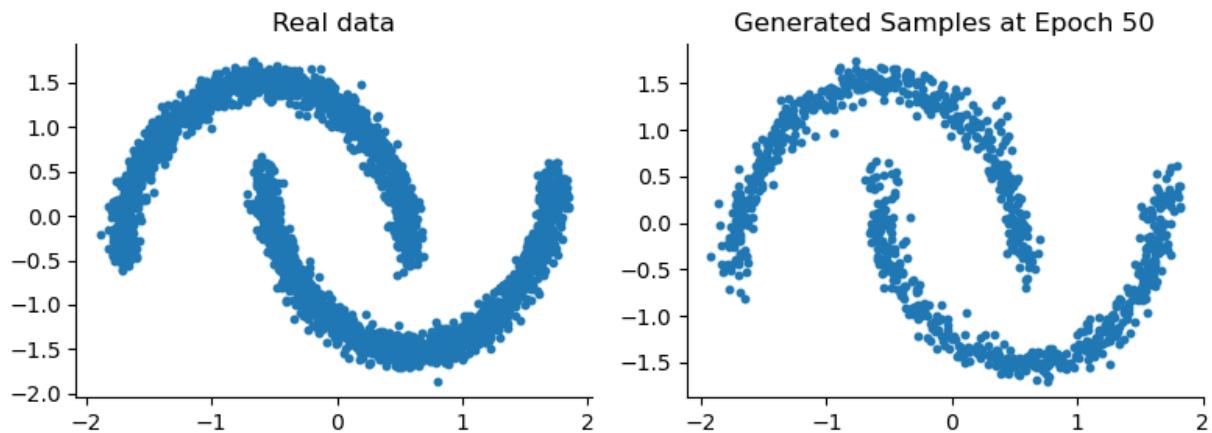


Problem 1: Flow models

4. Visualize 1000 samples drawn the model after it has been trained.



Problem 2: Generative adversarial networks

Q2) The minmax loss is given as:

$$L_a^{\text{minmax}}(\theta; \phi) = E_{z \sim N(0, I)} [\log (1 - D_\phi(h_\theta(z)))]$$

where, $h_\phi(x) = \sigma(h_\phi(x))$

and, $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

Let's compute the derivative of minmax loss w.r.t. θ

$$\begin{aligned} \frac{\partial}{\partial \theta} L_a^{\text{minmax}}(\theta; \phi) &= E_{z \sim N(0, I)} \left[\frac{\partial}{\partial \theta} \log (1 - \sigma(h_\phi(h_\theta(z)))) \right] \\ &= E_{z \sim N(0, I)} \left[\frac{1}{(1 - \sigma(h_\phi(h_\theta(z))))} \frac{\partial}{\partial \theta} (1 - \sigma(h_\phi(h_\theta(z)))) \right] \\ &= E_{z \sim N(0, I)} \left[\frac{-\sigma(h_\phi(h_\theta(z))) (1 - \sigma(h_\phi(h_\theta(z))))}{(1 - \sigma(h_\phi(h_\theta(z))))} \frac{\partial}{\partial \theta} h_\phi(h_\theta(z)) \right] \\ &= E_{z \sim N(0, I)} \left[-\sigma(h_\phi(h_\theta(z))) \frac{\partial}{\partial \theta} h_\phi(h_\theta(z)) \right] \end{aligned}$$

When discriminator successfully identifies a fake sample $h_\theta(z)$, its output $\sigma(h_\phi(h_\theta(z))) \approx 0$, which causes the derivative of the loss function to be very close to

zero (i.e. vanishing gradient). This stops the network from learning anything.

