

Question 1

1)

As it is mentioned in Bishop chapter 9 section 3 the relationship between these two methods are below:

- the K-means method can cluster spherical data appropriately while the EM can work well on elliptical data either.
- k-means performs hard assignment of data point to cluster in which each data is associated uniquely with one cluster but EM method makes soft assignment based on posterior in which each cluster have wight (probability) that indicate responsibility of that cluster to create data X
- **Yes on the some limits these two algorithms work similar** and the limits are below: lets first get some intuition of whats going on:

Intuition

as we said K-means work fine on spherical data so we will assume that data are spherical (we will see it when we write equations we will choose Σ to be a constant times Identity matrix which means data are distributed spherical around the μ). so as we know when the variance of one Gaussian model goes to zero the model shrink around the mean so it is some how the hard assignment on data like k-means and the log likelihood function increases so if we choose variance small and the data be spherical the GMM (EM) algorithm will work similar to k-means.

limits

Assume Gaussian Mixture model in which the co-variance matrices are multiplication of a constant to Identity matrix so they are like $\Sigma = \epsilon I$ where ϵ can be presentation of radius of sphere around mean.

so distribution is like below:

$$p(x|\mu_k, \Sigma_k) = \frac{1}{2\pi\epsilon} e^{-\frac{1}{2\epsilon}|x-\mu_k|^2} \quad (1)$$

so the posterior distribution of responsibility of each component will be:

$$\gamma(Z_{n,k}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)} \quad (2)$$

by applying eq(1) we get:

$$\gamma(Z_{n,k}) = \frac{\pi_k e^{-\frac{1}{2\epsilon} |x_n - \mu_k|^2}}{\sum_j \pi_j e^{-\frac{1}{2\epsilon} |x_n - \mu_j|^2}} \quad (3)$$

as $\epsilon \rightarrow 0$ so probability obtained from normal distribution goes to zero except for the component that $x_n \rightarrow \mu_i$ (notice that in such case we assumed that there is a data that is equal to μ more safe to say that we choose closest data to mean as representation of mean so in such case the responsibility of that component goes to one so it's like hard assignment in k-means method.

take the x_j closest to μ_j among other μ so we have:

$$\gamma(Z_{n,k}) = \frac{\pi_k e^{-\frac{1}{2\epsilon} |x_n - \mu_k|^2}}{\sum_j \pi_j e^{-\frac{1}{2\epsilon} |x_n - \mu_j|^2}} \quad (4)$$

$$\gamma(Z_{n,j}) = \frac{\pi_j}{\sum_k \pi_k e^{\frac{1}{2\epsilon} |x_n - \mu_j|^2 - |x_n - \mu_k|^2}} = \frac{\pi_j}{\sum_k \pi_k e^{\frac{1}{2\epsilon} \Delta_{K,n}}} \quad (5)$$

for closest to μ_k to x_n we can assume $\Delta_{K,n}$ go to zero so term is one and for others the term is zero so we have:

$$\gamma(Z_{n,j}) = \frac{\pi_j}{\pi_j \sum_{k/j} \pi_k e^{\frac{1}{2\epsilon} \Delta_{K,n}}} \quad (6)$$

$$\epsilon \rightarrow 0 \implies \gamma(Z_{n,j}) = 1 \quad (7)$$

$$\gamma(Z_{n,k/j}) = 0 \quad (8)$$

it's like hard assignment in k-means method.

so $\gamma(Z_{n,k}) \rightarrow r_{n,k}$ as we know the likelihood function of EM-estimator is

$$E_Z[\ln P(X, Z | \mu, \Sigma, \Pi)] = \sum_{n=1}^{n=N} \sum_{k=1}^{k=K} \gamma(Z_{n,k}) \{ \ln \pi_k + \ln N(x_n | \mu, \Sigma) \} \quad (9)$$

where $\gamma(Z_{n,k}) = E[Z_{n,k}]$

by replacing and applying logarithm to parameters we have :

$$\text{MAXIMIZE } E_Z[\ln P(X, Z | \mu, \Sigma, \Pi)] = \quad (10)$$

$$\text{MAXIMIZE } \sum_{n=1}^{n=N} \sum_{k=1}^{k=K} \gamma(Z_{n,k}) \{ \ln \pi_k + \ln N(x_n | \mu, \Sigma) \} \quad (11)$$

$$\epsilon \rightarrow \infty \implies \quad (12)$$

$$\text{MAXIMIZE } \sum_{n=1}^{n=N} \sum_{k=1}^{k=K} r_{n,k} \{ \ln \pi_k + \ln N(x_n | \mu, \Sigma) \} \quad (13)$$

