

(١) سطح

(الف)

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

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

$$I(X;Y|Z) - I(X;Y) = H(X|Z) - H(X|Y;Z) - H(x) + H(X|Y)$$

١: $H(X|Z)=0$; if X is completely determined by Z . Thus if $x=t=z$

$$I(X;Y|Z) - I(X;Y) = H(X|Z) - H(X|Y;Z) - H(x) + H(X|Y) < 0 \quad \text{for all } x's$$

٢: X, Z, Y are independent and X is conditionally independent of Z given Y :

$$H(X|Z) = H(x)$$

$$H(X|Y) = H(x)$$

$$H(X|Y;Z) = H(X|Y) = H(x)$$

$$\Rightarrow I(X;Y|Z) = I(X;Y)$$

So let $Z=X+Y$ and X and Y independent. we have:

$$I(X;Y|Z) - I(X;Y) = H(X|Z) - H(X|Y;Z) - H(x) + H(X|Y) = H(Y|Z) > 0$$

So $I(X;Y|Z) > I(X;Y)$ iff $X \perp Y$ and $Z = X+Y$ and $X \perp Z$ → from two given conditions, this always satisfies.

Ex: Let $X, Y \sim \text{Ber}(\frac{1}{2})$: $H(X|Z) = \sum_3 H(X|Z=3) \cdot \sum_3 p(3) \sum_x p(x) \log p(X=3|Z=3)$

$$= -\left(p(3=1) \sum_{x \in \{0,1\}} p(x) \log p(X=x|Z=3) + p(3=2) \sum_{x \in \{0,1\}} p(x) \log p(X=x|Z=3)\right. \\ \left. + p(3=3) \sum_{x \in \{0,1\}} p(x) \log p(X=x|Z=3)\right)$$

$$= -\left(\frac{1}{4} \cdot 1 \cdot \log(1) + \frac{1}{2} \cdot 1 \cdot \log(\frac{1}{2}) + \frac{1}{4} \cdot 1 \cdot \log(1)\right) = -\frac{1}{2} \cdot 1 \cdot \log(\frac{1}{2}) = \frac{1}{2} \log(2) = \frac{1}{2}$$

$$1. H(X,Y,Z) - H(X,Y) = H(Z|Y,X) + H(Y|X) + H(x) - H(Y|X) - H(x) = H(Z|Y,X) \quad (\rightarrow)$$

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

$$\text{Define } A = Z|X \Rightarrow H(X,Y,Z) - H(X,Y) - H(X|Z) + H(x) = H(Z|X,Y) - H(Z|X) = H(A|Y) - H(A)$$

$$= -I(Y;A) > 0$$

$$2. I(X; Z|Y) = I(Z; Y|X) = I(X; Z) - I(X; Z|Y)$$

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

$$\bullet H(Y|X) = H(X|Y) - H(X) + H(Y) \rightsquigarrow H(\frac{z}{x} | \frac{y}{x}) = H(X|Z|Y) + H(Z|Y) - H(X)$$

$$\rightsquigarrow H(Z|Y) - H(X) - H(Z|Y) + H(X|Z) = H(X|Z) - H(X) \leq 0 \quad H(X|Z) = H(X) \text{ iff } X, Z \text{ are independent.}$$

$$11 \quad P(Z=A|Y=y) = (1 + \exp(-\frac{y-5}{2}))^{-1} \rightsquigarrow \log_2 P(Z=A|Y=y) = -\log_2 (1 + \exp(-\frac{y-5}{2})) = x \quad (c)$$

$$1 - P(Z=A|Y=y) = 1 - (1 + \exp(-\frac{y-5}{2}))^{-1} \rightsquigarrow \log_2 (1 - P(Z=A|Y=y)) = \log_2 \left(\frac{\exp(-\frac{y-5}{2})}{1 + \exp(-\frac{y-5}{2})} \right) = \frac{y-5}{2} \log_2 e + \log_2 (1 + \exp(-\frac{y-5}{2}))$$

$$P(Z) = E_y [Z, Y=y] = \sum_y P(Z|Y=y) P(Y=y) = \frac{1}{4} \sum_y P(Z|Y=y) = \frac{1}{4}, \quad H(Z) = 1$$

$$H(Z|Y) = \sum_y P(Y=y) \sum_z P(Z|Y=y) \log(P(Z|Y=y)) = \frac{1}{4} \left(x \log_2 x + \frac{1-x}{2} \log_2 \frac{1-x}{2} + \frac{1-x}{3} \log_2 \frac{1-x}{3} + \frac{1-x}{6} \log_2 \frac{1-x}{6} + \frac{x}{2} \log_2 \frac{x}{2} \right)$$

$$= \frac{1}{4} \left(x \log_2 x + \frac{1-x}{2} \log_2 \frac{1-x}{2} + \frac{1-x}{3} \log_2 \frac{1-x}{3} - \frac{1-x}{3} \log_2 \frac{1-x}{3} + \frac{1-x}{6} \log_2 \frac{1-x}{6} - \frac{1-x}{6} \log_2 \frac{1-x}{6} \right)$$

$$= \frac{1}{4} (x \log_2 x + \log_2 (1-x) - (1-x) - \frac{1-x}{6} \log_2 (3)) = \frac{1}{4} (x \log_2 x - \frac{y-5}{2} \log_2 e + \log_2 (e) + x - 1 + \frac{1-x}{6} \log_2 (3))$$

2) $I(X; Y)$ represents the mutual information between RV X and Y . Thus if $I(X; Y)$ is big/small, it means with observing X , we can have more/less information about Y which represent the value or utility we get by choosing route Z .

