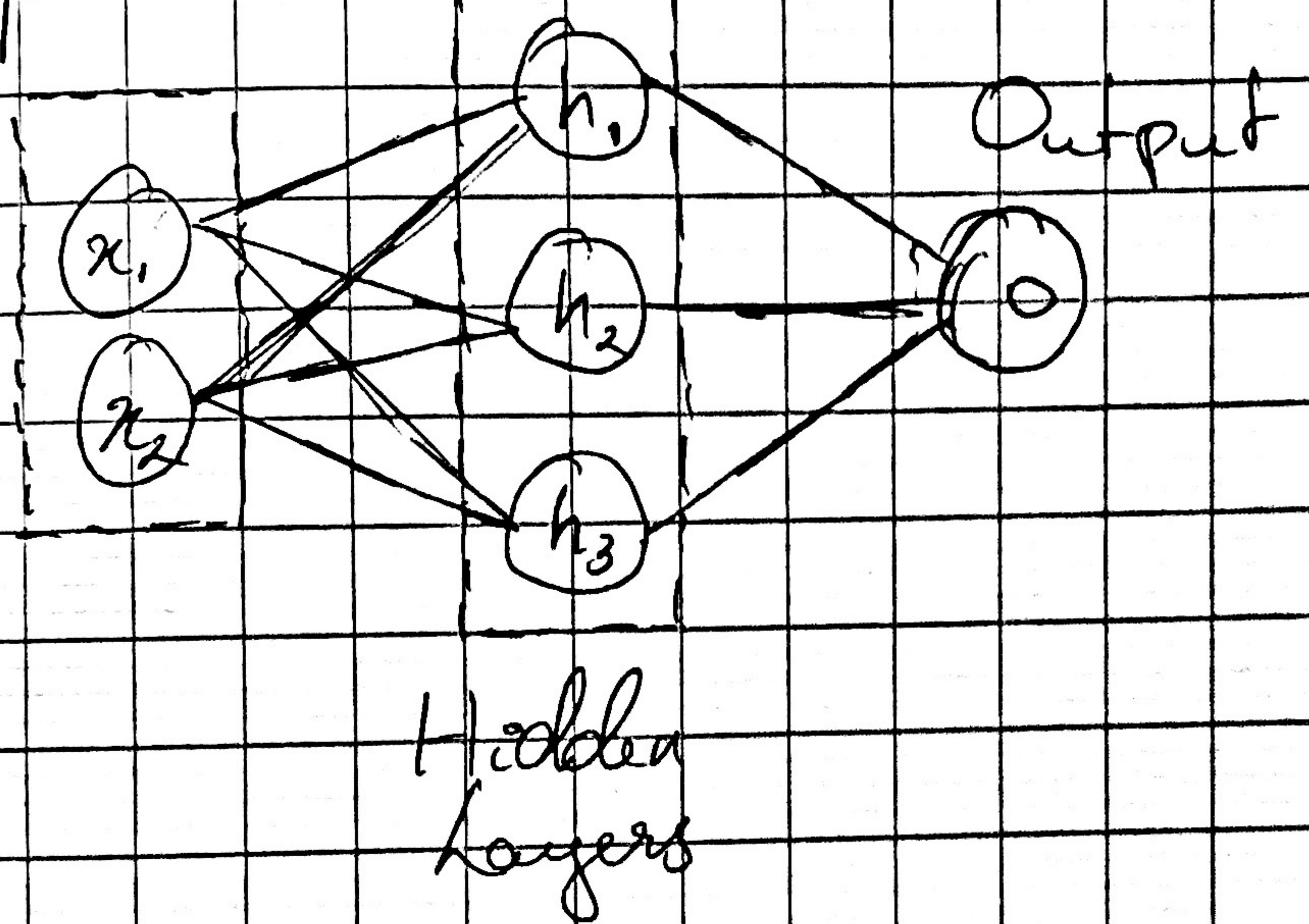


CS 229 Fall 2018

Problem Set #3: Deep learning &

Unsupervised learning

Inputs



Hidden
Layers

Output

$$h_j^{(i)} = g(w_{j,1}^{(i)} x_1 + w_{j,2}^{(i)} x_2 + \dots + w_{j,n}^{(i)})$$

$$o^{(i)} = g(w_{0,1}^{(i)} h_1 + w_{0,2}^{(i)} h_2 + w_{0,3}^{(i)} h_3 + w_0)$$

where g is the sigmoid function

$$\frac{\partial l}{\partial \omega_{1,2}^{(i)}} = \frac{\partial h_2}{\partial \omega_{1,2}^{(i)}} \cdot \frac{\partial o}{\partial h_2} \cdot \frac{\partial l}{\partial o}$$

$$= \sum_{i=0}^m x_1^{(i)} h_2^{(i)} (1 - h_2^{(i)}) \cdot w_2^{(i)} o^{(i)} (1 - o^{(i)}) \cdot \frac{2}{m} (o^{(i)} - y^{(i)})$$

$$\text{where } h_2^{(i)} = g(w_{1,2}^{(i)} x_1 + w_{2,2}^{(i)} x_2 + w_{0,2}^{(i)})$$

With the update rule of gradient descent

We get

$$\omega_{1,2}^{[1]} := \omega_{1,2}^{[1]} - \alpha \frac{\partial l}{\partial \omega_{1,2}^{[1]}}$$

where α is the learning rate.

- b) Each neuron in the hidden layer will check the position of the point $x^{(i)}$ with respect to each edge of the triangle.

$$\omega_{0,1}^{[1]} = 1 \quad \omega_{0,2}^{[1]} = -0,7 \quad \omega_{0,3}^{[1]} = -0,5$$

