



ALGORITMOS DE RECOMENDACIÓN

Clasificación Jerárquica

amazon

The Amazon logo, featuring the word "amazon" in a bold, black, sans-serif font. Below the text is a curved orange arrow that starts under the 'a' and points towards the 'n'.The Netflix logo, consisting of the word "NETFLIX" in a bold, red, sans-serif font, centered on a black rectangular background.

amazon

NETFLIX



PREMIO NETFLIX

- Competencia abierta que premió al mejor algoritmo que permitiera predecir con anticipación la evaluación de un usuario a una película nueva utilizando únicamente las evaluaciones anteriores y ninguna información adicional respecto al usuario o la película.

US\$ 1.000.000



amazon



29%

NETFLIX



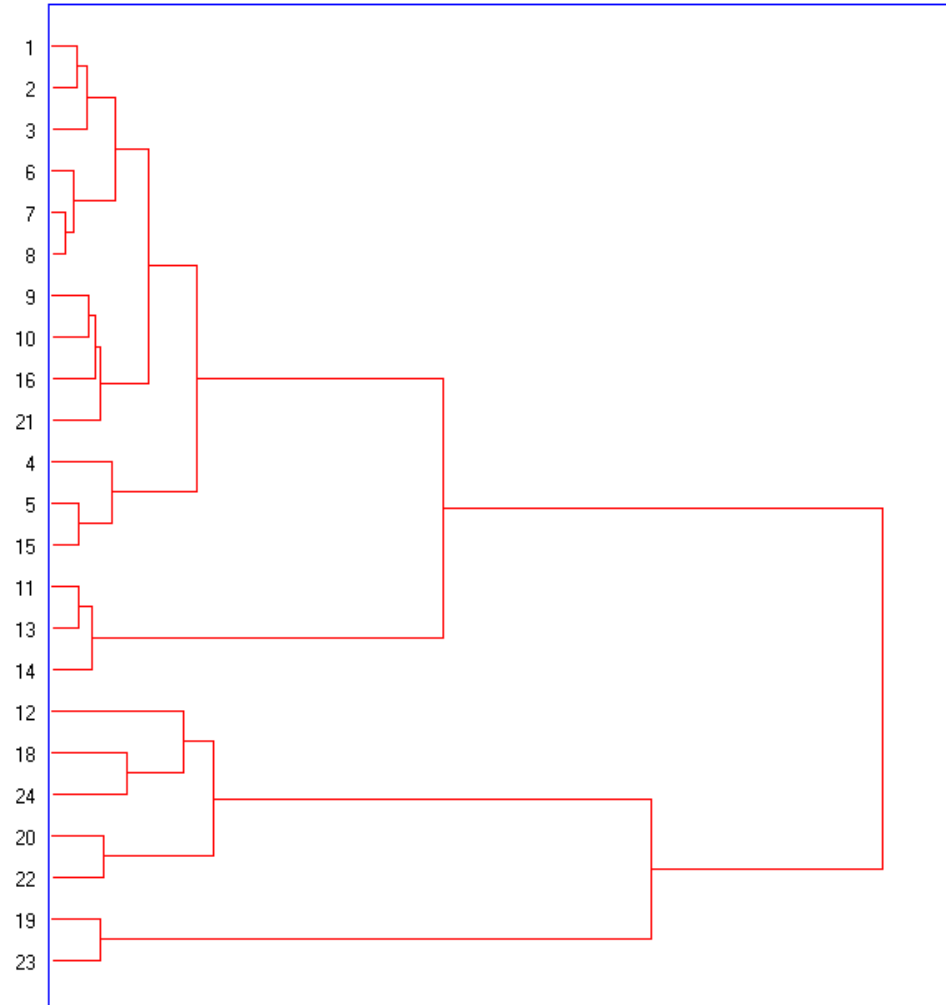
10%



...



CLASIFICACIÓN JERÁRQUICA



CLASIFICACIÓN AUTOMÁTICA

- “La clasificación automática tiene por objetivo reconocer grupos de individuos homogéneos, de tal forma que los grupos queden bien separados y bien diferenciados.”
- “Estos individuos pueden estar descritos por una tabla de datos de individuos por variables, con variables cuantitativas o cualitativas, o por una tabla de proximidades.”

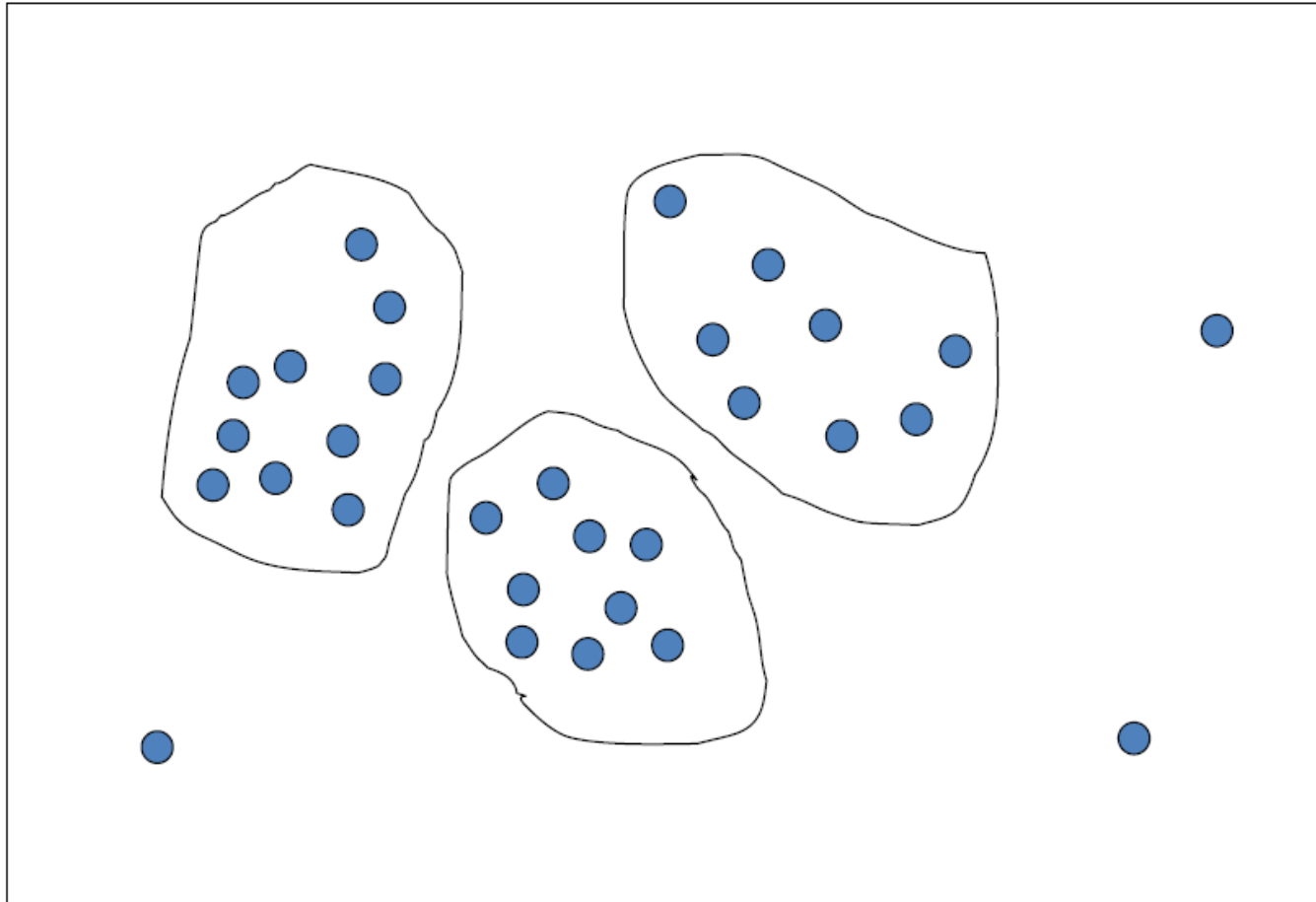


TAREA DE LA MINERÍA DE DATOS

- **“Clustering”**: *Es similar a la clasificación (discriminación), excepto que los grupos no son predefinidos. El objetivo es particionar o segmentar un conjunto de datos o individuos en grupos que pueden ser disjuntos o no. Los grupos se forman basados en la similaridad de los datos o individuos en ciertas variables. Como los grupos no son dados a priori el experto debe dar una interpretación de los grupos que se forman.*
- Métodos:
 1. Clasificación Jerárquica (grupos disjuntos).
 2. Nubes Dinámicas o k-means (grupos disjuntos).
 3. Clasificación Piramidal (grupos NO disjuntos).



ANÁLISIS DE CLÚSTER



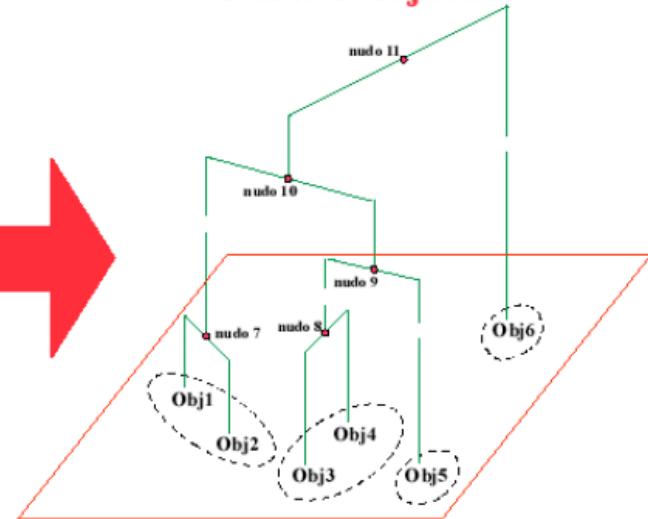
Clasificación Jerárquica

Tabla T(n,p)

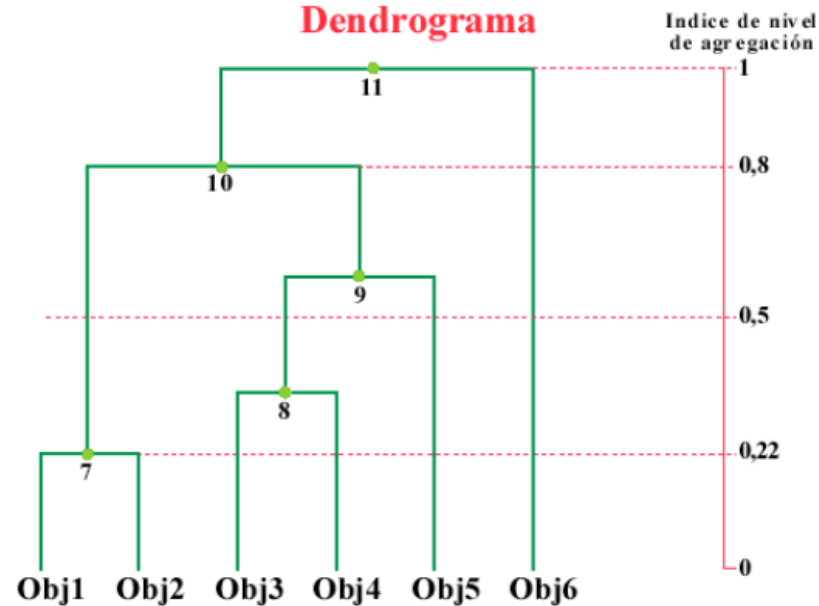
	Var.1	...	Var. j	...	Var. p
Obj.1					
Obj.2					
Obj.3			x_{ij}		
Obj.4					
Obj.5					
Obj.6					



Clases encajadas



Dendrograma



DEFINICIONES BÁSICAS

Sea X la matriz de datos cuyas n filas y p columnas, forman el conjunto del cual se busca una partición. Supondremos que X es una matriz de n individuos cada uno representado por p variables.

Tabla de Datos

$$\begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}$$



EJEMPLO: TABLA DE NOTAS ESCOLARES

	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7.0	6.5	9.2	8.6	8.0
Pedro	7.5	9.4	7.3	7.0	7.0
Inés	7.6	9.2	8.0	8.0	7.5
Luis	5.0	6.5	6.5	7.0	9.0
Andrés	6.0	6.0	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8.0	6.5
Carlos	6.3	6.4	8.2	9.0	7.2
José	7.9	9.7	7.5	8.0	6.0
Sonia	6.0	6.0	6.5	5.5	8.7
María	6.8	7.2	8.7	9.0	7.0



DISIMILITUD Y AGREGACIONES

Con el propósito de encontrar una clasificación de las filas o columnas de la matriz X , el primer problema a resolver es cómo cuantificar la similitud entre esos objetos o grupos de objetos.



