

1 Clustering

1.1 K Means

a. Definition

Given a set of observations (x_1, x_2, \dots, x_n) , where each observation is a dimensional real vector, k-means clustering aims to partition the n observations into k ($k \leq n$) sets $S = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within cluster sum of squares. Formally, the objective is to find:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^k \sum_{x \in S_i} \|x - u_i\|^2 = \underset{S}{\operatorname{argmin}} \sum_{i=1}^k |S_i| \operatorname{Var} S_i$$

b. Algorithm

- 1) Give the initial guess of k means m_1, \dots, m_k
- 2) Assign each observation to the cluster whose mean has the least squared Euclidean distance.
- 3) Calculate the new means to be the centroids of the observations in the new clusters.
- 4) $m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$

c. Time Complexity

$O(nkd)$, where n is the number of d dimensional vectors, k is the number of clusters and i is the number of iterations need till convergence.

2 Gaussian Mixture

a. Idea and Definition

1) In K means clustering, one sample point exclusively belongs to one cluster. In other words, we assign a sample point to a cluster with probability 1. In Mixture model, we assign sample point i to a cluster k with the probability r_{ik} , with

$$\sum_k r_{ik} = 1$$

The r_{ik} also follows the fact

$$\sum_i \sum_k r_{ik} = \sum_i 1 = N$$

By changing the order of summation

$$\sum_i \sum_k r_{ik} = \sum_k \sum_i r_{ik}$$

Define the weight of cluster: $w_k = \sum_i r_{ik} / N = \sum_k \omega_k * N = N$
So

$$\sum_k w_k = 1$$

We can also interpret w_k as a prior distribution of a sample point being assigned to cluster k.

2) And for each cluster k, we define the probability of having a sample point i at x_i use a normal distribution $N(x_i|u_k, \Sigma_k)$

Diagram:

3) The $r_{ik}\pi_k$ and $N(x_i|u_k, \Sigma_k)$ are connected with Bayesian rule

$$\begin{aligned} P(A|B) &= \frac{P(A)P(B|A)}{P(B)} \\ &= \frac{P(A)P(B|A)}{\sum_c P(C)|P(B|C)} \end{aligned}$$

According this rule, we have the following

$$\begin{aligned} &P(X_i = x_i \text{ and } X_i \text{ in cluster k}) \\ &= P(X_i \text{ in cluster k})P(X_i = x_i \text{ given } X_i \text{ in cluster k}) \\ &= P(X_i \text{ in cluster k} | X_i = x_i)P(X_i = x_i) \end{aligned}$$

So

$$\begin{aligned} &P(X_i \text{ in cluster k} | X_i = x_i) \\ &= P(X_i \text{ in cluster k})P(X_i = x_i \text{ given } X_i \text{ in cluster k}) / (X_i = x_i) \end{aligned}$$

Namely,

$$r_{ik} = \frac{\pi_k N(x_i|u_k, \Sigma_k)}{\sum_j \pi_j N(x_i|u_j, \Sigma_j)}$$

4) Our goal is the find u_k, Σ_k, w_k .

b. Cost function and Minimization

For a given point x_i , the likelihood function is

$$p(x_i) = \sum_k \pi_k N(x_i|u_k, \Sigma_k)$$

The likelihood function for the whole sample is

$$\Pi_{i=1}^N p(x_i) = \Pi_{i=1}^N \sum_k \pi_k N(x_i|u_k, \Sigma_k)$$

The goal is to minimize the negative of Log Likelihood

$$L = - \sum_{i=1}^N \ln \left(\sum_k \pi_k N(x_i|u_k, \Sigma_k) \right)$$

1) Take the derivative with respect to u_k

$$dL/du_k = \sum_i \frac{\pi_k N(x_i|u_k, \Sigma_k)}{\sum_j \pi_j N(x_i|u_j, \Sigma_j)} \Sigma^{-1}(x_i - u_k)$$

