

K-MEDIAS, Y AMIGOS

ALAN REYES-FIGUEROA

INTRODUCCIÓN A LA CIENCIA DE DATOS

(AULA 20) 04.ABRIL.2022

Suponga que tenemos un conjunto de datos $\mathbb{X} = \{\mathbf{x}_i\}_{i=1}^n$, con $\mathbf{x}_i \in \mathbb{R}^d$, y se quiere agrupar estos elementos en k conjuntos o categorías distintas.

- Definimos distancias entre observaciones $d(\mathbf{x}_i, \mathbf{x}_j)$ (e.g. distancia euclídeana).
- Por otro lado, elegimos un representante para cada uno de los k grupos o categorías a construir (e.g. el centroide $\mathbf{c}_k \in \mathbb{R}^d$ del grupo).

El método de **k -medias** (ó *k-means*) elige al azar k representantes, y asigna cada dato \mathbf{x}_i al representante más cercano, según la distancia $d(\cdot, \cdot)$.

Algoritmo: (K-medias)

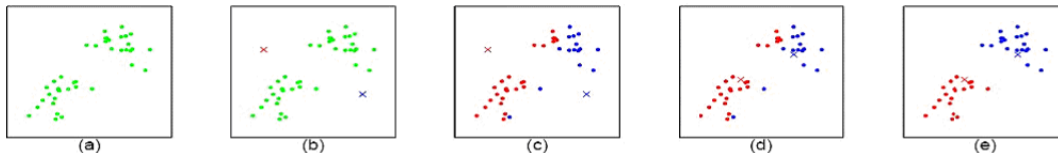
1. Para cada $g = 1, 2, 3, \dots, k$, elegimos de manera aleatoria el centroide $\mathbf{c}_g \in \mathbb{R}^d$ del grupo g .
2. Repetir hasta convergencia:
 - Asignar cada \mathbf{x}_i al representante más cercano según la métrica d , esto es:

$$h(\mathbf{x}_i) = g(i) \in C, \quad \text{con } g(i) = \operatorname{argmin}_{1 \leq g \leq k} d(\mathbf{c}_g, \mathbf{x}_i).$$

- Recalcular los representantes como los centroides de los datos asociados a cada grupo

$$\mathbf{c}_g = \frac{1}{|\{i : g(i) = g\}|} \sum_{i:g(i)=g} \mathbf{x}_i, \quad g = 1, 2, \dots, k.$$

K-medias



- (a) Conjunto de datos \mathbb{X} .
- (b) Se eligen aleatoriamente los centroides \mathbf{c}_g , $k = 2$.
- (c) Cada \mathbf{x}_i se asocia con su centroide más cercano.
- (d) Se recalculan los centroides \mathbf{c}_g , $k = 2$.
- (e) Cada \mathbf{x}_i se etiqueta con su centroide más cercano.
- (f) Repetir (b) a (e) hasta convergencia ...

K-medias

K-medias minimiza una función de costo de una manera particular:
Sea \mathbf{c}_g el representante del grupo g , $g(i)$ el grupo de \mathbf{x}_i , definimos

$$J(\mathbb{X}) = \min_g \min_{\mathbf{c}_g} \sum_g \sum_{i:g(i)=g} \|\mathbf{x}_i - \mathbf{c}_g\|^2 \quad (1)$$

En k-medias J se minimiza en dos pasos desacoplados:

- Fijando g , y minimizando sobre \mathbf{c}_g .
- Fijando \mathbf{c}_g , y minimizando sobre g .

Interpretación:

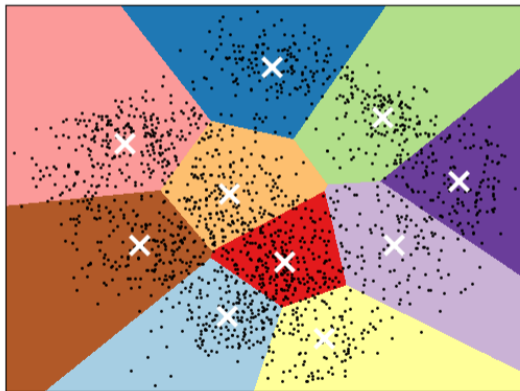
$$\|\mathbf{x}_i - \mathbf{c}_g\|^2 = \|\mathbf{x}_i - \text{decoder}(\text{encoder}(\mathbf{x}_i))\|^2. \quad (2)$$

Se trata de minimizar la varianza dentro de cada grupo. Entre mayor es k menor la varianza (2).

