



Módulo VI - Aprendizaje NO supervisado.

Clase 22: Clustering - Métricas



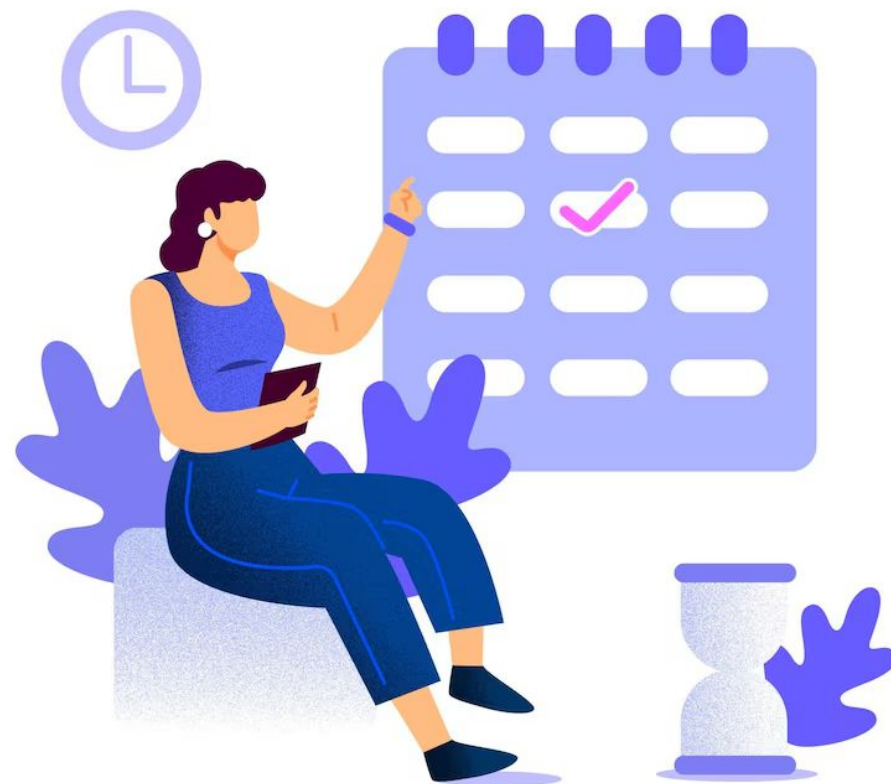


¿Ponemos a grabar el
taller?

¿QUÉ VAMOS A VER HOY?



