

→ refresher<sup>2,2</sup>

# GANs (Generative Adversarial Networks) (Ian Goodfellow + Yoshua Bengio)

↳ (Conflict or opposition)

GAN are deep neural net arch comprised of two neural networks, competing one against the other (i.e., why adversarial).

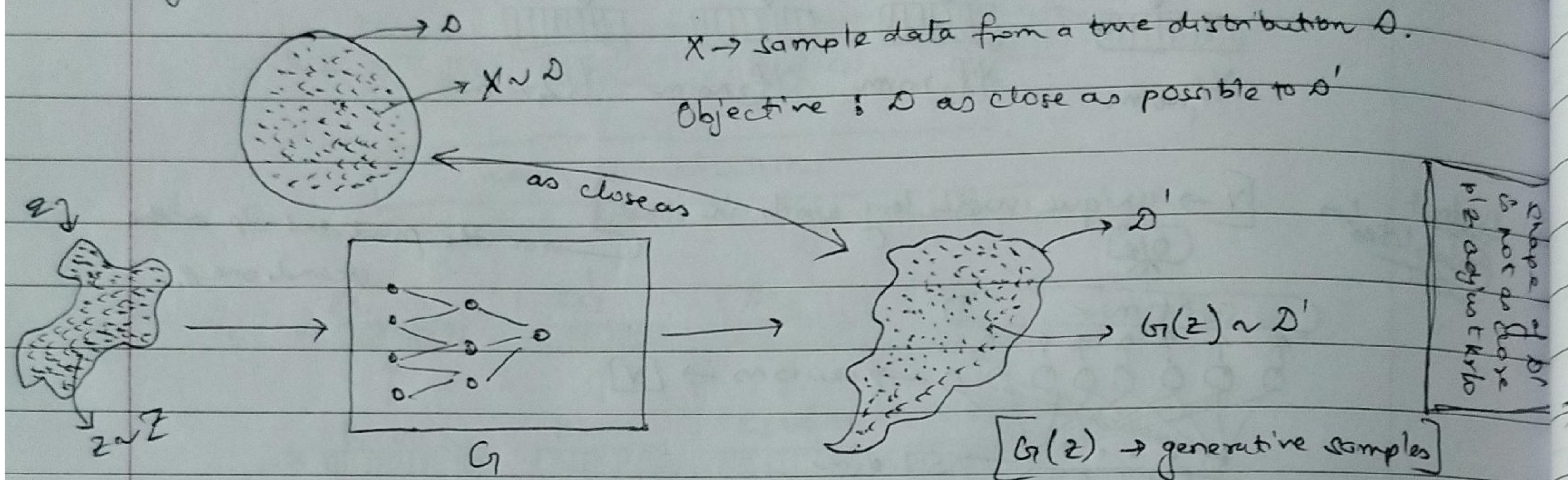
→ GAN are neural networks that are trained in an adversarial manner, to generate data mimicking some distribution.

→ Two class of models in machine learning? -

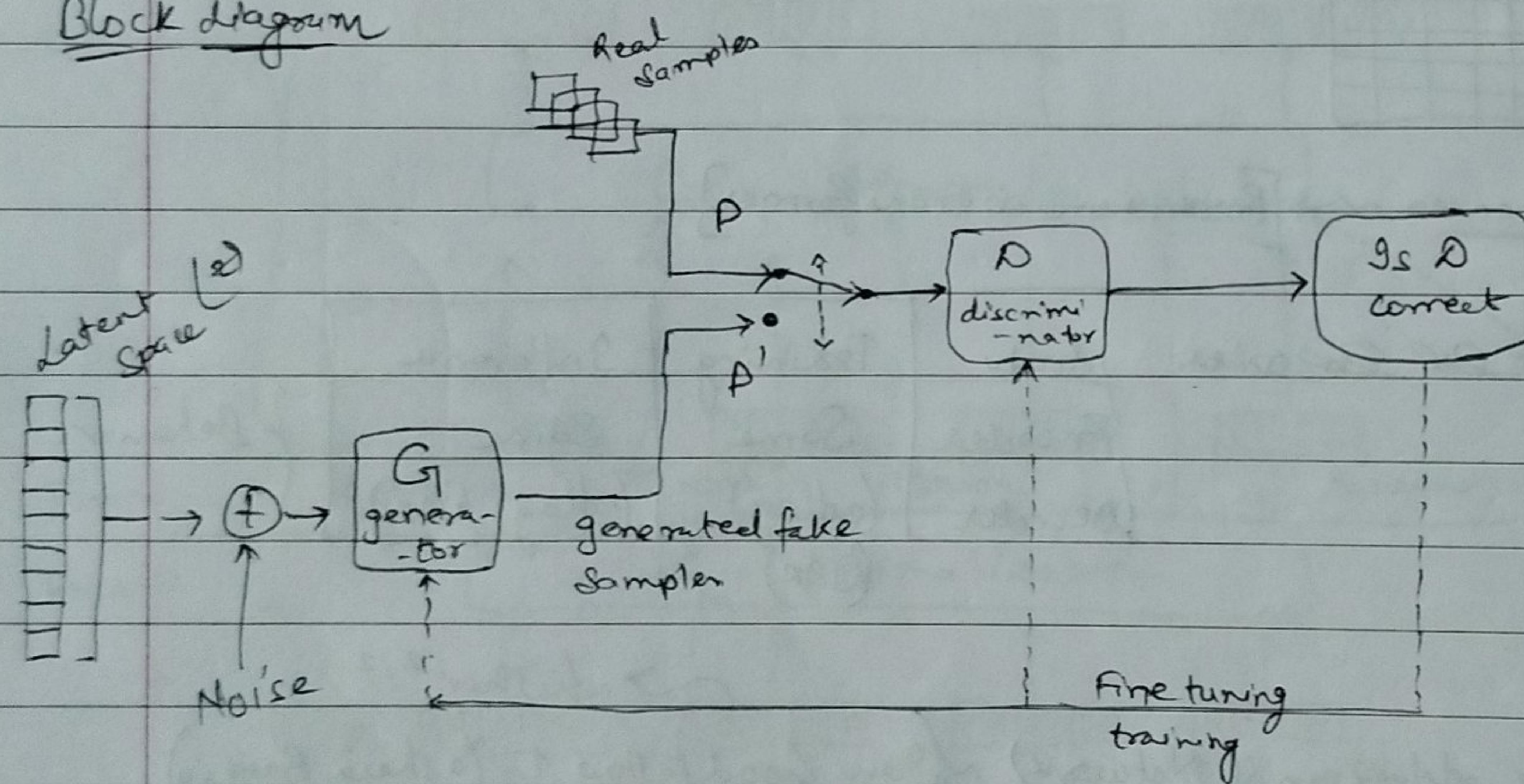
(a) Discriminative model:  $\rightarrow$  it is the one that discriminates b/w two different classes of data.

