

Objective: Develop a Mathematical Foundation for Gen AI.

Family of DGMs we'll look at:

1. Generative Adversarial Network (GAN)
2. Variational Auto-Encoder (VAE)
3. Denoising Diffusion Probabilistic Models (DDPM)
4. Auto-Regressive Models (AR)  
large language Models (LLM)
5. State-space Models (SSM)  
S4, Mamba
6. RL-Based Alignment for LLMs.  
RLHF, PPO, DPO.

Generative Models:

Examples: ChatGPT, Claude, Gemini, etc...  
Conditional Text Generators.

DALL-E, Stable Diffusion.

Conditional Image Generators.

Speech Generators : text  $\rightarrow$  wav

Mathematical Formulation of the Problem:

Starting point: Data.

Data  $D = \{x_1, x_2, \dots, x_n\}$  iid  $P_x$  (Unknown dist.)

$x_i \in \mathbb{R}^d$ , d-dimensionality of the data.

$X$ : Random Variable with a distribution

Suppose  $D = \{x_1, x_2, \dots, x_n\}_{n=1000}$ , then  $P_x$ .

$x_i, x_j$  are statistically independent and are sampled from the same distribution.  $\rightarrow$  Assumptions.

$x_i \perp\!\!\!\perp x_j, x_i \sim P_x$

Why?? For Mathematical Ease.

$x_i$ : Instances of a Vector-Valued Random Variable of size 'd'.

Goal: Estimate  $P_x$  & learn to sample from it.

Not all models explicitly figure out  $P_x$ . They may generate  $P_x$  in implicit manner.

But almost all of them learn to sample from the underlying distribution.

## General Principle of Gen Models:

↑ Density function.

- i) Assume a parametric family on  $P_x$ , denoted by  $P_\theta$ .

$P_\theta$ : Represented using Deep NNs.  
(Model)

- ii) Define & compute a divergence metric between the  $P_\theta$  &  $P_x$  → We dunno  $P_x$  tho??

- iii) Solve an optimization problem over the parameters of  $P_\theta$ , to minimize the above divergence metric.

The task of Generative modelling is 2-fold.

- a) Estimate the dist.
- b) Sample from it.

Using the above 3-step process, we accomplished

a). what about b)?

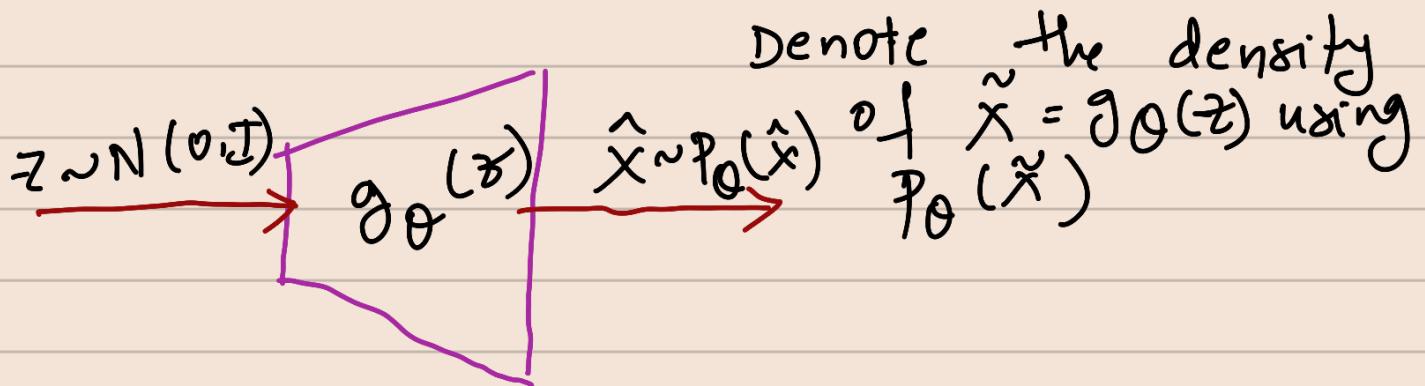
Example:

Some arbitrary but known distribution  $z \sim N(0, I)$

Suppose  $g_\theta(z) : z \rightarrow x$ , then

$\tilde{x} = g_\theta(z)$  has a different distribution than that of  $z$  & the distribution of  $\tilde{x}$  depends on the function  $g_\theta()$ .

