

# Errors of measurement in scientometrics: Classification schemes and document types in citation and publication rankings

Nicolas Robinson-Garcia<sup>1</sup>, Benjamín Vargas-Quesada<sup>1</sup>, Daniel Torres-Salinas<sup>1</sup>, Zaida Chinchilla-Rodríguez<sup>2</sup> and Juan Gorraiz<sup>3</sup>

<sup>1</sup> *elrobin@ugr.es; benjamin@ugr.es; torressalinas@ugr.es*

Unit for Computational Humanities and Social Sciences (U-CHASS), department of Information and Communication, University of Granada, Granada (Spain)

<sup>2</sup> *zaida.chinchilla@csic.es*

Consejo Superior de Investigaciones Científicas (CSIC), Instituto de Políticas y Bienes Públicos (IPP), Madrid, Spain

<sup>3</sup> *juan.gorraiz@univie.ac.at*

University of Vienna, Vienna University Library, Dept of Bibliometrics, Boltzmanngasse 5, A-1090 Vienna (Austria)

## Abstract

This research article delves into methodological challenges in scientometrics, focusing on errors stemming from the selection of classification schemes and document types. Employing two case studies, we examine the impact of these methodological choices on publication and citation rankings of institutions. We compute seven bibliometric indicators for over 8,434 institutions using 23 different classification schemes derived from Clarivate's InCites suite, as well as including all document types versus only citable items. Given the critical role university rankings play in research management and their methodological controversies, our goal is to propose a methodology that incorporates uncertainty levels when reporting bibliometric performance in professional practice. We then delve into differences in error estimates within research fields as well as between institutions from different geographic regions. The findings underscore the importance of responsible metric use in research evaluation, providing valuable insights for both bibliometricians and consumers of such data.

**Keywords** Responsible metrics; institutions rankings; citation indicators; publication counts; classifications of science; professional bibliometrics

## Introduction

### *General Context*

Errors constitute an inherent and inevitable aspect of the scientific process. Achieving perfect accuracy is unattainable, as tools for measurement will always include some level of uncertainty (Scuro 2004). According to the Joint Committee for Guides in Metrology (BIPM et al. 2008), uncertainty is defined as a 'parameter associated with the result of a measurement, that characterizes the dispersion of the values that could reasonably be attributed to the measurand' (p. 2), where the measurand refers to the object being measured.

In the field of scientometrics, uncertainty and errors in measurement are usually overlooked. This is problematic for various reasons. First, neglecting uncertainty leads to the misuse of bibliometric indicators, which are then employed to legitimize decisions in conditions of limited trust or political controversy (Ràfols et al. 2016). The professionalization of bibliometrics in academic libraries (Gorraiz et al. 2020; Gumpenberger et al. 2012) and Higher Education planning and research administration (Cox et al. 2019) has elevated bibliometric reporting to a valuable resource for decision-making. The widespread of metrics either explicitly or implicitly in research assessment exercises to support and judge individuals, departments or institutions (e.g., Hammarfelt and Rushforth 2017; Moed 2008) has created a landscape of tools and commercial solutions designed to respond to such demand. Bibliometric suites such as InCites

1 (from Clarivate) or Scival (Elsevier) offer a battery of research indicators based respectively on  
2 bibliographic data from Web of Science and Scopus aimed at responding at this demand. These  
3 indicators are not only built based on different sets of publications, but their definition also  
4 differs leading sometimes to contradictory results for monitoring a common object (Robinson-  
5 Garcia et al. 2020).

6  
7 Second, bibliometric indicators are usually reported with high levels of precision, being the  
8 more evident example the inclusion of up to three decimals of the Journal Impact Factor. This  
9 is particularly worrying given the fact that even Garfield himself considered the impact factor  
10 accurate only up to one decimal place (Bensman 2007). In later years he admitted that the only  
11 reason ISI calculated the impact factors reported in the JCRs out to three decimal places was to  
12 avoid the large number of ties that would have resulted in listing many journals alphabetically  
13 in the impact factor rankings. Any analysis carried out under such a premise would de facto  
14 lose their validity, and we should not forget that this error has a profound effect particularly at  
15 the lower frequencies on ordinal rankings by the impact factor, on which most journal  
16 evaluations are based (Schloegl and Gorraiz 2010).

17  
18 This sense of false precision is also present in league tables and rankings in which variations in  
19 positions may be the result of noise rather than improvement or decay, affecting the prestige of  
20 those being portrayed in such tables (Bastedo and Bowman 2010; Gadd et al. 2021). While  
21 there have been some efforts to recognize this uncertainty, such as the inclusion of stability  
22 intervals in the Leiden Ranking and the introduction of position intervals in the Shanghai  
23 Ranking below a certain threshold, these measures are, at best, modest. Among the many causes  
24 for error or uncertainty in league tables we identify four types: 1) those derived from the  
25 misassignment of research outputs (Waltman et al. 2012), 2) those derived from unequal  
26 coverage of fields and locations in the database (Hicks 1999; Rafols et al. 2019; van Leeuwen  
27 et al. 2001), 3) those inherent to the metadata such as incompleteness or low quality metadata  
28 (Franceschini et al. 2015, 2016; Guerrero-Bote et al. 2021; Selivanova et al. 2019), and 4) those  
29 derived from methodological choices. By the latter we refer to choices which can affect the  
30 results without having *a priori* criteria set to justify such choices.

### 31 *Objectives of the study*

32 Our ultimate goal is to propose a methodological framework for evaluating measurements of  
33 error and variability when reporting bibliometric indicators in professional practice. This is  
34 particularly important if we aim to advocate and promote a responsible use of metrics in  
35 research evaluation in times when their use are more questioned than ever (Torres-Salinas et al.  
36 2023). Both, bibliometricians and producers of bibliometric data and indicators, have the  
37 responsibility of bridging towards professionals and scientists (Leydesdorff et al. 2016)  
38 consuming this data and promoting good practices on the use of such metrics.

39  
40 In this paper we showcase an specific case study in which errors can be quantified, derived  
41 from methodological choices. Specifically, we will focus on two case studies:

- 42  
43 • **Diverging classification schemes.** Here errors arise from the selection of a given  
44 classification scheme over another one when producing field normalized indicators  
45 (Ruiz-Castillo and Waltman 2015). In this study we will focus on this latter aspect of  
46 indicator variability by computing the same set of indicators on the same set of  
47 publications using up to 23 different classification schemes. Differences are due to how  
48 publications are categorized differently according to each classification. Also, they can  
49 be due to the inclusion or exclusion of some of the records due to their unfitnes

1 regarding the classification scheme (e.g., the Essential Science Indicators classification  
2 scheme does not consider the fields of Arts and Humanities).

- 3 • **Selection of document types.** The second case study is derived from the definition and  
4 selection of document types included in each analysis (Moed and Van Leeuwen 1995).  
5 Here we will focus on estimating errors derived from including either all document  
6 types in an analysis or only citable items.

7  
8 In both cases we will be using publication and citation rankings of institutions as a means to  
9 showcase how different indicators are affected by these errors and how it affects the positioning  
10 of institutions in these rankings. It is important to emphasize that the indicators used in our  
11 study, including those available through InCites, are part of a larger methodological framework  
12 focused on evaluating measurement errors and variability. This framework aims to improve the  
13 accuracy and reliability of bibliometric analyses, particularly when applied to institutional  
14 rankings, which use in in research management is specially controversial (Gadd 2020).

### 15 *Structure of the paper*

16  
17 The paper is structured as follows. Next, we review literature related to methodological  
18 challenges imposed by the selection of classification schemes and definition of document types.  
19 Then we define our methodological design for measuring errors in scientometrics. Here we  
20 build on a previous study (Robinson-Garcia et al. 2023) to compute error estimates of 7  
21 bibliometric indicators by using up to 23 different classification schemes. We use the InCites  
22 bibliometric suite to calculate these indicators at the institutional level, analyzing a total of  
23 8,433 institutions. Third, we report average relative error estimates when considering different  
24 classification schemes at an aggregated level. We analyze differences in errors when accounting  
25 for all document types versus when considering only citable documents (articles, reviews and  
26 letters). Furthermore, we look into differences in errors when focusing on specific research  
27 fields as well as on geographic regions. We conclude by reporting the implications our findings  
28 have for professional practice and the use of scientometric indicators for decision-making.

