

→ We will skip the introduction part as you can find it anywhere on the internet.

AIM: Given some data, say images, a trained GAN will mimick the data and produce realistic looking images.

NOTE:

- ① Input to a trained GAN is just a random noise vector of any dimension (which should be the same dimensions the Generator was trained with).
- ② The Generator never really 'sees' the actual / real data for generating images. Thus, training the Generator is often slower and/or harder than training the Discriminator.
- ③ The Discriminator is just a binary image classifier, which is simultaneously trained from scratch along with the Generator.
- ④ A useful tip for implementing GANs is to **NOT** to assume that the Generator model is not enough capable / complex as compared to the Discriminator if you see the Discriminator loss going towards 0 and the Generator loss sky-rocketing during the initial training epochs.

MATHEMATICS:

1. If X is a Random Variable such that $X: \Omega \rightarrow \mathbb{R}$,
(Discrete)

Expectation of X , denoted by $E(X)$ is given by

$$E(X) = \sum_{x \in R(X)} x P(X=x)$$

$$E(X) = \sum_{x \in R(X)} x f_X(x)$$

Provided that $E(X)$ is absolutely convergent, i.e.

$$\sum_{x \in R(X)} |x| f_X(x) < \infty$$

where $f_X(x)$ is the Joint PMF of X .

(All thanks to the Probability & Statistics course!)

2. Log loss / Binary Cross Entropy loss :

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Click the image to get a detailed explanation

This is exactly the same loss function that we use for Logistic Regression.

If we carefully see the behavior of the BCE loss,

When $y_i = 0$:

$$\text{Loss}_0 = -\frac{1}{N} \sum_{i=1}^N \log(1 - p(y_i))$$

Since we want to make $p(y_i) = 0$ as well, thus we aim at making $\text{Loss}_0 \rightarrow 0$ as our model slowly trains.

If our model makes a terrible prediction of $p(y_i) = 1$, then $\log(1 - p(y_i)) \rightarrow (\log(0) = -\infty)$

Thus, our aim is to maximize Loss_0 .

When $y_i = 1$:

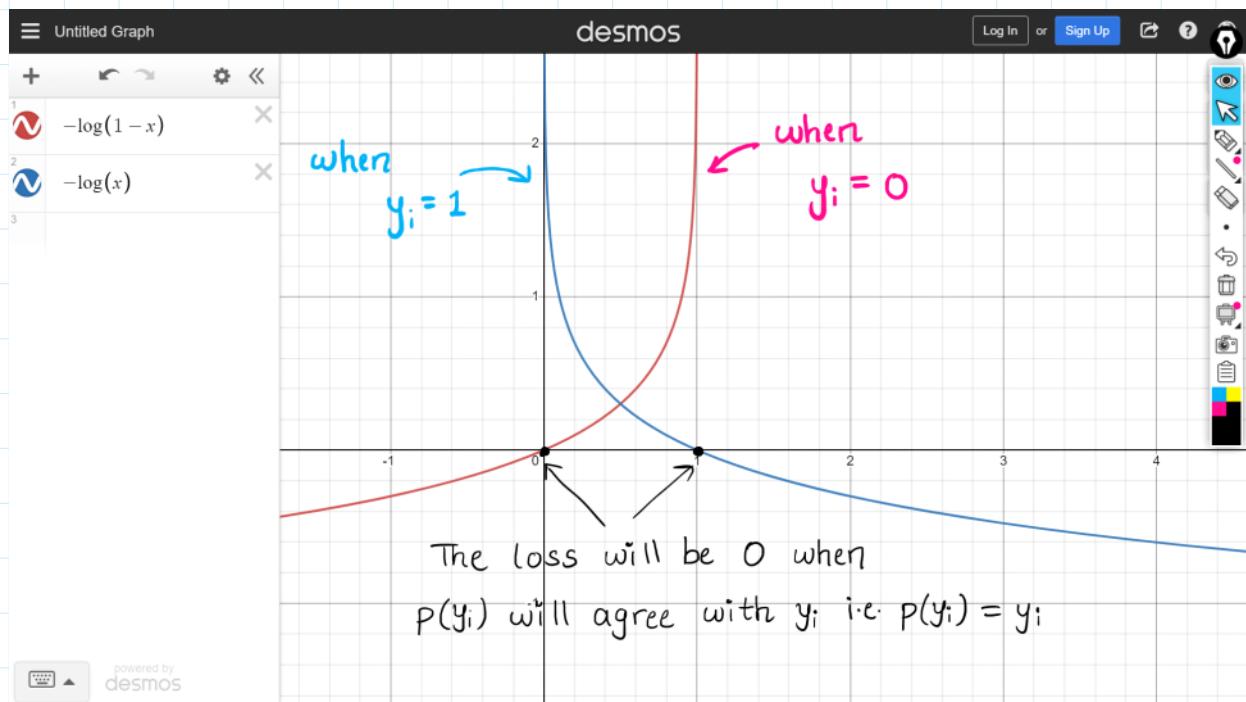
$$\text{Loss}_1 = -\frac{1}{N} \sum_{i=1}^N \log(p(y_i))$$