2. Implement and train a non-saturating GAN on Fashion MNIST for one epoch.



Problem 3: Divergence minimization

Q3) From the given hint:

$$f(t) = -p_{\text{data}}(x) \log t - p_{\theta}(x) \log(1-t) \quad \text{--- (i)}$$

In order to find the t which minimizes $f(t)$, we differentiate $f(t)$ with respect to 't' and then set it to zero. i.e.

$$\frac{df(t)}{dt} = 0$$

$$\Rightarrow \frac{-p_{\text{data}}(x)}{t} - \frac{p_{\theta}(x)}{(1-t)} (-1) = 0$$

$$\Rightarrow \frac{(1-t)p_{\text{data}}(x) + t p_{\theta}(x)}{t(1-t)} = 0$$

$$\Rightarrow -p_{\text{data}}(x) + t(p_{\text{data}}(x) + p_{\theta}(x)) = 0$$

$$\Rightarrow t = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\theta}(x)}$$

Hence, for $f = D_\theta(x)$, the function L_D is minimized when $D_\theta = D^*$

$$\text{where, } D^* = \frac{P_{\text{data}}(x)}{P_\theta(x) + P_{\text{data}}(x)} \quad \text{--- (i)}$$

$$\Rightarrow D_\theta(x) = \sigma(h_\theta(x))$$

$$\text{where, } \sigma(h_\theta(x)) = \frac{1}{1 + e^{-h_\theta(x)}}$$

$$\Rightarrow D_\theta(x) = \frac{1}{1 + e^{-h_\theta(x)}} \quad \text{--- (ii)}$$

Using equation (i), when $D_\theta(x) = D^*(x)$, eqn (ii) becomes

$$\frac{P_{\text{data}}(x)}{P_\theta(x) + P_{\text{data}}(x)} = \frac{1}{1 + e^{-h_\theta(x)}}$$

$$\Rightarrow e^{-h_\theta(x)} = \frac{P_\theta(x) + P_{\text{data}}(x)}{P_{\text{data}}(x)} - 1$$

$$\Rightarrow e^{-h_{\phi}(x)} = \frac{p_{\theta}(x)}{p_{\text{data}}(x)} + \frac{p_{\text{data}}(x)}{p_{\theta}(x)} - 1 = \frac{p_{\theta}(x)}{p_{\text{data}}(x)}$$

taking log on both sides of eq $\hat{=}$ (iii):

$$-h_{\phi}(x) = \log \frac{p_{\theta}(x)}{p_{\text{data}}(x)} = -\log \frac{p_{\text{data}}(x)}{p_{\theta}(x)}$$

$$\Rightarrow h_{\phi}(x) = \log \frac{p_{\text{data}}(x)}{p_{\theta}(x)}$$

3)

Generator Loss:

$$L_G(\theta; \phi) = E_{x \sim p_{\theta}(x)} [\log (1 - D_{\phi}(x))]$$

$$-E_{x \sim p_{\theta}(x)} [\log D_{\phi}(x)]$$

$$= E_{x \sim p_{\theta}(x)} \left[\log \frac{(1 - D_{\phi}(x))}{D_{\phi}(x)} \right] - \textcircled{i}$$

when, $D_\theta(x) = D^*(x)$, $D_\theta(x) = \frac{p_{\text{data}}(x)}{p_\theta(x) + p_{\text{data}}(x)}$ — (ii)

Substituting eqn (ii) in (i)

$$L_G(\theta; \phi) = E_{x \sim p_\theta(x)} \left[\log \frac{1 - \left(\frac{p_{\text{data}}(x)}{p_\theta(x) + p_{\text{data}}(x)} \right)}{\left(\frac{p_{\text{data}}(x)}{p_\theta(x) + p_{\text{data}}(x)} \right)} \right]$$

$$= E_{x \sim p_\theta(x)} \left[\log \frac{\left(\frac{p_\theta(x)}{p_\theta(x) + p_{\text{data}}(x)} \right)}{\left(\frac{p_{\text{data}}(x)}{p_\theta(x) + p_{\text{data}}(x)} \right)} \right]$$

$$= E_{x \sim p_\theta(x)} \left[\log \frac{p_\theta(x)}{p_{\text{data}}(x)} \right]$$

$$\Rightarrow L_G(\theta; \phi) = KL(p_\theta(x) || p_{\text{data}}(x))$$

4) The negative log-likelihood can be rewritten as

$$-E_{x \sim p_{\text{data}}(x)} [\log p_{\theta}(x)] = -E_{x \sim p_{\text{data}}(x)} [\log p_{\theta}(x)] \\ + E_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x)] \\ - E_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x)]$$

$$= E_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x) - \log p_{\theta}(x)]$$

$$- E_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x)]$$

$$= E_{x \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(x)}{p_{\theta}(x)} \right] - E_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x)]$$

$$= KL(p_{\text{data}}(x) \parallel p_{\theta}(x)) - E_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x)]$$