## 29 **Literature review**

### 30 *Selection of field classifications of science*

31 The selection of an appropriate classification scheme is a matter of concern in the field of  
32 scientometrics that has been discussed extensively in the literature using both quantitative and  
33 qualitative data (Gómez et al. 1996; Janssens et al. 2009; Minguillo 2010; Perianes-Rodriguez  
34 and Ruiz-Castillo 2018; Shu et al. 2019; Sugimoto and Weingart 2015). It is problematic for  
35 various reasons. First, in relation to the level of analysis at which the field delineation is made.  
36 Here, Gómez et al. (1996) point at four potential levels: document level, journal level, affiliation  
37 level and author level. Determining the level at which field delineation is done is crucial in  
38 terms of interpretability (Robinson-Garcia and Calero-Medina 2014) and accuracy (Shu et al.  
39 2019) as there is a problem of attribution especially when embedded in evaluation practices  
40 (Hansson et al. 2017). This issue is elegantly illustrated by Shu et al. (2019), who applied the  
41 Chinese Library Classification system to a set of publications using two different levels of  
42 analysis: at the journal level and the article level. They reported differences when rankings both:  
43 institutions and authors in terms of productivity, although this influence was mitigated at the  
44 institutional level.

45  
46 Second, there are many classification systems available databases (Gusenbauer 2022), each data  
47 sources has its own delineation of fields with *ad hoc* proposals (Gómez-Núñez et al. 2014;

1 Muñoz-Écija et al. 2019). But these classifications are not standardized, dealing to contradicting  
2 and non-comparable results. This is an important point in bibliometric studies in order to allow  
3 comparisons. The evaluation of the same content illustrates the structural differences in the  
4 databases, which is informative for interpreting bibliometric analyses (Stahlschmidt and  
5 Stephen 2022). A comparative study between the subject categories classification system in  
6 Web of Science and the All Science Journal Classification (ASJC) in Scopus, indicated that  
7 both classification were “too lenient in assigning journals to categories” (Wang and Waltman  
8 2016, p. 359). That is, they tended to include journals in multiple categories regardless how  
9 well connected they were in terms of citations with such categories. Another example is the  
10 classification proposed by Thijs et al. (2015), which combined hierarchical clustering and  
11 bibliographic coupling to create a 24 field classification system. This classification seemed to  
12 deviate greatly from the 22 fields from Essential Science Indicators included in Web of Science.  
13 These irregularities and discrepancies can lead to inconsistent outcomes (Reuven and Rosenfeld  
14 2022) hindering the interpretation of bibliometric indicators. As a means to improve these  
15 classification systems, some authors have suggested the use of hybrid methods, that is, refine  
16 journal-level classifications with paper level citation clustering and text mining to improve  
17 these classifications (e.g., Janssens et al. 2009).

18  
19 Third, there is the issue of balancing between accuracy, granularity and interpretability (Börner  
20 et al. 2012). In this sense, there are many multi-level classification systems which try to offer  
21 combinations by which users can zoom in or panning out. For instance, the publication level  
22 classification implemented by the CWTS (Waltman and van Eck 2012) is a three-tier  
23 classification system which goes spans from broad areas to micro topics based on direct citation  
24 networks between papers. Another example is the Australian and New Zealand Standard  
25 Research Classification (Australian Bureau of Statistics 2008) used by Dimensions, which also  
26 includes a three-level hierarchical classification system. The level of granularity of a  
27 classification will depend on the purpose for which the classification system is used. While  
28 broad fields may be desirable when reporting findings, more accuracy can lead to more robust  
29 indicators when normalizing by field (Ruiz-Castillo and Waltman 2015). In this sense, it is  
30 important to note that different indicators will show different levels of variability when moving  
31 from one classification system to another (Perianes-Rodriguez and Ruiz-Castillo 2018).

### 32 33 *Definition of document types*

34  
35 The definition and inclusion of document types will have a vital influence on the outcome of a  
36 bibliometric report. Furthermore, their definition and typology will vary depending on the  
37 database used as a data source. The underlying principle behind this distinction of documents  
38 is that different types of documents serve different functions in the scientific system and hence  
39 are read and cited differently, leading to differences in citation distributions (Lundberg 2007;  
40 Moed and Van Leeuwen 1995). Here we observe inconsistencies between databases, finding  
41 that a paper categorized as ‘Article’ in Web of Science (WoS) might be defined as a ‘Review’  
42 in Scopus, while Dimensions makes no distinction whatsoever (Visser et al. 2021).

43  
44 Scientometric studies have historically distinguished between citable and non-citable items  
45 (Heneberg 2014). But databases often misclassify records, being letters and reviews the most  
46 affected by these inaccuracies (Donner 2017). Gorraiz and Schloegl (2008) reported that there  
47 was a difference of over 10% between the sum of articles and reviews reported in Web of  
48 Science and Scopus. On a different study, Haunschild and Bornmann (2022) compared the  
49 scores that result from different normalization procedures, which have been performed based

1 on three different approaches of handling the document types. At least two of these approaches  
 2 are in use in popular university rankings. They showed that field-normalized scores strongly  
 3 depend on the choice of different document types and that the results on the aggregated level  
 4 (country, institution) are not supported by results on the level of individual publications.  
 5

## 6 **Data and methods**

### 7 *Data collection*

8 Data was retrieved using the bibliometric suite InCites, which includes the Web of Science  
 9 Core Collection including the Emerging Sources Citation Index (ESCI), for the period 1980-  
 10 2022. In this study we worked at the institutional level, using their “Organizations” option, that  
 11 is, their disambiguated list of institutions (more on issues for disambiguating institutional names  
 12 in Waltman et al. 2012). We include a total of 8,434 institutions which have published at least  
 13 1,000 documents (all document types) in the study period. InCites provides for each unit under  
 14 analysis a battery of bibliometric indicators which can be computed according to some  
 15 customizable parameters. One of them being the election of a given classification system. We  
 16 select 23 of the classification schemes currently available in InCites, including in here multi-  
 17 level classification systems. Table 1 provides a description of each of them, indicating the  
 18 number of categories per level and the assignment method Web of Science uses to create such  
 19 tables. Except of the Citation Topics classification, which follow the methodology developed  
 20 by Waltman and van Eck (2012), all schemes follow either a journal based method or aggregate  
 21 Web of Science subject categories (which also are journal based).  
 22

23 **Table 1. Description of the 23 classification schemes employed for the analysis\***

Acronym	Denomination	Description	Levels and categories	Method
ANVUR	ANVUR Category Schema	Official academic fields and disciplines list for Italian Universities Research and Teaching.	(1) 17 broad categories	Category-to-category mapping (WoS)
FOR	Australia ERA FOR	Revised Australian and New Zealand Standard Research Classification (ANZSRC 2020)	(1) 24 FoR2 (2) 212 FoR4	Journal mapping
CAPEB	CAPEB Brazil	Classification created by the Foundation CAPEB, linked to the Ministry of Education (Brazil)	(1) Capes 9 (2) Capes 49 (3) Capes 121	Category-to-category mapping (WoS)
CHINA	China SCADC Subject Categories	State Council Academic Degree Committee (SCADC) and Ministry of Education of China	(1) Broader 13 (2) Granular 96	Journal and other sources mapping
SHANGHAI	Shanghai Ranking Global Ranking of Subjects	Rankings of universities in 54 subjects across, Natural, , Life, Medical, and Social Sciences...	(1) 54 academic subjects	Category-to-category mapping (WoS)
TOPICS	Citation Topics	Algorithmically derived citation clusters (using an algorithm developed by CWTS, Leiden)	(1) Macro 10 (2) Meso 326	algorithmically on citation relationships
ESI	Essential Science Indicators Research Areas	All documents from Science Citation Index Expanded and Social Science Citation Index	(1) 22 broad categories	Journal mapping
FAPESP	FAPESP Brazil	Created by the São Paulo Research Foundation	(1) 9 High Level (2) 72 Detailed categories	Category-to-category mapping (WoS)
GIIP	Institutional Profiles Research Areas	Clarivate Analytics has been profiling the world’s leading universities and research institutions	(1) 6 broad academic fields	Category-to-category mapping (WoS)
KAKEN	KAKEN Category Schema (10 and 66)	From Japan called the Kakenhi Program (Grants-in-Aid for Scientific Research).	(1) 10 L2 (2) 66 L3	Category-to-category mapping (WoS)
OECD	OECD Category Schema	Revised Field of Science and Technology (FOS) Classification of the Frascati Manual.	(1) 42 fields	Category-to-category mapping (WoS)
PL19	PL19 Category Schema	The Polish PL19 category schema is used for annual evaluation exercise	(1) 44 underlying categories	Journal mapping
RIS	Research and Innovation Strategies for Specialization	The Research and Innovation Strategies for Smart Specialization (RIS3) for Latvia	(1) 7 specialization fields	Category-to-category mapping (WoS)
UKREF14	UK RAE Units of Assessment 2018	UK 2014 Research Assessment Exercise (RAE) Units of Assessment (UoA)	(1) 36 categories	Category-to-category mapping (WoS)

