

Introducción al *machine learning*

MACHINE LEARNING

A scene from Toy Story featuring Woody and Buzz Lightyear. Woody, on the left, is a brown cowboy doll with a white bandana and a drawstring belt. Buzz, on the right, is a green space ranger with a purple vest and a red button on his chest. Buzz is in mid-air, performing a celebratory dance move with his arms raised and legs spread. The background shows a wooden floor and a doorway.

MACHINE LEARNING EVERYWHERE

¿Qué NO es el *machine
learning*?



No es Hal-9000

A close-up, low-angle shot of a T-800 endoskeleton's head. The metallic, segmented faceplate reflects bright light, and its glowing red eyes are prominent. The mechanical nature of the head is visible through the joints and plates.

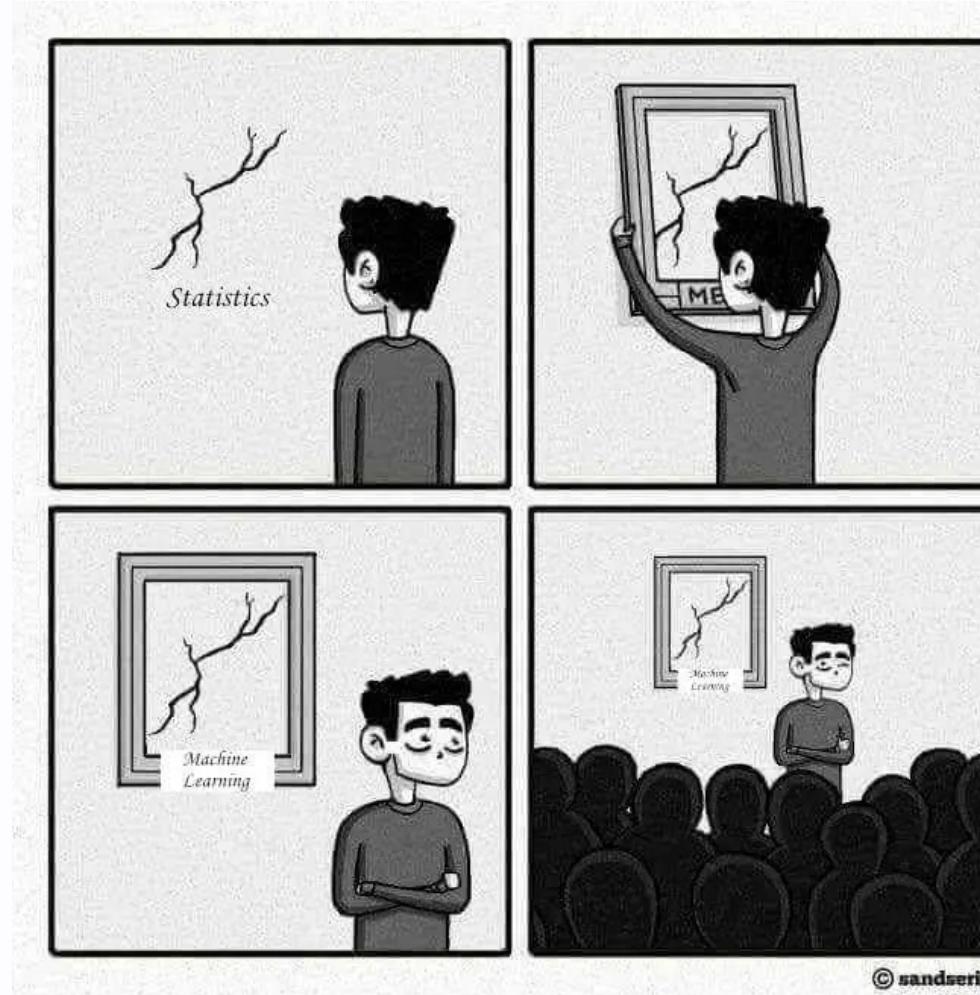
No es T-800



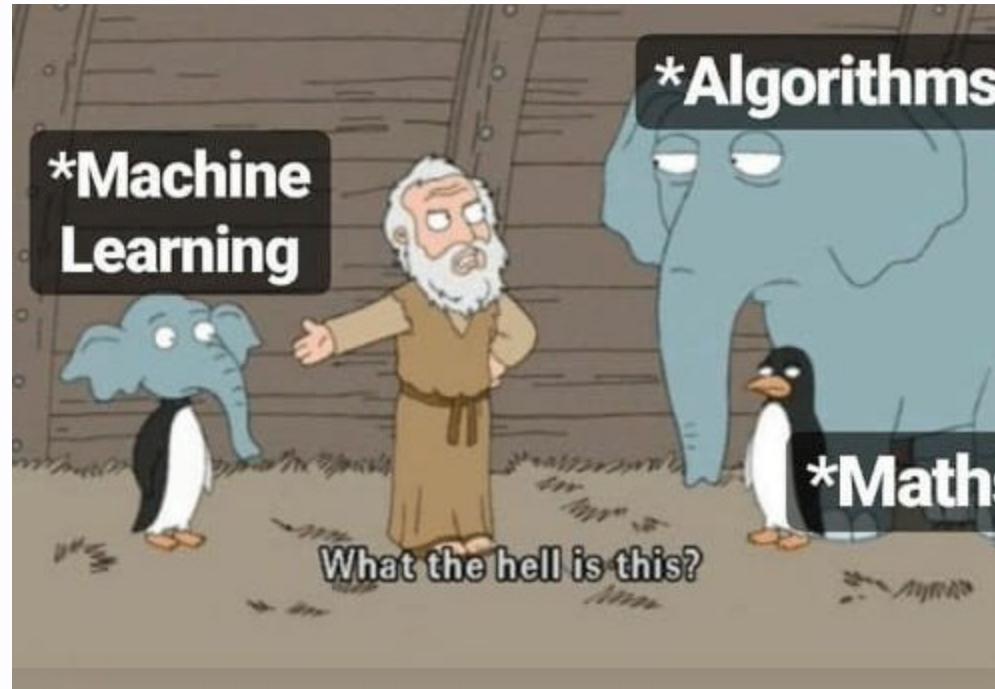
**No es una alternativa a los
seres humanos**

¿Qué es el *machine learning*?

El machine learning es estadística disfrazada



Machine learning = algorítmica + matemáticas



¿Machine learning?



Una mejor
denominación para
machine learning es
learning from data.



354 K/s 36 5:07 PM

← Tweet

 Chet Haase
@chethaase

A Machine Learning algorithm walks into a bar.
The bartender asks, "What'll you have?"
The algorithm says, "What's everyone else having?"

6:54 PM · 01 Nov 17

3,749 Retweets 7,605 Likes

1. Sin **datos** no hay **aprendizaje**
2. Sólo se **aprende** lo que está en los **datos**

**¿Por qué queremos necesitamos
el *machine learning*?**

2019 This Is What Happens In An Internet Minute



