXtam2,Xtam2[j])/PItam2[j])
deno<-deno+
((M.diseno[i,j]-PIcondi.tam2[i]*PIcondi.tam2[j])/(M.diseno[i,j]) )*
(Delta(QgorroXtam2,Xtam2[i])/PItam2[i])*
(Delta(QgorroXtam2,Xtam2[j])/PItam2[j])
}
a.opt.tam2.I<-(QgorroXtam2/QgorroYtam2)*(fx/fy)*(nume/deno)
b.tam2.I<-(fx/fy)*(nume/deno)
QgorroR.alfal <- QgorroYtam2*(QgorroXtam1/QgorroXtam2)
QgorroR.opt.tam2.I <- QgorroYtam2*(QgorroXtam1/QgorroXtam2)^(a.opt.tam2.I)
QgorroDif.tam2.I <- QgorroYtam2+b.tam2.I*(QgorroXtam1-QgorroXtam2)
listadedatos<-list(QgorroYtam2=QgorroYtam2,QgorroXtam2=QgorroXtam2,
QgorroR.alfal=QgorroR.alfal,
QgorroR.opt.tam2.I=QgorroR.opt.tam2.I, QgorroDif.tam2.I=QgorroDif.tam2.I)
listadedatos
}

Simula.bifasico.IC<-function(NombreP,Poblacion,Tipo.Muestreo,Tipo.Fd,variableY,
variableX,variableZ, tam1, tam2, replicas, beta)
# Realiza simulaciones de varios estimadores en M.bifásico.Los argumentos son:
# "NombreP": Nombre de Población. Sirve para darle nombre al fichero.
# "Población": Contiene datos (objeto de tipo matriz o dataframe) donde cada
#columna es una variable.
# "Tipo.Muestreo": Tipo=1: M.M, Tipo=5: M.Z, Tipo=8: M.P
# Tipo.Fd: Tipo de estimación para la F.d. Tipo.Fd = 1: H-T, Tipo.Fd = 2: Hayek
# "variableY": Indicamos un NÚMERO: Columna donde están los datos de Y.
# "variableX": Indicamos un NÚMERO: Columna donde están los datos de X.
# "variableZ": Indicamos un NÚMERO: Columna donde están los datos de Z.
# "tam1", "tam2": n' y n.
# "replicas": Número de muestras independientes a extraer.
# "beta": Orden del cuantil.
{
N<-length(Poblacion[,1])
variables<-length(Poblacion[1,])
datosx<-Poblacion[,variableX]
datosy<-Poblacion[,variableY]
datosz<-Poblacion[,variableZ]
salidaEstimador <- paste("C:\\Bifasico\\IC\\estI_P_",NombreP,"TM_",
Tipo.Muestreo,"TFd_",Tipo.Fd,"tam_",tam1,"y",tam2,"beta_",beta,".txt",sep = "")
#AHORA EMPIEZA EL PROCESO DE REPLICACIÓN #
Qy<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,vector=datosy)

```

```

    Qx<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,vectory=datosx)
    ecmDirecto <- 0 ; sesDirecto <- 0
    ecmR.alfal <- 0 ; sesR.alfal <- 0
    ecmR.opt.tam2.I<- 0 ; sesR.opt.tam2.I <- 0
    ecmDif.tam2.I<- 0 ; sesDif.tam2.I <- 0
    Matriz.Qgorros<-matrix(nro=replicas, nco=4)
    for(b in 1:replicas)
        {
            datos <- Datos.X.HT.ICB(Tipo.Muestreo,Tipo.Fd,datosx,datosy,
datosz, N, tam1, tam2, beta,Qy,Qx)
            Xtam2 <- datos$Xtam2
            Ytam2 <- datos$Ytam2
            Ztam2 <- datos$Ztam2
            X <- datos$X
            QgorroXtam1<-datos$QgorroXtam1
            PItam1<-datos$PItam1
            PItam2<-datos$PItam2
            M.diseno<-datos$M.diseno
            fx<-datos$fx
            fy<-datos$fy
            PIcondi.tam2<-datos$PIcondi.tam2
# ESTIMADORES
            estimadores<-est.IC.Bifasico(Tipo.Muestreo,Tipo.Fd,Xtam2,Ytam2,Ztam2,
                PItam1,PItam2,QgorroXtam1,beta,N,M.diseno,fx, fy,PIcondi.tam2 )
            length(Ytam2)->tam2.P
            Matriz.Qgorros[b,1]<-estimadores$QgorroYtam2
            Matriz.Qgorros[b,2]<-estimadores$QgorroXtam2
            Matriz.Qgorros[b,3]<-QgorroXtam1
            Matriz.Qgorros[b,4]<-tam2.P
### 1 ### El estimador directo para comparar:
                estDirecto <- estimadores$QgorroYtam2
                ecmDirecto <- ecmDirecto + (estDirecto - Qy)^2
                sesDirecto <- sesDirecto + abs(estDirecto - Qy)
### 2 ### Estimador de razón con alfa=1
                estR.alfal <- estimadores$QgorroR.alfal
                ecmR.alfal <- ecmR.alfal + (estR.alfal - Qy)^2
                sesR.alfal <- sesR.alfal + abs(estR.alfal - Qy)
## 3 ### Estimador de razón con alfa optimo version=tam2.B
                estR.opt.tam2.I <- estimadores$QgorroR.opt.tam2.I
                ecmR.opt.tam2.I <- ecmR.opt.tam2.I + (estR.opt.tam2.I - Qy)^2
                sesR.opt.tam2.I <- sesR.opt.tam2.I + abs(estR.opt.tam2.I - Qy)
### 4 ### Estimador de diferencia version=tam1.I
                estDif.tam2.I <- estimadores$QgorroDif.tam2.I
                ecmDif.tam2.I <- ecmDif.tam2.I + (estDif.tam2.I - Qy)^2
                sesDif.tam2.I <- sesDif.tam2.I+abs(estDif.tam2.I - Qy)
### Se escribe el valor de los estimadores
                linea <- paste(
b,tam2.P,Qy,estDirecto,estR.alfal,estR.opt.tam2.I,estDif.tam2.I,sep = "\t")
                write(linea, file = salidaEstimador, ncolumns = 1, append = T)
            } # end for (b in 1:replicas)
                ### AHORA EMPIEZA EL PROCESO CON COVARIANZAS ###
                cov.Qxy<-cov(Matriz.Qgorros[,1],Matriz.Qgorros[,2])
                var.Qx<-var(Matriz.Qgorros[,2])
                b.dif<-cov.Qxy/var.Qx
                salidaError <- paste("C:\\Bifasico\\IC\\errP_", NombreP, "TM_", Tipo.Muestreo,
                    "TFd_",Tipo.Fd, "tam1_", tam1,"beta_", beta, ".txt", sep = "")
                salidaEstimador<-paste("C:\\Bifasico\\IC\\estCovP_",NombreP,"TM_",
                    Tipo.Muestreo,"TFd_",Tipo.Fd,"tam_",tam1,"y",tam2,"beta_",beta,".txt",sep = "")
                ecmR.opt.tam2.C<- 0 ; sesR.opt.tam2.C <- 0
                ecmDif.C<- 0 ; sesDif.C <- 0
                for (b in 1: replicas)
                    {

```

```

a.opt.tam2.C<-(Matriz.Qgorros[b,2]/Matriz.Qgorros[b,1])*b.dif
### 2 ### Estimador de diferencia, version: .C
estDif.C <- Matriz.Qgorros[b,1] +
          b.dif*(Matriz.Qgorros[b,3]-Matriz.Qgorros[b,2])
ecmDif.C <- ecmDif.C + (estDif.C - Qy)^2
sesDif.C <- sesDif.C + abs(estDif.C - Qy)
### 1 ### Estimador de razón con alfa optimo. version: tam2.C
estR.opt.tam2.C <- Matriz.Qgorros[b,1]*
(Matriz.Qgorros[b,3]/Matriz.Qgorros[b,2])^(a.opt.tam2.C)
ecmR.opt.tam2.C <- ecmR.opt.tam2.C + (estR.opt.tam2.C - Qy)^2
sesR.opt.tam2.C <- sesR.opt.tam2.C + abs(estR.opt.tam2.C - Qy)
### Se escribe el valor de los estimadores
linea <- paste(b, estDif.C, estR.opt.tam2.C, sep = "\t")
write(linea, file = salidaEstimador, ncolumns = 1, append = T)
} # end for (b in 1:replicas)
#### SE ESCRIBEN LOS RESULTADOS DEL PROCESO DE SIMULACIÓN ###
### Errores Cuadráticos Medios y sus Ratios
ecmDirecto <- (ecmDirecto/replicas)
RecmDirecto<- 1
RecmR.alfal<-(ecmR.alfal/replicas)/ecmDirecto
RecmR.opt.tam2.I<-(ecmR.opt.tam2.I/replicas)/ecmDirecto
RecmDif.tam2.I<-(ecmDif.tam2.I/replicas)/ecmDirecto
RecmDif.C<-(ecmDif.C/replicas)/ecmDirecto
RecmR.opt.tam2.C<-(ecmR.opt.tam2.C/replicas)/ecmDirecto
## Sesgos absolutos relativos
sesDirecto <- sesDirecto/(replicas * Qy)
RsesR.alfal <- sesR.alfal/(replicas * Qy)
RsesR.opt.tam2.I <- sesR.opt.tam2.I/(replicas * Qy)
RsesDif.tam2.I<- sesDif.tam2.I/(replicas * Qy)
RsesDif.C<- sesDif.C/(replicas * Qy)
RsesR.opt.tam2.C <- sesR.opt.tam2.C/(replicas * Qy)
### Se escribe la matriz de resultados
mean(Matriz.Qgorros[,4])>mean.tam2
linea<-paste(tam2, mean.tam2, sesDirecto, ecmDirecto, RecmDirecto, RsesR.alfal,
RecmR.alfal,RsesR.opt.tam2.I,RecmR.opt.tam2.I,RsesDif.tam2.I,RecmDif.tam2.I,
RsesDif.C,RecmDif.C,RsesR.opt.tam2.C, RecmR.opt.tam2.C, sep = "\t")
write(linea, file = salidaError, ncolumns = 1, append = T)
## Ejemplo
#Simula.bifasico.IC(c("Fam1500YX1X2"), Fam1500, Tipo.Muestreo=5, Tipo.Fd=1,
#variableY=1,variableX=2,variableZ=3,tam1=150, tam2=50, replicas=10, beta=0.75)
}

```

C.3.4. Muestreo bifásico aplicado a la estratificación

```

# Diseños muestrales asumiendo muestreo bifásico aplicado a la estratificación.

muestra.bifasicoST.M.M<-funcion(Nh, tam1, tam2)
# Extracción de muestras en Muestreo bifasico aplicado a la estratificación y
#cálculo de probabilidades de inclusión: M.M = Mas.Mas = "TIPO=1".
# "Nh": Tamaños poblacionales de los estratos.
# Muestra 1(s´): tamaño n´= "tam1". Mediante MAS.
# Se asocian las unidades a los estratos (Nh: tamaños de los estratos).
# De aqui se obtiene: s´1,...,s´H, de tamaños: n´1,...,n´h,...,n´H.
# "tam2": tamaño total de la submuestra = n
# Mediante afijacion proporcional se divide tam2 en estratos: tam2h.
# De cada s´h se extraen muestras de tamaño "tam2h" Mediante MAS

```

```

# Nota1: Puesto muestra 1 es MAS, PItam1 son siempre iguales.
# Nota2: La v.auxiliar se usa en la segunda etapa. En MAS no es necesaria.
{ N<-sum(Nh)
  ceros<-1
  while (ceros>=1)
  {
    muestral <- sample(N,tam1)
    muestral <- sort(muestral)
    # tam1h: contendrá tamaños de estratos de donde se van a extraer muestras
    tam1h<-divideestratos(muestral,Nh)
    Trues<-tam1h==0
    ceros<-sum(Trues)
  }
  muestra2.prop<-vector(len=tam2)
  Eti.muestra2.prop<-vector(len=tam2)
  PItam2.prop<-vector(len=tam2)
  tam2h<-afijacion.proporcional(tam2,tam1h)
  L<-length(Nh)
### Cálculo de las probabilidades de Inclusión de la muestra 1
### De primer y segundo orden
  PItam1<-vector(len=tam1)
  Pij<-(tam1/N)*((tam1-1)/(N-1))
  PijTam1<-matrix(Pij, nr=tam1, nco=tam1)
  for (i in 1:tam1)
  { PItam1[i]<-tam1/N
    PijTam1[i,i]<-tam1/N # En la diagonal se ponen las de primer orden
  }
  ## EMPIEZA PROCESO DE SELECCION EN SEGUNDA ETAPA ##
PijCondi.tam2.prop<-matrix(nr=tam2,nc=tam2)
PIcondi.tam2.prop<-c()
  for (h in 1:L)
  { ## for #1
    if (h !=1 )
    { a<-sum(tam1h[1:(h-1)])+1
      a2<-sum(tam2h[1:(h-1)])+1
    }
    else
    { a<-1
      a2<-1
    }
    b<-sum(tam1h[1:h])
    b2<- sum(tam2h[1:h])
    Eti.muestra2.prop[a2:b2] <- sample(tam1h[h], tam2h[h])
    muestrah<-muestral[a:b]
    muestra2.prop[a2:b2] <- muestrah[Eti.muestra2.prop[a2:b2]]
### Cálculo de las probabilidades de Inclusión en la segunda fase
#### De la Muestra 2 de primer orden
    PI.Partel<-PItam1[1:tam2h[h]]
    PIcondi.tam2<-vector(len=tam2h[h])
    for (i in 1:tam2h[h]) PIcondi.tam2[i]<- tam2h[h]/tam1h[h]
    PIcondi.tam2.prop<-c(PIcondi.tam2.prop,PIcondi.tam2)
    PItam2.prop[a2:b2]<-PI.Partel*PIcondi.tam2
#### De la Muestra 2 de segundo orden ( matriz Pij.cond)
## Procedimiento: Obtener matriz Pij.cond y multiplicarlas compo. a componente
  for (l in 1:L)
  { ## for-l
    if (h==1)
      PijCondi.tam2.prop[a2:b2,a2:b2]<-
(tam2h[h]/tam1h[h])*((tam2h[h]-1)/(tam1h[h]-1))
    else ## else-l
    {

```

```

        if (l !=1 )
            al2<-sum(tam2h[1:(l-1)])+1
        else
            al2<-1
        bl2<- sum(tam2h[1:l])
        PijCondi.tam2.prop[a2:b2,al2:bl2]<-(tam2h[h]/tam1h[h])*(tam2h[l]/tam1h[l])
    } ## else-1
} ## for-1
## Se incluyen las de primer orden en la diagonal principal:
for (i in a2:b2) PijCondi.tam2.prop[i,i]<-tam2h[h]/tam1h[h]
} ## for #1
#Multiplicacion de matrices para obtener prob. de Inclusión en la segunda fase
PijTam2.prop<- PijTam1[1:tam2,1:tam2] * PijCondi.tam2.prop
# 1:tam2,1:tam2 ->Puesto que son todas iguales.
# " * " -> Producto componente a componente.
#### Para otros estimadores:
##### 1. Para estimador Directo y de Singh-Joader-Tracy (SJT):
# Se extrae otra muestra2 bajo muestreo bifásico. Muestral es la misma.
Eti.muestra2.dir <- sample(tam1,tam2)
muestra2.dir <- muestral[Eti.muestra2.dir]
### Calculo de las probabilidades de Inclusión
#### De primer orden de la Muestra 2
PI.Partel<-PItam1[1:tam2]
PIcondi.tam2.dir<-vector(len=tam2)
for (i in 1:tam2) PIcondi.tam2.dir[i]<- tam2/tam1
PItam2.dir<-PI.Partel*PIcondi.tam2.dir
#### De segundo orden de la Muestra 2
Pij<-(tam2/tam1)*((tam2-1)/(tam1-1))
PijCondi.tam2.dir<-matrix(Pij, nr=tam2,nc=tam2)
for (i in 1:tam2) PijCondi.tam2.dir[i,i]<-tam2/tam1
PijTam2.dir<- PijTam1[1:tam2,1:tam2] * PijCondi.tam2.dir
##### 2. Estimador de Silva y Skinner (ps):
ceros<-1
while (ceros>=1)
{ muestra2.ps <-sample(N,tam2)
muestra2.ps<-sort(muestra2.ps)
### Ahora se pos-estratifica: (mediante funcion divideestratos)
tam2h.ps<-divideestratos(muestra2.ps, Nh)
Trues<- tam2h.ps==0
ceros<-sum(Trues)
}
Pi.tam2.ps<-vector(len=tam2)
Pij<-(tam2/N)*((tam2-1)/(N-1))
Pij.tam2.ps<-matrix(Pij, nc=tam2, nr=tam2)
for (i in 1:tam2)
{ Pi.tam2.ps[i]<- tam2/N
Pij.tam2.ps[i,i]<-tam2/N ## En diagonal se ponen las 1er. orden
}
resultado <- list(
tam1h=tam1h,tam2h=tam2h,muestralh=muestral,muestra2.prop=muestra2.prop,
Eti.muestra2.prop=Eti.muestra2.prop, PITam1=PITam1,PITam2.prop = PITam2.prop,
PIcondi.tam2.prop=PIcondi.tam2.prop,
PijCondi.tam2.prop=PijCondi.tam2.prop,muestra2.dir=muestra2.dir,
PITam2.dir=PITam2.dir, PijTam1=PijTam1, PijTam2.prop=PijTam2.prop,
PijTam2.dir=PijTam2.dir,muestral=muestral, PijCondi.tam2.dir=PijCondi.tam2.dir,
PIcondi.tam2.dir=PIcondi.tam2.dir, Eti.muestra2.dir= Eti.muestra2.dir,
Pi.tam2.ps= Pi.tam2.ps, Pij.tam2.ps= Pij.tam2.ps,
tam2h.ps= tam2h.ps, muestra2.ps=muestra2.ps)
# EJEMPLO: muestra.bifasicoST.M.M(Fam1500.Nh, tam1=100, tam2=50)
resultado
}

```

