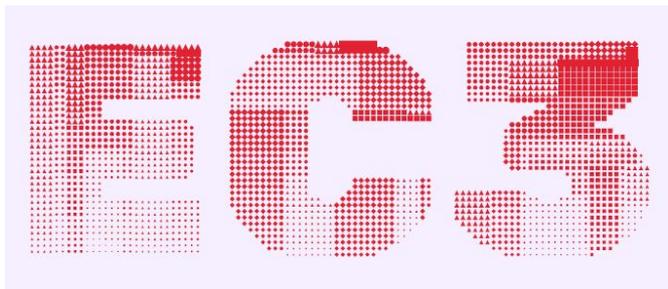


# Introducción al BIG DATA

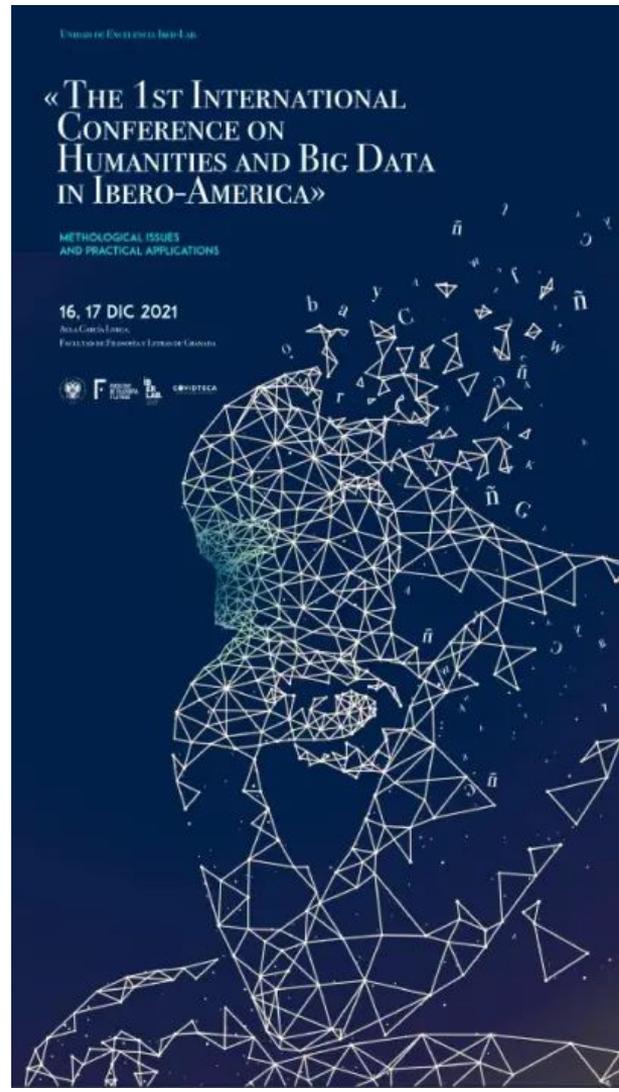


# Una introducción al Big Data



**Daniel Torres-Salinas**

**Sara Mariottini**





# Introducción y aplicaciones

*El big data sugiere un conocimiento absoluto. Las cosas revelan sus correlaciones secretas. Todo se vuelve calculable, predecible y controlable. Se anuncia toda una nueva era del saber*

Han, Byung-Chul

~~Generamos, recolectamos y almacenamos datos~~

Generamos, recolectan y almacenan nuestros datos

**Sensores**

RFID, GPS, Sensores, pulseras, drones...

**Internet**

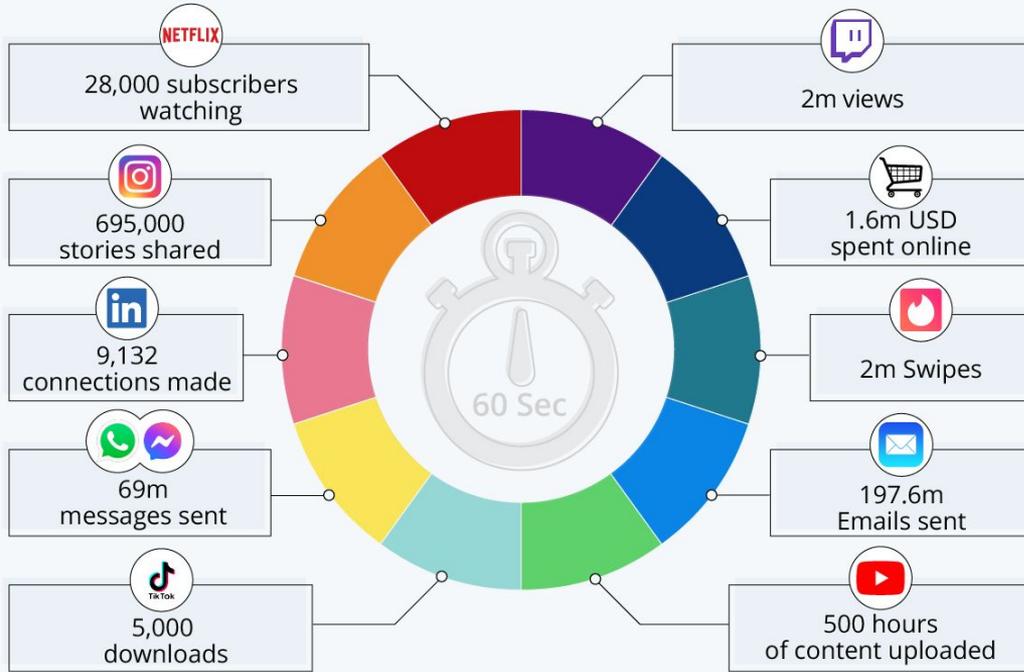
Click, mails, chats, búsquedas, tweets, transacciones...

**Real world**

Vídeo vigilancia, Cash Flows, Tráfico, smart grid, seguros, ...

# A Minute on the Internet in 2021

Estimated amount of data created on the internet in one minute



Source: Lori Lewis via AllAccess

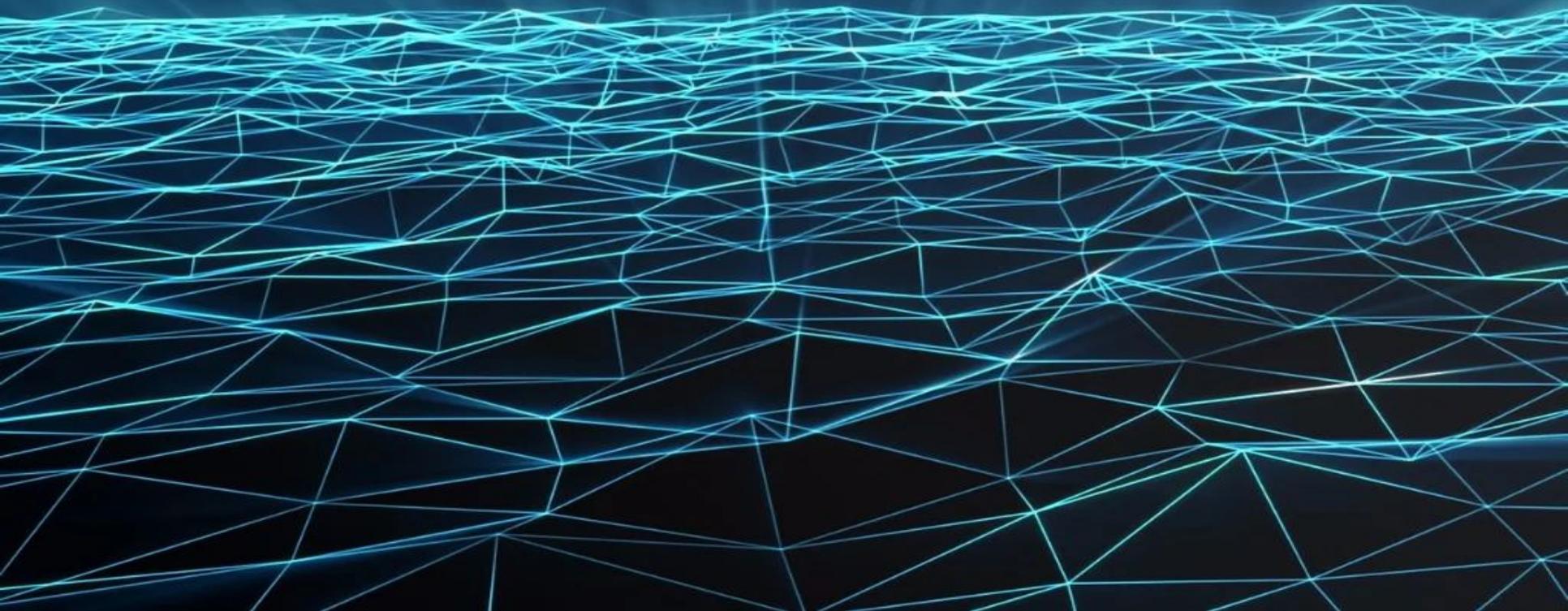
# Internet de la personas



# The Internet of Things: How Interconnectivity Is Changing Our World

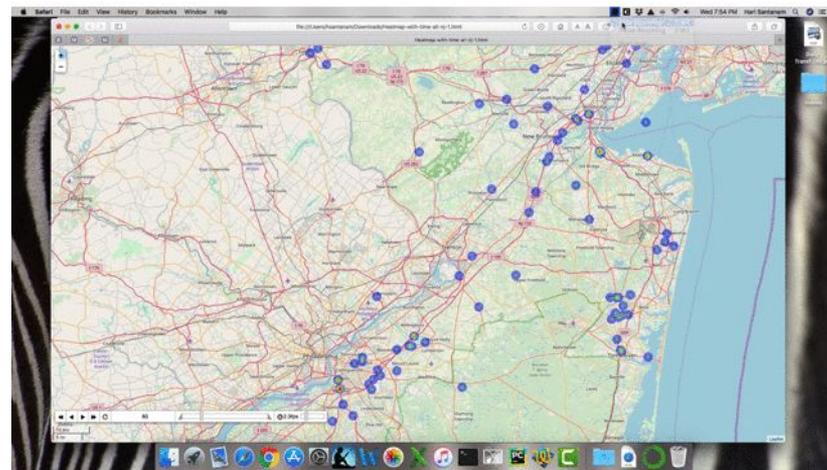
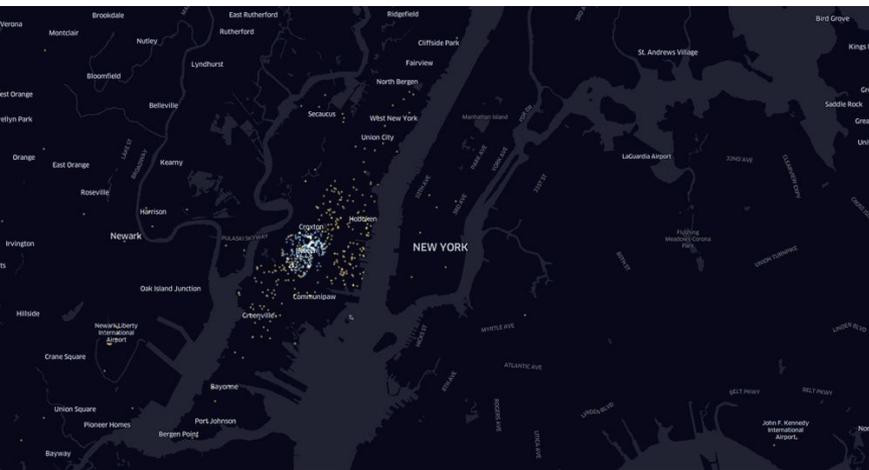