<b>UKREF21</b>	<b>UK REF Units of Assessment 2021</b>	UK 2021 Research Assessment Exercise (RAE) Units of Assessment (UoA)	(1) 36 units of assessment	Category-to-category mapping (WoS)
<b>WOS</b>	<b>Web of Science Research Areas</b>	The Web of Science schema comprises approximately 250 subject areas	(1) 250	---

\* Information extracted from <https://incites.help.clarivate.com/Content/Research-Areas/research-areas.htm> The data presented reproduce higher levels of categorization. Specifically, certain categories (FOR, CAPES, CHINA, TOPICS, FAPESP, KAKEN) encompass multiple levels, which are not detailed in this table. The total of 16 listed categories, when including the sublevels of the mentioned categories, sums up to the 23 classification schemes analysed in this paper.

1  
2 For each institution we focused on seven different indicators: total number of publications,  
3 times cited, the Category Normalized Citation Indicator, top 1% most cited papers, top 10%  
4 most cited papers, average percentile, and H-Index. Each indicator is computed 46 times;  
5 according to the 23 different classification schemes and including all document types versus  
6 only citable items.

### 7 *Calculation of errors*

8 The calculation of the errors followed the standard definition used in the experimental sciences,  
9 where the absolute error of a measurement is the difference between the measured value and  
10 the true value. In cases where the true value is unknown, it is replaced with the mean value  
11 obtained after multiple iterations of the measurement. In our case, each indicator is calculated  
12 a total of 23 times for each classification scheme and twice when looking into document types.  
13

14 To accurately determine the absolute error for each institution, we first calculate the mean value  
15 of the repeated measurements of each indicator. The absolute error for a given institution,  $\Delta x_i$   
16 is then computed as the difference between its measured value and the reference mean value.  
17 Mathematically, the absolute error of an institution can be defined as:

$$18 \Delta x_{\text{mean}} = (\text{ABS}(\Delta x_1) + \text{ABS}(\Delta x_2) + \dots + \text{ABS}(\Delta x_n))/n$$

19  
20  
21 Where  $x_i$  corresponds to an institution and  $\text{ABS}(\Delta x_i)$  is its absolute error.  
22

23 Let's consider we want to compute the absolute error after retrieving the number of publications  
24 produced by three departments from three different bibliometric databases. The measurements  
25 obtained for from each database are shown in Table 2.  
26

27 **Table 2. Mock up example of number of publications obtained for three institutions from three**  
28 **different databases**

	<b>Database I</b>	<b>Database II</b>	<b>Database III</b>
<b>Department A</b>	250	260	240
<b>Department B</b>	300	310	290
<b>Department C</b>	280	270	290

29  
30 For each research department, we first calculate the mean number of publications across the  
31 three databases:  
32

- 33 • Department A: 250
- 34 • Department B: 300
- 35 • Department C: 280

36  
37 Next, we calculate the absolute error for each database by comparing the measured values with  
38 the mean value for each department and dividing by the number of measurements. Hence, the  
39 absolute error for Department A would be:

$$\Delta x_{\text{mean}} = \frac{|250 - 250| + |260 - 250| + |240 - 250|}{3} = 6.67$$

Following the calculation of absolute errors, we also compute relative and percentage errors, facilitating comparisons between indicators across different classification schemes. The relative error is obtained by dividing the absolute error by the mean value, while the percentage error is calculated by multiplying the relative error by 100. Following the example of Department A, we now show its relative and percentage errors:

$$\text{Relative Error} = \frac{6.67}{250} \approx 0.0267$$

$$\text{Percentage Error} = 0.0267 \times 100 \approx 2.67\%$$

To account for the variability observed in repeated measurements, we calculated the confidence intervals associated with each error measurement. These intervals were derived from the standard deviation of the repeated measurements and provide a range within which the true value is expected to lie with a specified level of confidence.

## Results

### *General overview of error estimates by classification scheme and document types*

In table 3 we include an overview of our final dataset. For each classification scheme we include the total number of institutions covered when considering all document types and when considering only citable documents. As observed, only the WOS classification includes the 8,434 institutions originally included in our dataset when considering all document types. Simply by filtering to citable documents, we lose up to 559 institutions in the best of cases (WOS). The classification scheme with the lowest coverage is FOR1 which includes 42% of the institutions and 18% of documents, a share that increases up to 20% when considering only citable documents.

**Table 3. Total institutions and publications per classification scheme and document types**

Classification scheme	Institutions		Records	
	All docs.	Citable only	All docs.	Citable only
WOS	8,434	7,875	63,908,002	51,399,082
KAKENL2	8,433	7,874	63,895,799	51,389,373
KAKENL3	8,432	7,874	63,895,799	51,389,373
UKREF21	8,432	7,874	63,895,799	51,389,373
RIS	8,432	7,874	63,895,799	51,389,373
OCDE	8,432	7,874	63,895,799	51,389,373
GIPP	8,432	7,874	63,895,799	51,389,373
FAPESP	8,432	7,874	63,895,799	51,389,373
CAPES9	8,432	7,874	63,895,799	51,389,373
CAPES49	8,432	7,874	63,895,799	51,389,373
ANVUR	8,431	7,874	63,888,303	51,389,373
UKREF14	8,414	7,865	63,671,768	51,241,279
CAPES121	8,380	7,837	63,011,376	50,799,943

<b>SHANGHAI</b>	8,231	7,737	59,830,551	48,881,578
<b>CHINA BROAD</b>	8,141	7,610	58,977,093	47,077,808
<b>CHINA NARROW</b>	8,141	7,610	58,977,087	47,077,802
<b>TOPICSMACRO</b>	7,931	7,717	53,728,425	49,126,966
<b>TOPICSMESO</b>	7,931	7,717	53,728,425	49,126,966
<b>TOPICSMICRO</b>	7,931	7,717	53,728,425	49,126,966
<b>PL19</b>	7,746	7,187	51,317,825	40,102,542
<b>ESI</b>	7,497	6,943	49,844,239	39,158,270
<b>FOR2</b>	6,560	5,959	36,711,808	27,927,846
<b>FOR1</b>	3,549	3,348	11,671,816	10,107,716

1 Next, we report the average estimated error for each classification resulting from considering  
2 either all document types or only citable documents (Table 4). As observed, the lowest relative  
3 error is reported for the H-Index and the total number of citations, in both case below 1% for  
4 all classification schemes. We observe error estimates ranging between 0.7 for the CNCI to  
5 5.3% for the average percentile indicators. On the other extreme we find that number of  
6 documents is the indicator with the largest estimated error (7.4% on average). Furthermore, we  
7 observe different levels of variability in the error by indicator. While in the cases of times cited  
8 and H-Index these are below 0.1%, in the case of top 1% highly cited papers, we observe a  
9 variability of 8.3 between the largest estimated error (11% for TOPICSMICRO) and the lowest  
10 value (2.7% for TOPICSMACRO). A similar case we observe again in the number of  
11 publications, where there is a variability of up to 6.1% between FOR2 (9.6%) and the three  
12 TOPICS schemes (3.5%).