ALGUNOS ÍNDICES DE DISIMILITUD

Un índice de disimilitud entre objetos pertenecientes a un conjunto I , es una función d tal que:

$$d : I \times I \longrightarrow [0, +\infty[$$

y

$$d(x, y) = d(y, x) \text{ para todo } x, y \in I.$$



EJEMPLOS DE ÍNDICES

- Distancia Euclidiana: Sean $\mathbf{p} = (p_1, p_2, \dots, p_n)$ y $\mathbf{q} = (q_1, q_2, \dots, q_n)$ dos columnas de la tabla. La distancia Euclidiana entre ellas está definida por:

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



Tabla de Datos

	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7	6.5	9.2	8.6	8
Pedro	7.5	9.4	7.3	7	7
Inés	7.6	9.2	8	8	7.5
Luis	5	6.5	6.5	7	9
Andrés	6	6	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8	6.5
Carlos	6.3	6.4	8.2	9	7.2
José	7.9	9.7	7.5	8	6
Sonía	6	6	6.5	5.5	8.7
María	6.8	7.2	8.7	9	7

Distancia Lucía-Pedro

0.25	8.41	3.61	2.56	1	3.9787
------	------	------	------	---	--------



DISIMILITUD NOTAS ESCOLARES

Matriz de Distancias										
	Lucía	Pedro	Inés	Luis	Andrés	Ana	Carlos	José	Sonía	María
Lucía	0	3,98	3,11	3,85	1,947	3,89	1,517	4,28	4,32	1,39
Pedro		0	1,34	4,39	4,214	1,24	3,91	1,51	4,43	3,36
Inés			0	4,42	3,7	1,135782	3,265	1,69	4,77	2,53
Luis				0	3,072	1,89	3,439	5,45	1,89	4,07
Andrés					0	4,2	0,656	4,46	3,9	1,73
Ana						0	3,772	0,56	5,36	3
Carlos							0	4,05	4,2	1,09
José								0	5,64	3,3
Sonía									0	4,7
María										0



EJEMPLOS DE ÍNDICES

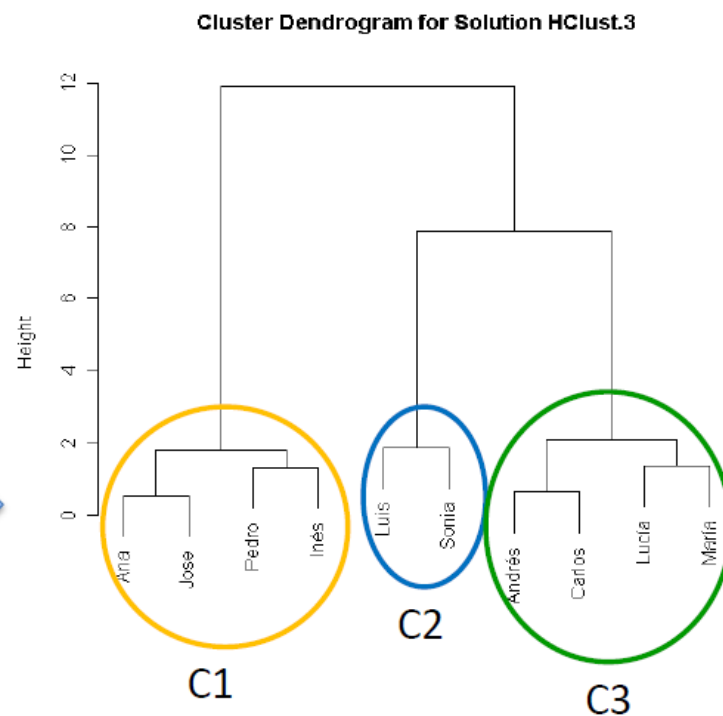
- Distancia Euclídea de las varianzas: cuando las variables tienen varianzas muy desiguales, la distancia entre filas puede depender más de la estructura de varianzas que de la estructura de correlaciones. Para corregir este efecto se utiliza el índice:

$$d(x_i, x_s) = \sqrt{\sum_{j=1}^p \frac{1}{\sigma_j^2} (x_{ij} - x_{sj})^2}$$



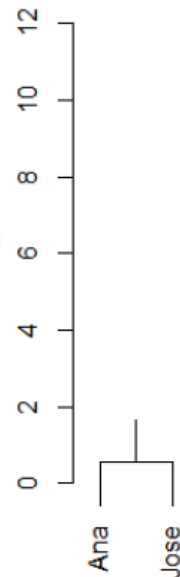
¿CÓMO SE CONSTRUYE EL ÁRBOL?

Análisis de los Clústeres					
	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7	6,5	9,2	8,6	8
Pedro	7,5	9,4	7,3	7	7
Inés	7,6	9,2	8	8	7,5
Luis	5	6,5	6,5	7	9
Andrés	6	6	7,8	8,9	7,3
Ana	7,8	9,6	7,7	8	6,5
Carlos	6,3	6,4	8,2	9	7,2
José	7,9	9,7	7,5	8	6
Sonía	6	6	6,5	5,5	8,7
María	6,8	7,2	8,7	9	7

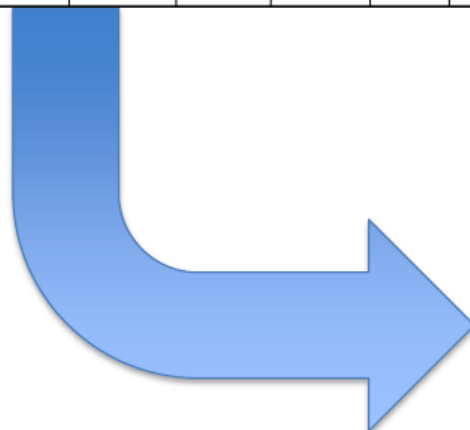


Matriz de Distancias										
	Lucía	Pedro	Inés	Luis	Andrés	Ana	Carlos	José	Sonía	María
Lucía	0	3,98	3,11	3,85	1,95	3,89	1,517	4,28	4,32	1,39
Pedro		0	1,34	4,39	4,21	1,24	3,91	1,51	4,43	3,36
Inés			0	4,42	3,7	1,14	3,265	1,69	4,77	2,53
Luis				0	3,07	1,89	3,439	5,45	1,89	4,07
Andrés					0	4,2	0,656	4,46	3,9	1,73
Ana						0	3,772	0,56	5,36	3
Carlos							0	4,05	4,2	1,09
José								0	5,64	3,3
Sonía									0	4,7
María										0

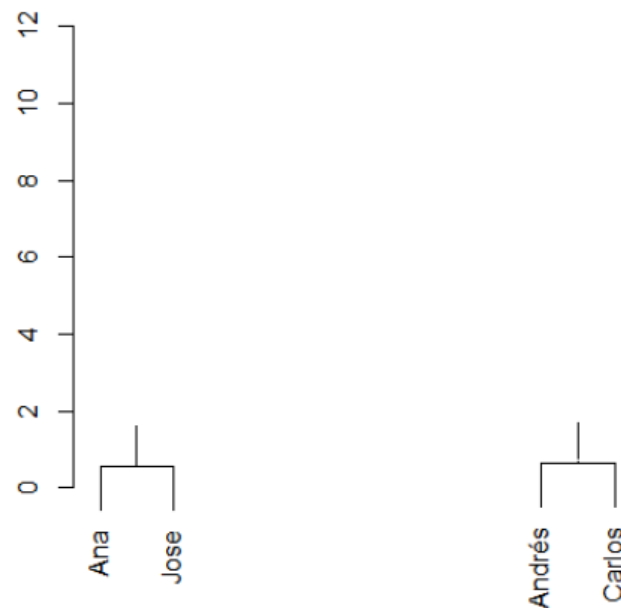
Cluster Dendrogram for Solution HClust.1



Matriz de Distancias										
	Lucía	Pedro	Inés	Luis	Andrés	Ana	Carlos	José	Sonía	María
Lucía	0	3,98	3,11	3,85	1,947	3,89	1,517	4,28	4,32	1,39
Pedro		0	1,34	4,39	4,214	1,24	3,91	1,51	4,43	3,36
Inés			0	4,42	3,7	1,14	3,265	1,69	4,77	2,53
Luis				0	3,072	1,89	3,439	5,45	1,89	4,07
Andrés					0	4,2	0,656	4,46	3,9	1,73
Ana						0	3,772	0,56	5,36	3
Carlos							0	4,05	4,2	1,09
José								0	5,64	3,3
Sonía									0	4,7
María										0



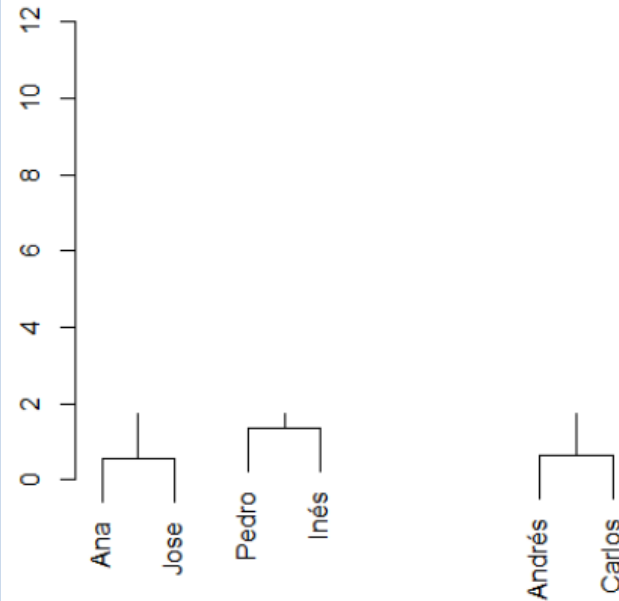
Cluster Dendrogram for Solution HClust.1



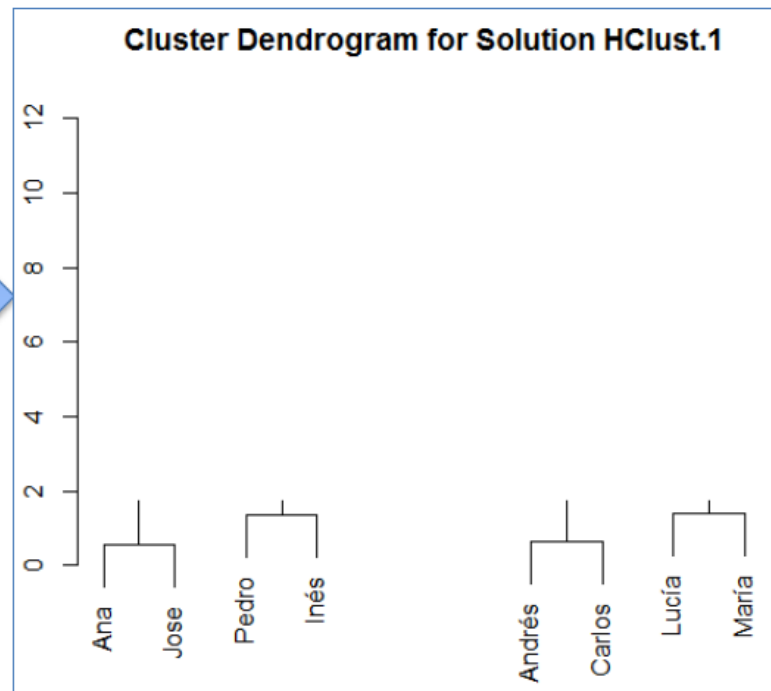
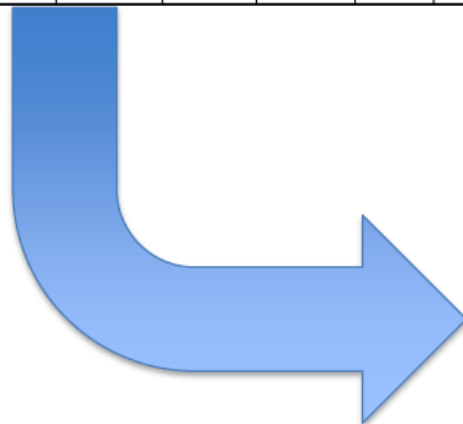
Matriz de Distancias

	Lucía	Pedro	Inés	Luis	Andrés	Ana	Carlos	José	Sonía	María
Lucía	0	3,98	3,11	3,85	1,947	3,89	1,517	4,28	4,32	1,39
Pedro		0	1,34	4,39	4,214	1,24	3,91	1,51	4,43	3,36
Inés			0	4,42	3,7	1,14	3,265	1,69	4,77	2,53
Luis				0	3,072	1,89	3,439	5,45	1,89	4,07
Andrés					0	4,2	0,656	4,46	3,9	1,73
Ana						0	3,772	0,56	5,36	3
Carlos							0	4,05	4,2	1,09
José								0	5,64	3,3
Sonía									0	4,7
María										0

