

D. Juan V. ...

Tipos y métodos de Machine Learning

Equipo Innovación Copec



Roberto Muñoz

Senior Data Scientist
Digital Data Analytics
EY



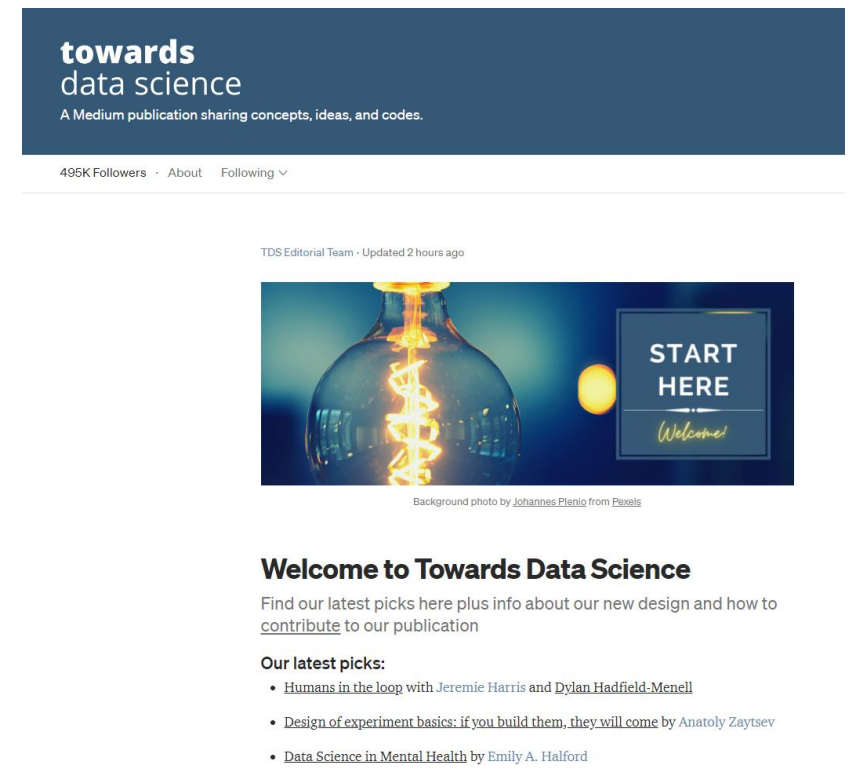
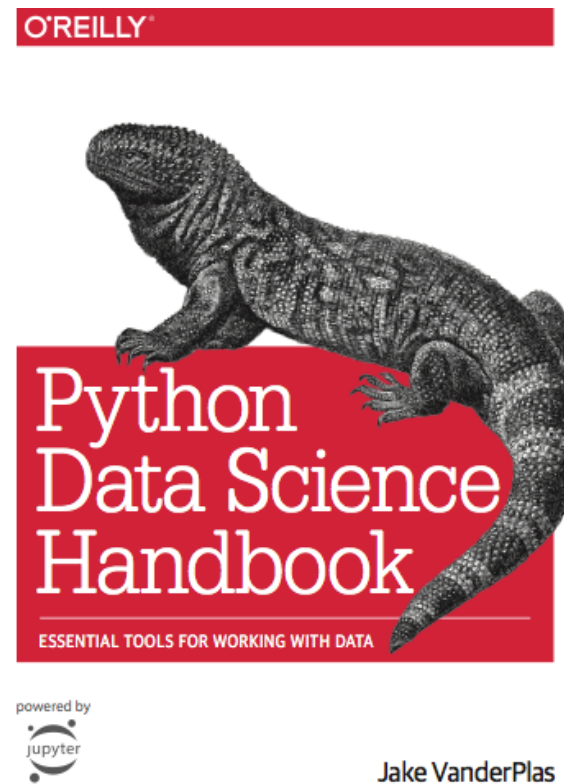
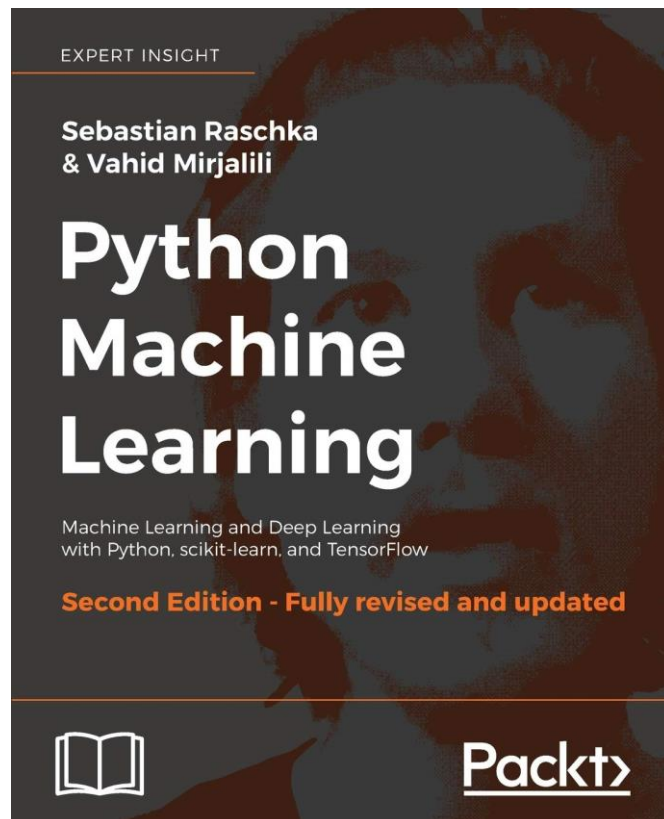
■ The better the question. ■ The better the answer. ■ The better the world works.



