

Problem 1: (a).

$$(1) \quad \mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{x}^{(i)} + \mathbf{b}^{[1]} : \quad \mathbf{z}_{1,1}^{[1]} = W_{1,1}^{[1]} x_1^{(i)} + b_{0,1}^{[1]}$$

$1 \times 3 \quad 3 \times 2 \quad 2 \times 1 \quad 3 \times 1$

$$(2) \quad a^{[1]} = g(\mathbf{z}^{[1]}) = \frac{1}{1 + e^{-\mathbf{z}^{[1]}}} \quad a_1^{[1]} = \frac{1}{1 + e^{-z_1^{[1]}}}$$

$$(3) \quad \mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$1 \times 1 \quad 1 \times 3 \quad 3 \times 1 \quad 1 \times 1$

$$(4) \quad a^{[2]} = g(\mathbf{z}^{[2]}) = \frac{1}{1 + e^{-\mathbf{z}^{[2]}}}$$

$$\frac{\partial}{\partial \mathbf{z}^{[2]}} : g(\mathbf{z}^{[2]}) (1 - g(\mathbf{z}^{[2]}))$$

$$\frac{\partial \mathbf{z}^{[2]}}{\partial a_2^{[1]}} : W_2^{[2]}$$

$$\frac{\partial a_{(2)}^{[1]}}{\partial z_{(2)}^{[1]}} : g(\mathbf{z}_2^{[1]}) (1 - g(\mathbf{z}_2^{[1]}))$$

$$\frac{\partial \mathbf{z}_{(2)}^{[1]}}{\partial W_{1,2}^{[1]}} : x_1^{(i)}$$

$$\frac{\partial \mathcal{L}}{\partial W_{1,2}^{[1]}} : g(\mathbf{z}^{[2]}) (1 - g(\mathbf{z}^{[2]})) W_2^{[2]} \cdot g(\mathbf{z}_2^{[1]}) (1 - g(\mathbf{z}_2^{[1]})) \cdot x_1$$

$$W_{1,2}^{[1]} := W_{1,2}^{[1]} - \eta \frac{\partial \mathcal{L}}{\partial W_{1,2}^{[1]}}$$

(b). It is possible.

The dataset can be perfectly separated by a triangle. For ~~each~~ ^{1st} hidden layer, it determines if the point is on the left side of $x_1 = 0.5$.

The 2nd hidden layer determines $x_2 \geq 0.5$. The third $x + y - 4 \geq 0$.

If any of the above is true, then $y^{(i)} = 1$.

(c). It is not possible.

When combine all the neurons with linear activation function. The final result is linear.

There is no way to separate the dataset perfectly with a single line in 2-dimensions.

Problem 2

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

$$-D_{KL} = -\sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} = \sum_{x \in X} P(x) \log \frac{Q(x)}{P(x)} \quad \text{Let } P(x) \leq 1 \leq \log \sum_{x \in X} P(x) \cdot \frac{Q(x)}{P(x)}$$

$$\leq \log 1$$

$$-D_{KL} \leq 0$$

$$D_{KL} \geq 0.$$

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

$$D_{KL}(P||Q) = \sum P(x) \log \frac{P(x)}{Q(x)} = \sum P(x) \cdot 0 = 0.$$

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

if $D_{KL}(P||Q)=0$. then $D_{KL}(P||Q) = \sum P(x) \log \frac{P(x)}{Q(x)} = 0 = -D_{KL} = \sum P(x) \log \frac{Q(x)}{P(x)}$

$$\sum P(x) \log \frac{Q(x)}{P(x)} \leq \log \sum Q(x) = 0$$

We know that $\log x$ is ~~convex~~ strictly concave, and $\sum P(x) \log \frac{Q(x)}{P(x)} = \log \sum P(x) \cdot \frac{Q(x)}{P(x)}$

Then $\frac{Q(x)}{P(x)}$ is a constant. Suppose $\frac{Q(x)}{P(x)} = c$. =0.

$$\text{Then } \sum P(x) \log \frac{P(x)}{Q(x)} = \sum P(x) \cdot \log c = \log(c) \cdot \sum P(x) = \log(c) = 0.$$

Therefore $c=1$. thus ~~$P(x)$ and $Q(x)$~~ $P=Q$.

(b).

