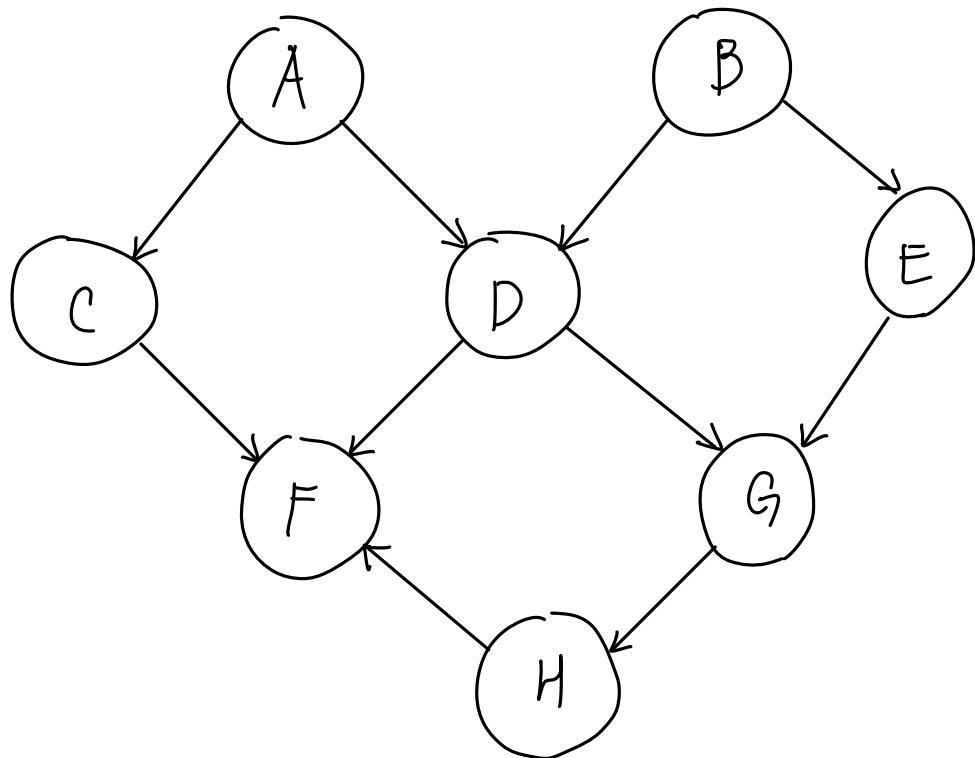


[Bayesian Network]

20170243 김재윤

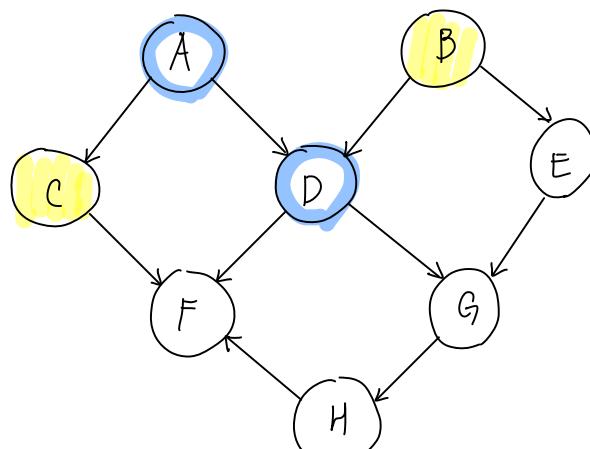


$(X \perp Y | Z)$: if X and Y are (conditionally) independent given Z .

(a) every pair of (marginally) independent variables.

(A, B) (A, E) (B, C) (C, E)

(b) $(B \perp C \mid A, D)$



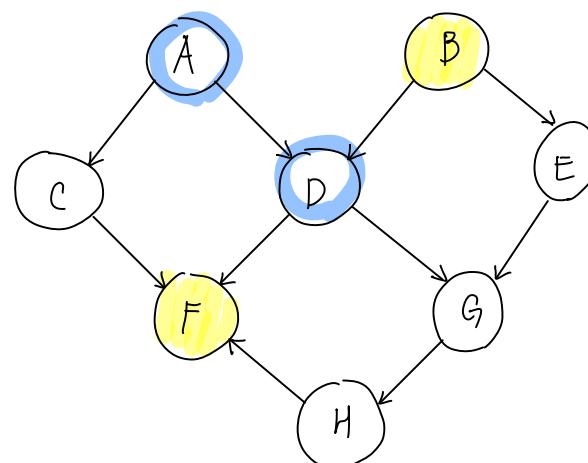
To check the conditional independence in directed graph, we should consider all paths from B to C and determine whether the path is blocked by A and D according to D-Separation.

