

1. (a)

$$h_2^{(i)} = \sigma(W_{1,2}^{[1]} x_1^{(i)} + W_{2,2}^{[1]} x_2^{(i)})$$

$$o^{(i)} = \sigma(W_1^{[2]} h_1^{(i)} + W_2^{[2]} h_2^{(i)} + W_3^{[2]} h_3^{(i)})$$

$$\frac{\partial L}{\partial W_{1,2}^{[1]}} = \sum_{i=1}^m \frac{\partial L}{\partial o^{(i)}} \frac{\partial o^{(i)}}{\partial h_2^{(i)}} \frac{\partial h_2^{(i)}}{\partial W_{1,2}^{[1]}}$$

$$= \frac{1}{m} \sum_{i=1}^m 2(o^{(i)} - y^{(i)}) o^{(i)}(1-o^{(i)}) \cdot W_1^{[2]} h_1^{(i)} (1-h_2^{(i)}) x_1^{(i)}$$

$$\text{so } W_{1,2}^{[1]} := W_{1,2}^{[1]} - \alpha \frac{2}{m} W_1^{[2]} \sum_{i=1}^m (o^{(i)} - y^{(i)}) o^{(i)}(1-o^{(i)})$$

$$h_2^{(i)} (1-h_2^{(i)}) x_1^{(i)}$$

$$\text{where } h_2^{(i)} = \sigma(W_{1,2}^{[1]} x_1^{(i)} + W_{2,2}^{[1]} x_2^{(i)})$$

(b)

Yes, it's possible. All the data of class 0 is inside a triangle, and all the examples outside the triangle is 1. The hidden layer consists of 3 neurons, each of which can be seened as if the point is on the right side of the triangle. The 3 sides is approximately  $x_1 = 0.5$ ,  $x_2 = 0.5$ ,  $x_1 + x_2 = 4$ .

(c) No, if  $f(x) = x$ ,

$\vec{h} = W^{[1]} \vec{x}$  is a linear transformation, but the dataset is not linear separable.

2.(a)

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}$$

$$= E_{x \sim P(x)} \left[ \log \frac{P(x)}{Q(x)} \right]$$

$$= E_{x \sim P(x)} \left[ -\log \frac{Q(x)}{P(x)} \right]$$

$$\geq -\log \left[ E_{x \sim P(x)} \frac{Q(x)}{P(x)} \right]$$

$$= -\log \sum_{x \in X} P(x) \frac{Q(x)}{P(x)}$$

$$= -\log \sum_{x \in X} Q(x) = -\log 1 = 0$$

if and only if  $\log \frac{P(x)}{Q(x)}$  is const,

$$D_{KL}(P||Q) = 0$$

$$\text{which implies } \frac{P(x)}{Q(x)} = C$$

$$\text{because } \sum_{x \in X} P(x) = \sum_{x \in X} Q(x) = 1$$

so it implies  $P(x) = Q(x)$

(b)  $D_{KL}(P(X, Y) || Q(X, Y))$

$$= \sum_{X, Y} P(x, y) \log \frac{P(x, y)}{Q(x, y)}$$

$$= \sum_{X, Y} P(y|X) P(x) \log \frac{P(y|X) P(x)}{Q(y|x) Q(x)}$$

$$= \sum_X P(x) \left( \sum_Y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right) + \sum_{X, Y} P(y|x) P(x) \log \frac{P(x)}{Q(x)}$$

$$= D_{KL}(P(X) || Q(X))$$

$$+ \sum_X P(x) \log \frac{P(x)}{Q(x)}$$

$$= D_{KL}(P(X) || Q(X)) +$$

$$D_{KL}(P(Y|X) || Q(Y|X))$$

$$(C) \arg \min_{\hat{P}} D_{KL}(\hat{P} || P_{\theta})$$

$$= \arg \min_{\theta} \sum_{x \in X} \hat{P}(x) \log \frac{\hat{P}(x)}{P_{\theta}(x)}$$

$$= \arg \max_{\theta} \sum_{x \in X} \hat{P}(x) \log P_{\theta}(x)$$

$$= \arg \max_{\theta} \sum_{x \in X} \sum_{i=1}^m \mathbb{I}_{\{X^{(i)}=x\}} \log P_{\theta}(x)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log P_{\theta}(x^{(i)})$$

