# Data Mining (Minería de Datos)

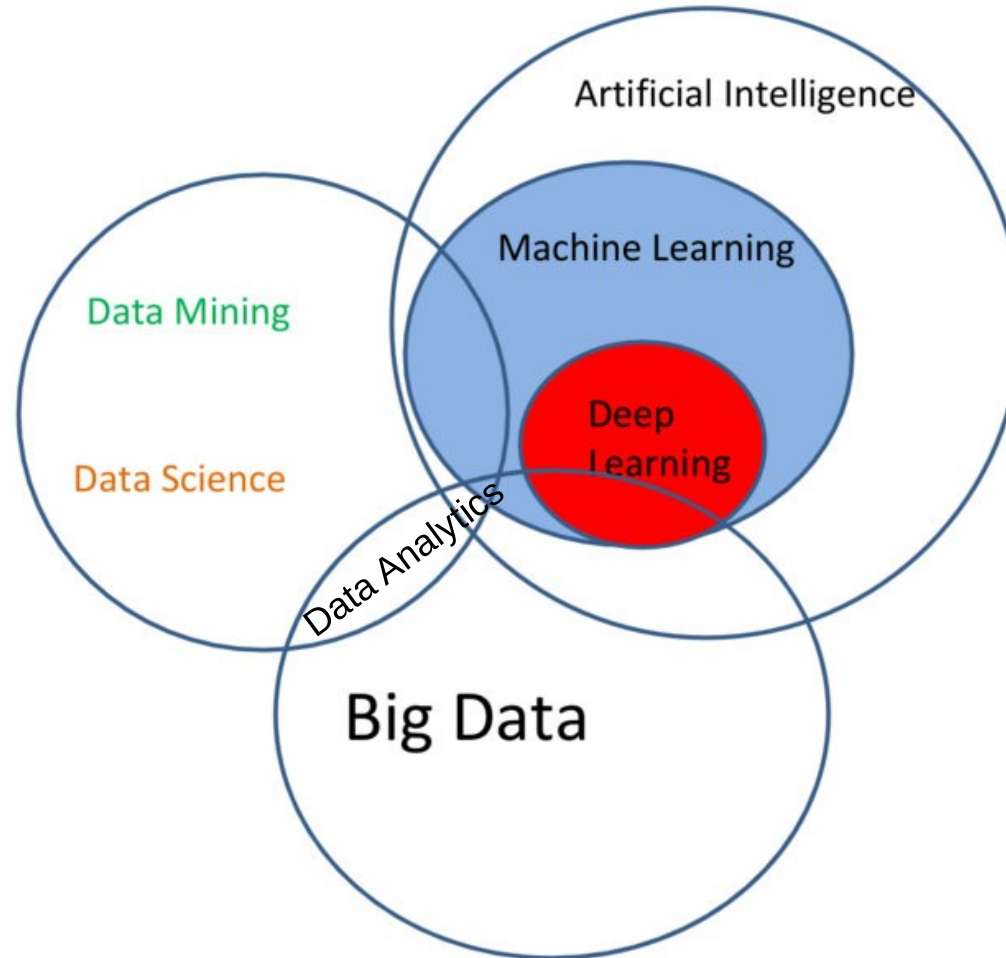# INTRODUCTION AND HISTORICAL PERSPECTIVE



**Sixto Herrera**

**Grupo de Meteorología**

**Univ. de Cantabria – CSIC**
**MACC / IFCA**

Master Universitario Oficial **Data Science**
con el apoyo del

UC
UNIVERSIDAD DE CANTABRIA

UIMP
Universidad Internacional
Menéndez Pelayo

CSIC
Consejo Superior de Investigaciones Científicas

**INTRO:** | **TERMINOLOGY & TRENDS** | 2

**Data Mining (DM)** can be defined as the process that starting from apparently unstructured data tries to extract knowledge and/or unknown interesting patterns.

**Machine Learning (ML)** relates with the study, design and development of the algorithms that give computers the capability to learn without being explicitly programmed (definition of A.Samuel).

Master Universitario Oficial **Data Science**

UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:** | **TERMINOLOGY & TRENDS** | 3

**Data Mining (DM)** can be defined as the process that starting from apparently unstructured data tries to extract knowledge and/or unknown interesting patterns.

**Machine Learning (ML)** relates with the study, design and development of the algorithms that give computers the capability to learn without being explicitly programmed (definition of A.Samuel).
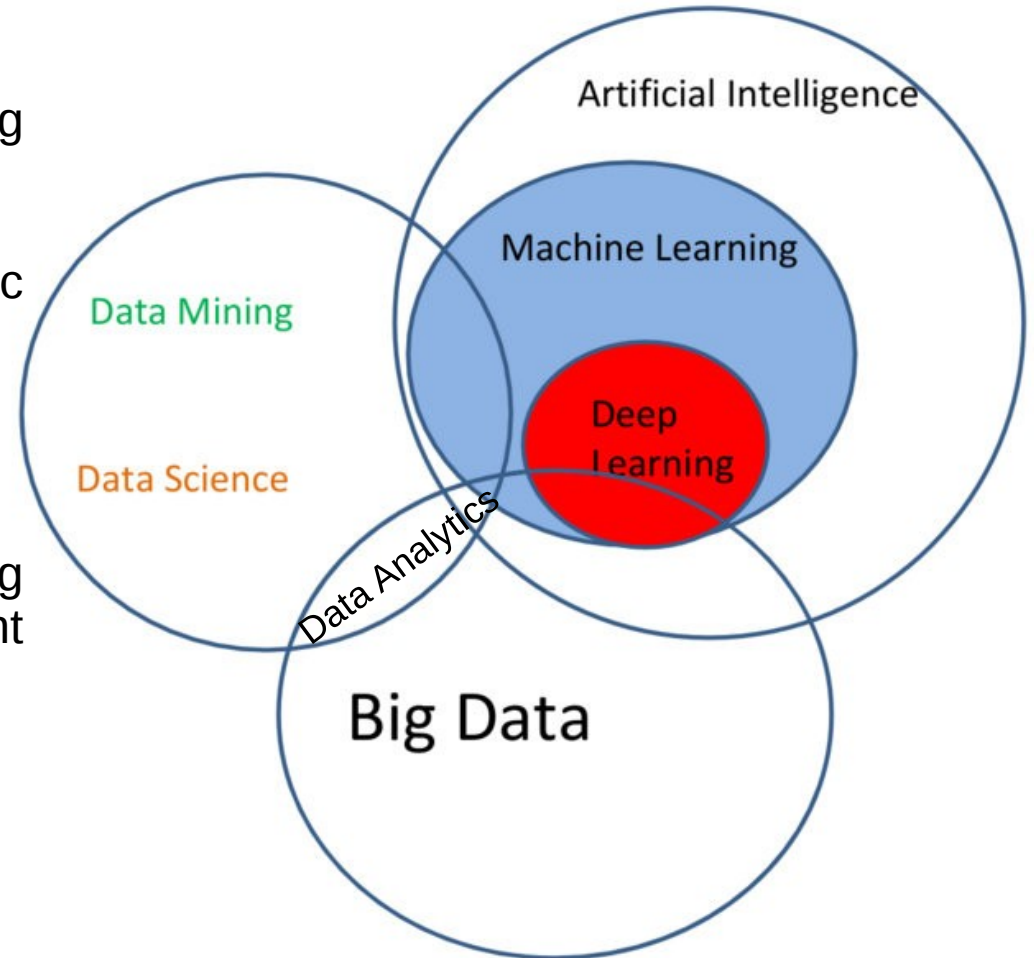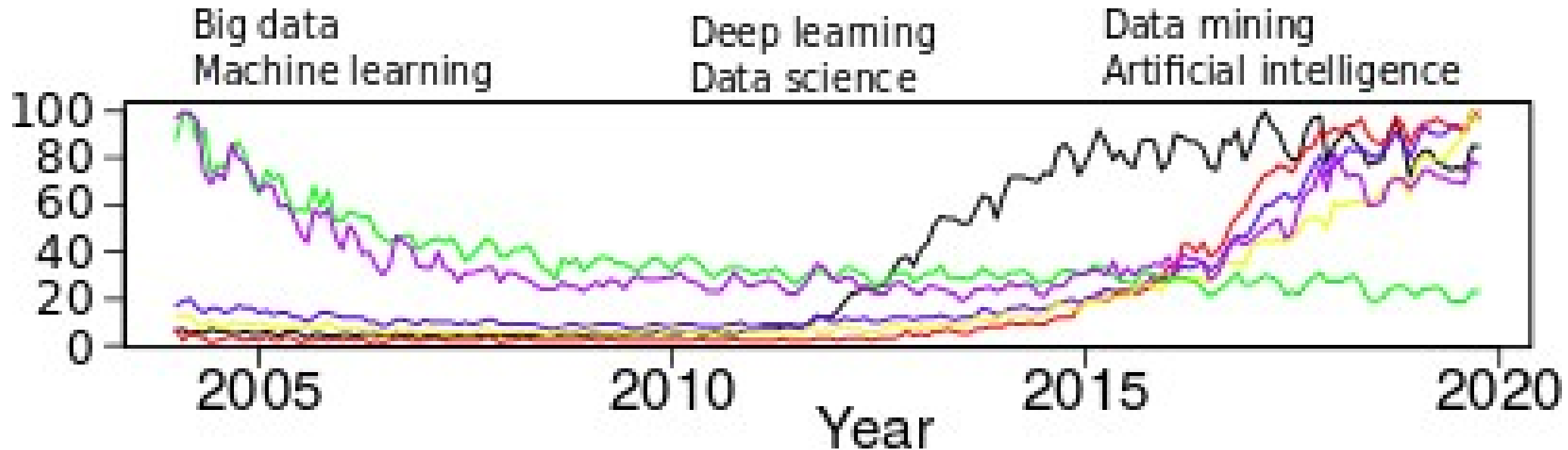
During this process machine learning algorithms are used (A. Flag).
ML Techniques → Generic
Data Mining → Understand some specific domain.

