

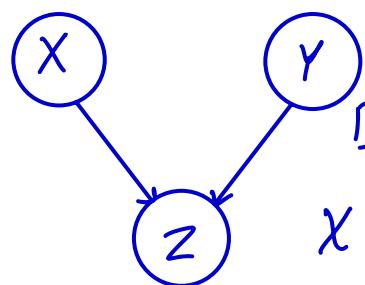
Q1

A) 1) $I(X;Y|Z) < I(X;Y)$

Let X and Y be random variables such that $I(X;Y) > 0$.
And X and Y deterministic functions of Z , therefore we have
 $I(X;Y|Z) = 0$.

2) $I(X;Y|Z) > I(X;Y)$

consider the following bayesian network. X and Y



are independent, giving $I(X;Y) = 0$.

But given Z , we have either X or Y . If X doesn't happen, the Y must happen and vice versa. Hence there is a conditional dependence between X and Y , given Z . Therefore $I(X;Y|Z) > 0 = I(X;Y)$

B) 1) $H(X, Y, Z) - H(X, Y) \leq H(X, Z) - H(X)$

$$H(X, Y, Z) = H(X) + H(Y|X) + H(Z|X, Y)$$

$$H(X, Y) = H(X) + H(Y|X)$$

$$\Rightarrow LHS := H(Z|X, Y)$$

$$RHS := H(Z|X)$$

Conditioning reduces entropy :

$$H(Z|X,Y) \leq H(Z|X)$$

Equality holds when $Z|X$ is independent of Y .

2) $I(X;Z|Y) \geq I(Z;Y|X) - I(Z;Y) + I(X;Z)$

$$I(Z;Y|X) - I(Z;Y) + I(X;Z)$$

$$= H(Z|X) - H(Z|Y,X) - H(Z) + H(Z|Y)$$

$$+ H(Z) - H(Z|X)$$

$$= H(Z|Y) - H(Z|X,Y)$$

$$= I(X,Z|Y)$$

C) ii First, we have to find the distribution of Z . Let Y be a uniform discrete random variable with values $\{4,5,6\}$. Thus

$$P(Z=A) = \sum_y P(Y=y) P(Z=A|Y=y)$$

$$= \frac{1}{3} \left(\frac{1}{1+e^{-0.5 \times (-1)}} + \frac{1}{1+e^{-0.5 \times 0}} + \frac{1}{1+e^{-0.5 \times 1}} \right)$$

$$= \frac{1}{2}$$

$$P(Z=B) = \sum_y P(Y=y) P(Z=B|Y=y) = \sum_y P(Y=y) \frac{1-P(Z=A|Y=y)}{2}$$

$$= \frac{1-P(Z=A)}{2}$$

$$= \frac{1 - \frac{1}{2}}{2} = \frac{1}{4}$$

$$P(Z=C) = \frac{1 - \frac{1}{2}}{6} = \frac{1}{12}$$

$$P(Z=D) = \frac{1 - \frac{1}{2}}{3} = \frac{1}{6}$$

$$I(Z; Y) = H(Z) - H(Z|Y)$$

$$H(Z) = - \sum_{z \in \{A, B, C, D\}} P(Z=z) \log P(Z=z) = 1.73$$

$$H(Z|Y) = - \sum_y P(Y=y) H(Z|Y=y)$$

$$= 1.70$$

$$\Rightarrow I(Z; Y) = 0.03$$

$I(Z; Y)$ is a measure of how much information we would have about Z , given we know Y . i.e. how much information does the reward give us about the action taken. If $I(Z; Y)$ is high, it means they are highly related and dependent and vice versa.

$$D) D_{KL}(p(x,y,z) \parallel p(x)p(y)p(z)) = -H(X,Y,Z) + H(X) + H(Y) + H(Z)$$

$$\begin{aligned} D_{KL}(p(x,y,z) \parallel p(x)p(y)p(z)) &= \sum p(x,y,z) \log \frac{p(x,y,z)}{p(x)p(y)p(z)} \\ &= \sum p(x,y,z) \log p(x,y,z) \\ &\quad - \sum p(x,y,z) \log [p(x)p(y)p(z)] \\ &= -H(X,Y,Z) - \sum p(x,y,z) \log p(x) \\ &\quad - \sum p(x,y,z) \log p(y) - \sum p(x,y,z) \log p(z) \\ &= -H(X,Y,Z) + H(X) + H(Y) + H(Z) \end{aligned}$$

$$E) I(X_1;X_3) + I(X_2;X_4) \leq I(X_1;X_4) + I(X_2;X_3)$$

Data processing inequality: $X \rightarrow Y \rightarrow Z \Rightarrow I(X;Z) \leq I(X;Y)$

$$I(X_1;X_3) + I(X_2;X_4) = H(X_1) - H(X_1|X_3) + H(X_2) - H(X_2|X_4)$$

$$I(X_1;X_4) + I(X_2;X_3) = H(X_1) - H(X_1|X_4) + H(X_2) - H(X_2|X_3)$$

So we have to prove:

$$H(X_1|X_3) + H(X_2|X_4) \geq I(X_1|X_4) + I(X_2|X_3)$$

$$H(X_1, X_2 | X_3) = H(X_1 | X_3) + H(X_2 | X_1, X_3)$$

$$\Rightarrow H(X_1 | X_3) = H(X_1, X_2 | X_3) - H(X_2 | X_1, X_3)$$

$$H(X_2 | X_4) = H(X_2, X_1 | X_4) - H(X_1 | X_2, X_4)$$

$$H(X_1|X_4) = H(X_1, X_2|X_4) - H(X_2|X_1, X_4)$$

$$H(X_2|X_3) = H(X_2, X_1|X_3) - H(X_1|X_2, X_3)$$

$$H(X_1, X_2|X_3) - H(X_2|X_1, X_3) + \cancel{H(X_2, X_1|X_4)} - \cancel{H(X_1|X_2, X_4)}$$

$$\cancel{H(X_1, X_2|X_4)} - \cancel{H(X_2|X_1, X_4)} + \cancel{H(X_2, X_1|X_3)} - \cancel{H(X_1|X_2, X_3)}$$

