



# Universidad Nacional del Altiplano

## Escuela de Posgrado

### Doctorado en Ciencias de la Computación

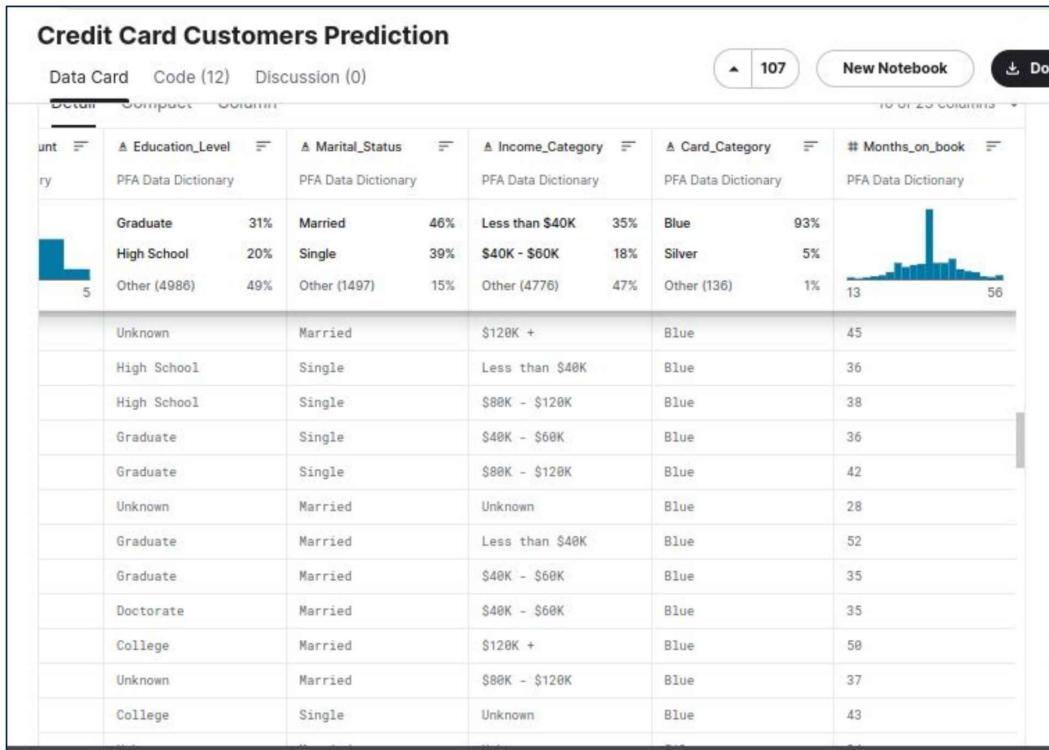
## Minería de Datos

Prof. Dr. Ivar Vargas Belizario

ivargasbelizario@gmail.com

2024 - I

## Structured data



# Structured data

High-dimensional Data  
Curse of Dimensionality

$X^{n \times m}$

	F1	F2	...	Fm	Label
n	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

3

# 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

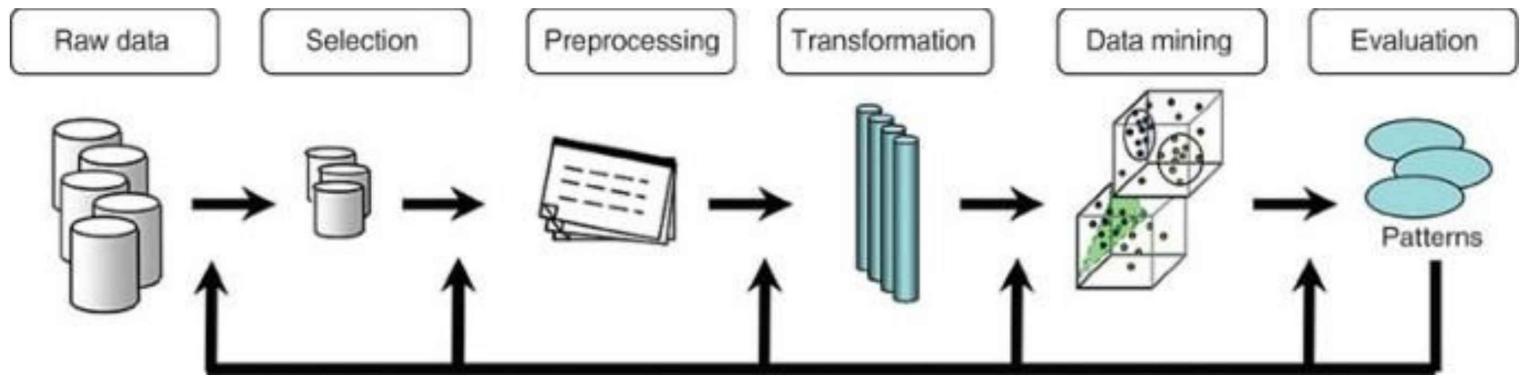
### 2. Complejidad:

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

Ivar Vargas Belizario

4

# Structured data (pipeline)

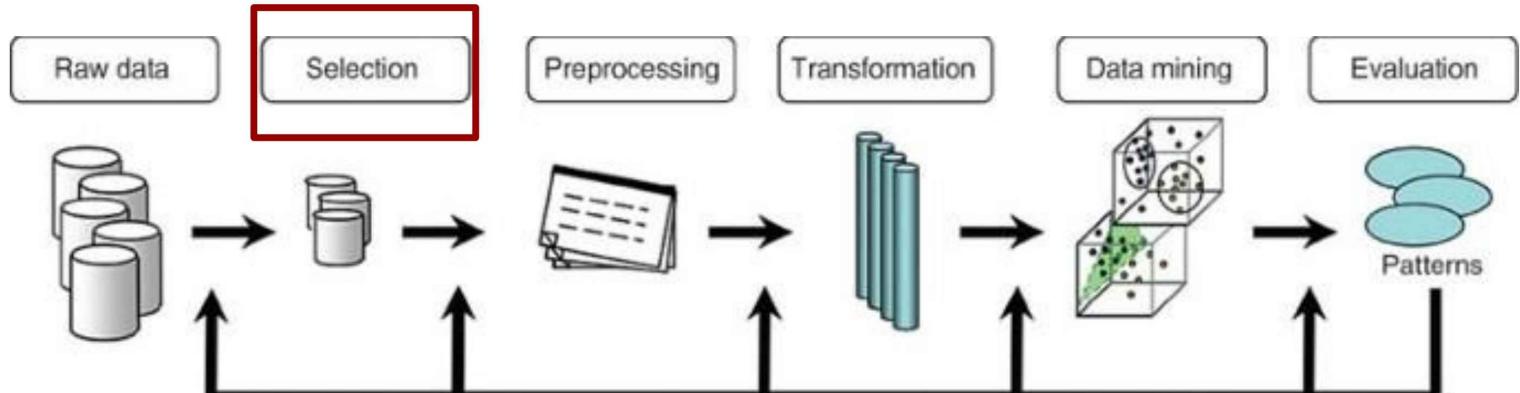


[4] [https://doi.org/10.1007/978-1-4899-7993-3\\_1134-2](https://doi.org/10.1007/978-1-4899-7993-3_1134-2)