Suppose  $g_\theta(z)$  is a neural network



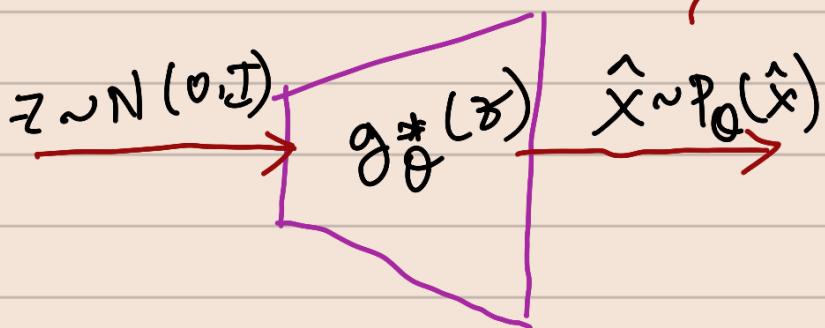
Suppose  $D(P_x || P_\theta)$  denote a divergence measure between  $P_x$  &  $P_\theta$ ,

Next step: solve the following optimization problem

$$\theta^* = \underset{\theta}{\operatorname{argmin}} D(P_x || P_\theta)$$

Upon solving the above optimization

problem, the distribution  $p_x$  is implicitly estimated by  $g_\theta(z)$  & one can sample from  $p_x$  using  $g_\theta(z)$ .



Push-forward  
Methods

A sample from  $z \sim N(0, I)$ , passed through  $g_\theta(z)$  would produce a sample from  $p_Q^*(\hat{x})$ , which is closer to  $p_x$ , we end up sampling from  $p_x$ .

General principle followed by all Generative Models

The Questions:

- ① How to compute the divergence metric without knowing  $p_\theta$  &  $p_x$ .
- ② What should be the choice of the divergence metric  $\mathcal{D}$ ?

③ How to choose the  $g(\cdot)$ , in turn  $P_\theta$ ?

④ How to solve the optimization problem of minimizing the divergence metric?

## Variational Divergence Minimization

Define divergence metrics between distributions

f-divergence:

Given 2 probability distribution functions with the corresponding density functions denoted by  $P_x \in P_\theta$ , the f-divergence between them is defined as follows:

$$D_f(P_x \parallel P_\theta) = \int_X P_\theta(x) f\left(\frac{P_x(x)}{P_\theta(x)}\right) dx.$$

$f(u): \mathbb{R}^+ \rightarrow \mathbb{R}$ , convex, left-semicontinuous  
 $f(1) = 0$

$X$ : space on which the  $P_x \in P_\theta$  are supported

## Properties of $f$ -Divergence:

- ①  $D_f(\cdot) \geq 0$  for any choice of  $f(\cdot)$
- (ii)  $D_f(P_x || P_\theta) = 0$ , if  $P_x = P_\theta$

## Examples of $f$ -divergence:

- a)  $f(u) = u \log u$  : KL-Divergence.

$$D_f(P_x || P_\theta) = \int_X P_\theta(x) \frac{P_x(x) \cdot \log\left(\frac{P_x(x)}{P_\theta(x)}\right)}{P_\theta(x)} dx$$

$$= D_{KL}$$

$$\underbrace{D_{KL}(P_x || P_\theta)}_{\text{Forward KL}} \neq \underbrace{D_{KL}(P_\theta || P_x)}_{\text{Reverse KL}}$$

KL Divergence  
is not symmetric.

Different choices of ' $f$ ' functions result in different divergence metrics with their own properties, which necessitates the need to look at swathes of divergence metrics.

b)  $f(u) = \frac{1}{2} \left( u \log u - (u+1) \log\left(\frac{u+1}{2}\right) \right)$  : JS-Divergence.

↳ Used famously in GANs

c)  $f(u) = \frac{1}{2}|u-1|$  : Total Variational Distance.

# Algorithm for F-Divergence Minimization.

Objective: Algorithm to minimize  $D_f$  b/w

$P_x \leq P_\theta$ , without knowing both, but having samples from both.

Samples from  $P_x$ : dataset D

Samples from  $P_\theta$ : outputs of  $\mathcal{G}_\theta(z)$  for different  $z$ .

Without knowing  $P_\theta(x) \leq P_x(x)$ , solving a high-dimensional integral is infeasible.

Key Idea: Integrals involving density functions can be approximated using samples from the distribution.

Suppose we want to compute

$$I = \int_X h(x) \cdot P_x(x) \cdot dx \text{ where } h(x) \text{ is a function \& } P_x(x) \text{ the density fn.}$$

We have samples drawn iid from  $P_x$   
 $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} P_x$

$$I = \mathbb{E}_{P_x}[h(x)] \quad \begin{array}{l} \text{[From the law of]} \\ \text{Unconscious Statistician]}\end{array}$$

## Law of Large Numbers:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n h(x_i) \approx \mathbb{E}_{P_X} [h(x)]$$

$x_i \sim \text{iid } P_X$

If the  $f$ -divergence metric can be expressed in terms of expectation of functions wrt  $P_X \not\ll P_0$ , then one can compute & optimize them.

