

Aprendizaje no Supervisado. Técnicas de Agrupamiento.

Minería de Datos

José T. Palma

Departamento de Ingeniería de la Información y las Comunicaciones
Universidad de Murcia

DIIC, UMU, 2018



Contenidos de la presentación

- 1 Introducción
- 2 Agrupamiento Jerárquico
- 3 Agrupamiento Particional
 - Método de las K-medias
 - K-medoides
 - DBSCAN
- 4 Evaluación de los agrupamientos
- 5 Resumen

Introducción

- En el caso de la clasificación no supervisada no tenemos ninguna información acerca de la organización de los elementos en grupos o clases.
- Por lo tanto, el objetivo consiste en encontrar dicha organización en base a la relación de proximidad de los elementos.
- No existe información previa sobre dicha organización y la interpretación de las clases y los grupos obtenidos hay que hacerla a posteriori.
- Para ello se aplican técnicas de agrupamiento o “clustering”:
 - Identificar distintos grupos de clientes en un banco para personalizar las ofertas de productos financieros.
 - Identificar distintos subgrupos en un tipo determinado de cáncer para ajustar los tratamientos.

Introducción

- La idea básica consiste en crear grupos que contengan elementos parecidos entre si y que elementos dispares se coloquen en grupos diferentes.
- Una técnica de clustering es capaz de describir la estructura subyacente a un conjunto de datos analizando las similitudes y diferencias (p. e., distancias) entre los elementos del conjunto.
- El objetivo final es obtener un conjunto de clases o grupos:
 - Cuando estos grupos son disjuntos y cubren todo el conjunto de elementos se dice que el agrupamiento es “particional”.
 - En algunos casos lo que interesa es una jerarquía de agrupamientos particionales anidados. En este caso tenemos un agrupamiento jerárquico que se suele representar mediante un dendograma.

Distancia y similaridad

- El concepto de similaridad y distancia es clave en las técnicas de agrupamiento ya que definen la lente que le vamos a dar al ordenador para que busque la estructura de los datos.
- Supongamos que tenemos m elementos recogidos en un conjunto $\Omega = \{1, 2, 3, \dots, m\}$

Definición (Distancia)

Una medida de distancia sobre el conjunto Ω es una función d tal que:

$$\begin{aligned} d : \Omega \times \Omega &\rightarrow \mathbb{R} \\ (i, j) &\rightarrow d(i, j) = d_{ij} \end{aligned}$$

Introducción

Distancia: propiedades

- 1 $d(i, j) \geq 0, \forall i, j \in \Omega$
- 2 $d(i, i) = 0, \forall i \in \Omega$
- 3 $d(i, j) = d(j, i) \forall i, j \in \Omega$

- Cuando además se cumple la propiedad de desigualdad triangular:

$$d(i, j) \leq d(i, k) + d(k, j) \forall i, j, k \in \Omega$$

diremos que la distancia es métrica y (Ω, d) forma un espacio métrico.

Distancia y similaridad

Definición (Similaridad)

Una medida de similaridad sobre el conjunto Ω es una función s tal que:

$$\begin{aligned}d : \Omega \times \Omega &\rightarrow \mathbb{R} \\(i, j) &\rightarrow s(i, j) = s_{ij}\end{aligned}$$

tal que:

- ❶ $s(i, j) \in [0, 1] \forall i, j \in \Omega$
- ❷ $1 = s(i, i) \geq s(i, j), \forall i, j \in \Omega$
- ❸ $s(i, j) = s(j, i) \forall i, j \in \Omega$

Distancia y similaridad

- Obviamente los conceptos de distancia y similaridad están relacionados: a mayor distancia menor similaridad.
- Existen diferentes formas para relacionar ambas medidas:
 - Transformación de Gower: $d_{ij}^2 = s_{ii} + s_{jj} - 2s_{ij}$
 - Distancia complemento: $d_{ij} = 1 - s_{ij}$
 - Raíz del complemento del cuadrado: $d_{ij} = \sqrt{1 - s_{ij}^2}$

Distancia y similaridad

- Generalmente cada elemento del conjunto Ω tendrá asociada una variable y podrá ser representado mediante el punto $x = \{x_1, x_2, \dots, x_n\}$ en el espacio \mathbb{R} .
- Dependiendo de la naturaleza de las variables, se deberán utilizar diferentes tipos de distancias y similaridades.
- A continuación comentaremos las más habituales.

Distancia para variables continuas I

- Sean $x = \{x_1, x_2, \dots, x_n\}$ e $y = \{y_1, y_2, \dots, y_n\}$ dos elementos del conjunto Ω , en el que todas las variables son continuas, las medidas de distancia más utilizadas son:

| Función de Distancia | Fórmula |
|----------------------|---|
| Euclidea | $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$ |
| Manhattan | $d(x, y) = \sum_{i=1}^n x_i - y_i $ |
| Norma del supremo | $d(x, y) = \sup_{i \in \{1, 2, \dots, n\}} x_i - y_i $ |
| Minkosky | $d(x, y) = \sqrt[p]{\sum_{i=1}^n (x_i - y_i)^p} \quad p > 0$ |
| Mahanalobis | $d(x, y) = [(x - y)^T \Sigma^{-1} (x - y)]^{1/2}$ Σ Matriz de covarianzas |

Distancia para variables continuas II

- Para evitar que unas variables dominen sobre otras, las variables continuas se suelen normalizar.
- Para la normalización se suele utilizar el z-score.
- Sea v_{ij} una el valor de la variable j en el objeto i , el valor normalizado z_{ij} se calcula de la siguiente forma:

$$z_{ij} = \frac{v_{ij} - \bar{v}_j}{s_j}$$

con

$$\bar{v}_j = \frac{1}{n} \sum_{i=1}^n v_{ij} \quad \text{y} \quad s_j = \sqrt{\frac{\sum_{i=1}^n (v_{ij} - \bar{v}_j)^2}{n - 1}}$$

