

# MAT-269: Sesión 22, Análisis de Conglomerados I

Felipe Osorio

[fosorios.mat.utfsm.cl](mailto:fosorios.mat.utfsm.cl)

Departamento de Matemática, UTFSM



## Objetivo:

El análisis de conglomerados intenta descubrir grupos (o cluster) de observaciones que son homogéneas dentro de cada grupo.

## Problema:

Dividir el análisis en dos pasos fundamentales.

- ▶ Elección de la medida de proximidad (similaridad).
- ▶ Selección del algoritmo de construcción de grupos.

Nos concentraremos en tres tipos de procedimiento de agrupamiento:

- ▶ Métodos jerárquicos aglomerativos.
- ▶ Métodos tipo  $K$ -means.
- ▶ Métodos de clasificación ML.



Estas técnicas operan sobre una matriz  $D = (d_{ij}) \in \mathbb{R}^{n \times n}$  de distancias<sup>1</sup> entre los puntos de  $\mathbf{X} \in \mathbb{R}^{n \times p}$ ,

$$D = \begin{pmatrix} d_{11} & d_{12} & \dots & d_{1n} \\ d_{21} & d_{22} & \dots & d_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{n1} & d_{n2} & \dots & d_{nn} \end{pmatrix}.$$

Por ejemplo, podríamos usar la distancia Euclidiana,

$$d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2 = \left\{ \sum_{k=1}^p (x_{ik} - x_{jk})^2 \right\}^{1/2}, \quad i, j = 1, \dots, n.$$

Note que, si  $d_{ij}$  es una distancia, entonces  $d'_{ij} = \max_{ij} \{d_{ij}\} - d_{ij}$  es una medida de proximidad.

---

<sup>1</sup> $D$  es construída usando medidas de similaridad o de disimilaridad



## Tipo de distancias:

- ▶ Norma Euclidiana con un métrica  $\mathbf{A} > \mathbf{0}$ ,

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_A = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j)},$$

es usual tomar  $\mathbf{A} = \mathbf{S}^{-1}$  o bien  $\mathbf{A} = \text{diag}(s_{11}^{-1}, \dots, s_{pp}^{-1})$ .

- ▶ Métrica de Minkowski

$$d(\mathbf{x}_i, \mathbf{x}_j) = \left\{ \sum_{k=1}^p |x_{ik} - x_{jk}|^m \right\}^{1/m}.$$

- ▶ Métrica Canberra

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^p \frac{|x_{ik} - x_{jk}|}{(x_{ik} + x_{jk})}.$$

- ▶ Coeficiente de Czekanowski

$$d(\mathbf{x}_i, \mathbf{x}_j) = 1 - 2 \frac{\sum_{k=1}^p \min(x_{ik}, x_{jk})}{\sum_{k=1}^p (x_{ik} + x_{jk})}.$$





Suponga dos objetos o grupos  $P$  y  $Q$ , y sea

$$n_P = \sum_{i=1}^n I(\mathbf{x}_i \in P),$$

el número de objetos en  $P$ , y análogamente para  $n_Q$ . Considere los siguientes procedimientos para agrupar las observaciones:

▶ single linkage:

$$d(P, Q) = \min_{i \in P, i \in Q} \{d_{ij}\},$$

▶ complete linkage:

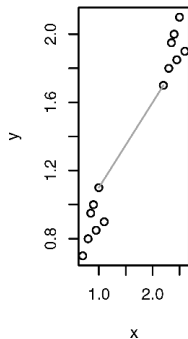
$$d(P, Q) = \max_{i \in P, i \in Q} \{d_{ij}\},$$

▶ average linkage:

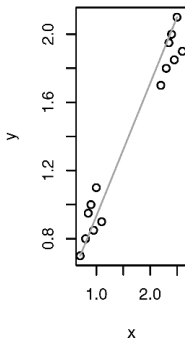
$$d(P, Q) = \frac{1}{n_P n_Q} \sum_{i \in P} \sum_{i \in Q} d_{ij}.$$