```

muestra.bifasicoST.M.Z<-function(Nh, x, tam1, tam2)
# Extracción de muestras en Muestreo bifasico aplicado a la estratificación y
#cálculo de probabilidades de inclusión: M.Z = Mas.Midzuno = TIPO=5.
{
  N1<-length(x)
  N<-sum(Nh)
  if (N != N1) stop("La longitud del vector x es dintinta a N=suma de Nh")
ceros<-1
  while (ceros>=1)
  {
    muestral <- sample(N,tam1)
    muestral <- sort(muestral)
# tam1h: contendrá tamaños de estratos de donde se van a extraer muestras
    tam1h<-divideestratos(muestral,Nh)
    Trues<-tam1h==0
    ceros<-sum(Trues)
  }
  datos.x.muestral<-x[muestral]
  muestra2.prop<-vector(len=tam2)
  Eti.muestra2.prop<-vector(len=tam2)
  PItam2.prop<-vector(len=tam2)
  tam2h<-afijacion.proporcional(tam2,tam1h)
  L<-length(Nh)
### Cálculo de las probabilidades de Inclusión
# Muestra 1- Primer orden y segundo orden.
  PItam1<-vector(len=tam1)
  Pij<-(tam1/N)*((tam1-1)/(N-1))
  PijTam1<-matrix(Pij, nr=tam1, nco=tam1)
  for (i in 1:tam1)
  {
    PItam1[i]<-tam1/N
    PijTam1[i,i]<-tam1/N ## En la diagonal se ponen las de 1er. orden
  }
### EMPIEZA PROCESO DE SELECCION EN SEGUNDA ETAPA ###
  datos.x.tam2<-vector(len=tam2)
  PijCondi.tam2.prop<-matrix(nr=tam2,nc=tam2)
  PIcondi.tam2.prop<-c()
  for (h in 1:L)
  { ## for #1
    if (h !=1 )
    {
      a<-sum(tam1h[1:(h-1)])+1
      a2<-sum(tam2h[1:(h-1)])+1
    }
    else
    {
      a<-1
      a2<-1
    }
    b<-sum(tam1h[1:h])
    b2<- sum(tam2h[1:h])
    Eti.muestra2.prop[a2:b2]<- Etiquetas.midzuno(tam2h[h],datos.x.muestral[a:b])
    muestrah<-muestral[a:b]
    muestra2.prop[a2:b2] <- muestrah[Eti.muestra2.prop[a2:b2]]
    datos.x.tam2[a2:b2]<-x[muestra2.prop[a2:b2]]
  } ## for #1
### Cálculo de las probabilidades de Inclusión en la segunda fase
  for (h in 1:L)
  { ## for #2
    if (h !=1 )
    {
      a<-sum(tam1h[1:(h-1)])+1
      a2<-sum(tam2h[1:(h-1)])+1
    }
    else
    {
      a<-1
      a2<-1
    }
  }
}

```

```

    }
    b<-sum(tamlh[1:h])
    b2<- sum(tam2h[1:h])
#### De la Muestra 2 de primer orden
    PI.Partel<-PItam1[1:tam2h[h]]
    PIcondi.tam2<-vector(len=tam2h[h])
    sum(datos.x.muestral[a:b])>TotalXtam1
    for (i in 1:tam2h[h]) PIcondi.tam2[i]<- ( (tam1h[h]-tam2h[h])/
(tamlh[h]-1) )*(datos.x.tam2[i]/TotalXtam1)+(tam2h[h]-1)/(tamlh[h]-1)
    PIcondi.tam2.prop<-c(PIcondi.tam2.prop,PIcondi.tam2)
    PITam2.prop[a2:b2]<-PI.Partel*PIcondi.tam2
#### De la Muestra 2 de segundo orden ( matriz Pij.cond)
    alfa<-datos.x.tam2/TotalXtam1 # dimension=tam2
    for (l in 1:L)
    {
      ## for-1
      if (h==1)
        for (i in a2:b2)
          for (j in a2:b2)
            PijCondi.tam2.prop[i,j]<- ((tam2h[h]-1)/(tamlh[h]-1))* ( ( (tam1h[h]-
tam2h[h])/((tamlh[h]-2) )*(alfa[i]+alfa[j]))+( (tam2h[h]-2)/(tamlh[h]-2) ) )
            else ## else-1
            {if (l !=1 )
              { al<-sum(tamlh[1:(l-1)])+1
                al2<-sum(tam2h[1:(l-1)])+1
              }
            else
              { al<-1
                al2<-1
              }
            bl<-sum(tamlh[1:l])
            bl2<- sum(tam2h[1:l])
            sum(datos.x.muestral[al:bl])>TotalXtam1.l
            for (i in a2:b2)
              for (j in al2:bl2)
                PijCondi.tam2.prop[i,j]<-(((tam1h[h]-tam2h[h])/((tam1h[h]-1)))*(datos.x.tam2[i] /
TotalXtam1)+(tam2h[h]-1)/(tam1h[h]-1))*(((tam1h[l]-tam2h[l])/((tam1h[l]-1)) *
(datos.x.tam2[j] / TotalXtam1.l)+ (tam2h[l]-1)/(tam1h[l]-1) )
            } ## else-1
          } ## for-1
## Se incluyen las de primer orden en la diagonal principal:
      for (i in a2:b2) PijCondi.tam2.prop[i,i]<- ( (tam1h[h]-tam2h[h])/
(tamlh[h]-1) )*(datos.x.tam2[i]/TotalXtam1)+(tam2h[h]-1)/(tam1h[h]-1)
    } ## for #2
# Multiplicacion de matrices para obtener prob. de Inclusión en la segunda fase
    PijTam2.prop<- PijTam1[1:tam2,1:tam2] * PijCondi.tam2.prop
    # 1:tam2,1:tam2 ->Puesto que son todas iguales en la primera etapa.
    # " * " -> Producto componente a componente.
#### Para otros estimadores:
##### 1. Para el estimador Directo y Singh-Joader-Tracy:
    Eti.muestra2.dir <- Etiquetas.midzuno(tam2,datos.x.muestral)
    muestra2.dir <- muestral[Eti.muestra2.dir]
## Cálculo de las probabilidades de Inclusión
##### De primer orden de la Muestra 2
    datos.x.tam2<-x[muestra2.dir]
    PI.Partel<-PItam1[1:tam2]
    PIcondi.tam2.dir<-vector(len=tam2)
    sum(datos.x.muestral)>TotalXtam1
    for (i in 1:tam2) PIcondi.tam2.dir[i]<- (
(tam1-tam2)/(tam1-1) )*(datos.x.tam2[i]/TotalXtam1)+tam2-1)/(tam1-1)
    PITam2.dir<-PI.Partel*PIcondi.tam2.dir
##### De segundo orden de la Muestra 2

```

```

        PijCondi.tam2.dir<-matrix(nr=tam2,nc=tam2)
alfa<-datos.x.tam2/TotalXtam1
  for (i in 1:tam2)
for (j in 1:tam2)
  {
    if (i ==j)
PijCondi.tam2.dir[i,i]<-
      ( (tam1-tam2)/(tam1-1) )*(alfa[i])+ (tam2-1)/(tam1-1)
    else
PijCondi.tam2.dir[i,j]<-((tam2-1)/(tam1-1))*(((tam1-tam2)/(tam1-2) ) *
(alfa[i]+alfa[j])+ ( tam2-2)/(tam1-2) ) )
  }
  PijTam2.dir<- PijTam1[1:tam2,1:tam2] * PijCondi.tam2.dir
#### 2. Silva and Skinner estimator (ps):
ceros<-1
while (ceros>=1)
{ muestra2.ps <-Etiquetas.midzuno(tam2,x)
  muestra2.ps<-sort(muestra2.ps)
  ### Ahora se pos-estratifica: (mediante funcion divideestratos)
  tam2h.ps<-divideestratos(muestra2.ps, Nh)
  Trues<- tam2h.ps==0
  ceros<-sum(Trues)
}
### Cálculo de las probabilidades de Inclusión
#### De primer orden
TotalX<-sum(x)
Xtam2.ps<-x[muestra2.ps]
Pi.tam2.ps<-vector(len=tam2)
for (i in 1:tam2)
  Pi.tam2.ps[i]<-((N-tam2)/(N-1) )*(Xtam2.ps[i]/TotalX)+(tam2-1)/(N-1)
#### De segundo orden
Pij.tam2.ps<-matrix(nc=tam2, nr=tam2)
alfa<-Xtam2.ps/TotalX
  for (i in 1:tam2)
  for (j in 1:tam2)
  if (i==j)
Pij.tam2.ps[i,i]<-((N-tam2)/(N-1))*(Xtam2.ps[i]/TotalX)+(tam2-1)/(N-1)
  else
    Pij.tam2.ps[i,j]<-(( tam2-1)/(N-1))*
      ( ( (N-tam2)/(N-2) )*(alfa[i]+alfa[j])+ ( tam2-2)/(N-2) ) )
  resultado <- list(
tam1h=tam1h,tam2h=tam2h,muestralh=muestral,
muestra2.prop=muestra2.prop, Eti.muestra2.prop=Eti.muestra2.prop,
PItam1=PItam1,PItam2.prop = PITam2.prop,
PIcondi.tam2.prop=PIcondi.tam2.prop,
PijCondi.tam2.prop=PijCondi.tam2.prop, muestra2.dir=muestra2.dir,
PItam2.dir=PItam2.dir,
PijTam1= PijTam1, PijTam2.prop=PijTam2.prop, PijTam2.dir=PijTam2.dir,
muestral=muestral, PijCondi.tam2.dir=PijCondi.tam2.dir,
PIcondi.tam2.dir=PIcondi.tam2.dir,Eti.muestra2.dir= Eti.muestra2.dir,
Pi.tam2.ps= Pi.tam2.ps, Pij.tam2.ps= Pij.tam2.ps, tam2h.ps= tam2h.ps,
muestra2.ps=muestra2.ps)
#### EJEMPLO
#muestra.bifasicoST.M.Z(Fam1500.Nh, Fam1500[,3], tam1=100, tam2=50)
  resultado
}

muestra.bifasicoST.M.P<-function(Nh, x, tam1, tam2)
# Extracción de muestras en Muestreo bifasico aplicado a la estratificación y
#cálculo de probabilidades de inclusión: M.P= Mas.Poisson = TIPO=8.
{
  N1<-length(x)

```

```

        N<-sum(Nh)
        if (N != N1) stop("La longitud del vector x es diferente a N=suma de Nh")
ceros<-1
while (ceros>=1)
  {
    muestral <- sample(N,tam1)
    muestral <- sort(muestral)
### tam1h: contendrá tamaños de estratos de donde se van a extraer muestras
    tam1h<-divideestratos(muestral,Nh)
    Trues<-tam1h==0
    ceros<-sum(Trues)
  }
  datos.x.muestral<-x[muestral]
  tam2h<-afijacion.proporcional(tam2,tam1h)
  L<-length(Nh)
### Cálculo de las probabilidades de Inclusión
### Muestra 1- Primer orden y segundo orden.
  PItam1<-vector(len=tam1)
  Pij<-(tam1/N)*((tam1-1)/(N-1))
  PijTam1<-matrix(Pij, nr=tam1, nco=tam1)
  for (i in 1:tam1)
    {
      PItam1[i]<-tam1/N
      PijTam1[i,i]<-tam1/N ## En la diagonal se ponen las de 1er. orden
    }
## EMPIEZA PROCESO DE SELECCION EN SEGUNDA ETAPA ####
  Eti.muestra2.prop<-c()
  muestra2.prop<-c()
  PItam2.prop<-c()
  tam2h.P<-vector(len=L)
  PIcondi.tam2.prop<-c()
  datos.x.tam2<-c()
  for (h in 1:L)
  { ## for #1
    if (h !=1 )
      {
        a<-sum(tam1h[1:(h-1)])+1
        a2<-sum(tam2h[1:(h-1)])+1      }
    else
      { a<-1 ; a2<-1 }
    b<-sum(tam1h[1:h])
    b2<- sum(tam2h[1:h])
    len<-0
    while (len==0)
      {
        muestra2h.prov <- Etiquetas.Poisson(tam2h[h],datos.x.muestral[a:b])
        len<-length(muestra2h.prov)
      }
    Eti.muestra2.prop<-c(Eti.muestra2.prop,muestra2h.prov )
    muestrah<-muestral[a:b]
    muestra2h.prov <- muestrah[muestra2h.prov]
    muestra2.prop<-c(muestra2.prop,muestra2h.prov)
    tam2h.P[h]<-length(muestra2h.prov)
### Cálculo de las probabilidades de Inclusión en la segunda fase
    datos.x.tam2<-c(datos.x.tam2, x[muestra2h.prov])
    datos.x.tam2h<-x[muestra2h.prov]
#### De la Muestra 2
    PI.Partel<-PItam1[1:tam2h.P[h]]
    PIcondi.tam2<-vector(len=tam2h.P[h])
    sum(datos.x.muestral[a:b])>TotalXtam1
    for (i in 1:tam2h.P[h])
  PIcondi.tam2[i]<- tam2h.P[h]*datos.x.tam2h[i]/TotalXtam1
  PItam2.prop<- c(PItam2.prop, PI.Partel*PIcondi.tam2)
  PIcondi.tam2.prop<-c(PIcondi.tam2.prop,PIcondi.tam2 )

```

```

    } ## for #1
    tam2<-sum(tam2h.P)
    PijCondi.tam2.prop<-matrix(nr=tam2,nc=tam2)
#### De la Muestra 2 de segundo orden ( matriz Pij.cond)
    for (h in 1:L)
    { ## for #1
        if (h !=1 )
        {
            a<-sum(tamlh[1:(h-1)])+1
            a2<-sum(tam2h.P[1:(h-1)])+1
        }
        else
        {
            a<-1
            a2<-1
        }
        b<-sum(tamlh[1:h])
        b2<- sum(tam2h.P[1:h])
    }
    for (l in 1:L)
    {
        ## for-1
        if (h==1)
            for (i in a2:b2)
                for (j in a2:b2)
                    PijCondi.tam2.prop[i,j]<- PIcondi.tam2.prop[i]*PIcondi.tam2.prop[j]
        else ## else-1
        {if (l !=1 )
            {
                al<-sum(tamlh[1:(l-1)])+1
                al2<-sum(tam2h.P[1:(l-1)])+1
            }
            else
            {
                al<-1
                al2<-1
            }
            bl<-sum(tamlh[1:l])
            bl2<- sum(tam2h.P[1:l])
            for (i in a2:b2)
                for (j in al2:bl2)
                    PijCondi.tam2.prop[i,j]<- PIcondi.tam2.prop[i]*PIcondi.tam2.prop[j]
        } ## else-1
    } ## for-1
## Se incluyen las de primer orden en la diagonal principal:
    for (i in a2:b2) PijCondi.tam2.prop[i,i]<-PIcondi.tam2.prop[i]
} ## for #2
# Multiplicacion de matrices para obtener prob. de Inclusión en la segunda fase
PijTam2.prop<- PijTam1[1:tam2,1:tam2] * PijCondi.tam2.prop
# 1:tam2,1:tam2 ->Puesto que son todas iguales en la primera etapa.
#### Para otros estimadores:
#### 1. Para estimador Directo y Singh-Joader-Tracy:
Eti.muestra2.dir <- Etiquetas.Poisson(tam2,datos.x.muestral)
muestra2.dir <- muestral[Eti.muestra2.dir]
tam2<-length(Eti.muestra2.dir)
### Cálculo de las probabilidades de Inclusión
#### De primer orden de la Muestra 2
datos.x.tam2<-x[muestra2.dir]
PI.Partel<-PItam1[1:tam2]
PIcondi.tam2.dir<-vector(len=tam2)
sum(datos.x.muestral)->TotalXtam1
for (i in 1:tam2)
PIcondi.tam2.dir[i]<- tam2*datos.x.tam2[i]/TotalXtam1
PItam2.dir<-PI.Partel*PIcondi.tam2.dir
#### De segundo orden de la Muestra 2
PijCondi.tam2.dir<-matrix(nr=tam2,nc=tam2)
alfa<-datos.x.tam2/TotalXtam1

```