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

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

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

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

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

$$= \sum_{y \in Y} P(y|x) D_{KL}(P(X) \| Q(X)) + D_{KL}(P(Y|X) \| Q(Y|X))$$

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

Q.E.D.

$$c). \arg \min_{\theta} D_{KL}(\hat{P} \| P_{\theta}) = \arg \min_{\theta} \sum_i^m \hat{P}(x^{(i)}) \log \frac{\hat{P}(x^{(i)})}{P_{\theta}(x^{(i)})}$$

Since \hat{P} is uniformly distributed, $\hat{P}(x^{(i)}) = \frac{1}{m}$, then the equation is.

$$\arg \min_{\theta} \sum_i^m \frac{1}{m} \log \frac{1}{m P_{\theta}(x^{(i)})}$$

$$= \arg \min_{\theta} \frac{1}{m} \sum_i^m -\log(m P_{\theta}(x^{(i)}))$$

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

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

$$= \arg \max_{\theta} \sum_i^m \log P_{\theta}(x^{(i)}) \quad \text{which is the maximum likelihood.}$$

Problem 3

$$\begin{aligned}
 (a). \quad E_{y \sim P(y; \theta)} [\nabla_{\theta'} \log P(y; \theta') | \theta' = \theta] &= \int_{-\infty}^{\infty} P(y; \theta) \nabla_{\theta'} \log P(y; \theta') |_{\theta' = \theta} dy \\
 &= \int_{-\infty}^{\infty} P(y; \theta) \cdot \frac{1}{P(y; \theta') |_{\theta' = \theta}} \cdot \nabla_{\theta'} P(y; \theta') |_{\theta' = \theta} dy \\
 &= \int_{-\infty}^{\infty} \nabla_{\theta'} P(y; \theta') |_{\theta' = \theta} dy \\
 &= \nabla_{\theta'} \left(\int_{-\infty}^{\infty} P(y; \theta') dy \right) \\
 &= \nabla_{\theta'} 1 \\
 &= 0.
 \end{aligned}$$

$$\begin{aligned}
 (b). \quad I(\theta) &= \text{Cov}_{y \sim P(y; \theta)} [\nabla_{\theta'} \log P(y; \theta') |_{\theta' = \theta}] \\
 &= E_{y \sim P(y; \theta)} [\nabla_{\theta'} \log P(y; \theta') \nabla_{\theta'} \log P(y; \theta')^T |_{\theta' = \theta}] - \left(E_{y \sim P(y; \theta)} [\nabla_{\theta'} \log P(y; \theta') |_{\theta' = \theta}] \right)^T \\
 &= E_{y \sim P(y; \theta)} [\nabla_{\theta'} \log P(y; \theta') \nabla_{\theta'} \log P(y; \theta')^T |_{\theta' = \theta}]
 \end{aligned}$$

$$(c). \quad I(\theta) = E [\nabla_{\theta'} \log P(y) \nabla_{\theta'} \log P(y)^T | \theta' = \theta].$$

For i, j entry of the covariance matrix, $\left(\frac{1}{P(y; \theta')} \cdot \nabla_{\theta'_i} P(y; \theta') \right) \cdot \left(\frac{1}{P(y; \theta')} \cdot \nabla_{\theta'_j} P(y; \theta') \right) |_{\theta' = \theta}$

$$= \frac{1}{(P(y; \theta'))^2} \nabla_{\theta'_i} P(y; \theta') \nabla_{\theta'_j} P(y; \theta') |_{\theta' = \theta}.$$

for $E_{y \sim P(y; \theta)} [-\nabla_{\theta'}^2 \log P(y; \theta') |_{\theta' = \theta}]$ the i, j entry of the ~~result~~ result matrix is.

$$\begin{aligned}
 -\nabla_{\theta'_j} (\nabla_{\theta'_i} \log P(y; \theta') |_{\theta' = \theta}) &= -\nabla_{\theta'_j} \left(\frac{1}{P(y; \theta')} \nabla_{\theta'_i} P(y; \theta') |_{\theta' = \theta} \right) \\
 &= - \left((-1) \frac{1}{P(y; \theta')^2} \cdot \nabla_{\theta'_j} P(y; \theta') \cdot \nabla_{\theta'_i} P(y; \theta') + \frac{1}{P(y; \theta')} \nabla_{\theta'_i} \nabla_{\theta'_j} P(y; \theta') |_{\theta' = \theta} \right) \\
 &= \underbrace{\frac{1}{P(y; \theta')^2} \cdot \nabla_{\theta'_i} P(y; \theta') \nabla_{\theta'_j} P(y; \theta')}_A + \underbrace{\frac{1}{P(y; \theta')} \nabla_{\theta'_i} \nabla_{\theta'_j} P(y; \theta') |_{\theta' = \theta}}_B
 \end{aligned}$$

Prob 3. (c) continue.