(b) Generative model :- A generative model  $G$  to be trained on training data  $X$   
assumed form some true distribution  $D$  is the one which, given some standard random  
distribution  $Z$  produces a distribution  $D'$  which is close to  $D$  according to some  
closeness metric, mathematically,

$\tilde{z} \sim \mathcal{E}$  maps to a sample  $G(\tilde{z}) \sim \mathcal{D}'$



## Block diagram

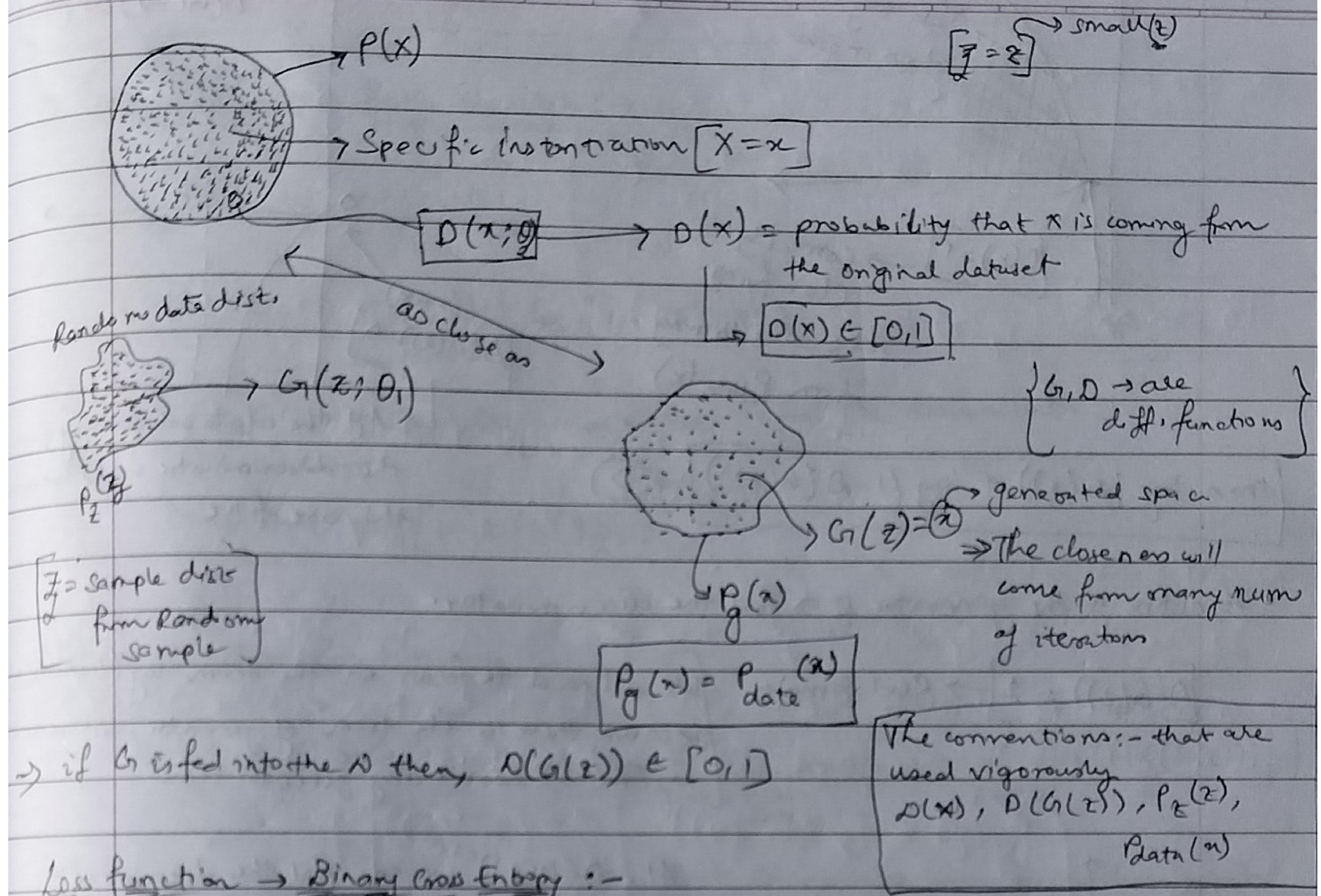


While training, we will not train our discriminator ( $w, b$ ), we can't make our discriminator too strong to beat.

- Loss function of GAN

Discriminator:- Role is to distinguish b/w the actual & fake data.

Generator:- Role is to create data in such a way so that it can feed the generator.



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

The label for the data coming from  $\text{data}(x)$  &  $y = 0(x)$  is putting this down will be:-

$$\lambda(\varphi(x), 1) = \log(\varphi(x)) \quad \text{---} \textcircled{5}$$

for data coming from generator, the label by  $y = 0$ , &  $\hat{g} = \rho(G(x))$  be in  
that case.

