

**Seminario Permanente de Formación en  
Inteligencia Artificial aplicada a  
Defensa  
SIADEF**

**Sesión 7: Aprendizaje no  
supervisado**

Cristóbal J. Carmona <[ccarmona@ujaen.es](mailto:ccarmona@ujaen.es)>  
Pedro González <[pglez@ujaen.es](mailto:pglez@ujaen.es)>



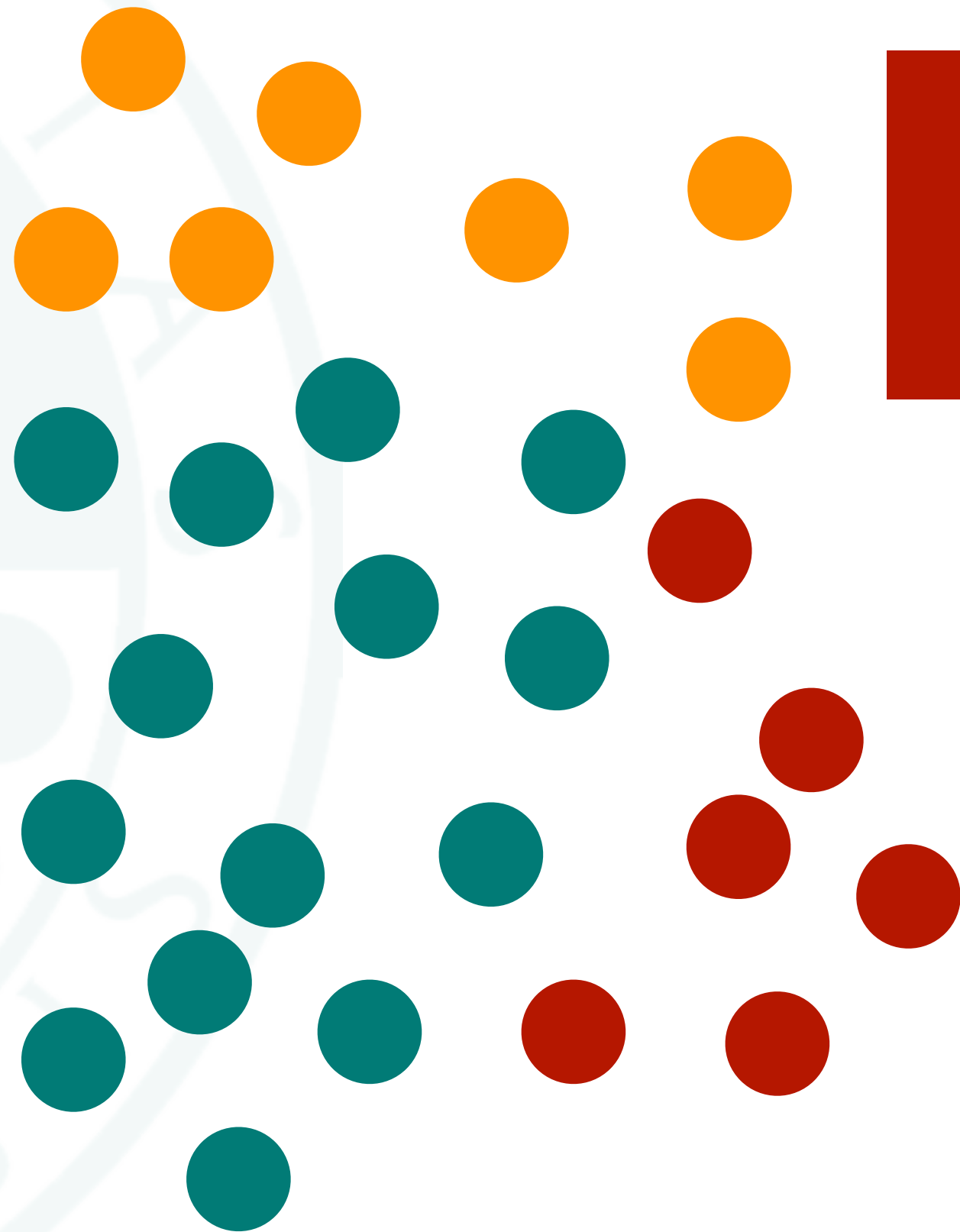


**Aprendizaje  
automático**



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

**Aprendizaje supervisado**



**Aprendizaje automático**

**Aprendizaje no supervisado**





## Aprendizaje supervisado



## Aprendizaje no supervisado







## Aprendizaje supervisado

Clasificación

Regresión

Series temporales

enfoque predictivo

## Aprendizaje no supervisado

Asociación

Agrupamiento

enfoque descriptivo

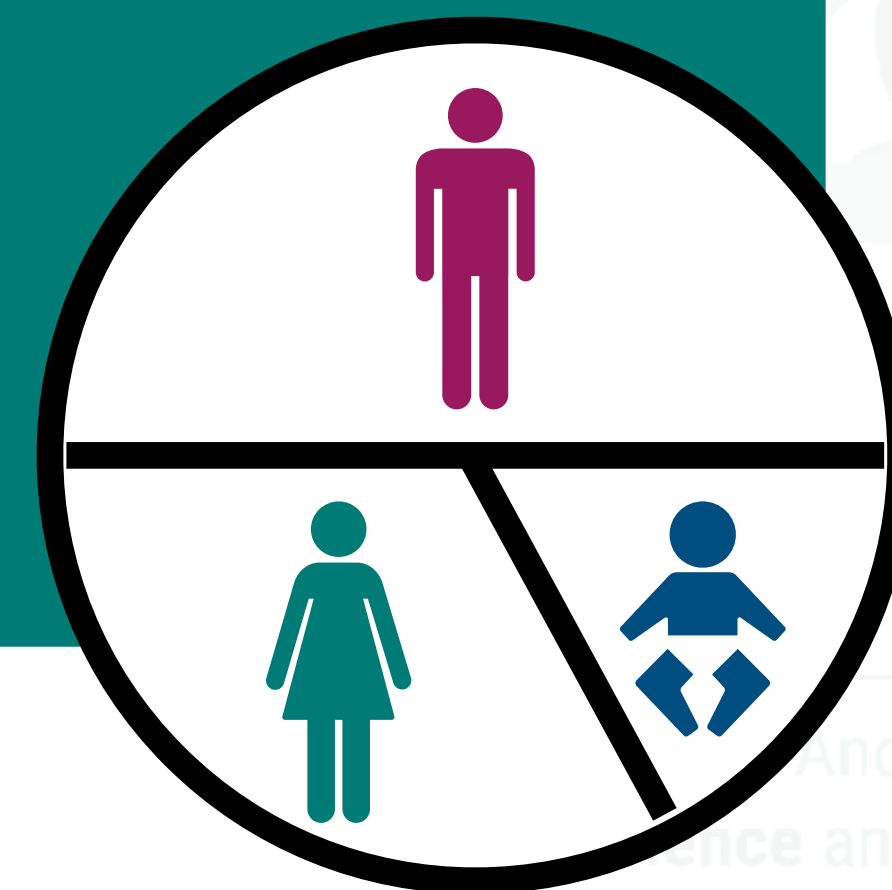
## Asociación

conceptos básicos  
ejemplos  
librerías

$X \rightarrow Y$

## Agrupamiento

conceptos básicos  
ejemplos  
paquetes en R



# Asociación

# Asociación

motivación

Para muchas empresas es un factor de éxito **comprender y actuar sobre los patrones de los clientes** utilizando los datos de las transacciones realizadas.

Sería muy interesante disponer de una herramienta que permita registrar y analizar las **transacciones** con los clientes para **revelar información valiosa sobre su comportamiento** respecto a la compra o consumo de artículos.

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Asociación

## definición

Dada una base de **datos de transacciones**, donde cada transacción es una lista de artículos (comprados por un cliente en la misma visita), encontrar **todas** las reglas que co-relacionen la presencia de un conjunto de artículos con otro conjunto de artículos

- **Ejemplos:**

- ▶ Los estudiantes que cursan Inteligencia Artificial tienden a cursar también Minería de Datos.
- ▶ El 98% de los clientes que compran neumáticos y accesorios para el automóvil, también adquieren servicios (cambio de neumáticos, ...).

# Asociación

conceptos básicos

**Itemset (Conjunto):** colección de uno o más items

- ▶ Ejemplo: {Milk, Bread, Diaper}
- ▶ k-itemset: itemset con k items

**Support count ( $\sigma$ ) (conteo soporte):** frecuencia de ocurrencia del itemset

- ▶ Ejemplo:  $\sigma$  ({Milk, Bread, Diaper}) = 2

**Support (Soporte):** fracción de transacciones que contienen un itemset

- ▶ Ejemplo:  $s$  ({Milk, Bread, Diaper}) = 2/5

**Frequent Itemset (Conjunto frecuente):** un itemset con soporte  $\geq \text{minsup}$ .  
En los conjuntos frecuentes los artículos están correlacionados positivamente

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

# Asociación

conceptos básicos

**Regla de asociación:**  $X \longrightarrow Y$ , donde  $X$  (antecedente) e  $Y$  (consecuente) son itemsets

▶ Ejemplo:  $\{\text{Milk, Diaper}\} \longrightarrow \{\text{Beer}\}$  (“ $\longrightarrow$ ” es co-ocurrencia, no causalidad)

**Tipos de reglas de asociación:**

▶ Asociaciones **Booleanas vs Cuantitativas** (tipo de los valores)

compra ( $x, \text{"PC"}$ )  $\longrightarrow$  compra ( $x, \text{"impresora"}$ ) vs. ingresos ( $x, \text{"2K..48K"}$ )  $\longrightarrow$  compra ( $x, \text{"PC"}$ )

▶ Asociaciones **Unidimensionales vs. Multidimensionales**

$A \longrightarrow B$  vs.  $A \& B \& N \longrightarrow D$

▶ Análisis con **distintos niveles** de abstracción:

Edad ( $x, \text{"30..39"}$ )  $\longrightarrow$  compra ( $x, \text{"Tablet"}$ ) vs. Edad ( $x, \text{"30..39"}$ )  $\longrightarrow$  compra ( $x, \text{"Tablet Galaxy"}$ )



# Asociación

conceptos básicos

## Métricas de evaluación de reglas

- ▶ **Soporte (s)**: fracción de las transacciones que contienen tanto a X como a Y
- ▶ **Confianza (c)**: frecuencia con la que los items de Y aparecen en transacciones que contienen los items de X

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Ejemplo:  $\{ \text{Milk}, \text{Diaper} \} \Rightarrow \text{Beer}$

$$s = \frac{\sigma(\text{Milk}, \text{Diaper}, \text{Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk}, \text{Diaper}, \text{Beer})}{\sigma(\text{Milk}, \text{Diaper})} = \frac{2}{3} = 0.67$$



# Asociación

conceptos básicos

## Minería de reglas de asociación:

- ▶ Dado un conjunto  $T$  de transacciones, el objetivo de la minería de reglas de asociación es encontrar **TODAS** las reglas que cumplan:
  - ▶ Soporte  $\geq \text{minsup}$
  - ▶ Confianza  $\geq \text{minconf}$

*Nota: encontrar estas reglas no significa que deba existir una relación entre antecedente y consecuente. Por lo tanto, un experto debería evaluarlas.*

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Asociación

algoritmos de descubrimiento de reglas de asociación

**Objetivo:** encontrar todas las reglas que superen los umbrales de soporte y confianza

## Algoritmos:

- ▶ Apriori y AprioriTid (Agrawal & Srikant, 1994)
- ▶ Opus (Webb, 1996)
- ▶ Direct Hasing and Pruning (DHP) (Adamo, 2001)
- ▶ Dynamic Set Counting (DIC) (Adamo, 2001)
- ▶ Charm (Zaki & Hsiao, 2002)
- ▶ FP-growth (Han, Pei & Yin, 1999)
- ▶ Closet (Pei, Han & Mao, 2000)
- ▶ .....

# Asociación

algoritmos de descubrimiento de reglas de asociación

**Los algoritmos deben generar siempre el mismo conocimiento**

## ¿Qué los hace diferentes?

- ▶ Forma en que los datos son cargados en memoria.
- ▶ Tiempo de procesamiento.
- ▶ Tipos de atributos (numéricos, categóricos).
- ▶ Forma en que generan los itemsets.
- ▶ Estructura de datos utilizada.



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Asociación

algoritmos de descubrimiento de reglas de asociación

Enfoque habitual: descomponer el problema en **dos pasos**

## 1. Generar itemsets frecuentes

- ▶ Genera todos los conjuntos frecuentes con **soporte**  $>$  **minsup**.

## 2. Utilizar itemsets frecuentes para generar reglas con fuerte asociación

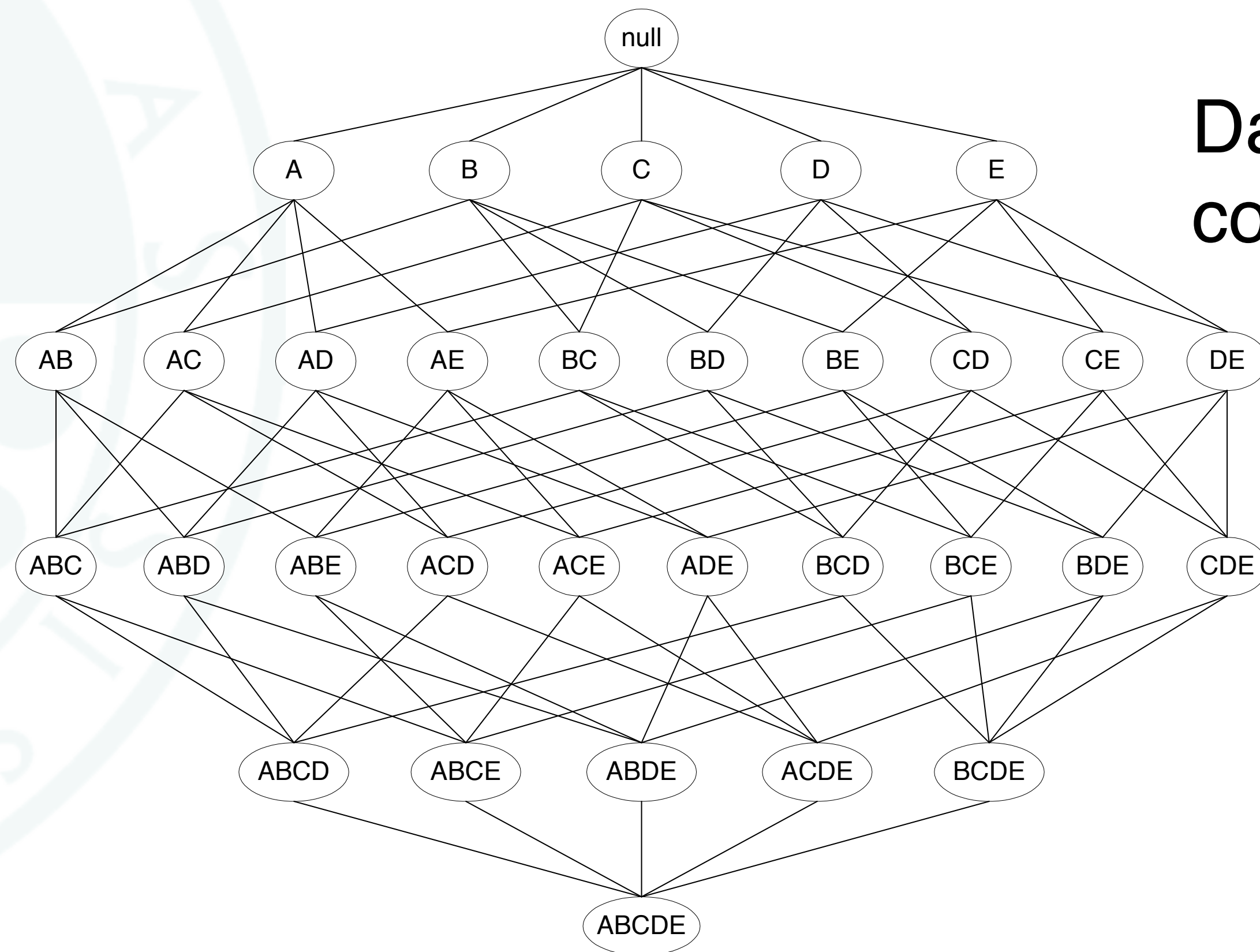
- ▶ Generar reglas de alta confianza (**confianza**  $>$  **minconf**) de cada conjunto frecuente (que cada regla es una partición binaria de un conjunto frecuente).



# Asociación

algoritmos de descubrimiento de reglas de asociación

La generación de conjuntos frecuentes es **costosa**



Dados  $d$  items, hay  $2^d$  posibles conjuntos (itemsets) candidatos

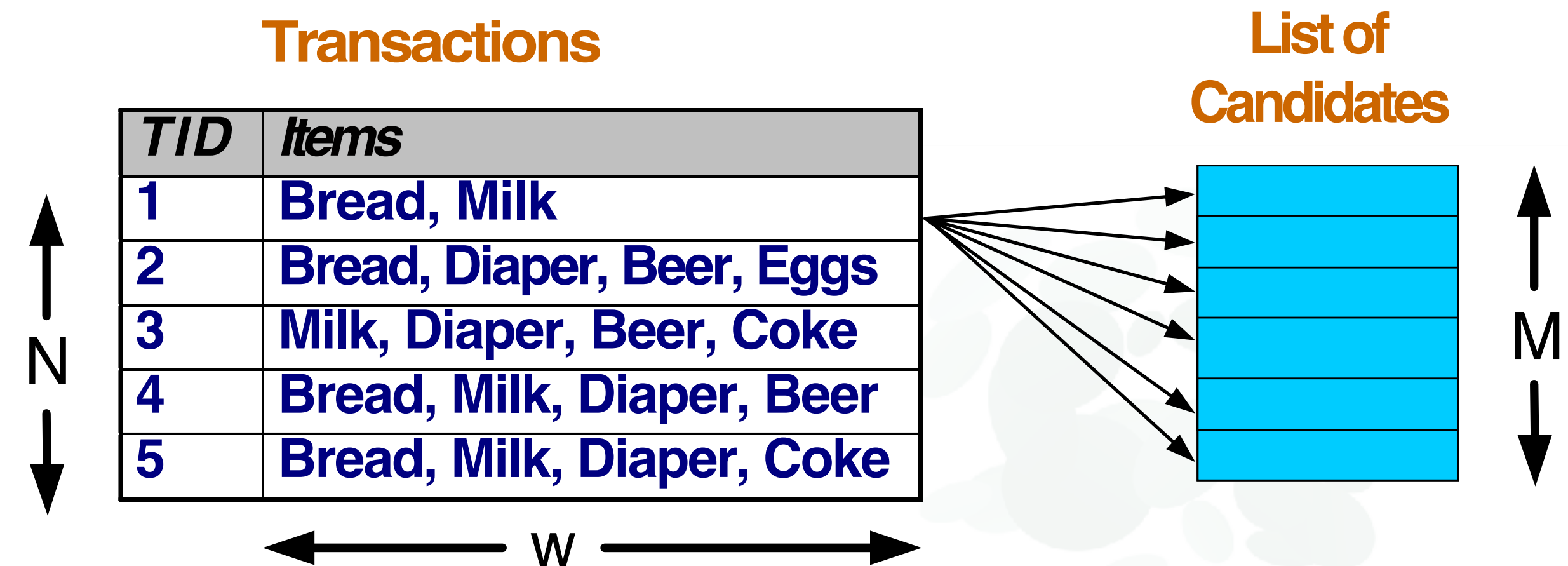


# Asociación

algoritmos de descubrimiento de reglas de asociación

La generación de conjuntos frecuentes es **costosa**

- Enfoque de **fuerza bruta**:
  - ▶ Cada Itemset del grafo es un conjunto frecuente **candidato**.
  - ▶ Calcular el soporte de cada candidato recorriendo la BD.
  - ▶ Comprobar cada transacción contra todos los candidatos.
  - ▶ Muy costoso:  $\sim O(NMw)$



# Asociación

algoritmos de descubrimiento de reglas de asociación

## Estrategias de generación de conjuntos frecuentes

- Reducir el **número de candidatos** (M)
  - ▶ Búsqueda completa:  $M=2^d$
  - ▶ Utilizar técnicas de poda para reducir M
- Reducir el **número de transacciones** (N)
  - ▶ Reducir el tamaño de N conforme aumenta el tamaño del Itemset
  - ▶ Lo utilizan los algoritmos DHP y de minado vertical.
- Reducir el **número de comparaciones** (NM)
  - ▶ Usar estructuras de datos eficientes para almacenar los candidatos o las transacciones.
  - ▶ No es necesario comprobar cada candidato contra cada transacción.