A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devices – are hard to comprehend, particularly when looked at in the context of one day

 500m

tweets are sent
every day

Twitter

294bn

billion emails are sent

320bn

emails to be sent
each day by 2021

306bn

emails to be sent
each day by 2020

3.9bn

people use emails



ACCUMULATED DIGITAL UNIVERSE OF DATA

4.4ZB

PwC

44ZB

2020

4PB

of data created by
Facebook, including

350m photos

100m hours of video
watch time

Facebook Research



4TB

of data produced by a connected car

Intel



DEMYSTIFYING DATA UNITS

From the more familiar 'bit' or 'megabyte', larger units of measurement are more frequently being used to explain the masses of data

Unit	Value	Size
b bit	0 or 1	1/8 of a byte
B byte	8 bits	1 byte
KB kilobyte	1,000 bytes	1,000 bytes
MB megabyte	1,000 ² bytes	1,000,000 bytes
GB gigabyte	1,000 ³ bytes	1,000,000,000 bytes
TB terabyte	1,000 ⁴ bytes	1,000,000,000,000 bytes
PB petabyte	1,000 ⁵ bytes	1,000,000,000,000,000 bytes
EB exabyte	1,000 ⁶ bytes	1,000,000,000,000,000,000 bytes
ZB zettabyte	1,000 ⁷ bytes	1,000,000,000,000,000,000,000 bytes
YB yottabyte	1,000 ⁸ bytes	1,000,000,000,000,000,000,000 bytes

*A lowercase "b" is used as an abbreviation for bits, while an uppercase "B" represents bytes.

463EB

of data will be created every day by 2025

IDC

95m

photos and videos are
shared on Instagram

Instagram Business

28PB

to be generated from wearable
devices by 2020

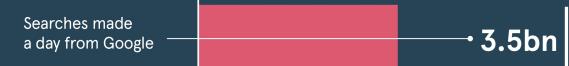
Statista



65bn

messages sent over WhatsApp and
two billion minutes of voice and
video calls made