# DATA LAKE



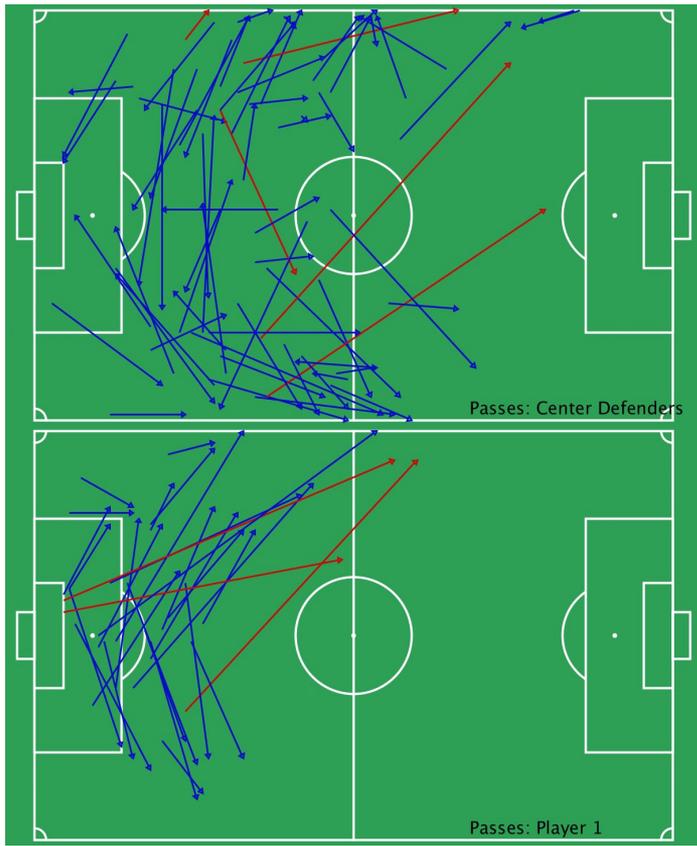


# Tesla's autonom ous-car use of Big Data



**Engineering Intelligence  
Through Data  
Visualization at Uber**

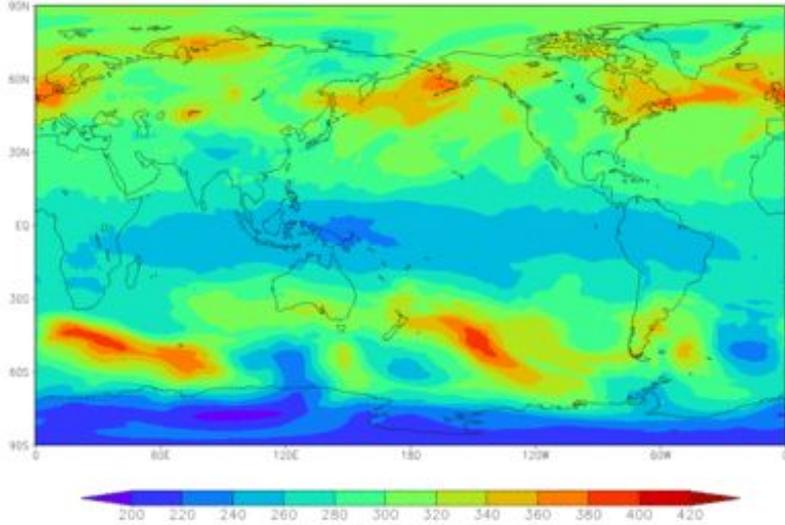
**Real world data science  
project: traffic accident  
analysis**



## **Data Mining for Strategy development in Football**

Broadcast tracking y equipos  
fantasmas, últimas noticias del  
Big Data en el fútbol

GFS Entire Atmosphere Total Ozone [Dobson]  
00Z12JUL2012+000Hrs



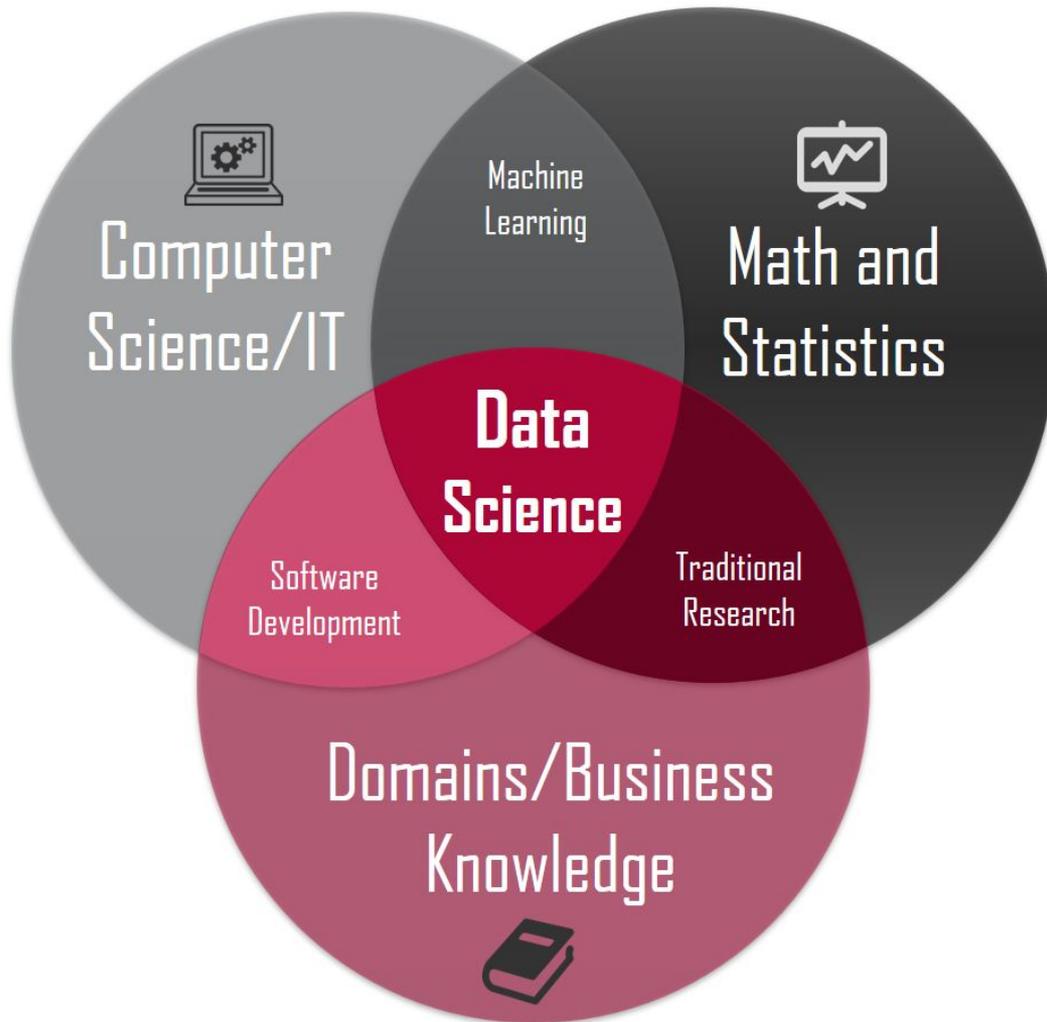
# Introducing Weacast

# Applying Data to the Next Phase of the COVID-19 Pandemic





**Elementos del big data**



**Podemos hablar de un campo multidisciplinar que usa métodos científicos , procesos y algoritmos para extraer conocimiento de datos estructurados y no estructurados**

