

In The Name of God

Sharif University of Technology
Electrical Engineering Department

Deep Generative Models

Assignment 3

Fall 2024

Instructor: Dr. S. Amini

Due on Azar 11, 1403 at 23:55

**Matin M.babaei 400102114****1 Probability Review**Consider the transformation $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$.

$$X = f(Z) = [z_1 \ z_2 e^{z_1} \ (z_3 e^{-z_1} + z_1^2)^{\frac{1}{2}}]$$

Is this an invertible transformation? What's the Jacobian of this transformation? Now assume we have a Multivariate Normal base distribution $p(Z) \sim N_3(0, I)$ on the domain of f . What's $p(x)$ at $x = [0 \ 1 \ 1/3]^T$ ($p(X)$ and $p(Z)$ are probability density functions)

If invertible, there should be a $z = f^{-1}(x)$:

Given Z - we have:

$$\begin{cases} x_1 = z_1 \quad (i) \\ x_2 = z_2 e^{z_1} \quad (ii) \\ x_3 = (z_3 e^{-z_1} + z_1^2)^{\frac{1}{2}} \quad (iii) \end{cases}$$

Now, if we are given X , can we obtain Z ? we'll see,

$$\begin{aligned} (x_1, x_2, x_3) &\xrightarrow{(i)} z_1 = x_1 \\ &\xrightarrow{(ii)} z_2 = \frac{x_2}{e^{x_1}} = \frac{x_2}{e^{z_1}} \\ &\xrightarrow{(iii)} x_3^2 = z_3 e^{-z_1} + z_1^2 \Rightarrow x_3^2 - x_1^2 = z_3 e^{-x_1} \Rightarrow z_3 = \frac{x_3^2 - x_1^2}{e^{-x_1}} \end{aligned}$$

Since Z was thoroughly determined having X , we can introduce this inverse transformation.

$$z = f^{-1}(x) = \begin{cases} z_1 = x_1 \\ z_2 = \frac{x_2}{e^{x_1}} = e^{-x_1} x_2 \\ z_3 = \frac{x_3^2 - x_1^2}{e^{-x_1}} = e^{x_1} (x_3^2 - x_1^2) \end{cases}$$

Thus, the transformation is invertible, using the function defined above.

$$J_f(z) = \frac{\partial \mathbf{x}}{\partial z} = \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \frac{\partial x_1}{\partial z_r} & \frac{\partial x_1}{\partial z_p} \\ \frac{\partial x_r}{\partial z_1} & \frac{\partial x_r}{\partial z_r} & \frac{\partial x_r}{\partial z_p} \\ \frac{\partial x_p}{\partial z_1} & \frac{\partial x_p}{\partial z_r} & \frac{\partial x_p}{\partial z_p} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ z_r e^{z_1} & e^{z_1} & 0 \\ \frac{1}{\mu} (-z_p e^{-z_1} + z_r) (z_r e^{-z_1} + z_p)^{-\frac{r}{\mu}} & 0 & \frac{1}{\mu} (e^{-z_1}) (z_r e^{-z_1} + z_p)^{-\frac{r}{\mu}} \end{bmatrix}$$

$$\alpha_1 = z_1, \quad \alpha_r = z_r e^{z_1}, \quad \alpha_p = (z_r e^{-z_1} + z_p)^{\frac{1}{\mu}}$$

$$\frac{\partial x_1}{\partial z_1} = 1, \quad \frac{\partial x_1}{\partial z_r} = \frac{\partial x_1}{\partial z_p} = 0, \quad \frac{\partial x_r}{\partial z_1} = z_r e^{z_1}, \quad \frac{\partial x_r}{\partial z_r} = e^{z_1}, \quad \frac{\partial x_r}{\partial z_p} = 0$$

$$\frac{\partial x_p}{\partial z_1} = \frac{1}{\mu} (-z_p e^{-z_1} + z_r) (z_r e^{-z_1} + z_p)^{-\frac{r}{\mu}}, \quad \frac{\partial x_p}{\partial z_r} = 0, \quad \frac{\partial x_p}{\partial z_p} = \frac{1}{\mu} (e^{-z_1}) (z_r e^{-z_1} + z_p)^{-\frac{r}{\mu}}$$