Facebook



Smart insights



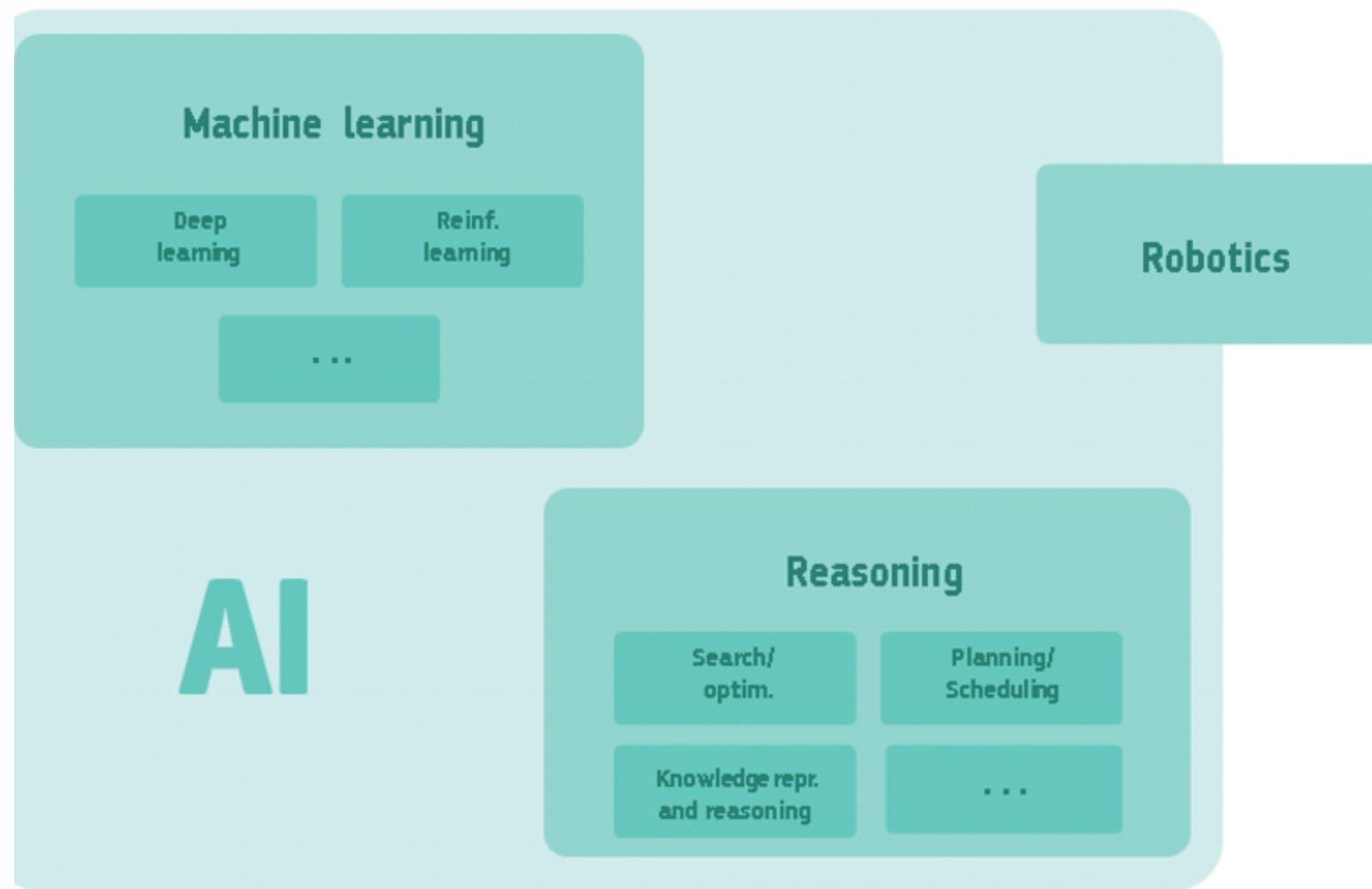
¿Por qué?

- Necesitamos obtener conocimiento de grandes volúmenes de datos
- Los ordenadores son muy buenos calculando:
 - Velocidad de ejecución
 - No incorporan sesgos
 - No se cansan, trabajando 24/7
 - Precisión en los cálculos