$$\text{We know } E_{y \sim P(y; \theta)} [A \cdot B] = E_{y \sim P(y; \theta)} [A] \cdot E_{y \sim P(y; \theta)} [B].$$

$$\text{Also, } E_{y \sim P(y; \theta)} [B] = \int_{-\infty}^{\infty} P(y; \theta) \cdot \frac{1}{P(y; \theta)} \nabla_{\theta_i} \nabla_{\theta_j} P(y; \theta') dy.$$

$$= \int_{-\infty}^{\infty} \nabla_{\theta_i} \nabla_{\theta_j} P(y; \theta') dy$$

$$= \nabla_{\theta_i} \nabla_{\theta_j} \int_{-\infty}^{\infty} P(y; \theta') dy$$

$$= \nabla_{\theta_i} \nabla_{\theta_j} 1$$

$$= 0.$$

$$\begin{aligned} \text{Thus, the original } i, j \text{ entry} &= E_{y \sim P(y; \theta)} \left[\frac{1}{P(y; \theta)^2} \nabla_{\theta_i} P(y; \theta) \nabla_{\theta_j} P(y; \theta) \right] \\ &= E_{y \sim P(y; \theta)} \left[\nabla_{\theta_i} \log P(y; \theta') \nabla_{\theta_j} \log P(y; \theta')^T \right]_{\theta' = \theta} \\ &= I(\theta). \end{aligned}$$

$$(d). D_{KL}(P_{\theta} \| P_{\theta+d}) = f(\hat{\theta}) \approx D_{KL}(P_{\theta} \| P_{\theta}) + (\hat{\theta} - \theta)^T \nabla_{\theta} D_{KL}(P_{\theta} \| P_{\theta})|_{\theta' = \theta} + \frac{1}{2} (\hat{\theta} - \theta)^T \nabla_{\theta}^2 D_{KL}(P_{\theta} \| P_{\theta})|_{\theta' = \theta} (\hat{\theta} - \theta).$$

$$= 0 + d^T \nabla_{\theta} P(y; \theta) \log \frac{P(y; \theta)}{P(y; \theta')} |_{\theta' = \theta} + \frac{1}{2} d^T \nabla_{\theta}^2 P(y; \theta) \log \frac{P(y; \theta)}{P(y; \theta')} |_{\theta' = \theta} d$$

$$= d^T \int_{-\infty}^{\infty} \cancel{P(y; \theta)} \cdot \frac{\cancel{P(y; \theta')}}{\cancel{P(y; \theta)}} \cdot (-1) \cdot \frac{\cancel{P(y; \theta)}}{\cancel{P(y; \theta)}} \cdot \nabla_{\theta} P(y; \theta') |_{\theta' = \theta} + \frac{1}{2} d^T \nabla_{\theta} \left(\cancel{P(y; \theta)} \cdot \frac{\cancel{P(y; \theta')}}{\cancel{P(y; \theta)}} \cdot (-1) \cdot \frac{\cancel{P(y; \theta)}}{\cancel{P(y; \theta)}} \cdot \nabla_{\theta} P(y; \theta') \right) |_{\theta' = \theta} d.$$

$$= \underbrace{d^T (-1) \int_{-\infty}^{\infty} \cancel{P(y; \theta)} \cdot \nabla_{\theta} P(y; \theta') |_{\theta' = \theta} dy}_A + \underbrace{\frac{1}{2} d^T \int_{-\infty}^{\infty} \cancel{P(y; \theta)} \cdot \nabla_{\theta} \left(\cancel{P(y; \theta)} \cdot \frac{\cancel{P(y; \theta')}}{\cancel{P(y; \theta)}} \cdot (-1) \cdot \frac{\cancel{P(y; \theta)}}{\cancel{P(y; \theta)}} \cdot \nabla_{\theta} P(y; \theta') \right) |_{\theta' = \theta} dy}_B.$$

$$A = d^T (-1) \nabla_{\theta'} \int_{-\infty}^{\infty} P(y; \theta') dy |_{\theta' = \theta}$$

$$= d^T (-1) \nabla_{\theta'} (1) |_{\theta' = \theta}$$

$$= 0.$$

$$B = \frac{1}{2} d^T \int_{-\infty}^{\infty} \nabla_{\theta'} \left(P(y; \theta) \cdot \frac{P(y; \theta')}{P(y; \theta)} \cdot P(y; \theta) \cdot (-1) \cdot \frac{1}{P(y; \theta')} \cdot \nabla_{\theta'} P(y; \theta') \right)_{\theta'=\theta} dy$$

$$= \frac{1}{2} d^T (-1) \int_{-\infty}^{\infty} \nabla_{\theta'} P(y; \theta) \cdot \frac{1}{P(y; \theta')} \nabla_{\theta'} P(y; \theta')_{\theta'=\theta} dy.$$

$$= \frac{1}{2} d^T (-1) \int_{-\infty}^{\infty} P(y; \theta) \left((-1) \frac{1}{P(y; \theta)^2} \nabla_{\theta'} P(y; \theta')^* \cdot \nabla_{\theta'} P(y; \theta')^T + \frac{1}{P(y; \theta)} \nabla_{\theta'}^2 P(y; \theta') \right)_{\theta'=\theta} dy$$