Building a better
working world

Bibliografía

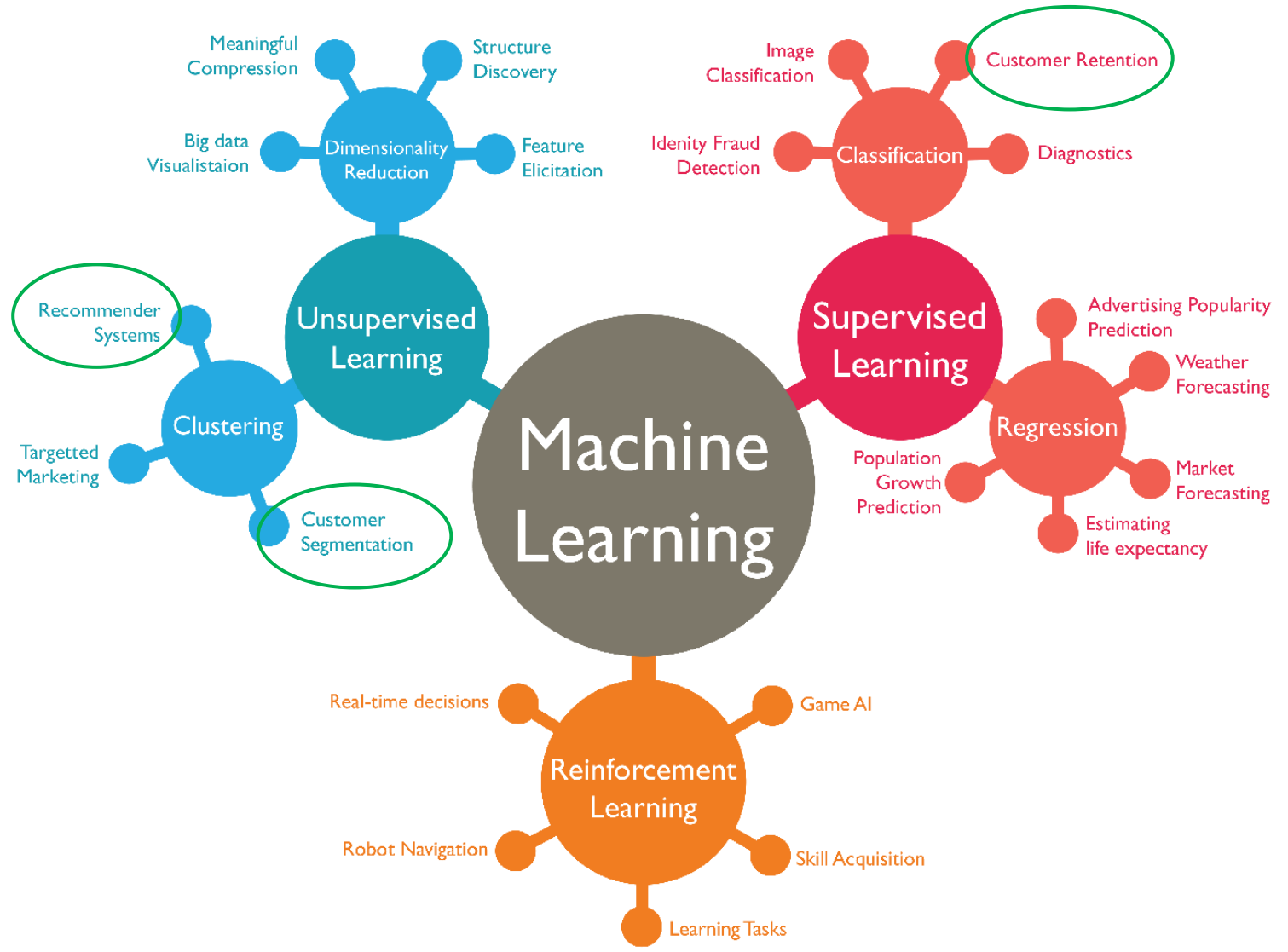
Libros y medios digitales



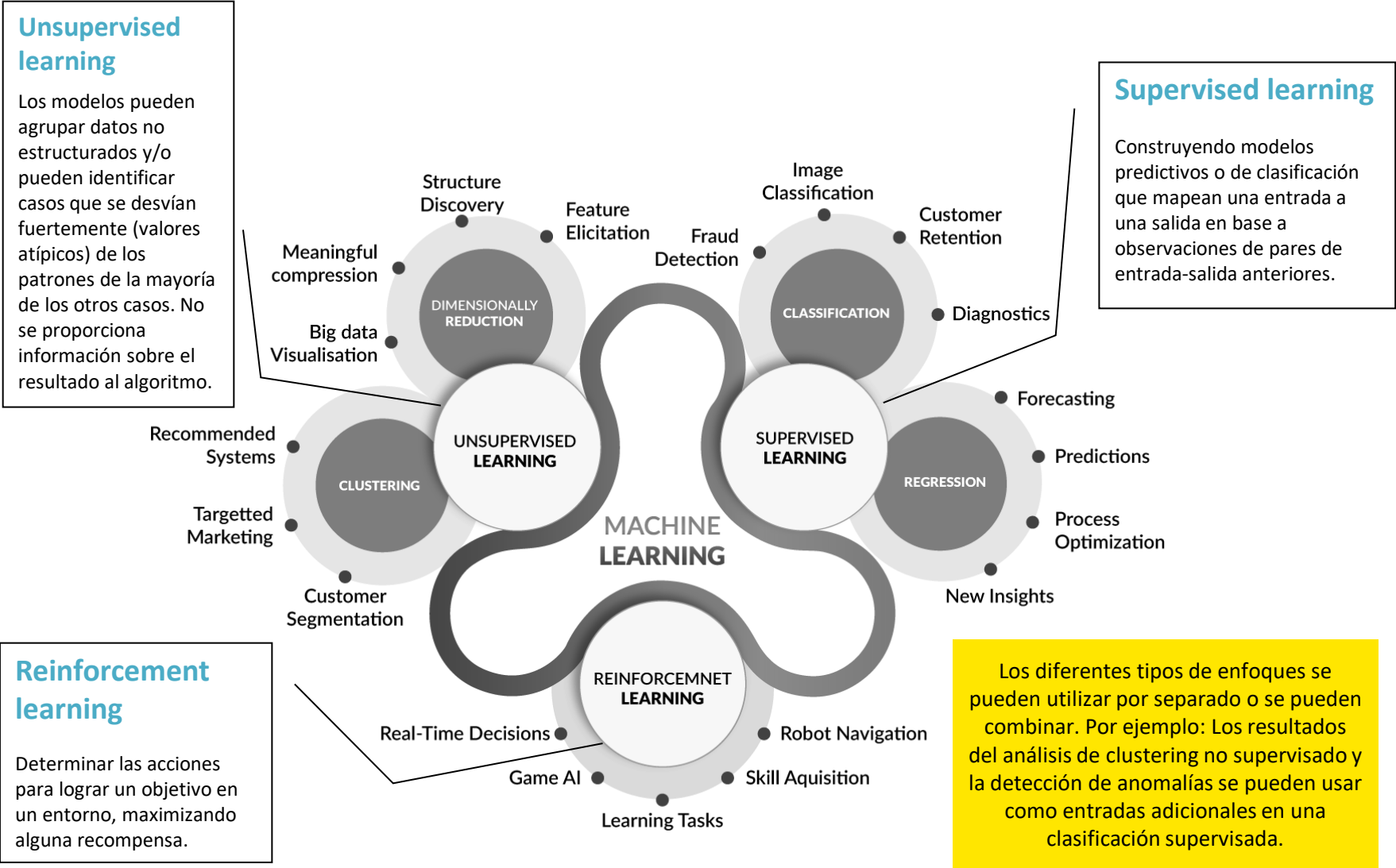
<https://jakevdp.github.io/PythonDataScienceHandbook/index.html>