$$D_{KL}(P(X, Y, Z) || P(X) P(Y) P(Z)) = \int_x \int_y \int_z P(x, y, z) \log \left(\frac{P(x, y, z)}{P(x) P(y) P(z)} \right) dx dy dz$$

$$= \int_x \int_y \int_z P(x, y, z) \log P(x, y, z) dx dy dz - \int_x \int_y \int_z P(x, y, z) \log P(x) dx dy dz - \int_x \int_y \int_z P(x, y, z) \log P(y) dx dy dz$$

$$= -H(X, Y, Z) - \int_x \int_y \int_z P(x, y, z) \log P(x) dx dy dz - \int_y \int_z P(x, y, z) \log P(y) dy dz$$

$$- \int_z \int_x P(x, z) \int_y P(x, y, z) dy dz dx$$

$$\int_x \int_y \int_z P(x, y, z) dx dy dz = P(X)$$

$$H(x, y, z) = \int p(x) \log p(x) dx + \int p(y) \log p(y) dy + \int p(z) \log p(z) dz$$

$$= H(x) + H(y) + H(z)$$

$D_{KL}(p(x, y, z) || p(x)p(y)p(z)) = 0$ iff $p(x, y, z) = p(x)p(y)p(z)$ $\Leftrightarrow x, y, z$ are independent.

or

$$H(x, y, z) = H(Z|x, y) + H(y|x) + H(x) = H(z) + H(y) + H(x) \quad \text{iff } z \text{ is independent from } y \text{ and } \begin{cases} y \perp\!\!\!\perp z \\ y \text{ is independent from } x \end{cases}$$

(\rightarrow)

$$\begin{aligned} & I(X_1; X_4) + I(X_2; X_3) - I(X_1; X_3) - I(X_2; X_4) \\ &= H(X_1) - H(X_1|X_4) + H(X_2) - H(X_2|X_3) - H(X_1) + H(X_1|X_3) - H(X_2) + H(X_2|X_4) \\ &= H(X_1|X_3) + H(X_2|X_4) - H(X_3|X_3) - H(X_4|X_1) \end{aligned}$$

from bayes rule we have $H(X_1, X_2|X_3) = H(X_1|X_3) + H(X_2|X_3, X_1)$

thus we have:

$$\begin{aligned} & H(X_1, X_2|X_3) - H(X_2|X_1, X_3) + H(X_2, X_1|X_4) - H(X_1|X_2, X_4) + H(X_1|X_2, X_3) - H(X_2, X_1|X_3) + H(X_2|X_1, X_4) \\ &= H(X_1|X_2, X_3) + H(X_2|X_1, X_4) - H(X_1|X_2, X_4) - H(X_2|X_1, X_3) \end{aligned}$$

we know that in markov process, $H(X_i|X_{i+1}, \dots) = 0$. Thus:

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

also: $I(X_2; X_3|X_1, X_4) = 0$ iff $X_2 \perp\!\!\!\perp X_3 | X_1, X_4$

$f(x)$ is convex iff for $0 < \lambda \leq 1$ and $u_1, u_2 \in \mathbb{R}$, $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$

or by first order condition: $f(u_2) - f(u_1) + \nabla f(u_1)(u_2 - u_1)$

if u_1 is a critical point for $f(u)$, then $\nabla f(u_1) = 0$. Thus: $f(u_2) \geq f(u_1)$ and u_1 is a global minimum of f ,

in reverse, if $f(u_1) = 0$ if u_1 is a global minima, because $f(x_2) = f(x_1) + \nabla f(x_1)(f(x_2) - f(x_1)) + O(n^2)$

then if $\nabla f(x_1) \neq 0$, let $u_2 = u_1 + \nabla f(x_1) \cdot u_1$ where $u \ll 1$, we have $f(x_2) < f(x_1)$ and this is a conflict

C

$$L(u, y, z, \lambda_1, \lambda_2) = f(u, y, z) - \lambda_1 g_1(u, y, z) - \lambda_2 g_2(u, y, z) = u^2 + y^2 + z^2 - \lambda_1(z^2 - u^2 - y^2) - \lambda_2(z - u - y - 1)$$

$$\begin{cases} \frac{\partial L}{\partial u} = 2u + 2\lambda_1 u + \lambda_2 = 0 \\ \frac{\partial L}{\partial y} = 2y + 2\lambda_1 y + \lambda_2 = 0 \end{cases} \xrightarrow{\# \cdot \frac{\partial L}{\partial y}} 2(u+y) + 2\lambda_1(u+y) = 0 \Rightarrow \begin{cases} u=y \\ 1+\lambda_1=0 \quad \lambda_1 \neq 0 \end{cases}$$

$$\begin{cases} \frac{\partial L}{\partial z} = 2z - 2\lambda_1 z - \lambda_2 = 0 \\ \frac{\partial L}{\partial \lambda_1} = z^2 - u^2 - y^2 = 0 \Rightarrow z^2 - 2u^2 = 0 \\ \frac{\partial L}{\partial \lambda_2} = z - u - y - 1 = 0 \Rightarrow z - 2u - 1 = 0 \end{cases} \Rightarrow z^2 - 2u^2 = 2u^2 \Rightarrow 4u^2 + 4u + 1 = 2u^2 \Rightarrow 2u^2 + 4u + 1 = 0 \Rightarrow u = -1 \pm \frac{\sqrt{2}}{2} \Rightarrow z = \pm \sqrt{2} - 1$$