13

14

**Table 4. Average relative error by classification scheme considering document type**

Classification scheme	Docs	+/-	Times Cited	+/-	CNCI	+/-	Top 1%	+/-	Top 10%	+/-	Avg. percentile	+/-	H-Index	+/-
WOS	8.0	6.3	0.8	0.6	2.2	2.0	4.5	3.7	2.7	2.3	4.3	3.6	0.3	0.3
KAKENL2	8.0	6.3	0.8	0.6	3.0	3.0	4.9	4.1	2.7	2.3	4.2	3.6	0.3	0.3
KAKENL3	8.0	6.3	0.8	0.6	2.8	2.7	4.5	3.8	2.7	2.3	4.2	3.6	0.3	0.3
UKREF21	8.0	6.3	0.8	0.6	2.7	2.6	4.7	3.9	2.9	2.4	4.3	3.6	0.3	0.3
RIS	8.0	6.3	0.8	0.6	3.3	3.2	5.0	4.3	2.7	2.2	4.2	3.6	0.3	0.3
OCDE	8.0	6.3	0.8	0.6	2.8	2.7	4.8	4.0	2.9	2.4	4.3	3.6	0.3	0.3
GIPP	8.0	6.3	0.8	0.6	3.0	2.9	5.0	4.2	2.9	2.3	4.4	3.7	0.3	0.3
FAPESP	8.0	6.3	0.8	0.6	2.9	2.8	4.8	3.9	2.8	2.4	4.3	3.6	0.3	0.3
CAPES9	8.0	6.3	0.8	0.6	3.0	2.9	5.1	4.3	2.9	2.4	4.3	3.6	0.3	0.3
CAPES49	8.0	6.3	0.8	0.6	2.7	2.6	4.6	3.9	2.7	2.3	4.3	3.6	0.3	0.3
ANWUR	8.0	6.3	0.8	0.6	3.0	2.9	4.8	4.0	2.9	2.5	4.3	3.6	0.3	0.3
UKREF14	8.0	6.3	0.8	0.6	3.0	2.6	4.6	3.9	2.9	2.4	4.3	3.6	0.3	0.3
CAPES121	7.9	6.3	0.8	0.6	3.0	2.1	4.5	3.7	2.7	2.3	4.3	3.6	0.3	0.3
SHANGHAI	7.7	6.0	0.8	0.6	3.0	2.5	4.5	3.8	2.7	2.2	3.9	3.3	0.3	0.3
CHINA BROAD	8.2	6.3	0.8	0.6	3.0	2.9	4.9	4.2	2.9	2.4	4.4	3.7	0.3	0.3
CHINA NARROW	8.2	6.3	0.8	0.6	2.7	2.6	4.6	3.7	3.1	2.6	4.5	3.7	0.3	0.3
TOPICSMACRO	3.5	2.6	0.8	0.6	1.0	0.9	2.7	2.4	1.3	1.0	0.7	0.5	0.3	0.3
TOPICSMESO	3.5	2.6	0.8	0.6	0.8	0.7	4.3	3.5	1.2	1.0	0.7	0.5	0.3	0.3
TOPICSMICRO	3.5	2.6	0.8	0.6	0.7	0.6	11.0	6.2	1.6	1.1	0.6	0.5	0.3	0.3
PL19	8.6	6.6	0.8	0.6	3.4	3.4	5.1	4.3	2.8	2.3	4.6	3.8	0.3	0.3



ESI	8.7	6.5	0.8	0.6	3.0	2.8	4.8	4.0	3.0	2.4	4.5	3.7	0.3	0.3
FOR2	9.6	7.0	0.8	0.6	2.8	2.7	4.8	3.9	2.8	2.2	5.3	4.0	0.3	0.3
FOR1	4.7	3.4	0.8	0.6	1.6	1.4	3.5	2.8	1.7	1.3	2.0	1.7	0.3	0.3

1  
2 Given differences on the institutional coverage observed in Table 3, we present error estimates  
3 only considering institutions which are present in all schemes. To show how the exclusion of  
4 schemes with lower coverage affect error estimates, we present the world average relative error  
5 when considering the 23 classification schemes and considering only the top 13 with the largest  
6 institutional coverage (Table 5). Here we observe relatively small differences between focusing  
7 on all document types or citable documents. When considering the 23 classification schemes,  
8 we observe that the largest estimated errors are found for number of documents (12.2-12.8%)  
9 and top 1% most highly cited papers (10.0-11.2%). When reducing the number of  
10 classifications, we find how the error estimates are reduced drastically for size dependent  
11 indicators (docs, times cited and H-Index) while they are very similar for non-size dependent  
12 indicators (CNCI, Top 1%, Top 10% and avg. percentile). However, while in the case of the  
13 CNCI and the average percentile, the error is reduced for 13 classification schemes, it increases  
14 for the Top 1% and Top 10% indicators.

15 **Table 5. World Percentage errors considering 23 and 13 classification schemes by document**  
16 **type selection**

17

Indicators	23 schemes		13 schemes	
	All docs	Citable only	All docs	Citable only
<b>Docs</b>	12.8	12.2	0.2	0.2
+/-	1.9	1.9	0.2	0.2
<b>Times Cited</b>	9.2	9.2	0.2	0.2
+/-	0.4	0.4	0.2	0.2
<b>CNCI</b>	4.8	4.6	3.6	3.4
+/-	1.7	1.4	1.8	1.5
<b>Top 1%</b>	11.2	10.0	10.0	11.8
+/-	4.1	3.9	4.8	6.6
<b>Top 10%</b>	5.8	5.8	5.3	6.0
+/-	1.8	1.7	1.8	2.2
<b>Avg. Percentile</b>	3.1	2.4	2.2	2.3
+/-	1.1	0.7	0.7	0.7
<b>H-Index</b>	4.1	4.1	0.1	0.1
+/-	0.4	0.4	0.1	0.1

18 *Field differences in error estimates for document types: Macro Topic and ESI fields*

19 As a means to deepen on how the choice of including all document types or only those defined  
20 as citable, in tables 6 and 7 we look into differences in error estimates by field. To do so, we  
21 use TOPICSMACRO scheme formed by 10 major fields and the 22 ESI fields respectively.

22

23 In the case of the macro topics, we observe that the largest errors can be found in the fields of  
24 Arts & Humanities (between 0.6% for H-Index and 8.1% for number of documents), followed  
25 by Clinical & Life Sciences (between 0.4% for the H-Index and 6.0% for number of  
26 documents). Interestingly, the greatest variability in percentage error is found in the number of  
27 documents, for which Arts & Humanities is the field with the largest average error estimate,

1 while Electrical Engineering, Electronics & Computer Science has an average percentage error  
 2 of 1%. The rest of the patterns between indicators hold to what we observed in Table 4.

3  
 4 By using the citation topics, we considered that their calculation is only possible for  
 5 publications with cited references. Therefore, all the document types usually non-including  
 6 cited references (like e.g., meeting abstracts) are automatically excluded. That is the reason why  
 7 the differences between using all document types and only citable items will be lower than  
 8 expected when considering other classifications elaborated on journal level.

9  
 10 **Table 6. Percentage errors from document type selection by Macro Topics**

TOPICSMACRO	Docs	+/-	Times Cited	+/-	CNCI	+/-	Top 1%	+/-	Top 10%	+/-	Avg. percentile	+/-	H-Index	+/-
Agriculture, Environment & Ecology	2.3	1.1	1.0	0.5	0.7	0.5	2.1	1.7	0.9	0.7	0.4	0.3	0.297	0.3
Arts & Humanities	8.1	2.2	3.0	0.8	1.7	1.0	5.2	3.8	2.5	1.2	2.8	1.0	0.6	0.7
Chemistry	1.9	1.1	1.0	0.6	0.5	0.4	1.8	1.7	0.7	0.6	0.3	0.2	0.2	0.3
Clinical & Life Sciences	6.0	2.8	1.4	0.6	1.6	1.3	3.5	2.7	1.7	1.2	1.0	0.6	0.4	0.3
Earth Sciences	2.3	1.0	0.7	0.3	0.6	0.4	1.8	1.4	0.7	0.4	0.4	0.2	0.2	0.2
Electrical Engineering, Electronics & Computer Science	1.0	0.5	0.6	0.4	0.4	0.3	1.2	1.2	0.4	0.3	0.1	0.1	0.2	0.3
Engineering & Materials Science	1.0	0.6	0.4	0.3	0.4	0.3	1.6	1.5	0.6	0.5	0.2	0.1	0.1	0.2
Mathematics	1.3	0.7	0.7	0.5	0.5	0.4	1.6	1.4	0.5	0.4	0.2	0.2	0.3	0.4
Physics	1.4	0.6	0.5	0.3	0.5	0.3	1.3	1.1	0.5	0.4	0.2	0.2	0.2	0.2
Social Sciences	3.4	1.6	1.3	0.6	1.0	0.8	2.4	1.9	0.9	0.6	0.7	0.4	0.4	0.4

11  
 12 In the case of the ESI fields (Table 7) a different pattern is observed. Here it is the  
 13 Multidisciplinary category the one accounting for the largest errors in all indicators with notable  
 14 differences with respect to the rest of the categories. The exception is found in the average  
 15 percentile with other fields exhibit greater percentage errors (e.g., Clinical Medicine, Social  
 16 Sciences, general). The other exhibiting a large, estimated error is Clinical Medicine, where the  
 17 error in terms of number of documents is just above 20%.