<https://towardsdatascience.com/>

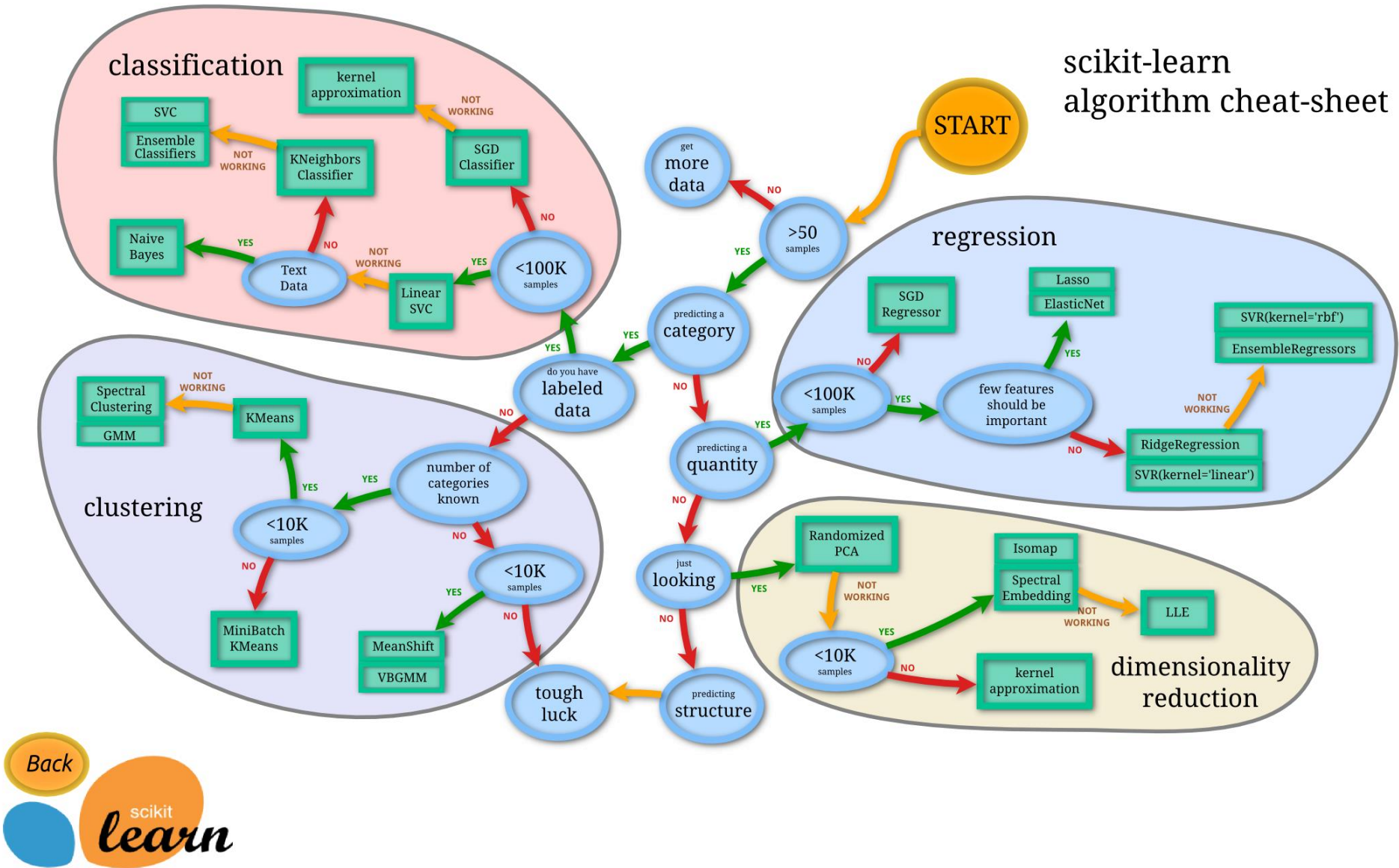
Tipos de aprendizaje



Tipos de aprendizaje



Scikit-learn methods



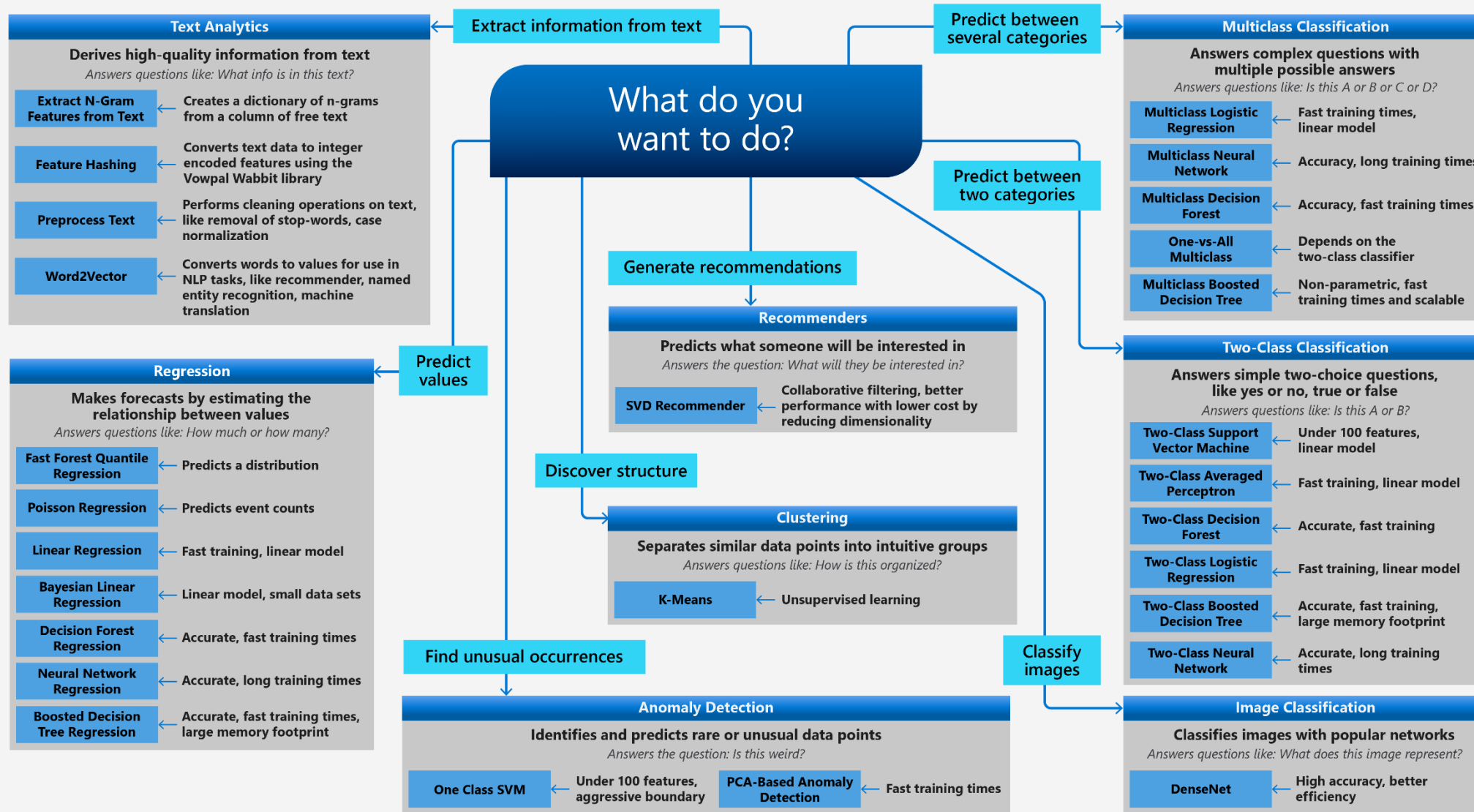


Microsoft Azure Machine Learning Algorithm Cheat Sheet

This cheat sheet helps you choose the best machine learning algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the goal you want to achieve with your data.