```

    for (i in 1:tam2)
  for (j in 1:tam2)
    {
      if (i ==j)
        PijCondi.tam2.dir[i,i]<- PIcondi.tam2.dir[i]
      else
        PijCondi.tam2.dir[i,j]<- PIcondi.tam2.dir[i]* PIcondi.tam2.dir[j]
    }
  PijTam2.dir<- PijTam1[1:tam2,1:tam2] * PijCondi.tam2.dir
#### 2. Silva and Skinner estimator (ps):
ceros<-1
while (ceros>=1)
{ muestra2.ps <-Etiquetas.Poisson(tam2,x)
  muestra2.ps<-sort(muestra2.ps)
### Ahora se pos-estratifica: (mediante funcion divideestratos)
  tam2h.ps<-divideestratos(muestra2.ps, Nh)
  Trues<- tam2h.ps==0
  ceros<-sum(Trues)
}
tam2<-length(muestra2.ps)
### Cálculo de las probabilidades de Inclusión
#### De primer orden
TotalX<-sum(x)
Xtam2.ps<-x[muestra2.ps]
Pi.tam2.ps<-vector(len=tam2)
for (i in 1:tam2)
  Pi.tam2.ps[i]<- tam2*Xtam2.ps[i]/TotalX
#### De segundo orden
Pij.tam2.ps<-matrix(nc=tam2, nr=tam2)
for (i in 1:tam2)
  for (j in 1:tam2)
    if (i==j)
      Pij.tam2.ps[i,i]<- Pi.tam2.ps[i]
    else
      Pij.tam2.ps[i,j]<-Pi.tam2.ps[i]*Pi.tam2.ps[j]

resultado <- list(
tam1h=tam1h, tam2h=tam2h.P, muestral1h=muestral1, muestra2.prop = muestra2.prop,
Eti.muestra2.prop=Eti.muestra2.prop,PItam1=PItam1,PItam2.prop = PITam2.prop,
PIcondi.tam2.prop=PIcondi.tam2.prop,PijCondi.tam2.prop=PijCondi.tam2.prop,
muestra2.dir=muestra2.dir,PItam2.dir=PItam2.dir, PijTam1=PijTam1,
PijTam2.prop=PijTam2.prop, PijTam2.dir=PijTam2.dir,
muestral1=muestral1, PijCondi.tam2.dir=PijCondi.tam2.dir,
PIcondi.tam2.dir=PIcondi.tam2.dir, Eti.muestra2.dir= Eti.muestra2.dir,
Pi.tam2.ps= Pi.tam2.ps, Pij.tam2.ps= Pij.tam2.ps      tam2h.ps= tam2h.ps,
muestra2.ps=muestra2.ps)
#### EJEMPLO
#muestra.bifasicoST.M.P(Fam1500.Nh, Fam1500[,3], tam1=100, tam2=50)
  resultado
}

muestra.bifasicoST<-function(TIPO,Nh, z, tam1, tam2)
# Para seleccionar un tipo de muestreo y llevarlo a la práctica. Se han
#implementado los siguientes TIPOS DE MUESTREOS:
#TIPO=1: M.M (Mas.Mas), TIPO=5: M.Z (Mas.Midzuno), TIPO=8: M.P (Mas.Poisson)
{
switch(TIPO,
  resultado<- muestra.bifasicoST.M.M(Nh, tam1, tam2),
  resultado<-Introduce.otro.muestreo,
  resultado<- Introduce.otro.muestreo,
  resultado<- Introduce.otro.muestreo,
  resultado<- muestra.bifasicoST.M.Z(Nh, z, tam1, tam2),

```

```

        resultado<- Introduce.otro.muestreo,
        resultado<- Introduce.otro.muestreo,
        resultado<- muestra.bifasicoST.M.P(Nh, z, tam1, tam2)
    )
resultado
}

DatosST<-function(Tipo.Muestreo,Tipo.Fd,datosy,datosx,Nh,tam1,tam2,beta,Qy)
# Estimadores y sus varianzas en M.bifásico aplicado a la estratificación
{
    N<-sum(Nh)
## 1. Captura de informacion: Tamaños muestrales, etiquetas y Prob. inclusión
    muestras <- muestra.bifasicoST(Tipo.Muestreo, Nh, datosx, tam1, tam2)
## Para estimador propuesto (1)
    tam1h <-muestras$tam1h
    tam2h <-muestras$tam2h
    muestralh <-muestras$muestralh
    muestra2.prop <-muestras$muestra2.prop
    Eti.muestra2.prop <-muestras$Eti.muestra2.prop
    PItam1 <-muestras$PItam1
    PItam2.prop <-muestras$PItam2.prop
    PIcondi.tam2.prop <-muestras$PIcondi.tam2.prop
    PijCondi.tam2.prop <-muestras$PijCondi.tam2.prop
    tam2.prop<-length(PItam2.prop)
    PijTam1<-muestras$PijTam1
    PijTam2.prop<-muestras$PijTam2.prop
## Para estimadores directo (2) y de Singh (3)
    muestra2.dir <- muestras$muestra2.dir
    PItam2.dir <-muestras$PItam2.dir
    PijTam2.dir <- muestras$PijTam2.dir
    PijCondi.tam2.dir <- muestras$PijCondi.tam2.dir
    muestral <- muestras$muestral
    PIcondi.tam2.dir <- muestras$PIcondi.tam2.dir
    Eti.muestra2.dir <- muestras$Eti.muestra2.dir
## Para estimador de Silva y Skinner (4)
    tam2h.ps <-muestras$tam2h.ps
    muestra2.ps <-muestras$muestra2.ps
    Pi.tam2.ps <-muestras$Pi.tam2.ps
    Pij.tam2.ps <-muestras$Pij.tam2.ps
### Aquí se pone el nuevo tam2 en el caso de M. de Poisson
    length(muestra2.prop)->tam2.prop
#### 2. Obtencion de estimadores.
## Para estimador propuesto (1)
    Ytam2.prop<-datosy[muestra2.prop]
    Q.prop<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PItam2.prop,Ytam2.prop)
## Para estimadores directo (2) y de Singh (3)
    # Directo
    Ytam2.dir<-datosy[muestra2.dir]
    Q.dir<-cuantil.HT.Rapido(Tipo.Fd,beta,N,PItam2.dir,Ytam2.dir)
    # Singh
    Xtam2.dir<-datosx[muestra2.dir]
    Xtam1<-datosx[muestral]
    Qx.1<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=PItam1,vector=Xtam1)
    Q.sjt<- cuantil.SJT(beta, Qx.1, Ytam2.dir, Xtam1, Xtam2.dir)

    tam2.dir<-length(Ytam2.dir)      ## Cambia del resto en M.Poisson
## Para calcular p11:
    Qx.2<-cuantil.HT.Rapido(TipoFd=1,beta,tam1,
        Pi=rep(tam2.dir/tam1,ti=tam2.dir), vector=Xtam2.dir)
    Qy.2<-cuantil.HT.Rapido(TipoFd=1,beta, tam1,
        Pi=rep(tam2.dir/tam1,ti=tam2.dir),vector=Ytam2.dir)

```

```

      Trues.x2<- Xtam2.dir<=Qx.2
      p11<- sum( Ytam2.dir[Trues.x2]<=Qy.2 ) / tam2.dir
## Para estimador de Silva y Skinner (4)
      Ytam2.ps<-datosy[muestra2.ps]
      Q.ps<-cuantil.PS(beta,Nh,nh=tam2h.ps,Pi=Pi.tam2.ps,vectory=Ytam2.ps)
#### 3. Obtencion de varianzas.
# "Estimacion" de la densidad de y
      fy<-densidad(datosy, Qy)
## Para estimador directo (2)
##### Suma 1.
PijTam1.v2<-PijTam1[Eti.muestra2.dir,Eti.muestra2.dir]
PITam1.v2<-PITam1[Eti.muestra2.dir]
col.PiTam1.v2<-matrix(PITam1.v2,nc=1,nr=tam2.dir)
matrix1.Tam21<- (PijTam1.v2- col.PiTam1.v2**PITam1.v2)/PijTam2.dir
matrix2.Tam21<- (Delta(Q.dir,Ytam2.dir) / col.PiTam1.v2 ) **
                 (Delta(Q.dir,Ytam2.dir) / PITam1.v2 )
suma.Tam21<-sum(matrix1.Tam21*matrix2.Tam21)
##### Suma 2.
col.PIcondi.tam2.dir<-matrix(PIcondi.tam2.dir,nc=1,nr=tam2.dir)
col.PITam2.dir<-matrix(PITam2.dir,nc=1,nr=tam2.dir)
matrix1.Tam2 <- (PijCondi.tam2.dir- col.PIcondi.tam2.dir **
                PIcondi.tam2.dir)/PijCondi.tam2.dir
matrix2.Tam2<- (Delta(Q.dir,Ytam2.dir) / col.PITam2.dir ) **
                (Delta(Q.dir,Ytam2.dir) / PITam2.dir )
suma.Tam2<-sum(matrix1.Tam2*matrix2.Tam2)
var.est.dir<-(1/fy^2)*(1/N^2)*( suma.Tam21+ suma.Tam2)
## Para estimador propuesto (1)
#### Suma 1
## En M. Poisson
PijTam1.v2<-PijTam1[1:tam2.prop, 1:tam2.prop]
PITam1.v2<-PITam1[1:tam2.prop]
col.PiTam1.v2<-matrix(PITam1.v2,nc=1,nr=tam2.prop)
matrix1.Tam21<- (PijTam1.v2- col.PiTam1.v2**PITam1.v2)/PijTam2.prop
matrix2.Tam21<- (Delta(Q.prop,Ytam2.prop) / col.PiTam1.v2 )**
                 (Delta(Q.prop,Ytam2.prop) / PITam1.v2 )
suma.Tam21.prop<-sum(matrix1.Tam21*matrix2.Tam21)
### Suma 2 # Se suman las unidades del mismo estrato.
L<-length(Nh)
suma.Tam2.prop<-0
col.PIcondi.tam2.prop<-matrix(PIcondi.tam2.prop,nc=1,nr=tam2.prop)
col.PITam2.prop<-matrix(PITam2.prop,nc=1,nr=tam2.prop)
Delta.tam2<-(PijCondi.tam2.prop- (col.PIcondi.tam2.prop**PIcondi.tam2.prop))
for (h in 1:L)
  { if (h !=1 )      a2<-sum(tam2h[1:(h-1)])+1
    else             a2<-1
    b2<- sum(tam2h[1:h])
    matrix1.Tam2.prop <-Delta.tam2[a2:b2,a2:b2]/PijCondi.tam2.prop[a2:b2,a2:b2]
    col.h<-matrix( (Delta(Q.prop,Ytam2.prop[a2:b2])/col.PITam2.prop[a2:b2,1]),
                  nr=tam2h[h],nc=1)
    matrix2.Tam2.prop<-
col.h**(Delta(Q.prop,Ytam2.prop[a2:b2])/PITam2.prop[a2:b2])
    suma.Tam2.prop<-suma.Tam2.prop+sum(matrix1.Tam2.prop*matrix2.Tam2.prop)
  }
var.est.prop<-(1/fy^2)*(1/N^2)*( suma.Tam21.prop + suma.Tam2.prop)
## Para estimador Singh (3)
var.est.sjt<-
(beta*(1-beta)/fy^2)*((1/tam2.dir)-(1/N)*(p11/(beta*(1-beta))-1)^2)
#### Para estimador Silva y Skinner (PS) (4)
Salida<-F.distribucion.PS(Nh,nh=tam2h.ps,t=Q.ps,Pi.tam2.ps,vectory=Ytam2.ps)
Fd.g<-Salida$Fd.g
Nh.gorro<-Salida$Nh.gorro

```

```

vector.Fd.g<-c()
vector.Nh<-c()
vector.Nh.gorro<-c()
for (h in 1:L)
  { vector.Fd.g<- c(vector.Fd.g, rep(Fd.g[h],ti=tam2h.ps[h]) )
    vector.Nh<-c(vector.Nh,rep(Nh[h],tim=tam2h.ps[h]) )
    vector.Nh.gorro<-c(vector.Nh.gorro, rep(Nh.gorro[h], tim=tam2h.ps[h]) )
  }
vector.a<-Delta(Q.ps,Ytam2.ps) - vector.Fd.g
var.est.ps1<- (1/fy^2)*varianza.PS(N,vector.a, Pi.tam2.ps, Pij.tam2.ps)
vector.a.2<- vector.a*vector.Nh/vector.Nh.gorro
var.est.ps2<- (1/fy^2)*varianza.PS(N,vector.a.2, Pi.tam2.ps, Pij.tam2.ps)
### Resultado del Programa
resultado<-list(Q.prop=Q.prop,Q.dir=Q.dir,Q.ps=Q.ps,Q.sjt=Q.sjt,var.est.prop =
var.est.prop,var.est.dir=var.est.dir,var.est.ps1=var.est.ps1,var.est.ps2=
var.est.ps2,var.est.sjt=var.est.sjt,PITam2.prop=PITam2.prop)
resultado
## EJEMPLO
#DatosST(Tipo.Muestreo=1,Tipo.Fd=1,datosy=Fam1500.st[,1],datosx=Fam1500.st[,1],
#Nh=Fam1500.Nh, tam1=100, tam2=50, beta=0.5, Qy=median(Fam1500.st[,1]))
}

Simula.bifasicoST<-function(NombreP,Poblacion,Nh,Tipo.Muestreo,Tipo.Fd,
variableY, variableZ, tam1, tam2, replicas, beta)
# Realiza la simulación en muestreo bifásico aplicado a la estratificación
{
  N<-sum(Nh)
  variables<-length(Poblacion[1,])
  datosx<-Poblacion[,variableZ]
  datosy<-Poblacion[,variableY]
  salidaError <- paste("C:\\BifasicoST\\errP_",NombreP, "TM_", Tipo.Muestreo,
    "TFd_",Tipo.Fd, "tam1_", tam1,"beta_", beta, ".txt", sep = "")
  salidaEstimador <- paste("C:\\BifasicoST\\estP_",NombreP,"TM_",Tipo.Muestreo,
    "TFd_",Tipo.Fd,"tam_", tam1,"y", tam2,"beta_", beta, ".txt", sep = "")
  salidaBoxEst <- paste("C:\\BifasicoST\\BoxEstP_",NombreP,"TM_",Tipo.Muestreo,
    "TFd_",Tipo.Fd,"tam_", tam1,"y", tam2,"beta_", beta, ".txt", sep = "")
  salidaBoxVar <- paste("C:\\BifasicoST\\BoxVarP_",NombreP,"TM_",Tipo.Muestreo,
    "TFd_",Tipo.Fd,"tam_", tam1,"y", tam2,"beta_", beta, ".txt", sep = "")
  # AHORA EMPIEZA EL PROCESO DE REPLICACIÓN #
  Qy<-cuantil.HT.Rapido(TipoFd=1,beta,N,Pi=1,vectory=datosy)
  ecm.est.dir <- 0 ; ses.est.dir <- 0 ; Coverage.dir <-0 ; values.var.dir<-c()
    sum1.dir<-0 ; sum2.dir<-0 ; len.dir<-0
  ecm.est.prop <-0 ; ses.est.prop <-0 ; Coverage.prop<-0 ; values.var.prop<-c()
    sum1.prop<-0 ; sum2.prop<-0 ; len.prop<-0
  ecm.est.ps <- 0 ; ses.est.ps <- 0 ; Coverage.ps1 <-0 ; Coverage.ps2<-0
  values.var.ps1<-c() ; values.var.ps2<-c() ; sum1.ps<-0 ; sum2.ps<-0
    len.ps1<-0 ; len.ps2<-0
  ecm.est.sjt <- 0 ; ses.est.sjt <- 0 ; Coverage.sjt<-0 ; values.var.sjt<-c()
    sum1.sjt<-0 ; sum2.sjt<-0 ; len.sjt<-0
    sum.PITam2<-0
  for(b in 1:replicas)
  {
  datos<- DatosST(Tipo.Muestreo,Tipo.Fd, datosy, datosx, Nh, tam1, tam2, beta,Qy)
    Q.prop<-datos$Q.prop
    Q.dir<-datos$Q.dir
    Q.ps<-datos$Q.ps
    Q.sjt<-datos$Q.sjt
    var.est.dir <-datos$var.est.dir
    var.est.prop <-datos$var.est.prop
    var.est.ps1 <-datos$var.est.ps1
    var.est.ps2 <-datos$var.est.ps2
    var.est.sjt <-datos$var.est.sjt
  }
}

```