Similaridad para variables binarias

- Sean $x = \{x_1, x_2, \dots, x_n\}$ e $y = \{y_1, y_2, \dots, y_n\}$ dos elementos del conjunto Ω , en el que todas las variables son binarias
- En este caso es más fácil calcular primero la similitud, para después transformarla en distancia.
- Para ello se calcula la matriz de confusión para calcular las coincidencias entre las n variables:

| | | x_i | | |
|-------|---|-------|-----|-----|
| | | 1 | 0 | |
| y_i | 1 | a | b | a+b |
| | 0 | c | d | c+d |
| | | a+c | b+d | n |

Similaridad para variables binarias

| Función de Distancia | Fórmula |
|----------------------|-----------------------------------|
| Índice de Acuerdo | $s(x, y) = \frac{a + d}{n}$ |
| Jaccard | $s(x, y) = \frac{a}{a + b + c}$ |
| Russel-Roo | $s(x, y) = \frac{a}{n}$ |
| Czekanowski | $s(x, y) = \frac{2a}{2a + b + c}$ |

Otro tipo de variables

- Para el caso de variables cualitativas o nominales existen dos posibilidades:
 - Contando las coincidencias:

$$d(x, y) = \frac{p - m}{p}$$

siendo p el número total de variables y m el número de coincidencias.

- Creando un atributo binario para cada uno de los posibles valores y calcular la similaridad como se describió anteriormente.

Otro tipo de variables

- Para el caso de variables ordinales (en este caso el orden si es importante) estas se tratan como numéricas después calcular su correspondencia al intervalo $[0, 1]$:

$$z_{i,k} = \frac{r_{i,k} - 1}{M_k - 1}$$

siendo :

- $z_{i,n}$ el valor estandarizado para el el objeto i de la variable x_k ,
- $r_{i,k}$ el valor antes de la transformación y
- M_k el límite superior del dominio de la variable x_k (se asume que el límite inferior es 1).

Similaridad para variables mixtas

- Supongamos que las n variables son mixtas: n_1 variables son cuantitativas, n_2 son binarias y n_3 son cualitativas.
- En este caso se puede utilizar la distancia de Gower que se calcula con $d(x, y) = \sqrt{1 - s(x, y)}$ con

$$s(x, y) = \frac{\sum_{l=1}^{n_1} (1 - |x_l - y_l|/R_l) + a + \alpha}{n_1 + (n_2 - d) + n_3}$$

- donde:
 - x_1, \dots, x_{n_1} e y_1, \dots, y_{n_1} representan los valores observados en las variables cuantitativas.
 - R_l es el rango de la l -ésima variable cuantitativa.
 - a y d son los valores de la matriz de confusión para las coincidencias de en las variables binarias.
 - α el número de coincidencias entre las variables cualitativas.

Similaridad para variables mixtas

- También se puede utilizar cualquier medida de agregación de las distancias/similaridades de las variables independientes.

$$d(x, y) = \sum_{l=1}^n \omega_l d(x_l, y_l) \quad \text{con} \quad \sum_{l=1}^n \omega_l = 1$$

- donde:
 - $d(x_l, y_l)$ se corresponden con las distancias de cada una de las variables.
 - ω_l es el peso asociado a cada una de las variables y se tiene que cumplir .

Agrupamiento Jerárquico

- Un agrupamiento jerárquico es una sucesión de particiones “anidadas”:
 - Cada grupo de elementos pertenecientes a una partición está totalmente incluido en alguna partición de nivel superior.
 - Esta estructura tiene una representación gráfica muy intuitiva denominada “dendograma”.
 - El dendograma representa cómo se van uniendo los distintos elementos en grupos.
 - Es un árbol binario en el que los elementos individuales se encuentran en los nodos hojas y los nodos intermedios representan diferentes agrupaciones de elementos.

Agrupamiento Jerárquico

- Existen dos tipos de técnicas para construir un dendrograma:
 - Las técnicas **aglomerativas** generan nuevos clusters uniendo clusters similares.
 - Se parte de una partición inicial en la que cada elemento forma un cluster.
 - Se van uniendo de dos en dos aquellos clusters que están más próximos.
 - Finaliza el proceso cuando todos los elementos están ubicados en un único cluster
 - Las técnicas **divisivas** los nuevos clusters se generan dividiendo clusters.
 - Se parte de un único cluster que contiene todos los elementos.
 - Se va dividiendo dicho cluster hasta alcanzar una partición en la que todos los clusters contiene un único elemento.

Agrupamiento Jerárquico

- Las técnicas aglomerativas son más eficientes que las divisivas.
- Las técnicas divisivas tienen la ventaja de que parten de la información global que hay en los datos y no tienen porque llegar hasta los clusters de tamaño 1.
- Sin embargo, las técnicas divisivas son muy lentas y, sólo se utilizan en el caso de que existan pocos datos.
- Esto hace que los métodos más utilizados sean los aglomerativos, que son los que vamos a analizar a continuación.

Agrupamiento Jerárquico

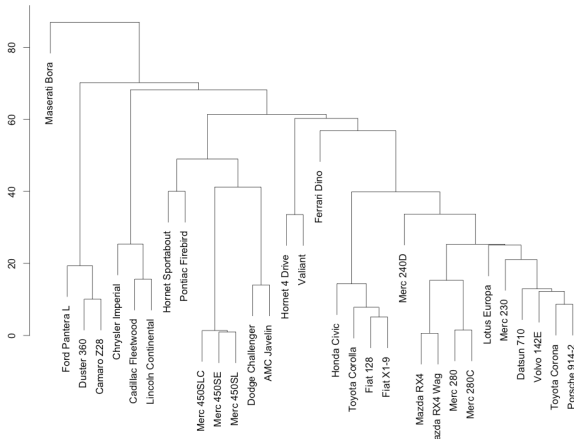


Figura: Dendrograma

Agrupamiento Jerárquico: Ejemplo

| | E_1 | E_2 | E_3 | E_4 | E_5 | E_6 | E_7 |
|-------|-------|-------|-------|-------|-------|----------|-------|
| E_1 | 0 | 10 | 10 | 7 | 6 | 13 | 8 |
| E_2 | | 0 | 4 | 7 | 8 | 2 | 8 |
| E_3 | | | 0 | 9 | 12 | 3 | 8 |
| E_4 | | | | 0 | 5 | 5 | 5 |
| E_5 | | | | | 0 | 6 | 6 |
| E_6 | | | | | | 0 | 9 |
| E_7 | | | | | | | 0 |

