

$$\begin{aligned}
1. (a) \quad I(X, Y) &= \sum_y \sum_x p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \\
&= \sum_y \sum_x p(x, y) \left[ \log \frac{p(x, y)}{p(x)} - \log p(y) \right] \\
&= \sum_y \sum_x p(x, y) [\log p(y|x) - \log p(y)] \\
&= \sum_y \sum_x p(x, y) (-\log p(y)) + \sum_y \sum_x p(x, y) \log p(y|x) \\
&= - \sum_y \log p(y) \left( \sum_x p(x, y) \right) - \left[ - \left( \sum_y \sum_x p(y|x) p(x) \log p(y|x) \right) \right] \\
&= - \sum_y \log p(y) \cdot p(y) - \left[ - \sum_x p(x) \sum_y p(y|x) \log p(y|x) \right] \\
&= -H(Y) - \sum_x p(x) H(Y|X=x) = H(Y) - H(Y|X)
\end{aligned}$$

Since  $I(X, Y)$  is symmetric for  $X, Y$ , similarly we can get  $I(X, Y) = H(X) - H(X|Y)$

$$\begin{aligned}
(b) \quad \text{From (a) we know, } I(X, Y) &= H(X) - H(X|Y), \Rightarrow \\
I(X, Y) &= H(X) + \sum_x p(x) \sum_y p(y|x) \log p(y|x). \text{ Since } X=f(Y), \text{ then} \\
p(y|x) &= 0 \text{ or } 1, \text{ for each case, } \sum_x p(x) \sum_y p(y|x) \log p(y|x) = 0, \text{ then} \\
I(X, Y) &= H(X) + 0 = H(X). \text{ Similarly, } I(X, Y) = H(Y)
\end{aligned}$$

$$\begin{aligned}
(c) \quad \min_{\theta} D_{KL}(\hat{p} \| q) &= H(\hat{p}, q) - H(\hat{p}) \propto H(\hat{p}, q) \\
&\quad (\text{Noticing only } q \text{ contains } \theta \text{ here})
\end{aligned}$$

$$\begin{aligned}
H(\hat{p}, q) &= - \sum_x \hat{p}(x) \log q(x|\theta) \\
&= - \sum_x \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{x=x_i\}} \log q(x|\theta) \\
&= - \frac{1}{N} \sum_{i=1}^N \log(q(x_i|\theta))
\end{aligned}$$

$$\begin{aligned}
\text{Thus } \min_{\theta} D_{KL}(\hat{p} \| q) &\Leftrightarrow \min_{\theta} - \frac{1}{N} \sum_{i=1}^N \log q(x_i|\theta) \\
&\Leftrightarrow \max_{\theta} \sum_{i=1}^N \log q(x_i|\theta) \\
&\Leftrightarrow \max_{\theta} \prod_{i=1}^N q(x_i|\theta) \rightarrow \text{likelihood}
\end{aligned}$$

Thus the minimum Kullback-Leibler divergence is obtained by the maximum likelihood.

(d)

For continuous variable,

$$\text{we have } H(x) = - \int p(x) \ln p(x) dx,$$

From the question, it has three constraint:

$$\text{s.t. } \begin{cases} \int p(x) dx = 1 \\ \int x p(x) dx = \mu \\ \int (x-\mu)^2 p(x) dx = \sigma^2 \end{cases}$$

Use Lagrange multipliers to max  $H(x)$ :

$$L = - \int p(x) \ln p(x) dx + \lambda_1 \left( \int p(x) dx - 1 \right) + \lambda_2 \left( \int x p(x) dx - \mu \right) + \lambda_3 \left( \int (x-\mu)^2 p(x) dx - \sigma^2 \right)$$

$$\frac{df}{dp(x)} = - \ln p(x) - 1 + \lambda_1 + \lambda_2 x + \lambda_3 (x-\mu)^2 = 0$$

$$\Rightarrow p(x) = e^{-1+\lambda_1+\lambda_2 x+\lambda_3 (x-\mu)^2}$$

$$\text{with } \begin{cases} \int e^{-1+\lambda_1+\lambda_2 x+\lambda_3 (x-\mu)^2} dx = 1 \\ \int x e^{-1+\lambda_1+\lambda_2 x+\lambda_3 (x-\mu)^2} dx = \mu \\ \int (x-\mu)^2 e^{-1+\lambda_1+\lambda_2 x+\lambda_3 (x-\mu)^2} dx = \sigma^2 \end{cases} \Rightarrow p(x) = \frac{1}{(\sqrt{2\pi}\sigma)^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

 $\therefore p = N(\mu, \sigma^2)$  is the solution of max  $H(x)$ 
 $\Leftrightarrow H(q) \leq H(p)$  for  $q$  be any probability density with mean  $\mu$  and variance  $\sigma^2$ .