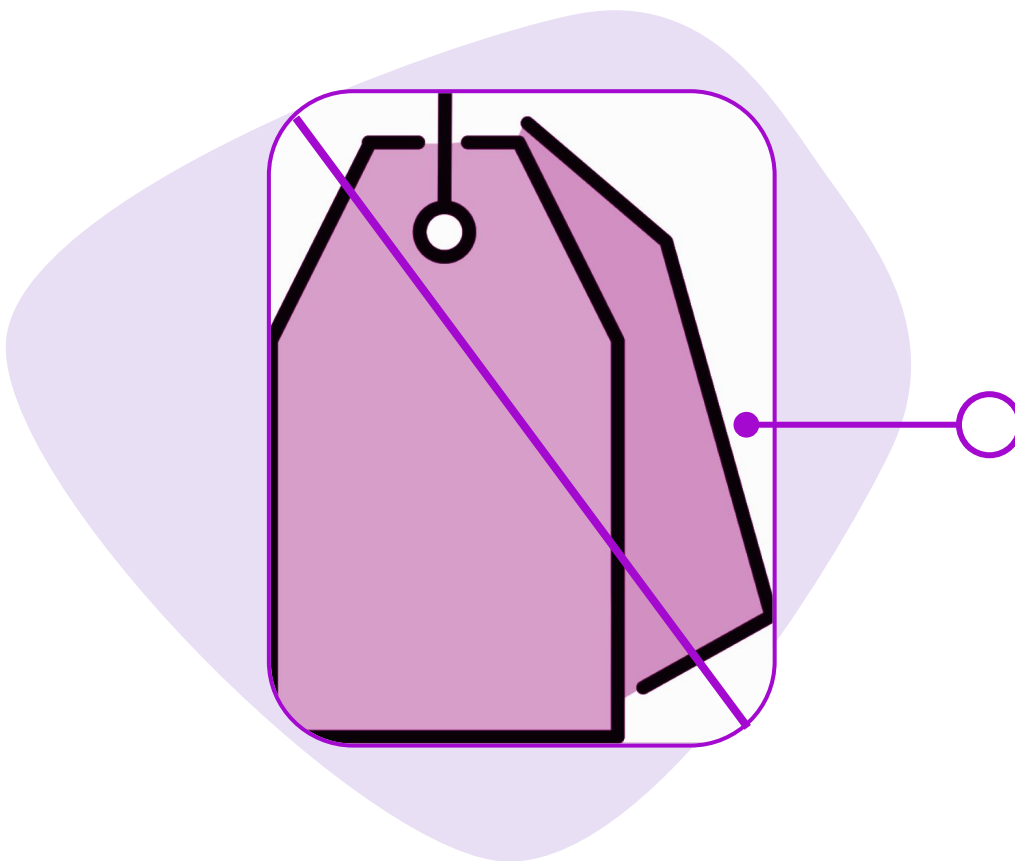
- Métricas para Clustering
 - Método Elbow
 - Método Silhouette





Métricas

Clustering



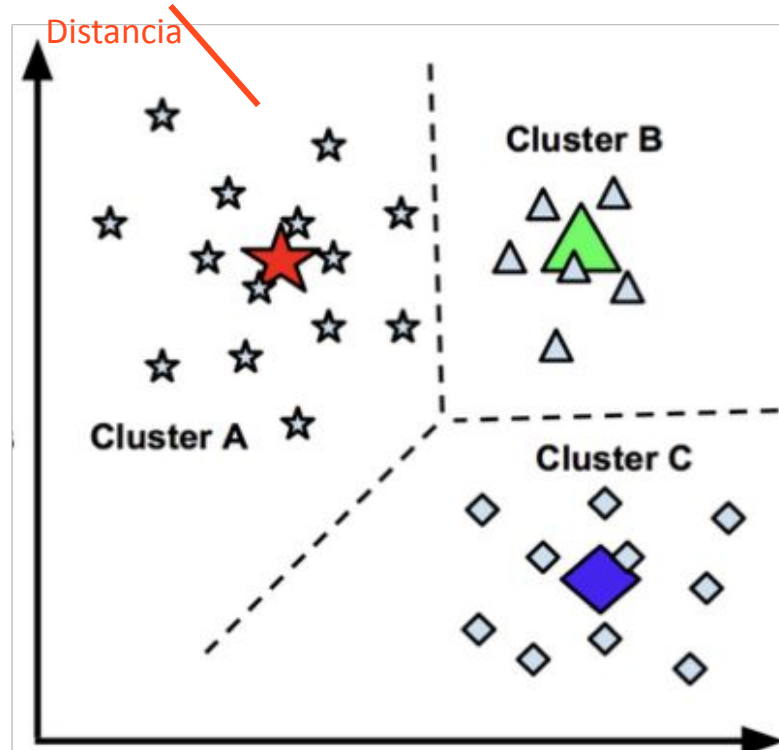
Los datos **NO** están
etiquetados.

No hay un valor
esperado.