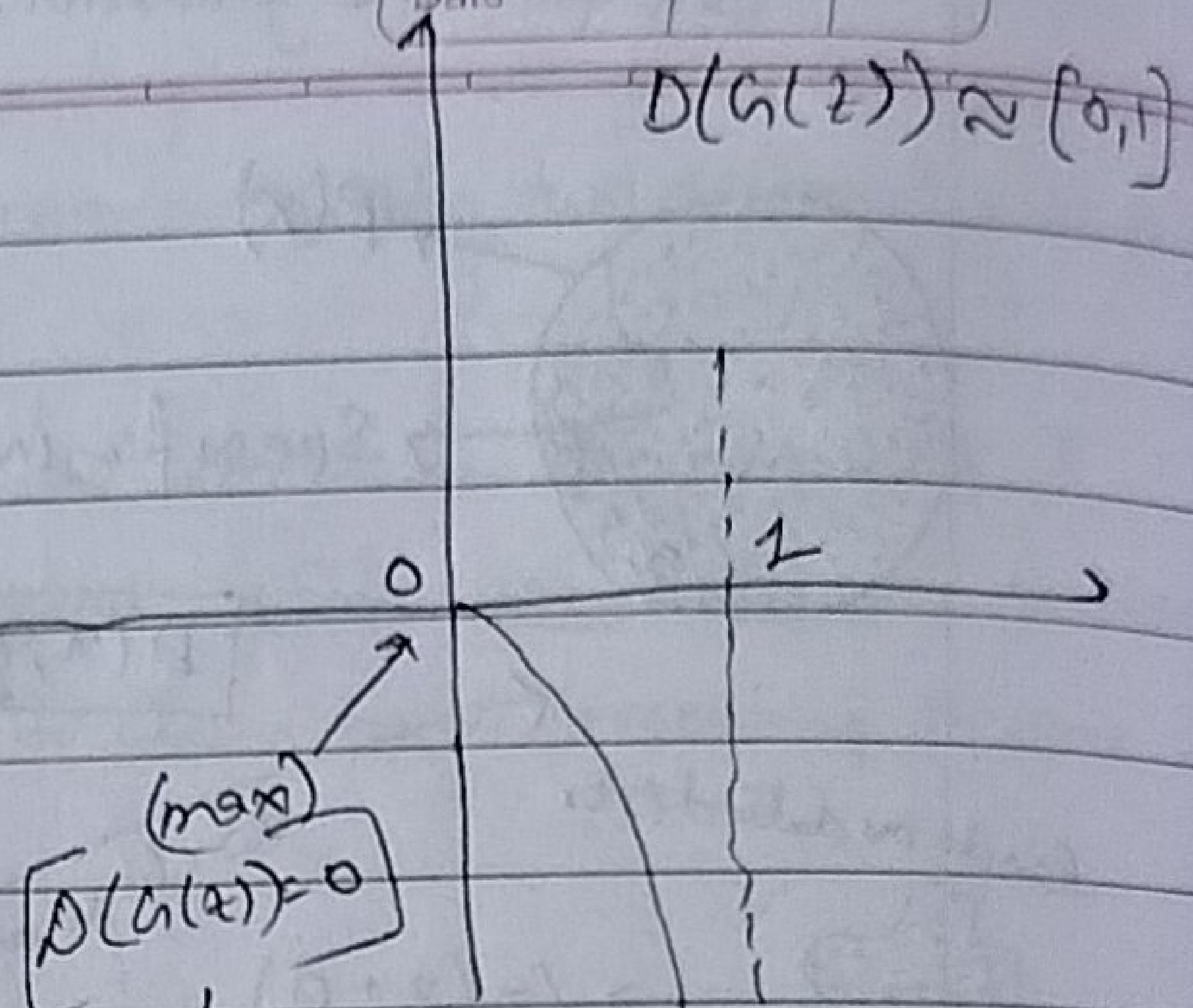
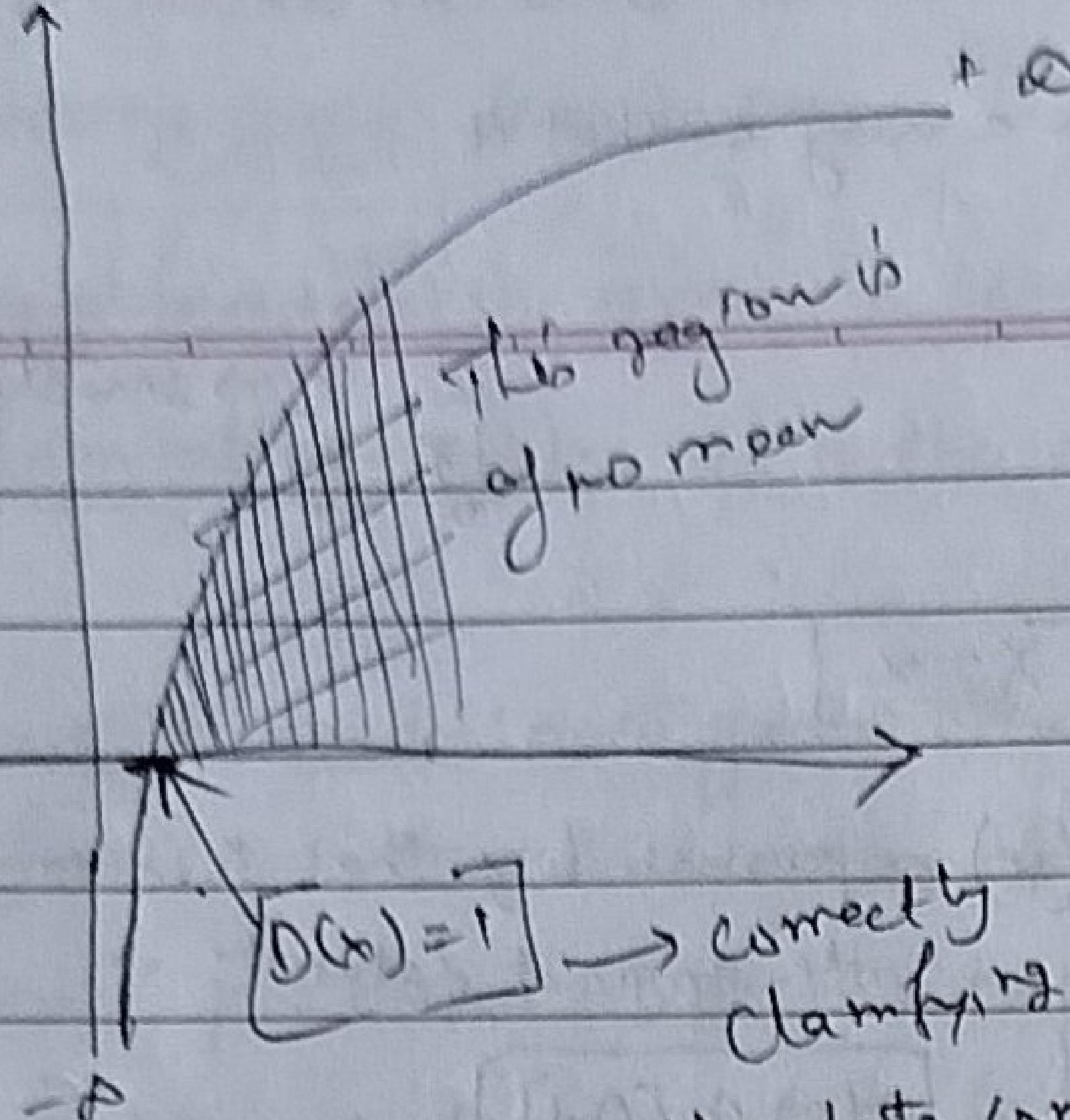
$$L(\delta G_1(z), 0) = (1 - 0) \log(1 - \delta(G_1(z))) = \log(1 - \delta(G_1(z))) \rightarrow B$$

- Objective of the discriminator is to correctly classify fake vs the Real dataset.  
For this  $A$  &  $B$  should be maximized