We need to show :

$$H(X_2|X_1, X_3) \leq H(X_2|X_1, X_4)$$

We know :

$$H(X_2|X_1, X_3) = H(X_2|X_1, X_3, X_4) \quad (\text{Markov property})$$

And conditioning reduces entropy, thus:

$$H(X_2|X_1, X_4) \geq H(X_2|X_1, X_3, X_4)$$

So we are done.

Q2

A) By first order condition, f is convex iff

$\text{dom } f$ is convex, $f(y) \geq f(x) + \nabla f(x)^T(y-x)$

Let for some $x_0 \in \text{dom } f$; $\nabla f(x_0) = 0$. Then by FOC we have :

$\forall y \in \text{dom } f$; $f(y) \geq f(x_0) + \nabla f(x_0)^T(y-x_0) = f(x_0)$. Thus x_0 is a minimum point.

$$B) \quad \min \quad x^2 + y^2 + z^2$$

$$\text{s.t.} \quad z^2 = x^2 + y^2$$

$$z = x + y + 1$$

$$z^2 = x^2 + y^2 \Rightarrow x^2 + y^2 + 1 + 2xy + 2x + 2y = x^2 + y^2$$

$$\Rightarrow y = -\frac{x+1}{x+1}$$

$$\Rightarrow \frac{1}{2(x+1)^2} (4x^4 + 8x^3 + 8x^2 + 4x + 1)$$

$$\frac{d}{dx} \frac{1}{(x+1)^3} (4x^4 + 12x^3 + 12x^2 + 6x + 1) = 0$$

$$\Rightarrow \lambda = 0, -1 \pm \frac{\sqrt{2}}{2}, \cancel{-\frac{1}{2}} \pm \cancel{\frac{i}{2}}$$

$$\Rightarrow \min = 6 - 4\sqrt{2} \quad \text{for } x = \frac{\sqrt{2}}{2} - 1$$

$$\max = +\infty$$

c) prove $L(\mu, \lambda) \leq \min_x f(x)$

Primal : $\min f(x)$

s.t. $g_i(x) \leq 0 \quad i=1, \dots, m$

$h_i(x) = 0 \quad i=1, \dots, p$

dual : $\max L(\mu, \lambda)$

s.t. $\lambda \geq 0$

$$L(\mu, \lambda) = \inf_x L(x, \mu, \lambda)$$

$$L^*(x, \mu, \lambda) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{i=1}^p \mu_i h_i(x)$$

consider any feasible point in the primal problem and the dual problem.

e.g.: $\lambda_i; g_i(x) \leq 0 \quad \forall i=1, \dots, m$ and $h_i(x) = 0 \quad \forall i=1, \dots, p$

$(\lambda, \mu) : \lambda \geq 0$

$$\begin{aligned}\Rightarrow \mathcal{L}(x, \lambda, \mu) &= f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{i=1}^p \mu_i h_i(x) \\ &\leq f(x)\end{aligned}$$

$$\Rightarrow L(\lambda, \mu) = \inf_x \mathcal{L}(x, \lambda, \mu) \leq f(x) \quad \forall x \text{ in dom primal}$$

$$\Rightarrow L(\lambda, \mu) \leq \min_x f(x)$$

Consider the primal problem

$$\begin{array}{ll}\text{min} & e^{-x} \\ \text{s.t.} & x^2/y \leq 0\end{array}$$

$$D = \{(x, y), y > 0\}$$

thus $p^* = 1$. The lagrangian is $L(x, y, \lambda) = e^{-x} + \lambda x^2/y$ thus the dual function is:

$$g(\lambda) = \inf_{x, y} (e^{-x} + \lambda x^2/y) = 0$$

so we have

$$\begin{array}{ll}d^* = & \max \quad 0 \\ & \text{s.t. } \lambda \geq 0\end{array}$$

Thus $d^* = 0$. so $p^* - d^* = 1$ and strong duality does not hold

$$D) \quad \min -\sum \log(\alpha_i + x_i)$$

$$\text{s.t. } x \geq 0$$

$$1^T x = 1$$

$$\mathcal{L}(x, \lambda, v) = -\sum \log(\alpha_i + x_i) - \lambda^T x + v 1^T x - v$$

to solve the problem, we find the variables such that they satisfy KKT conditions.

$$x \geq 0, 1^T x = 1, \lambda \geq 0, \lambda_i x_i = 0 \quad i=1, \dots, n$$

$$\frac{\partial \mathcal{L}}{\partial x_i} = \frac{1}{\alpha_i + x_i} - \lambda_i + v = 0$$

$$\Rightarrow \lambda_i = v - \frac{1}{\alpha_i + x_i} \Rightarrow x_i(v - \frac{1}{\alpha_i + x_i}) = 0$$

$$, v \geq \frac{1}{\alpha_i + x_i}$$

$$\Rightarrow \text{if } v < \frac{1}{\alpha_i} \Rightarrow x_i > 0 \Rightarrow \lambda_i = 0$$

$$\Rightarrow v = \frac{1}{\alpha_i + x_i} \Rightarrow x_i = \frac{1}{v} - \alpha_i$$

if $v \geq \frac{1}{\alpha_i}$, then if $x > 0$, we have $v > \frac{1}{\alpha_i + x_i}$. But by the complementary slackness we have $x_i(v - \frac{1}{\alpha_i + x_i}) = 0$, which is contradiction. Therefore we have

$$x_i = \begin{cases} \frac{1}{v} - \alpha_i & , \text{ if } v < \frac{1}{\alpha_i} \text{ or } \min\{0, \frac{1}{v} - \alpha_i\} \\ 0 & , \text{ o.w.} \end{cases}$$

and we have $\sum \min\{0, \frac{1}{v} - \alpha_i\} = 1$, which gives the value of v .