**¿CÓMO SABEMOS QUE NUESTRO MODELO
FUNCIONA CORRECTAMENTE?**

Clustering

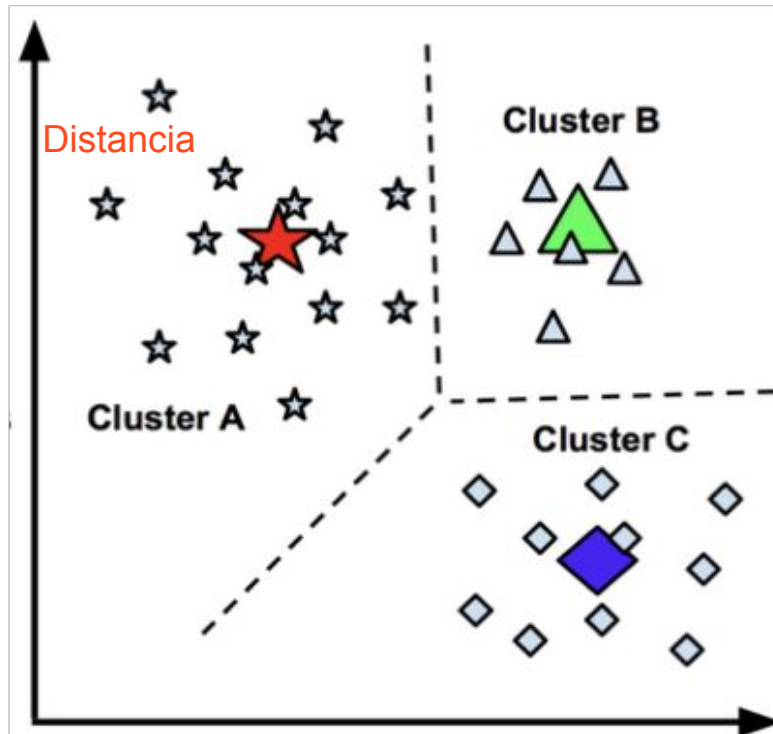
Debemos encontrar una medida para la **validación e interpretación de clusters en un dataset.**



Podemos medir cuál es la distancia media de cada dato al centroide más cercano.

Clustering

Debemos encontrar una medida para la **validación e interpretación de clusters en un dataset**.



$$d(i) = \|\mathbf{X}_i - \mathbf{C}_j\|^2$$

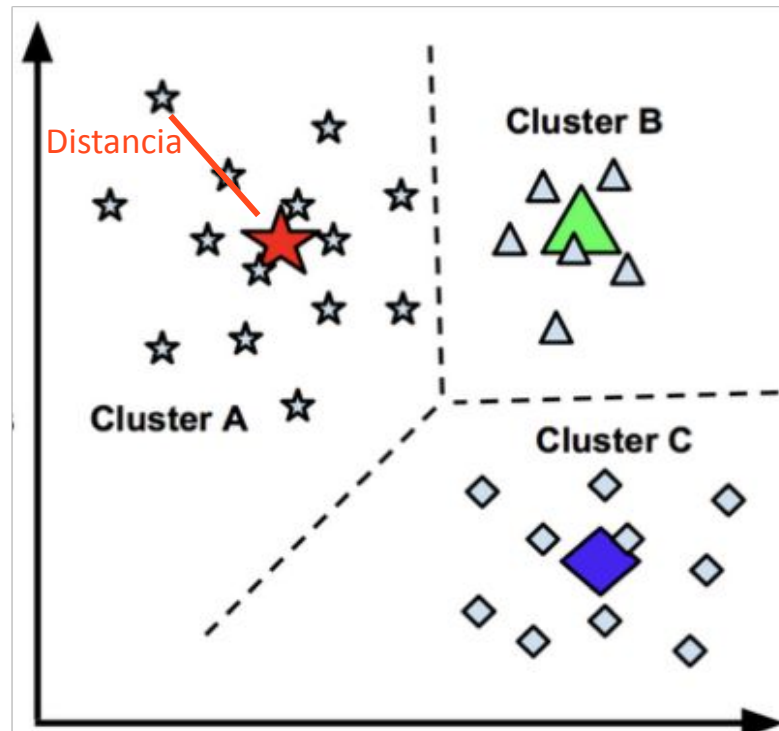
Distancia para
cada punto

Posición del
punto i

Posición del centroide
más cercano

Clustering (K-Means)

Buscamos una medida para evaluar qué tan buena resulta la **designación de clusters** en **K-Means**.