2.

$$\begin{aligned}
 (a) \quad \text{Dirichlet}(p|\alpha) &= \frac{T(\sum_{k=1}^m \alpha_k)}{\prod_{k=1}^m T(\alpha_k)} \prod_{k=1}^m p_k^{\alpha_k-1} \\
 &= \exp\left(\log \frac{T(\sum_{k=1}^m \alpha_k)}{\prod_{k=1}^m T(\alpha_k)} \prod_{k=1}^m p_k^{\alpha_k-1}\right) \\
 &= \exp\left[\log T(\sum_{k=1}^m \alpha_k) + \sum_{k=1}^m (\alpha_k-1) \log p_k - \sum_{k=1}^m \log T(\alpha_k)\right] \\
 &= \exp\left[\sum_{k=1}^m (\alpha_k-1) \log p_k - \left(\sum_{k=1}^m \log T(\alpha_k) - \log T(\sum_{k=1}^m \alpha_k)\right)\right]
 \end{aligned}$$

$$\text{Thus, set } \eta(\alpha) = \begin{pmatrix} \alpha_1-1 \\ \alpha_2-1 \\ \vdots \\ \alpha_m-1 \end{pmatrix}, \quad T(p) = \begin{pmatrix} \log p_1 \\ \log p_2 \\ \vdots \\ \log p_m \end{pmatrix},$$

$$A(\alpha) = \sum_{k=1}^m \log T(\alpha_k) - \log T(\sum_{k=1}^m \alpha_k). \text{ Then } D(p|\alpha) = \exp[\eta(\alpha)^T T(p) - A(\alpha)]$$

$$(b) \quad \text{The Dirichlet log-likelihood function } F(\alpha) = \log p(D|\alpha) = \sum_{j=1}^N \log D(p^j|\alpha)$$

$$\begin{aligned}
 \therefore F(\alpha) &= \sum_{j=1}^N \left[ \sum_{k=1}^m (\alpha_k-1) \log p_k^j + A(\alpha) \right] \\
 &= \sum_{k=1}^m (\alpha_k-1) \sum_{j=1}^N \log p_k^j + N A(\alpha) \\
 &= N \sum_{k=1}^m (\alpha_k-1) \frac{1}{N} \sum_{j=1}^N \log p_k^j + N \left[ \log T(\sum_{k=1}^m \alpha_k) - \sum_{k=1}^m \log T(\alpha_k) \right] \\
 &= N \left[ \sum_{k=1}^m (\alpha_k-1) \hat{t}_k + \log T(\sum_{k=1}^m \alpha_k) - \sum_{k=1}^m \log T(\alpha_k) \right], \quad \hat{t}_k = \frac{1}{N} \sum_{j=1}^N \log p_k^j
 \end{aligned}$$

$$\begin{aligned}
 (c) \quad \frac{\partial F}{\partial \alpha_k} &= N \left( \hat{t}_k + \frac{\partial \log T(\sum_{k=1}^m \alpha_k)}{\partial \sum_{k=1}^m \alpha_k} \cdot 1 - \frac{\partial \log T(\alpha_k)}{\partial \alpha_k} \right) \\
 &= N \left( \hat{t}_k + \bar{\Psi}(\sum_{k=1}^m \alpha_k) - \bar{\Psi}(\alpha_k) \right), \quad k=1, \dots, m
 \end{aligned}$$

$$\therefore \frac{\partial F}{\partial \alpha} = \left( \frac{\partial F}{\partial \alpha_1}, \frac{\partial F}{\partial \alpha_2}, \dots, \frac{\partial F}{\partial \alpha_m} \right)$$

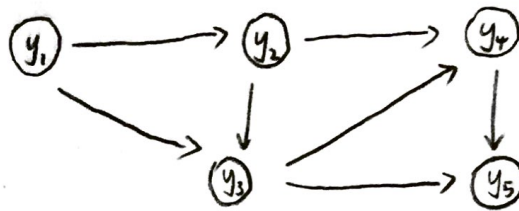
$$(d) \quad \frac{\partial^2 F}{\partial \alpha_k^2} = N \left( \bar{\Psi}'(\sum_{k=1}^m \alpha_k) - \bar{\Psi}'(\alpha_k) \right), \quad \frac{\partial^2 F}{\partial \alpha_k \partial \alpha_j} = N \bar{\Psi}'(\sum_{k=1}^m \alpha_k)$$

$$\text{Thus set } q_{kk} = -N \bar{\Psi}'(\alpha_k), \quad C = N \bar{\Psi}'(\sum_{k=1}^m \alpha_k).$$