*Constant with respect
to θ*

The GAN generator objective in the previous part is $L_G(\theta; \phi) = KL(p_\theta(x) || p_{\text{data}}(x))$ and we know that the KL-divergence in general is non-symmetric.

$$\text{i.e. } KL(p_\theta(x) || p_{\text{data}}(x)) \neq KL(p_{\text{data}}(x) || p_\theta(x))$$

Hence, the VAE decoder trained with negative ELBO and a GAN generator trained with L_G are not learning the same objective.

Problem 4: Conditional GAN with projection discriminator

(Q4) From problem 3.2, we get

$$h_{\phi}^*(x, y) = \log \frac{p_{\text{data}}(x, y)}{p_{\phi}(x, y)}$$

$$= \log \frac{p_{\text{data}}(x|y) p_{\text{data}}(y)}{p_{\phi}(x|y) p_{\phi}(y)}$$

$$= \log \frac{p_{\text{data}}(x|y)}{p_{\phi}(x|y)} + \log \frac{p_{\text{data}}(y)}{p_{\phi}(y)} \quad \text{--- (1)}$$

For simplicity, we assume that $p_{\text{data}}(x) = p_0(x) = \frac{1}{m}$

$$\Rightarrow \log \frac{p_{\text{data}}(y)}{p_0(y)} = 0$$

Equation (i) then becomes:

$$h_{\varphi}(x, y) = \log \frac{p_{\text{data}}(x|y)}{p_0(x|y)}$$

$$= \log \frac{\mathcal{N}(\varphi(x) | \mu_j, I)}{\mathcal{N}(\varphi(x) | \hat{\mu}_j, I)} \quad -\text{(ii)}$$

We know that,

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

for $\Sigma = I$, we get,

$$\mathcal{N}(\varphi(x) | \mu_j, I) = \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}\|\varphi(x) - \mu_j\|^2\right) \quad -\text{(iii)}$$

and,

$$\mathcal{N}(\varphi(x) | \hat{\mu}_j, I) = \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}\|\varphi(x) - \hat{\mu}_j\|^2\right) \quad -\text{(iv)}$$

Substituting equations (iii) and (iv) in (ii)

$$h_{\theta}^*(x, y) = \log \frac{\exp(-\frac{1}{2} \|\varphi(x) - \mu_y\|^2)}{\exp(-\frac{1}{2} \|\varphi(x) - \hat{\mu}_y\|^2)}$$

$$\begin{aligned} &= -\frac{1}{2} \|\varphi(x) - \mu_y\|^2 + \frac{1}{2} \|\varphi(x) - \hat{\mu}_y\|^2 \\ &= -\frac{1}{2} \|\varphi(x)\|^2 - \frac{1}{2} \|\mu_y\|^2 - \mu_y^T \varphi(x) \\ &\quad + \frac{1}{2} \|\varphi(x)\|^2 + \frac{1}{2} \|\hat{\mu}_y\|^2 + \hat{\mu}_y^T \varphi(x) \\ &= \frac{1}{2} (\|\hat{\mu}_y\|^2 - \|\mu_y\|^2) + (\hat{\mu}_y^T - \mu_y^T) \varphi(x) \end{aligned}$$

Let y^T be a one-hot vector denoting the class y , then $(\hat{\mu}_y^T - \mu_y^T)$ can be rewritten as:

$$(\hat{\mu}_y^T - \mu_y^T) = y^T \begin{bmatrix} \hat{\mu}_1^T - \mu_1^T \\ \hat{\mu}_2^T - \mu_1^T \\ \vdots \\ \hat{\mu}_y^T - \mu_y^T \\ \vdots \\ \hat{\mu}_m^T - \mu_m^T \end{bmatrix}$$

