

Data Mining

Unit 1. Structured data mining

Prof. Dr. Ivar Vargas Belizario

ivargasbelizario@gmail.com

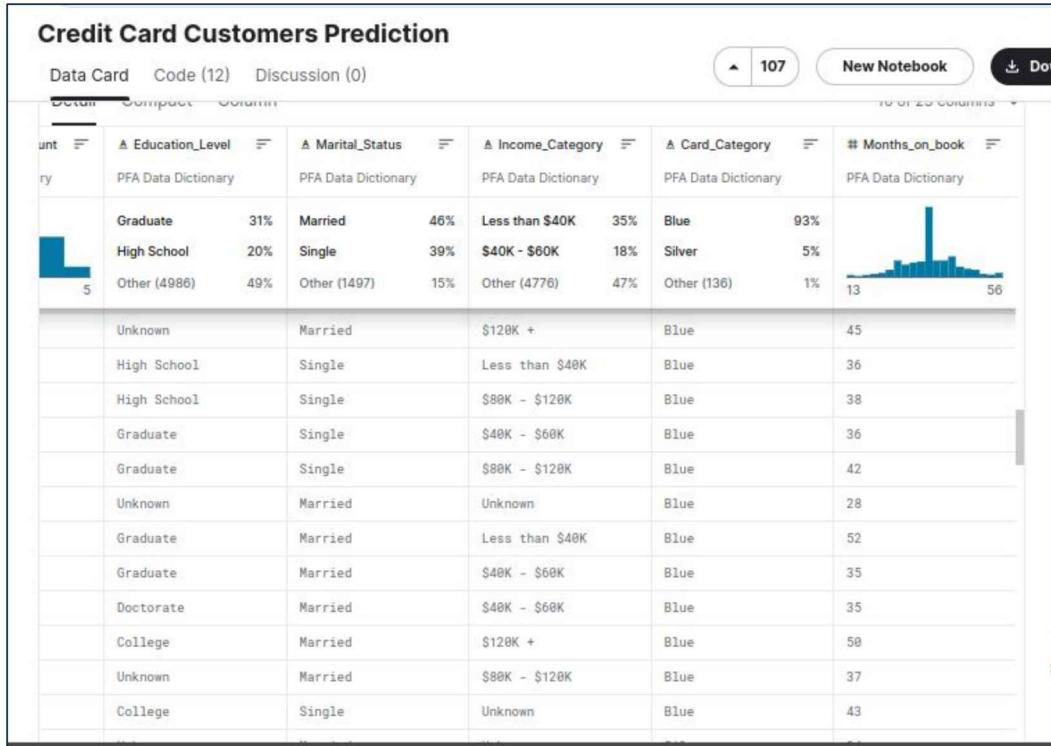
2025 - I



Contenido

- Structured data
- Dimensionality reduction
- Clustering
- Classification
- Regression
- Examples

Structured data



Ivar Vargas Belizario

3

Structured data

High-dimensional Data
Curse of Dimensionality

$X^{n \times m}$

y

n

m

F1	F2	...	Fm	Label
0.3	0.5	0.8	0.4	0.5
0.2	0.2	0.9	0.3	0.4
...
0.1	0.3	0.7	0.2	0.3

- Label:
- Discrete or
 - continuous

n: instances
m: features or attributes

Ivar Vargas Belizario

4

Structured data



Características:

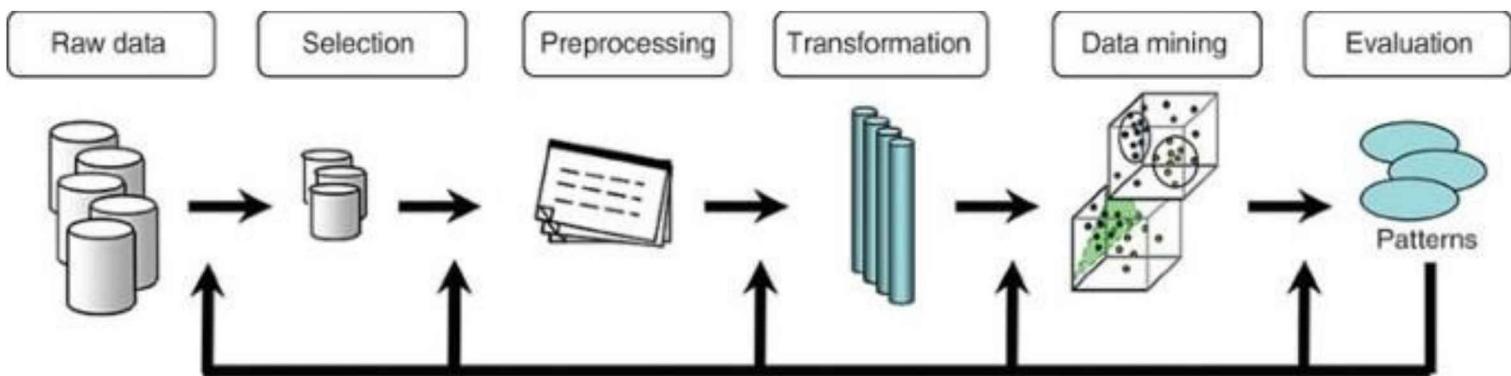
1. Grande tamaño de los datos:

- 1 Byte = 8 bits
- 1 Kilobyte (KB) = 1024 bytes
- 1 Megabyte (MB) = 1024 kilobytes
- 1 Gigabyte (GB) = 1024 megabytes
- 1 Terabyte (TB) = 1024 gigabytes
- 1 Petabyte (PB) = 1024 terabytes
- 1 Exabyte (EB) = 1024 petabytes
- 1 Zettabyte (ZB) = 1024 exabytes (2016-> tráfico en internet)
- 1 Yottabyte (YB) = 1024 zettabytes

2. Complejidad:

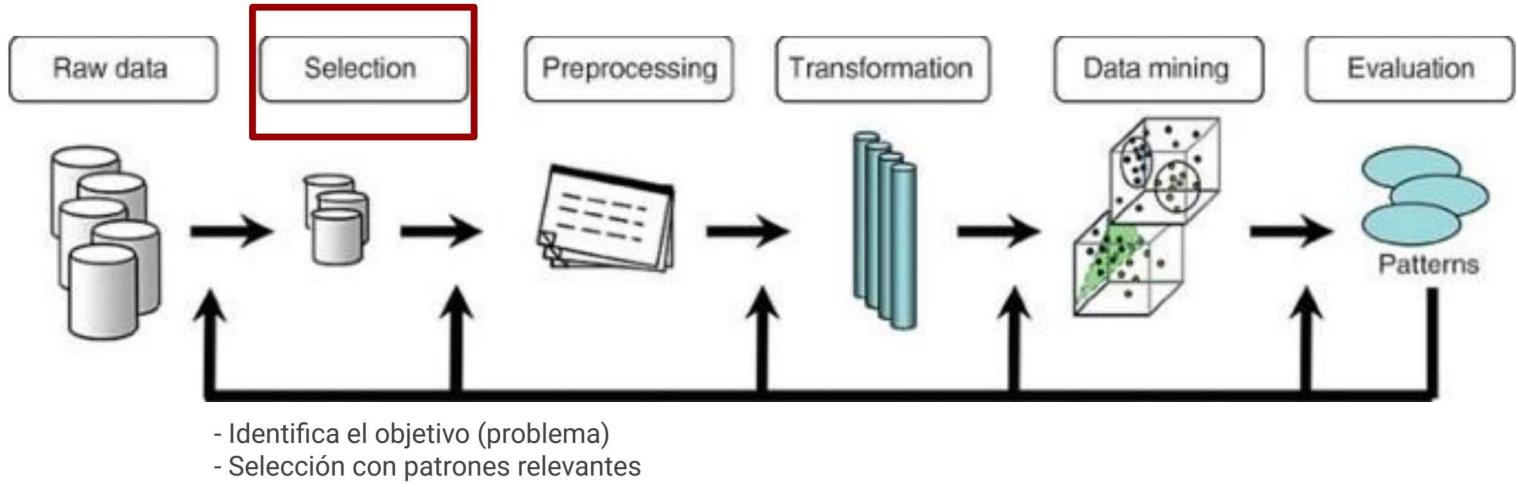
- No estructurados: Imágenes, videos, audios, etc.
- Sufren variación en el tiempo

Structured data (pipeline)



