

Problem 5: Justification of the K-means Algorithm (10 pts)

Let $x_1, \dots, x_n \in \mathbb{R}^p$ denote the expression levels of n genes in p samples, with x_{ij} indicating the expression of gene i in sample j . Let C_1, \dots, C_K denote the K non-overlapping clusters, each containing a subset of $\{1, \dots, n\}$, with $\cup_{k=1}^K C_k = \{1, \dots, n\}$. Let $|C_k|$ denote the size of cluster k and $m_k = (m_{k1}, \dots, m_{kp})'$ be the center of cluster k . The objective function to minimize is

$$f(C_1, \dots, C_K, m_1, \dots, m_K) = \sum_{k=1}^K \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - m_{kj})^2 = \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - m_{kj})^2.$$

1. In the first step of the $(t+1)$ -th iteration of the algorithm ($t = 0, 1, \dots$), given the clusters from the t -th iteration $C_1^{(t)}, \dots, C_K^{(t)}$. Show that updating the cluster centers as

$$m_{kj}^{(t+1)} = \frac{1}{|C_k^{(t)}|} \sum_{i \in C_k^{(t)}} x_{ij}, \quad j = 1, \dots, p$$

satisfies that

$$f(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t+1)}, \dots, m_K^{(t+1)}) \leq f(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t)}, \dots, m_K^{(t)}).$$

- ① From the problem statement

$$F(C_1, \dots, C_K, m_1, \dots, m_K) = \sum_{k=1}^K \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - m_{kj})^2 = \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - m_{kj})^2$$

Thus,

$$F(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t)}, \dots, m_K^{(t)}) = \sum_{k=1}^K \sum_{i \in C_k^{(t)}} \sum_{j=1}^p (x_{ij} - m_{kj}^{(t)})^2 = \sum_{k=1}^K \frac{1}{|C_k^{(t)}|} \sum_{i \in C_k^{(t)}} \sum_{j=1}^p (x_{ij} - m_{kj}^{(t)})^2$$

Rearranging the order of the summation

$$= \sum_{k=1}^K \sum_{j=1}^p \sum_{i \in C_k^{(t)}} (x_{ij} - m_{kj}^{(t)})^2$$

To minimize the sum of the squared distances between a group of points & its cluster center at time t , we must minimize:

$$\sum_{i \in C_k^{(t)}} (x_{ij} - m_{kj}^{(t)})^2$$

We take the partial derivative w.r.t. our current cluster center $m_{kj}^{(t)}$ & set it equal to 0

$$\frac{\partial}{\partial m_{kj}^{(t)}} \sum_{i \in C_k^{(t)}} (x_{ij} - m_{kj}^{(t)})^2 = 0$$

We move differentiation inside as a linear operator

chain rule

$$\begin{aligned} \sum_{i \in C_k^{(t)}} \frac{\partial}{\partial m_{kj}^{(t)}} (x_{ij} - m_{kj}^{(t)})^2 &= \sum_{i \in C_k^{(t)}} 2(x_{ij} - m_{kj}^{(t)}) (-1) \\ &\quad \text{Move constants outside} \\ &= -2 \sum_{i \in C_k^{(t)}} (x_{ij} - m_{kj}^{(t)}) = -2 \sum_{i \in C_k^{(t)}} x_{ij} + 2 \sum_{i \in C_k^{(t)}} m_{kj}^{(t)} \\ &\quad \text{Distribute the summation} \end{aligned}$$

$$2 \sum_{i \in C_k^{(t)}} m_{kj}^{(t)} = 2 \sum_{i \in C_k^{(t)}} x_{ij}$$

$C_k^{(t)}$ times

$$\text{Doesn't depend on } i, \text{ so } \sum_{i \in C_k^{(t)}} m_{kj}^{(t)} = m_{kj}^{(t)} + \dots + m_{kj}^{(t)}$$

$$|C_k^{(t)}| m_{kj}^{(t)} = \sum_{i \in C_k^{(t)}} x_{ij}$$

$$m_{kj}^{(t)} = \frac{1}{|C_k^{(t)}|} \sum_{i \in C_k^{(t)}} x_{ij}$$

Thus,

$$m_{kj}^{(t+1)} = \frac{1}{|C_k^{(t)}|} \sum_{i \in C_k^{(t)}} x_{ij}$$

minimizes the sum of the squared distances between a group of points & its cluster center & consequently $\mathcal{F}(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t)}, \dots, m_K^{(t)})$. $m_{kj}^{(t+1)}$ satisfies

$$\mathcal{F}(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t+1)}, \dots, m_K^{(t+1)}) \leq \mathcal{F}(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t)}, \dots, m_K^{(t)})$$