$$A \rightarrow \lceil \log(\rho(x)) \rceil$$

$$\log(1 - \delta(a_1(z))) \leq B$$

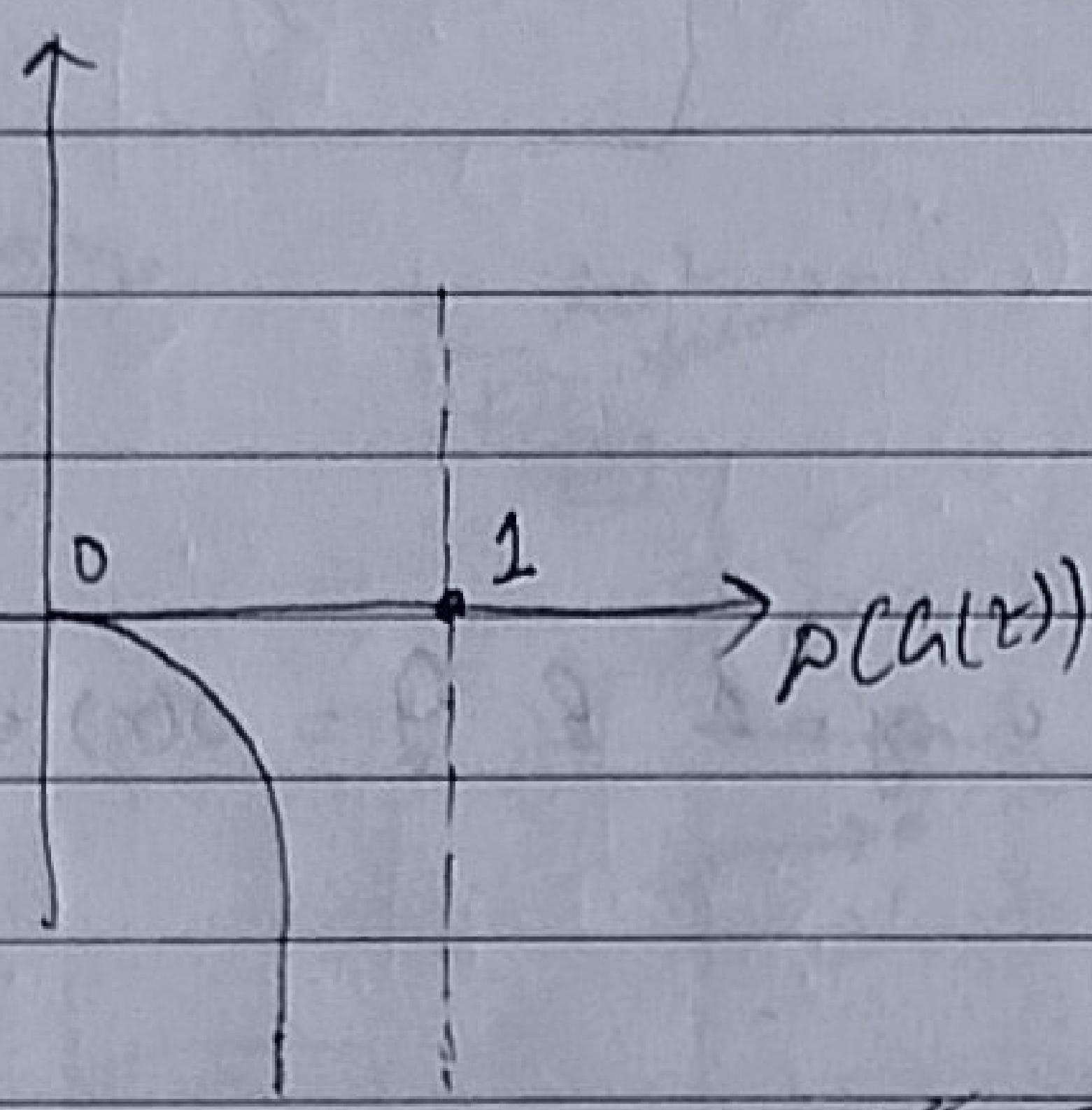
we need to maximize these two terms



$$\boxed{\max \log(D(x)) + \log(1 - D(G(z)))} = \text{⑤}$$

- Objective of the generator is to fool the discriminator.

$[D(G(z)) = 1] \rightarrow$  Real samples. We all need to focus on the term  $\log(1 - D(G(z)))$  because there is the term of generator, whereas the  $\log(D(x))$  hasn't any generator term



At,  $D(G(z)) = 1$  we can fool the discriminator.

We know, at  $D(G(z)) = 1$ , the discriminator is being fooled by generator. So at point (1) the value of  $D(G(z))$  is  $-2$  that is very less than 0. So we have to min the objective function

$$\min \boxed{[\log(D(x)) + \log(1 - D(G(z)))]}$$

→ Here  $(D(x))$  has no role however we are writing just to complete it. //

∴ by minimizing we are forcing  $[\log(1 - D(G(z)))]$  to reach at the value of 1, By this it can fool the discriminator.

$$\min_G \max_D \left\{ \left[ \log(D(x)) + \log(1 - D(G(z))) \right] \right\} \rightarrow \text{Objective Function}$$

↳ for only 1 instance of  $x$ .

So for all data, we need the Expected Value the final Objective func is

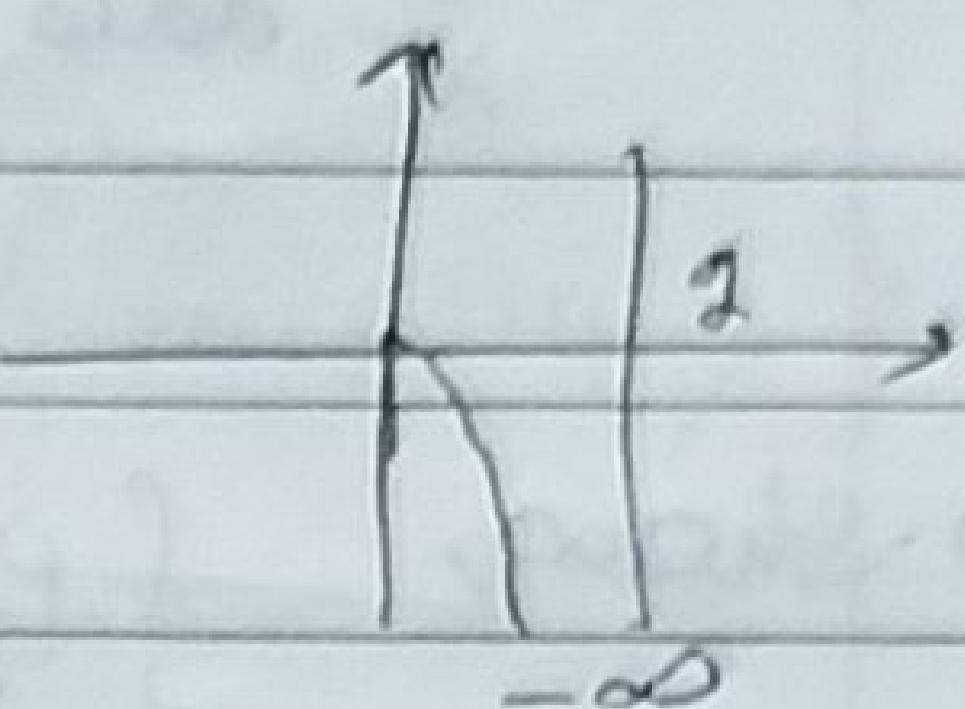
$$\min_{G_1} \max_D V(D, G_1) = \left[ \sum_{x \sim p_{\text{data}}(x)} E[\log D(x)] + \sum_{z \sim p_z(z)} E[\log(1 - D(G_1(z)))] \right]$$

Page No.

Date

④ ~~Generator~~ ~~Discriminator from scratch :-~~

$$y=0, \hat{y}=D(G_1(z))$$



$$L(D(G_1(z)), 0) = (1) \log [1 - D(G_1(z))]$$

$D(G_1(z)) = 1 \rightarrow$  our main goal to fool the discriminator

$$\min \log [1 - D(G_1(z))]$$

# Understanding Loss function :-

Ⓐ

$$\min_{G_1} \max_D V(D, G_1) = \min_{G_1} \max_D \left[ \sum_{x \sim p_{\text{data}}(x)} E[\log D(x)] + \sum_{z \sim p_z(z)} E[\log(1 - D(G_1(z)))] \right]$$

The major thing is :- How to find the best discriminator among all the discriminator that are possible for this optimization.