Agrupamiento Jerárquico: Ejemplo

| | E_1 | E_3 | E_4 | E_5 | E_7 | (E_2, E_6) |
|--------------|-------|-------|-------|-------|-------|--------------|
| E_1 | 0 | 10 | 7 | 6 | 8 | 10 |
| E_3 | | 0 | 9 | 12 | 8 | 3 |
| E_4 | | | 0 | 5 | 5 | 5 |
| E_5 | | | | 0 | 6 | 6 |
| E_7 | | | | | 0 | 8 |
| (E_2, E_6) | | | | | | 0 |

Agrupamiento Jerárquico: Ejemplo

| | E_1 | E_4 | E_5 | E_7 | $((E_2, E_6), E_3)$ |
|---------------------|-------|-------|----------|-------|---------------------|
| E_1 | 0 | 7 | 6 | 8 | 10 |
| E_4 | | 0 | 5 | 5 | 5 |
| E_5 | | | 0 | 6 | 6 |
| E_7 | | | | 0 | 8 |
| $((E_2, E_6), E_3)$ | | | | | 0 |

Agrupamiento Jerárquico: Ejemplo

| | E_1 | (E_4, E_5) | E_7 | $((E_2, E_6), E_3)$ |
|---------------------|-------|--------------|----------|---------------------|
| E_1 | 0 | 6 | 8 | 10 |
| (E_4, E_5) | | 0 | 5 | 5 |
| E_7 | | | 0 | 8 |
| $((E_2, E_6), E_3)$ | | | | 0 |

Agrupamiento Jerárquico: Ejemplo

| | E_1 | $((E_4, E_5), E_7)$ | $((E_2, E_6), E_3)$ |
|---------------------|-------|---------------------|---------------------|
| E_1 | 0 | 6 | 10 |
| $((E_4, E_5), E_7)$ | | 0 | 5 |
| $((E_2, E_6), E_3)$ | | | 0 |

Agrupamiento Jerárquico: Ejemplo

$$\begin{array}{rcl}
 & & E_1 \quad ((E_4, E_5), E_7), ((E_2, E_6), E_3)) \\
 & & 0 \qquad \qquad \qquad 6 \\
 E_1 & & \\
 ((E_4, E_5), E_7), ((E_2, E_6), E_3)) & & 0
 \end{array}$$

Agrupamiento Jerárquico: Ejemplo

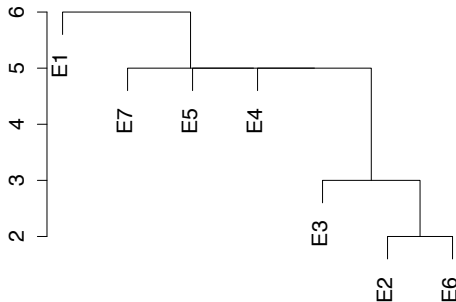


Figura: Dendrograma resultante del ejemplo.

Agrupamiento Jerárquico: Distancia entre clusters I

- Como se puede ver, la clave de un agrupamiento jerárquico reside en la forma en que se defina la distancia entre dos clusters.
- Sean A y B dos clusters, la distancia entre ambos, $d(A, B)$ se puede definir de diferentes formas:
 - **Método de enlace simple** (single link). En este caso la distancia $d(A, B)$ se calcula como la distancia mínima entre los elementos de ambos clusters:

$$d(A, B) = \min_{x \in A, y \in B} d(x, y)$$

Agrupamiento Jerárquico: Distancia entre clusters II

- **Método del enlace completo** (complete link). En este caso la distancia $d(A, B)$ se calcula como la distancia máxima entre los elementos de ambos clusters:

$$d(A, B) = \max_{x \in A, y \in B} d(x, y)$$

- **Método del enlace promedio** (average link, o UPGMA). En este caso la distancia $d(A, B)$ se calcula como el promedio de la distancia entre cada par de elementos de ambos clusters:

$$d(A, B) = \frac{1}{|A||B|} \sum_{x \in A, y \in B} d(x, y)$$

Agrupamiento Jerárquico: Distancia entre clusters III

- **El método del centroide** (UPGMC). En este caso la distancia $d(A, B)$ se calcula como la distancia entre los centroides de cada grupo. El centroide del cluster A se calcula de la siguiente forma

$$\bar{x} = \frac{1}{|A|} \sum_{x \in A} x$$

- **El método del Ward**. En este caso se trata de fusionar aquellos clusters de tal forma que en el nuevo cluster la suma de las distancias de los elementos al centroide sea menor.

Agrupamiento Jerárquico: Distancia entre clusters IV

- Existen versiones ponderadas de los métodos **promedio** y **centroide** (WPGMA y WPGMC respectivamente) que intentan compensar el hecho de fusionar cluster de tamaños muy dispares. Estos métodos se deberían utilizar cuando se sospeche de que el tamaño de los distintos clusters va a ser muy dispares.

Agrupamiento Jerárquico: Distancia entre clusters I

- En [Lance and Williams, 1967] y [Jain and Dubes, 1988] se puede encontrar una expresión genérica para todos los tipos de distancias analizadas.
- Supongamos vamos a fusionar los clusters C_i y C_j . Para poder seguir con el proceso primero debemos de calcular la distancia del nuevo cluster al resto de clusters C_k :

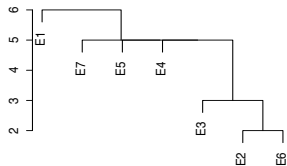
$$d(C_i \cup C_j, C_k) = \alpha_i d(C_i, C_k) + \alpha_j d(C_j, C_k) + \beta d(C_i, C_j) + \gamma |d(C_i, C_k) - d(C_j, C_k)|$$

donde:

