

MAT-269: Análisis de Conglomerados

Felipe Osorio

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 distancias¹ 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.

¹ 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

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



Example:

Para los datos de Iris, tenemos:

$$D = \begin{pmatrix} 0.00 & & & & & & & & \\ 0.54 & 0.00 & & & & & & & \\ 0.51 & 0.30 & 0.00 & & & & & & \\ 0.65 & 0.33 & 0.24 & 0.00 & & & & & \\ 0.14 & 0.61 & 0.51 & 0.65 & 0.00 & & & & \\ 0.62 & 1.09 & 1.09 & 1.17 & 0.62 & 0.00 & & & \\ 0.52 & 0.51 & 0.26 & 0.33 & 0.46 & 0.99 & 0.00 & & \\ 0.17 & 0.42 & 0.41 & 0.50 & 0.22 & 0.70 & 0.42 & 0.00 & \\ 0.92 & 0.51 & 0.44 & 0.30 & 0.92 & 1.46 & 0.55 & 0.79 & 0.00 \\ \vdots & \vdots & & & & & & & \ddots \end{pmatrix}$$

En este caso, tenemos que D es una matriz simétrica 150×150 .



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:

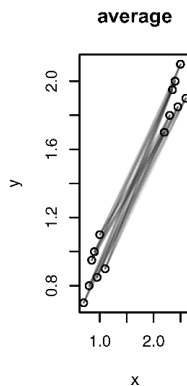
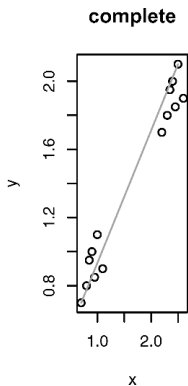
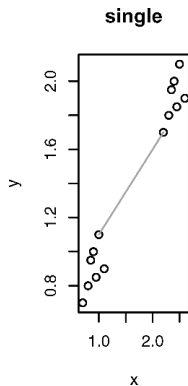
$$d(P, Q) = \max_{i \in P, i \in Q} \{d_{ij}\},$$

► average linkage:

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



Análisis de Conglomerados



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 δ_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	δ_1	δ_2	δ_3	δ_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 $\mathbf{X} = (\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top)^\top$.

```
1 begin
2   Construir la partición más fina.
3   Calcular la matriz de distancias  $\mathbf{D}$ .
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  $\mathbf{D}$ .
8   until todos los grupos están aglomerados en  $\mathbf{X}$ 
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

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

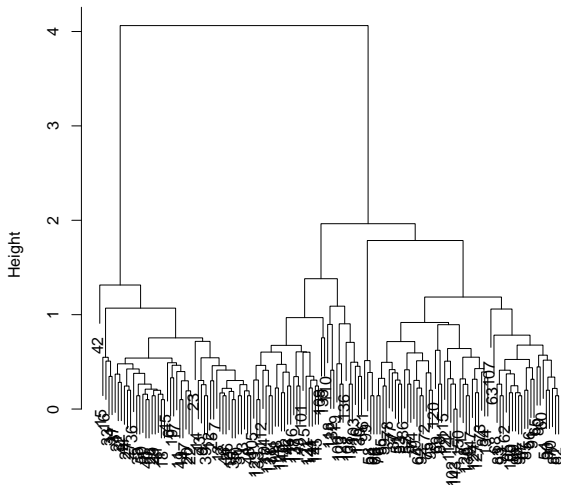
La menor distancia en \mathbf{D} se encuentra entre los grupos P y Q . De esta forma estos grupos se deben combinar en $P + Q = \{\mathbf{x}_1, \mathbf{x}_2\}$. Usando single linkage, obtenemos

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

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 $\mathbf{X} = (\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top)^\top$.

```
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
        grupo cuyo centroide (media) sea más cercano.
5     Recalcular centroides.
6   until no se pueda hacer más asignaciones.
7 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 μ_j son

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

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:

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



Una manera de caracterizar este tipo de funciones de densidad es asumiendo una **mezcla discreta de densidades**, como:

$$f(\mathbf{x}; \boldsymbol{\pi}, \boldsymbol{\theta}) = \sum_{j=1}^G \pi_j f_j(\mathbf{x}; \boldsymbol{\theta}_j),$$

donde \mathbf{x} es un vector aleatorio p -dimensional, $\boldsymbol{\pi} = (\pi_1, \dots, \pi_G)^\top$, y $\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_G^\top)^\top$, con π_j siendo las **proporciones de la mezcla** y f_j las **densidades** que componen la mezcla. Además,

$$\sum_{j=1}^G \pi_j = 1.$$

Una vez estimados los parámetros de la mezcla, las observaciones pueden ser asociadas con un particular cluster en base de la **probabilidad posterior estimada**,

$$\hat{P}(\text{cluster } j | \mathbf{x}_i) = \frac{\hat{\pi}_j f_j(\mathbf{x}_i; \hat{\boldsymbol{\theta}}_j)}{f(\mathbf{x}_i; \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}})}, \quad j = 1, \dots, G. \quad (1)$$



Dada una muestra de observaciones $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ desde una mezcla discreta de densidades, tenemos la función de log-verosimilitud

$$\ell_n(\boldsymbol{\pi}, \boldsymbol{\theta}) = \sum_{i=1}^n \log f(\mathbf{x}_i; \boldsymbol{\pi}, \boldsymbol{\theta}).$$

En el caso de que el j -ésimo componente siga una distribución normal multivariada con media $\boldsymbol{\mu}_j$ y covarianza $\boldsymbol{\Sigma}_j$, se puede mostrar que

$$\begin{aligned}\hat{\pi}_j &= \frac{1}{n} \sum_{i=1}^n \hat{P}(j|\mathbf{x}_i), \\ \hat{\boldsymbol{\mu}}_j &= \frac{1}{n\hat{\pi}_j} \sum_{i=1}^n \hat{P}(j|\mathbf{x}_i) \mathbf{x}_i, \\ \hat{\boldsymbol{\Sigma}}_j &= \frac{1}{n} \sum_{i=1}^n \hat{P}(j|\mathbf{x}_i) (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_j)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_j)^\top,\end{aligned}$$

donde $\hat{P}(j|\mathbf{x}_i)$ son las probabilidades estimadas definidas en (1). Este procedimiento es un caso particular del [algoritmo EM](#) para estimación ML en mezclas discretas.