Expressing  $D_f$  in terms of Expectations over  $P_X \not\ll P_0$ .

$$D_f(P_X \| P_0) = \int_x P_0(x) f\left(\frac{P_X(x)}{P_0(x)}\right) dx$$

Even though this expression looks really similar to the  $\int P_0(x)h(x).dx$ , you cannot directly rewrite that expression in terms of expectation precisely because of the arguments of  $f$ . They aren't simple "x" but instead are the ratio of probability functions.

So, we somehow have to decouple the ratio of PDFs from  $f$ .

SLIGHT DETOUR

# Conjugate function for a Convex fn:

If  $f(u)$  is a convex function, then there exists a conjugate function

$$f^*(t) := \sup_{u \in \text{dom}(f)} \{ut - f(u)\}$$

Point-wise defn

↳ Basically, we are lower bounding  $f(u)$  at every point  $t$ .

## Properties of Conjugate:

i)  $f^*(t)$  is also convex.

ii)  $[f^*(t)]^* = f(u) \Rightarrow$

$$f(u) = \sup_{t \in \text{dom}(f^*)} \{-tu - f^*(t)\}.$$

DETOUR DONE.

$$D_f(P_x \parallel P_\theta) = \int_X P_\theta(x) f\left(\underbrace{\frac{P_x(x)}{P_\theta(x)}}_u\right) dx$$

$$= \int_X P_\theta(x) f(u) dx.$$

$$f(u) = \sup_t \{tu - f^*(t)\}.$$

$$D_f(P_X || P_\theta) = \int_x p_\theta(x) \cdot \sup_t (t u - f^*(t)) \cdot dx.$$

$$= \int_x p_\theta(x) \sup_t \left\{ t \left[ \frac{p_X(x)}{p_\theta(x)} \right] - f^*(t) \right\} \rho \cdot dx.$$

$\hookrightarrow$  cannot directly pull out the  $\sup$ !!

$$= \sup_{T(x) \in \mathcal{T}} \int_x p_\theta(x) \left\{ T(x) \frac{p_X(x)}{p_\theta(x)} - f^*(T(x)) \right\} \rho \cdot dx$$

∴ The inner optimization problem involves  $x$  & the solution for it is dependent (a fn of  $x$ )

$$\begin{aligned} \mathcal{T}: X &\longrightarrow \text{dom } f^* \\ T(x) &\in \mathcal{T} \end{aligned}$$

Space of functions containing solutions for the inner optimization problem

But,

$$D_f \geq \sup_{T(x) \in \mathcal{T}} \int_x p_\theta(x) \left\{ T(x) \frac{p_X(x)}{p_\theta(x)} - f^*(T(x)) \right\} \rho \cdot dx$$

because the space of functions  $\mathcal{T}$ , that we are optimizing over may not contain the optimal  $T^*(x)$ , that is the soln for the inner optimization problem.

$$\geq \sup_{T(x)} \left[ \int_x P_x(x) T(x) \cdot d\pi - \int_x P_\theta(x) f^*(T(x)) \cdot d\pi \right]$$

$$\geq \sup_{T(x)} \left[ E_{P_x} T(x) - E_{P_\theta} f^*(T(x)) \right]$$

18/05/25, 02hs

## Realization of VDM

Given Data  $D = \{x_1, x_2, \dots, x_n\} \stackrel{\text{iid}}{\sim} P_x$

Goal:  
 $\theta^* = \arg \min_{\theta} D_f(P_x \parallel P_\theta)$

$$D_f \geq \sup_{T(x) \in \mathcal{T}} \left[ E_{P_x} T(x) - E_{P_\theta} f^*(T(x)) \right]$$

We set out to solve —  $\theta^* = \arg \min_{\theta} D_f(P_x \parallel P_\theta)$

but we are relegated to solve  $\theta^* \approx \arg \min_{\theta} [ \text{lower bound on } D_f ]$   
 and that's the best we can do.

$$= \arg \min_{\theta} \left[ \sup_{T(x)} \left( E_{P_x} T(x) - E_{P_\theta} f^*(T(x)) \right) \right]$$

wrt the parameters of  $g_\theta(z)$  ↪ wrt a class of functions  $T(x) \in \mathcal{T}$

