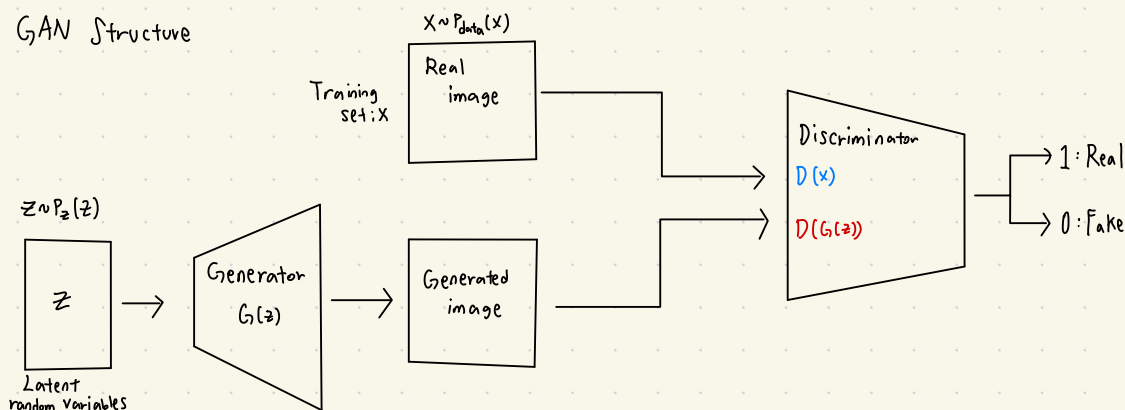


# GAN Structure



## GAN: Two Player game

- Generator: try to fool the discriminator by generating real-looking images
- Discriminator: try to distinguish between real and fake images
- GAN eventually minimizes the distance between the real data distribution and the model distribution.

Objective function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim P_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim P_z(z)} [\log (1 - D(G(z)))]$$

↗ Jensen-Shannon Divergence

Objective function of GANs is actually equal to  $\min_{G,D} \text{JSD}(P_{data} || P_{gen})$

$$\star \text{JSD}(P || Q) = \frac{1}{2} \text{KL}(P || M) + \frac{1}{2} \text{KL}(Q || M)$$

$$\text{where } M = \frac{1}{2}(P + Q)$$

Proof:

1) Fix  $G$ , and obtain optimal  $D^*$

$$D^*(x) = \arg\max_D V(D) = \mathbb{E}_{x \sim P_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim P_z(z)} [\log (1 - D(G(z)))]$$

$$= \mathbb{E}_{x \sim P_{data}(x)} [\log D(x)] + \mathbb{E}_{x \sim P_G(x)} [\log (1 - D(x))]$$

\*  $z \rightarrow G(z) \rightarrow x \rightarrow D(x)$

$$= \int P_{data}(x) \log D(x) dx + \int P_G(x) \log (1 - D(x)) dx$$

$$= \int P_{data}(x) \log D(x) + P_G(x) \log (1 - D(x)) dx$$

$$D^*(x) = \operatorname{Argmax}_D V(D) = \operatorname{argmax}_D P_{\text{data}}(x) \log(D(x)) + P_g(x) \log(1-D(x))$$

Substitute  $a = P_{\text{data}}(x)$ ,  $y = D(x)$ ,  $b = P_g(x)$

Then,  $a \log(y) + b \log(1-y)$

Differentiate w.r.t  $D(x)$ , then

$$\frac{a}{y} + \frac{-b}{1-y} = \frac{a-(a+b)y}{y(1-y)}$$

To get  $D^*$ , we need  $\frac{a-(a+b)y}{y(1-y)} = 0$ . And when  $y = \frac{a}{a+b}$ ,  $\frac{a-(a+b)y}{y(1-y)} = 0$

