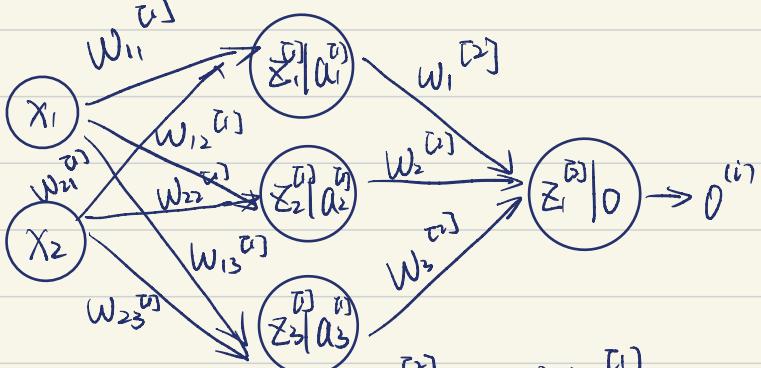


Question 1:

a.



$$\frac{\partial L}{\partial w_{12}^{(1)}} = \frac{\partial L}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial a_2^{(1)}} \cdot \frac{\partial a_2^{(1)}}{\partial z_2^{(1)}} \cdot \frac{\partial z_2^{(1)}}{\partial w_{12}^{(1)}}$$

$$\therefore \frac{\partial L}{\partial w^{(2)}} = \frac{\partial}{\partial w^{(2)}} \left( \frac{1}{m} \sum_{i=1}^m (O^{(ii)} - y^{(ii)})^2 \right)$$

$$= \frac{\partial}{\partial w^{(2)}} \frac{1}{m} \sum_{i=1}^m \left( b(w^{(2)} a^{(1)}) - y^{(i)} \right)^2 = \frac{1}{m} \sum_{i=1}^m 2 O^{(i)} \cdot O^{(i)} \cdot (1 - O^{(i)}) \cdot a^{(1)} - 2 y^{(i)} O^{(i)} \cdot (1 - O^{(i)}) \cdot a^{(1)}$$

$$= \frac{2}{m} \sum_{i=1}^m O^{(i)} (1 - O^{(i)}) \cdot a^{(1)} (O^{(i)} - y^{(i)})$$

Forward Propagation:

$$z_2^{(1)} = w_{12}^{(1)} x_1 + w_{22}^{(1)} x_2$$

$$a_2^{(1)} = g(z_2^{(1)})$$

$$z_1^{(2)} = w_{12}^{(2)} a_1^{(1)} + w_{22}^{(2)} a_2^{(1)} + w_{32}^{(2)} a_3^{(1)}$$

$$O = g(z_1^{(2)})$$

$$L = \frac{1}{m} \sum_{i=1}^m (O^{(i)} - y^{(i)})^2$$

$$a^{[1]} = \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

$$\therefore \frac{\partial L}{\partial z^{[2]}} = \frac{2}{m} \sum_{i=1}^m o^{(i)} (1 - o^{(i)}) (o^{(i)} - y^{(i)})$$

$$\therefore \frac{\partial L}{\partial w_1^{[1]}} = \frac{2}{m} \sum_{i=1}^m \left( o^{(i)} (1 - o^{(i)}) (o^{(i)} - y^{(i)}) \right) \cdot w_2^{[2]} \cdot a_2^{[1]} (1 - a_2^{[1]}) \cdot x_1^{(i)}$$

(b) The decision boundary of "0":  $x_1 = 0.25, x_2 = 0.25, x_1 + x_2 = 4.5$

$$z_1^{[1]} = w_{01}^{[1]} + w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2$$

want: when  $x_1 > 0.25, z_1^{[1]} < 0$

$$0.25 w_{11}^{[1]} + w_{01}^{[1]} = 0 \Rightarrow w_{11}^{[1]} = -1, w_{01}^{[1]} = 0.25$$

$$\text{same for } z_2^{[1]}: 0.25 w_{22}^{[1]} + w_{02}^{[1]} = 0 \Rightarrow w_{22}^{[1]} = -1, w_{02}^{[1]} = 0.25$$

$$z_3^{[1]} = w_{03}^{[1]} + w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2$$

$$\Rightarrow w_{03}^{[1]} = 4.5, w_{13}^{[1]} = 1, w_{23}^{[1]} = 1$$

$$\Rightarrow W^{[1]} = \begin{pmatrix} 0.25 & -1 & 0 \\ 0.25 & 0 & -1 \\ -4.5 & 1 & 1 \end{pmatrix}$$

for output layer, if 0, then  $W_{01}^{[2]} + 0 < 0 \Rightarrow W_{01}^{[2]} = -1$

$$\Rightarrow W^{[2]} = (-1 \quad 1 \quad 1 \quad 1)$$

(c) The problem is not linearly separable. Therefore, it can not make the loss 0.