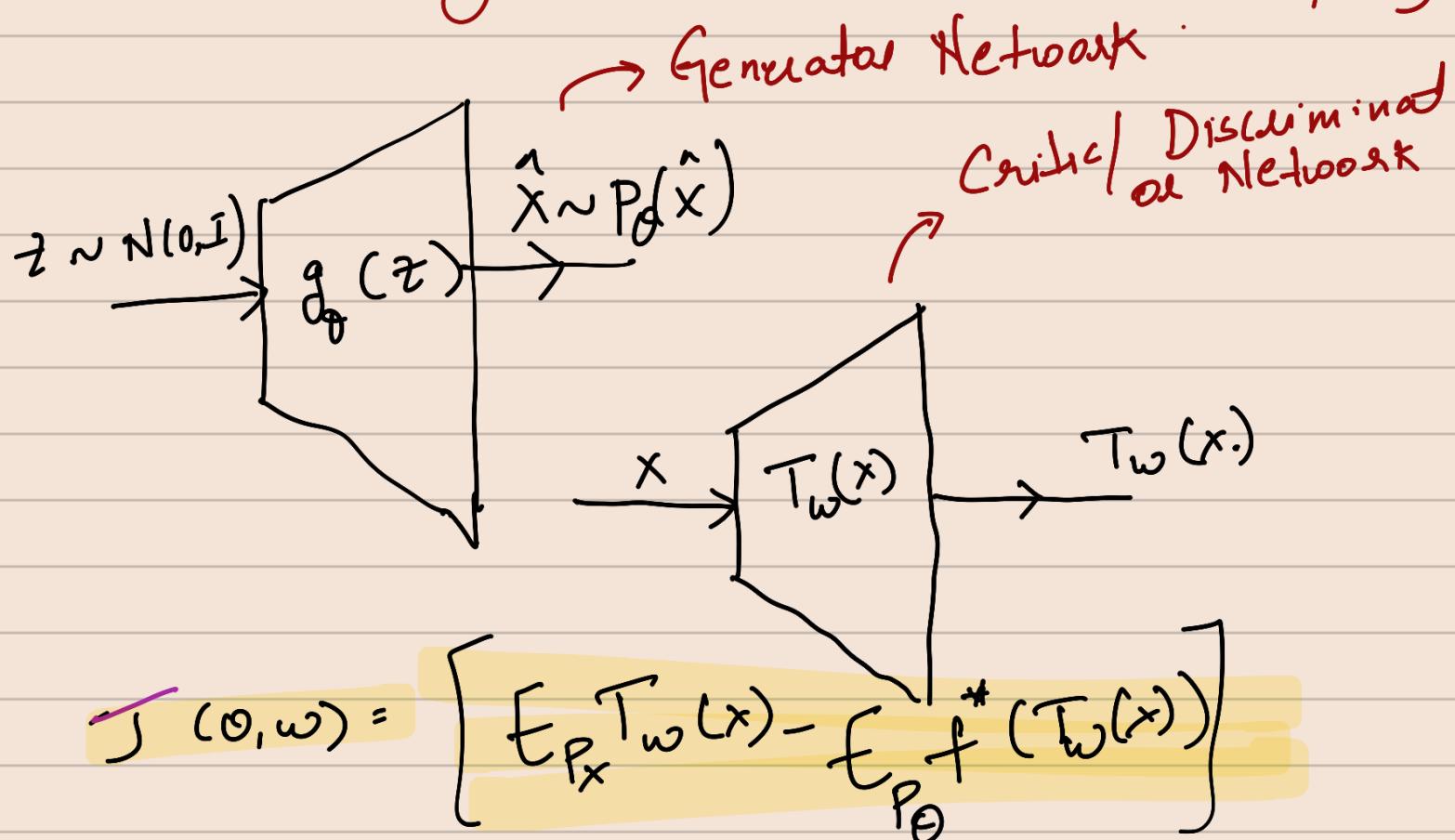
This, cannot be done analytically by optimizing over space of non-trivial functions.

In practice, we represent  $\mathcal{T}$  via neural networks:  $T_w(x)$ , where  $w$  are the parameters of the neural network.

With this, the objective will become:

$$\theta^*, w^* = \underset{\theta}{\text{arg min}} \max_w \left[ \mathbb{E}_{P_X} T_w(x) - \mathbb{E}_{P_\theta} f^*(T_w(x)) \right]$$

Implementing VDM for Generative Sampling.



$$\theta^*, w^* = \arg \min_{\theta} \max_w J(\theta, w)$$

↳ Saddle-point optimization

Problem where we have alternate minimization & maximization.

