

# Científico de Datos

Aliados:



**Microsoft**

Vigilada Mineducación



# APRENDIZAJE NO SUPERVISADO

Aliados:



**Microsoft**

Vigilada Mineducación

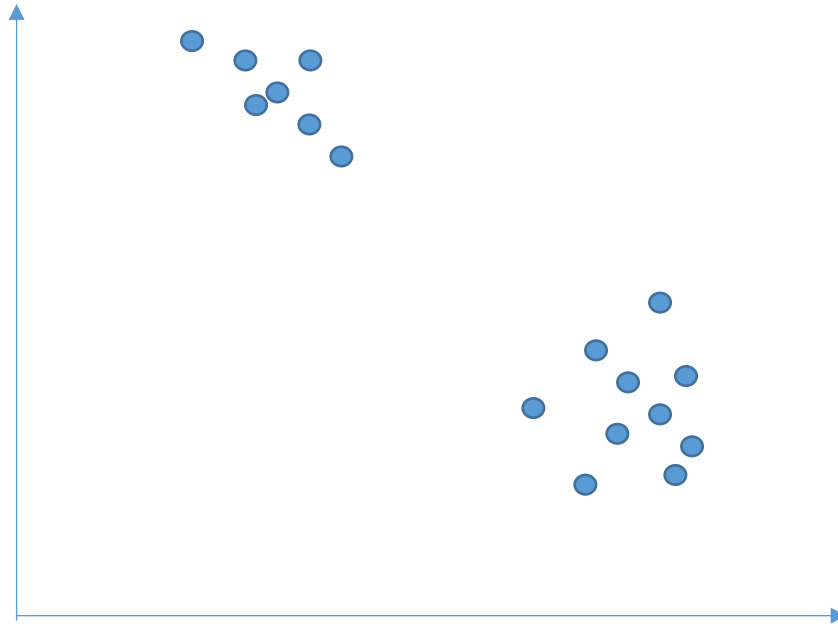


# DEFINTION

- Encontrar etiquetas (labels) a objetos sin etiquetas
- Es el proceso de particionar un conjunto de objetos en subconjuntos.

