



K-means clustering - 01

tds sheets | clustering

CONTEXT

We consider that we have N data points in a simple d -dimensional Euclidean space

$$\{x_1, x_2, \dots, x_N\}$$

We want both to identify clusters among these data and to get the centers of each cluster. For a given clustering with K clusters, we define

$$c : \{1, 2, \dots, N\} \rightarrow \{1, 2, \dots, K\} \quad \text{and} \quad \alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)$$

respectively **the mapping that assigns to each point a cluster** and **the centers of the clusters**. The criteria we want to use to assess the quality of our clustering is the **inertia**

$$I(c, \alpha) = \sum_{k=1}^K \sum_{n, c(n)=k} \|x_n - \alpha_k\|^2 = \sum_{n=1}^N \|x_n - \alpha_{c(n)}\|^2$$

inertia sum over the K clusters sum of squared distances to center within cluster k sum over the N points squared distance between point n and its cluster's center

that is the sum of squared distance between points and the center of their respective cluster. We want the inertia to be as small as possible.

INTUITION

As indicated by its name, **K-means assumes there are K clusters in the data**.

For these K unknown clusters, **if we knew the centers we could assign each of our points to a cluster** (choosing the cluster whose center is the closest).

In the other way, **if we knew the points assigned to each cluster we could easily compute the centers** (taking the mean of the points belonging to each cluster).

The idea behind K-means is to use an EM algorithm (Expectation-Maximisation) that will alternatively:

- ➡ fix cluster's centers and define new clusters based on it \textcircled{c} $\textcircled{\alpha}$
- ➡ fix clusters and compute new cluster's centers based on it \textcircled{c} $\textcircled{\alpha}$

REMARKS

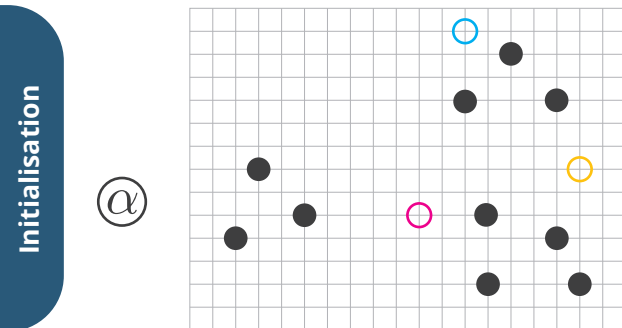
With this algorithm, we are ensured that the **inertia will decrease** at each step. However, **we can be stuck into a local minimum** instead of reaching the global minimum. As for any algorithm that can reach local optimum, the **initialisation has a huge importance** in the final result.



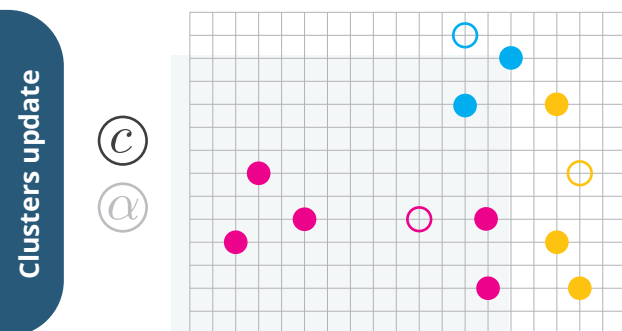
K-means clustering - 02

tds sheets | clustering

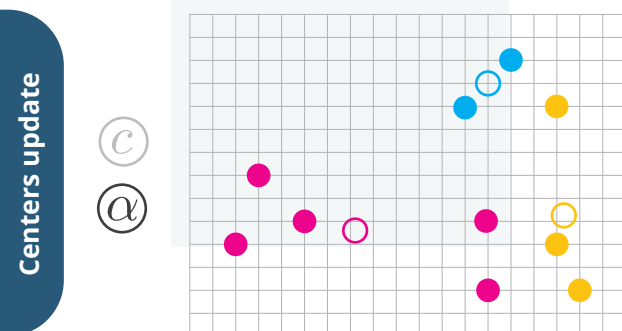
EXAMPLE WITH K=3



- ➔ centers are randomly initialised $\alpha^{(0)}$
- ➔ no clusters are defined yet —
- ➔ inertia can't be computed for now —

