

k-means

Reviews

Unsurprised learning

$$x_1, x_2, \dots, x_n \quad x_i \in \mathbb{R}^d$$

Goal: Assign labels to each data point y_1, y_2, \dots, y_n $y_i \in \{1, \dots, k\}$
 $k = \text{number of clusters}$

• Distance-based clustering



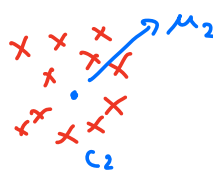
Given centers \rightarrow Assignment is easy



Define: $\mu = \frac{1}{n} \sum x_i$ (Centroid)

$$F(c_1, c_2, \dots, c_k) = \sum_{j=1}^k \left(\sum_{x_i \in C_j} \|x_i - \mu_j\|_2^2 \right)$$

k -cluster Variance in cluster j



If we have centroid, it's quite easy to cluster. But given the dataset, how to choose k suitable centroid?
 \rightarrow Lloyd's algorithm.

Lloyd's algorithm (k-means)

• Initialize: pick random centers $\mu_1, \mu_2, \dots, \mu_k$

① Assignment $z_i = \arg \min_{j=1, \dots, k} \|x_i - \mu_j\|$

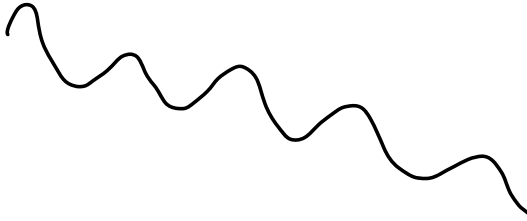
• ② Update centers $\mu_i = \frac{1}{n} \sum_{x_i \in C_i} x_i$

Step 1

Calculate the distance between each point and each centroid, and the point belongs to the group of certain centroid once it has the closest distance!

Step 2. Update the centers using existing points belongs to current center.

not a convex optimization problem.



Some questions:

- ① How to initialize?
- ② Do we always find optimal clustering?
- ③ Convergence?
- ④ For what kind of data is this suited?

Question 1

$$F(C_1, C_2, \dots, C_k) = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|_2^2,$$

$\mu_1, \dots, \mu_k \equiv$ centroids of C_1, \dots, C_k respectively

• For a fixed K , prove/argue that F achieve a minimum value.

Solution: • Finite ways to assign points into cluster C_1, \dots, C_k

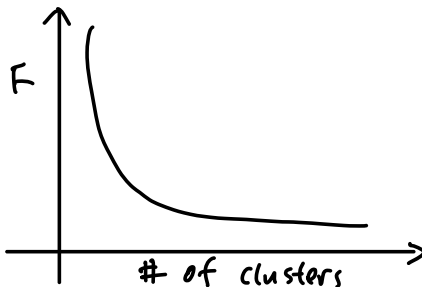
• Among these clusterings, pick the one with smallest value for F

Question 2

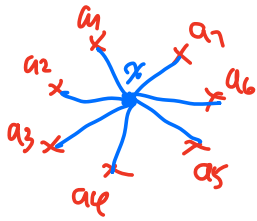
What is the minimum value if $k=n$? ($n \equiv$ number of point)

Solution: Zero. Each point is its own cluster.

Elbow method



Lemma : Let $\{a_1, \dots, a_n\}$ be a set of n points. Let x be an arbitrary point.



Find the blue point, which has the minimum distances among all red points. \rightarrow the blue point must be the centroid!

Prove:
$$\sum_{i=1}^n \|a_i - x\|^2 = \sum_{i=1}^n \|a_i - c\|^2 + n \|c - x\|^2$$

 where c is the centroid of $\{a_1, \dots, a_n\}$

Proof:
$$\sum_{i=1}^n \|a_i - x\|^2$$

$\Rightarrow \sum_{i=1}^n \|a_i - x + (-c)\|^2$

$\Rightarrow \sum_{i=1}^n \underbrace{\|a_i - c\|}_A + \underbrace{\|c - x\|}_B$

$\Rightarrow \sum_{i=1}^n \|a_i - c\|^2 + \|c - x\|^2 + 2 \underbrace{(a_i - c)^T (c - x)}$

$= n \|c - x\|^2 + \sum_{i=1}^n \|a_i - c\|^2 + \underbrace{\sum_{i=1}^n 2(a_i - c)^T (c - x)}_0$ $c = \frac{1}{n} \sum a_i$

$$\begin{aligned} \sum_{i=1}^n (a_i - c) &= \sum_{i=1}^n \left(a_i - \frac{1}{n} \sum_{j=1}^n a_j \right) \\ &= \sum_{i=1}^n a_i - \sum_{i=1}^n \left(\frac{1}{n} \sum_{j=1}^n a_j \right) \\ &= n c - n c = 0 \end{aligned}$$

Conclusion: $\sum_{i=1}^n \|a_i - x\|^2 \geq \sum_{i=1}^n \|a_i - c\|^2$

Example

N points in five tight clusters.

