

K-means, K-Medoid, Mean-shift Clustering

Pattern Recognition, Jul-Nov 2019

Indian Institute of Technology Madras

September 6, 2019

K Means Clustering

- 1 Randomly initialize centroids of K Gaussian
- 2 Assign data-points to to C_1 , C_2 or C_3 according to each one's distance
- 3 Reestimate means.
- 4 Repeat 2 & 3 until convergence. (after that there is no point in changing the clusters)

$$\begin{aligned}\vec{\mu}_k &= [\mu_{k1}, \mu_{k2} \dots \mu_{kd}]^t \\ &= \text{mean of } k^{th} \text{ cluster}\end{aligned}$$

Our goal is :

- ▶ To find an assignment of data-points to clusters
- ▶ To find μ_k such that sum of distances of each data point closest to the $\vec{\mu}_k$ is minimum

Let us define γ_{nk} for each \vec{x}_n such that $\gamma_{nk} \in \{0, 1\}$

$\gamma_{nk} = 1$, if \vec{x}_n belongs to k^{th} cluster

$\gamma_{nj} = 0$, if $j \neq k$

Distortion

$$D = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \|\vec{x}_n - \vec{\mu}_k\|^2 \quad \text{objective is to minimise } D.$$

But we need to address few issues before moving forward.

- ▶ γ_{nk} is not known
- ▶ More variables than equations

$$\begin{aligned}
D &= \mathbf{E} [d(\vec{x}_n, \vec{\mu}_k)] \\
&= \sum_{k=1}^K \int_{\vec{x}_n \in C_k} d(\vec{x}_n, \vec{\mu}_{k(\text{old})}) p(\vec{x}_n, \vec{\mu}_{k(\text{old})}) d\vec{x} \\
&= \sum_{k=1}^K P(\mu_k) \int_{\vec{x}_n \in C_k} d(\vec{x}_n, \vec{\mu}_k) p(\vec{x}_n / \vec{\mu}_k) d\vec{x} \\
&= \sum_{k=1}^K D_k
\end{aligned}$$

$q(\vec{x}_n) = \vec{\mu}_k$ if $\vec{x}_n \in C_k$ also called as vector quantization
 happens if $d(\vec{x}_n, \vec{\mu}_k) \leq d(\vec{x}_n, \vec{\mu}_j) \quad k \neq j \quad 1 \leq k, j \leq K$

Derivating w.r.t μ_k

$$\nabla_{\mu_k} D = \frac{1}{N} \sum_{\vec{x}_n \in C_k} \nabla_k (\vec{x}_n - \vec{\mu}_k)^t (C_k^{-1}) (\vec{x}_n - \vec{\mu}_k)$$

$$\sum_{\vec{x}_n \in C_k} \vec{x}_n = N_k \vec{\mu}_k$$

$$\vec{\mu}_k = \frac{1}{N_k} \sum_{\vec{x}_n \in C_k} \vec{x}_n$$

$$C_k = \frac{1}{N_k} \sum_{\vec{x}_n \in C_k} (\vec{x}_n - \vec{\mu}_k)(\vec{x}_n - \vec{\mu}_k)^t$$

K-Means Clustering

We need to prove that $D_{i+1} \leq D_i$ where D_{i+1} is the Distortion in $(i+1)^{th}$ iteration
 Let $\vec{Z} \in \{\vec{z}_1, \vec{z}_2 \dots \vec{z}_k\}$ and $L(\vec{x}_i, \vec{z}_k)$ be an assignment that minimises distortion.

$$D_i = E [L(\vec{x}_n, \vec{z}_k)], n = 1, 2 \dots N$$

\vec{Z}_k set of prototypes, $k = 1, 2 \dots K$

The $(i+1)^{th}$ distortion is getting computed using z_k of the i^{th} iteration

Now, Using Euclidean Mean

$$E[\vec{z}] = \frac{1}{2} \sum_{n=1}^N ((\vec{x}_n - \vec{z}_{nk} S_{ni}(\vec{z})))^2$$

$S_{ni}(\vec{z})$ is the closest centroid to example \vec{x}_n in i^{th} iteration.

But we don't have $S_i(\vec{z})$ so we define an auxiliary function:

$$Q(\vec{z}, \vec{z}') = \frac{1}{2} \sum_{n=1}^N (\vec{x}_n - \vec{z}'_{S_i(\vec{z})})^2$$

Now derivating $Q(\vec{z}, \vec{z}')$ w.r.t \vec{z}'_k and equating it to 0.

$$\vec{z}'_k = \frac{1}{N_k} \sum_{\gamma_{nk}=S_i(\vec{z})} \vec{x}_n \gamma_{nk} \quad (N_k = \text{no. of pts. assigned to } k^{\text{th}} \text{ cluster})$$

To prove: $E[\vec{z}'] - E[\vec{z}] \leq 0$

$$\text{LHS} = E[\vec{z}'] - Q(\vec{z}', \vec{z}) + Q(\vec{z}', \vec{z}) - Q(\vec{z}, \vec{z}) \quad (\text{Writing } E[\vec{z}] = Q(\vec{z}, \vec{z}))$$

We can see that in LHS, $Q(\vec{z}', \vec{z}) - Q(\vec{z}, \vec{z}) \leq 0$ because using \vec{z} we are creating \vec{z}' i.e \vec{z}' is the best assignment obtained by minimizing the auxiliary function. Hence $\text{LHS} \leq 0$. Hence proved.

K-Means Vs K-medoid

- ▶ Centroid is an actual vector coming from training data.
- ▶ We first compute k-means centroid, which might not be an actual datapoint. So, the nearest data point to this centroid is called k-medoid centroid.
- ▶ k-Medoid is computationally expensive
- ▶ k-Medoid is more robust to outliers – in that outliers are not taken into account during the computation of the medoid.
- ▶ In K-Medoid it is possible to have a cluster with one point, while in K-Means it is possible to have a cluster with no points – why?

K-Means for representing the distribution of variables for a class

- ▶ The training data of a class is used to find K prototypes for every class.
- ▶ During testing, given a set of $\mathcal{D}_{test} = \{\bar{x}_{t1}, \dots \bar{x}_{tN_t}\}$ test vectors corresponding to a given class, the Distortion is computed with the prototypes of every class, and the class which yields smallest distortion is identified as the class corresponding to the set of test vectors.
- ▶ Voting can also be used – although leads to quantization.
- ▶ K-Means clustering is the pre-cursor for GMMs to be discussed later.

Mean shift clustering

The mean update equation for mean shift clustering is given by

$$m(\bar{x}) = \frac{\sum_{\bar{x}_i \in N(x_i)} K(\bar{x} - \bar{x}_i) \cdot x_i}{K(x - \bar{x}_i)}$$

where k is a kernel given by

$$K(x_i - \bar{x}) = \frac{e^{-||x_i - \bar{x}||^2}}{h^2}$$

here h is band width parameter

Drawbacks of K Means clustering algorithm

- ▶ The clustering is hard, i.e., a point can belong to only one cluster
- ▶ Solution: Soft clustering of points
- ▶ Each cluster is replaced by a Gaussian
- ▶ Every point belong to all the clusters to some extent
- ▶ The belongingness of a point \bar{x} to a cluster k is given by Gaussian PDF.

$$p(\bar{x}|z_k = 1) = \sum_{k=1}^K \pi_k N(\bar{x}|\bar{\mu}_k, \Sigma_k)$$

where $\sum_{k=1}^K \pi_k = 1$