Ivar Vargas Belizario

5

# Structured data (pipeline)

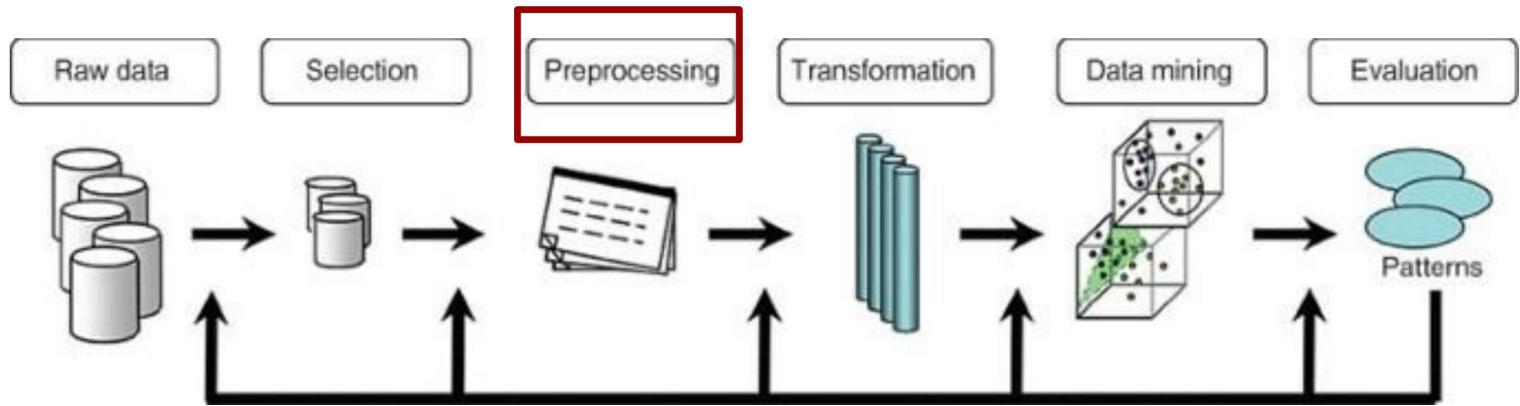


- Identifica el objetivo (problema)
- Selección con patrones relevantes

Ivar Vargas Belizario

6

# Structured data (pipeline)

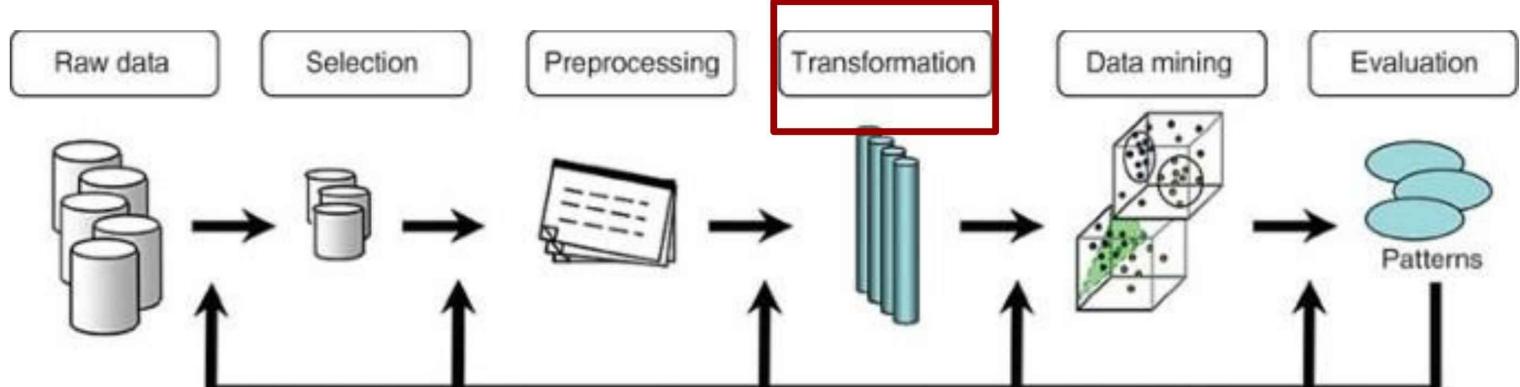


- Limpieza
- Aumentar la calidad de los datos:
  - Completar datos faltantes
  - Eliminar instancias duplicadas

Ivar Vargas Belizario

7

# Structured data (pipeline)

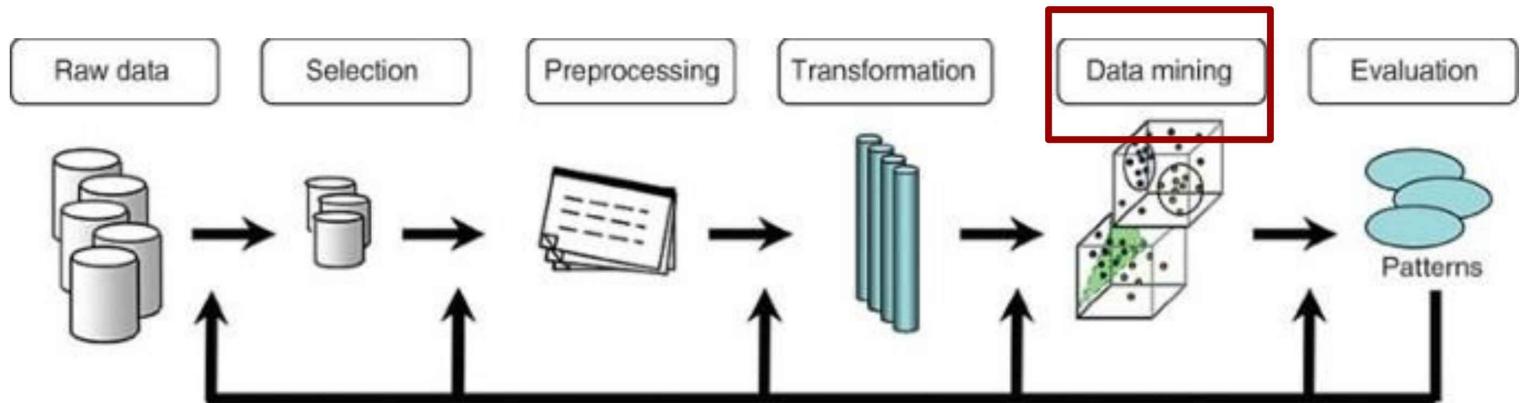


- Eliminar atributos correlacionados e irrelevantes.
- Se crea una nueva definición de atributos com mayor significancia.

Ivar Vargas Belizario

8

# Structured data (pipeline)

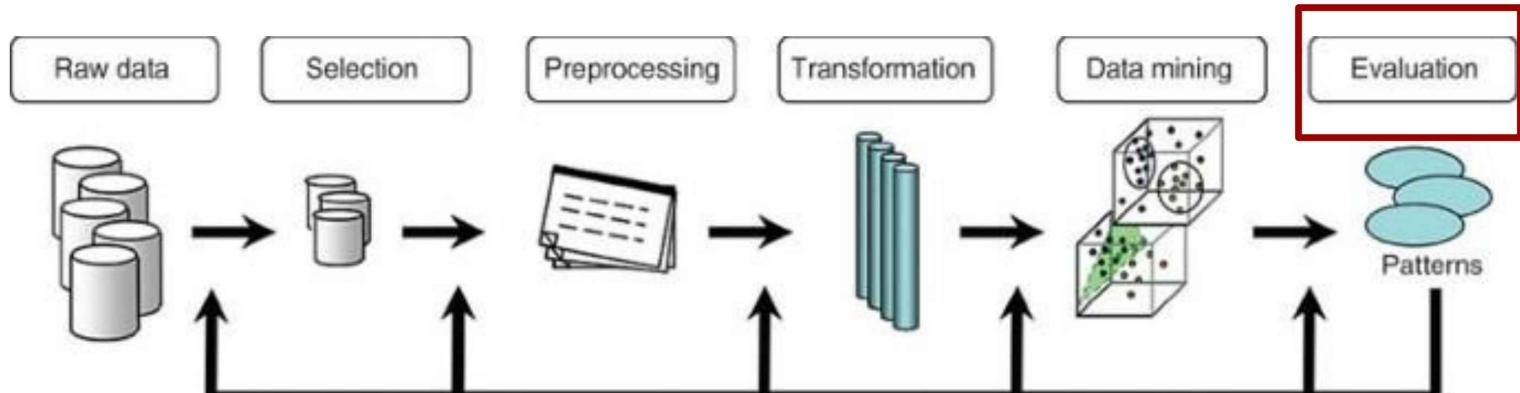


- Selecciona el tipo de algoritmo de minería de datos identificado en la etapa de **Selection**.
- Aprender patrones según:
  - Clasificación,
  - Agrupamiento
  - Regresión,
  - Asociación