$$P(x) @ x = [a, 1, \frac{1}{\mu}]^T$$

$$P(x) = P(z) \cdot |\det(J_f(z))|^{-1}$$

$$z_1 = \alpha_1 = a, \quad z_r = \frac{\alpha_r}{e^{z_1}} = 1, \quad z_p = e^{z_1} (\alpha_p - z_r) = (\frac{1}{\mu})^{\frac{r}{\mu}} = \frac{1}{\mu^r}$$

$$P(z_p) = (1)^{-\frac{r}{\mu}} |\sum \exp(-\frac{1}{\mu} (z_p - \mu)^T \Sigma^{-1} (z_p - \mu))|$$

$$d = r, \quad \Sigma = I, \quad \mu = 0, \quad \|z\|_p^r = z_1^r + z_r^r + z_p^r = a + 1 + (\frac{1}{\mu})^{\frac{r}{\mu}}$$

$$P(z_p) = \frac{1}{(1)^{\frac{r}{\mu}} (1)^{\frac{1}{\mu}}} \exp(-\frac{1}{\mu} z^T z) = (1)^{-\frac{r}{\mu}} \exp(-\frac{1}{\mu} \|z\|_p^r)$$

$$\rightarrow P(z) = (1)^{-\frac{r}{\mu}} \exp(-\frac{1}{\mu} (\frac{\|z\|_p^r}{\sqrt{r}}))$$

$$J_f(z) = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ -\frac{1}{\mu} & 0 & \mu \end{bmatrix}$$

$$|\det(J_f(z))| = 1 \cdot (\mu - 1) = \mu$$

$$P(x) = P(z) |\det(J_f(z))|^{-1} = (1)^{-\frac{r}{\mu}} \exp(-\frac{1}{\mu} (\frac{\|z\|_p^r}{\sqrt{r}})) \frac{1}{\mu}$$

2 VP-NFs are not universal

Remember a Normalizing Flow formulated as:

$$p_\theta(\mathbf{x}) = p(z = f_\theta^{-1}(\mathbf{x})) \left| \det \frac{\partial f_\theta^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right|$$

In a Volume-Preserving Normalizing Flow (VP-NF), the determinant of the Jacobian matrix of each transformation is constant, specifically:

$$\left| \det \frac{\partial f_\theta(\mathbf{z})}{\partial \mathbf{z}} \right| = 1$$

This implies that the transformation preserves volume in the latent space, meaning the model can rearrange the probability mass but cannot scale it.

A set of probability distributions \mathcal{P} is called a distributional universal approximator if for every possible target distribution $p(x)$ there is a sequence of distributions $p_n(x) \in \mathcal{P}$ such that $\lim_{n \rightarrow \infty} p_n(x) \sim p(x)$ with a chosen type of convergence.

We want to prove: The family of Normalizing Flows with constant Jacobian determinant $\left| \det \frac{\partial f_\theta(\mathbf{z})}{\partial \mathbf{z}} \right| = \text{const}$ with latent $z \sim \mathcal{N}(0, 1)$ is not a universal distribution approximator under KL divergence. We do this by finding a lower bound for the KL objective. Use the Pinsker's inequality as below and the counter example given:

$$\delta(p_\theta, q) = \sup_{E: \text{ a measurable event}} \{ |p_\theta(E) - q(E)| \} \leq \sqrt{\frac{1}{2} D_{\text{KL}}(p\|q)}$$

counter example:

$$p(x, y) = \begin{cases} 0.9 & \text{if } (x, y) \in [-0.5, 0.5] \times [-0.5, 0.5], \\ 0.9 - k \cdot (|x| - 0.5) & \text{if } |x| \in [0.5, 0.9k + 0.5] \text{ and } |y| \in [0, |x|], \\ 0.9 - k \cdot (|y| - 0.5) & \text{if } |y| \in [0.5, 0.9k + 0.5] \text{ and } |x| \in [0, |y|], \\ 0 & \text{otherwise.} \end{cases}$$

$$A = \{ (x, y) \in \mathbb{R}^2 : p_\theta(x, y) \geq 0.9 - \epsilon \}$$

$$B = [-0.5, 0.5] \times [-0.5, 0.5]$$

$$\bar{A} = B \setminus A$$

2.1 Part (a)

Prove there exist an event E such that $|p_\theta(E) - p(E)| > 0$ when $A = \emptyset$, and show the lower bound for KL divergence can be positive for some $\epsilon > 0$, using Pinkster's inequality. (hint: prove \bar{A} can be such E)