# Asociación

algoritmo Apriori

Primer algoritmo de descubrimiento de reglas de asociación

Utiliza estrategias para optimizar el número de conjuntos frecuentes candidatos y de posibles reglas a generar

R. Agrawal, R. Srikant - Fast Algorithms for Mining Association Rules - Proc. of the 20th International Conference on Very Large Databases, Santiago, Chile, Sept. 1994.

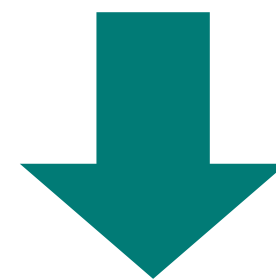


# Asociación

algoritmo Apriori

## Reducción de conjuntos candidatos: principio Apriori

- Si un itemset es frecuente, TODOS sus subconjuntos serán frecuentes.
  - El soporte de un itemset nunca es superior al soporte de sus subconjuntos.
  - Esto se conoce como la propiedad de **anti-monotonía** del soporte.

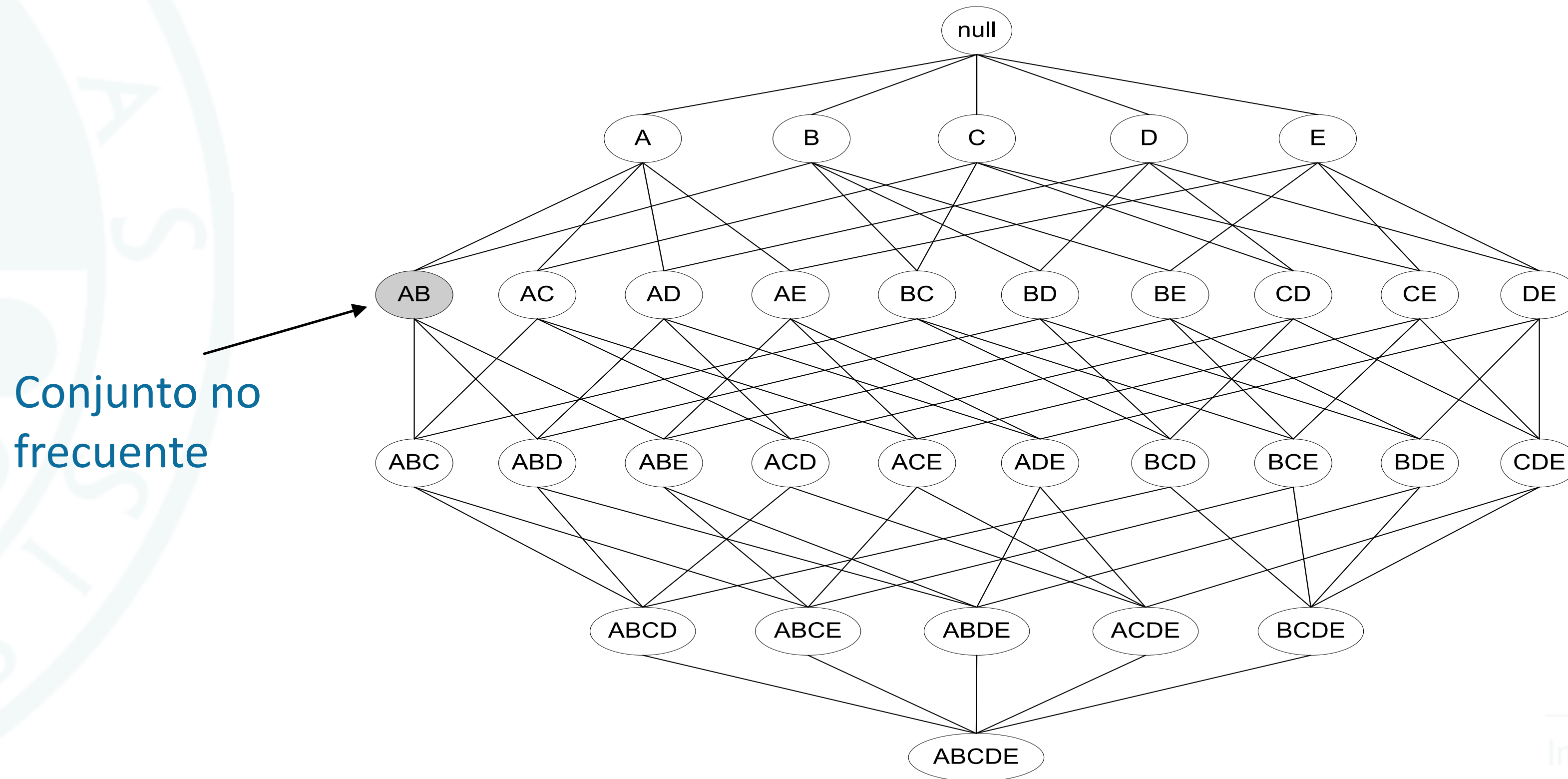


**Si un itemset NO es frecuente, ningún superconjunto suyo es frecuente**

# Asociación

algoritmo Apriori

Reducción de conjuntos candidatos: **principio Apriori**



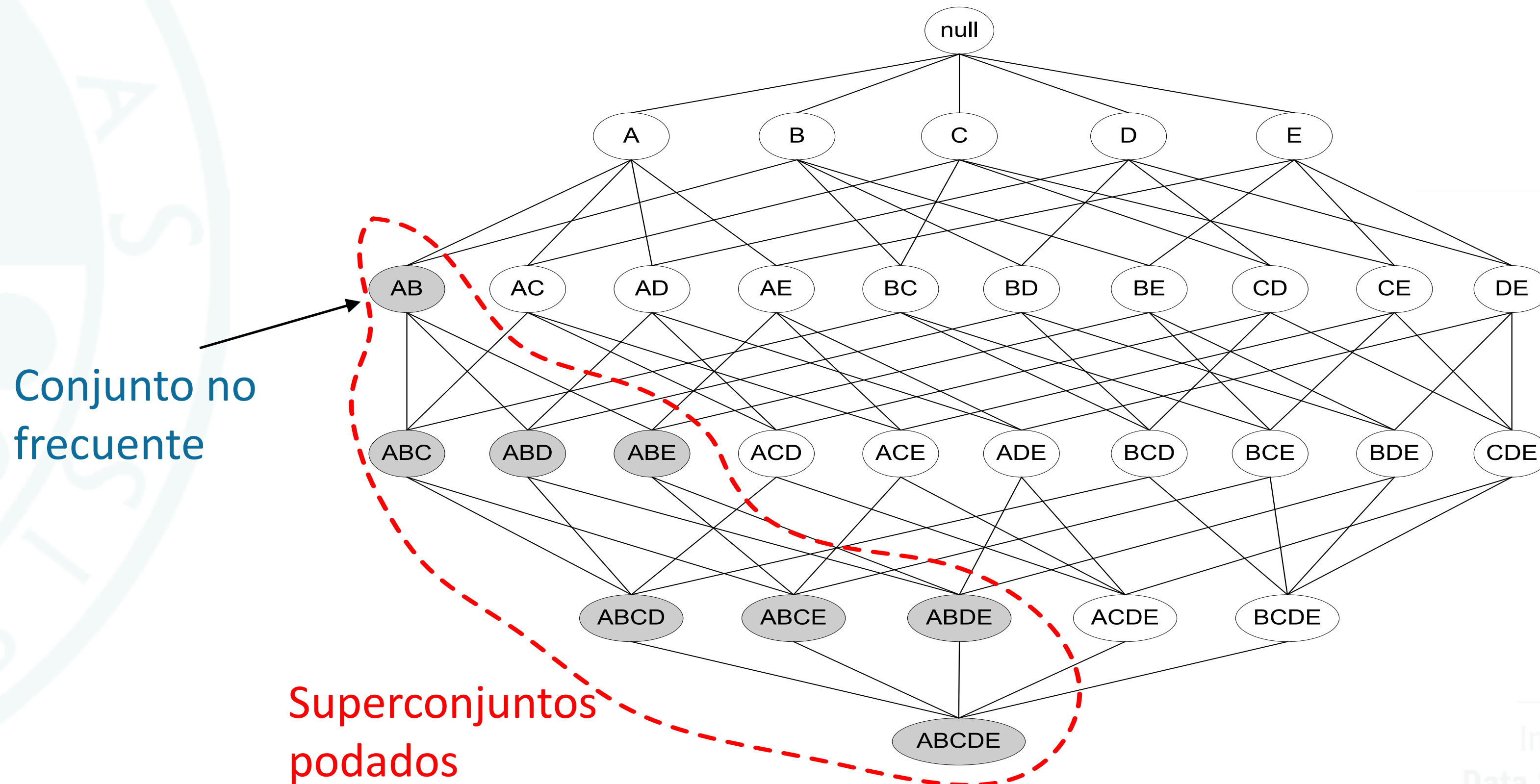
DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Asociación

algoritmo Apriori

Reducción de conjuntos candidatos: **principio Apriori**



DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Asociación

algoritmo Apriori

## Reducción de conjuntos candidatos: principio Apriori

Item	Count
<b>Bread</b>	<b>4</b>
<b>Coke</b>	<b>2</b>
<b>Milk</b>	<b>4</b>
<b>Beer</b>	<b>3</b>
<b>Diaper</b>	<b>4</b>
<b>Eggs</b>	<b>1</b>

Items (1-itemsets)



Itemset	Count
{ <b>Bread,Milk</b> }	<b>3</b>
{ <b>Bread,Beer</b> }	<b>2</b>
{ <b>Bread,Diaper</b> }	<b>3</b>
{ <b>Milk,Beer</b> }	<b>2</b>
{ <b>Milk,Diaper</b> }	<b>3</b>
{ <b>Beer,Diaper</b> }	<b>3</b>

Pairs (2-itemsets)

(No es necesario generar candidatos en que aparezca Coke o Eggs)



Triplets (3-itemsets)

Itemset	Count
{ <b>Bread,Milk,Diaper</b> }	<b>3</b>

Soporte mínimo = 3

Si se consideran todos los subconjuntos,

$$C_{6,1} + C_{6,2} + C_{6,3} = 41$$

Con poda basada en soporte,

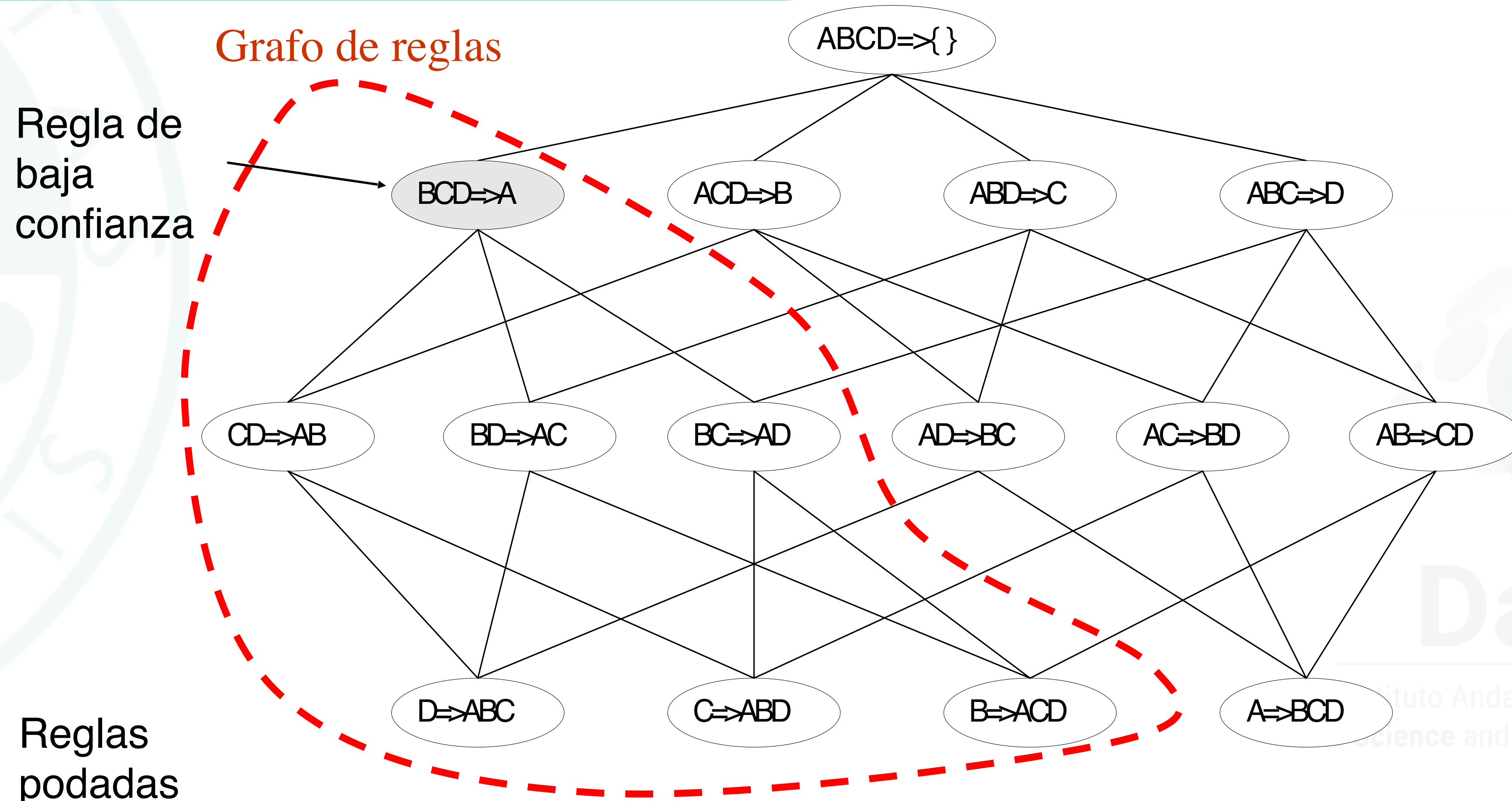
$$6 + 6 + 1 = 13$$



# Asociación

algoritmo Apriori

## Reducción de reglas candidatas

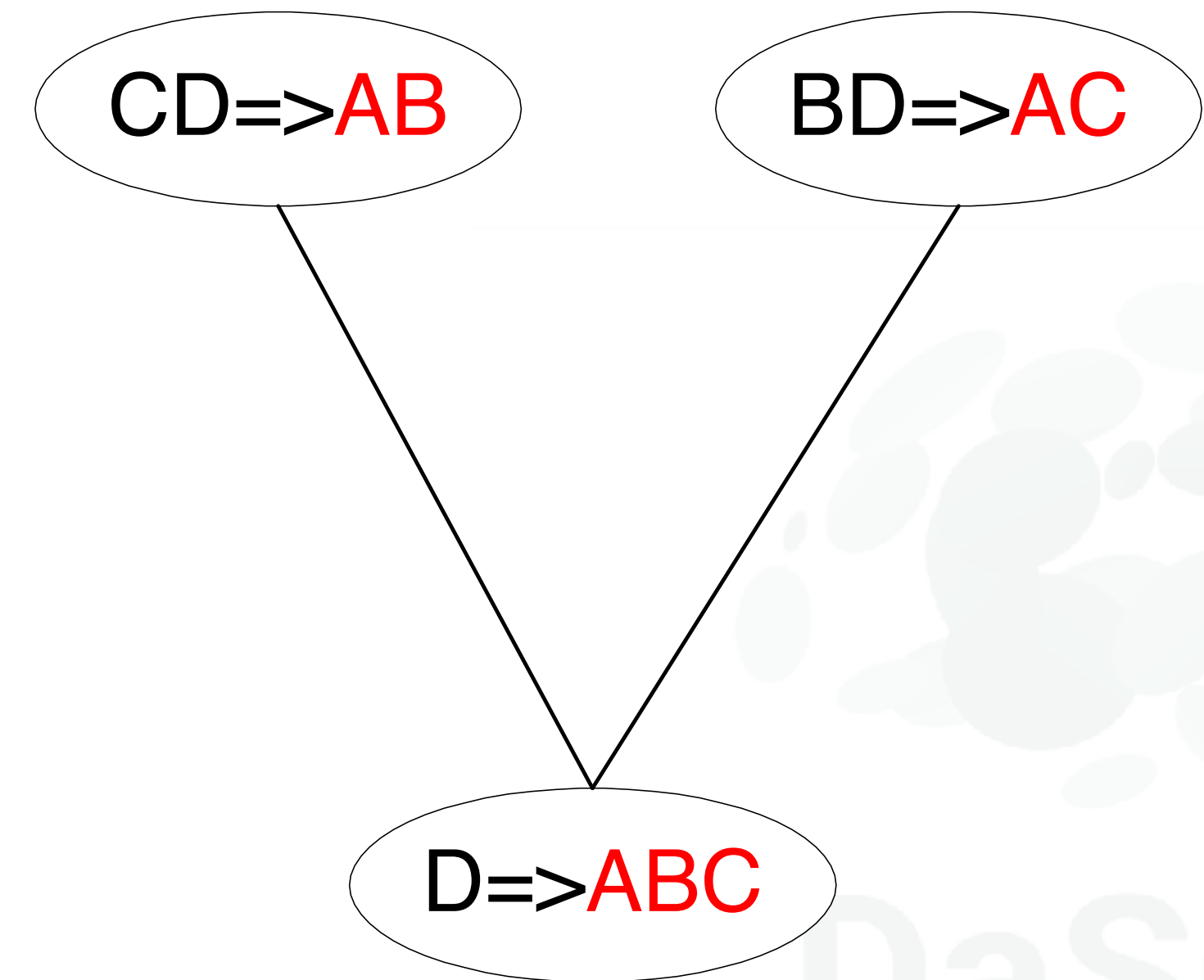


# Asociación

algoritmo Apriori

## Reducción de reglas candidatas

- La regla candidata se genera uniendo dos reglas que compartan el mismo prefijo en el consecuente de la regla
- $\text{join}(CD \Rightarrow AB, BD \Rightarrow AC)$  produciría la regla candidata  $D \Rightarrow ABC$
- Podar la regla  $D \Rightarrow ABC$  si su subconjunto  $AD \Rightarrow BC$  no tiene alta confianza



# Asociación

algoritmo FP-Growth

Permite extraer reglas de asociación a partir de itemsets frecuentes sin necesidad de generar candidatos para cada tamaño.