```

PITam2.prop<- datos$PITam2.prop
sum.PITam2<-sum.PITam2+ (1/N)*sum(1/PITam2.prop)
### 1 ### El estimador directo para comparar:
    ecm.est.dir <- ecm.est.dir + (Q.dir - Qy)^2
    ses.est.dir <- ses.est.dir + (Q.dir - Qy)
    Box.dir<-(Q.dir - Qy) /Qy
    Estado<-abs((Q.dir-Qy)/sqrt(abs(var.est.dir)))<=1.96
if (is.na(Estado)==TRUE) valor.Cove<-0
    else
    if (Estado==TRUE) valor.Cove<-1 else valor.Cove<-0
Coverage.dir<-Coverage.dir+ valor.Cove
values.var.dir<-c(values.var.dir,var.est.dir )
sum1.dir<- sum1.dir+Q.dir
sum2.dir<-sum2.dir+Q.dir^2
len.dir<-len.dir+sqrt(abs(var.est.dir))
### 2 ### Estimador estratificado (propuesto)
    ecm.est.prop <- ecm.est.prop + (Q.prop - Qy)^2
    ses.est.prop <- ses.est.prop + (Q.prop - Qy)
    Box.prop<-(Q.prop - Qy) /Qy
    Estado<-abs((Q.prop-Qy)/sqrt(abs(var.est.prop)))<=1.96
if (is.na(Estado)==TRUE) valor.Cove<-0
    else
    if (Estado==TRUE) valor.Cove<-1 else valor.Cove<-0
Coverage.prop<-Coverage.prop+ valor.Cove
values.var.prop<-c(values.var.prop,var.est.prop)
sum1.prop<- sum1.prop+Q.prop
sum2.prop<-sum2.prop+Q.prop^2
len.prop<-len.prop+sqrt(abs(var.est.prop))
### 3 ### Estimador SJT
    ecm.est.sjt <- ecm.est.sjt + (Q.sjt - Qy)^2
    ses.est.sjt <- ses.est.sjt + (Q.sjt - Qy)
    Box.sjt<-(Q.sjt - Qy) /Qy
    Estado<-abs((Q.sjt-Qy)/sqrt(abs(var.est.sjt)))<=1.96
if (is.na(Estado)==TRUE) valor.Cove<-0
    else
    if (Estado==TRUE) valor.Cove<-1 else valor.Cove<-0
Coverage.sjt<-Coverage.sjt+ valor.Cove
values.var.sjt<-c(values.var.sjt,var.est.sjt)
sum1.sjt<- sum1.sjt+Q.sjt
sum2.sjt<-sum2.sjt+Q.sjt^2
len.sjt<-len.sjt+sqrt(abs(var.est.sjt))
### 4 ### Estimador Silva Skinner (PS)
    ecm.est.ps <- ecm.est.ps + (Q.ps - Qy)^2
    ses.est.ps <- ses.est.ps + (Q.ps - Qy)
    Box.ps<-(Q.ps - Qy) /Qy
    Estado1<-abs((Q.ps-Qy)/sqrt(abs(var.est.ps1)))<=1.96
    Estado2<-abs((Q.ps-Qy)/sqrt(abs(var.est.ps2)))<=1.96
if (is.na(Estado1)==TRUE) valor.Cove1<-0
    else
    if (Estado1==TRUE) valor.Cove1<-1 else valor.Cove1<-0
if (is.na(Estado2)==TRUE) valor.Cove2<-0
    else
    if (Estado2==TRUE) valor.Cove2<-1 else valor.Cove2<-0
Coverage.ps1<-Coverage.ps1+ valor.Cove1
Coverage.ps2<-Coverage.ps2+ valor.Cove2
values.var.ps1<-c(values.var.ps1,var.est.ps1)
values.var.ps2<-c(values.var.ps2,var.est.ps2)
sum1.ps<- sum1.ps+Q.ps
sum2.ps<-sum2.ps+Q.ps^2
len.ps1<-len.ps1+sqrt(abs(var.est.ps1))
len.ps2<-len.ps2+sqrt(abs(var.est.ps2))

```

```

### Se escribe el valor de los estimadores
  linea <- paste(b, Qy, Q.dir, Q.prop, Q.sjt, Q.ps, var.est.dir,
var.est.prop, var.est.sjt, var.est.ps1, var.est.ps2, sep = "\t")
  write(linea, file = salidaEstimador, ncolumns = 1, append = T)
  linea <- paste(b, Box.dir, Box.prop, Box.sjt, Box.ps, sep = "\t")
  write(linea, file = salidaBoxEst, ncolumns = 1, append = T)
} # end replicas
E.Pitam2<- (1/replicas)*sum.PITam2
##### Obtencion de longitud media
  E.len.dir <-(1/replicas)*2*1.96*len.dir
  E.len.prop<-(1/replicas)*2*1.96*len.prop
  E.len.sjt <-(1/replicas)*2*1.96*len.sjt
  E.len.ps1<-(1/replicas)*2*1.96*len.ps1
  E.len.ps2<-(1/replicas)*2*1.96*len.ps2
### Obtencion de ECM, RB (SR), RE (ER) y RRMSE de varianzas
  ## Estimador Directo
var.dir<-(1/replicas)*sum2.dir - ((1/replicas)*sum1.dir)^2
MSE.var.dir<- (1/replicas)*sum((values.var.dir-var.dir)^2)
RB.var.dir<- 100*(1/replicas)*(1/var.dir)*sum(values.var.dir-var.dir)
RRMSE.var.dir<- 100*sqrt(MSE.var.dir)/var.dir
RE.var.dir<-1
  ## Estimador propuesto
var.prop<-(1/replicas)*sum2.prop - ((1/replicas)*sum1.prop)^2
MSE.var.prop<- (1/replicas)*sum((values.var.prop-var.prop)^2)
RB.var.prop<- 100*(1/replicas)*(1/var.prop)*sum(values.var.prop-var.prop)
RRMSE.var.prop<- 100*sqrt(MSE.var.prop)/var.prop
RE.var.prop<- MSE.var.prop/MSE.var.dir
  ## Estimador de SJT
var.sjt<-(1/replicas)*sum2.sjt - ((1/replicas)*sum1.sjt)^2
MSE.var.sjt<- (1/replicas)*sum((values.var.sjt-var.sjt)^2)
RB.var.sjt<- 100*(1/replicas)*(1/var.sjt)*sum(values.var.sjt-var.sjt)
RRMSE.var.sjt<- 100*sqrt(MSE.var.sjt)/var.sjt
RE.var.sjt<- MSE.var.sjt/MSE.var.dir
  ## Estimadores PS
var.ps<-(1/replicas)*sum2.ps - ((1/replicas)*sum1.ps)^2
MSE.var.ps1<- (1/replicas)*sum((values.var.ps1-var.ps)^2)
MSE.var.ps2<- (1/replicas)*sum((values.var.ps2-var.ps)^2)
RB.var.ps1<- 100*(1/replicas)*(1/var.ps)*sum(values.var.ps1-var.ps)
RB.var.ps2<- 100*(1/replicas)*(1/var.ps)*sum(values.var.ps2-var.ps)
RRMSE.var.ps1<- 100*sqrt(MSE.var.ps1)/var.ps
RRMSE.var.ps2<- 100*sqrt(MSE.var.ps2)/var.ps
RE.var.ps1<- MSE.var.ps1/MSE.var.dir
RE.var.ps2<- MSE.var.ps2/MSE.var.dir
### Obtencion de ECM, RB, RE y RRMSE de estimadores
  ## Estimador Directo
MSE.est.dir<-(1/replicas)*ecm.est.dir
RB.est.dir<-100*(1/replicas)*(1/Qy)*ses.est.dir
RRMSE.est.dir<-100*sqrt(MSE.est.dir)/Qy
RE.est.dir<-1
  ## Estimador propuesto
MSE.est.prop<-(1/replicas)*ecm.est.prop
RB.est.prop<-100*(1/replicas)*(1/Qy)*ses.est.prop
RRMSE.est.prop<-100*sqrt(MSE.est.prop)/Qy
RE.est.prop<-MSE.est.prop/MSE.est.dir
  ## Estimador de SJT
MSE.est.sjt<-(1/replicas)*ecm.est.sjt
RB.est.sjt<-100*(1/replicas)*(1/Qy)*ses.est.sjt
RRMSE.est.sjt<-100*sqrt(MSE.est.sjt)/Qy
RE.est.sjt<-MSE.est.sjt/MSE.est.dir
  ## Estimador PS
MSE.est.ps<-(1/replicas)*ecm.est.ps

```

```

RB.est.ps<-100*(1/replicas)*(1/Qy)*ses.est.ps
RRMSE.est.ps<-100*sqrt(MSE.est.ps)/Qy
RE.est.ps<-MSE.est.ps/MSE.est.dir
### Obtención de Coverage.
    Cove.dir<-100*(1/replicas)* Coverage.dir
    Cove.prop<-100*(1/replicas)* Coverage.prop
    Cove.sjt<-100*(1/replicas)* Coverage.sjt
    Cove.ps1<-100*(1/replicas)* Coverage.ps1
    Cove.ps2<-100*(1/replicas)* Coverage.ps2
### Se escribe la matriz de resultados
linea <- paste(round(E.Pitam2,4), tam2,
    round(RB.est.dir,4),    round(RRMSE.est.dir,4),    round( RE.est.dir,4),
    round(RB.est.prop,4),  round(RRMSE.est.prop,4),    round(RE.est.prop,4),
    round(RB.est.sjt,4) ,  round(RRMSE.est.sjt,4),    round( RE.est.sjt,4),
    round(RB.est.ps,4),    round(RRMSE.est.ps,4),    round( RE.est.ps,4),
    round(RB.var.dir,4),   round(RRMSE.var.dir,4),    round(RE.var.dir,4),
round(Cove.dir,4),    round( E.len.dir,4),
    round(RB.var.prop,4),  round(RRMSE.var.prop,4),    round(RE.var.prop,4),
round(Cove.prop,4),  round( E.len.prop,4),
    round( RB.var.sjt,4),  round(RRMSE.var.sjt,4),    round( RE.var.sjt,4),
round(Cove.sjt,4),  round( E.len.sjt,4),
    round(RB.var.ps1,4),  round(RRMSE.var.ps1,4),    round(RE.var.ps1,4),
round( Cove.ps1,4),  round(E.len.ps1,4),
    round(RB.var.ps2,4),  round(RRMSE.var.ps2,4),    round(RE.var.ps2,4),
round( Cove.ps2,4),  round(E.len.ps2,4),    sep = "\t")
write(linea, file = salidaError, ncolumns = 1, append = T)
##### Box-varianzas
    for (b in 1:replicas)
    {
        Box.dir<- (values.var.dir[b] -var.dir )/var.dir
        Box.prop<- (values.var.prop[b]-var.prop)/var.prop
        Box.sjt<- (values.var.sjt[b] -var.sjt )/var.sjt
        Box.ps1<- (values.var.ps1[b] -var.ps )/var.ps
        Box.ps2<- (values.var.ps2[b] -var.ps )/var.ps
        linea <- paste(b, Box.dir, Box.prop, Box.sjt, Box.ps1, Box.ps2, sep = "\t")
        write(linea, file = salidaBoxVar, ncolumns = 1, append = T)
    }
##### EJEMPLO
#Simula.bifasicoST(NombreP=c("Fam1500"), Población=Fam1500.st, Nh=Fam1500.Nh,
#Tipo.Muestreo=5, Tipo.Fd=1, variableY=1, #variableZ=2, tam1=150, tam2=50,
#replicas=100, beta=0.5)
}

```

C.3.5. Estimación modelo-asistida usando verosimilitud empírica

```
## ESTIMADORES DE CALIBRACION ##
```

```
cuantil.cal.1va.3ptos
```

```

<-function(beta,puntos.t,to=vector.t,s=etiquetas,d=1/Pi,y,x,q=rep(1,N) )
# Estimadores de calibración. 1 variable auxiliar y 3 puntos.
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<- calibracionfd2(puntos.t, to, s, d, y, x, q )
if (Fd.n>beta)

```

```

    { while(B-A != 1)
      { t<-datos[M]
        Fd<- calibracionfd2(puntos.t, to, s, d, y, x, q )
        if (Fd>beta)
          { B<-M
            M<-floor((A+B)/2)
            if (A==B) B<-A+1
            t<-datos[B]
          }
        else if (Fd<beta)
          { A<-M
            M<-floor((A+B)/2)
            if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
          }
        else {B<-2 ; A<-1 ; t<-datos[M]}
      }
    }
  }
else { t<-datos[n] }
t
}

```

cuantil.cal.2va.3ptos<-

```

function(beta, puntos.t,to=c(valores.t,1),s=etiquetas,d=1/Pi,y,X=x2,q=rep(1,N))
# Estimadores de calibración. 2 variables auxiliares y 3 puntos.

```

```

{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
Fd.n<- calibracionfd1(puntos.t, to, s, d, y, X, q)
if (Fd.n>beta)
  { while(B-A != 1)
    { t<-datos[M]
      Fd<- calibracionfd1(puntos.t, to, s, d, y, X, q)
      if (Fd>beta)
        { B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        { A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
else { t<-datos[n] }
t
}

```

ESTIMADORES DE VEROSIMILITUD EMPÍRICA

cuantil.MA.RX.1pto<-function(beta, N, t0, Pi, muestray, muestrax,x)

```

# Estimadores modelo-asistidos con coeficiente R(ver Rao et al, 1990) y 1 punto

```

```

{
length(muestray)->n

```

```

sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
resultados<- F.distribucion.ve.RX.vec(N,t,t0, Pi,muestray,muestrax, x)
Estado.NR<-resultados$Estado.NR
if (Estado.NR==0) t<-0
else
{
Fd.n<-resultados$Fd
if (Fd.n>beta)
{
while(B-A != 1)
{
t<-datos[M]
resultados<- F.distribucion.ve.RX.vec(N,t,t0, Pi,muestray,muestrax, x)
Estado.NR<-resultados$Estado.NR
Fd<-resultados$Fd
if (Fd>beta)
{
B<-M
M<-floor((A+B)/2)
if (A==B) B<-A+1
t<-datos[B]
}
else if (Fd<beta)
{
A<-M
M<-floor((A+B)/2)
if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
}
}
else {B<-2 ; A<-1 ; t<-datos[M]}
}
}
else { t<-datos[n] }
}
list(t=t, Estado.NR=Estado.NR )
}

```

```

cuantil.MA.BX.lpto<-function(beta, N, t0, Pi, muestray, muestrax,x)
# Estimadores modelo-asistidos con coeficiente B de regresión y 1 punto
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
resultados<-F.distribucion.ve.BX.vec(N,t, t0, Pi,muestray,muestrax, x)
Estado.NR<-resultados$Estado.NR
if (Estado.NR==0) t<-0
else
{
Fd.n<- resultados$Fd
if (Fd.n>beta)
{
while(B-A != 1)
{
t<-datos[M]
resultados<- F.distribucion.ve.BX.vec(N, t, t0, Pi,muestray,muestrax, x)
Estado.NR<-resultados$Estado.NR

```

```

      Fd<-resultados$Fd
      if (Fd>beta)
        {
          B<-M
          M<-floor((A+B)/2)
          if (A==B) B<-A+1
          t<-datos[B]
        }
      else if (Fd<beta)
        {
          A<-M
          M<-floor((A+B)/2)
          if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
      else {B<-2 ; A<-1 ; t<-datos[M]}
    }
  }
else { t<-datos[n] }
}
list(t=t, Estado.NR=Estado.NR)
}

```

```

cuantil.MA.RX.3ptos<-function(beta, N, valores.t, Pi, muestray, muestrax,x)
# Estimadores modelo-asistidos con coeficiente R y 3 puntos.
{
  length(muestray)->n
  sort(muestray)->datos
  t<-datos[n]
  A<-1
  B<-n
  M<-floor(n/2)
  resultados<- F.distribucion.RX.Tt.vec(N,t,valores.t, Pi,muestray,muestrax, x)
  Fd.n<- resultados$Fd
  if (Fd.n>beta)
    { while(B-A != 1)
      {
        t<-datos[M]
resultados<- F.distribucion.RX.Tt.vec(N,t,valores.t, Pi,muestray,muestrax, x)
        Fd<-resultados$Fd
        if (Fd>beta)
          {
            B<-M
            M<-floor((A+B)/2)
            if (A==B) B<-A+1
            t<-datos[B]
          }
        else if (Fd<beta)
          {
            A<-M
            M<-floor((A+B)/2)
            if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
          }
        else {B<-2 ; A<-1 ; t<-datos[M]}
      }
    }
  }
else { t<-datos[n] }
t
}

```

```

cuantil.MA.BX.3ptos<-function(beta, N, valores.t, Pi, muestray, muestrax,x)
# Estimadores modelo-asistidos con coeficiente B de regresión y 3 puntos.
{
  length(muestray)->n
  sort(muestray)->datos

```

```

t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
resultados<- F.distribucion.BX.T.vec(N, t, valores.t, Pi, muestray, muestrax, x)
Fd.n<-resultados$Fd
if (Fd.n>beta)
  {
    while(B-A != 1)
      {
        t<-datos[M]
resultados<- F.distribucion.BX.T.vec(N, t, valores.t, Pi, muestray, muestrax, x)
Fd<-resultados$Fd
if (Fd>beta)
  {
    B<-M
    M<-floor((A+B)/2)
    if (A==B) B<-A+1
    t<-datos[B]
  }
else if (Fd<beta)
  {
    A<-M
    M<-floor((A+B)/2)
    if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
  }
else {B<-2 ; A<-1 ; t<-datos[M]}
      }
  }
else { t<-datos[n] }
t
}

```