2.2 Part(b)

Using the change of variables formula, create the event C in latent space resulting from A . Compare its volume with A . Then compute its volume and provide a simple upper bound for it. (hint: volume here is a 2D circle because of the normal distribution)

2.3 Part(c)

Now show that there is an event E such that $|p_\theta(E) - p(E)| > 0$ when $A \neq \emptyset$, and show the lower bound for KL divergence can be positive for some $\epsilon > 0$, using Pinkster's inequality

2.4 Part(d)

Now briefly explain why this means VP-NFs can't approximate every distribution.

2.1)

Note that $A = \emptyset \Rightarrow$ this means:

i) $\forall (x, y) \in \mathbb{R}^2, P_\theta(x, y) < 0.9 - \epsilon$

ii) $\bar{A} = B \setminus A \rightarrow \bar{A} = B \setminus \emptyset = B$

So, using (i) $\bar{A} = [-\frac{1}{r}, \frac{1}{r}] \times [-\frac{1}{r}, \frac{1}{r}] \rightarrow P(\bar{A}) = 19$ (Counter example)

and using (iii), $\forall (m, y) \in \mathbb{R}^2: P_\theta(m, y) < 19 - \epsilon \Rightarrow \forall (m, y) \in \bar{A}: P_\theta(m, y) < 19 - \epsilon$

so taking $E := \bar{A}$, we would have:

$|P_\theta(E) - P(E)| > |19 - (19 - \epsilon)| = \epsilon \rightarrow$ So for $E := \bar{A}$, there is

an event which $|P_\theta(E) - P(E)| > 0$ when $A = \emptyset$ ($E \neq \bar{A}$)

Now, let's see.

we know: $|P_\theta(E) - P(E)| > 0$

Also: $|P_\theta(E) - P(E)| < \delta(P_\theta, \rho) = \sup_E \{|P_\theta(E) - P(E)|\} \leq \sqrt{\frac{1}{r} D_{KL}(P_\theta \| P)}$

So we infer: $\delta(P_\theta, \rho) > \epsilon \Rightarrow \frac{1}{r} D_{KL}(P_\theta \| P) > \epsilon \Rightarrow D_{KL} \geq r \delta^2(P_\theta, \rho)$

2.2) $z \xrightarrow{f_\theta} x: P_\theta(m) = P_\theta(f_\theta^{-1}(m)) \cdot \frac{1}{|\det(J_f)|}$ (J_f symmetric)

Note that: $z \sim \mathcal{N}(0, I)$

$\forall m \in A: P_\theta(m) \geq 19 - \epsilon \Rightarrow P_z(2) \geq \frac{1}{|\det(J_f)|} \geq 19 - \epsilon$

$\rightarrow P_z(2) \geq |\det(J_f)| (19 - \epsilon) \rightarrow \log P_z(2) \geq \log |\det(J_f)| + \log(19 - \epsilon)$

$\rightarrow \log \left(\frac{1}{r^n} \exp(-\frac{1}{r} z^T z) \right) = -\log(rn) - \frac{1}{r} \|z\|^2 \geq \log |\det(J_f)| + \log(19 - \epsilon)$

$\rightarrow \|z\|^2 \leq -r(\log(rn) + \log |\det(J_f)| + \log(19 - \epsilon)) = -rn$

$m \in A: \log \left(\frac{1}{r^n} \exp(-\frac{1}{r} z^T z) \right) \geq \log(19 - \epsilon) \Rightarrow \log |\det(J_f)| \leq \log(19 - \epsilon)$