Blueprint for adversarial training

Any saddle-point optimization problem is also a adversial optimization problem.

19 May 2025, W2L6:

We know, by construction,  $T(\cdot) : X \rightarrow \text{dom } f^*$

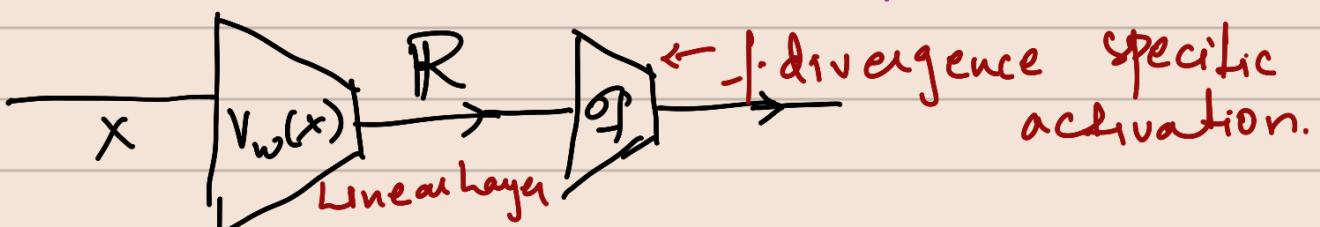
In practice, we do this by:

$T_w(x)$  is represented as  $\sigma_f(V_w(x))$ , so that the range of  $T$  network corresponds to domain of  $f^*$ .

Depending on the choice of  $f^*$ , we need to tweak the 'T' network

where  $\sigma_f$  is a  $\gamma$ -divergence specific activation  
 $V_w(x) : X \rightarrow \mathbb{R}$ ,  $\sigma_f(v) : \mathbb{R} \rightarrow \text{dom } f^*$

This means the T network that we'll approximate using a neural network is represented as a composition of  $\gamma$ 's.



$$J(\theta, \omega) = \mathbb{E}_{P_X} [g_f(V_\omega(x))] - \mathbb{E}_{P_\theta} [f^*(g_f(V_\omega(x)))]$$

# Generative Adversarial Networks

A special case of VDM algorithms:

The choice of  $f$ -divergence:  $\text{wlog } u = (u+1)\log(u/(u+1))$

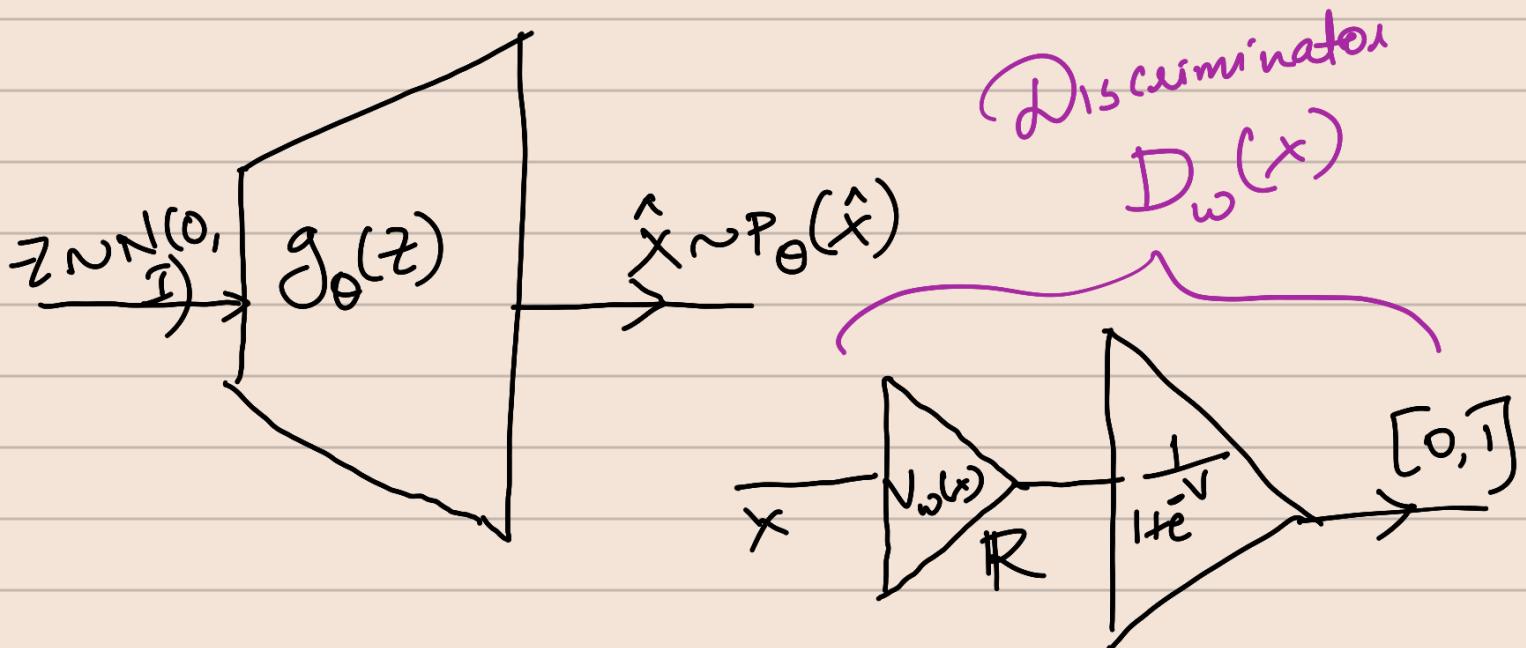
$$f^*(t) = -\log(1-e^t), \quad \text{dom } f^* = \bar{\mathbb{R}}$$

