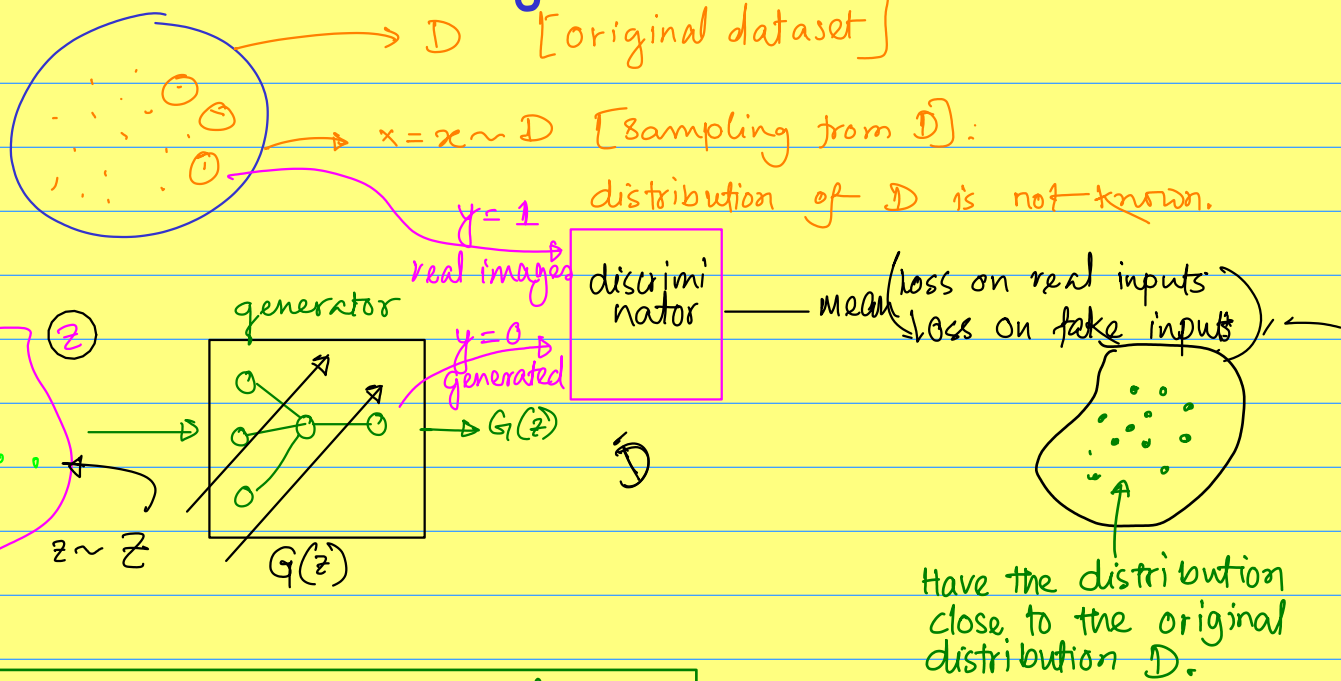