$$= \frac{1}{2} d^T \left[\underbrace{\int_{-\infty}^{\infty} P(y; \theta) \cdot \frac{1}{P(y; \theta)^2} \nabla_{\theta'} P(y; \theta')^* \nabla_{\theta'} P(y; \theta')^T_{\theta'=\theta} dy}_C + \underbrace{\int_{-\infty}^{\infty} \nabla_{\theta'}^2 P(y; \theta)_{\theta'=\theta} dy}_D \right] d$$

$$D = \nabla_{\theta'}^2 \int_{-\infty}^{\infty} P(y; \theta') dy \big|_{\theta'=\theta} = \nabla_{\theta'}^2 (1) \big|_{\theta'=\theta} = 0$$

$$C = \int_{-\infty}^{\infty} P(y; \theta) \cdot \frac{\nabla_{\theta'} P(y; \theta')}{P(y; \theta)} \cdot \frac{\nabla_{\theta'} P(y; \theta')^T}{P(y; \theta)} \big|_{\theta'=\theta} dy$$

$$= E_{y \sim P(y; \theta)} \left[\nabla_{\theta'} \log P(y; \theta') \nabla_{\theta'} \log P(y; \theta')^T \big|_{\theta'=\theta} \right]$$

$$= I(\theta)$$

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

Q.E.D.

$$(e). \mathcal{L}(d, \lambda) = \ell(\theta + d) - \lambda [\text{D}_{KL}(P_\theta \| P_{\theta+d}) - C]$$

$$= \ell(\theta) + d^T \nabla_{\theta'} \ell(\theta')|_{\theta'=\theta} - \lambda \left(\frac{1}{2} d^T I(\theta) d - C \right)$$

$$= \log P(y, \theta) + d^T \frac{\nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}}{P(y; \theta)} - \lambda \left(\frac{1}{2} d^T I(\theta) d - C \right)$$

$$\nabla_d \mathcal{L}(d, \lambda) = \frac{\nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}}{P(y; \theta)} - \lambda I(\theta) d = 0$$

$$\lambda I(\theta) d = \frac{\nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}}{P(y; \theta)}$$

$$d = \frac{1}{\lambda} I^{-1}(\theta) \cdot \frac{\nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}}{P(y; \theta)}$$

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

$$d^T I(\theta) d = 2C$$

$$\left(\frac{1}{\lambda} I^{-1}(\theta) \cdot \frac{\nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}}{P(y; \theta)} \right)^T I(\theta) \left(\frac{1}{\lambda} I^{-1}(\theta) \cdot \frac{\nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}}{P(y; \theta)} \right) = 2C$$

$$\frac{1}{\lambda^2 P(y; \theta)^2} \cdot \nabla_{\theta'} P(y; \theta')^T|_{\theta'=\theta} I^{-1}(\theta) I(\theta) I^{-1}(\theta) \nabla_{\theta'} P(y; \theta')|_{\theta'=\theta} = 2C$$

$$\lambda^2 = \frac{1}{2C} \cdot \frac{1}{P(y; \theta)^2} \nabla_{\theta'} P(y; \theta')^T|_{\theta'=\theta} I^{-1}(\theta) \nabla_{\theta'} P(y; \theta')|_{\theta'=\theta}$$

$$\lambda = \frac{1}{P(y; \theta)} \sqrt{\frac{1}{2C} \nabla_{\theta'} P(y; \theta')^T I^{-1}(\theta) \nabla_{\theta'} P(y; \theta')}$$

$$\lambda = \frac{1}{\sqrt{2C}} \sqrt{\nabla_{\theta} \ell(\theta) I^{-1}(\theta) \nabla_{\theta} \ell(\theta)}$$

$$d^* = \sqrt{\frac{2C}{\nabla_{\theta} \ell(\theta) I^{-1}(\theta) \nabla_{\theta} \ell(\theta)}} I^{-1}(\theta) \nabla_{\theta} \ell(\theta) = \sqrt{\frac{2C}{\nabla_{\theta} P(y; \theta) I^{-1}(\theta) \nabla_{\theta} P(y; \theta)}} I^{-1}(\theta) \nabla_{\theta} P(y; \theta)$$

Q. (f). In Newton's method, the update is $-H^{-1} \nabla_{\theta} \ell(\theta)$. and we know from log likelihood is concave, then eventually $\ell(\theta)$ will converge to the maximum, which is H is P.S.D. Since natural parameter also has a maximum that ~~max~~ maximize the $\log P(y; \theta)$ and minimize D.K.L. ~~to~~ the ~~update~~ update is $\frac{1}{\lambda} I^{-1}(\theta) \nabla_{\theta} \ell(\theta) = \frac{1}{\lambda} (-I^{-1}(\theta)) \nabla_{\theta} \ell(\theta)$. we know λ is positive, which does not contribute the direction of ~~X~~.