Path that passes through A: A is tail-to-tail case \rightarrow block.
ex: B-D-A-C

Path that passes through D: D is head-to-tail case \rightarrow block
ex: B-D-F-C

Path that passes through F: F is head-to-head case & not observed \rightarrow block
ex: B-E-G-H-F-C
 \therefore independent

(c) $(B \not\perp F \mid A, D)$



Path that passes through D: D is head-to-tail case \rightarrow block
ex: B-D-F

Path that passes E, G, H: we can reach from B to F
ex: B-E-G-H-F
 \therefore not independent

(d)

① We assume that each variable can take one of 10 values. Then, if we do not use the knowledge encoded in graph, the number of parameters we need to model the full joint distribution is $10^8 - 1$. The reason why we subtract 1 is that the sum of all probabilities is 1, so we can get it by subtracting the rest from 1. So we will skip the last entry, so we don't need it.

② We should first joint probability with the factorization. Then, we can write the factorization form as follow. Then, the number of parameters we have to need is decreasing compared to ①.

$$\begin{aligned} & P(A, B, C, D, E, F, G, H) \\ &= P(A)P(B)P(C|A)P(D|A, B)P(E|B)P(F|C, D, H)P(G|D, E)P(H|G) \end{aligned}$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $10^0 \quad 10^0 \quad 10^1 \quad 10^2 \quad 10^1 \quad 10^3 \quad 10^2 \quad 10^1$
 $'' \quad '' \quad '' \quad '' \quad '' \quad '' \quad '' \quad ''$
 $| \quad | \quad | \quad | \quad | \quad | \quad | \quad |$
 $10 \quad 10 \quad 10 \quad 100 \quad 10 \quad 1000 \quad 100 \quad 10$

$$1 + 1 + 10 + 100 + 10 + 1000 + 100 + 10 = 1232$$

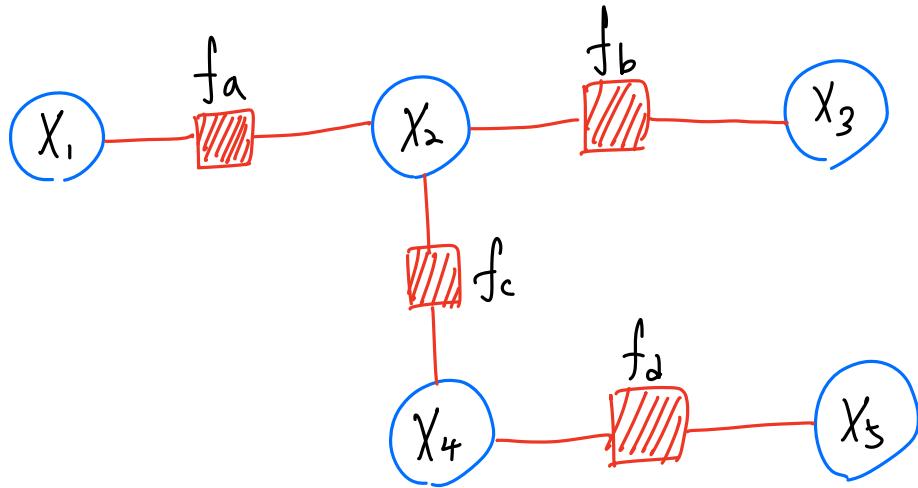
$$\therefore (10 - 1) \times 1232 = 1088$$

2. [Belief Propagation]