2. In the second step of the $(t+1)$ -th iteration of the algorithm, given the cluster centers from the first step $m_1^{(t+1)}, \dots, m_K^{(t+1)}$. Show that if we update the cluster membership of gene i as

$$c(i)^{(t+1)} = \arg \min_{k \in \{1, \dots, K\}} \sum_{j=1}^p (x_{ij} - m_{kj}^{(t+1)})^2,$$

the resulting updated clusters

$$C_k^{(t+1)} = \{i : c(i)^{(t+1)} = k\}, k = 1, \dots, K$$

satisfy that

$$f(C_1^{(t+1)}, \dots, C_K^{(t+1)}, m_1^{(t+1)}, \dots, m_K^{(t+1)}) \leq f(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t+1)}, \dots, m_K^{(t+1)}) .$$

- ② From problem S.1, we have shown that the new cluster centers are computed as

$$m_{kj}^{(t+1)} = \frac{1}{|C_k^{(t+1)}|} \sum_{i \in C_k^{(t+1)}} x_{ij}$$

Thus, the new loss function becomes

$$\mathcal{F}(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t+1)}, \dots, m_K^{(t+1)}) = \sum_{k=1}^K \sum_{i \in C_k^{(t+1)}} \sum_{j=1}^p (x_{ij} - m_{kj}^{(t+1)})^2$$

From the problem statement,

$$c(i)^{(t+1)} = \arg \min_{k \in \{1, \dots, K\}} \sum_{j=1}^p (x_{ij} - m_{kj}^{(t+1)})^2$$

Thus,

$$\sum_{j=1}^p (x_{c(i)^{(t+1)} j} - m_{kj}^{(t+1)})^2 \leq \sum_{j=1}^p (x_{ij} - m_{kj}^{(t+1)})^2$$

as $c(i)^{(t+1)}$ is the argument that minimizes $\sum_{j=1}^p (x_{ij} - m_{kj}^{(t+1)})^2$

And the resulting updated clusters

$$C_k^{(t+1)} = \{i : c(i)^{(t+1)} = k\}, k = 1, \dots, K$$

satisfy

$$\mathcal{F}(C_1^{(t+1)}, \dots, C_K^{(t+1)}, m_1^{(t+1)}, \dots, m_K^{(t+1)}) \leq \mathcal{F}(C_1^{(t)}, \dots, C_K^{(t)}, m_1^{(t+1)}, \dots, m_K^{(t+1)})$$

Problem 6: Practice of the Hierarchical Clustering (10 pts)

Suppose that we have four observations, for which we have a dissimilarity matrix

$$\begin{bmatrix} & 0.3 & 0.4 & 0.7 \\ 0.3 & & 0.5 & 0.8 \\ 0.4 & 0.5 & & 0.9 \\ 0.7 & 0.8 & 0.9 & \end{bmatrix}.$$

For instance, the dissimilarity between the first and second observations is 0.3, and the dissimilarity between the second and fourth observations is 0.8.

- On the basis of this dissimilarity matrix, sketch the dendrogram that results from hierarchical clustering on these four observations using complete linkage. Be sure to indicate on the plot the height at which each merging occurs, as well as the observations corresponding to each leaf in the dendrogram. Cut the dendrogram to obtain two clusters. Which observations are in each cluster?

$$\begin{array}{c} x_1 \quad x_2 \quad x_3 \quad x_4 \\ \left[\begin{array}{cccc} & 0.3 & 0.4 & 0.7 \\ 0.3 & & 0.5 & 0.8 \\ 0.4 & 0.5 & & 0.9 \\ 0.7 & 0.8 & 0.9 & \end{array} \right] \\ x_2 \\ x_3 \\ x_4 \end{array}$$

① Complete linkage is given as

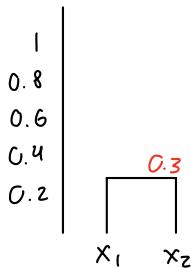
$$d_{cl}(A, C) = \max_{x_i \in A, x_j \in C} d_{ij}$$

$$\begin{aligned} d(x_1, x_2) &= 0.3 & d(x_2, x_3) &= 0.5 & d(x_3, x_4) &= 0.9 \\ d(x_1, x_3) &= 0.4 & d(x_2, x_4) &= 0.8 & \\ d(x_1, x_4) &= 0.7 \end{aligned}$$