$$E) E[x] = \int_{-\infty}^{+\infty} x P(x) dx = \mu$$

$$\text{Var}(X) = E[(X - \mu)^2] = \int_{-\infty}^{+\infty} (x - \mu)^2 P(x) dx = \sigma^2$$

$$F(-\infty \leq x \leq +\infty) = \int_{-\infty}^{+\infty} P(x) dx = 1$$

Since we are integrating over all of the real numbers, the mean of the distribution has no effect on the differential entropy. Thus we only take the variance and the probability sum constraints into account.

$$\text{minimize} \quad \int_{-\infty}^{+\infty} P(x) \log P(x) dx$$

$$\text{subject to} \quad \int_{-\infty}^{+\infty} (x - \mu)^2 P(x) dx = \sigma^2$$

$$\int_{-\infty}^{+\infty} P(x) dx = 1$$

$$\Rightarrow \mathcal{L}(P(x), v_1, v_2) = - \int_{-\infty}^{+\infty} P(x) \log P(x) dx + v_1 (1 - \int_{-\infty}^{+\infty} P(x) dx) \\ + v_2 (\sigma^2 - \int_{-\infty}^{+\infty} (x-\mu)^2 P(x) dx)$$

By the KKT conditions, for the optimal y , the small variation of \mathcal{L} , i.e. $\delta \mathcal{L}$ must be zero.

$$\delta \mathcal{L} = \int_{-\infty}^{+\infty} \delta P(x) \left(-\frac{\partial P(x) \log P(x)}{\partial P(x)} + v_1 \frac{\partial P(x)}{\partial P(x)} + v_2 (x-\mu)^2 \frac{\partial P(x)}{\partial P(x)} \right) \\ = \int_{-\infty}^{+\infty} \delta P(x) (-\log P(x) - 1 + v_1 + v_2 (x-\mu)^2) = 0$$

By the Euler-Lagrange equation (or since the above equality has to hold for any small $\delta P(x)$), regardless of the sign):

$$-\log P(x) - 1 + v_1 + v_2 (x-\mu)^2 = 0 \implies P(x) = e^{-v_1 - 1 - v_2 (x-\mu)^2}$$

Using the constraints we get:

$$\int_{-\infty}^{+\infty} e^{-v_1 - 1 - v_2 (x-\mu)^2} dx = e^{-v_1 - 1} \left[\frac{\sqrt{\pi}}{2\sqrt{v_2}} \operatorname{erf}(\sqrt{v_2}(x-\mu)) \right]_{-\infty}^{+\infty} \\ = e^{-v_1 - 1} \times \frac{\sqrt{\pi}}{2\sqrt{v_2}} \times 2 = 1 \\ e^{-v_1 - 1} = \frac{\sqrt{v_2}}{\sqrt{\pi}}$$

$$\begin{aligned}
 \frac{\sqrt{\nu_2}}{\sqrt{\pi}} \int_{-\infty}^{+\infty} (x-\mu)^2 e^{-\nu_2(x-\mu)^2} dx = \sigma^2 &= \left(\frac{(x-\mu) e^{-\nu_2(x-\mu)^2}}{2a} \right) \Big|_{-\infty}^{+\infty} \\
 &\quad - \frac{\sqrt{\pi}}{4\nu_2} \left. \text{erf}(\sqrt{\nu_2}(x-\mu)^2) \right|_{-\infty}^{+\infty} \Big| \frac{\sqrt{\nu_2}}{\sqrt{\pi}} \\
 &= \frac{\sqrt{\pi}}{4\sqrt{\nu_2}} \times 2 \times \frac{\sqrt{\nu_2}}{\sqrt{\pi}} = \sigma^2
 \end{aligned}$$

$$\Rightarrow \frac{1}{2\nu_2} = \sigma^2 \Rightarrow \nu_2 = \frac{1}{2\sigma^2}$$

$$\Rightarrow P(n) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The objective functional and the constraints of the problem are all convex, hence the KKT conditions give the optimal solution, since we have strong duality.

A3

A) To prove the stationary distribution of the given Markov chain is μ , we must show that if the transition probability matrix is P , then $P\mu = \mu$. Now we show that the SD is a uniform distribution.

For an arbitrary node x we have

$$\mu(x) \times \left(1 - \frac{|N(x)|}{M}\right) + \sum_{y \in N(x)} \mu(y) \frac{1}{|N(x)|} = \mu(x) - \frac{\mu(x) |N(x)|}{M} + \frac{\sum_{y \in N(x)} \mu(y)}{M}$$
$$= \mu(x)$$

Now consider the node with the highest probability. It's probability is the average of the probability of the nodes connected to it. Therefore the probabilities of this node and its neighbours are equal. We apply the same argument to the neighbours of the so-called node. Since the Markov chain is irreducible, it is connected. Hence the probability of all nodes is equal. Therefore the stationary distribution is uniform.

(If the graph is directed, the in-degree and out-degree of each node has to be equal, and we can reason as above for this case.)