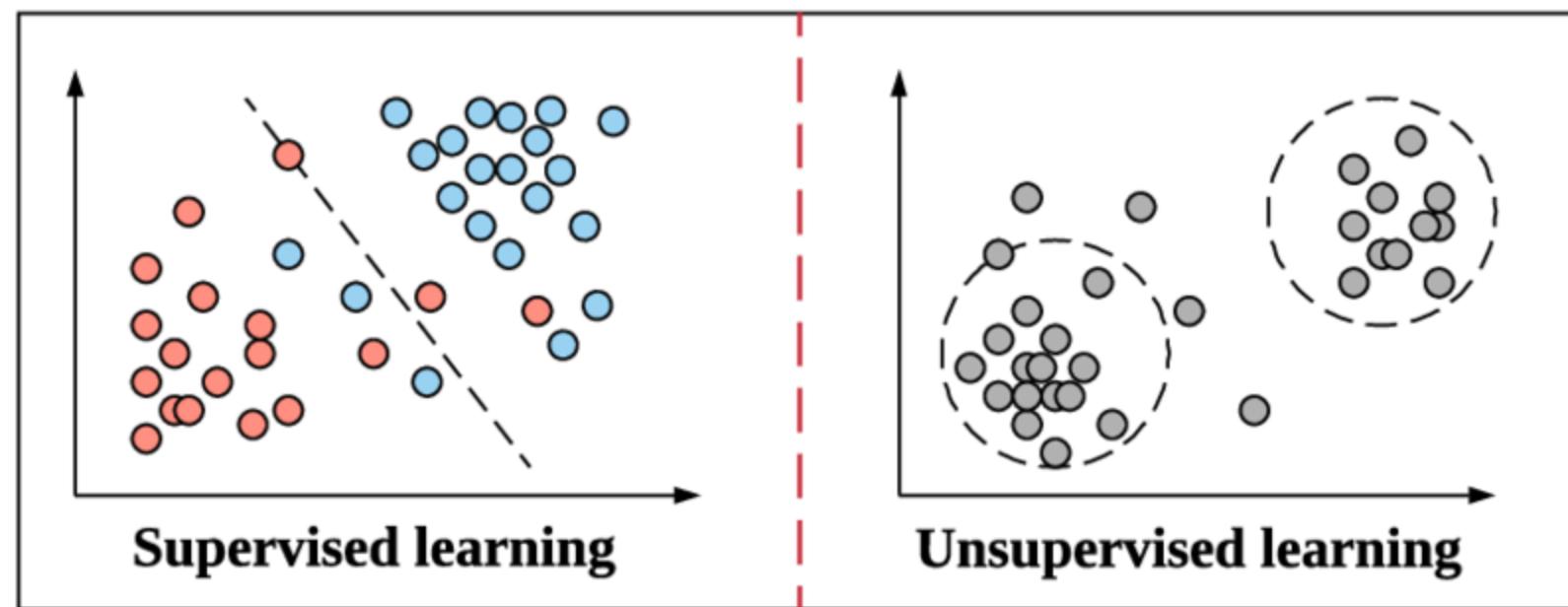
Algunas aplicaciones

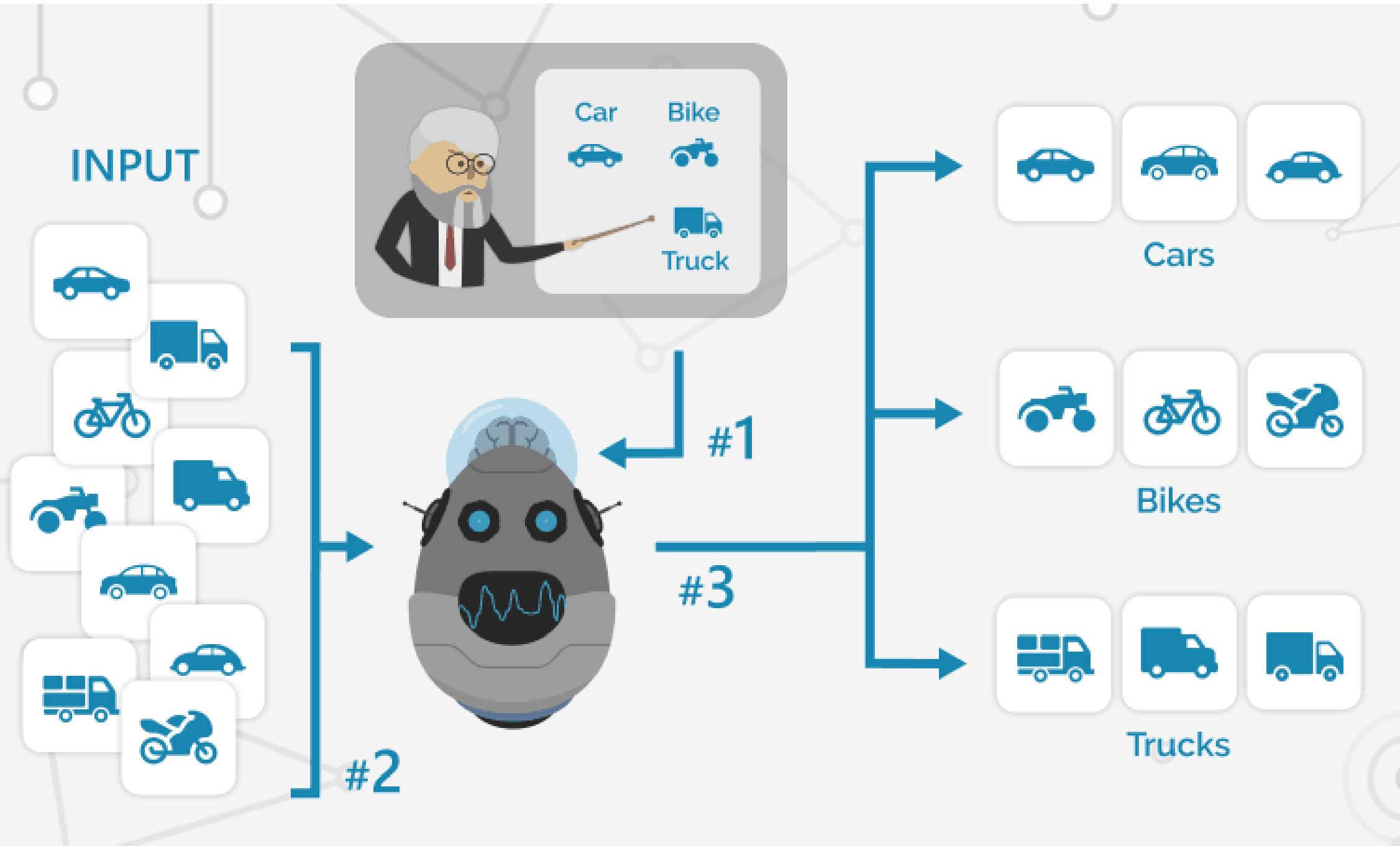
- Visión por computador
 - Conducción autónoma, detección de cáncer, ...
- Reconocimiento del habla
 - Siri, Alexa, OK Google, ...
- Procesamiento del lenguaje natural
 - Traducción automática, chat bots, ...
- Detección de patrones
 - Sistemas de recomendación, detección de fraude, análisis de comportamiento en webs, ...