For hierarchical clustering, we start with n datapoints

- Find the two closest data points
- Merge
- Repeat

x_1 & x_2 are the closest datapoints with the lowest dissimilarity / highest similarity of 0.3

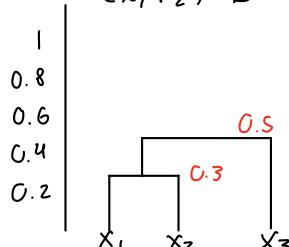


Since this is complete linkage, we find the max distance between clusters & points

$$\begin{aligned} d_{cl}((x_1, x_2), (x_3)) &= \max(d(x_1, x_3), d(x_2, x_3)) \\ &= \max(0.4, 0.5) \\ &= 0.5 \end{aligned}$$

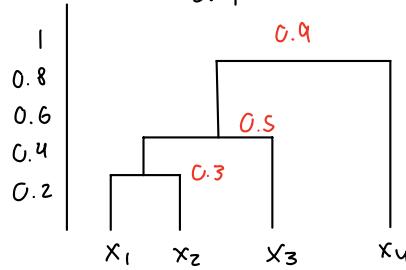
$$\begin{aligned} d_{cl}((x_1, x_2), (x_4)) &= \max(d(x_1, x_4), d(x_2, x_4)) \\ &= \max(0.7, 0.8) \\ &= 0.8 \end{aligned}$$

The closest point to (x_1, x_2) is x_3 with a dissimilarity of 0.5



We calculate the distance of x_u to the rest of the points, taking the max as at splitting threshold

$$\begin{aligned} d_{\text{cl}}((x_1, x_2, x_3), (x_u)) &= \max(d(x_1, x_u), d(x_2, x_u), d(x_3, x_u)) \\ &= \max(0.7, 0.8, 0.9) \\ &= 0.9 \end{aligned}$$



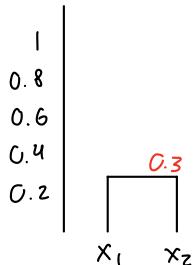
If we want to cut the dendrogram into two clusters, x_1, x_2 , & x_3 belong in one cluster while x_u is its own cluster.

2. Repeat 1, this time using single linkage.

② Single linkage is given as

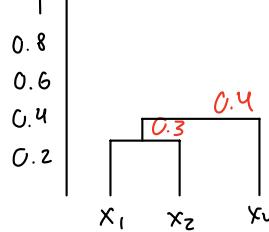
$$d_{\text{SL}}(A, C) = \min_{x_i \in A, x_j \in C} d_{ij}$$

x_1 & x_2 are the closest datapoints with the lowest dissimilarity/highest similarity of 0.3



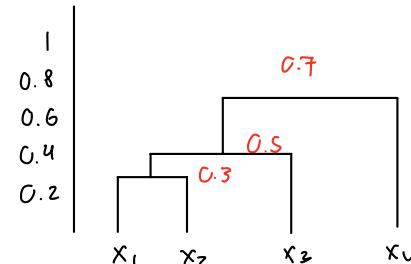
Since this is single linkage, we find the min distance between clusters & points

$$\begin{aligned} d_{\text{cl}}((x_1, x_2), (x_3)) &= \min(d(x_1, x_3), d(x_2, x_3)) \\ &= \min(0.4, 0.5) \\ &= 0.4 \end{aligned} \quad \begin{aligned} d_{\text{cl}}((x_1, x_2), (x_u)) &= \min(d(x_1, x_u), d(x_2, x_u)) \\ &= \min(0.7, 0.8) \\ &= 0.7 \end{aligned}$$



We calculate the distance of x_u to the rest of the points, taking the min as at splitting threshold

$$\begin{aligned} d_{\text{cl}}((x_1, x_2, x_3), (x_u)) &= \min(d(x_1, x_u), d(x_2, x_u), d(x_3, x_u)) \\ &= \min(0.7, 0.8, 0.9) \\ &= 0.7 \end{aligned}$$



Problem 7: EM Algorithm for the Gaussian Mixture Model (20 pts)

In the following Gaussian Mixture Model

$$X_i | Z_i = 0 \sim N(\mu_0, \sigma_1^2);$$

$$X_i | Z_i = 1 \sim N(\mu_1, \sigma_2^2);$$

$$Z_i \sim \text{Bernoulli}(\gamma), i = 1, \dots, n,$$

where X_i 's are observable random variables, and Z_i 's are hidden random variables.