```

cuantil.MC.1va.1pto<-function(beta, N, t0, Pi, muestray, muestrax,x)
# Estimadores Modelo-Calibrados. 1 variable auxiliar y 1 punto.
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
resultados<- F.distribucion.Mve.2t.vec(N,t, t0, Pi, muestray, muestrax, x)
Estado.NR<-resultados$Estado.NR
if (Estado.NR==0) t<-0
else
{
Fd.n<-resultados$Fd
if (Fd.n>beta)
  {
    while(B-A != 1)
      {
        t<-datos[M]
resultados<- F.distribucion.Mve.2t.vec(N,t, t0, Pi, muestray, muestrax, x)
Estado.NR<-resultados$Estado.NR
Fd<- resultados$Fd
if (Fd>beta)
  {
    B<-M
    M<-floor((A+B)/2)
    if (A==B) B<-A+1
    t<-datos[B]
  }
      }
  }
}
}

```

```

        else if (Fd<beta)
        {
            A<-M
            M<-floor((A+B)/2)
            if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
        }
        else {B<-2 ; A<-1 ; t<-datos[M]}
    }
}
else { t<-datos[n] }
}
list( t=t, Estado.NR=Estado.NR )
}

cuantil.MC.2va.1pto<-function(beta, N, t0, Pi, muestray, muestrax2,x2)
# Estimadores Modelo-Calibrados. 2 variables auxiliares y 1 punto.
{
length(muestray)->n
sort(muestray)->datos
t<-datos[n]
A<-1
B<-n
M<-floor(n/2)
resultados<- F.distribucion.Mve.Jvar.vec(N, t,t0, Pi,muestray,muestrax2, x2)
Estado.NR<-resultados$Estado.NR
if (Estado.NR==0) t<-0
else
{
Fd.n<-resultados$Fd
if (Fd.n>beta)
{
while(B-A != 1)
{
t<-datos[M]
resultados<- F.distribucion.Mve.Jvar.vec(N, t,t0, Pi,muestray,muestrax2, x2)
Estado.NR<-resultados$Estado.NR
Fd<-resultados$Fd
if (Fd>beta)
{
B<-M
M<-floor((A+B)/2)
if (A==B) B<-A+1
t<-datos[B]
}
else if (Fd<beta)
{
A<-M
M<-floor((A+B)/2)
if (A==B) { B<-A+1; t<-datos[M] } else t<-datos[B]
}
else {B<-2 ; A<-1 ; t<-datos[M]}
}
}
}
else { t<-datos[n] }
}
list(t=t, Estado.NR=Estado.NR)
}

```

Obtención de cuantiles y diversas medidas de pobreza

```

Quantiles.MA.Bootstrap<-function(n,y,x, vector.beta)
# Obtiene cuantiles y sus varianzas mediante Bootstrap asumiendo MAS
# "n": Tamaño de la muestra
# "t": Punto en el que se calcula la Fd.

```

```

# "y": Característica de interés
# "x": Variable auxiliar. Cada fila es una variable.
# "t0": cuantil para utilizar en las restricciones.
{
##### Preparacion de datos y obtencion de probabilidades de inclusion.
  N<-length(y)
  k<-N%n
  if (k != 0) stop("Introduce un valor de 'n' que sea divisor de 'N' ")
Estado<-0
while (Estado==0)
{
  vector.Est<-c()
  sample(N,n)->etiquetas
  muestray<-y[etiquetas]
  is.vector(x)->valor.X
  if (valor.X==T) # Si solo se tiene 1 v.a.:
  {
    Jvar<-1
    muestrax<-x[etiquetas]
    else
    {
      Jvar<-length(x[,1])
      muestrax2<-x[,etiquetas]
      muestrax<-muestrax2[1,] # muestrax no es una matriz, es un vector.
    }
    if (Jvar!=1)
    {
      x2<-x
      x<-x[1,] # x es un vector.
    }
    Prob.inclusion.Mas(N,n)->prob
    Pi<-prob$Pi
    Pij<-prob$Pij
# Obtención de valores t0 para usar en estimadores M.Calibrados y M.A con 1pto.
len.beta<-length(vector.beta)
## Para una variable
  muestraxB<-as.matrix(muestrax)
  muestraxB<-t(muestraxB)
  datosx<-as.matrix(x)
  datosx<-t(datosx)
### Rx
  Rgorro<-matrix(nr=1,nc=1)
  Rgorro[1,1]<-sum(muestray/Pi)/sum(muestraxB[1,]/Pi)
  g.N.Rx<-t(Rgorro)%*%datosx
  vector.t0Rx<-c()
for (p in 1:len.beta)
vector.t0Rx<-c(vector.t0Rx, cuantil(1,vector.beta[p],N,Pi=1,muestray=g.N.Rx) )
### Bx
  betagorrox<-matrix(nr=1,nc=1)
  betagorrox[1,1]<-sum((1/Pi)*muestraxB*muestray)/sum((1/Pi)*muestraxB^2)
  g.N.Bx<-t(betagorrox)%*%datosx
  vector.t0Bx<-c()
  for (p in 1:len.beta) vector.t0Bx<-c(vector.t0Bx, cuantil(1,
vector.beta[p], N, Pi=1,muestray=g.N.Bx) )
## Si hay dos variables auxiliares
if (Jvar !=1)
{
### Bx2
  xs<-t(muestrax2)*(1/Pi)
  td<-muestrax2%*%xs
  ts<-solve(td)
}
}

```

```

        r<-muestrax2%*((1/Pi)*muestray)
        betagorrox2<-ts%*(r)
        g.N.Bx2<-t(betagorrox2)%*%x2
        vector.t0Bx2<-c()
        for (p in 1:len.beta) vector.t0Bx2<-c(vector.t0Bx2, cuantil(1,
vector.beta[p], N, Pi=1,muestray=g.N.Bx2) )
    }
# Para los modelo-Calibrados se usa Qy(b).est: Se estima a partir de la muestra
    vector.t0MC<-c()
    for (p in 1:len.beta) vector.t0MC<-c(vector.t0MC, cuantil(1,
vector.beta[p], N, Pi=Pi, muestray=muestray) )
# Mmatrices de cuantiles y Varianzas. Por columnas se ponen los estimadores.
    valores.t<-c(0.25,0.5,0.75) ## Para usar en las restricciones
    Matriz.Cuantiles <-matrix(nr=len.beta, nc=13)
    Matriz.Varianzas <-matrix(nr=len.beta, nc=13)
    Matriz.Varianzas.Formulas<-matrix(nr=len.beta, nc=13)
    Matriz.Delta<-(Pij-as.matrix(Pi)%*%Pi)/Pij
    for (be in 1:len.beta)
    {
    beta<-vector.beta[be]
    t0.Rx<-vector.t0Rx[be]
    t0.Bx<-vector.t0Bx[be] #t0=Qg(beta)
    if (Jvar != 1) t0.Bx2<-vector.t0Bx2[be]
    t0.MC<-vector.t0MC[be] # t0=Qy(beta)
    cuantil.HK(TipoFd=2,beta,N,Pi,muestray) ->Matriz.Cuantiles[be,1]
    cuantil.r(beta, N, Pi, muestray, muestrax,x )->Matriz.Cuantiles[be,2]
    cuantil.dl(beta,N,Pi,muestray, muestrax,x) ->Matriz.Cuantiles[be,3]
    cuantil.CD(beta,N,etiquetas,muestray, muestrax,x) ->Matriz.Cuantiles[be,4]
    cuantil.RKM(beta,N,Pi, Pij, muestray, muestrax,x) ->Matriz.Cuantiles[be,5]
    cuantil.MA.RX.1pto(beta, N, t0.Rx, Pi, muestray, muestrax,x) ->resultado
    vector.Est<-c(vector.Est, resultado$Estado.NR)
                                resultado$t->Matriz.Cuantiles[be,6]
    cuantil.MA.BX.1pto(beta, N, t0.Bx, Pi, muestray, muestrax,x) ->resultado
    vector.Est<-c(vector.Est, resultado$Estado.NR)
                                resultado$t->Matriz.Cuantiles[be,7]

    if (Jvar != 1)
    {
    cuantil.MA.BX.1pto(beta, N, t0.Bx2, Pi, muestray, muestrax2,x2) ->resultado
    vector.Est<-c(vector.Est, resultado$Estado.NR)
                                resultado$t->tMatriz.Cuantiles[be,8]
    }
    else
                                Matriz.Cuantiles[be,7] ->Matriz.Cuantiles[be,8]
    cuantil.MC.1va.1pto(beta, N, t0.MC, Pi, muestray, muestrax,x) ->resultado
    vector.Est<-c(vector.Est, resultado$Estado.NR)
    resultado$t->Matriz.Cuantiles[be,9]
    if (Jvar != 1)
    {
    cuantil.MC.2va.1pto(beta, N, t0.MC, Pi, muestray, muestrax2,x2) ->resultado
    vector.Est<-c(vector.Est, resultado$Estado.NR)
                                resultado$t->Matriz.Cuantiles[be,10]
    }
    else
                                Matriz.Cuantiles[be,9] ->Matriz.Cuantiles[be,10]
    cuantil.MA.RX.3ptos(beta,N,valores.t,Pi,muestray,muestrax,x)
->Matriz.Cuantiles[be,11]
    cuantil.MA.BX.3ptos(beta, N, valores.t, Pi, muestray, muestrax,x)
->Matriz.Cuantiles[be,12]
    if (Jvar != 1)
    cuantil.MA.BX.3ptos(beta, N, valores.t, Pi, muestray, muestrax2,x2)
->Matriz.Cuantiles[be,13]

```

```

else
    Matriz.Quantiles[be,12] ->Matriz.Quantiles[be,13]
} # for (be in 1:len.beta)
## Estimacion de varianzas mediante Bootstrap
valor.k<-N/n
# Datos Poblacionales Bootstrap
y.B<- rep(muestray, valor.k)
x.B<- rep(muestrax, valor.k)
if (Jvar != 1)
{
vectorx2.B<- rep(muestrax2, valor.k)
x2.B<-matrix(vectorx2.B, nr=Jvar)
}
Sum.ArrayMatrizFd.i      <- matrix(0, nr=len.beta, nc=13)
Sum.Cuad.ArrayMatrizFd.i<- matrix(0, nr=len.beta, nc=13)
for (k in 1:valor.k)
{
## Obtenemos muestras Bootstrap
# Datos Muestrales Bootstrap
sample(N,n)->etiquetas.B
muestray.B<-y.B[etiquetas.B]
muestrax.B<-x.B[etiquetas.B]
if (Jvar != 1)
{
muestrax2.B<-x2.B[,etiquetas]
}
# Valores t0 (Bootstrap) para usar en estimadores M.Calibrados y M.A con 1pto.
## Para una variable
muestraxB.B<-as.matrix(muestrax.B)
muestraxB.B<-t(muestraxB.B)
datosx.B<-as.matrix(x.B)
datosx.B<-t(datosx.B)
### Rx
Rgorro.B<-matrix(nr=1,nc=1)
for (j in 1:1)      Rgorro.B[j,1]<-sum(muestray.B/Pi)/sum(muestraxB.B[j,]/Pi)
g.N.Rx.B<-t(Rgorro.B)%*%datosx.B
vector.t0Rx.B<-c()
for (p in 1:len.beta)  vector.t0Rx.B<-c(vector.t0Rx.B,  cuantil(1,
vector.beta[p], N, Pi=1,muestray=g.N.Rx.B) )
### Bx
betagorrox.B<-matrix(nr=1,nc=1)
betagorrox.B[1,1]<-sum((1/Pi)*muestraxB.B*muestray.B)/sum((1/Pi)*muestraxB.B^2)
g.N.Bx.B<-t(betagorrox.B)%*%datosx.B
vector.t0Bx.B<-c()
for (p in 1:len.beta)  vector.t0Bx.B<-c(vector.t0Bx.B,  cuantil(1,
vector.beta[p], N, Pi=1,muestray=g.N.Bx.B) )
## Si hay dos variables auxiliares
if (Jvar !=1)
{
### Bx2
xs<-t(muestrax2.B)*(1/Pi)
td<-muestrax2.B%*%xs
ts<-solve(td)
r<-muestrax2.B%*%((1/Pi)*muestray.B)
betagorrox2.B<-ts%*%(r)
g.N.Bx2.B<-t(betagorrox2.B)%*%x2.B
vector.t0Bx2.B<-c()
for (p in 1:len.beta)  vector.t0Bx2.B<-c(vector.t0Bx2.B,  cuantil(1,
vector.beta[p], N, Pi=1,muestray=g.N.Bx2.B) )
}
# Para los modelo-Calibrados se usa Qy(b).est: Se estima a partir de la muestra

```

```

vector.t0MC.B<-c()
for (p in 1:len.beta) vector.t0MC.B<-c(vector.t0MC.B, cuantil(1,
vector.beta[p], N, Pi=Pi, muestray=muestray.B) )
#### Fin de cálculo valores t0
for (be in 1:len.beta)
{
beta<-vector.beta[be]
t0.Rx.B <-vector.t0Rx.B[be]
t0.Bx.B <-vector.t0Bx.B[be]
if (Jvar != 1) t0.Bx2.B <-vector.t0Bx2.B[be]
t0.MC.B <-vector.t0MC.B[be]
## Calculo de estimadores
cuantil.HK(TipoFd=2,beta,N,Pi, muestray.B) ->valor1
Sum.ArrayMatrizFd.i[be,1]<- Sum.ArrayMatrizFd.i[be,1] + valor1
Sum.Cuad.ArrayMatrizFd.i[be,1] <- Sum.Cuad.ArrayMatrizFd.i[be,1] + valor1^2
cuantil.r (beta, N, Pi, muestray.B, muestrax.B, x.B ) ->valor2
Sum.ArrayMatrizFd.i[be,2]<- Sum.ArrayMatrizFd.i[be,2] + valor2
Sum.Cuad.ArrayMatrizFd.i[be,2] <- Sum.Cuad.ArrayMatrizFd.i[be,2] + valor2^2
cuantil.dl (beta,N,Pi, muestray.B, muestrax.B, x.B) ->valor3
Sum.ArrayMatrizFd.i[be,3]<- Sum.ArrayMatrizFd.i[be,3] + valor3
Sum.Cuad.ArrayMatrizFd.i[be,3] <- Sum.Cuad.ArrayMatrizFd.i[be,3] + valor3^2
cuantil.CD (beta,N,etiquetas.B,muestray.B, muestrax.B, x.B) ->valor4
Sum.ArrayMatrizFd.i[be,4]<- Sum.ArrayMatrizFd.i[be,4] + valor4
Sum.Cuad.ArrayMatrizFd.i[be,4] <- Sum.Cuad.ArrayMatrizFd.i[be,4] + valor4^2
cuantil.RKM(beta,N,Pi, Pij, muestray.B, muestrax.B, x.B) ->valor5
Sum.ArrayMatrizFd.i[be,5]<- Sum.ArrayMatrizFd.i[be,5] + valor5
Sum.Cuad.ArrayMatrizFd.i[be,5] <- Sum.Cuad.ArrayMatrizFd.i[be,5] + valor5^2
#### Se usan los nuevos valores t0:
cuantil.MA.RX.1pto(beta,N,t0.Rx.B, Pi, muestray.B, muestrax.B,x.B)->resultado
vector.Est<-c(vector.Est, resultado$Estado.NR)

resultado$t->valor6
Sum.ArrayMatrizFd.i[be,6]<- Sum.ArrayMatrizFd.i[be,6] + valor6
Sum.Cuad.ArrayMatrizFd.i[be,6] <- Sum.Cuad.ArrayMatrizFd.i[be,6] + valor6^2
cuantil.MA.BX.1pto(beta,N,t0.Bx.B, Pi, muestray.B, muestrax.B,x.B)-> resultado
vector.Est<-c(vector.Est, resultado$Estado.NR)
resultado$t->valor7
Sum.ArrayMatrizFd.i[be,7]<- Sum.ArrayMatrizFd.i[be,7] + valor7
Sum.Cuad.ArrayMatrizFd.i[be,7] <- Sum.Cuad.ArrayMatrizFd.i[be,7] + valor7^2
if (Jvar != 1)
{
cuantil.MA.BX.1pto(beta,N,t0.Bx2.B,Pi,muestray.B, muestrax2.B,x2.B)->resultado
vector.Est<-c(vector.Est, resultado$Estado.NR)
resultado$t->valor8
Sum.ArrayMatrizFd.i[be,8]<- Sum.ArrayMatrizFd.i[be,8] + valor8
Sum.Cuad.ArrayMatrizFd.i[be,8] <- Sum.Cuad.ArrayMatrizFd.i[be,8] + valor8^2
}
else
{
Sum.ArrayMatrizFd.i[be,8]<- Sum.ArrayMatrizFd.i[be,7]
Sum.Cuad.ArrayMatrizFd.i[be,8] <- Sum.Cuad.ArrayMatrizFd.i[be,7]
}
cuantil.MC.1va.1pto(beta,N,t0.MC.B,Pi,muestray.B,muestrax.B,x.B) ->resultado
vector.Est<-c(vector.Est, resultado$Estado.NR)
resultado$t->valor9
Sum.ArrayMatrizFd.i[be,9]<- Sum.ArrayMatrizFd.i[be,9] + valor9
Sum.Cuad.ArrayMatrizFd.i[be,9] <- Sum.Cuad.ArrayMatrizFd.i[be,9] + valor9^2
if (Jvar != 1)
{
cuantil.MC.2va.1pto(beta,N,t0.MC.B,Pi,muestray.B,muestrax2.B,x2.B)->resultado
vector.Est<-c(vector.Est, resultado$Estado.NR)
}
}

```

