

Name: 费祥 Student ID: 120090414

T₁.

A. According to the assumption, variances are shared spherical, $\Sigma = \sigma^2 I$, and we have $\pi_k = \frac{1}{K}$

$$\begin{aligned} \underset{z, \mu}{\operatorname{argmax}} \log P(D, z | \pi, \mu, \sigma) &= \underset{z, \mu}{\operatorname{argmax}} \sum_i \log \pi(z_i) P_{\mu, \sigma}(x_i | z_i) \\ &= \underset{z, \mu}{\operatorname{argmax}} \sum_i \left(\log \frac{1}{K} + \log N(x_i; \mu_{z_i}, \sigma) \right) = \underset{z, \mu}{\operatorname{argmin}} \sum_i \frac{\|\mu_{z_i} - x_i\|^2}{2\sigma^2} \\ &= \underset{z, \mu}{\operatorname{argmin}} \sum_i \|\mu_{z_i} - x_i\|^2 = \underset{z, \mu}{\operatorname{argmin}} \sum_i \sum_k r_{ik} \|\mu_k - x_i\|^2 \end{aligned}$$

\therefore the above is the same as k-means

\therefore Q.E.D.

B. according to the objective function, we don't need to compute the covariance matrices.

T₂.

(i) (11.87) is the term $\sum_z q_i(z) \log \frac{P(x_i, z_i | \theta)}{q_i(z)}$ from the left hand of (11.85), and the lower bound in (11.85) is a sum over i of this term.

Then, since joint probability $P(x_i, z_i | \theta) = P(z_i | x_i, \theta)P(x_i | \theta)$, we can reach (11.88)

$$(ii) Q(\theta^t, \theta^t) = \sum_i L(\theta^t, q_i^t)$$

Since $q_i^t(z_i) = P(z_i | x_i, \theta^t)$ \therefore KL divergence is 0

$$\therefore L(\theta^t, q_i^t) = \log P(x_i | \theta^t) \Rightarrow Q(\theta^t, q_i^t) = \sum_i \log P(x_i | \theta^t)$$

$$\text{and we have } L(\theta^t) = \sum_i \log \left[\sum_z P(x_i, z_i | \theta^t) \right] = \sum_i \log P(x_i | \theta^t)$$

$$\therefore Q(\theta^t, \theta^t) = \sum_i \log P(x_i | \theta^t) = L(\theta^t)$$

(iii) since $\log P(x_i|\theta) \geq L(\theta, q_i)$ always holds

$$\therefore L(\theta^{t+1}) \geq Q(\theta^{t+1}, \theta^t)$$

and since $\theta^{t+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta^t)$

$$\therefore Q(\theta^{t+1}, \theta^t) \geq Q(\theta^t, \theta^t)$$

and here we have $Q(\theta^t, \theta^t) = L(\theta^t)$

\therefore even though the hidden variable z do not appear, we still know that $q_i^t(z_i) = P(z_i|x_i, \theta^t)$ since $Q(\theta^t, \theta^t) = L(\theta^t)$
 \therefore they still have the right values in iteration t ,

T₃,

To explain why M-step works, we first need to know what M-step does.

In fact, our final purpose is to increase the log-likelihood to a maximum. but for the original optimized problem, there is no closed form solution since:

$$\log P(D|\theta) = \sum_{n=1}^N \log P(X^{(n)}|\theta) = \sum_{n=1}^N \log \left(\sum_{z^{(n)}} P(z^{(n)}, X^{(n)}|\theta) \right)$$

so there is a summation inside the log, it is hard to optimize. therefore, we use the auxiliary distribution of hidden variables $q_n(z^{(n)})$, and we have :

$$\ln P(D|\theta) = L(\theta, \theta) + KL(q(z)||P(z|D; \theta))$$

$L(\theta, \theta)$ is the lower bound, and in fact it is $Q(\theta; \theta_{old})$. and it is much easier to optimize, so we can get the θ_{new} in M-step.

Then, we just need to illustrate that this optimization can increase the original log-likelihood.

In E-step, we get the q_i^t that make $L(\theta, \theta)$ and $L(\theta; X, Y)$ tightest, and since $q_i^t(z) = P(z_i|x_i, \theta^t)$, here, θ^t is old, θ^{t+1} is θ_{new} .

So the KL divergence is 0 $\Rightarrow L(\theta^t) = Q(\theta^t, \theta^t)$

and since $\log P(x_i | \theta) \geq L(\theta, q_i)$ always holds

$$\therefore L(\theta^{t+1}) \geq Q(\theta^{t+1}, \theta^t)$$

and since $\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^t)$

$$\therefore Q(\theta^{t+1}, \theta^t) \geq Q(\theta^t, \theta^t) = L(\theta^t)$$

$$\therefore L(\theta^{t+1}) \geq L(\theta^t) \Rightarrow L(\theta_{\text{new}}) \geq L(\theta_{\text{old}})$$

so the likelihood increase.

Therefore, M-step works.