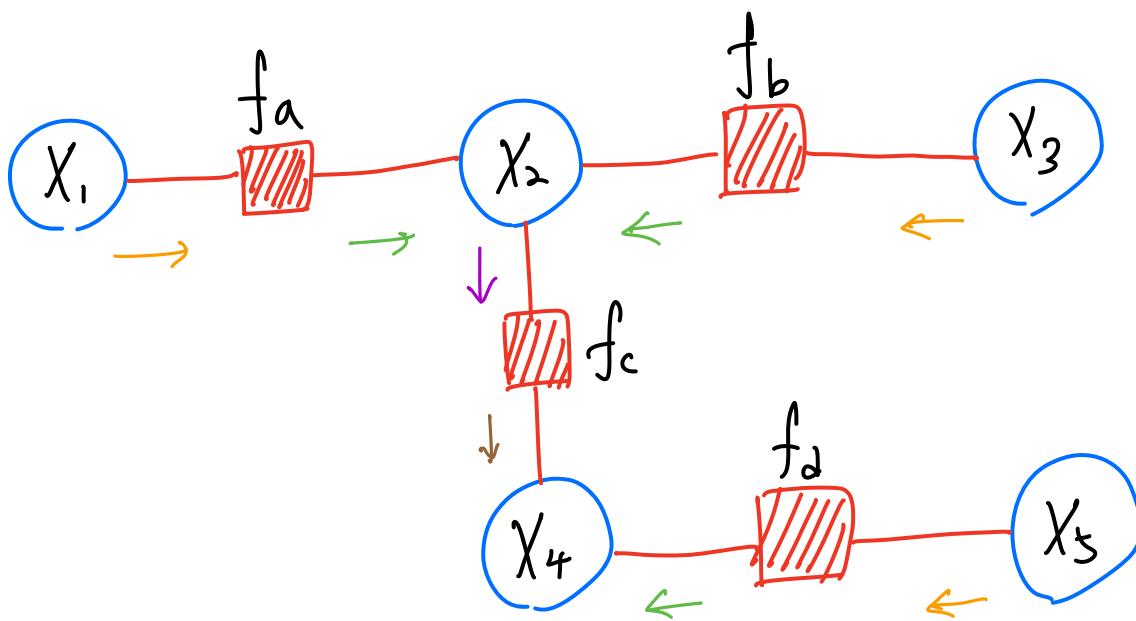
binary random variables $X_1, \dots, X_5 \in \{0,1\}$

$$p(X) = p(X_1, X_2, X_3, X_4, X_5) = \frac{1}{Z} f_a(X_1, X_2) f_b(X_2, X_3) f_c(X_2, X_4) f_d(X_4, X_5)$$

(a) corresponding factor graph.



(b) sum-product belief propagation \rightarrow compute $p(X_4)$



$$V_{1 \rightarrow a}(x_1) = 1$$

$$V_{3 \rightarrow b}(x_3) = 1$$

$$V_{5 \rightarrow d}(x_5) = 1$$

$$M_{a \rightarrow 2}(x_2) = \sum_{x_1} f_a(x_1, x_2) \cdot V_{1 \rightarrow a} = \sum_{x_1} f_a(x_1, x_2)$$

$$M_{b \rightarrow 2}(x_2) = \sum_{x_3} f_b(x_2, x_3) \cdot V_{3 \rightarrow b} = \sum_{x_3} f_b(x_2, x_3)$$

$$\underbrace{M_{d \rightarrow 4}(x_4)}_{\sim} = \sum_{x_5} f_d(x_4, x_5) \cdot V_{5 \rightarrow b} = \sum_{x_5} f_d(x_4, x_5)$$

$$V_{2 \rightarrow c}(x_2) = M_{a \rightarrow 2} \cdot M_{b \rightarrow 2} = \left(\sum_{x_1} f_a(x_1, x_2) \right) \left(\sum_{x_3} f_b(x_2, x_3) \right)$$

$$M_{c \rightarrow 4}(x_4) = \sum_{x_2} f_c(x_2, x_4) \cdot V_{2 \rightarrow c}$$