[4] https://doi.org/10.1007/978-1-4899-7993-3_1134-2

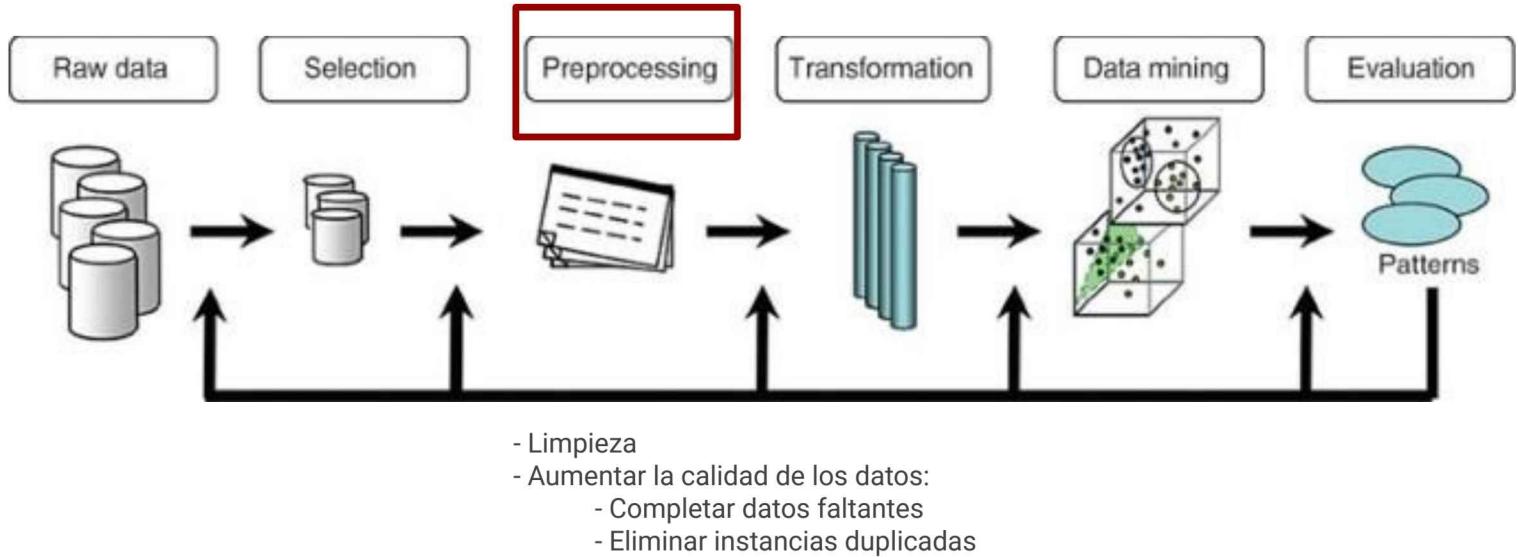
Structured data (pipeline)



7

Ivar Vargas Belizario

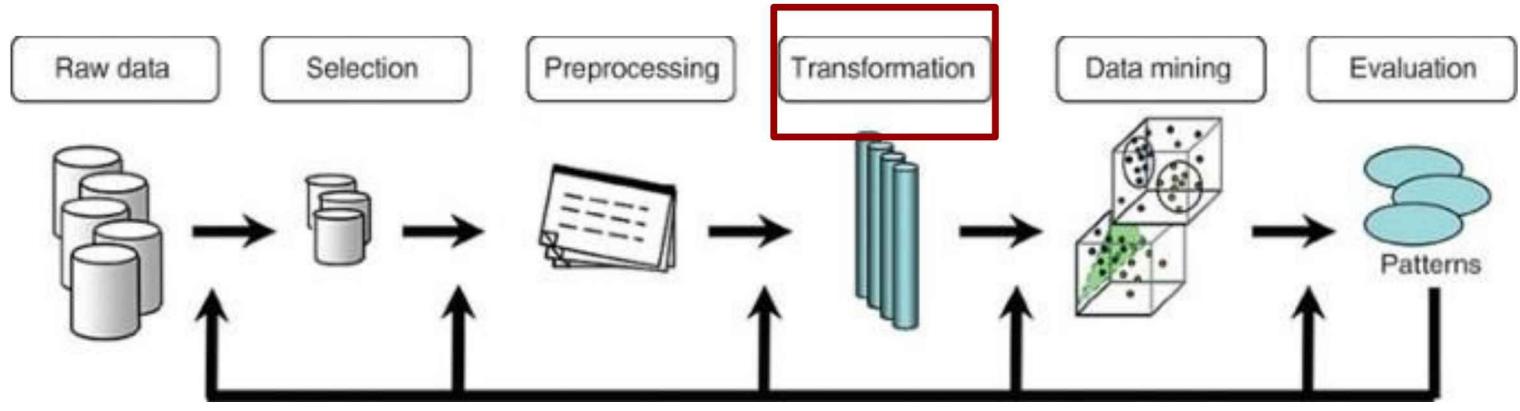
Structured data (pipeline)



8

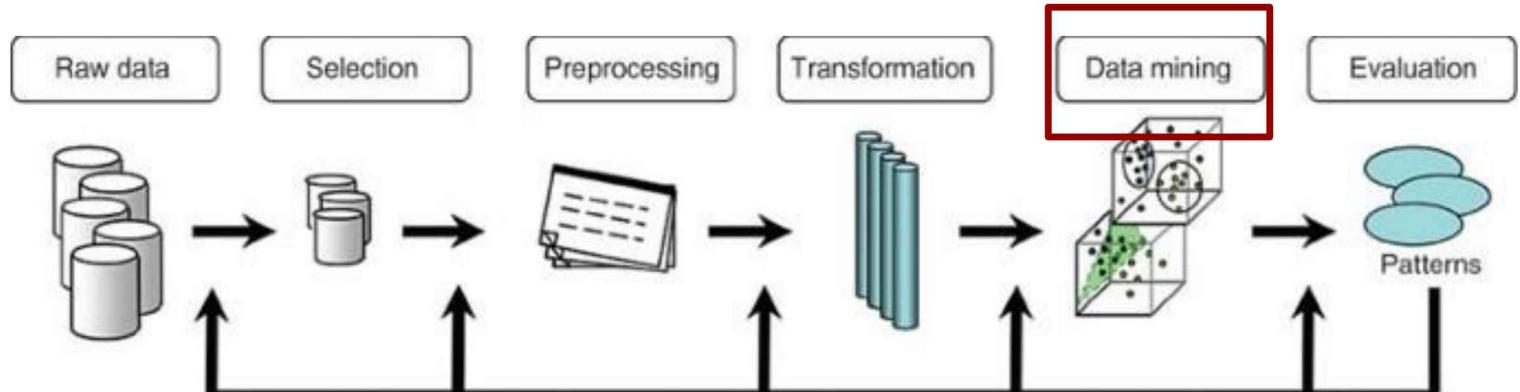
Ivar Vargas Belizario

Structured data (pipeline)



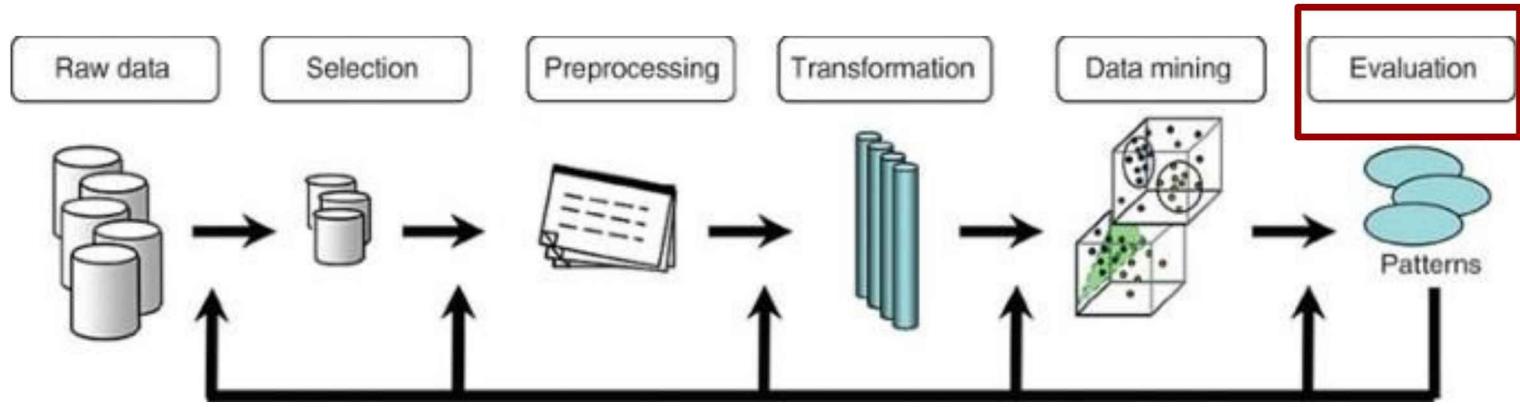
- Eliminar atributos correlacionados e irrelevantes.
- Se crea una nueva definición de atributos con mayor significancia.