Roberto Muñoz

7

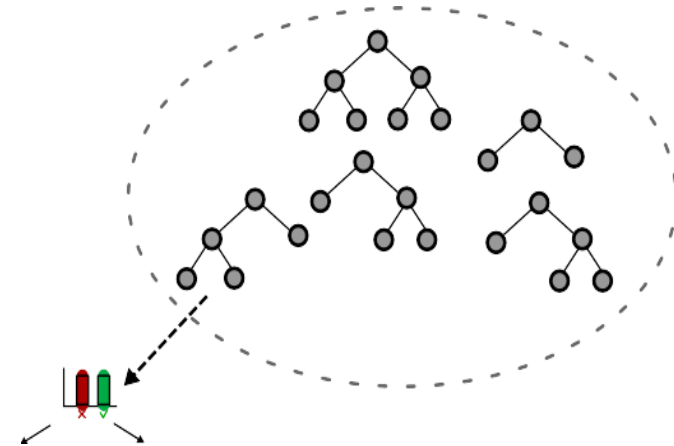
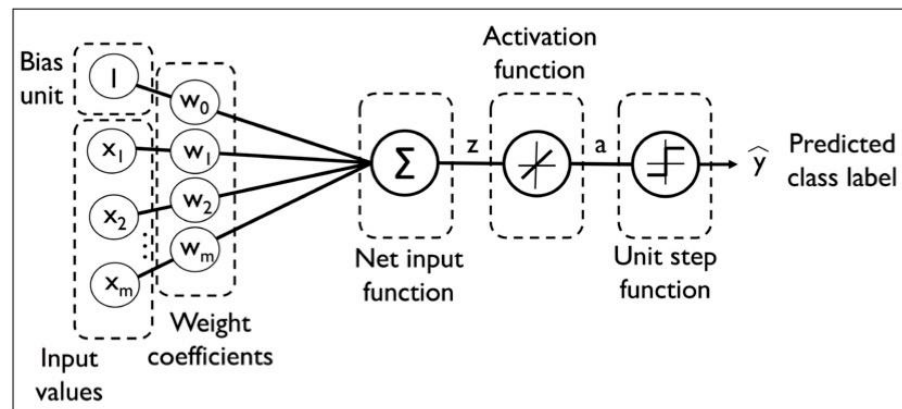


Supervised learning



Clasificación

- Los métodos más usados para resolver los problemas de clasificación en ML son
 - Support Vector Machine (SVM)
 - Arbol de decisión
 - Random forest
 - Deep Neural Networks



Atributos				Atributo objetivo
Nombre	Saldo	Edad	Empleo	Pérdida
Mike	\$200,000	42	no	yes
Mary	\$35,000	33	yes	no
Claudio	\$115,000	40	no	no
Robert	\$29,000	23	yes	yes
Dora	\$72,000	31	no	no

This is one row (example).
Feature vector is: **<Claudio,115000,40,no>**
Class label (value of Target attribute) is **no**

Clasificación

1)



→ perro

6)



→ gato

2)



→ perro

7)



→ perro

3)



→ gato

8)



→ gato

4)



→ perro

9)



→ perro

5)



→ perro

10)



→ gato

Matriz de confusión

id	observado	prediccion
1	perro	perro
2	perro	perro
3	perro	gato
4	perro	perro
5	perro	perro
6	perro	gato
7	perro	perro
8	gato	gato
9	gato	perro
10	gato	gato

		True/Actual	
		Positive (🐶)	Negative
Predicted	Positive (🐶)	5 (TP)	1 (FP)
	Negative	2 (FN)	2 (TN)

Matriz de confusión

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

Inglés

		Predicción	
		Positivos	Negativos
Observación	Positivos	Verdaderos Positivos (VP)	Falsos Negativos (FN)
	Negativos	Falsos Positivos (FP)	Verdaderos Negativos (VN)

Español

Recall, Precision, Accuracy

$$\text{Recall} = \frac{tp}{tp + fn}$$

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn}$$

- **Entrenamiento:** Ajustar los parámetros del algoritmo de forma tal de que se minimicen la cantidad de predicciones que no correspondan a la etiqueta original.
 - **Recall:** Porcentaje de clasificados correctamente como positivos sobre todos los que realmente eran positivos.
 - **Precision:** Porcentaje de clasificados correctamente como positivos sobre todos los clasificados como positivos.
 - **Accuracy:** Porcentaje de clasificados correctamente.
- **Recall:** Sensibilidad
 - **Precision:** Precisión
 - **Accuracy:** Exactitud

Unsupervised learning



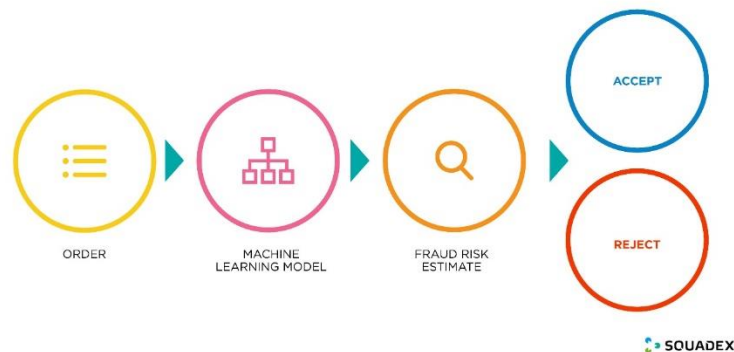
Créditos: Interstellar

Casos de uso

• Detección de fraudes

Los fraudes con tarjeta de crédito aumentaron un 104% entre Q1 2019 y Q2 2020.

Usar modelo para identificar operaciones fraudulentas.