Emplea una estructura de árbol (Frequent Pattern Tree) que almacena la información de las transacciones, comprimiéndola hasta 200 veces.

Después separa la estructura asociada a cada patrón frecuente para analizarlo de forma separada y concatenar los resultados.

En la mayoría de casos, FP-Growth es más rápido que Apriori.

Han, J.; Pei, J.; Yin, Y. Mining frequent patterns without candidate generation. In Proc. 2000 ACM SIGMOD Int. Conf. Management of Data (SIGMOD'00), 1–12. 2000.



# Asociación

evaluación de los patrones de asociación

Algoritmos de reglas de asociación tienen a **generar demasiadas reglas**:

- ▶ Muchas de ellas no son interesantes o son redundantes.
- ▶ Aplicar **postprocesamiento** para seleccionar reglas valiosas.
- ▶ Ejemplo:  $\{A,B,C\} \rightarrow \{D\}$  y  $\{A,B\} \rightarrow \{D\}$ 
  - ▶ Son redundantes si tienen el mismo soporte y confianza.

Se pueden utilizar medidas de Interés para podar o calificar los patrones obtenidos (originalmente sólo se usaba soporte y confianza).



# Asociación

evaluación de los patrones de asociación

La confianza no siempre es buena para buscar reglas interesantes:

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Regla: Tea  $\rightarrow$  Coffee

Confianza =  $P(\text{Coffee}|\text{Tea}) = 0.75$

- ▶ Por la alta confianza (0.75) parece una buena regla
- ▶ Pero es engañoso, porque  $P(\text{Coffee}) = 0.9$ 
  - ▶  $P(\text{Coffee}|\text{Tea}) = 0.75$
  - ▶  $P(\text{Coffee}|\overline{\text{Tea}}) = 0.9375$

# Asociación

evaluación patrones asociación

## Medidas de interés

- **Medidas objetivas:**
  - ▶ Cataloga los patrones a partir de estadísticas calculadas sobre los datos
  - ▶ Ejemplo: 21 medidas de asociación (soporte, confianza, Laplace, Gini, información mutua, Jaccard, etc):

#	Measure	Formula
1	$\phi$ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's ( $\lambda$ )	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio ( $\alpha$ )	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's Q	$\frac{P(A,B)P(\bar{A}\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A}\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
5	Yule's Y	$\frac{\sqrt{P(A,B)P(\bar{A}\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A}\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha - 1}}{\sqrt{\alpha + 1}}$
6	Kappa ( $\kappa$ )	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information ( $M$ )	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$
8	J-Measure ( $J$ )	$\max \left( P(A, B) \log \left( \frac{P(B A)}{P(B)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{B} \bar{A})}{P(\bar{B})} \right), \right. \\ \left. P(A, B) \log \left( \frac{P(A B)}{P(A)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{A} \bar{B})}{P(\bar{A})} \right) \right)$
9	Gini index ( $G$ )	$\max \left( P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] \right. \\ \left. - P(B)^2 - P(\bar{B})^2, \right. \\ \left. P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] \right. \\ \left. - P(A)^2 - P(\bar{A})^2 \right)$
10	Support ( $s$ )	$P(A, B)$
11	Confidence ( $c$ )	$\max(P(B A), P(A B))$
12	Laplace ( $L$ )	$\max \left( \frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2} \right)$
13	Conviction ( $V$ )	$\max \left( \frac{P(A)P(\bar{B})}{P(\bar{A}B)}, \frac{P(B)P(\bar{A})}{P(\bar{B}A)} \right)$
14	Interest ( $I$ )	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine ( $IS$ )	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's ( $PS$ )	$P(A, B) - P(A)P(B)$
17	Certainty factor ( $F$ )	$\max \left( \frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)} \right)$
18	Added Value ( $AV$ )	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength ( $S$ )	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
20	Jaccard ( $\zeta$ )	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Klosgen ( $K$ )	$\sqrt{P(\bar{A}, \bar{B}) \max(P(B A) - P(B), P(A B) - P(A))}$

# Asociación

evaluación patrones asociación

## Medidas de interés

- **Medidas subjetivas:**
  - ▶ Cataloga los patrones a partir de la interpretación del usuario:
    - Un patrón es interesante subjetivamente si contradice las expectativas del usuario (Silberschatz & Tuzhilin).
    - Un patrón es interesante subjetivamente si permite actuar (Silberschatz & Tuzhilin).

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Asociación

## aplicaciones

- ▶ Recomendación de productos
- ▶ Recomendaciones de medios digitales
- ▶ Diagnóstico médico
- ▶ Optimización de contenidos
- ▶ Bioinformática
- ▶ Minería Web
- ▶ Análisis de datos científicos



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Asociación

aplicación: análisis de la cesta de la compra

Un aplicación muy conocida es el **Análisis de la cesta de la compra** (*Market Basket Analysis*), normalmente asociado con ventas de productos.



iz Interuniversitario en  
Data Science and Computational Intelligence

# Asociación

aplicación: análisis de la cesta de la compra

Se utiliza para descubrir los patrones de co-ocurrencia de grupos específicos de artículos o items.

- ▶ Si alguien compra un grupo de artículos, es más probable (o menos) que también compre otro grupo de artículos
- ▶ Se trata de identificar grupos de artículos que son adquiridos en conjunto.

Intenta identificar reglas de la forma:

- ▶ {fideos, queso rallado} → {salsa}

Julander, C. (1992) Basket Analysis: A New Way of Analysing Scanner Data. International Journal of Retail & Distribution Management 20 (7), 10-18.



# Asociación

aplicación: análisis de la cesta de la compra

## Ejemplo: Compra de productos de alimentación

- Tenemos 10.000 recibos de compras de clientes.
  - ▶ Cada recibo es una transacción.
  - ▶ Cada línea de cada recibo es un artículo (no se repiten).
- Para guardar toda esa información, creamos una única tabla que contiene una fila por cada artículo de cada transacción:

Transaction	Items
A0001	citrus fruit
A0001	margarine
A0001	ready soups
A0001	semi-finished bread
A0002	coffee
A0002	tropical fruit

# Asociación

aplicación: análisis de la cesta de la compra

## Ejemplo: Compra de productos de alimentación

- Para usar esta información con los algoritmos de análisis de la cesta de la compra, primero debemos transformarlo a una **representación binaria** (0,1) indicando si un artículo concreto se compró en una transacción específica.
- Tendremos una única fila por transacción, y tantas columnas como artículos:

Transaction	Items
A0001	citrus fruit
A0001	margarine
A0001	ready soups
A0001	semi-finished bread
A0002	coffee
A0002	tropical fruit



Transaction	citrus fruit	margarine	ready soups	semi-finished bread	coffee	tropical fruit
A0001	1	1	1	1	0	0
A0002	0	0	0	0	1	1
A0003	0	0	0	0	0	0
A0004	0	0	0	0	0	0



# Asociación

aplicación: análisis de la cesta de la compra

## Ejemplo: Compra de productos de alimentación

- Hay que establecer los **parámetros** para la detección de patrones.
  - **Soporte mínimo:** 0.001 (muchos recibos y productos)
  - **Confianza mínima:** 0.70.
  - **Longitud de la regla:** máximo 3 elementos (máximo 2 en el antecedente)

Rules	Support	Confidence	Lift
{liquor, red/blush wine} => {bottled beer}	0.002	0.90	11.24
{cereals, yogurt} => {whole milk}	0.002	0.81	3.17
{butter, jam} => {whole milk}	0.001	0.83	3.26
{chocolate, pickled vegetables} => {whole milk}	0.001	0.86	3.35
{grapes, onions} => {other vegetables}	0.001	0.92	4.74
{hard cheese, oil} => {other vegetables}	0.001	0.92	4.74

# Asociación

aplicación: análisis de la cesta de la compra

## Ejemplo: Compra de productos de alimentación

- **Ajuste** de los parámetros del algoritmo:
  - ▶ La distribución de artículos por transacción será distinta de unos negocios a otros, siendo necesarios valores distintos de soporte y confianza.
  - ▶ Para determinar qué funciona mejor hay que experimentar con los parámetros:
    - Umbrales bajos —> más reglas: más difícil identificar las que tienen mayor impacto.
  - ▶ No hay certeza: una opción es comenzar experimentando con valores bajos de los parámetros e irlos subiendo en función de los resultados.

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Asociación

aplicación: análisis de la cesta de la compra

## Ejemplo: Compra de productos de alimentación

- Para aprovechar este conocimiento adicional, nos centramos en 3 artículos específicos: "Yogurt", "Tropical Fruit" y "Bottled Beer".
  - Para ello usamos los mismos umbrales y especificamos que aparezcan estos **términos en el consecuente**:

Rules	Support	Confidence	Lift
{pip fruit, sausage, sliced cheese} => {yogurt}	0.001	0.86	6.14
{butter, cream cheese , root vegetables} => {yogurt}	0.001	0.91	6.52
{butter, margarine, tropical fruit} => {yogurt}	0.001	0.85	6.07
{butter, curd, other vegetables, tropical fruit} => {yogurt}	0.001	0.83	5.97
{liquor, red/blush wine} => {bottled beer}	0.002	0.90	11.24
{citrus fruit, fruit/vegetable juice, grapes} => {tropical fruit}	0.001	0.85	8.06
{ham, other vegetables, pip fruit, yogurt} => {tropical fruit}	0.001	0.83	7.94



# Asociación

aplicación: análisis de la cesta de la compra

## Ejemplo: Compra de productos de alimentación

- También podemos utilizar estos **términos como antecedente** en la generación del conjunto de reglas:

Rules	Support	Confidence	Lift
{yogurt} => {whole milk}	0.056	0.40	1.57
{tropical fruit} => {other vegetables}	0.036	0.34	1.77
{yogurt} => {rolls/buns}	0.034	0.25	1.34
{tropical fruit} => {rolls/buns}	0.025	0.23	1.27
{bottled beer} => {soda}	0.017	0.21	1.21
{bottled beer} => {bottled water}	0.016	0.20	1.77
{tropical fruit} => {pip fruit}	0.020	0.19	2.57
{tropical fruit} => {citrus fruit}	0.020	0.19	2.29

# Asociación

aplicación: análisis de la cesta de la compra

## Uso del análisis para la toma de decisiones de negocio

- Antes de usar los datos para tomar ninguna decisión de negocio, es importante dar un paso atrás y recordar que:
  - ▶ La salida del análisis refleja la frecuencia de **co-ocurrencia** de los artículos en las transacciones. Esto es función tanto de la fortaleza de la asociación entre artículos como de la forma en que el negocio los ha presentado al cliente.
    - ▶ Es decir: los artículos pueden aparecer juntos no porque estén naturalmente conectados, sino porque los encargados de la organización los han presentado juntos.

# Asociación

aplicación: análisis de la cesta de la compra

## Uso del análisis para la toma de decisiones de negocio

- Los resultados del análisis de la cesta de la compra se pueden usar para realizar campañas de **marketing dirigido**:
  - ▶ Por ejemplo, para cada regla, elegimos algunos de los productos comprados con márgenes y retornos altos, y mandamos correos personalizados a los clientes.
  - ▶ La forma de usar el análisis es importante para el propio análisis:
    - Para alimentar a un sistema automático de recomendaciones, nos interesará obtener un conjunto amplio de reglas.
    - Para experimentar, tiene más sentido obtener pocas reglas de gran valor, y actuar sólo sobre ellas.



# Asociación

reglas de asociación en R

## Paquetes en R para descubrimiento de reglas de asociación

Paquete	Descripción
arules	algoritmos apriori y eclat
arulesViz	visualizaciones para exploración de reglas
rCBA	algoritmo fpgrowth

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Asociación

## referencias bibliográficas

- ▶ Adamo, J.M. Data Mining For Association Rules and Sequential Patterns Sequential and Parallel Algorithms. Springer. 2001
- ▶ Agrawal, R; Srikant, R. Fast Algorithms for Mining Association Rules - Proc. of the 20th International Conference on Very Large Databases, Santiago, Chile, Sept. 1994.
- ▶ Han, J.; Pei, J.; Yin, Y. Mining frequent patterns without candidate generation. In Proc. 2000 ACM SIGMOD Int. Conf. Management of Data (SIGMOD'00), 1–12. 2000.
- ▶ Julander, C. Basket Analysis: A New Way of Analysing Scanner Data. International Journal of Retail & Distribution Management 20 (7), 10-18, 1992.
- ▶ Pei, J; Han, J; Mao, R. CLOSET: An efficient algorithm for mining frequent closed itemsets. In DMKD 2000, pp. 11—20, 2000
- ▶ Silberschatz, A; Tuzhilin, A. On subjective measures of interestingness in knowledge Discovery. Proceedings of the First International Conference on Knowledge Discovery and Data Mining, 275–281. 1995.
- ▶ Webb, G.I. OPUS: An efficient admissible algorithm for unordered search. Journal of Artificial Intelligence Research, 3:45--83, 1996.
- ▶ Zaki, M.J; Hsiao, C.J. CHARM: An Efficient Algorithm for Closed Itemset Mining. 2002
- ▶ Zhang, C; Zhang, S. Association rule mining: models and algorithms. Springer. 2002.

# Agrupamiento



# Agrupamiento

motivación

**Crear grupos de instancias o ejemplos de similares características**

Segmentación de clientes de una empresa



# Agrupamiento

¿qué es un cluster? ¿qué es clustering?

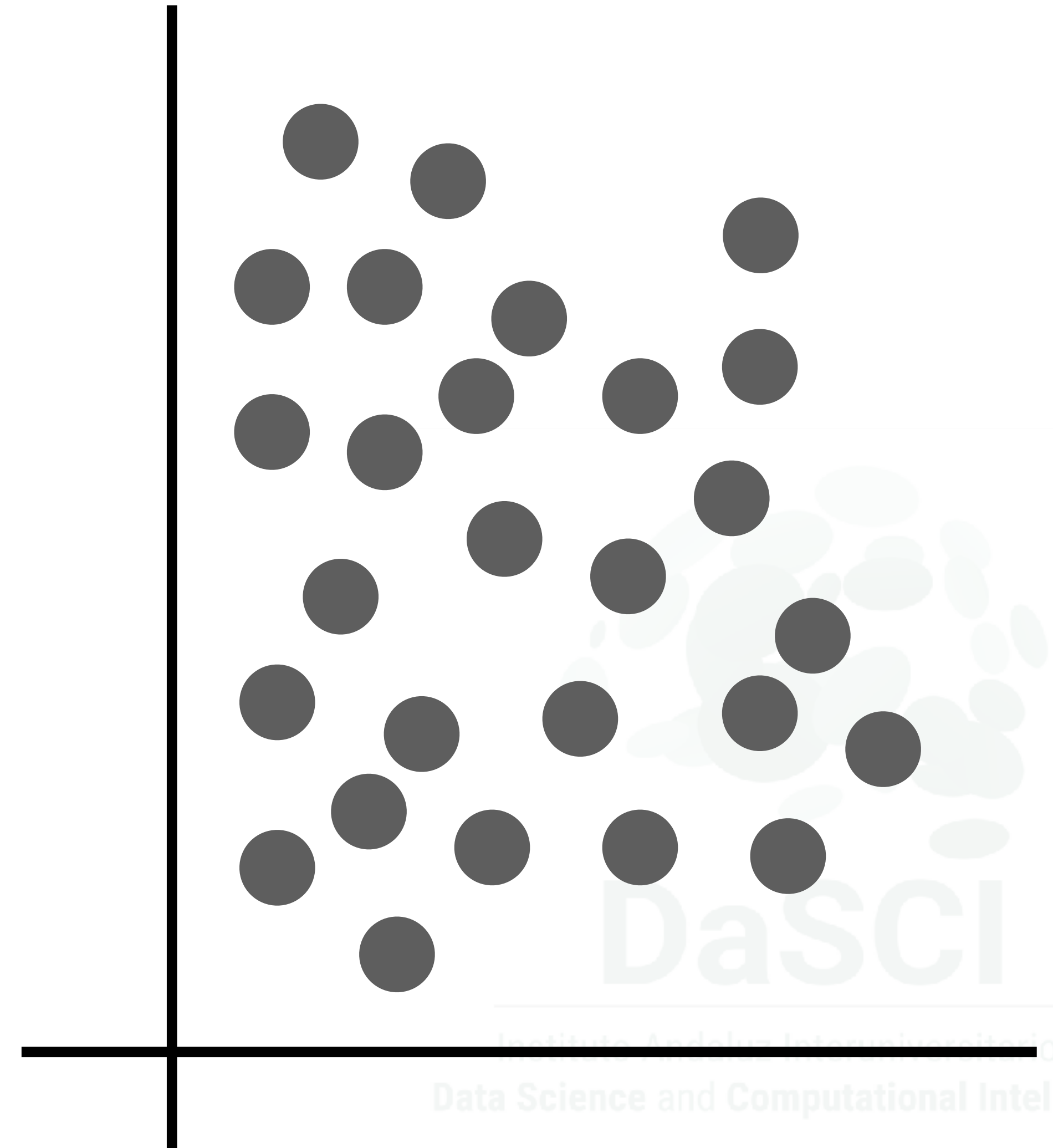
Por definición, **un *cluster* es un grupo o conjunto de objetos.**

- Similares a cualquier otro incluido en el mismo *cluster*.
- Distintos a los objetos incluidos en otros grupos.

