

Sistemas Inteligentes

Escuela Técnica Superior de Informática

Universitat Politècnica de València

Tema B2T4:

Aprendizaje no supervisado: algoritmo k-medias.

Índice

- 1 Introducción ▷ 2
- 2 Agrupamientos particionales ▷ 4
- 3 Algoritmo C-Medias ▷ 9

Índice

- 1 *Introducción* ▷ 2
- 2 Agrupamientos particionales ▷ 4
- 3 Algoritmo C-Medias ▷ 9

Clustering

Clustering es un problema generalmente mal definido: no existe una definición precisa y universalmente aceptada.

Según [Anderberg, 1973], el objetivo del clustering es:

Agrupar objetos en clases tales que los de una misma clase presenten un alto grado de **asociación natural** entre sí, mientras que las clases sean relativamente distintas unas de otras.

Dicho de otra manera:

Encontrar **agrupamientos naturales** en un conjunto de objetos, de forma que la descripción de éstos se realice en términos de clases o grupos de objetos con fuertes semejanzas internas.

Dos tipos de clustering: **Particional** y **Jerárquico**.

Índice

- 1 Introducción ▷ 2
- 2 *Agrupamientos particionales* ▷ 4
- 3 Algoritmo C-Medias ▷ 9

Clustering particional

Problema genérico:

Asumimos disponible una **función criterio** J para evaluar la calidad de cualquier partición de N datos en C clases. De este modo, el problema del clustering puede verse como uno de búsqueda del tipo:

$$\Pi^* = \arg \min_{\Pi = \{X_1, \dots, X_C\}} J(\Pi) \quad (1)$$

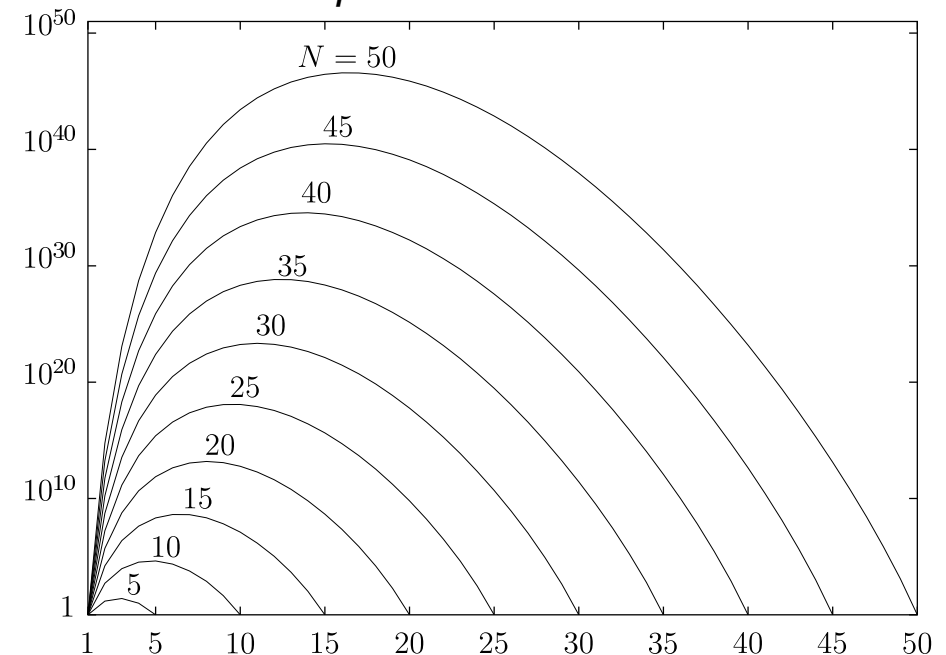
Dificultad:

El número de particiones a explorar es muy elevado incluso para valores pequeños de N y C (ver dcha.). No es factible buscar soluciones globalmente óptimas mediante técnicas de enumeración completa (explícita o implícita) salvo en casos particulares.

Solución:

Soluciones subóptimas obtenidas mediante algoritmos aproximados.