$$\omega_{1,1}^{[1]} = -1 \quad \omega_{1,2}^{[1]} = 1 \quad \omega_{1,3}^{[1]} = 0$$

$$\omega_{2,1}^{[1]} = -1 \quad \omega_{2,2}^{[1]} = 0 \quad \omega_{2,3}^{[1]} = 1$$

$$\omega_0^{[2]} = -3$$

With the output layer

$$\omega_1^{[2]} = 1$$

$$\omega_2^{[2]} = 1$$

$$\omega_3^{[2]} = 1$$

We can verify if all the

three conditions were met.

c) With the hidden layer we don't have any information gain.

2 KL divergence and Maximum likelihood

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

a) Prove that $\forall P, Q, D_{KL}(P \parallel Q) \geq 0$ and

$$D_{KL}(P \parallel Q) = 0 \text{ if and only if } P = Q$$

$$\Rightarrow D_{KL}(P \parallel Q) = \sum_{x \in X} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

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

Jensen's Inequality

$$\geq - \sum_{x \in X} \log(Q(x))$$

$$> -\log(1) \\ \geq 0$$

$$D_{KL}(P \parallel Q) = 0$$

According to the inequality of Jensen

$$-\sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right) = \sum_x \log\left(\frac{P(x)^2}{Q(x)}\right)$$

If and only if $\frac{P(x)}{Q(x)}$ is constant

and since $P(x) > 0$ $D_{KL}(P \parallel Q) = 0$ if

and only if $P(x) = Q(x)$.