Critical points on $y=0$ and $y \geq 0$ for f : $f(-1 - \frac{\sqrt{2}}{2}, -1 - \frac{\sqrt{2}}{2}, \sqrt{2} - 1) = 2(\sqrt{2} - 1)^2 = 2(3 - 2\sqrt{2}) = 6 - 4\sqrt{2}$ minimum

$$f(-1 + \frac{\sqrt{2}}{2}, -1 + \frac{\sqrt{2}}{2}, \sqrt{2} - 1) = 2(-\sqrt{2} - 1)^2 = 2(3 + 2\sqrt{2}) = 6 + 4\sqrt{2}$$
 maximum

C

$$L(\mu, \lambda) = \inf_u \{f(u) + \mu^T g(u) + \lambda^T h(u)\}$$

let a^* be the feasible point; then we have: $h(a^*) = 0$, $\mu^T g(a^*) \leq 0 \Rightarrow L(\mu, \lambda) = \inf_u \{f(u) + \mu^T g(u) + \lambda^T h(u)\} \leq f(a^*) + \mu^T g(a^*) + \lambda^T h(a^*) \leq f(a^*)$

considering the point that a^* is any feasible point, then if a_0 is optimal solution for $\min_u f(u)$, a_0 is also a feasible point. Then we have: $L(\mu, \lambda) = \inf_u \{f(u) + \mu^T g(u) + \lambda^T h(u)\} \leq f(a_0) = \min_u f(u)$

$$\text{let } \min_{\substack{\text{s.t.} \\ a_i \geq 0}} e^{-a_i} \Rightarrow L(\lambda) = \inf_{a_i \geq 0} \left\{ e^{-a_i} + \lambda \frac{a_i^2}{2} \right\} = \begin{cases} 0 & \lambda > 0 \\ -\infty & \lambda \leq 0 \end{cases} \Rightarrow d^* = \max \{0 : \lambda > 0\}$$

however: $e^{-d^*} \cdot e^{2d^*} = 1$ Thus $\max L(\lambda) = 0 < \min_u f(u) = 1$

D

$$\text{minimize } -\sum_i \log(a_i + x_i) \quad (\text{Also known as water filling problem})$$

s.t. $x_i \geq 0$ $\sum_i x_i = 1$

first, we clarify our feasible set. x_i is optimal if $x_i \geq 0$ and $\sum x_i = 1$. If we form the lagrangian/dual problem, we have:

$$L(p, \lambda) = -\sum_{i=1}^n \log(x_i + \alpha_i) + \lambda(\mathbf{1}^\top \mathbf{x} - 1) + \mu^\top \mathbf{x} = -\sum_{i=1}^n \log(x_i + \alpha_i) + \lambda(\sum_{i=1}^n x_i - 1) + \sum_{i=1}^n \lambda_i x_i$$

$$\frac{\partial L}{\partial x_i} = -\frac{1}{x_i + \alpha_i} + \lambda_i + \mu_i = 0 \Rightarrow \mu_i = \frac{1}{x_i + \alpha_i} - \lambda_i$$

From KKT conditions, we know that $\mu_i^\top \mathbf{1} = 0$. Thus: $\mu_i (\frac{1}{x_i + \alpha_i} - \lambda_i) = 0$. Also from the fact that $\mu_i \geq 0$, we have: $\lambda_i \geq \frac{1}{x_i + \alpha_i}$. Let's break it into 2 conditions:

$$\alpha_i > 0 \Rightarrow \lambda_i = \frac{1}{x_i + \alpha_i} \Rightarrow x_i = \frac{1}{\lambda_i} - \alpha_i \quad \text{or} \quad \alpha_i = 0 \Rightarrow x_i^* = \max\{0, \frac{1}{\lambda_i} - \alpha_i\} \quad \text{such that } \sum_i x_i^* = 1$$

first we form the lagrangian:

$$\begin{aligned} L(p, \lambda, \mu, \sigma) &= f(x) - \lambda^\top g(x) = -\int p(x) \log p(x) - \lambda_0 (\int p(x) dx - 1) - \lambda_1 (\int (x - \mu)^2 p(x) dx - \sigma^2) \\ &= -\int p(x) [\log p(x) - \lambda_0 - \lambda_1 (x - \mu)^2] dx - \lambda_0 - \lambda_1 \sigma^2 \end{aligned}$$

Using Calculus of Variations, we can find critical points of $L(p, \lambda, \mu, \sigma^2)$ by solving following system of equations:

$$\frac{d}{du} \frac{\partial L}{\partial p} - \frac{\partial L}{\partial p} = 0 \quad ((\frac{\partial L}{\partial p} = 0) \Rightarrow \frac{\partial L}{\partial p} = -\log p(x) - 1 - \lambda_0 - \lambda_1 (x - \mu)^2 = 0 \Rightarrow p(x) = e^{-1-\lambda_0-\lambda_1(x-\mu)^2})$$

now, since $p(x)$ is a probability distribution, we have: $\int p(x) dx = 1 \Rightarrow \int e^{-1-\lambda_0-\lambda_1(x-\mu)^2} dx = 1 \Rightarrow e^{-\lambda_0} \int e^{-\lambda_1(x-\mu)^2} dx = 1$

$$\text{Let } t = \sqrt{\lambda_1}(x - \mu) \Rightarrow dt = \sqrt{\lambda_1} dx \Rightarrow \int e^{-\lambda_1(x-\mu)^2} dx = \frac{1}{\sqrt{\lambda_1}} \int e^{-t^2} dt = \frac{\sqrt{\pi}}{\sqrt{\lambda_1}} \Rightarrow e^{-\lambda_0} \cdot \frac{\sqrt{\pi}}{\sqrt{\lambda_1}} \Rightarrow \lambda_0 = -1 + \frac{1}{2} \log \pi - \frac{1}{2} \log \lambda_1$$