While DM may utilize machine learning techniques, it may also drive the advancement of ML techniques/algorithms (P. Anantharam).
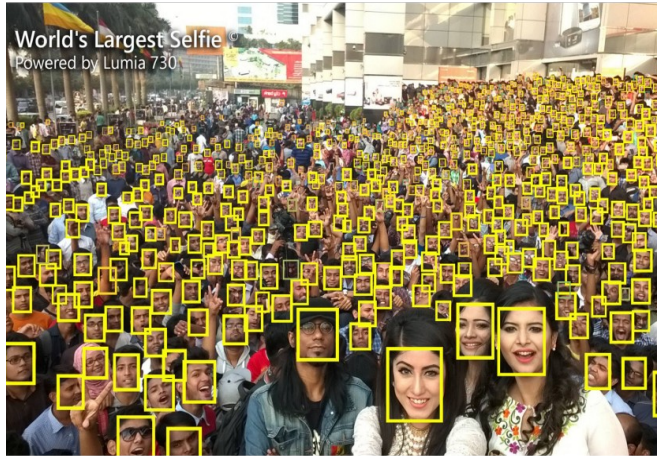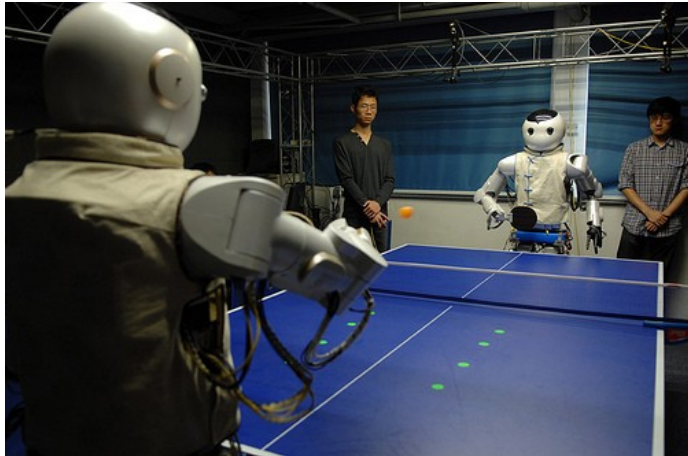


Artificial Intelligence

Machine Learning

Data Mining

Data Science

Deep Learning

Data Analytics

Big Data

Master Universitario Oficial **Data Science**
UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**   **TERMINOLOGY & TRENDS**   4