**Clustering es la técnica de segmentar o agrupar una población heterogénea en un número de subgrupos homogéneos o clusters.**

# Agrupamiento

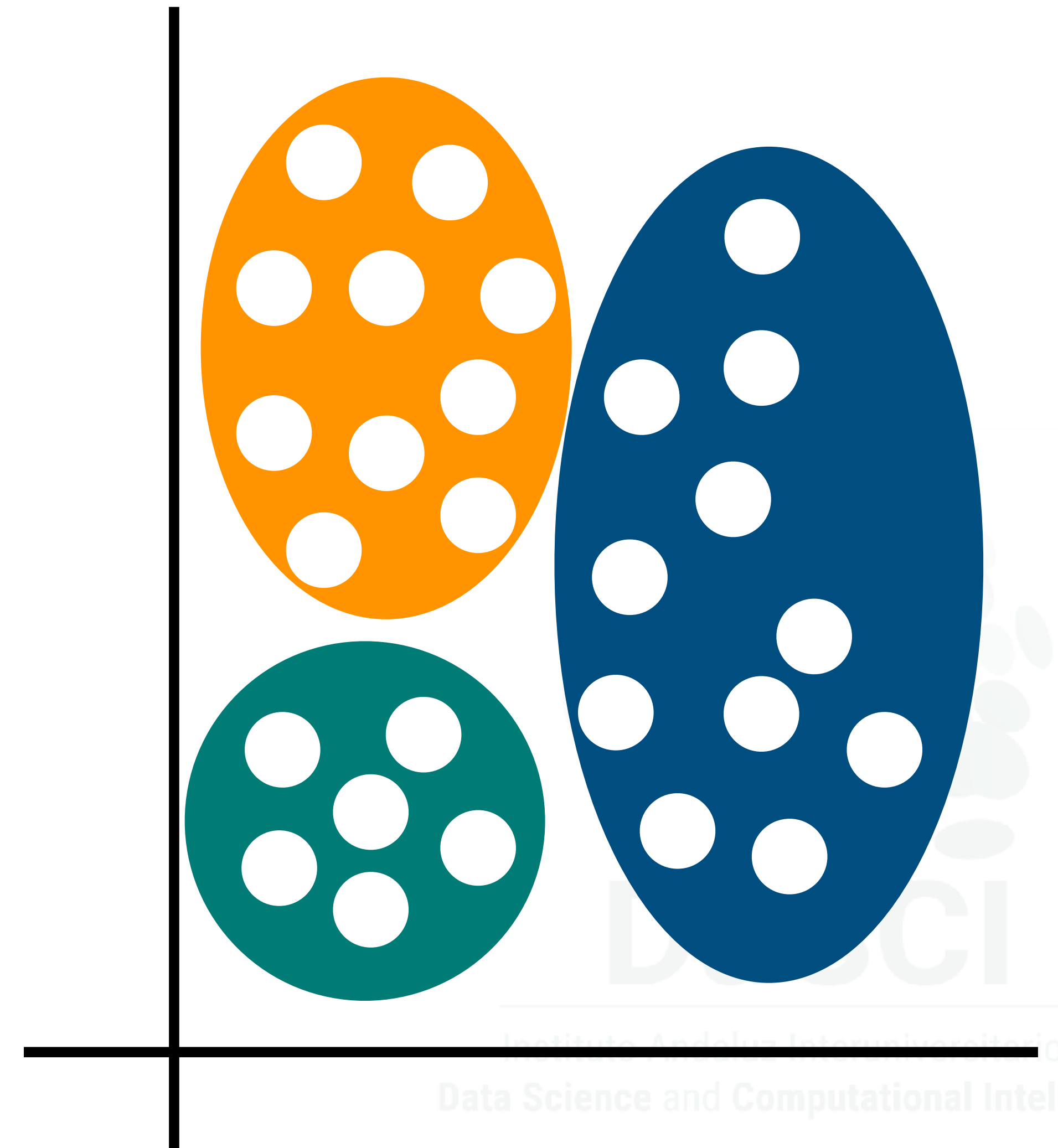
clasificación versus clustering





# Agrupamiento

clasificación versus clustering



# Agrupamiento

aplicaciones típicas

- Para tareas de preprocesamiento antes de aplicar otra técnica de descubrimiento del conocimiento
- Técnica de descubrimiento del conocimiento para obtener información acerca de la distribución de los datos

# Agrupamiento

problemas reales

- Dificultad en los manejos de outliers
- En bases de datos dinámicas implica que la pertenencia a clusters varía en el tiempo
- Interpretar el significado de cada cluster puede ser difícil
- **NO tenemos una única solución para un problema. El número de clusters es difícil de determinar**

Datos o instancias con valores atípicos respecto del conjunto completo



# Agrupamiento

aplicaciones



Segmentación de clientes de una empresa para usos publicitarios



amazon

# Agrupamiento

aplicaciones



Identificación de áreas con usos similares a partir de observaciones en cultivos



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Agrupamiento

aplicaciones



Agrupar asegurados con similares características para ofertar otros productos que estos clientes ya tienen contratados



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Agrupamiento

aplicaciones



Grupos de viviendas de acuerdo a su tipo, valor o situación geográfica



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Agrupamiento

aplicaciones



Agrupar documentos, analizar ficheros .log, patrones similares de comportamiento en clientes



Instituto Andaluz Interuniversitario en Data Science and Computational Intelligence

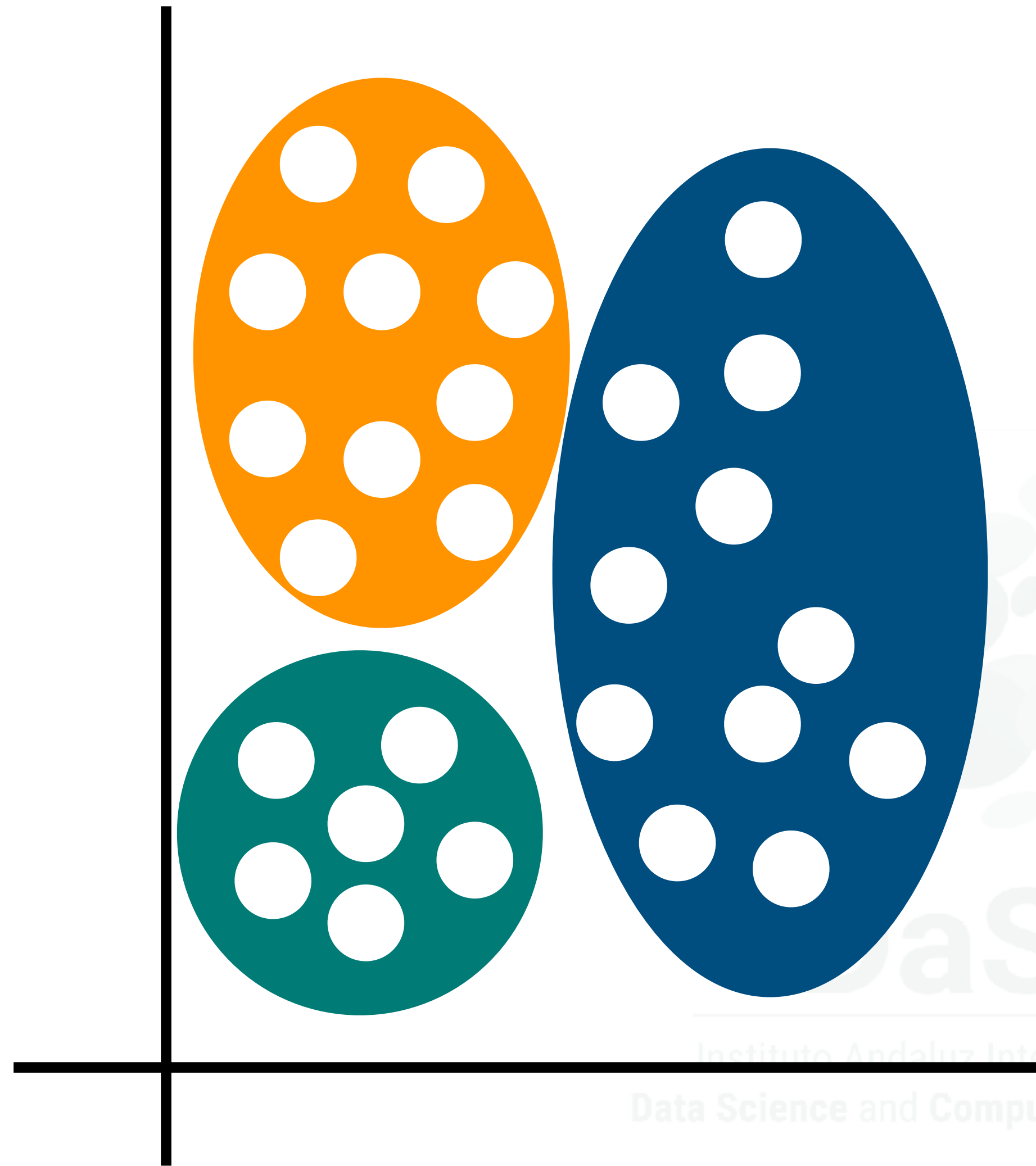
# Agrupamiento

bondad de un método de clustering

Dos conceptos clave

**MINIMIZAR** la  
similaridad intra-cluster

**MAXIMIZAR** la  
similaridad inter-cluster





# Agrupamiento

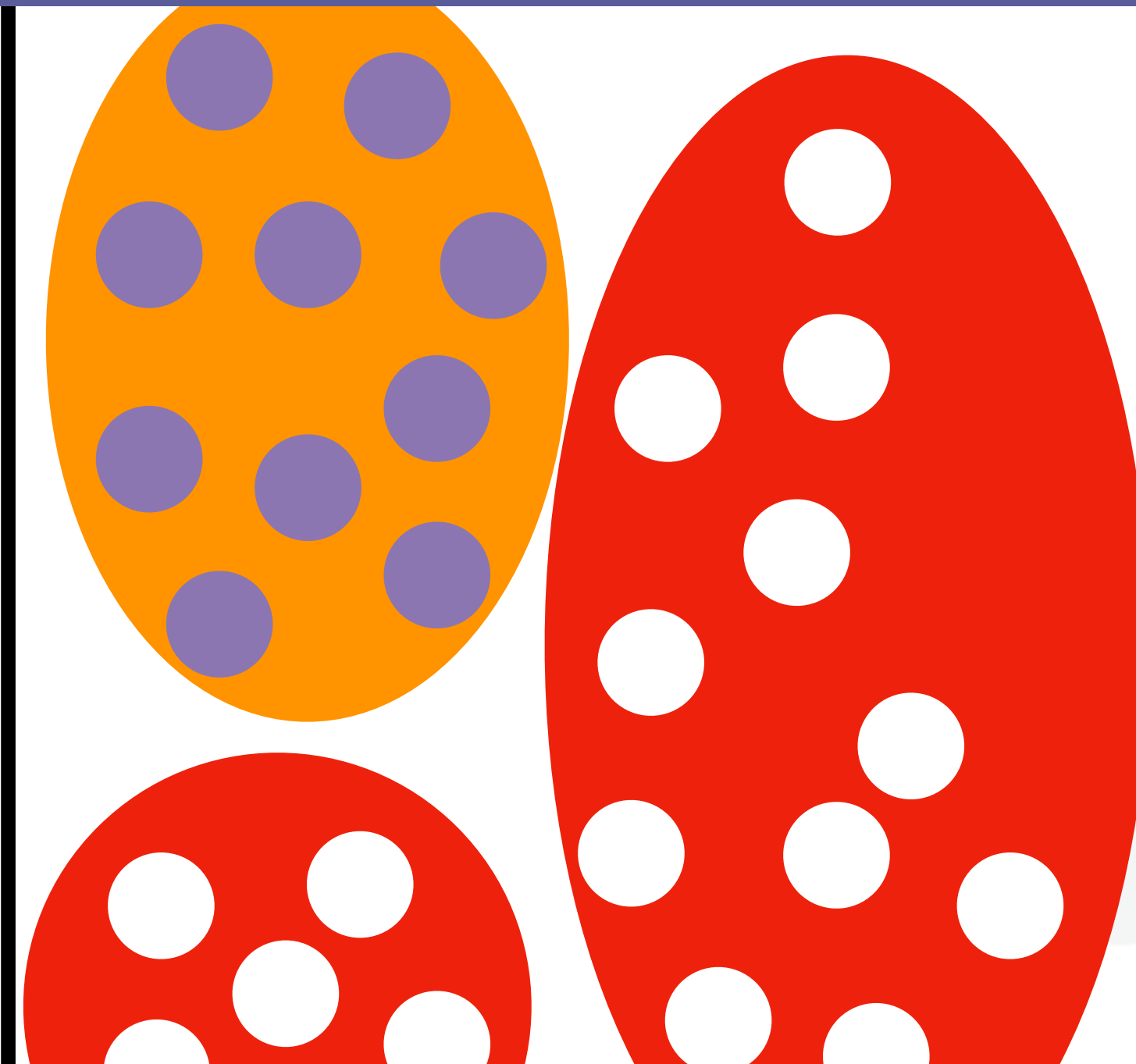
bondad de un método de clusterización

Dos conceptos clave

**MINIMIZAR la  
similaridad intra-cluster**

**MAXIMIZAR la  
similaridad inter-cluster**

Los elementos de un mismo cluster  
deben ser lo más similares posibles



Los elementos de distintos clusters  
deben ser lo menos similares posibles

# Agrupamiento

bondad de un método de clustering

Las medidas de similaridad/disimilaridad se asocian normalmente a una medida de distancia:

- Son funciones muy sensibles al tipo de variable usadas (intervalos, binarias, booleanas, categóricas, ordinales)
- Posible asignar peso a variables por criterios del experto o la propia aplicación
- Complicado de dar definiciones para términos de "suficientemente similar"

# Agrupamiento

subjetividad



**¿Cómo agruparíamos a este grupo de personas?**

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Agrupamiento

subjetividad

**¿Cómo agruparíais a este grupo de personas?  
Hombres/Mujeres**



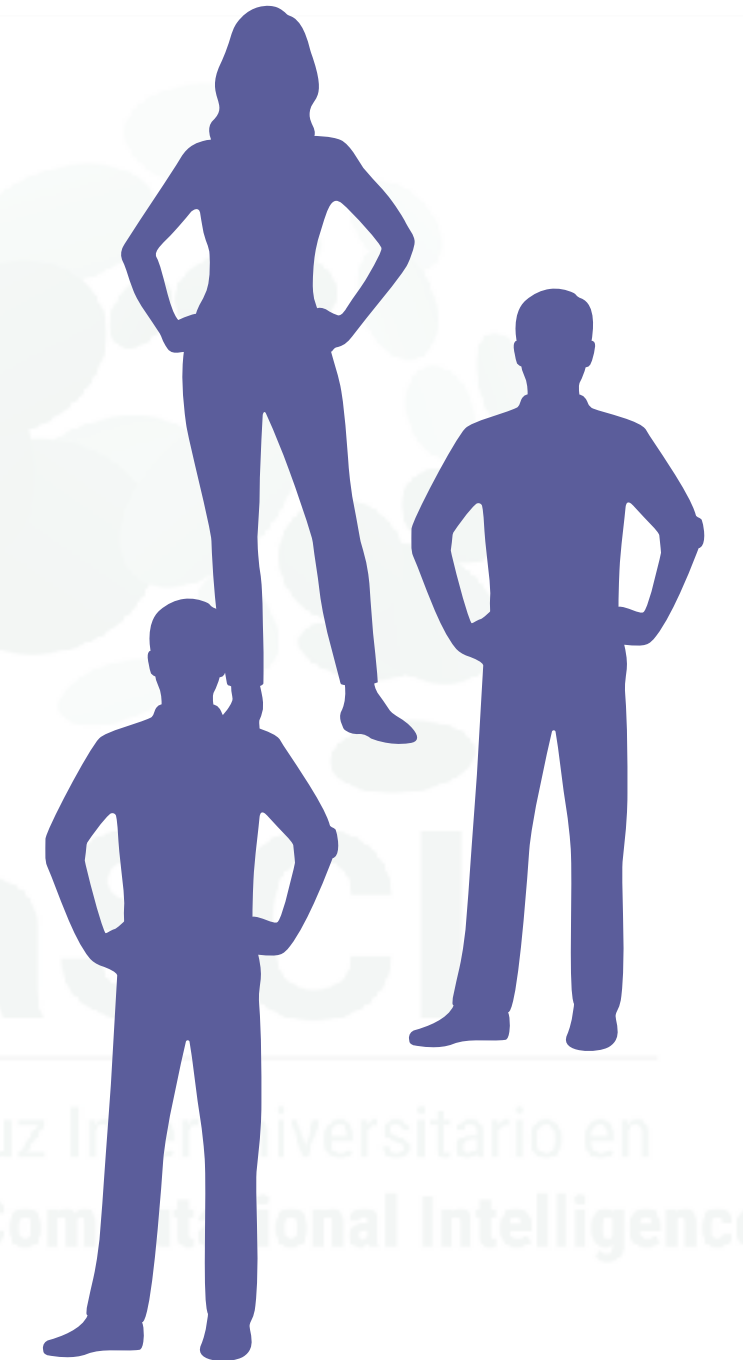
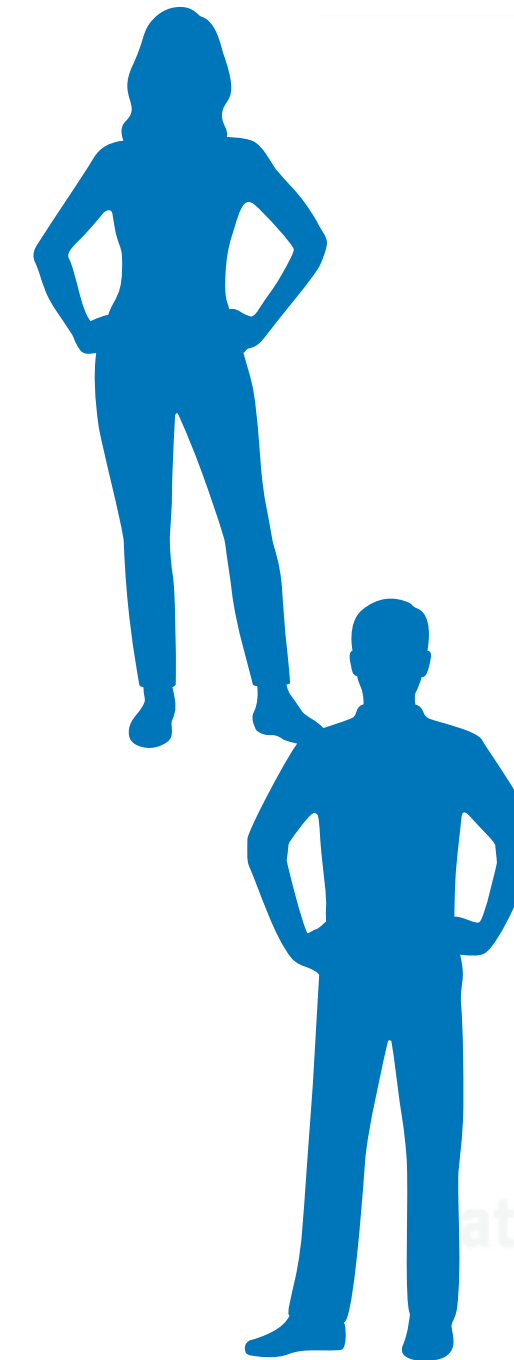
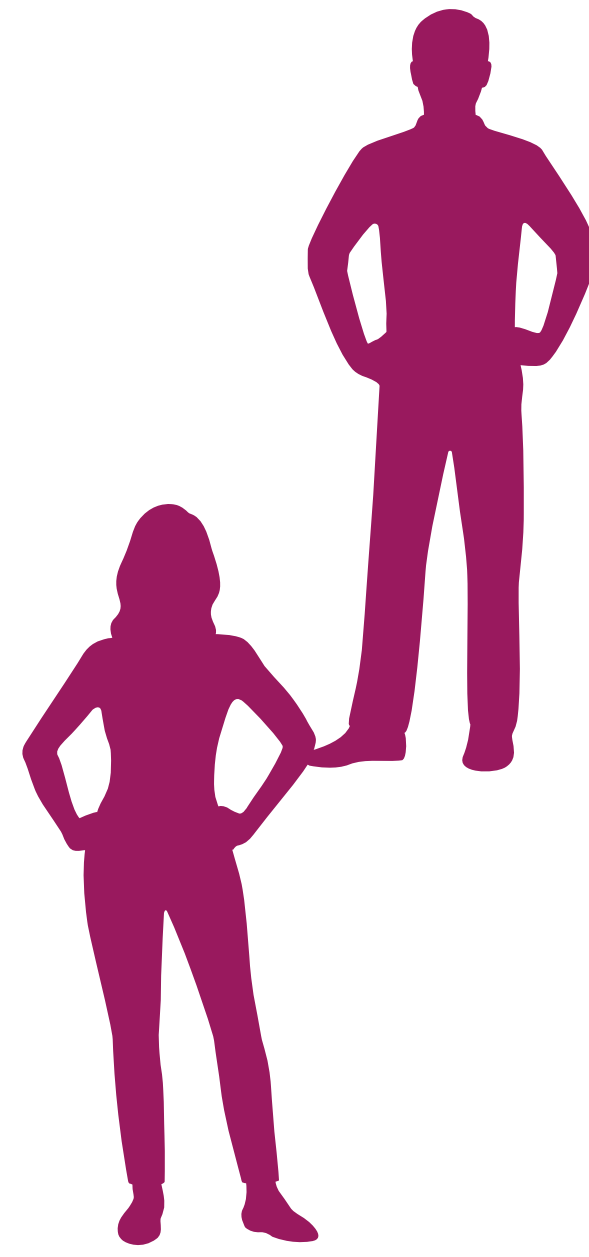
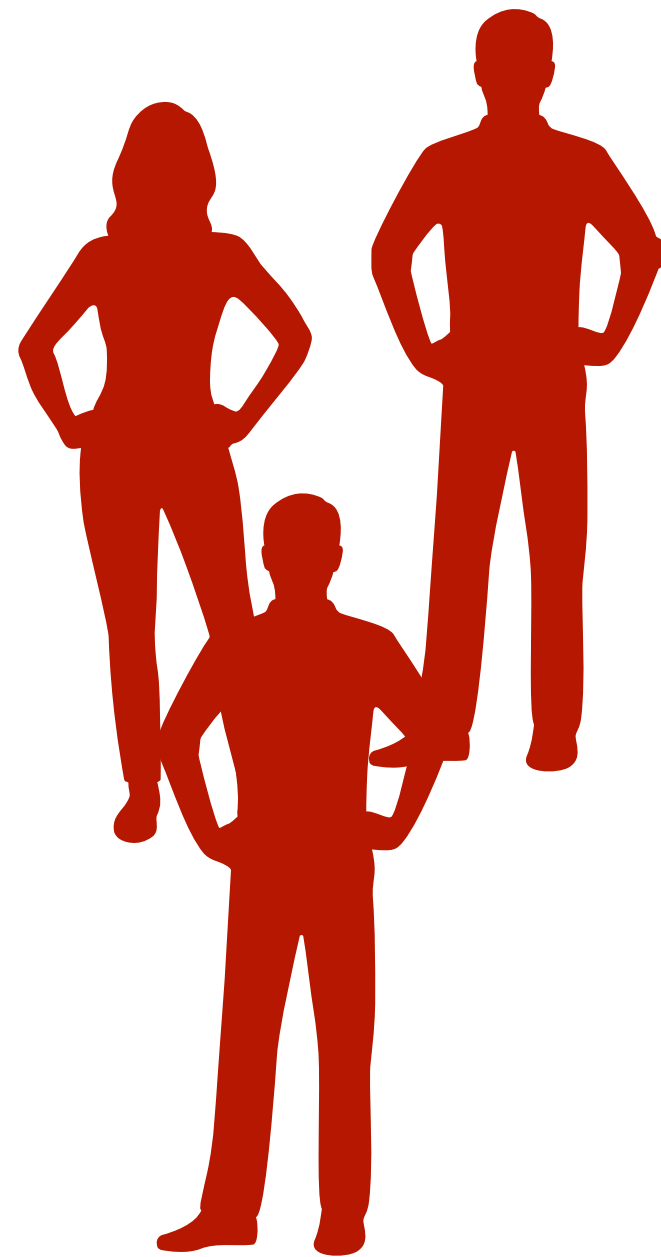