Ivar Vargas Belizario

9

# Structured data (pipeline)

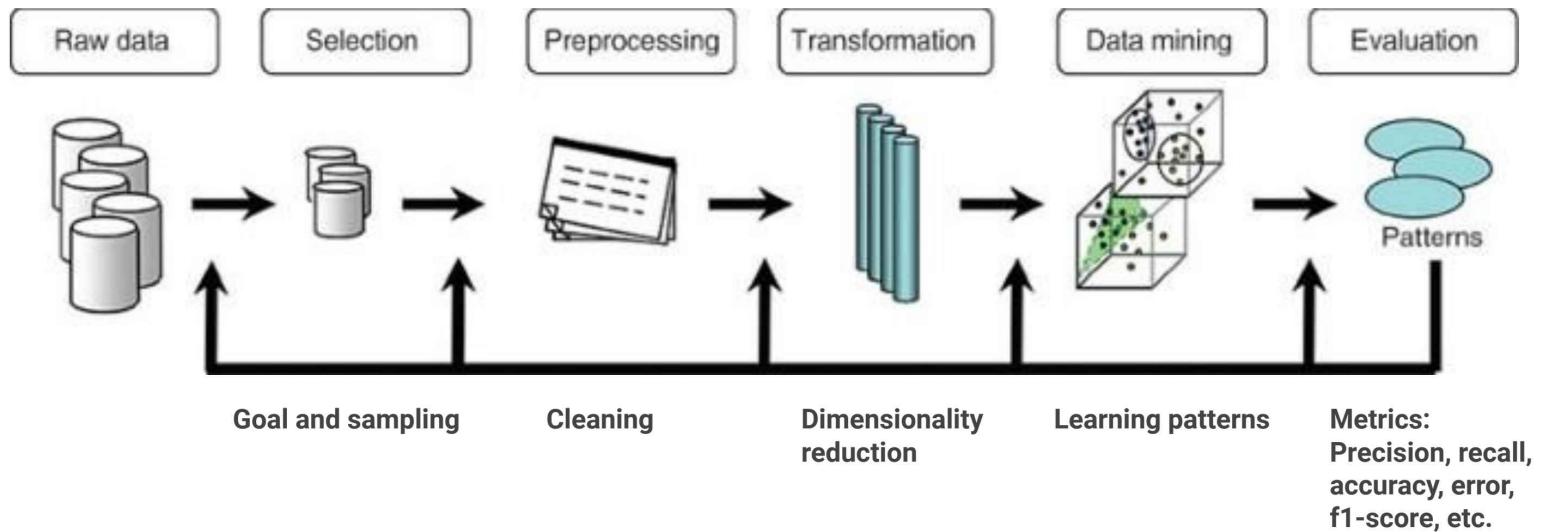


- Evalúa lo aprendido

Ivar Vargas Belizario

10

# Structured data



Ivar Vargas Belizario

11

## Dimensionality reduction

$X^{n \times m}$					$y$	$X^{n \times p}$					$y$
$\brace{F_1, F_2, \dots, F_m}$					$\brace{Label}$	$\brace{F_1, F_2, \dots, F_p}$					$\brace{Label}$
0.3	0.5	0.8	0.4	0.5	n	0.3	0.5	0.8	0.4	0.5	...
0.2	0.2	0.9	0.3	0.4		0.2	0.2	0.9	0.3	0.4	...
...	...	...	...	...		...	...	...	...	...	...
0.1	0.3	0.7	0.2	0.3		0.1	0.3	0.7	0.2	0.3	

donde,  $p < m$

Ivar Vargas Belizario

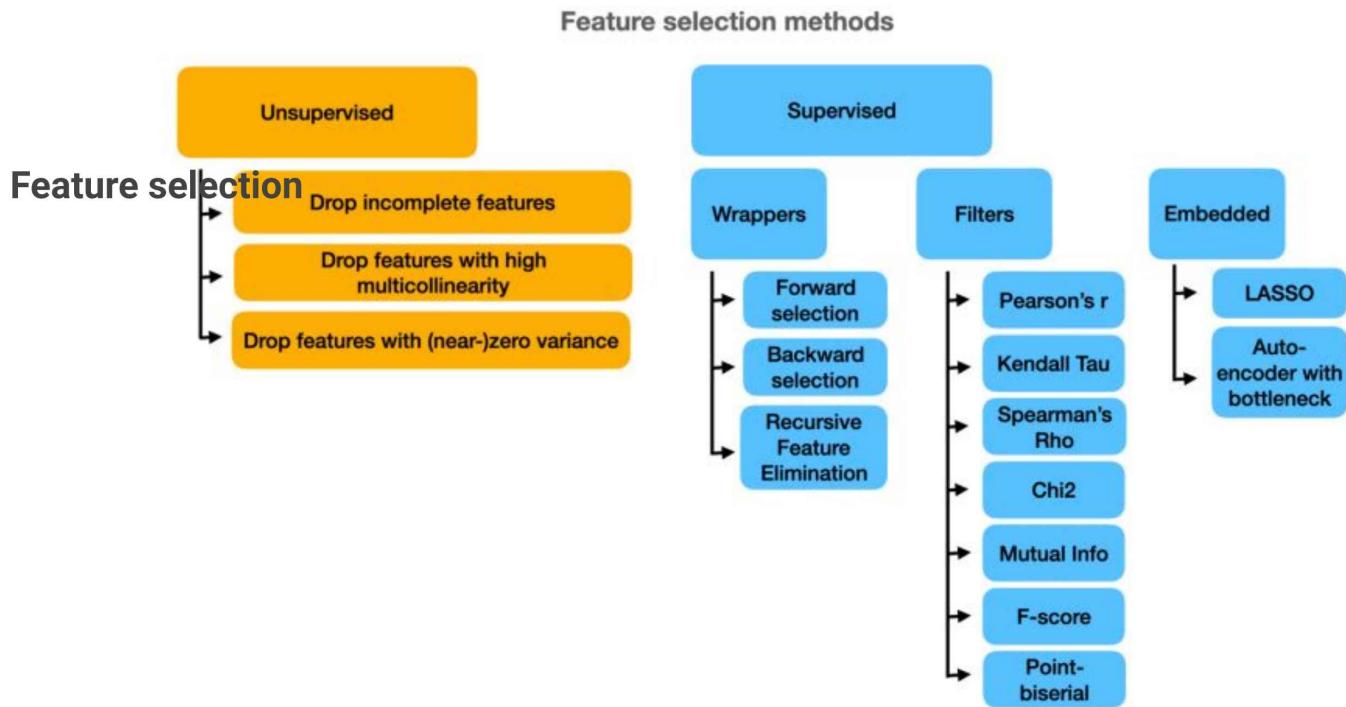
12

# 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

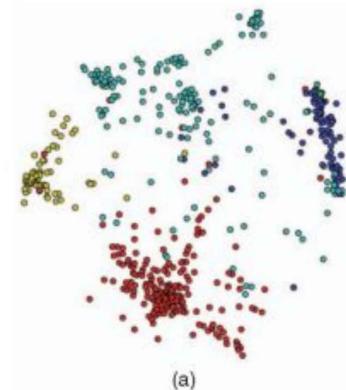


Fig. 2. Projection of a document collection composed of scientific papers in four different areas (colors indicate the areas). (a) Whole map. (b) Zoomed part.