Big data
Machine learning

Deep learning
Data science
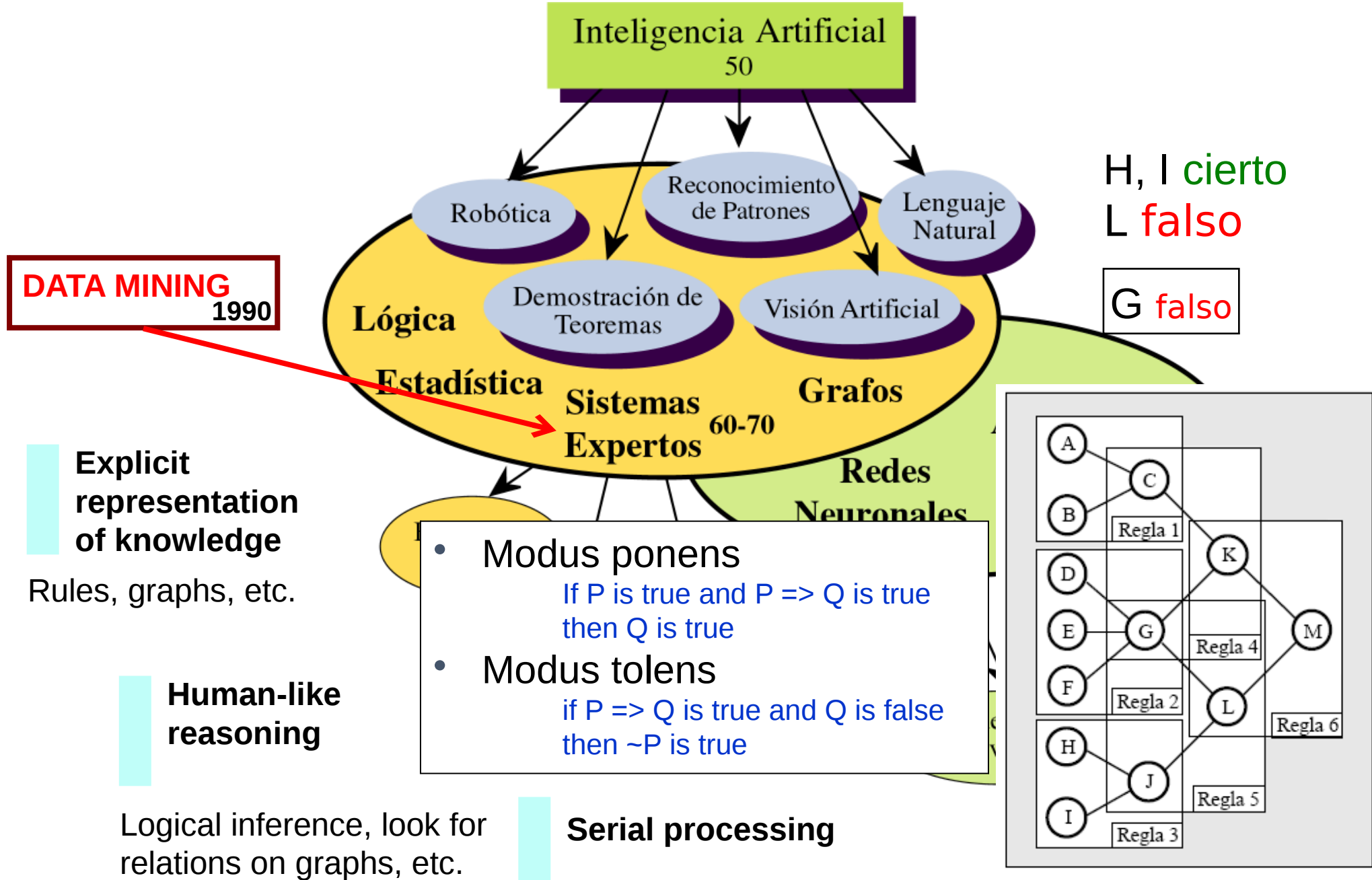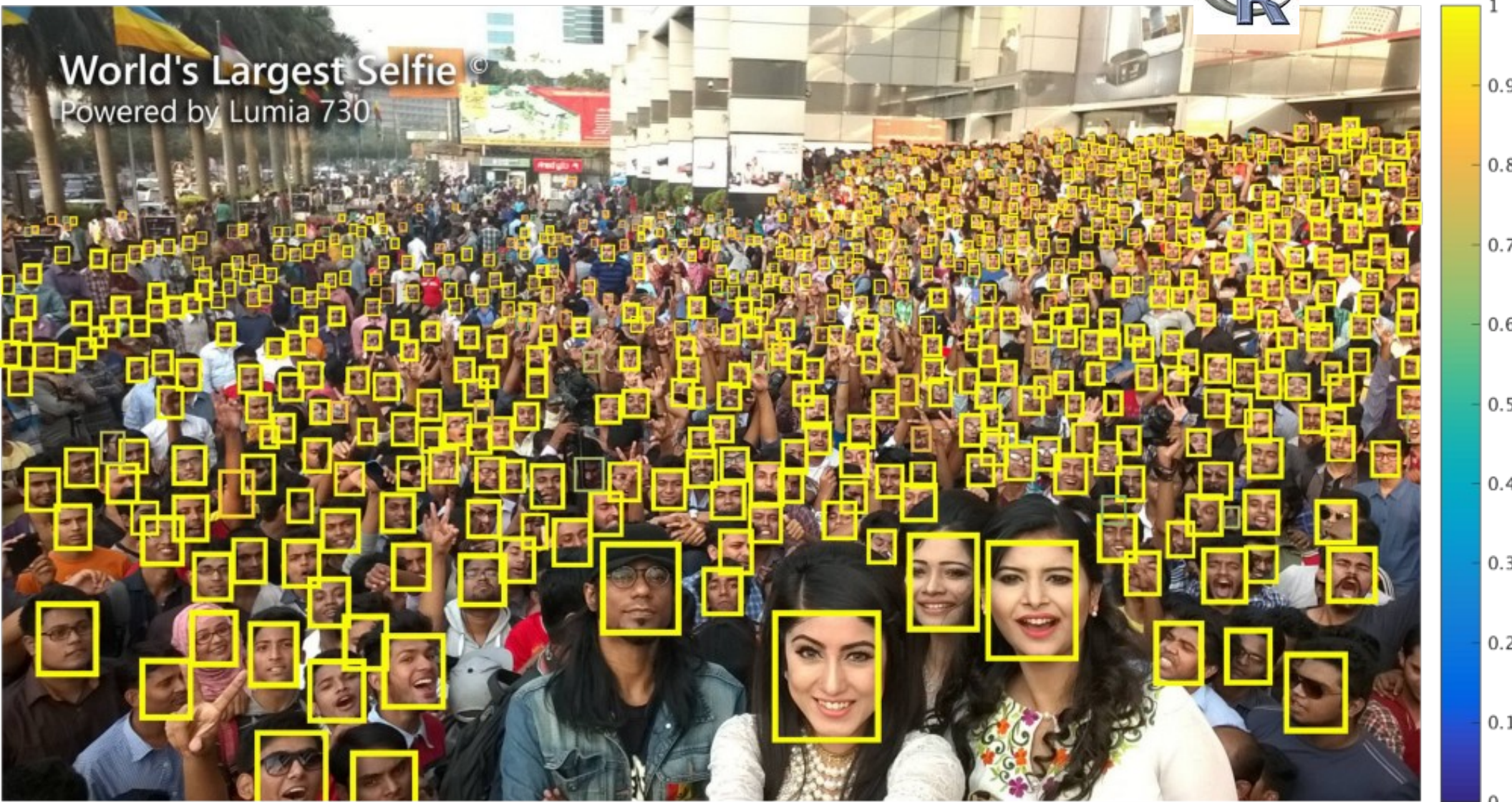
Data mining
Artificial intelligence

```
library(gtrendsR)
library(reshape2)
google.trends = gtrends(c("big data"), gprop = "web", time = "all")[[1]]
google.trends = dcast(google.trends, date ~ keyword + geo, value.var = "hits")
rownames(google.trends) = google.trends$date
plot(google.trends, type = "l")
google.trends = gtrends(c("machine learning"), gprop = "web", time = "all")[[1]]
google.trends = dcast(google.trends, date ~ keyword + geo, value.var = "hits")
rownames(google.trends) = google.trends$date
lines(google.trends, col = "blue")
                    ## Reproducir la figura anterior:
```

https://www.displayr.com/extracting-google-trends-data-in-r

Master Universitario Oficial **Data Science**
UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

INTRO:          **TERMINOLOGY & TRENDS**          5

Inteligencia Artificial 50
- Robótica
- Reconocimiento de Patrones
- Lenguaje Natural
- Demostración de Teoremas
- Visión Artificial
- Lógica
- Estadística
- Sistemas Expertos 60-70
- Grafos

World's Largest Selfie
Powered by Lumia 730

1996
Games
Kasparov
Deep Blue
Subscribe

Master Universitario Oficial **Data Science**
UC UNIVERSIDAD DE CANTABRIA
UIMP Universidad Internacional Menéndez Pelayo
con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:** **ARTIFICIAL INTELLIGENCE** 6

# Inteligencia Artificial
## 50

Robótica

Reconocimiento de Patrones

Lenguaje Natural

Demostración de Teoremas

Visión Artificial

**H, I** cierto
**L** falso

**G** falso

**DATA MINING**
**1990**

Lógica

Estadística

**Sistemas Expertos** 60-70

**Grafos**

**Redes Neuronales**

**Explicit representation of knowledge**

Rules, graphs, etc.

- Modus ponens
  If P is true and P => Q is true then Q is true
- Modus tolens
  if P => Q is true and Q is false then ~P is true

**Human-like reasoning**

Logical inference, look for relations on graphs, etc.

**Serial processing**

A
B
C
Regla 1
D
E
G
F
Regla 2
H
J
I
Regla 3
K
M
Regla 4
L
Regla 5
Regla 6

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA
UIMP Universidad Internacional Menéndez Pelayo
CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**     **ML Techniques**     7

World's Largest Selfie ©
Powered by Lumia 730

...evelop a face detector (Tiny Face Detector) that can find ~800 faces out of ~1000 reportedly preser...
...aking use of novel characterization of scale, resolution, and context to find small objects.

Master Universitario Oficial **Data Science**
con el apoyo del

UC — Universidad de Cantabria
UIMP — Universidad Internacional Menéndez Pelayo
CSIC — Consejo Superior de Investigaciones Científicas

INTRO:     AN EXAMPLE IN R     8

# Using Machine Learning to Explore Neural Network Architecture

Wednesday, May 17, 2017

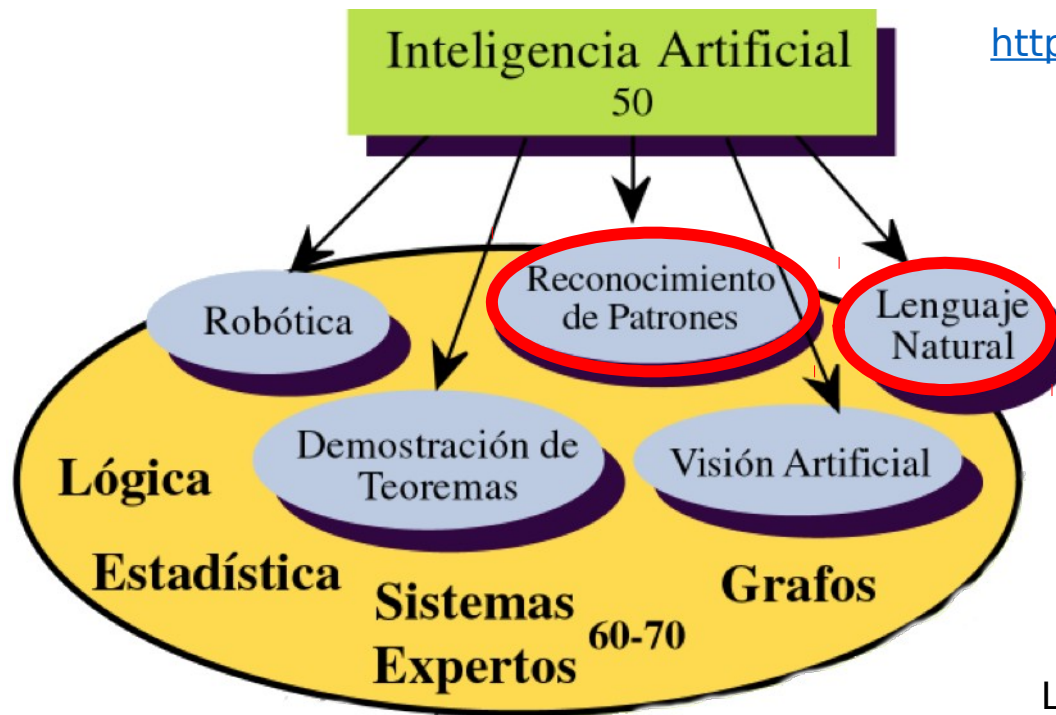Posted by Quoc Le & Barret Zoph, Research Scientists, Google Brain team