now π_k is obtained under condition mentioned maximization is equal to : (14)

$$\text{minimize } \sum_{n=1}^{n=N} \sum_{k=1}^{k=K} r_{n,k} \{ |x_n - \mu_k|^2 \} \quad (15)$$

eq(15) is exactly the objective function in k-means method ✓

$$P(x_b | x_a) = \frac{P(x_a, x_b)}{P(x_a)} = \frac{P(x)}{P(x_a)} = \frac{\sum_{k=1}^K \pi_k P(x|k)}{P(x_a)} = \sum_{k=1}^K \frac{\pi_k}{P(x_a)} P(x|k) \quad (1)$$

$$= \sum_{k=1}^K \frac{\pi_k}{P(x_a)} P(x_a, x_b | k) = \sum_{k=1}^K \frac{\pi_k}{P(x_a)} P(x_b | x_a, k) P(x_a | k) = \sum_{k=1}^K \frac{\pi_k P(x_a | k)}{P(x_a)} P(x_b | x_a, k)$$

$$P(x_a | k) = \sum_{x_b} P(x_a, x_b | k)$$

$$P(x_a) = \sum_{x_b} P(x_a, x_b) = \sum_{x_b} \sum_{k=1}^K \pi_k P(x_a, x_b | k) = \sum_{k=1}^K \sum_{x_b} \pi_k P(x_a, x_b | k) =$$

$$\sum_{k=1}^K \pi_k \sum_{x_b} P(x_a, x_b | k) = \sum_{k=1}^K \pi_k P(x_a | k)$$

$$P(x_b | x_a, k) = \frac{P(x_a, x_b | k)}{P(x_a | k)}$$

$$l(\theta) = \sum_{i=1}^N \sum_{k=1}^K \gamma_k^i \log(\pi_k |\Sigma|^{-\frac{1}{V}} (\gamma \pi)^{-\frac{K}{V}} \exp(-\frac{1}{V} (x_i - \mu'_k)^T \Sigma^{-1} (x_i - \mu'_k))) = \quad (2)$$

$$\sum_{i=1}^N \sum_{k=1}^K \gamma_k^i (\log(\pi_k) - \frac{1}{V} \log(|\Sigma|) - \frac{K}{V} \log(\gamma \pi) - \frac{1}{V} (x_i - \mu'_k)^T \Sigma^{-1} (x_i - \mu'_k))$$

$$\Rightarrow \frac{\partial l(\theta)}{\partial \Sigma} = \sum_{i=1}^N \sum_{k=1}^K \gamma_k^i \left(-\frac{1}{V} \Sigma^{-1} + \frac{1}{V} \Sigma^{-1} (x_i - \mu'_k) (x_i - \mu'_k)^T \Sigma^{-1} \right) = 0$$

$$\Rightarrow -N \Sigma^{-1} + \sum_{i=1}^N \sum_{k=1}^K \gamma_k^i \Sigma^{-1} (x_i - \mu'_k) (x_i - \mu'_k)^T \Sigma^{-1} = 0 \Rightarrow N \Sigma = \sum_{i=1}^N \sum_{k=1}^K \gamma_k^i (x_i - \mu'_k) (x_i - \mu'_k)^T$$

$$\Rightarrow \Sigma = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \gamma_k^i (x_i - \mu'_k) (x_i - \mu'_k)^T$$

μ'_j, π'_j مانند حالت ماکزیمم لایکلیت برای مقادیر حساب می شود.

میزر Z را طبق مدل Mixture تعریف کنیم: احتمال اینکه داده نام برده شده از دسته jام باشد. $P(z_j^{(i)} = 1 | x_i, \theta) = \gamma_j^i$

توزیع پواسون به صورت اربده است. $Poisson(X=x|\theta) = \frac{\theta^x e^{-\theta}}{x!}$

E-step

$P(z|x, \theta)$

$$\gamma_j^i = P(z_j^{(i)} = 1 | x_i, \theta) = \frac{P(x_i | z_j^{(i)} = 1, \theta) P(z_j^{(i)} = 1)}{\sum_k P(x_i | z_k^{(i)} = 1, \theta) P(z_k^{(i)} = 1)} = \frac{\pi_j \theta_j^{x_i} e^{-\theta_j}}{\sum_k \pi_k \theta_k^{x_i} e^{-\theta_k}}$$