```

resultado$t->valor10
Sum.ArrayMatrizFd.i[be,10]<- Sum.ArrayMatrizFd.i[be,10] + valor10
Sum.Cuad.ArrayMatrizFd.i[be,10] <- Sum.Cuad.ArrayMatrizFd.i[be,10] + valor10^2
}
else
{
Sum.ArrayMatrizFd.i[be,10]<- Sum.ArrayMatrizFd.i[be,9]
Sum.Cuad.ArrayMatrizFd.i[be,10] <- Sum.Cuad.ArrayMatrizFd.i[be,9]
}
cuantil.MA.RX.3ptos(beta,N,valores.t, Pi, muestray.B, muestrax.B,x.B)->valor11
Sum.ArrayMatrizFd.i[be,11]<- Sum.ArrayMatrizFd.i[be,11] + valor11
Sum.Cuad.ArrayMatrizFd.i[be,11] <- Sum.Cuad.ArrayMatrizFd.i[be,11] +
valor11^2
cuantil.MA.BX.3ptos(beta,N,valores.t, Pi, muestray.B, muestrax.B,x.B)->valor12
Sum.ArrayMatrizFd.i[be,12]<- Sum.ArrayMatrizFd.i[be,12] + valor12
Sum.Cuad.ArrayMatrizFd.i[be,12] <- Sum.Cuad.ArrayMatrizFd.i[be,12] + valor12^2
if (Jvar != 1)
{
cuantil.MA.BX.3ptos(beta,N,valores.t,Pi,muestray.B, muestrax2.B,x2.B)->valor13
Sum.ArrayMatrizFd.i[be,13]<- Sum.ArrayMatrizFd.i[be,13] + valor13
Sum.Cuad.ArrayMatrizFd.i[be,13] <- Sum.Cuad.ArrayMatrizFd.i[be,13] + valor13^2
}
else
{
Sum.ArrayMatrizFd.i[be,13]<- Sum.ArrayMatrizFd.i[be,12]
Sum.Cuad.ArrayMatrizFd.i[be,13] <- Sum.Cuad.ArrayMatrizFd.i[be,12]
}
} # for (be in 1:len.beta)
} # end for (k in 1:valor.k)
## Estimadores de las varianzas mediante Bootstrap
Matriz.Varianzas <-
(1/(valor.k-1))*Sum.Cuad.ArrayMatrizFd.i -
(valor.k/(valor.k-1))* ( (1/valor.k)* Sum.ArrayMatrizFd.i )^2
# Estimacion de varianzas mediante formulas (falta multiplicar las densidades)
for (be in 1:len.beta)
{
### Estimador de tipo Hajek (=HT pq es MAS)
Delta(Matriz.Quantiles[be,1], muestray)->vectorDy.HK
cociente<-vectorDy.HK/Pi
Matriz.HK<-as.matrix(cociente)%*%cociente
Matriz.Varianzas.Formulas[be,1]<-(1/N^2)*sum(Matriz.Delta*Matriz.HK)
#### Estimador de tipo Razon
Delta(Matriz.Quantiles[be,2], muestray)->vectorDy.Razon
Delta(Matriz.Quantiles[be,2], Rgorro*muestrax)->vectorDRx.Razon
vector.a<-
vectorDy.Razon-(sum(vectorDy.Razon)/sum(vectorDRx.Razon))*vectorDRx.Razon
Matriz.Varianzas.Formulas[be,2]<-varianza.PS(N,vector.a, Pi, Pij)
#### Estimador de tipo Diferencia
Delta(Matriz.Quantiles[be,3], muestray)->vectorDy.Difel
Delta(Matriz.Quantiles[be,3], Rgorro*muestrax)->vectorDRx.Difel
vector.a<-vectorDy.Difel-vectorDRx.Difel
Matriz.Varianzas.Formulas[be,3]<-varianza.PS(N,vector.a, Pi, Pij)
### Estimador de CD.
Matriz.Varianzas.Formulas[be,4]<-Matriz.Varianzas[be,4]
#### Estimador de tipo RAO, KOVAR, MANTEL (dm)
qn1<-(muestray-Rgorro*muestrax)/sqrt(muestrax)
qn2<-(Matriz.Quantiles[be,5]-Rgorro*muestrax)/sqrt(muestrax)
Matriz.Gi<-matrix(nr=n,nc=n)
Matriz.Gj<-matrix(nr=n,nc=n)
for (i in 1:n)
for (j in 1:n)

```

```

{
  if (i == j)
  {
    vector.Pijk<- matrix( (n-1)/(N-1), nr=n, nc=1 )
    vector.Pijk<-as.vector(vector.Pijk)
    vector.Pijk[i]<-1
  }
  else
  {
    vector.Pijk<- matrix( (n-2)/(N-2), nr=n, nc=1 )
    vector.Pijk<- as.vector(vector.Pijk)
    vector.Pijk[i]<-1
    vector.Pijk[j]<-1
  }
  sum(Pij[i,j]/vector.Pijk)->suma.1
  sum( (Pij[i,j]/vector.Pijk)*Delta(qn2[i],qn1) )->suma.2
  Matriz.Gi[i,j]<-(1/suma.1)*suma.2
  sum( (Pij[i,j]/vector.Pijk)*Delta(qn2[j],qn1) )->suma.2
  Matriz.Gj[i,j]<-(1/suma.1)*suma.2
}
Delta(Matriz.Quantiles[be,5], muestray)->vectorDy.RKM
suma<-0
Indices<-1:n
for (j in 2:n)
{
  Ind.i<- Indices < j
  Ind.i <- Indices[Ind.i]
  for (i in Ind.i)
    suma<-suma+ ( (Pi[i]*Pi[j]-Pij[i,j])/Pij[i,j] ) *
      ( (vectorDy.RKM[i]-Matriz.Gi[i,j])/Pi[i]-
        (vectorDy.RKM[j]-Matriz.Gj[i,j])/Pi[j]
        )^2
  }
Matriz.Varianzas.Formulas[be,5]<-(1/N^2)*suma
### Estimadores Modelo Asistidos:
Matriz.Varianzas.Formulas[be,6]<- Matriz.Varianzas[be,6]
Matriz.Varianzas.Formulas[be,7]<- Matriz.Varianzas[be,7]
Matriz.Varianzas.Formulas[be,8]<- Matriz.Varianzas[be,8]
#### Estimadores de Calibracion
Datos.MC<-
F.distribucion.Mve.2t.vec(N,t=Matriz.Quantiles[be,9],t0=vector.t0MC[be],
Pi,muestray,muestrax, datosx)
vector.Est<-c(vector.Est, Datos.MC$Estado.NR)
Fd.MC<-Datos.MC$Fd
vector.g<-Datos.MC$vector.g
Delta(Matriz.Quantiles[be,9], muestray)->vectorDy.MC1
B.gorro<-sum( (vector.g-mean(vector.g) ) *vectorDy.MC1)/sum((vector.g-
mean(vector.g))^2)
vector.a<-vectorDy.MC1-Fd.MC- ( vector.g-mean(vector.g) ) *B.gorro
Matriz.Varianzas.Formulas[be,9]<-varianza.PS(N,vector.a, Pi, Pij)

  if (Jvar != 1)
  {
Datos.MC<-
F.distribucion.Mve.Jvar.vec(N,t=Matriz.Quantiles[be,10],t0=vector.t0MC[be],
Pi,muestray,muestrax2, x2)
vector.Est<-c(vector.Est, Datos.MC$Estado.NR)
Fd.MC<-Datos.MC$Fd
vector.g<-Datos.MC$vector.g
Delta(Matriz.Quantiles[be,10], muestray)->vectorDy.MC1

```

```

B.gorro<-sum( (vector.g-mean(vector.g) ) *vectorDy.MC1)/sum((vector.g-
mean(vector.g))^2)
vector.a<-vectorDy.MC1-Fd.MC- ( vector.g-mean(vector.g) )*B.gorro
Matriz.Varianzas.Formulas[be,10]<-varianza.PS(N,vector.a, Pi, Pij)
}
else
  Matriz.Varianzas.Formulas[be,9] ->Matriz.Varianzas.Formulas[be,10]
### Estimadores Modelo Asistidos:
Matriz.Varianzas.Formulas[be,11]<- Matriz.Varianzas[be,11]
Matriz.Varianzas.Formulas[be,12]<- Matriz.Varianzas[be,12]
Matriz.Varianzas.Formulas[be,13]<- Matriz.Varianzas[be,13]
} # for (be in 1:len.beta)
Estado<-min(vector.Est)
} # end while (Estado==0)
list(Matriz.Quantiles=Matriz.Quantiles, Matriz.Varianzas=Matriz.Varianzas,
muestray=muestray, Matriz.Varianzas.Formulas=Matriz.Varianzas.Formulas)
# EJEMPLO:
# Quantiles.MA.Bootstrap(n=40,y=Fam1500[,1],x=t(Fam1500[,2:3]),vector.beta=0.5)
# Quantiles.MA.Bootstrap(n=40,y=Murthy5[,1],x=Murthy5[,2], vector.beta=0.5)
}

simula.Quantiles.MA.Bootstrap
<-function(NombreP,B,y,x,n,vector.beta1, vector.beta2)
# Simulaciones de líneas de pobreza bajo distintos tipos de estimadores.
# "vector.beta1": Indices de cuantiles que se tomara para el estudio conjunto
# "vector.beta2": Cuantiles que no entran en el estudio conjunto
{
  vector.beta<-c(vector.beta1, vector.beta2)
  N<-length(y)
  k<-N%n
  if (k != 0) stop("Introduce un valor de 'n' que sea divisor de 'N' ")

  len.beta<-length(vector.beta)
# Cuantiles verdaderos de la variable principal y Densidades verdaderas.
  vector.Qy<-c()
  fy.ver<- c()
  for (p in 1:len.beta)
  {
    vector.Qy<-c(vector.Qy, cuantil(1, vector.beta[p], N, Pi=1, muestray=y) )
    fy.ver<- c(fy.ver, densidad(y, vector.Qy[p]) )
  }
# Indices de estimadores con varianzas:
  indices.for<-c(1,2,3,5,9,10)
#A) ##### Creando archivos (En XLS y en TXT)
##### Estimadores
  lineaEst<-paste("b", "Qy", "HK", "Razon", "Dif1", "CD",
"RKM", "MA.Rx1p", "MA.Bx1p1v", "MA.Bx1p2v", "MC1", "MC2", "MA.RX3p",
"MA.Bx3p1v", "MA.Bx3p2v", sep="\t")
  OutQuantilesTXT<-c()
  OutQuantilesXLS<-c()
  for (i in 1:len.beta)
  {
    OutQuantilesTXT<-c(OutQuantilesTXT,
paste("C:\\Pemle\\Quantiles\\est",NombreP,"QuantMAB",B,"n",n,"beta",vector.beta
[i],".txt", sep=""))
    OutQuantilesXLS<-c(OutQuantilesXLS,
paste("C:\\Pemle\\Quantiles\\est",NombreP,"QuantMAB",B,"n",n,"beta",vector.beta
[i],".xls", sep=""))
    write(lineaEst, file=OutQuantilesXLS[i], ncolumns=1, append=T)
  }
##### BOX-PLOT

```

```

OutQuantilesBoxTXT<-c()
OutQuantilesBoxXLS<-c()
OutVarJackBoxTXT<-c()
OutVarJackBoxXLS<-c()
OutVarVerBoxTXT<-c()
OutVarVerBoxXLS<-c()
OutVarEstBoxTXT<-c()
OutVarEstBoxXLS<-c()
  for (i in 1:len.beta)
  {
    ##Estimadores
    OutQuantilesBoxTXT<-c(OutQuantilesBoxTXT,
paste("C:\\Pemle\\Quantiles\\estBox",NombreP,"QuantMAB",B,"n",n,"beta",vector.b
eta[i],".txt", sep="" )
    OutQuantilesBoxXLS<-c(OutQuantilesBoxXLS,
paste("C:\\Pemle\\Quantiles\\estBox",NombreP,"QuantMAB",B,"n",n,"beta",vector.b
eta[i],".xls", sep="" )
    write(lineaEst, file=OutQuantilesBoxXLS[i], ncolumns=1, append=T)
    ## Varianzas Bootstrap
    OutVarJackBoxTXT<-c(OutVarJackBoxTXT,
paste("C:\\Pemle\\Quantiles\\Bootstrap\\estBox",NombreP,"QuantMAB",B,"n",n,"bet
a",vector.beta[i],".txt", sep="" )
    OutVarJackBoxXLS<-c(OutVarJackBoxXLS,
paste("C:\\Pemle\\Quantiles\\Bootstrap\\estBox",NombreP,"QuantMAB",B,"n",n,"bet
a",vector.beta[i],".xls", sep="" )
    write(lineaEst, file=OutVarJackBoxXLS[i], ncolumns=1, append=T)
    ## Densidades Verdaderas
    OutVarVerBoxTXT<-c(OutVarVerBoxTXT,
paste("C:\\Pemle\\Quantiles\\densidadVer\\estBox",NombreP,"QuantMAB",B,"n",n,"b
eta",vector.beta[i],".txt", sep="" )
    OutVarVerBoxXLS<-c(OutVarVerBoxXLS,

paste("C:\\Pemle\\Quantiles\\densidadVer\\estBox",NombreP,"QuantMAB",B,"n",n,"b
eta",vector.beta[i],".xls", sep="" )
    write(lineaEst, file=OutVarVerBoxXLS[i], ncolumns=1, append=T)
    ## Densidades Estimadores
    OutVarEstBoxTXT<-c(OutVarEstBoxTXT,
paste("C:\\Pemle\\Quantiles\\densidadEst\\estBox",NombreP,"QuantMAB",B,"n",n,"b
eta",vector.beta[i],".txt", sep="" )
    OutVarEstBoxXLS<-c(OutVarEstBoxXLS,
paste("C:\\Pemle\\Quantiles\\densidadEst\\estBox",NombreP,"QuantMAB",B,"n",n,"b
eta",vector.beta[i],".xls", sep="" )
    write(lineaEst, file=OutVarEstBoxXLS[i], ncolumns=1, append=T)
  }
### Errores estimadores
  ## Archivos en blanco
  OutErroresEstTXT<-c()
  OutErroresEstXLS<-c()
  ### Linea para identificar columnas
  lineal<-paste("n", "Qy", "Denominador",
"RE.HK", "RE.Razon", "RE.Dif1", "RE.CD",
"RE.RKM", "RE.MA.Rx1p", "RE.MA.BX1p1v", "RE.MA.Bx1p2v",
"RE.MC1", "RE.MC2", "RE.MA.RX3p", "RE.MA.Bx3p1v",
"RE.MA.Bx3p2v",
"Rba.HK", "Rba.Razon", "Rba.Dif1", "Rba.CD",
"Rba.RKM", "Rba.MA.Rx1p", "Rba.MA.BX1p1v", "Rba.MA.Bx1p2v",
"Rba.MC1", "Rba.MC2", "Rba.MA.RX3p", "Rba.MA.Bx3p1v",
"Rba.MA.Bx3p2v",
"Rbr.HK", "Rbr.Razon", "Rbr.Dif1", "Rbr.CD",
"Rbr.KM", "Rbr.MA.Rx1p", "Rbr.MA.BX1p1v", "Rbr.MA.Bx1p2v",

```