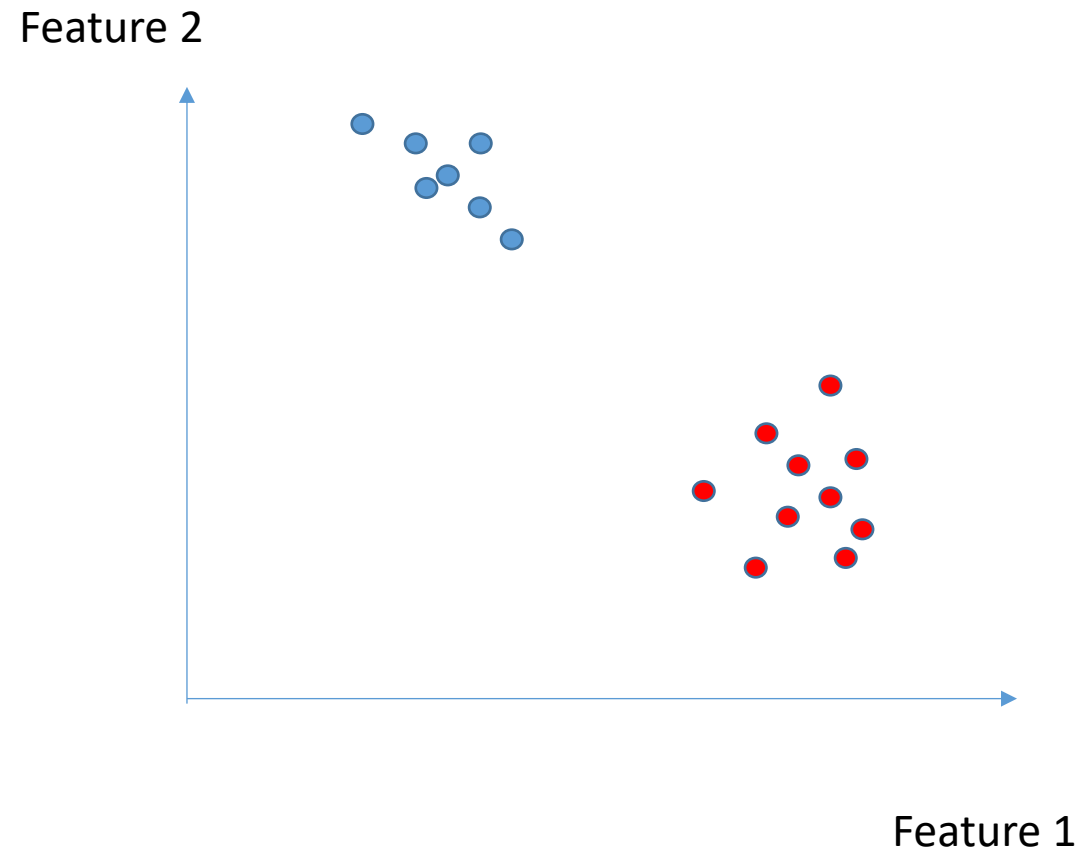
# DEFINTION

Feature 2



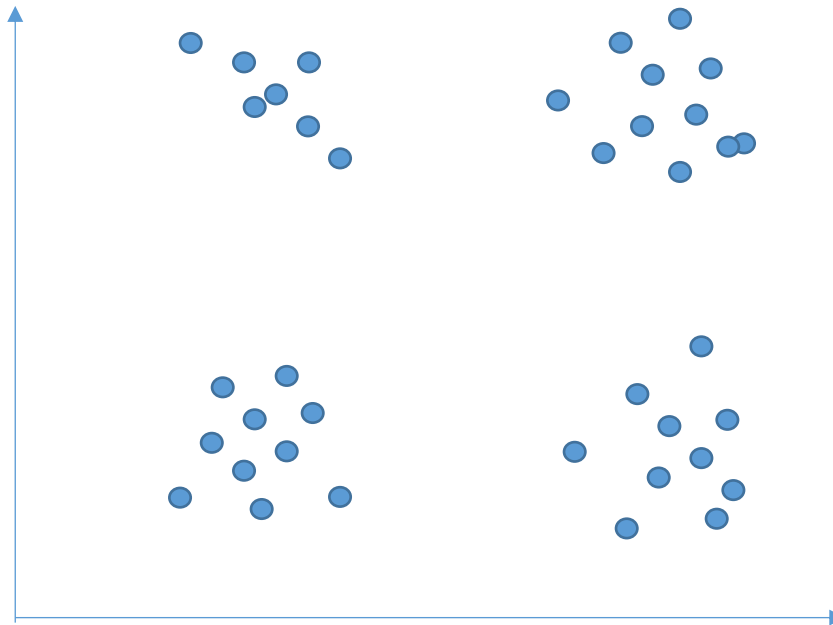
Featue 1

# DEFINTION



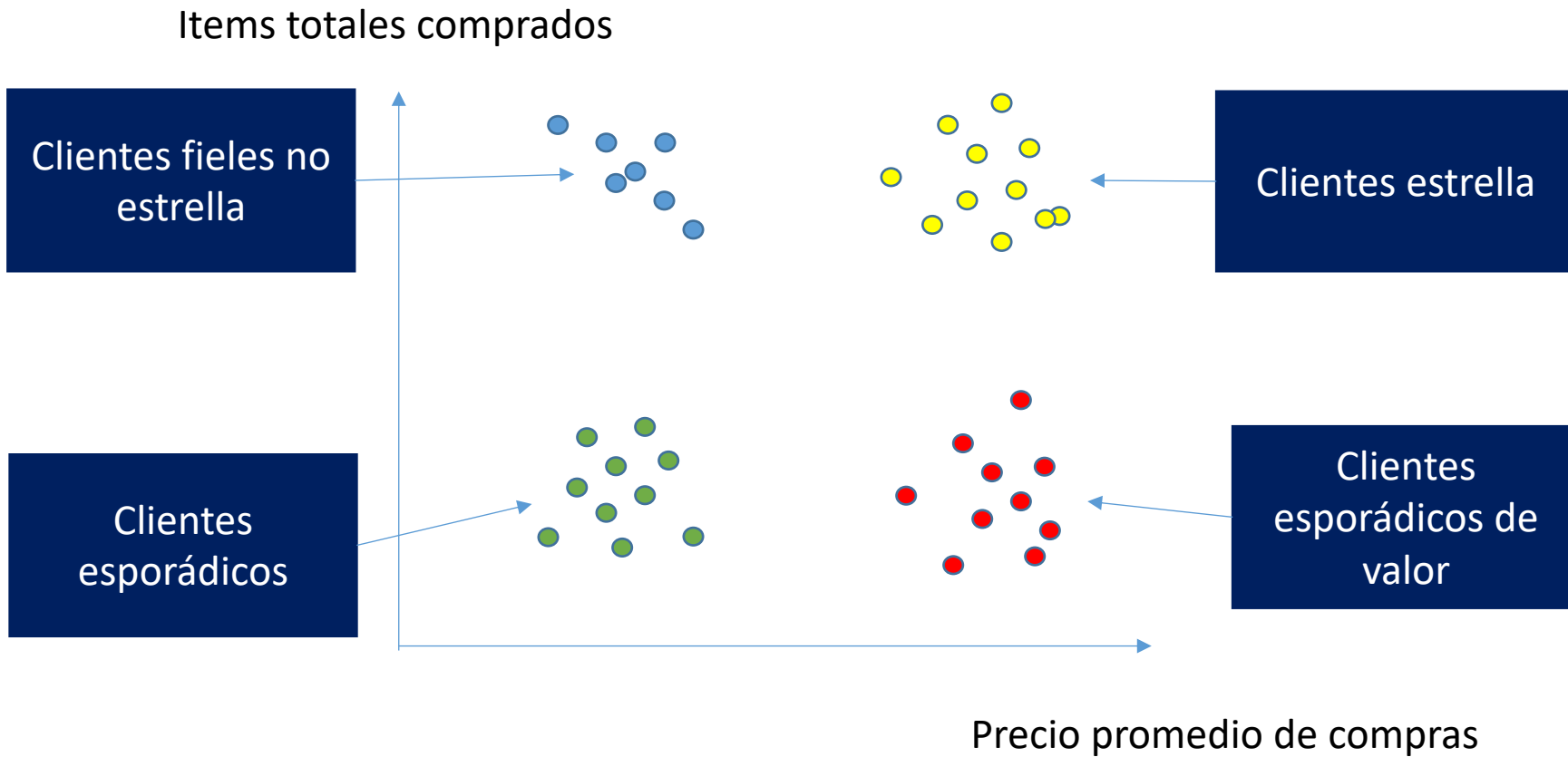
# DEFINTION

Items totales comprados



Precio promedio de compras

# DEFINTION



# DEFINITION





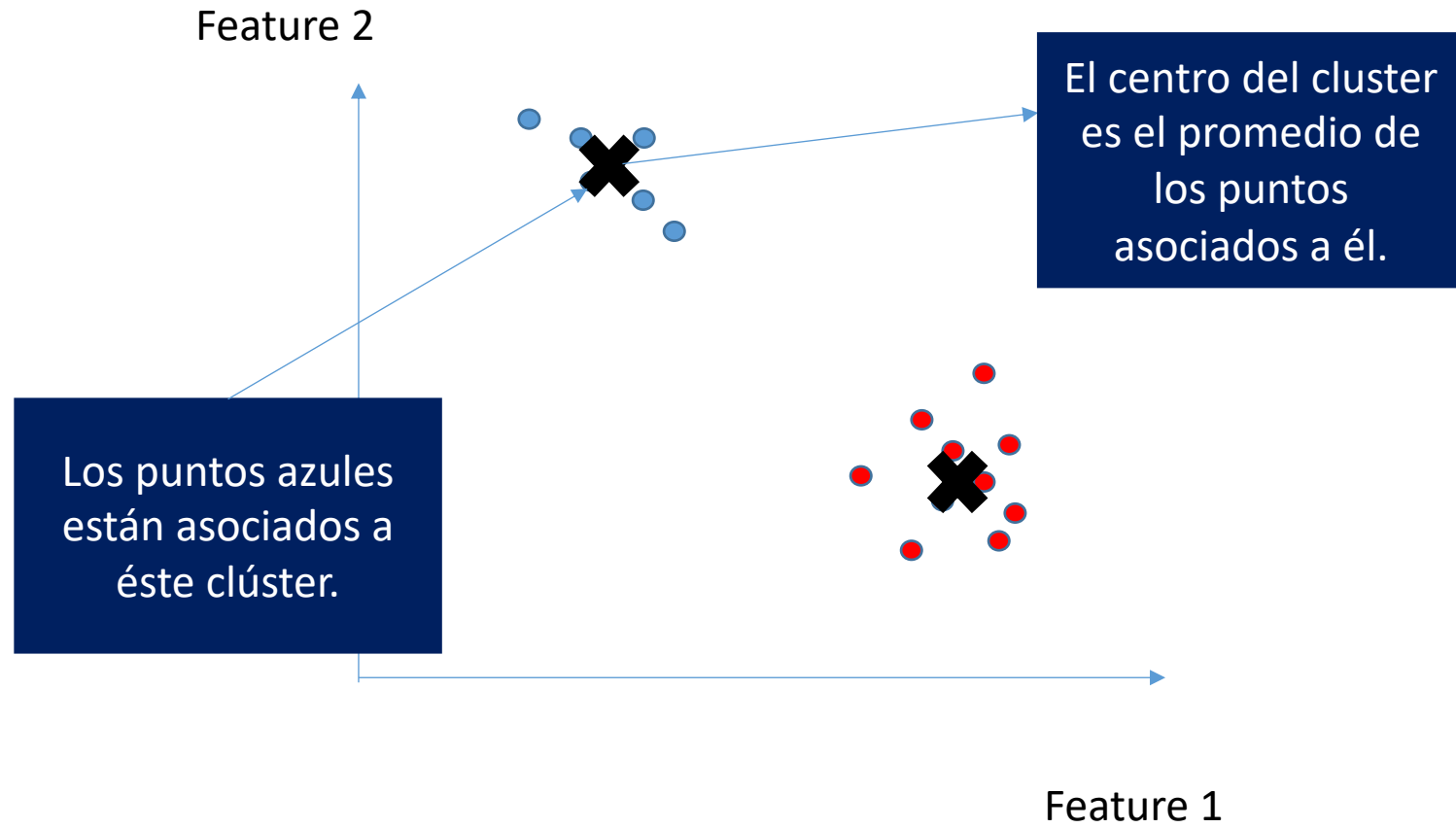
# DEFINTION

Method	General Characteristics
Partitioning methods	<ul style="list-style-type: none"><li>– Find mutually exclusive clusters of spherical shape</li><li>– Distance-based</li><li>– May use mean or medoid (etc.) to represent cluster center</li><li>– Effective for small- to medium-size data sets</li></ul>
Hierarchical methods	<ul style="list-style-type: none"><li>– Clustering is a hierarchical decomposition (i.e., multiple levels)</li><li>– Cannot correct erroneous merges or splits</li><li>– May incorporate other techniques like microclustering or consider object “linkages”</li></ul>