# Agrupamiento

subjetividad

¿Cómo agruparíais a este grupo de personas?

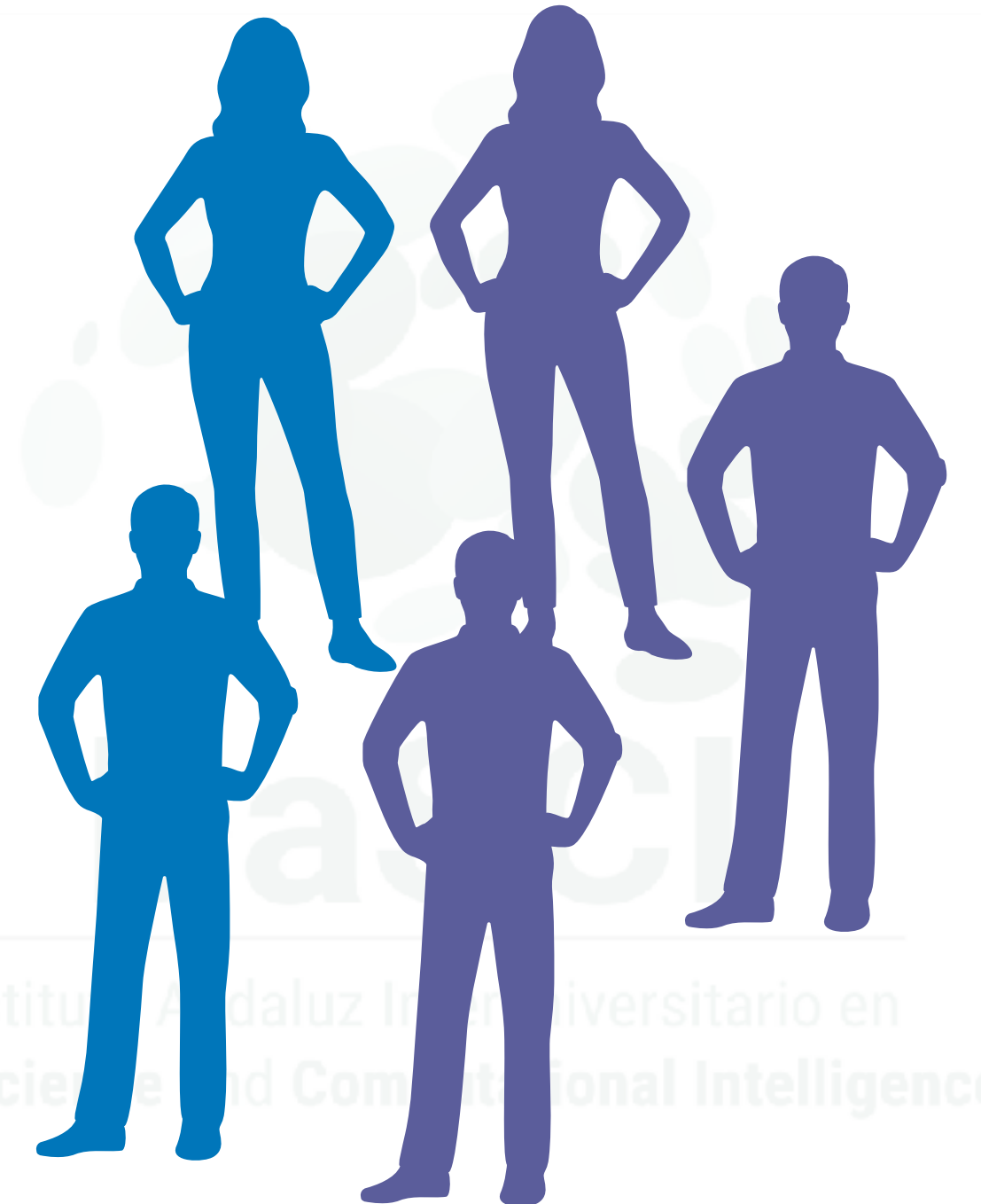
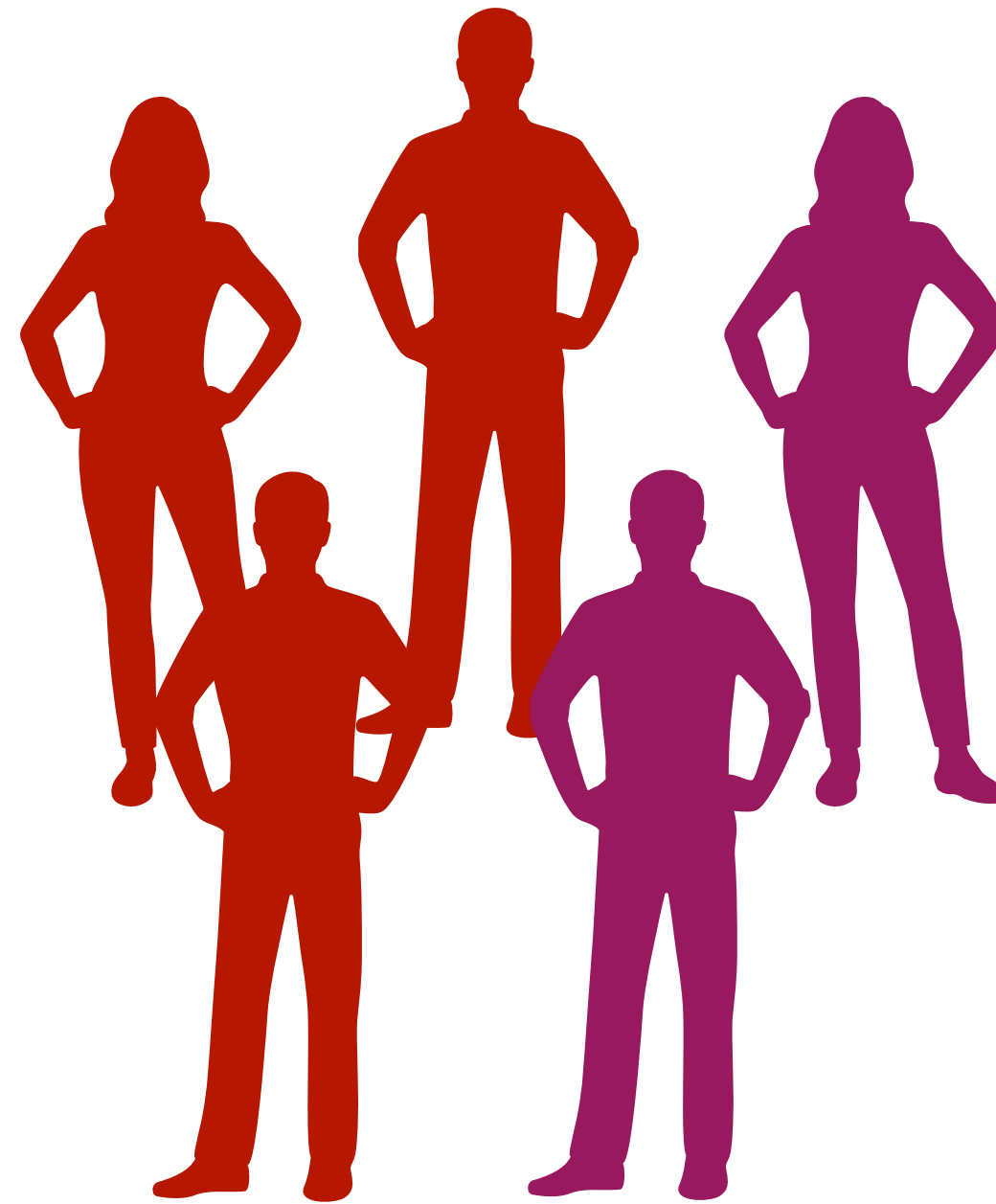
Raza



# Agrupamiento

subjetividad

**¿Cómo agruparíais a este grupo de personas?  
Por número de clusters, por ejemplo, TRES**



# Agrupamiento

propiedades deseables

escalables

insensible al orden de registros de entrada

tratar distintos tipos de variables

descubrir clusters formas arbitrarias

ruido y outliers

alta dimensionalidad

Incorporar restricciones del usuario

resultados interpretables

requisitos mínimos del problema

# Medidas de distancia y similitud

¿Qué es la similitud? Se define por la RAE como semejanza (cualidad de semejante (que semeja o se parece a alguien o algo (tener determinada apariencia o aspecto)))





# Medidas de distancia y similitud

¿Qué es la similitud? Se define por la RAE como semejanza (cualidad de semejante (que semeja o se parece a alguien o algo (tener determinada apariencia o aspecto)))



# Medidas de distancia y similitud

id	sexo	nacimiento	categoría	salario	experiencia	minoría
25	Mujer	6-10-36	Administ	18750	54	No
26	Mujer	26-9-65	Administ	38550	22	No
27	Hombre	6-10-60	Administ	27450	173	Sí
28	Hombre	21-1-51	Seguridad	24300	191	Sí
29	Hombre	1-9-50	Seguridad	30750	209	Sí
30	Mujer	25-7-46	Directivo	68750	38	No
31	Hombre	18-7-59	Administ	19650	229	Sí
32	Hombre	6-9-58	Directivo	59375	6	No

# Medidas de distancia y similitud



id	sexo	nacimiento	categoria	salario	experiencia	minoría
25	Mujer	6-10-36	Administ	18750	54	No
26	Mujer	26-9-65	Administ	38550	22	No
27	Hombre	6-10-60	Administ	27450	173	Sí
28	Hombre	21-1-51	Seguridad	24300	191	Sí
29	Hombre	1-9-50	Seguridad	30750	209	Sí
30	Mujer	25-7-46	Directivo	68750	38	No
31	Hombre	18-7-59	Administ	19650	229	Sí
32	Hombre	6-9-58	Directivo	59375	6	No

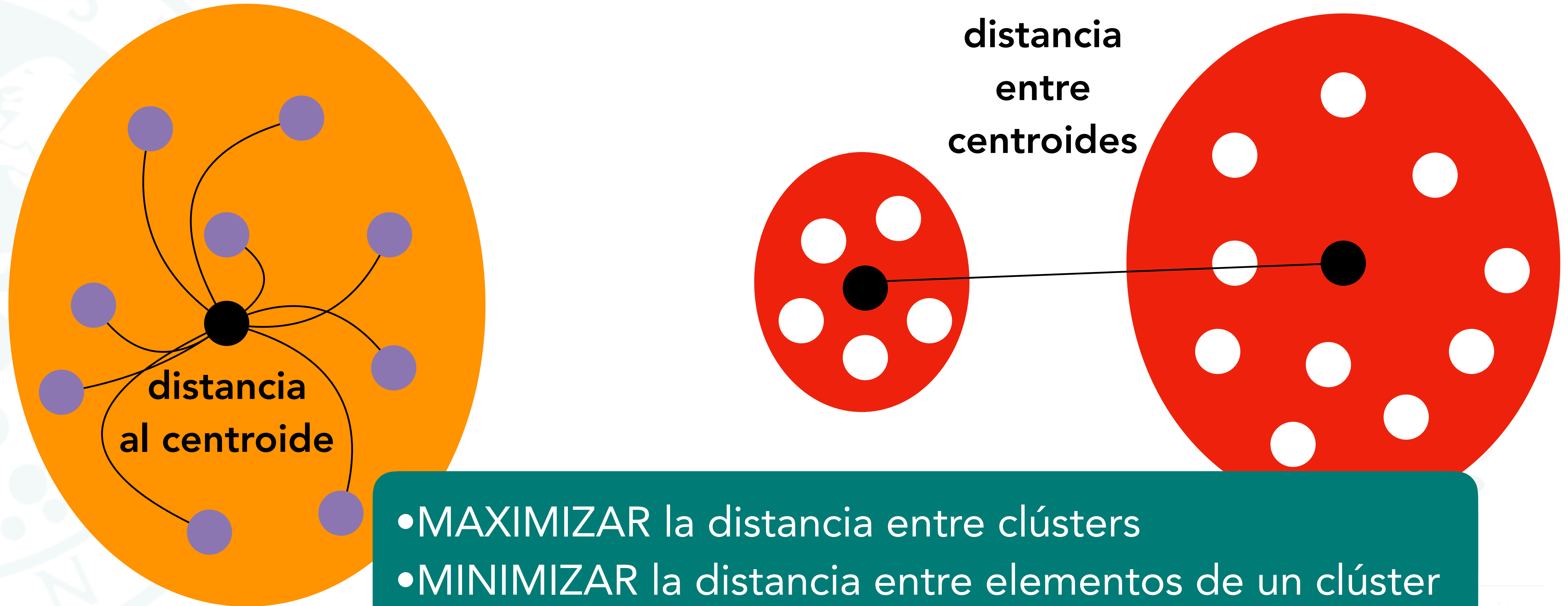
# Medidas de distancia y similitud

id	sexo	nacimiento	categoría	salario	experiencia	minoría
25	Mujer	6-10-36	Administrativa	18750	54	No
26	Mujer	26-9-65	Administrativa	38550	22	No
27	Hombre	6-10-60	Administrativa	27450	173	Sí
31	Hombre	18-7-59	Administrativa	19650	229	Sí

- Dentro de un mismo grupo tendríamos variables que nos podrían dividir el grupo en otros subgrupos o categorías
- No es necesario utilizar todos los atributos para agrupar objetos o elementos



# Medidas de distancia y similitud



- MAXIMIZAR la distancia entre clústers
- MINIMIZAR la distancia entre elementos de un clúster

# Medidas de distancia y similitud

La definición de la medida de distancia suele depender del tipo de variable:

- Variables intervalares
- Variables binarias o booleanas (verdadero/falso)
- Variables nominales o categóricas
- Variables ordinales



# Medidas de distancia y similitud

variables numéricas

Son las más sencillas de calcular

$$\text{Distancia } (X, Y) = A(X) - A(Y)$$

Cuando tenemos  $n$  dimensiones se utiliza la famosa

**distancia Euclidea**

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_i - q_i)^2 + \dots + (p_n - q_n)^2} = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}.$$

# Medidas de distancia y similitud

variables nominales

Son también muy sencillas de calcular porque la distancia se fija a 1 si los valores son diferentes, y a 0 si son iguales

Una variación sería ponderar las variables según su importancia en el problema

$$d(i, j) > d(i, k)$$



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Medidas de distancia y similitud

variables nominales

Son también muy sencillas de calcular porque la distancia se fija a 1 si los valores son diferentes, y a 0 si son iguales

Una variación sería ponderar las variables según su importancia en el problema

$$d(i, j) > d(i, k)$$

**el elemento  $i$  es más parecido a  $k$  que al elemento  $j$**

# Medidas de distancia y similitud

variables continuas

Para evitar que unas variables dominen sobre otras, los valores de los atributos se "normalizan" a priori:

- Desviación absoluta media

$$s_f = \frac{1}{n} (|x_{1f} - m_f| + \dots + |x_{nf} - m_f|)$$

$$m_f = \frac{1}{n} (x_{1f} + \dots + x_{nf})$$

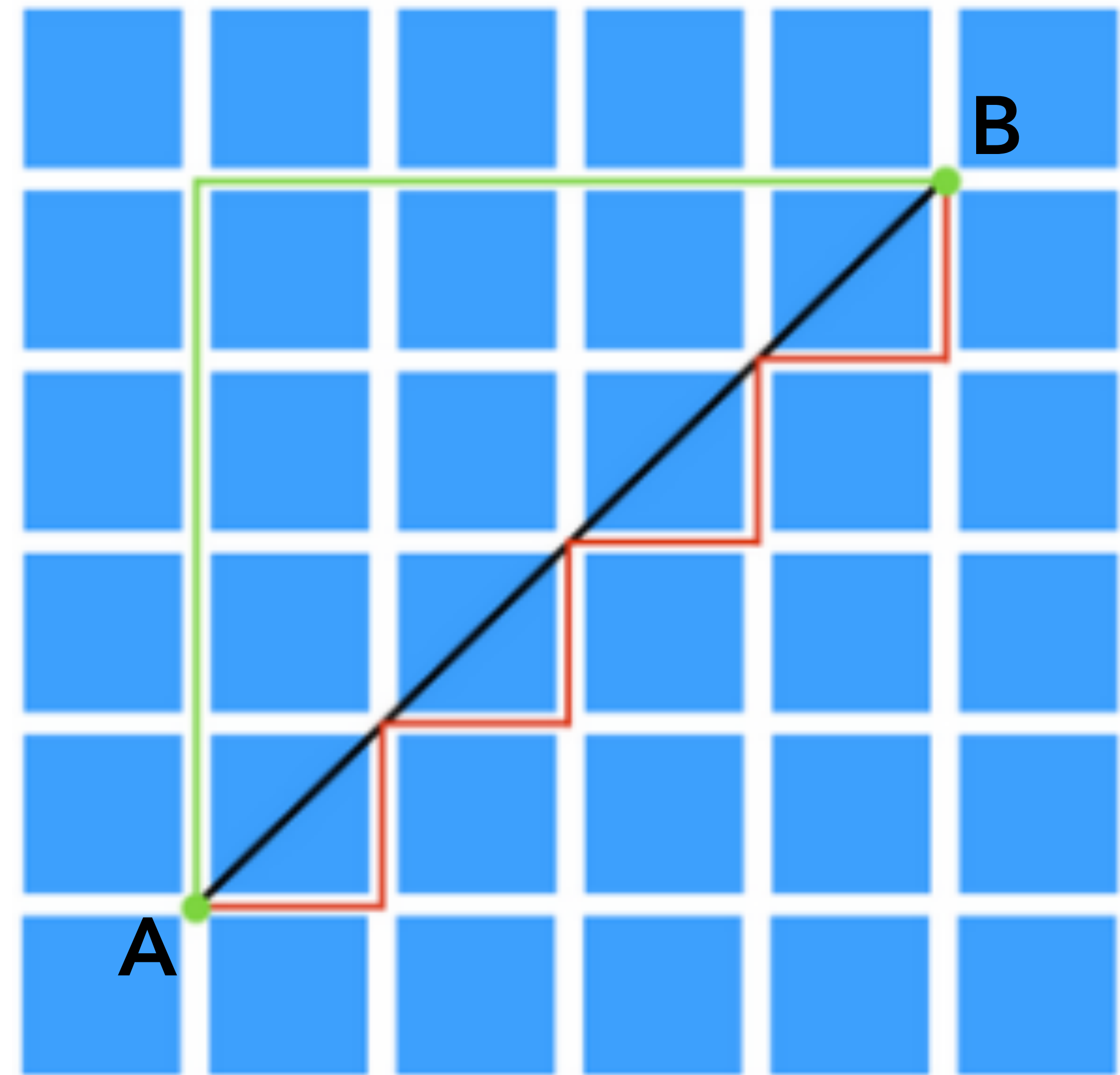
- z-score (medida estandarizada)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

# Medidas de distancia y similitud

distancia de Minkowski

- La distancia entre dos puntos (AB) se calcula mediante la distancia **Euclidea** (línea negra)
- Si consideramos un mapa cuadrículado como el de las calles de la isla de Manhattan hace que el camino más corto posible en taxi sea exactamente la distancia de **Manhattan** (línea roja y verde)



# Medidas de distancia y similitud

distancia de Minkowski

$$d_r(x, y) = \left( \sum_{j=1}^J |x_j - y_j|^r \right)^{\frac{1}{r}}$$

- Si  $r=1$  es Distancia de Manhattan

$$d_1(x, y) = \sum_{j=1}^J |x_j - y_j|$$

- Si  $r=2$  es Distancia Euclidea

$$d_2(x, y) = \sqrt{\sum_{j=1}^J |x_j - y_j|^2}$$

- Si  $r=\infty$  es Distancia de Chebyshev



# Medidas de distancia y similitud

otras distancias

Distancia de Mahalanobis

- Correlación entre variables y no depende de la escala de medida

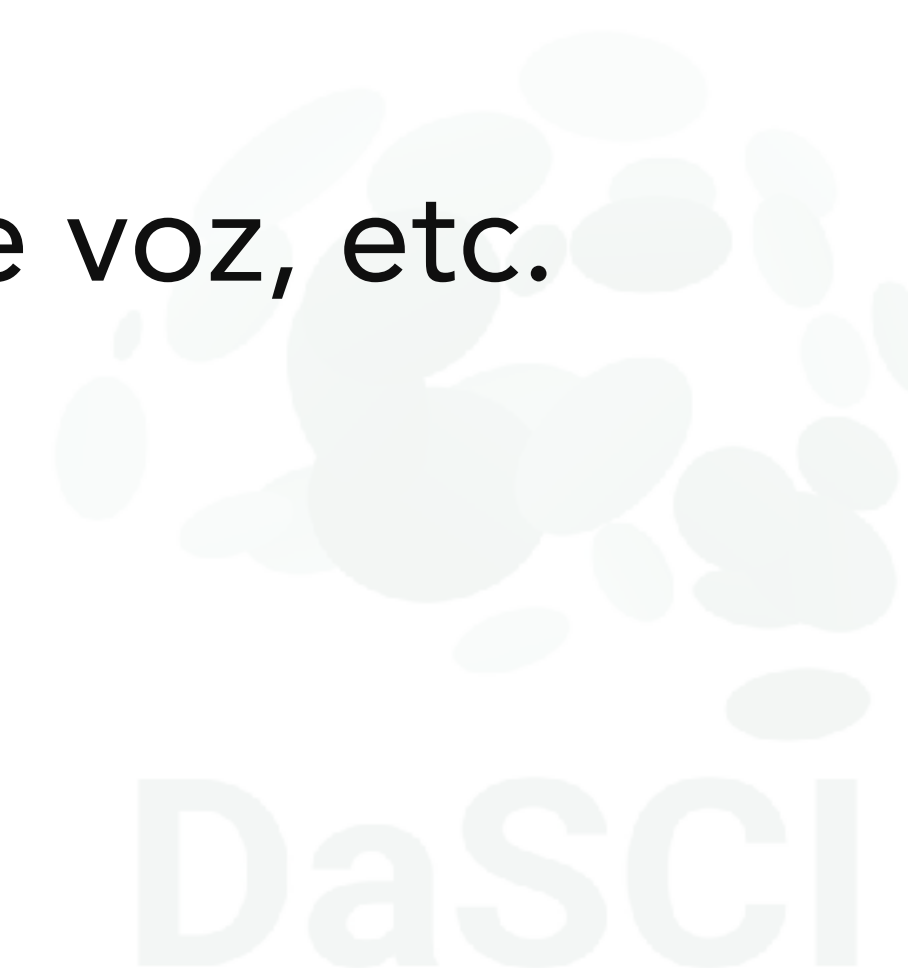
Distancia de Levenshtein

- Empleada en correctores ortográficos, reconocimiento de voz, etc.

Modelos basados en teorías de conjuntos

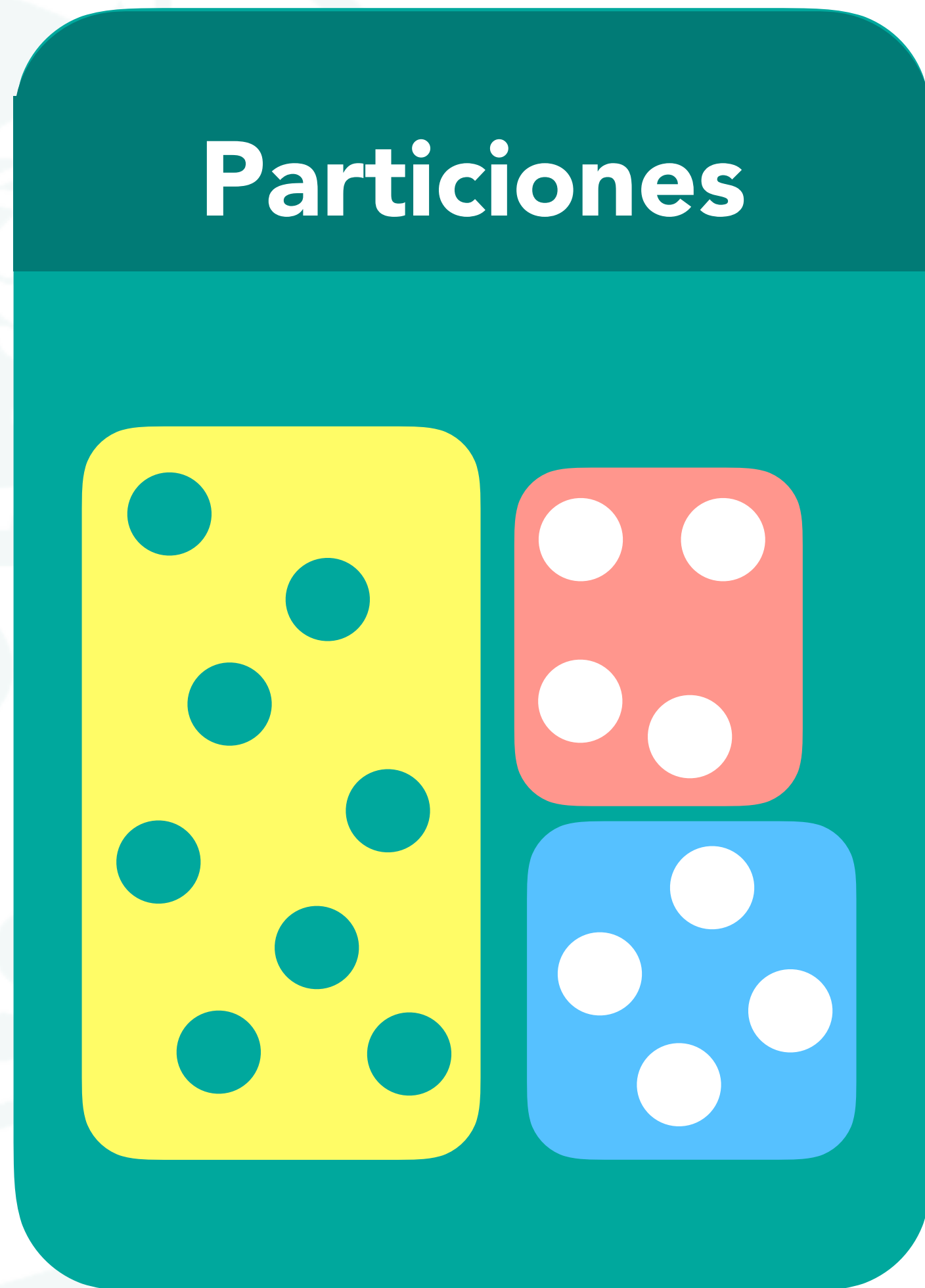
Distancia de vecinos compartidos

Medidas de correlación

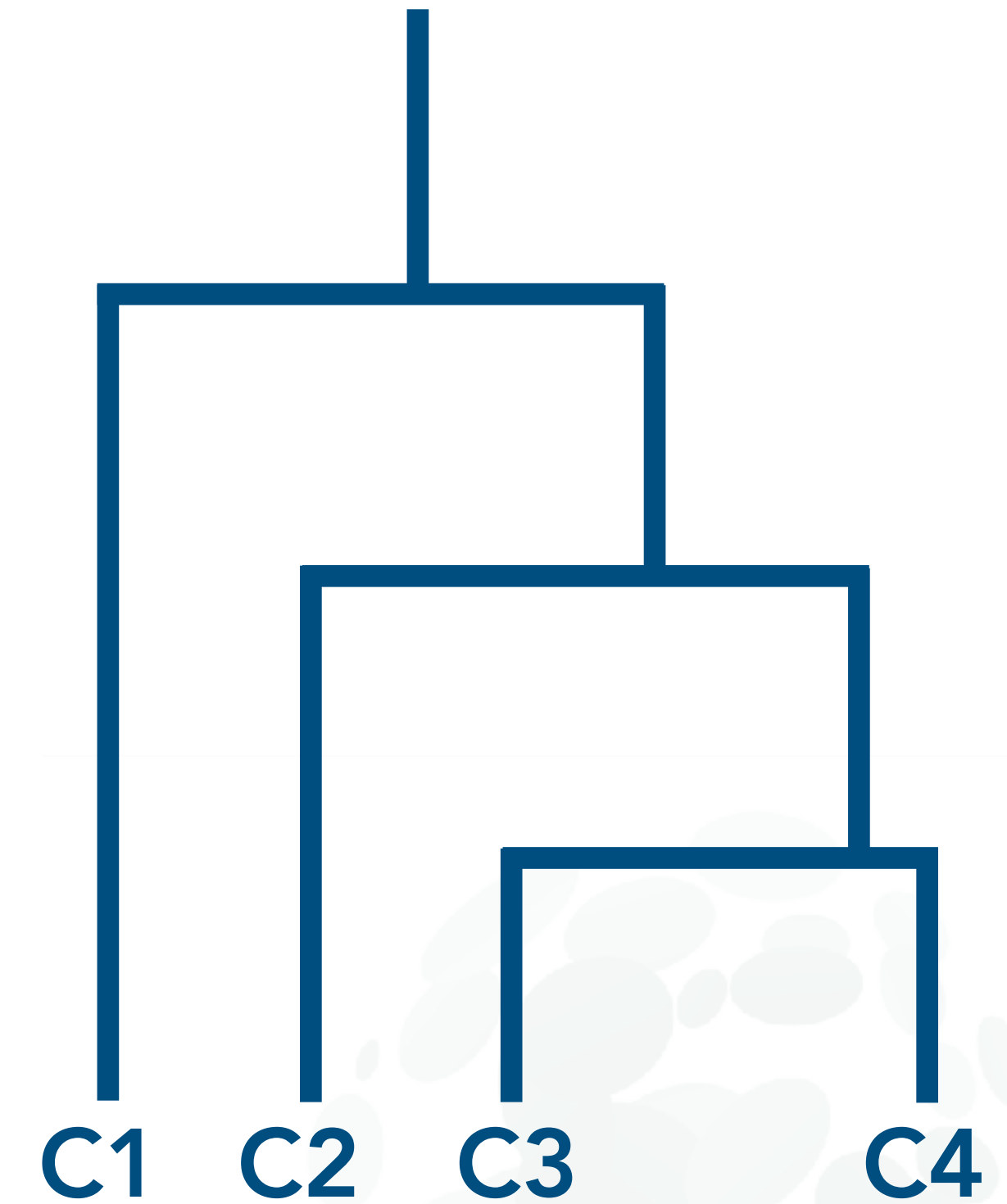


Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Clasificación de algoritmos de clustering



## DENDOGRAMA



**Aglomerativo**

**Divisivo**

# Clasificación de algoritmos de clustering

## Exclusiva versus no exclusiva

- No exclusivos los elementos pueden pertenecer a varios grupos, y pueden representar múltiples "fronteras"

## Difuso versus no difuso

- A un elemento se le asocia un peso entre 0 y 1, y todos los pesos del cluster deben sumar 1.
- El agrupamiento probabilístico tiene características similares

## Parcial versus completa

## Heterogéneo versus homogéneo

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Método basado en particionamiento clásico

## k-Means

Se basa en la idea de ir moviendo entre clusters hasta que se alcanza el número de clusters deseado

Cada cluster se representa por el centro del mismo, también conocido como CENTROIDE

Parámetro de entrada ( $k$ ) que es el número de clusters

OBJETIVO: Minimizar distancia euclídea total entre cada punto y su representante de cluster más cercano de forma iterativa



# Método basado en particionamiento clásico

---

## Algorithm 2.1 The k-means algorithm

---

**Input:** Dataset  $D$ , number clusters  $k$

**Output:** Set of cluster representatives  $C$ , cluster membership vector  $\mathbf{m}$

/\* Initialize cluster representatives  $C$  \*/

Randomly choose  $k$  data points from  $D$

5: Use these  $k$  points as initial set of cluster representatives  $C$

**repeat**

/\* Data Assignment \*/

Reassign points in  $D$  to closest cluster mean

Update  $\mathbf{m}$  such that  $m_i$  is cluster ID of  $i$ th point in  $D$

10: /\* Relocation of means \*/

Update  $C$  such that  $c_j$  is mean of points in  $j$ th cluster

**until** convergence of objective function  $\sum_{i=1}^N (\operatorname{argmin}_j \|\mathbf{x}_i - \mathbf{c}_j\|_2^2)$

