

1. A Simple Neural Network.

(a)

$$z^{(1)} = W^{(1)}x + b^{(1)}$$

$$h = \sigma(z^{(1)})$$

$$z^{(2)} = W^{(2)}h + b^{(2)}$$

$$o = \sigma(z^{(2)})$$

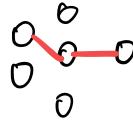
$$l = \frac{1}{m} \sum_{i=1}^m (o^{(i)} - y^{(i)})^2 = \frac{1}{m} \sum J^{(i)}$$

$$\frac{\partial J}{\partial w_{i,j}^{(1)}} = \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial h} \cdot \frac{\partial h}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial w_{i,j}^{(1)}}$$

$$= \sigma(o)y \cdot o(1-o) \cdot W_2^{(2)} \cdot h_2(l-h_2) \cdot x_i$$

Sigmoid derivative:

$$z' = z(1-z)$$



$$h_2 = W_{1,2}^{(2)} x_1 + W_{2,2}^{(2)} x_2 + b^{(2)}$$

so the updating rule of $w_{i,j}^{(1)}$ is:

$$w_{i,j}^{(1)} := w_{i,j}^{(1)} - \alpha \cdot \frac{1}{m} \sum_{i=1}^m \sigma(\sigma(o)y) \sigma(1-\sigma) \cdot W_2^{(2)} \cdot h_2(l-h_2) \cdot x_i^{(i)}$$

(b) Yes. Check the plot of the dataset, we can see the classifying boundary is a triangle, the three edges of triangle can be composed of three linear classifier.

Each neuron in the hidden layer can perform as a binary classifier. so it's possible to classify the dataset 100% right.

(c) It's not possible. If the activation function is linear,

$$z^{(2)} = W^{(2)}(W^{(1)}x + b^{(1)}) + b^{(2)}$$

$$= W^{(2)}W^{(1)}x + W^{(2)}b^{(1)} + b^{(2)}$$

so the result is still linear. However, the dataset is not linearly separable.

2. KL Divergence and Maximum Likelihood

(a) Non-negativity: $\log \frac{P(x)}{Q(x)}$ is a convex function,

$$\begin{aligned} D_{KL}(P||Q) &= \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \\ &= - \sum_{x \in X} P(x) \log \frac{Q(x)}{P(x)} \\ &= E \left[-\log \frac{Q(x)}{P(x)} \right] \end{aligned}$$

For a convex function f , according to Jensen's Inequality:

$$\begin{aligned} E \left[-\log \frac{Q(x)}{P(x)} \right] &\geq -\log \left(E \left[\frac{Q(x)}{P(x)} \right] \right) \\ &= -\log \left(\sum_{x \in X} P(x) \frac{Q(x)}{P(x)} \right) \\ &= -\log \left(\sum_{x \in X} Q(x) \right) \\ &= 0. \end{aligned}$$

so $D_{KL}(P||Q) \geq 0$.

If $P = Q$, then $D_{KL}(P||Q) = \sum_{x \in X} P(x) \log 1 = 0$.

If $D_{KL}(P||Q) = 0$, then $\sum_{x \in X} P(x) \frac{P(x)}{\log Q(x)} = 0$.

In the problem setting, $P(x) > 0$ for $\forall x$. $\log \frac{P(x)}{Q(x)} \geq 0$. $P = Q$.

So, $D_{KL}(P||Q) = 0 \iff P = Q$.

(b) Chain Rule for KL Divergence.

$$\begin{aligned} D_{KL}(P(X, Y) || Q(X, Y)) &= \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \frac{P(x, y)}{Q(x, y)} \\ &= \sum_{x \in X} \sum_{y \in Y} P(x) \cdot P(y|x) \cdot \log \frac{P(x)P(y|x)}{Q(x)Q(y|x)} \\ &= \sum_{x \in X} \sum_{y \in Y} P(x)P(y|x) \left(\log \frac{P(x)}{Q(x)} + \log \frac{P(y|x)}{Q(y|x)} \right) \\ &= \sum_{x \in X} \sum_{y \in Y} P(x)P(y|x) \log \frac{P(x)}{Q(x)} + \sum_{x \in X} \sum_{y \in Y} P(x)P(y|x) \frac{P(y|x)}{Q(y|x)} \\ &= \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \left(\sum_{y \in Y} P(y|x) \right) + \sum_{x \in X} P(x) \left(\sum_{y \in Y} P(y|x) \frac{P(y|x)}{Q(y|x)} \right) \\ &= D_{KL}(P(X) || Q(X)) + D_{KL}(P(Y|X) || Q(Y|X)) \end{aligned}$$