Given observed data points x_1, \dots, x_n , derive the EM algorithm for estimating $\mu_0, \mu_1, \sigma_1^2, \sigma_2^2$ and γ in the following steps .

1. Write down the complete log-likelihood $\ell(\mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma)$ in terms of x_1, \dots, x_n and Z_1, \dots, Z_n .

The problem statement defines

$$X_i | Z_i = 0 \sim N(\mu_0, \sigma_1^2)$$

$$X_i | Z_i = 1 \sim N(\mu_1, \sigma_2^2)$$

$$Z_i \sim \text{Bernoulli}(\gamma), i = 1, \dots, n$$

Where γ is the probability of $Z_i = 1$

To find the log likelihood, we first find the complete likelihood

Note: Probability chain rule
 $P(A, B, C) = P(A|B, C)P(B, C)$

$$L(x_1, z_1, x_2, z_2, \dots, x_n, z_n; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma) = \prod_{i=1}^n P(X_i | Z_i; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma) P(Z_i; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma)$$

$$= \prod_{i=1}^n P(X_i | Z_i; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma) P(Z_i; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma)$$

Note : The PDF of $X \sim N(\mu, \sigma^2)$ & $Z \sim \text{Bernoulli}(\gamma)$ is given as

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}} (1-\gamma)^{1-z_i} \cdot \left[\frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}} \right]^{z_i} \gamma^{z_i} \right]$$

$$= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}} (1-\gamma)^{1-z_i} \cdot \left[\frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^{z_i} \right]$$

$$\ell_c(x_1, z_1, x_2, z_2, \dots, x_n, z_n; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma) = \log(L(x_1, z_1, x_2, z_2, \dots, x_n, z_n; \mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma))$$

$$= \log \left(\prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}} (1-\gamma)^{1-z_i} \cdot \left[\frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^{z_i} \right] \right)$$

$$= \sum_{i=1}^n \log \left(\left[\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}} (1-\gamma)^{1-z_i} \right] + \log \left(\left[\frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^{z_i} \right) \right)$$

$$= \sum_{i=1}^n (1-z_i) \log \left(\left[\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}} (1-\gamma)^{1-z_i} \right] \right) + z_i \log \left(\left[\frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^{z_i} \right)$$

$$\log \left(\left[\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}} (1-\gamma)^{1-z_i} \right] \right) = \log \left(\frac{1}{\sqrt{2\pi}\sigma_1} \right) + \log \left(\frac{-\frac{(x_i-\mu_0)^2}{2\sigma_1^2}}{e} \right) + \log(1-\gamma)$$

$$= \log \left(\frac{1}{\sqrt{2\pi}\sigma_1} \right) - \frac{(x_i-\mu_0)^2}{2\sigma_1^2} + \log(1-\gamma)$$

$$\log \left(\left[\frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right] \right) = \log \left(\frac{1}{\sqrt{2\pi}\sigma_2} \right) + \log \left(\frac{-\frac{(x_i-\mu_1)^2}{2\sigma_2^2}}{e} \right) + \log(\gamma)$$

$$= \log \left(\frac{1}{\sqrt{2\pi}\sigma_2} \right) - \frac{(x_i-\mu_1)^2}{2\sigma_2^2} + \log(\gamma)$$

$$\ell_c(\mu_0, \mu_1, \sigma_1^2, \sigma_2^2, \gamma) = \sum_{i=1}^n (1-z_i) \left(\log \left(\frac{1}{\sqrt{2\pi}\sigma_1} \right) - \frac{(x_i-\mu_0)^2}{2\sigma_1^2} + \log(1-\gamma) \right) + z_i \left(\log \left(\frac{1}{\sqrt{2\pi}\sigma_2} \right) - \frac{(x_i-\mu_1)^2}{2\sigma_2^2} + \log(\gamma) \right)$$

2. In the E-step of the $(t+1)$ -th iteration ($t = 0, 1, 2, \dots$), derive the conditional expectation of Z_i given x_i and the current parameter estimates $(\hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)})$:

$$\tau_i^{(t+1)} = E \left[Z_i | x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)} \right].$$

$$\tau_i^{(t+1)} = E[Z_i | x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}]$$

Note: The expectation of a discrete random variable X is given as
 $E[X] = \sum_{x \in \mathcal{X}} x p(x)$

$$\begin{aligned} \tau_i^{(t+1)} &= 1 \cdot P(Z_i = 1 | x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) + \\ &\quad 0 \cdot P(Z_i = 0 | x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) \end{aligned}$$