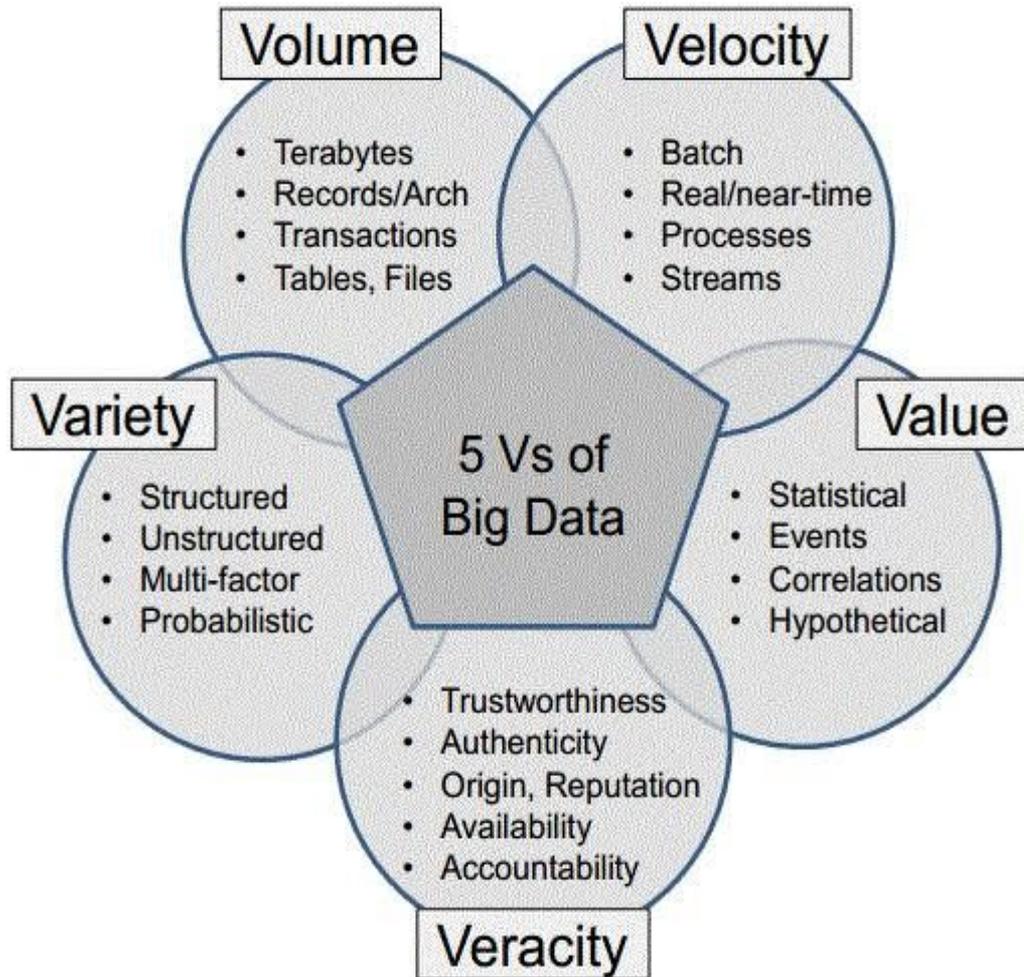
Estadística  
Bussines Analytics  
Bussines Intelligence  
Gestión de Base de datos  
Visualización  
Machine Learning  
Data Mining  
Artificial intelligence  
Modelos y predicciones

# Una definición como otra cualquiera

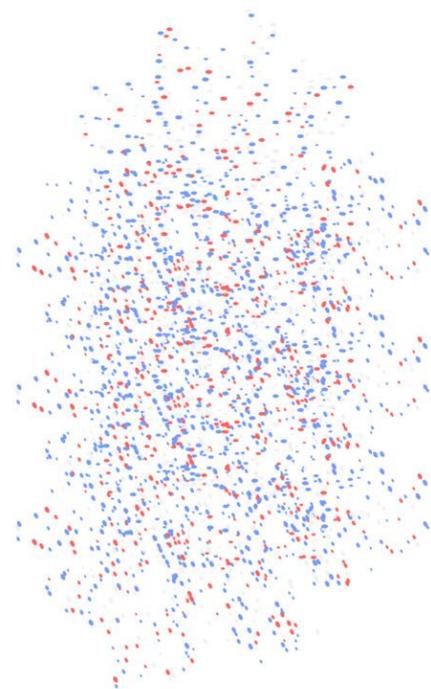
Large amounts of different types of data produced from various types of sources, such as people, machines or sensors. This data includes climate information, satellite imagery, digital pictures and videos, transition records or GPS signals. Big Data may involve personal data: that is, any information relating to an individual, and can be anything from a name, a photo, an email address, bank details, posts on social networking websites, medical information, or a computer IP address

# Una definición como otra cualquiera

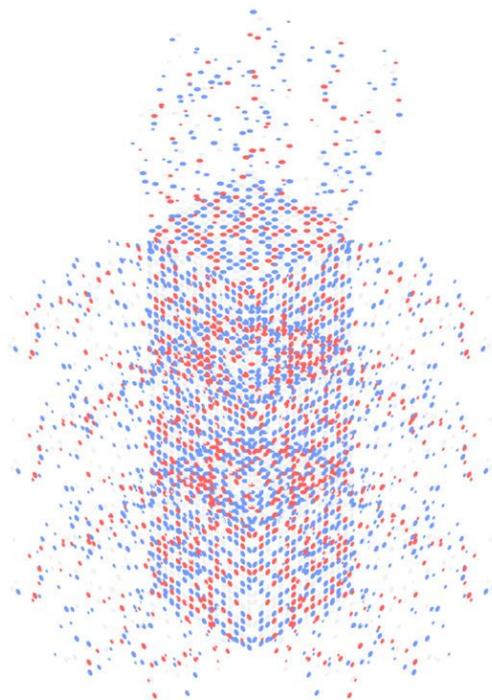
El big data trabaja con datos gigantescos, han de ser, capturados y almacenados y con capacidad para ser permanente actualizados y posteriormente analizados. Asimismo estos procesos de análisis, normalmente relacionados con la visualización y predicción, implican la puesta en marcha de métodos que nos permitan extraer el valor de los datos. Estas técnicas de análisis están sobre todo orientadas a tres objetivos principales, como son la búsqueda de patrones, la identificación de asociaciones y elaboración de modelos de pronóstico.



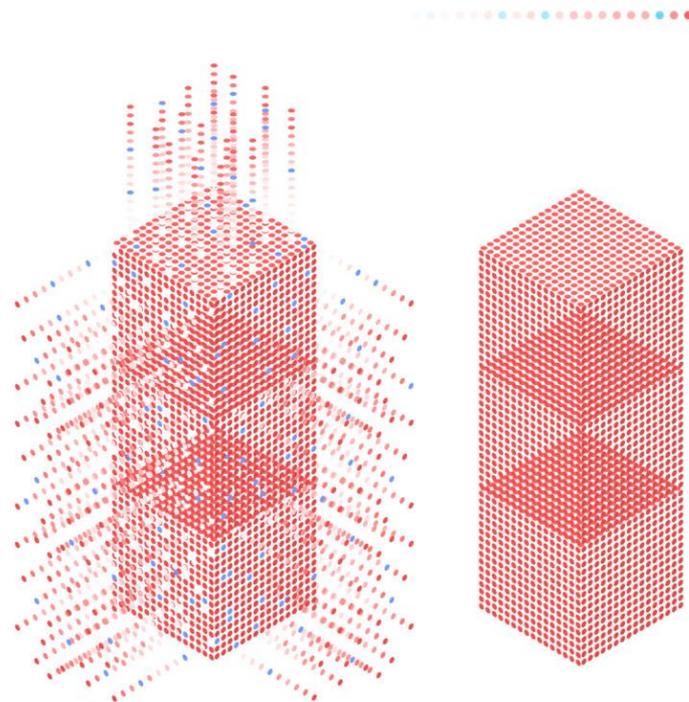
Unstructured data



Semi-structured data



Structured data



## Unstructured data

The university has 5600 students.  
John's ID is number 1, he is 18 years old and already holds a B.Sc. degree.  
David's ID is number 2, he is 31 years old and holds a Ph.D. degree. Robert's ID is number 3, he is 51 years old and also holds the same degree as David, a Ph.D. degree.

## Semi-structured data

```
<University>
  <Student ID="1">
    <Name>John</Name>
    <Age>18</Age>
    <Degree>B.Sc.</Degree>
  </Student>
  <Student ID="2">
    <Name>David</Name>
    <Age>31</Age>
    <Degree>Ph.D. </Degree>
  </Student>
  ....
</University>
```

## Structured data

ID	Name	Age	Degree
1	John	18	B.Sc.
2	David	31	Ph.D.
3	Robert	51	Ph.D.
4	Rick	26	M.Sc.
5	Michael	19	B.Sc.

¿Cómo analizamos estos datos?

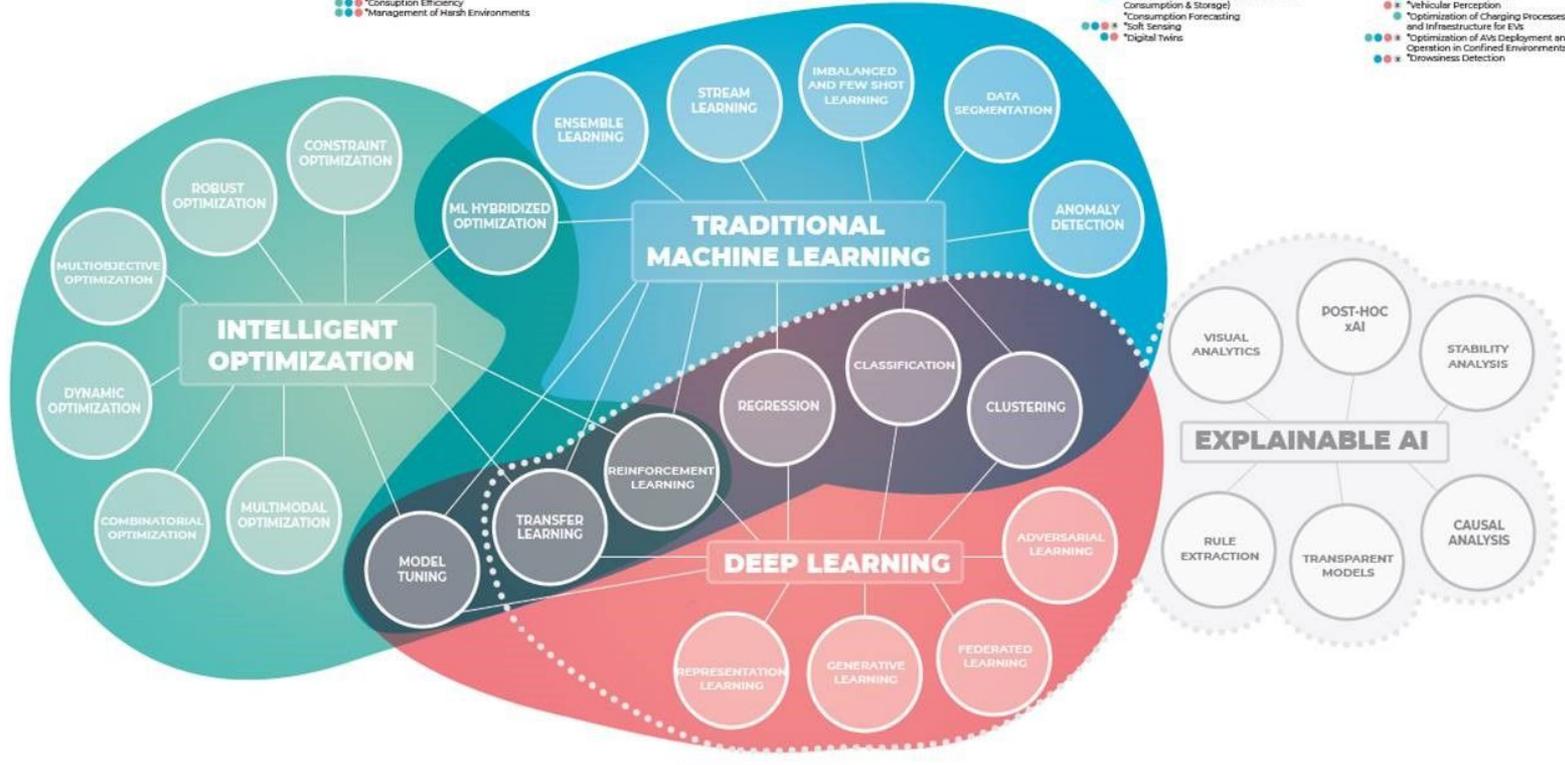
- INDUSTRY**
- Production Scheduling
  - Inventory Management
  - Fault Diagnosis
  - Tool Prediction
  - Demand Forecasting
  - Supply Chain Logistics
  - Maintenance Planning
  - Waste/Scrap management
  - Security (including inspection)
  - Truck/Container Load Optimization
  - Warehouse Optimization
  - Consumption Efficiency
  - Management of Harsh Environments

