

CLUSTERING

Dato un set S , un numero K , l'idea è trovare le K partizioni di S : S_1, \dots, S_K che sono **omogenee** e **ben separate**. È un problema **non supervisionato**.

Dati i punti $p_1, \dots, p_n \in \mathbb{R}^m$

Dada la funzione distanza $d: \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$, dove va definita una distanza fra due vettori: $d(x, y) = \|x - y\|_{1/2/\infty}$

Definiamo per ogni cluster un centroide $x_j \in \mathbb{R}^m$.

Il problema di trovare i migliori cluster:

$$\begin{cases} \min \sum_{i=1}^n \min_{j=1 \dots K} d(p_i, x_j) \\ x_j \in \mathbb{R}^m \quad \forall j=1 \dots K \end{cases}$$

Nota 2: $\|\cdot\|_2$

$$\begin{cases} \min \sum_{i=1}^n \min_{j=1 \dots K} \|p_i - x_j\|_2^2 \\ x_j \in \mathbb{R}^m \quad \forall j=1 \dots K \end{cases}$$

Per trovare l'ottimo per **$K=1$** , calcoliamo $\nabla f(x) = 0$.

Notiamo che per $K=1$ è inutile l'operazione di $\min_{j=1 \dots K} (\cdot)$:

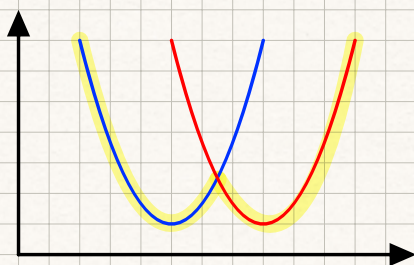
$$\min \sum_{i=1}^n \|p_i - x\|_2^2 = \sum_{i=1}^n (x - p_i)^T (x - p_i)$$

$$\nabla \left| \sum_{i=1}^n (x - p_i)^T (x - p_i) \right| = 2ex - 2 \sum_{i=1}^n p_i = 0 \Rightarrow x = \underbrace{\frac{\sum_{i=1}^n p_i}{n}}_{\text{media}}$$

Nota: questo problema è convesso solo per $K=1$. Per $K > 1$ il problema diventa **non-convesso** e **non differenziabile**.

Perché è non convesso se $K > 1$?

Trovare il miglior centroide per un singolo cluster equivale a minimizzare un quadrato, convesso per definizione ($\min \|p_i - x_i\|^2$). Per più cluster, diciamo 2, dobbiamo minimizzare contemporaneamente più parabole:



Fissando p_i , e per i centroidi $x_j : j=1 \dots K$, possiamo scrivere che:

$$\min_{j=1 \dots K} \|p_i - x_j\|_2^2 = \begin{cases} \min \sum_{j=1}^K \alpha_{ij} \|p_i - x_j\|_2^2 \\ \sum_{j=1}^K \alpha_{ij} = 1 \\ \alpha_{ij} \geq 0 \quad \forall j=1 \dots K \end{cases}$$

Dove $\alpha_{ij} = |\emptyset, \dots, \emptyset, 1, \emptyset, \dots, \emptyset|$ ovvero con 1 al j -esimo elemento. La soluzione ottima dell'ultimo sistema è:

$$\alpha_{ij}^* = \begin{cases} 1 & \|p_i - x_j\|_2 = \min_{h=1 \dots K} \|p_i - x_h\|_2 \\ \emptyset & \text{altrimenti} \end{cases}$$

TEOREMA

Il problema con $\|\cdot\|_2$, equivale a risolvere il seguente sistema **non-convesso** ma **differentenziabile**:

$$\left\{ \begin{array}{l} \min_{x, \alpha} f(x, \alpha) := \sum_{i=1}^e \sum_{j=1}^k \alpha_{ij} \|p_i - x_j\|_2^2 \\ \sum_{j=1}^k \alpha_{ij} = 1 \quad \forall i=1 \dots e \\ \alpha_{ij} \geq 0 \quad \forall i=1 \dots e \text{ e } \forall j=1 \dots k \\ x_j \in \mathbb{R}^n \quad \forall j=1 \dots k \end{array} \right.$$

Dove bisogna cercare un solo minimo, ma essendo differentenziabile possiamo usare il KKT, o altri metodi di risoluzione.