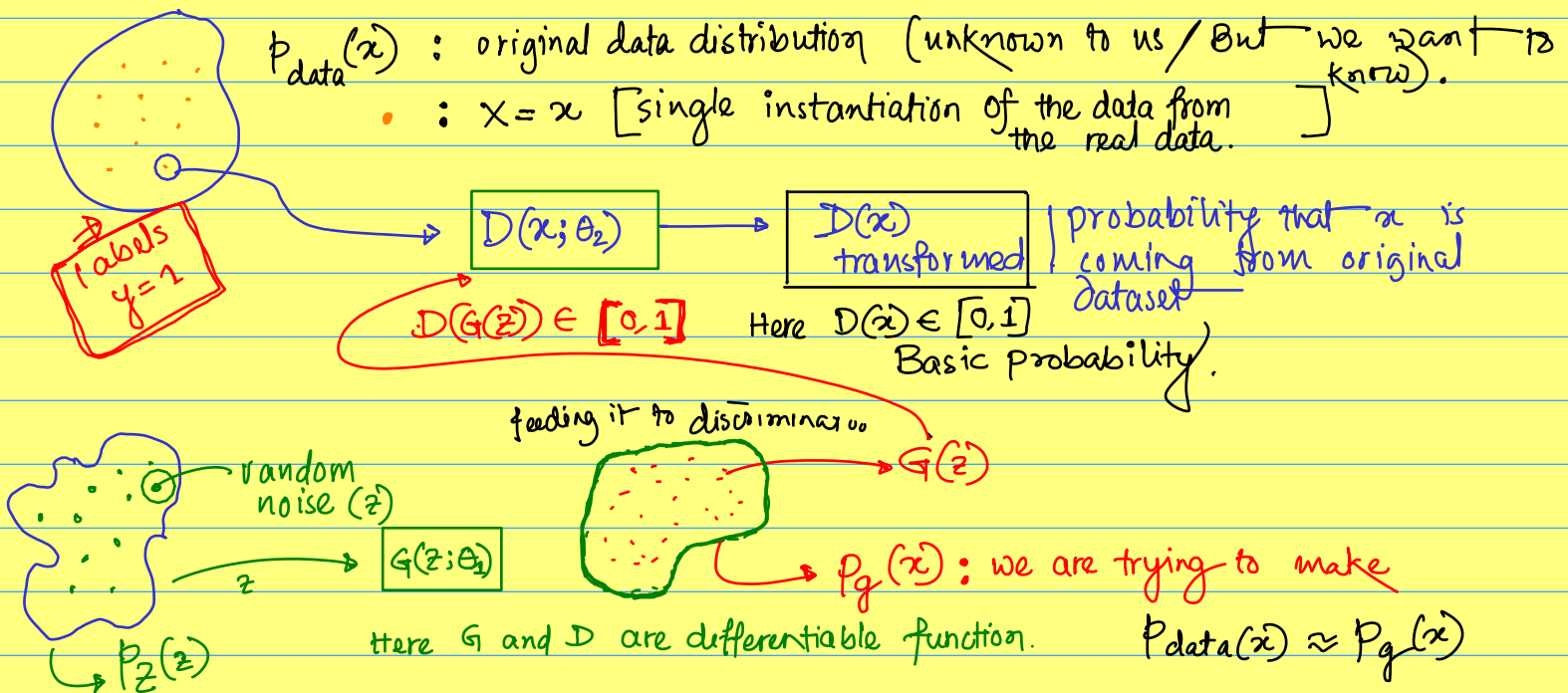