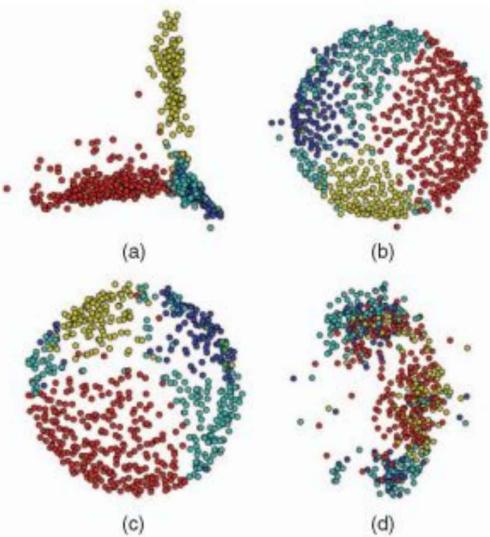


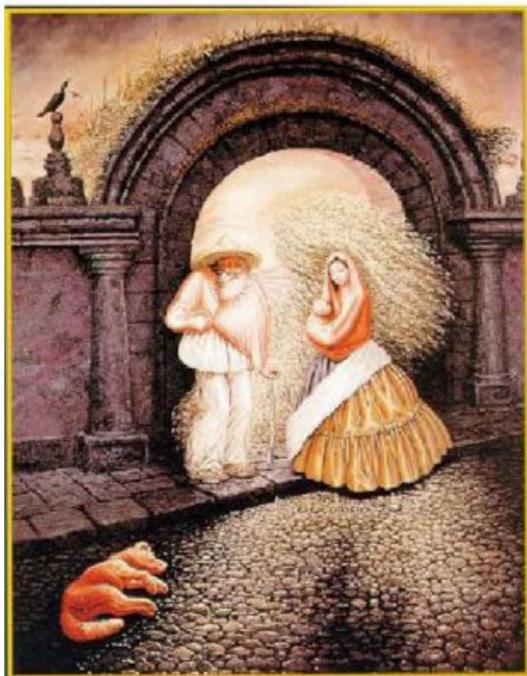
Fig. 11. Examples of projections generated using different techniques for the same data set used in the LSP projection presented in Fig. 2a. (a) PCA. (b) Sammon's mapping. (c) Original FDP model. (d) Approximated FDP model.