For Generalized linear model, $\ell(\theta) = \arg \max_{\theta} P(y; \theta) = b(y) \exp(\theta^T T(y) - a(\theta))$

From problem set 1 question 4, we know $H_{ij} = \frac{\partial^2}{\partial \eta^2} a(\eta)$ where η is the natural parameter

In natural gradient, the parameter is θ . And,

$$I(\theta)_{ij} = \int_{-\infty}^{\infty} P(y; \theta) - \nabla_{\theta}^2 \log P(y; \theta) dy \quad \text{where let } a(\theta) = \log P(y; \theta)$$

$$= \int_{-\infty}^{\infty} P(y; \theta) - \nabla_{\theta}^2 \log (b(y) \exp(\theta^T T(y) - a(\theta))) dy.$$

$$= \int_{-\infty}^{\infty} P(y; \theta) - \nabla_{\theta}^2 (\log b(y) + \theta^T T(y) - a(\theta)) dy.$$

$$= \int_{-\infty}^{\infty} P(y; \theta) - \nabla_{\theta} (0 + T(y) - \nabla_{\theta} a(\theta)) dy.$$

$$= \int_{-\infty}^{\infty} P(y; \theta) + \nabla_{\theta}^2 a(\theta) dy.$$

$$= \nabla_{\theta}^2 a(\theta) \int_{-\infty}^{\infty} P(y; \theta) dy$$

$$= \nabla_{\theta}^2 a(\theta) \quad \text{which is} \quad \text{Thus, } I(\theta) = H$$

Thus, $I(\theta)^{-1} \nabla_{\theta} \log P(y; \theta) = H^{-1} \nabla_{\theta} \log P(y; \theta)$ which indicates the same direction.

Since $\frac{1}{\lambda}$ is a positive scalar that does not contribute to the direction.

(f) Answer:

Newton's Method: $\theta \leftarrow \theta := \theta - H^{-1} \nabla_{\theta} l(\theta).$

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

$$= E_{y \sim p(y; \theta)} [-\nabla_{\theta}^2 l(\theta)]$$

$$= -E_{y \sim p(y; \theta)} [\nabla_{\theta}^2 l(\theta)]$$

$$\theta := \theta + \hat{d}$$

$$= \theta + \frac{1}{\lambda} I(\theta)^{-1} \nabla_{\theta} l(\theta).$$

$$= \theta + \frac{1}{\lambda} E_{y \sim p(y; \theta)} [H]^{-1} \nabla_{\theta} l(\theta).$$

(4)-(c).

$$l_{\text{semi-sup}}(\theta^{(t+1)}) = l_{\text{unsup}}(\theta^{(t+1)}) + 2l_{\text{sup}}(\theta^{(t+1)}).$$

expanding

$$= \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} + 2 \sum_{i=1}^m \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)})$$

Jensen's inequality

$$\geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} + 2 \sum_{i=1}^m \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)})$$

$\theta^{(t+1)}$ is the arg max

$$\geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} + 2 \sum_{i=1}^m \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t)})$$

$$= l_{\text{unsup}}(\theta) + l_{\text{sup}}(\theta)$$

$$= l_{\text{semi-sup}}(\theta).$$

(b). latent variables: $W_j^{(i)}$ for $\forall i \in \{1, \dots, m\}$ and $\forall j \in \mathbb{R}^n$

$$W_j^{(i)} := Q_i^{(t)}(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma).$$

$$= \frac{P(z^{(i)}, x^{(i)}; \phi, \mu, \Sigma)}{P(x^{(i)})}$$

$$= \frac{P(x^{(i)} | z^{(i)}) \cdot P(z^{(i)})}{\sum_{z^{(i)}} P(x^{(i)} | z^{(i)}) \cdot P(z^{(i)})} \quad \text{with } \mu, \Sigma, \phi$$

~~(c) latent variables parameters: $\mu^{(t+1)}, \Sigma^{(t+1)}, \phi^{(t+1)}$~~

~~$$l_{\text{semi-sup}}(\theta) = \sum_{i=1}^m \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \mu, \Sigma, \phi)}{Q_i^{(t)}(z^{(i)})} + 2 \sum_{i=1}^m \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \mu, \Sigma, \phi)$$~~

~~$$W_j^{(i)} := \mathbb{1}\{z^{(i)} = j\} \text{ for } \forall i \in \{1, \dots, m\} \quad \forall j \in \mathbb{R}^n$$~~

(c). $\mathcal{L}_{\text{semi-sup}}(\theta) = \mathcal{L}_{\text{unsup}}(\theta) + \alpha \mathcal{L}_{\text{sup}}(\theta)$