At Google, we have successfully applied deep learning models to many applications, from image recognition to speech recognition to machine translation. Typically, our machine learning models are painstakingly designed by a team of engineers and scientists. This process of manually designing machine learning models is difficult because the search space of all possible models can be combinatorially large — a typical 10-layer network can have ~$10^{10}$ candidate networks! For this reason, the process of designing networks often takes a significant amount of time and experimentation by those with significant machine learning expertise.



Un ejemplo de cómo identifica imágenes NASNet (Google Research)

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA
UIMP Universidad Internacional Menéndez Pelayo
CSIC Consejo Superior de Investigaciones Científicas

**INTRO:** **ImageNet and AutoML** 9

(a)


(b)

Overview of Natural Language Processing(NLP) with R and OpenNLP

## Inteligencia Artificial
### 50

Robótica

Reconocimiento de Patrones

Lenguaje Natural

Demostración de Teoremas

Visión Artificial

Lógica

Estadística

Sistemas Expertos 60-70

Grafos

60000+10000 images 28x28
Labeled as {0,...,9}



Lineal: 10%. k-NN: 3%. SVM: 1%. Deep: 0.3%




World's Largest Selfie
Powered by Lumia 730


1996
Games
Kasparov
Deep Blue
Subscribe

Master Universitario Oficial **Data Science**

UC UNIVERSIDAD DE CANTABRIA  UIMP Universidad Internacional Menéndez Pelayo  con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**    **ARTIFICIAL INTELLIGENCE**    10

ImageNet is an image database organized according to the (nouns of the) WordNet hierarchy, in which each node of the hierarchy is depicted by an average of over five hundred images.
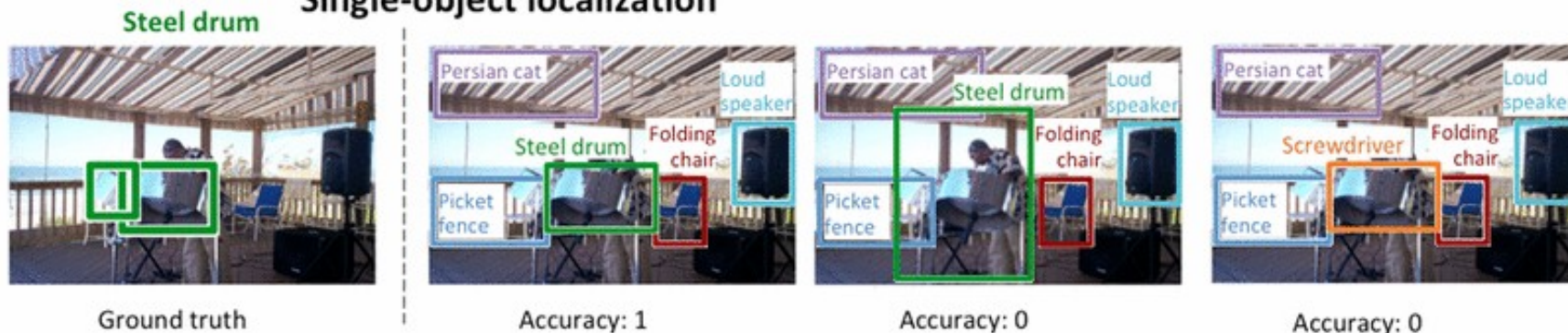
#synsets: 21841
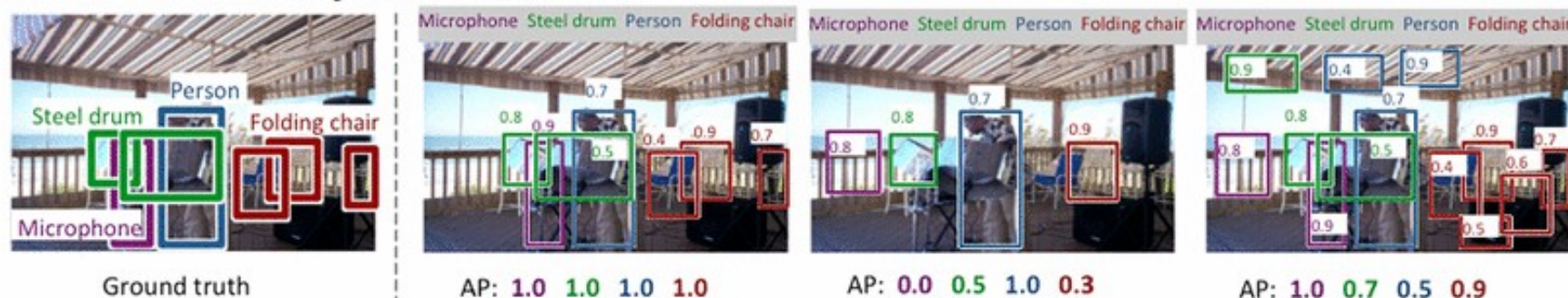#images: 14197122

150 GB   [kaggle]



David G. Lowe, **Distinctive Image Features from Scale-Invariant Keypoints.** *International Journal of Computer Vision, 2004.*

**Single-object localization**

Steel drum — Ground truth | Accuracy: 1 | Accuracy: 0 | Accuracy: 0

**Object detection**

Ground truth — AP: 1.0 1.0 1.0 1.0 | AP: 0.0 0.5 1.0 0.3 | AP: 1.0 0.7 0.5 0.9

Validation: top-5 error rate

2017 video included

Inception-v3: 3.46% top-5 and 17.3% top-1 (25 million parameters). [Inception In kaggle]

O. Russakovsk (2015) ImageNet Large Scale Visual Recognition Challenge, International Journal of Computer Vision, 115, 211–252

Master Universitario Oficial **Data Science**

UC UNIVERSIDAD DE CANTABRIA | UIMP Universidad Internacional Menéndez Pelayo | con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**    **ILSVRC Challenge**    12

**Nuevos DATA-Paradigmas driven**

**Statistical Inspiration**

**STATISTICAL LEARNING** 2000

**DEEP LEARNING** 2010

**Inspiración Biológica**

**Data driven using abstract representations**

Kernels, neural network, etc.

**Optimization-based reasoning (error function).**

Empirical risk, gradient descend, etc.

**Parallel processing**

HPC, GPUs, cloud, etc.