- HEALTH**
- Patient Characterization
  - Clinical Diagnostic Support
  - Resource Planning Management System
  - Predictive Models for Ovocytary Quality
  - Cognitive Status Characterization

- CYBERSECURITY**
- Synergy of ML & Blockchain in Federated Systems
  - ML in Adversarial Settings
  - Intrusion Detection System
- E-SERVICES**
- Market Basket Analysis
  - Preference / Recommendation System

- ENERGY**
- Analytics tools:
- Hybridization of Digital Twins and Data
  - Renewable Energy Forecasting
  - HVAC Optimization
  - Demand Side Management
  - Load Patterns from AMIs
  - Submetering
  - Tool estimation for medium voltage cables
  - Battery Life Prediction (Prognosis)
  - Microgrid Scheduling
  - Design of Energy Assets (in Production, Consumption & Storage)
  - Consumption Forecasting
  - Self Sensing
  - Digital Twins

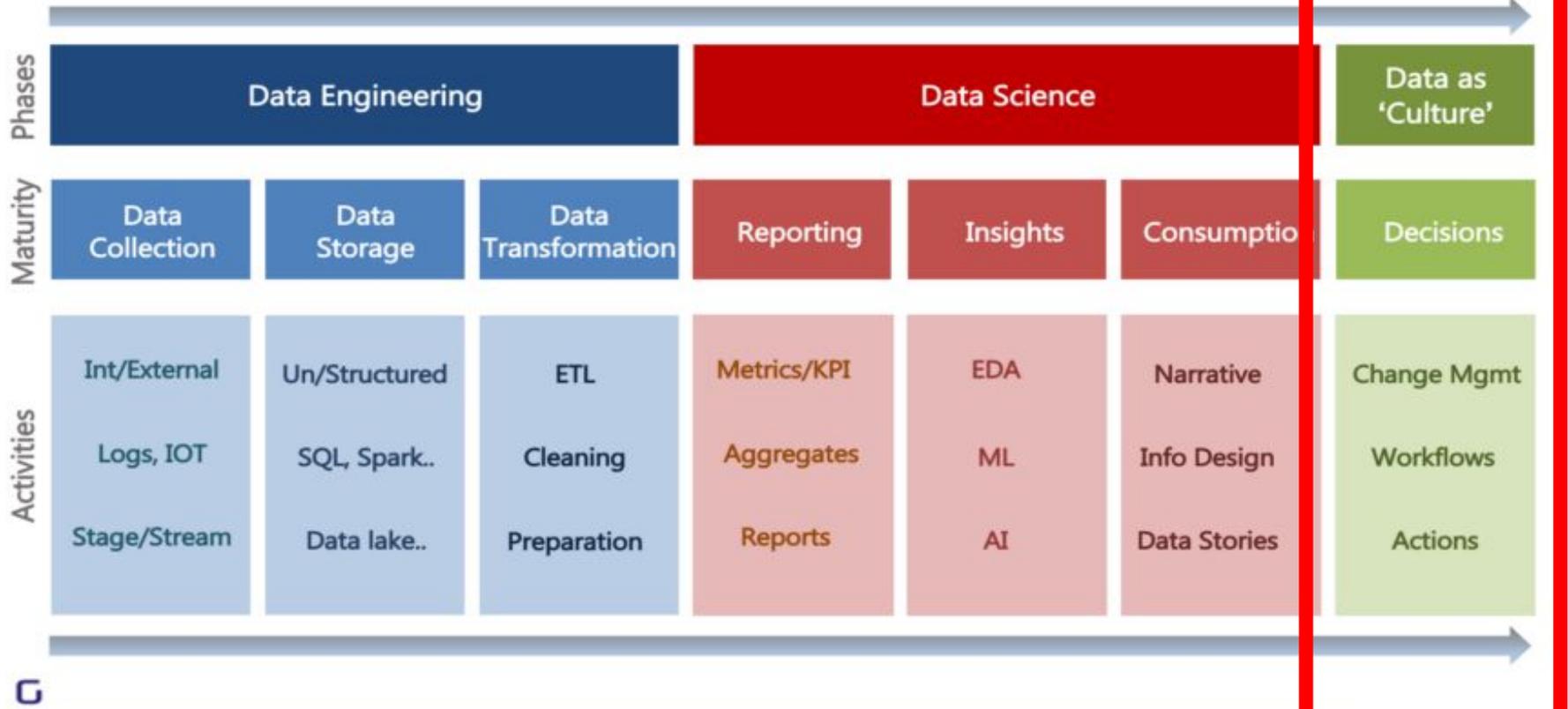
- MOBILITY**
- Multimodal Route Planning (Bike, car, public transportation...)
  - Traffic Forecasting
  - Operational Efficiency in Airports
  - Vehicular Communication
  - Accident Prediction
  - Infrastructure Allocation
  - Last-mile Delivery Scheduling
  - Optimum Car-sharing & Car-Pooling
  - Mobility Characterization
  - Turner Characterization
  - Vehicular Perception
  - Optimization of Changing Processes and Infrastructure for EVs
  - Optimization of AVs Deployment and Operation in Confined Environments
  - Towniness Detection



**BIG DATA**

- DISTRIBUTED OPTIMIZATION
- BIG DATA ML
- GREEN COMPUTING
- DISTRIBUTED LEARNING AND OPTIMIZATION
- DATA/JOBS AS/OPS
- CONTAINERS

# MATURITY LEVELS WITH DATA



HUMANISTA

??

Pic: Scaling the 7 levels of maturity with data

## THE 5 ROLES AND SKILLS NEEDED IN EVERY DATA SCIENCE TEAM



### Data Translator



- Domain expertise
- Business analysis
- Team leaders

### Data Scientist



- Statistics, ML, AI
- Identify insights
- Scripting skills

**HUMANISTA**

### Information Designer



- Information design
- Interaction design
- Visual design

### ML Engineer



- Software engineering
- Front/back-end coding
- DevOps

### Data Science Manager



- Project management
- Business analysis
- Change management



**las humanidades y las sociales**

# THE SPATIAL HUMANITIES

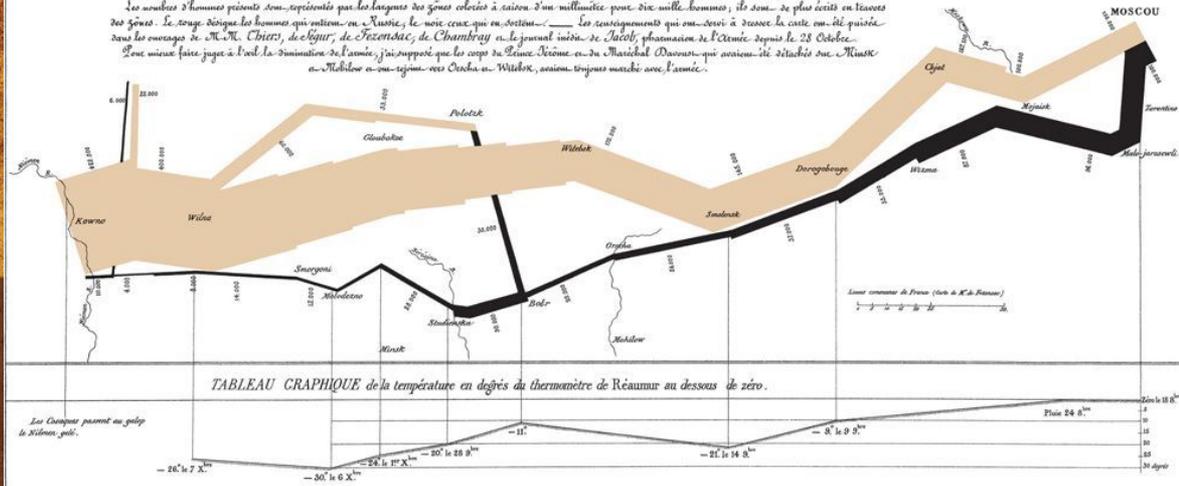
## GIS AND THE FUTURE OF HUMANITIES SCHOLARSHIP

EDITED BY DAVID J. BODENHAMER, JOHN CORRIGAN, AND TREVOR M. HARRIS

### Carte Figurative des aspects successifs en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.

Ordonné par M. Minaud, Inspecteur Général des Ponts et Chaussées en retraite. Paris, le 20 Novembre 1869.

Les marches d'hommes présents sont représentés par les lignes des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en lettres des zones. Le tracé désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. de Chateaubriand, de Ségur, de Ségur, de Chateaubriand et le journal inédit de Napoléon, par le Baron de Ségur le 28 Octobre. Les unités sont jugées à l'ordinaire de l'armée, j'ai supposé que les corps de Blücher, de Wittgenstein, de Dorsow, qui avaient été attachés sur Minsk et Mollat, n'ont rejoint nos Corps à Minsk, ainsi qu'on le voit sur la carte.



Atlas par Boyer, à Paris, 1789, page 21 et 22.

Map Lib. Reprint in Dresden.