18 **Table 7. Percentage errors from document type selection by ESI fields**

ESI	Docs	+/-	Times Cited	+/-	CNCI	+/-	Top 1%	+/-	Top 10%	+/-	Avg. percentile	+/-	H-Index	+/-
Agricultural Sciences	3.5	2.7	0.4	0.3	1.8	1.5	3.0	2.3	1.3	1.0	1.8	1.7	0.1	0.2
Biology & Biochemistry	11.5	5.5	1.3	0.610	3.0	2.2	6.5	4.7	3.9	2.9	6.9	3.8	0.4	0.3
Chemistry	5.4	4.6	0.7	0.6	3.2	3.1	5.5	5.4	3.3	3.6	3.8	4.0	0.2	0.3
Clinical Medicine	20.5	5.5	1.8	0.6	6.5	4.6	8.0	4.7	6.5	3.5	10.6	3.2	0.4	0.3
Computer Science	3.1	1.2	0.8	0.5	1.4	1.1	3.0	2.2	1.1	0.8	0.8	0.4	0.3	0.3
Economics & Business	6.6	2.9	1.2	0.5	3.0	2.3	4.8	3.7	1.8	1.1	3.5	2.0	0.4	0.3
Engineering	1.9	1.1	0.5	0.3	1.7	1.5	2.9	2.6	1.3	1.0	0.5	0.4	0.1	0.2
Environment/Ecology	2.2	0.9	1.1	0.5	1.1	0.7	2.1	1.5	1.0	0.7	0.3	0.2	0.5	0.4
Geosciences	3.5	1.4	0.7	0.3	1.2	0.9	3.2	2.8	1.1	0.9	1.1	0.7	0.2	0.3
Immunology	13.7	3.1	1.8	0.5	3.1	2.1	5.5	3.4	4.6	2.3	6.9	2.2	0.5	0.3
Materials Science	1.1	0.6	0.3	0.2	1.1	1.0	2.3	2.4	0.8	0.6	0.3	0.2	0.1	0.1
Mathematics	1.3	0.6	0.3	0.2	0.6	0.5	1.9	1.6	0.6	0.4	0.5	0.3	0.1	0.2
Microbiology	4.1	1.3	1.5	0.7	1.1	0.8	3.1	2.3	1.3	0.8	1.2	0.7	0.4	0.3

<b>Molecular Biology &amp; Genetics</b>	9.1	2.8	0.9	0.3	2.8	1.7	5.0	2.4	3.9	1.8	5.3	1.9	0.3	0.2
<b>Multidisciplinary</b>	23.0	22.6	3.1	1.9	12.6	7.0	19.9	15.1	13.1	9.4	8.8	9.8	1.1	0.8
<b>Neuroscience &amp; Behavior</b>	14.6	3.8	1.2	0.4	3.4	2.5	5.4	3.3	3.6	3.0	8.4	2.4	0.2	0.2
<b>Pharmacology &amp; Toxicology</b>	11.8	4.5	1.6	0.7	2.8	2.0	5.2	3.5	4.1	2.4	6.4	2.8	0.4	0.4
<b>Physics</b>	1.6	0.7	0.7	0.4	0.6	0.4	1.5	1.1	0.6	0.4	0.3	0.2	0.2	0.2
<b>Plant &amp; Animal Sciences</b>	5.2	2.8	1.0	0.4	2.0	1.4	4.2	3.4	2.0	1.6	2.2	1.8	0.2	0.3
<b>Psychiatry/Psychology</b>	12.5	2.8	1.3	0.5	3.3	2.4	4.8	3.2	3.4	1.7	7.5	1.9	0.3	0.3
<b>Social Sciences, general</b>	15.1	5.6	1.5	0.5	4.1	3.0	5.8	3.9	4.1	2.4	9.8	4.6	0.4	0.3
<b>Space Science</b>	1.3	0.4	0.4	0.2	0.5	0.4	1.2	1.0	0.5	0.4	0.3	0.1	0.1	0.1

1  
2 *Regional differences in error estimates for classification schemes: United States vs. South*  
3 *America*

4 Finally, we delve into regional differences in order to understand how homogeneous or  
5 heterogeneous is the effect of using different classification schemes in institutions located in  
6 different parts of the world. As an illustrative example, in Table 8 we report the average  
7 percentage errors of institutions located in the United States and in South America. Again, we  
8 report errors considering all document types or citable documents, as well as including the 23  
9 classification schemes or only the 13 classifications with the largest institutional coverage.  
10 While the general pattern of errors is similar to that observed in table 5, we do observe  
11 differences between institutions of these two regions. Overall, we observe that differences of  
12 error are always lower than 1% with some exceptions. The largest difference is that observed  
13 for the Top 1% most highly cited publications, where differences of error are above 4%  
14 (favoring US institutions) when considering the 23 classification schemes. These differences  
15 are below 1% when considering only 13 schemes. The other exception is for all document types  
16 and Top 10%, where the difference of error is just above 1%, being larger for South American  
17 institutions.  
18

19 **Table 8. Percentage errors from classification schemes for institutions located in the United**  
20 **States vs. institutions located in South America**

Indicators	UNITED STATES				SOUTH AMERICA			
	23 schemes		13 schemes		23 schemes		13 schemes	
	All docs	Citable only	All docs	Citable only	All docs	Citable only	All docs	Citable only
<b>Docs</b>	13.8	11.8	0.2	0.3	12.4	12.6	0.1	0.1
+/-	1.6	1.5	0.2	0.2	1.0	1.3	0.1	0.1
<b>Times Cited</b>	9.3	9.3	0.3	0.3	9.3	9.2	0.2	0.2
+/-	0.3	0.3	0.2	0.2	0.3	0.3	0.1	0.1
<b>CNCI</b>	5.4	4.9	4.5	3.6	5.1	5.0	2.0	2.0
+/-	1.9	1.4	2.3	1.6	1.2	1.4	0.8	0.6
<b>Top 1%</b>	8.6	7.6	9.2	8.8	13.2	12.2	9.3	9.7
+/-	2.7	2.4	4.3	4.1	3.6	3.4	3.8	4.1
<b>Top 10%</b>	5.0	5.2	5.1	5.2	6.0	6.0	5.0	5.3
+/-	1.2	1.3	1.9	1.8	1.2	1.3	1.4	1.5
<b>Avg. Percentile</b>	4.1	2.0	1.9	2.0	3.2	2.9	2.4	2.5
+/-	1.1	0.5	0.5	0.5	0.8	0.6	0.5	0.5

H-Index	4.3	4.3	0.1	0.1	4.0	4.0	0.1	0.1
+/-	0.3	0.3	0.1	0.1	0.3	0.3	0.1	0.1

1

## 2 **Conclusions & Discussion**

3 In bibliometric practices as well as in many bibliometric studies, several choices always have  
4 to be made which naturally can seriously affect the results obtained and their interpretations.  
5 In this article, we have assessed the magnitude of these errors occurring in two recurring  
6 situations: 1) considering the document type and 2) considering different classification  
7 schemes. Moreover, we think that a science that is mainly based on statistics should indicate  
8 the validity of the retrieved values, both of its numbers and of its decimals, in order to reveal  
9 their significance. The purpose of this paper is to reveal bibliometricians which errors will be  
10 derived from their document type and classification schemes decisions.