Inteligencia Artificial 50

Robótica

Reconocimiento de Patrones

Lenguaje Natural

Demostración de Teoremas

Visión Artificial

Lógica

Estadística

Sistemas Expertos 60-70

Grafos

Algoritmos Evolutivos 80

Redes Neuronales 80

Basados en Reglas

Basados en Probabilidad

Fuzzy Sets

Perceptrones Multicapa

Algoritmos Genéticos

Redes no Supervisadas

Estrategias Evolutivas

the non trivial extraction of implicit, previously unknown, and potentially useful information from data

*W. Frawley and G. **Piatetsky-Shapiro** and C. Matheus, Knowledge Discovery in Databases: An Overview.*
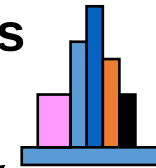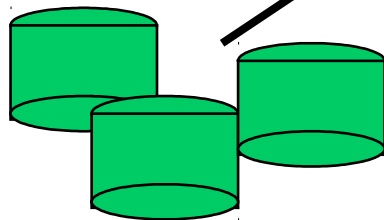
*AI Magazine, Fall **1992**, 213-228.*

**Knowledge Patterns**

**AI & Machine learning**

**DATA MINING**

**Task-relevant Data**

**Data Warehouse**

**Data Cleaning**

**Databases**

*Data Mining: Practical Machine Learning Tools and Techniques with Java Impleméntatenos*

*Ian H. Witten, Eibe Frank (**1999**)*

**Machine Learning and Data Mining Open Soure Tools in Java**

`http://www.cs.waikato.ac.nz/~ml/weka/`

WEKA The University of Waikato

KDnuggets

the non trivial extraction of implicit, previously unknown, and potentially useful information from data

*W. Frawley and G.* **Piatetsky-Shapiro** *and C. Matheus, Knowledge Discovery in Databases: An Overview.*
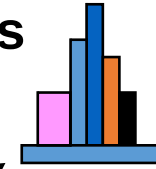
*AI Magazine, Fall* **1992***, 213-228.*

**Knowledge Patterns**

**AI & Machine learning**

**DATA MINING**

**Task-relevant Data**

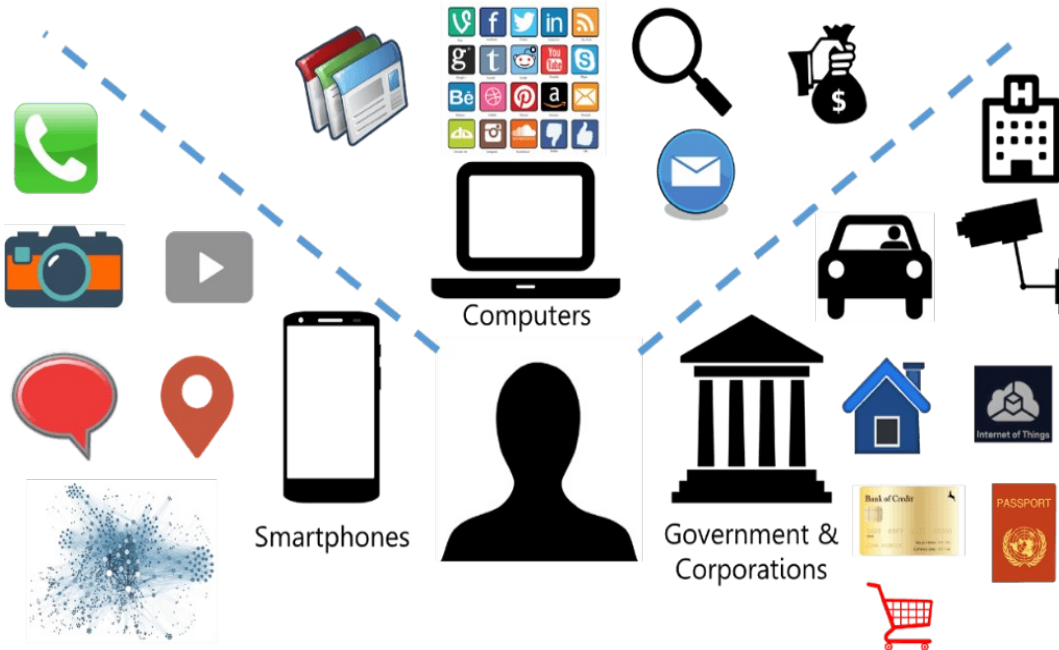**Data Warehouse**

**Data Cleaning**

**Databases**

The essence of machine learning:

- A pattern exists.

- We cannot pin it down mathematically.

- We have data on it.

the non trivial extraction of implicit, previously unknown, and potentially useful information from data

*S. Bryson et al., Visually exploring gigabyte data sets in real time.*
*Communications of the ACM, 42, 82-90,*
*Aug.* **1999**

**Knowledge Patterns**

**AI & Machine learning**

**DATA ANALYTICS**

**Data Discovery, Cleaning and Reduction**

**Modeling and Prediction**

Volume
Data Size
Data Complexity
Speed of Change
Velocity
Data Sources
Variety

**Big data**
**(integration of heterogeneous real-time sources)**

Computers
Smartphones
Government & Corporations

**2011**

WINTER 2011    VOL.52 NO.2
**MITSloan**
Management Review

Steve LaValle, Eric Lesser, Rebecca Shockley, Michael S. Hopkins and Nina Kruschwitz

Big Data, Analytics and the Path From Insights to Value

**2014**

Raghupathi and Raghupathi *Health Information Science and Systems* 2014, 2:3
http://www.hissjournal.com/content/2/1/3

HEALTH INFORMATION SCIENCE AND SYSTEMS

**REVIEW**                                    **Open Access**

Big data analytics in healthcare: promise and potential

Wullianallur Raghupathi[1*] and Viju Raghupathi[2]

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA    UIMP Universidad Internacional Menéndez Pelayo    CSIC CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

**INTRO:**    **BIG DATA & DATA ANALYTICS**    16
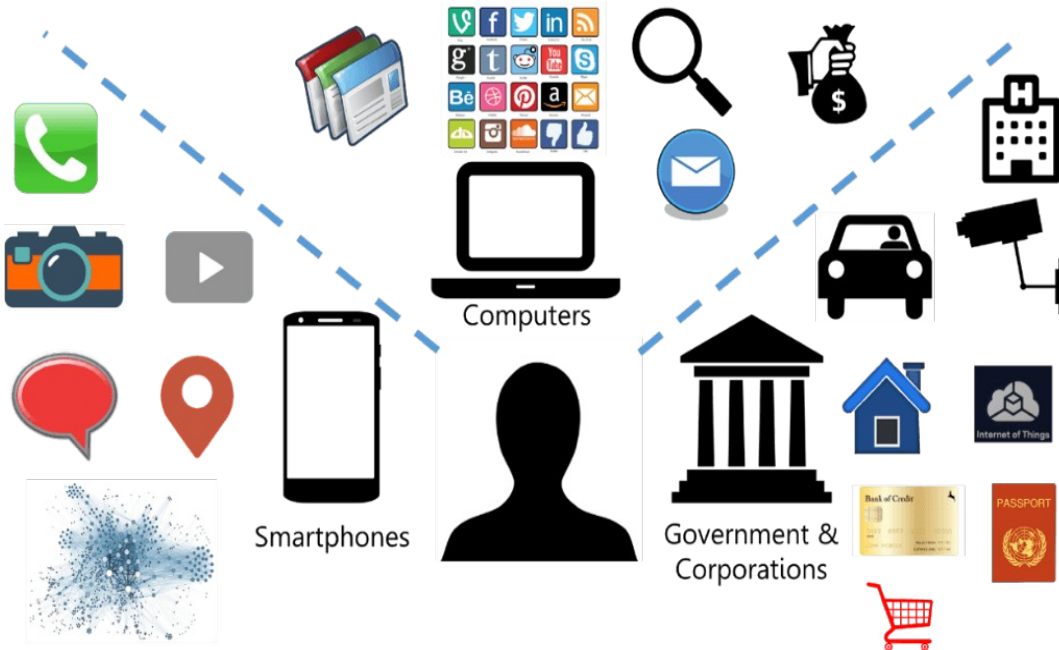
the non trivial extraction of implicit, previously unknown, and potentially useful information from data

*S. Bryson et al., Visually exploring gigabyte data sets in real time.*
*Communications of the ACM, 42, 82-90,*
*Aug.* **1999**

**Knowledge Patterns**

**AI & Machine learning**

**DATA ANALYTICS**

**Data Discovery, Cleaning and Reduction**

**Modeling and Prediction**

Volume
Data Size

Data Complexity

Speed of Change

**Velocity**

Data Sources

**Variety**

**Big data**
**(integration of heterogeneous real-time sources)**

Computers

Smartphones

Government & Corporations

Bank of Credit

PASSPORT

Internet of Things

Domain-specific Problems

**DataLabs**

Weather

Crime reports

Transportation

Social media

Linked open urban data

Energy

Real estate

Maps

Policies

Master Universitario Oficial **Data Science**

UC UNIVERSIDAD DE CANTABRIA    UIMP Universidad Internacional Menéndez Pelayo    con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