(c) KL and Maximum likelihood.

$$\begin{aligned} \text{proof: } \arg\min_{\theta} D_{KL}(\hat{P} || P_{\theta}) &= \arg\min_{\theta} \sum_{x \in X} \hat{P}(x) \log \frac{\hat{P}(x)}{P_{\theta}(x)} \\ &= \arg\min_{\theta} \sum_{x \in X} \hat{P}(x) \log \hat{P}(x) - \hat{P}(x) \log P_{\theta}(x) \\ &= \arg\min_{\theta} \sum_{x \in X} -\hat{P}(x) \log P_{\theta}(x) \\ &= \arg\max_{\theta} \sum_{x \in X} \hat{P}(x) \log P_{\theta}(x) \\ &= \arg\max_{\theta} \sum_{x \in X} \left(\frac{1}{m} \sum_{t=1}^m \delta(x^{(t)} = x) \right) \log P_{\theta}(x) \\ &= \arg\max_{\theta} \sum_{t=1}^m \log P_{\theta}(x^{(t)}) \end{aligned}$$

3. KL Divergence, Fisher Information, and the Natural Parameter

KL divergence is invariant to model parameterization⁽ⁱ⁾, but the gradient w.r.t. the model parameter is not invariant to model parameterization.

proof: (i) $y = mx + b$, $p(x) dx = p(y) dy$.

$$\begin{aligned} D_{KL}(P(x) || Q(x)) &= \int_a^b p(x) \log \frac{P(x)}{Q(x)} dx \\ &= \int_a^b p(y) \log \frac{P(y)}{Q(y)} \frac{dy}{dx} dy \\ &= \int_a^b p(y) \log \left(\frac{P(y)}{Q(y)} \right) dy = D_{KL}(P(y) || Q(y)) \end{aligned}$$

(ii) e.g. $D_{KL}(\hat{P}(x) || P_{\theta}(x))$ is will change when updating parameter θ .

(a) Score function:

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

(b) Fisher Information.

X is real random vector,

$$\begin{aligned}
 \text{cov}(x, x) &= E[(x - E[x])(x - E[x])^T] \\
 &= E[xx^T] - E[x]E[x]^T
 \end{aligned}$$

$$\begin{aligned}
 I(\theta) &= \text{cov}_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta') |_{\theta'=\theta}] \\
 &= E_{y \sim p(y; \theta)} [\nabla_{\theta} \log p(y; \theta) (\nabla_{\theta} \log p(y; \theta))^T |_{\theta'=\theta}] \\
 &\quad - E_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta') |_{\theta'=\theta}] E_{y \sim p(y; \theta)} [\nabla_{\theta} \log p(y; \theta') |_{\theta'=\theta}]^T \\
 &= E_{y \sim p(y; \theta)} [\nabla_{\theta} \log p(y; \theta) (\nabla_{\theta} \log p(y; \theta))^T |_{\theta'=\theta}]
 \end{aligned}$$

(c) Fisher Information (alternative form)

$$\frac{\partial \log p(y; \theta)}{\partial \theta_i} = \frac{1}{p(y; \theta)} \frac{\partial p(y; \theta)}{\partial \theta_i}$$

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

$$= E_{y \sim p(y; \theta)} \left[\frac{\partial \log p(y; \theta')}{\partial \theta_i} \frac{\partial \log p(y; \theta')^T}{\partial \theta_j} \right]$$

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

$$\frac{\partial^2 \log p(y; \theta)}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_j} \left(\frac{1}{p(y; \theta)} \cdot \frac{\partial p(y; \theta)}{\partial \theta_i} \right)$$

$$= -\frac{1}{p(y; \theta)^2} \cdot \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} + \frac{1}{p(y; \theta)} \cdot \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j}$$

$$\begin{aligned}
 E_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta') |_{\theta'=\theta}]_{ij} &= E_{y \sim p(y)} \left[\frac{1}{p(y; \theta)^2} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} - \frac{1}{p(y; \theta)} \cdot \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \right]_{ij} \\
 &= E_{y \sim p(y)} \left[\frac{1}{p(y; \theta)^2} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \right]_{ij} - E_{y \sim p(y)} \left[\frac{1}{p(y; \theta)} \cdot \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \right]_{ij} \\
 &= E_{y \sim p(y)} \left[\frac{1}{p(y; \theta)^2} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \right]_{ij} - \int_{-\infty}^{\infty} p(y; \theta) \frac{1}{p(y; \theta)} \cdot \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} dy
 \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{y \sim p(y)} \left[\frac{1}{p(y_i; \theta)^2} \frac{\partial^2 p(y_i; \theta)}{\partial \theta_i \partial \theta_j} \right]_{ij} - \frac{\delta^2}{2} \int_{-\infty}^{\infty} p(y_i; \theta) dy \\
&= \mathbb{E}_{y \sim p(y)} \left[\frac{1}{p(y_i; \theta)^2} \frac{\partial^2 p(y_i; \theta)}{\partial \theta_i \partial \theta_j} \right]_{ij} = I(\theta)_{ij}
\end{aligned}$$