Structured data (pipeline)



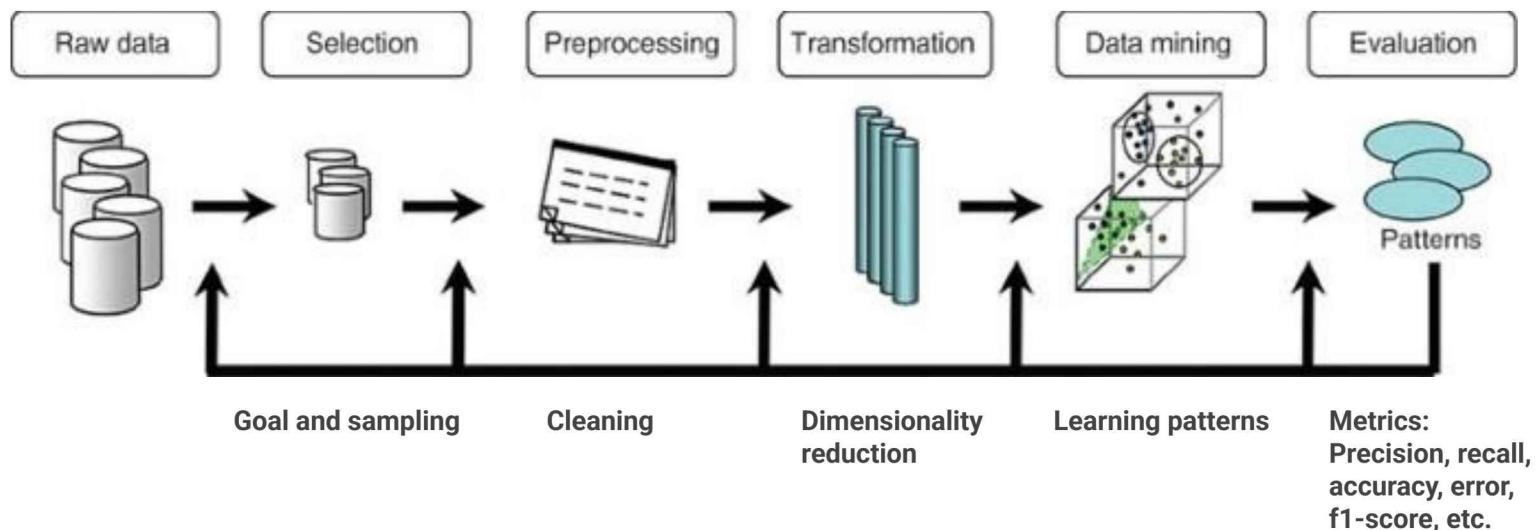
- Selecciona el tipo de algoritmo de minería de datos identificado en la etapa de **Selection**.
- Aprender patrones según:
 - Clasificación,
 - Agrupamiento
 - Regresión,
 - Asociación

Structured data (pipeline)

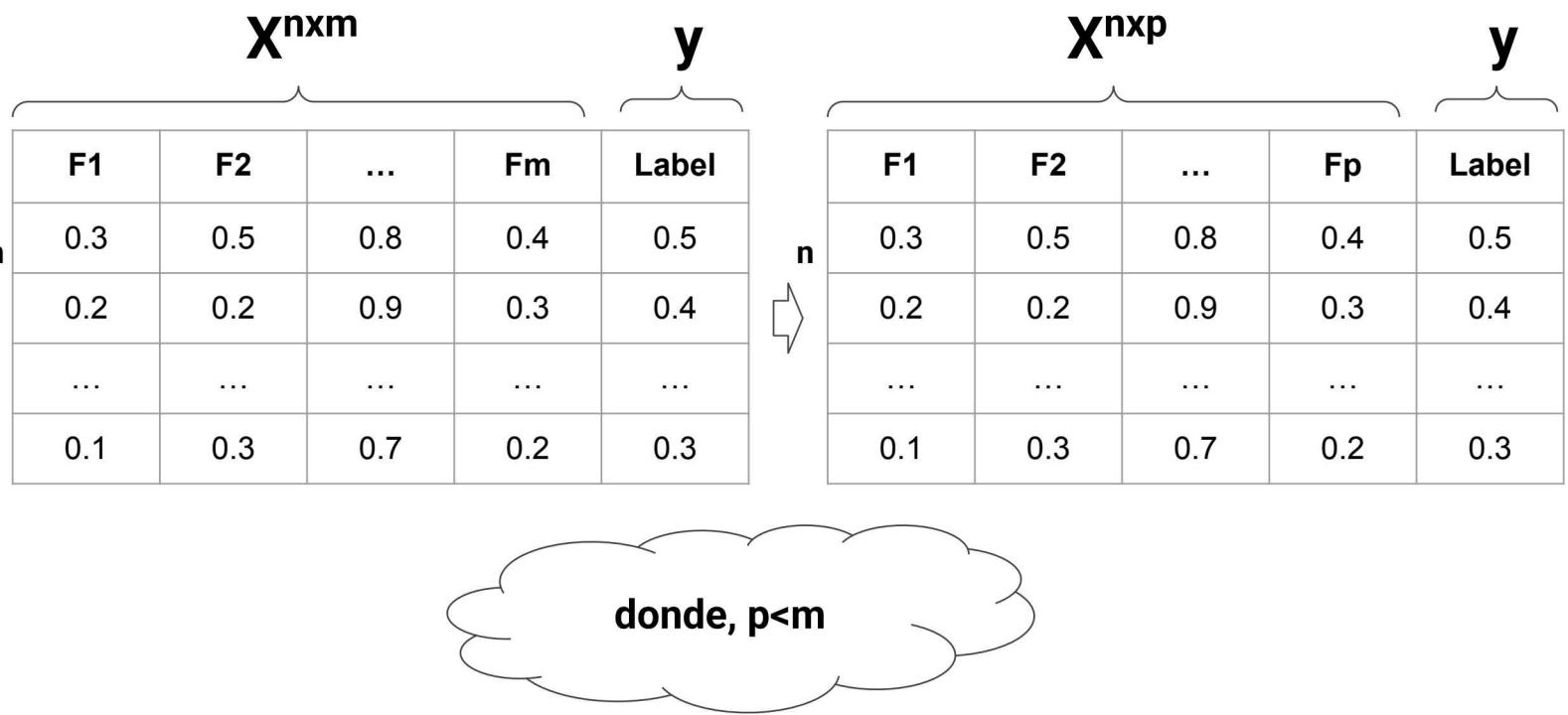


- Evalúa lo aprendido

Structured data



Dimensionality reduction

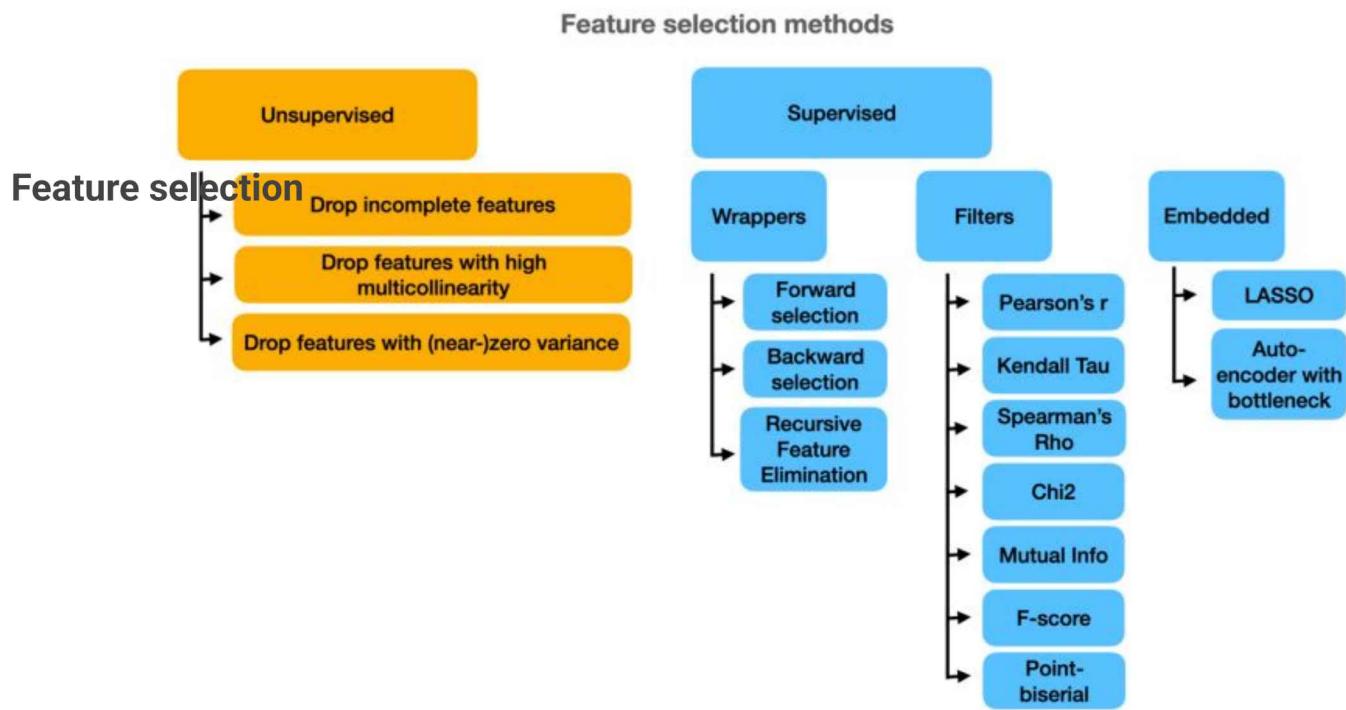


Dimensionality reduction

Tipos

1. **Feature selection:** Los datos no varian, se seleccionan los mejores atributos manteniendo su integridad.
 - i. Feature importance
 - ii. Based in correlation matrix
2. **Feature extraction:** Los datos sí varían por fusión o transformación de características.
 - i. PCA
 - ii. T-SNE
 - iii. UMAP
 - iv. LSP

Dimensionality reduction



Dimensionality reduction

Feature extraction

- PCA
- T-SNE
- UMAP
- LSP

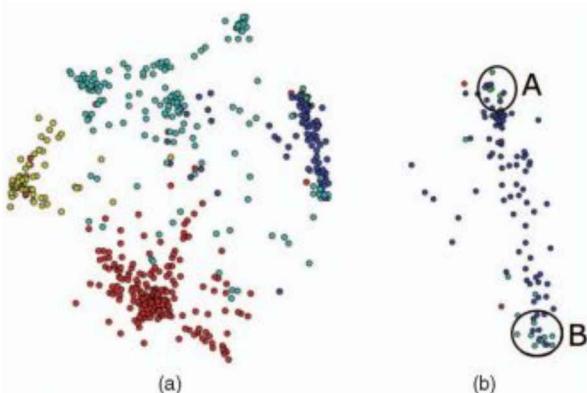


Fig. 2. Projection of a document collection composed of scientific papers in four different areas (colors indicate the areas). (a) Whole map. (b) Zoomed part.