[5] <https://doi.org/10.1109/TVCG.2007.70443>

Ivar Vargas Belizario

15

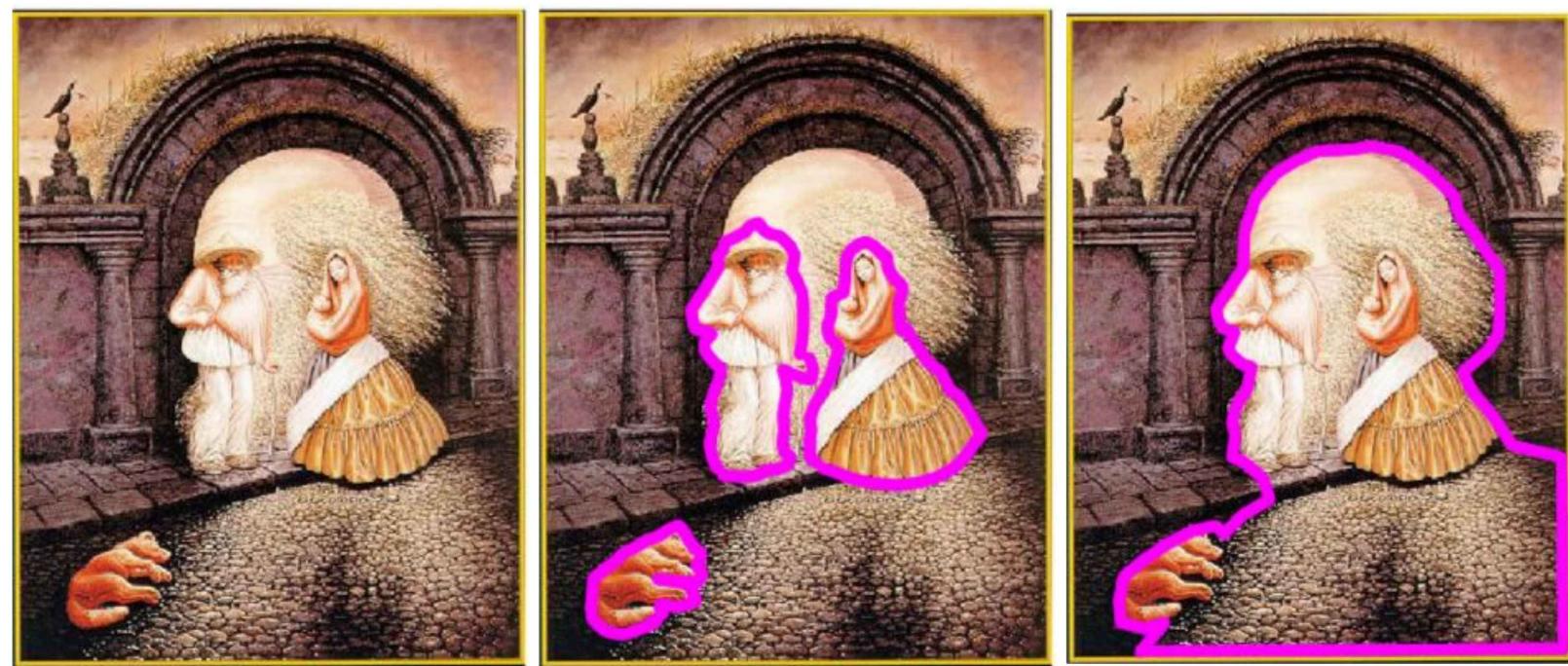
## Clustering



Ivar Vargas Belizario

16

# Clustering



Ivar Vargas Belizario

17

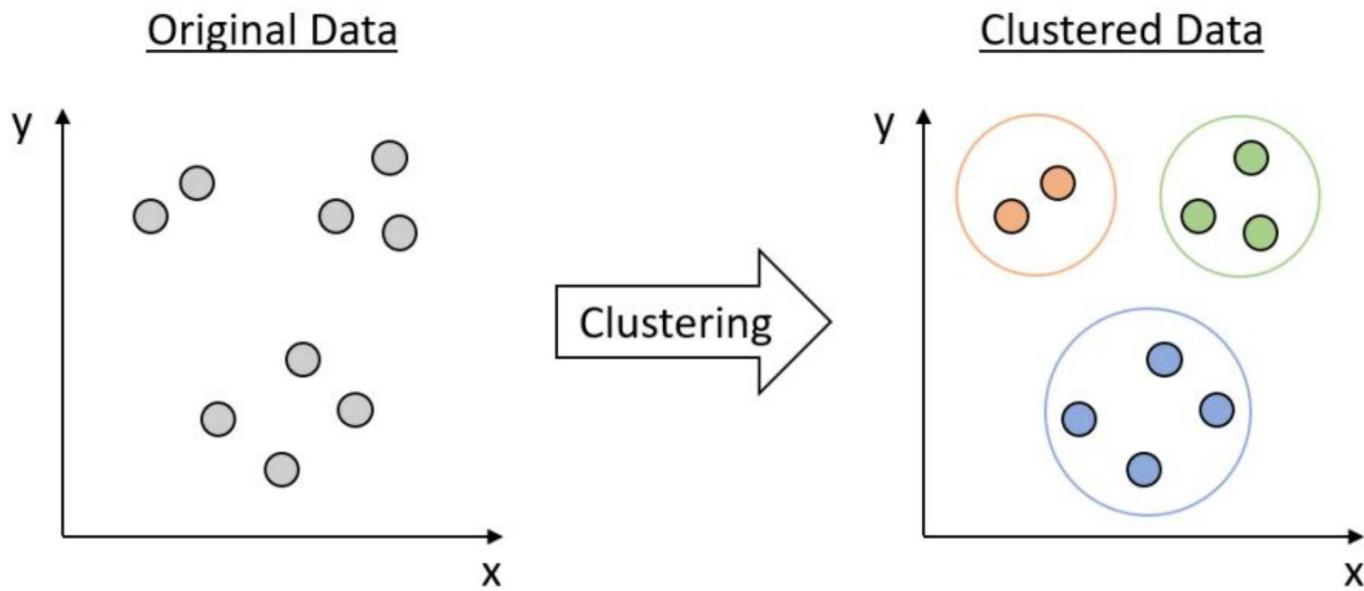
# Clustering



Ivar Vargas Belizario

18

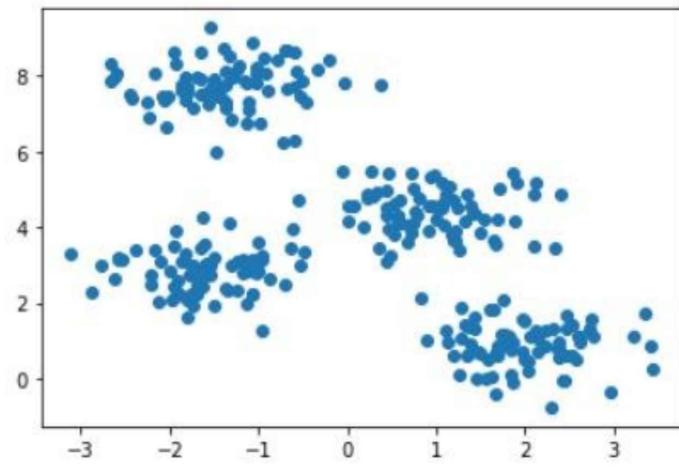
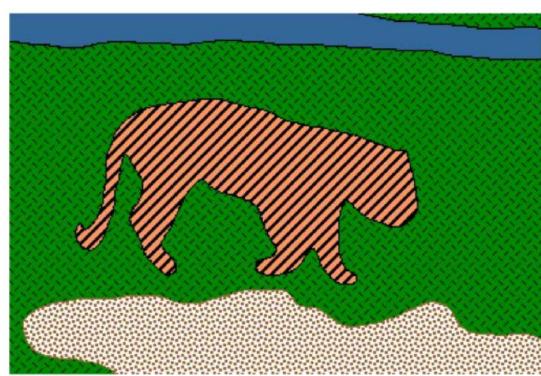
# Clustering



# Clustering

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

# Clustering - K-means clustering



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

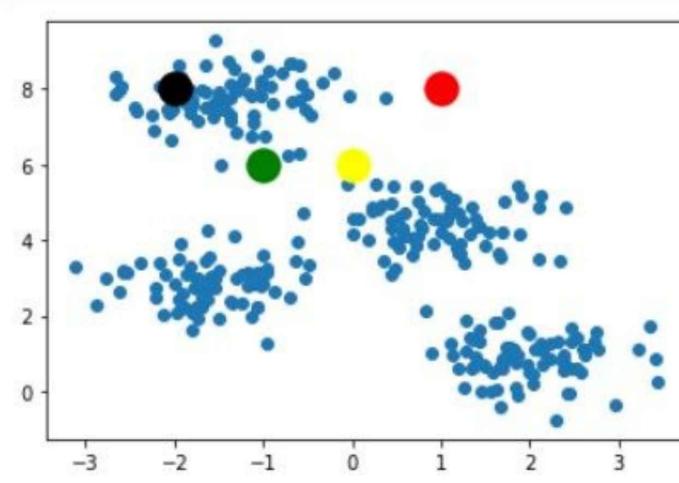
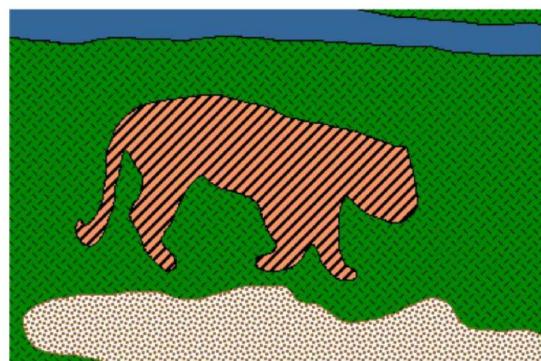
number of clusters      number of cases  
case  $i$       centroid for cluster  $j$

Distance function

Ivar Vargas Belizario

21

# Clustering - K-means clustering



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

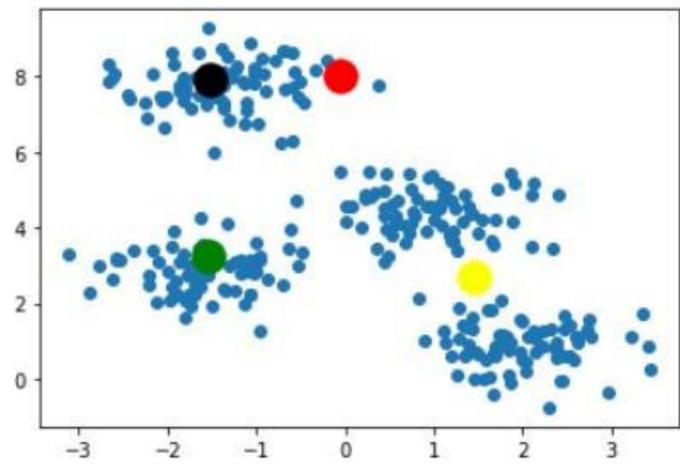
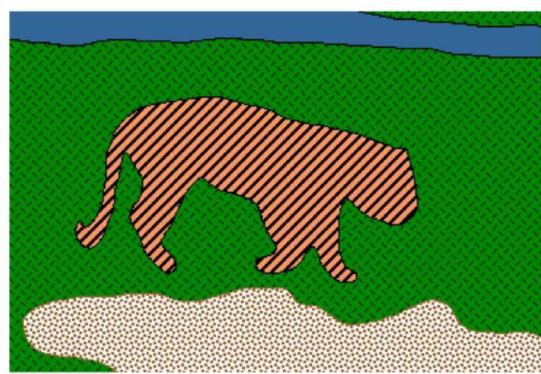
number of clusters      number of cases  
case  $i$       centroid for cluster  $j$

Distance function

Ivar Vargas Belizario

22

# Clustering - K-means clustering



number of clusters      number of cases      centroid for cluster  $j$   
case  $i$

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

Distance function

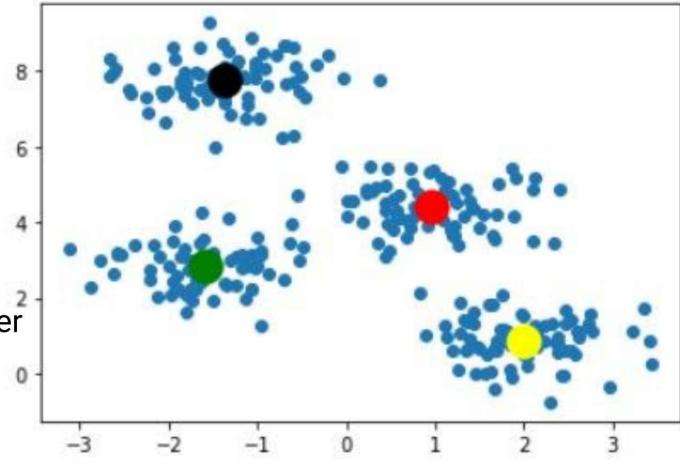
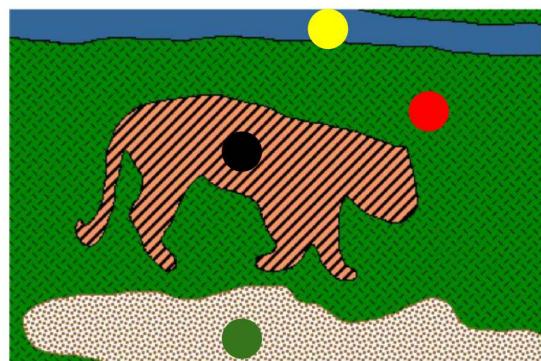
Ivar Vargas Belizario

23

# Clustering - K-means clustering



$k$  = is indicated by the user



number of clusters      number of cases      centroid for cluster  $j$   
case  $i$