According to paper:- Finding the best discriminator:-

Proposition:- For  $G_1$  fixed, the optimal discriminator  $D$  is :-

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

Proof:- The training criterion for the discriminator  $D$ , given any generator  $G_1$  is to maximize eq:- Ⓐ

Hence, the optimal discriminator for a given  $G_1$  is denoted as  $D_G^*$ .

$$\therefore D_G^* = \underset{D}{\operatorname{argmax}} V(D, G_1)$$

Please Note:-

$$E_{p(x)}[x] = \int x p_x(x) dx$$

-①

This is imp.

What value  
D will max

$$D_G^* = \underset{D}{\operatorname{argmax}}$$

$$\left\{ \sum_{x \sim p_{\text{data}}(x)} E[\log(D(x))] + \sum_{z \sim p_z(z)} E[\log(1 - D(G_1(z)))] \right\}$$

②

the argument/  
that's why argmax

$$= \underset{D}{\operatorname{argmax}} V(G_1, D)$$

Applying ① to ② we get

$$V(G, D) = \int_{\text{data}} p(x) \log(D(x)) dx + \int_{\mathbb{Z}} p_z(z) \log(1 - D(G(z))) dz \quad \text{Page No. } \boxed{2} \quad \text{Date } \boxed{\phantom{00}}$$

In the paper it is written as:-

$$= \int_{\text{data}} p(x) \log(D(x)) dx + \int_{\mathbb{X}} p_g(x) \log(1 - D(x)) dx \quad ??$$

To show,  $\int_{\mathbb{Z}} p_z(z) \log(1 - D(G(z))) dz = \int_{\mathbb{X}} p_g(x) \log(1 - D(x)) dx.$

To do this we need to Recall some concepts in probability.

Recall Probability:-

If the probability density function of a random variable  $X$  is given as  $p_X(x)$  it is possible to calculate the probability density function of some variable  $Y = g(x)$  this is called as "change of variable" & is defined as

$$p_Y(y) = p_X(g^{-1}(y)) \left| \frac{d}{dy} (g^{-1}(y)) \right|$$

$$\xrightarrow[X, p_X(x)]{G} \bullet Y = g(X) \\ p_Y(y) ??$$

$G$  is the transformation function  $\times$

On our case

$$\xrightarrow[Z, p_Z(z)]{G} \bullet X = G(Z)$$

$$p_g(x) = p_z(g^{-1}(x)) \left| \frac{d}{dx} (g^{-1}(x)) \right|$$

mapping all the  $x$  are replaced by  $z$ , and all the  $y$  are replaced by  $x$

To prove/show:-  $\int_{\mathbb{Z}} p_z(z) \log(1 - D(G(z))) dz = \int_{\mathbb{X}} p_g(x) \log(1 - D(x)) dx$

$$\xrightarrow[Z, p_Z(z)]{G} \bullet X = G(Z)$$

assume  $G$  is invertible  $\Rightarrow Z = G^{-1}(X)$

lets start from LHS

$$\text{or } \int_z P_2(z) \log(1 - \sigma(G^{-1}(z))) dz = \int_x P_2(G^{-1}(x)) \log(1 - \sigma(x)) dG^{-1}(x)$$

Page No.

Date

$$= \int_x P_2(G^{-1}(x)) \log(1 - \sigma(x)) \frac{dG^{-1}(x)}{dx} dx$$

Using relation  $P_g(x) = P_2(G^{-1}(x)) \frac{d}{dx} G^{-1}(x)$  {derived earlier}

$$= \int_x P_2(G^{-1}(x)) \log(1 - \sigma(x)) \frac{dG^{-1}(x)}{dx} dx \quad \{ \text{rewriting} \}$$

$$= \int_x \left\{ P_2(G^{-1}(x)) \frac{dG^{-1}(x)}{dx} \right\} \log(1 - \sigma(x)) dx$$

$$= \int_x P_g(x) \log(1 - \sigma(x)) dx \quad \text{--- (3) Proved}$$

So, as we derived this which has the same equation in the paper

$$i. V(G, D) = \int_x \text{data} P(x) \log(P(x)) dx + \int_x P_g(x) \log(1 - \sigma(x)) dx$$

Clubbing both integrals now,

$$V(G, D) = \int_x \left[ P(\text{data}) \log(P(x)) + P_g(x) \log(1 - \sigma(x)) \right] dx$$

The optimal  $D^*$  for a given  $G$  is obtained by maximizing the argument  $V(G, D)$ .

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

$$\text{So, } \frac{d}{dD(x)} \left[ P_{\text{data}}(x) \log(P(x)) + P_g(x) \log(1 - \sigma(x)) \right] = 0$$

$$\Rightarrow \left\{ \frac{P_{\text{data}}(x)}{P(x)} - \frac{P_g(x)}{1 - \sigma(x)} = 0 \right\} \quad \text{--- (4)}$$

$$D_G^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \quad \boxed{\text{proved}} \quad (\text{optimal value})$$

We got the optimal one, now we need to get that whether it is maximum or not!!! to get that we need to do one more derivation

To prove it is a maximum value, we differentiate eqn ④

$$-\frac{P_{\text{data}}(x)}{(D(x))^2} - \frac{P_g(x)}{(1-D(x))^2} < 0. \text{ Hence maximum}$$

value is achieved for  $D$  given  $G$  as,

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

optimal as well as the  
maximum value

→ we have calculated the value of  $D$ , tho it is not practically possible but mathematically it is important. Why not practically? Because the  $P_{\text{data}}(x)$  is non-appron.

↓  
not based on theory, but  
rather on exp & observation

} It is mathematically important  
because its existence enables to  
prove an optimal value of  
generator. //

The role of generator is to reverse of that of  $D$ . i.e. min so the optimal  $G$  that minimize the loss function occurs when  $D = D_G^*$ . So we get optimal  $G^*$  as

$$G^* = \underset{G}{\operatorname{argmin}} V(D_G^*, G)$$

at this point, we must show that the optimization problem stated in (A) has a unique solution  $G^*$  & this solution satisfies  $P_g = P_{\text{data}}$ .

$$\text{So, our optimal } D_G^* = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}$$

By substituting optimal value of  $D_G^*$  in  $G^* = \underset{G}{\operatorname{argmin}} V(D_G^*, G)$  we obtain.

$$= \underset{G}{\operatorname{argmin}} \left\{ \int_x \left\{ P_{\text{data}}(x) \log \left( \frac{P_{\text{data}}^*(x)}{P_{\text{data}}(x) + P_g(x)} \right) + P_g(x) \log \left( 1 - \frac{P_{\text{data}}^*(x)}{P_{\text{data}}(x) + P_g(x)} \right) \right\} dx \right\}$$

$$= \underset{G}{\operatorname{argmin}} \left[ \int_x \left\{ P_{\text{data}}(x) \log \left( \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \right) + P_g(x) \log \left( \frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)} \right) \right\} dx \right]$$

add & subtract  $(\log 2) P_{\text{data}}(x)$  &  $(\log 2) P_g(x)$  in above eqn