Cluster Dendrogram for Solution HClust.1

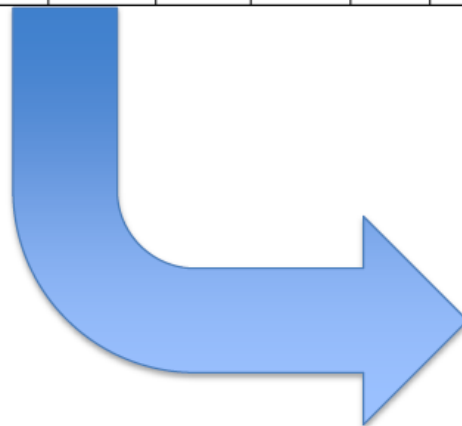


Matriz de Distancias										
	Lucía	Pedro	Inés	Luis	Andrés	Ana	Carlos	José	Sonía	María
Lucía	0	3,98	3,11	3,85	1,947	3,89	1,517	4,28	4,32	1,39
Pedro		0	1,34	4,39	4,214	1,24	3,91	1,51	4,43	3,36
Inés			0	4,42	3,7	1,14	3,265	1,69	4,77	2,53
Luis				0	3,072	1,89	3,439	5,45	1,89	4,07
Andrés					0	4,2	0,656	4,46	3,9	1,73
Ana						0	3,772	0,56	5,36	3
Carlos							0	4,05	4,2	1,09
José								0	5,64	3,3
Sonía									0	4,7
María										0

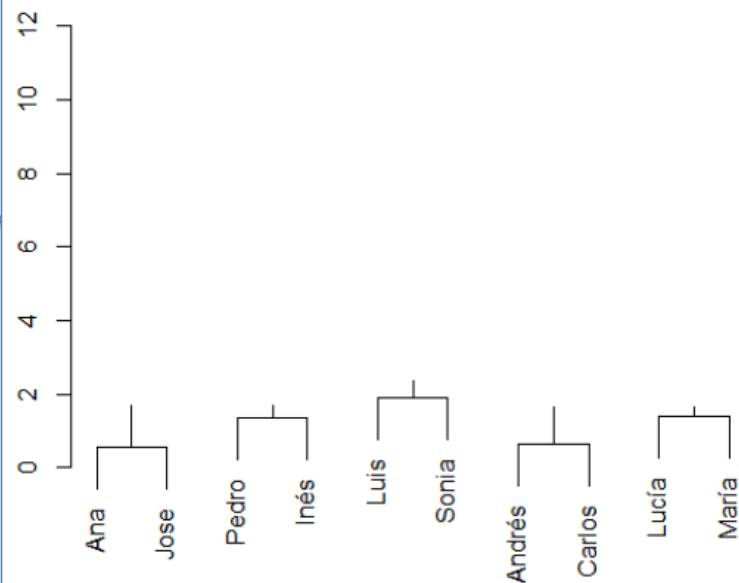


Matriz de Distancias

	Lucía	Pedro	Inés	Luis	Andrés	Ana	Carlos	José	Sonía	María
Lucía	0	3,98	3,11	3,85	1,947	3,89	1,517	4,28	4,32	1,39
Pedro		0	1,34	4,39	4,214	1,24	3,91	1,51	4,43	3,36
Inés			0	4,42	3,7	1,14	3,265	1,69	4,77	2,53
Luis				0	3,072	1,89	3,439	5,45	1,89	4,07
Andrés					0	4,2	0,656	4,46	3,9	1,73
Ana						0	3,772	0,56	5,36	3
Carlos							0	4,05	4,2	1,09
José								0	5,64	3,3
Sonía									0	4,7
María										0



Cluster Dendrogram for Solution HClust.1



ÍNDICES DE AGREGACIÓN

- Permiten cuantificar la similitud entre grupos de objetos del conjunto a clasificar.

Una agregación es una función tal que:

$$\delta : P(I) \times P(I) \longrightarrow [0, +\infty[$$

$$\delta(x, x) = 0 \quad \forall x \in P(I)$$

$$\delta(x, y) = \delta(y, x),$$



EJEMPLOS

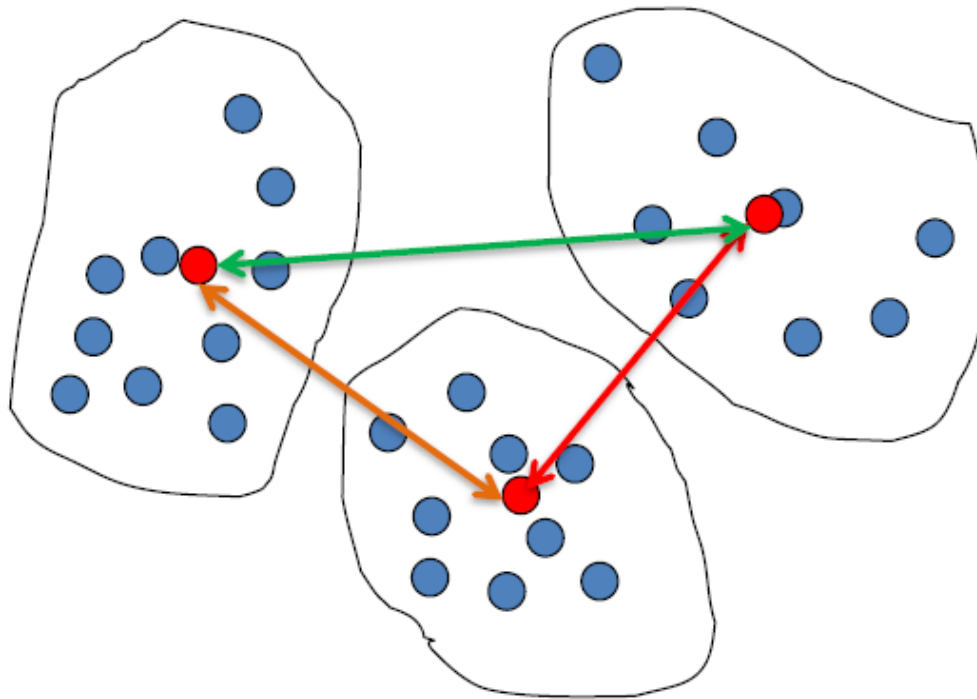
- Agregación de *Ward*:

$$\delta_w(x, y) = \frac{|x| \cdot |y|}{|x| + |y|} \|g_x - g_y\|^2$$

g_x es el baricentro de x .

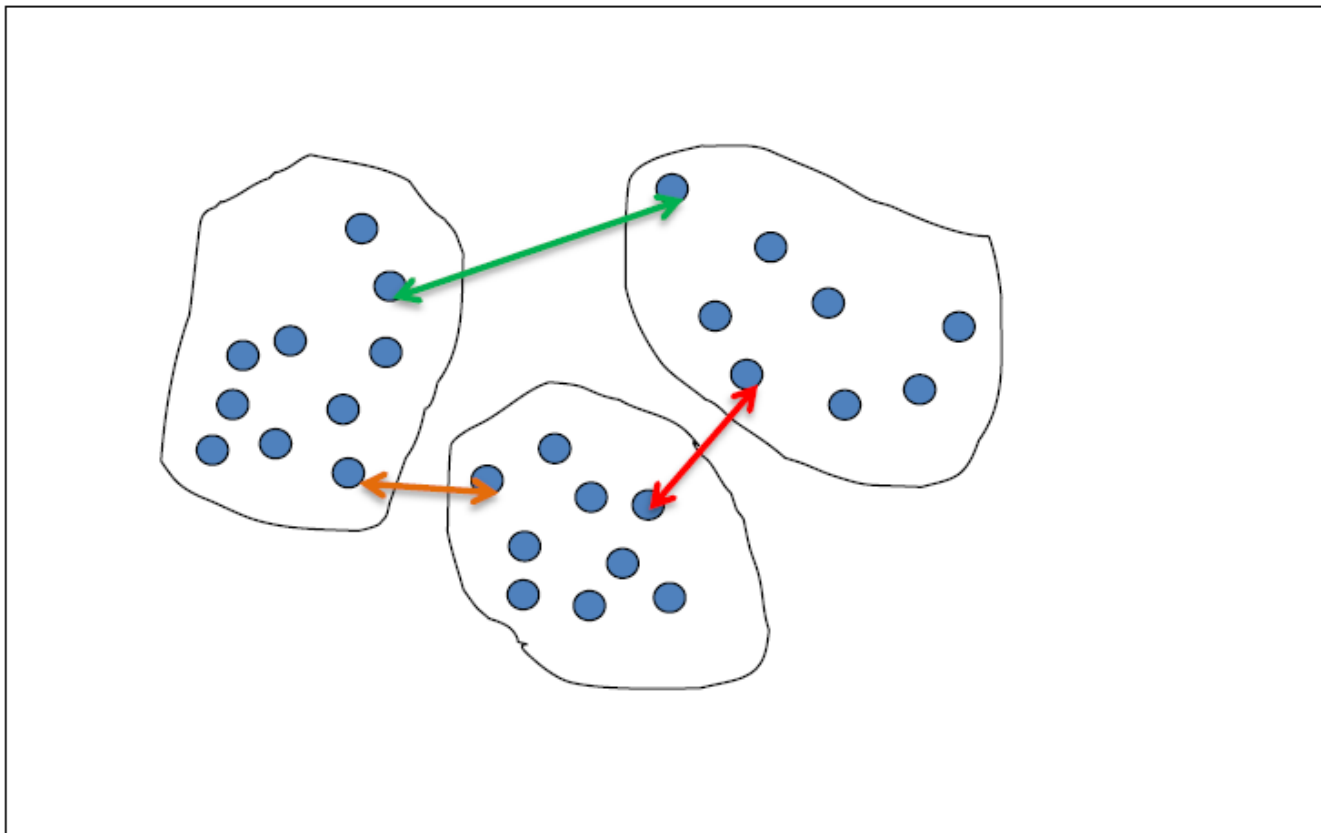


AGREGACIÓN DE WARD



AGREGACIÓN DEL SALTO MÍNIMO

$$\delta_{\min}(x, y) = \min \{d(h, k) \mid h \in x \text{ y } k \in y\}.$$



OTROS EJEMPLOS

- Agregación del salto máximo:

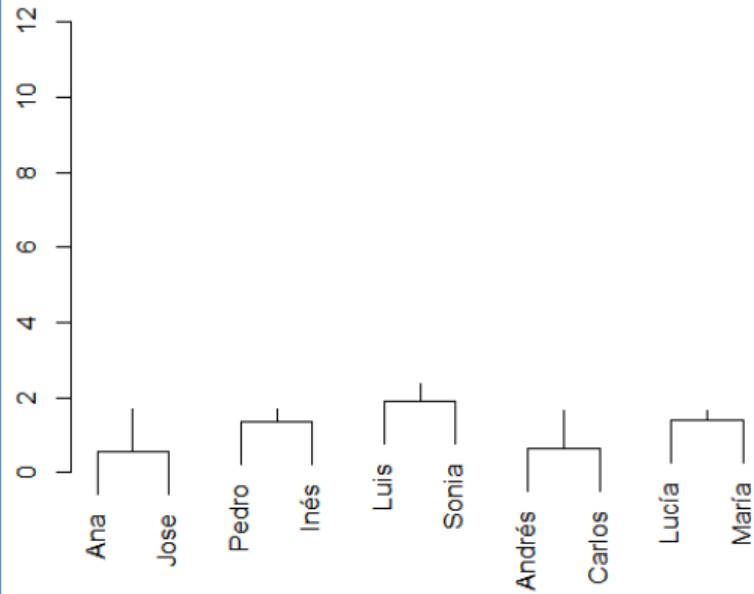
$$\delta_{\max}(x, y) = \max \{d(h, k) \mid h \in x \text{ y } k \in y\}$$

- Agregación del promedio de las disimilitudes:

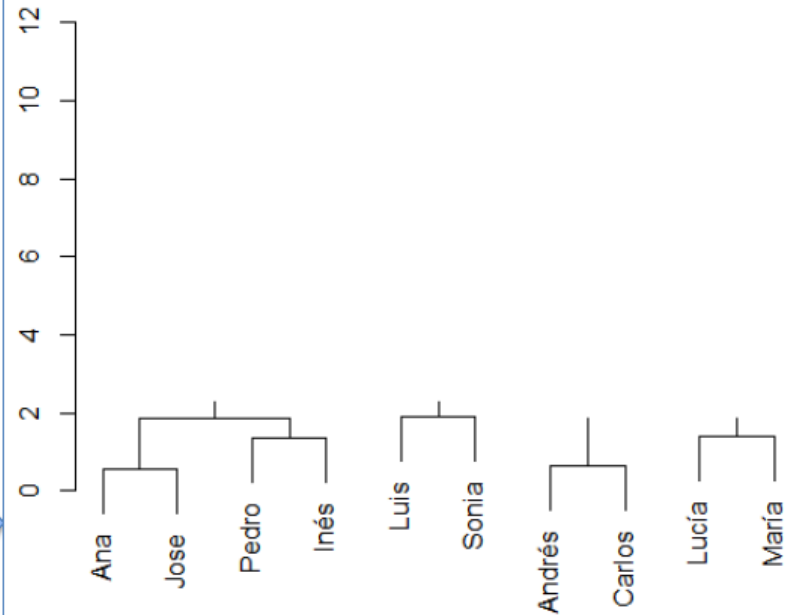
$$\delta_{\text{prom}}(x, y) = \frac{1}{|x| + |y|} \sum \{d(h, k) \mid h \in x \text{ y } k \in y\}$$



