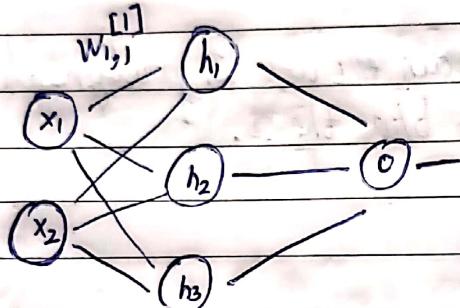


PS-3

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

$$(a) \text{ let } l^{(i)} = (o^{(i)} - y^{(i)})^2$$



$$m^{(i)} = \begin{pmatrix} x_1^{(i)} \\ x_2^{(i)} \end{pmatrix}_{2 \times 1}$$

$$\text{let } w^{[1]} = \begin{pmatrix} w_{1,1}^{[1]} & w_{2,1}^{[1]} \\ w_{1,2}^{[1]} & w_{2,2}^{[1]} \\ w_{1,3}^{[1]} & w_{2,3}^{[1]} \end{pmatrix}_{3 \times 2}$$

$$\text{and } b^{[1]} = \begin{pmatrix} w_0^{[1]} \\ w_0^{[1]} \\ w_0^{[1]} \end{pmatrix}_{3 \times 1}$$

$$\Rightarrow z^{[1](i)} = w^{[1]} m^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$\text{let } w^{[2]} = (w_1^{[2]} \quad w_2^{[2]} \quad w_3^{[2]})_{1 \times 3}$$

$$\text{and } b^{[2]} = (w_0^{[2]})$$

$$\Rightarrow O^{(i)} = a^{[2]^{(i)}} = \sigma(z^{[2]^{(i)}})$$

$$z^{[2]^{(i)}} = W^{[2]} a^{[1]^{(i)}} + b^{[2]}$$

$$\frac{\partial L}{\partial O^{(i)}} = 2(O^{(i)} - y^{(i)})$$

$$\frac{\partial L}{\partial a^{[2]^{(i)}}} = 2(a^{[2]^{(i)}} - y^{(i)})$$

$$\frac{\partial L}{\partial z^{[2]^{(i)}}} = \frac{\partial L}{\partial a^{[2]^{(i)}}} \times \frac{\partial a^{[2]^{(i)}}}{\partial z^{[2]^{(i)}}}$$

$$\frac{\partial L}{\partial z^{[2]^{(i)}}} = 2(a^{[2]^{(i)}} - y^{(i)}) * a^{[2]^{(i)}} (1 - a^{[2]^{(i)}})$$

$$\begin{aligned} \frac{\partial L}{\partial a^{[1]^{(i)}}} &= \frac{\partial L}{\partial z^{[2]^{(i)}}} \times \frac{\partial z^{[2]^{(i)}}}{\partial a^{[1]^{(i)}}} \\ &= \underbrace{\frac{\partial L}{\partial z^{[2]^{(i)}}}}_{1 \times 1} \times \underbrace{(W^{[2]})^T}_{3 \times 1} = (W^{[2]})^T \times \frac{\partial L}{\partial z^{[2]^{(i)}}} \end{aligned}$$

$$\begin{aligned} \left(\frac{\partial L}{\partial z^{[2]^{(i)}}} \right) \frac{\partial L}{\partial z^{[1]^{(i)}}} &= \frac{\partial L}{\partial a^{[1]^{(i)}}} \times \frac{\partial a^{[1]^{(i)}}}{\partial z^{[1]^{(i)}}} \\ &= \underbrace{\frac{\partial L}{\partial a^{[1]^{(i)}}}}_{3 \times 1} * \underbrace{a^{[1]^{(i)}} (1 - a^{[1]^{(i)}})}_{3 \times 1} \end{aligned}$$

↙ this is okay!

$$\begin{aligned} \frac{\partial L}{\partial W^{[1]}} &= \frac{\partial L}{\partial z^{[1]^{(i)}}} \times \frac{\partial z^{[1]^{(i)}}}{\partial W^{[1]}} \\ &= \underbrace{\frac{\partial L}{\partial z^{[1]^{(i)}}}}_{3 \times 1} \times \underbrace{(M^{(i)})^T}_{1 \times 2} \end{aligned}$$