El *machine learning* es una disciplina de la **inteligencia artificial**



Dos enfoques:



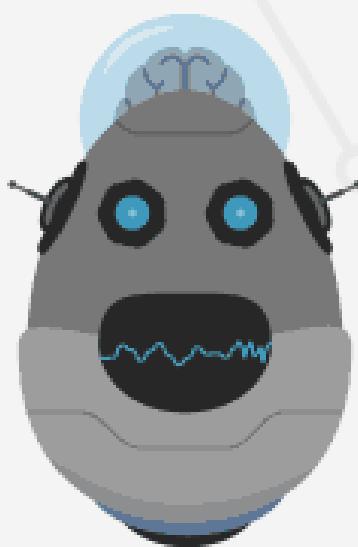


INPUT



#1

"OK, let me try..."



#2

Medium size, 4 wheels, glass...



Small size, 2 wheels...



Big size, 4+ wheels, glass...



Classical Machine Learning

Task Driven
↓
Supervised Learning
(Pre Categorized Data)

Classification
(Divide the socks by Color)
Eg. Identity Fraud Detection

Regression
(Divide the Ties by Length)
Eg. Market Forecasting

Data Driven
↓
Unsupervised Learning
(Unlabelled Data)

Clustering
(Divide by Similarity)
Eg. Targeted Marketing

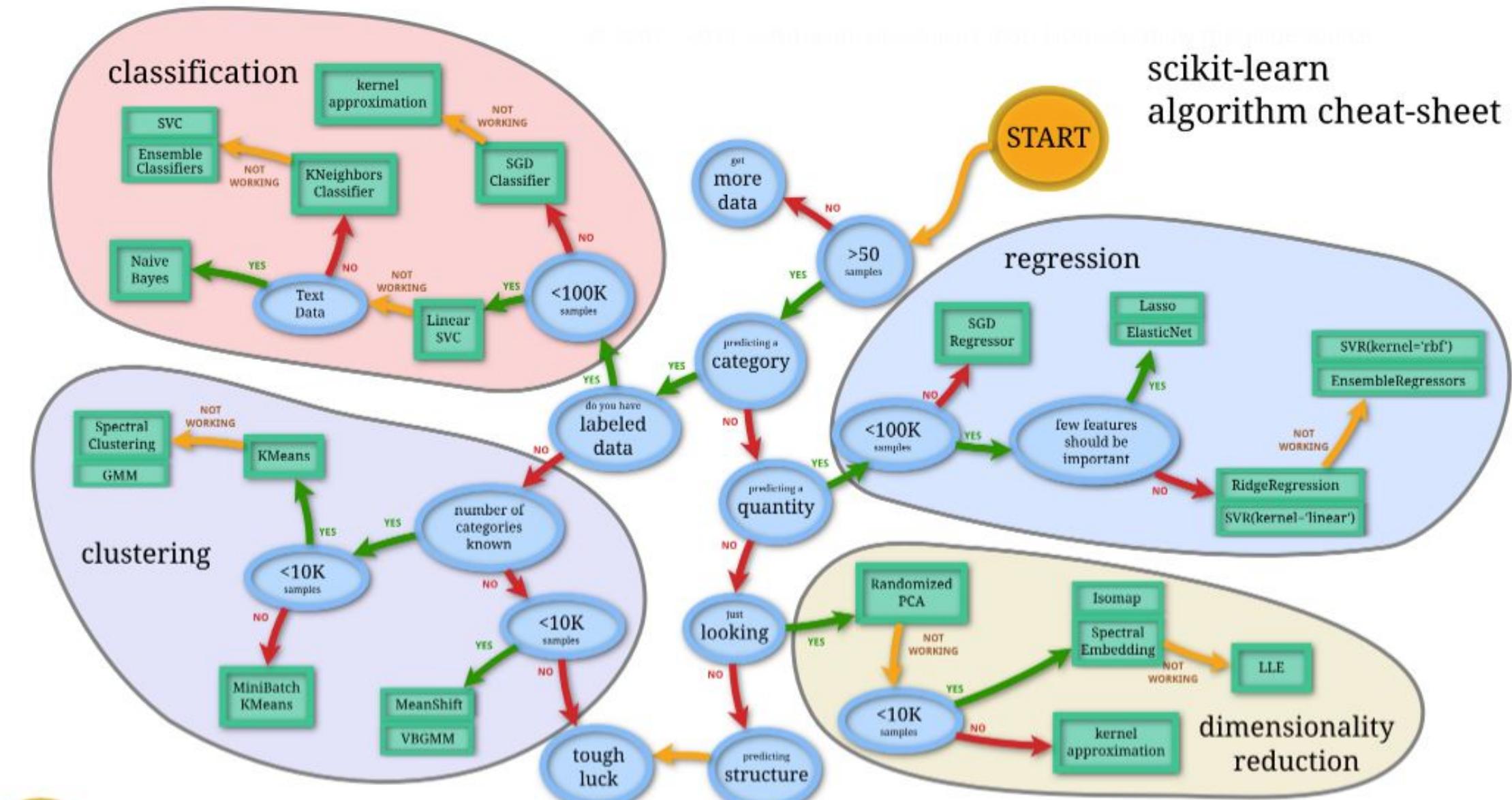
Association
(Identify Sequences)
Eg. Customer Recommendation