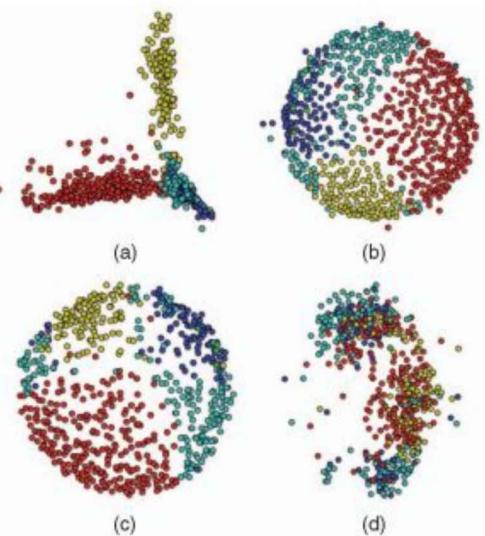
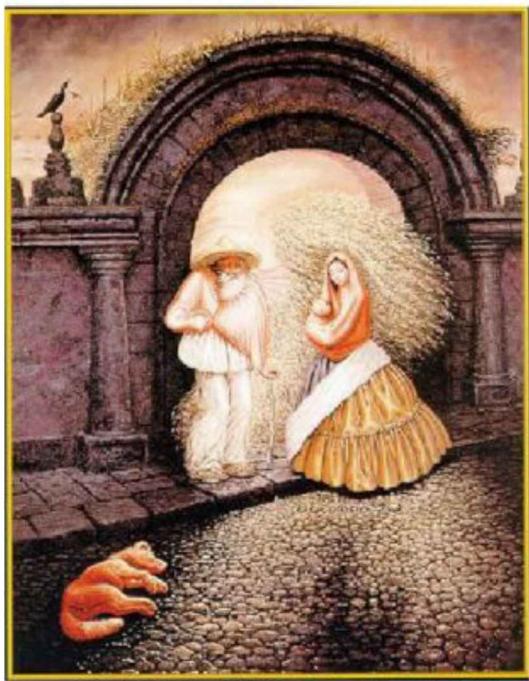


Fig. 11. Examples of projections generated using different techniques for the same data set used in the LSP projection presented in Fig. 2a. (a) PCA. (b) Sammon's mapping. (c) Original FDP model. (d) Approximated FDP model.

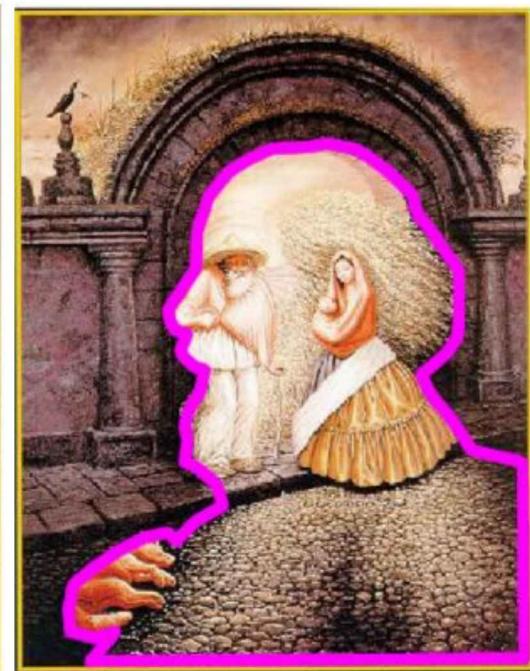
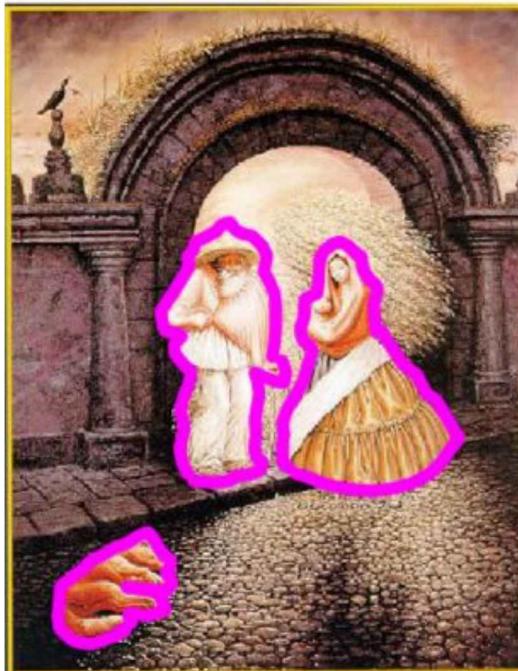
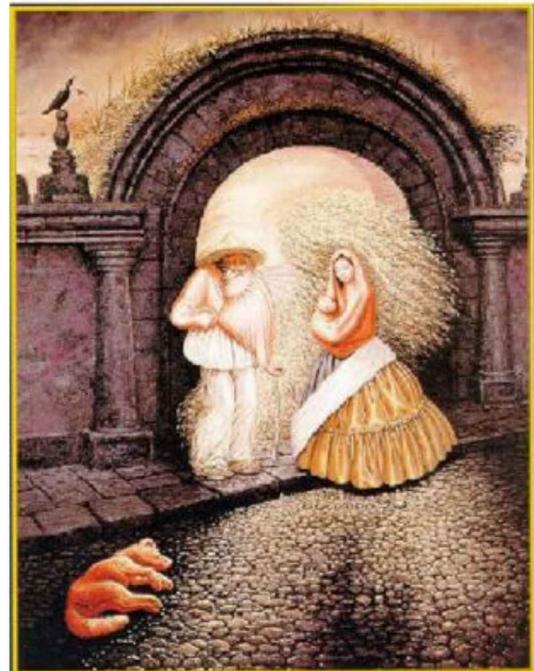
Clustering