from second constraint, we have:

$$\int (x - \mu)^2 p(x) dx = \int (x - \mu)^2 e^{-1-\lambda_0-\lambda_1(x-\mu)^2} dx - \sigma^2$$

$$\text{again, let } t = \sqrt{\lambda_1}(x - \mu). \text{ we have: } \sigma^2 = \int \frac{t^2}{\lambda_1} e^{-1-\lambda_0-t^2} \frac{dt}{\sqrt{\lambda_1}} = \frac{e^{-\lambda_0}}{\lambda_1 \sqrt{\lambda_1}} \int t^2 e^{-t^2} dt = \frac{e^{-\lambda_0}}{\lambda_1 \sqrt{\lambda_1}} \cdot \frac{\sqrt{\pi}}{2} (e^{-1-\lambda_0} \cdot \frac{\sqrt{\lambda_1}}{\sqrt{\pi}}) \Rightarrow \sigma^2 = \frac{\sqrt{\lambda_1} \times \sqrt{\pi}}{2 \lambda_1 \sqrt{\lambda_1} \times \sqrt{\pi}} \cdot \frac{1}{2} \Rightarrow \lambda_1 = \frac{1}{2\sigma^2}$$

$$\text{Thus we have: } p(x) = e^{-\lambda_0} e^{-\lambda_1(x-\mu)^2} \cdot \frac{\sqrt{\lambda_1}}{\sqrt{\pi}} e^{-\frac{\lambda_1(x-\mu)^2}{2\sigma^2}} \sim p(x) \sim N(\mu, \sigma^2)$$

The chain is A-periodic and irreducible, thus we have $\pi = \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_m \end{bmatrix}$

$$\lim_{n \rightarrow \infty} P^{(n)} = \begin{bmatrix} \cdots & \pi & \cdots \\ \cdots & \pi & \cdots \\ \vdots & & \vdots \\ \cdots & \pi & \cdots \end{bmatrix}$$

using the probability of transition $\begin{cases} p_{ii} = \frac{1}{m} |N(i)| \\ p_{ij} = \begin{cases} \frac{1}{m} & j \in N(i) \\ 0 & j \notin N(i) \end{cases} \end{cases}$

we know that $M \geq \arg \max_{x \in \Omega} |N(x)|$; Thus $p_{ij} \geq p_{ji} = \frac{1}{m}$. Thus P is symmetric.

using the fact that for a node i , we have $\sum_{j \in \Omega} p_{ij} = |N(i)| \times \frac{1}{m} + 1 - \frac{|N(i)|}{m} = 1$: column and rows of P sum to be 1.

we know that if P 's columns sum to 1, then the stationary distribution π defines as $\pi_i = \frac{1}{m}$.

Main problem in this counting is that we are counting one sample multiple times, i.e., the sampling is uniform; however we can take a lot more information via a better sampling. Also, there is a 2^n multiplier when there is no need to count the possible outcomes of X in the numerator when we have uniform sampling; i.e. a better way could be initializing $X = x_1, \dots, x_n$ 2^n times non-uniformly and then check the constraints; however given approach in C is better and uses almost all of driven data.

1. let $x_1 = (0, 0, \dots, 0)$. then we have:

$$P_{11} = \frac{\sum_{i=1}^n I(a_i \geq b)}{n}, P_{12} = 1 - P_{11}, P_{22} = \frac{1}{n} + \frac{\sum_{i=1}^n I(a_i < c_i \leq b)}{n-1} \Rightarrow P_{22} = 1 - P_{12}$$

from part 1, we know that if a markov chain is A-periodic and irreducible, it has a stationary distribution π over nodes.

Given Markov chain is obviously irreducible, and it's A-periodic (since biggest period has length 1). Thus the Markov-Process has a stationary distribution π which tells us after infinitely step Random walk, which nodes are most possible to be found our agent at.

2. we just have to check when the MC stops, i.e., after $T \rightarrow \infty$ passes, MC stops at X_Z , where Z is all possible X 's where $\sum a_i x_{i,t} b_i$ constrain is satisfied. Thus Z is all possible valid outcomes.

$$p(s|same) = \alpha_{ss} e^{s\lambda_s}, p(s|different) = \alpha_{ds} e^{-s\lambda_d}, s \in [0, 1]$$

The probability distributions seem rational, because if two pictures are same, we expect a higher score proportional to s . i.e., $dP(s|same) = \lambda_s \alpha_{ss} s e^{s\lambda_s} > 0$ and vice versa for different pictures: $dP(s|different) = -\lambda_d \alpha_{ds} s e^{-s\lambda_d} < 0$
 α_{ss} and α_{ds} are normalizing constants to make $p(s|same)$ and $p(s|different)$ a probability distribution over some interval s . for example, for $s \in [0, 1]$; we have:

$$\int p(s|same) ds = \int \alpha_{ss} e^{s\lambda_s} ds = \frac{\alpha_{ss}}{\lambda_s} e^{\lambda_s s} \Big|_0^1 = \alpha_{ss} \frac{e^{\lambda_s} - 1}{\lambda_s} = 1 \Rightarrow \alpha_{ss} = \frac{\lambda_s}{e^{\lambda_s} - 1}$$

$$\int p(s|different) ds = \int \alpha_{ds} e^{-s\lambda_d} ds = \frac{\alpha_{ds}}{\lambda_d} e^{-\lambda_d s} \Big|_0^1 = \alpha_{ds} \frac{1 - e^{-\lambda_d}}{\lambda_d} = 1 \Rightarrow \alpha_{ds} = \frac{\lambda_d}{e^{-\lambda_d} - 1}$$