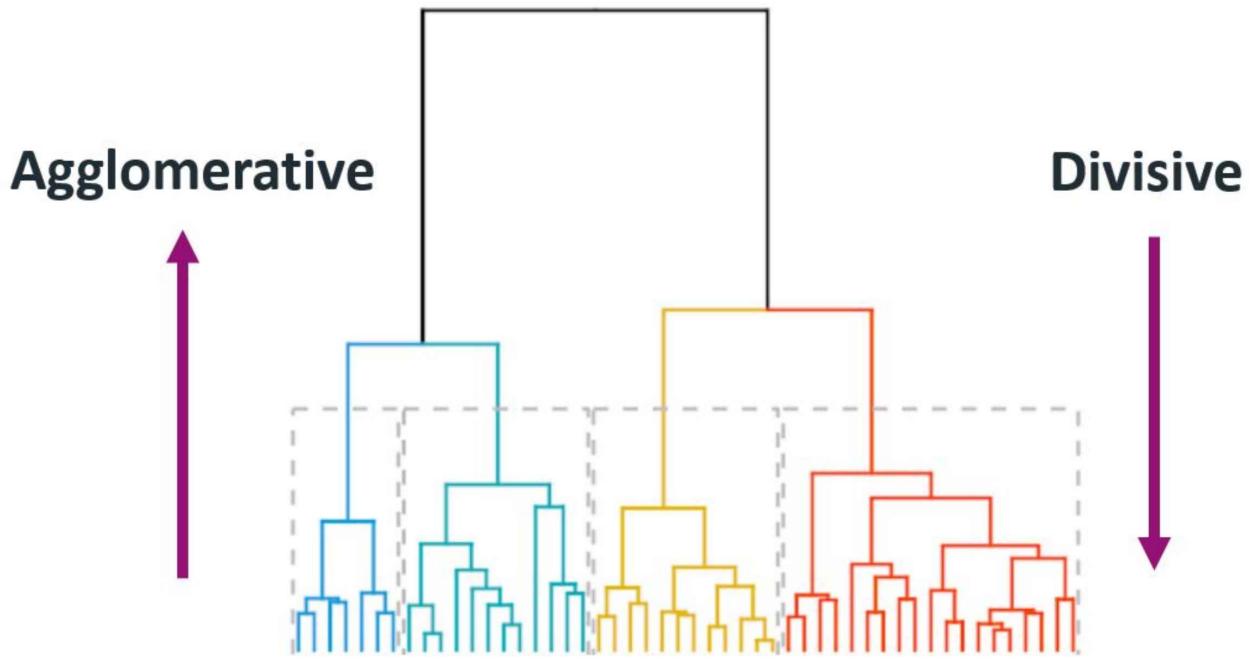
$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