Density-based methods	<ul style="list-style-type: none"><li>– Can find arbitrarily shaped clusters</li><li>– Clusters are dense regions of objects in space that are separated by low-density regions</li><li>– Cluster density: Each point must have a minimum number of points within its “neighborhood”</li><li>– May filter out outliers</li></ul>

# K- MEANS

- Método de clustering duro (hard) porque cada dato pertenece sólo a un cluster dado.
- Dos conceptos: **centroides y puntos.**

# K-MEANS



# K-MEANS

Inicializar centros  $k = \{\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k\}$

1. Para cada centro identificamos qué datos tiene asociados.
2. Calculamos el nuevo valor del centro como la media de todos los datos.
3. Evaluamos criterio de convergencia.  
Si no se cumple, volver al paso 1.

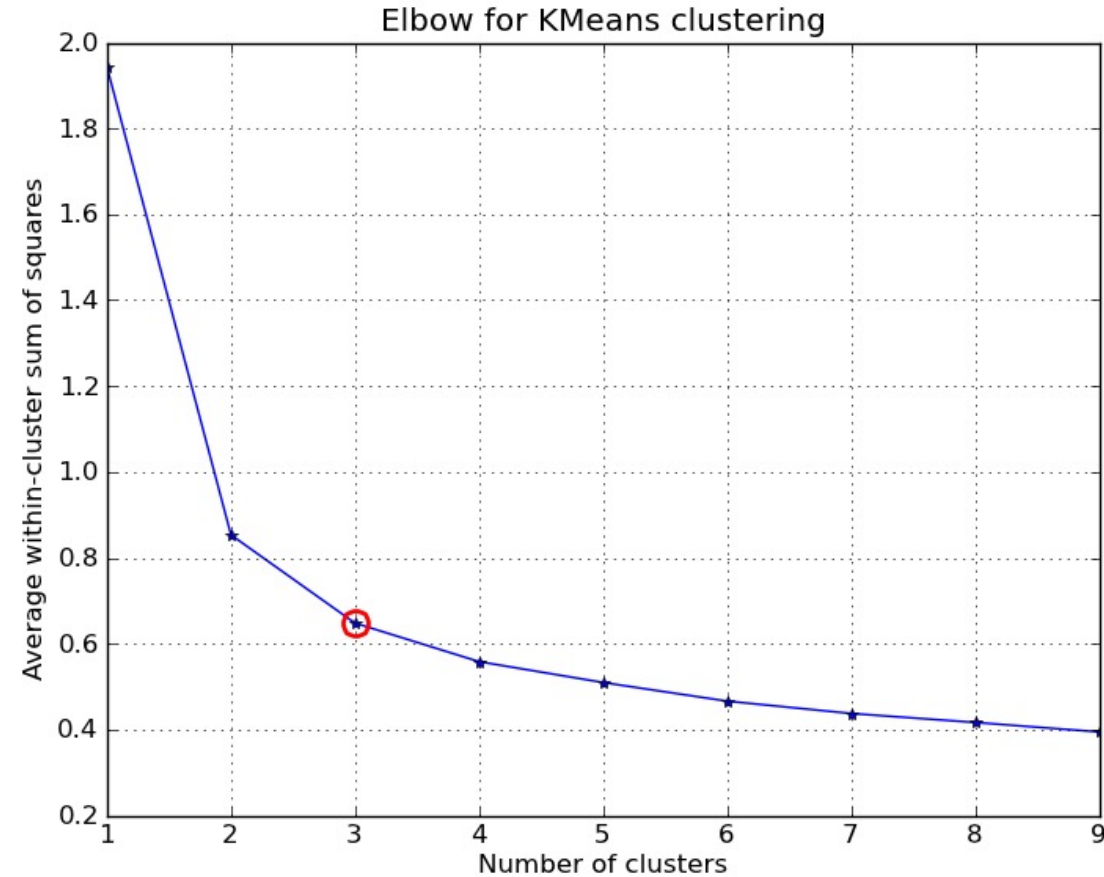
# CONSIDERACIONES

- Asume clusters circulares
- Se debe determinar K: hay algunas heurísticas para elegir K.
- Se deben inicializar centros. Esto hace que pueda caer en mínimos locales. Se aconseja inicializar varias veces.
- Sensible a outliers.
- Problemas con muchas dimensiones.
- ¡Converge!