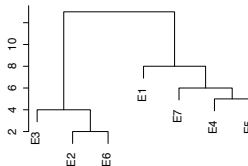
Agrupamiento Jerárquico: Distancia entre clusters II

| Método | α_i | α_j | β | γ |
|------------------|-------------------------------------|-------------------------------------|----------------------------------|----------|
| Simple | $1/2$ | $1/2$ | 0 | $-1/2$ |
| Completo | $1/2$ | $1/2$ | 0 | $1/2$ |
| Promedio (UPGMA) | $\frac{n_i}{n_i + n_j}$ | $\frac{n_j}{n_i + n_j}$ | 0 | 0 |
| WPGMA | $1/2$ | $1/2$ | 0 | 0 |
| Centroide(UPGMC) | $\frac{n_i}{n_i + n_j}$ | $\frac{n_j}{n_i + n_j}$ | $\frac{-n_i n_j}{(n_i + n_j)^2}$ | 0 |
| WPGMC | $1/2$ | $1/2$ | $-1/4$ | 0 |
| Ward | $\frac{n_i + n_k}{n_i + n_j + n_k}$ | $\frac{n_j + n_k}{n_i + n_j + n_k}$ | $\frac{-n_k}{n_i + n_j + n_k}$ | 0 |

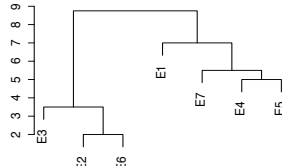
Agrupamiento Jerárquico: Ejemplo



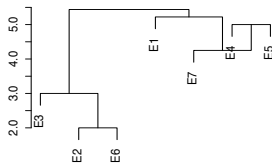
Enlace simple



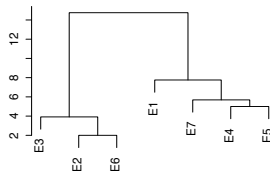
Enlace completo



Enlace Promedio



Centroides



Ward

Agrupamiento particional

- El objetivo de un agrupamiento particional es, dado un conjunto n elementos representados en un espacio d —dimensional:
 - encontrar una partición del mismo en k subconjuntos.
 - los elementos dentro de un grupo se tiene que parecer más entre sí que a los elementos de otros grupos
- El número k de subgrupos puede ser conocido apriori o no,
 - En la mayoría de las técnicas ese dato es un parámetro.

Agrupamiento particional

- Es necesario un criterio para medir la coherencia de cada grupo, así como la de entre grupos.
- Existen dos tipos de criterios:
 - Los métodos basados en **criterios globales** representan cada grupo mediante un prototipo, asignando cada elemento al grupo del prototipo más cercano. **K-medias** y **K-medoides** son técnicas que se corresponden con este tipo de criterio.
 - Los métodos basados en **criterios locales** forman grupos utilizando la estructura local de los datos, por ejemplo, identificando regiones de alta densidad de puntos. Uno de los métodos más conocidos dentro de este enfoque es **DBSCAN**.

Agrupamiento en base a encoders

- Supongamos que tenemos un conjunto de n elementos $X = \{x_1, x_2, x_3, \dots, x_n\}$.
- El resultado de aplicar una técnica de agrupamiento que define k clusters ($k < N$) se puede definir mediante un encoder C :

Encoder

$$C(i) = k \Leftrightarrow x_i \in k$$

es decir, el encoder C nos indica a qué cluster pertenece cada elemento.

- Por lo tanto, el objetivo de una técnica de agrupamiento debe ser encontrar el encoder $C^*(i)$ que optimice algún criterio determinado.

Distancia intra e intercluster I

- Sea C un encoder de k clusters sobre el conjunto $X = \{x_1, x_2, x_3, \dots, x_n\}$.
- Hay que tener en cuenta que la separación total entre los puntos en X siempre es la misma:

$$T = \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N d_{ii'} = \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \left(\sum_{i': C(i')=k} d_{ii'} + \sum_{i': C(i') \neq k} d_{ii'} \right)$$

con lo que en realidad tenemos $T = W(C) + B(C)$

Distancia intra e intercluster II

- Es decir, la distancia total en entre los puntos de un conjunto dividido en k clusters se puede calcular mediante la suma de la distancia entre los puntos de distintos clusters (intercluster, $W(C)$) y la distancia entre los puntos de cada cluster (intracluster, $B(C)$).

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \sum_{i': C(i') \neq k} d(x_i, x_{i'})$$

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \sum_{i': C(i')=k} d(x_i, x_{i'})$$

Distancia intra e intercluster III

- Por lo tanto, para obtener un buen agrupamiento (encoder) podemos:
 - Maximizar la distancia intercluster $B(C)$, buscar clusters lo más separados entre sí.
 - Minimizar la distancia intracluster $W(C)$, buscar clusters lo más compactos posibles.
- Ambas opciones son equivalentes ya que $W(C) = T - B(C)$

Método K-medias

- El algoritmo iterativo **K-medias** es el más popular que se puede aplicar a variables numéricas y la distancia euclídea:

$$d(x_i, x_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = \|x_i - x_{i'}\|^2$$

- Por tanto, la distancia intracluster puede escribirse como

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i: C(i)=k} \sum_{i': C(i')=k} \|x_i, x_{i'}\|^2 = \sum_{k=1}^K \sum_{i: C(i)=k} \|x_i - \bar{x}_k\|^2$$

donde \bar{x}_k es el centroide

- Como podemos ver, en este caso, la distancia intracluster coincide con la suma del cuadrado de los errores (SSE).

Método K-medias

- Por lo tanto, ahora nuestro problema de optimización es

$$C^* = \min_C \sum_{k=1}^K \sum_{i: C(i)=k} ||x_i - \bar{x}_k||^2$$

- El algoritmo K-means es un algoritmo que intenta resolver este problema siguiendo un esquema de ascensión de colinas por la máxima pendiente.
 - Después de cada iteración no se puede volver atrás y probar otros centroides.

Algoritmo K-medias

- 1 Comenzar con alguna de las dos configuraciones iniciales:
 - Si se inicializan aleatoriamente los representantes, m_i de cada uno de los k clusters ir al paso 2.
 - Si se parte de una partición aleatoria del conjunto de entrada en k grupos ir al paso 3.