# Understanding Loss Function of GAN.



## Basic Conventions to understand loss function:



## Loss function (Binary Cross Entropy)

$$L(\hat{y}, y) = [y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

### Loss contribution from real images:

Label for the Original data coming from  $p_{data}(x)$  will be  $y=1$  and for the output produced by the discriminator  $D(x)$  the labels will be  $\hat{y} = D(x) = 1$

So the amount of loss from the real images will be

$$L(\hat{y}, y) = L(D(x), 1) = [1 \cdot \log D(x) + 0 \cdot \log (1-D(x))] = \log(D(x))$$

## ⊕ Loss contribution from the generated images:

For data coming from the generator  $G(z)$ , the ground truth is  $y=0$  and prediction  $\hat{y} = D(G(z))$ . Therefore

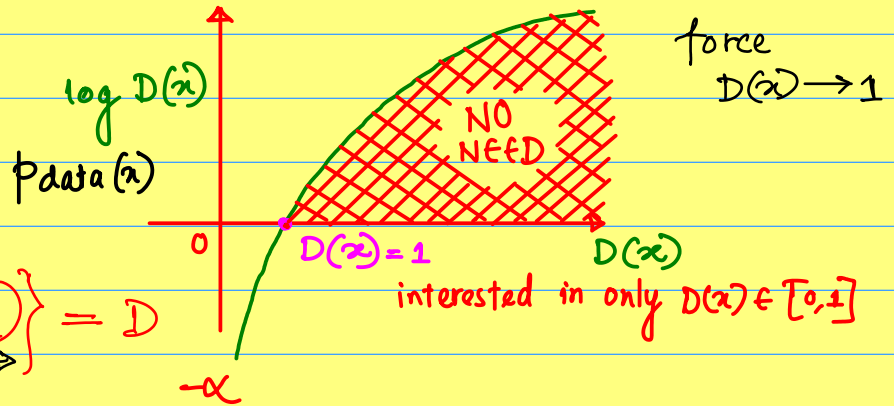
$$L(D(G(z)), 0) = [(1-0) \log(1 - D(G(z)))] \\ = \log(1 - D(G(z))) \quad \text{--- (B)}$$

⊕ Objective of the Discriminator: We want to maximize both the A and B

(A)  $[\log(D(x))]$

(B)  $[\log(1 - D(G(z)))]$

log plot:



Therefore the objective function

$$\max \left\{ \log(D(x)) + \log(1 - D(G(z))) \right\} = D$$

⊕ Objective of the generator: fool the discriminator

$$D(G(z)) \rightarrow 1 \quad \text{[tentative move]}$$

(target of discriminator)

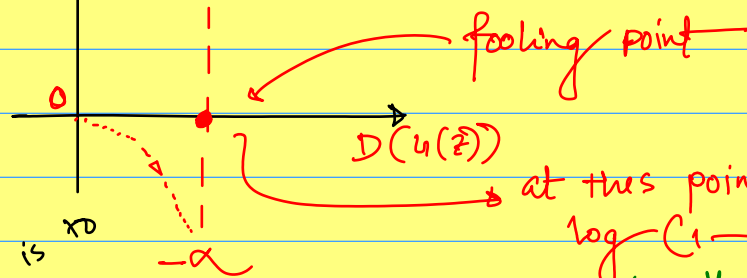
and force  $D(G(z)) \rightarrow 0$

Not at all dependent on the generator.

Look at the second term.

$$\log(1 - D(G(z)))$$

if  $D(G(z)) = 1$  then it is fooled.



So here for generator the aim is min

$$\log(D(x)) + \log(1 - D(G(z)))$$

there is no gen term here.

produce data to the original data distribution.

Therefore combining both of the optimization task we get —

$$\min_G \max_D \left\{ \log(D(x)) + \log(1 - D(G(z))) \right\} \quad \text{--- (C)}$$

Important:  $\oplus$  when generator is considered <sup>(minimized)</sup> only the  $\log(1 - D(G(z)))$  term will be considered. i.e. only the parameters of  $G(z)$  will be updated.

$\oplus$  when the discriminator is <sup>(maximized)</sup> optimized, then both of the term will be considered. i.e. Both parameters of  $D$  and  $G$  will be updated.

The above equation (C) is only for a single value or point. But for all the points in the dataset of  $n$  points. It will be something like

$$\min_G \max_D \frac{1}{m} \sum_{i=1}^m \log(D(x_i)) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z_i)))$$

which can be re written as the expectation value

$$\min_G \max_D V(D, G) = \underbrace{\min_G \max_D}_{\substack{\min \max \\ G \quad D}} \left\{ E_{x \sim P_{\text{data}}(x)} \log(D(x)) + E_{z \sim P_z(z)} \log(1 - D(G(z))) \right\}$$

## Intuitive understanding of loss function

# Finding the best discriminator:

Proposition: For a fixed  $G$ , the optimal discriminator  $D$  is

$$D_G^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} \quad \text{optimal/best.}$$

Proof: the training criterion for the discriminator  $D$  given any generator  $G$  is to maximize the equation (A) which is

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

if  $G$  is given, that means the generator is given i.e. fixed, then the goal is to maximize the remaining discriminator.

Hence the optimal discriminator is defined as

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

here please note that  $E_{P(x)}[x] = \int x p_x(x) dx$  — (1)

Expectation of the variable  $x$  having the probability density function  $P(x)$ .