We obtain:-

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

Page No.  
Date

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

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

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

$$\textcircled{*} \quad \int_x P(x) dx = 1, \quad \log \left( \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} / 2 \right), \quad \log \left( \frac{P_g(x)}{P_g(x) + P_{\text{data}}(x)} / 2 \right)$$

$$G^* = -\log_2 (1+1) + \int_x P_{\text{data}}(x) \log \left( \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} / 2 \right) +$$

$$\int_x P_g(x) \log \left( \frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)} / 2 \right).$$

KL Divergence

$$G^* = \underset{G_1}{\text{argmin}} \left\{ -\log 4 + \text{KL} \left[ P_{\text{data}}(x) \parallel \frac{P_{\text{data}} + P_g}{2} \right] + \text{KL} \left[ P_g(x) \parallel \frac{P_{\text{data}} + P_g}{2} \right] \right\}$$

must be wondering there is no  $G_1$  at RHS. But  $P_g$  is dependent on  $G_1$  (prob. derivation)

$$G^* = \underset{G}{\operatorname{argmin}} \left\{ -\log 4 + 2 \text{JSD} \left( P_{\text{data}}^{(n)} \parallel P_g^{(n)} \right) \right\}$$

Jensen Shannon Divergence

We know,  $P_g(x) = \int \int \int \dots \int \underbrace{P_g(g^{-1}(x))}_{d^n} dx$ ,  $P_g^{(n)}$  independent on  $x$

Here Jensen Shannon Divergence is:

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

$\therefore$  As we are minimizing that term if  $P_g(x) = P_{\text{data}}^{(n)}$  then the JSD will be 0 (whole term) //

min value achieved by generator will be  $\rightarrow [-\log(4)]$

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

if  $P_g(x) = P_{\text{data}}^{(n)}$

So, JSD is 0 only when  $P_g(x) = P_{\text{data}}^{(n)}$  which minimizes the argument & the value obtain is  $-\log 4$  //

Theorem: The global minimum of the criterion  $G^* = \underset{G}{\operatorname{argmin}} \mathcal{V}(D_g^*, G)$  is achieved if & only if  $P_g = P_{\text{data}}$ . At that point  $G^*$  value is  $-\log 4$

Proof: The goal of GAN process is to have  $P_g(x) = P_{\text{data}}^{(n)}$ , so substituting these values for optional  $D$  results in  $D_g^* = 1/2$

NOTE - While we were deriving we got the solution  $P_g(x) = P_{\text{data}}^{(n)}$ , We didn't said that directly //

Let's go with Goodfellow idea

Recalling,

$$V(G, D) = \underset{G}{\operatorname{argmin}} \int_x \left[ P_{\text{data}}(x) \log(D(x)) + P_g(x) \log(1 - D(x)) \right] dx$$

Page No.	
Date	

Substituting,  $D(x) = 1/2$  we get,

$$V(G, D_G^*) = \int_x \left[ P_{\text{data}}(x) \log \frac{1}{2} + P_g(x) \log \frac{1}{2} \right] dx = -\log 2 \int_x [P_{\text{data}}(x) + P_g(x)] dx$$

$$V(G, D_G^*) = -2 \log 2 = -\log 4 \quad [\because P_{\text{data}}(x) = P_g(x) \& \int_p p(x) dx = 1]$$

## # Optimization of Loss function

Let's Recall the loss function:

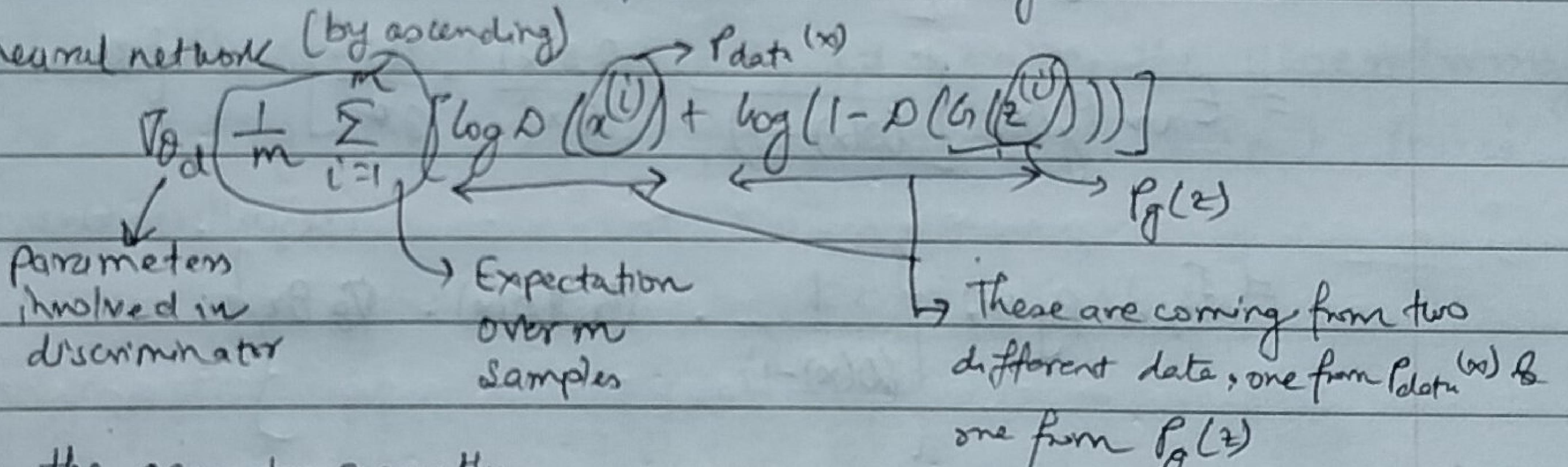
$$L_2 = \min_G \max_P \left\{ \underset{x \sim p(x)}{\mathbb{E}} [\log D(x)] + \underset{z \sim p(z)}{\mathbb{E}} [\log (1 - D(G(z)))] \right\}$$

In paper  $\Rightarrow$  Algorithm 1  $\rightarrow$  they showed that they are using Mini-batch Gradient Descent

(1) Taking  $m$  fake sample &  $m$  real samples ( $y=1$ )

(2) Update the discriminator, to update the W & B of the discriminator network

i.e. the neural network (by ascending)



(3) Train the generator separately:

→ Again, take  $m$  noise samples from  $P_g(z)$

→ Update the generator (by descending) using stochastic gradient

$$\nabla \theta_g \left( \frac{1}{m} \sum_{i=1}^m (1 - D(G(z^{(i)}))) \right)$$

↓ Parameters involved in Generator      ↓ Expectation over(m) samples.

$$\max_P \left\{ \underset{x \sim p(x)}{\mathbb{E}} [\log D(x)] + \underset{z \sim p(z)}{\mathbb{E}} [\log (1 - D(G(z)))] \right\}$$

$$\min_G \left\{ \underset{z \sim p(z)}{\mathbb{E}} [\log (1 - D(G(z)))] \right\}$$

