

Problem 1

(a) - For $D_{KL}(Q||P) = \infty$, we can have $P(x) = 0$ where $Q(x) \neq 0$.

- for example we can have P to be uniform random variable on $[0, 1]$ and Q be a uniform random variable on $[-1, 1]$.

- for $x \in [-1, 0)$, $Q(x)$ is divided by 0 hence $D_{KL}(Q||P) = \infty$.

$$- D_{KL}(P||Q) = E_{x \sim P(x)} \left[\log \frac{P(x)}{Q(x)} \right]$$

$$= \int_0^1 \log \frac{P(x)}{Q(x)} dx$$

$$= \int_0^1 \log \frac{1}{(1/2)} dx$$

$$= \log 2 \neq \infty$$

$$\textcircled{b} \quad \text{Forward KL} \Rightarrow D_{KL}(P||Q) = E \left[\log \frac{P(x)}{Q(x)} \right]$$

$$\text{Reverse KL} \Rightarrow D_{KL}(Q||P) = E \left[\log \frac{Q(x)}{P(x)} \right]$$

For Plot A :-

- Q avoids being tiny where P is large, which is not very costly for forward KL.
- Q establishes "average" balance between the modes of P where P is small. This is costly in reverse KL
- Hence we can conclude that 'A' corresponds to minimizing forward KL

For Plot B :-

- where P is large, Q is extremely small, causing high forward KL
- where Q is large, P is also notably large, leading to not very costly reverse KL
- we can conclude that B minimizes reverse KL.

(C) For the ELBO

$$\begin{aligned} \mathcal{L}(x; \theta, \phi) &= \mathbb{E}_{q_{\phi}(z|x)} [\log p(x|z; \theta)] \\ &\quad - D_{KL}(q_{\phi}(z|x) || p(z)) \end{aligned}$$

- made up of two terms
 - reconstruction \approx regularization (KL)

- The regularization/KL term measures the divergence b/w the learned latent distribution $q_{\phi}(z|x)$ and prior distribution $p(z)$ which is typically gaussian.
- The asymmetry here means the VAE is more penalized when it proposes a distribution $q_{\phi}(z|x)$ that places probability mass in regions where $p(z)$ does not. This means
 - "penalizing" for proposing unlikely/unexpected latent space and is less concerned about "missing out" on regions prior considers likely
- Implications:- ① regularization pushes posterior to match prior, preventing overfitting & ensuring well-structured latent space
- Asymmetry ensures the encoder does not produce representations that are too unusual. Even if it helps reconstruction

Problem 2

(a) Mixture of Bernoullis

$$\rightarrow (i) \quad x \in \{0,1\}^D \rightarrow \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$P(x_d=1) = p_d \quad i \quad p \in (0,1)^D \rightarrow \begin{bmatrix} P(x_0=1) \\ P(x_1=1) \\ \vdots \\ P(x_D=1) \end{bmatrix}$$

$$P(x|p) = \prod_{d=1}^D p_d^{x_d} (1-p_d)^{(1-x_d)}$$

$$\rightarrow (ii) \text{ mixture of } k \text{ Bernoullis}$$

$x^{(i)}$ is drawn from some vector of Bernoulli random variables with parameters $p^{(k)}$ distribution $\pi(k)$

$$x^{(i)} \text{ drawn from } \begin{bmatrix} x_0 \\ \vdots \\ x_m \end{bmatrix} \sim \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix}$$

mixture with distribution $\pi(k)$

Let A_k be an event of drawing/sampling $x=x^{(i)}$ from the mixture. Then,

$$\begin{aligned} P(x|p, \pi) &= \sum P(x, A_k | p, \pi) \\ &= \sum P(x_i | A_k, p, \pi) \times P(A_k | p, \pi) \end{aligned}$$

① prob of observing $x=x^{(i)}$ given it's drawn from $p^k \rightarrow P(x|p)$ because knowing x was sampled from p^k makes π irrelevant