B) The problem with this method is that each valid vector in the sample contains more than one valid state in itself. But we count only one of them. (The same problem arises with FVMC algorithm in estimating $\langle Q \rangle$ values.) Therefore this method has high variance.

- C) 1) In the described Markov process, the transition is only possible adjacent valid state. Since we have $\sum_{i=1}^n a_i < b$, every state is valid. Thus this Markov chain is irreducible. Therefore by part (A), this Markov chain has a uniform stationary distribution.
- 2) The transition probability of going to each valid state from one is $\frac{1}{n}$. For a given state x we have:

$$P_{x,y} = \begin{cases} \frac{1}{n} & y=x, y \in N(x) \\ 1 - \frac{|N(x)|}{n} & y \neq x \\ 0 & y \notin N(x) \end{cases}$$

where $N(x)$ is the set of states which are reachable from x . If we reduce the whole process to the valid states, the the Markov process is irreducible, thus by part A, it has a uniform distribution. Thus if we run the process for infinite steps, the probability of each state would be equal, and so we should visit all valid states. so we only need to count the newly visited states.

Q4

- A) s is the reward. If the images are the same, the higher the value of s , the higher its probability (given that $\lambda_s > 0$). And if the images are different, it's vice versa. Which it seems sensible. ($\lambda_d > 0$)
 The α_s and α_d coefficients are for normalizations, so the given functions are probabilistic.
- B) We would like to compute $P(\text{same}_j | s_1, \dots, s_n)$, which by the bayes rule we have ..

$$P(\text{same}_j | s_1, \dots, s_n) = \frac{P(s_1, \dots, s_n | \text{same}_j) P(\text{same}_j)}{P(s_1, \dots, s_n)}$$

$$= \frac{P(s_1, \dots, s_n | \text{same}_j)}{\sum_i P(s_1, \dots, s_n | \text{same}_i) P(\text{same}_i)} \times P(\text{same}_j)$$

since the prior for all the images is equal, then $\forall i \in \{1, \dots, n\} P(\text{same}_i) = \frac{1}{n}$
 and $P(\text{different}_j) = \frac{n-1}{n}$. Thus we have:

$$P(\text{same}_j | s_1, \dots, s_n) = \frac{P(s_1, \dots, s_n | \text{same}_j)}{\sum_i P(s_1, \dots, s_n | \text{same}_i)}$$

$$P(s_1, \dots, s_n | \text{same}_j) = \prod_k P(s_k | \text{same}_j) = P(s_j | \text{same}_j) \prod_{k \neq j} P(s_k | \text{different}_k)$$

$$= \alpha_{ss} \times \alpha_{ds} e^{\lambda_s s_j} \prod_{k \neq j} e^{-\lambda_d s_k}$$

$$\Rightarrow P(\text{same}_j | s_1, \dots, s_n) = \frac{e^{(\lambda_s + \lambda_d)s_j}}{\sum_i e^{(\lambda_s + \lambda_d)s_i}}$$

C) We want the probability of the case which the maximum reward is for a correct recognition. In other words, if M_s and M_d are the maximum values

of the scores for the same and different cases, we want to compute $P(M_s > M_d)$. By the equality of priors, we expect that N_s images belong to the "same" category, and $N - N_s$ to the other one. So if we suppose that the scores for each category are sorted, by the order stats, if $X_1 \leq \dots \leq X_K$ are K i.i.d random variables with CDF of F_X , then

$$F_{X_r}(x) = \sum_{j=r}^K \binom{K}{j} F_X(x)^j (1 - F_X(x))^{n-j}$$

$$f_{X_r}(x) = \frac{K!}{(r-1)!(K-r)!} f_X(x) F_X(x)^{r-1} (1 - F_X(x))^{n-r}$$

Therefore we have

$$F_D = F_{d_{(N-N_s)}} \text{ and } F_D(n) = F_{\text{diff}}(x)^{N - N_s}$$

$$f_D(n) = (N - N_s) f_D(n) F_{\text{diff}}(x)^{N - N_s - 1} \rightarrow f_d(x) = P(x | \text{different})$$

$$F_S = F_{s_{N_s}} \text{ and } F_D(n) = F_{\text{same}}(x)^{N_s}$$

$$f_S(n) = N_s f_S(n) F_{\text{same}}(x)^{N_s - 1} \rightarrow f_s(x) = P(x | \text{same})$$

Now we need to compute $P(M_s \geq M_d)$, which is

$$F_D(M_d) \left(1 - \int_0^{M_d} f_S(n) dn\right)$$

(Q5)

A) $\theta = (\mu, \sigma^2)$, $P_\theta(X) = \mathcal{N}(X; \mu, \sigma^2)$

Let the sequence of samples be (x_1, x_2, \dots, x_n) . Then and x_i s are i.i.d.

$$P_\theta(X) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

For simplicity, we use the log-likelihood function.

$$L_\theta(X) = \log P_\theta(X) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

To find the MLE, we derivate L_θ w.r.t θ and set it to zero.

$$\theta = (\mu, \sigma^2)$$

$$\begin{aligned} \frac{\partial}{\partial \mu} L_\theta(X) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \\ \Rightarrow \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n x_i \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} L_\theta(X) &= -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = 0 \\ \Rightarrow \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

For $\hat{\mu}$ we have :

$$\begin{aligned} E[\hat{\mu}] &= E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i] \\ &= \frac{1}{n} \sum_{i=1}^n \mu = \mu \end{aligned}$$

Thus MLE for μ is unbiased.