$$g_f(v) = -\log(1+e^{-v})$$

An fully similar  
to JS Div.

$$\overline{J}_{\text{Gan}}(\theta, \omega) = \mathbb{E}_{P_X} [\log D_\omega(x)] + \mathbb{E}_{P_\theta} [\log (1-D_\omega(x))]$$

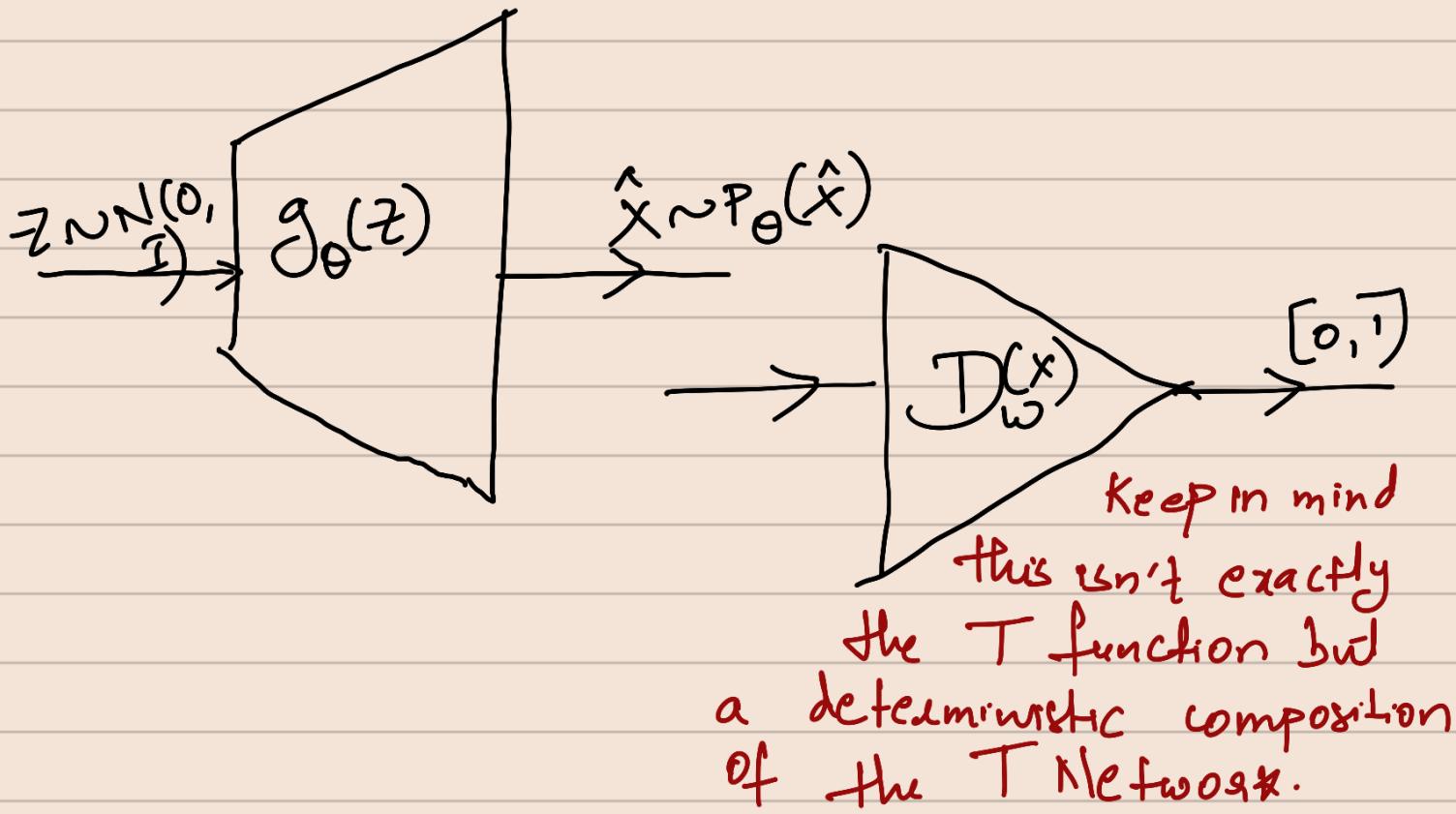
where  $D_\omega(x) = \frac{1}{1+\bar{\rho}^{-V_\omega(x)}}$