②. $P(A_K | p, \pi)$ - measures prob of sampling $n(i)$
nothing but $\pi^{(K)}$

Hence

$$P(n^{(i)} | p, \pi) = P(x_i, A_K | p, \pi) = \sum_k \pi_k P(x_i | p^{(k)})$$

→ $x = \{x_i^{(i)}\}_{i=1 \dots n} \rightarrow$ drawing from entire dataset

$x^{(i)}$ is individual observation in x

log likelihood is calculated by summing the
log prob of each individual observation $x^{(i)}$
assuming the observations are independent

$$\log L(p, \pi) = \sum_{i=1}^N \log P(x_i^{(i)} | p, \pi)$$

(b) Expectation step

(i) Let $A_k^{(i)}$ be the event that $x^{(i)}$ was drawn from $p^{(k)}$

$$P(z^{(i)} | \pi) = \prod_{k=1}^K \pi_k^{z_k^{(i)}}$$

$$\begin{aligned} P(x^{(i)} | z^{(i)}, p, \pi) &= \prod_k P(x^{(i)} | z^{(i)}, p, \pi, A_k^{(i)}) \\ &= \prod_{k=1}^K [P(x^{(i)} | p^{(k)})]^{z_k^{(i)}} \end{aligned}$$

(ii) likelihood of data and latent variable

$$P(z, x | \pi, p) = \prod_{i=1}^N P(x^{(i)}, z^{(i)} | \pi, p)$$

$$= \prod_{i=1}^N P(x^{(i)} | z^{(i)}, \pi, p) P(z^{(i)} | \pi)$$

$$= \prod_{i=1}^N \left[\prod_{k=1}^K \left[P(x^{(i)} | p^{(k)}) \right]^{z_k^{(i)}} \right] \left[\prod_{k=1}^K \pi_k^{z_k^{(i)}} \right]$$

$$= \prod_{i=1}^N \left[\prod_{k=1}^K \left[\pi_k P(x^{(i)} | p^{(k)}) \right]^{z_k^{(i)}} \right]$$

(c)

$$\begin{aligned} \eta(z_k^{(i)}) &= E[z_k^{(i)} | x^{(i)}, \pi, p] \\ &= P(z_k^{(i)} = 1 | x^{(i)}, \pi, p) \\ &= \frac{P(x^{(i)}) \cdot P(z_k^{(i)} = 1 | \pi, p)}{\sum_j P(x^{(i)}) \prod_{d=1}^D (p_d)^{z_d^{(i)}} (1-p_d)^{1-z_d^{(i)}}} \end{aligned}$$

Computing log likelihood:-

$$\begin{aligned} \log P(Z, D | \pi, p) &= \sum_{i=1}^N \left[\sum_{k=1}^K z_k^{(i)} \left(\log P(z_k^{(i)} | p^{(i)}) + \log \pi_k \right) \right. \\ &= \sum_{i=1}^N \sum_{k=1}^K z_k^{(i)} \left[\log \pi_k + \sum_{d=1}^D (z_d^{(i)} \log p_d + \right. \\ &\quad \left. \left. + (1-z_d^{(i)}) \log (1-p_d^{(i)}) \right) \right] \end{aligned}$$

Taking the expected value \mathbb{E} replacing

$$E[z_k^{(i)}] = \eta(z_k^{(i)})$$
 gives the solution

(i) setting the derivative of

$$\frac{d}{d p_d^{(k)}} E[\log P(z, d | \pi, p)] = 0.$$

$$\sum_{i=1}^N \eta(z_k^{(i)}) \left[\frac{\pi_d^{(i)}}{p_d^{(k)}} - \frac{1-\pi_d^{(i)}}{1-p_d^{(k)}} \right] = 0$$

$$\sum_{i=1}^N \eta(z_k^{(i)}) \left[\pi_d^{(i)} (1-p_d^{(k)}) - (1-\pi_d^{(i)}) p_d^{(k)} \right] = 0$$

$$\Rightarrow p_d^{(k)} = \frac{\sum_{i=1}^N \eta(z_k^{(i)}) \pi_d^{(i)}}{\sum_{i=1}^N \eta(z_k^{(i)})} = \frac{\sum_{i=1}^N \eta(z_k^{(i)}) \pi_d^{(i)}}{N_k}$$