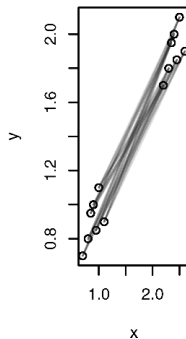
single



complete



average



## Análisis de Conglomerados

Suponga dos objetos o grupos  $P$  y  $Q$  que están unidos, y deseamos calcular la distancia entre este nuevo grupo  $P + Q$  con un grupo  $R$ , digamos:

$$d(R, P + Q) = \delta_1 d(R, P) + \delta_2 d(R, Q) + \delta_3 d(P, Q) \\ + \delta_4 |d(R, P) - d(R, Q)|,$$

donde diferentes elecciones de las ponderaciones  $\delta_i$ 's da origen a distintos tipos de algoritmos aglomerativos.

Sea

$$n_P = \sum_{i=1}^n I(\mathbf{x}_i \in P),$$

el número de objetos en  $P$ , y análogamente para  $n_Q$  y  $n_R$ . Por ejemplo,

Linkage	$\delta_1$	$\delta_2$	$\delta_3$	$\delta_4$
single	1/2	1/2	0	-1/2
complete	1/2	1/2	0	1/2
average	1/2	1/2	0	0
median	1/2	1/2	-1/4	0
centroid	$\frac{n_P}{n_P+n_Q}$	$\frac{n_Q}{n_P+n_Q}$	$-\frac{n_P n_Q}{(n_P+n_Q)^2}$	0





---

## Algoritmo 1: Método Jerárquico Aglomerativo.

---

Entrada: Matriz de datos  $X = (\mathbf{x}_1^T, \dots, \mathbf{x}_n^T)^T$ .

```
1 begin
2   Construir la partición más fina.
3   Calcular la matriz de distancias  $D$ .
4   do
5     Hallar dos grupos con la distancia más cercana.
6     Agrupar dos grupos en un único grupo.
7     Calcular la distancia entre los nuevos grupos y obtener una matriz
       reducida  $D$ .
8   until todos los grupos están aglomerados en  $X$ 
9 end
```

---



## Ejemplo:

Considere

$$\mathbf{x}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{x}_3 = \begin{pmatrix} 5 \\ 5 \end{pmatrix}$$

El algoritmo inicia con  $K = 3$  grupos,  $P = \{\mathbf{x}_1\}$ ,  $Q = \{\mathbf{x}_2\}$ ,  $R = \{\mathbf{x}_3\}$ . La matriz de distancias  $\mathbf{D}$  es dada por

$$\mathbf{D} = \begin{pmatrix} 0 & 1 & 50 \\ 1 & 0 & 41 \\ 50 & 41 & 0 \end{pmatrix}.$$

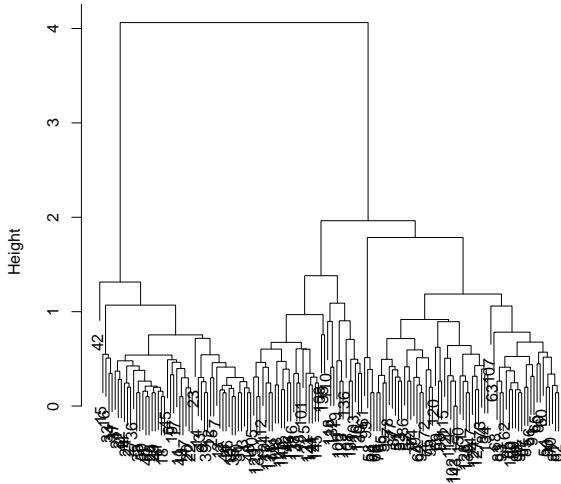
La menor distancia en  $\mathbf{D}$  se encuentra entre los grupos  $P$  y  $Q$ . De esta forma estos grupos se deben combinar en  $P + Q = \{\mathbf{x}_1, \mathbf{x}_2\}$ . Usando single linkage, obtenemos