Since we want to make $p(y_i) = 1$ as well, thus we aim at making $\text{Loss}_1 \rightarrow 0$ as our model slowly trains.

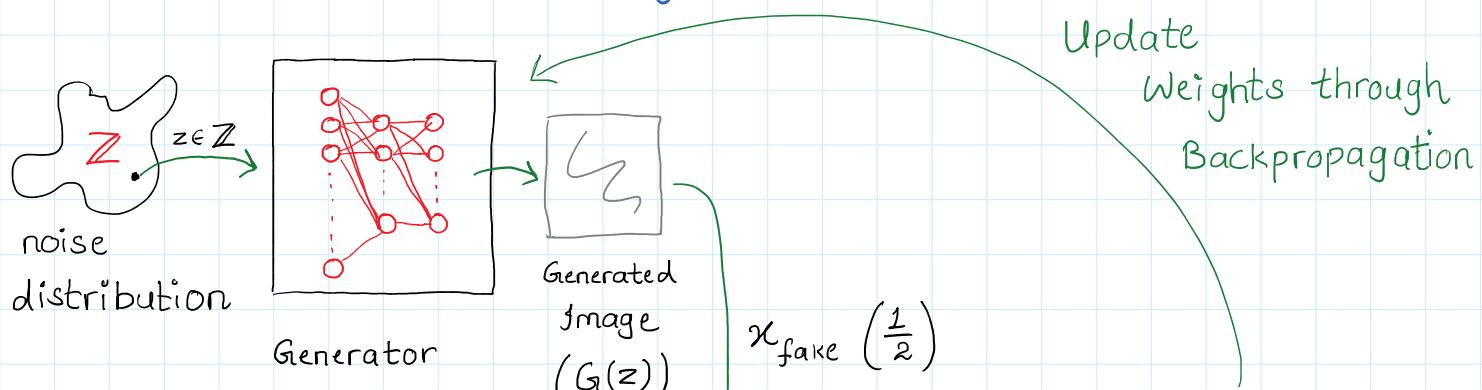
If our model makes a terrible prediction of $p(y_i) = 0$, then $\log(p(y_i)) \rightarrow (\log(0) = -\infty)$