(d) Approximating D_{KL} with Fisher Information:

$$\begin{aligned}
\log p(y; \tilde{\theta}) &\approx \log p(y; \theta) + (\tilde{\theta} - \theta)^T \nabla_{\theta} \log p(y; \theta')|_{\theta'=\theta} + \frac{1}{2} (\tilde{\theta} - \theta)^T (\nabla_{\theta}^2 \log p(y; \theta'))|_{\theta=\theta} \\
\tilde{\theta} = \theta + d &= \log p(y; \theta) + d^T \nabla_{\theta} \log p(y; \theta')|_{\theta=\theta} + \frac{1}{2} d^T (\nabla_{\theta}^2 \log p(y; \theta'))|_{\theta=\theta} d.
\end{aligned}$$

$$\mathbb{E}_{y \sim p(y; \theta)} [\log p(y; \tilde{\theta})] = \mathbb{E}_{y \sim p(y; \theta)} [\log p(y; \theta)] + d + \frac{1}{2} d^T \mathbb{E}_{y \sim p(y; \theta)} [\nabla_{\theta}^2 \log p(y; \theta)] d$$

$$\begin{aligned}
D_{KL}(P(\theta) || P(\theta+d)) &= \int_{-\infty}^{\infty} p(y; \theta) \log \frac{p(y; \theta)}{p(y; \tilde{\theta})} \\
&= \mathbb{E}_{y \sim p(y; \theta)} [\log p(y; \theta)] - \mathbb{E}_{y \sim p(y; \theta)} [\log p(y; \tilde{\theta})] \\
&= \frac{1}{2} d^T I(\theta) d
\end{aligned}$$

(e) Natural Gradient.

$$d^* = \underset{d}{\operatorname{argmax}} l(\theta+d) \text{ subject to } D_{KL}(p||p_{\theta+d}) = c.$$

construct Lagrangian

$$\begin{aligned}
L(d, \lambda) &= l(\theta+d) - \lambda [D_{KL}(p||p_{\theta+d}) - c] \\
&\approx \log p(y; \theta) + d^T \nabla_{\theta} \log p(y; \theta)|_{\theta=\theta} - \lambda \left[\frac{1}{2} d^T I(\theta) d - c \right]
\end{aligned}$$

$$\begin{aligned}
\nabla_d L(d, \lambda) &\approx \nabla_{\theta} \log p(y; \theta)|_{\theta=\theta} - \lambda I(\theta) d \\
&= \frac{\nabla_{\theta} p(y; \theta)|_{\theta=\theta}}{p(y; \theta)} - \lambda I(\theta) d = 0. \quad (a)
\end{aligned}$$

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

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

$$\begin{aligned}
\text{plug } \tilde{d} \text{ into (b)} &= -\frac{1}{2} \frac{1}{\lambda} \frac{\nabla_{\theta} p(y; \theta)|_{\theta=\theta}^T (I(\theta)^{-1})^T I(\theta) \cdot \frac{1}{\lambda} I(\theta)^{-1} \nabla_{\theta} p(y; \theta)|_{\theta=\theta}}{p(y; \theta)} + c \\
&= -\frac{1}{2 \lambda^2 p(y; \theta)^2} \nabla_{\theta} p(y; \theta)|_{\theta=\theta}^T (I(\theta)^{-1})^T \nabla_{\theta} p(y; \theta)|_{\theta=\theta} + c \\
&= 0
\end{aligned}$$

$$\lambda = \sqrt{\frac{1}{2 c p(y; \theta)^2} \cdot \nabla_{\theta} p(y; \theta)|_{\theta=\theta}^T (I(\theta)^{-1})^T \nabla_{\theta} p(y; \theta)|_{\theta=\theta}}$$

plug λ back to \mathcal{J} .

$$\begin{aligned} d^* &= \sqrt{\frac{2c p(y; \theta)^2}{\nabla_{\theta'} p(y; \theta)|_{\theta=0}^T (\mathcal{I}(\theta))^T \nabla_{\theta'} p(y; \theta)|_{\theta=0}}} \mathcal{I}(\theta)^T \frac{\nabla_{\theta'} p(y; \theta)|_{\theta=0}}{p(y; \theta)} \\ &= \sqrt{\frac{2c}{\nabla_{\theta'} p(y; \theta)|_{\theta=0}^T (\mathcal{I}(\theta))^T \nabla_{\theta'} p(y; \theta)|_{\theta=0}}} \mathcal{I}(\theta)^T \nabla_{\theta'} p(y; \theta)|_{\theta=0} \end{aligned}$$

(f) Relation with Newton's Method.