(ii) we need to only minimize $\sum_{i=1}^N \sum_{k=1}^K \eta(z_k^{(i)}) \log \pi_k$

since the rest are not a function of π

In order to keep π a distribution, we require

$$\sum_k \pi_k = 1 \text{ let } \lambda \text{ be the dual variable for this}$$

constraint.

$$L(\pi, \lambda) = -\sum_{i=1}^N \sum_{k=1}^K \eta(z_k^{(i)}) \log \pi_k + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

Taking the derivative w.r.t π_k :-

$$\frac{d}{d\pi_k} \mathcal{L}(\pi, \lambda) = -\sum_{i=1}^N \frac{\eta(z_k^{(i)})}{\pi_k} + \lambda = 0$$

$$\Rightarrow \pi_k = \frac{\sum_{i=1}^N \eta(z_k^{(i)})}{\lambda} = \frac{N_k}{\lambda}$$

Solving for λ

$$\mathcal{L}(\lambda) = -\sum_{i=1}^{N_K} \eta(z_k^{(i)}) (\log N_k - \log \lambda) + \sum_{k=1}^K (N_k - \lambda)$$

taking derivative w.r.t λ

$$\frac{1}{\lambda} \sum_{i=1}^N \sum_{k=1}^K (\eta(z_k^{(i)}) - 1) = 0.$$

$$\lambda = \sum_{i=1}^N \sum_{k=1}^K \eta(z_k^{(i)}) = \sum_{k=1}^K N_k.$$

Optimal $\boxed{\pi_k = \frac{N_k}{\sum_{k=1}^K N_k}}$

Problem 3

→ To generate samples,

- we start by sampling a mixture component / cluster from the distribution $\text{mix-}\pi$, means choosing a cluster ' k ' with prob $\text{mix-}\pi[k]$
- after choosing cluster, we generate a new image by sampling each pixel based on Bernoulli distribution parameterized by $p[k]$.
- once all pixels are sampled, we reshape the flattened binary image data.

→ Plausibility

- largely depends on how well the EM algo has clustered the data & how accurately the Bernoulli params p represent each cluster's typical image.

For $k=5$; EM algo might group diverse digits into same cluster, leading to muddled representations shown in the images attached.

For $k=20$; more clusters → can better segregate different digits & variations of digits. This results in more recognizable digits looking more like MNIST samples as shown in images attached.

Problem 4

→ NICE:-

→ volume-preserving transformations & invertible mappings, tries to learn a more structured latent space.

→ visualised samples do not show clear digit shapes, there are occasional unusual digit forms, as shown in the image

→ VAE

→ has ability to produce more smoother & continuous latent space, helpful in generative tasks.

→ the images generated are more consistent with MNIST data compared to NICE

→ VAE regularizes latent space using probabilistic framework, ensuring that samples generated from any point in latent space produces valid outputs as shown in images

→ SM generation quality largely depends on the number of clusters. Given the nature of MNIST where the goal is to generate clear & recognizable digits, VAE is most suitable because of the ability to generate images from smooth & continuous latent space.

Problem 5

(a) Cauchy distribution:-

$$f_X(u) = \frac{1}{\pi(1+u^2)} \quad \text{for } -\infty < u < \infty$$

$$\text{Let } Y = \frac{1}{X},$$

Consider the CDF of Y , for $v > 0$.

$$\begin{aligned} F_Y(v) &= P\{Y \leq v\} = P\left\{\frac{1}{X} \leq v\right\} \\ &= P\left\{\frac{1}{v} \leq X, X > 0\right\} + P\left\{\frac{1}{v} \geq X, X < 0\right\} \\ &= P\left\{\frac{1}{v} \leq X\right\} + P\{X < 0\} \\ &= 1 - P\left\{\frac{1}{v} \geq X\right\} + P\{X < 0\} \\ &= 1 - F_X\left(\frac{1}{v}\right) + F_X(0) \end{aligned}$$

Taking the derivative w.r.t v .

