# Generative Adversial Network (GAN)

## Adversial Learing

Generator



$$z \to P_z(z) \to \Delta \quad [G] \to \Delta \quad x \sim P_g(x)$$

(Noise)

dataset $\quad x \sim P_{data}(x)$

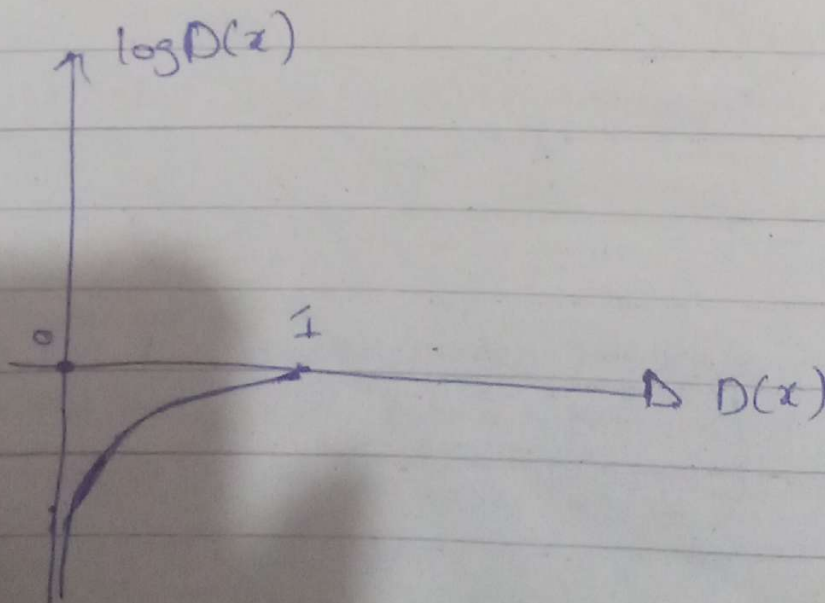Discriminator

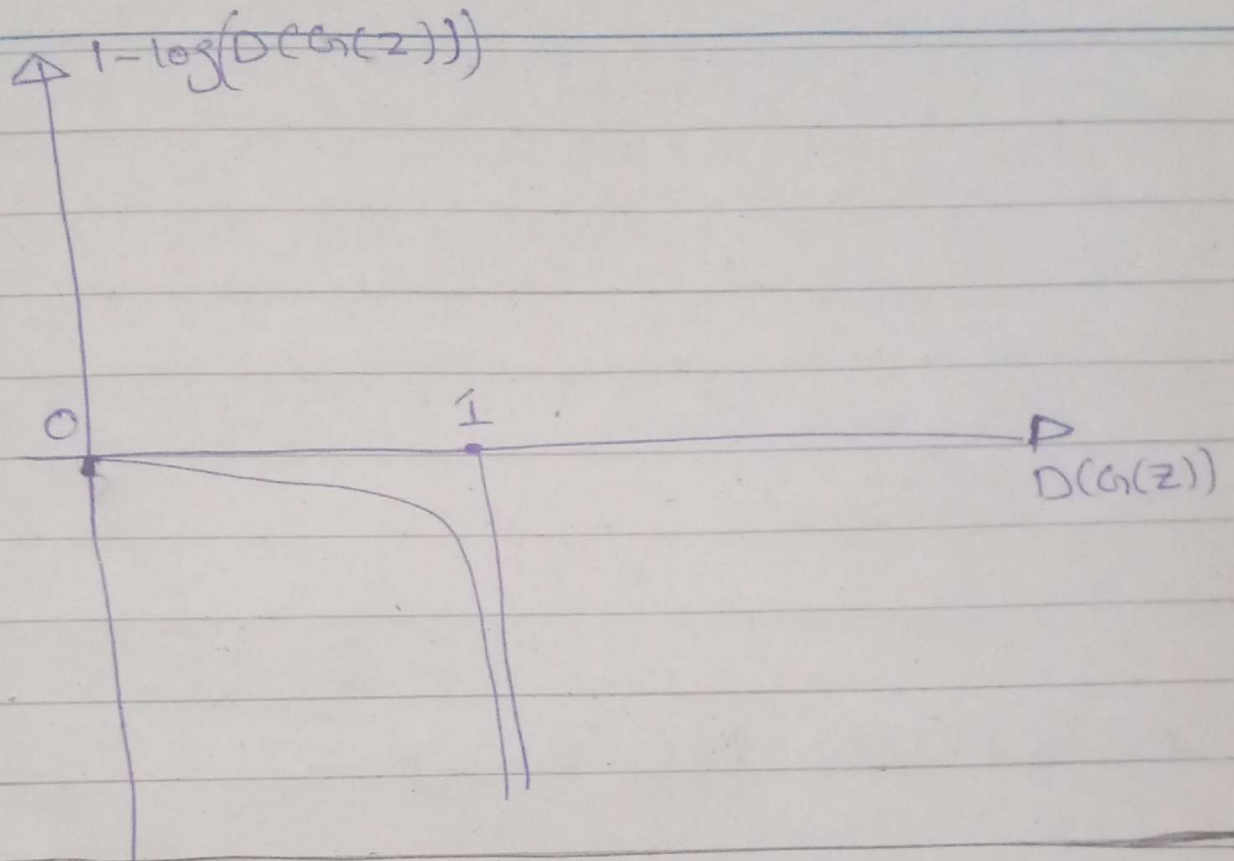## Adversial Learning

$$\min_G \max_D V(D, G)$$

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



$\log D(x)$

$0$

$1$

$\Delta \; D(x)$

De convolution → ~~Auto~~

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



Generator's Task is to fool the Discriminator

Training GAN's.