3. (a)

$$\begin{aligned} & \int_{-\infty}^{+\infty} p(y; \theta) \nabla_{\theta} \log P(y; \theta) dy \\ &= \int_{-\infty}^{+\infty} P(y; \theta) \frac{\partial_p P(y; \theta)}{P(y; \theta)} dy \\ &= \int_{-\infty}^{+\infty} \nabla_{\theta} P(y; \theta) dy \\ &= \nabla_{\theta} \int_{-\infty}^{+\infty} P(y; \theta) dy = \nabla_{\theta} I = 0 \end{aligned}$$

$$(b) \text{Conf} \hat{x} \in E[(x - \bar{x})(x - \bar{x})^T]$$

$$\text{so } I(\theta) = E_{y \sim p(y; \theta)} [(\nabla_{\theta} \log(y; \theta) - E(\nabla_{\theta} \log(y; \theta)))^T] = \frac{1}{2} d^T I(\theta) d$$

$$(\nabla_{\theta} \log(y; \theta) - E(\nabla_{\theta} \log(y; \theta)))^T]$$

$$= E_{y \sim p(y; \theta)} [\nabla_{\theta} \log(y; \theta) \nabla_{\theta} \log(y; \theta)^T]$$

$$(c) I(\theta)_{(i,j)} = E_{y \sim p(y; \theta)} [\nabla_{\theta} \log(y; \theta)_i \nabla_{\theta} \log(y; \theta)_j]$$

$$= \int_{-\infty}^{+\infty} \frac{1}{p(y; \theta)} \frac{\partial P}{\partial \theta_i} \frac{\partial P}{\partial \theta_j} dy$$

$$\int_{-\infty}^{+\infty} p(y; \theta) \left( -\frac{\partial^2 \log P(y; \theta)}{\partial \theta_i \partial \theta_j} \right) dy$$

$$= \int_{-\infty}^{+\infty} p(y; \theta) \left( -\frac{\partial}{\partial \theta_i} \left( \frac{1}{p(y; \theta)} \frac{\partial P}{\partial \theta_j} \right) \right) dy$$

$$= \int_{-\infty}^{+\infty} p(y; \theta) \left( \frac{1}{p(y; \theta)} \frac{\partial P}{\partial \theta_i} \frac{\partial P}{\partial \theta_j} - \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \frac{1}{p(y; \theta)} \right) dy$$

$$= I(\theta)_{ij} - \int_{-\infty}^{+\infty} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} dy = I(\theta)_{ij} - \frac{\partial^2}{\partial \theta_i \partial \theta_j} \int_{-\infty}^{+\infty} p(y; \theta) dy$$

$$= I(\theta)_{ij}$$

$$E_{y \sim p(y; \theta)} [-\nabla_{\theta}^2 \log P(y; \theta)]_{\theta=\theta_0} = I(\theta)$$

$$(d) \log P(y; \theta) = \log P(y; \theta_0)$$

$$+ d^T \nabla_{\theta} (\log(y; \theta'))|_{\theta=\theta_0}$$

$$+ \frac{1}{2} d^T (\nabla_{\theta}^2 f(\theta')|_{\theta=\theta_0}) d$$

$$E_{y \sim p(y; \theta_0)} [\log P(y; \theta')] = E_{y \sim p(y; \theta_0)} [\log P(y; \theta)]$$

$$+ d^T E_{y \sim p(y; \theta_0)} [\nabla_{\theta} \log(y; \theta')|_{\theta=\theta_0}]$$

$$- \frac{1}{2} d^T I(\theta) d$$

$$\text{so } E_{y \sim p(y; \theta)} [\log P(y; \theta)] - E_{y \sim p(y; \theta_0)} [\log P(y; \theta)]$$