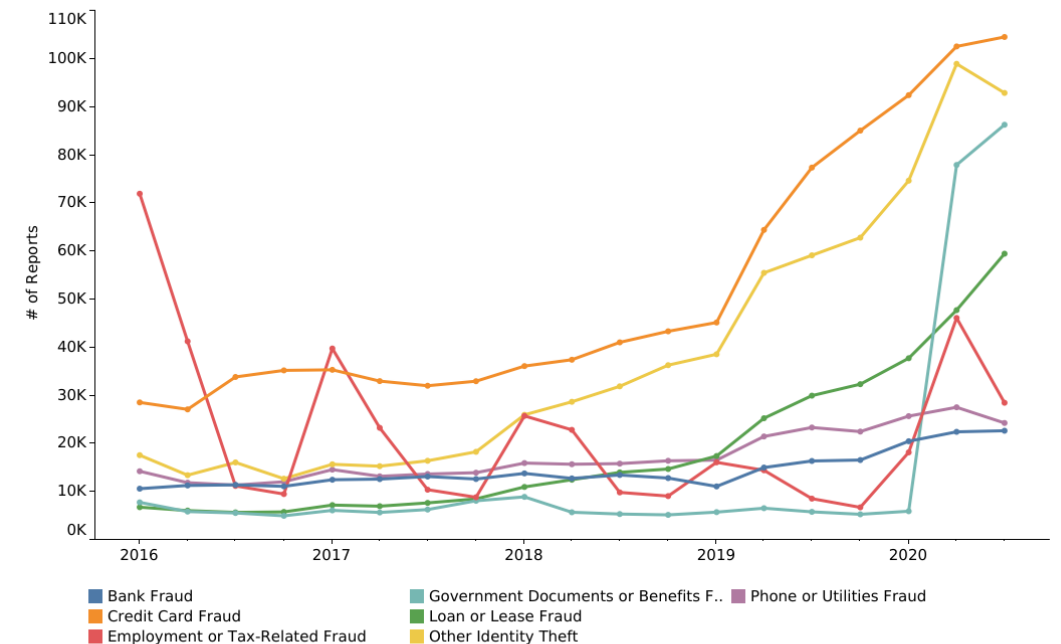
FTC CONSUMER SENTINEL NETWORK

Published October 16, 2020
(data as of September 30, 2020)

Compare Identity Theft Report..

Date Range
All values

Theft Type
All



Consumers can report multiple types of identity theft.

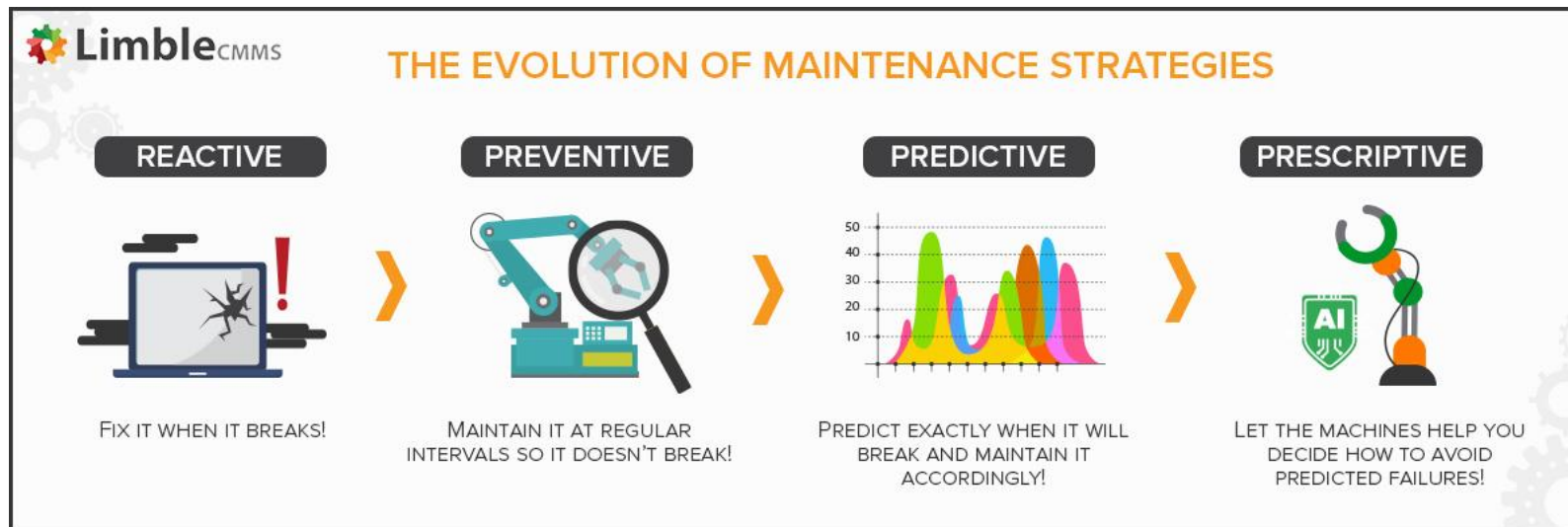
FEDERAL TRADE COMMISSION · ftc.gov/exploredata

Casos de uso

- **Mantenición preventiva de fallas**

Las industrias de las telecomunicaciones y manufactura están constantemente recolectando datos de sus operaciones. Cuentan con máquinas equipadas con múltiples sensores.

Usar modelos para identificar fallas de manera temprana y hacer mantenciones a la maquinaria.