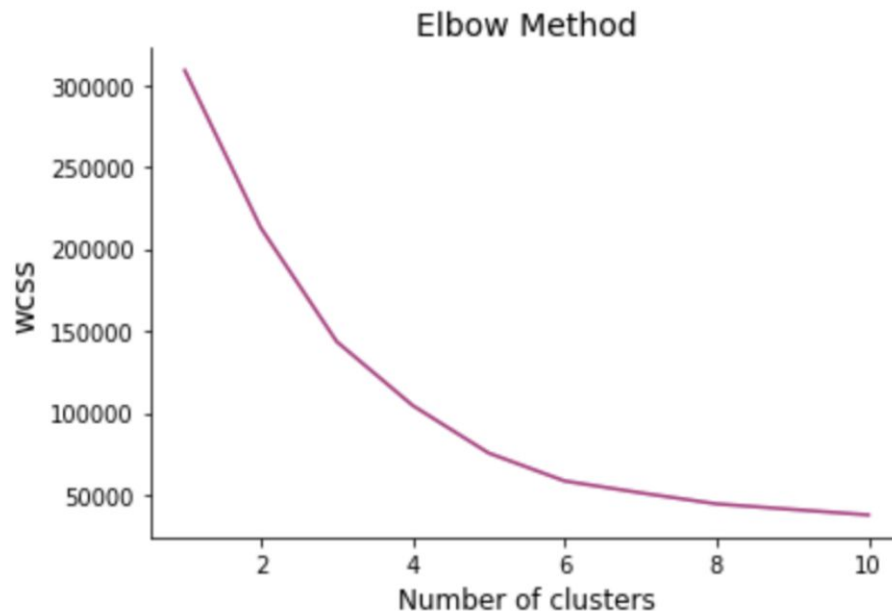
$$D = \frac{1}{N} \sum_{i=1}^N d(i)$$

Distancia
media total



Método Elbow

Buscamos una medida para evaluar qué tan buena resulta la **designación de clusters** en **K-Means**.



Se gráfica distancia (o inercia en sklearn) vs. k y se encuentra donde está el codo

Se busca dónde está el **codo** de la curva. El valor de la distancia o inercia siempre descende cuando aumenta el número de clusters.

Método Silhouette

Medida para la validación e interpretación de clusters en un dataset (para cualquier método de clustering).

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

Se mide qué tan parecidos son los datos con respecto a otros de su propio cluster (cohesión) en comparación con qué tan parecidos son a los datos en otros clusters (separación)

s(i): Valor de silhouette para el dato i.

a(i): Distancia media del dato i con el resto de su cluster

b(i): Distancia media del dato i con el cluster más cercano

Método Silhouette

Medida para la validación e interpretación de clusters en un dataset (para cualquier método de clustering).

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

Nos da una medida para cada punto de qué tan bien están ubicados en sus clusters. Para una medida de **todo el conjunto**, se toma la **media** de todos los valores.

s(i): Valor de silhouette para el dato i.