K-medias

Al final, el algoritmo de k -means induce una partición de Voronoi sobre el espacio \mathbb{R}^d :

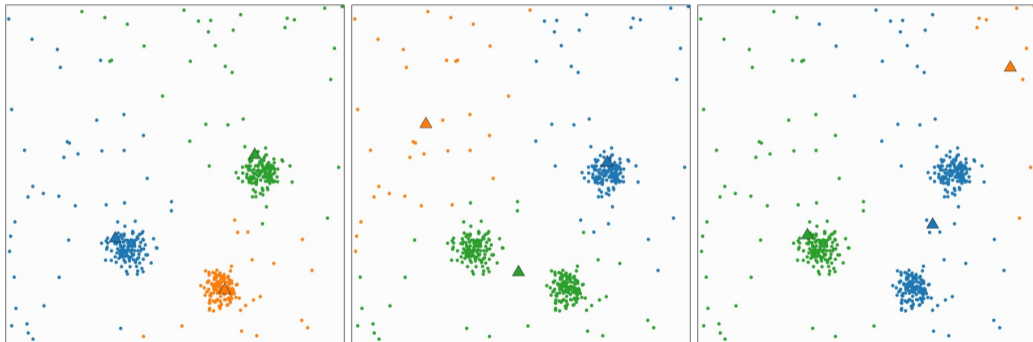
K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross



K-medias

<https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html>

Como la elección inicial de los c_g es aleatoria, es posible obtener diferentes resultados:



Distintos resultados con *k*-means en el mismo conjunto de datos

Observaciones:

- Se sugiere replicar k -means varias veces (con diferentes inicializaciones). Esto permite que aparezcan diferentes estructuras de agrupamiento.
- El diagrama de Voronoi es útil, ya pensando en k -means como un método predictivo.

Existen diversas variantes:

- k -medianas.
- k -medioides.
- Fuzzy k -means.

K-medianas

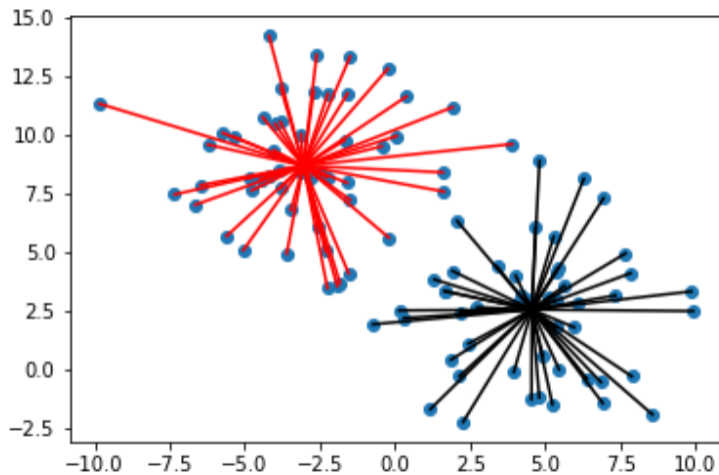
Es una variante de k -medias. En lugar de utilizar el centroide \mathbf{c}_g como representante de cada grupo g , se utiliza la mediana.

Esta mediana se calcula, componente a componente, usando la distancia Manhattan (norma $\|\cdot\|_1$). De este modo, cada uno de los atributos o variables es una observación dentro del conjunto de datos.