Volume of $C := V(C)$

$\rightarrow m \in A: V(C) = \pi \|z\|^2 \leq -rn\pi$

$V(A) = |\det(J_f)| V(C) \rightarrow V(A) \leq |\det(J_f)| (-rn\pi)$

$$\rightarrow U(A) \leq |\det J_f|(-\Gamma \pi) = \Gamma \pi / |\det J_f| \left(\log \frac{1}{\Gamma \pi / |\det J_f| (19 - \epsilon)} \right)$$

$$\log(19) \leq \frac{1}{e} : \frac{1}{19} \log(19) \leq \frac{1}{e}$$

$$\text{put } \alpha \leftarrow \frac{1}{\Gamma \pi / |\det J_f| (19 - \epsilon)} \text{ then:}$$

$$\Gamma \pi / |\det J_f| (19 - \epsilon) \left(\log \frac{1}{\Gamma \pi / |\det J_f| (19 - \epsilon)} \right) \leq \frac{1}{e}$$

So:

$$U(A) \leq \Gamma \pi / |\det J_f| \left(\log \frac{1}{\Gamma \pi / |\det J_f| (19 - \epsilon)} \right) \leq \frac{1}{e} \frac{1}{19 - \epsilon}$$

$$\rightarrow U(A) \leq \frac{1}{|\det J_f|} \frac{1}{e} \frac{1}{19 - \epsilon}$$

2.3)

$A \neq \emptyset$:

$$i) P_\theta(A) = U(A)(19 - \epsilon) \stackrel{(2.2)}{\leq} (19 - \epsilon) \frac{1}{(19 - \epsilon) - \epsilon} = \frac{1}{e} \approx 0.37$$

$$ii) P(A) \geq P(B) = 19$$

$$\text{So } |P_\theta(A) - P(A)| \geq 19 - \frac{1}{e} \approx 0.37$$

Assume $E = A$:

$$f(P_\theta, P) = \sup_E \{ |P_\theta(A) - P(A)| \} = 19 - \frac{1}{e} \approx 0.37$$

using Pinsker's Inequality: $\sup_E |P_\theta(A) - P(A)| \geq \epsilon > 0$

$$f(P_\theta, P) \leq \sqrt{\frac{1}{2} D_{KL}(P, P_\theta)} \rightarrow D_{KL}(P, P_\theta) \geq \epsilon^2 (P, P_\theta) = \epsilon^2 (19 - \frac{1}{e})^2$$

2.4) we saw either $A = \emptyset$ or $A \neq \emptyset$, $D_{KL}(P, P_\theta) \geq \epsilon^2 (P, P_\theta) > 0$

And $f(P, P_\theta) > 0 \rightarrow \text{UP-NFs can not approximate every distribution.}$

3 Optimality of GAN Framework

In class, we covered that the optimal discriminator satisfies the following condition

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \quad (1)$$

for arbitrary values of \mathbf{x} . In this problem, you will prove the optimality of this discriminator. To be more specific, for a fixed generator θ , show that $D_\phi^*(\mathbf{x})$ minimizes the following loss function:

$$\mathcal{L}(\phi; \theta) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log (1 - D_\phi(\mathbf{x}))] \quad (2)$$

we only need to consider what happens at each input \mathbf{x} pointwise, since our discriminator is perfectly expressive.

The relevant part of the objective for the discriminator value at a particular \mathbf{x} is

$$\mathcal{L}(\phi; \theta) = \min_{D_\phi(\mathbf{x})} - (p_{\text{data}}(\mathbf{x}) \log D_\phi(\mathbf{x}) + p_\theta(\mathbf{x}) \log (1 - D_\phi(\mathbf{x})))$$

