

1 Conceptual Review of K-means and mixture modeling

1. What is the benefit of using k-means algorithm when solving a partitioning problem?

Brute-forcing (exhaustively enumerating all partitions) is NP-hard. Using K-means can solve the problem in $O(lKN)$, where $l = \#$ iterations, $K = \#$ cluster centers, and $N = \#$ points.

2. Recap the steps to k-means algorithm

- 1: Fixing the cluster centers, assign points to nearest clusters

- 2: Given the point assignments, re-estimate cluster centers

Termination: No points change clusters in next iteration

3. What is K-medoids, and how is it different from k-means? Discuss their pros and cons.

K-medoids: Cluster centers are estimated using the median point instead of mean

Benefit: Less sensitive to outliers, its "visualization" can be easier interpreted

Con: More work to compute medoid than mean

4. What is a key difference between mixture modeling and k-means?

K-means is hard assignment (each point belongs to only one cluster) while mixture modeling is soft assignment (calculates probability that a point belongs to a cluster).

2 Expectation Maximization (EM) with Gaussian Mixture Models (GMM)

Let z be a multinomial random latent variable with components z_1, z_2, \dots, z_k , where each component takes on 0 or 1 *i.e.* $P(z_j = 1)$ is the probability that a point comes from gaussian distribution j .

Let $\lambda = \mu_1, \mu_2, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k, \pi_1, \dots, \pi_k$ where $\pi_j = P(z_j=1)$.

The log likelihood $\ell(\lambda|x_1, x_2, \dots, x_m) = \sum_{i=1}^m \log P(x_i|\lambda) = \sum_{i=1}^m \log \sum_{j=1}^k \pi_j \mathcal{N}(x_i|\mu_j, \Sigma_j)$.
[Refer to EM lecture slide 17 for breaking down $P(x_i|\lambda)$]

(a) E-step: Calculate the posterior probability $P(z_j = 1|x_i, \lambda) \forall i, j$.

$$\begin{aligned} P(z_j = 1|x_i, \lambda) &= \frac{p(x_i|z_j=1, \mu_j, \Sigma_j) p(z_j=1|\pi_j)}{p(x_i|\lambda)} && \text{[Bayes Rule]} \\ &= \frac{\mathcal{N}(x_i|\mu_j, \Sigma_j) \pi_j}{\sum_{l=1}^k \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l)} && \text{[Marginalization for denominator]} \end{aligned}$$

(Note: In lecture, Pat removed the denominator and represented the proportional probability with the numerator)

(b) M-step: Apply MLE and update the parameters $\pi_j, \mu_j, \Sigma_j \forall j$.

$$\begin{aligned} \frac{\partial \ell}{\partial \mu_j} &= \frac{\partial \ell}{\partial \mu_j} \sum_{i=1}^m \log \sum_{l=1}^k \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l) && \text{[Log likelihood function]} \\ &= \sum_{i=1}^m \frac{1}{\sum_{l=1}^k \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l)} \frac{\partial \ell}{\partial \mu_j} \sum_{l=1}^k \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l) && \text{[Differentiation rule: } \frac{\partial}{\partial x} \ln(u(x)) = \frac{1}{u(x)} * u'(x)\text{]} \\ &= \sum_{i=1}^m \frac{1}{\sum_{l=1}^k \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l)} \frac{\partial \ell}{\partial \mu_j} \pi_j \mathcal{N}(x_i|\mu_j, \Sigma_j) && \text{[Eliminating terms with no } u_j\text{]} \\ &= \sum_{i=1}^m \frac{\mathcal{N}(x_i|\mu_j, \Sigma_j) \pi_j}{\sum_{l=1}^k \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l)} \frac{\partial \ell}{\partial \mu_j} \frac{(x_i - \mu_j)^2}{2\Sigma_j} && \text{[Exponential rule: } \frac{\partial}{\partial x} e^{u(x)} = e^{u(x)} * u'(x)\text{]} \\ &= \sum_{i=1}^m P(z_j = 1|x_i, \lambda) \frac{\partial \ell}{\partial \mu_j} \frac{(x_i - \mu_j)^2}{2\Sigma_j} && \text{[Substitute from E-step]} \\ &= \sum_{i=1}^m P(z_j = 1|x_i, \lambda) \Sigma_j^{-1} (x_i - \mu_j) && \text{[Derivative of log gaussian density function]} \end{aligned}$$

Setting this to 0, you get: $\mu_j = \frac{\sum_{i=1}^m P(z_j=1|x_i, \lambda) x_i}{\sum_{i=1}^m P(z_j=1|x_i, \lambda)}$

Similar calculation produces Σ_j and π_j .