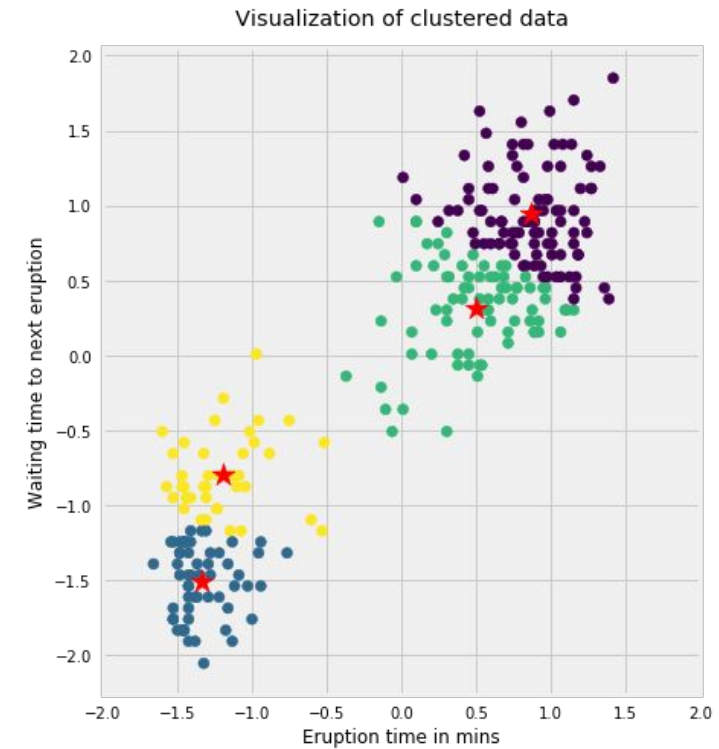
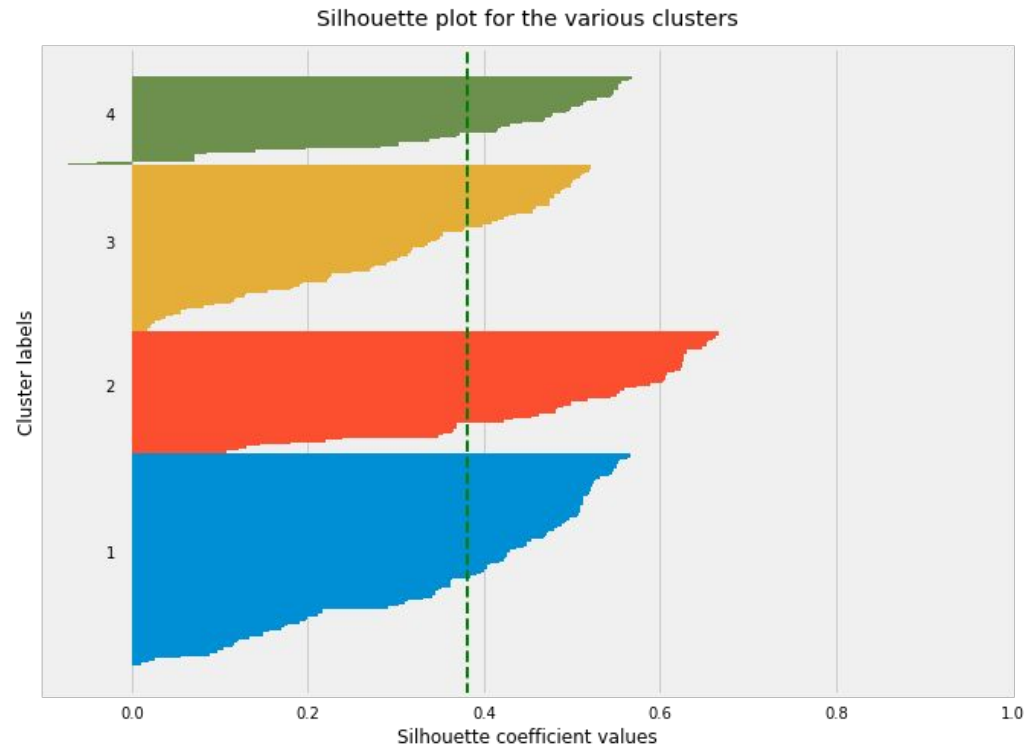
a(i): Distancia media del dato i con el resto de su cluster

b(i): Distancia media del dato i con el cluster más cercano

Método Silhouette

Se espera que si están bien asignados y armados los clusters, el perfil de todos los datos sea parejo.