$\|\hat{M}_y\|^2 - \|M_y\|^2$ can also be rewritten as:

$$\|\hat{M}_y\|^2 - \|M_y\|^2 = y^T \begin{bmatrix} \|\hat{M}_1\|^2 - \|M_1\|^2 \\ \|\hat{M}_2\|^2 - \|M_2\|^2 \\ \vdots \\ \|\hat{M}_y\|^2 - \|M_y\|^2 \\ \vdots \\ \|\hat{M}_m\|^2 - \|M_m\|^2 \end{bmatrix}$$

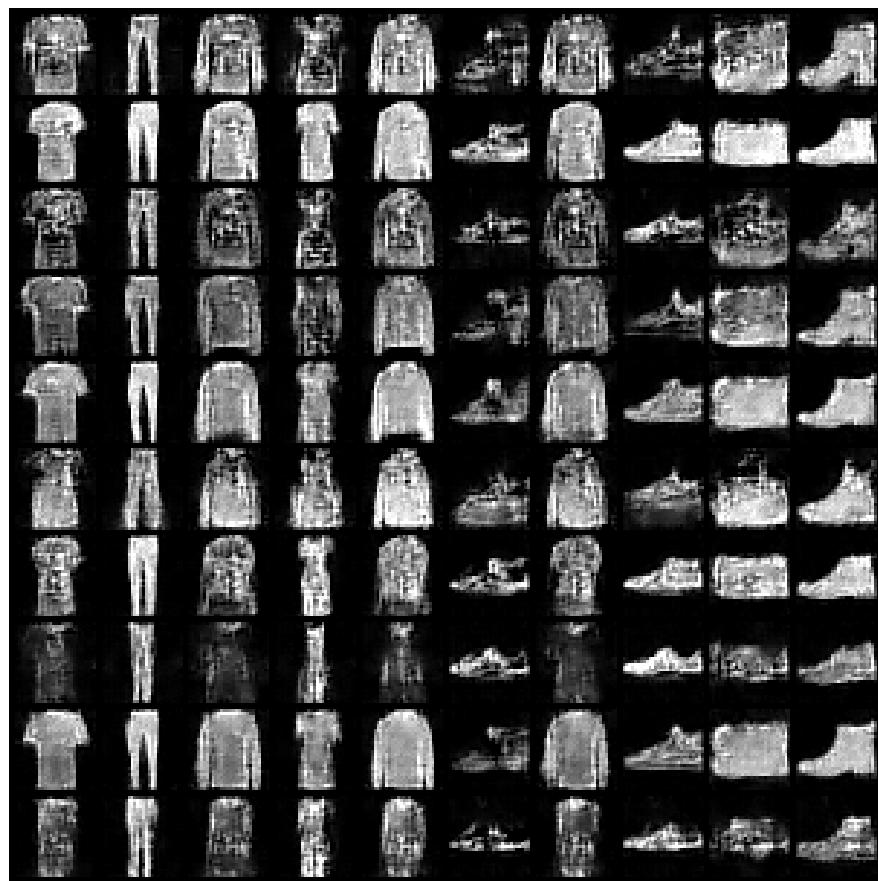
Hence, equation ① can be rewritten as

$$h_{\varphi}^*(x, y) = y^T \left(\begin{bmatrix} M_1^T - \hat{M}_1^T \\ M_2^T - \hat{M}_2^T \\ \vdots \\ M_y^T - \hat{M}_y^T \\ \vdots \\ M_m^T - \hat{M}_m^T \end{bmatrix} \varphi(x) + \begin{bmatrix} \|\hat{M}_1\|^2 - \|M_1\|^2 \\ \|\hat{M}_2\|^2 - \|M_2\|^2 \\ \vdots \\ \|\hat{M}_y\|^2 - \|M_y\|^2 \\ \vdots \\ \|\hat{M}_m\|^2 - \|M_m\|^2 \end{bmatrix} \right)$$