Note: Bayes Rule & the law of total probability states

$$P(B_j | A) = \frac{P(A | B_j) P(B_j)}{\sum_{i=1}^k P(A | B_i) P(B_i)}$$

$$\begin{aligned} &P(Z_i = 1 | x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) \\ &= \sum_{z_i \in \{0, 1\}} \frac{P(X_i | Z_i = z_i, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = z_i; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)})}{P(X_i | Z_i = z_i, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = z_i; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)})} \\ &= \frac{P(X_i | Z_i = 1, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = 1; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)})}{P(X_i | Z_i = 0, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = 0; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)})} \\ &\quad + P(X_i | Z_i = 0, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = 0; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) \end{aligned}$$

From problem 5.1

$$\begin{aligned} &P(X_i | Z_i = 0, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = 0; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) \\ &= \left[\frac{1}{\sqrt{2\pi} \sigma_1} \cdot e^{-\frac{(x_i - \mu_0)^2}{2\sigma_1^2}} \cdot (1-\gamma) \right]^{1-0} \cdot \left[\frac{1}{\sqrt{2\pi} \sigma_2} \cdot e^{-\frac{(x_i - \mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^0 \end{aligned}$$

$$\begin{aligned} &P(X_i | Z_i = 1, x_i; \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) P(Z_i = 1; x_i, \hat{\mu}_0^{(t)}, \hat{\mu}_1^{(t)}, (\hat{\sigma}_1^{(t)})^2, (\hat{\sigma}_2^{(t)})^2, \hat{\gamma}^{(t)}) \\ &= \left[\frac{1}{\sqrt{2\pi} \sigma_1} \cdot e^{-\frac{(x_i - \mu_0)^2}{2\sigma_1^2}} \cdot (1-\gamma) \right]^{1-1} \cdot \left[\frac{1}{\sqrt{2\pi} \sigma_2} \cdot e^{-\frac{(x_i - \mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^1 \\ &= \left[\frac{1}{\sqrt{2\pi} \sigma_2} \cdot e^{-\frac{(x_i - \mu_1)^2}{2\sigma_2^2}} \cdot \gamma \right]^1 \end{aligned}$$

$$\begin{aligned} \tau_i^{(t+1)} &= \frac{\frac{1}{\sqrt{2\pi} \sigma_2} \cdot e^{-\frac{(x_i - \mu_1)^2}{2\sigma_2^2}} \cdot \gamma}{\frac{1}{\sqrt{2\pi} \sigma_2} \cdot e^{-\frac{(x_i - \mu_1)^2}{2\sigma_2^2}} \cdot \gamma + \frac{1}{\sqrt{2\pi} \sigma_1} \cdot e^{-\frac{(x_i - \mu_0)^2}{2\sigma_1^2}} \cdot (1-\gamma)} \end{aligned}$$

3. In the M-step of the $(t+1)$ -th iteration, derive the updated parameter estimates based on x_1, \dots, x_n and $\tau_1^{(t+1)}, \dots, \tau_n^{(t+1)}$.

$$\left(\hat{\mu}_0^{(t+1)}, \hat{\mu}_1^{(t+1)}, (\hat{\sigma}_1^{(t+1)})^2, (\hat{\sigma}_2^{(t+1)})^2, \hat{\gamma}^{(t+1)} \right).$$

Using the formulas given in the notes sheet & as derived in class,

$$\begin{aligned}\hat{\gamma}^{(t+1)} &= \frac{1}{n} \sum_{i=1}^n \tau_i^{(t)} \\ \hat{\mu}_0^{(t+1)} &= \frac{\sum_{i=1}^n x_i (1 - \tau_i^{(t)})}{\sum_{i=1}^n (1 - \tau_i^{(t)})} & \hat{\mu}_1^{(t+1)} &= \frac{\sum_{i=1}^n x_i \tau_i^{(t)}}{\sum_{i=1}^n \tau_i^{(t)}} \\ (\hat{\sigma}_1^2)^{(t+1)} &= \frac{\sum_{i=1}^n (x_i - \hat{\mu}_0^{(t+1)})^2 (1 - \tau_i^{(t+1)})}{\sum_{i=1}^n (1 - \tau_i^{(t+1)})} & (\hat{\sigma}_2^2)^{(t+1)} &= \frac{\sum_{i=1}^n (x_i - \hat{\mu}_0^{(t+1)})^2 \tau_i^{(t+1)}}{\sum_{i=1}^n \tau_i^{(t+1)}}\end{aligned}$$