For $\hat{\sigma}^2$ we have :

$$\begin{aligned} E[\hat{\sigma}^2] &= E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2\right] = E\left[\frac{1}{n} \sum_{i=1}^n ((x_i - \mu) - (\hat{\mu} - \mu))^2\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n ((x_i - \mu)^2 - 2(x_i - \mu)(\hat{\mu} - \mu) + (\hat{\mu} - \mu)^2)\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - \frac{2}{n}(\hat{\mu} - \mu) \sum_{i=1}^n (x_i - \mu) + (\hat{\mu} - \mu)^2\right] \\ &= \frac{1}{n} \sum_{i=1}^n E[(x_i - \mu)^2] - E[2(\hat{\mu} - \mu)^2] + E[(\hat{\mu} - \mu)^2] \\ &= \sigma^2 - E[(\hat{\mu} - \mu)^2] = (1 - \frac{1}{n})\sigma^2 < \sigma^2 \end{aligned}$$

$$E\left[\left(\frac{1}{n} \sum_{i=1}^n x_i - \mu\right)^2\right] = E\left[\frac{1}{n^2} \sum_{i=1}^n (x_i - \mu)^2\right] = \frac{1}{n^2} \sum_{i=1}^n E[(x_i - \mu)^2] = \frac{\sigma^2}{n}$$

Thus MLE for variance is biased.

B) Let L be the likelihood function. To prove the invariance of MLE, we first define the induced likelihood function. For the function $T(\theta)$, the induced likelihood function L^* is defined as

$$L^*(\eta | X) = \sup_{\theta: T(\theta)=\eta} L(\theta | X)$$

and by definition, the value $\hat{\eta}$ that maximizes L^* is the

MLE for $\eta = T(\theta)$. We need to prove $L^*(\hat{\eta}|X) = L^*(T(\hat{\theta})|X)$.
we have :

$$\begin{aligned} L^*(\hat{\eta}|X) &= \sup_{\eta} L^*(\eta|X) = \sup_{\eta} \sup_{\theta: T(\theta)=\eta} L(\theta|X) \\ &= \sup_{\theta} L(\theta|X) \\ &= L^*(\hat{\theta}|X) \end{aligned}$$

And

$$L(\hat{\theta}|X) = \sup_{\theta: T(\theta)=T(\hat{\theta})} L(\theta) = L^*(T(\hat{\theta})|X)$$

thus

$$L^*(\hat{\eta}|X) = L^*(T(\hat{\theta})|X)$$

And we have :

$$T_1(\mu) = \frac{1}{\mu^2 + 1} \implies T(\hat{\mu}) = \frac{1}{\frac{1}{n} \sum x_i + 1} = \frac{n^2}{\sum x_i + n^2}$$

$$T_2(\sigma^2) = \sqrt{\sigma^2} \implies T(\hat{\sigma}^2) = \sqrt{\frac{1}{n} \sum (x_i - \hat{\mu})^2}$$

C) 1) Consider the estimator $\hat{\mu}_{IS} = \frac{1}{n} \sum_{i=1}^n w_i f(x_i)$. So

$$E_q[\hat{\mu}] = E_q \left[\frac{1}{n} \sum_{i=1}^n w_i f(x_i) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n E_q [w(x_i) f(x_i)]$$

For the random variable x we have

$$\begin{aligned} E_q [w(x) f(x)] &= \int w(x) f(x) q(x) dx \\ &= \int \frac{p(x)}{q(x)} q(x) f(x) dx \\ &= \int p(x) f(x) dx \\ &= E_p [f(x)] \end{aligned}$$

Therefore

$$E_q [\hat{\mu}_{IS}] = \frac{1}{n} \sum_{i=1}^n E_q [w(x_i) f(x_i)] = E_p [f(x)]$$

And thus the IS estimator is unbiased.

2) For the variance we have

$$\begin{aligned} \text{Var}(\hat{\mu}_{IS}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n w(x_i) f(x_i)\right) \\ &= \frac{1}{n} \text{Var}(w(x) f(x)) \quad x_i \text{ s are i.i.d.} \\ &= \frac{1}{n} \left(\int \frac{f(x)^2 p(x)^2}{q(x)} dx - \mu^2 \right) \end{aligned}$$

Let $\bar{q}(x) = \frac{|f(x)| p(x)}{c}$, where c is a constant making \bar{q} a distribution

i.e. $c = \int |f(x)| p(x)$. Then for any arbitrary distribution we have:

$$\begin{aligned} \sigma_{\bar{q}}^2 &= \int \frac{|f(x)|^2 p(x)^2}{\bar{q}(x)} dx - \mu^2 \\ &= c \int |f(x)| p(x) dx - \mu^2 \end{aligned}$$

$$\begin{aligned}
&= \left(\int |f(x)| p(x) dx \right)^2 - \mu^2 \\
&= \left(\int \frac{|f(x)| p(x)}{q(x)} q(x) dx \right)^2 - \mu^2 \\
&\leq \int \frac{f^2(x) p^2(x)}{q^2(x)} q(x) dx - \mu^2 \quad (\text{Cauchy-Schwarz}) \\
&= \sigma_q^2
\end{aligned}$$

If the function f is non-negative, then the variance would be zero, otherwise it has no certain value.