~~SE~~

$$\theta^{(t+1)} := \arg \max_{\theta} \sum_{i=1}^m \left(\sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z^{(i)})} \right) + \alpha \left(\sum_{i=1}^{\tilde{m}} \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta) \right)$$

where $Q_i^{(t)}(z^{(i)}) = P(z^{(i)}=j | x^{(i)}; \phi, \mu, \Sigma)$ and

$$P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma) = P(x^{(i)} | z^{(i)}; \phi, \mu, \Sigma) \cdot P(z^{(i)}=j; \phi).$$

$$= \frac{1}{(2\pi)^{p/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j$$

$$\nabla_{\mu} \cdot \theta^{(t+1)} = \nabla_{\mu} \sum_{i=1}^m \sum_{z^{(i)}} -w_j^{(i)} \cdot \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)$$

$$= -\frac{1}{2} \sum_{i=1}^m w_j^{(i)} \left(x^{(i)T} \Sigma_j^{-1} x^{(i)} - x^{(i)T} \Sigma_j^{-1} \mu_j - \mu_j^T \Sigma_j^{-1} x^{(i)} + \mu_j^T \Sigma_j^{-1} \mu_j \right)$$

$$= -\frac{1}{2} \sum_{i=1}^m w_j^{(i)} \left(-x^{(i)T} \Sigma_j^{-1} - \Sigma_j^{-1} x^{(i)} + 2 \Sigma_j^{-1} \mu_j \right)$$

$$= -\frac{1}{2} \sum_{i=1}^m w_j^{(i)} \left(-2 \Sigma_j^{-1} x^{(i)} + 2 \Sigma_j^{-1} \mu_j \right)$$

$$= \sum_{i=1}^m w_j^{(i)} \left(\Sigma_j^{-1} x^{(i)} - \Sigma_j^{-1} \mu_j \right)$$

~~Handwritten scribbles and crossed-out equations.~~

$$= \Sigma_j^{-1} \sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)$$

Handwritten derivation showing the simplification of the gradient expression:

$$\mu_j = \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

c). continued.

$$\nabla_{\mu_k} \propto \sum_{i=1}^{\tilde{N}} \log(P(\tilde{x}^{(i)} | \hat{z}^{(i)}; \mu_k, \Sigma_k) P(\hat{z}^{(i)}; \phi))$$

$$= \propto \sum_{i=1}^{\tilde{N}} \nabla_{\mu_k} \log\left(\frac{1}{(2\pi)^{n/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(\tilde{x}^{(i)} - \mu_k)^T \Sigma_k^{-1} (\tilde{x}^{(i)} - \mu_k)\right)\right) \cdot \prod_{j=1}^K \phi_j^{1\{\hat{z}^{(i)}=j\}}$$

$$= \propto \sum_{i=1}^{\tilde{N}} \nabla_{\mu_k} \left(-\frac{1}{2}(\tilde{x}^{(i)} - \mu_k)^T \Sigma_k^{-1} (\tilde{x}^{(i)} - \mu_k)\right)$$

$$= \propto \sum_{i=1}^{\tilde{N}} (-1) \Sigma_k^{-1} (\tilde{x}^{(i)} - \mu_k) \cdot (-1)$$

$$= \propto \sum_{i=1}^{\tilde{N}} \Sigma_k^{-1} (\tilde{x}^{(i)} - \mu_k) = \underbrace{\propto \sum_{i=1}^{\tilde{N}} \Sigma_k^{-1} (\tilde{x}^{(i)} 1\{\hat{z}^{(i)}=1\} - \mu_k 1\{\hat{z}^{(i)}=1\})}$$

Combine the two terms. $\nabla_{\mu_k} \ell_{\text{semi-sup}}(0) = \sum_{i=1}^m \Sigma_k^{-1} w_k^{(i)} (\tilde{x}^{(i)} - \mu_k) + \propto \sum_{i=1}^{\tilde{N}} \Sigma_k^{-1} (\tilde{x}^{(i)} - \mu_k) = 0$
 $\tilde{x}^{(i)}$ belongs to μ_k

$$\Sigma_k^{-1} \left[\sum_{i=1}^m w_k^{(i)} (\tilde{x}^{(i)} - \mu_k) + \propto \sum_{i=1}^{\tilde{N}} (\tilde{x}^{(i)} - \mu_k) \right] = 0.$$

$$\sum_{i=1}^m w_k^{(i)} \tilde{x}^{(i)} - \mu_k \sum_{i=1}^m w_k^{(i)} + \propto \sum_{i=1}^{\tilde{N}} \tilde{x}^{(i)} - \propto \mu_k \cdot \text{count}(\tilde{x}^{(i)}) =$$