## \* Vanishing Gradient problem

As the weights and bias of Discriminator and generator have ~~page~~ together individual.  
so, we will represent the parameters of discriminator as  $\phi$

and the parameters of the generator as  $\theta$ :

and by derivatives it (the loss func) gets to the 0, so let's see: →

$$L = \min_{G} \max_{\theta} E_{x \sim P(x) \text{ data}} \left[ \log D_{\theta}(x) \right] + E_{z \sim p(z)} \left[ \log (1 - D_{\theta}(G_{\theta}(z))) \right]$$

depends on generator only

$$\nabla_{\theta} L = \nabla_{\theta} \left[ E_{z \sim p(z)} \log (1 - D_{\theta}(G_{\theta}(z))) \right]$$

$p_z(z) \rightarrow$  Random distribution

$$= E_{z \sim p(z)} \nabla_{\theta} \left[ \log (1 - D_{\theta}(G_{\theta}(z))) \right]$$

$$= E_{z \sim p(z)} \left[ \frac{\partial G_{\theta}(z)}{\partial \theta} \cdot \frac{\partial D_{\theta}(G_{\theta}(z))}{\partial G_{\theta}(z)} \left\{ \frac{1}{D(G_{\theta}(z)) - 1} \right\} \right]$$

$$= \cancel{\nabla_{\theta} D_{\theta}(z)} \quad z = G_{\theta}(z)$$

$$P_g(x)$$

$$= E_{x \sim P_g(x)} \left[ \frac{\partial x}{\partial \theta} \frac{\partial D(x)}{\partial x} \cdot \frac{1}{(D(x)-1)} \right]$$

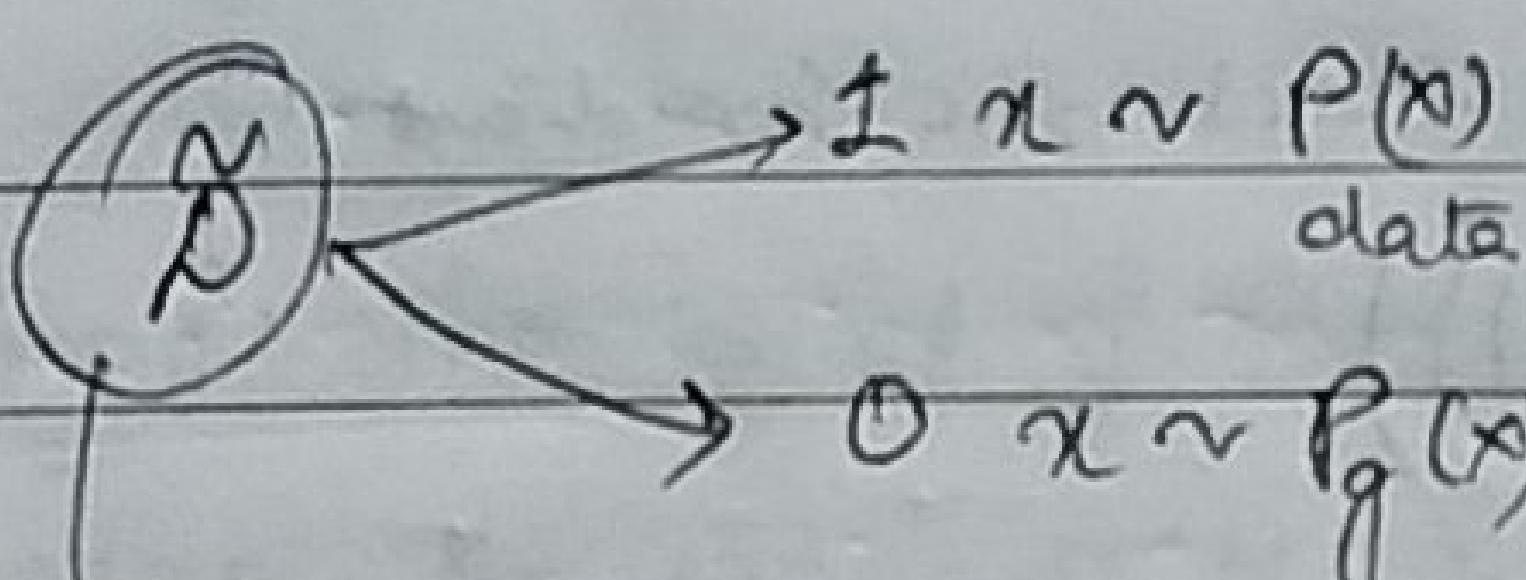
$$= E_{x \sim P_g(x)} \left[ \frac{1}{(D(x)-1)} \cdot \frac{\partial D(x)}{\partial x} \cdot \frac{\partial x}{\partial \theta} \right] \quad \{ \text{Rewriting} \}$$

$$\nabla_x = \frac{\partial}{\partial x}$$

$$\nabla_{\theta} = \frac{\partial}{\partial \theta}$$

$$= E_{x \sim P_g(x)} \left[ \frac{1}{(D(x)-1)} \cdot \nabla_x D(x) \cdot \nabla_{\theta} \nabla_x \right]$$

We are training the generator, after the training of the discriminator.



This is the perfect discriminator

if  $D$  tends to  $\delta$  &  $x \sim P_g$  (fake sample)

then,

$$\lim_{D \rightarrow \delta} D(x) = 0$$

$$\lim_{D \rightarrow \delta} \nabla_x D(x) = 0$$

Let's put this into our eqn

$$\nabla_{\theta} L = \lim_{D \rightarrow \delta} E_{x \sim P_g(x)} \left\{ \frac{\nabla_x D(x) \cdot \nabla_{\theta} \nabla_x}{(D(x)-1)} \right\} = 0 \quad (\text{Vanishing Gradient})$$

This is more of the mathematical Reasoning for our discussion of Vanishing Gradients.  
As we seen it turn out to be 0, which will not help in the back propagation  
on the w & b of the generator and the generator will produce

Page No.