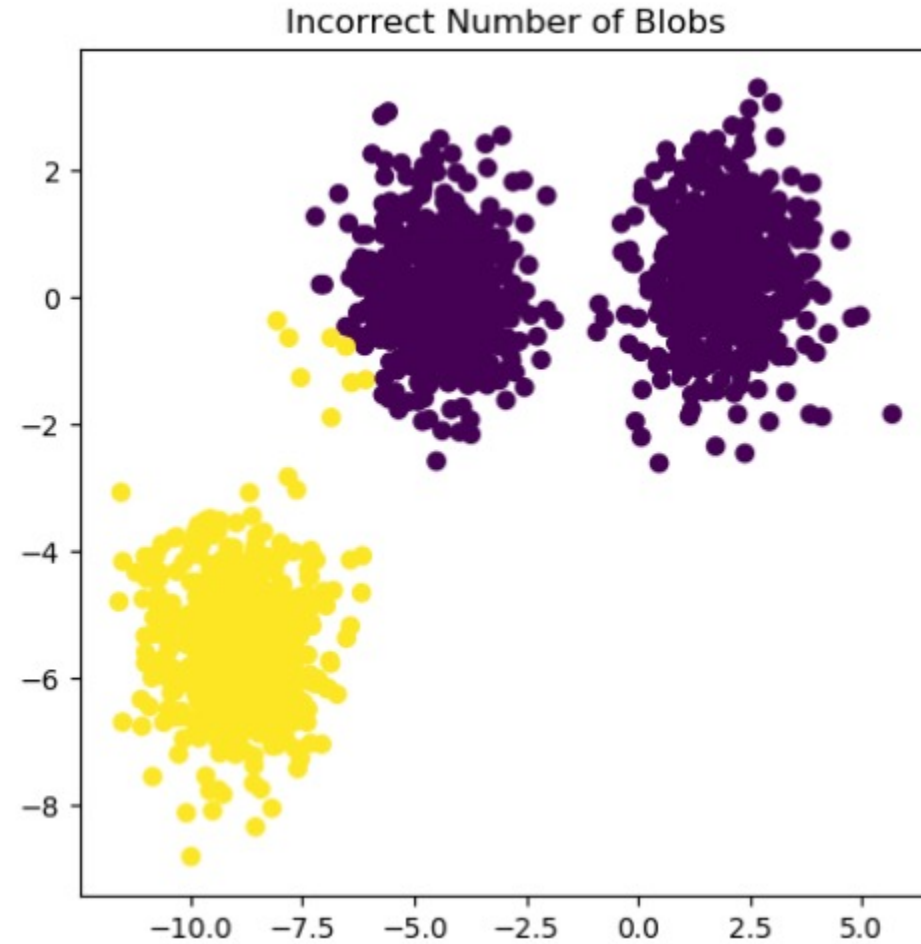
# CÓMO ELEGIR EL NÚMERO DE CLUSTERS (K)

- Heurística: Criterio del codo (elbow criterion).
- Correr el algoritmo K-Means de 1 hasta un valor máximo. Evaluar la calidad del clustering.
- Seleccionar k cuando la calidad del clustering comience a disminuir poco.