$$\begin{aligned}
 J_{GAN}(\theta, \omega) &= E_{P_X} T_\omega(x) - E_{P_\theta} f^*(T_\omega(x)) - \\
 &: E_{P_X} \sigma_f(V_\omega(x)) - E_{P_\theta} f^*(\sigma_f(V_\omega(x))) \\
 &: E_{P_X} -\log(1+e^{-V_\omega(x)}) - E_{P_\theta} f^*(-\log(1+e^{-V_\omega(x)})) \\
 &= E_{P_X} \log\left(\frac{1}{1+e^{-V_\omega(x)}}\right) - E_{P_\theta} f^*\left(\log\left(\frac{1}{1+e^{-V_\omega(x)}}\right)\right) \\
 &: E_{P_X} \log\left(\frac{1}{1+e^{-V_\omega(x)}}\right) + E_{P_\theta} \log\left(\frac{1}{1+e^{-\log\left(\frac{1}{1+e^{-V_\omega(x)}}\right)}}\right) \\
 &: E_{P_X} \log\left(\frac{1}{1+e^{-V_\omega(x)}}\right) + E_{P_\theta} \log\left(1 - \frac{1}{1+e^{-V_\omega(x)}}\right) \\
 &: E_{P_X} \log D_\omega(x) + E_{P_\theta} \log(1 - D_\omega(x))
 \end{aligned}$$

Derivation of  $J_{GAN}$  for  $\int := u \log u - (u+1) \log(u+1)$

# GAN ARCHITECTURE



$$J_{GAN}(\theta, w) = \underset{x \sim P_x}{\mathbb{E}} \log D_w(x) + \underset{\hat{x} \sim P_\theta}{\mathbb{E}} \log(1 - D_w(\hat{x}))$$

21 May 2025

## IMPLEMENTATION OF GAN IN PRACTICE

Input:  $D = \{x_1, x_2, \dots, x_n\} \stackrel{iid}{\sim} P_x$

$$\omega^* = \arg \max_w \left[ \underset{P_x}{\mathbb{E}} (\log D_w(x)) + \underset{P_\theta}{\mathbb{E}} (\log (1 - D_w(\hat{x}))) \right]$$

$$\approx \underset{\omega}{\operatorname{argmax}} \left[ \frac{1}{B_1} \sum_{j=1}^{B_1} \log D_{\omega}(x_i) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_{\omega}(\hat{x}_j)) \right]$$

$$\begin{aligned} x_1, \dots, x_{B_1} &\sim P_X \\ \hat{x}_1, \dots, \hat{x}_{B_2} &\sim P_{\theta} \end{aligned}$$

$$\omega^{t+1} \leftarrow \omega^t + \alpha_1 \nabla J_{GAN}(\theta, \omega); \text{ One gradient step. thru Discriminator.}$$

$\theta$  is kept a constant.

$$\theta^* = \underset{\theta}{\operatorname{argmin}} J_{GAN}(\theta, \omega)$$

$$\approx \underset{\theta}{\operatorname{argmin}} \left[ \frac{1}{B_1} \sum_{j=1}^{B_1} \log D_{\omega}(x_i) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_{\omega}(\hat{x}_j)) \right]$$

Independent of  $\theta$ .