Ivar Vargas Belizario

17

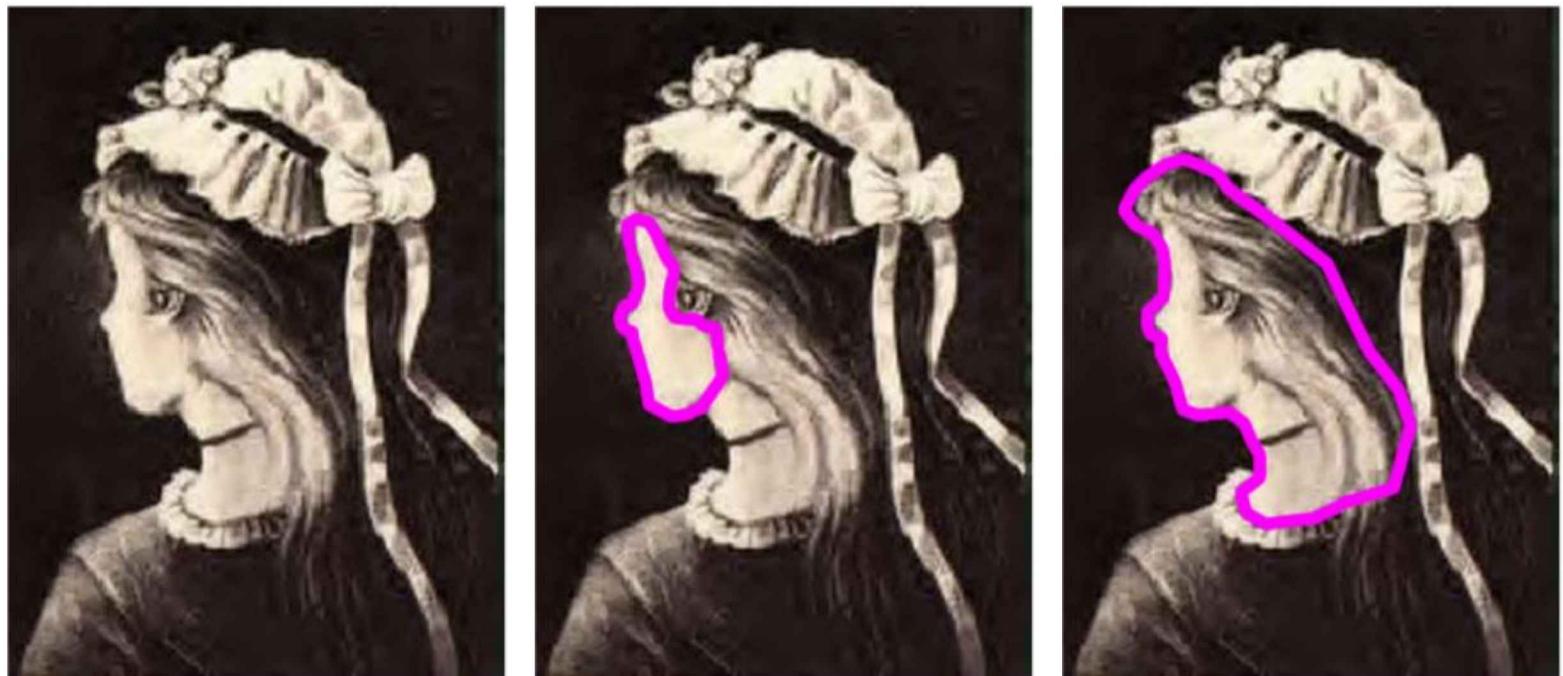
Clustering



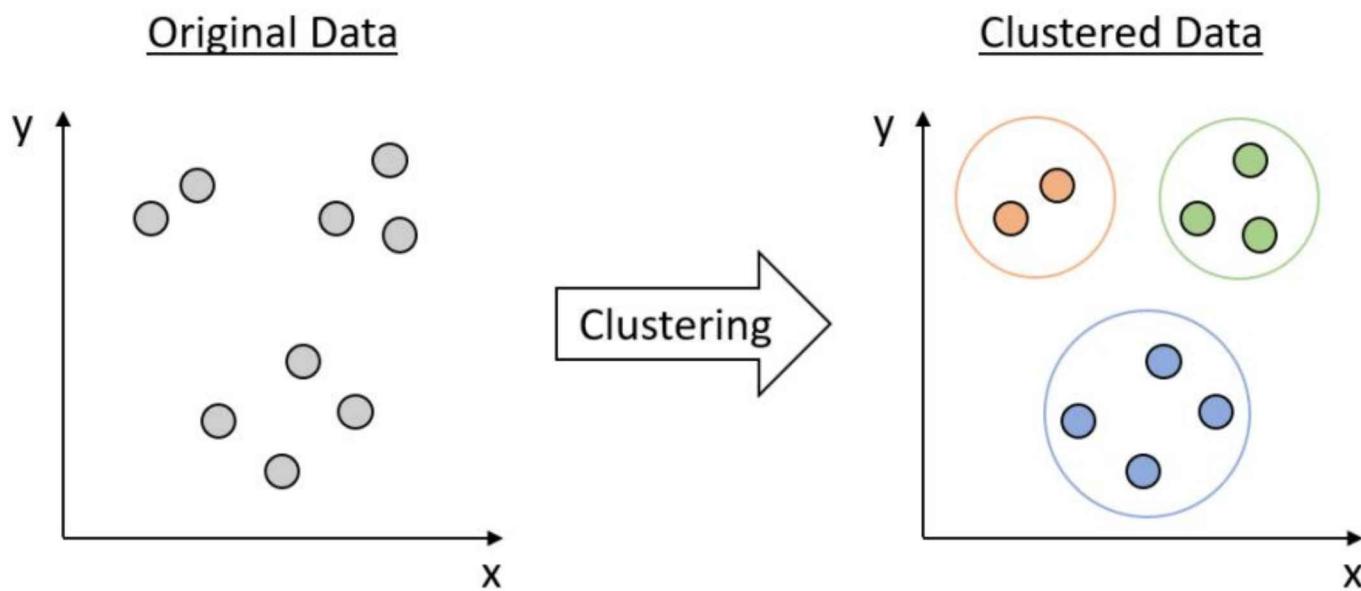
Ivar Vargas Belizario

18

Clustering



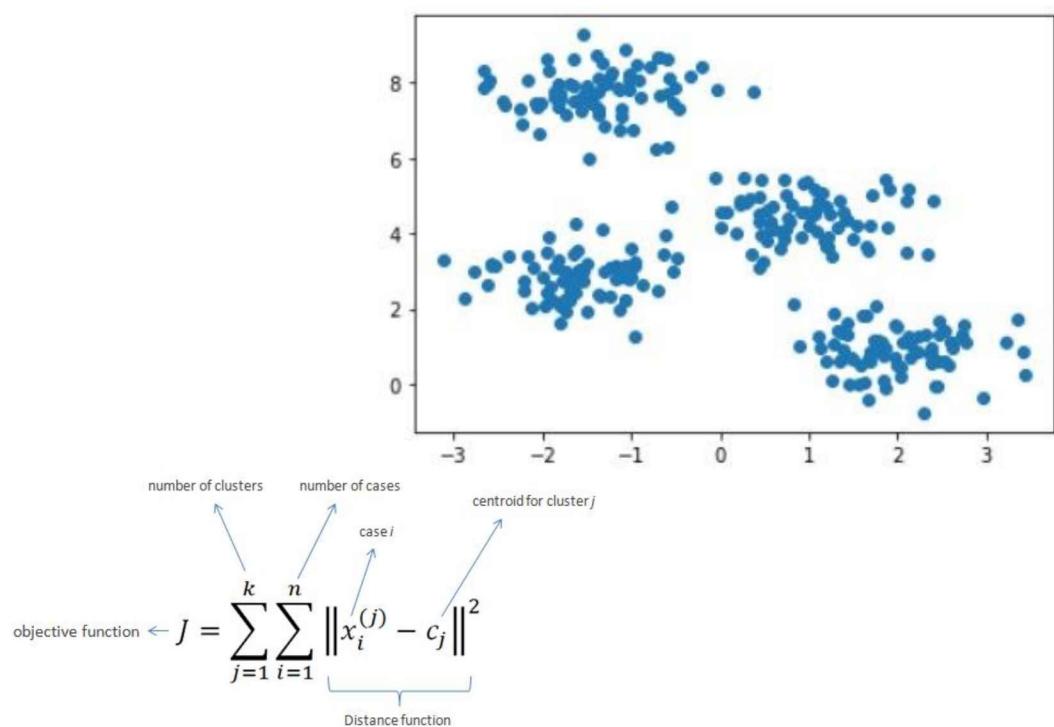
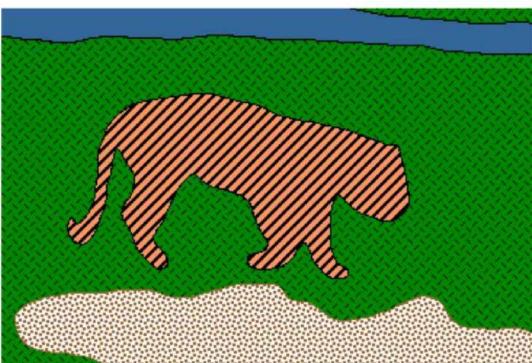
Clustering



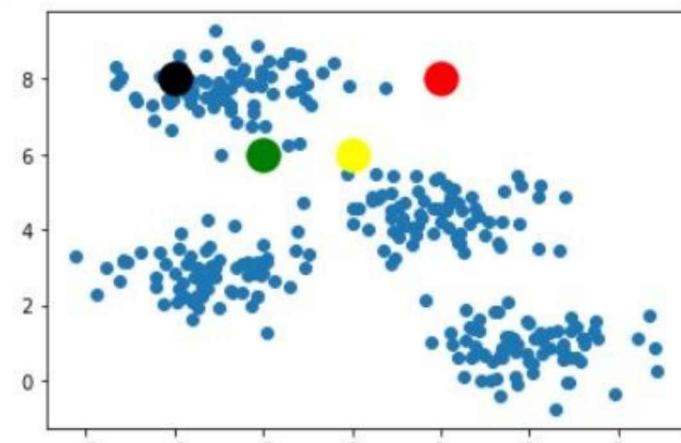
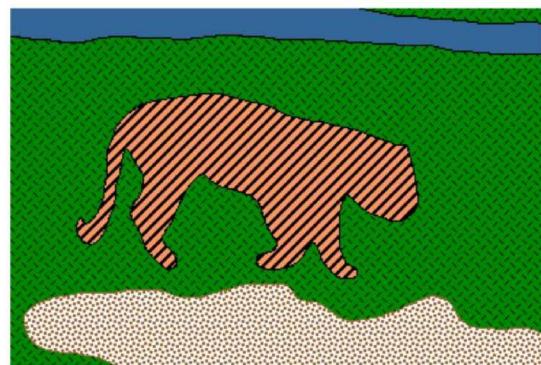
Clustering

- K-means clustering
- Hierarchical clustering
- Spectral clustering
- Mean shift algorithm
- etc.