# CÓMO ELEGIR EL NÚMERO DE CLUSTERS (K)

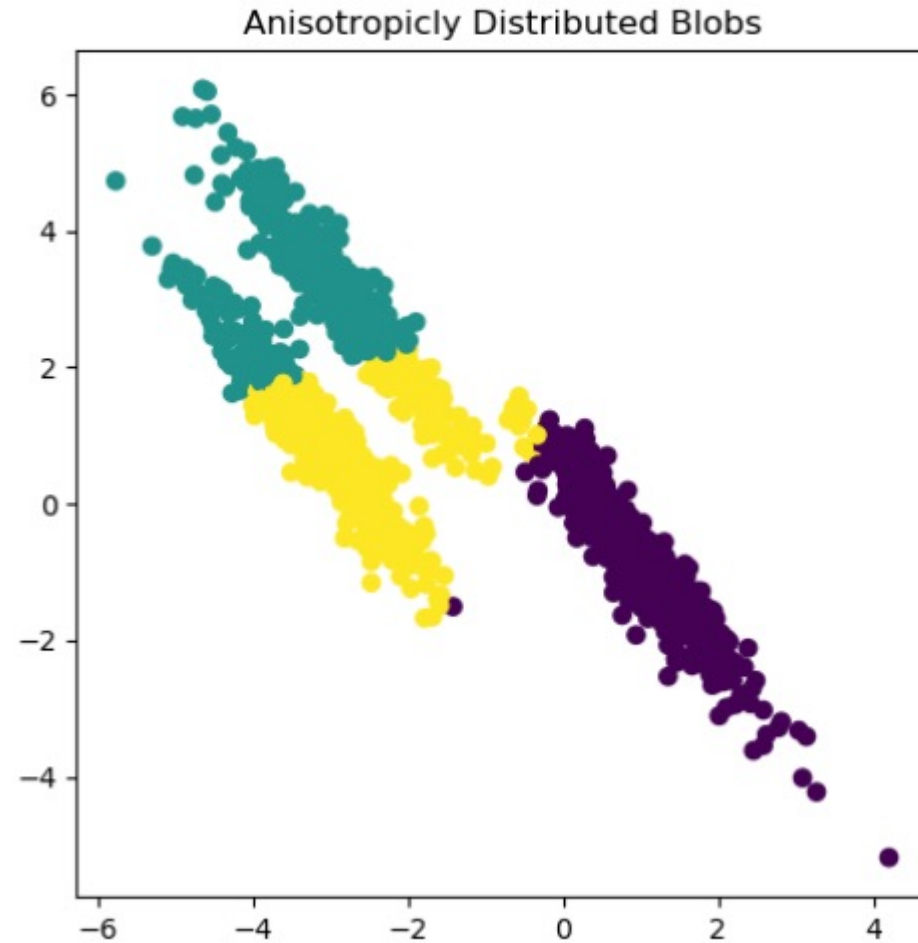


# NÚMERO INCORRECTO DE CLUSTERS (K)

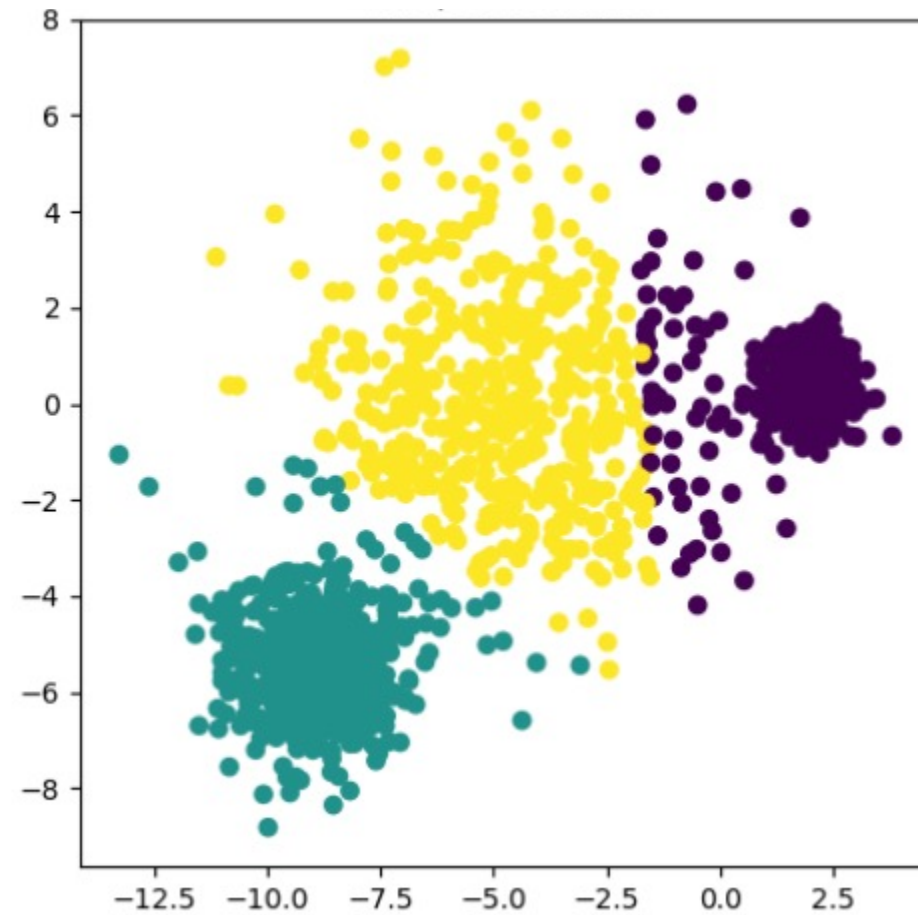




# DISTRIBUCIÓN NO ADECUADA PARA K-MEANS



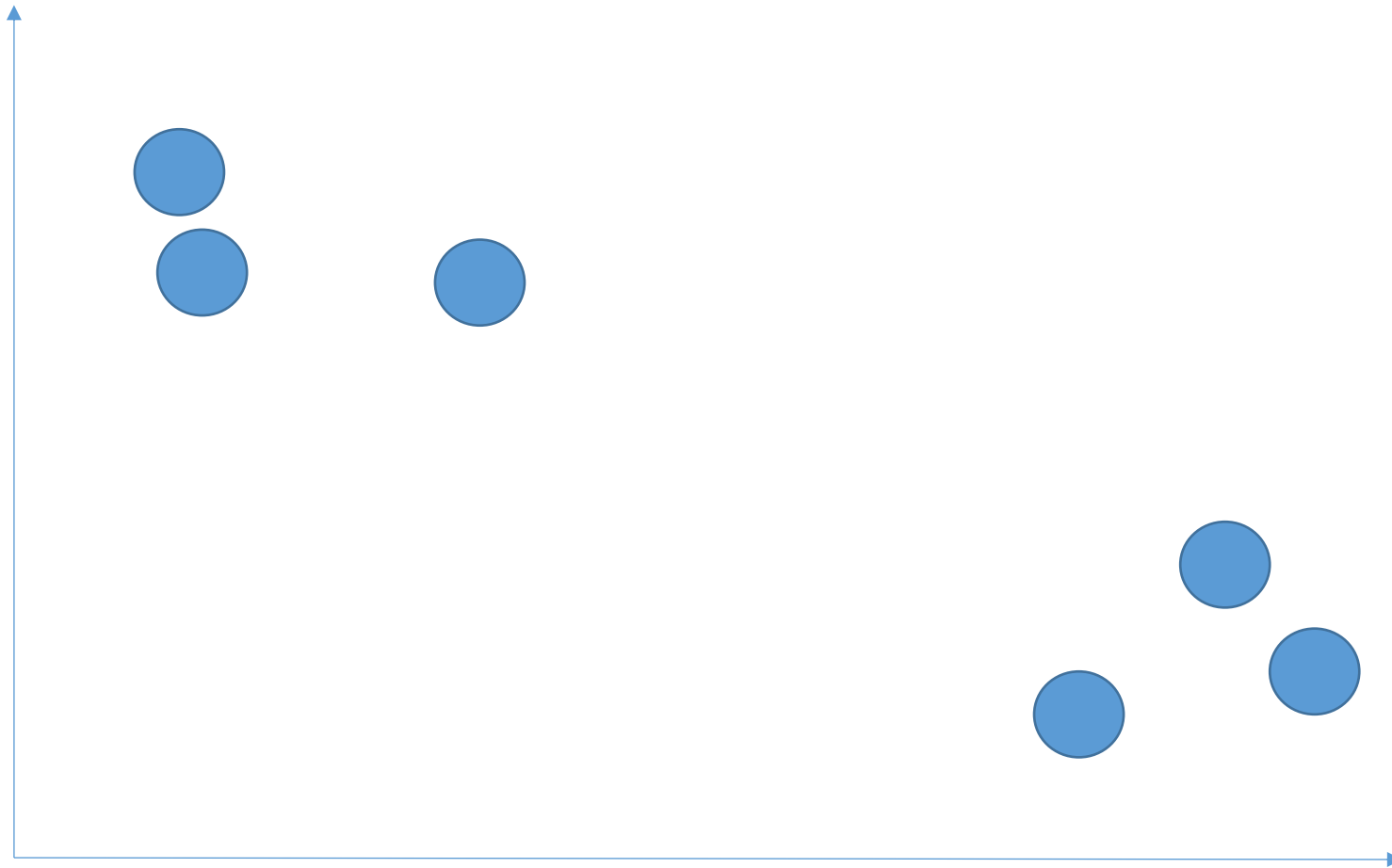
# K-MEANS



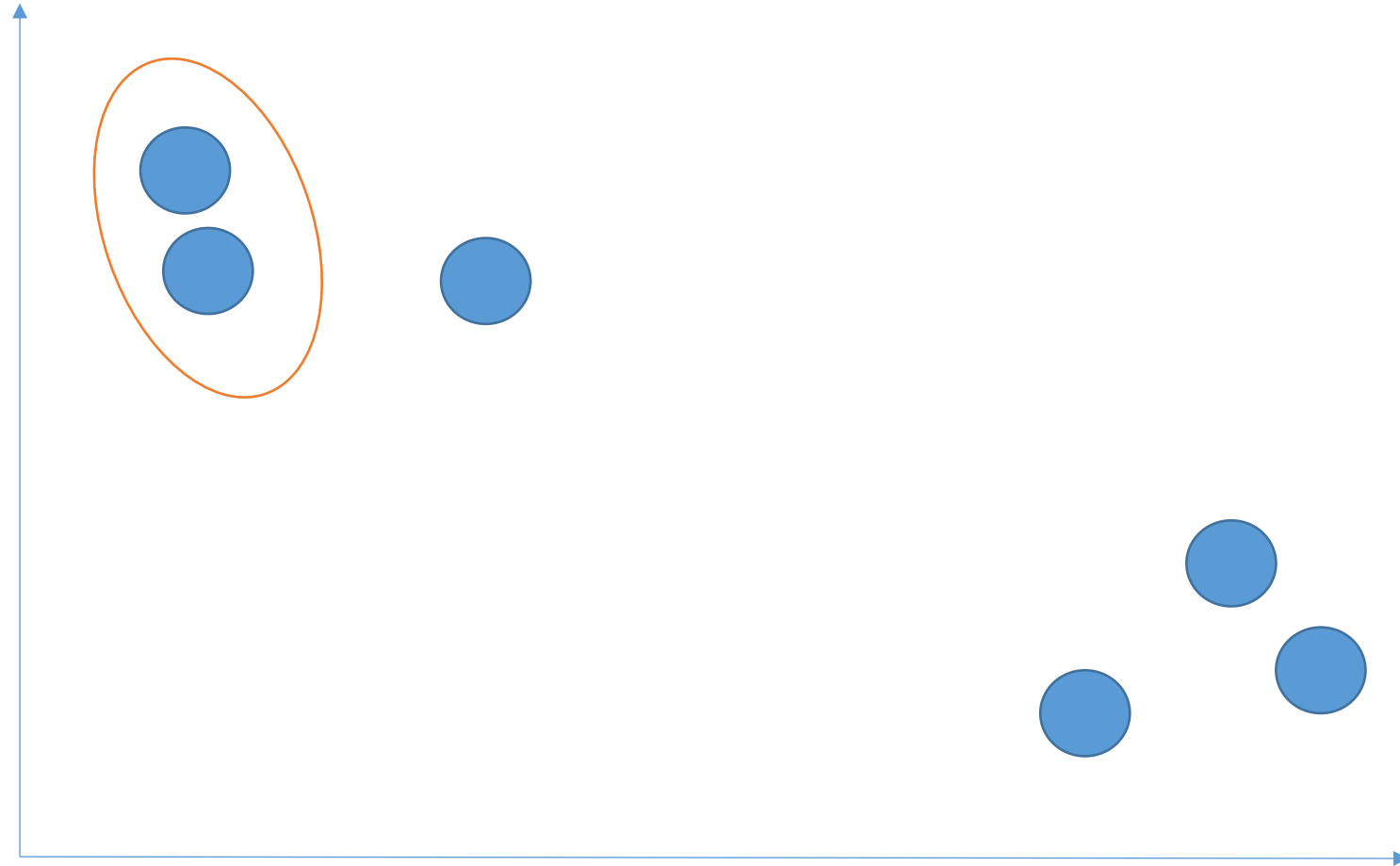
# Fuzzy C-Means

- Permite tener un grado de pertenencia a cada clúster por muestra. El grado de pertenencia es un número entre 0 y 1.
- Tiene problemas similares al K-Means: se debe determinar el número de clúster y tiene problemas en altas dimensionalidades.
- Es muy conveniente cuando no se desea asignar cada muestra a un clúster específico.