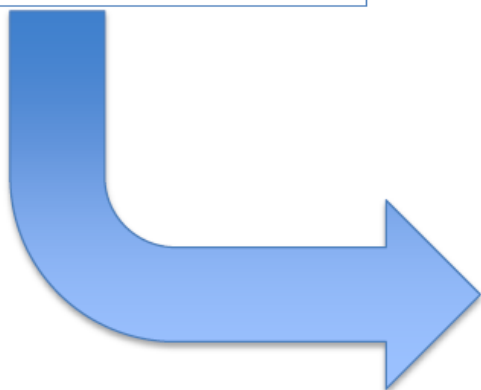
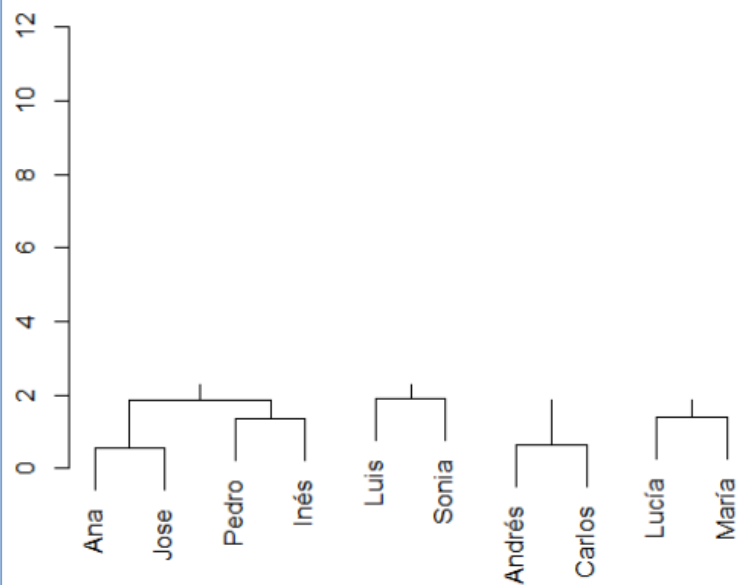
Cluster Dendrogram for Solution HClust.1



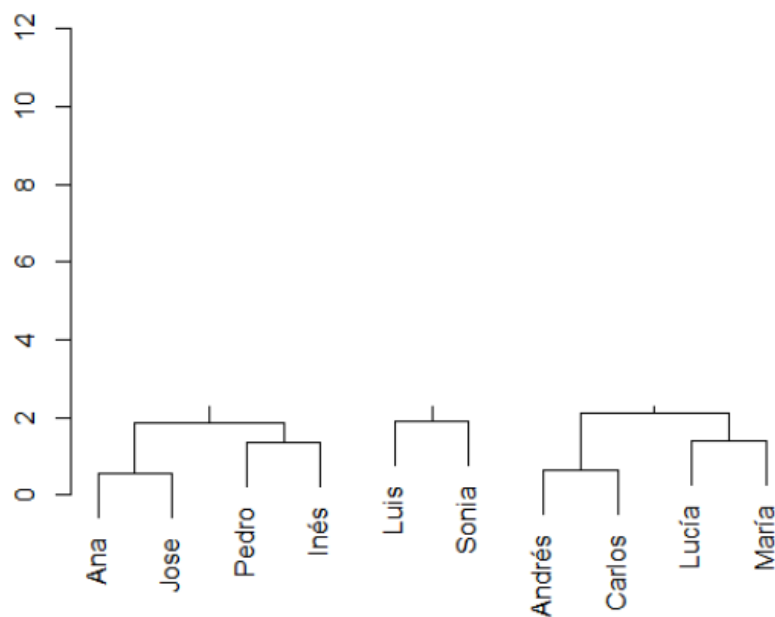
Cluster Dendrogram for Solution HClust.1



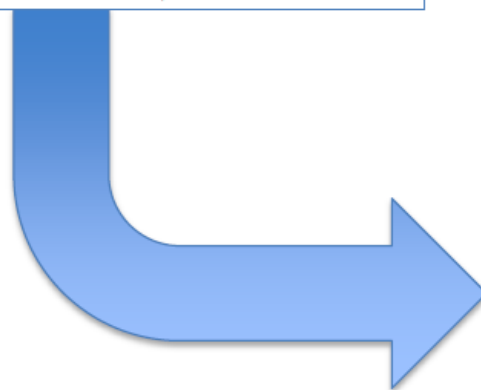
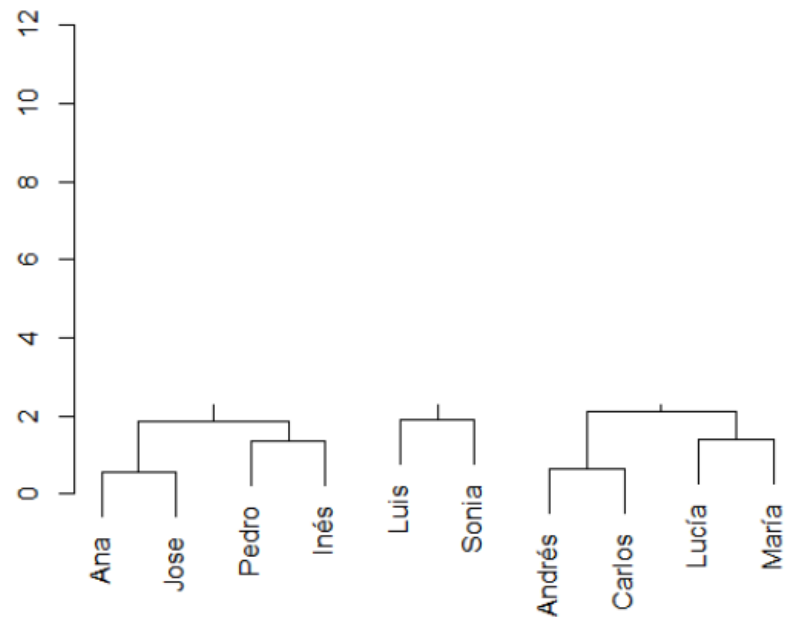
Cluster Dendrogram for Solution HClust.1



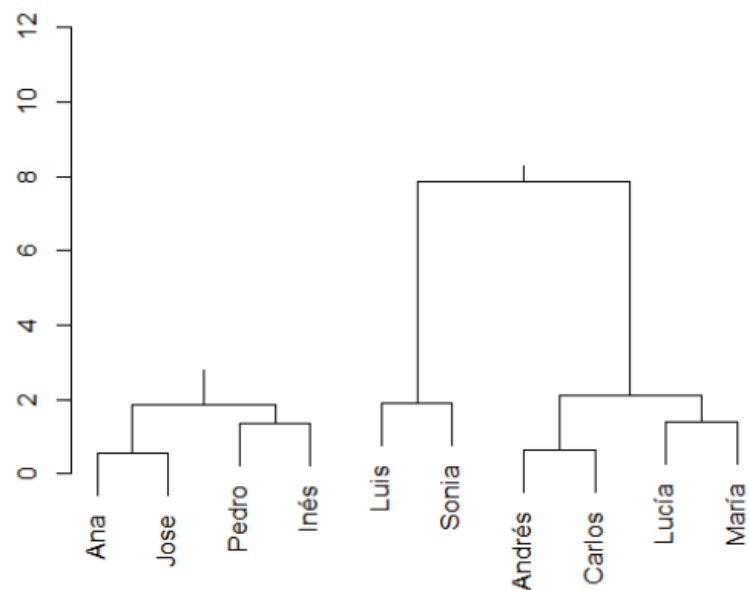
Cluster Dendrogram for Solution HClust.1



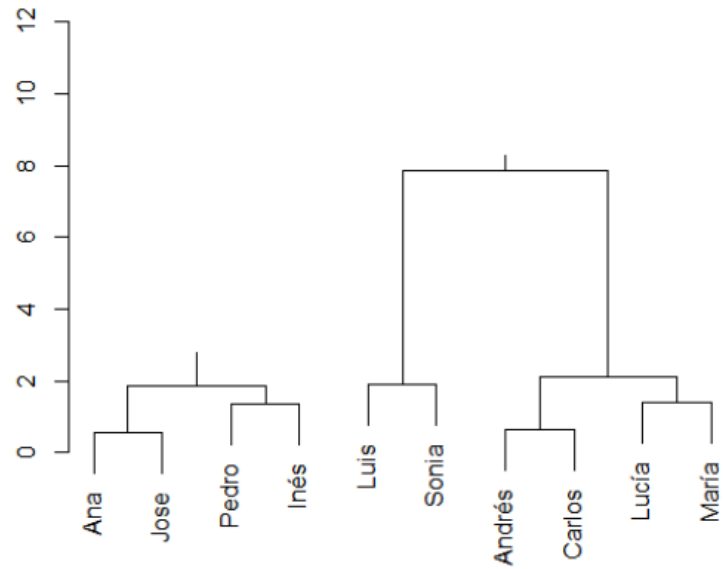
Cluster Dendrogram for Solution HClust.1



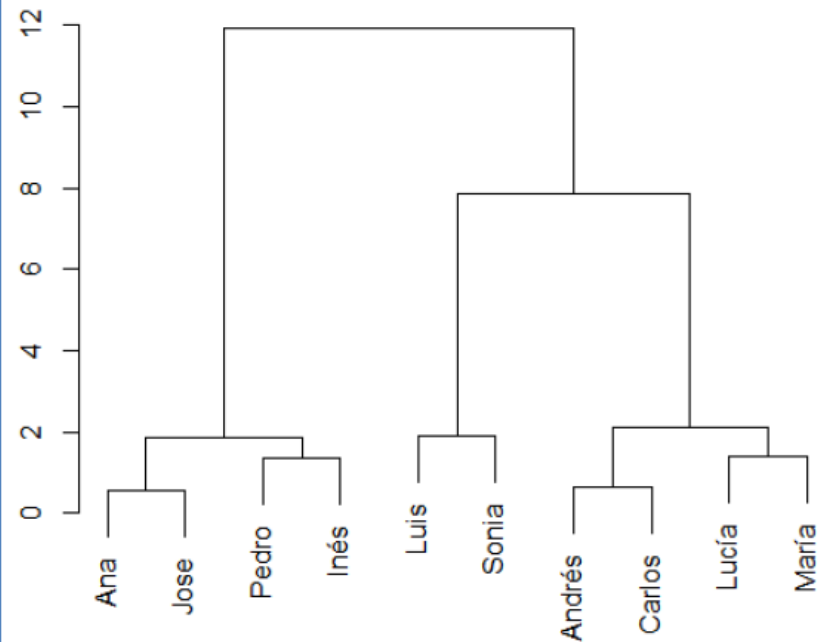
Cluster Dendrogram for Solution HClust.1



Cluster Dendrogram for Solution HClust.1



Cluster Dendrogram for Solution HClust.1



ALGORITMO

1. Inicialización: Se define P como el conjunto que contiene las clases conformadas por sólo un elemento.
2. Formación de nuevos nodos: Se funcionan los dos nodos más cercanos en sentido de la agregación elegida.
3. Actualización de P : Se agrega a P el nuevo nodo y se eliminan de él los nodos que lo conforman.
4. Test: Se detiene el algoritmo cuando el cardinal de P es mayor o igual a dos.



EJEMPLO

- Supongamos que tenemos los siguientes valores de disimilitud $s(\Omega = \{x_1, x_2, x_3, x_4\})$:

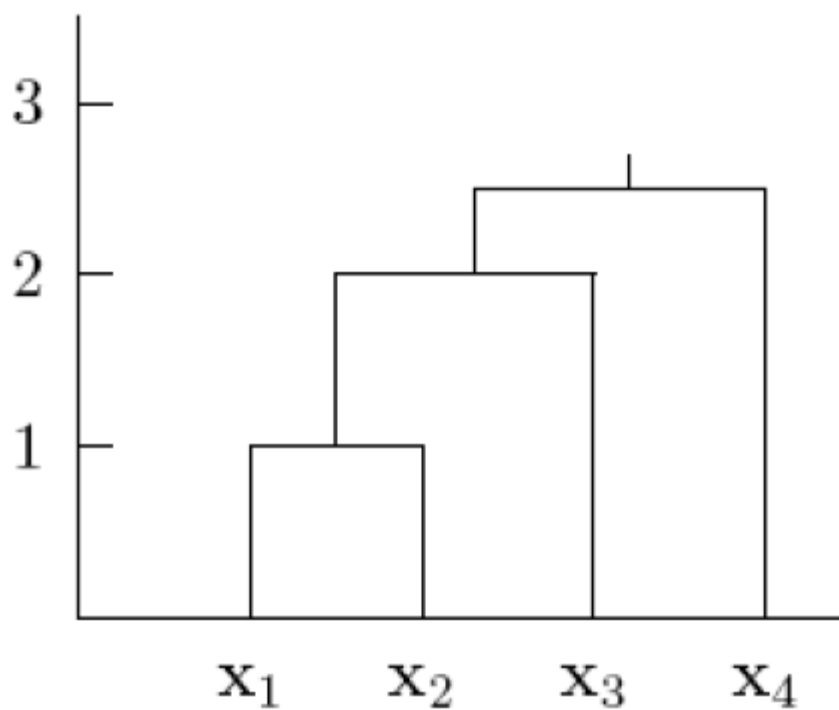
	x_1	x_2	x_3	x_4
x_1	0	1	3	5.5
x_2		0	2	4.5
x_3			0	2.5
x_4				0

- Vemos claramente que la disimilitud mínima se alcanza para la distancia entre x_1 y x_2 . Por lo tanto agregamos estos datos y utilizando la agregación del salto mínimo obtenemos:

	$\{x_1, x_2\}$	x_3	x_4
$\{x_1, x_2\}$	0	2	4.5
x_3		0	2.5
x_4			0



	$\{x_1, x_2, x_3\}$	x_4
$\{x_1, x_2, x_3\}$	0	2.5
x_4		0



TAREA

- Repetir el ejemplo anterior para la agregación del salto máximo.
- ¿Se obtiene el mismo dendograma?