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

# Diffusion Models:-

## Recap of VAE

$$\epsilon \sim N(0,1)$$

$$Z = E(z) + V(z) \cdot \epsilon$$

We sequentially add noise until we completely distorted the image.
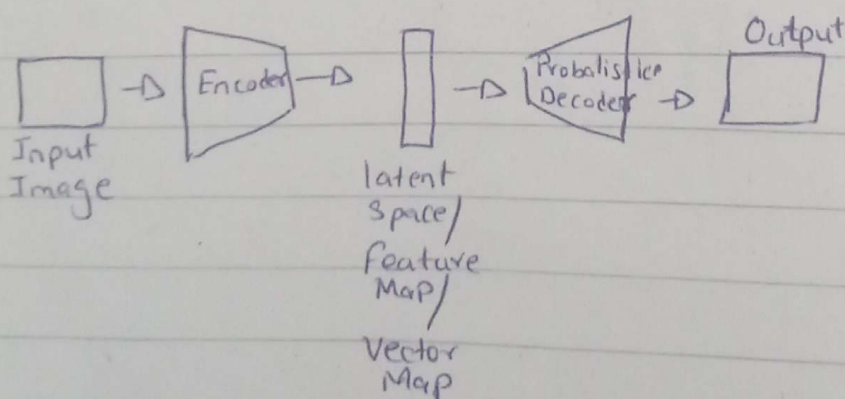
(Noise is always Gaussian)

(Gaussian distribution is the most closet to Natural)

$$z_1 = x_0 + n \quad (\text{Forward})$$
$$x_0 = x_1 - \hat{n} \quad (\text{Backward})$$

## VAE's is related to Diffussion Model



Variational auto encoder

So Encoder in VAE is similiar to Forward pass in D.M. & Decoder in VAE is similiar to Backward pass in D.M.