Now let $p(x)$ be a distribution concentrated in region A of the domain and $q(x)$ have almost zero probability there. Now let $f(x)$ have high values at A. Then the variance of $\hat{\mu}_{IS}$ can be arbitrarily large.

$$3) \hat{\mu}_{N-IS} = \frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i}$$

$$(ct) C = \sum_{i=1}^n w(x_i)$$

$$\Rightarrow E_q[\hat{\mu}_{N-IS}] = E_q \left[\frac{\frac{1}{n} \sum_{i=1}^n w(x_i) f(x_i)}{\frac{1}{n} \sum_{i=1}^n w(x_i)} \right] \neq \frac{E_p[f(X)]}{E_q[w(X)]} = E_p[f(X)] = \mu$$

$$E_q[w(X)] = \int \frac{p(x)}{q(x)} q(x) dx = 1$$

Thus $\hat{\mu}_{N-IS}$ is biased

But by the strong law of large numbers we have:

$$\lim_{n \rightarrow \infty} \hat{\mu}_{N-IS} = \lim_{n \rightarrow \infty} \frac{\frac{1}{n} \sum_{i=1}^n w(x_i) f(x_i)}{\frac{1}{n} \sum_{i=1}^n w(x_i)} = \frac{E_q[w(X)f(X)]}{E[w(X)]} = \mu$$

4) We have $\hat{\mu}_{N-IS} = \frac{\sum w_i f(x_i)}{\sum w_i}$. So it is a weighted average of $f(x)$. Let's denote $\hat{\mu}_{N-IS}$ by $\hat{\mu}$. Then $m \leq \hat{\mu} \leq M$. Now let's define $g(t)$ as:

$$\begin{aligned} g(t) &= E[(\hat{\mu} - t)^2] \\ \Rightarrow g'(t) &= \frac{d}{dt} E[(\hat{\mu} - t)^2] = E[\frac{d}{dt}(\hat{\mu} - t)^2] \\ &= E[-2\hat{\mu} + 2t] = 2E[t] - 2E[\hat{\mu}] \\ &= 2t - 2E[\hat{\mu}] = 0 \Rightarrow t = E[\hat{\mu}] \end{aligned}$$

And $g''(t) = 2 > 0 \Rightarrow g$ is convex $\Rightarrow E[\hat{\mu}] = \underset{t}{\operatorname{argmin}} g(t)$.

Therefore for any point e.g. $t = \frac{M+m}{2}$ we have

$$\text{var}(\hat{\mu}) = g(E[\hat{\mu}]) \leq g\left(\frac{M+m}{2}\right)$$

Now we have:

$$g\left(\frac{M+m}{2}\right) = E\left[(\hat{\mu} - \frac{M+m}{2})^2\right] = \frac{1}{4} E[((\hat{\mu} - M) + (\hat{\mu} - m))^2]$$

And $m \leq \hat{\mu} \leq M$, thus $\hat{\mu} - M \leq 0$ and $\hat{\mu} - m \geq 0$, so

$$((\hat{\mu} - M) + (\hat{\mu} - m))^2 \leq ((\hat{\mu} - m) - (\hat{\mu} - M))^2 = (M - m)^2$$

And since E is linear, then

$$\text{var}(\hat{\mu}_{N-IS}) \leq g\left(\frac{M+m}{2}\right) = \frac{1}{4} E[(\hat{\mu} - M) + (\hat{\mu} - m)]^2 \leq \frac{(M-m)^2}{4}$$

$$5) 1. \int_2^\infty e^{-\frac{x^2}{2}} dx \longleftrightarrow \int_{\frac{1}{2t^2}}^{\frac{1}{2}} e^{-\frac{1}{2t^2}} dt$$

Now we have to filter the samples through the interval $[0, \frac{1}{2}]$.

$$S_{\text{filtered}} = \{0.26, 0.41, 0.40, 0.10\}$$

$$H = \left(\frac{1}{2} - 0\right) \sum_{t \in S_{\text{filtered}}} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \times \frac{1}{t^2} e^{-\frac{1}{2t^2}} = 0.108$$

2. Now we apply IS. Take the sample from $g = N(3, 1)$, by just shifting the S sample by 3. Then the shifted filtered sample is :

$$S' = \{3.26, 2.73, 3.41, 3.4, 2.98, 3.1, 2.72\}$$

$$\Rightarrow H = \sum_{x \in S'} \frac{\frac{1}{|S'|}}{\frac{1}{\sqrt{2\pi(1)^2}} e^{-\frac{(x-3)^2}{2}}} = 0.03$$

The actual value of H is 0.05. Thus this estimator is better. It is better, since it gives more weight to the samples in the interval and reduces the probability of outlier samples.

(Q6)

A) $-D_{KL}(p||q) = - \int p(x) \log \frac{p(x)}{q(x)} dx$
 $= \int p(x) \log \frac{q(x)}{p(x)} dx$

$\leq \log \int p(x) \frac{q(x)}{p(x)} dx$

$= \log \int q(x) dx$

$= \log 1$

$\Rightarrow D_{KL}(p||q) \geq 0$

Again by the Jensen, the equality holds iff $\frac{p(x)}{q(x)}$ is constant, and
since $\int p(x) dx = \int q(x) dx = 1$, the ratio is 1. Thus $D_{KL}(p||q) = 0 \Leftrightarrow p = q$.

B) $D_{KL}(q(z) || P(z|x)) = E_q \left[\log \frac{q(z)}{P(z|x)} \right]$
 $= E_q \left[\log q(z) - \log P(z|x) \right]$
 $= E_q [\log q(z)] - E_q [\log P(z|x)]$
 $= E_q [\log q(z)] - E_q [\log P(z,x)] + E[\log P(x)]$
 $= \log P(x) - (E_q [\log P(x,z)] - E[\log q(z)])$

(I have problem with this. It has to be x , not X , i.e. a random variable)