EJEMPLO

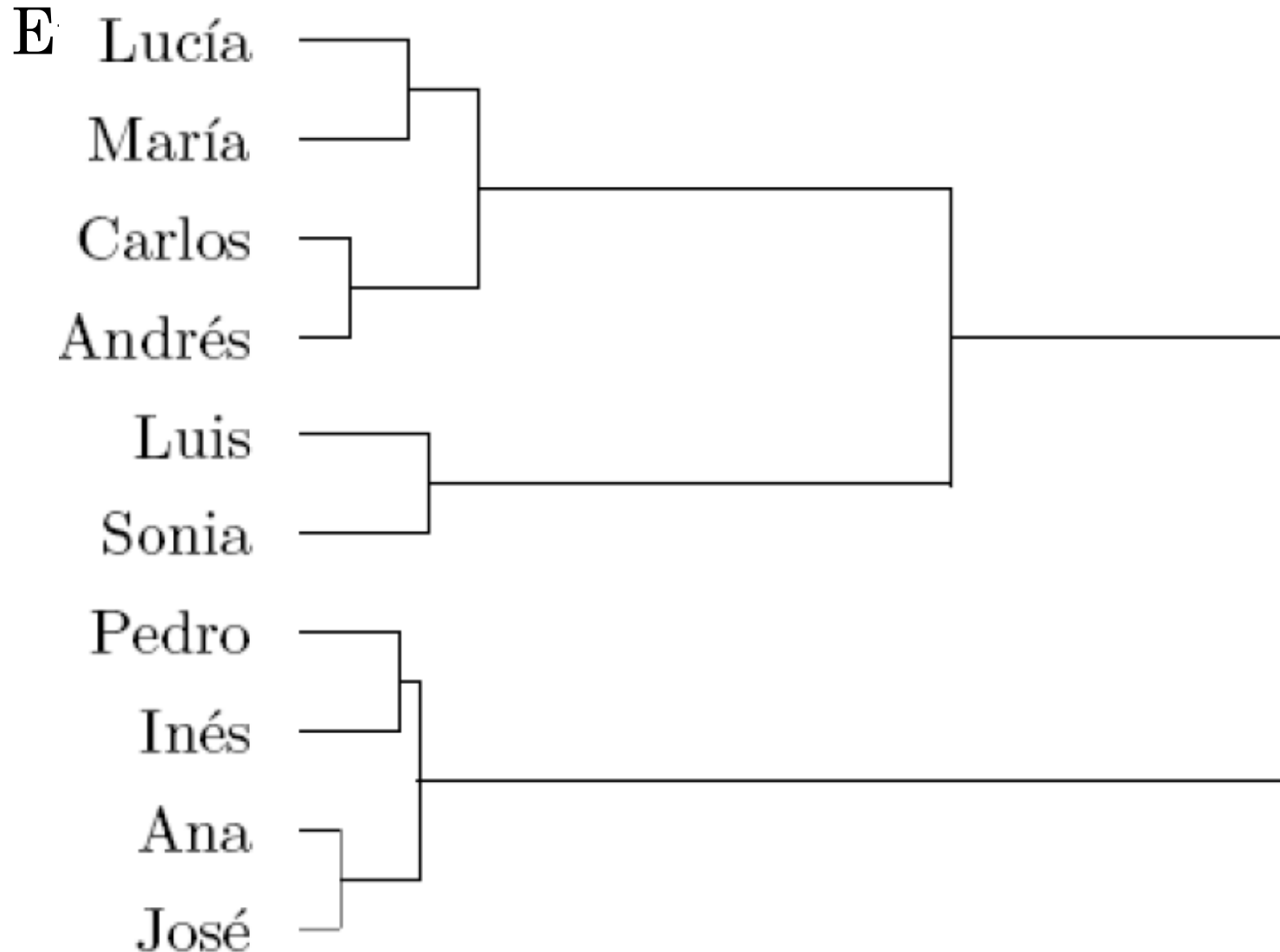
- Consideremos la tabla de datos siguiente que contiene notas obtenidas por diez estudiantes en cinco materias. Todas las notas están en escala de 1 a 10.

Estudiante	Matemáticas	Ciencias	Español	Historia	Ed. Física
Lucía	7.0	6.5	9.2	8.6	8.0
Pedro	7.5	9.4	7.3	7.0	7.0
Inés	7.6	9.2	8.0	8.0	7.5
Luis	5.0	6.5	6.5	7.0	9.0
Andrés	6.0	6.0	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8.0	6.5
Carlos	6.3	6.4	8.2	9.0	7.2
José	7.9	9.7	7.5	8.0	6.0
Sonia	6.0	6.0	6.5	5.5	8.7
María	6.8	7.2	8.7	9.0	7.0

Tabla 3.1: Tabla de datos de las notas escolares.



- La clasificación jerárquica usando la agregación de *Ward* con la distancia



INTERPRETACIÓN

Análisis de los Clústeres

	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7	6.5	9.2	8.6	8
Pedro	7.5	9.4	7.3	7	7
Inés	7.6	9.2	8	8	7.5
Luis	5	6.5	6.5	7	9
Andrés	6	6	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8	6.5
Carlos	6.3	6.4	8.2	9	7.2
José	7.9	9.7	7.5	8	6
Sonía	6	6	6.5	5.5	8.7
María	6.8	7.2	8.7	9	7

Centro Gravedad C1={Pedro, Inés, Ana, José}

Matemáticas	Ciencias	Español	Historia	EdFísica
7.7	9.475	7.625	7.75	6.75

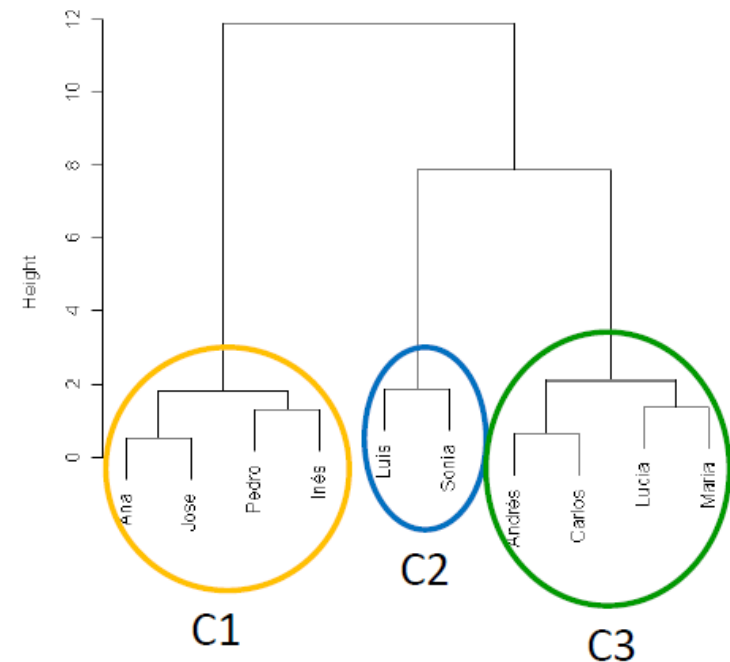
Centro Gravedad C2={Luis, Sonía}

Matemáticas	Ciencias	Español	Historia	EdFísica
5.5	6.25	6.5	6.25	8.85

Centro Gravedad C3={Lucía, Andrés, Carlos, María}

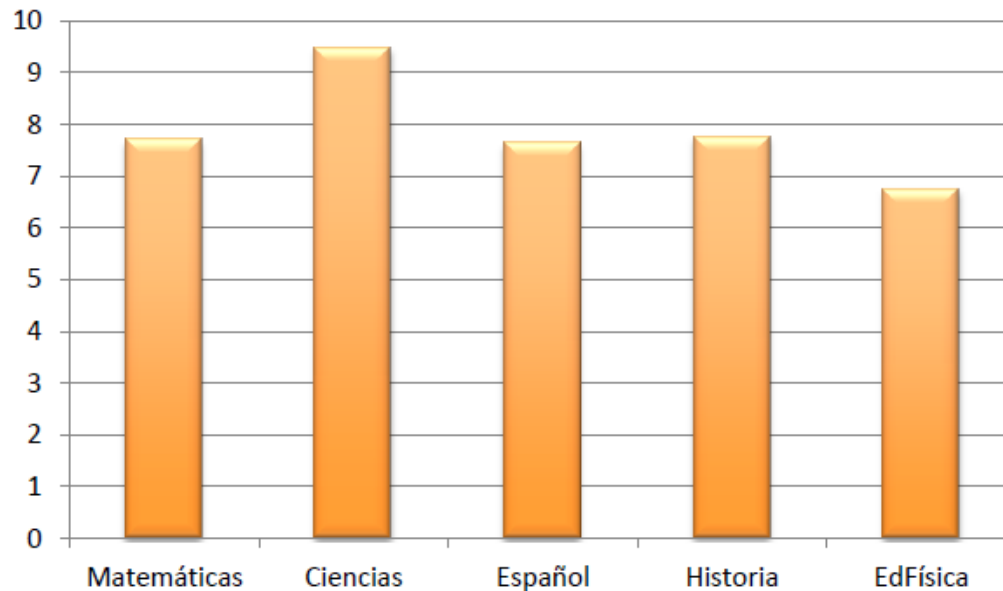
Matemáticas	Ciencias	Español	Historia	EdFísica
6.525	6.525	8.475	8.875	7.375

Cluster Dendrogram for Solution HClust.3



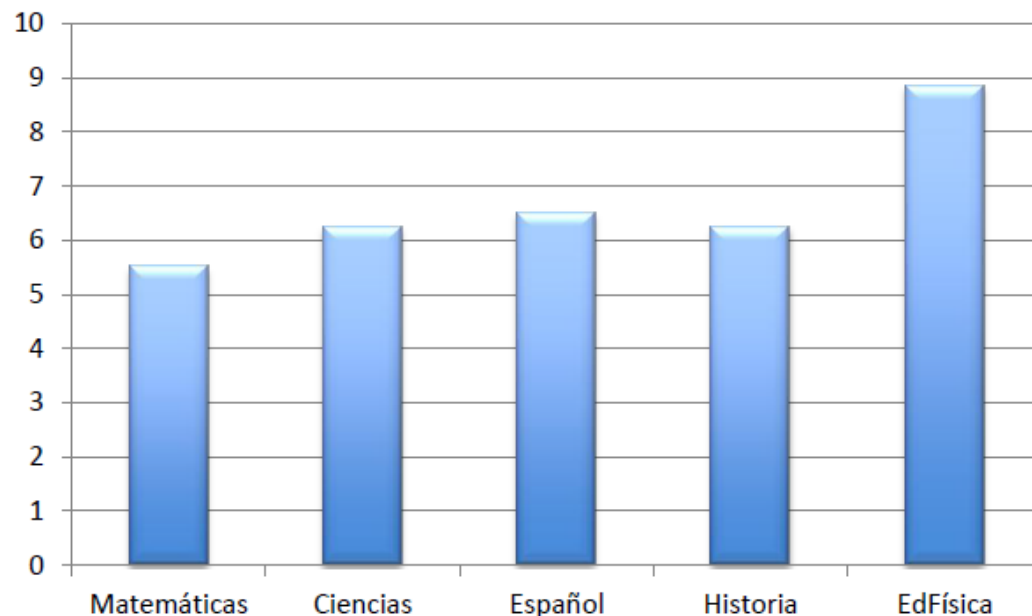
INTERPRETACIÓN HORIZONTAL

- El primer grupo está conformado por los estudiantes buenos en Ciencias, Matemáticas y con rendimiento promedio en las demás materias.