Nota:

Nonostante abbiamo definito $\alpha_{ij} = 1$ o $\alpha_{ij} = 0$, nel sistema non c'è nessuna condizione che esplicitamente dice che sia così. Abbiamo solo la condizione che $\sum_{j=1}^k \alpha_{ij} = 1$ e $\alpha \geq 0$. Quindi una soluzione ammissibile può essere $\alpha_i = (\frac{1}{2}, \frac{1}{2})$ per $k=2$, fissato un punto i . Questo può non succedere perché fissato un punto, il problema diventa convesso su una regione ammissibile che è un poliedro, per cui l'ottimo è necessariamente su un vertice del poliedro, per cui $\alpha_{ij} = (1, 0)$ o $(0, 1)$.

K-MEANS

L'algoritmo di k-means si basa sulla seguente proprietà del sistema definito sopra nel teorema, ovvero che se fissiamo il centroide x_i , possiamo dividere il problema in k sotto problemi di tipo **LP**, nella forma:

$$\begin{cases} \min \sum_{j=1}^k \alpha_{ij} \|p_i - x_j\|_2^2 \\ \sum_{j=1}^k \alpha_{ij} = 1 \\ \alpha_{ij} \geq 0 \quad \forall j=1 \dots k \end{cases}$$

Dove la soluzione ottima è

$$\alpha_{ij}^* = \begin{cases} 1 & \text{se } j : \|p_i - x_j\|_2 = \min_{h=1 \dots k} \|p_i - x_h\|_2 \\ 0 & \text{altrimenti} \end{cases}$$

SS invece fissiamo α_{ij} , il problema si riduce in k sottoproblemi di tipo **QP** **convessi**

$$\begin{cases} \min \sum_{i=1}^e \alpha_{ij} \|p_i - x_j\|_2^2 \\ x_j \in \mathbb{R}^n \end{cases}$$

Dove l'ottimo è: $x_j = \frac{\sum_{i=1}^e \alpha_{ij} p_i}{\sum_{i=1}^e \alpha_{ij}}$] **media**

Objective Function

$$f(x, \alpha) = \sum_{i=1}^e \sum_{j=1}^k \alpha_{ij} \|p_i - x_j\|_2^2$$

TEOREMA

L'algoritmo di K-means si ferma dopo un numero finito di iterazioni, a una soluzione (x^*, α^*) del KKT del problema

$$\left\{ \begin{array}{l} \min_{x, \alpha} f(x, \alpha) := \sum_{i=1}^e \sum_{j=1}^k \alpha_{ij} \|p_i - x_j\|_2^2 \\ \sum_{j=1}^k \alpha_{ij} = 1 \quad \forall i=1 \dots e \\ \alpha_{ij} \geq 0 \quad \forall i=1 \dots e \quad \text{e} \quad \forall j=1 \dots k \\ x_j \in \mathbb{R}^n \quad \forall j=1 \dots k \end{array} \right.$$

Tale che:

$$f(x^*, \alpha^*) \leq f(x^*, \alpha) \quad \forall \alpha \geq 0: \sum_{j=1}^k \alpha_{ij} = 1 \quad \forall i=1 \dots e$$

$$f(x^*, \alpha^*) \leq f(x, \alpha^*) \quad \forall x \in \mathbb{R}^{kn}$$

Note: K-means non garantisce di trovare l'ottimo globale, perché il problema è convesso in una sola direzione e non in tutte.

CLUSTER CON $\|\cdot\|_1$

δ^- è lo stesso sistema, ma con la norma 1:

$$\begin{cases} \min \sum_{i=1}^e \min_{j=1 \dots k} \|p_i - x_j\|_1 \\ x_j \in \mathbb{R}^m \quad \forall j=1 \dots k \end{cases}$$

Prendiamo ora e numeri reali $a_1 < \dots < a_e$, e vogliamo cercare la soluzione ottima del seguente sistema, con $k=1$.

$$\begin{cases} \min_{x \in \mathbb{R}} \sum_{i=1}^e |x - a_i| = f(x) \end{cases}$$

δ^- è il **mediano** fra $(a_1 \dots a_e)$: $\begin{cases} a_{e+1/2} & \text{e dispari} \\ \frac{a_{e/2} + a_{1+e/2}}{2} & \text{e pari} \end{cases}$