A b

$$\Rightarrow h_{\varphi}^*(x, y) = y^T (A \varphi(x) + b)$$

2. Implement and train a conditional GAN on Fashion MNIST for one epoch.



Problem 5: Wasserstein GAN

$$\text{Q5) } \Downarrow \quad p_{\theta}(x) = N(x | \theta, \varepsilon^2)$$

$$p_{\text{data}}(x) = N(x | \theta_0, \varepsilon^2)$$

$$\Rightarrow KL(p_{\theta}(x) || p_{\text{data}}(x)) = E_{x \sim p_{\theta}(x)} \left[\log \frac{p_{\theta}(x)}{p_{\text{data}}(x)} \right]$$

$$= E_{x \sim N(\theta, \varepsilon^2)} \left[\log \frac{N(x | \theta, \varepsilon^2)}{N(x | \theta_0, \varepsilon^2)} \right]$$

We know that,

$$N(x; \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

$$\Rightarrow KL(p_{\theta}(x) || p_{\text{data}}(x)) =$$

$$E_{x \sim N(\theta, \varepsilon^2)} \left[\log \frac{\frac{1}{(2\pi)^{m/2} \varepsilon^{1/2}} \exp \left(-\frac{1}{2\varepsilon^2} (x-\theta)^2 \right)}{\frac{1}{(2\pi)^{m/2} \varepsilon^{1/2}} \exp \left(-\frac{1}{2\varepsilon^2} (x-\theta_0)^2 \right)} \right]$$

$$= E_{x \sim N(\theta, \varepsilon^2)} \left[\log \frac{\exp \left(-\frac{1}{2\varepsilon^2} (x-\theta)^2 \right)}{\exp \left(-\frac{1}{2\varepsilon^2} (x-\theta_0)^2 \right)} \right]$$

$$\begin{aligned}
&= E_{x \sim N(\theta, \varepsilon^2)} \left[-\frac{1}{2\varepsilon^2} (x-\theta)^2 + \frac{1}{2\varepsilon^2} (x-\theta_0)^2 \right] \\
&= E_{x \sim N(\theta, \varepsilon^2)} \left[\frac{1}{2\varepsilon^2} (-x^2 - \theta^2 + 2x\theta + x^2 + \theta_0^2 - 2x\theta_0) \right] \\
&= E_{x \sim N(\theta, \varepsilon^2)} \left[\frac{1}{2\varepsilon^2} (\theta_0^2 - \theta^2 + 2x(\theta - \theta_0)) \right] \\
&= \frac{1}{2\varepsilon^2} (\theta_0^2 - \theta^2) + \frac{1}{2\varepsilon^2} E_{x \sim N(\theta, \varepsilon^2)} [2x(\theta - \theta_0)] \\
&= \frac{1}{2\varepsilon^2} (\theta_0^2 - \theta^2 + 2\theta(\theta - \theta_0)) \\
&= \frac{1}{2\varepsilon^2} (\theta_0^2 - \theta^2 + 2\theta^2 - 2\theta\theta_0) \\
&= \frac{1}{2\varepsilon^2} (\theta - \theta_0)^2 \\
\Rightarrow \text{KL}(p_\theta(x) || p_{\text{data}}(x)) &= \frac{(\theta - \theta_0)^2}{2\varepsilon^2}
\end{aligned}$$

2) If $\theta \neq \theta_0$ and $\varepsilon \rightarrow 0$, the KL-divergence and its derivative will go to infinity. This will result in a very unstable generator training since the gradients obtained by the optimizer will be large.

3) - the loss function, $L_D(\phi; \theta)$, is unbounded.

As $\varepsilon \rightarrow 0$, if $\theta \neq \theta_0$, there will not be an overlap between the distributions.

In the process of minimizing L_D , the optimizer can set $D_\phi(\theta) = -\infty$ and $D_\phi(\theta_0) = \infty$, thereby L_D approaching $-\infty$. Hence, there is no discriminator D_ϕ that minimizes this new objective.

$$L_D = D_\phi(\theta) - D_\phi(\theta_0)$$

derivative of D_ϕ (same as the slope) is restricted to be between -1 and $+1$

When D_ϕ is set to be either -1 or $+1$, that is when we achieve the smallest $D_\phi(\theta)$

This can be done by connecting two points at $(\theta_0, k) \neq (0, k - |\theta - \theta_0|)$

5. Implement and train WGAN-GP for one epoch on Fashion MNIST.