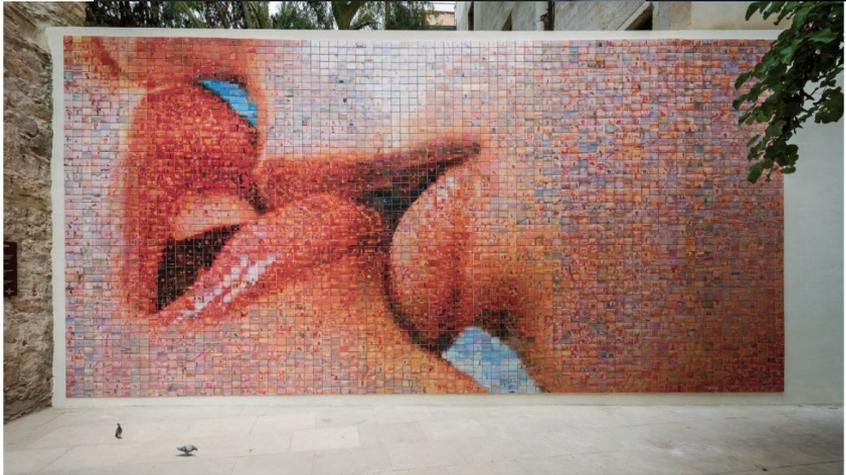
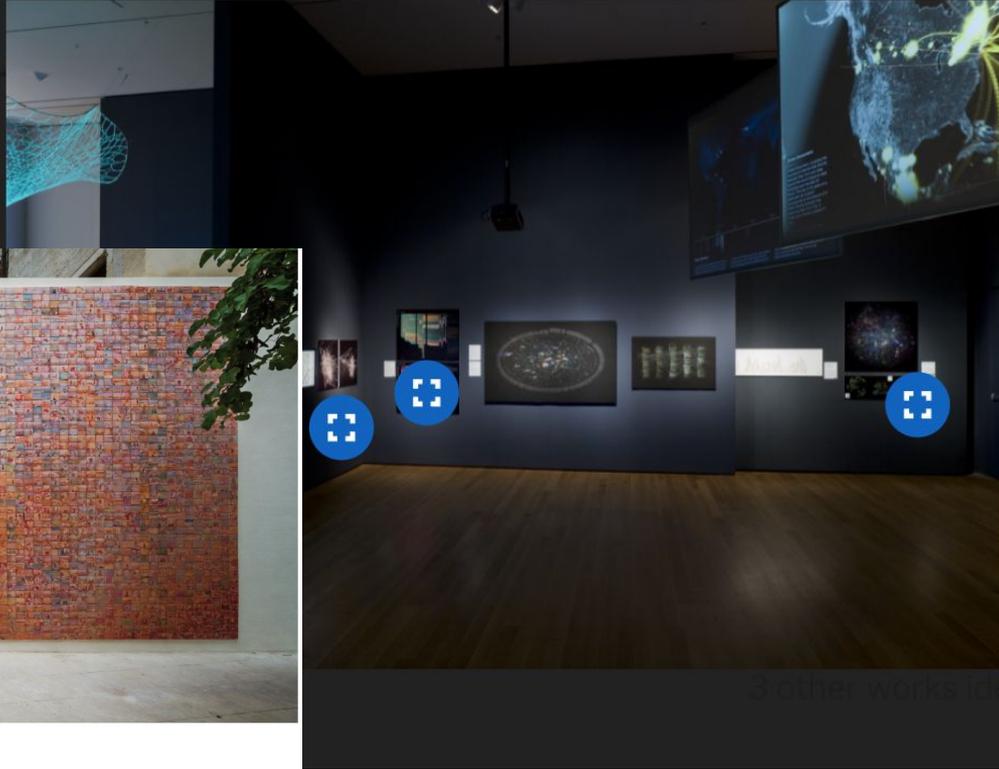


# Design and the Elastic Mind

Timed tickets    Member

Plan your visit    WI

We have identified these



Cultura

Art Data: cuando el Big Data emerge dentro del mundo del arte

A Collection of  
Feb 15, 2014-16  
Other works by  
Other work by

# MoMA

[Plan your visit](#)

[What's on](#)

[Art and artists](#)

[Store](#)



# Applied Design

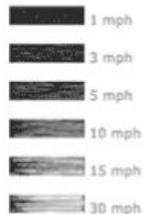
March 20, 2013

11:59 am EST

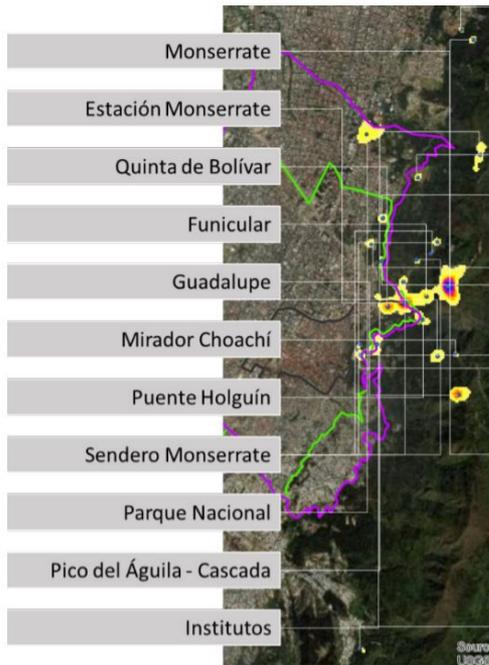
(time of forecast download)

top speed: 31.5 mph

average: 10.5 mph



# Big Data Meets Digital Cultural Heritage: Design and Implementation of SCRABS, A Smart Context-awaRe Browsing Assistant for Cultural Environments



Una propuesta de valoración de patrimonio urbano a la luz de Panoramio, una fuente Big-Data. El caso del centro de Bogotá

*An Urban Heritage Assessment with Panoramio, a Big-Data source. A Study of Downtown Bogotá*

Search: Sherlock Holmes, Frankenstein



1800 - 2019

English (2019)

Case-Insensitive

Smoothing of 5

### Choose corpus

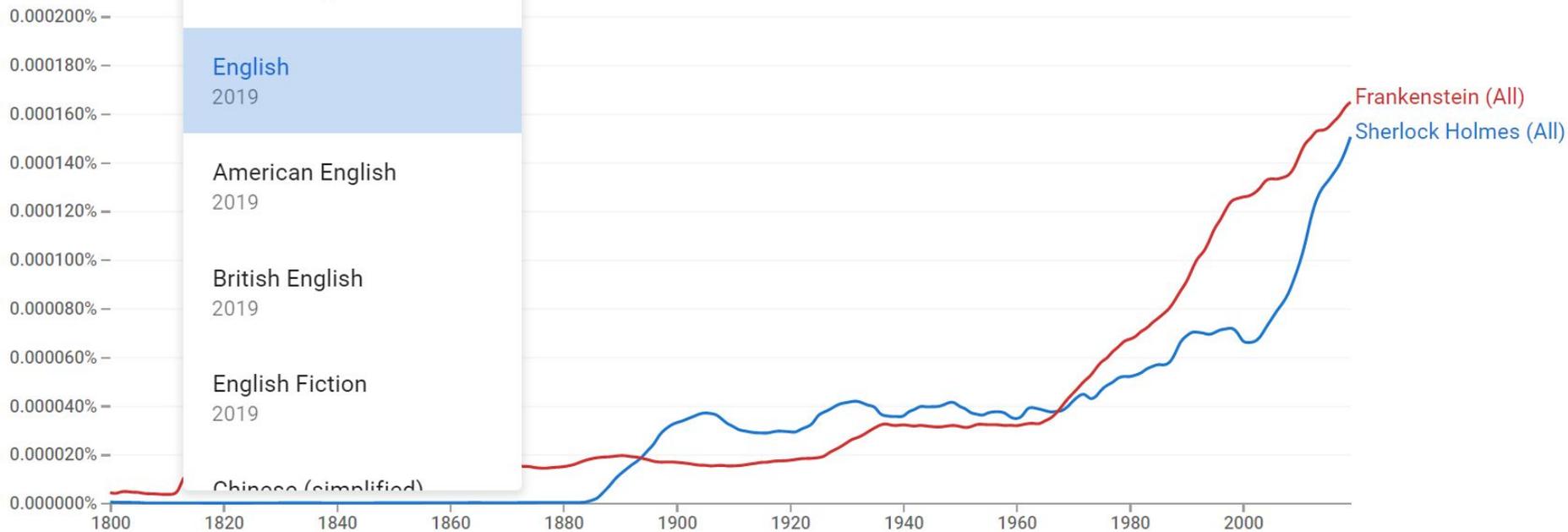
English  
2019

American English  
2019

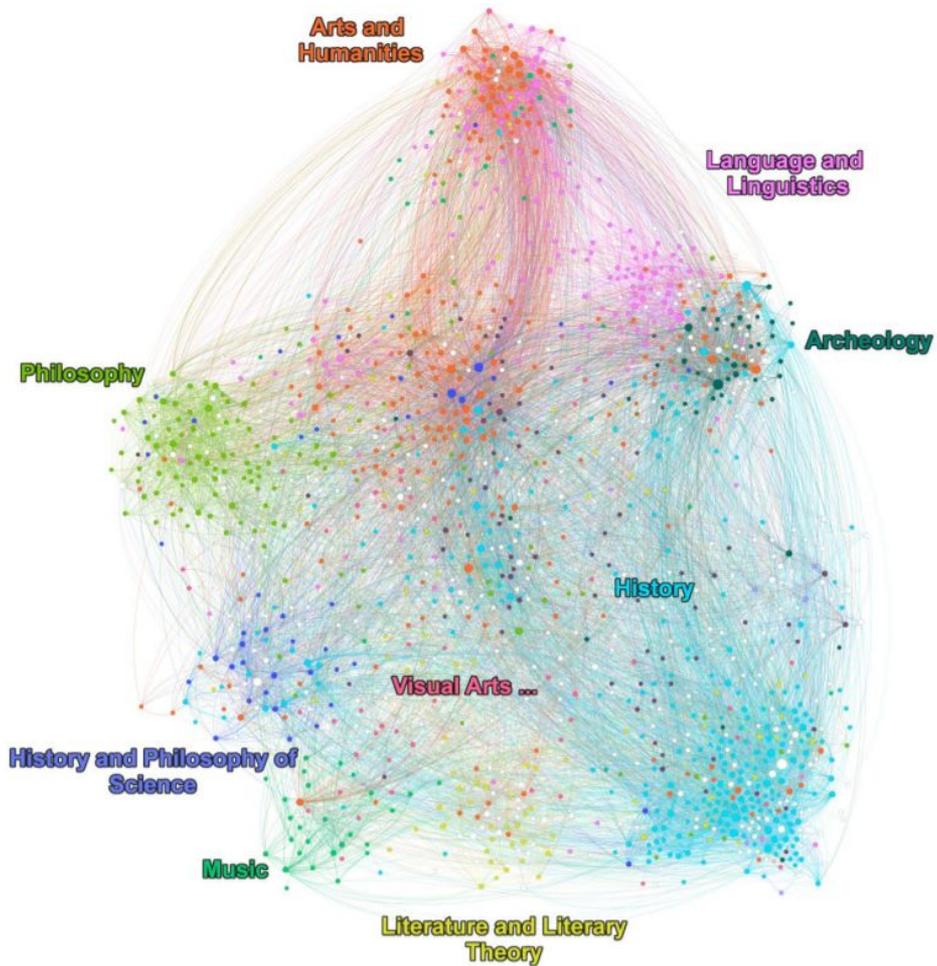
British English  
2019

English Fiction  
2019

Chinese (simplified)