Therefore the optimal discriminator will be

$$D_G^* = \arg \max_D \left\{ E_{x \sim P_{data}(x)} \log(D(x)) + E_{z \sim P_z(z)} [\log(1 - D(G(z)))] \right\}$$

$$= \arg \max_D \{ V(D, G) \} \quad \text{what value of } D \text{ will maximize } V(D, G). \quad \text{--- (2)}$$

# Expectation of any random variable  $x$  with probability density function pdf of  $P_x(x)$  is denoted as

$$E_{P(x)}[x] = \int x \cdot P_x(x) dx \quad \text{--- (1)}$$

Putting the equation (1) in the equation (2) we get

$$V(G, D) = \int_x P_{data}(x) \cdot \log(D(x)) dx + \int_z P_z(z) \log(1 - D(G(z))) dz$$

But the original paper directly jumped from the above equation to the following equation

$$V(G, D) = \int_{\mathbf{z}} P_{\mathbf{z}}(\mathbf{z}) \log D(\mathbf{z}) d\mathbf{z} + \int_{\mathbf{x}} P_{\mathbf{g}}(\mathbf{x}) \cdot \log(1 - D(\mathbf{x})) d\mathbf{x}.$$

Therefore we need to find out — how the following is possible

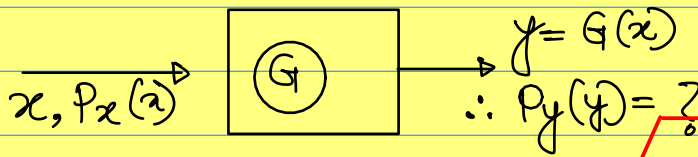
$$\int_{\mathbf{z}} P_{\mathbf{z}}(\mathbf{z}) \log D(G(\mathbf{z})) d\mathbf{z} = \int_{\mathbf{x}} P_{\mathbf{g}}(\mathbf{x}) \cdot \log(1 - D(\mathbf{x})) d\mathbf{x}$$

Probability Concept: If the probability density function (PDF) of a random variable  $\mathbf{x}$  is given as  $P_{\mathbf{x}}(\mathbf{x})$ , it is possible to calculate the probability density function of some other variable  $\mathbf{y} = G(\mathbf{x})$ . This is called "change of variable" and is defined as

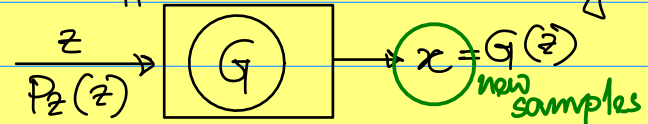
$$P_{\mathbf{y}}(\mathbf{y}) = P_{\mathbf{x}}(G^{-1}(\mathbf{y})) \left| \frac{d}{dy} (G^{-1}(\mathbf{y})) \right|$$

important.

Here  $\mathbf{y} = G(\mathbf{x}) = \mathbf{x}^2$   
 $= \mathbf{x}^2 + 4$   
 $= \mathbf{x}^3 + \mathbf{x} + 1$   
 it can be any function.



Now in our case the generator takes the input  $\mathbf{z}$ . Therefore the system is like the following.



Therefore,

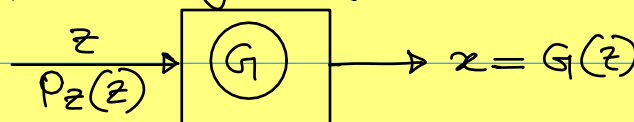
$$P_{\mathbf{g}}(\mathbf{x}) = P_{\mathbf{z}}(G^{-1}(\mathbf{x})) \cdot \frac{d}{dx} (G^{-1}(\mathbf{x})).$$

our case

We need to show that

$$\int_{\mathbf{z}} P_{\mathbf{z}}(\mathbf{z}) \cdot \log(1 - D(G(\mathbf{z}))) d\mathbf{z} = \int_{\mathbf{x}} P_{\mathbf{g}}(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x}$$