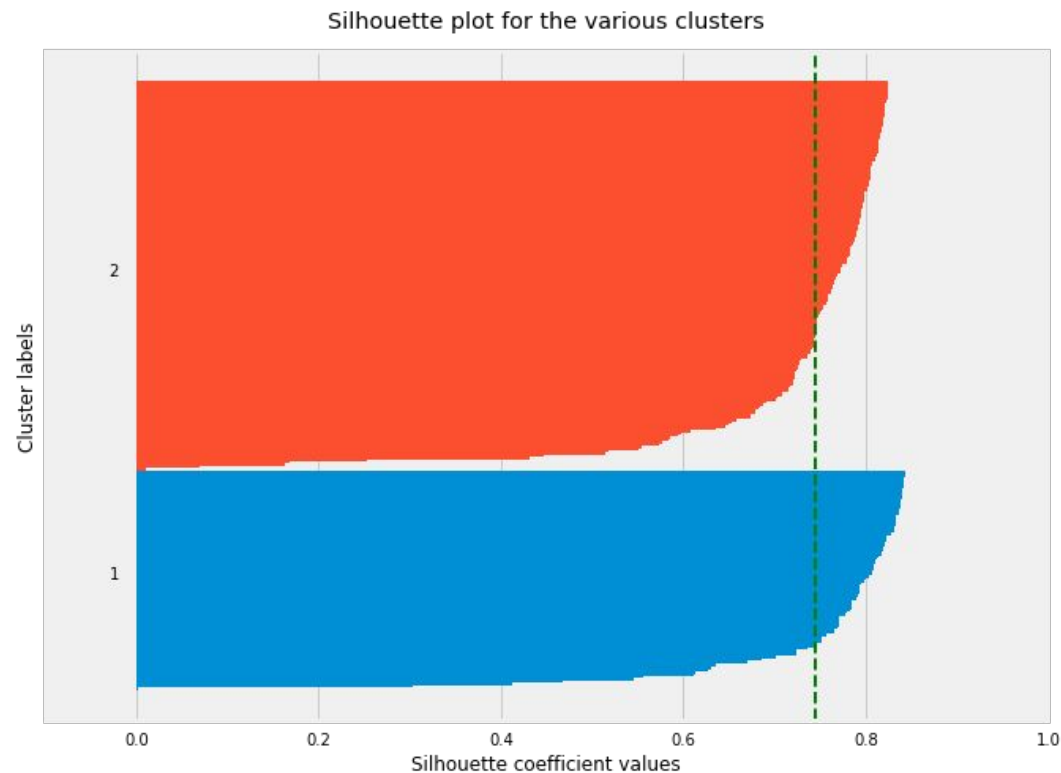
Silhouette analysis using $k = 4$



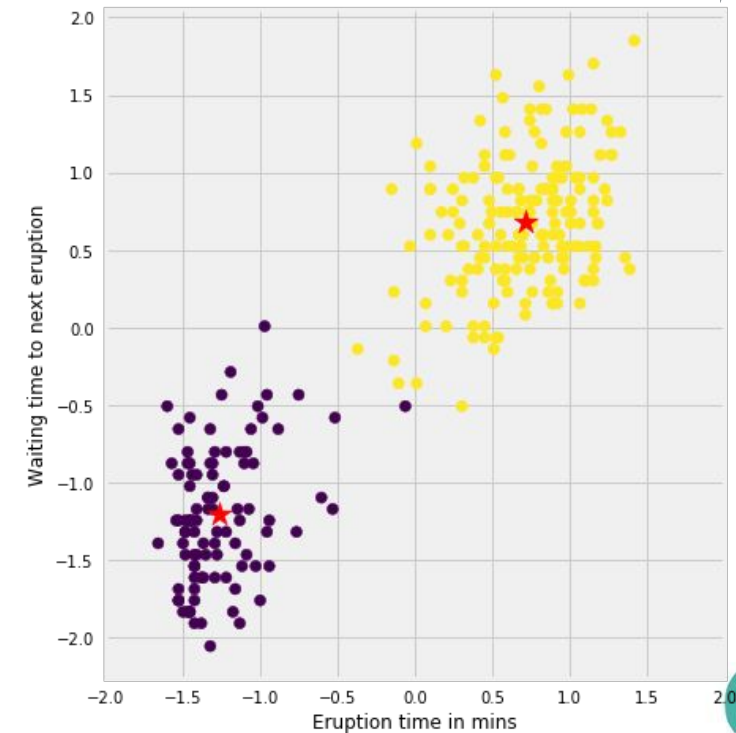
Método Silhouette

Se espera que si están bien asignados y armados los clusters, el perfil de todos los datos sea parejo.

Silhouette analysis using $k = 2$



Visualization of clustered data



Sección Práctica

Trabajamos con métricas en la Notebook 21

TRABAJAMOS EN SALAS



Sección Práctica

Trabajamos con métricas en la Notebook 21

TRABAJAMOS EN SALAS