- 2 Calcular el encoder (distribuir los elementos entre los k clusters de acuerdo a los representantes m_i).

$$C(i) = \arg \min_{1 \leq k \leq K} \|x_i - m_k\|^2$$

- 3 Calcular los nuevos $m_k, i = 1, \dots, K$ como el centroide de todos los puntos $x \in X$ tales que $C(x) = k$.
- 4 Si los nuevos $m_k, i = 1, \dots, K$ no se han estabilizado, volver al paso 4. Si no, fin

Método K-medias: Ejemplo

Ejemplo K-medias con 3 clusters

Método K-medias: Ejemplo

Ejemplo K-medias con 4 clusters

Método K-medias: Consideraciones

- El algoritmo K-medias se aplica en el caso de que todos los atributos sean reales.
- En principio hemos utilizado la distancia euclídea:
 - Los representantes se corresponden la media aritmética de los elementos del cluster.
 - Lo que se está minimizando es la desviación respecto a los representantes de cada cluster. Es decir, la distancia intracluster.
 - Sin embargo, amplifica el efecto de los outliers.
- Se pueden considerar otras medidas de distancia.
- También se pueden utilizar medidas de similitud con lo que en vez de minimizar habría que maximizar.

Método K-medias: Problemas I

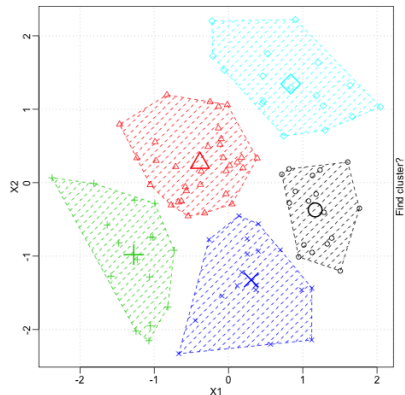
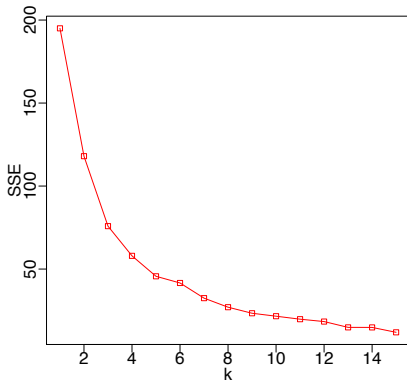
- Es muy sensible a la elección de los centroides
 - Se pueden realizar varias ejecuciones con diferentes centroides iniciales y comparar resultados.
 - El método más simple es escoger, de entre las N instancias del conjunto de observaciones, k instancias para inicializar los m_i
 - Una posibilidad algo más elaborada es usar el vector de medias \bar{x} de todo el cjto. de datos. Para obtener cada m_i sumamos o restamos valores aleatorios a cada componente de \bar{x}
 - Hacer un análisis PCA (Principal Component Analysis), dividir el rango de la primera variable generada por PCA en k intervalos iguales y generar los centroides a partir de la media de dichos intervalos
 - Usar clustering jerárquico

Método K-medias: Problemas II

- K-medias es muy sensible al número de clusters, k , y hay que elegirlos a priori.
 - Se puede usar un método jerárquico para estimar el valor de k .
 - Se puede aplicar K-medias para distintos valores de k y comprobar cuando no hay mejoras significativas del SEE.

Método K-medias: Problemas III

- En nuestro ejemplo, parece que el corte puede ser con $k = 5$



Método K-medias: Problemas IV

- Al usar la media para calcular los centroides el método es sensible a los outliers
 - Utilizar las medianas.
 - Eliminar los outliers (en algunos casos pueden ser de interés).
 - Utilizar K-medoides: el representante tiene que ser el elemento más representativo del cluster.
- Para manejar datos no numéricos se requiere la redefinición de la función de distancia, trabajar con la moda y no la media.
- K-medias no funciona bien cuando los clusters son de: distinto tamaño, diferente densidad y no convexos.
 - Esto requiere una revisión posterior de los resultados y hacer varias pruebas para distintos valores de K.

Método K-medoides

- Para evitar la sensibilidad del método K-medias a los outliers, K-medoides elige como representante de cada cluster a un punto del mismo considerado más representativo, la **mediana**.
- Al no basarse en los centroides (valores medios), no hace falta la definición de una función de distancia ya que puede operar directamente con la matriz de distancias (o similaridad).
- El proceso es idéntico al método K-medias sólo que el cálculo de los centroides se sustituye por el de los medoides.
- Es más costoso computacionalmente que el método K-medias, además de necesitar también saber a priori el número de grupos

Algoritmo K-medoides

- 1 Comenzar con alguna de las dos configuraciones iniciales:
 - Si escogen aleatoriamente k elementos como representantes m_i de los k clusters ir al paso 2.
 - Si se parte de una partición aleatoria del conjunto de entrada en k grupos ir al paso 3.

- 2 Calcular el encoder (distribuir los elementos entre los k clusters de acuerdo a los representantes m_i).

$$C(i) = \arg \min_{1 \leq k \leq K} D(x_i, m_k)$$

- 3 Calcular los nuevos $m_k, i = 1, \dots, K$ como medoides de cada cluster:

$$m_k = \arg \min_{i: C(i)=k} \sum_{i': C(i')=k} D(x_i, x_{i'})$$

- 4 Si los nuevos $m_k, i = 1, \dots, K$ no se han estabilizado, volver al paso 2. Si no, fin

Método K-medoides I

- Al utilizar las medianas en vez valores medios, se seleccionan como representantes elementos del conjunto de datos:
 - Esto permite que el método sea más **robusto** y no se ve tan afectado por la presencia de outliers
 - También se gana en **interpretabilidad**.
- El problema que plantea es el alto coste computacional.
 - Funciona bien para conjunto pequeños de datos y pocos clusters (100 elementos y 5 clusters aproximadamente).

Método K-medoides II

- Existen varias implementaciones:
 - **PAM** (Partition Around Medoids) consiste en la implementación de las ideas anteriormente propuestas.
 - **CLARA** (Clustering LARge Applications), intenta reducir la carga computacional de PAM seleccionando los medoides de una muestra aleatoria y significativa de los datos y después aplica PAM. Básicamente realiza varios muestreos y da como resultado el mejor clustering.
 - **CLARANS** que se diferencia del anterior en que la búsqueda de los medoides se aproxima como un proceso de búsqueda, realizando un muestreo cada vez que se calcula un medoide.

Método DBSCAN

- El método **DBSCAN** (Density-based spatial clustering of applications with noise) es un método basado en criterios locales que se apoya en el concepto de densidad de los puntos.
 - Se consideran clusters aquellas regiones del espacio con una alta densidad de puntos.
 - Las regiones con una baja densidad se podrán corresponder con puntos que no están asociados a la mayoría de los datos.