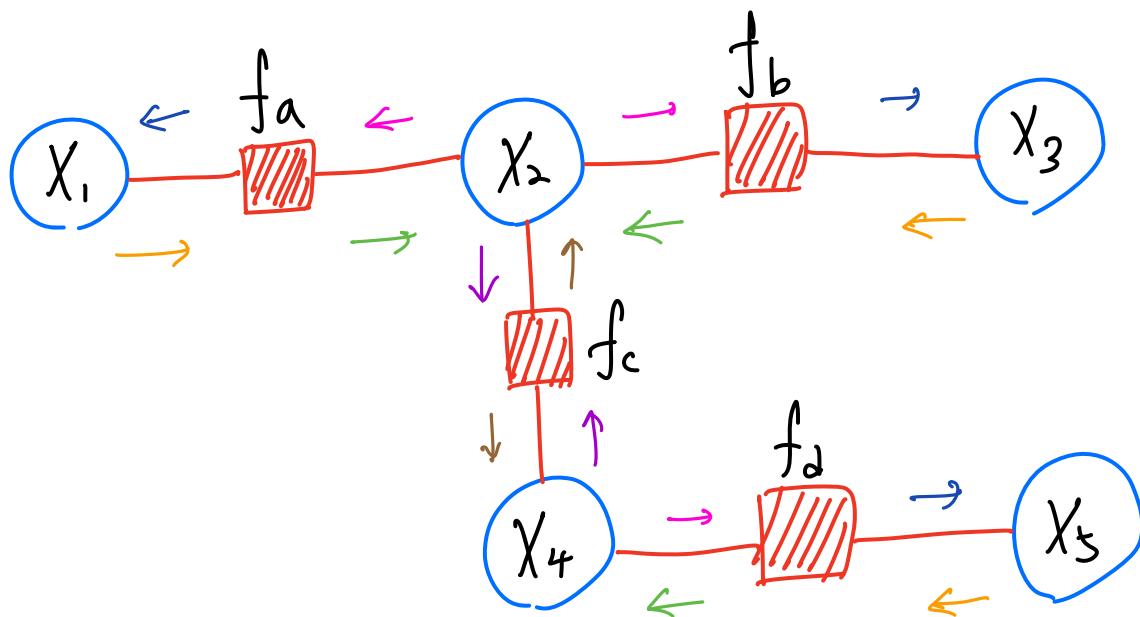
$$\underbrace{\quad}_{\sim} = \sum_{x_2} f_c(x_2, x_4) \left(\sum_{x_1} f_a(x_1, x_2) \sum_{x_3} f_b(x_2, x_3) \right)$$

$$\therefore P(x_4) \propto M_{d \rightarrow 4} \cdot M_{c \rightarrow 4}$$

$$= \underbrace{\frac{1}{2} \sum_{x_5} f_d(x_4, x_5)}_{M_{d \rightarrow 4}} \underbrace{\sum_{x_2} f_c(x_2, x_4) \left(\sum_{x_1} f_a(x_1, x_2) \sum_{x_3} f_b(x_2, x_3) \right)}_{M_{c \rightarrow 4}}$$

$$(c) f_u(x, y) = x + y \quad \text{for every } u \in \{a, b, c, d\}$$

max-sum BP algorithm \rightarrow (i) optimal configuration X^* maximizing the joint distribution
(ii) maximum of the joint distribution.



$$V_{1 \rightarrow a}(x_1) = 0$$

$$V_{2 \rightarrow b}(x_3) = 0$$

$$V_{5 \rightarrow d}(x_5) = 0$$

$$\begin{aligned} M_{a \rightarrow 2}(x_2) &= \max_{x_1} \{ \log f_a(x_1, x_2) + V_{1 \rightarrow a} \} \\ &= \max_{x_1} \{ \log (x_1 + x_2) \} = \log (1 + x_2) \end{aligned}$$

$$M_{b \rightarrow 2}(x_2) = \max_{x_3} \left\{ \log f_b(x_2, x_3) + V_{3 \rightarrow b} \right\}$$

$$= \max_{x_3} \left\{ \log (x_2 + x_3) \right\} = \log (x_2 + 1)$$

$$M_{d \rightarrow 4}(x_4) = \max_{x_5} \left\{ \log f_d(x_4, x_5) + V_{5 \rightarrow d} \right\}$$

$$= \max_{x_5} \left\{ \log (x_4 + x_5) \right\} = \log (x_4 + 1)$$

$$V_{2 \rightarrow c}(x_2) = M_{a \rightarrow 2} + M_{b \rightarrow 2}$$

$$= \log(1 + x_2) + \log(x_2 + 1)$$

$$= \log(1 + x_2)(x_2 + 1)$$

$$V_{4 \rightarrow c}(x_4) = M_{d \rightarrow 4}$$

$$= \log(x_4 + 1)$$

$$M_{c \rightarrow 4}(x_4) = \max_{x_2} \left\{ \log f_c(x_2, x_4) + V_{2 \rightarrow c}(x_2) \right\}$$