We found that the term

$$\frac{\pi_k N(x_i|u_k, \Sigma_k)}{\sum_j \pi_j N(x_i|u_j, \Sigma_j)}$$

is exactly r_{ik}

Let the derivative equal to zero, we have

$$u_k = \frac{1}{N_k} \sum_i r_{ik} x_i \quad (N_k = \sum_i r_{ik})$$

2) Taking the derivative with respect to Σ_k gives

$$\Sigma_k = 1/N_k \sum_i r_{ik} (x_i - u_k)(x_i - u_k)^T$$

3) Taking the derivative with respect to π_k gives

$$\pi_k = \frac{N_k}{N}$$

We see $u_k, \Sigma_k, w_k, r_{ik}$ are mutually dependent, therefore we need to solve this iteratively.

c. Algorithm

1) Initialize cluster prior assignment $\pi_k = P(z_i = k)$

2) Given an observation x_i from cluster k , calculate $P(x_i|z = k, u_k, \Sigma_k) = N(x_i|u_k, \Sigma_k)$

3) E step

Given an observation x_i , calculate r_{ik}

$$r_{ik} = \frac{\pi_k N(x_i|u_k, \Sigma_k)}{\sum_j \pi_j N(x_i|u_j, \Sigma_j)}$$

4) M step

$$N_k = \sum_i^N r_{ik}$$

$$\hat{u}_k = \sum_i^N \frac{r_{ik}}{N_k} x_i$$

$$\hat{\Sigma}_k = \sum_i^N \frac{r_{ik}}{N_k} (x_i - \hat{u}_k)(x_i - \hat{u}_k)^T$$

$$\pi_k = \frac{N_k}{N}$$

5) Repeat 3) and 4)

d. Connection to K means

In order to easily see how Gaussian mixture clustering relates to K means, we need to introduce another latent variable Z and consider the log likelihood function of the complete data set (X, Z) . We discussed the probability to assign a sample point to cluster k as π_k , now we denote this assignment using an indicator random variable $Z^{(i)} = Z_k = (z_{k1}, z_{k2}, \dots, z_{kK})^T$ Where

$$\begin{aligned} z_{kj} &= 1 \text{ when } j = k \\ &= 0 \text{ otherwise} \end{aligned}$$

In other words, only Z_k is a K dimensional vector with only k th component being 1, other components are zero. For example, $Z_1 = (1, 0, 0, \dots, 0)^T$, $Z_k = (0, 0, 0, \dots, 1 \text{ (kth element)}, \dots, 0)^T$
And

$$p(Z^{(i)} = Z_k) = \pi_k = \prod_j \pi_j^{z_{kj}}$$

We rewrite the likelihood function given X and Z

$$L = \prod_{i=1}^N \prod_k^K \prod_{j=1}^K \pi_j^{z_{kj}} N(x_i | u_j, \Sigma_j)^{z_{kj}}$$

If we let $\pi_k = 1/K$, and the covariance matrix $= \sigma^2 \mathbf{I}$

$$\text{Log} L = \sum_i \sum_k (\log \pi_k - \frac{1}{2} \frac{1}{\sigma^2} \|x_i - u_k\|^2)$$

This reduces to the K means cost function

Now let σ goes to 0

$$\begin{aligned} r_{ik} &= \frac{\pi_k N(x_i | u_k, \Sigma_k)}{\sum_j \pi_j N(x_i | u_j, \Sigma_j)} \\ &= \frac{\pi_k \exp(-\frac{(x_i - u_k)^2}{2\sigma^2})}{\sum_j \pi_j \exp(-\frac{(x_i - u_j)^2}{2\sigma^2})} \end{aligned}$$

When σ goes to zero, the exponential term that decays the slowest survives, and the term that decays the slowest is the one that minimize $\|x_i - u_k\|$. Let u^* be the u_k that minimize $\|x_i - u_k\|$, then

$$\begin{aligned} r_{ik} &= 1 \text{ for } k = k^* \\ &= 0 \text{ otherwise} \end{aligned}$$

This reduces to the k means where a sample point i is solely assigned to a cluster k .