```

"RBr.MC1", "RBr.MC2", "RBr.MA.RX3p", "RBr.MA.Bx3p1v",
"RBr.MA.Bx3p2v",
"RRMSE.HK", "RRMSE.Razon", "RRMSE.Dif1", "RRMSE.CD",
"RRMSE.KM", "RRMSE.MA.Rx1p", "RRMSE.MA.BX1p1v",
"RRMSE.MA.Bx1p2v",
"RRMSE.MC1", "RRMSE.MC2", "RRMSE.MA.RX3p",
"RRMSE.MA.Bx3p1v", "RRMSE.MA.Bx3p2v",
sep="\t")
### Escribiendo la linea anterior.
for (i in 1:len.beta)
{
OutErroresEstTXT<-c(OutErroresEstTXT,
paste("C:\\Pemle\\Quantiles\\errEst",NombreP,"QuantMAB",B,"beta",vector.beta[i]
, ".txt", sep=""))
OutErroresEstXLS<-c(OutErroresEstXLS,
paste("C:\\Pemle\\Quantiles\\errEst",NombreP,"QuantMAB",B,"beta",vector.beta[i]
, ".xls", sep=""))
write(linea1, file=OutErroresEstXLS[i], ncolumns=1, append=T)
}
#### Errores Varianzas
## Archivos en blanco
OutErroresVarTXT<-c()
OutErroresVarXLS<-c()
OutErroresVarVerTXT<-c()
OutErroresVarVerXLS<-c()
OutErroresVarEstTXT<-c()
OutErroresVarEstXLS<-c()
linea2<-paste("n", "Qy", "Denominador",
"RE.HK", "RE.Razon", "RE.Dif1", "RE.CD",
"RE.RKM", "RE.MA.Rx1p", "RE.MA.BX1p1v", "RE.MA.Bx1p2v",
"RE.MC1", "RE.MC2", "RE.MA.RX3p", "RE.MA.Bx3p1v",
"RE.MA.Bx3p2v",
"RBa.HK", "RBa.Razon", "RBa.Dif1", "RBa.CD",
"RBa.RKM", "RBa.MA.Rx1p", "RBa.MA.BX1p1v", "RBa.MA.Bx1p2v",
"RBa.MC1", "RBa.MC2", "RBa.MA.RX3p", "RBa.MA.Bx3p1v",
"RBa.MA.Bx3p2v",
"RBr.HK", "RBr.Razon", "RBr.Dif1", "RBr.CD",
"RBr.KM", "RBr.MA.Rx1p", "RBr.MA.BX1p1v", "RBr.MA.Bx1p2v",
"RBr.MC1", "RBr.MC2", "RBr.MA.RX3p", "RBr.MA.Bx3p1v",
"RBr.MA.Bx3p2v",
"RRMSE.HK", "RRMSE.Razon", "RRMSE.Dif1", "RRMSE.CD",
"RRMSE.KM", "RRMSE.MA.Rx1p", "RRMSE.MA.BX1p1v",
"RRMSE.MA.Bx1p2v",
"RRMSE.MC1", "RRMSE.MC2", "RRMSE.MA.RX3p",
"RRMSE.MA.Bx3p1v", "RRMSE.MA.Bx3p2v",
"Cove.HK", "Cove.Razon", "Cove.Dif1", "Cove.CD",
"Cove.KM", "Cove.MA.Rx1p", "Cove.MA.BX1p1v", "Cove.MA.Bx1p2v",
"Cove.MC1", "Cove.MC2", "Cove.MA.RX3p", "Cove.MA.Bx3p1v",
"Cove.MA.Bx3p2v",
"Len.HK", "Len.Razon", "Len.Dif1", "Len.CD",
"Len.KM", "Len.MA.Rx1p", "Len.MA.BX1p1v", "Len.MA.Bx1p2v",
"Len.MC1", "Len.MC2", "Len.MA.RX3p", "Len.MA.Bx3p1v",
"Len.MA.Bx3p2v",
sep="\t")
## Escribiendo la linea anterior.
### Para Bootstrap
for (i in 1:len.beta)
{
OutErroresVarTXT<-c(OutErroresVarTXT,
paste("C:\\Pemle\\Quantiles\\Bootstrap\\errVar",NombreP,"QuantMAB",B,"beta",vec
tor.beta[i], ".txt", sep=""))
}

```

```

        OutErroresVarXLS<-c(OutErroresVarXLS,
paste("C:\\Pemle\\Quantiles\\Bootstrap\\errVar",NombreP,"QuantMAB",B,"beta",vector.beta[i],".xls", sep="" )
        write(linea2, file=OutErroresVarXLS[i], ncolumns=1, append=T)
    }
### Para Formulas con fy.ver
for (i in 1:len.beta)
{
    OutErroresVarVerTXT<-c(OutErroresVarVerTXT,
paste("C:\\Pemle\\Quantiles\\densidadVer\\errVar",NombreP,"QuantMAB",B,"beta",vector.beta[i],".txt", sep="" )
        OutErroresVarVerXLS<-c(OutErroresVarVerXLS,
paste("C:\\Pemle\\Quantiles\\densidadVer\\errVar",NombreP,"QuantMAB",B,"beta",vector.beta[i],".xls", sep="" )
        write(linea2, file=OutErroresVarVerXLS[i], ncolumns=1, append=T)
    }
### Para Formulas con fy.est
for (i in 1:len.beta)
{
    OutErroresVarEstTXT<-c(OutErroresVarEstTXT,
paste("C:\\Pemle\\Quantiles\\densidadEst\\errVar",NombreP,"QuantMAB",B,"beta",vector.beta[i],".txt", sep="" )
        OutErroresVarEstXLS<-c(OutErroresVarEstXLS,
paste("C:\\Pemle\\Quantiles\\densidadEst\\errVar",NombreP,"QuantMAB",B,"beta",vector.beta[i],".xls", sep="" )
        write(linea2, file=OutErroresVarEstXLS[i], ncolumns=1, append=T)
    }
#### Medidas AV.
##### Estimadores
    OutErroresAVestTXT<-c()
    OutErroresAVestXLS<-c()

lineal<-paste("n",
"RE.HK", "RE.Razon", "RE.Dif1", "RE.CD",
"RE.RKM", "RE.MA.Rx1p", "RE.MA.BX1p1v", "RE.MA.Bx1p2v",
"RE.MC1", "RE.MC2", "RE.MA.RX3p", "RE.MA.Bx3p1v",
"RE.MA.Bx3p2v",
"RBa.HK", "RBa.Razon", "RBa.Dif1", "RBa.CD",
"RBa.RKM", "RBa.MA.Rx1p", "RBa.MA.BX1p1v", "RBa.MA.Bx1p2v",
"RBa.MC1", "RBa.MC2", "RBa.MA.RX3p", "RBa.MA.Bx3p1v",
"RBa.MA.Bx3p2v",
"RBr.HK", "RBr.Razon", "RBr.Dif1", "RBr.CD",
"RBr.KM", "RBr.MA.Rx1p", "RBr.MA.BX1p1v", "RBr.MA.Bx1p2v",
"RBr.MC1", "RBr.MC2", "RBr.MA.RX3p", "RBr.MA.Bx3p1v",
"RBr.MA.Bx3p2v",
"RRMSE.HK", "RRMSE.Razon", "RRMSE.Dif1", "RRMSE.CD",
"RRMSE.KM", "RRMSE.MA.Rx1p", "RRMSE.MA.BX1p1v",
"RRMSE.MA.Bx1p2v",
"RRMSE.MC1", "RRMSE.MC2", "RRMSE.MA.RX3p",
"RRMSE.MA.Bx3p1v", "RRMSE.MA.Bx3p2v",
sep="\t")
    OutErroresAVestTXT<-
paste("C:\\Pemle\\Quantiles\\errAVest",NombreP,"QuantMAB",B, ".txt", sep="" )
    OutErroresAVestXLS<-
paste("C:\\Pemle\\Quantiles\\errAVest",NombreP,"QuantMAB",B, ".xls", sep="" )
    write(lineal, file=OutErroresAVestXLS, ncolumns=1, append=T)
##### Varianzas
# Bootstrap
    OutErroresAVvarTXT<-c()
    OutErroresAVvarXLS<-c()
# Formulas y fy.ver

```

```

OutErroresAVvarVerTXT<-c()
OutErroresAVvarVerXLS<-c()
# Formulas y fy.est
OutErroresAVvarEstTXT<-c()
OutErroresAVvarEstXLS<-c()
linea2<-paste("n",
"RE.HK", "RE.Razon", "RE.Dif1", "RE.CD",
"RE.RKM", "RE.MA.Rx1p", "RE.MA.BX1p1v", "RE.MA.Bx1p2v",
"RE.MC1", "RE.MC2", "RE.MA.RX3p", "RE.MA.Bx3p1v",
"RE.MA.Bx3p2v",
"Rba.HK", "Rba.Razon", "Rba.Dif1", "Rba.CD",
"Rba.RKM", "Rba.MA.Rx1p", "Rba.MA.BX1p1v", "Rba.MA.Bx1p2v",
"Rba.MC1", "Rba.MC2", "Rba.MA.RX3p", "Rba.MA.Bx3p1v",
"Rba.MA.Bx3p2v",
"RBr.HK", "RBr.Razon", "RBr.Dif1", "RBr.CD",
"RBr.KM", "RBr.MA.Rx1p", "RBr.MA.BX1p1v", "RBr.MA.Bx1p2v",
"RBr.MC1", "RBr.MC2", "RBr.MA.RX3p", "RBr.MA.Bx3p1v",
"RBr.MA.Bx3p2v",
"RRMSE.HK", "RRMSE.Razon", "RRMSE.Dif1", "RRMSE.CD",
"RRMSE.KM", "RRMSE.MA.Rx1p", "RRMSE.MA.BX1p1v",
"RRMSE.MA.Bx1p2v",
"RRMSE.MC1", "RRMSE.MC2", "RRMSE.MA.RX3p",
"RRMSE.MA.Bx3p1v", "RRMSE.MA.Bx3p2v",
"Cove.HK", "Cove.Razon", "Cove.Dif1", "Cove.CD",
"Cove.KM", "Cove.MA.Rx1p", "Cove.MA.BX1p1v", "Cove.MA.Bx1p2v",
"Cove.MC1", "Cove.MC2", "Cove.MA.RX3p", "Cove.MA.Bx3p1v",
"Cove.MA.Bx3p2v",
sep="\t")
OutErroresAVvarTXT<-
paste("C:\\Pemle\\Quantiles\\Bootstrap\\errAVvar",NombreP,"QuantMAB",B,".txt",
sep=" ")
OutErroresAVvarXLS<-
paste("C:\\Pemle\\Quantiles\\Bootstrap\\errAVvar",NombreP,"QuantMAB",B,".xls",
sep=" ")
write(linea2, file=OutErroresAVvarXLS, ncolumns=1, append=T)
OutErroresAVvarVerTXT<-
paste("C:\\Pemle\\Quantiles\\densidadVer\\errAVvar",NombreP,"QuantMAB",B,".txt",
, sep=" ")
OutErroresAVvarVerXLS<-
paste("C:\\Pemle\\Quantiles\\densidadVer\\errAVvar",NombreP,"QuantMAB",B,".xls",
, sep=" ")
write(linea2, file=OutErroresAVvarVerXLS, ncolumns=1, append=T)
OutErroresAVvarEstTXT<-
paste("C:\\Pemle\\Quantiles\\densidadEst\\errAVvar",NombreP,"QuantMAB",B,".txt",
, sep=" ")
OutErroresAVvarEstXLS<-
paste("C:\\Pemle\\Quantiles\\densidadEst\\errAVvar",NombreP,"QuantMAB",B,".xls",
, sep=" ")
write(linea2, file=OutErroresAVvarEstXLS, ncolumns=1, append=T)
#### Definicion de Matrices y Arrays
Array.Quantiles<- array(0, dim=c(len.beta, 13, B))
Sum.Quantiles <-matrix(0, nr=len.beta, nc=13)
Sum.Cuad.Quantiles<-matrix(0, nr=len.beta, nc=13)
Array.Varianzas<- array(0, dim=c(len.beta, 13, B))
Array.Var.For.fy.ver<- array(0, dim=c(len.beta, 13, B))
Array.Var.For.fy.est<- array(0, dim=c(len.beta, 13, B))
Array.Box.Quantiles<- array(0, dim=c(len.beta, 13, B))
Array.Box.Varianzas<- array(0, dim=c(len.beta, 13, B))
Array.Box.Var.For.ver <- array(0, dim=c(len.beta, 13, B))
Array.Box.Var.For.est <- array(0, dim=c(len.beta, 13, B))
M.ECM <-matrix(0, nr=len.beta, nc=13)

```

```

M.ses <-matrix(0, nr=len.beta, nc=13)
M.sesRe<-matrix(0, nr=len.beta, nc=13)
M.ECM.var <-matrix(0, nr=len.beta, nc=13)
M.ses.var <-matrix(0, nr=len.beta, nc=13)
M.sesRe.var <-matrix(0, nr=len.beta, nc=13)
M.ECM.var.ver <-matrix(0, nr=len.beta, nc=13)
M.ses.var.ver <-matrix(0, nr=len.beta, nc=13)
M.sesRe.var.ver <-matrix(0, nr=len.beta, nc=13)
M.ECM.var.est <-matrix(0, nr=len.beta, nc=13)
M.ses.var.est <-matrix(0, nr=len.beta, nc=13)
M.sesRe.var.est <-matrix(0, nr=len.beta, nc=13)
Coverage <-matrix(0, nr=len.beta, nc=13)
Coverage.ver <-matrix(0, nr=len.beta, nc=13)
Coverage.est <-matrix(0, nr=len.beta, nc=13)
len.Int.Boot <-matrix(0, nr=len.beta, nc=13)
len.Int.ver <-matrix(0, nr=len.beta, nc=13)
len.Int.est <-matrix(0, nr=len.beta, nc=13)
# Errores de estimadores y obtencion de array para varianzas y Box-PLot.
for (b in 1:B)
{
  Quantiles.MA.Bootstrap(n,y,x, vector.beta)->Datos
  Matriz.Quantiles<- Datos$Matriz.Quantiles
  Matriz.Varianzas<- Datos$Matriz.Varianzas
  muestray <- Datos$muestray
  Matriz.Varianzas.Formulas<-Datos$Matriz.Varianzas.Formulas
  fy.est<-matrix(nr=len.beta, nc=6)
  for(i in 1:len.beta) ### Calculo de medidas y densidad estimada
  {
    Array.Quantiles[i,,b]<- Matriz.Quantiles[i,]
    Sum.Quantiles[i,] <- Sum.Quantiles[i,]+Matriz.Quantiles[i,]
    Sum.Cuad.Quantiles[i,]<- Sum.Cuad.Quantiles[i,]+ Matriz.Quantiles[i,]^2
    Array.Varianzas[i,,b]<- Matriz.Varianzas[i,]
    Array.Var.For.fy.ver[i,,b]<-Matriz.Varianzas.Formulas[i,]
    Array.Var.For.fy.est[i,,b]<-Matriz.Varianzas.Formulas[i,]
    k<-0
    for (j in indices.for)
    {
      k<-k+1
      Array.Var.For.fy.ver[i,j,b]<-(1/fy.ver[i]^2)*Array.Var.For.fy.ver[i,j,b]
      fy.est[i,k]<-densidad(muestray,Matriz.Quantiles[i,j])
      Array.Var.For.fy.est[i,j,b]<-(1/fy.est[i,k]^2)*Array.Var.For.fy.est[i,j,b]
    }
    M.ECM[i,] <-M.ECM[i,] + (Matriz.Quantiles[i,]-vector.Qy[i])^2
    M.ses[i,] <-M.ses[i,] + abs(Matriz.Quantiles[i,]-vector.Qy[i])
    M.sesRe[i,]<-M.sesRe[i,]+ (Matriz.Quantiles[i,]-vector.Qy[i])
  Estado<-
  abs((Matriz.Quantiles[i,]-vector.Qy[i])/sqrt(Matriz.Varianzas[i,]) )<=1.96
  Coverage[i,]<-Coverage[i,]+ Estado
  Estado.ver<-abs( (Matriz.Quantiles[i,]-
vector.Qy[i])/sqrt(Array.Var.For.fy.ver[i,,b]) )<=1.96
  Coverage.ver[i,]<-Coverage.ver[i,]+ Estado.ver
  Estado.est<-abs( (Matriz.Quantiles[i,]-
vector.Qy[i])/sqrt(Array.Var.For.fy.est[i,,b]) )<=1.96
  Coverage.est[i,]<-Coverage.est[i,]+ Estado.est
  len.Int.Boot[i,]<- len.Int.Boot[i,] + 2*1.96*sqrt(Matriz.Varianzas[i,])
  len.Int.ver[i,] <- len.Int.ver[i,] + 2*1.96*sqrt(Array.Var.For.fy.ver[i,,b])
  len.Int.est[i,] <- len.Int.est[i,] + 2*1.96*sqrt(Array.Var.For.fy.est[i,,b])
  } # end for(i in 1:len.beta)
##### Preparando y escribiendo resultados
##### Estimadores
  for (i in 1:len.beta)

```