Let's optimize this:

$$\frac{\partial \mathcal{L}}{\partial D_\phi(\mathbf{x})} = -p_{\text{data}}(\mathbf{x}) \frac{1}{D_\phi(\mathbf{x})} + p_\theta(\mathbf{x}) \cdot \frac{1}{1 - D_\phi(\mathbf{x})}$$

$$\Rightarrow p_\theta(\mathbf{x}) D_\phi^*(\mathbf{x}) = p_{\text{data}}(\mathbf{x}) (1 - D_\phi^*(\mathbf{x}))$$

$$\Rightarrow D_\phi^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_\theta(\mathbf{x})}$$

$$\frac{\partial}{\partial D_\phi(\mathbf{x})} \mathcal{L}(\phi, \theta) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{D_\phi(\mathbf{x})} \right] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\frac{1}{1 - D_\phi(\mathbf{x})} \right]$$

$$\begin{aligned} \frac{\partial^2}{\partial D_\phi^*(\mathbf{x})} \mathcal{L}(\phi, \theta) &= -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{D_\phi^*(\mathbf{x})} \right] - \mathbb{E}_{\mathbf{x} \sim p_\theta} \left[\frac{1}{(1 - D_\phi^*(\mathbf{x}))^2} \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{D_\phi^*(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta} \left[\frac{1}{(1 - D_\phi^*(\mathbf{x}))^2} \right] > 0 \end{aligned}$$

$\rightarrow \mathcal{L}(\phi, \theta)$ is convex w.r.t D_ϕ , thus D_ϕ^* is the global minima.

$$\mathcal{L}(D_\phi^*, \theta) < \mathcal{L}(D_\phi, \theta)$$

4 Perfect Discriminator, Perfect Generation

Recall that when training GAN, the gradient signals for θ come from the term

$$\mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} [\log (1 - D_\phi(\mathbf{x}))], \quad (3)$$

which can also be written as

$$L_G(\theta; \phi) = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log (1 - \sigma(h_\phi(G_\theta(\mathbf{z}))))], \quad (4)$$

where $\sigma(\cdot)$ denotes the sigmoid function, $h_\phi(\cdot)$ denotes the logits, and G_θ is the generator. Now assume that our discriminator is perfect, i.e. $\nabla_\theta L_G(\theta; \phi) \rightarrow 0$, show that the gradient of θ would vanish to zero in this case, i.e. $\nabla_\theta L_G(\theta; \phi) \rightarrow 0$.

$$\mathcal{L}_G(\theta; \phi) = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log (1 - \sigma(h_\phi(G_\theta(\mathbf{z}))))]$$

Let's compute gradient:

$$\nabla_\theta \mathcal{L}_G(\theta; \phi) = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_\theta \log (1 - \sigma(h_\phi(G_\theta(\mathbf{z}))))]$$

we shall use chain rule:

$$\nabla_\theta \log (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) = \frac{\nabla_\theta (1 - \sigma(h_\phi(G_\theta(\mathbf{z}))))}{1 - \sigma(h_\phi(G_\theta(\mathbf{z})))}$$

$$\nabla_\theta (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) = -\nabla_\theta \sigma(h_\phi(G_\theta(\mathbf{z})))$$

For Sigmoid, Gradient is $\nabla \sigma = \sigma(1 - \sigma)$

$$\nabla_\theta \sigma(h_\phi(G_\theta(\mathbf{z}))) = \sigma(h_\phi(G_\theta(\mathbf{z}))) / (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) \nabla_\theta h_\phi(G_\theta(\mathbf{z}))$$

$$\text{So we have: } \nabla_\theta \log (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) = -\frac{\sigma(h_\phi(G_\theta(\mathbf{z}))) (1 - \sigma(h_\phi(G_\theta(\mathbf{z}))))}{1 - \sigma(h_\phi(G_\theta(\mathbf{z})))} \nabla_\theta h_\phi(G_\theta(\mathbf{z})) \quad (< \infty)$$

$$\text{So, } \nabla_\theta \mathcal{L}_G(\theta; \phi) = -\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma(h_\phi(G_\theta(\mathbf{z}))) \nabla_\theta h_\phi(G_\theta(\mathbf{z}))]$$

If the discriminator is perfect, we are told that $\sigma(h_\phi(G_\theta(\mathbf{z}))) \rightarrow 0$

So, in spite of $\nabla_\theta h_\phi(G_\theta(\mathbf{z}))$, when $\sigma(h_\phi(G_\theta(\mathbf{z}))) \rightarrow 0$, the

entire expression vanishes: $\nabla_\theta \mathcal{L}_G(\theta; \phi) \rightarrow 0$