Since here the system is as follows:



$$\begin{aligned} \mathbf{x} &= G(\mathbf{z}) \\ \mathbf{z} &= G^{-1}(\mathbf{x}). \end{aligned}$$

Now if we assume that  $G$  is invertible, then it is true that  $\mathbf{z} = G^{-1}(\mathbf{x})$  [inverse of a function]

Therefore from the left hand side we get

$$\int_{\mathbf{z}} P_{\mathbf{z}}(\mathbf{z}) \cdot \log(1 - D(G(\mathbf{z}))) d\mathbf{z}$$

$$= \int_{\mathbf{x}} P_{\mathbf{z}}(G^{-1}(\mathbf{x})) \log(1 - D(\mathbf{x})) dG^{-1}(\mathbf{x})$$

$$= \int_{\mathbf{x}} P_{\mathbf{z}}(G^{-1}(\mathbf{x})) \cdot \log(1 - D(\mathbf{x})) \frac{dG^{-1}(\mathbf{x})}{d\mathbf{x}} d\mathbf{x}.$$

using the relation  $P_g(\mathbf{x}) = P_{\mathbf{z}}(G^{-1}(\mathbf{x})) \cdot \frac{d}{d\mathbf{x}} G^{-1}(\mathbf{x})$ . [derived earlier].

$$= \int_{\mathbf{x}} P_{\mathbf{z}}(G^{-1}(\mathbf{x})) \cdot \log(1 - D(\mathbf{x})) \cdot \frac{dG^{-1}(\mathbf{x})}{d\mathbf{x}} d\mathbf{x}$$

$$= \int_{\mathbf{x}} \boxed{P_{\mathbf{z}}(G^{-1}(\mathbf{x})) \cdot \frac{d}{d\mathbf{x}} G^{-1}(\mathbf{x})} \cdot \log(1 - D(\mathbf{x})) d\mathbf{x}$$

$P_g(\mathbf{x})$

$$= \int_{\mathbf{x}} P_g(\mathbf{x}) \cdot \log(1 - D(\mathbf{x})) d\mathbf{x} \quad \checkmark \text{ ————— } (3)$$

Now the whole loss function looks like the following:

$$V(G, D) = \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \cdot \log D(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{z}} P_{\mathbf{z}}(\mathbf{z}) \cdot \log(1 - D(G(\mathbf{z}))) d\mathbf{z}$$

$$= \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \cdot \log D(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x}} P_g(\mathbf{x}) \cdot \log(1 - D(\mathbf{x})) d\mathbf{x}$$

$$= \int_{\mathbf{x}} [P_{\text{data}}(\mathbf{x}) \cdot \log(D(\mathbf{x})) + P_g(\mathbf{x}) \cdot \log(1 - D(\mathbf{x}))] d\mathbf{x}$$

Thus the optimal  $D^*$  for a given  $G$  is obtained by maximizing  $\underline{V(G, D)}$  from above expression.

So we will find the maximum value of the integrand and choose the value at the maximum value to be the optimal value of  $D$  for a given  $G$  i.e.  $D_G^*$ . Therefore,

$$\frac{d}{dD(\mathbf{x})} [P_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) + P_g(\mathbf{x}) \log(1 - D(\mathbf{x}))] = 0$$

But here

$$\frac{d}{dD(x)} = \frac{P_{data}(x)}{\underbrace{D(x)}_{D_G^*}} - \frac{P_g(x)}{1-D(x)} = 0.$$

Formula:

$$\frac{d}{dx} \log(x) = \frac{1}{x}$$

$$\therefore \frac{P_{data}(x)}{D(x)} = \frac{P_g(x)}{1-D(x)}$$

$$\Rightarrow P_{data}(x) \cdot (1-D(x)) = P_g(x) \cdot D(x).$$

$$\Rightarrow P_{data}(x) - P_{data}(x) \cdot D(x) = P_g(x) \cdot D(x).$$

$$\Rightarrow P_{data}(x) = P_{data}(x) D(x) + P_g(x) D(x).$$

$$\Rightarrow D(x) \cdot (P_{data}(x) + P_g(x)) = P_{data}(x).$$

$$\Rightarrow \underbrace{D_G^*(x)} = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \quad \text{--- (Proved)}$$

To prove that its the maximum value we differentiate the following equation again

$$\frac{P_{data}(x)}{D(x)} - \frac{P_g(x)}{1-D(x)} = 0$$

$$\Rightarrow -\frac{P_{data}(x)}{D(x)^2} - \frac{P_g(x)}{(1-D(x))^2} < 0$$

Hence it was in fact the maximum value of the function.