Algoritmo: (K-medianas)

1. Para cada $j = 1, 2, 3, \dots, k$, elegimos de manera aleatoria el representante $\mathbf{c}_j \in \mathbb{R}^d$.
2. Repetir hasta convergencia:
 - Asignar cada \mathbf{x}_i al representante más cercano según la métrica d .
 - Recalcular los representantes como las medianas, componente a componente, de los datos asociados a cada grupo.

K-medianas



K-medioides

Es otra variante de k -medias. En lugar de utilizar el centroide \mathbf{c}_g como representante de cada grupo g , se utiliza el dato más cercano al centroide \mathbf{c}_g :

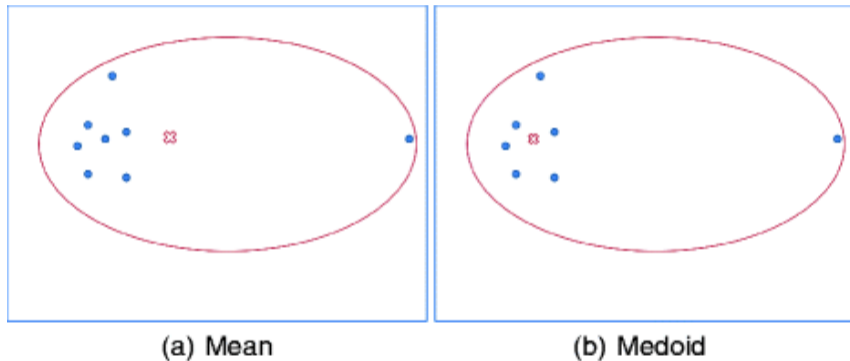
$$\mathbf{c}_g = \mathbf{x}_j, \quad \text{con } j = \operatorname{argmin}_\ell \left\| \mathbf{x}_\ell - \frac{1}{|g(i)|} \sum_{i:g(i)=g} \mathbf{x}_i \right\|^2$$

De este modo, cada uno de los representante se elige dentro del mismo conjunto de datos.

Algoritmo: (K-medioides)

1. Para cada $j = 1, 2, 3, \dots, k$, elegimos de manera aleatoria el representante $\mathbf{c}_j \in \mathbb{R}^d$.
2. Repetir hasta convergencia:
 - Asignar cada \mathbf{x}_i al representante más cercano según la métrica d .
 - Recalcular los representantes como el dato más cercano a \mathbf{c}_g .

K-medioides



Agrupamiento difuso:

El agrupamiento difuso (*fuzzy clustering*) es una clase de algoritmos de agrupamiento donde, en lugar de asignar un único grupo a cada dato \mathbf{x}_i , cada elemento tiene un grado de pertenencia (difuso) a cada uno de los grupos.

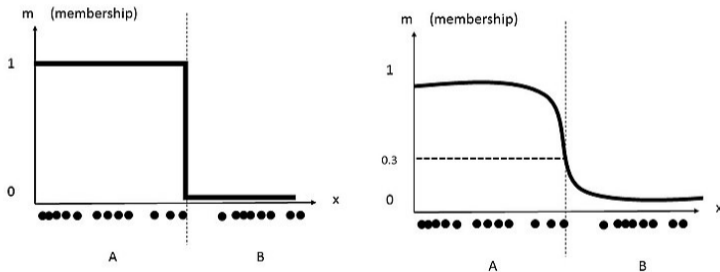
- surge de la necesidad de resolver una deficiencia del agrupamiento exclusivo (agrupación inequívoca).
- implementaciones a partir del surgimiento de la lógica difusa (Zadeh, 1965).
- Se representa la similitud entre un elemento \mathbf{x}_i y un grupo g por una función, llamada función de pertenencia $\mathbf{w}_i : \mathbf{x}_i \rightarrow [0, 1]^k$, que toma valores entre cero y uno.