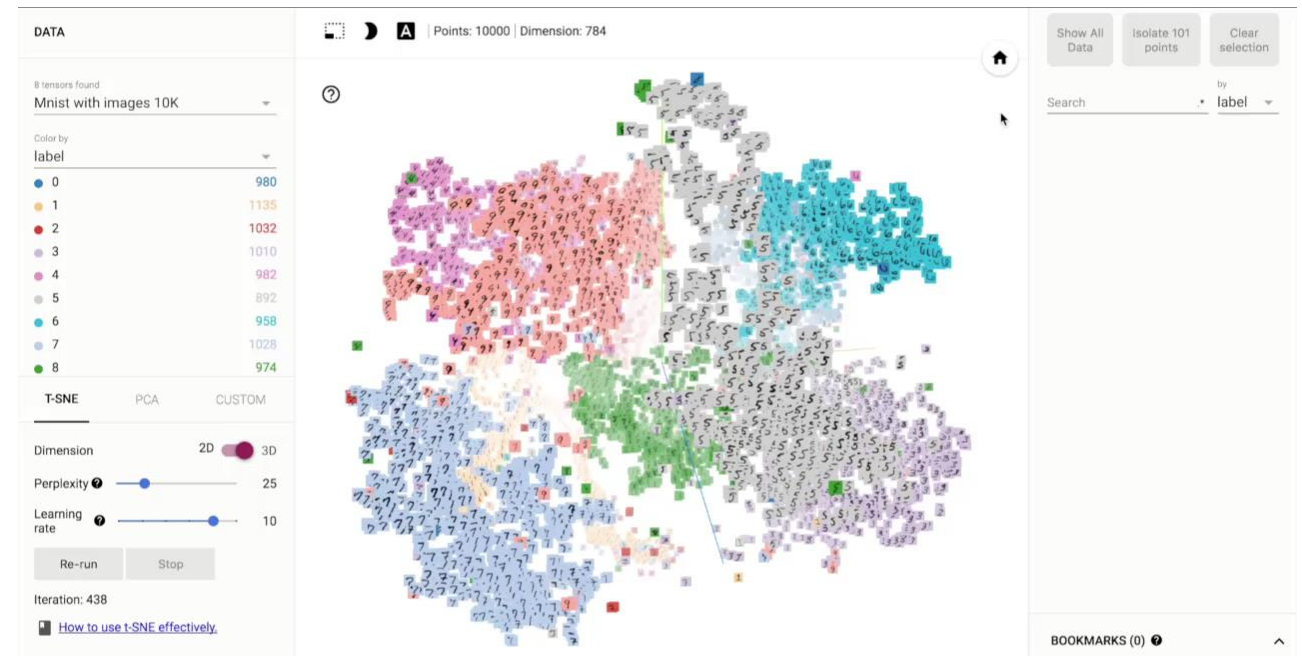
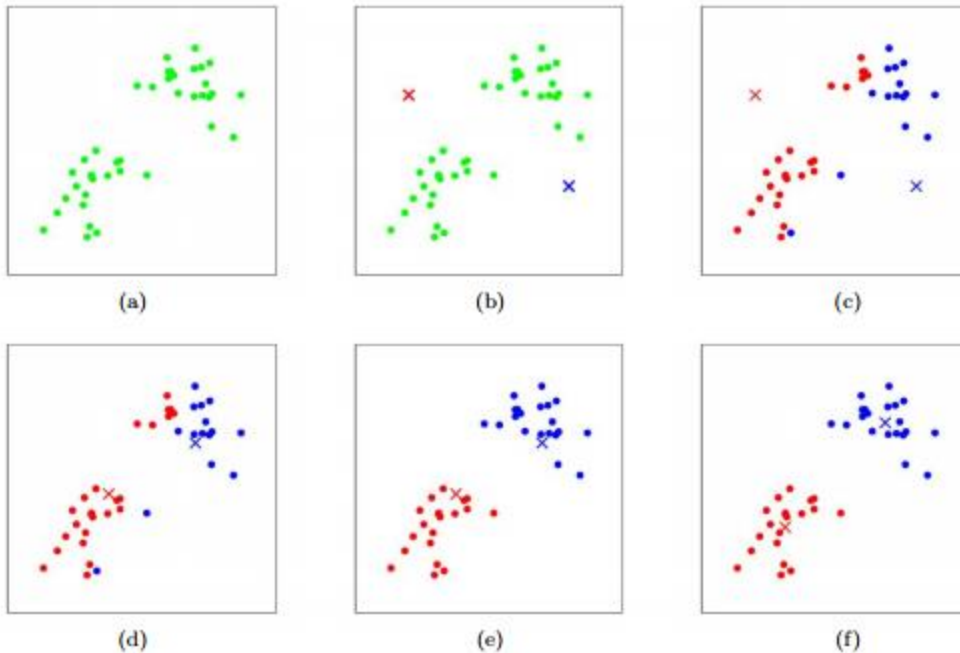


Tipos de ML

- Aprendizaje no supervisado

Entrenar un modelo usando datos que no han sido clasificados previamente.

El sistema debe poder reconocer patrones y generar sus propias etiquetas. El modelo debe clasificar los nuevos datos de entrada.