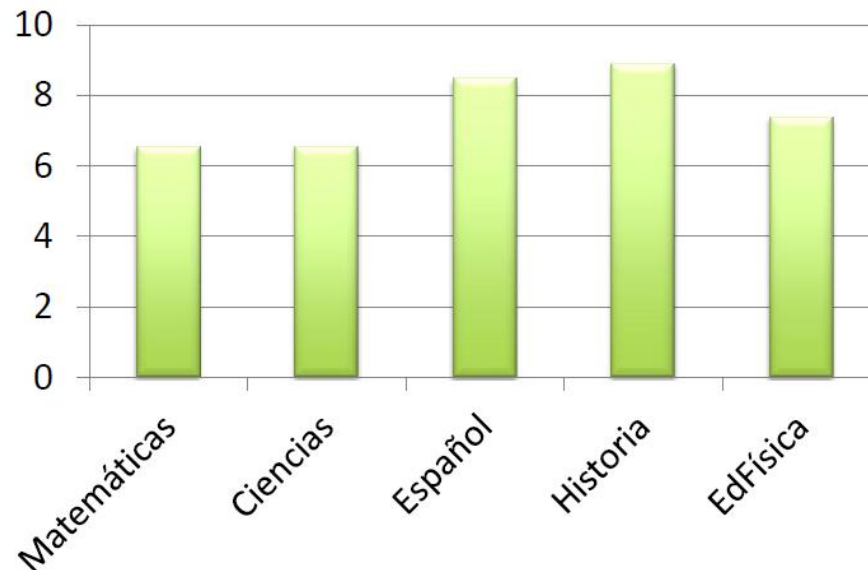
INTERPRETACIÓN HORIZONTAL

- El segundo grupo está conformado por los estudiantes buenos en Educación Física y que tienen un rendimiento de regular a malo en las demás materias.



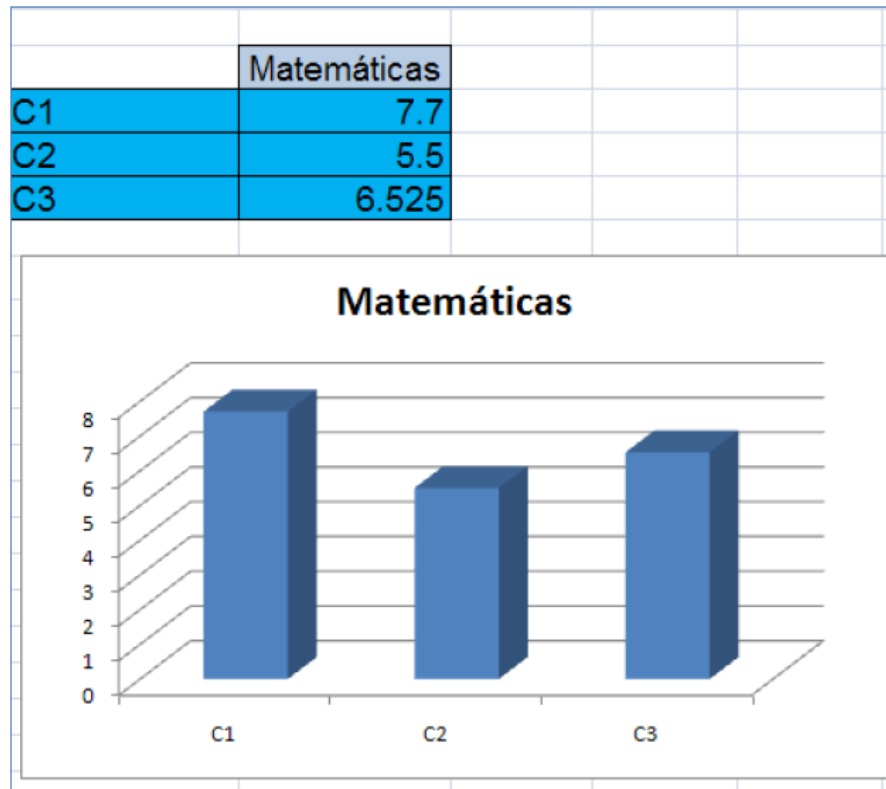
INTERPRETACIÓN HORIZONTAL

- El tercer grupo está conformado por aquellos estudiantes buenos en Español e Historia y con un rendimiento promedio en las demás asignaturas.



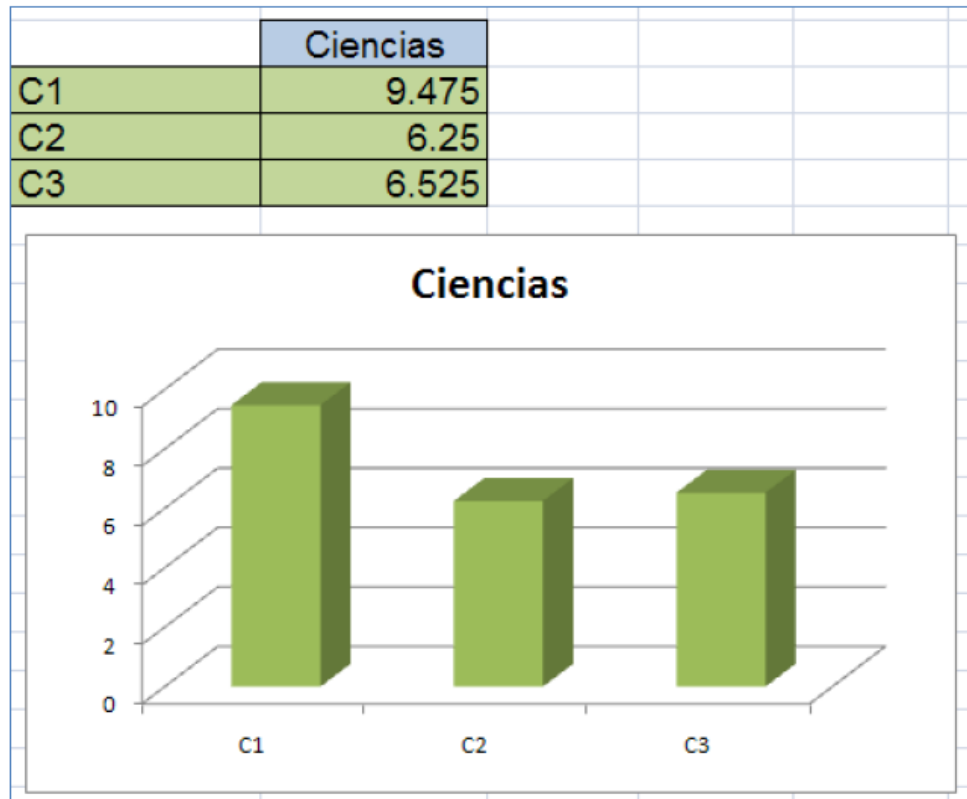
INTERPRETACIÓN VERTICAL

- El primer grupo es el que tiene mejores resultados en matemática.



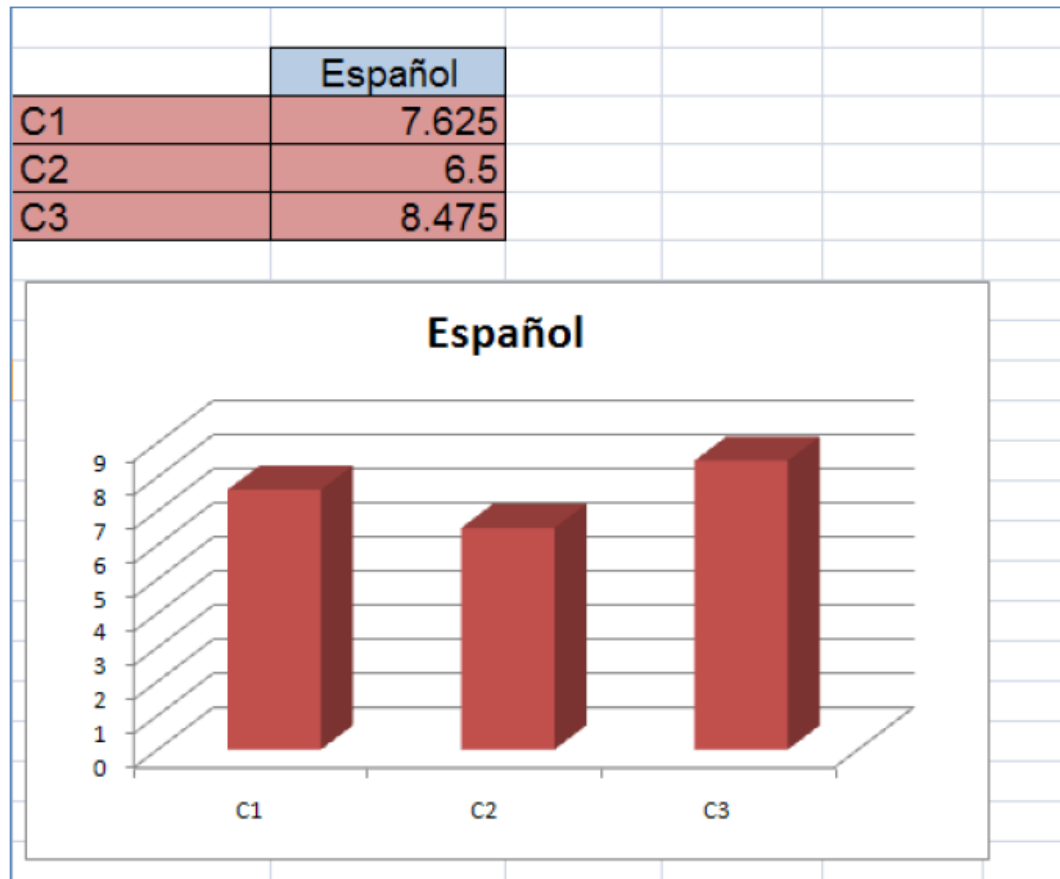
INTERPRETACIÓN VERTICAL

- El primer grupo es el con mejores resultados en Ciencias.



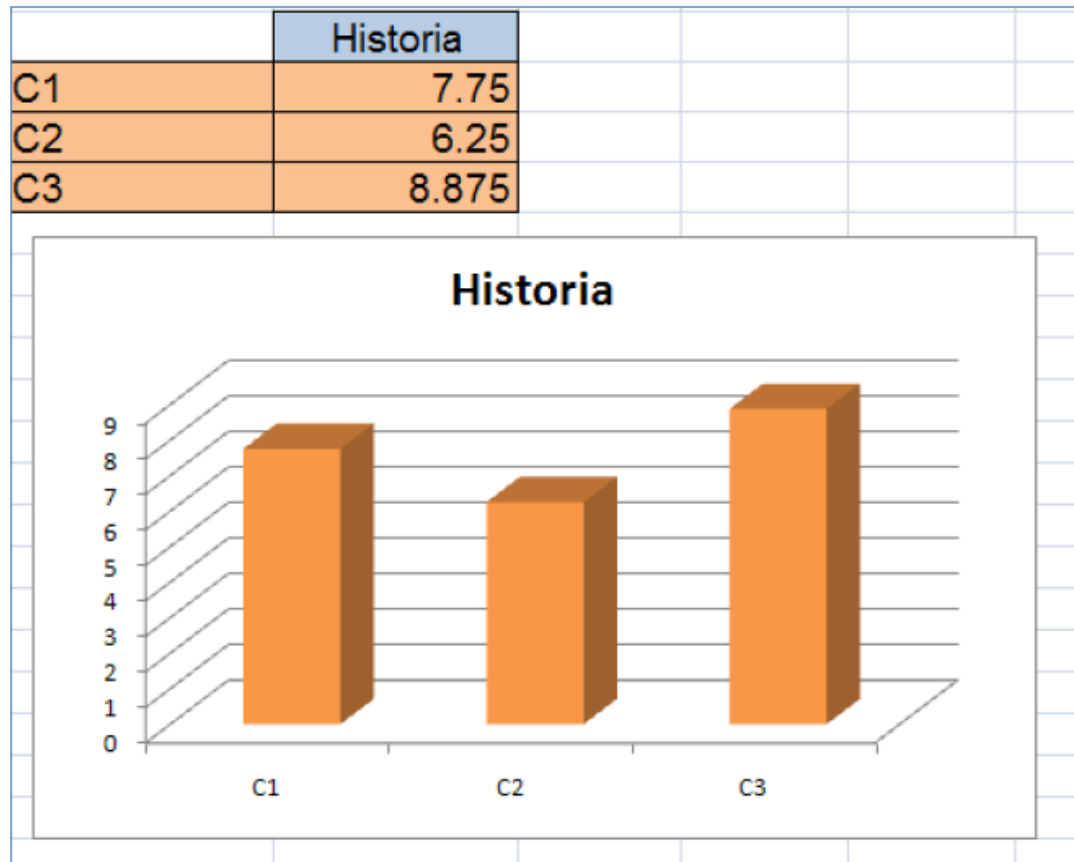
INTERPRETACIÓN VERTICAL

- El tercer grupo es el mejor en Español.