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

$$\begin{aligned} \text{Natural gradient: } \mathcal{I}(\theta) &= \mathbb{E}_{y \sim p(y; \theta)} [\nabla_{\theta'}^2 \log p(y; \theta')|_{\theta'=\theta}] \\ &= \mathbb{E}_{y \sim p(y; \theta)} [H] \\ &= -\mathbb{E}_{y \sim p(y; \theta)} [H] \\ \theta &:= \theta + \mathcal{J} \\ &= \theta + \frac{1}{\lambda} \mathcal{I}(\theta)^T \nabla_{\theta} l(\theta) \\ &= \theta - \frac{1}{\lambda} \mathbb{E}_{y \sim p(y; \theta)} [H]^{-1} \nabla_{\theta} l(\theta) \end{aligned}$$

They're the same.

4. Semi-supervised EM.

$$\begin{aligned} l_{\text{semi-sup}}(\theta) &= l_{\text{unsup}}(\theta) + \alpha l_{\text{sup}}(\theta) \\ &= \sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) + \alpha \sum_{i=1}^N \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta) \end{aligned}$$

$$\begin{aligned} \sum_v \log p(x^{(i)}; \theta) &= \sum_v \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) \\ &= \sum_v \log \sum_{z^{(i)}} Q_v(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_v(z^{(i)})} \\ \text{Jensen's Inequality} &\geq \sum_v \sum_{z^{(i)}} Q_v(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_v(z^{(i)})} \end{aligned}$$

E-step:

$$\text{for each } i \in \{1, \dots, m\}, \text{ set} \\ Q_v^{(t)}(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta^{(t)})$$

M-step:

$$\theta^{(t+1)} := \operatorname{argmax}_{\theta} \left[\sum_{i=1}^m \left(\sum_{z^{(i)}} Q_v^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_v^{(t)}(z^{(i)})} \right) + \alpha \left(\sum_{i=1}^N \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta) \right) \right]$$

(a) Convergence.

$$\begin{aligned}
 l_{\text{semi-sup}}(\theta^{(t+1)}) &= l_{\text{unsup}}(\theta^{(t+1)}) + \alpha l_{\text{sup}}(\theta^{(t+1)}) \\
 &\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i(z^{(i)})} + \alpha \sum_{i=1}^m \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)}) \\
 &\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^t)}{Q_i(z^{(i)})} + \alpha \sum_{i=1}^m \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^t) \\
 &= l_{\text{unsup}}(\theta^{(t)}) + \alpha l_{\text{sup}}(\theta^{(t)}) \\
 &= l_{\text{semi-sup}}(\theta^{(t)})
 \end{aligned}$$

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

(b) In E step, we need to reestimate $w_j^{(i)}$

For each i, j set

$$\begin{aligned}
 w_j^{(i)} &:= p(z^{(i)}=j | x^{(i)}; \phi, \mu, \Sigma) \\
 &= \frac{p(x^{(i)} | z^{(i)}=j; \mu, \Sigma) p(z^{(i)}=j | \phi)}{\sum_{l=1}^k p(x^{(i)} | z^{(i)}=l; \mu, \Sigma) p(z^{(i)}=l | \phi)} \\
 &= \frac{\frac{1}{(2\pi)^{k/2} |\Sigma_j|^{1/2}} \exp(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)) \phi_j}{\sum_{l=1}^k \frac{1}{(2\pi)^{k/2} |\Sigma_l|^{1/2}} \exp(-\frac{1}{2} (x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l)) \phi_l}.
 \end{aligned}$$

(c) parameter ϕ, μ, Σ need to be updated.