↙ done!

$$\therefore \frac{\partial l^{(i)}}{\partial w^{[2]}} = (w^{[2]})^T \times \left(2(a^{[2](i)} - y^{(i)}) * a^{[2](i)} (1 - a^{[2](i)}) \right)$$

$\underbrace{a^{[1](i)} (1 - a^{[1](i)})}_{3 \times 1} \quad \underbrace{x (M^{(i)})^T}_{1 \times 2}$

$$= \begin{pmatrix} w_1^{[2]} \\ w_2^{[2]} \\ w_3^{[2]} \end{pmatrix} * \begin{pmatrix} w_{0,1}^{[1]} + w_{1,1}^{[1]} M_1^{(i)} + w_{2,1}^{[1]} M_2^{(i)} \\ w_{0,2}^{[1]} + w_{1,2}^{[1]} M_1^{(i)} + w_{2,2}^{[1]} M_2^{(i)} \\ w_{0,3}^{[1]} + w_{1,3}^{[1]} M_1^{(i)} + w_{2,3}^{[1]} M_2^{(i)} \end{pmatrix}$$

$$\times (1 \times 1) \times (M_1^{(i)} \quad M_2^{(i)})$$

finally we need $\frac{\partial l^{(i)}}{\partial w_{1,2}^{[2]}}$

again rearranging,

~~$$\frac{\partial l^{(i)}}{\partial w^{[2]}} = (w^{[2]})^T \times 2(a^{[2](i)} - y^{(i)}) * (a^{[1](i)} (1 - a^{[1](i)}))$$~~

$$\frac{\partial l^{(i)}}{\partial w^{[2]}} = (w^{[2], T}) \times \left(2(o^{(i)} - y^{(i)}) * o^{(i)} (1 - o^{(i)}) \right)$$

$\hookrightarrow 1^{\text{st}} \text{ col } m \text{ of this}$

$\hookrightarrow 2^{\text{nd}} \text{ row of this}$

$$\Rightarrow \frac{\partial l^{(i)}}{\partial w_{1,2}^{[2]}} = 2 w_2^{[2]} (o^{(i)} - y^{(i)}) o^{(i)} (1 - o^{(i)})$$

$\times (w_{0,2}^{[1]} + w_{1,2}^{[1]} M_1^{(i)} + w_{2,2}^{[1]} M_2^{(i)}) \times M_2^{(i)}$

to stop the formula from becoming too long,
 keep $(a^{[1][i]})_2 = \sigma(W_{0,2}^{[1]} + W_{1,2}^{[1]} M_1^{[i]} + W_{2,2}^{[1]} M_2^{[i]})$

$$\therefore \Delta^{(i)} = 2 W_{2,2}^{[2]} \times (\alpha^{(i)} - y^{(i)}) (0^{(i)}) (1 - \alpha^{(i)}) \\ \times (a_2^{[1][i]}) (1 - a_2^{[1][i]}) \quad XM^{(i)}$$

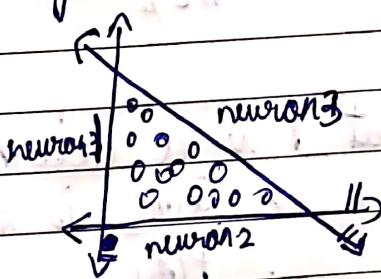
\therefore over all m e.g.,

$$W_{1,2}^{[1]} := W_{1,2}^{[1]} - \alpha \times \frac{2}{m} \sum_{i=1}^m ()$$

(b) Yes, it is possible to get 100% accuracy
 on the data.

possible :: as given in the hint, the 3 neurons
 can learn a Δ boundary s.t. the region
 of a^3 overlap with o_3 in the data set.

something like,



$$x_1 = 0.5$$

neuron 1 can learn the line - $x_1 = 0.5$

neuron 2 " " " " - $x_2 = 0.5$

" 3 " " " " - $x_1 + x_2 = 4$

(then o/p layer weights will need to invert)
 o/p from neuron 1 & 2. This way we will have
 a classifier detecting pts. inside o_3 outside the Δ .