$$b) D_{KL}(P(X) \parallel Q(X)) + D_{KL}(P(Y|X) \parallel Q(Y|X))$$

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

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

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

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

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

c) Prove that $\underset{\hat{P}}{\operatorname{argmin}} D_{KL}(\hat{P} || P_0) = \underset{\hat{P}}{\operatorname{argmax}} \sum_{i=1}^m \log P_0(x^{(i)})$

$$\underset{\hat{P}}{\operatorname{argmin}} D_{KL}(\hat{P} || P_0) = \sum_{x \in X} P(x) \log \left(\frac{\hat{P}(x)}{P_0(x)} \right)$$

$$= \underset{\hat{P}}{\operatorname{argmin}} \sum_{x \in X} \hat{P}(x) \log (\hat{P}(x)) - \sum_{x \in X} \hat{P}(x) \log (P_0(x))$$

$$= \underset{\hat{P}}{\operatorname{argmin}} - \sum_{x \in X} \hat{P}(x) \log (P_0(x))$$

$$= \underset{\hat{P}}{\operatorname{argmax}} \sum_{x \in X} \frac{1}{m} \sum_{i=1}^m \delta_{\{x^{(i)} = x\}} \log (P_0(x))$$

$$= \underset{\hat{P}}{\operatorname{argmax}} \sum_{i=1}^m \log P_0(x^{(i)})$$

3) KL Divergence, Fisher Information and the Natural Gradient

$$D_{KL}(p||q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

$$= E_{x \sim p(x)} [\log(p(x))] - E_{x \sim p(x)} [\log(q(x))]$$

a) The score function of $p(y; \theta)$ is $\nabla_{\theta} \log p(y; \theta)$

Show that the expected value of the score function is 0

$$E_{y \sim p(y; \theta)} [\nabla_{\theta} \log p(y; \theta) | \theta' = \theta] = 0$$

$$* E_{y \sim p(y; \theta)} [\nabla_{\theta} \log p(y; \theta) | \theta' = \theta] = \int_{-\infty}^{+\infty} p(y; \theta) \nabla_{\theta} \log p(y; \theta) dy$$

$$= \int p(y; \theta) \times \frac{1}{p(y; \theta)} \times \nabla_{\theta} p(y; \theta) dy$$

$$= \int \nabla_{\theta} p(y; \theta) dy = \nabla_{\theta} \int p(y; \theta) dy = \nabla_{\theta} (1) = 0$$

b) Fisher Information

$$I(\theta) = \text{Cov}_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta') | \theta' = \theta]$$

Show that the Fisher Information can be written as:

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

$$\times \text{Cov}_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta') | \theta' = \theta]$$

$$= E(\nabla_{\theta} \log p(y; \theta) - E(\nabla_{\theta} \log p(y; \theta)))$$

$$\nabla_{\theta} \log p(y; \theta) - E(\nabla_{\theta} \log p(y; \theta))$$

$$= E (\nabla_{\theta} \log(p(y; \theta)) \nabla_{\theta} \log(p(y; \theta)))^T$$

c) $E_{y \sim p(y; \theta)} [-\nabla_{\theta}^2 \log p(y; \theta)]_{ij}$

$$= E_{y \sim p(y; \theta)} \left[\frac{1}{p(y; \theta)} \frac{\partial p(y; \theta)}{\partial \theta_i \partial \theta_j} - \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \right]$$

$$= E_{y \sim p(y; \theta)} \left[\frac{1}{p(y; \theta)^2} \frac{\partial p(y; \theta)}{\partial \theta_i \partial \theta_j} \right]$$

d) Approximation of D_{KL} with Fisher Information

$$D_{KL}(p \parallel p_{\theta}) = E_{y \sim p(y; \theta)} [\log p(y; \theta)] - E_{y \sim p(y; \theta)} [\log(p(y; \bar{\theta}))]$$

Taylor Series of $\log(p(y; \theta))$

$$\Rightarrow \log(p(y; \bar{\theta})) = \log p(y; \theta) + d^T \nabla_{\theta} \log p(y; \theta) +$$

$$\frac{1}{2} d^T \nabla_{\theta}^2 \log p(y; \theta) d$$

$$= \log p(y; \theta) + d^T \frac{\nabla_{\theta} p(y; \theta)}{p(y; \theta)} + \frac{1}{2} d^T \nabla_{\theta}^2 \log p(y; \theta) d$$

$$\Rightarrow E_{y \sim p(y; \theta)} [\log p(y; \bar{\theta})] = -E_{y \sim p(y; \theta)} [\log p(y; \theta)] + E_{y \sim p(y; \theta)} [$$

$$-\frac{1}{2} d^T \nabla_{\theta}^2 \log p(y; \theta) d]$$

$$\Rightarrow D_{KL} = E_{y \sim p(y; \theta)} \left[-\frac{1}{2} d^T \nabla_{\theta}^2 \log p(y; \theta) d \right]$$