Método DBSCAN: Conceptos previos I

- Para describir el algoritmo necesitamos definir los siguientes conceptos:

ϵ -vecindad

Sean x_k un elemento del conjunto de datos $X = \{x_1, \dots, x_n\}$ y $\epsilon > 0$ con $\epsilon \in \mathbb{R}$, la ϵ -vecindad del punto x_k , $N_\epsilon(x_k)$, se define como:

$$N_\epsilon(x_k) = \{x \in X \mid d(x, x_k) \leq \epsilon\}$$

es decir, la ϵ -vecindad de un punto incluye todos aquellos puntos que están a una distancia menor o igual que ϵ .

- La geometría de la ϵ -vecindad vendrá determinada por la medida de distancia utilizada.

Método DBSCAN: Conceptos previos II

- Para determinar cuando una ϵ -vecindad tiene una densidad alta se define el parámetro *MinPts*, de tal forma que si:

$$|N_{\epsilon}(x_k)| \geq MinPts$$

diremos que la ϵ -vecindad del punto x_k es alta y x_k es un **punto núcleo**.

- Todo aquél punto dentro de la ϵ -vecindad de un punto núcleo se denomina **punto frontera**. Un punto frontera puede pertenecer a dos ϵ -vecindades distintas.
- El resto de puntos se consideran **puntos ruido**.

Método DBSCAN: Conceptos previos III

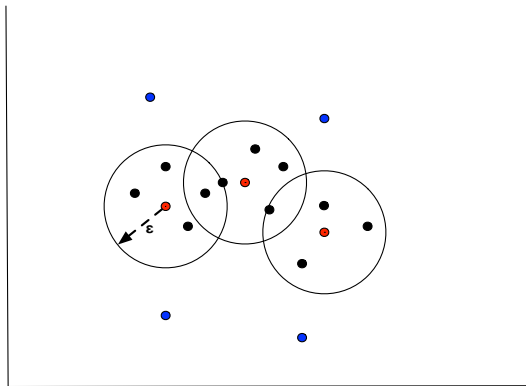


Figura: Concepto de ϵ — vecindad y puntos núcleo (rojo), puntos frontera (negro) y puntos ruido (azul).

Método DBSCAN: Conceptos previos IV

Densidad alcanzable directa

Sean p y q dos elementos del conjunto X , decimos que q es directamente densidad alcanzable desde p si:

- 1 $q \in N_{\epsilon}(p)$
- 2 p es un punto núcleo.

Método DBSCAN: Conceptos previos V

Densidad alcanzable

Sean p y q dos elementos del conjunto X , decimos que p es densidad alcanzable desde q , si existe una cadena de puntos p_1, p_2, \dots, p_n tal que:

- 1 $p_1 = q$ y $p_n = p$
- 2 $\forall i \in \{1, \dots, n\}$ p_{i+1} es directamente densidad alcanzable desde p_i

Es transitiva pero no simétrica

Método DBSCAN: Conceptos previos VI

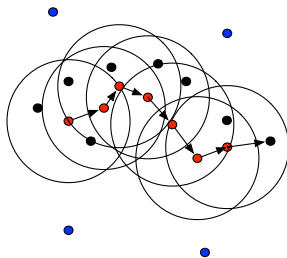


Figura: Concepto de densidad alcanzable

Método DBSCAN: Conceptos previos VII

Puntos densamente conectados

Sean p y q dos elementos del conjunto X , decimos que p y q están densamente conectados si:

- Si son directamente alcanzables desde un mismo punto o

La conectividad densa es simétrica.

Método DBSCAN: Conceptos previos VIII

Cluster

Sean un conjunto X y los parámetros ϵ y $MinPts$, un cluster C es un subconjunto de X que satisface los siguientes criterios

- 1 **Maximalidad:** $\forall p, q$ si $p \in C$ y q es densidad alcanzable desde p , entonces $q \in C$
- 2 **Conectividad:**
 $\forall p, q \in C$, p y q están densamente conectados

Algoritmo DBSCAN I

- La idea básica del método DBSCAN consiste en crear clusters con todos los puntos que son densidad alcanzable.
 - 1 Se especifican los parámetros ϵ y $MinPts$.
 - 2 Seleccionar arbitrariamente un punto, x_k .
 - 3 Encontrar todos aquellos puntos densidad alcanzables desde x_k .
 - 4 Si x_k es un punto núcleo se forma un cluster y se intenta expandir añadiendo todos los puntos densidad alcanzable a otros puntos núcleos e incluyendo también los puntos en su ϵ -vecindad.
 - 5 Si x_k es un punto frontera se procede con el siguiente punto.
 - 6 En otro caso el punto se etiqueta como ruido y se desecha.

Algoritmo DBSCAN II

- Los pasos 2-5 se repiten hasta que todos los puntos han sido visitados o añadidos a algún cluster.
- Básicamente se añaden al cluster todos aquellos puntos densamente conectados desde los puntos del cluster. Esto permite una gran cantidad de geometrías para los clusters.
- Sin embargo hay tres parámetros que influyen notablemente en el método:
 - La función de distancia elegida, que definirá la geometría de la ϵ -vecindad.
 - Valores altos para ϵ requieren valores altos para *MinPts*.
 - Un valor bajo para ϵ dará lugar a un número alto de clusters pequeños. A medida que se vaya aumentando dicho valor se irán produciendo un número más pequeño de clusters, pero aumentará el número de puntos ruido.

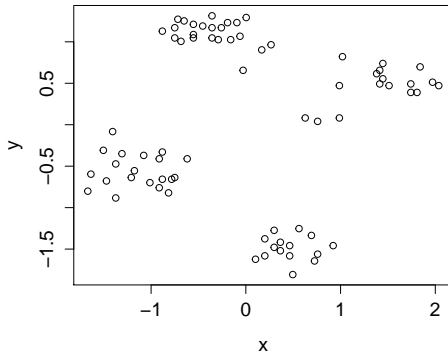
Método DBSCAN: Selección de parámetros I

- El parámetro $MinPts$ se suele fijar a $MinPts = d + 1$, siendo d el número de dimensiones (algunos autores) utilizan $MinPts = 2d - 1$.
 - El valor $MinPts = 1$ no tiene sentido.
 - El valor $MinPts = 2$ equivale a un agrupamiento jerárquico de enlace simple cortado a la altura ϵ .
 - Los valores grandes de $MinPts$ son generalmente mejores para datos con ruido.
 - A medida que el conjunto de datos sea mayor $MinPts$ debe ser mayor.