(click on line/label for focus, right click to expand/contract wildcards)



Regular article

# Mapping the backbone of the Humanities through the eyes of Wikipedia

Daniel Torres-Salinas , Esteban Romero-Frías , Wenceslao Arroyo-Machado  

[Show more](#) 





Mentioned by

- 666 news outlets
- 129 blogs
- 11 policy sources
- 137644 tweeters
- 28 Facebook pages
- 18 Wikipedia pages

SUMMARY

News

Blogs

Policy documents

Twitter

Facebook

More...

**Title** Report 9: Impact of non-pharmaceutical interventions (NPIs) to reduce COVID19 mortality and healthcare demand

**DOI** 10.25561/77482 [↗](#)

**Authors** Ferguson, N, Laydon, D, Nedjati Gilani, G, Imai, N, Ainslie, K, Baguelin, M, Bhatia, S, Boonyasiri... [\[show\]](#)

[↗ View on publisher site](#)

[✉ Alert me about new mentions](#)



**Total mentions**  
197.3 million



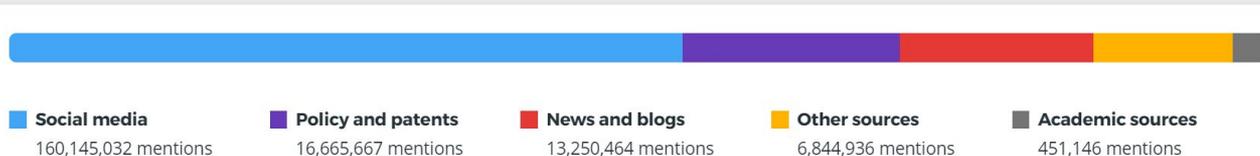
**Outputs with attention**  
19.7 million



**Total outputs tracked**  
35.7 million

[≡+ SAVE SEARCH](#)

Attention breakdown



コンピューター

Arts & Humanities

ARTS & HUMANITIES TOTAL MENTIONS / AUTORES Y PAPERS

Twitter	Wikipedia	News	Policy
16.5K	84	1.6K	6
9.5K	44	887	4



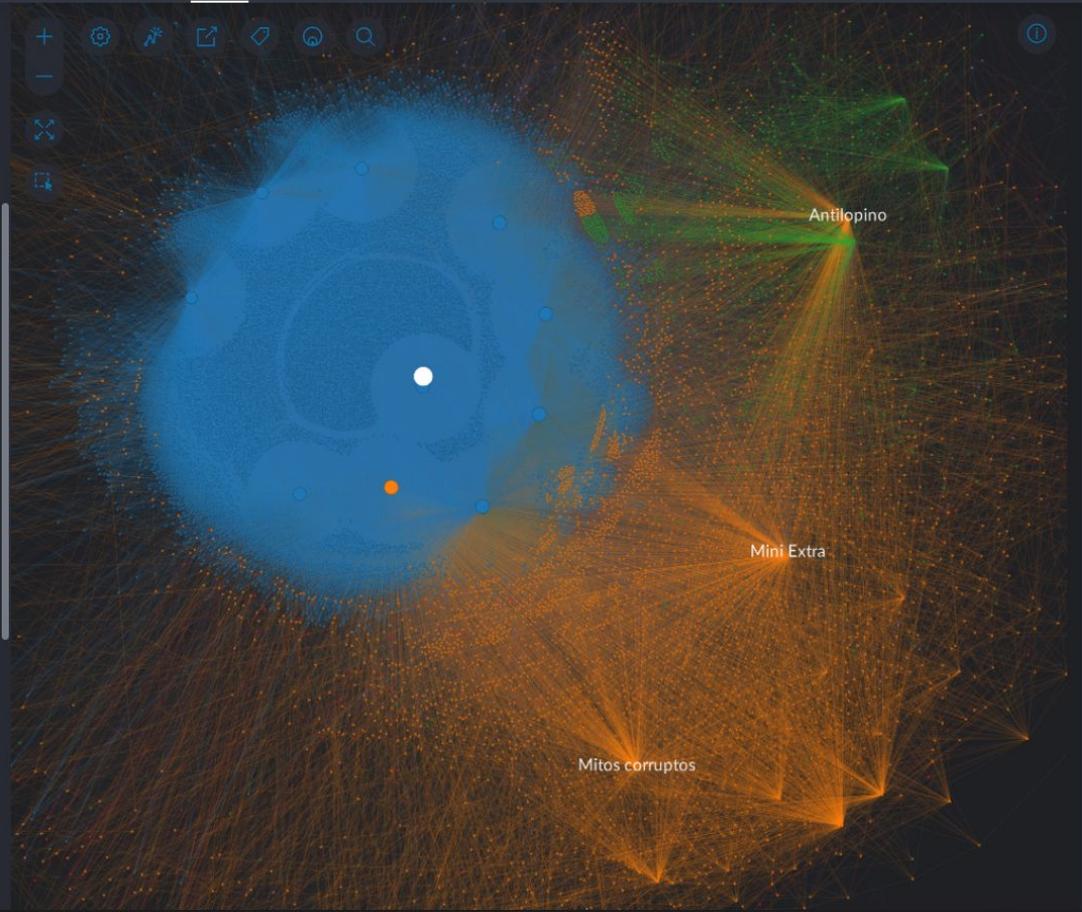
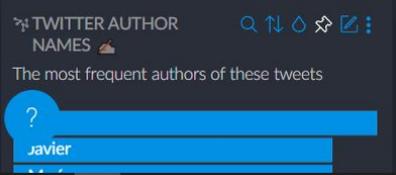
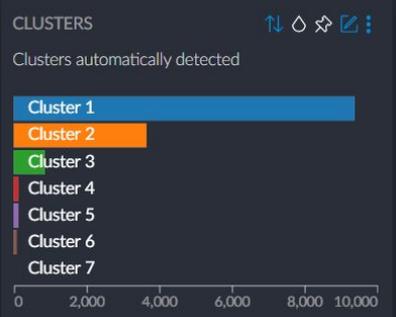
## Top 3 Personas InfluScieners

Autor/a	InfluRatio	Wikipedia	Twitter	News	Policy
 <b>1.</b> Jacob Morales UNIV LAS PALMAS GRAN CANARIA	<b>616</b>	0	227	56	0
 <b>2.</b> Diego Garate UNIV CANTABRIA	<b>563</b>	0	369	34	0

# Borges a través de las redes sociales

100% 14,138      CLEAR

+ New Segmentation



Search in 46 variables

Auto: Absolute

DATE

The date tweets were tweeted



DATE RANGE



FAVORITES COUNT

Count of favorites a user has

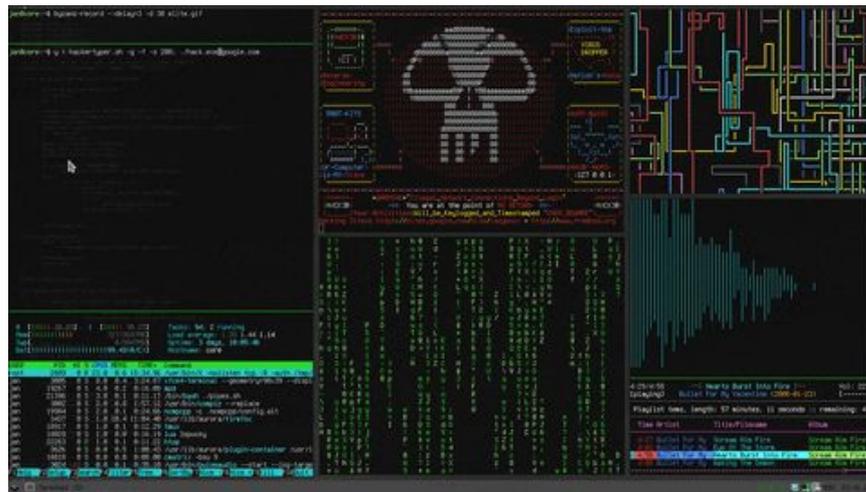


**Lado oscuro del big data**

FACEBOOK, CUATRO AÑOS EN EL CENTRO DE LA POLÉMICA >

## Caso #3: Cambridge Analytica, la gran fuga de datos

La información de más de 50 millones de usuarios de Facebook fueron utilizados sin consentimiento desde 2014 para comercializarlos ilegalmente con terceros



Cambridge Analytica, pero con sede en Estados Unidos, comenzaba a acaparar los titulares. A mediados del mes, The New York Times y The Observer revelan que en 2014 la compañía se hizo con una base de datos de Facebook para supuesto uso académico pero la explotó sin permiso para elaborar estrategias electorales. Entre ellas, prestó sus servicios a la candidatura del republicano Donald Trump, que el 8 noviembre ganó la presidencia frente a la demócrata Hillary Clinton.