Fuzzy K-medias

Básicamente, a cada dato \mathbf{x}_i , el clasificador difuso asigna un vector de coeficientes

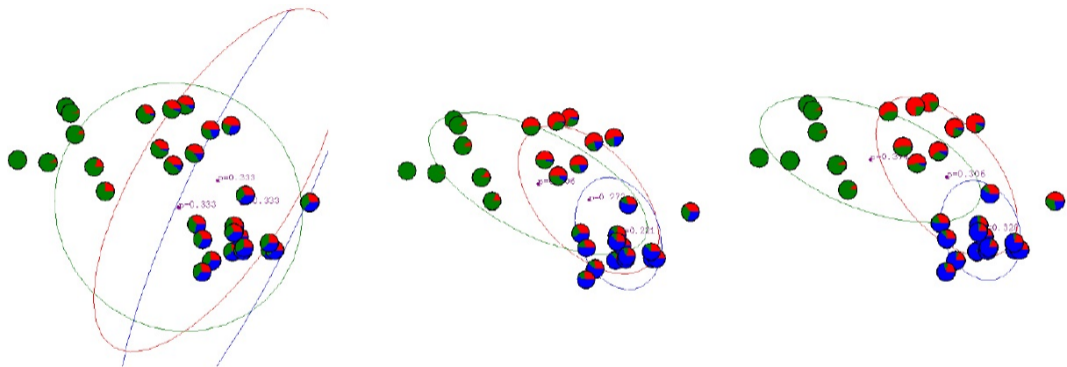
$$h(\mathbf{x}_i) = \mathbf{w}_i = (w_{i1}, w_{i2}, \dots, w_{ik}),$$

con $w_{ij} = \mathbb{P}(\mathbf{x}_i \in g_j)$.



Diferencias entre (a) *hard-clustering*, y (b) *fuzzy-clustering*.

Fuzzy K-medias



Segmentación o clasificación obtenida con *fuzzy-clustering*.

Fuzzy K-medias

Para cada dato $\mathbf{x}_i \in \mathbb{R}^d$, consideramos un vector de coeficientes $\mathbf{w}_i = (w_{i1}, \dots, w_{ik}) \in \mathbb{R}^k$, con $w_{ij} \geq 0, \forall j = 1, 2, \dots, k$.

Con este esquema, el centroide de un grupo g se calcula como el promedio ponderado de sus elementos:

$$\mathbf{c}_j = \frac{\sum_{i:g(i)=j} w_{ij}^m \mathbf{x}_i}{\sum_{i:g(i)=j} w_{ij}^m},$$

$m > 0$ es un hiperparámetro que controla el suavizamiento. A mayor m , mayor difusividad.

Los pesos w_{ij} se recalculan como

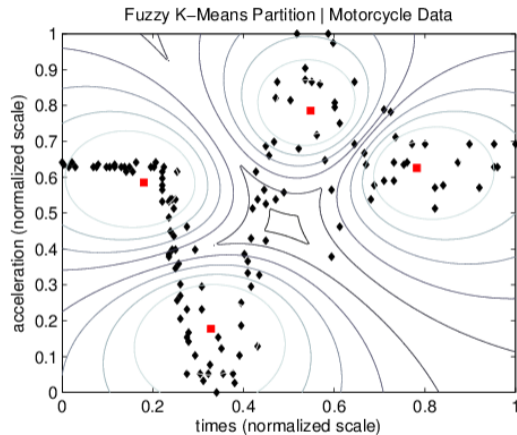
$$w_{ij} = \frac{1}{\sum_{\ell=1}^k \left(\frac{\|\mathbf{x}_i - \mathbf{c}_j\|}{\|\mathbf{x}_i - \mathbf{c}_\ell\|} \right)^{\frac{2}{m-1}}}.$$

Fuzzy K-medias

Algoritmo: (Fuzzy k -medias) Dado un conjunto de datos $\mathbb{X} = \{\mathbf{x}_i\}$, y un conjunto de clases $C = \{g_1, g_2, \dots, g_k\}$, el algoritmo construye una matriz $\mathbf{W} = (w_{ij}) \in \mathbb{R}^{n \times k}$ con los grados de pertenencia.