Tipos de tareas en Aprendizaje no supervisado

- Clustering o Agrupación

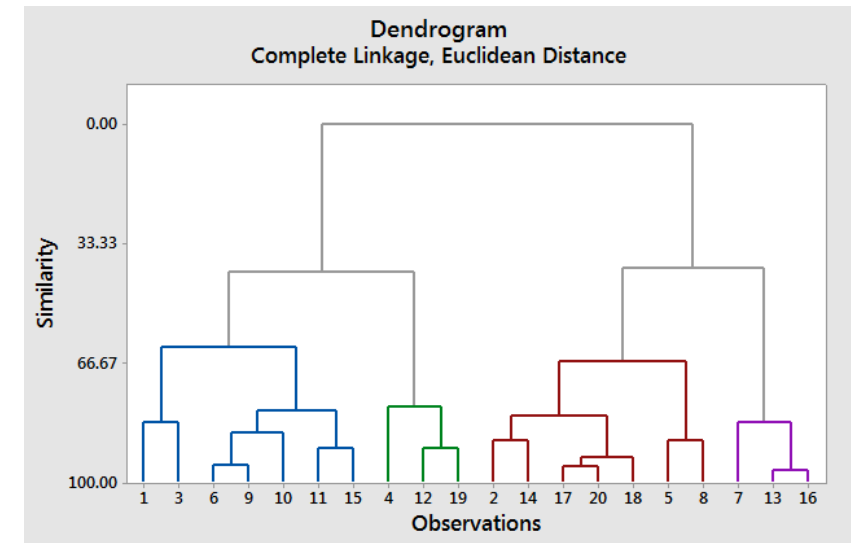
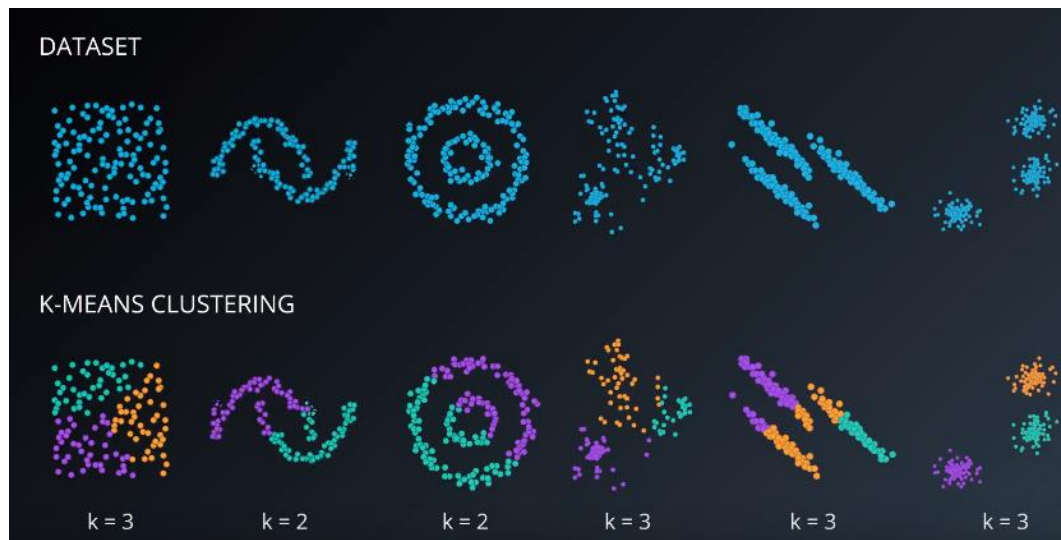
Encontrar diferentes grupos o segmentos presentes en los datos. Usado principalmente en segmentación de clientes y sistemas de recomendación.

- Detección de anomalías

Analizar el comportamiento y los patrones regulares en los datos y detectar desviaciones respecto al comportamiento normal. Usado principalmente en análisis de series de tiempo y análisis de imágenes.

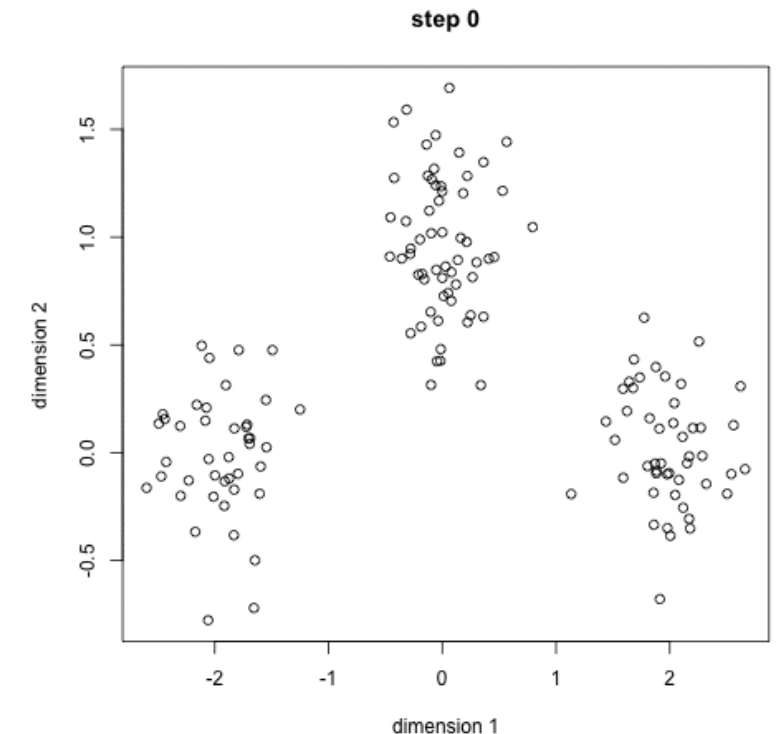
Métodos de clustering

- K-means
- Mean-shift
- Density-Based Spatial Clustering of Applications with Noise
- Gaussian mixture models
- Clustering jerquizado



Método de K-means

- K-Means tiene como objetivo encontrar y agrupar en clases los puntos de datos que tienen una alta similitud entre ellos.
1. Primero, necesitamos elegir k , el número de clusters que queremos que nos encuentren.
 2. Luego, el algoritmo seleccionará aleatoriamente los centroides de cada grupo.
 3. Se asignará cada punto de datos al centroide más cercano (utilizando la distancia euclídea).
 4. Se calculará la inercia del conglomerado.
 5. Los nuevos centroides se calcularán como la media de los puntos que pertenecen al centroide del paso anterior. En otras palabras, calculando el error cuadrático mínimo de los puntos de datos al centro de cada cluster, moviendo el centro hacia ese punto.
 6. Volver al paso 3.



The background is a complex, abstract digital network. It features a dense web of thin, reddish-brown lines connecting numerous small, glowing nodes. The nodes are primarily white and yellow, with some appearing as bright cyan or blue spheres. A prominent, bright yellow-orange glow emanates from a central node, radiating outwards. The overall color palette is dark blue and black, with the network elements providing a high-contrast, futuristic aesthetic. Faint, vertical columns of binary code (0s and 1s) are visible in the background, adding to the digital theme.

Taller