INTRO:    **BIG DATA & DATA ANALYTICS**    17

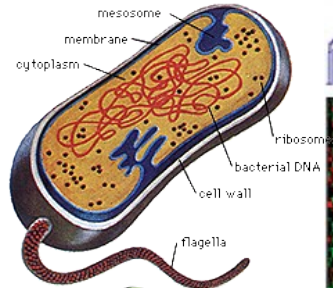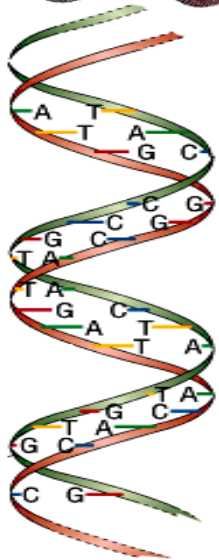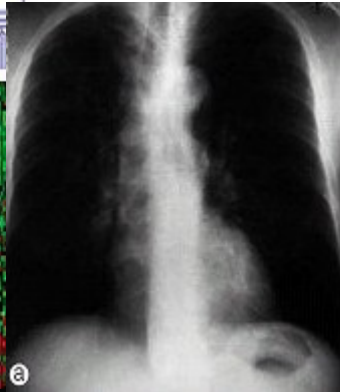| Financiero Seguros | Comercio y marketing | Industria y gestión empresarial | Tecnologías información y comunicaciones | Sanitario y farmacéutico | Meteorología, clima y medio ambiente |
|---|---|---|---|---|---|



El SNS genera **5 millones** de altas hospitalarias al año almacenando datos sobre diagnósticos y procedimientos asociados a cada paciente que contienen información necesaria para la **gestión** del SNS.

http://icmbd.es/

Master Universitario Oficial **Data Science** con el apoyo del

UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   CSIC CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

**INTRO:**   **BIG DATA in BIO AND HEALTH**   18

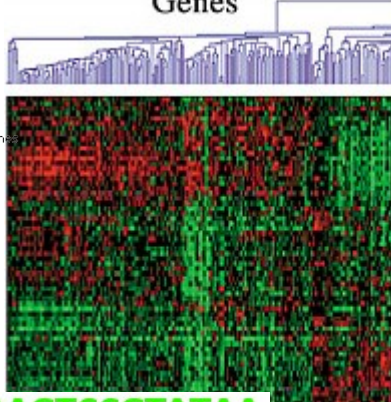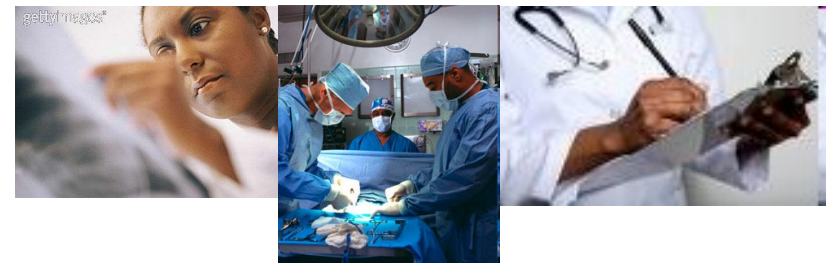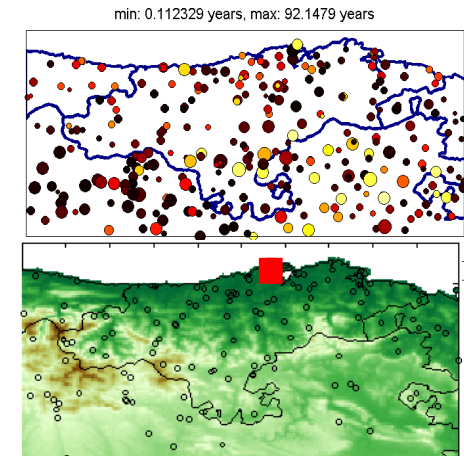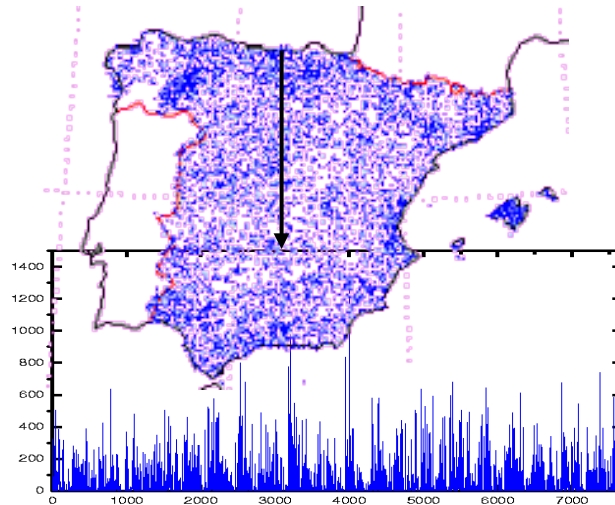| Financiero Seguros | Comercio y marketing | Industria y gestión empresarial | Tecnologías información y comunicaciones | Sanitario y farmacéutico | Meteorología, clima y medio ambiente |
|---|---|---|---|---|---|

Las observaciones y simulaciones globales y regionales del clima generan **cientos de TB** de información heterogénea necesaria para la predicción meteorológica.

```
86010150000000001086010150000000010
86010210000000001086010210000000010
860103                          0000010
860104                          0000010
860105                          0000010
860106                          0000010
860107                          0000010
860108                          0000010
860109                          0001000
860110                          1001100
86011100100000000086011100100000000
```

min: 0.112329 years, max: 92.1479 years

http://www.meteo.unican.es/downscaling

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   CSIC CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

**INTRO:**     **BIG DATA in METEROLOGY**     19

Startups Using Big Data

A key factor for the quick growth of data science is the efficient frameworks (and infrastructures) available:



https://www.kdnuggets.com/2017/09/datacamp-keras-cheat-sheet-deep-learning-python.html
https://project.inria.fr/deeplearning/files/2016/05/DLFrameworks.pdf
https://www.kdnuggets.com/2018/09/deep-learning-framework-power-scores-2018.html
https://github.com/amueller/scipy_2015_sklearn_tutorial

Master Universitario Oficial **Data Science**
UC UNIVERSIDAD DE CANTABRIA    UIMP Universidad Internacional Menéndez Pelayo    con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**    **SOFTWARE & FRAMEWORKS**    21

- Clearly defined business problem
- Set success criteria
- Define clear data science objectives

- Understand data points and constraints
- Formulate data analytics strategy
- Perform required transformation

- Experiment with multiple models
- Choose the most optimal model
- Create a feedback loop

**Define Business Problem** → **Map to Machine Learning Problem** → **Data Preparation** → **Exploratory Data Analysis** → **Modeling** → **Evaluation**

- Break business problems to data science problems
- Identify Machine Learning Problem categories

- Perform statistical and visual analysis
- Discover and handle outliers/errors
- Shortlist predictive modeling techniques

**80% of work**          **20% of work**

Master Universitario Oficial **Data Science**
UC UNIVERSIDAD DE CANTABRIA    UIMP Universidad Internacional Menéndez Pelayo    con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**          **DATA ANALYTICS**          22

- Clearly defined business problem
- Set success criteria
- Define clear data science objectives