$$D_{KL}(P_{\theta} || P_{\theta+\Delta}) = \frac{1}{2} d^T I(\theta) d$$

$$(e) \frac{1}{2} d^T I(\theta) d = c$$

$$d^* = \arg \max_d d^T \nabla_{\theta} l(\theta')|_{\theta=\theta_0}$$

$$L(d, \lambda) = d^T \nabla_{\theta} l(\theta')|_{\theta=\theta_0} - \lambda (\frac{1}{2} d^T I d - c)$$

$$\nabla_d L = \nabla_{\theta} l(\theta')|_{\theta=\theta_0} - \lambda I d = 0$$

$$d = \frac{1}{\lambda} I^{-1}(\theta) \nabla_{\theta} l(\theta')|_{\theta=\theta_0}$$

$$\nabla_{\lambda} L = -\frac{1}{2} d^T I(\theta) d + c = 0$$

$$2\lambda^2 c = \nabla_{\lambda} L(\theta')|_{\theta=\theta_0}^T I^{-1}(\theta) \nabla_{\theta} l(\theta')|_{\theta=\theta_0}$$

$$\lambda = \sqrt{\frac{\nabla_{\lambda} L(\theta')|_{\theta=\theta_0}^T I^{-1}(\theta) \nabla_{\theta} l(\theta')|_{\theta=\theta_0}}{2c}}$$

$$\text{so } d^* = \sqrt{\frac{2c}{P_{\theta} P_{\theta'} I^{-1}(\theta) P_{\theta'}}} I^{-1}(\theta) \nabla_{\theta} l(\theta)$$

$$(f) \text{Newton's method: } \theta := \theta - H^{-1} \nabla_{\theta} L(\theta)$$

Nature gradient:

$$\theta - \theta + \frac{1}{\lambda} I^{-1}(\theta) \nabla_{\theta} L(\theta) = \theta - \frac{1}{\lambda} E(H^{-1}) \nabla_{\theta} L(\theta)$$

4. (1)

$$l_{\text{semi-sup}}(\theta^{(t+1)}) \geq \sum_{j=1}^m \sum_{z^{(i)}} Q^{(i)} P(x^{(i)}, z^{(i)}; \theta^{(t+1)}) \\ + \alpha \sum_{i=1}^m \log P(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta^{(t+1)})$$

$$\geq \sum_{i=1}^m \sum_{z^{(i)}} Q^{(i)} \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q^{(i)}} + \alpha \sum_{i=1}^m \log P(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta^{(t+1)})$$

$$= l_{\text{semi-sup}}(\theta^{(t+1)})$$

$$(b) w_j^{(i)} = P(z=j | \mathbf{x}) = \frac{P(x|z=j) P(z=j)}{p(x)}$$

$$= \frac{P(x|z=j) p(z=j)}{\sum_{l=1}^k P(x|z=l) P(z=l)} \\ = \frac{\frac{1}{|\Sigma_j|} \exp(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j))}{\sum_{l=1}^k \frac{1}{|\Sigma_l|} \exp(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_l)^T \Sigma_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l))} \phi_j$$

$$(c) ELBO(x^{(i)}, Q^{(i)}, \theta) =$$

$$\sum_{j=1}^m \left( \sum_{i=1}^k w_j^{(i)} \log \frac{P(x^{(i)}|z^{(i)}=j) P(z^{(i)}_j)}{w_j^{(i)}} \right)$$

$$+ \alpha \sum_{i=1}^m \log P(\bar{x}^{(i)} | \bar{z}^{(i)}) P(\bar{z}^{(i)})$$

$$= \sum_{j=1}^m \left( \sum_{i=1}^k w_j^{(i)} \frac{\frac{1}{(2\pi)^d} \exp(-\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j))}{w_j^{(i)}} \right)$$

$$+ \alpha \sum_{i=1}^m \log \frac{1}{(2\pi)^d |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)\right)$$

$$\nabla_{M_j} ELBO_m = \nabla_{M_j} \sum_{i=1}^m \frac{1}{2} w_j^{(i)} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)$$