11

12 In the first case, we have seen how the choice of what type of documents should be included in  
13 the analysis, can influence not only the indicators of publication activity, where they are  
14 manifest, but also those of impact and especially those of normalized impact, such as the CNCI,  
15 and the Top 1% and Top 10%, which are the most commonly used in bibliometric practices  
16 (Moed 2017).

17

18 Generally, in this case, the decision falls between choosing the “citable publications” or all  
19 types of documents that are available. This difference was already introduced by Garfield, when  
20 he introduced his measure of the Journal Impact Factor by considering only the citable items  
21 (articles, reviews and proceedings) in the denominator, while in the numerator including the  
22 citations to all document types.

23

24 Our study shows that this decision can severely distort the results in the Essential Science  
25 Indicators (ESI) Category “Multidisciplinary”. As it is well known, this category appears as  
26 well in the ESI classification scheme as in the Journal Citation Reports, while in InCites all  
27 these publications are reallocated to categories that are more precise. Anyhow, our results show  
28 the big differences bibliometricians may confront with dealing with journals and journal impact  
29 measures assigned to this category.

30

31 Other categories are also affected by the document type decision, especially *Clinical Medicine*,  
32 *Immunology* and *Pharmacology & Toxicology*, where percentage errors higher than 4% are  
33 reported for the calculation of the top 10% most cited, that is commonly used as a measure for  
34 academic excellence. These are mainly due to the effect of the *Meeting Abstracts*, as it has  
35 already been often reported (Gorraiz et al. 2016). This also corroborated by the results obtained  
36 when using the Citation Topics – Macrotopics. In this case, the percentage errors are  
37 considerably diminished, because Meeting Abstracts lack of references and cannot be  
38 considered in this scheme. *Editorial Materials* and *Book Chapters* are the other document types  
39 responsible for the differences in the impact measures, also in other categories like *Social*  
40 *Sciences, general, Biology & Biochemistry* and *Psychiatry/Psychology*. Interestingly, *Physics*,  
41 *Space Physics* and *Mathematics* are almost not affected, and the only consideration of citable  
42 items is a sound decision.

43

44 Therefore, in our analyses we are considering errors due to three different decision-making  
45 processes: 1) to select a classification scheme, 2) to select a classification schema based on  
46 journal level versus one based on document level; and 3) to use a classification schema on  
47 different aggregation levels. And for the bibliometrician community it is crucial to be able to

1 estimate the effect that these decisions are expected to have on the values calculated for their  
2 impact measures.

3  
4 In this our first error analysis study, we have considered all the classification schemes available  
5 nowadays in InCites and calculated the errors due to their selection. In a first instance, we  
6 wanted to extend our work to the classifications available in SciVal and compare our results.  
7 But to our great surprise, we discovered that the selection of the classification did not change  
8 the impact results at the working aggregation level (meso-level) in that analytical tool, so they  
9 could not be incorporated into our study.

10  
11 Our results show the alterations that the impact measures would undergo in the case of using  
12 another classification scheme. Of the four normalized impact indicators, the top 1% shows the  
13 largest margin of error ranging from 3 to 17%. This is mainly due to the short number of  
14 publications contributing to this percentile. The measure of excellence (Top 10% most cited)  
15 fluctuates between 5 and 8%, and the average percentile varying between approximately 1.5  
16 and 3%. The value of the CNCI can vary from 1.5 to 5%.

17  
18 Differences in the values of the H-index are only reported when including classification  
19 schemes reducing considerably the number of publications considered (10 from 23 Schemes).  
20 Furthermore, topological factors, would only increase slightly these errors or deviations due to  
21 the selection of a classification scheme. In a study case, the differences between North-  
22 American and South-American institutions, the discrepancies between the percentage errors  
23 were only slightly higher for the percentile indicators for the South-American organizations  
24 when considering only the citable items. Finally, our results show that the decision made by for  
25 using all types of documents or only the citable ones hardly alter the divergences resulting from  
26 the use of one or the other classification system (see table 5) at the meso level.

## 27 28 **Limitations and further research**

29  
30 All the analyses in this study have been carried out at the meso-level, which is the most  
31 significant for this purpose. All analyses have been performed for the period 1980-2022 to  
32 increase the significance of the results. For shorter periods, the errors, especially for indicators  
33 that are not standardized, such as the number of publications and the number of citations, would  
34 naturally be much higher and should therefore be recalculated for each specific situation in  
35 further studies

36  
37 At the macro level very similar results are expected, while at the micro level, i.e., in the  
38 evaluation of individuals, the analysis is much more problematic due to the small number of  
39 publications available, and the great diversity of the cases and criteria to be considered (such as  
40 gender, age of career, etc.). For example, Åström, Hammarfelt and Hansson (2017) discuss how  
41 scientific publications can be categorized in different fields depending on the unit of assessment  
42 being evaluated: the publication, the individual or the institution. They found variations in terms  
43 of purpose of categorization as well as purpose of evaluation, i.e., the definition and function  
44 of the publications depending on whether it is situated in a context of scholarly communication  
45 or a context of research evaluation. The raising questions such as on what levels the distinctions  
46 are made, and in terms of on what principles the categories are being defined. The varying  
47 functions of the boundary object becomes critical when contextualized within the concept of  
48 infrastructures (publication databases, citation indices, evaluation systems and classification  
49 systems). Therefore, it is always advisable to perform the measurements on a case-by-case basis  
50 and with different data sources and purpose of use.

1  
2 One of the most subtle and critical problems of bibliometrics is classifications. As is well  
3 known, there is no standard, and each database and even each nation or continent uses its own  
4 schema for their evaluation systems. That is why Clarivate, in a big effort, tried to collect the  
5 most used ones on different aggregation levels (macro, meso and micro) in its flagship product  
6 “InCites” and to use them for the calculation of the most common bibliometric indicators (see  
7 Table 1). Besides, classification systems are usually created at the level of journals (Pudovkin  
8 & Garfield, 2002) but also at the paper level (Waltman & van Eck 2012; Rivest et al 2021).  
9 Comparisons of these two levels of aggregations, journal classification versus paper  
10 classification using the same classification scheme and the same dataset revealed that almost  
11 half of the papers could be misclassified in journal classification systems (Shu et al 2019). When  
12 comparing rankings of the most productive institutions and authors, classification of papers has  
13 less influence on rankings at the institutional level than at the individual level (Shu et al 2020),  
14 which has implications for bibliometric evaluation. At this point, it is important to emphasize  
15 that InCites also includes other classifications based not only on journal level (e.g. ESI and  
16 WoS subject categorization) but also most recent ones based on document level like the  
17 classification system builds on “Citation Topics”, algorithmically derived citation clusters using  
18 an algorithm developed by CWTS, Leiden<sup>1</sup>. Further studies can bring light on differences using  
19 journal or article level when calculating errors.

20  
21 Another limitation of this study is that we have performed all the analyses on a single data  
22 source, the Web of Science Core Collection. Therefore, it will be necessary to perform future  
23 analyses comparing the results in different sources, such as WoS CC, Scopus, Dimensions and  
24 even other Open sources, such as OpenAlex, Crossref or Lens. We are fully convinced that  
25 these studies will be of great help to all those involved in providing bibliometric services to be  
26 able to argue, justify and foresee the effects of the decisions they have had to make in carrying  
27 out their analyses.

28  
29 **Conflict of interest:** Nicolas Robinson-Garcia is Associate Editor of *Scientometrics*.

30 **Acknowledgments:** The authors acknowledge funding from the Spanish Ministry of Science (COMPARE project  
31 PID2020-117007RA-I00; and RESPONSIBLE project PID2021-128429NB-I00). Nicolas Robinson-Garcia is  
32 funded by a Ramón y Cajal grant from the Spanish Ministry of Science and Innovation (REF: RYC2019-027886-  
33 I).