Clustering - K-means clustering



Clustering - K-means clustering



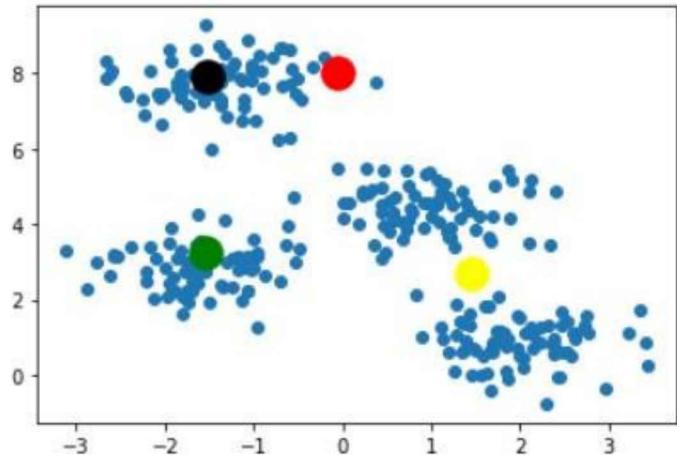
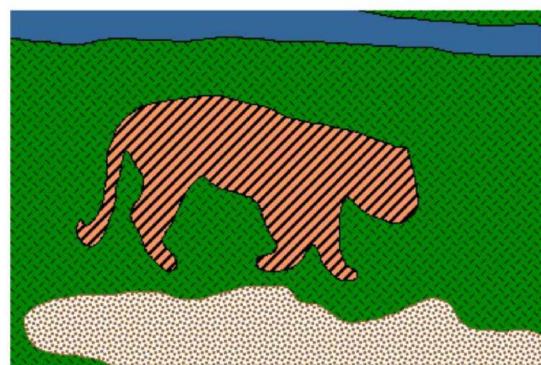
number of clusters number of cases centroid for cluster j

case i

Distance function

$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

Clustering - K-means clustering



number of clusters number of cases centroid for cluster j

case i

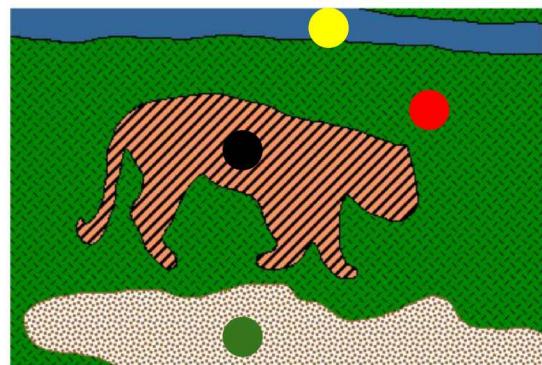
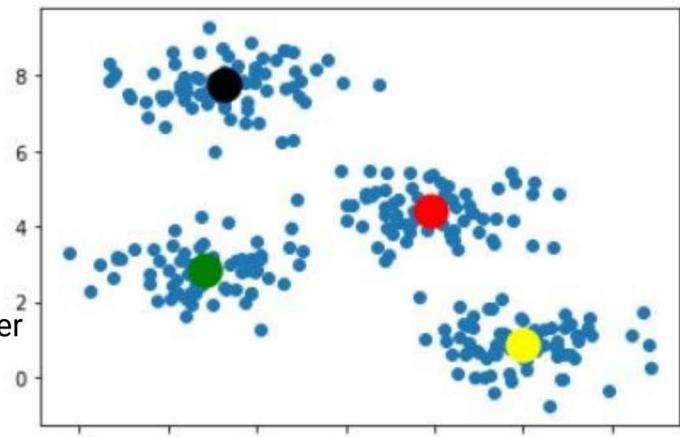
Distance function

$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

Clustering - K-means clustering



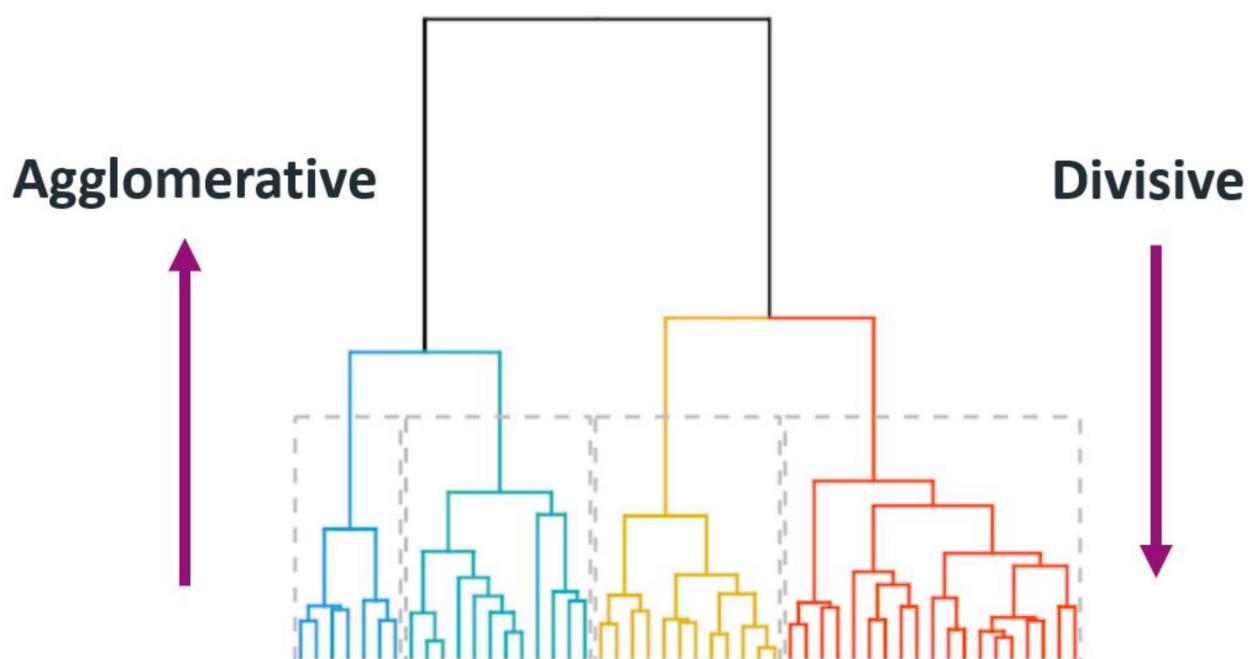
k = is indicated by the user



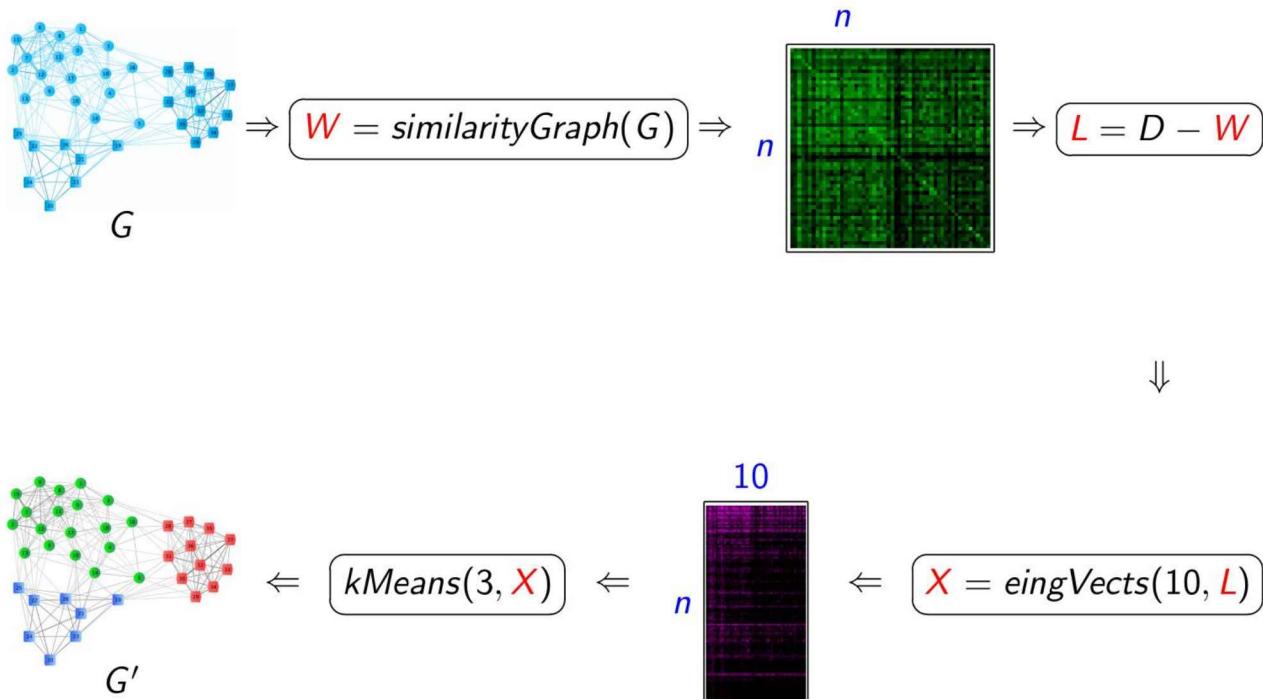
$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