$$\begin{aligned} f_Y(v) &= \frac{d}{dv} F_Y(v) = -\frac{1}{v^2} f_X\left(\frac{1}{v}\right) \\ &= \frac{1}{v^2} \cdot \frac{1}{\pi\left(1+\left(\frac{1}{v}\right)^2\right)} = \frac{1}{\pi(1+v^2)} \end{aligned}$$

similarly for $v < 0$,

$$\begin{aligned}F_Y(v) &= P(Y \leq v) = P\left\{\frac{1}{X} \leq v^2\right\} \\&= P\left\{\frac{1}{v^2} \geq X, X > 0\right\} + P\left\{\frac{1}{v} \leq X, X < 0\right\} \\&\Rightarrow 0 + P\left\{\frac{1}{v} \leq X < 0\right\} = F_X(0) - F_X\left(\frac{1}{v}\right)\end{aligned}$$

Taking the derivative w.r.t. v .

$$\begin{aligned}f_Y(v) &= \frac{d}{dv} F_Y(v) = -\left(\frac{-1}{v^2}\right) f_X\left(\frac{1}{v}\right) \\&= \frac{1}{\pi(1+v^2)}\end{aligned}$$

Therefore $Y = \frac{1}{X}$ also has Cauchy distribution

⑥ $Z = e^Y$ where Y is normally distributed with parameters μ & σ .

CDF of Z , in terms of CDF of Y

For $z > 0$,

$$P(Z \leq z) = P(e^Y \leq z) = P(Y \leq \ln(z))$$

This is the CDF of Y evaluated at $\ln(z)$, so:

$$F_Z(z) = F_Y(\ln(z))$$

PDF of Z by differentiating

$$f_Z(z) = \frac{d}{dz} F_Y(\ln(z))$$

$$= f_Y(\ln(z)) \left(\frac{1}{z} \right)$$

using the PDF of Y

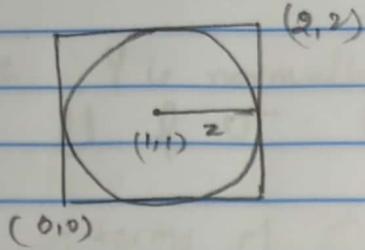
$$f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

we get

$$f_Z(z) = \frac{1}{\sqrt{2\pi}\sigma z} e^{-\frac{(\ln(z)-\mu)^2}{2\sigma^2}}$$

Q this is the density of e^Y
where $z > 0$

Q

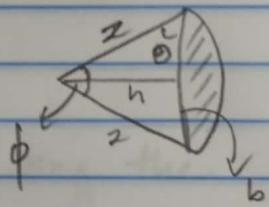
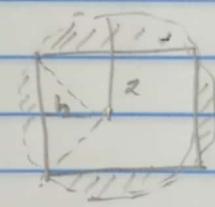


(i) for $0 \leq z \leq 1$, $CDF = \frac{1}{4} \text{ area of circle} = \frac{1}{4} \pi z^2$

$1/\mu \rightarrow$ normalizing factor.

' for $1 \leq z \leq \sqrt{2}$

(ii) $CDF = \frac{1}{4} (\text{area of outer circle}) - \text{shaded region}$



$$\begin{aligned} \text{area of } A^c &= \frac{1}{2} b \cdot h. \\ &= \frac{1}{2} (\sqrt{z^2-1}) \cdot (1) \\ &= \underline{\underline{(\sqrt{z^2-1})}} \end{aligned} \quad \left. \begin{array}{l} \left(\frac{b^2}{2} + h^2 = z^2 \right. \\ b = \sqrt{z^2-h^2} \\ b = \sqrt{z^2-1} \end{array} \right.$$

$$\begin{aligned} \text{area of the sector} &= \frac{\phi}{2\pi} \cdot (\pi z^2) \\ &= \frac{\phi z^2}{2} \\ &= z^2 \tan^{-1} \sqrt{z^2-1} \end{aligned}$$

$$\begin{aligned} \tan\left(\frac{\phi}{2}\right) &= \frac{b/2}{h} \\ &= \sqrt{z^2-1} \\ \phi &= 2 \left(\tan^{-1} \sqrt{z^2-1} \right) \end{aligned}$$