## 34 **References**

- 35 Åström, F., Hammarfelt, B., & Hansson, J. (2017). Scientific publications as boundary objects:  
36 theorising the intersection of classification and research evaluation. *Information Research*,  
37 22(1), CoLIS paper 1623. <http://InformationR.net/ir/22-1/colis/colis1623.html>  
38 Australian Bureau of Statistics. (2008, March 31). Australian and New Zealand Standard  
39 Research Classification (ANZSRC). c=AU; o=Commonwealth of Australia; ou=Australian  
40 Bureau of Statistics.  
41 <https://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage/1297.02008?OpenDocument>.  
42 Accessed 30 November 2023  
43 Bastedo, M. N., & Bowman, N. A. (2010). U.S. News & World Report College Rankings:  
44 Modeling Institutional Effects on Organizational Reputation. *American Journal of Education*,  
45 116(2), 163–183. <https://doi.org/10.1086/649436>

---

<sup>1</sup> Traag, V.A., Waltman, L. & van Eck, N.J. From Louvain to Leiden: guaranteeing well-connected communities. *Sci Rep* 9, 5233 (2019). <https://doi.org/10.1038/s41598-019-41695-z>

- 1 Bensman, S. J. (2007). Garfield and the impact factor. *Annual Review of Information Science*  
2 *and Technology*, 41(1), 93–155. <https://doi.org/10.1002/aris.2007.1440410110>
- 3 BIPM, IEC, IFCC, ILAC, ISO, IUPAC, et al. (2008). Evaluation of measurement data —  
4 Supplement 2 to the “Guide to the expression of uncertainty in measurement” — Extension  
5 to any number of output quantities.  
6 [https://www.bipm.org/documents/20126/2071204/JCGM\\_102\\_2011\\_E.pdf/6a3281aa-1397-](https://www.bipm.org/documents/20126/2071204/JCGM_102_2011_E.pdf/6a3281aa-1397-d703-d7a1-a8d58c9bf2a5)  
7 [d703-d7a1-a8d58c9bf2a5](https://www.bipm.org/documents/20126/2071204/JCGM_102_2011_E.pdf/6a3281aa-1397-d703-d7a1-a8d58c9bf2a5)
- 8 Börner, K., Klavans, R., Patek, M., Zoss, A. M., Biberstine, J. R., Light, R. P., et al. (2012).  
9 Design and Update of a Classification System: The UCSD Map of Science. *PLOS ONE*, 7(7),  
10 e39464. <https://doi.org/10.1371/journal.pone.0039464>
- 11 Cox, A., Gadd, E., Petersohn, S., & Sbaiffi, L. (2019). Competencies for bibliometrics. *Journal*  
12 *of Librarianship and Information Science*, 51(3), 746–762.  
13 <https://doi.org/10.1177/0961000617728111>
- 14 Donner, P. (2017). Document type assignment accuracy in the journal citation index data of  
15 Web of Science. *Scientometrics*, 113(1), 219–236. <https://doi.org/10.1007/s11192-017-2483->  
16 [y](https://doi.org/10.1007/s11192-017-2483-y)
- 17 Franceschini, F., Maisano, D., & Mastrogiacomo, L. (2015). Errors in DOI indexing by  
18 bibliometric databases. *Scientometrics*, 102(3), 2181–2186. <https://doi.org/10.1007/s11192->  
19 [014-1503-4](https://doi.org/10.1007/s11192-014-1503-4)
- 20 Franceschini, F., Maisano, D., & Mastrogiacomo, L. (2016). Empirical analysis and  
21 classification of database errors in Scopus and Web of Science. *Journal of Informetrics*, 10(4),  
22 933–953. <https://doi.org/10.1016/j.joi.2016.07.003>
- 23 Gadd, E. (2020). University rankings need a rethink. *Nature*, 587(7835), 523–523.  
24 <https://doi.org/10.1038/d41586-020-03312-2>
- 25 Gadd, E., Holmes, R., & Shearer, J. (2021). Developing a Method for Evaluating Global  
26 University Rankings, 3(1), 2. <https://doi.org/10.29024/sar.31>
- 27 Gómez, I., Bordons, M., Fernández, M. T., & Méndez, A. (1996). Coping with the problem of  
28 subject classification diversity. *Scientometrics*, 35(2), 223–235.  
29 <https://doi.org/10.1007/BF02018480>
- 30 Gómez-Núñez, A. J., Batagelj, V., Vargas-Quesada, B., Moya-Anegón, F., & Chinchilla-  
31 Rodríguez, Z. (2014). Optimizing SCImago Journal & Country Rank classification by  
32 community detection. *Journal of Informetrics*, 8(2), 369–383.  
33 <https://doi.org/10.1016/j.joi.2014.01.011>
- 34 Gorraiz, J., Melero-Fuentes, D., Gumpenberger, C., & Valderrama-Zurián, J.-C. (2016).  
35 Availability of digital object identifiers (DOIs) in Web of Science and Scopus. *Journal of*  
36 *Informetrics*, 10(1), 98–109. <https://doi.org/10/f8d2bp>
- 37 Gorraiz, J., & Schloegl, C. (2008). A bibliometric analysis of pharmacology and pharmacy  
38 journals: Scopus versus Web of Science. *Journal of Information Science*, 34(5),  
39 715–725. <https://doi.org/10.1177/0165551507086991>
- 40 Gorraiz, J., Wieland, M., Ulrych, U., & Gumpenberger, C. (2020). De Profundis: A Decade of  
41 Bibliometric Services Under Scrutiny. In C. Daraio & W. Glänzel (Eds.), *Evaluative*  
42 *Informetrics: The Art of Metrics-Based Research Assessment : Festschrift in Honour of Henk*  
43 *F. Moed* (pp. 233–260). Cham: Springer International Publishing.  
44 [https://doi.org/10.1007/978-3-030-47665-6\\_11](https://doi.org/10.1007/978-3-030-47665-6_11)
- 45 Guerrero-Bote, V. P., Chinchilla-Rodríguez, Z., Mendoza, A., & de Moya-Anegón, F. (2021).  
46 Comparative Analysis of the Bibliographic Data Sources Dimensions and Scopus: An  
47 Approach at the Country and Institutional Levels. *Frontiers in Research Metrics and*  
48 *Analytics*, 5. <https://www.frontiersin.org/articles/10.3389/frma.2020.593494>. Accessed 5  
49 December 2023

- 1 Gumpenberger, C., Wieland, M., & Gorraiz, J. (2012). Bibliometric practices and activities at  
2 the University of Vienna. *Library Management*, 33(3), 174–183.  
3 <https://doi.org/10.1108/01435121211217199>
- 4 Gusenbauer, M. (2022). Search where you will find most: Comparing the disciplinary coverage  
5 of 56 bibliographic databases. *Scientometrics*, 127(5), 2683–2745.  
6 <https://doi.org/10.1007/s11192-022-04289-7>
- 7 Hammarfelt, B., & Rushforth, A. D. (2017). Indicators as judgment devices: An empirical study  
8 of citizen bibliometrics in research evaluation. *Research Evaluation*, 26(3), 169–180.  
9 <https://doi.org/10.1093/reseval/rvx018>
- 10 Haunschild, R., & Bornmann, L. (2022). Relevance of document types in the scores' calculation  
11 of a specific field-normalized indicator: Are the scores strongly dependent on or nearly  
12 independent of the document type handling? *Scientometrics*, 127(8), 4419–4438.  
13 <https://doi.org/10.1007/s11192-022-04446-y>
- 14 Heneberg, P. (2014). Parallel worlds of citable documents and others: Inflated commissioned  
15 opinion articles enhance scientometric indicators. *Journal of the Association for Information  
16 Science and Technology*, 65(3), 635–643. <https://doi.org/10.1002/asi.22997>
- 17 Hicks, D. (1999). The difficulty of achieving full coverage of international social science  
18 literature and the bibliometric consequences. *Scientometrics*, 44(2), 193–215.
- 19 Janssens, F., Zhang, L., Moor, B. D., & Glänzel, W. (2009). Hybrid clustering for validation  
20 and improvement of subject-classification schemes. *Information Processing & Management*,  
21 45(6), 683–702. <https://doi.org/10.1016/j.ipm.2009.06.003>
- 22 Leydesdorff, L., Wouters, P., & Bornmann, L. (2016). Professional and citizen bibliometrics:  
23 complementarities and ambivalences in the development and use of indicators—a state-of-  
24 the-art report. *Scientometrics*, 109(3), 2129–2150. [https://doi.org/10.1007/s11192-016-2150-](https://doi.org/10.1007/s11192-016-2150-8)  
25 8
- 26 Lundberg, J. (2007). Lifting the crown—citation z-score. *Journal of Informetrics*, 1(2), 145–  
27 154. <https://doi.org/10.1016/j.joi.2006.09.007>
- 28 Minguillo, D. (2010). Toward a new way of mapping scientific fields: Authors' competence for  
29 publishing in scholarly journals. *Journal of the American Society for Information Science and  
30 Technology*, 61(4), 772–786. <https://doi.org/10.1002/asi.21282>
- 31 Moed, H. F., & Van Leeuwen, Th. N. (1995). Improving the accuracy of institute for scientific  
32 information's journal impact factors. *Journal of the American Society for Information Science*,  
33 46(6), 461–467. [https://doi.org/10.1002/\(SICI\)1097-4571\(199507\)46:6<461::AID-](https://doi.org/10.1002/(SICI)1097-4571(199507)46:6<461::AID-ASI5>3.0.CO;2-G)  
34 ASI5>3.0.CO;2-G
- 35 Moed, Henk F. (2008). UK Research Assessment Exercises: Informed judgments on research  
36 quality or quantity? *Scientometrics*, 74(1), 153–161. [https://doi.org/10.1007/s11192-008-](https://doi.org/10.1007/s11192-008-0108-1)  
37 0108-1
- 38 Moed, Henk F. (2017). *Applied Evaluative Informetrics*. Cham: Springer.
- 39 Muñoz-Écija, T., Vargas-Quesada, B., & Chinchilla Rodríguez, Z. (2019). Coping with  
40 methods for delineating emerging fields: Nanoscience and nanotechnology as a case study.  
41 *Journal of Informetrics*, 13(4), 100976. <https://doi.org/10.1016/j.joi.2019.100976>
- 42 Perianes-Rodríguez, A., & Ruiz-Castillo, J. (2018). The impact of classification systems in the  
43 evaluation of the research performance of the Leiden Ranking universities. *Journal of the  
44 Association for Information Science and Technology*, 69(8), 1046–1053.  
45 <https://doi.org/10.1002/asi.24017>
- 46 Rafols, I., Ciarli, T., & Chavarro, D. (2019). Under-reporting research relevant to local needs  
47 in the South: Database biases in rice research. In R. Arvanitis & D. O'Brien (Eds.), *The  
48 Transformation of Research in the South Policies and outcomes*. Éditions des archives  
49 contemporaines. <https://digital.csic.es/handle/10261/226953>. Accessed 28 November 2023