(c) Not possible, the whole network becomes a single linear unit & we can learn only one line & dataset is not linearly separable.

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

$$\begin{aligned} P(m) > 0 \forall m \\ k \cdot \log 0 = 0. \end{aligned} \quad \hookrightarrow D_{KL}(P(x) || Q(x))$$

(a)

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

$$= - \sum_{m \in X} P(m) \times (\log(Q(m)))$$

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

note that P is a deterministic fn
for an m $\Rightarrow \frac{Q(m)}{P(m)}$ is 0

& so if $X \sim P$ i.e. frequency
of $m \in X = P(m) \quad \forall m \in X$

$$D_{KL}(P || Q) = E[-\log(Q(m)/P(m))]$$

$$\geq -\log E[Q(m)/P(m)]$$

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

$$D_{KL}(P||Q) \geq -\log \left(\sum_{m \in X} Q(m) \right)$$

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

$$D_{KL}(P||Q) \geq 0, P=Q \text{ iff } D_{KL}(P||Q)=0$$

$\therefore P \text{ & } Q \text{ are pmfs of } X$

$$\begin{aligned} E[f(X)] \\ = f(E[X]) \\ = 0 \end{aligned}$$

~~$$D_{KL}(P(X|Y)||Q(X|Y)) = \sum_y \sum_m P(y) P(m|y) \log \frac{P(m|y)}{Q(m|y)}$$~~

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

$$= \sum_m P(m) \sum_y P(y|m) \log \frac{P(m)P(y|m)}{Q(m)Q(y|m)}$$

$$= \sum_m \left(P(m) \sum_y P(y|m) \log \frac{P(m)}{Q(m)} \right) + D_{KL}(P(Y|X)||Q(Y|X))$$

$$= \sum_m P(m) \log \left(\frac{P(m)}{Q(m)} \times \sum_y P(y|m) \right) + D_{KL}(\quad)$$

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

(c) clearly, $\arg \max_{\theta} \sum_{i=1}^m P_\theta(m^{(i)})$

$$(c) \text{ clearly, } \arg \max_{\theta} \sum_{i=1}^m \log P_\theta(m^{(i)}) = \sum_{i=1}^m \frac{1}{m} \log \left(\frac{P_\theta(m^{(i)})}{1/m} \right)$$

\because we are dividing & subtracting constants

$$\Rightarrow \arg \max_{\theta} \sum_{i=1}^m \log P_\theta(m^{(i)}) = \arg \max_{\theta} \sum_{i=1}^m \hat{P}(m^{(i)}) \log \frac{P_\theta(m^{(i)})}{\hat{P}(m^{(i)})}$$

$$= \arg \max_{\theta} - D_{KL}(\hat{P} || P_\theta)$$

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

~~the remark is also~~

v good remark.

Q5 (a) Code

(b) 16 colors \Rightarrow 4 bits to completely describe a color.

So $512 \times 512 \times \frac{1}{2}$ size from $512 \times 512 \times 3$

by a factor of 6.

PS - 3 (Contd.)

(a) Score $f^n = \nabla_{\theta} \log \{p(y; \theta)\}$

$$E[\text{Score } f^n] = E_{y \sim p(y; \theta)} [\nabla_{\theta} \log \{p(y; \theta)\}]$$

$$= \int (p(y; \theta)) \nabla_{\theta} (\log p(y; \theta')) dy$$

$$= \int_{-\infty}^{\infty} p(y; \theta) \times \frac{1}{p(y; \theta)} \nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta} dy$$

$$= \nabla_{\theta'} \left(\int_{-\infty}^{\infty} p(y; \theta') dy \right) \Big|_{\theta'=\theta}$$

$$= \nabla_{\theta'} (1) \Big|_{\theta'=\theta}$$

$$= 0$$

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

↳ Covariance matrix of the score f^n .