Let A_j is the picture j is for person j . Then we have:

$$p(A_j|s_j) = \frac{p(s_j|A_j) p(A_j)}{\sum p(s_i|A_i) p(A_i)} = \frac{p(s_j|A_j) p(A_j)}{\sum_i p(s_i|A_i) p(A_i)}$$

Let's assume the condition of truly assigning pictures to someone is uniform and uncorrelated, i.e. $\forall i, j : p(A_i) = p(A_j)$

$$\text{thus: } p(A_j|s_j) = \frac{p(s_j|A_j) p(A_j)}{\sum_i p(s_i|A_i) p(A_i)} = \frac{p(s_j|A_j)}{\sum_i p(s_i|A_i)} = \frac{\alpha_{ss} e^{s_j \lambda_s}}{\sum_i \alpha_{ss} e^{s_i \lambda_s}} = \frac{(n-1)\alpha_{ss} e^{-s_j \lambda_s} + \alpha_{ss} e^{s_j \lambda_s}}{e^{\lambda_s} - 1} = \frac{(n-1)\lambda_d e^{-s_j \lambda_d} + \lambda_d e^{s_j \lambda_d}}{e^{\lambda_d} - 1}$$

$$p(A|\max) = \frac{p(\max|A_j) p(A_j)}{p(\max|A_j) p(A_j) + p(\max|A^c) p(A^c)} = \frac{p(\max|A)}{p(\max|A) + p(\max|A^c)}$$

So; if we have a rating policy; we should build it based on two different condition. one, takes the hypothesis that given picture for everyone is right; i.e. $p(\max|A)$. And the other one takes the hypothesis that our indicator thinks that picture is not mapped quite correctly. Then, if we take the assumption that $p(A) = p(A^c)$, we will have a short formula like given above. Secondly, if we consider the case where A_i 's are independent; we will have: $p(A|\max) = \frac{p(\max|A)}{(n-1)p(\max|A^c), p(\max|A)}$ which almost a common sense.

$$\text{Then we have: } P_{s_{\max}}(\max) = p(A_j, i_j \leq \max) p(\max) = \left(\int_0^{\max} p(x) dx \right)^n p(\max)$$

$$\Rightarrow \text{from part b: } P_{s_{\max}} = \int_0^{\max} \frac{1}{n} p(x|A_j) + p(x|A^c) \frac{n-1}{n} dx$$

$$\text{also note that } p(A) = \int_0^1 p(A|\max) p(\max) d\max = \int_{-\infty}^1 \frac{\alpha_{ss} e^{s\lambda_s}}{-\alpha_{ss} e^{s\lambda_s} + (n-1)\lambda_d e^{-s\lambda_d}} (e^{s\lambda_s} - e^{-(n-1)\lambda_d e^{-s\lambda_d}}) \left(\frac{\alpha_{ss}}{\lambda_s} e^{s\lambda_s} - \frac{\lambda_d}{e^{-\lambda_d} - 1} \right) ds$$

$$= \frac{1}{n^2} \int_0^1 \alpha_{ss} e^{s\lambda_s} \left(\frac{\alpha_{ss}}{\lambda_s} (e^{s\lambda_s} - 1) + (n-1) \frac{\lambda_d}{e^{-\lambda_d} - 1} (e^{-s\lambda_d} - 1) \right)^{n-1} ds$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(X_{1:n} | \theta) = \arg \max_{\theta} L(\theta; X) \propto P(X_i; \mu, \sigma^2) = \frac{1}{6\sqrt{2\pi}} \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right]$$

$$L(\mu, \sigma) = \prod_{i=1}^n P(X_i | \mu, \sigma^2) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right]$$

$$\log L(\mu, \sigma) = -n \log \sigma - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2$$

$$\begin{aligned} \frac{\partial \log(L(\mu, \sigma))}{\partial \mu} &= 0 \rightarrow \frac{-2 \sum (x_i - \mu)}{2\sigma^2} = 0 \Rightarrow \hat{\mu} = \frac{\sum x_i}{n} \\ \frac{\partial \log(L(\mu, \sigma))}{\partial \sigma} &= 0 \rightarrow -\frac{n}{\sigma} + \frac{\sum (x_i - \mu)^2}{\sigma^3} = 0 \Rightarrow \hat{\sigma}^2 = \frac{\sum (x_i - \bar{x})^2}{n} = \frac{1}{n} \sum x_i^2 - \frac{1}{n^2} \sum_i \sum_j x_i x_j \end{aligned}$$

$$\text{Let } \delta_i = \mu - x_i \text{ then } \hat{\sigma}^2 = \frac{1}{n} \sum_i (\mu - \delta_i)^2 - \frac{1}{n^2} \sum_i \sum_j (\mu - \delta_i)(\mu - \delta_j)$$

$$\begin{aligned} \mathbb{E}[\hat{\sigma}^2] &= \frac{1}{n} \sum_i \mathbb{E}[(\mu - \delta_i)^2] - \frac{1}{n^2} \sum_i \sum_j \mathbb{E}(\mu - \delta_i)(\mu - \delta_j) \\ &= \frac{1}{n} \sum_i \left(\mathbb{E}[\mu^2] - 2\mu \mathbb{E}[\delta_i] + \mathbb{E}[\delta_i^2] \right) - \frac{1}{n^2} \sum_i \sum_j \left(\mathbb{E}[\mu^2] - \mu \mathbb{E}[\delta_i + \delta_j] + \mathbb{E}[\delta_i \delta_j] \right) \\ &= \frac{1}{n} (n\mu^2 + n\sigma^2) - \frac{1}{n^2} (n^2\mu - n\sigma^2) = \mu^2 + \sigma^2 - \mu^2 - \frac{1}{n} \sigma^2 = \frac{n-1}{n} \sigma^2 \\ &\text{biased!} \quad \checkmark \end{aligned}$$