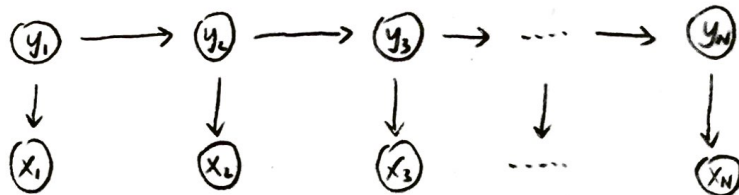
$$\text{then } \nabla_{\alpha}^T F(\alpha) = \left( \frac{\partial F}{\partial \alpha_k \alpha_j} \right)_{k,j=1,\dots,m} = (Q + C) \mathbf{1}^T$$

$$\begin{aligned}
 (e) \quad \alpha^{\text{new}} &= \alpha^{\text{old}} - [H_F(\alpha^{\text{old}})]^{-1} \nabla F(\alpha^{\text{old}}) \\
 &= \alpha^{\text{old}} - (Q + C \mathbf{1} \mathbf{1}^T)^{-1} \cdot \nabla F(\alpha^{\text{old}}) \\
 &= \alpha^{\text{old}} - \left( Q^{-1} - \frac{Q^{-1} C \mathbf{1} \mathbf{1}^T Q^{-1}}{1 + \mathbf{1}^T Q^{-1} C \mathbf{1}} \right) \nabla F(\alpha^{\text{old}})
 \end{aligned}$$

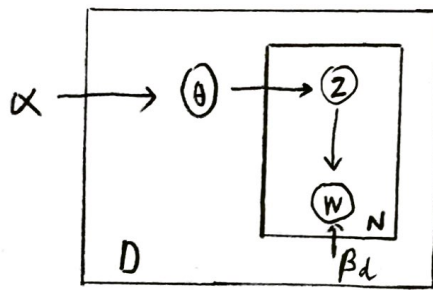
3. (a) (i)  $P(y_1, y_2, y_3, y_4, y_5) = P(y_1)P(y_2) \prod_{k=3}^5 P(y_k | y_{k-1}, y_{k-2})$



(ii)  $P(x_1, \dots, x_N, y_1, \dots, y_N) = P(y_1) \prod_{k=2}^N P(y_k | y_{k-1}) \prod_{k=1}^N P(x_k | y_k)$



(b)



(c) (1) :  $P(r, \theta, \phi, z) = \frac{1}{N!} [P(r_i) P(\theta | r_i) \prod_{j=1}^N P(\phi_{ij}) P(z_{ij} | \phi_{ij})]$

(2) :  $P(z, w) = \frac{1}{N!} [P(z_i) \prod_{j=1}^N P(w_{ij} | z_i)]$

$$5. \quad (a) \quad f(y|x; \theta) = \prod_{k=1}^K \pi_k \phi(y; w_k^T x + b_k, \sigma_k^2) \\
\Rightarrow f(y|x; \theta) = \prod_{i=1}^N \prod_{k=1}^K \pi_k \phi(y_i; w_k^T x_i + b_k, \sigma_k^2) \\
\therefore \log f(y|x; \theta) = \sum_{i=1}^N \log \sum_{k=1}^K \pi_k \phi(y_i; w_k^T x_i + b_k, \sigma_k^2)$$

$$(b) \quad f(y|x, z=k; \theta) = \phi(y; w_k^T x + b_k, \sigma_k^2), \\
f(y, z|x; \theta) = \pi_z \phi(y; w_z^T x + b_z, \sigma_z^2), \text{ Let } \Delta_{ik} = 1_{\{z_i=k\}} \\
\text{Then } \log f(y, z|x; \theta) = \sum_{k=1}^K \sum_{i=1}^N 1_{\{z_i=k\}} \log \pi_k \phi(y_i; w_k^T x_i + b_k, \sigma_k^2) \\
= \sum_{k=1}^K \sum_{i=1}^N \Delta_{ik} \log \pi_k \phi(y_i; w_k^T x_i + b_k, \sigma_k^2)$$

$$(c) \quad Q(\theta, \theta^{old}) = E_z [\log f(y, z|x; \theta) | y, x; \theta^{old}]$$

$$\text{Let } Y_{ik} = p(z_i=k | y_i, x_i; \theta^{old})$$

$$Y_{ik} = \frac{f(z_i, y_i | x_i; \theta^{old})}{f(y_i | x_i; \theta^{old})} = \frac{\pi_k^{old} \phi(y_i, w_k^{old} x_i + b_k^{old}, \sigma_k^{old})}{\sum_{k=1}^K \pi_k^{old} \phi(y_i, w_k^{old} x_i + b_k^{old}, \sigma_k^{old})}$$