Número de particiones en función de C para varios N



Clustering particional: Criterio “suma de errores cuadráticos” (SEC)

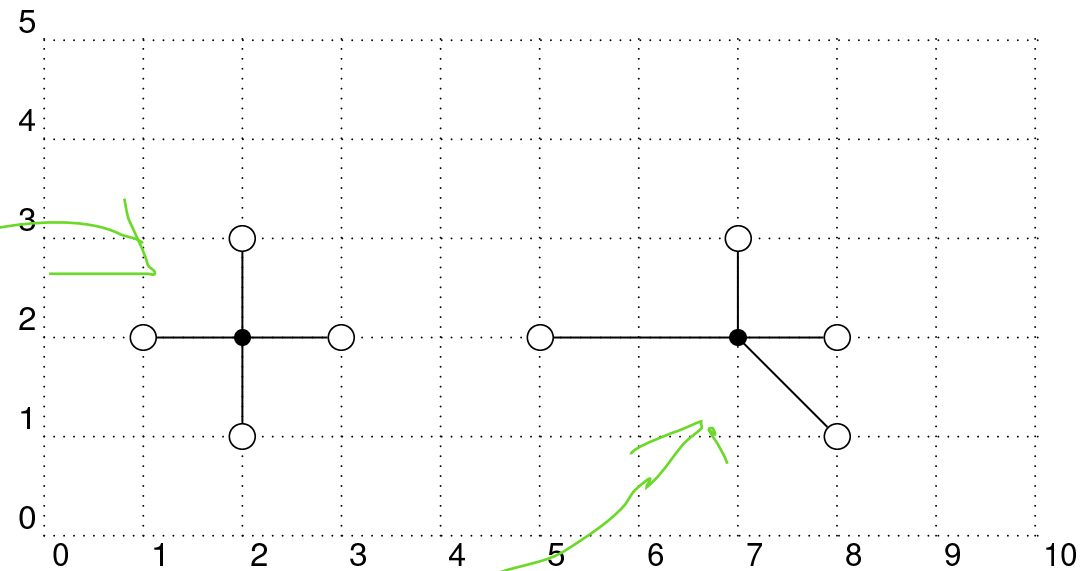
La SEC de una partición de N datos en C clusters, $\Pi = \{X_1, \dots, X_C\}$, es:

$$J(X_1, \dots, X_C) = \sum_c J_c, \quad J_c = \sum_{x \in X_c} \|x - m_c\|^2, \quad m_c = \frac{1}{|X_c|} \sum_{x \in X_c} x \quad (2)$$

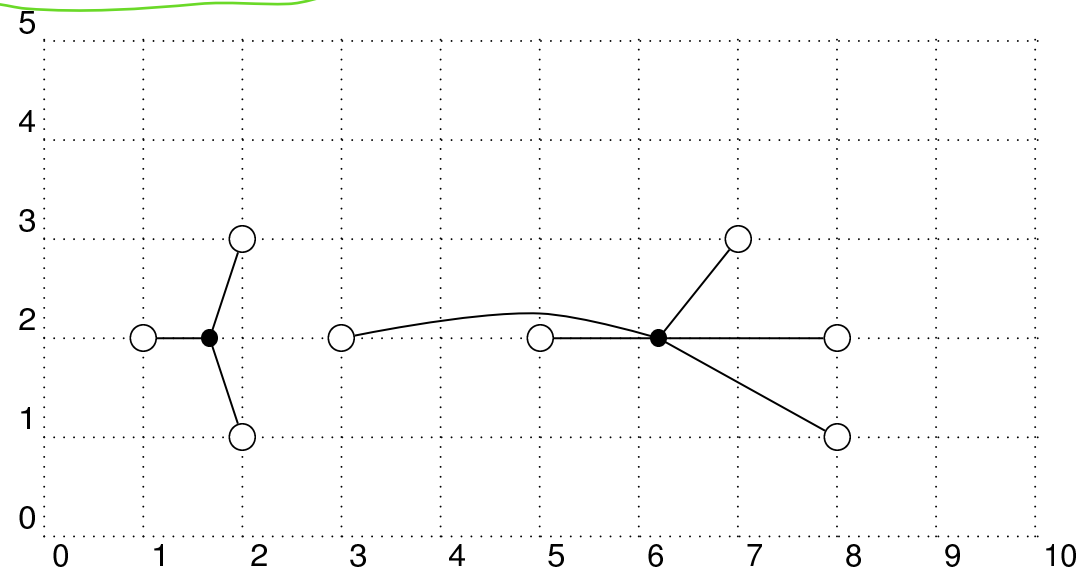
Interpretación:

- En cada cluster X_c su *media*, m_c , se interpreta como el “prototipo natural” de X_c . Cada dato $x \in X_c$, se interpreta como una “versión distorsionada” de m_c y la distorsión de x se caracteriza por el *vector error* $x - m_c$.
- Como su nombre indica, el criterio SEC mide la suma (o media) de los cuadrados de las magnitudes de estos vectores error y, obviamente, es un criterio a minimizar.
- La *media* de cada cluster es el punto que representa los datos del cluster con menor SEC.

Ejemplo de clustering particional



SEC = 12



SEC > 12

$$SEC = 1^2 + 1^2 + 1^2 + 1^2$$

+

$$2^2 + 1^2 + 1^2 + \sqrt{2}^2$$

=

$$4 + 8 = 12$$

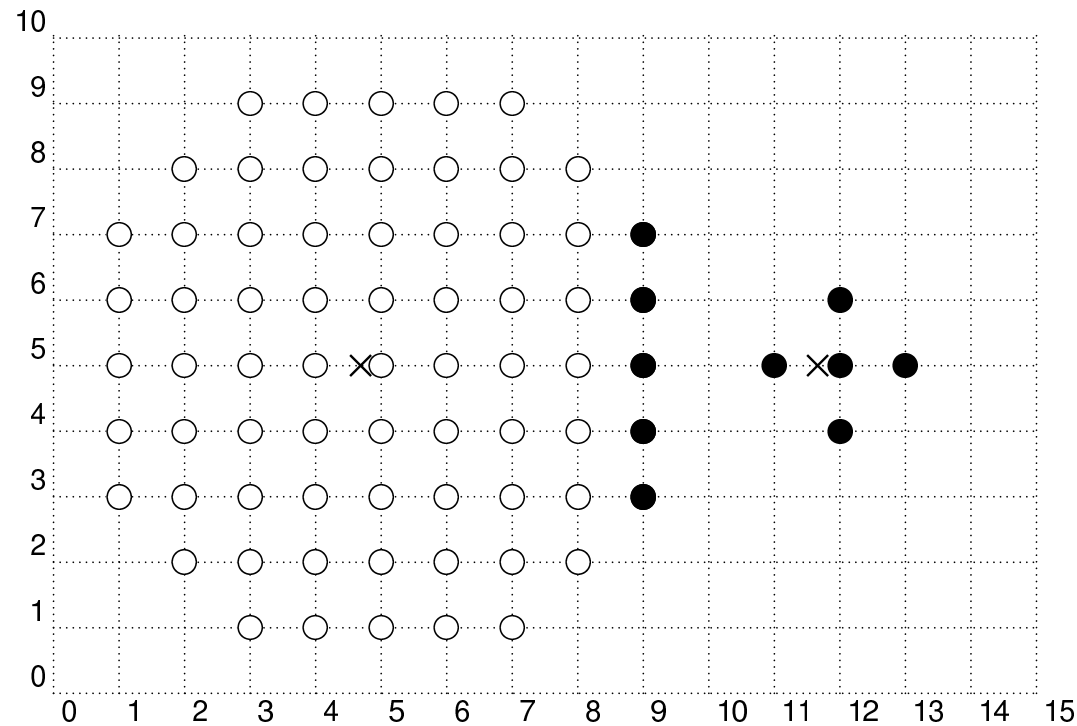
$$\vec{m}_1 = \frac{1}{4} \left(\begin{pmatrix} 1 \\ 2 \end{pmatrix} + \begin{pmatrix} 2 \\ 1 \end{pmatrix} + \begin{pmatrix} 2 \\ 3 \end{pmatrix} + \begin{pmatrix} 3 \\ 2 \end{pmatrix} \right) = \frac{1}{4} \begin{pmatrix} 8 \\ 8 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

$$J(x_1) = \underbrace{(1-2)^2}_{x_1} + \underbrace{(2-2)^2}_{x_1} + \underbrace{(2-2)^2}_{x_1} + \underbrace{(3-2)^2 + (2-2)^2}_{x_1}$$

Bondad del criterio SEC

El criterio SEC es apropiado sólo si los datos forman ***clusters hiperesféricos de tamaño similar***.