$$= \max_{x_2} \left\{ \log (x_2 + x_4) + \log (1 + x_2)(x_2 + 1) \right\}$$

$$= \max_{x_2} \left\{ \log (x_2 + x_4)(1 + x_2)(x_2 + 1) \right\}$$

$$= \log(1 + x_4)(1 + 1)(1 + 1)$$

$$M_{c \rightarrow 2}(x_2) = \max_{x_4} \left\{ \log f_c(x_2, x_4) + V_{4 \rightarrow c}(x_4) \right\}$$

$$= \max_{x_4} \left\{ \log (x_2 + x_4) + \log (x_4 + 1) \right\}$$

$$= \max_{x_4} \left\{ \log (x_2 + x_4)(x_4 + 1) \right\}$$

$$= \log(x_2 + 1)(1 + 1)$$

$$V_{4 \rightarrow d} (x_4) = M_{C \rightarrow 4} (x_4)$$

$$= \log (1+x_4)(1+1)(1+1)$$

$$V_{2 \rightarrow a} (x_2) = M_{C \rightarrow 2} (x_2)$$

$$= \log (x_2+1)(1+1)$$

$$V_{2 \rightarrow b} (x_2) = M_{C \rightarrow 2} (x_2)$$

$$= \log (x_2+1)(1+1)$$

$$M_{d \rightarrow 5} (x_5) = \max_{x_4} \{ \log f_d(x_4, x_5) + V_{4 \rightarrow d} \}$$

$$= \max_{x_4} \{ \log (x_4+x_5) + \log (1+x_4)(1+1)(1+1) \}$$

$$= \log (1+x_5)(1+1)(1+1)(1+1)$$

$$M_{a \rightarrow 1} (x_1) = \max_{x_2} \{ \log f_a(x_1, x_2) + V_{2 \rightarrow a} \}$$

$$= \max_{x_2} \{ \log (x_1+x_2) + \log (x_2+1)(1+1) \}$$

$$= \log (x_1+1)(1+1)(1+1)$$

$$M_{b \rightarrow 3} (x_3) = \max_{x_2} \{ \log f_b (x_2, x_3) + V_{2 \rightarrow b} \}$$

$$= \max_{x_2} \{ \log (x_2+x_3) + \log (x_2+1)(1+1) \}$$

$$= \log (1+x_3)(1+1)(1+1)$$

$$\begin{aligned}\log p(x_4) &= \max_{x_4} (\mu_{c \rightarrow 4} + \mu_{d \rightarrow 4}) \\ &= \max_{x_4} (\log(1+x_4)(1+1)(1+1) + \log(1+x_4)) \\ &= \log 2^4 = \log 16\end{aligned}$$

And, we need normalization term $(\frac{1}{Z})$, then

$$P_{\max} = \frac{16}{Z}, X^* = (x_1^*, x_2^*, x_3^*, x_4^*, x_5^*) = (1, 1, 1, 1, 1)$$

3 [Gaussian Mixture Model & Expectation - Maximization]

Gaussian Mixture Model (GMM) with K components.

each component $k \in [K] = \{1, 2, \dots, K\}$

mean μ_k , variance σ_k^2 , mixing parameter $\pi_k \rightarrow \theta = (\mu_k, \sigma_k^2, \pi_k)_{k \in [K]}$

dataset $X = \{x_i\}_{i \in [n]}, x_i \in \mathbb{R}$

latent variables $Z = \{z_i\}_{i \in [n]}, z_i = k$

$$\mathcal{N}(x_i | \mu_k, \sigma_k^2) := \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right)$$

(a)

① log-likelihood of the data

$$\begin{aligned} \log p(X|\theta) &= \log \prod_{i=1}^n p(x_i|\theta) = \sum_{i=1}^n \log p(x_i|\theta) \\ &= \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \sigma_k^2) \end{aligned}$$