- 1 Ràfols, I., Molas-Gallart, J., Chavarro, D. A., & Robinson-Garcia, N. (2016). *On the*  
2 *Dominance of Quantitative Evaluation in 'Peripheral' Countries: Auditing Research with*  
3 *Technologies of Distance* (SSRN Scholarly Paper No. ID 2818335). Rochester, NY: Social  
4 Science Research Network. <https://papers.ssrn.com/abstract=2818335>. Accessed 7 January  
5 2019
- 6 Robinson-Garcia, N., & Calero-Medina, C. (2014). What do university rankings by fields rank?  
7 Exploring discrepancies between the organizational structure of universities and bibliometric  
8 classifications. *Scientometrics*, 98(3), 1955–1970. [https://doi.org/10.1007/s11192-013-1157-](https://doi.org/10.1007/s11192-013-1157-7)  
9 7
- 10 Robinson-Garcia, N., Torres-Salinas, D., Vargas-Quesada, B., Chinchilla-Rodríguez, Z., &  
11 Gorraiz, J. (2023). Errors of measurement in scientometrics: Identification and calculation of  
12 systematic errors. In *Proceedings of ISSI 2023 – the 19th International Conference of the*  
13 *International Society for Scientometrics and Informetrics* (Vol. 2, pp. 387–393). Presented at  
14 the ISSI 2023, Bloomington, IN (United States). <https://doi.org/10.5281/zenodo.8428899>
- 15 Robinson-Garcia, N., Van Leeuwen, Th. N., & Torres-Salinas, D. (2020). Measuring Open  
16 Access uptake: Data sources, expectations, and misconceptions. *Scholarly Assessment*  
17 *Reports*. <https://doi.org/10.5281/zenodo.4071143>
- 18 Ruiz-Castillo, J., & Waltman, L. (2015). Field-normalized citation impact indicators using  
19 algorithmically constructed classification systems of science. *Journal of Informetrics*, 9(1),  
20 102–117. <https://doi.org/10.1016/j.joi.2014.11.010>
- 21 Schloegl, C., & Gorraiz, J. (2010). Comparison of citation and usage indicators: the case of  
22 oncology journals. *Scientometrics*, 82(3), 567–580. [https://doi.org/10.1007/s11192-010-](https://doi.org/10.1007/s11192-010-0172-1)  
23 0172-1
- 24 Scuro, S. R. (2004). Introduction to error theory. *Visual Physics Laboratory, Texas A&M*  
25 *University, College Station, TX, 77843*.  
26 <http://web.ist.utl.pt/~mcasquilho/compute/errtheory/basics/ScuroErrTheo.pdf>. Accessed 28  
27 November 2023
- 28 Selivanova, I. V., Kosyakov, D. V., & Guskov, A. E. (2019). The Impact of Errors in the Scopus  
29 Database on the Research Assessment. *Scientific and Technical Information Processing*,  
30 46(3), 204–212. <https://doi.org/10.3103/S0147688219030109>
- 31 Shu, F., Julien, C.-A., Zhang, L., Qiu, J., Zhang, J., & Larivière, V. (2019). Comparing journal  
32 and paper level classifications of science. *Journal of Informetrics*, 13(1), 202–225.  
33 <https://doi.org/10.1016/j.joi.2018.12.005>
- 34 Stahlschmidt, S., & Stephen, D. (2022). From indexation policies through citation networks to  
35 normalized citation impacts: Web of Science, Scopus, and Dimensions as varying resonance  
36 chambers. *Scientometrics*, 127(5), 2413–2431. <https://doi.org/10.1007/s11192-022-04309-6>
- 37 Sugimoto, C. R., & Weingart, S. (2015). The kaleidoscope of disciplinarity. *Journal of*  
38 *Documentation*, 71(4), 775–794. <https://doi.org/10.1108/JD-06-2014-0082>
- 39 Thijs, B., Zhang, L., & Glänzel, W. (2015). Bibliographic coupling and hierarchical clustering  
40 for the validation and improvement of subject-classification schemes. *Scientometrics*, 105(3),  
41 1453–1467. <https://doi.org/10.1007/s11192-015-1641-3>
- 42 Torres-Salinas, D., Arroyo-Machado, W., & Robinson-Garcia, N. (2023). Bibliometric  
43 denialism. *Scientometrics*, 128(9), 5357–5359. <https://doi.org/10.1007/s11192-023-04787-2>
- 44 van Leeuwen, T. N., Moed, H. F., Tijssen, R. J. W., Visser, M. S., & Raan, A. F. J. V. (2001).  
45 Language biases in the coverage of the Science Citation Index and its consequences for  
46 international comparisons of national research performance. *Scientometrics*, 51(1), 335–346.  
47 <https://doi.org/10.1023/A:1010549719484>
- 48 Visser, M., van Eck, N. J., & Waltman, L. (2021). Large-scale comparison of bibliographic data  
49 sources: Scopus, Web of Science, Dimensions, Crossref, and Microsoft Academic.  
50 *Quantitative Science Studies*, 1–22. [https://doi.org/10.1162/qss\\_a\\_00112](https://doi.org/10.1162/qss_a_00112)

- 1 Waltman, L., Calero-Medina, C., Kosten, J., Noyons, E. C. M., Tijssen, R. J. W., van Eck, N.  
2 J., et al. (2012). The Leiden ranking 2011/2012: Data collection, indicators, and interpretation.  
3 *Journal of the American Society for Information Science and Technology*, 63(12), 2419–2432.  
4 <https://doi.org/10.1002/asi.22708>
- 5 Waltman, L., & van Eck, N. J. (2012). A new methodology for constructing a publication-level  
6 classification system of science. *Journal of the American Society for Information Science and*  
7 *Technology*, 63(12), 2378–2392. <https://doi.org/10.1002/asi.22748>
- 8 Wang, Q., & Waltman, L. (2016). Large-scale analysis of the accuracy of the journal  
9 classification systems of Web of Science and Scopus. *Journal of Informetrics*, 10(2), 347–  
10 364. <https://doi.org/10.1016/j.joi.2016.02.003>
- 11