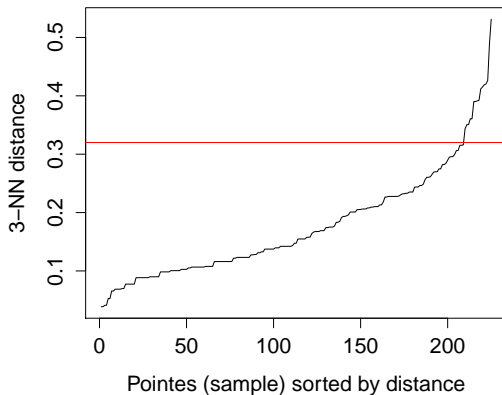
Método DBSCAN: Selección de parámetros II

- El valor ϵ puede ser elegido por medio de un gráfico de k-distancias con $k = \text{MinPts}$.
 - Se fija el valor de ϵ a la distancias en la que se muestre una fuerte curvatura (es decir, un codo).
 - A medida que ϵ se va haciendo más grande, el tamaño de los clusters obtenidos aumentará.

Método DBSCAN: Ejemplo I



Método DBSCAN: Ejemplo II



Método DBSCAN: Ejemplo III

Ejemplo DBSCAN con $\epsilon = 0,32$ y $MinPts = 3$.

Método DBSCAN: Conclusiones

● Ventajas:

- Los clusters pueden tener formas y tamaños arbitrarios.
- El número de clusters se determina automáticamente.
- Puede detectar y aislar el ruido, siendo robusto a los outliers.
- Se puede optimizar utilizando estructuras de datos para los índices (por ejemplo árboles K-D).

● Desventajas:

- No es enteramente determinista, un punto frontera puede pertenecer a dos clusters distintos.
 - Los parámetros necesarios pueden ser difíciles de encontrar.
 - Es muy sensible a los valores de dichos parámetros.
- OPTICS (Ordering points to identify the clustering structure) es una generalización de DBSCAN en la que sólo se fija el parámetro *MinPts*

Método OPTICS: Generalidades

- OPTICS sólo requiere el parámetro *MinPts*.
- No genera un conjunto de clusters
 - Ordena los elementos del conjunto de datos de tal forma que aquellos puntos cercanos son vecinos en dicha ordenación.
 - También se almacena la distancia que se necesita para que dichos puntos pertenezcan al mismo cluster.
- La información sobre la ordenación y la distancia es equivalente a un DBSCAN para distintos valores de ϵ .
- Por lo tanto, se puede utilizar tanto de forma automática como interactiva a la hora de encontrar un clustering en el conjunto de datos.

Evaluación de los agrupamientos I

- Una vez aplicada una técnica de agrupamiento concreta:
 - ¿Son los clusters generados un fiel reflejo de la verdadera naturaleza de los datos?
- Por regla general, la mayoría de las técnicas se ven influenciadas por dos parámetros:
 - La medida de distancia utilizada.
 - El número de clusters que la técnica concreta debe buscar.
- Esto nos lleva a que generalmente deberíamos elegir entre varias configuraciones posibles después de probar varias combinaciones de parámetros.
- Esta tarea se puede llevar a cabo por medio de una **medida de la calidad del agrupamiento**.

Evaluación de los agrupamientos II

- Una medida de calidad trata de evaluar cómo se de buena es la estructura revelada por el agrupamiento obtenido respecto a la estructura real que presentan los datos.
- De todas formas, hay que tener en cuenta que el éxito de la medida de calidad seleccionada depende de la técnica utilizada y la propias características de los datos.

Evaluación de los agrupamientos III

- Atendiendo al resultado de una técnica de agrupamiento, este debe satisfacer dos propiedades:
 - **Compactación:** nos indica cuán cerca están entre sí los elementos de un cluster. A mayor varianza entre los elementos de un cluster menos compacto será este y, al contrario, a menor varianza mas compacto será.
 - **Separabilidad:** nos indica cuán distintos son los clusters entre sí.
- Una forma intuitiva de medir estas características puede ser las distancias intra e intercluster.
 - Un agrupamiento compacto y separable se debe caracterizar por una distancia intracluster pequeña y una distancia intercluster grande.

Índice Silueta (Silhouette index) I

- Supongamos que la observación i pertenece al cluster C_i
- Para cada observación i se define el índice silueta $s(i)$ de la siguiente forma:
- Sea $a(i)$ la distancia media entre i y todos los elementos de su mismo cluster.

$$a(i) = \frac{1}{|C_i|} \sum_{\forall j \in C_i \wedge j \neq i} d(i, j)$$

- Si el cluster $C(i)$ tiene un sólo elemento entonces $s(i) = 0$

Índice Silueta (Silhouette index) II

- Sea $b(i)$ la distancia media entre i y todos los elementos del cluster más cercano

$$b(i) = \min_{k \neq i} (d(i, C_k)) \text{ con } d(i, C_k) = \frac{1}{|C_k|} \sum_{\forall j \in C_k} d(i, j)$$

- El índice silueta se define como:

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

Índice Silueta (Silhouette index) III

- El índice silueta varía entre $[-1, 1]$.
- Un $s(i)$ cercano a 1 indica que la observación está muy bien agrupada.
- Un $s(i)$ cercano a 0 indica que la observación está entre dos clusters.
- Un $s(i)$ negativo indica que la observación está mal agrupada.
- El coeficiente silueta para todo el agrupamiento sería la media de todos los índices siluetas

Índice Gap I

- El índice gap [Tibshirani *et al.*, 2001] compara la varianza total intra-cluster observada para diferentes valores de k con el valor esperado en una distribución uniforme de referencia.
- Supongamos que nuestros datos han sido agrupados en k clusters $\{C_1, C_2, \dots, C_k\}$, con $n_r = |C_r|$.