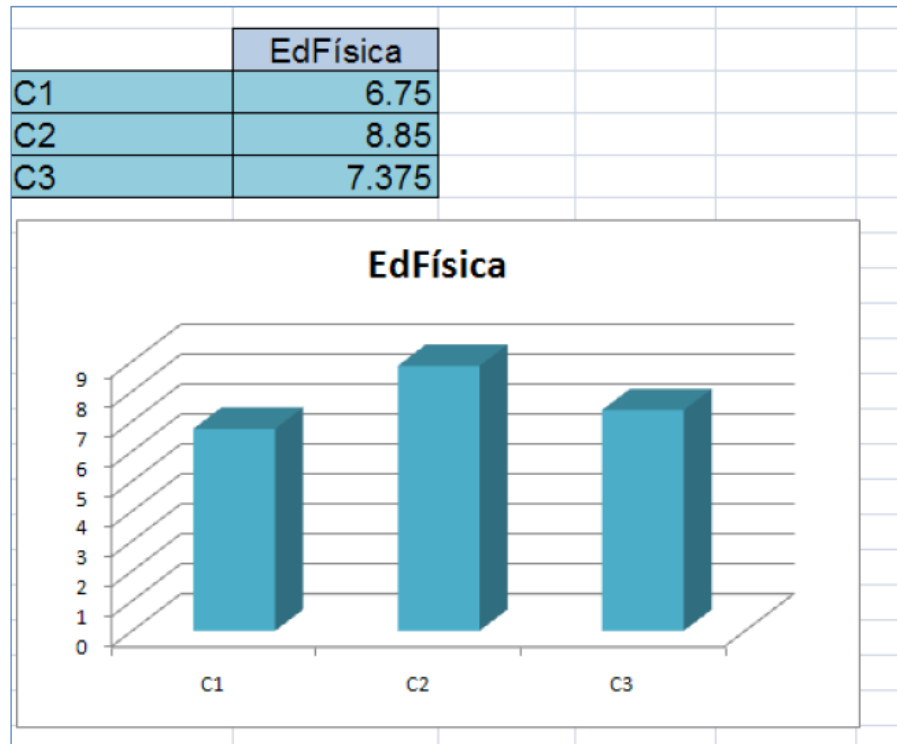
INTERPRETACIÓN VERTICAL

- El tercer grupo es el mejor en Historia.

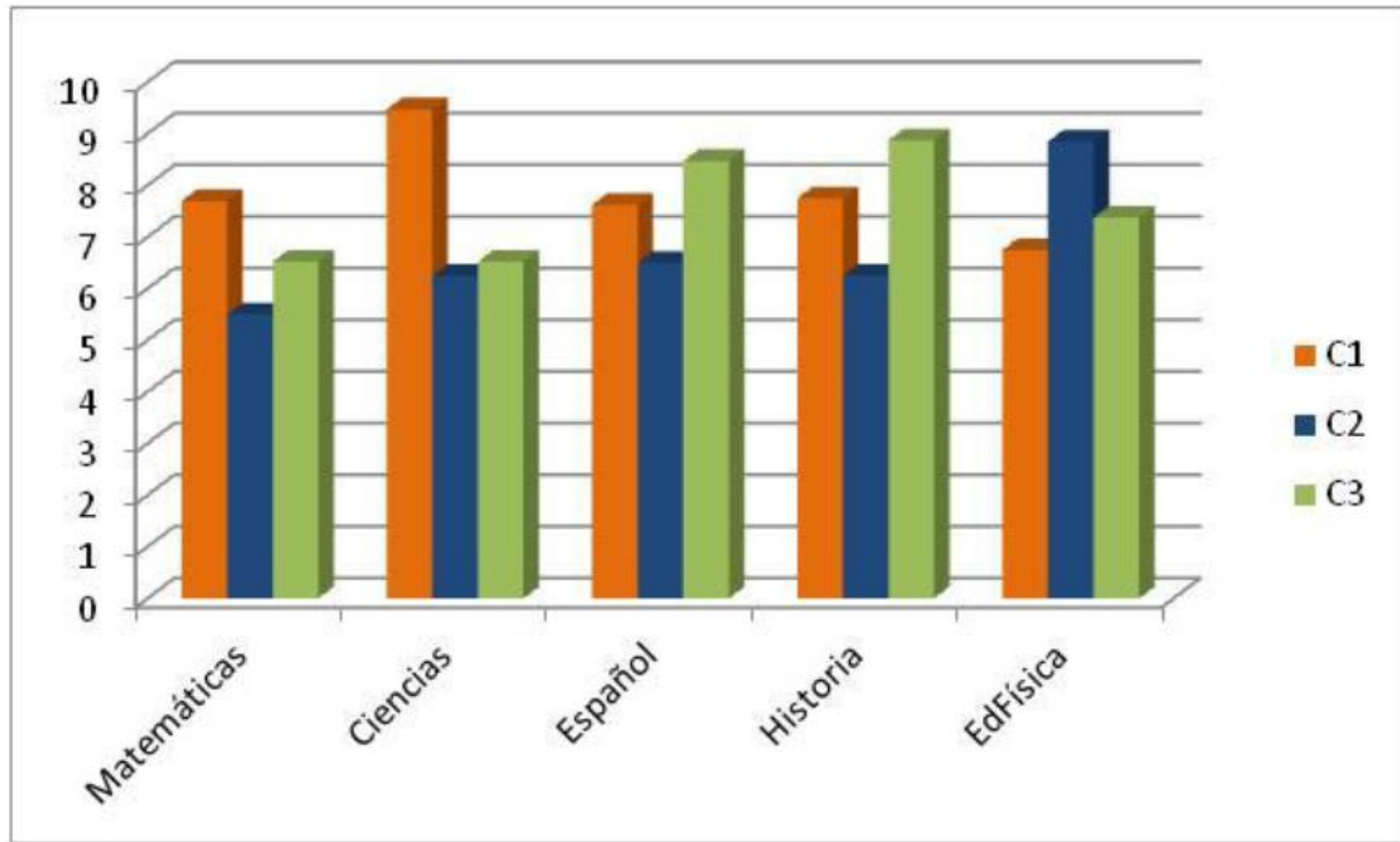


INTERPRETACIÓN VERTICAL

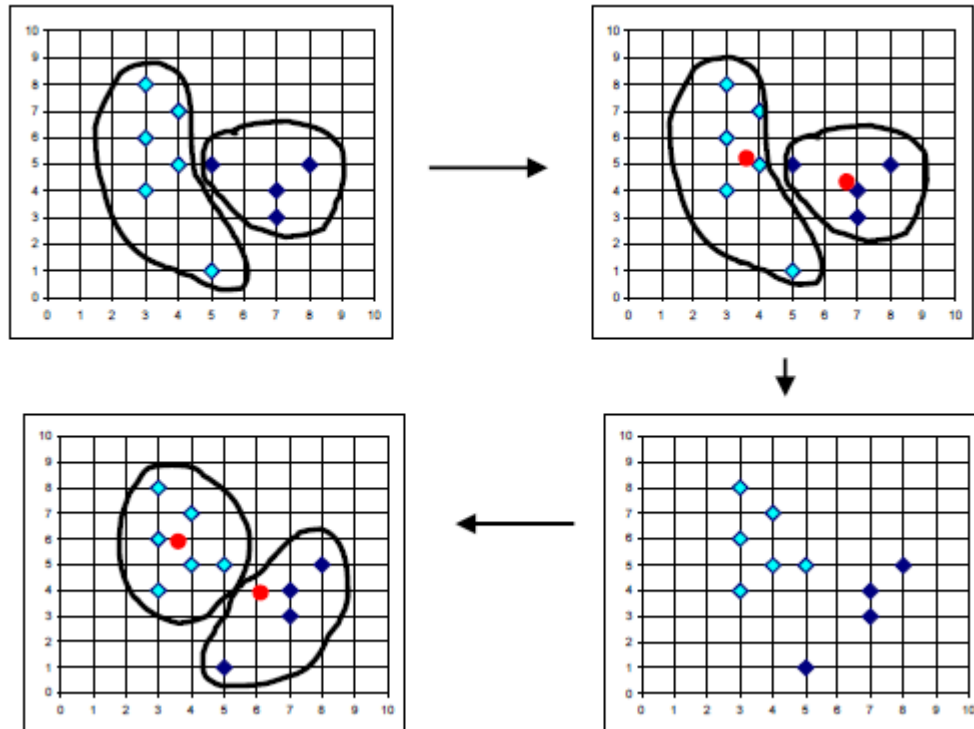
- El segundo grupo es el mejor en Educación Física.



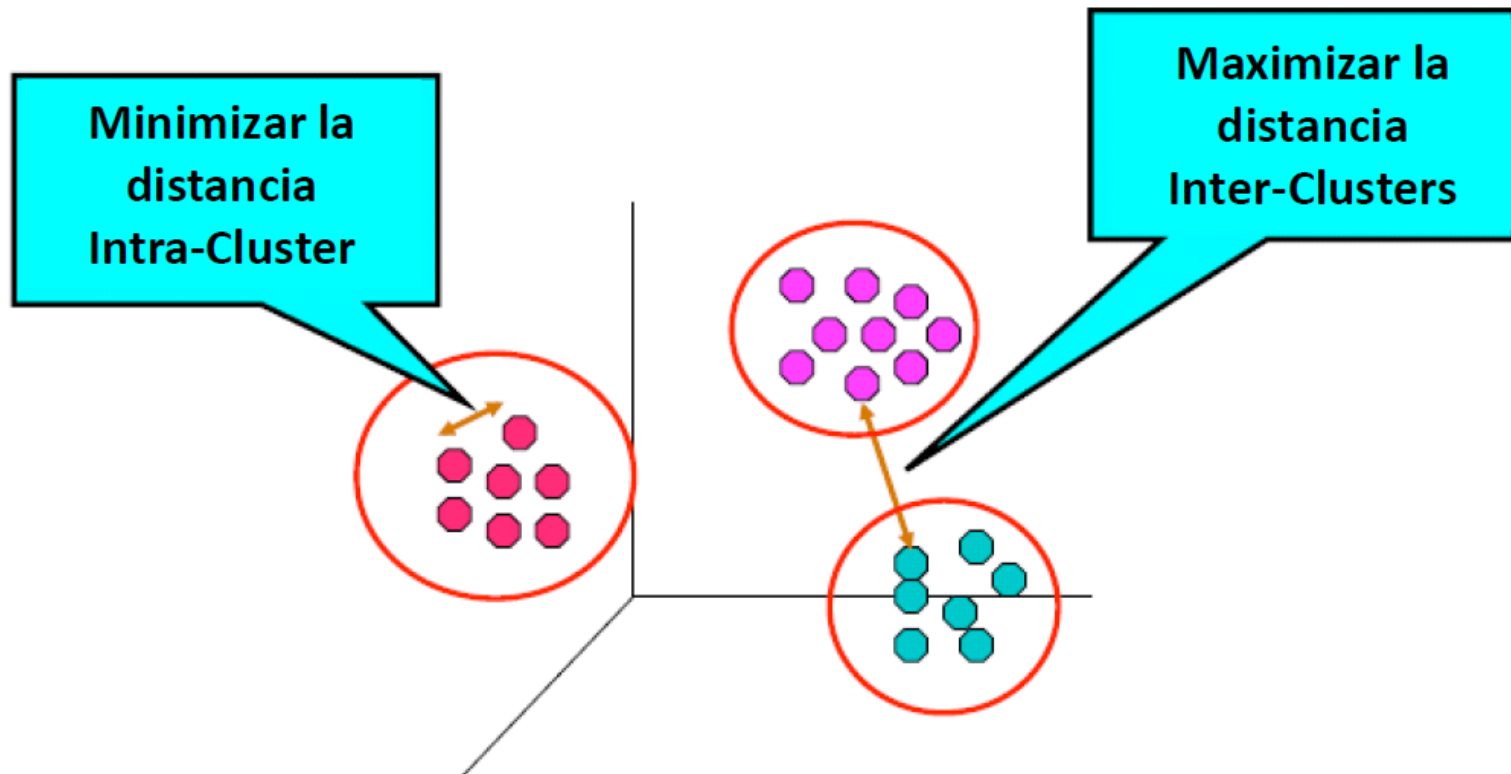
INTERPRETACIÓN HORIZONTAL-VERTICAL



MÉTODO K-MEANS



CRITERIO DE LA INERCIA



CENTRO DE GRAVEDAD

- El centro de gravedad de una clase C_K está dado por:

$$\mathbf{g}_k = \frac{1}{|C_k|} \sum_{i \in C_k} \mathbf{x}_i$$



Ejemplo: Estudiantes

Ver

NotasEscolaresExcelKMeans.xlsx

Análisis de los Clústeres					
	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7	6.5	9.2	8.6	8
Pedro	7.5	9.4	7.3	7	7
Inés	7.6	9.2	8	8	7.5
Luis	5	6.5	6.5	7	9
Andrés	6	6	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8	6.5
Carlos	6.3	6.4	8.2	9	7.2
José	7.9	9.7	7.5	8	6
Sonía	6	6	6.5	5.5	8.7
María	6.8	7.2	8.7	9	7
Centro Gravedad Total de la Nube de Puntos					
	Matemáticas	Ciencias	Español	Historia	EdFísica
	6.79	7.65	7.74	7.9	7.42
Centro Gravedad C1={Pedro,Inés,Ana,José}					
	Matemáticas	Ciencias	Español	Historia	EdFísica
	7.7	9.475	7.625	7.75	6.75
Centro Gravedad C2={Luis,Sonía}					
	Matemáticas	Ciencias	Español	Historia	EdFísica
	5.5	6.25	6.5	6.25	8.85
Centro Gravedad C3={Lucía,Andrés,Carlos,María}					
	Matemáticas	Ciencias	Español	Historia	EdFísica
	6.525	6.525	8.475	8.875	7.375

INERCIAS

- Inercia total:

$$I = \frac{1}{n} \sum_{i=1}^n ||\mathbf{x}_i - \mathbf{g}||^2$$

- Inercia inter-clases: Es la inercia de los centros de gravedad de cada clase respecto al centro de gravedad total.