In order to maximize log likelihood,

$$\begin{aligned}
 l_{\text{semi-sup}} &= l_{\text{unsup}} + \alpha l_{\text{sup}} \\
 &= \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} + \alpha \sum_{i=1}^m \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \phi, \mu, \Sigma) \\
 &= \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} + \alpha \sum_{i=1}^m \log p(\tilde{x}^{(i)} | z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi) \\
 &= \sum_{i=1}^m \sum_{j=1}^{|\Sigma|} w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{k/2} |\Sigma_j|^{1/2}} \exp(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)) \cdot \phi_j}{w_j^{(i)}} + \\
 &\quad \alpha \sum_{i=1}^m \sum_{j=1}^{|\Sigma|} \mathbb{I}_{\{z^{(i)}=j\}} \log \frac{1}{(2\pi)^{k/2} |\Sigma_j|^{1/2}} \exp(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)) \cdot \phi_j
 \end{aligned}$$

Similar to how we solve in unsupervised case, we construct the Lagrangian

$$\phi: \quad \mathcal{L}(\phi) = \sum_{i=1}^m \sum_{k=1}^K w_k^{(i)} \log \phi_k + \beta \left(\sum_{k=1}^K \phi_k - 1 \right) + \alpha \sum_{i=1}^m \sum_{k=1}^K 1\{\tilde{z}^{(i)}=k\} = l^T \log \phi$$

$$\nabla_{\phi} \mathcal{L}(\phi) = \sum_{i=1}^m \frac{w_i^{(i)}}{\phi_j} + \alpha \sum_{i=1}^m \frac{1\{\tilde{z}^{(i)}=j\}}{\phi_j} + \beta = 0, \quad \phi_j = \frac{\sum_{i=1}^m w_i^{(i)} + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\}}{\beta}$$

$$\sum_{k=1}^K \phi_k = \frac{\sum_{i=1}^m \sum_{k=1}^K w_k^{(i)} + \alpha \sum_{i=1}^m \sum_{k=1}^K 1\{\tilde{z}^{(i)}=k\}}{-\beta} = \frac{m + \alpha m}{-\beta} = 1. \quad \therefore \beta = -(m + \alpha m)$$

$$\phi_j = \frac{\sum_{i=1}^m w_i^{(i)} + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\}}{m + \alpha m} \quad (1)$$

$$\mu_j: \quad \nabla_{\mu_j} \text{loss semi-sup} = \nabla_{\mu_j} \text{loss sup} + \alpha \nabla_{\mu_j} l^{\text{sup}}$$

$$= \sum_{i=1}^m \left[\left(\sum_{j=1}^m w_j^{(i)} (x^{(i)} - \mu_j) \right) + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\} (\tilde{x}^{(i)} - \mu_j) \right]$$

$$= \sum_{i=1}^m \left[\left(\sum_{j=1}^m w_j^{(i)} x^{(i)} + \alpha \sum_{j=1}^m 1\{\tilde{z}^{(i)}=j\} \tilde{x}^{(i)} \right) - \left(\sum_{j=1}^m w_j^{(i)} + \alpha \sum_{j=1}^m 1\{\tilde{z}^{(i)}=j\} \right) \mu_j \right]$$

$$\mu_j = \frac{\left(\sum_{j=1}^m w_j^{(i)} x^{(i)} + \alpha \sum_{j=1}^m 1\{\tilde{z}^{(i)}=j\} \tilde{x}^{(i)} \right)}{\sum_{j=1}^m w_j^{(i)} + \alpha \sum_{j=1}^m 1\{\tilde{z}^{(i)}=j\}} \quad (2)$$

$$\Sigma_j^{-1}: \quad \nabla_{\Sigma_j} \text{loss semi-sup} = \nabla_{\Sigma_j} \text{loss sup} + \alpha \nabla_{\Sigma_j} l^{\text{sup}}$$

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

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

$$= -\frac{1}{2} \sum_{j=1}^m \left(\sum_{i=1}^m w_i^{(i)} + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\} \right)$$

$$+ \frac{1}{2} \sum_{j=1}^m \left(\sum_{i=1}^m w_i^{(i)} (x^{(i)} - \mu_j)^T + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\} (\tilde{x}^{(i)} - \mu_j)^T \right) \Sigma_j^{-1}$$

$$= 0$$

$$\Sigma_j^{-1} = \frac{\sum_{i=1}^m w_i^{(i)} (x^{(i)} - \mu_j)^T + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\} (\tilde{x}^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_i^{(i)} + \alpha \sum_{i=1}^m 1\{\tilde{z}^{(i)}=j\}}$$

(f) (i) Semi-supervised EM algorithm requires less iteration than unsupervised EM.

(ii) Semi-supervised EM is more stable. The labeled data can efficiently guide unlabeled data to converge to the same color after different initialization.

(iii) The overall quality of semi-sup EM is better than unsup EM.
EM sometimes failed to distinguish two close low-variance gaussian, but mix them into one gaussian instead.