Distance function

Clustering - Hierarchical clustering



Clustering - Spectral clustering



[6] <https://doi.org/10.1109/34.868688>

Classification

Classification

X					y
F1	F2	...	Fn	Label	
0.3	0.5	0.8	0.4	0	
0.2	0.2	0.9	0.3	1	
...	
0.1	0.3	0.7	0.2	0	

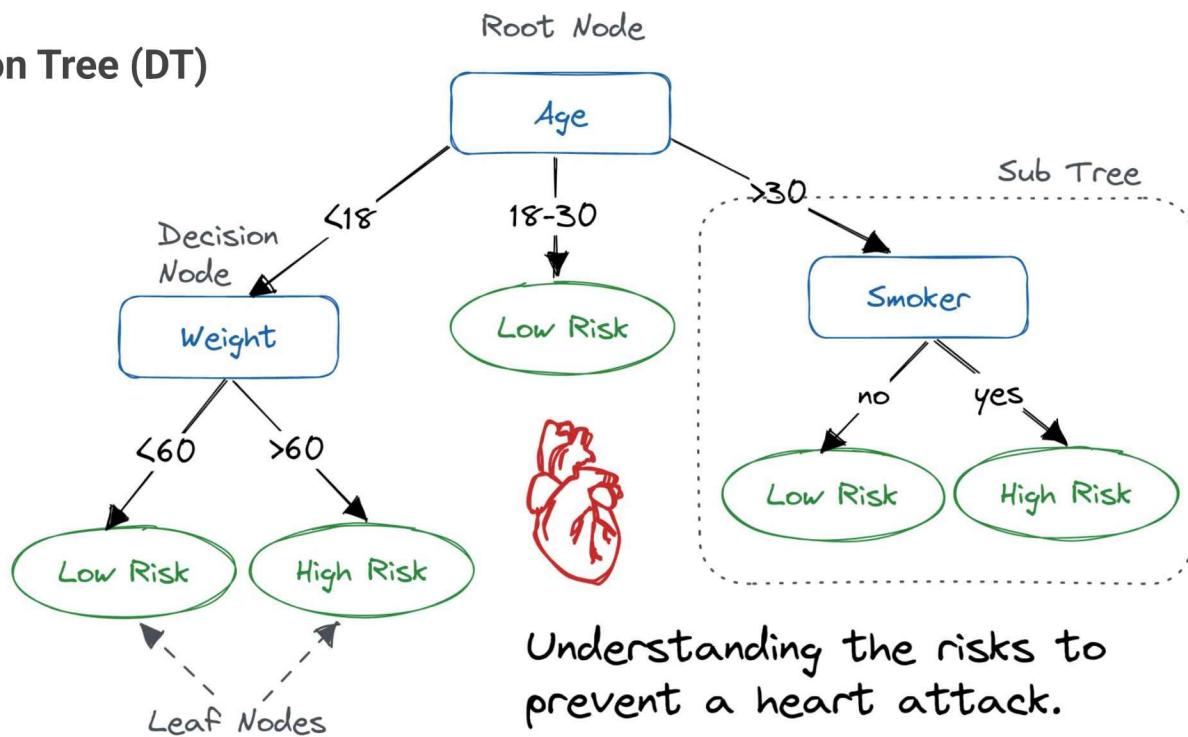
- y : son **datos discretos**, es decir números enteros positivos que representan objetos **finitos**.
 - 0: Cat; 1: Dog
 - 0: Cancer; 1: No cancer
 - 0: Virus; 1: Bacteria; 2: hongos

Classification

- Decision Tree (DT)
- Random Forest (RF)
- Multi-layer Perceptron (MLP)
- etc.

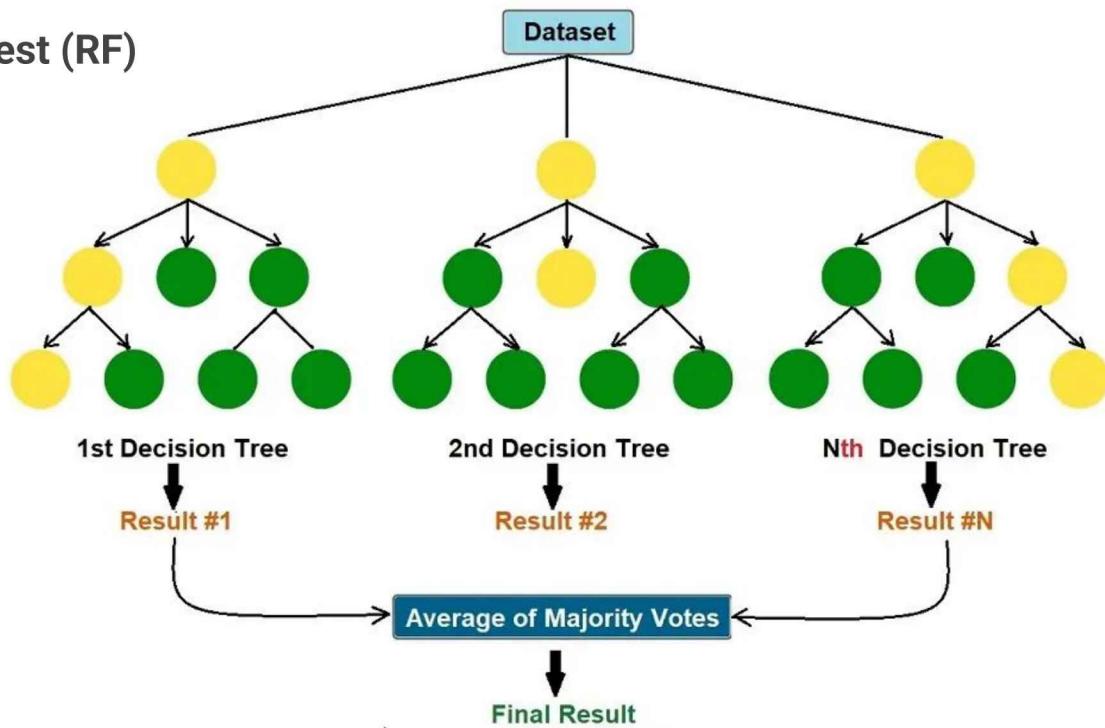
Classification

Decision Tree (DT)



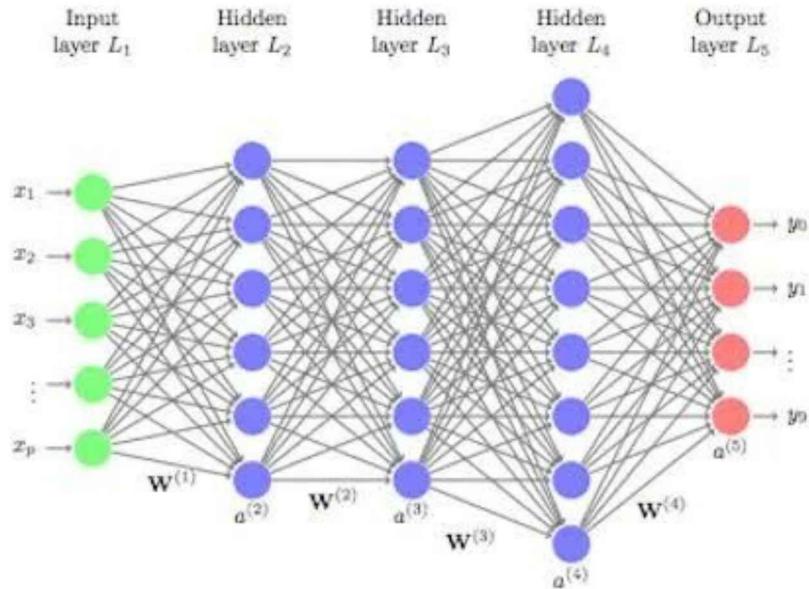
Classification

Random Forest (RF)



Classification

Multi-layer Perceptron (MLP)



Regression

Regression

F1	F2	...	Fn	Label
0.3	0.5	0.8	0.4	0.5
0.2	0.2	0.9	0.3	0.4
...
0.1	0.3	0.7	0.2	0.3

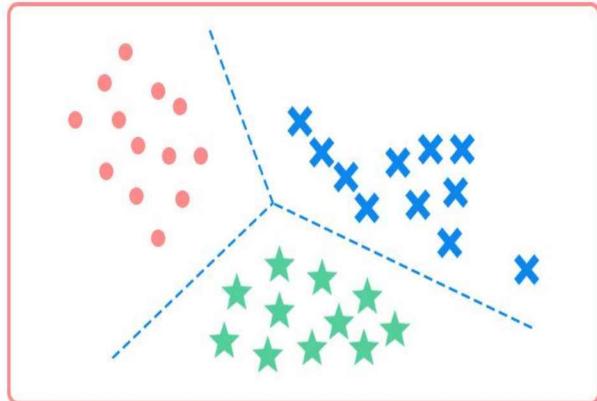
- Y: son datos **continuos** dentro de un **rango infinito** de valores.
Generalmente son representados por valores reales.
 - Temperatura
 - Consumo de energía
 - Latitud, Longitud

Regression

- Decision Tree Regression (DTR)
- Random Forest Regression(RFR)
- Multi-layer Perceptron Regression (MLPR)
- etc.