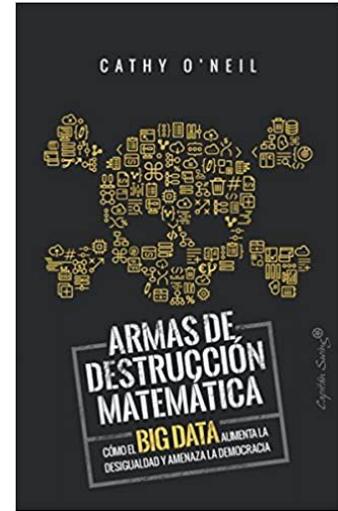


< NOTICIAS

July 24, 2019 1:05 pm

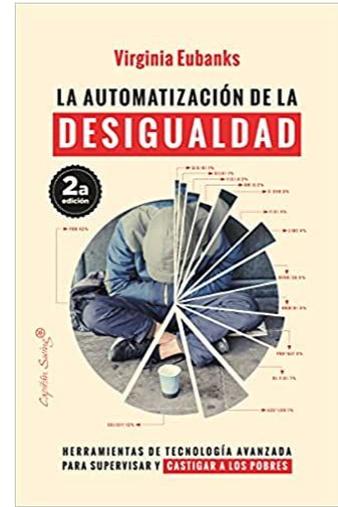
**“El gran hackeo”: Cambridge Analytica es sólo la punta del iceberg**

**AÑADIDOS  
RECIENTEMENTE**



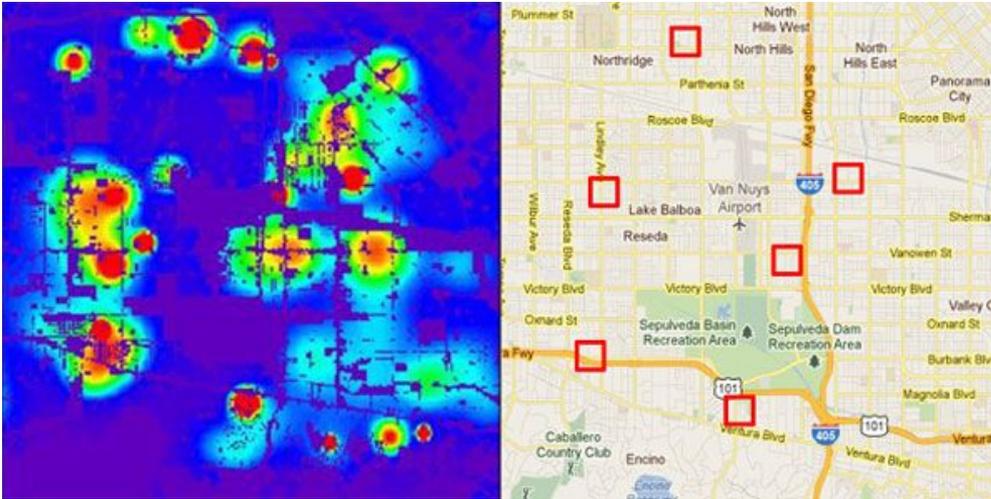
*“... Vivimos en la edad del algoritmo las decisiones no están hechas por humanos, sino por modelos matemáticos. Los modelos que se utilizan en la actualidad son opacos, no regulados e incontestables, incluso cuando están equivocados. Esto deriva en un refuerzo de la discriminación: si un estudiante pobre no puede obtener un préstamo porque un modelo de préstamo lo considera demasiado arriesgado (en virtud de su código postal), quedará excluido del tipo de educación que podría sacarlo de la pobreza...”*

**Los modelos apuntalan a los afortunados y castigan a los oprimidos**



**Virginia Eubanks investiga los impactos de la minería de datos, las políticas de algoritmos y los modelos de riesgo predictivo aplicados a las personas pobres y de clase trabajadora en Estados Unidos.**

*... En Los Ángeles, un algoritmo calcula la vulnerabilidad comparativa de decenas de miles de personas sin hogar con el fin de priorizarlas para un grupo insuficiente de recursos de vivienda. En Pittsburgh, una agencia de bienestar infantil utilizar un modelo estadístico para tratar de predecir qué niños podrían ser víctimas futuras de abuso o negligencia...*



# Did you know that police forces in the UK are trying to predict where crimes will happen and who will commit them?

The illustration shows a police officer in profile, wearing a yellow uniform jacket with a black collar and a black cap with a white checkered band. He is wearing glasses and pointing his right hand towards a smartphone held in his left hand. The background is a light blue map of a city street grid.

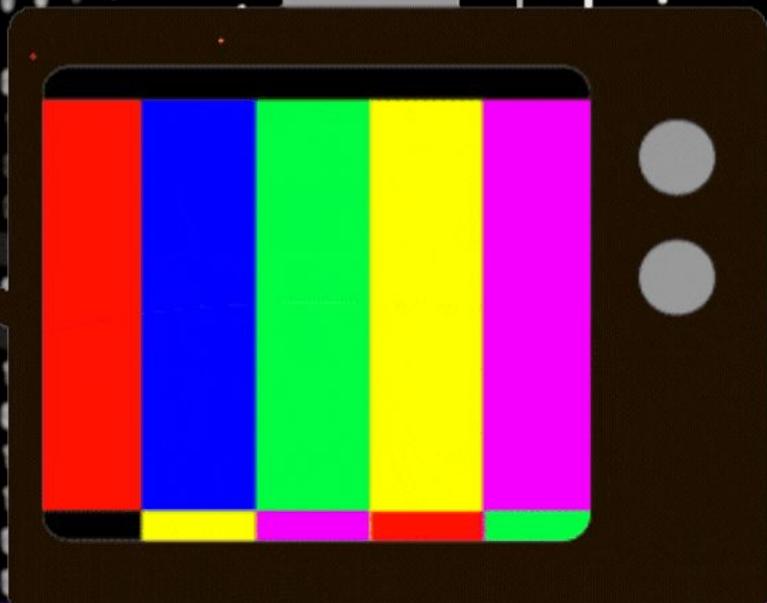
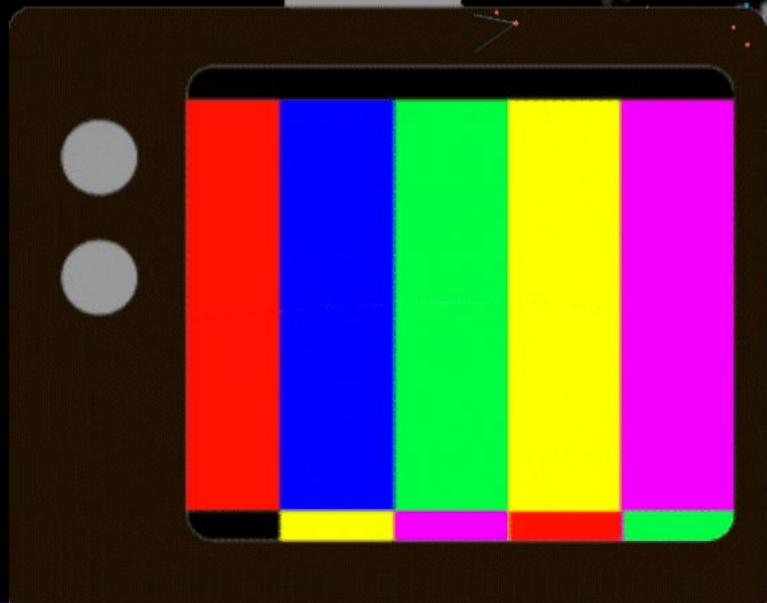
**PredPol**  
Predict Crime in **Real Time**™

PredPol provides targeted, real-time crime prediction designed for and successfully tested by officers in the field.

The inset map shows a street intersection with labels for Wilkes Cr, Pendegast Ave, and Emble Cr.

PIM

NETFLIX



# Style Is an Algorithm

No one is original anymore, not even you.

By Kyle Chayka | Apr 17, 2018, 10:00am EDT

*Illustrations by Momo Pixel*

☰ EL PAÍS

SUSCRÍBETE

INICIAR SESIÓN

EN PORTADA

## El gusto en la era del algoritmo

La prescripción artificial en plataformas digitales como Amazon, Netflix, Google o Facebook eleva el riesgo de homogeneizar la identidad y los hábitos de consumo cultural

Campus • NOTICIA 

# Humanidades contra la barbarie: "Dominar los algoritmos es como dominar las armas nucleares"

RUTH DÍAZ  | Madrid

8 MAY. 2019 | 02:20





# Métodos y técnicas



“ The big mystery of Big Data  
is causation versus  
correlation ”

Alexander Nix

empresario británico, ex director ejecutivo de  
Cambridge Analytica



1

## Clasificación de los métodos

Introducción

Sara Mariottini

---



Podemos agrupar los algoritmos de Machine Learning en tres macro grupos:

### Aprendizaje a partir de ejemplos.

Encontrar las dependencias deseadas utilizando un número "limitado" de observaciones. Este tipo de aprendizaje procede de la estadística. Existen dos subcategorías de este tipo de aprendizaje: el aprendizaje inductivo y el deductivo.

La mayoría de los enfoques de aprendizaje automático pertenecen a la categoría de aprendizaje inductivo.

- 1 Métodos basados en la supervisión humana en el proceso de aprendizaje de la Inteligencia Artificial: **Supervised learning (SL)**, **Unsupervised learning (UL)** **Semi-Supervised learning**, **Reinforcement learning (RL)**



2

Métodos basados en la generalización a partir de muestras de datos vista durante la fase de entrenamiento del modelo. Cuando la Inteligencia Artificial obtienen nuevas informaciones/instancias, esa viene **comparada** con la información almacenada en su **memoria**. Aquí no hay el manejo del **concept drift** de hecho se habla de una tipología de aprendizaje “perezoso”:

**Instance-based Learning, Model-based Learning.**

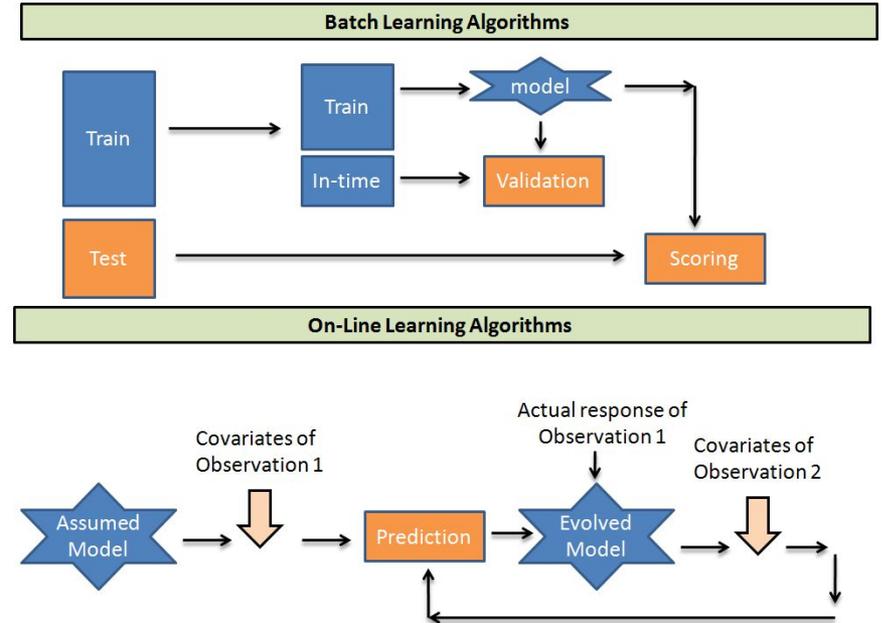


3

Métodos basados en la capacidad de aprender a partir de muestras de **datos incrementales** (ej.

algoritmos que manejan los cambios de conceptos):

**Batch Learning, Online Learning**



<https://www.analyticsvidhya.com/blog/2015/01/introduction-online-machine-learning-simplified-2/>



“Data is the new oil!”