Distance function

Ivar Vargas Belizario

24

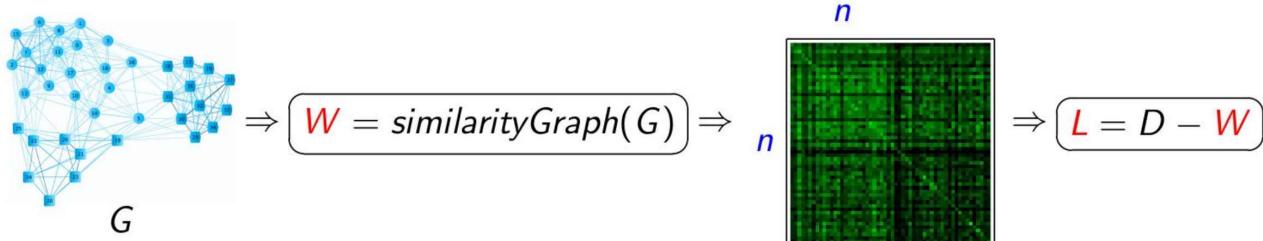
# Clustering - Hierarchical clustering



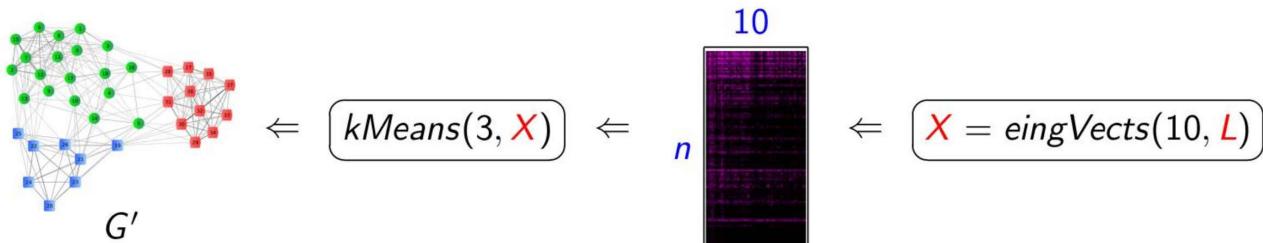
Ivar Vargas Belizario

25

# Clustering - Spectral clustering



⇓



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

Ivar Vargas Belizario

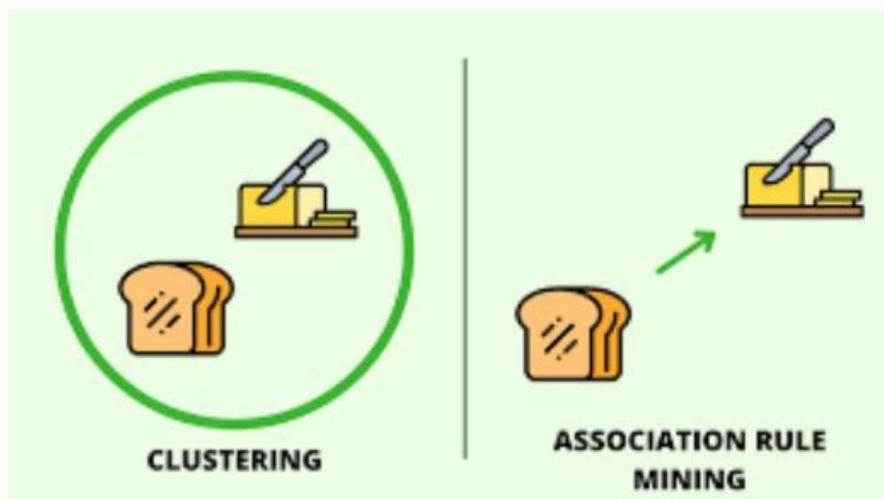
26

# Association

{Bread, Egg} => {Milk}

Itemset = {Bread, Egg, Milk}

# Association



# Classification

## Classification

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

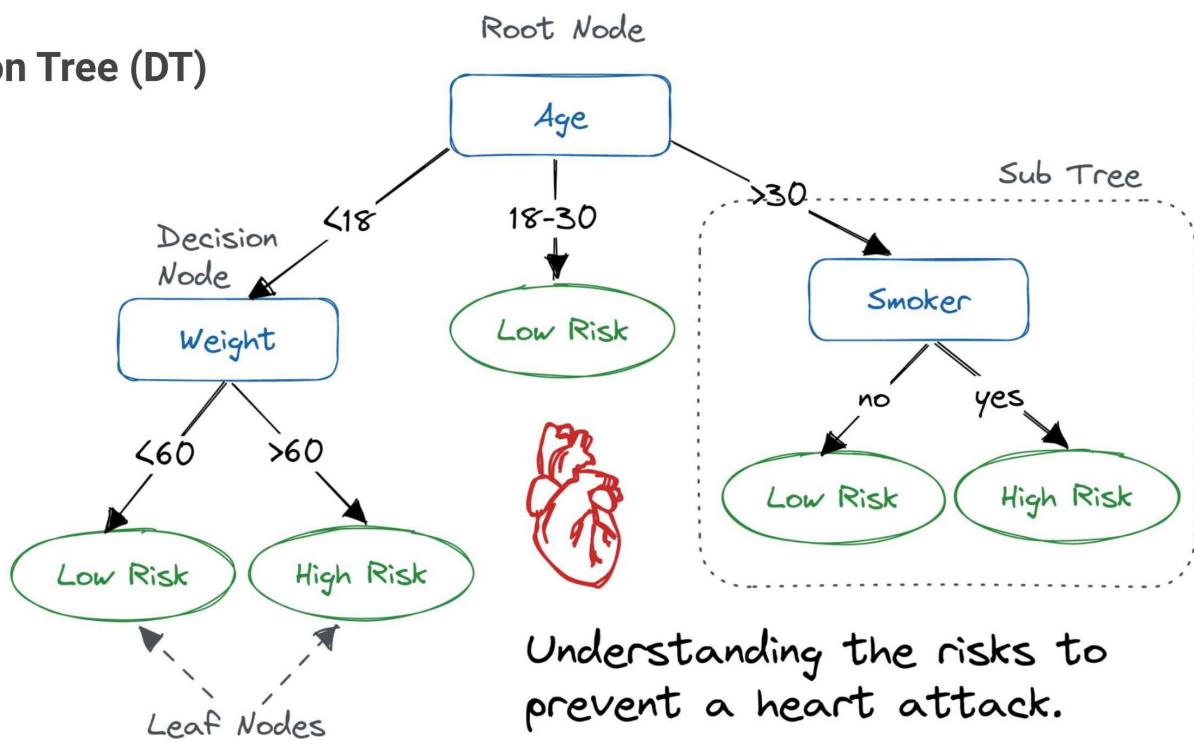
- $y$ : son **datos discretos**, es decir número enteros positivos que representan objetos **finitos**.
  - 0: Cat; 1: Dog
  - 0: Cancer; 1: No cancer
  - 0: Virus; 1: Bacteria; 2: hongos

# Classification

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

# Classification

## Decision Tree (DT)

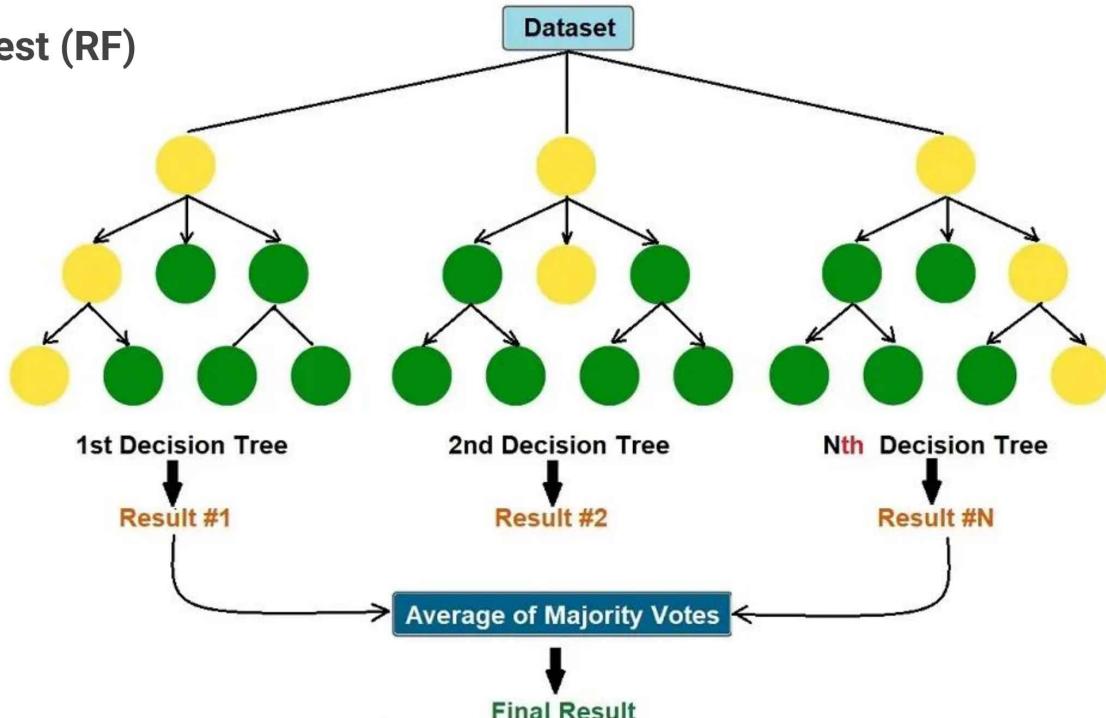


Ivar Vargas Belizario

31

## Classification

## Random Forest (RF)

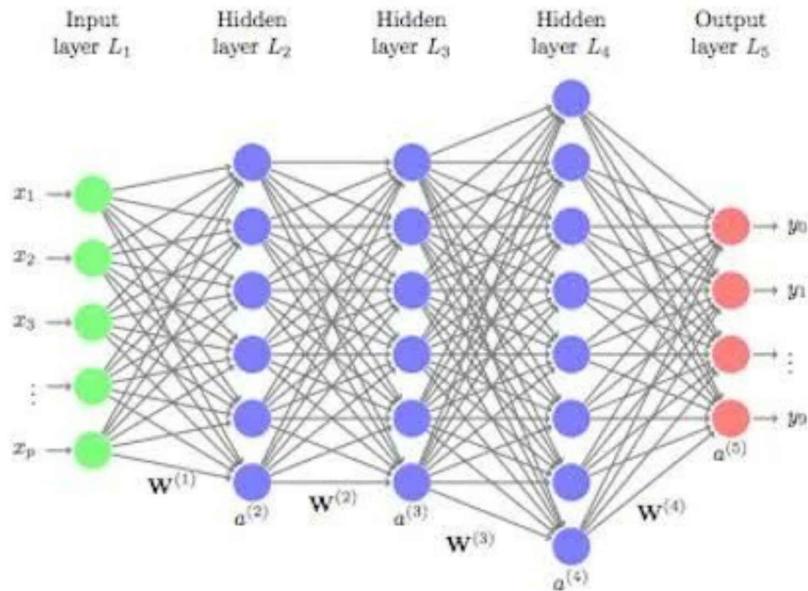


Ivar Vargas Belizario

32

# Classification

## Multi-layer Perceptron (MLP)



# Regression

## Regression

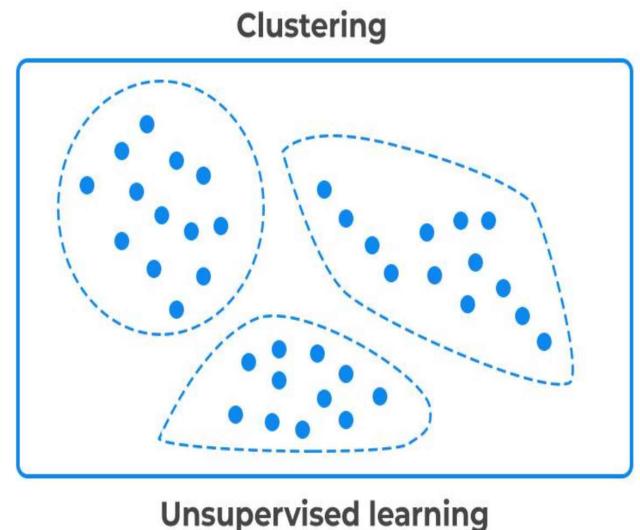
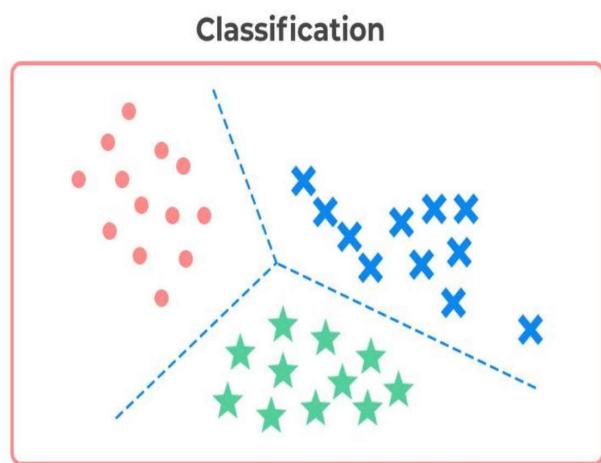
X				y
F1	F2	...	Fn	Label
0.3	0.5	0.8	0.4	0.5
0.2	0.2	0.9	0.3	0.4
...	...	...	...	...
0.1	0.3	0.7	0.2	0.3

- 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.