$\theta \mapsto \hat{\theta}(\theta) : \theta \mapsto \arg \max_{\theta} P(X | \theta = \hat{\theta}) = \arg \max_{\theta} P(X | \hat{L}(\theta))$, where $\hat{L}(\theta) = L(\hat{\theta})$. Thus we have:

$$\begin{cases} \hat{\mu}_s = \hat{L}(\hat{\theta}) = \frac{1}{1 + \hat{\theta}^2} = \frac{1}{1 + (\frac{\sum x_i}{n})^2} \\ \hat{\theta} = \sqrt{\frac{\sum (x_i - \hat{\mu}_s)^2}{n}} \end{cases}$$

$$\begin{aligned} 1. \hat{\mu}_{ss} &= \frac{1}{n} \sum_{i=1}^n f_{(x_i)} \frac{P(x_i)}{q_{(x_i)}} \rightarrow \mathbb{E}_q[\hat{\mu}_{ss}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{x \sim q} \left[f_{(x_i)} \frac{P(x_i)}{q_{(x_i)}} \right] = \frac{1}{n} \sum_{i=1}^n \int_{x \sim q} f_{(x_i)} \frac{P(x_i)}{q_{(x_i)}} q_{(x_i)} dx_i = \frac{1}{n} \sum_{i=1}^n \int_{x \sim q} f_{(x_i)} P(x_i) dx_i \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[f_{(x_i)}] \end{aligned}$$

$$\rightarrow \hat{\mu}_{ss} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{x \sim p}[f_{(x_i)}] = \frac{1}{n} \sum_{i=1}^n \mu = \mu \Rightarrow s \text{ is unbiased}$$

$$\begin{aligned} 2. \text{Var}(\hat{\mu}_{ss}) &= \text{Var}(\frac{1}{n} \sum_{i=1}^n g_{(x_i)}) \text{ where } g_{(x_i)} = f_{(x_i)} \frac{P(x_i)}{q_{(x_i)}} \Rightarrow \hat{\sigma}_{ss}^2 = \frac{1}{n^2} \text{Var}(\sum_i g_{(x_i)}) = \frac{1}{n} \text{Var}(g_{(x_i)}) + \frac{1}{n} \left[\mathbb{E}[g_{(x_i)}]^2 - \hat{\mu}_{ss}^2 \right] \\ &= \frac{1}{n} \left[\int f_{(x_i)}^2 \frac{P(x_i)^2}{q_{(x_i)}} dx_i - \hat{\mu}_{ss}^2 \right] = \frac{1}{n} \left[\int f_{(x_i)}^2 \frac{P(x_i)^2}{q_{(x_i)}} dx_i - \mu^2 \right] \end{aligned}$$

for minimum $\text{Var}(\hat{\mu}_{ss})$, we can choose $f_{(x)} = \frac{q(x)}{P(x)} \times \mu$ which gives us zero variance.

Also, let's assume $f(x) = \frac{1}{p(x)q(x)}$. Then $\sigma^2 = \infty$ if $q(x)=0$ for some x .

$$3. \hat{\mu}_{N-SS} = \frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i} = \frac{\sum_{i=1}^n \frac{p_{xi}}{q_{xi}} f(x_i)}{\sum_{i=1}^n \frac{p_{xi}}{q_{xi}}}$$

for $n \rightarrow \infty$, like $n=1$ we have: $E[\hat{\mu}_{N-SS}] = E\left[\frac{\frac{p_{x1}}{q_{x1}} f(x_1)}{\frac{p_{x1}}{q_{x1}}}\right] = E[f(x_1)]$ which is not necessarily M .

for $n \rightarrow \infty$, $\lim_{n \rightarrow \infty} E[\hat{\mu}_{N-SS}] = \lim_{n \rightarrow \infty} E\left[\frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i}\right] = E\left[\frac{\sum_{i=1}^n w_i}{\sum_{i=1}^n w_i} \cdot \text{cor}\left(\frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i}, \frac{n}{\sum_{i=1}^n w_i}\right)\right]$

for $f(x)=1$ we have:

$$\lim_{n \rightarrow \infty} E[\hat{\mu}_{N-SS}] = M \times 1 + \lim_{n \rightarrow \infty} \text{cor}\left(\frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i}, \frac{n}{\sum_{i=1}^n w_i}\right) = M$$

So $\hat{\mu}_{N-SS}$ is unbiased for $n \rightarrow \infty$

$$4. m < f(x) \leq M \Rightarrow \begin{cases} \hat{\mu}_{N-SS} = \frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i} \leq \frac{\sum_{i=1}^n w_i M}{\sum_{i=1}^n w_i} = M \\ \hat{\mu}_{N-SS} = \frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i} \geq \frac{\sum_{i=1}^n w_i m}{\sum_{i=1}^n w_i} = m \end{cases} \Rightarrow m < \hat{\mu}_{N-SS} \leq M$$

$$\sigma^2 = E[\hat{\mu}_{N-SS}^2] - M^2 \quad \text{also} \quad E[(M - \hat{\mu}_{N-SS})(\hat{\mu}_{N-SS} - m)] = -E[\hat{\mu}_{N-SS}] \cdot mM + (M+m)m \quad \Rightarrow E[\hat{\mu}_{N-SS}] \leq mM + (M+m)m$$