## Diffusion Models are trained on noise

$$x^n \longrightarrow \boxed{D.M} \longrightarrow \hat{n}$$

$n = noise$,

$$x_3 \longrightarrow \boxed{D.M} \longrightarrow \hat{n}$$

$$x_2 = x_3 - \hat{n}$$

## Forward Pass

$$q(x_t \mid x_{t-1}) = N(x_t; \underbrace{\sqrt{1-\beta_t} \, x_{t-1}}_{\text{Mean}}, \underbrace{\beta_t I}_{\substack{\text{Variance} \\ \text{(factor of} \\ \text{Identity} \\ \text{Matrix)}}})$$

$\underset{\text{output}}{\downarrow}$

e.g. $x_t = x_{t1} + x_{t2} + x_{t3}$

$\quad x_{t1} \quad x_{t2} \quad x_{t3}$

$$\begin{array}{c} x_{t1} \\ x_{t2} \\ x_{t3} \end{array} \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix} \right\} \text{Covariance}$$

forward pass is very simple:
Declare a distribution

Backward pass:

$$P_\theta(x_{t-1}|x_t) = \mathcal{N}\left(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I\right)$$

↓ Output    ↓ Mean    ↓ S.D/variance
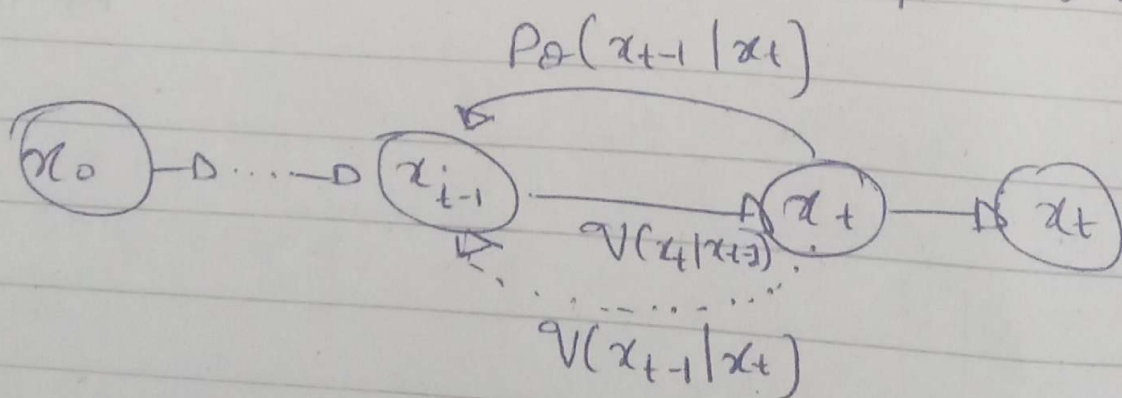
Identity matrix is used to simplify calculations, it reduces relations b/w random variables

Standard Deviation will remain same in FP & BP

· but we will try to approximate $\mu$ (Mean)

$$\Rightarrow |\mu_n - \mu_{\hat{n}}|$$

$$P_\theta(x_{t-1}|x_t)$$

$$x_0 - \square \cdots - \square \, x_{t-1} - \underset{V(x_t|x_{t-1})}{\text{———}} \, x_t - \square \, x_t$$

$$V(x_{t-1}|x_t)$$

A case study on how to solve
the problems

Multi model Learning

$$\epsilon \sim \mathcal{N}(0,1)$$

$$x_t = \sqrt{\overline{\alpha_t}}\, x_0 + (1-\overline{\alpha_t}) \cdot \epsilon$$

Mean of $q(x_{t-1} | x_t, x_0)$
(Backward Pass)

$$\mu_t = \underbrace{\left(\frac{1}{\sqrt{\alpha_t}}\right)}_{\text{Constant}} \left( x_t - \underbrace{\frac{1-\alpha_t}{\sqrt{1-\overline{\alpha_t}}}}_{\text{Constant}} \underbrace{\epsilon_t}_{\substack{\text{Actual} \\ \text{Noise}}} \right)$$

$$\boxed{\alpha = 1 - \beta} \longrightarrow \begin{bmatrix} \beta_1 = 0.002 \\ \beta_2 = 0.003 \\ \beta_3 = 0.007 \end{bmatrix}$$