Clive Humby

matemático británico y empresario de la ciencia de los datos



2

## Principales Técnicas de Machine Learning

Sara Mariottini

---



# Metodos y Tecnicas de Data Science

## Técnicas de Machine Learning



- Técnicas de Ciencia de datos más utilizadas en 2018/19 – Kdnuggets Poll.
- Los resultados están basados en 833 votantes
- En comparación con los años anteriores, el aumento más notable se produjo en el uso de diversas tecnologías de redes neuronales (Deep Learning)



# Técnicas de Machine Learning

## Supervised Learning – Clasificación

**Aprendizaje  
Automático  
Supervisado**



**Problemas de  
Clasificación**

**objetivo clave:** predecir las etiquetas de salida de los datos de entrada relacionados con lo que el modelo ha aprendido durante la fase de entrenamiento.

- **binaria** (dos clases, 0 y 1)
- **multiclase** (múltiples categorías)

### Algoritmos de clasificación más populares:

- **Decision trees**
- **K-Nearest Neighbors**
- **Support Vector Machines**
- **Deep Learning**

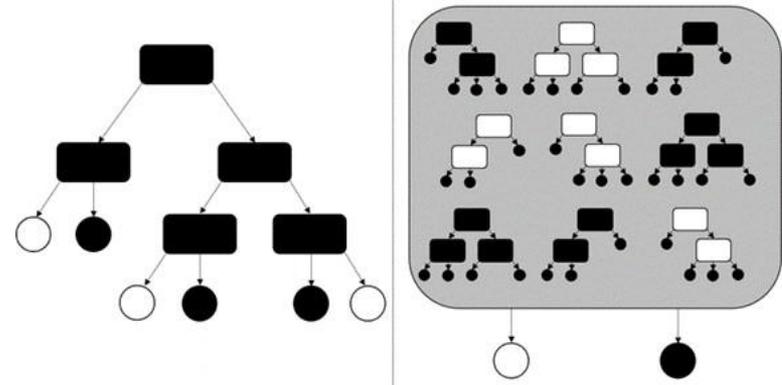


# Técnicas de Machine Learning

## Supervised Learning – Clasificación

### Decision Trees

- Modelo de aprendizaje supervisado jerárquico.
- Gráficos dirigidos / Diagrama de flujo - Lógica "booleana".
- Predice o clasifica el valor de una variable objetivo en función de varias variables de entrada.
- Se compone de **nodos** de decisión internos y **hojas** terminales: número menor de etapas define la región local en una serie de **divisiones secuenciales**. Pueden mostrarse gráficamente de forma que sean
- fáciles de interpretar para los no expertos.





# Técnicas de Machine Learning

## Supervised Learning – Clasificación

### K-Nearest Neighbors

#### Clasificadores basados en casos (k-NN)

*K-nearest neighbors* = *k-vecinos más cercanos*.

La clasificación consiste en encontrar los  $k$  vecinos más cercanos y se le asigna al nuevo dato la clase más común entre los  $k$  vecinos.

Cercanía  $\longrightarrow$  Medida de distancia

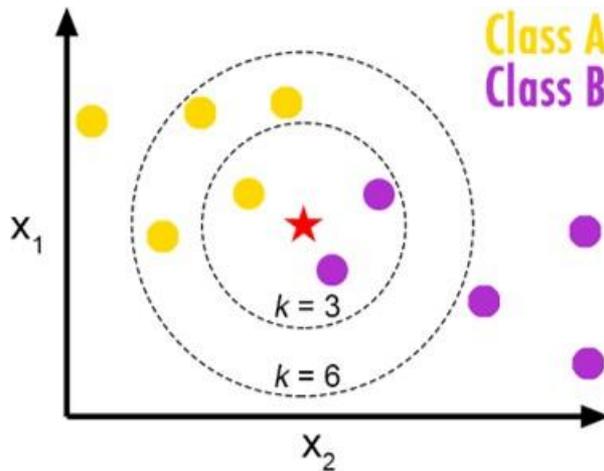
Ejemplo con 2 Atributos. Distancia Euclídea.

Datos de referencia, dataset de partida:

Nuevo punto que entra en el dataset:

$k=3$   $\longrightarrow$  Clase asignada: B

$k=6$   $\longrightarrow$  Clase asignada: A



Los valores de los atributos del  $i$ -ésimo ejemplo se presentan por el vector  $p$ -dimensional:

$$x_i = (x_{1i}, x_{2i}, \dots, x_{pi}) \in X$$

• Distancia Euclídea

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^p (x_{ri} - x_{rj})^2}$$

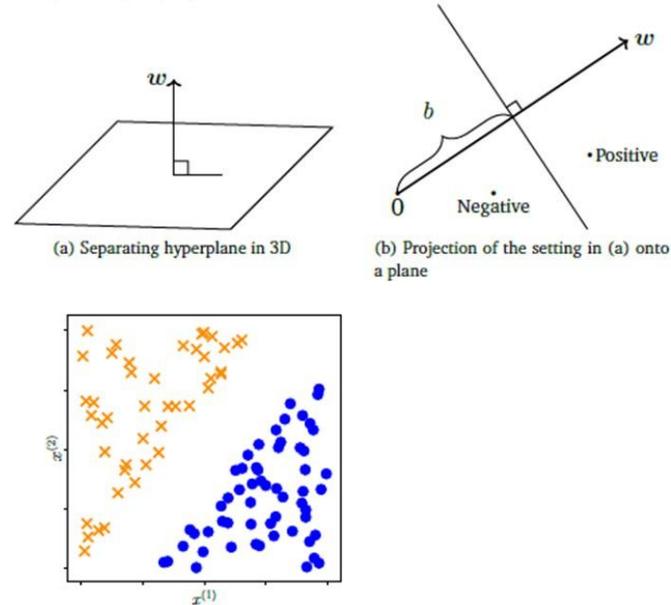


# Técnicas de Machine Learning

## Supervised Learning – Clasificación

### Support Vector Machines

- **Vectores de soporte:** puntos de datos más cercanos al hiperplano (dependen obviamente de la conformación del conjunto de datos).
- Representar la información del conjunto de datos en  $\{R\}^D$  y particionar este espacio de manera que sólo las observaciones con la misma etiqueta estén en la misma partición.
- **Objetivo:** localizar la función de clasificación más adecuada para separar las clases en los datos de entrenamiento.

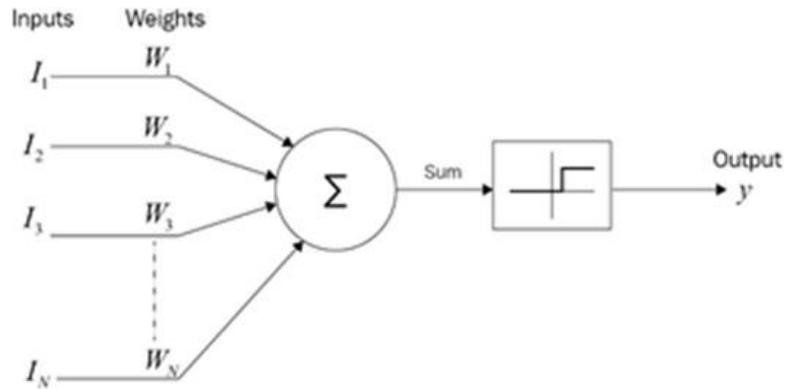




# Técnicas de Machine Learning

## Supervised learning - Deep Learning

### Deep Learning



Modelo red neuronal McCulloch-Pitts 1943.

- Subcampo del Machine Learning y se basa su proceso de aprendizaje en las redes neuronales artificiales (RNA).
- Inspirados en la profundidad arquitectónica del cerebro.
- Capa de entrada, una capa oculta y una capa de salida.
- Modelo puede aprender niveles de características de abstracción creciente.

### Sentiment analysis





# Técnicas de Machine Learning

## Unsupervised learning - Deep Learning

### Generative adversarial networks (GANs)



Figura 2.17.: CycleGAN aplicada a transferencia de estilo artístico

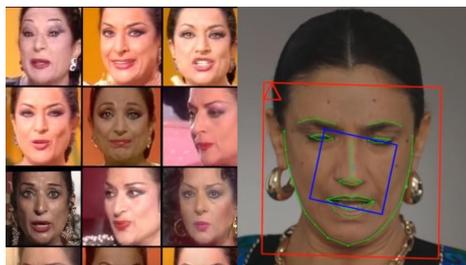
- Sistema de dos redes neuronales que compiten recíprocamente en una especie de juego de suma cero.
- Muy empleadas para la generación de imágenes “deep fakes”.



# Técnicas de Machine Learning

## Unsupervised learning - Deep Learning

### Campaña publicitaria viral de Cruzcampo (2021) - 'Deepfake' de Lola Flores



A woman with long, straight, light-colored hair and bangs is shown in profile, pointing her right index finger towards the right. She is in the foreground, looking towards a balcony in the distance. On the balcony, a small figure of a person is visible against the night sky. The scene is lit with a mix of blue and purple light, creating a cinematic atmosphere. The background shows a building with a balcony and a railing.

# Gracias

## Torres-Salinas & Mariottini