- Understand data points and constraints
- Formulate data analytics strategy
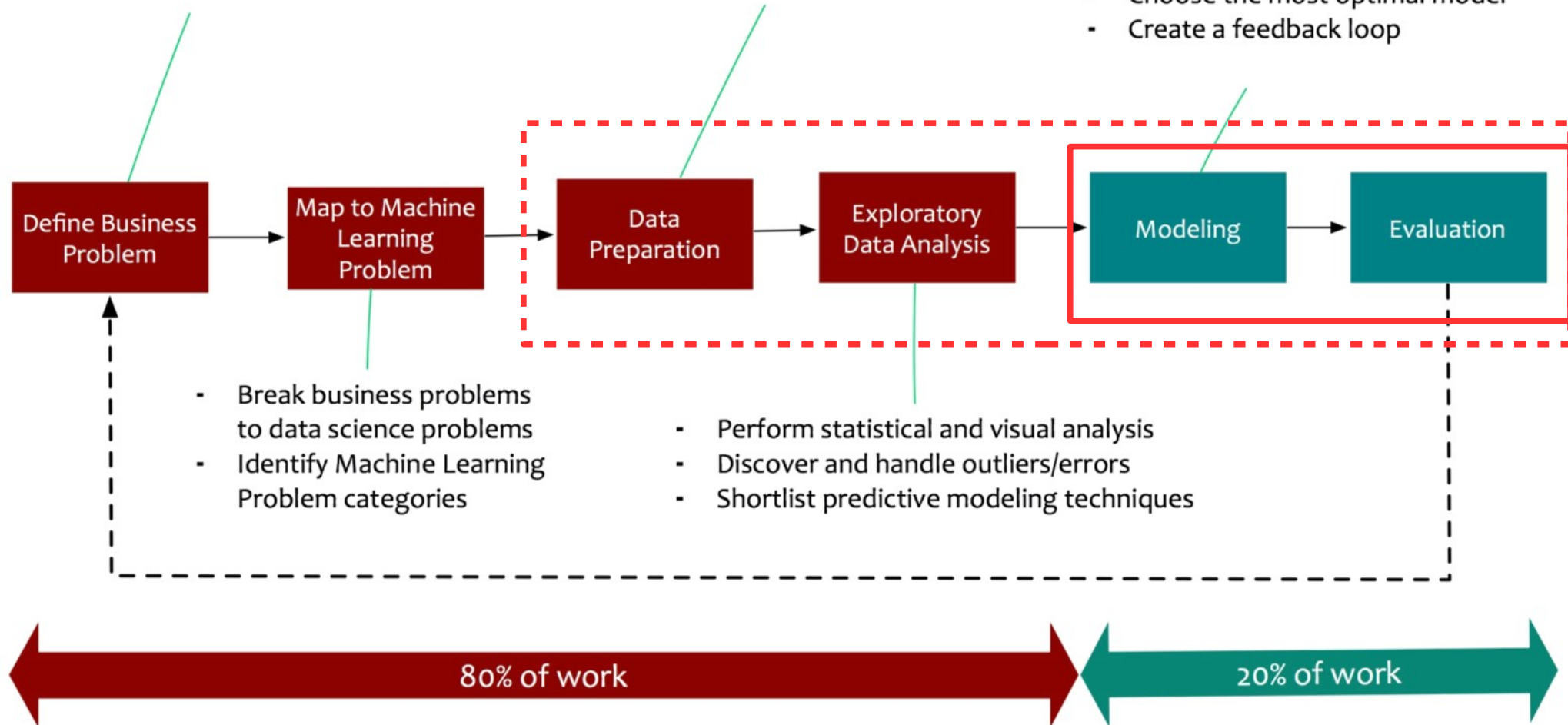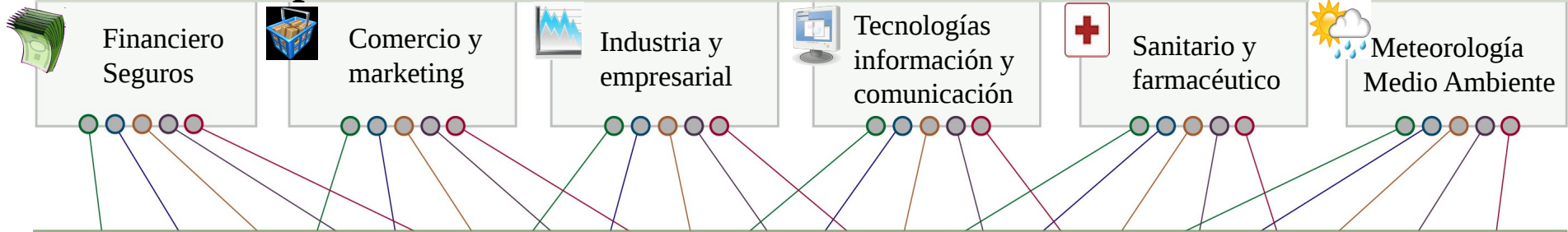- Perform required transformation

- Experiment with multiple models
- Choose the most optimal model
- Create a feedback loop

```
Define Business Problem → Map to Machine Learning Problem → Data Preparation → Exploratory Data Analysis → Modeling → Evaluation
```

- Break business problems to data science problems
- Identify Machine Learning Problem categories

- Perform statistical and visual analysis
- Discover and handle outliers/errors
- Shortlist predictive modeling techniques

80% of work

20% of work

Master Universitario Oficial **Data Science**
UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   con el apoyo del CSIC Consejo Superior de Investigaciones Científicas

**INTRO:** | **DATA ANALYTICS** | 23

# Sectores de aplicación

Financiero Seguros

Comercio y marketing

Industria y empresarial

Tecnologías información y comunicación

Sanitario y farmacéutico

Meteorología Medio Ambiente

## Proceso de Minería de Datos

Data Selection and Cleaning

Data Transformation
feature extraction

Data Modeling

Evaluation / Deployment

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA
UIMP Universidad Internacional Menéndez Pelayo
CSIC CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

**INTRO:** | **Data Mining: Transformation** | 24

**Sectores de aplicación**

| | | | | | |
|---|---|---|---|---|---|
| Financiero Seguros | Comercio y marketing | Industria y empresarial | Tecnologías información y comunicación | Sanitario y farmacéutico | Meteorología Medio Ambiente |

**Proceso de Minería de Datos**

Data Selection and Cleaning → Data Transformation *feature extraction* → Data Modeling → Evaluation / Deployment

**Problemas habituales**

| Descripción y visualización | Asociación | Segmentación | Clasificación | Predicción |
|---|---|---|---|---|

Machine learning develop methods for data modelling and prognosis.

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**   **Canonical Problems**   25

# Sectores de aplicación

| Financiero Seguros | Comercio y marketing | Industria y empresarial | Tecnologías información y comunicación | Sanitario y farmacéutico | Meteorología Medio Ambiente |
|---|---|---|---|---|---|

## Proceso de Minería de Datos

Data Selection and Cleaning  →  **Data Transformation feature extraction**  →  Data Modeling  →  Evaluation / Deployment

**WIREs DATA MINING AND KNOWLEDGE DISCOVERY** — WILEY
Explore this journal >

Advanced Review

## Data discretization: taxonomy and big data challenge

Knowledge-Based Systems
Volume 86, September 2015, Pages 33-45
ELSEVIER

## Recent advances and emerging challenges of feature selection in the context of big data
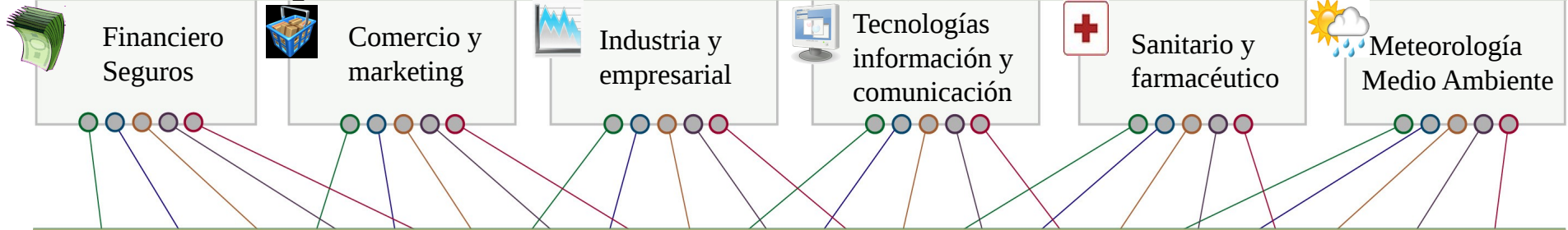
V. Bolón-Canedo, N. Sánchez-Maroño, A. Alonso-Betanzos

http://onlinelibrary.wiley.com/doi/10.1002/widm.1173/full

Master Universitario Oficial **Data Science** con el apoyo del
UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   CSIC