$$\begin{aligned} d(R, P + Q) &= \frac{1}{2}d(R, P) + \frac{1}{2}d(R, Q) - \frac{1}{2}|d(R, P) - d(R, Q)| \\ &= \frac{1}{2}d_{13} + \frac{1}{2}d_{23} - \frac{1}{2}|d_{13} - d_{23}| = \frac{50}{2} + \frac{41}{2} - \frac{|50-41|}{2} = 41. \end{aligned}$$

y la matriz de distancias *reducida* adopta la forma  $\mathbf{D}_* = \begin{pmatrix} 0 & 41 \\ 41 & 0 \end{pmatrix}$ . Detenemos el algoritmo uniendo los grupos  $R$  y  $P + Q$  para formar el cluster  $\mathbf{X}$ , la matriz de datos original.



# Análisis de Conglomerados



**K-means** busca particionar los  $n$  individuos en  $K$  grupos, digamos  $G_1, G_2, \dots, G_K$ . El tipo más común de algoritmo halla una partición que minimice la **suma de cuadrados dentro-de-grupo**,

$$WGSS = \sum_{j=1}^q \sum_{r=1}^K \sum_{i \in G_r} (x_{ij} - \bar{x}_j^{(r)})^2,$$

donde  $\bar{x}_j^{(r)} = \frac{1}{n_i} \sum_{i \in G_r} x_{ij}$ .

$n$	$k$	Num. de particiones posibles
15	3	2 375 101
20	4	45 232 115 901
25	8	690 223 721 118 368 580
100	5	$10^{68}$



---

## Algoritmo 2: Método $K$ -medias.

---

**Entrada:** Matriz de datos  $\mathbf{X} = (\mathbf{x}_1^T, \dots, \mathbf{x}_n^T)^T$ .

```
1 begin
2   Hallar una partición inicial de los individuos en los  $K$  grupos.
3   do
4     Proceder a través de la lista de elementos y asignar una observación al
5     grupo cuyo centroide (media) sea más cercano.
6     Recalcular centroides.
7   until no se pueda hacer más asignaciones.
8 end
```

---

### Observación:

El método de  $K$ -medias sufre principalmente de dos problemas:

- ▶ No es invariante a transformaciones de escala.
- ▶ Impone una estructura “esférica” a los datos.



El procedimiento de agrupamiento por ML es basado en asumir  $G$  subpoblaciones

$$f_j(\mathbf{x}; \boldsymbol{\theta}_j), \quad \boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_G^\top)^\top.$$

Además, se introduce un vector  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_n)^\top$  donde  $\gamma_i = k$  si  $\mathbf{x}_i$  pertenece a la  $k$ -ésima población.

De este modo, el problema de agrupamiento resulta de escoger  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_G^\top)^\top$  y  $\boldsymbol{\gamma}$  maximizando la verosimilitud:

$$L(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \prod_{i=1}^n f_{\gamma_i}(\mathbf{x}_i; \boldsymbol{\theta}_{\gamma_i}).$$



Bajo normalidad tenemos  $\theta_j = (\mu_j, \Sigma_j)$ ,  $j = 1, \dots, G$  y los MLE de  $\mu_j$  son

$$\bar{\mathbf{x}} = \frac{1}{n_j} \sum_{i \in A_j} \mathbf{x}_i,$$

con  $A_j = \{i : \gamma_i = j\}$  y  $n_j$  es el número de elementos de  $A_j$ . En este caso, la función de log-verosimilitud perfilada adopta la forma:

$$\ell(\boldsymbol{\theta}, \boldsymbol{\gamma}) = c - \frac{n}{2} \sum_{j=1}^G \left\{ \text{tr} \mathbf{S}_j \boldsymbol{\Sigma}_j^{-1} + \log |\boldsymbol{\Sigma}_j| \right\}.$$