let $\nabla_{\theta'} \log p(y; \theta') \Big|_{\theta'=\theta} = S$

$$I(\theta) = \text{Cov}_{y \sim p(y; \theta)} [S]$$

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

$$J(\theta) = E_{y \sim p(y; \theta)} [(s)(s)^T] \quad (\because E[s] = 0)$$

~~Ans~~

$$\Rightarrow J(\theta) = E_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta') \times \nabla_{\theta'} \log p(y; \theta')^T] \Big|_{\theta'=\theta}$$

$$(C) \left[\nabla_{\theta'} \log p(y; \theta') \nabla_{\theta'} \log p(y; \theta')^T \right]_{ij} \Big|_{\theta'=\theta} = \frac{\partial \log p(y; \theta')}{\partial \theta_i} \times \frac{\partial \log p(y; \theta')}{\partial \theta_j} \Big|_{\theta'=\theta}$$

$$= \frac{1}{(p(y; \theta))^2} \times \left(\frac{\partial p(y; \theta')}{\partial \theta_i} \times \frac{\partial p(y; \theta')}{\partial \theta_j} \right) \Big|_{\theta'=\theta}$$

~~(Ans)~~

$$= \frac{1}{p(y; \theta)^2} \times \frac{\partial^2 p(y; \theta')}{\partial \theta_i \partial \theta_j} \Big|_{\theta'=\theta}$$

by Clairaut's thm.

Now,

\rightarrow note that we cannot split this
because: log is not twice differentiable.

$$[-\nabla_{\theta'}^2 \log p(y; \theta')]_{ij} \Big|_{\theta'=\theta} = - \frac{\partial^2 \{\log p(y; \theta')\}}{\partial \theta_i \partial \theta_j} \Big|_{\theta'=\theta}$$

$$= - \frac{\partial}{\partial \theta_i} \left(\frac{1}{p(y; \theta')} \times \frac{\partial^2 p(y; \theta')}{\partial \theta_i \partial \theta_j} \right) \Big|_{\theta'=\theta}$$

~~$\frac{1}{p(y; \theta')} \times \frac{\partial^2 p(y; \theta')}{\partial \theta_i \partial \theta_j}$~~

$$= \left(-1 \times \frac{1}{p(y; \theta')^2} \times \frac{\partial^2 p(y; \theta')}{\partial \theta_i \partial \theta_j} - \frac{1}{p(y; \theta')} \times \frac{\partial^2 p(y; \theta')}{\partial \theta_i \partial \theta_j} \right) \Big|_{\theta'=\theta}$$

$$\Rightarrow \left[-\nabla_{\theta}^2 \log p(y; \theta') \right]_{ij} \Big|_{\theta'=\theta} = \frac{1}{p(y; \theta)^2} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j}$$

$$- \frac{1}{p(y; \theta)} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \Big|_{\theta'=\theta}$$

$$E_{y \sim p(y; \theta)} \left[-\nabla_{\theta}^2 \log p(y; \theta') \right]_{ij} \Big|_{\theta'=\theta} = E_{y \sim p(y; \theta)} \left[\frac{1}{p(y; \theta)^2} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \Big|_{\theta'=\theta} \right]$$

$$- \int_{-\infty}^{\infty} \frac{p(y; \theta)}{p(y; \theta)} \frac{\partial^2 p(y; \theta)}{\partial \theta_i \partial \theta_j} \Big|_{\theta'=\theta} dy$$

↳ this is = 0.

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

∴ proved.

(d) $D_{KL}(p_\theta || p_{\theta+d}) = \int_{-\infty}^{\infty} p(y; \theta) \log \frac{p(y; \theta)}{p(y; \theta+d)} dy$

if $\hat{\theta} = \theta + d$, then Taylor Series helps us approximate it-

$$f(\hat{\theta}) = f(\theta) + (\hat{\theta} - \theta)^T \nabla_{\theta} f(\theta') \Big|_{\theta'=\theta} + \frac{1}{2} (\hat{\theta} - \theta)^T \nabla_{\theta}^2 f(\theta') \Big|_{\theta'=\theta} (\hat{\theta} - \theta)$$