$\triangle$ we end up with

$$\mu_t = x_t - \text{noise}(\epsilon_t)$$

$\downarrow$

This is the noise
we added into the image

Actual Mean for the reverse process
is

$$\mu_t = x_t - \text{noise}(\epsilon_t)$$

**#** If you have 2 distributions "P" & "q" And both are ~~as~~ gaussian & you want to find KL-Divergence there is a generalized formula for it there.

$$D_{KL}(q||P) = \frac{1}{2}\left( tr\left(\underline{\Sigma_P^{-1}\Sigma_q}\right) + (\mu_P - \mu_q)^T \Sigma_P^{-1}(\mu_P - \mu_q)\right.$$

$$\left. -K + \ln\left(\frac{\det \Sigma_P}{\det \Sigma_q}\right)\right.$$

Covariance
of P & q
A

for q =

$$\tilde{\mathcal{N}}(x_t, z_0) = \frac{1}{\sqrt{\alpha_t}}\left(x_* - \frac{\beta_t}{\sqrt{1-\alpha_t}}\in_t\right)$$

for p =

Work with the transformer layers.

# ATTENTION MECHANISM:—

1) Draw back ⟶

~~Background~~ Very Data hungry (20,000 to 30,000)

RNN (Recap) ⟶

~~PIEAS~~

The final layers lose the information in the earlier words

feature ~~of~~ vector

Attention Mechanism : It is all about weighted ~~elements~~ average

Context — window

Cross — attention & self — attention is used in Transformers.

$$C_4 = W_1 [\quad] + W_2 [\quad] + W_3 [\quad]$$

و خ      = لا      بر

As it has more importance, it will have more value weighted Avg

Decoder = I am a _____

Cross-Attention Mechanism ↑

Self-attention → Every word will be checked for context

For self-attention we take concepts from databases. (Query, Key, Value)

Every word will be converted to 3

$$\dot{V} | \dot{E} = S_1 \dot{V} + S_2 \underset{V}{\overset{V}{E}} + S_3 V$$

(Query) لا | (Key) لا

(Query) لا | (Key)

لا

(Query) لا | (Key)

و خ

Paper is called ← CLIP
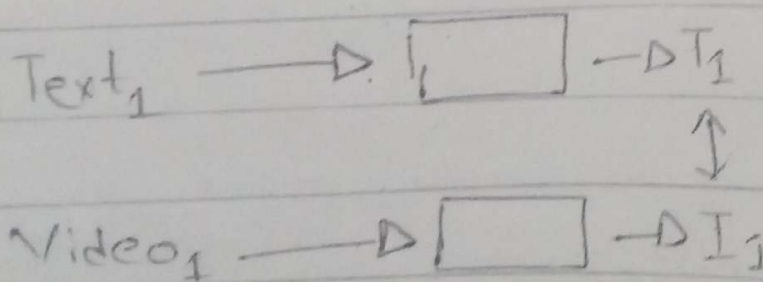
▷ Contrastive language Image
Pre la

Learning Transferable Visual Models from
Natural language Supervision

Berd used to extract features from
Text.
From Features we would used Vision
encoder's.

CLIP ( Research paper
reading is a must )

▷ Contrastive Learning Image Pre-training

$Text_1 \longrightarrow$ ▷ $\boxed{\phantom{I_1}}$ $\longrightarrow T_1$

$\updownarrow$

$Video_1 \longrightarrow$ ▷ $\boxed{\phantom{I_1}}$ $\longrightarrow I_1$

As these 2 belong to the same sample
we put the closer in the feature space.

ASR (ase Study $\rightarrow$ (NSF) $\rightarrow$ (Grant)

