

9. EM Algorithm

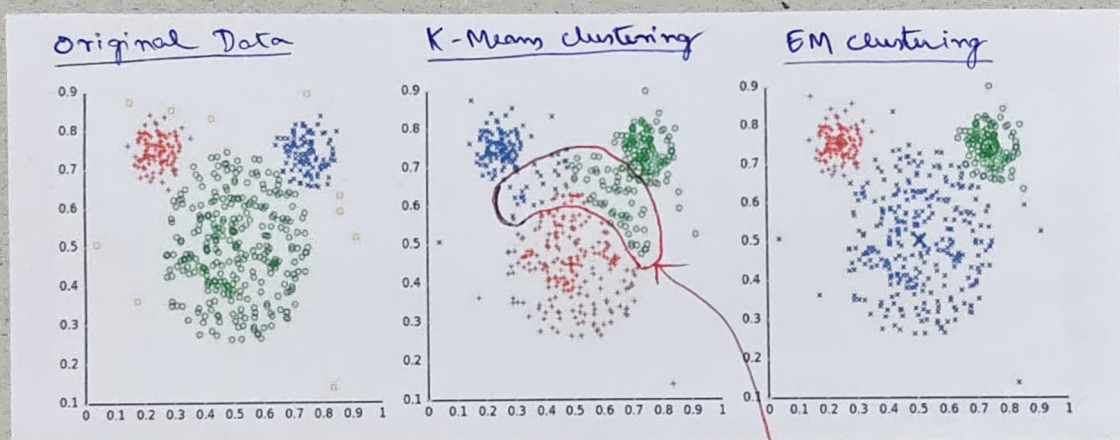
EM stands for Expectation Maximization. This can be used for Probabilistic clustering or Soft-clustering.

Previously we've seen the K-Means algorithm, which is used for Hard-clustering.

What is Probabilistic clustering?

Previously, we assigned each point to a Unique cluster. But the problem is, there might be some points which cannot be assigned to any particular cluster, because those points seemingly belong to either clusters.

In Probabilistic clustering, we calculate the probability that a data point belongs to a cluster! This is why, we also call this as Soft clustering.



Plot: Different Cluster Analysis results on "mouse" data set.

In K-Means clustering, these data points are misclassified.

The Probabilistic clustering is clearly able to recover the original clusters to a great degree. Therefore, EM algorithm can yield a much better clustering performance than the K-Means algorithm.

To understand EM algorithm, let us consider a Gaussian cluster model. (i) we assume that the clusters are generated by Gaussian distributions. With probability P_i (prior probability of i^{th} cluster), we generate a sample \bar{x} from Gaussian cluster i .

$$(i) \mathcal{N}(\bar{\mu}_i, \sigma^2 \mathbf{I})$$

Covariance Matrix
Centroid cluster

$$\Rightarrow \underbrace{\mathcal{N}(\bar{\mu}_1, \sigma^2 \mathbf{I}), \mathcal{N}(\bar{\mu}_2, \sigma^2 \mathbf{I}), \dots, \mathcal{N}(\bar{\mu}_K, \sigma^2 \mathbf{I})}_{\text{Assuming there are } K \text{ clusters}}$$

The Probability Density Function (PDF) is given as

$$f_x(\bar{x}) = \sum_{i=1}^K P_i \times \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{x} - \bar{\mu}_i\|^2}$$

where, P_1, P_2, \dots, P_K are Prior Probabilities of clusters.

This is termed as Gaussian Mixture model, as there are K Gaussians, one for each cluster.

Consider now M data points.

$$\bar{x}(1), \bar{x}(2), \dots, \bar{x}(M).$$

How do we cluster this data?

(i) We desire to estimate Centroids $\bar{\mu}_i$

As well as the cluster Assignments $\alpha_i(j)$

Here, performing direct ML estimation is mathematically intractable (NOT possible)

In this problem, there are 2 unknowns.

① Centroid $\bar{\mu}_i$

② Cluster assignment $\alpha_i(j)$

However, if cluster assignment $\alpha_i(j)$ is known, problem is simple. For example: we have $m=8$ data points, out of which

$$\bar{x}(1), \bar{x}(3), \bar{x}(5), \bar{x}(8) \in \text{cluster 1}$$

$$\bar{x}(2), \bar{x}(4), \bar{x}(6), \bar{x}(7) \in \text{cluster 2}$$

Then the centroids are

$$\hat{\mu}_1 = \frac{\bar{x}(1) + \bar{x}(3) + \bar{x}(5) + \bar{x}(8)}{4}$$

(Average of points in cluster 1)

$$\hat{\mu}_2 = \frac{\bar{x}(2) + \bar{x}(4) + \bar{x}(6) + \bar{x}(7)}{4}$$

(Average of points in cluster 2)

Recall, clustering problem is Unsupervised Learning (No Labels).
(i) there is NO cluster information.

Here, we introduce the concept of missing data or latent information.