Question 2:

$$MAP = \left( \frac{m}{l} \sum_{z^{(i)}} p(x^{(i)}, z^{(i)} | \theta) \right) p(\theta).$$

$$\begin{aligned}\Rightarrow l(\theta) &= \log \left( \frac{m}{T} \sum_{i=1}^m P(x^{(i)}, z^{(i)} | \theta) \right) \cdot P(\theta) \\ &= \log P(\theta) + \sum_{i=1}^m \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)} | \theta) \\ &= \log P(\theta) + \sum_{i=1}^m \sum_{z^{(i)}} \log P(x^{(i)}, z^{(i)} | \theta).\end{aligned}$$

$$\begin{aligned}\text{E-step: } l(\theta) &= \log P(\theta) + \sum_{i=1}^m \sum_{z^{(i)}} \log Q_i(z^{(i)}) \cdot \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \\ &= \log P(\theta) + \sum_{i=1}^m \log \bar{E}_{z^{(i)} \sim Q_i} \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \geq \log P(\theta) + \sum_{i=1}^m \bar{E}_{z^{(i)} \sim Q_i} \log \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \\ &\quad = \log P(\theta) + \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \cdot \log \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})}\end{aligned}$$

$$\text{Want: } \log \bar{E}_{z^{(i)} \sim Q_i} \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} = \bar{E}_{z^{(i)} \sim Q_i} \log \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})}$$

$$\Rightarrow \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \text{ is constant. } \Rightarrow Q_i(z^{(i)}) = \frac{1}{\lambda} P(x^{(i)}, z^{(i)} | \theta).$$

$$\therefore \sum_{z^{(i)}} Q_i(z^{(i)}) = 1 \Rightarrow \lambda = \sum_{z^{(i)}} P(x^{(i)}, z^{(i)} | \theta).$$

$$\Rightarrow Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)} = \frac{P(x^{(i)}, z^{(i)}; \theta)}{P(x^{(i)}; \theta)} = P(z^{(i)} | x^{(i)}; \theta)$$

M-step:  $\theta := \arg\max_{\theta} \sum_i^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} + \log P(\theta).$   
 Take the gradient of  $\theta$  and set to 0.

Proof of convergence: see lecture notes

Question 3:

$$(a) i. x^{pr} = y^{pr} + z^{pr} + \epsilon^{pr}$$

$y^{pr}, z^{pr}, \epsilon^{pr}$  are all Gaussian

$\Rightarrow x^{pr}$  is Gaussian

$\Rightarrow$  the mean for  $(x^{pr}, y^{pr}, z^{pr})$  is  $(\mu_p + \nu_r, \mu_p, \nu_r)^T$

and the covariance is

$$\begin{pmatrix} \hat{\sigma}_p^2 + \hat{\sigma}_r^2 + \hat{\sigma}^2 & \hat{\sigma}_p^2 & \hat{\sigma}_r^2 \\ \hat{\sigma}_p^2 & \hat{\sigma}_p^2 & 0 \\ \hat{\sigma}_r^2 & 0 & \hat{\sigma}_r^2 \end{pmatrix}$$

ii. E-step:  $\frac{P(x^{pr}, y^{pr}, z^{pr})}{Q_{pr}(y^{pr}, z^{pr})}$  is constant.

$$\therefore \sum_{r=1}^R \sum_{p=1}^P Q_{pr}(y^{pr}, z^{pr}) = 1.$$

$$\therefore Q_{pr} = \frac{P(x^{pr}, y^{pr}, z^{pr})}{\sum_{r=1}^R \sum_{p=1}^P P(x^{pr}, y^{pr}, z^{pr})} = \frac{P(x^{pr}, y^{pr}, z^{pr})}{P(x^{pr})} = P(y^{pr}, z^{pr} | x^{pr}).$$

$\Rightarrow Q_{pr}$  is also a Gaussian.