Since M, m are constant and $\hat{\mu}_{N-SS} - m$
both are positive values with positive expected value

$$\Rightarrow \sigma^2 \leq (M+m)/M - mM - M^2 = (M-m)(m-M) \leq \frac{(m+M)^2}{4}$$

5.

$$H = \int_{-\infty}^{\infty} \exp(-\frac{u^2}{2}) du = \sqrt{2\pi} (1 - \Phi(2)) \approx 0.0577$$

(الف)

$$\hat{\mu}_{MC} = \frac{1}{n} \sum \frac{f(x_i)}{p(x_i)} I\{u_i > 2\} = 0 \quad \text{because all } u_i's \text{ are less than 2}$$

(بـ)

$$\hat{\mu}_{SS} = \frac{1}{n} \sum_i \frac{\sqrt{2\pi} p(x_i)}{q(x_i)} \quad \text{where } p \sim N(0,1) \text{ and } q \sim N(3,1)$$

because samples now are from $N(3,1)$, we shift all of samples by 3. and then we have:

$$\hat{\mu}_{SS} = \frac{1}{|S_1|} \sum_{x \in S_1} \frac{\sqrt{2\pi} e^{-\frac{x^2}{2}}}{e^{-\frac{(x-3)^2}{2}}} = \frac{1}{|S_1|} \sum_{x \in S_1} e^{\frac{2-6x}{2}} \approx 0.0223$$

we see that SS estimator is a better estimator than MC; because of its low variance

$$D_{KL}(P||q) = \int p(x) \log \frac{p(x)}{q(x)} dx \rightarrow \text{Jensen inequality: } f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i) \text{ for } f \text{-convex}, \left\{ \begin{array}{l} \lambda_i \geq 0 \\ \sum_{i=1}^n \lambda_i = 1 \end{array} \right.$$

from Jensen ineq., we have: $\log(E_x[g(x)]) \geq E_x[\log(g(x))]$

$$\begin{aligned} -D_{KL}(P||q) &= -\int p(x) \left(-\log \frac{q(x)}{p(x)} \right) dx = \int p(x) \log \frac{p(x)}{q(x)} dx \leq \log \int p(x) \frac{p(x)}{q(x)} dx = \log 1 = 0 \\ \rightarrow D_{KL}(P||q) &\geq 0 \end{aligned}$$

Now, let $q(x) = p(x)$, we have: $D_{KL}(P||P) = \int p(x) \log \frac{p(x)}{p(x)} dx = 0$
So minimum of D_{KL} is when $q(x) = p(x)$ where KL divergence is zero.

$$p(z|x) = \frac{p(z,x)}{p(x)}$$

$$\begin{aligned} D_{KL}(q(z)||p(z|x)) &= \int q(z) \log \frac{q(z)p(x)}{p(z,x)} dz = \int q(z) \left[\log q(z) + \log p(x) - \log p(x,z) \right] dz \\ &= \int q(z) \left[\log q(z) - \log p(z,x) \right] dz + \int q(z) \underbrace{\log p(x)}_{I} dz \end{aligned}$$

$$I: \int q(z) \log p(x) dz = \log p(x) \int q(z) dz = \log p(x)$$

$$\begin{aligned} \rightarrow D_{KL}(q(z)||p(z|x)) &= \int q(z) \log q(z) - \log p(x,z) + \log p(x) \\ &\quad \underbrace{\vphantom{\int} \log p(x) - \log p(x,z)}_{E_q[\log q(z) - \log p(x,z)]} \end{aligned}$$

$$\begin{aligned} \rightarrow D_{KL}(q(z)||p(z|x)) &= \log p(x) + E_q[\log q(z) - \log p(x,z)] \\ &= \log p(x) - (E_q[\log p(x,z)] - E_q[\log q(z)]) \\ &= \log p(x) - L(q) \end{aligned}$$

$$L(Q) = E_q[\log p(x, z)] - E_q[\log Q(z)], Q(z) = \prod_{i=1}^n q_i(z)$$

$$p(z_1, z_2, \dots, z_n, x_1, x_2, \dots, x_m) = p(x, m) \prod_{i=1}^n p(z_i | z_1, z_2, \dots, z_{i-1}, x_1, \dots, x_m)$$

$$\approx L(Q) = \sum_{i=1}^n E_i[\log p(z_i | z_1, \dots, z_{i-1}, x)] - E_i[\log q(z_i)]$$

$$\approx L(q_k) = E[\log p(z_k | z_{-k}, x) - E_j[\log q(z_k)] + C$$

Define: $F = \int q(z_k) E_{-k}[\log p(z_k | z_{-k}, x)] dz_k - \int q(z_k) \log q(z_k) dz_k$

$$= \int q(z_k) (E_{-k}[\log p(z_k | z_{-k}, x)] - q(z_k)) dz_k$$

using calculus of variation, we have $\frac{d}{dz_k} \frac{\partial F}{\partial q(z_k)} = \frac{\partial F}{\partial q(z_k)}$ to maximize $L(q_k)$. Also it's obvious that if all of q_k 's are optimal for $L(q_k)$, the $L(q_k)$ is optimal.