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

A Semi-supervised EM

$$* \text{Lunsup}(\theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)$$

$$* \text{Lsup}(\theta) = \sum_{i=1}^m \log (p(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta))$$

$$* E\text{-step}: Q_i^{(t)}(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta^{(t)})$$

$$* M\text{-step}: \theta^{(t+1)} := \underset{\theta}{\operatorname{argmax}} \left[\sum_{i=1}^m \left(\sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z^{(i)})} + \alpha \left(\sum_{i=1}^m \log p(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta) \right) \right) \right]$$

a) $\text{Lsemi-sup}(\theta^{t+1}) \geq \text{Lsemi-sup}(\theta^{(t)})$

$$\begin{aligned} \text{Lsemi-sup}(\theta^{t+1}) &= \sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta^{t+1}) + \alpha \sum_{i=1}^m \log (p(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta^{t+1})) \\ &= \sum_{i=1}^m \log \sum_{z^{(i)}} Q_i^{(t+1)}(z^{(i)}) \underbrace{p(x^{(i)}, z^{(i)}; \theta^{t+1})}_{Q_i^{(t+1)}(z^{(i)})} + \alpha \sum_{i=1}^m \log p(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta^{t+1}) \end{aligned}$$

$$\text{Jensen's inequality} \geq \sum_{i=1}^m \left(\sum_{z^{(i)}} Q_i^{(t+1)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{t+1})}{Q_i^{(t+1)}(z^{(i)})} \right) + \alpha \sum_{i=1}^m \log p(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta^{t+1})$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \left(\sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \right) + \alpha \sum_{i=1}^m \log p(\bar{x}^{(i)}, \bar{z}^{(i)}; \theta^{(t)})$$

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

$$= \text{lower-sup}(\theta^{(t)}) + \alpha \text{lower-sup}(\theta^{(t)})$$

$$= \text{lower-sup}(\theta^{(t)})$$

b) Semi-supervised E-Step:

$$Q_j^{(i)} = p(z_j^{(i)} | x_j^{(i)}; \phi, \mu, \Sigma)$$

$$= \frac{p(x_j^{(i)} | z_j^{(i)}; \mu, \Sigma) p(z_j^{(i)}; \phi)}{\sum_{l=1}^k p(x_j^{(i)} | z_l^{(i)}; \mu, \Sigma) p(z_l^{(i)}; \phi)}$$

$$\sum_{l=1}^k p(x_j^{(i)} | z_l^{(i)}; \mu, \Sigma) p(z_l^{(i)}; \phi)$$

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

$$\sum_{l=1}^k \frac{1}{(2\pi)^{n/2} |\Sigma_l|^{1/2}} \exp\left(-\frac{1}{2} (x_j^{(i)} - \mu_l)^T \Sigma_l^{-1} (x_j^{(i)} - \mu_l)\right) \phi_l$$