$$\mu_k \sum_{i=1}^m w_k^{(i)} + \propto \mu_k \cdot \text{count}(\tilde{x}) = \sum_{i=1}^m w_k^{(i)} \tilde{x}^{(i)} + \propto \sum_{i=1}^{\tilde{N}} \tilde{x}^{(i)}$$

$$\mu_k = \frac{\sum_{i=1}^m w_k^{(i)} \tilde{x}^{(i)} + \propto \sum_{i=1}^{\tilde{N}} \tilde{x}^{(i)} 1\{\hat{z}^{(i)}=1\}}{\sum_{i=1}^m w_k^{(i)} + \propto \sum_{i=1}^{\tilde{N}} 1\{\hat{z}^{(i)}=1\}}$$

$$\nabla_{\phi} \ell_{\text{unsup}}(\theta) = \nabla_{\phi} \sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log \phi_j$$

$$\nabla_{\phi} \ell_{\text{sup}}(\theta) = \nabla_{\phi} \sum_{i=1}^{\tilde{m}} \log \phi_{z_j^{(i)}} = \nabla_{\phi} \sum_{i=1}^{\tilde{m}} \sum_{j=1}^K \mathbb{1}_{\{z_j^{(i)} = j\}} \log \phi_j$$

$$\nabla_{\phi} \propto \sum_{i=1}^{\tilde{m}} \log \prod_{j=1}^K \phi^{z_j^{(i)}} = \nabla_{\phi} \propto \sum_{i=1}^{\tilde{m}} \sum_{j=1}^K \mathbb{1}_{\{z_j^{(i)} = j\}} \log \phi_j$$

~~Thus~~

$$\nabla_{\phi} \left(\sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log \phi_j + \alpha \sum_{i=1}^{\tilde{m}} \sum_{j=1}^K z_j^{(i)} \log \phi_j \right)$$

We know $\sum_{j=1}^K \phi_j = 1$ thus $\sum_{j=1}^K \phi_j - 1 = 0$.

$$\mathcal{L}(\phi) = \sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log \phi_j + \alpha \sum_{i=1}^{\tilde{m}} \sum_{j=1}^K z_j^{(i)} \log \phi_j + \beta \left(\sum_{j=1}^K \phi_j - 1 \right).$$

$$\frac{\partial}{\partial \phi_j} \mathcal{L}(\phi) = \sum_{i=1}^m \frac{w_j^{(i)}}{\phi_j} + \alpha \sum_{i=1}^{\tilde{m}} \frac{z_j^{(i)}}{\phi_j} + \beta = 0.$$

$$\Rightarrow \frac{1}{\phi_j} \sum_{i=1}^m w_j^{(i)} + \frac{1}{\phi_j} \alpha \sum_{i=1}^{\tilde{m}} z_j^{(i)} + \beta = 0.$$

$$\phi_j = \frac{\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} z_j^{(i)}}{-\beta}.$$

We know $\sum_j \phi_j = 1$

$$\sum_j \frac{\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} z_j^{(i)}}{-\beta} = 1$$

$$\sum_j \sum_{i=1}^m w_j^{(i)} + \alpha \sum_j \sum_{i=1}^{\tilde{m}} z_j^{(i)} = -\beta$$

$$m + \alpha \tilde{m} = -\beta$$

Thus
$$\phi_j = \frac{\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} z_j^{(i)}}{m + \alpha \tilde{m}}$$

$$\nabla_{\Sigma_\ell^{(t+1)}} \ell_{\text{unsup}}(\theta) = \sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log \frac{1}{(2\pi)^{k/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j.$$

$$= \nabla_{\Sigma_\ell} \sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log \frac{1}{(2\pi)^{k/2} |\Sigma_j|^{1/2}} + w_j^{(i)} \left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right)$$

$$= \sum_{i=1}^m \sum_{j=1}^K w_j^{(i)} \log |\Sigma_j|^{-1/2} - \frac{1}{2} w_j^{(i)} \left((x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right)$$

$$= \sum_{i=1}^m -\frac{1}{2} w_{\ell}^{(i)} \Sigma_\ell^{-1} + \sum_{i=1}^m -\frac{1}{2} w_{\ell}^{(i)} \left(-\Sigma^{-1} \right) (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T \Sigma^{-1}$$

$$= \sum_{i=1}^m -\frac{1}{2} w_{\ell}^{(i)} \Sigma_\ell^{-1} + \sum_{i=1}^m \frac{1}{2} w_{\ell}^{(i)} \left(\Sigma^{-1} (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T \right) \Sigma^{-1}$$

$$\nabla_{\Sigma_\ell} \ell_{\text{sup}}(\theta) = \cancel{\frac{\alpha \tilde{m}}{2} \Sigma^{-1}} - \frac{\alpha \tilde{m}}{2} \Sigma^{-1} + \sum_{i=1}^{\tilde{m}} \frac{1}{2} \Sigma^{-1} (\tilde{x}^{(i)} - \mu_\ell) (\tilde{x}^{(i)} - \mu_\ell)^T \Sigma^{-1}$$