$$= \sum_{i=1}^m w_j^{(i)} \sum_j (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)$$

$$\nabla_{M_j} ELBO_{\text{sup}} = \sum_{j=1}^m |\{z^{(i)}=j\} \sum_j (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)$$

$$\nabla_{M_j} ELBO$$

$$= \sum_{j=1}^m \left( \left( \sum_{i=1}^k w_j^{(i)} x^{(i)} \right) + \alpha \sum_{i=1}^m \log P(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta) \right) - M_j \left( \sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^m |\{z^{(i)}=j\}| \right)$$

$$M_j = \frac{\sum_{i=1}^k w_j^{(i)} x^{(i)}}{\sum_{i=1}^k w_j^{(i)} + \alpha \sum_{i=1}^m |\{z^{(i)}=j\}|}$$

$$\nabla_{\Sigma_j} ELBO_m$$

$$= \sum_{i=1}^m \sum_{j=1}^m w_j^{(i)} \left( -\frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right)$$

$$= \sum_{i=1}^m w_j^{(i)} \left( -\frac{1}{2} |\Sigma_j|^{-1} + \frac{1}{2} \sum_j (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} \right)$$

$$\nabla_{\Sigma_j} ELBO_{\text{sup}}$$

$$= \sum_{i=1}^m \left( \sum_{j=1}^m w_j^{(i)} \left( -\frac{1}{2} |\Sigma_j|^{-1} + \frac{1}{2} \sum_j (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) \right)$$

$$\nabla_{\Sigma_j} ELBO =$$

$$-\frac{1}{2} \left( \sum_{j=1}^m w_j^{(i)} \sum_j |\Sigma_j|^{-1} + \alpha \sum_{i=1}^m |\{z^{(i)}=j\}| \sum_j |\Sigma_j|^{-1} \right) \\ + \frac{1}{2} \left( \sum_{j=1}^m w_j^{(i)} \sum_j (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} \right. \\ \left. + \sum_{j=1}^m |\{z^{(i)}=j\}| \sum_j (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} \right)$$

$$= 0 \\ \Sigma_j = \frac{\sum_{i=1}^m w_j^{(i)} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T + \alpha \sum_{i=1}^m |\{z^{(i)}=j\}| (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T}{\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^m |\{z^{(i)}=j\}|}$$

$$P_{\phi_j} \text{ELBO}_{un} = \sum_{i=1}^m \frac{w_j^{(i)}}{\phi_j}$$

$$P_{\phi_j} \text{ELBO}_{sup} = \sum_{i=1}^m I\{z^{(i)} = j\} \frac{1}{\phi_j}$$

$$\sum_{j=1}^k \phi_j = 1$$

$$\text{so } L(\phi, \lambda) = \text{ELBO} - \lambda \left( \sum_{j=1}^k \phi_j - 1 \right)$$

$$\nabla_{\phi_j} L = \frac{m w_j^{(i)}}{\phi_j} + \alpha \sum_{i=1}^m I\{z^{(i)} = j\} - \lambda = 0$$

$$\phi_j = \frac{\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^m I\{z^{(i)} = j\}}{\lambda}$$

$$\sum_{j=1}^k \phi_j = 1$$

$$\lambda = \sum_{j=1}^k \left( \frac{m}{\sum_{i=1}^m w_j^{(i)}} + \alpha \sum_{i=1}^m I\{z^{(i)} = j\} \right)$$

$$= m + \alpha \bar{m}$$

$$\phi_j = \frac{\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^m I\{z^{(i)} = j\}}{m + \alpha \bar{m}}$$

(f)

(i) for semi-supervised EM, it takes less time to converge

(ii) semi-supervised EM is more stable. For a unsupervised EM, it's possible to get different results with different initialization, but for semi-supervised EM, it always generate the same results

(iii) Semi-supervised EM has higher quality because it has lower variance for 3 distribution and higher variance for another distribution

5.(b)

$$\text{factor} = \frac{\log_2 16}{24} = \frac{4}{24} = \frac{1}{6}$$