---

## k-Means

Selección aleatoria de  $k$  centroides

Asignar puntos a cada centroide

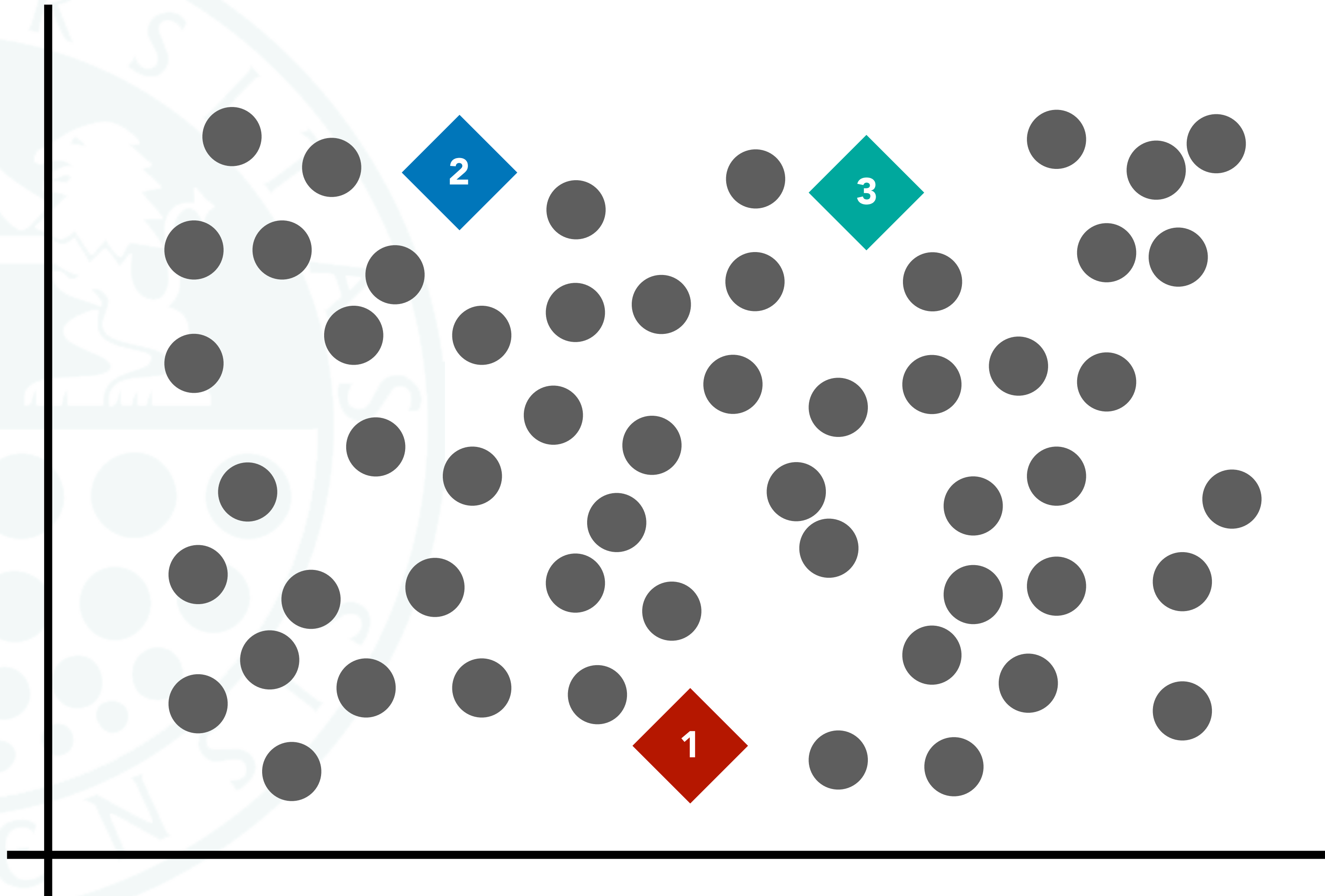
Re-colocación de medias

0 asignaciones

# Método basado en particionamiento clásico

## k-Means

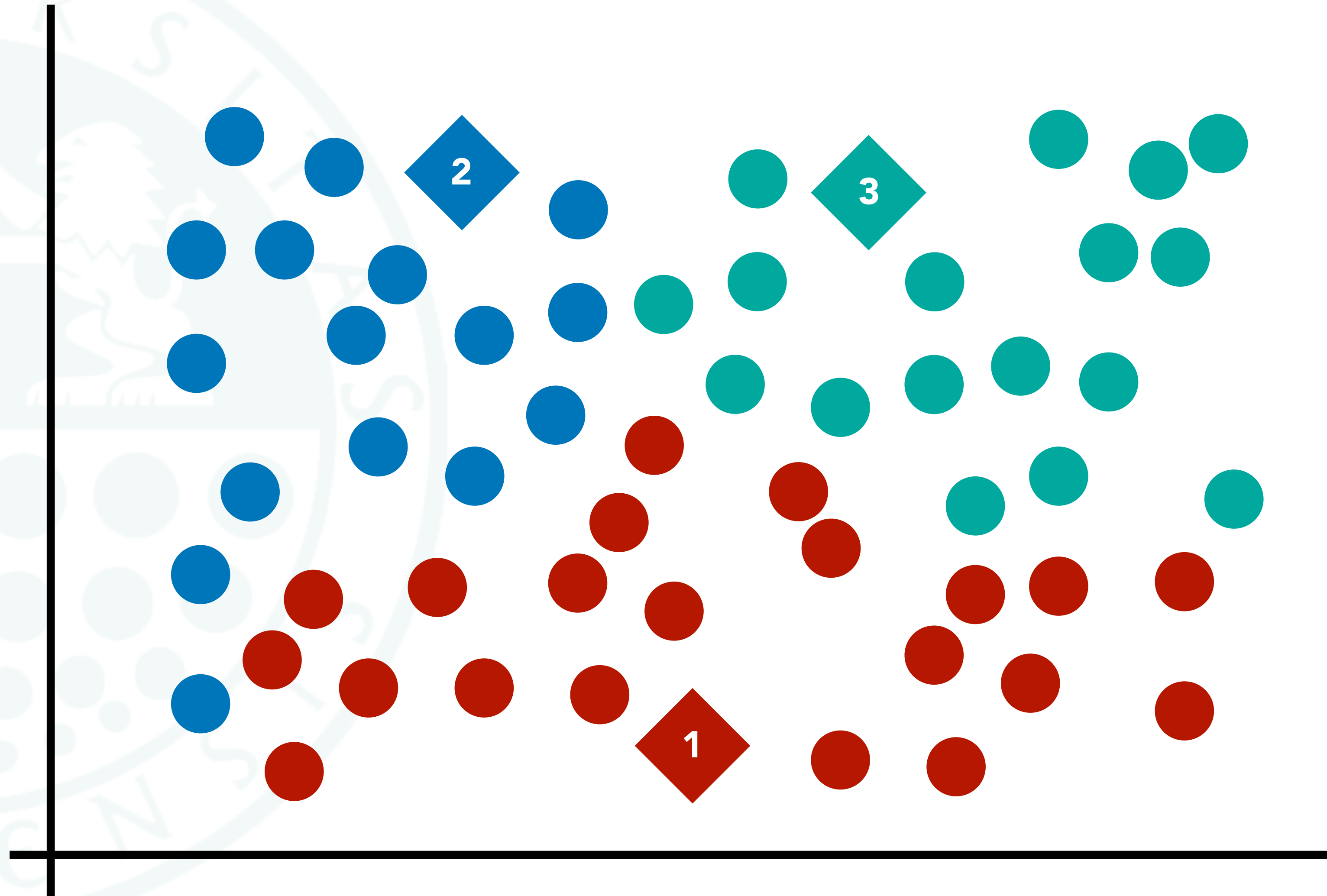
Selección aleatoria de  $k=3$  centroides



Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Método basado en particionamiento clásico



## k-Means

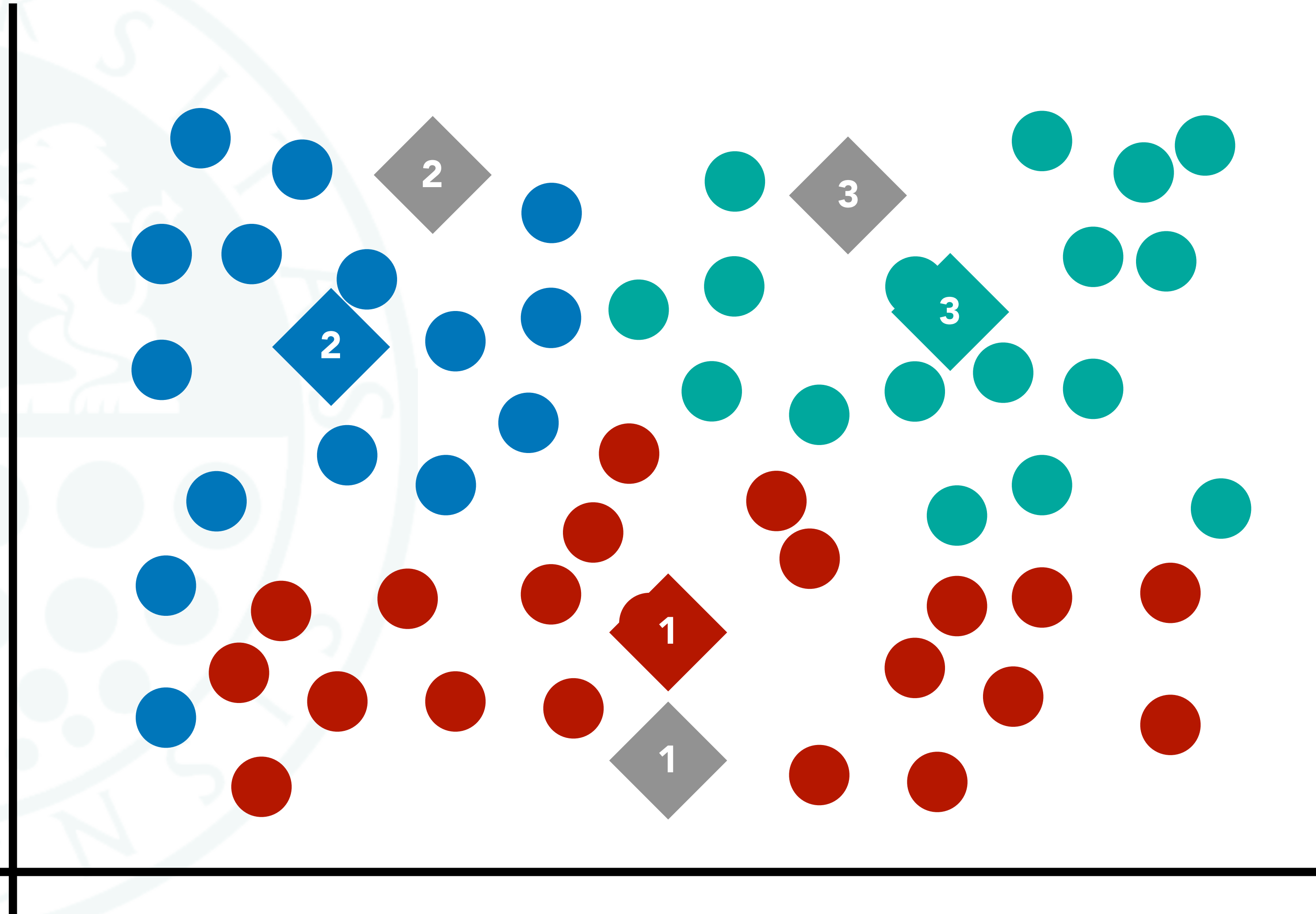
Selección aleatoria de  $k=3$  centroides

Asignar puntos a cada centroide

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence

# Método basado en particionamiento clásico



## k-Means

Selección aleatoria de  $k=3$  centroides

Asignar puntos a cada centroide

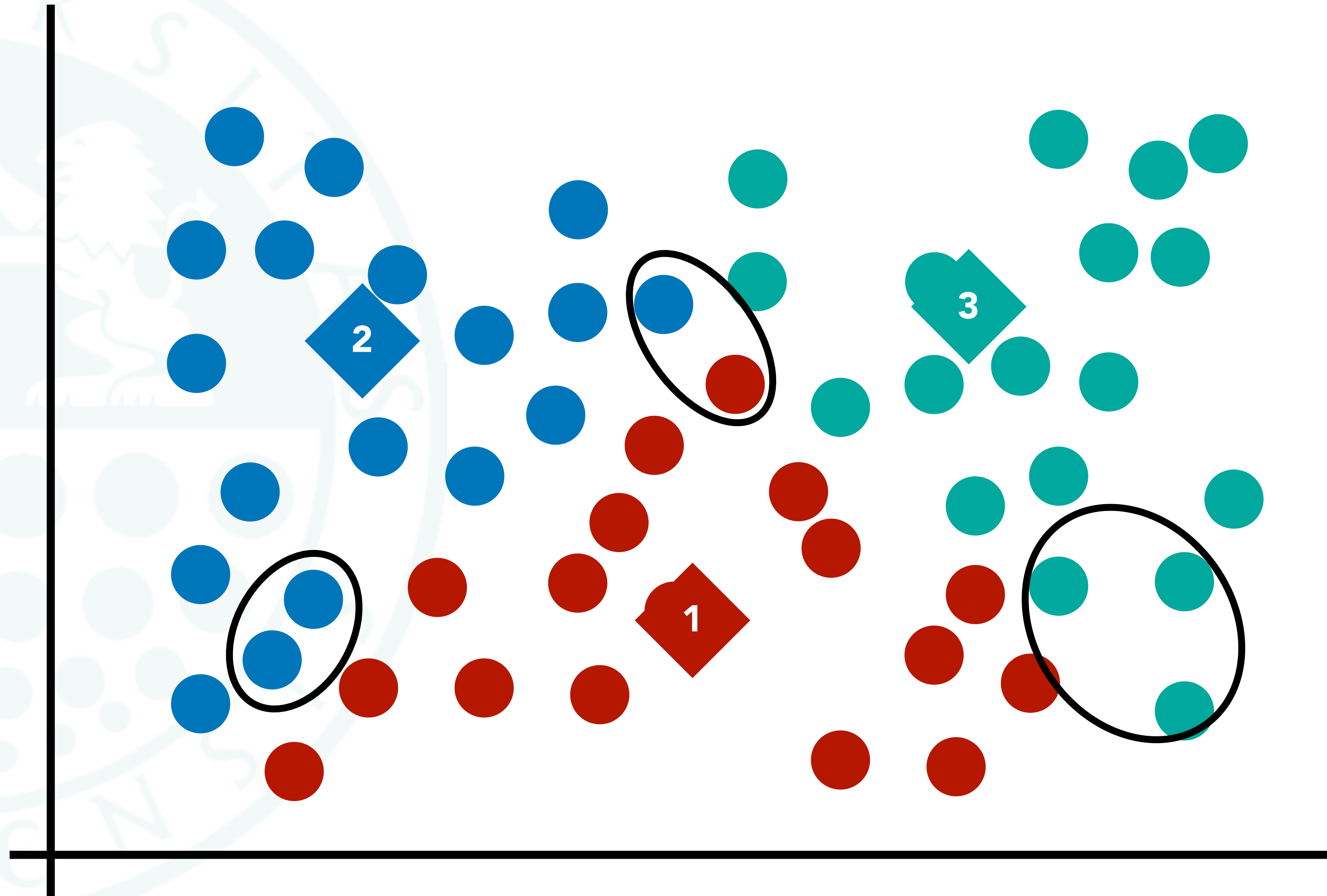
Re-colocación de medias

DaSCI

Instituto Andaluz Interuniversitario en  
Data Science and Computational Intelligence



# Método basado en particionamiento clásico

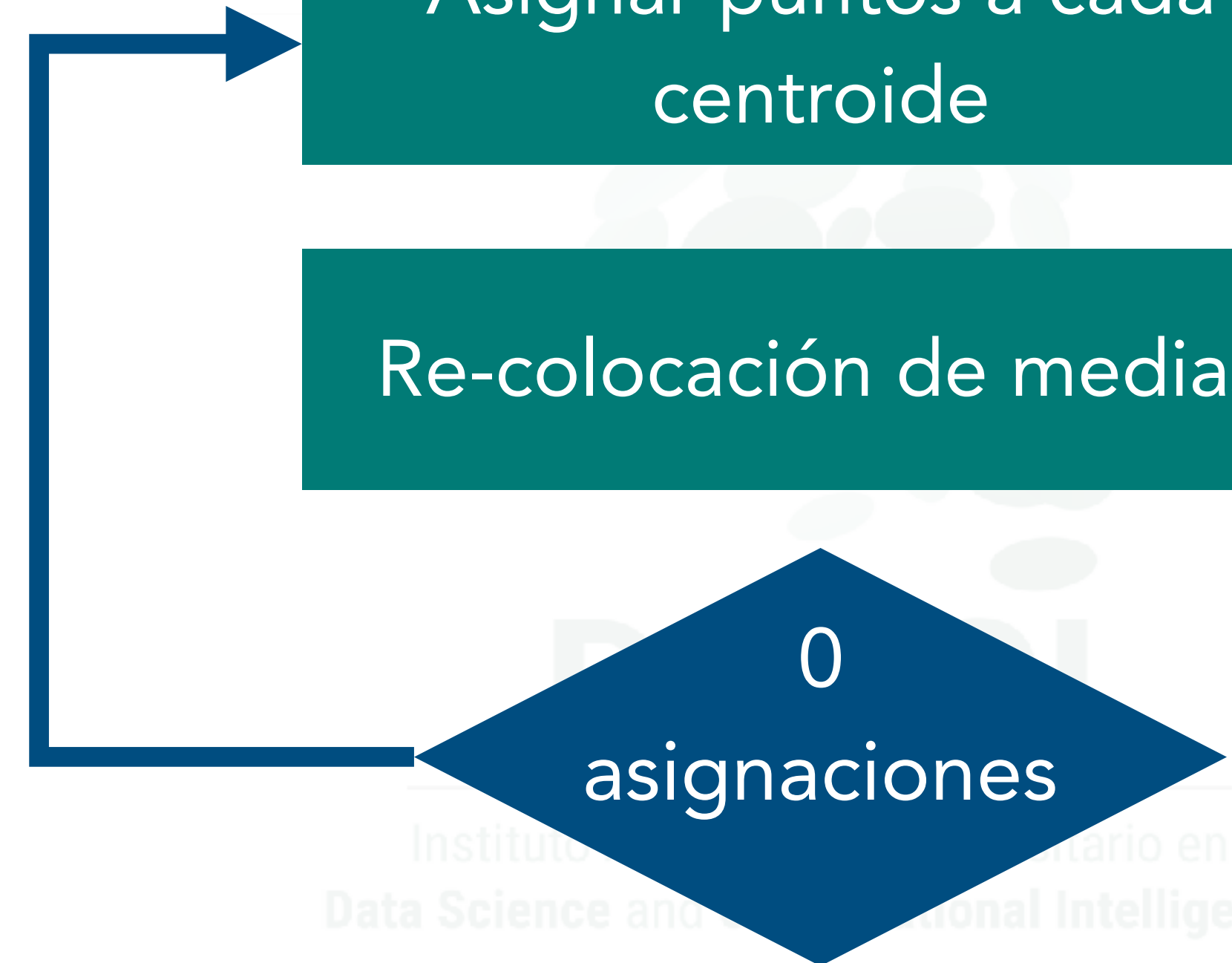


## k-Means

Selección aleatoria de  $k=3$  centroides

Asignar puntos a cada centroide

Re-colocación de medias



# Método basado en particionamiento clásico

## k-Means

### VENTAJAS

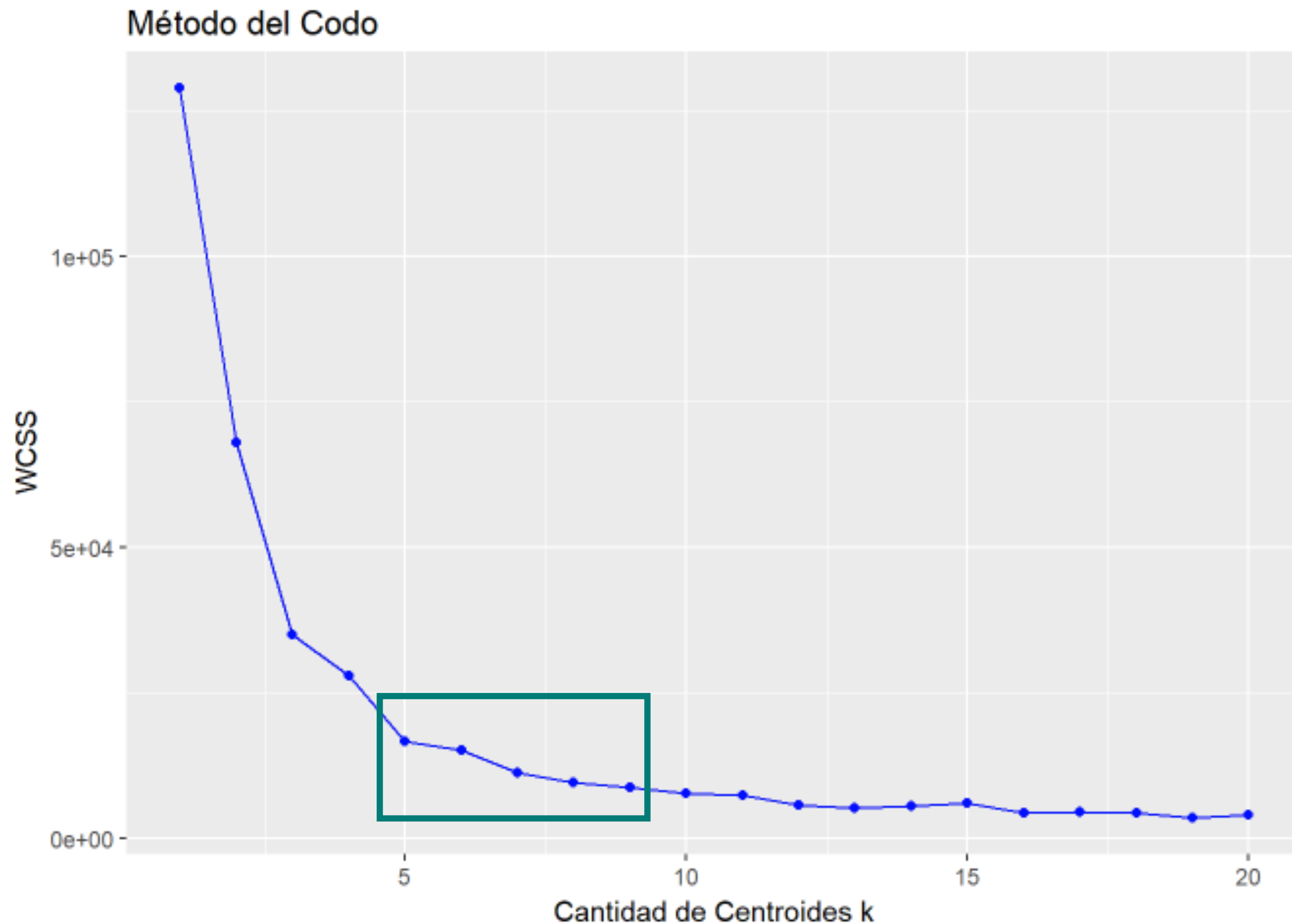
- Relativamente eficiente  
 $O(tkn)$  —  $t$  son iteraciones,  
 $k$  es clusters y  $n$  objetos
- Finaliza de forma frecuente  
en un óptimo local