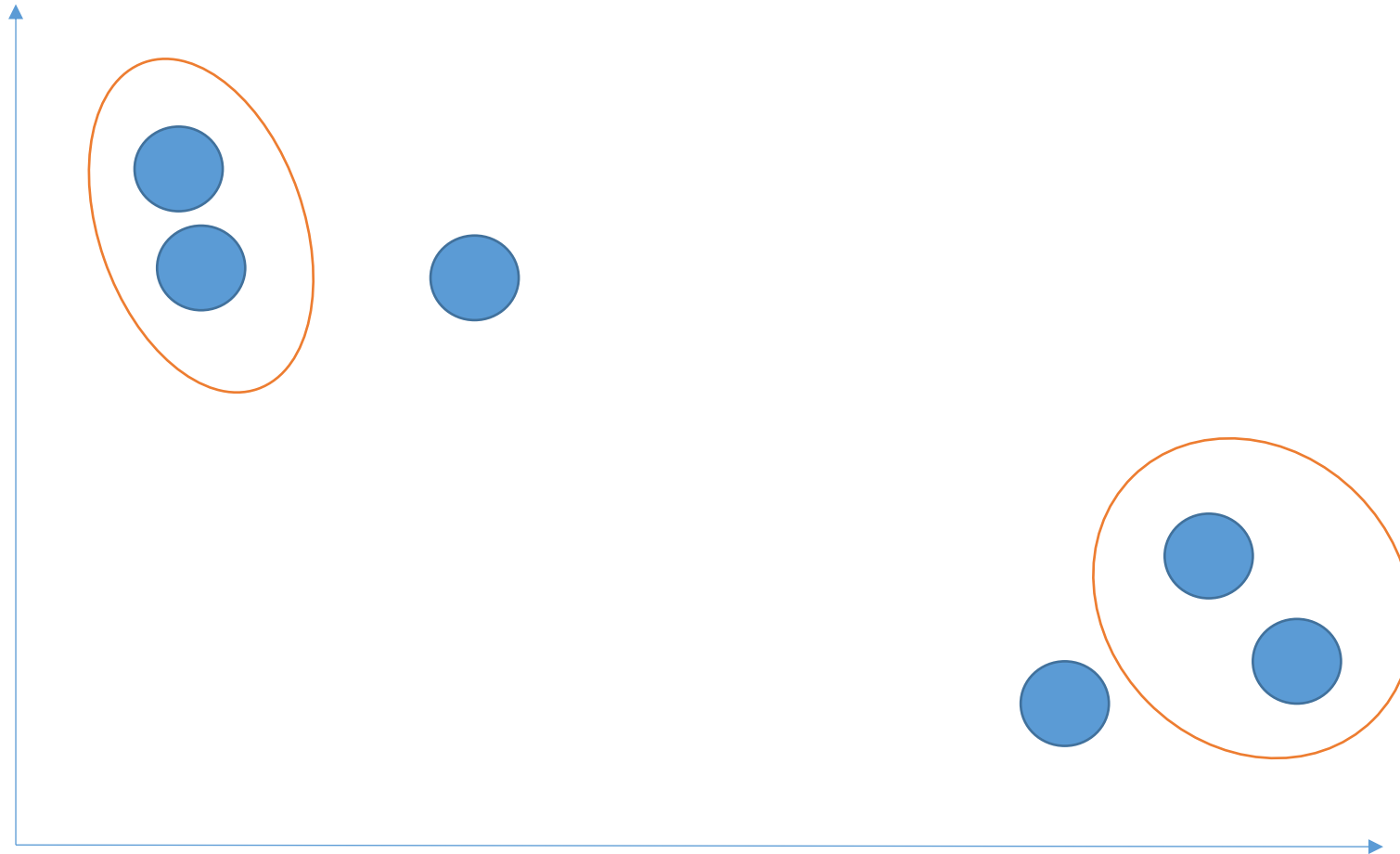
# CLUSTERING JERÁRQUICO



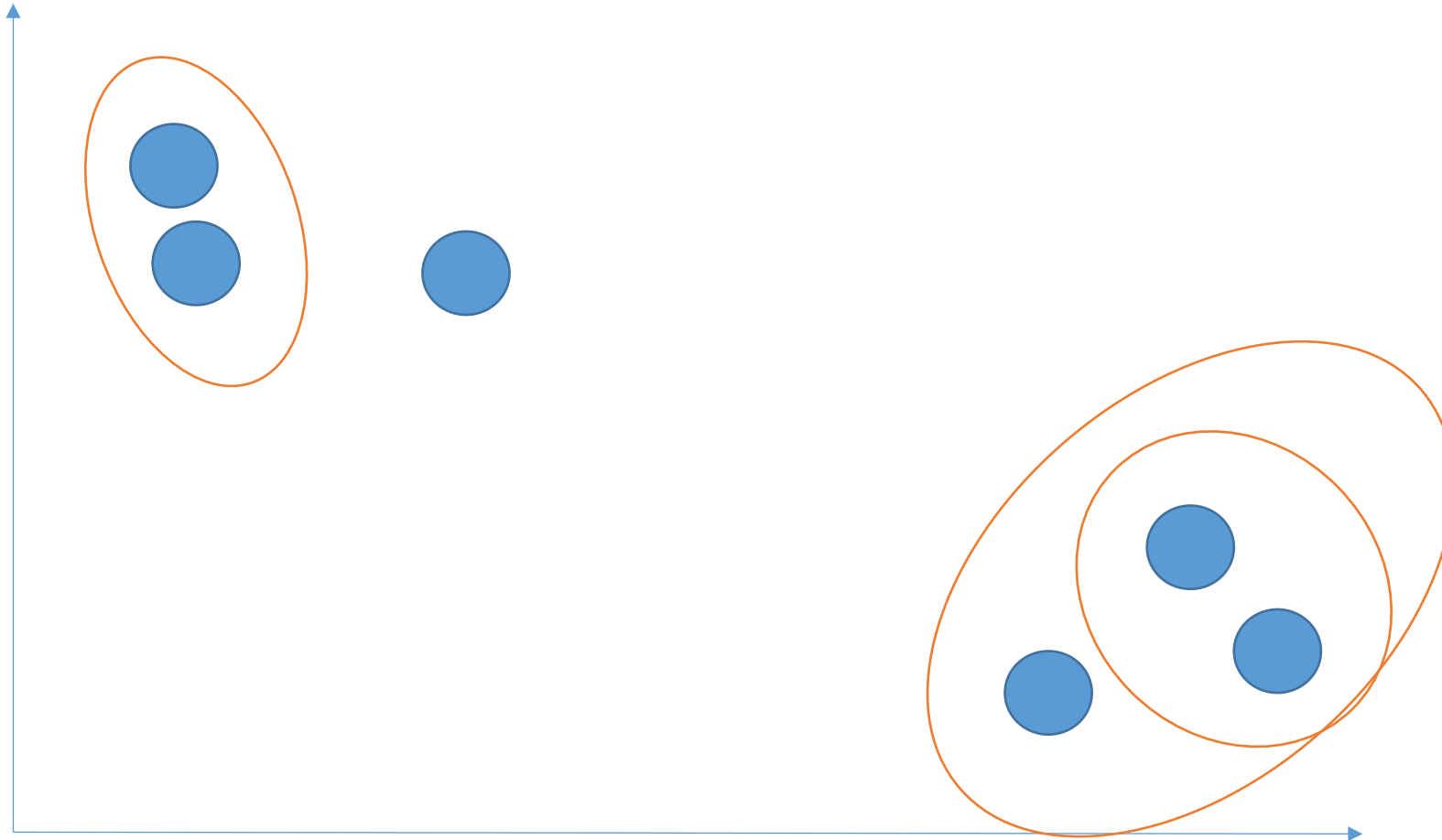
# CLUSTERING JERÁRQUICO



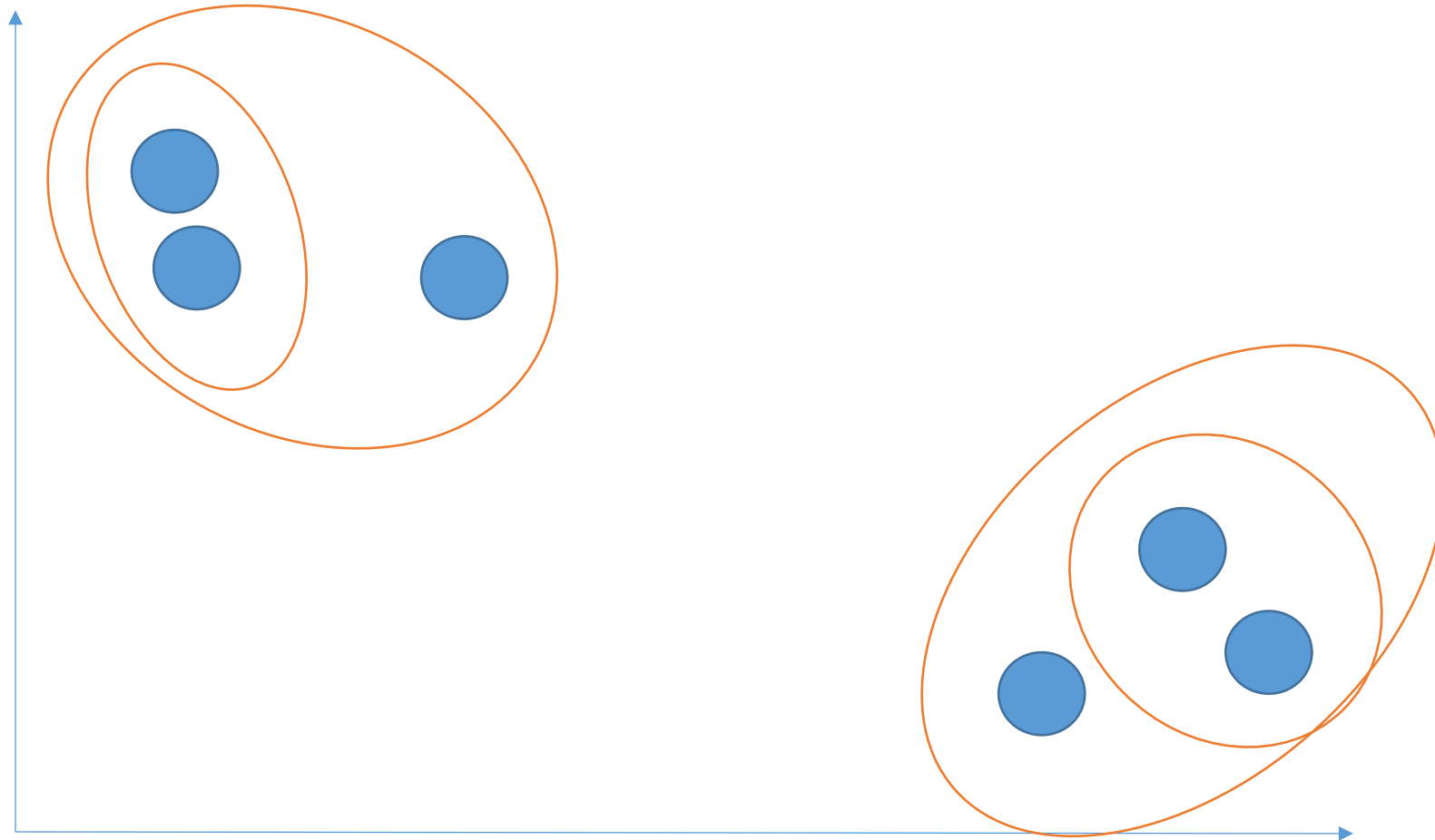
# CLUSTERING JERÁRQUICO



# CLUSTERING JERÁRQUICO



# CLUSTERING JERÁRQUICO





# CLUSTERING JERÁRQUICO

- Agrupa datos de manera consecutiva mientras construye una jerarquía.
- No hay que especificar un número de grupos deseados, pero hay que especificar un punto de corte.