Hence,  $D^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}$

$$\log\left(1 - \frac{2}{2+3}\right) = \log\left(\frac{3}{2+3}\right)?$$

$$\log\left(1 - \frac{2}{5}\right) = \log\left(\frac{3}{5}\right)$$

2) Fix  $D$  to  $D^*$ , train Generator.

$$\min_G \max_D V(D, G) = \min_G V(D^*, G)$$

$$\min_G V(D^*, G) = \mathbb{E}_{x \sim P_{\text{data}}(x)} [\log D^*(x)] + \mathbb{E}_{x \sim P_g(x)} [\log (1 - D^*(x))]$$

$$= \int P_{\text{data}}(x) \cdot \log\left(\frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx + \int P_g(x) \log\left(\frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx$$

$\uparrow = 2 \log 2 = \log 2 + \log 2$

Since  $1 - \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} = \frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)}$

$$= -\log 4 + \log 4 + \int P_{\text{data}}(x) \log\left(\frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx + \int P_g(x) \log\left(\frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx$$

$$= -\log 4 + \int P_{\text{data}}(x) \log\left(\frac{2 \cdot P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx + \int P_g(x) \log\left(\frac{2 \cdot P_g(x)}{P_{\text{data}}(x) + P_g(x)}\right) dx$$

$$= -\log 4 + KL(P_{\text{data}}(x) \parallel \frac{P_{\text{data}}(x) + P_g(x)}{2}) + KL(P_g(x) \parallel \frac{P_{\text{data}}(x) + P_g(x)}{2}), \quad \star KL(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$\star JSD(P \parallel Q) = \frac{1}{2} KL(P \parallel M) + \frac{1}{2} KL(Q \parallel M)$

$= -\log 4 + 2 JSD(P_{\text{data}} \parallel P_g)$

$\Rightarrow P_{\text{data}}(x) = P_g(x)$ , then  $D^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} = \frac{1}{2}$  Hence, objective function of GAN = JSD

① Training Discriminator:  $\min_G \max_D \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))]$

$\downarrow$   
 fixed  
 $\downarrow$   
 When  $D(x)=1$ ,  
 $\Rightarrow \log(1) : \text{maximum}$

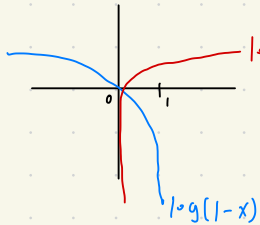
$\downarrow$   
 When  $D(G(z))=0$ ,  
 $\Rightarrow \log(1) : \text{maximum}$

② Training Generator:  $\min_G \max_D \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))] = \max_G \mathbb{E}_{z \sim p_z(z)} [\log (D(G(z)))]$

$\downarrow$   
 min  
 $G$

$\downarrow$   
 fixed  
 When  $D(G(z))=0$   
 $\Rightarrow \log(0) : \text{minimum}$

$\downarrow$   
 When  $D(G(z))=1$   
 $\Rightarrow \log(1) : \text{maximum}$



In most case, Discriminator > Generator at first, which lead  $D(G(z))=0$  more.  
 When  $x=0$ ,  $\log(x)$  has steeper slope than  $\log(1-x)$ , which enable Generator to be trained efficiently.

★ Binary Cross Entropy: Measuring the error of reconstruction in for example, an auto-encoder. Note that the targets 'y' should be numbers between 0 and 1.

$$l(x, y) = L = \{l_1, l_2, \dots, l_N\}^T = -W_N [y_n \cdot \log x_n + (1 - y_n) \log (1 - x_n)]$$

Then, for Discriminator,  $\rightarrow$  minimizing loss function

$$L_D = L_{D_{\text{real}}} + L_{D_{\text{fake}}} = l(D(x), \text{torch.ones\_like}(D(x))) + l(D(G(z)), \text{torch.zeros\_like}(D(G(z))))$$

For Generator,

$$L_{\text{gen}} = l(D(G(z)), \text{torch.ones\_like}(D(G(z))))$$