② Complete-data log-likelihood

$$\begin{aligned} \log p(X, Z|\theta) &= \sum_{i \in [n]} \log p(x_i, z_i|\theta) = \sum_{i \in [n]} \log \left(p(x_i|z_i, \theta) p(z_i|\theta) \right) \\ &= \sum_{i \in [n]} \log \left(\prod_{k \in [K]} \left(\pi_k \mathcal{N}(x_i | \mu_k, \sigma_k^2) \right)^{z_{ki}} \right) \\ &= \sum_{i \in [n]} \sum_{k \in [K]} z_{ki} \log (\pi_k \mathcal{N}(x_i | \mu_k, \sigma_k^2)) \end{aligned}$$

$$(b) \quad \theta^{(t)} = (\mu_k^{(t)}, \sigma_k^{(t)}, \pi_k^{(t)})_{k \in [K]}$$

posterior probability

$$\begin{aligned} & p(Z_i=k | x_i, \theta^{(t)}) \\ &= \mathbb{E}_{Z_i|x_i, \theta^{(t)}} [Z_{ik}] \\ &= \frac{\pi_k p(x_i | Z_i=k, \mu_k^{(t)}, \sigma_k^{(t)^2})}{\sum_{l \in [K]} \pi_l p(x_i | Z_i=l, \mu_l^{(t)}, \sigma_l^{(t)^2})} \\ &= r_{ik} = Z_{ik}^{(t)} \end{aligned}$$

(c)

$$\begin{aligned} & \mathbb{E}_{Z_i|x_i, \theta^{(t)}} [\log p(x_i, Z_i | \theta)] \\ &= \sum_{k \in [K]} \mathbb{E}_{Z_i|x_i, \theta^{(t)}} [Z_{ik} \log (\pi_k N(x_i | \mu_k^{(t)}, \sigma_k^{(t)^2}))] \\ &= \sum_{k \in [K]} r_{ik} \log (\pi_k N(x_i | \mu_k^{(t)}, \sigma_k^{(t)^2})) \\ &= \sum_{k \in [K]} Z_{ik}^{(t)} \log (\pi_k N(x_i | \mu_k^{(t)}, \sigma_k^{(t)^2})) \end{aligned}$$

(d) M-step to compute $\theta^{(t+1)}$ from $\theta^{(t)}$ and $Z_{ik}^{(t)}$'s

$$\theta^{(t+1)} = (\mu_k^{(t+1)}, \sigma_k^{(t+1)}, \pi_k^{(t+1)})_{k \in [K]} = \arg \max_{\theta} \sum_{i \in [n]} \mathbb{E}_{Z_i|x_i, \theta^{(t)}} [\log p(x_i, z_i | \theta)]$$

$$= \arg \max_{\theta} Q(\theta; \theta^{(t)})$$

① mean

$$\frac{\partial Q}{\partial \mu_k} = - \frac{1}{\sigma_k^{(t)2}} \sum_{i \in [n]} Z_{ik}^{(t)} (x_i - \mu_k^{(t)}) = 0$$

$$\Rightarrow \mu_k^{(t+1)} = \frac{\sum_{i \in [n]} Z_{ik}^{(t)} x_i}{\sum_{i \in [n]} Z_{ik}^{(t)}}$$

② variance

$$\frac{\partial Q}{\partial \sigma_k^2} = \sum_{i \in [n]} Z_{ik}^{(t)} \left[-\frac{D}{\sigma_k^{(t)2}} + \frac{1}{\sigma_k^{(t)3}} \|x_i - \mu_k^{(t)}\|^2 \right] = 0$$

$$\Rightarrow \sigma_k^{(t+1)} = \frac{1}{D} \frac{\sum_{i \in [n]} Z_{ik}^{(t)} \|x_i - \mu_k^{(t+1)}\|^2}{\sum_{i \in [n]} Z_{ik}^{(t)}}$$

③ mixing parameter

$\{\pi_k\}$ must verify $\sum_{k \in [K]} \pi_k = 1$. Hence, recalling the theory of constrained optimization, consider the Lagrangian

$$Q'(\theta, \lambda) = Q(\theta) + \lambda \left(1 - \sum_{k \in [K]} \pi_k \right)$$

$$\frac{\partial Q'(\theta, \lambda)}{\partial \pi_k} = \sum_{i \in [n]} \frac{z_{ik}^{(t)}}{\pi_k^{(t)}} - \lambda = 0$$

$$\Rightarrow \pi_k^{(t+1)} = \frac{1}{n} \sum_{i \in [n]} z_{ik}^{(t)}$$

(e) For every cluster $k \in [K]$, when we choose

$$\pi_k = \frac{1}{K}, \quad \sigma_k^2 = \sigma^2 \quad \text{and} \quad \tau^2 \text{ goes to } 0,$$

then EM algorithm for GMM will be same with K-means.