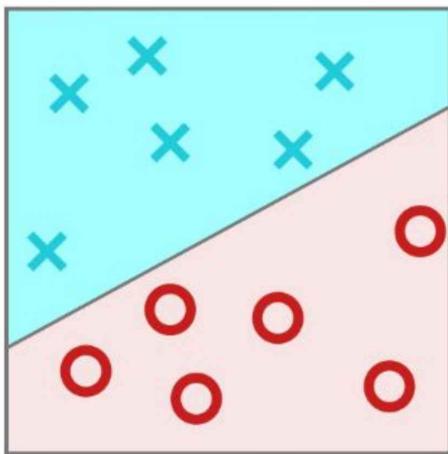
## Classification x Clustering



**Se determina la clase (discreto)**

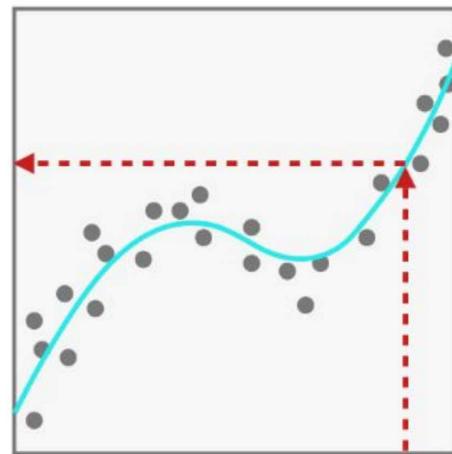
**Se determina el agrupamiento (discreto)**

# Classification x Regression



Classification

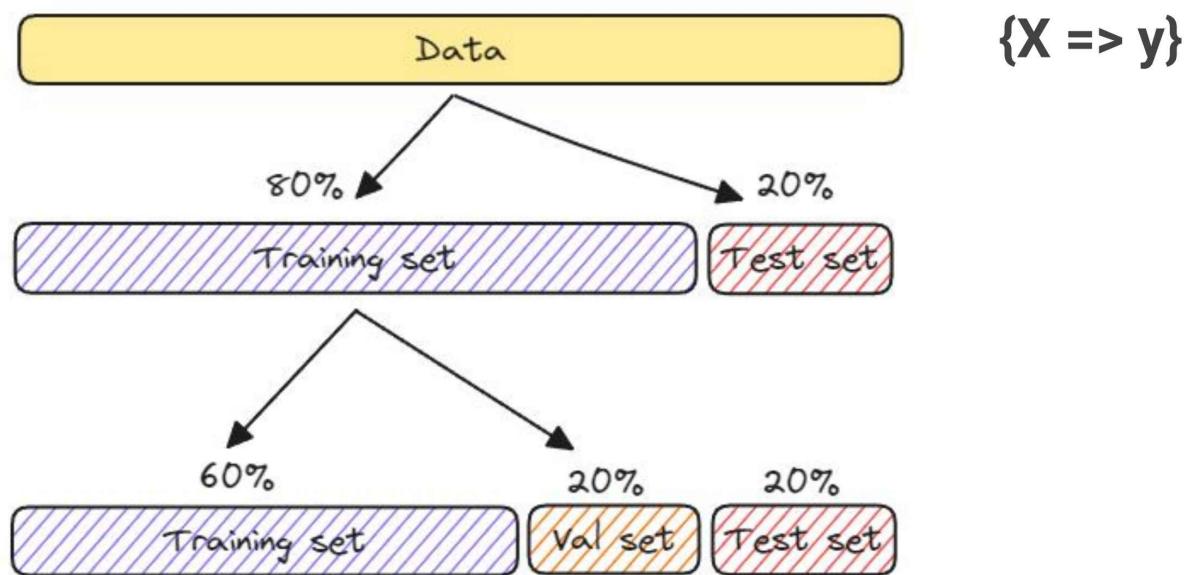
Se determina la clase (discreto)



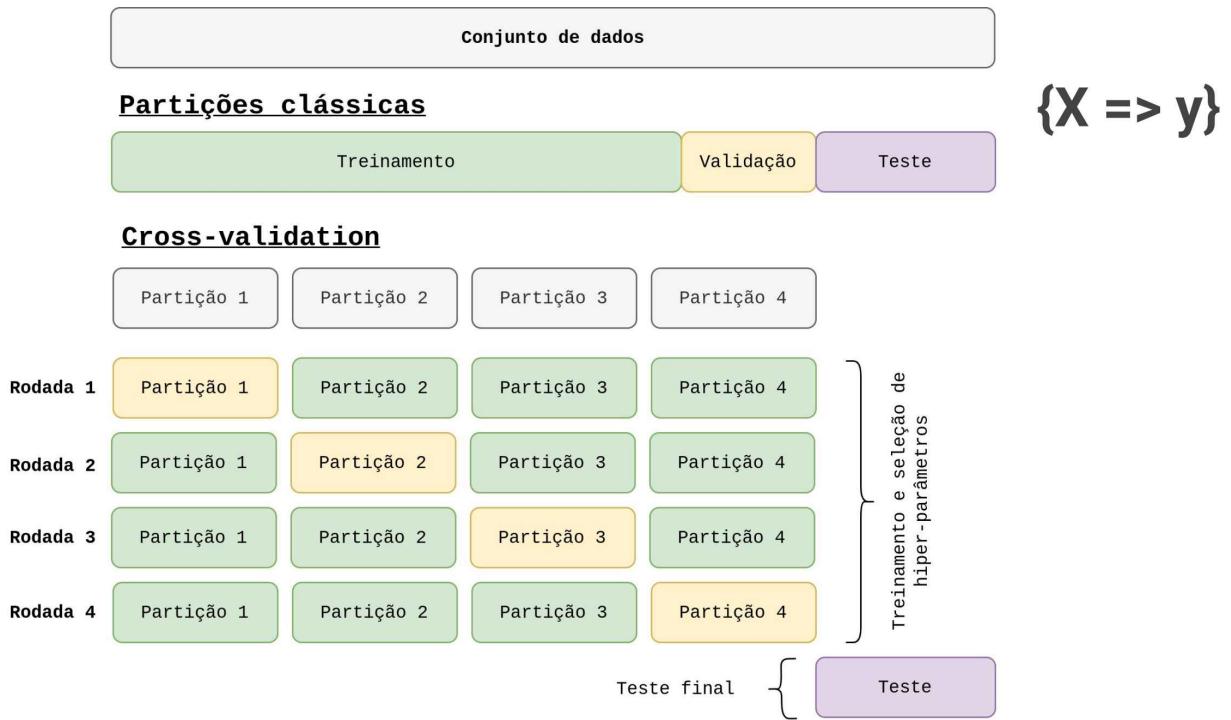
Regression

Se determina valores continuos

## Classification and Regression



# Classification and Regression



<https://computacao-inteligente.com.br/conceitos/avaliando-performance-cross-validation/>

Ivar Vargas Belizario

39

## Examples

Digit Recognizer (profesor)

<https://www.kaggle.com/competitions/digit-recognizer/>

Shelter Animal Outcomes (Trabajo Encargado 1 - TE1)

<https://www.kaggle.com/competitions/shelter-animal-outcomes/>

Ivar Vargas Belizario

40



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



# Minería de Datos

**Gracias**

Prof. Dr. Ivar Vargas Belizario

[ivargasbelizario@gmail.com](mailto:ivargasbelizario@gmail.com)

2024 - I