M-step

$$Q(\theta) = E_{Z \sim p(z|x, \theta)} [\log P(x, z|\theta)] = \sum_{i=1}^n \sum_{j=1}^J \gamma_j^i \log [\pi_j \theta_j^{x_i} e^{-\theta_j} \frac{1}{x_i!}]$$

$$\frac{\partial Q}{\partial \theta_k} = \sum_{i=1}^n \sum_{j=1}^J \gamma_j^i \frac{\partial}{\partial \theta_k} [\log \pi_j + x_i \log \theta_j - \theta_j - \log x_i!] = 0$$

نقطه عملی کام در اینجا نام برده است $\rightarrow \sum_{i=1}^n \gamma_k^i (\frac{x_i}{\theta_k} - 1) = 0 \rightarrow \theta_k \sum_{i=1}^n \gamma_k^i = \sum_{i=1}^n \gamma_k^i x_i$

$$\rightarrow \theta_k = \frac{\sum_{i=1}^n \gamma_k^i x_i}{\sum_{i=1}^n \gamma_k^i}$$

اینست مقدار تخمین برآورد ما

$L(\theta, \lambda) = Q(\theta) - \lambda (\sum_{j=1}^J \pi_j - 1)$

نقطه عملی

$$\frac{\partial L}{\partial \pi_k} = \sum_{i=1}^n \sum_{j=1}^J \gamma_j^i \frac{\partial}{\partial \pi_k} [\log \pi_j + x_i \log \theta_j - \theta_j - \log x_i!] - \lambda = 0$$

نقطه عملی کام در اینجا نام برده است $\rightarrow \sum_{i=1}^n \gamma_k^i \frac{1}{\pi_k} - \lambda = 0 \rightarrow \lambda = \frac{\sum_{i=1}^n \gamma_k^i}{\pi_k} \rightarrow \pi_k = \frac{\sum_{i=1}^n \gamma_k^i}{\lambda}$

$\sum_{k=1}^J \pi_k = 1 \rightarrow \sum_{k=1}^J \frac{\sum_{i=1}^n \gamma_k^i}{\lambda} = 1 \rightarrow \lambda = \frac{\sum_{k=1}^J \sum_{i=1}^n \gamma_k^i}{1} = \frac{\sum_{i=1}^n \sum_{k=1}^J \gamma_k^i}{1} = \frac{n}{1} = n$

$\rightarrow \lambda = n \rightarrow \pi_k = \frac{\sum_{i=1}^n \gamma_k^i}{n}$

Question 3)

1)

As it's binomial distribution it's easy to obtain:

$$\begin{aligned} I) \quad p(w_t, m_t | k) &= \binom{m_t}{w_t} p_k^{w_t} (1 - p_k)^{m_t - w_t} \\ II) \quad p(w_t, m_t | ct) &= \binom{m_t}{w_t} p_{ct}^{w_t} (1 - p_{ct})^{m_t - w_t} \end{aligned}$$

2)

As we know from previous question for E-step we have :

$$\begin{aligned} Q(k) = p(k|x, \theta) &\propto p(x|ct = k, \theta) p(ct = k) \propto \binom{m_t}{x} p_k^x (1 - p_k)^{m_t - x} \pi_k \\ \sum_K Q(k) = 1 &\implies \sum_K \binom{m_t}{x} p_k^x (1 - p_k)^{m_t - x} \pi_k = 1 \end{aligned}$$

As we know the equation above is propto so we need to normalize it that it adds up to one:

$$Q_t[k]^i = \frac{\binom{m_t}{x} p_k^x (1 - p_k)^{m_t - x} \pi_k}{\sum_K \binom{m_t}{x} p_j^x (1 - p_j)^{m_t - x} \pi_j}$$

3)

For writing M-step we have to maximize equation below:

$$\begin{aligned} L &= \sum_N \sum_K Q_t[k]^i \log(p(w_t | p_k)) = \sum_N \sum_K Q_t[k]^i \log\left(\binom{m_t}{w_t} p_k^{w_t} (1 - p_k)^{m_t - w_t}\right) = \\ &\quad \sum_N \sum_K Q_t[k]^i \left\{ \log\left(\binom{m_t}{w_t}\right) + w_t \log(p_k) + (m_t - w_t) \log(1 - p_k) \right\} \\ \frac{\partial L}{\partial p_j} &= \sum_N Q_t[j]^i \left\{ \frac{w_t}{p_j} - \frac{m_t - w_t}{1 - p_j} \right\} = \sum_N Q_t[j]^i \left\{ \frac{w_t - m_t p_j}{p_j - p_j^2} \right\} = 0 \implies \\ \sum_N Q_t[j]^i \{w_t - m_t p_j\} &= 0 \implies p_j = \frac{\sum_{t=1}^N Q_t[j]^i w_t}{\sum_{t=1}^N Q_t[j]^i m_t} \checkmark \end{aligned}$$