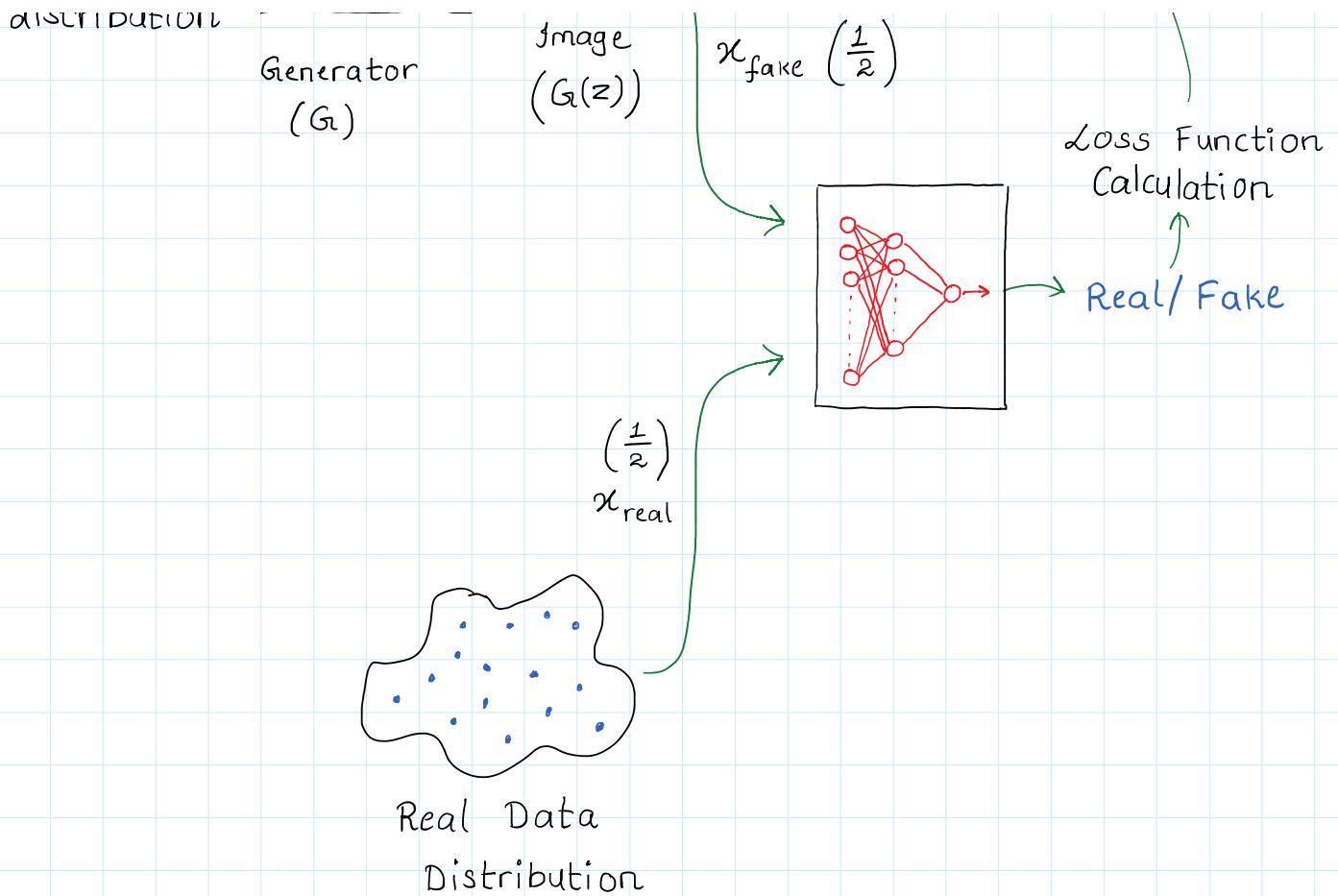
Thus, our aim is to maximize Loss_1 .

3. Getting a taste of the log loss / BCE loss.



Workflow of GAN Training:





The Original min-max Loss Function Derivation :

Given the BCE loss, removing the negative sign gives us :

$$\mathcal{L}(y, \hat{y}) = \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (\frac{1}{N} \text{ doesn't make a fundamental diff.})$$

we can replace \hat{y}_i with $D(x)$ (The discriminator's prediction of whether x is real or fake).

When we feed a real image, $y_i = 1$

$$\mathcal{L}(y_i, \hat{y}_i) = \log(D(x_i)) \quad \text{where } x \sim P_{\text{real}}(x)$$

$$\rightarrow \mathcal{L}(1, D(x_i)) = \log(D(x_i)) = \mathcal{L}_1 \quad - \textcircled{1}$$

When we feed a fake image, $y_i = 0$

$$\mathcal{L}(y_i, \hat{y}_i) = \log(1 - D(G(z_i))) \quad \text{where } z_i \sim \mathcal{Z}$$

$$\rightarrow \mathcal{L}(0, D(G(z_i))) = \log(1 - D(G(z_i))) = \mathcal{L}_0 \quad - (2)$$

The Discriminator's Story :

It wants to correctly classify the real and the fake images, thus :

$$\begin{aligned} D(x_i) &= 1 & \forall x_i \sim P_{\text{real}}(X), z_i \sim \mathcal{Z} \\ D(G(z_i)) &= 0 \end{aligned}$$

This means it's trying to make \mathcal{L}_0 and \mathcal{L}_1 converge to 0 (from ① & ②).

Initially, when the discriminator performs poorly :

$$\begin{aligned} D(G(z_i)) &= 1 \\ \text{or} & \quad \text{for some } z_i \sim \mathcal{Z} \text{ and/or some } x_i \sim P_{\text{real}}(X) \\ D(x_i) &= 0 \end{aligned}$$

This will cause \mathcal{L}_0 and/or \mathcal{L}_1 to tend to $-\infty$.