Therefore the optimal  $D$  for a given  $G$  will be

$$D_G^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \quad \checkmark$$

→ optimal.

→ maximal.

It is theoretically solvable but practically impossible to solve because the a priori  $P_{data}(x)$  is not possible to calculate deterministically.



The role of the generator is to reverse that of the discriminator. i.e. minimizing the objective function. So the optimal  $G$  that minimize the loss function occurs when  $D = D_G^*$ . Therefore we get the optimal  $G^*$  as  $G^* = \arg \min_G V(D_G^*, G)$

At this point we must show that the optimization problem stated in (A) has unique solution  $G^*$  and this solution satisfies the  $p_g = p_{data}$ .

$$\text{So from previous proof } D_G^* = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

Therefore substituting the optimal value of  $D_G^*$  in the equation (A) we get

$$G^* = \arg \min_G V(D_G^*, G) \\ = \arg \min_G \int_x [p_{data}(x) \cdot \log(D_G^*(x)) + p_g(x) \cdot \log(1 - D_G^*(x))] dx$$

Now again we have  $D_G^* = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$ ; Substituting this

$$\text{we get } G^* = \arg \min_G \int_x \left[ p_{data}(x) \cdot \log\left(\frac{p_{data}(x)}{p_{data}(x) + p_g(x)}\right) + p_g(x) \cdot \log\left(1 - \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}\right) \right] dx$$

$$= \arg \min_G \int_x \left[ p_{data}(x) \cdot \log\left(\frac{p_{data}(x)}{p_{data}(x) + p_g(x)}\right) + p_g(x) \cdot \log\left(\frac{p_{data}(x) + p_g(x) - p_{data}(x)}{p_{data}(x) + p_g(x)}\right) \right] dx$$

$$= \arg \min_G \int_x \left[ p_{data}(x) \cdot \log\left(\frac{p_{data}(x)}{p_{data}(x) + p_g(x)}\right) + p_g(x) \cdot \log\left(\frac{p_g(x)}{p_{data}(x) + p_g(x)}\right) \right] dx$$

Now we add and subtract  $(\log 2) p_{data}(x)$  and  $(\log 2) p_g(x)$  in the above equation, then we get

$$G^* = \arg \min_G \int_x \left[ (\log 2 - \log 2) p_{data}(x) + p_{data}(x) \log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} + (\log 2 - \log 2) p_g(x) + p_g(x) \log \frac{p_g(x)}{p_{data}(x) + p_g(x)} \right] dx$$



$$G^* = \arg \min_G \int_x \left[ -\log_2(P_{data}(x) + P_g(x)) + P_{data}(x) \left\{ \log_2 + \log\left(\frac{P_{data}(x)}{P_{data}(x) + P_g(x)}\right) \right\} + P_g(x) \cdot \left\{ \log \frac{P_g(x)}{P_{data}(x) + P_g(x)} + \log_2 \right\} \right] dx$$

$$G^* = \arg \min_G \log_2 \int_x (P_{data}(x) + P_g(x)) dx$$

$$+ \int_x P_{data}(x) \left[ \log_2 + \log\left(\frac{P_{data}(x)}{P_{data}(x) + P_g(x)}\right) \right] dx$$

$$+ \int_x P_g(x) \left[ \log_2 + \log\left(\frac{P_g(x)}{P_{data}(x) + P_g(x)}\right) \right] dx$$

$$G^* = \arg \min_G \log_2 (1+1)$$

$$+ \int_x P_{data}(x) \cdot \log\left(\frac{P_{data}(x)}{\frac{P_{data}(x) + P_g(x)}{2}}\right) dx$$

$$+ \int_x P_g(x) \cdot \log\left(\frac{P_g(x)}{\frac{P_{data}(x) + P_g(x)}{2}}\right) dx$$

Now for any PDF its integral is always 1.