Índice Gap II

- Sea D_r la suma de la distancia entre todos los elementos del cluster r :

$$D_r = \sum_{i, i' \in C_r} d_{ii'}$$

- Sea W_k la distancia intra-cluster

$$W_k = \sum_{r=1}^k \frac{1}{2n_r} D_r$$

Índice Gap III

- **Algoritmo:**

- 1 Calcular W_k para distintos valores de k
- 2 Generar B conjuntos de referencia usando un muestreo uniforme
- 3 Calcular la suma de la distancia intra-cluster, W_{bk}^* en cada uno de los B conjuntos y para distintos valores de k .
- 4 El índice Gap se calcula como

$$Gap(k) = \frac{1}{B} \sum_b \log(W_{kb}^*) - \log(W_k)$$

Índice Gap IV

- 5 Sea $\bar{l} = \frac{1}{B} \sum_b \log(W_{kb}^*)$, calcular las desviaciones estandares s_k como:

$$s_k = \sqrt{\frac{1}{B} \sum_b (\log(W_{kb}^*) - \bar{l})^2}$$

- 6 Se determina el valor óptimo de k como:

$$k_{opt} = \underset{k}{\text{mín}} (gap(k) \geq gap(k+1) - s_k)$$

- Existen otros criterios pero este criterio ha demostrado experimentalmente mejor comportamiento.

Índice Davies-Bouldin I

- Sea un agrupamiento de c clusters $\{C_1, C_2, \dots, C_c\}$.
- Se calcula la distancia intracluster para cada cluster

$$s_i = \frac{1}{|C_i|} \sum_{x \in C_i} d(x - \bar{x}_i)^2$$

- donde $\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_c\}$ son los centroides de cada cluster.
- Calculamos la distancia entre prototipos:

$$d_{ij} = d(\bar{x}_i - \bar{x}_j)^2$$

- Para cada cluster, se calcula el siguiente ratio;

$$r_i = \max_{j: j \neq i} \frac{s_i + s_j}{d_{ij}}$$

Índice Davies-Bouldin II

- Siendo el índice de Davies-Bouldin la media de dichos valores:

$$r = \frac{1}{c} \sum_{i=1}^c r_i$$

- Según este índice el valor óptimo para el número de clusters es aquel que hace mínimo el índice.
- Hay que tener en cuenta el valor de r mínimo se consigue con valores pequeños el numerador de r_i y valores grandes en el denominador.
 - Es decir favorece la creación de agrupamientos compactos y separados.

Índice de separación de Dunn

- Para definir este índice necesitamos primero definir el **diámetro de un cluster**:

$$\Delta(C_i) = \max_{x,y \in C_i} d(x - y)$$

- siendo la distancia intercluster

$$\delta(C_i, C_j) = \min_{x \in C_i, y \in C_j} d(x - y)$$

- de esta forma el **índice de separación de Dunn** se define como:

$$r = \min_k \min_{j,j=1..c} \frac{\delta(C_i, C_j)}{\max_k \Delta(C_k)}$$

Índice de validez I

- El **índice de validez** intenta evaluar la validez de un agrupamiento buscando un compromiso entre separabilidad y compactación.
- Para ello necesitamos calcular la distancia intracluster media:

$$\bar{s} = \frac{1}{c} \sum_{i=1}^c \frac{s_i}{s}$$

- donde s es la distancia intracluster considerado el conjunto de datos como un único cluster.

Índice de validez II

- La medida de separación del agrupamiento se puede calcular como

$$sep = \frac{D_{min}}{D_{max}} \sum_{i=1}^c \left(\sum_{j=1}^c d(\bar{x}_i - \bar{x}_j) \right)^{-1}$$

- con

$$D_{min} = \min_{i,j=1..c} d(\bar{x}_i - \bar{x}_j) \text{ y } D_{max} = \max_{i,j=1..c} d(\bar{x}_i - \bar{x}_j)$$

Índice de validez III

- de esta forma **índice de validez**, SD, es una combinación ponderada de la separación y la distancia intercluster media.

$$DS = \alpha \bar{s} + sep$$

- donde α se utiliza para balancear el peso de los componentes en el índice.
- El valor de c que haga mínimo el índice se puede tomar como una medida del número óptimo de clusters.

Resumen I

- En este capítulo hemos abordado las técnicas principales para aprendizaje no supervisado, centrándonos en las técnicas de agrupamiento o clustering.
- Se ha presentado el concepto de distancia y similaridad como elemento clave para definir los agrupamientos, así como diferentes formas de medirlos
- Primero hemos analizado las técnicas de agrupamiento jerárquico y sus distintas variantes dependiendo de cómo se calcule la distancia entre grupos.

Resumen II

- El segundo grupo de técnicas que se han analizado corresponde con las técnicas particionales, que se han dividido en dos grupos:
 - Técnicas basadas en criterios globales: K-medias y K-medoides.
 - Técnicas basadas en criterios locales como DBSCAN
- Por último, se han analizado algunas medidas para medir la calidad del agrupamiento obtenido.

Referencias I



Krzysztof J Cios, Witold Pedrycz, and Roman W Swiniarski.
Data mining methods for knowledge discovery, volume 458.
Springer Science & Business Media, 2012.



Anil K. Jain and Richard C. Dubes.
Algorithms for Clustering Data.
Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1988.



Alboukadel Kassambara.
Practical Guide to Cluster Analysis in R *Practical Guide to Cluster Analysis in R*.
STHDA, 2017.



Godfrey N Lance and William Thomas Williams.
A general theory of classificatory sorting strategies ii. clustering systems.
The computer journal, 10(3):271–277, 1967.



Roque Luis Marín Morales and José Tomás Palma Méndez, editors.
Inteligencia artificial: técnicas, métodos y aplicaciones.
McGraw-Hill, 2008.

Referencias II



Basilio Sierra Araujo.

Aprendizaje automático: conceptos básicos y avanzados: aspectos prácticos utilizando el software Weka.

Pearson Prentice Hall Madrid, 2006.



Robert Tibshirani, Guenther Walther, and Trevor Hastie.

Estimating the number of clusters in a dataset via the gap statistic.

Journal of the Royal Statistical Society B, 63:411–423, 2001.