O Semi-Supervised learning

M-Step:

$$L = L_{\text{unsup}} + L_{\text{sup}}$$

$$= \sum_{i=1}^m \sum_{j=1}^k \omega_j^{(i)} \log(p(x^{(i)} | z_j; \theta) p(z_j)) +$$

$$\alpha \sum_{i=1}^m \sum_{j=1}^k \mathbb{I}\{\bar{z}^{(i)} = j\} \log(p(\bar{z}^{(i)} | \bar{z}_j) p(\bar{z}_j)) + \beta (1 - \sum_{j=1}^k \phi_j)$$

$$\frac{\partial L}{\partial \phi_j} = \sum$$

$$\Rightarrow \frac{\partial L}{\partial \phi_j} = \sum_{i=1}^m \frac{\omega_j^{(i)}}{\phi_j} + \alpha \cdot \frac{\sum_{i=1}^m \mathbb{I}\{\bar{z}^{(i)} = j\}}{\phi_j} - \beta = 0$$

$$\Rightarrow \phi_j = \sum_{i=1}^m \omega_j^{(i)} + \alpha \sum_{i=1}^m \mathbb{I}\{\bar{z}^{(i)} = j\}$$

B

$$\sum_{j=1}^k \phi_j = 1 = \sum_{i=1}^m \sum_{j=1}^k \omega_j^{(i)} + \alpha \sum_{i=1}^m \sum_{j=1}^k 1\{\bar{z}^{(i)} = j\}$$

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

$$*\mu_j$$

$$L_{sup} = \sum_{i=1}^m \sum_{j=1}^k \omega_j^{(i)} \log(p(z_j) p(x^{(i)}|z_j; \theta))$$

$$= \sum_{i=1}^m \sum_{j=1}^k \omega_j^{(i)} \log(\phi_j) + \log(p(x^{(i)}|z_j; \theta))$$

$$\frac{\partial L_{sup}}{\partial \mu_j} = \sum_{i=1}^m \omega_j^{(i)} \sum_j (\bar{x}^{(i)} - \mu_j)$$

$$= D \frac{\partial L_{sup}}{\partial \mu_j} = \sum_{i=1}^m 1\{\bar{z}^{(i)} = j\} \sum_j (\bar{x}^{(i)} - \mu_j)$$

$$D \frac{\partial L_{semi-sup}}{\partial \mu_j} = \frac{\partial L_{sup}}{\partial \mu_j} + \frac{\partial L_{unsup}}{\partial \mu_j}$$

$$= \sum_{i=1}^m \left[\sum_{j=1}^k \omega_j^{(i)} x_j^{(i)} + \alpha \sum_{j=1}^k 1\{\bar{z}^{(i)} = j\} \bar{x}_j^{(i)} - \mu_j \left(\sum_{i=1}^m \omega_i^{(i)} + \alpha \sum_{i=1}^m 1\{\bar{z}^{(i)} = j\} \right) \right]$$

$$\Rightarrow \mu_j = \sum_{i=1}^m w_j^{(i)} x^{(i)} + \alpha \sum_{i=1}^m 1 \{ \bar{z}^{(i)} = j \} \bar{x}^{(i)}$$

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

Likewise:

$$\frac{\partial L_{\text{sup}}}{\partial \sum_j} = -\frac{1}{2} \sum_{i=1}^m w_j^{(i)} \sum_j^{-1} + \frac{1}{2} \sum_j \left(\sum_{i=1}^m w_j^{(i)} (\bar{x} - \mu_j) (\bar{x} - \mu_j)^T \right) \sum_j^{-1}$$

$$\frac{\partial L_{\text{sup}}}{\partial \sum_j} = -\frac{1}{2} \sum_{i=1}^m 1 \{ \bar{z}^{(i)} = j \} \sum_j^{-1} + \frac{1}{2} \sum_j \left(\sum_{i=1}^m 1 \{ \bar{z}^{(i)} = j \} \right) \\ (\bar{x} - \mu_j) (\bar{x} - \mu_j)^T \sum_j^{-1}$$

$$\Rightarrow \sum_j = \sum_{i=1}^m w_j^{(i)} (\bar{x} - \mu_j) (\bar{x} - \mu_j)^T + \alpha \sum_{i=1}^m 1 \{ \bar{z}^{(i)} = j \} \\ (\bar{x} - \mu_j) (\bar{x} - \mu_j)^T$$

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