Trabajamos en salas de zoom

Métricas

Trabajaremos con la Notebook 21

En los grupos establecidos,
ejercitamos como se evalúan
los modelos no supervisados



50 minutos de actividad



Descanso

Nos vemos en 10 minutos

Repasamos dudas

**Trabajamos con la
Notebook 21**

Revisamos los conceptos y el
código trabajados en la
notebook 21

Desafío 15 (continuación)

Para la siguiente repasar la notebook 21



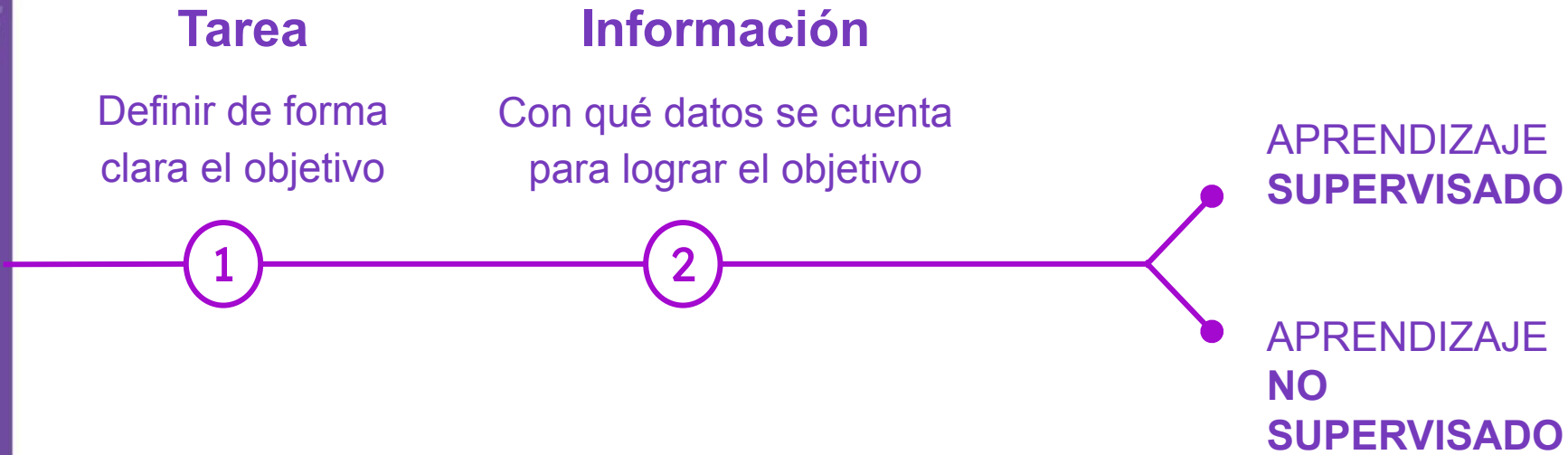


¿Alguna consulta?