now,

$$D_{KL}(p_\theta || p_{\theta+d}) = E_{y \sim p(y; \theta)} [p(y; \theta)] - E_{y \sim p(y; \theta)} [p(y; \hat{\theta})]$$

$$= E_{y \sim p(y; \theta)} [p(y; \theta)] - \left(E_{y \sim p(y; \theta)} [p(y; \theta)] + d^T E_{y \sim p(y; \theta)} [\nabla_{\theta} p(y; \theta)] \right)$$

$$\approx E_{y \sim p(y; \theta)} [p(y; \theta)] - E_{y \sim p(y; \theta)} [p(y; \theta)]$$

Applying the
Taylor series
expansion

$$\left. -E_{y \sim p(y; \theta)} [d^T \nabla_{\theta} p(y; \theta')] \right|_{\theta'=\theta}$$

$$\left. -\frac{1}{2} E_{y \sim p(y; \theta)} [d^T \nabla_{\theta'}^2 p(y; \theta')] \right|_{\theta'=\theta} d$$

is 0 ! part (a).

$$\approx -d^T E_{y \sim p(y; \theta)} [\nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta}]$$

$$-\frac{1}{2} d^T E_{y \sim p(y; \theta)} [\nabla_{\theta'}^2 p(y; \theta') \Big|_{\theta'=\theta}] d$$

$$\approx \frac{1}{2} d^T J(\theta) d$$

$$(e) \quad d^* = \arg \max_d l(\theta + d)$$

$$\text{Subject to } D_{KL}(p_\theta || p_{\theta+d}) = c$$

Lagrangian -

$$L(d, \lambda) = l(\theta) + d^T \nabla_{\theta} l(\theta) \Big|_{\theta=\theta} - \lambda \left[\frac{1}{2} d^T I(\theta) d - c \right]$$

$$\Rightarrow \mathcal{L}(d, \lambda) \approx \log p(y; \theta) + d^T \nabla_{\theta'} p(y; \theta')|_{\theta'=\theta}$$

$$-\lambda \left(\frac{1}{2} d^T \mathcal{I}(\theta) d - c \right)$$

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

$$\frac{p(y; \theta)}{p(y; \theta')}$$

$$-\lambda \left(\frac{1}{2} d^T \mathcal{I}(\theta) d - c \right)$$

Linear Systems -

$$\nabla_d \mathcal{L}(d, \lambda) = 0 + \nabla_{\theta'} p(y; \theta')|_{\theta'=\theta} - \frac{\lambda}{2} \mathcal{I}(\theta) d \neq 0$$

$$= 0$$

$$\Rightarrow \nabla_{\theta'} p(y; \theta')|_{\theta'=\theta} = p(y; \theta) \times \lambda \mathcal{I}(\theta) d$$

$\hookrightarrow (a)$

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

$$\Rightarrow c = \frac{1}{2} d^T \mathcal{I}(\theta) d \rightarrow (b)$$

from (a) -

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

from (b) A d -

$$\frac{1}{2} \times \frac{\lambda}{p(y; \theta)} \left(\mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0} \right)^T \mathcal{I}(\theta)$$

$$\times \left(\mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0} \right)$$

$$\frac{\lambda}{2p(y; \theta)^2} \times \left(\nabla_{\theta} p(y; \theta) \Big|_{\theta=0} \right)^T \mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0}$$

$$\lambda = \sqrt{\frac{1}{2c p(y; \theta)^2} \left(\nabla_{\theta} p(y; \theta) \Big|_{\theta=0} \right)^T \mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0}}$$

calculating d^* -

~~$$d^* = \sqrt{2c} \left(\nabla_{\theta} p(y; \theta) \Big|_{\theta=0} \right)^T \mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0}$$~~

$$d^* = \sqrt{\frac{2c}{\left(\nabla_{\theta} p(y; \theta) \Big|_{\theta=0} \right)^T \mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0}}} \mathcal{I}(\theta)^{-1} \nabla_{\theta} p(y; \theta) \Big|_{\theta=0}$$

(P) for GLMs,

$$(P(y; \theta))$$

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla_{\theta} l(\theta)$$

Newton's method.

$$J(\theta) = -E_{y \sim p(y; \theta)} [+ \nabla_{\theta}^2 p(y; \theta')]_{\theta'=\theta}$$

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

$$= -E_{y \sim p(y; \theta)} [H]$$

$$\Rightarrow \text{using } \hat{d} = \frac{1}{\lambda} J(\theta)^{-1} \nabla_{\theta} \log p(y; \theta') \Big|_{\theta'=\theta}$$

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

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

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

Same

PS-3 (contd.)