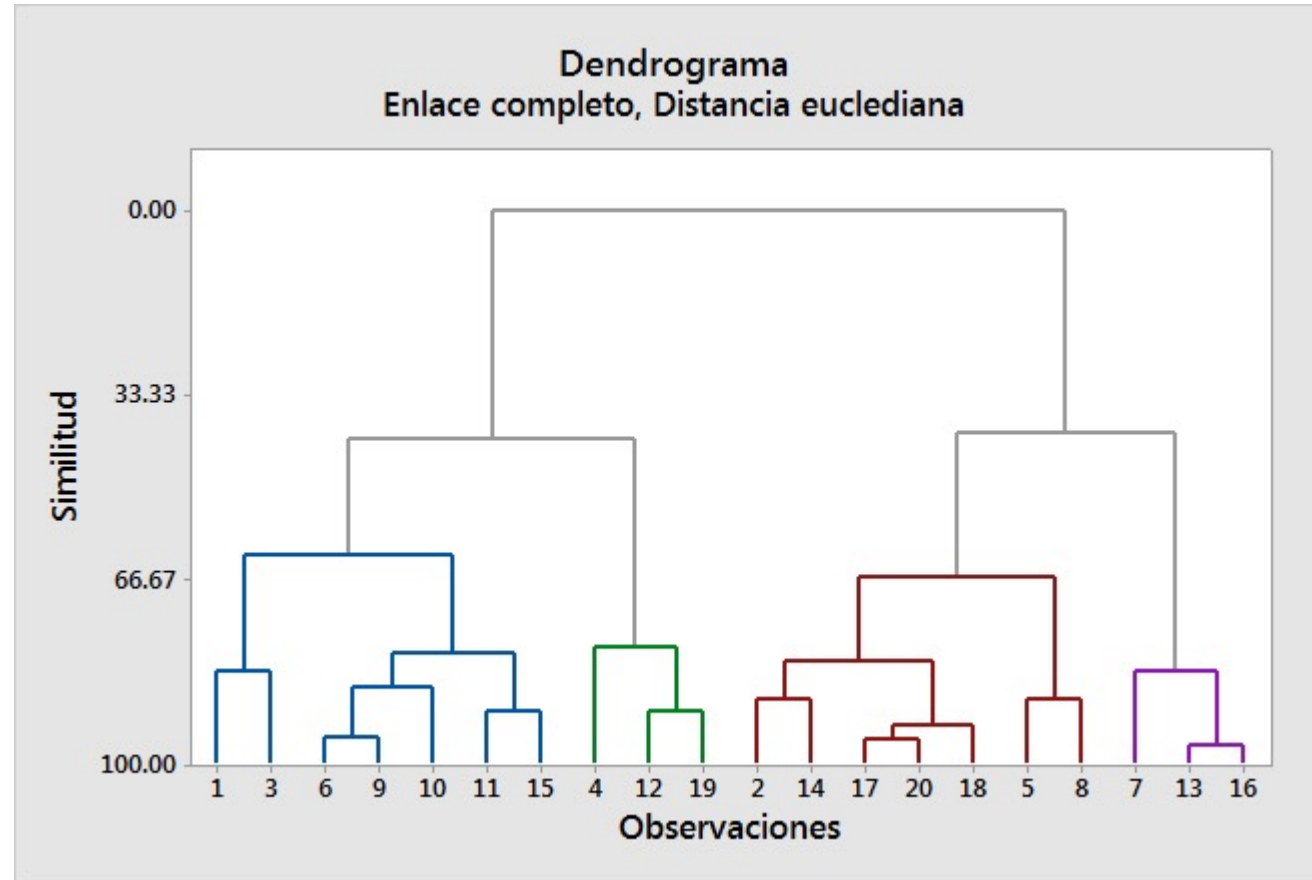
# CLUSTERING JERÁRQUICO

1. Empieza con clusters para cada dato:  $C_n = \{\{x_1\}, \{x_2\}, \dots, \{x_n\}\}$
2. Busca la mínima distancia (u otro criterio) entre los cluster.
3. Se unen los dos clusters más cercanos en un nuevo cluster.

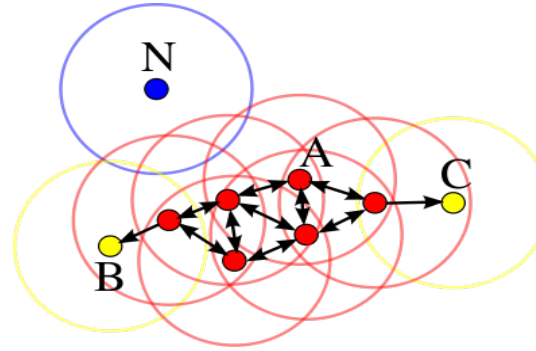
# CRÍTERIOS DE ENLACE

- Distancia mínima (Single Linkage)
- Distancia máxima (Complete Linkage)
- Distancia promedio (average linkage)

# CLUSTERING JERÁRQUICO



# DBSCAN



- No necesita especificar el número de clusters
- Permite encontrar clusters con formas geométricas arbitrarias
- Es robusto detectando outliers

# Índice Silhoutte

- **a** es el promedio de las disimilitudes (o distancias) de la observación  $i$  con las demás observaciones del cluster al que pertenece  $i$
- **b** es la distancia mínima a otro cluster que no es el mismo en el que está la observación  $i$ . Ese cluster es la segunda mejor opción para  $i$  y se lo denomina vecindad de  $i$ .

$$s(\mathbf{o}) = \frac{b(\mathbf{o}) - a(\mathbf{o})}{\max\{a(\mathbf{o}), b(\mathbf{o})\}}$$

- $s(i) \approx 1$ , la observación  $i$  está bien asignada a su *cluster*
- $s(i) \approx 0$ , la observación  $i$  está entre dos *cluster*
- $s(i) \approx -1$ , la observación  $i$  está mal asignada a su *cluster*

# ¡Gracias!

Aliados:



**Microsoft**

Vigilada Mineducación



Advanced analytics for business