$$\approx \frac{d}{dz_k} \frac{\partial F}{\partial q(z_k)} = \frac{\partial F}{\partial q(z_k)} = E_{-k}[\log p(z_k | z_{-k}, x)] - \log q(z_k) - 1 = 0 \Rightarrow \log q_k(z_k) = E_{-k}[\log p(z_k | z_{-k}, x)] + C$$

$$\ln q^*(z) = E_{\pi, \mu, \Lambda}[\ln p(x, z, \pi, \mu, \Lambda)] + C = E_{\pi, \mu, \Lambda}[\ln p(z | \pi) p(x | z, \mu, \Lambda)] + C$$

$$= E_\pi[\ln p(z | \pi)] + E_{\mu, \Lambda}[\ln p(x | z, \mu, \Lambda)] + C$$

$$= E_\pi[\ln \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}}] + E_{\mu, \Lambda}[\ln \frac{1}{(2\pi)^D |\Lambda_k|^{1/2}} \exp(-\frac{1}{2}(x - \mu_k)^T \Lambda_k^{-1} (x - \mu_k))] + C$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} [E_\pi[\pi_k] + \frac{1}{2} E[\ln |\Lambda_k|] - \frac{D}{2} \ln(2\pi) - \frac{1}{2} E_{\mu, \Lambda}[(x - \mu_k)^T \Lambda_k^{-1} (x - \mu_k)]] + C$$

Let $p_{nk} = \underbrace{\pi_k}_{\text{we have: } \ln q^*(z)} \Rightarrow \ln q^*(z) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} p_{nk} \Rightarrow q^*(z) \propto \prod_{n=1}^N \prod_{k=1}^K p_{nk}^{z_{nk}}$

Let $r_{nk} = \frac{p_{nk}}{\sum_j p_{nj}}$ which normalizes p_{nk} such that $q^*(z)$ becomes a categorical probability distribution

whit $E[z_{nk}] = r_{nk}$

Due to structural property of our model, we have: $q(\pi, \mu, \Lambda) = q(\pi) q(\mu, \Lambda)$

$$= q(\pi) \prod_{k=1}^K q(\mu_k, \Lambda_k)$$

Then:

$$\begin{aligned} \log q_{(z)}^+ &= E_z [\log p(z|y)] + C = E_z [\log p(y) + \log p(z|y)] + C = \log p(y) \cdot E_z [\log p(z|y)] + C \\ &= (\alpha_0 - 1) \sum_{k=1}^K \pi_k + \sum_{n=1}^N \sum_{k=1}^K \gamma_k E_z [z_{nk}] = \sum_{k=1}^K \pi_k (\alpha_0 + \sum_{n=1}^N \gamma_{nk} - 1) \end{aligned}$$

$\downarrow \sim \gamma_{nk}$

Thus $q_{(z)}^+$ has dirichlet distribution with $\alpha = \alpha_0 + \sum_{n=1}^N \gamma_{nk}$

$$\begin{aligned} \log q_{(\mu_k, \lambda_k)}^+ &= E_z [\log p(\mu_k, \lambda_k) p(z|\mu_k, \lambda_k)] + C = \log p(\mu_k, \lambda_k) + E_z [\log p(x|z, \mu_k, \lambda_k) p(z|\mu_k, \lambda_k)] + C \\ &= \log p(\mu_k, \lambda_k) + E_z [\log \prod_{n=1}^N N(x_n | \mu_k, \lambda_k^{-1})^{z_{nk}}] + C \\ &= \log p(\mu_k, \lambda_k) + E_z [\sum_{n=1}^N z_{nk} \log N(x_n | \mu_k, \lambda_k^{-1})] + C \\ &= \log p(\mu_k, \lambda_k) + \sum_{n=1}^N E_z [z_{nk}] \log N(x_n | \mu_k, \lambda_k^{-1}) + C \\ &= \log N(\mu_k | m_0, (\beta_0 \lambda_k)^{-1}) \ln(\lambda_k | w_0, v_0) + \sum_{n=1}^N \gamma_{nk} \log N(x_n | \mu_k, \lambda_k^{-1}), C \end{aligned}$$

Given Normal-Wishart properties, we have: $q_{(\mu_k, \lambda_k)}^+ = N(\mu_k | m_k, (\beta_k \lambda_k)^{-1}) \ln(\lambda_k | w_k, v_k)$

where: $N_k = \sum_{n=1}^N \gamma_{nk}$, $\bar{x}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} x_n$, $S_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (x_n - \bar{x}_k)(x_n - \bar{x}_k)^T$, $V_k = V_0 + N_k$

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

To calculate p_{hk} which gives us γ_{nk} , we have to calculate $E_{\mu_k, \lambda_k} [(X_n - \mu_k)^T \ln(\lambda_k | w_k, v_k)], E[\ln \lambda_k]$

$$E_{\mu_k, \lambda_k} [(X_n - \mu_k)^T \ln(\lambda_k | w_k, v_k)] = D \beta_k^{-1} + V_k (X_n - m_k)^T W_k (X_n - m_k)$$

$$E[\ln \lambda_k] = \ln \tilde{\lambda}_k, E[\ln \lambda_k] = \ln \tilde{\lambda}_k$$

to calculate γ_{nk} , we need rest of parameters, and vice-versa. To solve this problem we can use estimation-maximization to estimate γ_{nk} and then update μ_k, λ_k to maximize $\log q_{(z)}^+$, which gives us minimum D_{KL} .

$$Q(\theta | \theta^{(t)}) = E_{z \sim p(z|x, \theta^{(t)})} [\log(x, z|\theta)]$$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta | \theta^{(t)})$$

$$\text{where } \theta^{(t)} = \{\mu^{(t)}, \lambda^{(t)}, \pi^{(t)}\}$$