**INTRO:**   **Data Mining: Transformation**   26

# Sectores de aplicación

| | | | | | |
|---|---|---|---|---|---|
| Financiero Seguros | Comercio y marketing | Industria y empresarial | Tecnologías información y comunicación | Sanitario y farmacéutico | Meteorología Medio Ambiente |

## Proceso de Minería de Datos

Data Selection and Cleaning → Data Transformation *feature extraction* → **Data Modeling** → Evaluation / Deployment

## Simple Neural Network

## Deep Learning Neural Network

● Input Layer  ● Hidden Layer  ● Output Layer

$x_1$ $w_1$ → Σ → y  numeric

$x_2$ $w_2$

numeric or binary

$y = w_0 + w_1x_1 + w_2x_2$

$y = f(X,W) = X^T.W$

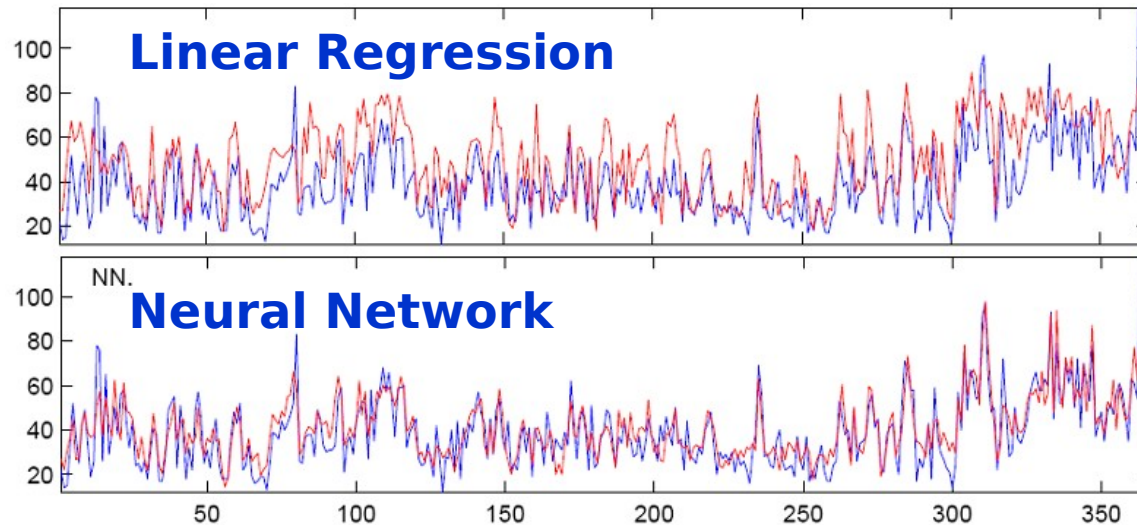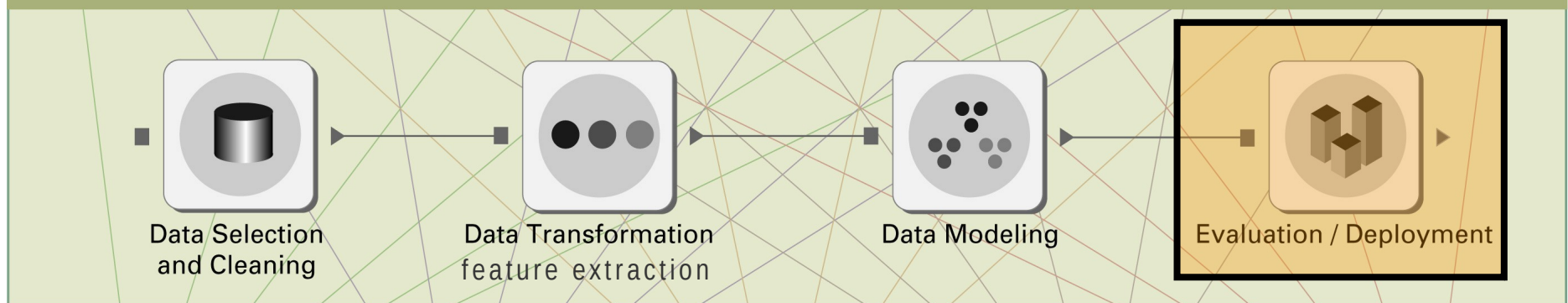**REGRESSION**

$$W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

Master Universitario Oficial **Data Science** con el apoyo del

UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**   **Data Mining: Data Modeling**   27

# Sectores de aplicación

| Financiero Seguros | Comercio y marketing | Industria y empresarial | Tecnologías información y comunicación | Sanitario y farmacéutico | Meteorología Medio Ambiente |

# Proceso de Minería de Datos

**Data Selection and Cleaning**

**Data Transformation** *feature extraction*

**Data Modeling**

**Evaluation / Deployment**

**Linear Regression**

**Wind Speed**

NN.

**Neural Network**

Master Universitario Oficial **Data Science**

UC UNIVERSIDAD DE CANTABRIA   UIMP Universidad Internacional Menéndez Pelayo   *con el apoyo del* CSIC CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

INTRO:     **Data Mining: Evaluation**     28

# Sectores de aplicación

| Financiero Seguros | Comercio y marketing | Industria y empresarial | Tecnologías información y comunicación | Sanitario y farmacéutico | Meteorología Medio Ambiente |

## Proceso de Minería de Datos

| Data Selection and Cleaning | Data Transformation *feature extraction* | Data Modeling | Evaluation / Deployment |



**K-FOLD STRATEGY**

TRAIN — TEST

**1** Set aside the test set and split the train set into k folds

FOLD 1  FOLD 2  FOLD 3  ...  FOLD K

**2** For each parameter combination

FOLD 1  OTHER FOLDS

Parameter (e.g., depth) A   5  15   Parameter B (e.g., n trees)
6  16

OTHER FOLDS  FOLD K

Compute metric  Average

METRIC 1
METRIC K

**3** Choose the parameter combinaison with the best metrics

A  6  14  B

Retrain model on all training data   Compute metric on test set

**HOLDOUT STRATEGY**

TRAIN  VALIDATION  TEST

**1** Split your data into train / validation / test

**2** For each parameter combination

TRAIN A MODEL   COMPUTE METRIC ON VALIDATION SET

Parameter (e.g., depth) A   5  15   Parameter B (e.g., n trees)
6  16

VALIDATION METRIC

**3** Choose the parameter combination with the best metric

A  6  14  B

Retrain model on all training data   Compute metric on test set

TEST METRIC (can compare with other models)

## Sectores de aplicación

| | | | | | |
|---|---|---|---|---|---|
| Financiero Seguros | Comercio y marketing | Industria y empresarial | Tecnologías información y comunicación | Sanitario y farmacéutico | **Meteorología** Medio Ambiente |

## Proceso de Minería de Datos



Data Selection and Cleaning → Data Transformation *feature extraction* → Data Modeling → Evaluation / Deployment

## Problemas habituales

| Descripción y visualización | Asociación | Segmentación | Clasificación | Predicción |
|---|---|---|---|---|

Machine learning develop methods for data modelling and prognosis.

Master Universitario Oficial **Data Science**
con el apoyo del
UC UNIVERSIDAD DE CANTABRIA  UIMP Universidad Internacional Menéndez Pelayo  CSIC Consejo Superior de Investigaciones Científicas

INTRO: | DATA ANALYTICS | 30

An Introduction to Statistical Learning: With Applications in R

James, G., Witten, D., Hastie, T., Tibshirani, R.

Springer (2013)

http://www-bcf.usc.edu/~gareth/ISL

[PDF]

require(ISLR)

The Elements of Statistical Learning

Trevor Hastie, Robert Tibshirani, Jerome Friedman

Springer (2nd ed. 2009, Corr. 9th printing 2017)

https://web.stanford.edu/~hastie/ElemStatLearn/

[PDF]

Gregory Piatetsky-Shapiro

https://www.kdnuggets.com

Master Universitario Oficial **Data Science**

con el apoyo del

UC UNIVERSIDAD DE CANTABRIA  UIMP Universidad Internacional Menéndez Pelayo  CSIC Consejo Superior de Investigaciones Científicas

**INTRO:**      **BIBLIOGRAPHY**      31