C) We have $L(q) = E_q [\log P(x,z)] - E_q [\log q(z)]$

Let's denote $E_{q_j, j \neq i}$ by E_{-i} . We can write

$$q(z) = \prod_{i=1}^n q(z_i) \Rightarrow \log q(z) = \sum_{i=1}^n \log q(z_i)$$

And we have

$$\begin{aligned} L(q) &= \int q(z) \log p(X, z) dz - \int q(z) \log q(z) dz \\ &= \int \prod_{i=1}^n q(z_i) \log(X, z) dz - \int \prod_{i=1}^n q(z_i) \sum_{i=1}^n \log q(z_i) dz \\ &= \int \prod_{i=1}^n q(z_i) (\log(X, z) - \sum_{i=1}^n \log q(z_i)) dz \\ &= \int (q(z_i) \prod_{j \neq i} q(z_j)) (\log(X, z) - (\log q(z_i) + \sum_{j \neq i} \log q(z_j))) dz \\ &= \int_{z_i} q(z_i) \int_{Z_{j \neq i}} \prod_{j \neq i} q(z_j) ((\log p(X, z) - (\log q(z_i) + \sum_{j \neq i} \log q(z_j))) dz_j) dz_i \\ &= \int_{z_i} q(z_i) \int_{Z_{j \neq i}} \prod_{j \neq i} q(z_j) (\log p(X, z)) dz_j dz_i \\ &\quad - \int_{z_i} q(z_i) \int_{Z_{j \neq i}} \prod_{j \neq i} q(z_j) (\log q(z_i) + \sum_{j \neq i} \log(q_j)) dz_j dz_i \\ &= \int_{z_i} q(z_i) E_{-i} [\log p(X, z)] dz_i - \int_{z_i} q(z_i) \log q(z_i) dz_i + C \end{aligned}$$

where $C = \int_{Z_{j \neq i}} \prod_{j \neq i} q(z_j) \sum_{j \neq i} \log(q_j) dz_{j \neq i}$. The Lagrangian of this optimization

problem is:

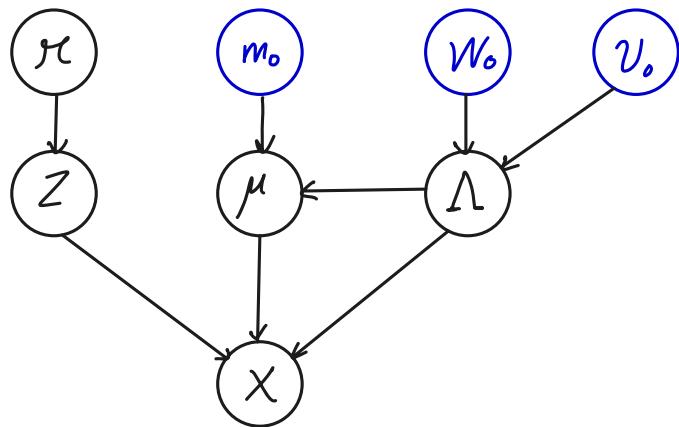
$$L(q) - \sum_{j=1}^n \lambda_j \int_{z_j} q(z_j) dz_j$$

$$\frac{\delta L(q)}{\delta q(z_i)} = \frac{\partial}{\partial q(z_i)} [q(z_i) [E_{-i} [\log p(X, z)] - \log q(z_i)] - \lambda_i q(z_i)]$$

$$= E_{z_i} [\log P(X, Z)] - \log q(z_i) - 1 - \lambda_i = 0$$

$$\Rightarrow \log q(z_i) = E_{z_i} [\log P(X, Z)] + \text{const} \quad \forall 1 \leq i \leq n$$

D) The diagram of this GMM can be represented as a PGM like :



In this setting, the posterior we want to estimate is q , and the actual joint distribution is $P(X, Z, \pi, \mu, \Lambda)$. Therefore the ELBO of interest is

$$\begin{aligned} L(q) &= E_q \left[\log \frac{P(X, Z, \pi, \mu, \Lambda)}{q(Z, \pi, \mu, \Lambda)} \right] \\ &= E_q \left[\log P(X, Z, \pi, \mu, \Lambda) - \log q(Z, \pi, \mu, \Lambda) \right] \end{aligned}$$

Now we have to compute $P(X, Z, \pi, \mu, \Lambda)$. By the represented PGM, which is a Bayes Net, we have

$$P(X, Z, \pi, \mu, \Lambda) = P(X|Z, \mu, \Lambda) P(Z|\pi) P(\pi) \\ P(\mu, \Lambda)$$

Thus we have

$$L(q) = E_q[\log P(X|Z, \mu, \Lambda)] + E_q[\log P(Z|\pi)] + E_q[\log P(\pi)] \\ + E_q[\log P(\mu, \Lambda)] - E_q[\log q(z)] - E_q[\log q(x, \mu, \Lambda)] \\ (\text{we have } q(z, \pi, \mu, \Lambda) = q(z) q(x, \mu, \Lambda))$$

by the previous part we have,

$$\log q^*(z) = E_{q(\pi, \mu, \Lambda)}[\log P(X, Z, \pi, \mu, \Lambda)] + \text{const} \\ = E_{q(\pi, \mu, \Lambda)}[\log(P(X|Z, \mu, \Lambda) P(Z|\pi) P(\pi) P(\mu, \Lambda))] + \text{const} \\ = E_{q(x, \mu, \Lambda)}[\log P(X|Z, \mu, \Lambda)] + E_{q(\pi, \mu, \Lambda)}[\log P(Z|\pi)] \\ + \underbrace{E_{q(\pi, \mu, \Lambda)}[\log P(\pi)] + E_{q(\pi, \mu, \Lambda)}[\log P(\mu, \Lambda)]}_{\text{const w.r.t } z} + \text{const} \\ = E_{q(x, \mu, \Lambda)}[\log P(X|Z, \mu, \Lambda) + \log P(Z|\pi)] + \text{const}$$

$$E_{q(\pi)}[\log P(Z|\pi)] = \sum_{n=1}^N \sum_{k=1}^K z_{nk} E_{q(\pi)}[\log \pi_k]$$

$$E_{q(\mu, \Lambda)}[\log P(X|Z, \mu, \Lambda)] = \sum_{n=1}^N \sum_{k=1}^K z_{nk} E_{q(\mu, \Lambda)}[\log N(x_n | \mu_k, \Lambda^{-1})]$$