Si los tamaños de los clusters son muy distintos, es posible la agrupación natural *no* tenga el mínimo SEC:



Índice

- 1 Introducción ▷ 2
- 2 Agrupamientos particionales ▷ 4
- 3 *Algoritmo C-Medias* ▷ 9

Cálculo incremental de la SEC al transferir x del cluster X_i al X_j

$$X'_i = X_i - \{x\}$$

$$X'_j = X_j + \{x\}$$

$$m'_i = m_i - \frac{x - m_i}{n_i - 1}$$

$$m'_j = m_j + \frac{x - m_j}{n_j + 1}$$

$$J'_i = J_i - \frac{n_i}{n_i - 1} \|x - m_i\|^2$$

$$J'_j = J_j + \frac{n_j}{n_j + 1} \|x - m_j\|^2$$

$$\Delta J = \frac{n_j}{n_j + 1} \|x - m_j\|^2 - \frac{n_i}{n_i - 1} \|x - m_i\|^2$$

La transferencia será provechosa si el incremento de SEC es negativo; es decir:

$$\frac{n_j}{n_j + 1} \|x - m_j\|^2 < \frac{n_i}{n_i - 1} \|x - m_i\|^2 \quad (3)$$

Estas ecuaciones permiten minimizar la SEC mediante refinamientos sucesivos a partir una partición inicial dada.

Optimización de la SEC: algoritmo C -medias

Algorithm C -means (versión “correcta” [Duda & Hart])

Input: X ; C ; $\Pi = \{X_1, \dots, X_C\}$;

Output: $\Pi^* = \{X_1, \dots, X_C\}$; $\mathbf{m}_1, \dots, \mathbf{m}_C$; J

for $c = 1$ **to** C **do** $\mathbf{m}_c = \frac{1}{n_c} \sum_{\mathbf{x} \in X_c} \mathbf{x}$ **endfor**

repeat

$transfers = false$

forall $\mathbf{x} \in X$ (let $i : \mathbf{x} \in X_i$) **do**

if $n_i > 1$ **then**

$$j^* = \arg \min_{j \neq i} \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2$$

$$\Delta J = \frac{n_{j^*}}{n_{j^*} + 1} \|\mathbf{x} - \mathbf{m}_{j^*}\|^2 - \frac{n_i}{n_i - 1} \|\mathbf{x} - \mathbf{m}_i\|^2$$

if $\Delta J < 0$ **then**

$transfers = true$

$$\mathbf{m}_i = \mathbf{m}_i - \frac{\mathbf{x} - \mathbf{m}_i}{n_i - 1} \quad \mathbf{m}_{j^*} = \mathbf{m}_{j^*} + \frac{\mathbf{x} - \mathbf{m}_{j^*}}{n_{j^*} + 1}$$

$$X_i = X_i - \{\mathbf{x}\} \quad X_{j^*} = X_{j^*} + \{\mathbf{x}\}$$

$$J = J + \Delta J$$

endif

endif

endforall

until $\neg transfers$ // Coste por iteración: $O(N \cdot C \cdot D)$, $N = |X|$, $D = \text{dimensión}$

Optimización de la SEC: otra versión de C -medias

Algorithm C -means (versión “popular”)

Input: X ; C ; $\Pi = \{X_1, \dots, X_C\}$;