Thus, in a nutshell, the Discriminator is trying to maximize both \mathcal{L}_0 and \mathcal{L}_1 .

$$D: \max_D \sum_{i=1}^N [\log(D(x_i)) + \log(1 - D(G(z_i)))] \quad - (3)$$

The Generator's Story :

It wants to fool the Discriminator, thus

$$D(G(z_i)) = 1 \quad \forall z_i \sim \mathcal{Z}$$

This means it wants to make \mathcal{L}_0 tend to $-\infty$. (From ②)

Initially, when the Generator performs poorly :

$$D(G(z_i)) = 0 \quad \text{for some } z_i \sim Z$$

This will cause L_0 to tend to 0.

Thus, in a nutshell, the Generator is trying to minimize L_0 .

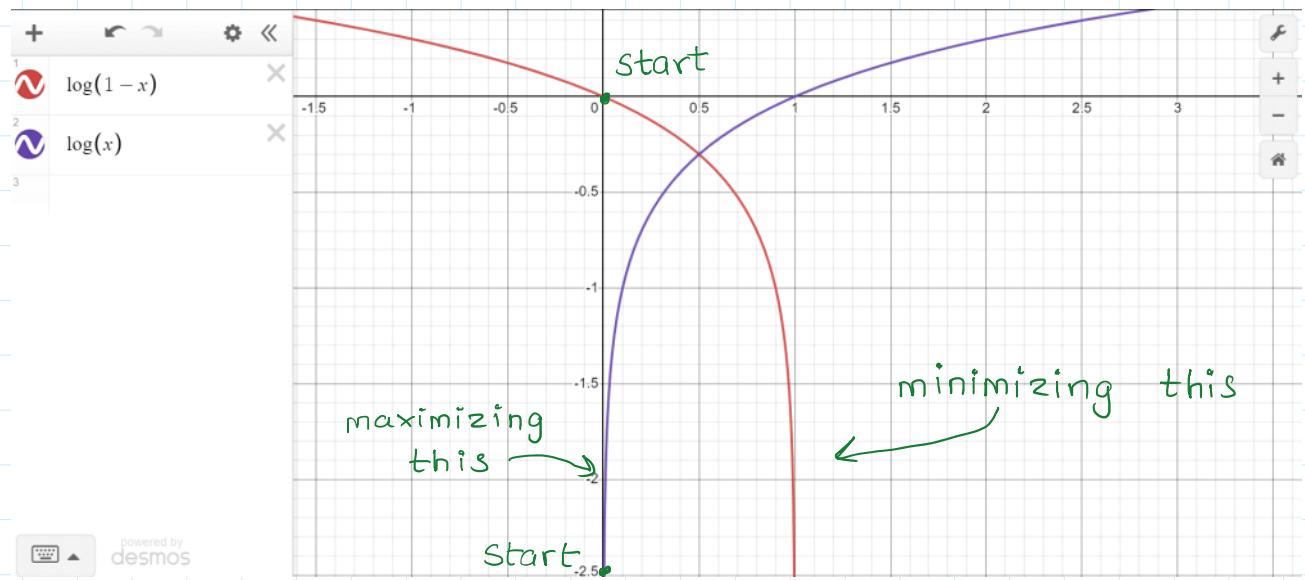
This is the same as saying that the Generator is trying to minimize both L_0 and L_1 as there is no role of the generator in the value of L_1 .

$$G: \min_G \sum_{i=1}^N [\log(1 - D(G(z_i)))] \quad (\text{In practice})$$

$$G: \min_G \sum_{i=1}^N [\log(D(x_i)) + \log(1 - D(G(z_i)))] \quad (\text{In theory})$$

- ④

A problem with this Generator Loss is saturating gradients during the early stages of Generator training.



So instead of minimizing $\log(1 - D(G(z^{(i)})))$, we will try maximizing $\log(D(G(z^{(i)})))$.

$$G: \sum_{i=1}^N \log(D(G(z^{(i)})))$$

Our aim : Maximize {G}

In practice, the Generator is only trained to

$$\max_G \sum_{i=1}^N \log(D(G(z^{(i)}))) \quad (\text{Non-saturating gradients!})$$

From ③ and ④,

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

I can also write this like :

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

(Not quite sure how and why) 