```

    {
      linea<-paste(b, vector.Qy[i], sep="\t")
      for (j in 1:13) linea<-paste(linea, Matriz.Quantiles[i,j],sep="\t")
      write(linea,file=OutQuantilesTXT[i], ncolumns=1, append=T)
      write(linea,file=OutQuantilesXLS[i], ncolumns=1, append=T)
    }
  } # for (b in 1:B)
# Varianzas reales (basado en la B replicaciones) y de los datos para BOX-PLOTS
Matriz.Varianzas.ver<-matrix(nr=len.beta, nc=13)
Matriz.Varianzas.ver<- (1/B)*Sum.Cuad.Quantiles - ( (1/B)* Sum.Quantiles )^2
  for (i in 1:len.beta)
    for (j in 1:13)
      {
Array.Box.Quantiles[i,j,<- (Array.Quantiles[i,j,]-vector.Qy[i])/vector.Qy[i]
  Array.Box.Varianzas[i,j,<- (Array.Varianzas[i,j,]-
Matriz.Varianzas.ver[i,j])/Matriz.Varianzas.ver[i,j]
  Array.Box.Var.For.ver[i,j,<- (Array.Var.For.fy.ver[i,j,]-
Matriz.Varianzas.ver[i,j])/Matriz.Varianzas.ver[i,j]
  Array.Box.Var.For.est[i,j,<- (Array.Var.For.fy.est[i,j,]-
Matriz.Varianzas.ver[i,j])/Matriz.Varianzas.ver[i,j]
      }
##### Escribiendo datos BOX-PLOTS
for (b in 1:B)
  {
    for (i in 1:len.beta)
      {
## Estimadores
        linea<-paste(b, vector.Qy[i], sep="\t")
        for (j in 1:13) linea<-paste(linea,
Array.Box.Quantiles[i,j,b],sep="\t")
          write(linea,file=OutQuantilesBoxTXT[i], ncolumns=1, append=T)
          write(linea,file=OutQuantilesBoxXLS[i], ncolumns=1, append=T)
### Bootstrap
        linea<-paste(b, vector.Qy[i], sep="\t")
        for (j in 1:13) linea<-paste(linea,
Array.Box.Varianzas[i,j,b],sep="\t")
          write(linea,file=OutVarJackBoxTXT[i], ncolumns=1, append=T)
          write(linea,file=OutVarJackBoxXLS[i], ncolumns=1, append=T)
### Formulas y fy.ver
        linea<-paste(b, vector.Qy[i], sep="\t")
        for (j in 1:13) linea<-paste(linea,
Array.Box.Var.For.ver[i,j,b],sep="\t")
          write(linea,file=OutVarVerBoxTXT[i], ncolumns=1, append=T)
          write(linea,file=OutVarVerBoxXLS[i], ncolumns=1, append=T)
### Formulas y fy.est
        linea<-paste(b, vector.Qy[i], sep="\t")
        for (j in 1:13) linea<-paste(linea,
Array.Box.Var.For.est[i,j,b],sep="\t")
          write(linea,file=OutVarEstBoxTXT[i], ncolumns=1, append=T)
          write(linea,file=OutVarEstBoxXLS[i], ncolumns=1, append=T)
      }
    } # for (b in 1:B)
##### Calculo de errores de varianzas
for (b in 1:B)
  {
    for(i in 1:len.beta)
      for (j in 1:13)
        {
          M.ECM.var[i,j] <-M.ECM.var[i,j] + (Array.Varianzas[i,j,b]-
Matriz.Varianzas.ver[i,j])^2
        }
      }
    }
  }

```

```

M.ses.var[i,j] <-M.ses.var[i,j] + abs(Array.Varianzas[i,j,b]-
Matriz.Varianzas.ver[i,j])
M.sesRe.var[i,j]<-M.sesRe.var[i,j] + (Array.Varianzas[i,j,b]-
Matriz.Varianzas.ver[i,j])
M.ECM.var.ver[i,j] <-M.ECM.var.ver[i,j] + (Array.Var.For.fy.ver[i,j,b]-
Matriz.Varianzas.ver[i,j])^2
M.ses.var.ver[i,j] <-M.ses.var.ver[i,j] + abs(Array.Var.For.fy.ver[i,j,b]-
Matriz.Varianzas.ver[i,j])
M.sesRe.var.ver[i,j]<-M.sesRe.var.ver[i,j] + (Array.Var.For.fy.ver[i,j,b]-
Matriz.Varianzas.ver[i,j])
M.ECM.var.est[i,j] <-M.ECM.var.est[i,j] + (Array.Var.For.fy.est[i,j,b]-
Matriz.Varianzas.ver[i,j])^2
M.ses.var.est[i,j] <-M.ses.var.est[i,j] + abs(Array.Var.For.fy.est[i,j,b]-
Matriz.Varianzas.ver[i,j])
M.sesRe.var.est[i,j]<-M.sesRe.var.est[i,j] + (Array.Var.For.fy.est[i,j,b]-
Matriz.Varianzas.ver[i,j])
}
} # for (b in 1:B)
#### Preparando y escribiendo resultados
##### Errores Estimadores
#### ECM:
RE<-matrix(nr=len.beta, nc=13)
Denominador<-M.ECM[,1]/B
for(i in 1:13) RE[,i]<-(M.ECM[,i]/B)/Denominador
##### Sesgos Absolutos
RBa<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBa[,i]<-M.ses[,i]/(B*vector.Qy)
##### Sesgos relativos
RBr<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBr[,i]<-100*M.sesRe[,i]/(B*vector.Qy) # En tanto por ciento.
##### RRMSE:
RRMSE<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RRMSE[,i]<-sqrt(M.ECM[,i]/B)/vector.Qy
for (i in 1:len.beta)
{
linea<-paste(n, vector.Qy[i], Denominador[i], sep="\t")
for (j in 1:13) linea<-paste(linea, RE[i,j],sep="\t")
for (j in 1:13) linea<-paste(linea, RBa[i,j],sep="\t")
for (j in 1:13) linea<-paste(linea, RBr[i,j],sep="\t")
for (j in 1:13) linea<-paste(linea, RRMSE[i,j],sep="\t")
write(linea,file=OutErroresEstTXT[i], ncolumns=1, append=T)
write(linea,file=OutErroresEstXLS[i], ncolumns=1, append=T)
}
##### Errores Varianzas y Coverage
#### ECM:
RE.var<-matrix(nr=len.beta, nc=13)
Denominador.var<-M.ECM.var[,1]/B
for(i in 1:13) RE.var[,i]<-(M.ECM.var[,i]/B)/Denominador.var
#### Sesgos Absolutos
RBa.var<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBa.var[,i]<-M.ses.var[,i]/(B*Matriz.Varianzas.ver[,i])
#### Sesgos relativos
RBr.var<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBr.var[,i]<-100*M.sesRe.var[,i]/(B*Matriz.Varianzas.ver[,i])
#### RRMSE:
RRMSE.var<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RRMSE.var[,i]<-sqrt(M.ECM.var[,i]/B)/Matriz.Varianzas.ver[,i]

Cove<-100*(1/B)*Coverage
len.Int.Boot<- (1/B)*len.Int.Boot
##### Errores Varianzas y Coverage (COOn Formulas y fy.ver)

```

```

#### ECM:
RE.var.ver<-matrix(nr=len.beta, nc=13)
Denominador.var.ver<-M.ECM.var.ver[,1]/B
for(i in 1:13) RE.var.ver[,i]<-(M.ECM.var.ver[,i]/B)/Denominador.var.ver
##### Sesgos Absolutos
RBa.var.ver<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBa.var.ver[,i]<-
M.ses.var.ver[,i]/(B*Matriz.Varianzas.ver[,i])
##### Sesgos relativos
RBr.var.ver<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBr.var.ver[,i]<-
100*M.sesRe.var.ver[,i]/(B*Matriz.Varianzas.ver[,i])
##### RRMSE:
RRMSE.var.ver<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RRMSE.var.ver[,i]<-
sqrt(M.ECM.var.ver[,i]/B)/Matriz.Varianzas.ver[,i]

Cove.ver<-100*(1/B)*Coverage.ver
len.Int.ver<- (1/B)*len.Int.ver
##### Errores Varianzas y Coverage (COon Formulas y fy.est)
#### ECM:
RE.var.est<-matrix(nr=len.beta, nc=13)
Denominador.var.est<-M.ECM.var.est[,1]/B
for(i in 1:13) RE.var.est[,i]<-(M.ECM.var.est[,i]/B)/Denominador.var.est
#### Sesgos Absolutos
RBa.var.est<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBa.var.est[,i]<-M.ses.var.est[,i]/(B*Matriz.Varianzas.ver[,i])
#### Sesgos relativos
RBr.var.est<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RBr.var.est[,i]<-
100*M.sesRe.var.est[,i]/(B*Matriz.Varianzas.ver[,i])
##### RRMSE:
RRMSE.var.est<-matrix(nr=len.beta, nc=13)
for(i in 1:13) RRMSE.var.est[,i]<-
sqrt(M.ECM.var.est[,i]/B)/Matriz.Varianzas.ver[,i]

Cove.est<-100*(1/B)*Coverage.est
len.Int.est<- (1/B)*len.Int.est
for (i in 1:len.beta)
{
  linea<-paste(n, vector.Qy[i], Denominador.var[i], sep="\t")
  for (j in 1:13) linea<-paste(linea, RE.var[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RBa.var[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RBr.var[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RRMSE.var[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, Cove[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, len.Int.Boot[i,j],sep="\t")
  write(linea,file=OutErroresVarTXT[i], ncolumns=1, append=T)
  write(linea,file=OutErroresVarXLS[i], ncolumns=1, append=T)
  linea<-paste(n, vector.Qy[i], Denominador.var.ver[i], sep="\t")
  for (j in 1:13) linea<-paste(linea, RE.var.ver[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RBa.var.ver[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RBr.var.ver[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RRMSE.var.ver[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, Cove.ver[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, len.Int.ver[i,j],sep="\t")
  write(linea,file=OutErroresVarVerTXT[i], ncolumns=1, append=T)
  write(linea,file=OutErroresVarVerXLS[i], ncolumns=1, append=T)
  linea<-paste(n, vector.Qy[i], Denominador.var.est[i], sep="\t")
  for (j in 1:13) linea<-paste(linea, RE.var.est[i,j],sep="\t")
  for (j in 1:13) linea<-paste(linea, RBa.var.est[i,j],sep="\t")
}

```

```

        for (j in 1:13) linea<-paste(linea, RBr.var.est[i,j],sep="\t")
        for (j in 1:13) linea<-paste(linea, RRMSE.var.est[i,j],sep="\t")
        for (j in 1:13) linea<-paste(linea, Cove.est[i,j],sep="\t")
        for (j in 1:13) linea<-paste(linea, len.Int.est[i,j],sep="\t")
        write(linea,file=OutErroresVarEstTXT[i], ncolumns=1, append=T)
        write(linea,file=OutErroresVarEstXLS[i], ncolumns=1, append=T)
    }
##### Estudio Conjunto: AVRBR, AVRRE, AVRRMSE, AVCoverage
length(vector.betal)->len.betal
# Paso 1: Restringirnos a los indices que componen este estudio
# Estimadores
RE.C <-RE[1:len.betal,]
RBa.C <-RBa[1:len.betal,]
RBr.C <-RBr[1:len.betal,]
RRMSE.C<-RRMSE[1:len.betal,]
# Varianzas Bootstrap
RE.var.C <-RE.var[1:len.betal,]
RBa.var.C <-RBa.var[1:len.betal,]
RBr.var.C <-RBr.var[1:len.betal,]
RRMSE.var.C<-RRMSE.var[1:len.betal,]

Cove.C <- Cove[1:len.betal,]
Len.C.Boot <- len.Int.Boot[1:len.betal,]
# Varianzas con Formulas y fy.ver
RE.var.C.ver <-RE.var.ver[1:len.betal,]
RBa.var.C.ver <-RBa.var.ver[1:len.betal,]
RBr.var.C.ver <-RBr.var.ver[1:len.betal,]
RRMSE.var.C.ver<-RRMSE.var.ver[1:len.betal,]

Cove.C.ver <- Cove.ver[1:len.betal,]
Len.C.ver <- len.Int.ver[1:len.betal,]

# Varanzas con Formulas y fy.est
RE.var.C.est <-RE.var.est[1:len.betal,]
RBa.var.C.est <-RBa.var.est[1:len.betal,]
RBr.var.C.est <-RBr.var.est[1:len.betal,]
RRMSE.var.C.est<-RRMSE.var.est[1:len.betal,]

Cove.C.est <- Cove.est[1:len.betal,]
Len.C.est <- len.Int.est[1:len.betal,]
# Paso 2: Creacion de los nuevos vectores
# Estimadores
AVRRE<- vector(len=13)
AVRBA<- vector(len=13)
AVRBR<- vector(len=13)
AVRRMSE<- vector(len=13)
# Varianzas Bootstrap
AVRRE.var<- vector(len=13)
AVRBA.var<- vector(len=13)
AVRBR.var<- vector(len=13)
AVRRMSE.var<- vector(len=13)

AVCoverage<- vector(len=13)
AVLen.Boot<- vector(len=13)
# Varianzas con formulas y fy.ver
AVRRE.var.ver<- vector(len=13)
AVRBA.var.ver<- vector(len=13)
AVRBR.var.ver<- vector(len=13)
AVRRMSE.var.ver<- vector(len=13)

AVCoverage.ver<- vector(len=13)

```

```

AVLen.ver<- vector(len=13)
# Varianzas con formulas y fy.est
AVRRE.var.est<- vector(len=13)
AVRBA.var.est<- vector(len=13)
AVRBr.var.est<- vector(len=13)
AVRRMSE.var.est<- vector(len=13)

AVCoverage.est<- vector(len=13)
AVLen.est<- vector(len=13)

# Paso 3: Rellenamos los vectores
for (j in 1:13)
{
  # Estimadores
  AVRRE[j]<-sqrt( mean(RE.C[,j]) )
  AVRBA[j]<-mean(RBa.C[,j])
  AVRBr[j]<-mean( abs(RBr.C[,j]) )
  AVRRMSE[j]<-mean(RRMSE.C[,j])
  # Varianzas Bootstrap
  AVRRE.var[j]<-sqrt( mean(RE.var.C[,j]) )
  AVRBA.var[j]<-mean(RBa.var.C[,j])
  AVRBr.var[j]<- mean( abs( RBr.var.C[,j] ) )
  AVRRMSE.var[j]<-mean(RRMSE.var.C[,j])

  AVCoverage[j]<- mean(Cove.C[,j])
  AVLen.Boot[j]<- mean(Len.C.Boot[,j])
  # Varianzas Formulas y fy.ver
  AVRRE.var.ver[j]<-sqrt( mean(RE.var.C.ver[,j]) )
  AVRBA.var.ver[j]<-mean(RBa.var.C.ver[,j])
  AVRBr.var.ver[j]<- mean( abs( RBr.var.C.ver[,j] ) )
  AVRRMSE.var.ver[j]<-mean(RRMSE.var.C.ver[,j])

  AVCoverage.ver[j]<- mean(Cove.C.ver[,j])
  AVLen.ver[j]<- mean(Len.C.ver[,j])
  # Varianzas Formulas y fy.est
  AVRRE.var.est[j]<-sqrt( mean(RE.var.C.est[,j]) )
  AVRBA.var.est[j]<-mean(RBa.var.C.est[,j])
  AVRBr.var.est[j]<- mean( abs( RBr.var.C.est[,j] ) )
  AVRRMSE.var.est[j]<-mean(RRMSE.var.C.est[,j])

  AVCoverage.est[j]<- mean(Cove.C.est[,j])
  AVLen.est[j]<- mean(Len.C.est[,j])
}
### Paso 4: Escribimos resultados
### De estimadores
linea<-paste(n, sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRE[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBA[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBr[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRMSE[j],sep="\t")
write(linea,file=OutErroresAVestTXT, ncolumns=1, append=T)
write(linea,file=OutErroresAVestXLS, ncolumns=1, append=T)
### De varianzas Bootstrap
linea<-paste(n, sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRE.var[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBA.var[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBr.var[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRMSE.var[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVCoverage[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVLen.Boot[j],sep="\t")
write(linea,file=OutErroresAVvarTXT, ncolumns=1, append=T)

```

```

write(linea,file=OutErroresAVvarXLS, ncolumns=1, append=T)
### De varianzas Formula y fy.ver
linea<-paste(n, sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRE.var.ver[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBa.var.ver[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBr.var.ver[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRMSE.var.ver[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVCoverage.ver[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVLen.ver[j],sep="\t")
write(linea,file=OutErroresAVvarVerTXT, ncolumns=1, append=T)
write(linea,file=OutErroresAVvarVerXLS, ncolumns=1, append=T)
### De varianzas Formula y fy.est
linea<-paste(n, sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRE.var.est[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBa.var.est[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRBr.var.est[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVRRMSE.var.est[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVCoverage.est[j],sep="\t")
for (j in 1:13) linea<-paste(linea, AVLen.est[j],sep="\t")
write(linea,file=OutErroresAVvarEstTXT, ncolumns=1, append=T)
write(linea,file=OutErroresAVvarEstXLS, ncolumns=1, append=T)
### EJEMPLOS
# simula.Quantiles.MA.Bootstrap(NombreP=c("Fam1500"), B=50, y=Fam1500[,1],
# x=Fam1500[,2], n=100, vector.beta1=seq(0.1, 0.9, by=0.2),
# vector.beta2=c(0.25,0.75) )
# simula.Quantiles.MA.Bootstrap(NombreP=c("Ecpf1997Trim1N9000"), B=500,
# y=Ecpf1997Trim1N9000[,1], x=Ecpf1997Trim1N9000[,2], n=100,
# vector.beta1=c(0.5), vector.beta2=c() )
}

```