$\downarrow$

Automatic Speech Recognition 03/05/2025
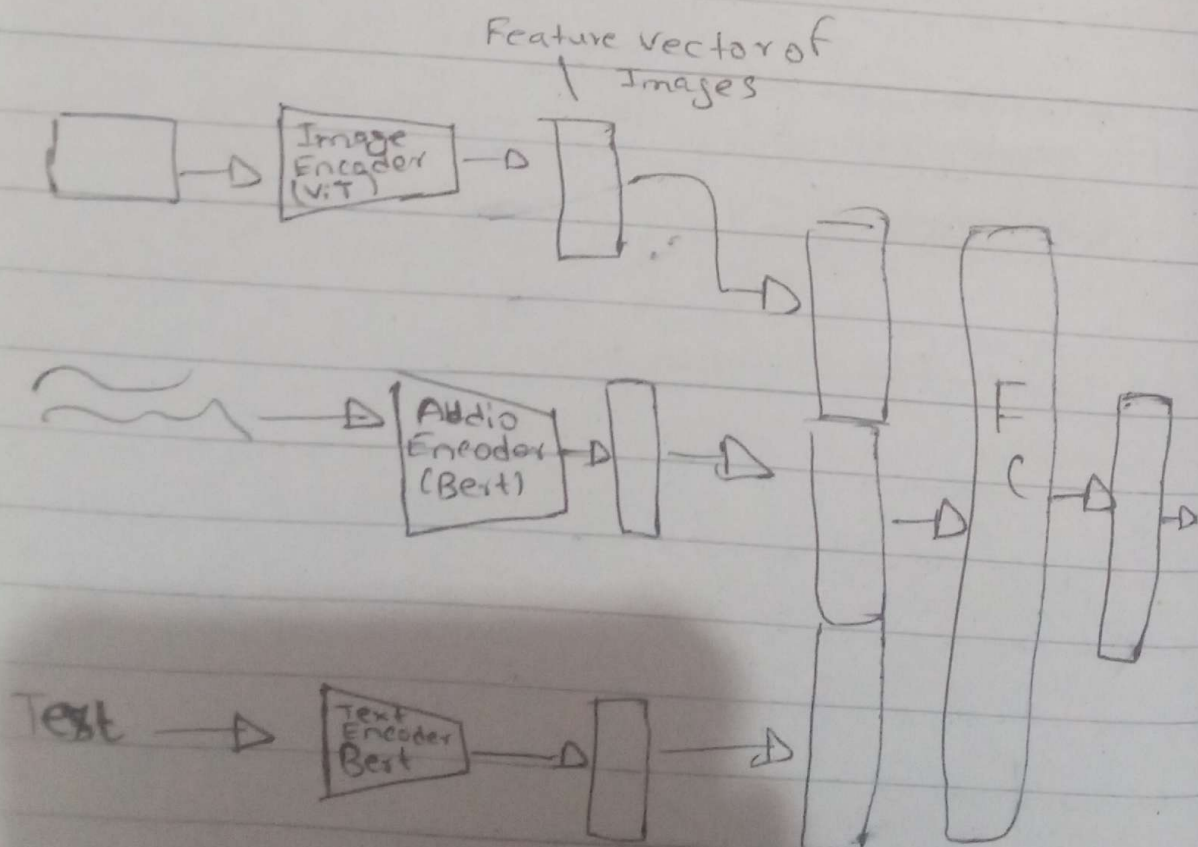
(ASD)

Kids with ~~erat~~ autism

Web-based Emotion recognition algo
that takes video as i/p

~~End goal of~~

Text encoder $\rightarrow$ bert
~~Img~~ Image Encoder $\rightarrow$ Vit
Audio encoder $\Rightarrow$ bert



Feature vector of Images

F
C

Audio is aligned with the video

Do frame by frame but it has less accuracy then video based.

3D - CNN for videos but has more paramet~~s~~ ~~but~~ web-based has limited parameters.

Department of Education
- It's very slow

What is the reason

↳ The Text is not bringing very big impact into the picture, remove it

~~the~~

By applying Fourier Transformation you can ~~app~~ convert $1D_\uparrow$ into $2D_\downarrow$
Signal     Signal

As Image encoder is 2D & audio encoder is 1D and eve end up with 2 encoders slowing us down, using Fourier transform

we end up with only 1 encoder speed is really good.

It is called spectrogram, when you end up converting 1D signal into 2D signal

CLIP

Contrastive learning, very import
used in computer vision

Masked auto encoding