Output: $\Pi^* = \{X_1, \dots, X_C\}$; $\mathbf{m}_1, \dots, \mathbf{m}_C$

repeat

transfers = false

for $c = 1$ **to** C **do** $\mathbf{m}_c = \frac{1}{n_c} \sum_{\mathbf{x} \in X_c} \mathbf{x}$ **endfor**

forall $\mathbf{x} \in X$ (let $i : \mathbf{x} \in X_i$) **do**

if $n_i > 1$ **then**

$j^* = \arg \min_{1 \leq j \leq C} d(\mathbf{x}, \mathbf{m}_j)$

if $j^* \neq i$ **then**

transfers = true

$X_i = X_i - \{\mathbf{x}\}$; $X_{j^*} = X_{j^*} + \{\mathbf{x}\}$

endif

endif

endforall

until $\neg \text{transfers}$

// Coste por iteración: $O(N \cdot C \cdot D)$, $N = |X|$, $D = \text{coste de } d(\cdot, \cdot)$

Optimalidad de los algoritmos *C*-medias

- Ninguna de las versiones del algoritmo *C-medias* garantiza la obtención de un mínimo global de la SEC
- La versión de Duda & Hart obtiene un *mínimo local*
- La versión “popular” no garantiza la minimización local en algunos casos

Ejemplo:

$$X = \{1, 3, 4.5\} \subset \mathbb{R}; \quad \Pi^0 = \{\{1, 3\}, \{4.5\}\}; \quad J^0 = 2.0$$

$$\text{C-medias “popular”}: \quad \Pi^* = \Pi^0; \quad J^* = J^0 = 2.0$$

$$\text{C-medias Duda \& Hart}: \quad \Pi^* = \Pi^1 = \{\{1\}, \{3, 4.5\}\}; \quad J^* = J^1 = 1.125$$

El criterio SEC y Cuantificación Vectorial

Los siguientes criterios a minimizar son equivalentes:

$$J(X_1, \dots, X_C) = \sum_c \sum_{\mathbf{x} \in X_c} \|\mathbf{x} - \mathbf{m}_c\|^2 \quad (4)$$

$$J(X_1, \dots, X_C; \mathbf{r}_1, \dots, \mathbf{r}_C) = \sum_c \sum_{\mathbf{x} \in X_c} \|\mathbf{x} - \mathbf{r}_c\|^2 \quad (5)$$

$$J(\mathbf{r}_1, \dots, \mathbf{r}_C) = \sum_{\mathbf{x}} \min_c \|\mathbf{x} - \mathbf{r}_c\|^2 \quad (6)$$

Justificación:

- (5) equivale a (4) pues, para toda partición X_1, \dots, X_C , los *representantes (de cluster)* $\mathbf{r}_1, \dots, \mathbf{r}_C$ que minimizan (5) son las medias de los clusters.
- (5) equivale a (6) ya que, para todo conjunto de representantes $\mathbf{r}_1, \dots, \mathbf{r}_C$, la partición que minimiza (5) es aquella en la que cada dato se asigna al cluster de su representante más próximo.
- (6) se conoce como el problema del *diseño de un cuantificador vectorial en teoría de la información*

Otra interpretación del criterio SEC

El criterio SEC se puede reescribir sin incluir las medias de los clusters:

$$J(X_1, \dots, X_C) = \frac{1}{2} \sum_c n_c \bar{s}_c \quad (7)$$

donde n_c es el número de datos en X_c y \bar{s}_c es la media de las distancias Euclídeas al cuadrado entre todos estos datos:

$$\bar{s}_c = \frac{1}{n_c^2} \sum_{\mathbf{x}, \mathbf{x}' \in X_c} \|\mathbf{x} - \mathbf{x}'\|^2 \quad (8)$$

Luego la SEC se puede interpretar como una suma ponderada de medias de distancias al cuadrado “intra-cluster”.

Con base en esta interpretación, podemos redefinir \bar{s}_c para obtener criterios parecidos a la SEC (válidos incluso con datos *no-vectoriales*):

$$\bar{s}_c = \frac{1}{n_c^2} \sum_{x, x' \in X_c} d(x, x') \quad \bar{s}_c = \frac{1}{n_c^2} \max_{x, x' \in X_c} d(x, x') \quad (9)$$