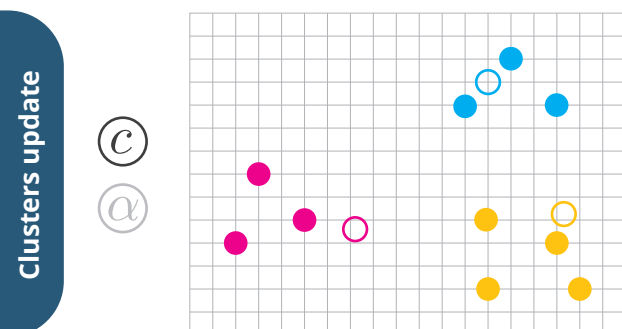


- ➔ centers are fixed $\alpha^{(0)}$
- ➔ clusters are updated — $\rightarrow c^{(1)}$
- ➔ inertia can be computed $I(c^{(1)}, \alpha^{(0)})$



- ➔ centers are updated $\alpha^{(0)} \rightarrow \alpha^{(1)}$
- ➔ clusters are fixed $c^{(1)}$
- ➔ inertia decrease $I(c^{(1)}, \alpha^{(1)}) \leq I(c^{(1)}, \alpha^{(0)})$

$$\begin{aligned} \forall k \quad C_k^{(1)} &= \{n, c^{(1)}(n) = k\} & \alpha_k^{(1)} &= \frac{1}{C_k^{(1)}} \sum_{n \in C_k^{(1)}} x_n \\ \Rightarrow \sum_{n \in C_k^{(1)}} \|x_n - \alpha_k^{(1)}\|^2 &\leq \sum_{n \in C_k^{(1)}} \|x_n - \alpha_k^{(0)}\|^2 \\ \Rightarrow I(c^{(1)}, \alpha^{(1)}) &\leq I(c^{(1)}, \alpha^{(0)}) \end{aligned}$$



- ➔ centers are fixed $\alpha^{(1)}$
- ➔ clusters are updated $c^{(1)} \rightarrow c^{(2)}$
- ➔ inertia decrease $I(c^{(2)}, \alpha^{(1)}) \leq I(c^{(1)}, \alpha^{(1)})$

$$\begin{aligned} \forall n \quad c^{(2)}(n) &= \arg \min_k \|x_n - \alpha_k^{(1)}\|^2 \\ \Rightarrow \|x_n - \alpha_{c^{(2)}(n)}^{(1)}\|^2 &\leq \|x_n - \alpha_{c^{(1)}(n)}^{(1)}\|^2 \\ \Rightarrow I(c^{(2)}, \alpha^{(1)}) &\leq I(c^{(1)}, \alpha^{(1)}) \end{aligned}$$



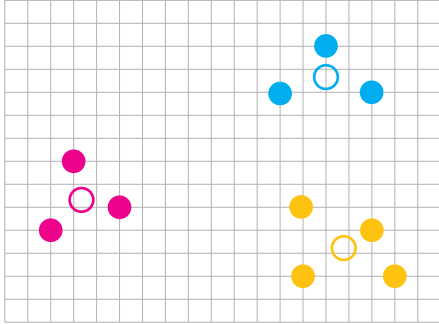
K-means clustering - 03

tds sheets | clustering

EXAMPLE WITH K=3 (SUITE)

Centers update

c
 α



- ➡ centers are updated
- ➡ clusters are fixed
- ➡ inertia decrease

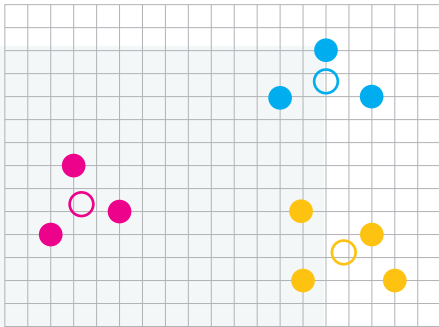
$$\alpha^{(1)} \rightarrow \alpha^{(2)}$$

$$c^{(2)}$$

$$I(c^{(2)}, \alpha^{(2)}) \leq I(c^{(2)}, \alpha^{(1)})$$

End

c
 α



- ➡ centers are fixed
- ➡ clusters are already up-to-date
- ➡ inertia has reached a local minimum $I(c^{(2)}, \alpha^{(2)})$

$$\alpha^{(2)}$$

$$c^{(2)}$$

$$I(c^{(2)}, \alpha^{(2)})$$

HOW TO CHOOSE K?

As K-means algorithm is defined for a given K, one main question is **"How to choose the value of K?"**. We can easily verify that **the higher K is, the lower optimal inertia will be**. So we can't just choose the value K that minimise inertia, it won't have any sense.

Instead, a good way to select a value of K is to use the **elbow method**. The idea is to plot **the graph of the best obtained inertia with respect to the value of K**. Then, if there is one, choose the value of K "at the elbow". If such elbow appears, it means indeed that after it, adding more clusters brings less value and so is not interesting.