1. Para cada $j = 1, 2, 3, \dots, k$, elegimos de manera aleatoria el centroide $\mathbf{c}_j \in \mathbb{R}^d$. Elegimos pesos aleatorios \mathbf{w}_j (e.g. $\mathbf{w}_j = (1, 1, \dots, 1)$, $\forall j$).
2. Repetir hasta convergencia:
 - Asignar cada \mathbf{x}_i al representante más cercano según la métrica d .
 - Recalcular los representantes como los centroides ponderados de los datos asociados a cada grupo.
 - Recalcular los grados de pertenencia w_{ij} .

Fuzzy K-medias



Curvas de nivel de la partición difusa.

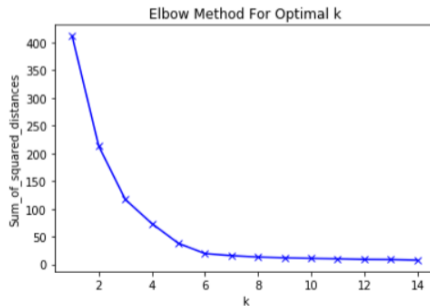
Segmentación de imágenes:

https://www.youtube.com/watch?v=yR7k19YBqiwab_channel=Computerphile



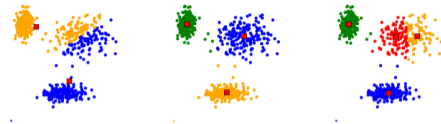
Observaciones:

- Como la distancia euclideana da igual peso a cada dimensión, mejor normalizar los datos (normalizar o estandarizar).
- Existen muchas heurísticas para elegir k , e.g. el método del “codo”: cómo cambia la suma de variaciones por grupo (1) en función de k .

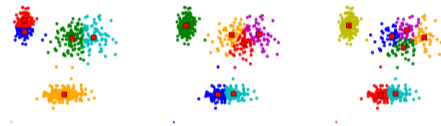


Observaciones

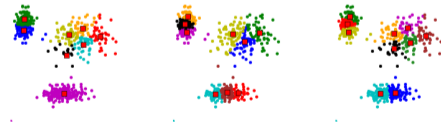
Centers = 2; FPC = 0.79 Centers = 3; FPC = 0.88 Centers = 4; FPC = 0.81



Centers = 5; FPC = 0.72 Centers = 6; FPC = 0.71 Centers = 7; FPC = 0.69



Centers = 8; FPC = 0.64 Centers = 9; FPC = 0.58 Centers = 10; FPC = 0.58

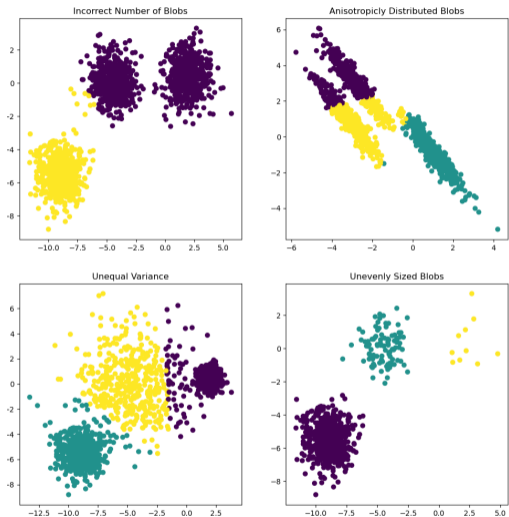


Particiones obtenidas con diferentes valores de k .

Observaciones

- En muchas situaciones, k -medias falla.

Solución:
transformar los datos o usar métodos donde $d(\cdot, \cdot)$ captura la forma local de los datos.



- Muchas veces se usa primero algún método de reducción de dimensión, cuando la dimensionalidad del espacio es muy alta. Las distancias (euclideanas) pierden poder discriminativo en dimensiones altas. Se llama la *maldición de la alta dimensionalidad*.
- Si los datos no son muy continuos, se prefiere tomar como representantes observaciones de la muestra (*k*-medianas, *k*-medioides).
- Miles de variantes! (también porque para conjuntos de datos grandes, el algoritmo básico es demasiado costoso).
- Es buena idea correr el algoritmo con diferentes puntos de arranque para evitar óptimos locales.