Dimensionality Reduction
(Wider Dependencies)
Eg. Big Data Visualization

Obj: Predictions & Predictive Models

Pattern/ Structure Recognition

scikit-learn algorithm cheat-sheet



Back

scikit
learn

Proyectos *machine learning*

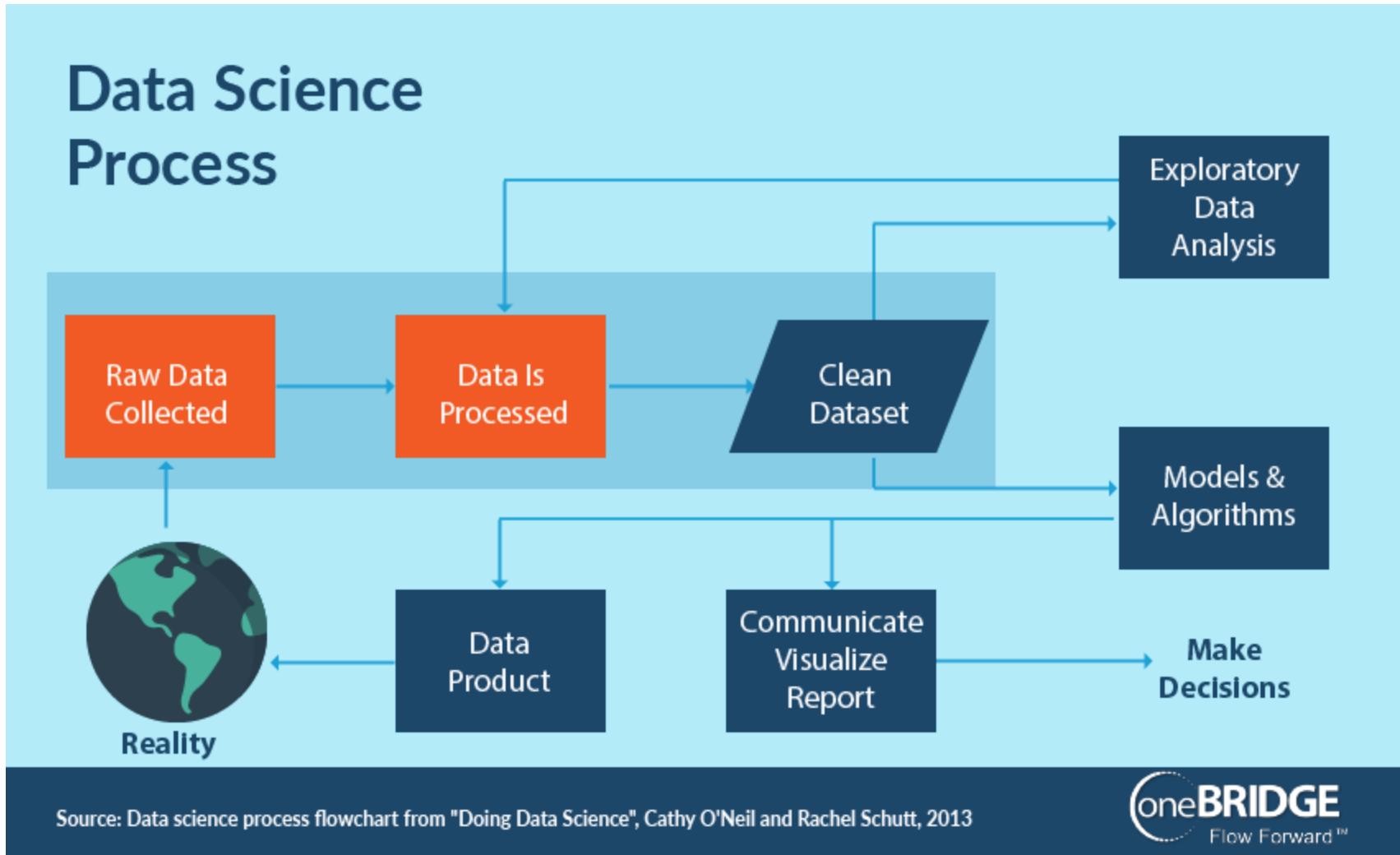
Nomenclatura de un *dataset*:

Feature 1	Feature 2	...	Feature N	Label
Observation 1,1	Observation 1,2	...	Observation 1,N	Label 1
Observation 2,1	Observation 2,2	...	Observation 2,N	Label 2
...
Observation M,1	Observation M,2	...	Observation M,N	Label M

Dataset de ejemplo:

Latitude	Longitud	Timestamp	Temperature
40,3214873	-3,8123123	1587980647	18,3
18,9230112	15,2394502	1587984567	6,4
...
-7,2321231	10,7234433	1587968742	25,7