Objective function:

$$f(x) = \begin{cases} -ex + \sum_{i=1}^e a_i & \text{se } x < a_1 \\ (2-e)x + \sum_{i=2}^e a_i - a_1 & \text{se } x \in [a_1, a_2] \\ \dots \\ (2r-e)x + \sum_{i=r+1}^e a_i - \sum_{i=1}^r a_i & \text{se } x \in [a_r, a_{r+1}] \\ \dots \\ (2-e)x + a_e + \sum_{i=1}^{e-1} a_i & \text{se } x \in [a_{e-1}, a_e] \\ ex - \sum_{i=1}^e a_i & \text{se } x > a_e \end{cases}$$

Questa funzione è **convessa** e **lineare a tratti**.

Esempio

$$f(x) = |x-2| + |x-4| + |x-10| \quad \text{con } \ell=3$$

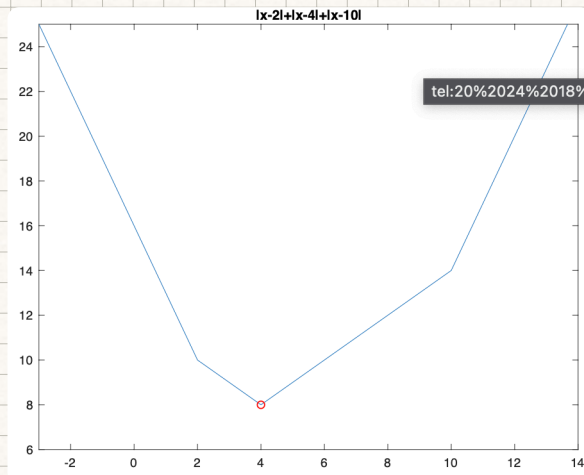
$$1) \text{ se } x < 2 : 2-x + 4-x + 10-x = -3x+16 = f(x)$$

$$\text{se } x \in [2, 4] : x-2 + 4-x + 10-x = -x+12 = f(x)$$

$$\text{se } x \in [4, 10] : x-2 + x-4 + 10-x = x+4 = f(x)$$

$$\text{se } x > 10 : x-2 + x-4 + x-10 = 3x-16 = f(x)$$

Usando la regola, non facciamo conti.



If $k > 1$ (at least two clusters), then the problem is **nonconvex and nonsmooth**:

$$\begin{cases} \min_x \sum_{i=1}^{\ell} \min_{j=1, \dots, k} \|p_i - x_j\|_1 \\ x_j \in \mathbb{R}^n \quad \forall j = 1, \dots, k \end{cases} \quad (8)$$

Theorem

Problem (8) is equivalent to the following problem:

$$\begin{cases} \min_{x, \alpha} \sum_{i=1}^{\ell} \sum_{j=1}^k \alpha_{ij} \|p_i - x_j\|_1 \\ \sum_{j=1}^k \alpha_{ij} = 1 \quad \forall i = 1, \dots, \ell \\ \alpha_{ij} \geq 0 \quad \forall i = 1, \dots, \ell, j = 1, \dots, k \\ x_j \in \mathbb{R}^n \quad \forall j = 1, \dots, k. \end{cases} \quad (9)$$

Il problema (9) del teorema è equivalente al seguente problema **bilineare** (ovvero differenziabile) e **non-convesso**

$$\left\{ \begin{array}{l} \min_{x, \alpha, u} \sum_{i=1}^{\ell} \sum_{j=1}^k \sum_{h=1}^n \alpha_{ij} u_{ijh} \\ u_{ijh} \geq (p_i)_h - (x_j)_h \quad \forall i = 1, \dots, \ell, j = 1, \dots, k, h = 1, \dots, n \\ u_{ijh} \geq (x_j)_h - (p_i)_h \quad \forall i = 1, \dots, \ell, j = 1, \dots, k, h = 1, \dots, n \\ \sum_{j=1}^k \alpha_{ij} = 1 \quad \forall i = 1, \dots, \ell \\ \alpha_{ij} \geq 0 \quad \forall i = 1, \dots, \ell, j = 1, \dots, k \\ x_j \in \mathbb{R}^n \quad \forall j = 1, \dots, k. \end{array} \right.$$