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

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

$$l_{\text{unsup}}(\theta) = \sum_{i=1}^m \log p(m^{(i)}; \theta)$$

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

for E-step -

~~$$l_{\text{semi-sup}}(\theta) = \sum_{i=1}^m \log Q_i(z^{(i)})$$~~

$$l_{\text{semi-sup}}(\theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} Q_i(z^{(i)}) \cdot \frac{p(m^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} + l_{\text{sup}}(\theta)$$

Now with the same steps as in the notes, we would apply the Jensen's inequality & try to make the lower bound tight at current θ . Note that to ensure "tightness" it only depends on the first term so we would derive

$$Q_i(z^{(i)}) = p(z^{(i)} | m^{(i)}; \theta)$$

Same as in the lecture notes.

now with Jensen's inequality,

$$l_{\text{semi-sup}}(\theta) \geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(m^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} + \alpha l_{\text{sup}}(\theta)$$

\Rightarrow for the N-step we will update

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

$$(a) l_{\text{semi-sup}}(\theta^{(t+1)}) = \sum_{i=1}^m \log \sum_{z^{(i)}} Q_i(z^{(i)}) \times \frac{p(m^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i(z^{(i)})}$$

$$+ \alpha \sum_{i=1}^{\tilde{m}} \log p(\tilde{m}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)})$$

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

$$+ \alpha \sum_{i=1}^{\tilde{m}} \log p(\tilde{m}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)})$$

(Jensen's inequality)

$$\text{now, } \therefore \theta^{(t+1)} := \arg \max_{\theta} (l_{\text{semi-sup}}(\theta^{(t)}))$$

$$l_{\text{semi-sup}}(\theta^{(t+1)}) \geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(m^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i(z^{(i)})}$$

$$+ \alpha \sum_{i=1}^{\tilde{m}} \log p(\tilde{m}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)})$$

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

Note that the first term in the last eqn. is exactly equal to $\text{lowerup}(\theta^{(t)}) \because$ of our choice of $Q_i(z^{(i)})$ to make the lowerbound tight.

$$(b) Q_i(z^{(i)}) = P(z^{(i)} | M^{(i)}; \theta)$$

from given info. \nearrow , i.e., E-step only involves latent variables

$$\text{here } \theta = \{\mu, \Sigma, \phi\}$$

$$\Rightarrow Q_i(z^{(i)}=j) = P(z^{(i)}=j | M^{(i)}; \mu, \Sigma, \phi)$$

$$= \underbrace{P(M^{(i)} | z^{(i)}=j; \mu, \Sigma)}_{\sum_{j=1}^k p(\alpha^{(i)} | z^{(i)}=j; \mu, \Sigma)} \times \underbrace{p(z^{(i)}=j | \phi)}$$

$$\sum_{j=1}^k p(\alpha^{(i)} | z^{(i)}=j; \mu, \Sigma) \times p(z^{(i)}=j | \phi)$$

$$\text{let } Q_i(z^{(i)}=j) = w_j^{(i)}.$$

(c) M-step - (parameters re-estimated include Σ, μ, ϕ)

$$l_{\text{semi-sup}}(\theta^{(t)}) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{p(M^{(i)} | z^{(i)}=j; \mu, \Sigma) \times p(z^{(i)}=j | \phi)}{w_j^{(i)}}$$

$$+ \alpha \sum_{i=1}^m \log \{ p(\tilde{M}^{(i)} | \tilde{Z}^{(i)}=j; \mu, \Sigma) \times p(\tilde{z}^{(i)}=j | \phi) \}$$