$$B(P) = \sum_{k=1}^K \frac{|C_k|}{n} ||\mathbf{g}_k - \mathbf{g}||^2$$



INERCIAS

- Inercia intra-clases: Es la inercia de los individuos de cada clase con respecto al centro de gravedad de la clase.

$$W(P) = \sum_{k=1}^K I(C_k) = \frac{1}{n} \sum_{k=1}^K \sum_{i \in C_k} \|\mathbf{x}_i - \mathbf{g}_k\|^2$$



TEOREMA DE FISHER

$$I = B(P) + W(P)$$



EJERCICIO: EJEMPLO ESTUDIANTES

Análisis de los Clústeres						Cálculo de I =Inercia Total			Cálculo de $B(P)$ =Inercia Inter-Clases		
	Matemáticas	Ciencias	Español	Historia	EdFísica						
Lucía	7	6.5	9.2	8.6	8	4.3246			4.6434		
Pedro	7.5	9.4	7.3	7	7	4.7466			9.9291		
Inés	7.6	9.2	8	8	7.5	3.1426			2.8287		
Luis	5	6.5	6.5	7	9	9.3706		$B(P)=$	4.9747		
Andrés	6	6	7.8	8.9	7.3	4.3646					
Ana	7.8	9.6	7.7	8	6.5	5.6806					
Carlos	6.3	6.4	8.2	9	7.2	3.2726					
José	7.9	9.7	7.5	8	6	7.5186					
Sonia	6	6	6.5	5.5	8.7	12.283					
María	6.8	7.2	8.7	9	7	2.5106					
						$I=$	5.7214				
Centro Gravedad Total de la Nube de Puntos									$W(P)=$	0.7468	
	Matemáticas	Ciencias	Español	Historia	EdFísica						
	6.79	7.65	7.74	7.9	7.42						
Centro Gravedad C1={Pedro, Inés, Ana, José}											
	Matemáticas	Ciencias	Español	Historia	EdFísica						
	7.7	9.475	7.625	7.75	6.75						
Centro Gravedad C2={Luis, Sonia}											
	Matemáticas	Ciencias	Español	Historia	EdFísica						
	5.5	6.25	6.5	6.25	8.85						
Centro Gravedad C3={Lucía, Andrés, Carlos, María}											
	Matemáticas	Ciencias	Español	Historia	EdFísica						
	6.525	6.525	8.475	8.875	7.375						

$$I=B(P)+W(P) \quad 5.7214$$



OBJETIVO

- Recordemos que la idea es generar clases lo más homogéneas posibles. Esto coincide con minimizar la inercia inter-clases $W(P)$ y maximizar la inercia inter-clases $B(P)$.
- Gracias al Teorema de Fisher, dado que la inercia total es constante, la optimización de cualquiera de estas inercias inmediatamente optimiza la otra.



EJERCICIO COMBINATORIO

- Calcule el número de particiones en dos clases de un conjunto de 60 elementos.
- Repita lo anterior para 100 elementos.



MÉTODO K-MEANS

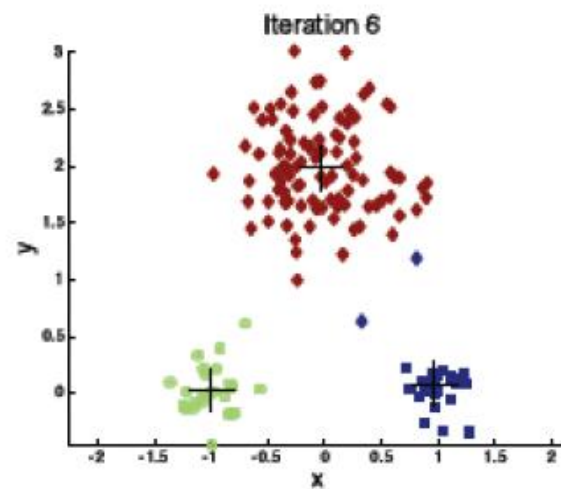
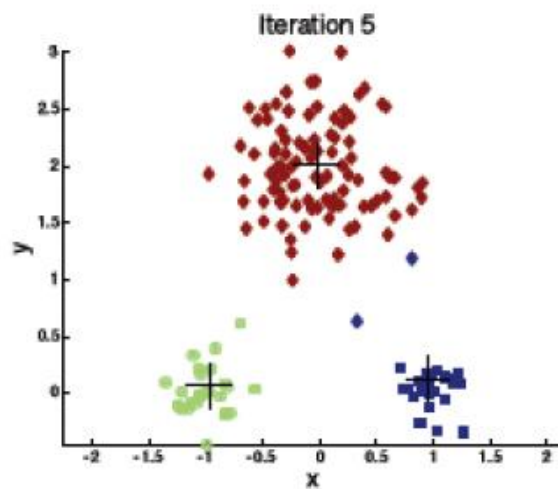
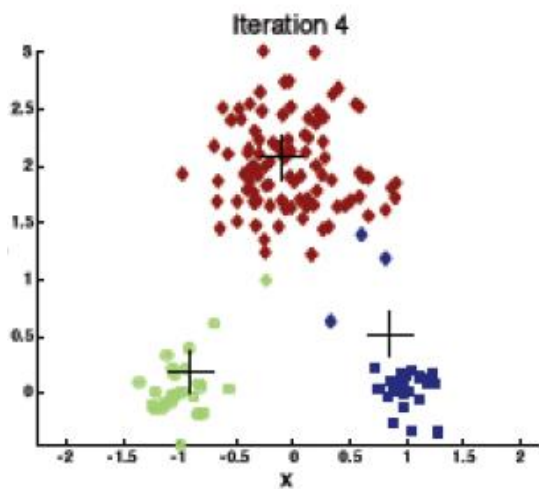
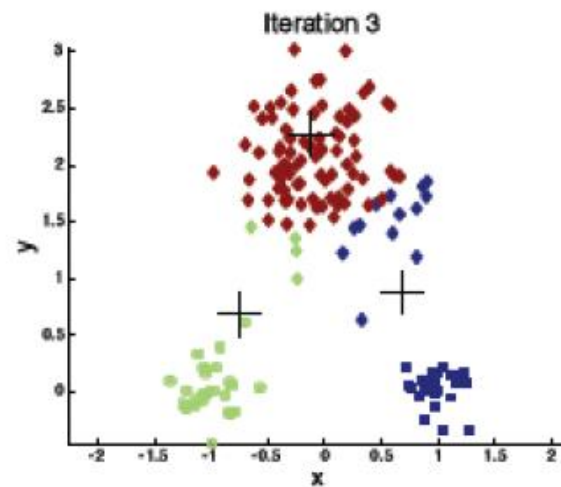
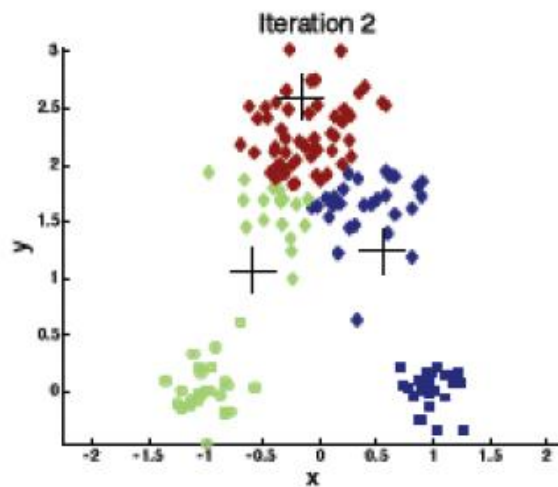
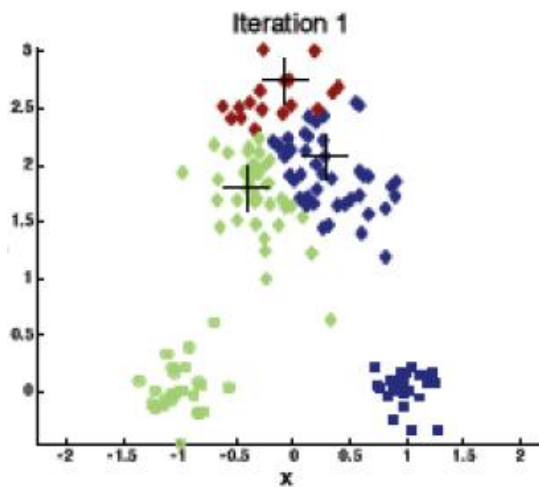
Objetivo: Encontrar una partición P del conjunto de puntos y representantes de cada una de las clases de modo que $W(P)$ sea mínima.



MÉTODO DE FORGY

- Consiste en un método de reasignación-recentraje que itera sucesivamente las dos operaciones siguientes hasta lograr convergencia:
 1. Representar una clase por su centro de gravedad;
 2. Asignar los objetos a la clase del centro de gravedad más cercano.



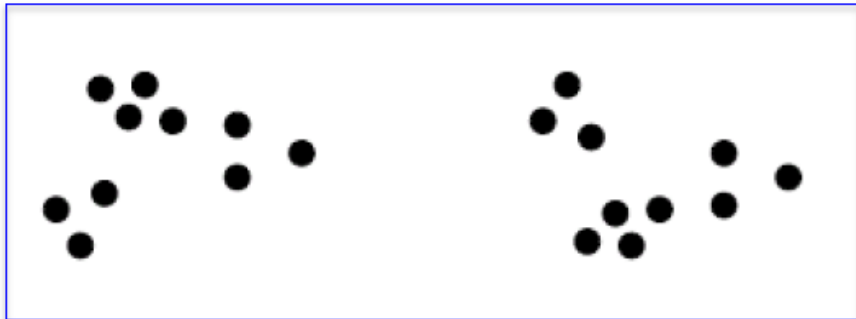


MÉTODO DE MCQUEEN

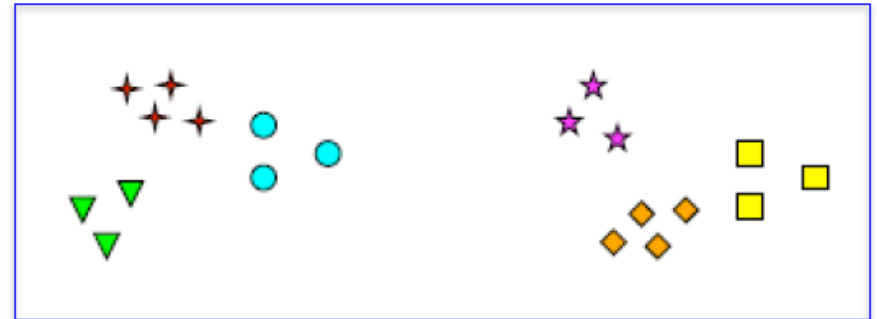
- Tal y como en el método de Forgy, las clases son representadas por su centro de gravedad y se examina a cada individuo de modo de asignarlo a la clase más cercana.
- La diferencia con el método de Forgy radica en que luego de asignar un individuo a una clase, el centro de ésta es re-calculado inmediatamente.



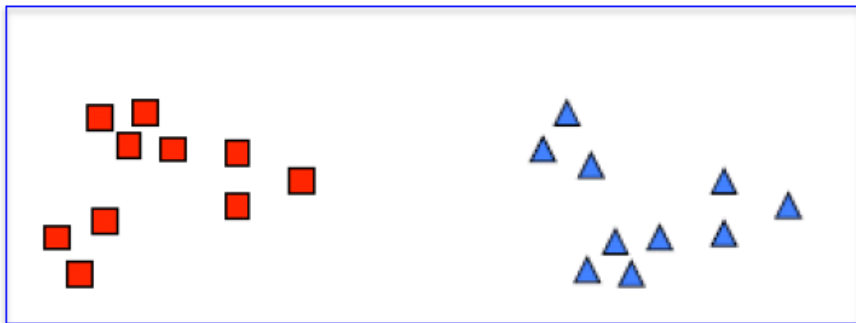
¿CUÁNTAS CLASES?



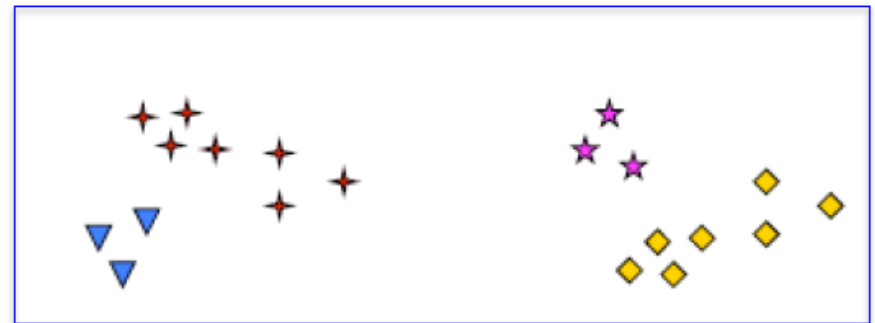
Datos originales



6 clústeres



2 clústeres



4 clústeres



IDEA

- Graficar la inercia intra-clases versus el número de clases. A esto se le llama “codo”.