$$\boldsymbol{\mu} = [\mu_p \ \nu_r] + \frac{x^{pr} - \mu_p - \nu_r}{\hat{\sigma}_p^2 + \hat{\sigma}_r^2 + \hat{\sigma}^2} \begin{bmatrix} \hat{\sigma}_p^2 & \hat{\sigma}_r^2 \\ \hat{\sigma}_r^2 & \hat{\sigma}_r^2 \end{bmatrix}$$

$$\boldsymbol{\Sigma} = [\hat{\sigma}_p^2 \ 0 \ 0 \ \hat{\sigma}_r^2] - \frac{1}{\hat{\sigma}_p^2 + \hat{\sigma}_r^2 + \hat{\sigma}^2} \begin{bmatrix} \hat{\sigma}_p^4 & \hat{\sigma}_p^2 \hat{\sigma}_r^2 & \hat{\sigma}_r^2 \hat{\sigma}_p^2 & \hat{\sigma}_r^4 \\ \hat{\sigma}_p^2 \hat{\sigma}_r^2 & \hat{\sigma}_r^4 & 0 & \hat{\sigma}_r^4 \\ \hat{\sigma}_r^2 \hat{\sigma}_p^2 & 0 & \hat{\sigma}_p^4 & \hat{\sigma}_p^4 \\ \hat{\sigma}_r^4 & \hat{\sigma}_r^4 & \hat{\sigma}_p^4 & \hat{\sigma}_r^4 \end{bmatrix}$$

$$(b) \text{ M-step: } \underset{M_p, O_p, V_r, J_r}{\operatorname{argmax}} \frac{P(x^{\text{pr}}, y^{\text{pr}}, z^{\text{pr}}; M_p, O_p, V_r, J_r)}{O_{\text{pr}}}$$

use gradient descent, set all the gradient to 0.

Question 4:

$$(a) \text{ KL}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = - \sum_x P(x) \log \frac{Q(x)}{P(x)} \stackrel{\text{Jensen inequality}}{\geq} -\log \sum_x P(x) \cdot \frac{Q(x)}{P(x)}$$

$$\geq -\log \sum_x Q(x) \geq -\log 1 \geq 0.$$

$$\text{If } P=Q, \text{ then } \text{KL}(P||Q) = \sum_x P(x) \log 1 = 0.$$

$$\text{If } \text{KL}(P||Q)=0, \text{ then } - \sum_x P(x) \log \frac{Q(x)}{P(x)} = -\log \sum_x P(x) \cdot \frac{Q(x)}{P(x)}$$

$$\Rightarrow \frac{Q(x)}{P(x)} \text{ is constant} \Rightarrow P=Q.$$

$$(b) \text{KL}(P(x) || Q(x)) + \text{KL}(P(Y|x) || Q(Y|x))$$

$$= \sum_x P(x) \log \frac{P(x)}{Q(x)} + \sum_x P(x) \left( \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right).$$

$$= \sum_x P(x) \left[ \log \frac{P(x)}{Q(x)} + \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right]$$

$$\therefore \sum_y P(y|x) = 1$$

$$\therefore = \sum_x P(x) \sum_y P(y|x) \left( \log \frac{P(x)}{Q(x)} + \log \frac{P(y|x)}{Q(y|x)} \right)$$

$$= \sum_x P(x) \sum_y P(y|x) \left( \log \frac{P(x) \cdot P(y|x)}{Q(x) \cdot Q(y|x)} \right)$$

$$= \sum_x P(x) \sum_y P(y|x) \log \frac{P(x,y)}{Q(x,y)} = \sum_x P(x,y) \log \frac{P(x,y)}{Q(x,y)} = \text{KL}(P(x,y) || Q(x,y))$$

$$\begin{aligned}
 (c) \quad KL(\hat{P} || P_{\theta}) &= \sum_{\hat{x}} \hat{P}(\hat{x}) \log \frac{\hat{P}(\hat{x})}{P_{\theta}(\hat{x})} = - \sum_{\hat{x}} \hat{P}(\hat{x}) \log \frac{P_{\theta}(\hat{x})}{\hat{P}(\hat{x})} \\
 &= -\frac{1}{m} \sum_{\hat{x}} \sum_{i=1}^m I\{x^{(i)} = \hat{x}\} \cdot \log \frac{P(x^{(i)}, \theta)}{\frac{1}{m} \sum_{i=1}^m I\{x^{(i)} = \hat{x}\}} \\
 &= -\frac{1}{m} \sum_{i=1}^m \log P_{\theta}(x^{(i)})
 \end{aligned}$$

Thus, minimizing  $KL(\hat{P} || P_{\theta})$  is equivalent to maximizing  $\sum_{i=1}^m \log P_{\theta}(x^{(i)}) = L(\theta)$ .