terms involving ϕ ,

$$\begin{aligned} L(\phi) = & \sum_{i=1}^m \sum_{j=1}^k \log w_j^{(i)} \log p(z^{(i)}=j | \phi) + \alpha \sum_{i=1}^m \log p(\tilde{z}^{(i)}=j | \phi) \\ & + P \left(\sum_{i=1}^k \phi_i - 1 \right) \end{aligned}$$

note that α is not a Lagrange multiplier
but β is. ($\because \sum_j \phi_j = 1$)

$$\nabla_{\phi} L(\phi) =$$

$$\frac{\partial L(\phi)}{\partial \phi_j} = \sum_{i=1}^m w_j^{(i)} x_i + \alpha \sum_{i=1}^m I\{z^{(i)} = j\} + \beta = 0$$

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

$$\phi_j = \left(\sum_{i=1}^m w_j^{(i)} + \alpha \sum_{i=1}^m I\{z^{(i)} = j\} \right) / (-\beta)$$

applying $\sum_{j=1}^k (\cdot)$ on both sides,

$$1 = \left(\sum_{i=1}^m (1) + \alpha \sum_{i=1}^m I\{z^{(i)} = j\} \right) / (-\beta)$$

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

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

terms involving μ_j ,

$$\begin{aligned}
 \mathcal{L}(\mu) &= \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} (\log p(m^{(i)} | z_j^{(i)}; \mu, \Sigma) - \log w_j^{(i)}) \\
 &\quad + \alpha \sum_{i=1}^m \log (p(\tilde{m}^{(i)} | \tilde{z}^{(i)}; \mu, \Sigma)) \\
 &= \sum_{i=1}^m \sum_{j=1}^k \left\{ w_j^{(i)} \log \left(\sum_{l=1}^k p(m^{(i)} | z_l^{(i)} = j; \mu, \Sigma) \right) \right. \\
 &\quad \left. + \alpha \sum_{i=1}^m \log (p(\tilde{m}^{(i)} | \tilde{z}^{(i)}; \mu, \Sigma)) \right\}
 \end{aligned}$$

($\frac{\partial \mathcal{L}}{\partial \mu_j}$)

$$\begin{aligned}
 \frac{\partial \mathcal{L}(\mu)}{\partial \mu_j} &= \left(\sum_{i=1}^m \right) \sum_{j=1}^k \left[\left(\sum_{i=1}^m w_j^{(i)} m^{(i)} - \mu_j \sum_{i=1}^m w_j^{(i)} \right) \right. \\
 &\quad \left. + \alpha \left(\sum_{i=1}^m I\{z^{(i)} = j\} \tilde{m}^{(i)} - \mu_j \sum_{i=1}^m I\{z^{(i)} = j\} \right) \right] \\
 &= 0
 \end{aligned}$$

$$\Rightarrow \mu_j = \frac{\sum_{i=1}^m w_j^{(i)} m^{(i)}}{\sum_{i=1}^m w_j^{(i)}} + \alpha \frac{\sum_{i=1}^m I\{z^{(i)} = j\} \tilde{m}^{(i)}}{\sum_{i=1}^m I\{z^{(i)} = j\}}$$

$$\text{similarly, } \Sigma_j = \frac{\sum_{i=1}^m w_j^{(i)} (m^{(i)} - \mu_j)(m^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}} + \alpha \frac{\sum_{i=1}^m I\{z^{(i)} = j\} \tilde{m}^{(i)} \tilde{m}^{(i)}^T}{\sum_{i=1}^m I\{z^{(i)} = j\}}$$

Code

Code

(SSEM)

(i) Semi-Supervised EM took much less iterations ~ 50

④ Un - " EM ~ 1000
(USEM)

(ii) SSEM was very stable, even with different initializations.
I got the same plots.

USEM was not very stable, although some pts. were
classified to the same class, the cluster itself changed

(iii) SSEM ~~was~~ is ~~completely~~ almost correct if not
completely correct, acc. to given info.

USEM is not very accurate to the give info.