### DESVENTAJAS

- Aplicable si se puede definir  
el concepto de media
- Necesidad de fijar clusters
- Débil en ruido/outliers
- Sólo indicado en clusters  
convexos

## Método basado en particionamiento clásico

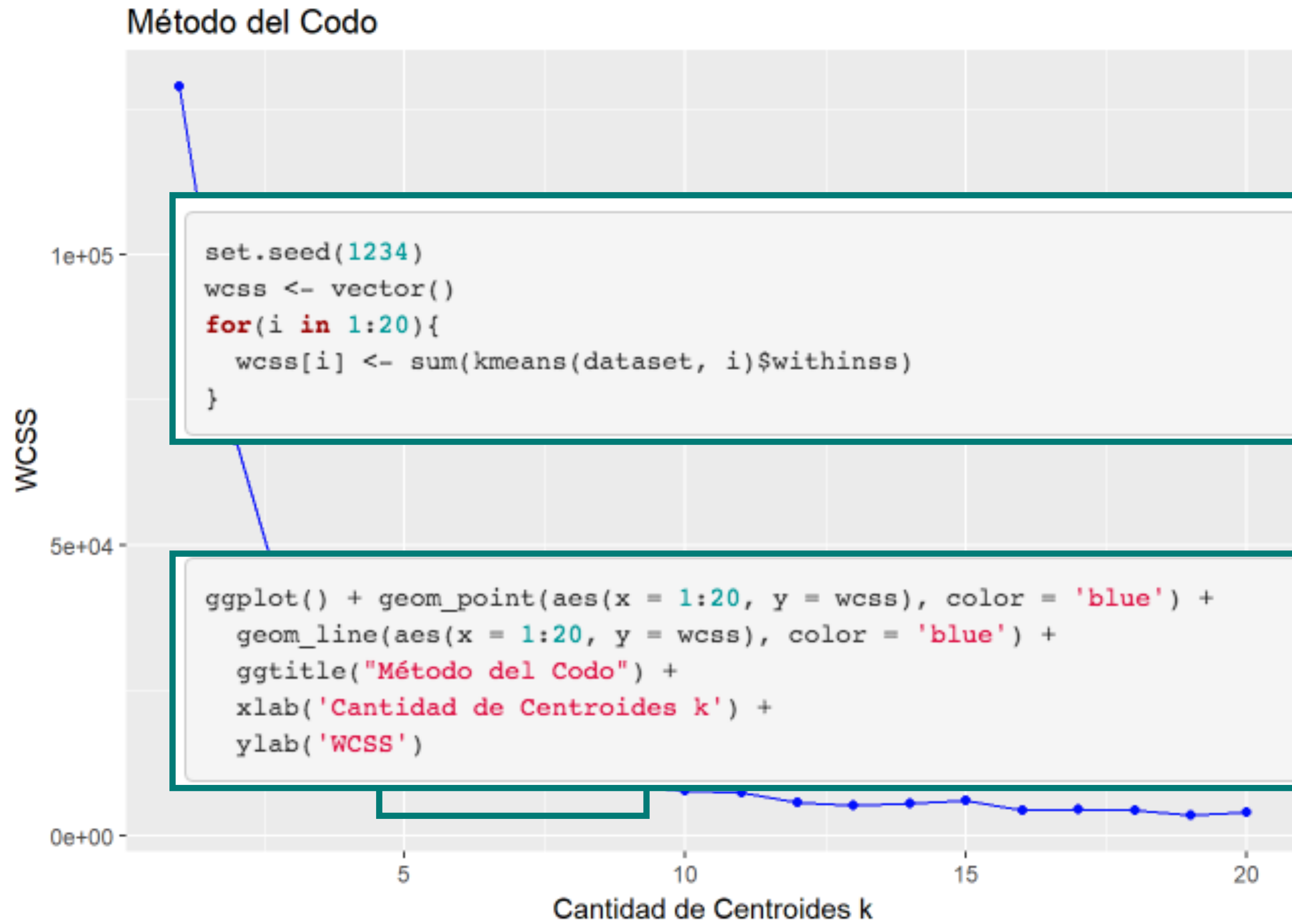
### k-Means Método del codo



- Al aumentar los centroides, disminuye el WCSS
- Se genera una forma de codo.
- El valor óptimo de k se escoge donde ya no se dejan de producir variaciones importantes del valor de WCSS (Within Clusters Summed Squares).
- A partir de  $k \geq 5$ , por lo que sería interesante evaluar los resultados con valores como 5, 6 y 7.

# Método basado en particionamiento clásico

## k-Means Método del codo



• Al aumentar los centroides,

donde ya no se dejan de producir

• A partir de  $k \geq 5$ , por lo que sería interesante evaluar los resultados con valores como 5, 6 y 7.



# Clustering en R

paquetes más citados

Nombre	Citas
apcluster	3682
akmedoids	1220
clustMixType	925
mclust	885
conclust	532
EMMIXgene	338

**View: Cluster**  
<https://cran.r-project.org/web/views/Cluster.html>

**Listado de paquetes**  
[\\_clustering\\_cran\\_r.xlsx](#)

# Clustering en R

paquete CLUSTERING

Es un nuevo paquete que hemos desarrollado que permite analizar datos no etiquetados (aprendizaje no supervisado) con distintos objetivos:

- Comparativa de algoritmos de la literatura
- Ranking de las variables que más influyen a la hora de agrupar los datos por diferentes medidas
  - Externas: Purity, Entropy, Recall, Precision, etc
  - Internas: Dunn, Connectivity y Silhoutte

## Clustering en R

paquete CLUSTERING

<https://cran.r-project.org/web/packages/Clustering/Clustering.pdf>

The screenshot displays the RStudio interface. The left pane shows an R script with the following code:

```
1 library(Clustering)
2
3 # Realizamos comparativa de métodos de clustering con
4 # nuestros datos.
5 comparativa <- clustering(df = Clustering::basketball, # datos
6                           packages = NULL,           # Paquetes de clustering a ejecutar: NULL = todos
7                           algorithm = NULL,         # Algoritmos de clustering a ejecutar: NULL = todos
8                           min = 3,                 # nº mínimo de clusters
9                           max = 6)                 # nº máximo de variables
10
11
12 # Obtener un resumen de los resultados (valores medios de métricas en cada
13 # algoritmo, etc.), mediante summary()
14 resumen <- summary(comparativa)
15 print(resumen)
16
17
18 # Determinar que número de clusters es el óptimo en nuestra comparativa.
19 # Por ejemplo, comparamos quien tiene mayor precisión
20 plot_clustering(comparativa, "precision")
21
22
23 # Tras la visualización de los resultados, queremos quedarnos únicamente con
24 # aquellos métodos cuya precisión sea mayor a 0.15 y cuyo recall sea mayor a 0.5:
25 resultadoFiltrado <- resultadoFiltrado <- comparativa[precision > 0.15 & recall > 0.5]
26
27
28 # Finalmente, exportamos nuestros resultados filtrados a LaTeX:
29 # Se guardará un fichero llamado "external_data.tex" en el directorio actual.
30 export_file_external(df = resultadoFiltrado)
31
```

The right pane shows the documentation for the 'Clustering' package version 1.6.1, titled "Execution of Multiple Clustering Algorithm". It includes a "DESCRIPTION file" link and a "Help Pages" section with links to various datasets and functions:

- [appClustering](#): Clustering GUI.
- [basketball](#): This data set contains a series of statistics (5 attributes) about 96 basketball players.
- [best\\_ranked\\_external\\_metrics](#): Best rated external metrics.
- [best\\_ranked\\_internal\\_metrics](#): Best rated internal metrics.
- [bolts](#): Data from an experiment on the affects of machine adjustments on the time to count bolts.
- [clustering](#): Clustering algorithm.
- [datasetTest](#): Dataset with training data. The data is obtained from running the clustering algorithm with the basketball dataset. Of all the packages we use the cluster package and all the metrics. The number of clusters is between 3 and 4.
- [evaluate\\_best\\_validation\\_external\\_by\\_metrics](#): Evaluation of the algorithms by measures of dissimilarity.
- [evaluate\\_best\\_validation\\_internal\\_by\\_metrics](#): Evaluation of the algorithms by measures of dissimilarity.



## Clustering en R

paquete CLUSTERING

<https://cran.r-project.org/web/packages/Clustering/Clustering.pdf>

The screenshot displays the RStudio interface. The top-left pane shows R code for evaluating clustering algorithms. The bottom-left pane shows the console output with various metrics. The bottom-right pane shows a plot of Precision vs. Clustering (number of clusters).

```
16
17
18 # Determinar que número de clusters es el óptimo en nuestra comparativa.
19 # Por ejemplo, comparamos quien tiene mayor precision
20 plot_clustering(comparativa, "precision")
21
22
23 # Tras la visualización de los resultados, queremos quedarnos únicamente con
24 # aquellos métodos cuya precisión sea mayor a 0.15 y cuyo recall sea mayor a 0.5:
25 resultadoFiltrado <- resultadoFiltrado <- comparativa[precision > 0.15 & recall > 0.5]
26 print(summary(resultadoFiltrado))
27
28
29 # Finalmente, exportamos nuestros resultados filtrados a LaTeX:
30 # Se guardara un fichero llamado "external_data.tex" en el directorio actual.
31 export_file_external(df = resultadoFiltrado)
32
23:1 (Top Level) =
```

Environment History Connections Build Git Tutorial

Data

Object	Class
comparativa	List of 5
resultadoFiltrado	List of 5
resumen	List of 5

Console

```
~/Escritorio/clustering/
Mean time for evaluation of external metrics:
[1] "6.0311"
```

Metric mean entropy:  
[1] "6.2538"

Metric mean variation\_information:  
[1] "5.364"

Metric mean precision:  
[1] "6.166"

Metric mean recall:  
[1] "6.5964"

Metric mean f\_measure:  
[1] "6.2583"

Metric mean fowlkes\_mallows\_index:  
[1] "6.3134"

Metric mean connectivity:  
[1] "22.08"

Metric mean dunn:  
[1] "6.0572"

Metric mean silhouette:  
[1] "6.184"

Mean time for evaluation of internal metrics:  
[1] "6.0822"

> |

Files Plots Packages Help Viewer

Precision

CLUSTERING

Algorithm

- agnes
- apclusterK
- clara
- daisy
- diana
- fanny
- fuzzy\_cm
- fuzzy\_gg
- fuzzy\_gk
- gama
- gmm
- hclust
- kmeans\_arma
- kmeans\_rcpp
- mini\_kmeans
- mona
- pam
- pvclust



## Clustering en R

paquete CLUSTERING

<https://cran.r-project.org/web/packages/Clustering/Clustering.pdf>

The screenshot displays the Clustering Shiny application interface. On the left is a control panel with the following sections:

- Do you want to use test data or a file directory?**
  - File Directory
  - Test Data
- Directory Datasets**
  - Path: `/home/agvico/R/x86_64-pc-linux-gnu-library/4.0/Clustering/shiny/`
- Dataset Test**
  - Selected: `Baseball`
- Packages**
  - 7 items selected
- Algorithms**
  - 18 items selected
- Number of Clustering**
  - Slider range: 1 to 10, with markers at 3 and 6.
- Metrics**
  - 9 items selected
- Do you want to show the variable?**
  - Yes
  - No

On the right, the results table is shown under the 'Summary' tab. It contains two tables of clustering results:

Algorithm	Distance	Clusters	Dataset	timeExternal	entropy	variation information	precision	recall	f measure	towlkes mallows index
1	apclusterK_euclidean	3	dataframe	4	4	5	2	2	2	2
2	apclusterK_euclidean	4	dataframe	5	2	3	2	2	2	2
3	apclusterK_euclidean	5	dataframe	4	2	4	2	2	2	2
4	apclusterK_euclidean	6	dataframe	3	2	3	2	4	2	2
5	apclusterK_manhattan	3	dataframe	3	4	5	2	2	2	2
6	apclusterK_manhattan	4	dataframe	4	4	3	2	2	2	2
7	apclusterK_manhattan	5	dataframe	3	2	3	2	4	2	2
8	apclusterK_manhattan	6	dataframe	3	2	3	2	4	2	2
9	apclusterK_minkowski	3	dataframe	4	4	5	2	2	2	2
10	apclusterK_minkowski	4	dataframe	5	2	3	2	2	2	2

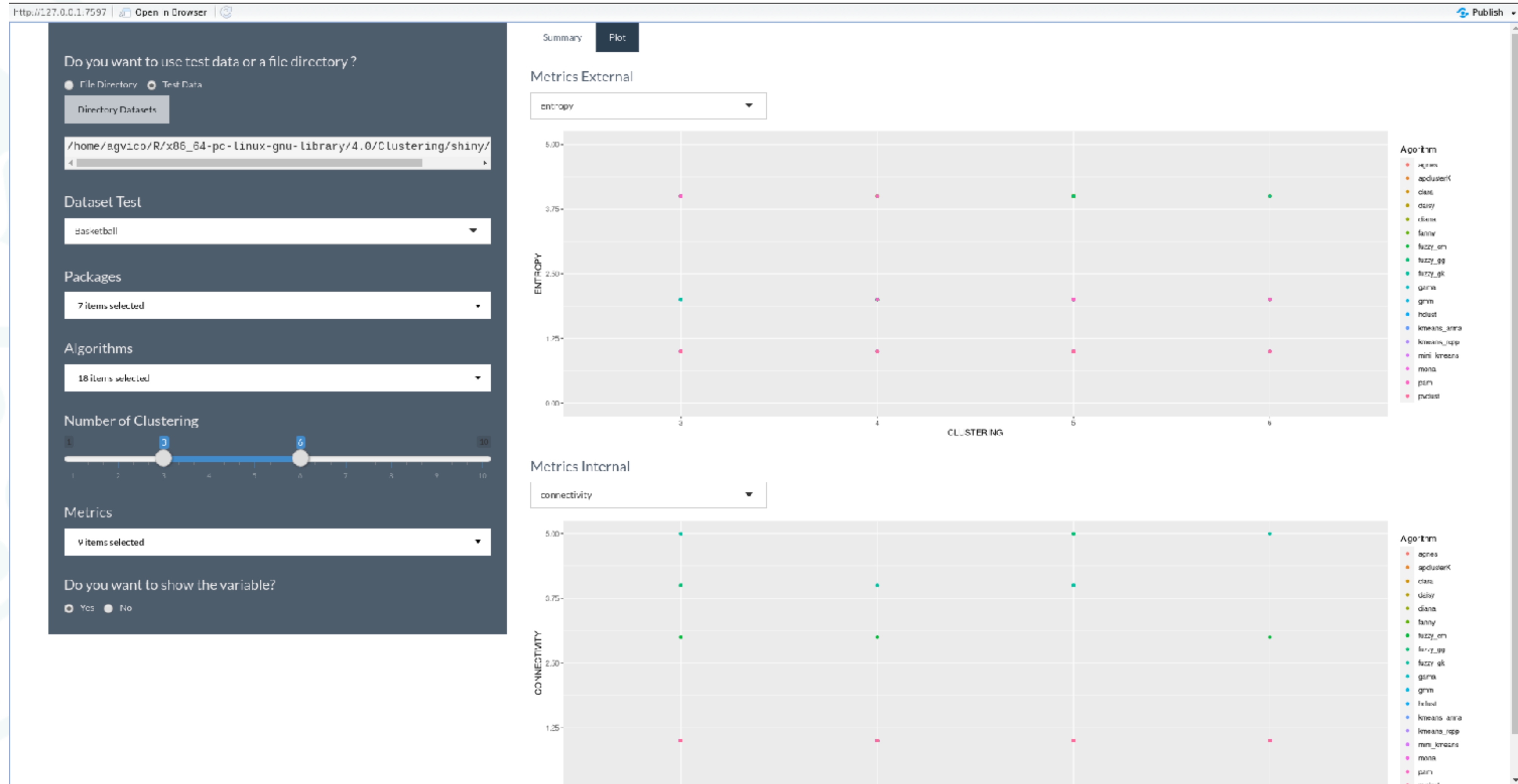
  

Algorithm	Distance	timeExternal	entropy	variation information	precision	recall	f measure	towlkes mallows index
1	agnes_euclidean	1	1	1	1	1	1	1
2	agnes_manhattan	1	1	1	1	1	1	1
3	apclusterK_euclidean	5	4	5	2	4	2	2
4	apclusterK_manhattan	4	4	5	2	4	2	2
5	apclusterK_minkowski	5	4	5	2	4	2	2
6	clara_euclidean	5	4	5	2	4	2	2
7	clara_manhattan	5	4	5	2	3	2	2
8	daisy_euclidean	5	4	1	2	3	2	2
9	daisy_power	5	4	1	2	5	2	2
10	daisy_manhattan	5	4	1	2	5	2	2

## Clustering en R

paquete CLUSTERING

<https://cran.r-project.org/web/packages/Clustering/Clustering.pdf>



# Clustering

referencias bibliográficas

- ▶ Pang-Ning Tan, Michael Steinbach & Vipin Kumar. Introduction to Data Mining, 2006.
- ▶ Jiawei Han. Data mining: concepts and techniques, 2006.
- ▶ Charu C. Aggrawal & Chandan K. Reddy (editors): Data Clustering: Algorithms and Applications. Champan & Hall / CRC Press, 2014.
- ▶ Lloyd, S. P. (1957). Least squares quantization in PCM. Technical Report RR-5497, Bell Lab, September 1957.
- ▶ MacQueen, J. B. (1967). Some methods for classification and analysis of multivariate observations. In L. M. Le Cam & J. Neyman (Eds.), Proceedings of the fifth Berkeley symposium on mathematical statistics and probability (Vol. 1, pp. 281–297). California: University of California Press.



**Seminario Permanente de Formación en  
Inteligencia Artificial aplicada a  
Defensa  
SIADEF**

**Sesión 7: Aprendizaje no  
supervisado**

Cristóbal J. Carmona <[ccarmona@ujaen.es](mailto:ccarmona@ujaen.es)>  
Pedro González <[pglez@ujaen.es](mailto:pglez@ujaen.es)>