$$\sum_{i=1}^m -\frac{1}{2} w_{\ell}^{(i)} + \sum_{i=1}^m \frac{1}{2} w_{\ell}^{(i)} \Sigma^{-1} (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T - \frac{\alpha \tilde{m}}{2} + \alpha \sum_{i=1}^{\tilde{m}} \frac{1}{2} \Sigma^{-1} (\tilde{x}^{(i)} - \mu_\ell) (\tilde{x}^{(i)} - \mu_\ell)^T =$$

$$\Sigma^{-1} \left(\sum_{i=1}^m \frac{1}{2} w_{\ell}^{(i)} (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T + \alpha \sum_{i=1}^{\tilde{m}} \frac{1}{2} (\tilde{x}^{(i)} - \mu_\ell) (\tilde{x}^{(i)} - \mu_\ell)^T \right) =$$

$$\sum_{i=1}^m \frac{1}{2} w_{\ell}^{(i)} + \frac{1}{2} \alpha \tilde{m}$$

$$\sum_{i=1}^m w_{\ell}^{(i)} (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T + \alpha \sum_{i=1}^{\tilde{m}} (\tilde{x}^{(i)} - \mu_\ell) (\tilde{x}^{(i)} - \mu_\ell)^T = (\Sigma) \left(\sum_{i=1}^m w_{\ell}^{(i)} + \alpha \tilde{m} \right) \mathbf{I}$$

$$\Sigma = \frac{\sum_{i=1}^m w_{\ell}^{(i)} (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T + \alpha \sum_{i=1}^{\tilde{m}} (\tilde{x}^{(i)} - \mu_\ell) (\tilde{x}^{(i)} - \mu_\ell)^T}{\left(\sum_{i=1}^m w_{\ell}^{(i)} + \alpha \tilde{m} \right) \mathbf{I}}$$

$$\Sigma_\ell = \left(\sum_{i=1}^m w_{\ell}^{(i)} (x^{(i)} - \mu_\ell) (x^{(i)} - \mu_\ell)^T + \alpha \sum_{i=1}^{\tilde{m}} (\tilde{x}^{(i)} - \mu_\ell) (\tilde{x}^{(i)} - \mu_\ell)^T \right) \left[\left(\sum_{i=1}^m w_{\ell}^{(i)} + \alpha \tilde{m} \right) \mathbf{I} \right]^{-1}$$

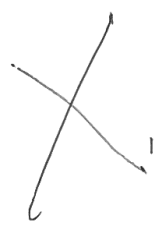
4). unsupervised semi-supervised

(i) 117, 119, 102 35, 34, 34.

(ii) not stable (different clusters) very stable.

(iii). not good quality. good quality
(not Gaussian shape) (Gaussian shape).

Problem 5

(b). $\frac{256 \times 256 \times 256}{16} \approx ? \text{ M}$ 

$3 \times 8 = 24$ bit. to store colors.

4 bits for 16 colors.

$$\frac{24}{4} = 6$$


4f). unsupervised semi-supervised

(i) 117, 119, 102 35, 34, 34.

(ii) not stable (different clusters) Very stable.

(iii). not good quality. good quality
(not Gaussian shape) (Gaussian shape).

Problem 5

(b). $\frac{256 \times 256 \times 256}{16} \approx ? \text{ M}$ 

$3 \times 8 = 24$ bit. to store colors.

4 bits for 16 colors.

$$\frac{24}{4} = 6$$

Problem 1 Rewrite

$$h_1 = \sigma(x_1 w_{1,1}^{[1]} + x_2 w_{2,1}^{[1]} + w_{0,1}^{[1]})$$

$$h_2 = \sigma(x_1 w_{1,2}^{[1]} + x_2 w_{2,2}^{[1]} + w_{0,2}^{[1]})$$

$$h_3 = \sigma(x_1 w_{1,3}^{[1]} + x_2 w_{2,3}^{[1]} + w_{0,3}^{[1]})$$

$$O = \sigma(h_1 w_1^{[2]} + h_2 w_2^{[2]} + h_3 w_3^{[2]} + w_0^{[2]})$$

$$\frac{\partial \mathcal{L}}{\partial w_{1,2}^{[1]}} = \frac{\partial \mathcal{L}}{\partial h_2} \frac{\partial h_2}{\partial w_{1,2}^{[1]}} = \frac{1}{m} \sum_{i=1}^m 2(O^{(i)} - y) \cdot \sigma^{(i)} \cdot (1 - O^{(i)}) \cdot w_2^{[2]} \cdot h_2 \cdot (1 - h_2) \cdot x_1$$