Thus :

$$\log q^*(z) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \gamma_{nk} + \text{const}$$

where

$$\begin{aligned} \gamma_{nk} &= E_{q(\pi)} [\log x_k] + \frac{1}{2} E_{q(\mu, \Lambda)} [\log |\Lambda_k|] - \frac{D}{2} \log(2\pi) \\ &\quad - \frac{1}{2} E_{q(\mu, \Lambda)} [(x_n - \mu_k)^T \Lambda_k (x_k - \mu_k)] \end{aligned}$$

$$\begin{aligned} (\text{Because } E[\log N(x_n | \mu_k, \Lambda_k^{-1})] &= E[-\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\Lambda_k| \\ &\quad - \frac{1}{2} (x_n - \mu_k)^T \Lambda_k (x_n - \mu_k)] \\ &= -\frac{D}{2} \log(2\pi) + \frac{1}{2} E[\log |\Lambda_k|] \\ &\quad - \frac{1}{2} E[(x_n - \mu_k)^T \Lambda_k (x_n - \mu_k)] \end{aligned}$$

Therefore $q^*(z) = C \prod_{n=1}^N \prod_{k=1}^K e^{z_{nk} \gamma_{nk}}$ where C is a normalizing constant.

Now we need to compute $q^*(\pi)$, $q^*(\mu)$ and $q^*(\Lambda)$. Similar as before,

$$\begin{aligned} \log q^*(\pi) &= E_{q(z, \mu, \Lambda)} [\log P(X, Z, \pi, \mu, \Lambda)] + \text{const} \\ &= E_{q(Z)} [\log P(Z|\mu) + \log P(\pi)] + \text{const} \quad (\text{Based on the Bayes Net}) \\ &= \log P(\pi) + E_{q(Z)} [\log P(Z|\pi)] + \text{const} \\ &= (\alpha_0 - 1) \log \sum_{k=1}^K \pi_k + C \sum_{n=1}^N \sum_{k=1}^K e^{\gamma_{nk}} \log x_k + \text{const} \end{aligned}$$

And

$$\begin{aligned}
 \log q^*(\mu, \Lambda) &= E_{q(Z, \pi)} [\log P(X, Z, \pi, \mu, \Lambda)] + \text{const} \\
 &= E_{q(Z, \pi)} [\log P(X|Z, \mu, \Lambda) + \log P(\mu, \Lambda)] + \text{const} \\
 &= \log P(\mu, \Lambda) + E_{q(Z)} [\log P(X|Z, \mu, \Lambda)] + \text{const}
 \end{aligned}$$

To normalize the proposed expression for each latent variable to form a distribution, since each vector is K-of-1, the sum of z_{nk} s over all values of X is 1. We have

$$\begin{aligned}
 q^*(Z) &= \prod_{n=1}^N \prod_{k=1}^K r_{nk}^{z_{nk}} \quad \text{where } r_{nk} = \frac{e^{\lambda_{nk}}}{\sum_{i=1}^K e^{\lambda_{ni}}} = E[z_{nk}] \\
 q^*(\pi) &= (\alpha_0 - 1) \log \sum_{k=1}^K \pi_k + \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log \pi_k + \text{const}
 \end{aligned}$$

And

$$\begin{aligned}
 \log q^*(\mu, \Lambda) &= \sum_{k=1}^K \log N(\mu_k | m_0, (\beta_0 \Lambda_k)^{-1}) \mathcal{W}(\Lambda_k | W_0, V_0) \\
 &\quad + \sum_{n=1}^N \sum_{k=1}^K E_{q(Z)} [z_{nk} (\frac{1}{2} \log |\Lambda_k| - \frac{D}{2} \log 2\pi - (\mu_n - \mu_k)^T \Lambda_k (K_n - \mu_k))]
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow \log q^*(\mu_k, \Lambda_k) &= \log N(\mu_k | m_0, (\beta_0 \Lambda_k)^{-1}) + \log \mathcal{W}(\Lambda_k | W_0, V_0) \\
 &\quad + \sum_{n=1}^N (\frac{1}{2} \log |\Lambda_k| - \frac{D}{2} \log 2\pi - \frac{D}{2}) E[z_{nk}]
 \end{aligned}$$

$$= \log(N(\mu_k | m_k, (\beta_k \Lambda_k)^{-1})) W(\Lambda_k | W_k, U_k)$$

where

$$\beta_k = \beta_0 + N_k , \quad N_k = \sum_{n=1}^N r_{nk}$$

$$m_k = \frac{\beta_0}{\beta_k} m_0 + \frac{N_k}{\beta_k} \bar{x}_k , \quad \bar{x}_k = \frac{\sum_{n=1}^N r_{nk} x_n}{N_k}$$

$$W_k^{-1} = W_0^{-1} + N_k S_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{x}_k - m_0)(\bar{x}_k - m_0)^T$$

$$U_k = V_0 + N_k$$

The analogy of EM algorithm to the process above is

- 1) The E step : computes the value of r_{nk} using the current values.
- 2) The M step : uses the new r_{nk} values to update the parameters.

I have used the Bishop's Pattern Recognition book for the update rules And deriving $q^*(\mu_k, \Lambda_k)$