Date

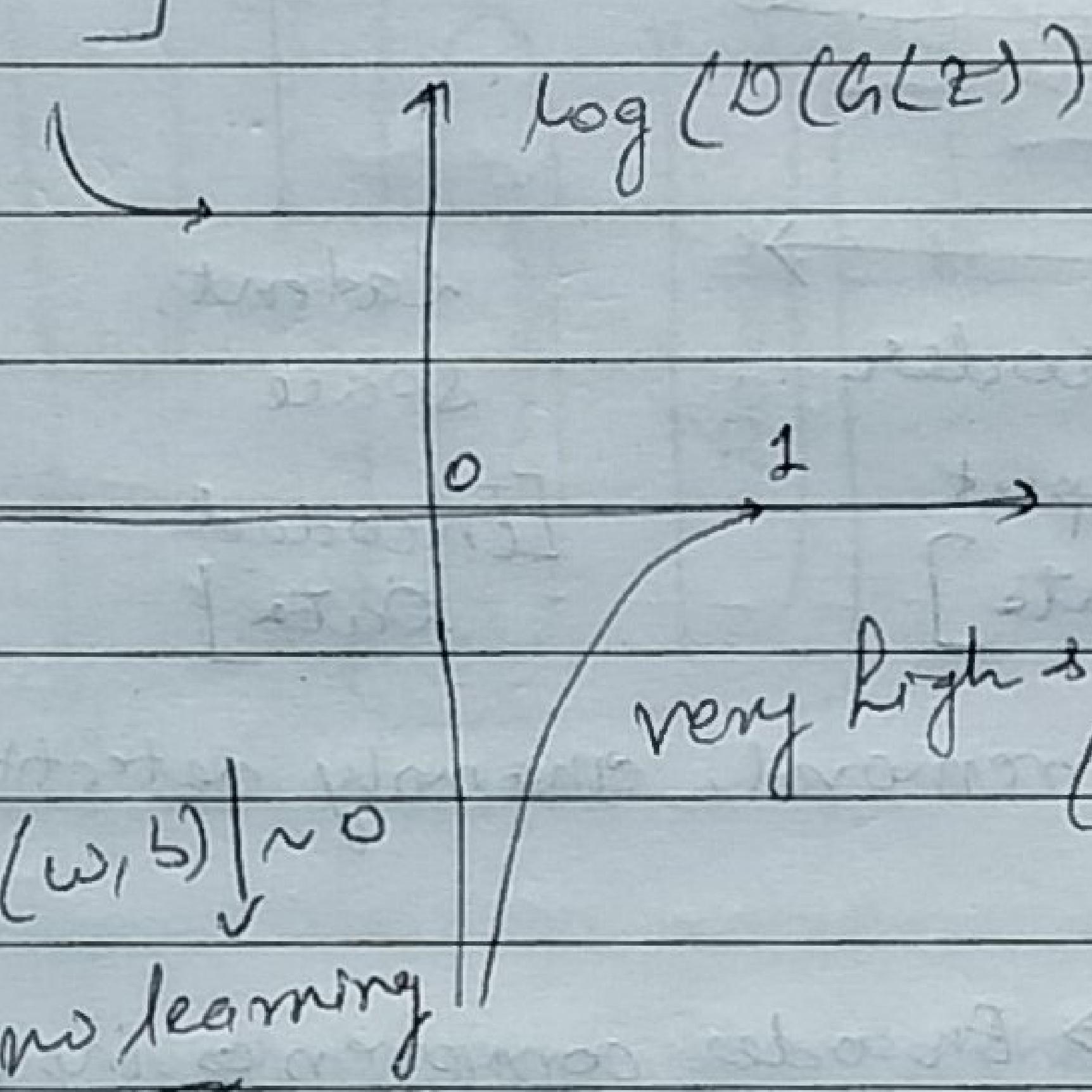
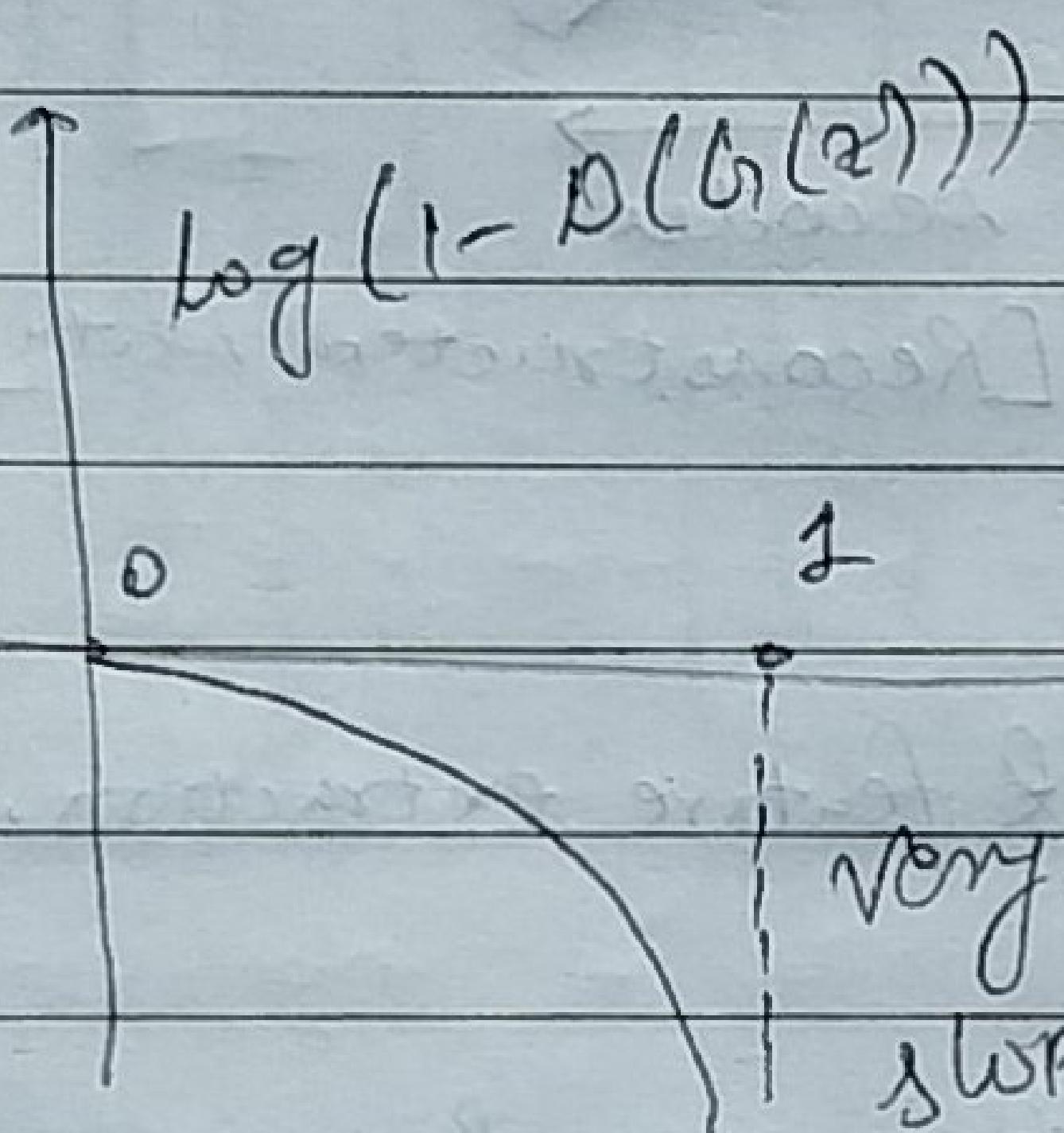
fake samples and it will not learn at all. (which are not close as real sample)

So to solve this just see the cost function :-

$$\Rightarrow \min_{G} E_{z \sim P_z(z)} [\log (1 - D_G G_G(z))]$$

To avoid this situation we will change the func

$$\max_{G} E_{z \sim P_z(z)} [\log (D_G G_G(z))]$$



very high slope  
(w, b) ↑, learning fast

- Mode Collapse:- During the training process, the generator may collapse to a setting where it always produces same output. This is called mode collapse.

[Code :[GitHub](#)]