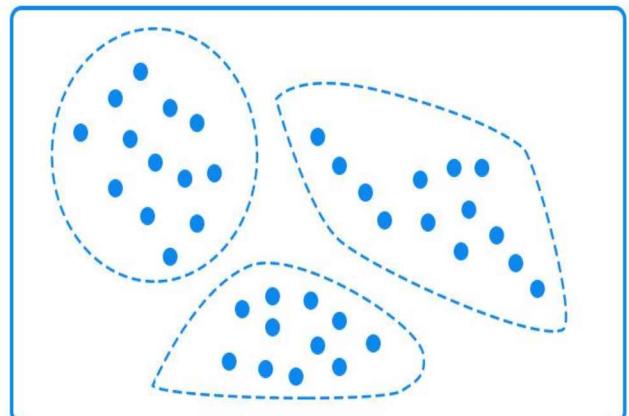
Classification x Clustering

Classification



Supervised learning

Clustering

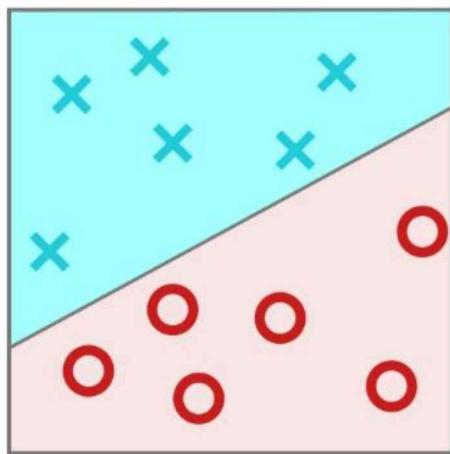


Unsupervised learning

Se determina la clase (discreto)

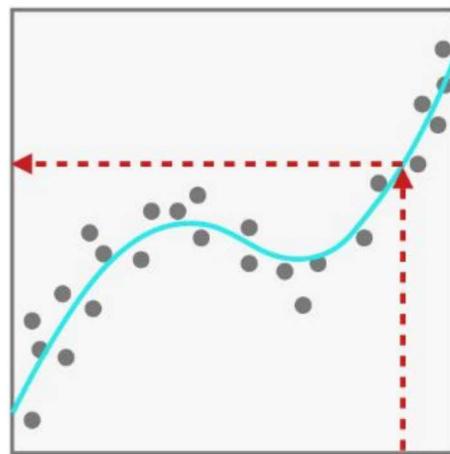
Se determina el agrupamiento (discreto)

Classification x Regression



Classification

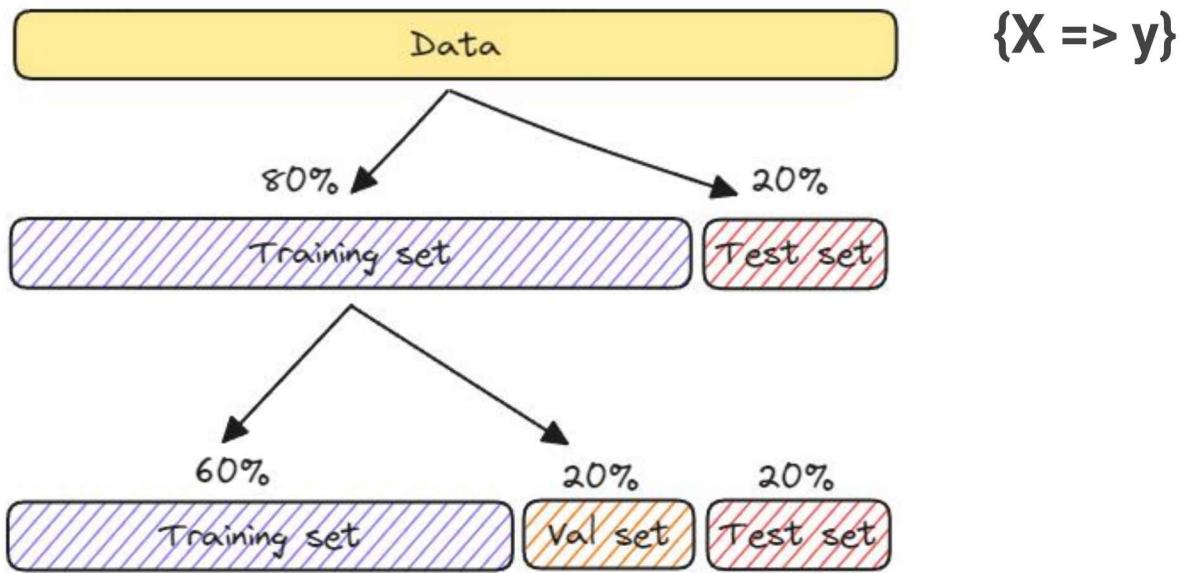
Se determina la clase (discreto)



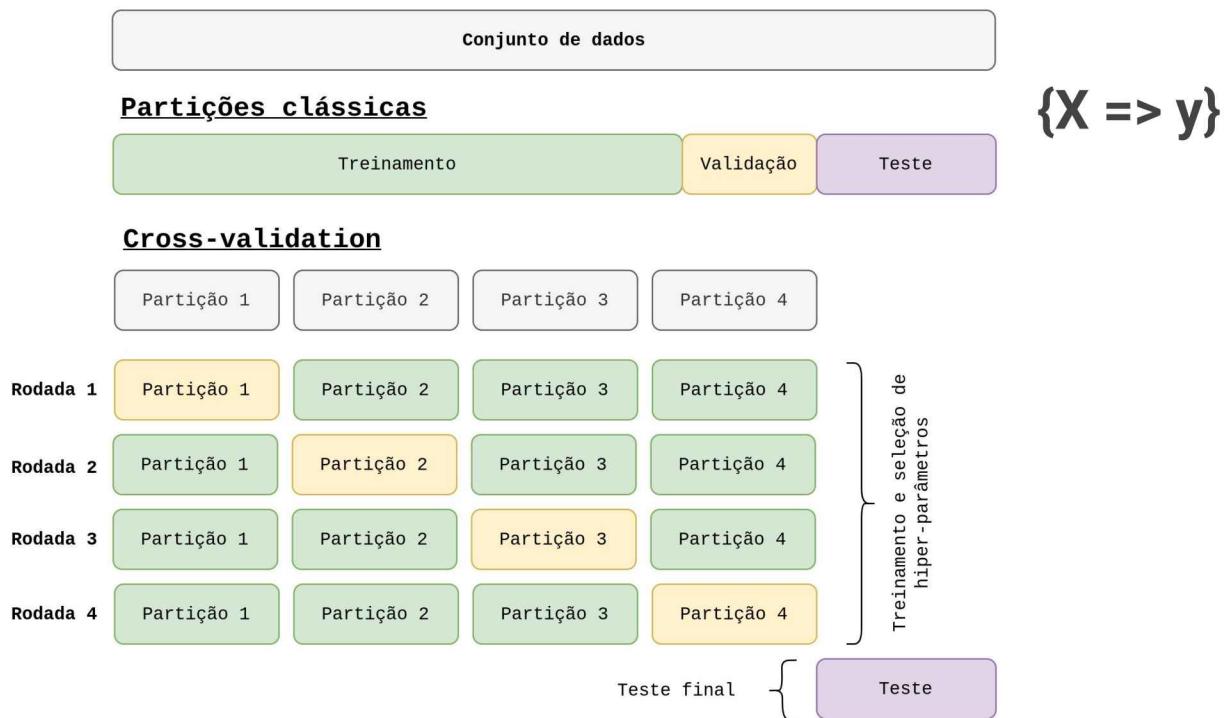
Regression

Se determina valores continuos

Classification and Regression

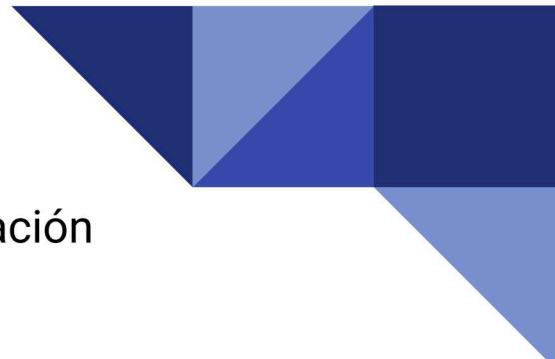


Classification and Regression





Universidad Nacional del Altiplano
Escuela de Posgrado
Doctorado en Ciencias de la Computación



Data Mining

Unit 1. Structured data mining

Gracias

Prof. Dr. Ivar Vargas Belizario

ivargasbelizario@gmail.com

2025 - I