Therefore

$$\int_x P_{data}(x) dx = 1$$

$$\text{and } \int_x P_g(x) dx = 1$$

From KL divergence equation we get

$$KL(P(x) || Q(x))$$

$$= \int_x P(x) \cdot \log \frac{P(x)}{Q(x)} dx$$

$$\therefore G^* = \arg \min_G \left\{ -\log 4 + KL\left(P_{data}(x) || \frac{P_{data}(x) + P_g(x)}{2}\right) + KL\left(P_g(x) || \frac{P_{data}(x) + P_g(x)}{2}\right) \right\}$$

One thing to look for here is that there is no term  $G$  on the right hand side to optimize. But we have  $P_g$  which is actually dependent on  $G$ .

Previously we used  $P_z(\bar{G}^{-1}(z))$

$$P_g(x)$$

$$x = G(z) \\ G^{-1}(z) = x$$

JSD = Jensen Shannon Divergence.

$$G^* = \arg \min_G \left\{ -\log 4 + 2 \text{JSD}(P_{data}(x) || P_g(x)) \right\}$$

Here we get

$$G^* = \arg \min_G \left\{ -\log 4 + 2 \text{JSD} \left( P_{\text{data}}(x) \parallel P_g(x) \right) \right\}$$

$$\text{where } \text{JSD} \left( P(x) \parallel Q(x) \right) = \frac{1}{2} \left[ \text{KL} \left( P(x) \parallel M(x) \right) + \text{KL} \left( Q(x) \parallel M(x) \right) \right]$$
$$\text{where } M(x) = \frac{P(x) + Q(x)}{2}$$

But the JSD term becomes 0 when  $P_g(x) = P_{\text{data}}(x)$ .

which is our ultimate goal but it is very hard to achieve.

Therefore in the above equation JSD is 0 only when  $P_g(x) = P_{\text{data}}(x)$ . which minimizes the argument and the value obtained is  $-\log 4$ .

$$G^* = \arg \min_G \left\{ -\log 4 + 2 \text{JSD} \left( P_g(x) \parallel P_{\text{data}}(x) \right) \right\}$$

one way to make it zero is to make  $P_g(x) = P_{\text{data}}(x)$ . the optimization process is forcing this to happen.

**Theorem:** The global minimum of the criterion  $G^* = \arg \min_G V(D_G^*, G)$  is achieved if and only if  $P_g(x) = P_{\text{data}}(x)$ . At that point  $G^*$  will have the value of  $-\log 4$

**Proof:** From previous proof we know  $D_G^* = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}$

Now if we put  $P_g(x) = P_{\text{data}}(x)$  then

$$D_G^* = \frac{1}{2}$$

Now we put it in the expression of  $G^* = \arg \min_G V(D_G^*, G)$

$$\therefore G^* = \arg \min_G V \left( \frac{1}{2}, G \right)$$

$$= \arg \min_G \left\{ \int_x \left[ P_{\text{data}}(x) \cdot \log(D(x)) + P_g(x) \cdot \log(1-D(x)) \right] dx \right\}$$
$$= \arg \min_G \left\{ \int_x \left[ P_{\text{data}}(x) \log\left(\frac{1}{2}\right) + P_g(x) \cdot \log\left(\frac{1}{2}\right) \right] dx \right\}$$

$$= \arg \min_G -\log(2) \left[ \int_x p_{data}(x) dx + \int_x p_g(x) dx \right]$$

$$= -\log(2) * [1 + 1]$$

$$= -\log 4$$

$$\int_x p_{data}(x) dx = 1 \text{ [PDF]}$$

$$\int_x p_g(x) dx = 1 \text{ [PDF]}$$