$$Q(\theta, \theta^{old}) = \sum_{i=1}^N \sum_{k=1}^K E_z [1_{\{z_i=k\}} \log \pi_k \phi(y, w_k^T x + b_k, \sigma_k^2)] \\
= \sum_{i=1}^N \sum_{k=1}^K Y_{ik} \log \pi_k \phi(y, w_k^T x + b_k, \sigma_k^2) \\
= \sum_{i=1}^N \sum_{k=1}^K Y_{ik} \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K Y_{ik} \log \phi(y, w_k^T x + b_k, \sigma_k^2)$$

$$(d) \quad \theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old}), \text{ s.t. } \sum_{k=1}^K \pi_k = 1, \forall \pi_k > 0 \text{ for } k=1, \dots, K \\
\text{let } f = Q(\theta, \theta^{old}) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right) - \sum_{k=1}^K \alpha_k \pi_k, \quad \alpha_k = 0, \forall k=1, \dots, K$$

$$\frac{\partial f}{\partial \pi_k} = \sum_{i=1}^N \frac{Y_{ik}}{\pi_k} + \lambda = 0 \Rightarrow \pi_k = \frac{\sum_{i=1}^N Y_{ik}}{-\lambda}$$

$$\frac{\partial f}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1 = 0 \Rightarrow \sum_{k=1}^K \sum_{i=1}^N Y_{ik} / -\lambda - 1 = 0 \Rightarrow \lambda = -N \Rightarrow \pi_k = \frac{\sum_{i=1}^N Y_{ik}}{N}$$

$$\frac{\partial f}{\partial w_k} = \sum_{i=1}^N Y_{ik} \left( \frac{\log \phi(y, w_k^T x + b_k, \sigma_k^2)}{\partial w_k} \right) = \sum_{i=1}^N Y_{ik} \frac{\partial \left( \log \frac{1}{\sqrt{2\pi}\sigma_k} + \frac{y_i - (w_k^T x + b_k)}{2\sigma_k^2} \right)}{\partial w_k} \\
= \sum_{i=1}^N Y_{ik} \frac{2(y_i - w_k^T x + b_k)}{2\sigma_k^2} (-x_i) = 0$$

$$\frac{\partial f}{\partial b_k} = \sum_{i=1}^N Y_{ik} \frac{\partial \left( \log \frac{1}{\sqrt{2\pi}\sigma_k} + \frac{y_i - (w_k^T x + b_k)}{2\sigma_k^2} \right)}{\partial b_k} = \sum_{i=1}^N Y_{ik} \frac{y_i - (w_k^T x + b_k)}{\sigma_k^2} \cdot (-1) = 0$$

$$\text{Set } \tilde{w}_k = \begin{bmatrix} w_k \\ b_k \end{bmatrix}, \quad \tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} = [x, 1], \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$\text{then we can combine } \frac{\partial f}{\partial w_k} \text{ \& } \frac{\partial f}{\partial b_k} \text{ to get } \sum_{i=1}^N Y_{ik} (y_i - \tilde{w}_k^T \tilde{x}_i) \tilde{x}_i = 0$$



Use closed form to represent the whole question, we can set the objective function as

$$\sum_{i=1}^N Y_{ik} (y_i - \tilde{w}_k^T \tilde{x}_i)^2 \text{ can be written as } (\tilde{x} \tilde{w}_k - y)^T R (\tilde{x} \tilde{w}_k - y)$$

$$\text{where } R = \text{diag}(Y_{1k}, Y_{2k}, \dots, Y_{Nk})$$

$$\therefore \sum_{i=1}^N Y_{ik} (y_i - \tilde{w}_k^T \tilde{x}_i) \tilde{x}_i = 2 \tilde{x}^T R \tilde{x} w - 2 \tilde{x}^T R y = 0$$

$$\Rightarrow \tilde{w}_k = (\tilde{x}^T R \tilde{x})^{-1} \tilde{x}^T R y$$

$$\frac{\partial f}{\partial \tilde{b}_k} = \frac{\partial \sum_{i=1}^N Y_{ik} \left( \log \frac{1}{\sqrt{2\pi \tilde{b}_k}} + \left( -\frac{1}{2\tilde{b}_k} (y_i - \tilde{w}_k^T \tilde{x}_i)^2 \right) \right)}{\partial \tilde{b}_k}$$

$$= \sum_{i=1}^N Y_{ik} + \frac{1}{\tilde{b}_k} \sum_{i=1}^N (y_i - \tilde{w}_k^T \tilde{x}_i)^2 Y_{ik} = 0$$

$$\Rightarrow \tilde{b}_k = \frac{1}{\sum_{i=1}^N Y_{ik}} \sum_{i=1}^N (y_i - \tilde{w}_k^T \tilde{x}_i)^2 Y_{ik}$$

$$= \frac{1}{\sum_{i=1}^N Y_{ik}} (\tilde{x} \tilde{w}_k - y)^T R (\tilde{x} \tilde{w}_k - y)$$