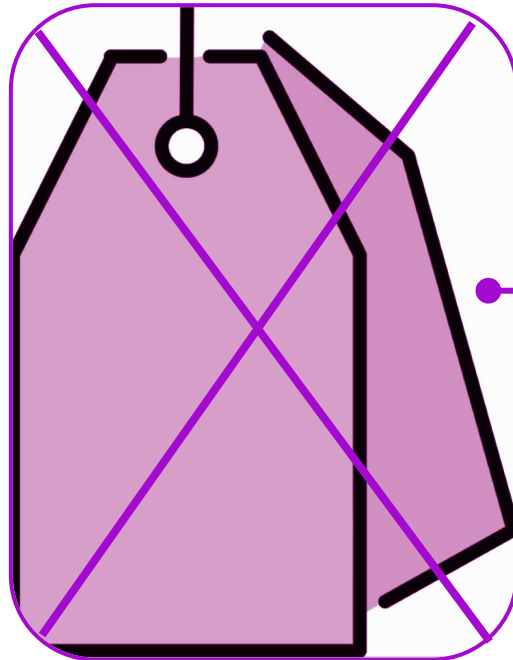


Repasemos

Eligiendo algoritmo

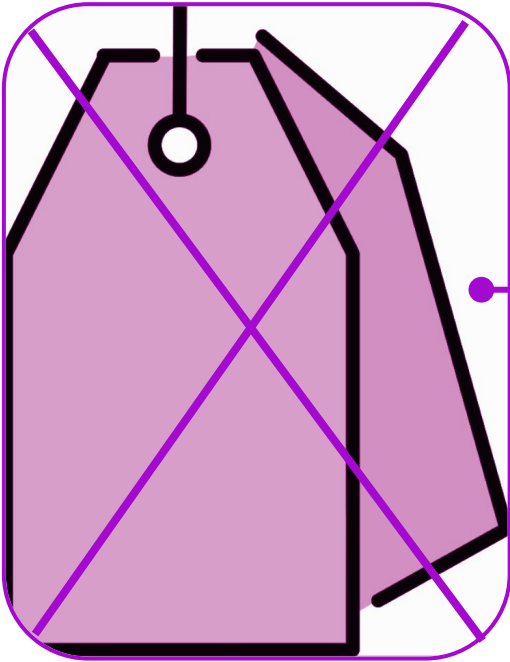


Aprendizaje No Supervisado



El algoritmo infiere patrones de un conjunto de datos que, a diferencia del aprendizaje supervisado, **no** están **etiquetados**. Puede utilizarse para descubrir la estructura subyacente de los datos

Aprendizaje No Supervisado



- Clustering
- Reducción de la Dimensionalidad



Reducción de la dimensionalidad

Reducción de la dimensionalidad

El objetivo consta de **reducir la cantidad de features** de un dataset, pero **reteniendo la mayor** cantidad de **información** posible.

Reducción de la dimensionalidad

Consiste en **reducir** la cantidad de **features** de un dataset, pero **reteniendo** la mayor cantidad de **información** posible.

- **Aplicaciones**
 - Reducir la complejidad del input en un modelo de regresión o clasificación
 - Visualización
 - Detectar features relevantes en datasets
- **Algoritmos**
 - **Principal Component Analysis**
 - Multidimensional Scaling
 - t-SNE: t-distributed Stochastic Neighbor Embedding
 - LDA: Linear Discriminant Analysis

Principal Component Analysis

Técnica utilizada para describir un conjunto de datos en términos de nuevas variables o componentes no correlacionados.

- El algoritmo encuentra nuevos componentes que describen los datos
- Los componentes se ordenan por la cantidad de varianza original que describen, reduciendo la dimensionalidad del conjunto de datos.
- Utiliza **Single Vector Decomposition (SVD)**

FUNDACIÓN
YPF

¡Muchas gracias!