\therefore area of shaded region

$$= \underline{\underline{(z^2 \tan^{-1} \sqrt{z^2-1})}} - \underline{\underline{(\sqrt{z^2-1})}}$$

CDF :-

$$\frac{F(z)}{z} = \frac{1}{4} \left(\pi z^2 - 4 \left(z^2 \tan^{-1}(\sqrt{z^2-1}) - \sqrt{z^2-1} \right) \right)$$

Differentiating to find PDF :-

(i) for $0 < z \leq 1$

$$f_z(z) = \frac{\pi z}{2}$$

for $1 < z \leq \sqrt{2}$

$$f_z(z) = \frac{1}{4} 2\pi z - 4 \left(2z \tan^{-1}(\sqrt{z^2-1}) + \frac{z^2-1}{4(\sqrt{z^2-1})^2} \cdot \frac{dz}{dz} - \frac{dz}{2\sqrt{z^2-1}} \right)$$

$$= \frac{1}{4} (2\pi z - 8z \tan^{-1} \sqrt{z^2-1})$$

$$= \frac{\pi z}{2} - 2z \tan^{-1} \sqrt{z^2-1}$$

Problem 6. Jacobian of Leaky flow

Given $f(z) = [\text{LReLU}[z_1], \text{LReLU}[z_2], \dots, \text{LReLU}[z_n]]$

The Jacobian matrix J is an $n \times n$ matrix where entry j_{ij} is $\frac{\partial f_i}{\partial z_j}$.

For Leaky ReLU

$$\text{LReLU}[z_i] \rightarrow \begin{cases} 0 \cdot z_i & \text{for } z_i < 0 \\ z_i & \text{for } z_i \geq 0 \end{cases}$$

So partial derivatives are:-

$$\frac{\partial \text{LReLU}(z_i)}{\partial z_i} \rightarrow \begin{cases} 0 \cdot 1 & \text{for } z_i < 0 \\ 1 & \text{for } z_i \geq 0. \end{cases}$$

The Jacobian matrix J is diagonal for elementwise operation because f_i depends only on z_i and not on other components of z .

$$J = \text{diag} \left(\frac{\partial \text{LReLU}(z_1)}{\partial z_1}, \dots, \frac{\partial \text{LReLU}(z_n)}{\partial z_n} \right)$$

Determinate of the Jacobian:-

$$\det(J) = \prod_{i=1}^n \frac{\partial \text{LReLU}(z_i)}{\partial z_i}$$

Inverse of the Jacobian

Inverse of diagonal matrix is also diagonal with inverse of each diag element

$$J^{-1} = \text{diag} \left(\frac{1}{\frac{\partial \text{ReLU}(z_i)}{\partial z_i}}, \dots, \frac{1}{\frac{\partial \text{ReLU}(z_n)}{\partial z_n}} \right)$$

Compute the inverse absolute determinate of Jacobian

Inverse absolute of \det of Jacobian

$$\det(J^{-1}) = \frac{1}{\prod_{i=1}^n \frac{\partial \text{ReLU}(z_i)}{\partial z_i}}$$

Problem :-

Difference :-

- Continuous flows operate over continuous state spaces, where any small change in parameters can cause change in model output. Discrete flows operate over finite set of states making modeling difficult.
- In continuous flows, the determinant of the Jacobian gives a measure of change in volume (density). whereas, in discrete space there is no volume/density change in conventional sense, hence Jacobian is not directly applicable.
- In discrete domain, a simple invertible function is a bitwise NOT operation. A discrete flow model could consist of a sequence of such invertible operations. For instance, given a binary input, applying a series of bitwise not operations & bitwise shifts could constitute the transformations of a flow model.

Each transformation in the sequence must have a computable inverse to adhere to the principles of flow-based models. This simplistic model demonstrates discrete flow models' essence - applying a series of invertible transformations to discrete data while ensuring each step's invertibility & tractability.