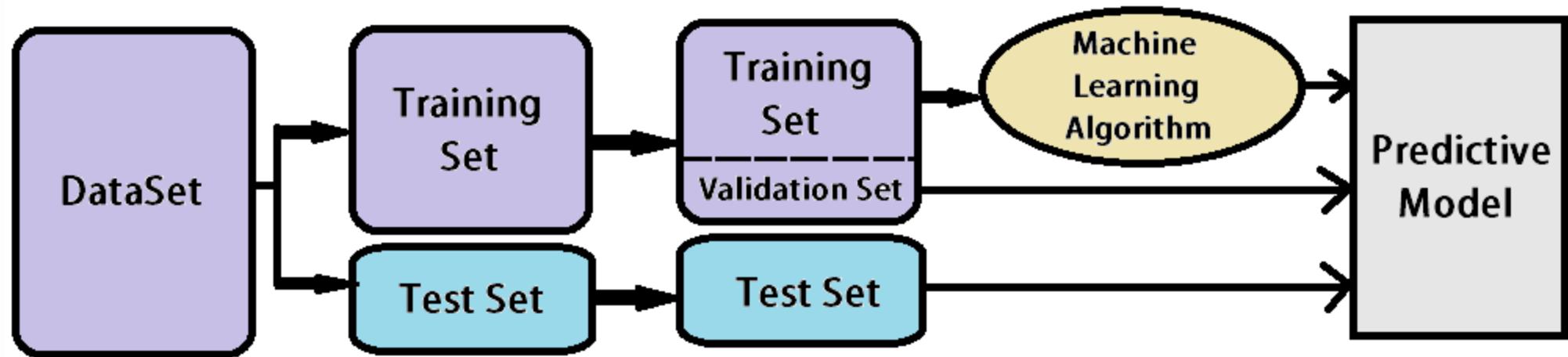
El proceso de la ciencia de datos



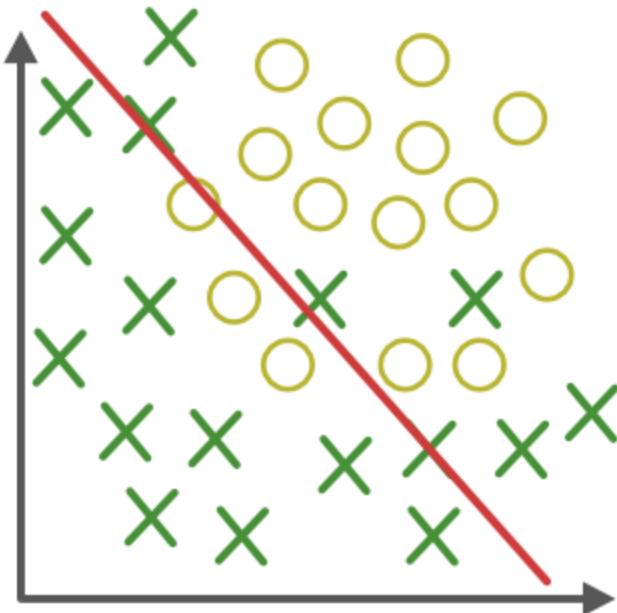
Machine learning paso a paso:



Proceso de evaluación:

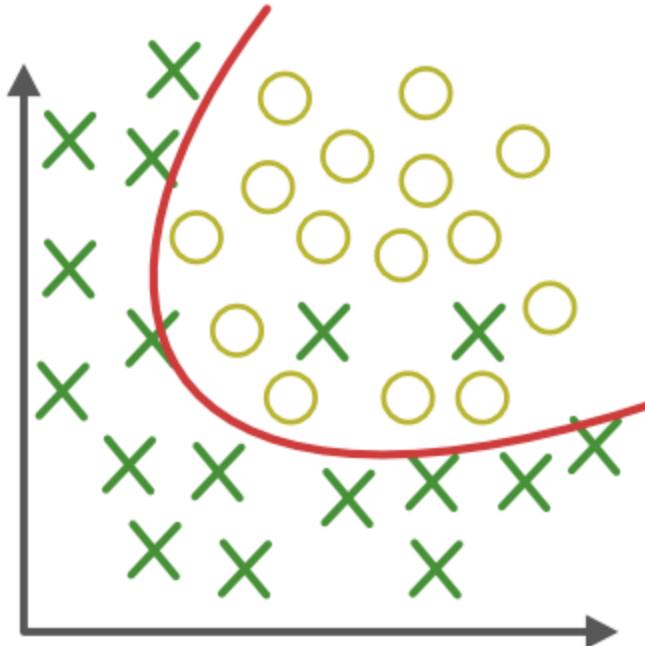


Conseguir el modelo perfecto:

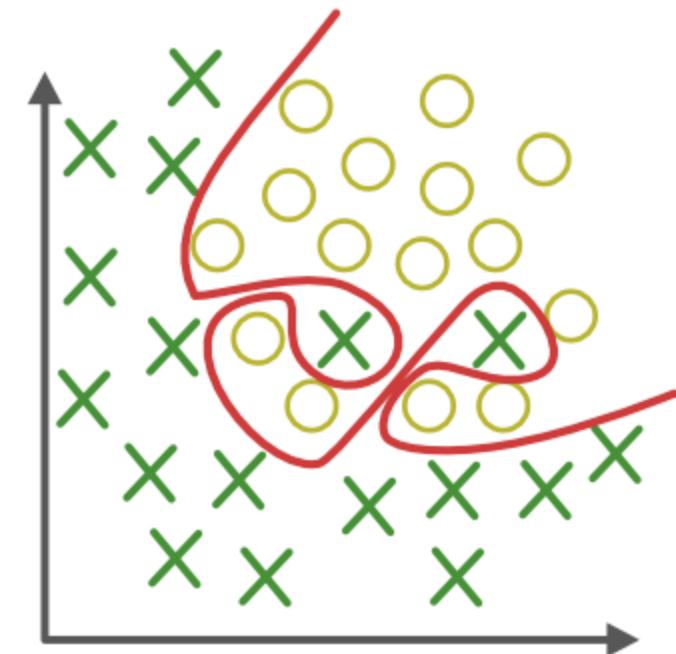


Under-fitting

(too simple to explain the variance)



Appropriate-fitting



Over-fitting

(forcefitting--too good to be true)