The cluster assignment variable

$$\alpha_i(j) = \begin{cases} 1 & \text{if } \bar{x}(j) \in \text{Cluster } i \\ 0 & \text{if } \bar{x}(j) \notin \text{Cluster } i \end{cases}$$

is missing here.

$$\bar{x}(1), \bar{x}(2), \dots, \bar{x}(m)$$

$$\alpha_i(j), \quad i=1, 2, \dots, K$$

$$j=1, 2, \dots, m$$

Complete Data

We have to estimate the Missing data as well as the Centroids. To do these, we have to first start with cost function, and that cost function comes from the likelihood.

The Likelihood of the complete data

$$\prod_{j=1}^M \prod_{i=1}^K \left(p_i \times \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{x}(j) - \bar{\mu}_i\|^2} \right) \alpha_i(j)$$

Product

Missing Data

This is the joint PDF of the points, considering the availability of the missing data, which can be treated as the Likelihood function.

The Log-Likelihood of the complete Data

$$\sum_{j=1}^M \sum_{i=1}^K \alpha_i(j) \left(p_i - \frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\bar{x}(j) - \bar{\mu}_i\|^2 \right)$$

Product becomes Sum for Log.

EM Algorithm

Week 6 : Session 2

The EM algorithm proceeds iteratively.

(i) in each iteration, we perform

Expectation step \rightarrow E-Step

Maximization step \rightarrow M-Step

Let us consider the centroids in $(l-1)$ iterations.

$$\mu_0^{(l-1)}, \mu_1^{(l-1)}, \dots, \mu_K^{(l-1)}$$

At the end of $(l-1)^{th}$ iteration, we have these centroids, which are nothing but cluster means. Recall, we have Gaussian clusters (i) Gaussian Mixture model. We are modeling the clusters as Gaussian. So, the centers are nothing but the mean of the different Gaussian components in the mixture.

The Expected value of the Log Likelihood in iteration l is,

$$\sum_{j=1}^M \sum_{i=1}^K \alpha_i^{(l)}(j) \left\{ \ln p_i - \frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\bar{\pi}(j) - \bar{\mu}_i\|^2 \right\}$$

$Q(\bar{\mu})$

$$\rightarrow E\{\alpha_i(j)\} = \alpha_i^{(l)}(j)$$

Expected value of $\alpha_i(j)$ in l^{th} iteration

How to calculate $\alpha_i^{(l)}(j)$?

Recall, $\alpha_i^{(l)}(j) = 1$ if $\bar{\pi}(j) \in C_i$

$$= \Pr(\bar{\pi}(j) \in C_i) \times 1$$

$$+ \Pr(\bar{\pi}(j) \notin C_i) \times 0$$

$$= \Pr(\bar{\pi}(j) \in C_i)$$

$$= \Pr(C_i | \bar{\pi}(j))$$

$$= \frac{\Pr(\bar{\pi}(j) | C_i) \cdot P(C_i)}{\sum_{k=1}^K \Pr(\bar{\pi}(j) | C_k) \cdot P(C_k)}$$

(Bayes Rule)

Prior Probability

$$= \frac{P_i \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\pi}(j) - \bar{\mu}_i^{(l-1)}\|^2}}{\sum_{k=1}^K P_k \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\pi}(j) - \bar{\mu}_k^{(l-1)}\|^2}}$$

$$= E\{\alpha_i(j)\}$$

M-Step

The M-Step is the Maximization step (ii) Maximize the Expected value of Likelihood. To determine $\bar{\mu}_j$, differentiate w.r.t. $\bar{\mu}_j$ and set equal to zero.

(ii) To maximize, take the Gradient of the Expected value of the Log Likelihood and set equal to zero.

$$\nabla_{\bar{\mu}_i} Q(\bar{\mu}) = 0$$

Solve for $\bar{\mu}_i^{(l)}$, we get

Centroid of i^{th} cluster in iteration l .

$$\bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{x}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}$$

Now, compare the expressions in K-Means and EM,

$$\text{EM: } \bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{x}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}$$

$$\text{K-Means: } \bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{x}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}$$

Same Expression!

Both are Identical!

Then what is the difference?

| <u>EM</u> | <u>K-Means</u> |
|--|---|
| <ul style="list-style-type: none"> EM performs Weighted average $\alpha_i^{(l)}(j)$ - Soft clustering $\alpha_i^{(l)}(j) \in [0, 1]$ (i) can take any value in the interval $[0, 1]$ Example: 0.95 | <ul style="list-style-type: none"> K-Means performs Hard average $\alpha_i^{(l)}(j)$ - Hard clustering $\alpha_i^{(l)}(j) \in \{0, 1\}$ (ii) can only take 2 possible values |

Therefore, in EM, $\alpha_i^{(l)}(j)$ denote probabilities.

Previously, we assumed the Prior probabilities to be known.
But here, we compute Prior probabilities.

The Prior probabilities p_i can be computed as follows.

$$p_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j)}{M}$$

Estimate of Prior probability in iteration l

Sum of all A posteriori probabilities of points in cluster i