Bayes' Theorem: $P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$

4.1

$\Rightarrow \ln P(\theta|X) = \ln P(X|\theta) + \ln P(\theta) - \ln P(X)$

$\Rightarrow \ln P(\theta|X) \propto \ln \left(\sum_z P(X,z|\theta) \right) + \ln P(\theta) = \ln \left(\sum_z P(X,z|\theta) P(\theta) \right)$

بین تنها تفاوت با حالت اصلی $P(\theta)$ است که اضافه شده است. پس مشابه مرحله E باید $P(z|X, \theta^{old})$ را حساب کنیم در M تابع را بر حسب θ بیفک کنیم.

E-Step: calculate $P(z|X, \theta^{old})$

M-Step: $\theta^{new} = \operatorname{argmax}_{\theta} \sum_z P(z|X, \theta^{old}) \ln [P(X,z|\theta) P(\theta)]$

$\xrightarrow{\text{marginalizing}} = \operatorname{argmax}_{\theta} \underbrace{\sum_z P(z|X, \theta^{old}) \ln P(X,z|\theta)}_{\text{همان تابع ثابت قبل}} + \ln P(\theta)$

4.2 ابتدا باید X (متغیر observed) ، Z (متغیر latent) را تعریف کنیم.

با توجه به داده های سوال، به این صورت تعریف می شوند:

$X \in \{A \text{ or } C, B, D\}$ $Z \in \{A, B, C, D\}$

با توجه به داده های سوال، X تنها ۳ حالت دارد در حالی که بر حسب اصلی ۴ حالت دارد.

حالا با توجه به قسمت اول سوال، معادلات EM را می نویسیم.

$$E\text{-Step: } P\{Z=A \mid X=A \text{ or } C\} = \frac{1}{1+2\theta}$$

$$P\{Z=C \mid X=A \text{ or } C\} = \frac{2\theta}{1+2\theta}$$

بغیر از حالت = با، $P\{Z \mid X\}$ برابر باشد یا یک است و به سادگی تعیین می شود.

$$M\text{-Step: } \theta^{\text{new}} = \arg \max_{\theta} \sum_{i=1}^n \sum_Z P\{Z \mid X^{(i)}; \theta\} (\ln P\{X^{(i)}, Z; \theta\} + \ln P\{\theta\})$$

$$= \arg \max_{\theta} (n_a + n_c) \left(\frac{1}{1+2\theta^{\text{old}}} \ln \frac{1}{3} + \frac{2\theta^{\text{old}}}{1+2\theta^{\text{old}}} \ln \frac{2}{3} \theta \right) + n_b \ln \frac{1}{3} (1-\theta)$$

$$+ n_d \ln \frac{1}{3} (1-\theta) + n(v_1-1) \ln \theta + n(v_2-1) \ln (1-\theta)$$

$$= \arg \max_{\theta} (n_a + n_c) \frac{2\theta^{\text{old}}}{1+2\theta^{\text{old}}} \ln \theta + (n_b + n_d) \ln (1-\theta) + n(v_1-1) \ln \theta + n(v_2-1) \ln (1-\theta)$$

$$\Rightarrow \frac{2\theta^{\text{old}} (n_a + n_c)}{(1+2\theta^{\text{old}}) \theta} + \frac{n(v_1-1)}{\theta} = \frac{n_b + n_d + n(v_2-1)}{1-\theta}$$

$$\Rightarrow \frac{2\theta^{\text{old}} (n_a + n_c) + n(1+2\theta^{\text{old}})(v_1-1)}{(1+2\theta^{\text{old}}) \theta} = \frac{n_b + n_d + n(v_2-1)}{1-\theta}$$

$$\Rightarrow \theta^{\text{new}} = \frac{\frac{2\theta^{\text{old}}}{1+2\theta^{\text{old}}} (n_a + n_c) + n(v_1-1)}{\frac{2\theta^{\text{old}}}{1+2\theta^{\text{old}}} (n_a + n_c) + n(v_1-1) + n_b + n_d + n(v_2-1)}$$

$$n_a + n_c = n - n_b - n_d \quad \checkmark \text{ این است}$$