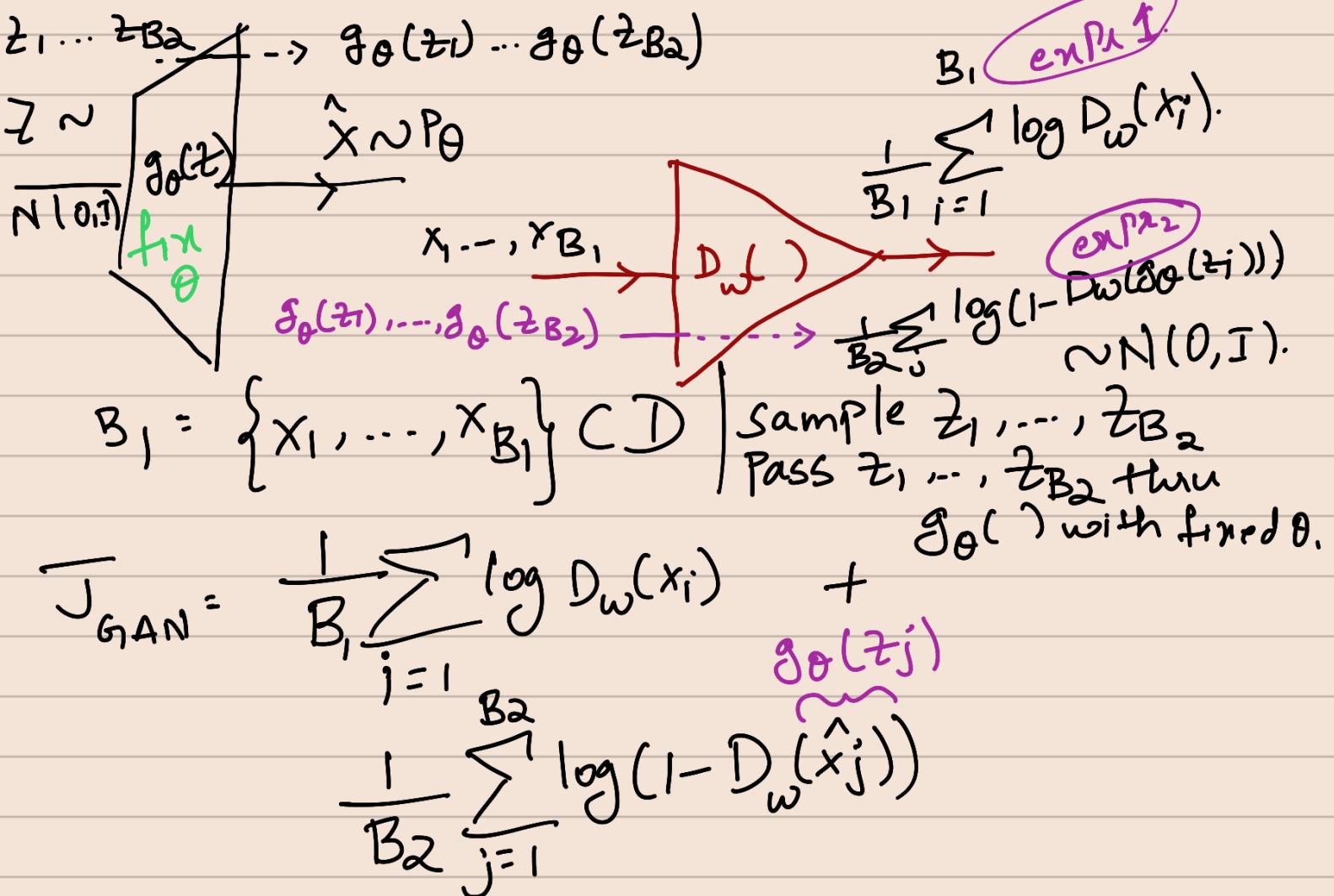
$$\approx \underset{\theta}{\operatorname{argmin}} \left[ \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_{\omega}(g_{\theta}(z_j))) \right]$$

$$\theta^{t+1} \leftarrow \theta^t - \alpha_2 \nabla J_{GAN}(\theta, \omega); \text{ 1 GD step thru the generator.}$$

$\omega$  is kept a constant

## TO TRAIN THE DISCRIMINATOR:

keep  $\theta$  constant



Once we have  $\text{exp1}$  &  $\text{exp2}$ , let us calculate the loss and backpropagate the gradients all the way to the inputs of the discriminator.

$$J \stackrel{P_\theta(x) = P_X(x)}{\rightarrow} E_{P_X}[\log(D(x))] + E_{P_\theta}[\log(1 - D(x))] = J_{GAN}.$$

The optimization over  $D(x)$  is treated independently - for each 'x', allowing us to find optimal  $D^*(x)$ . by maximizing point-wise loss.

$$-P(D(x)) = P_x(x) \log(P_D(x)) + P_\theta(x) \log(1-P_D(x))$$

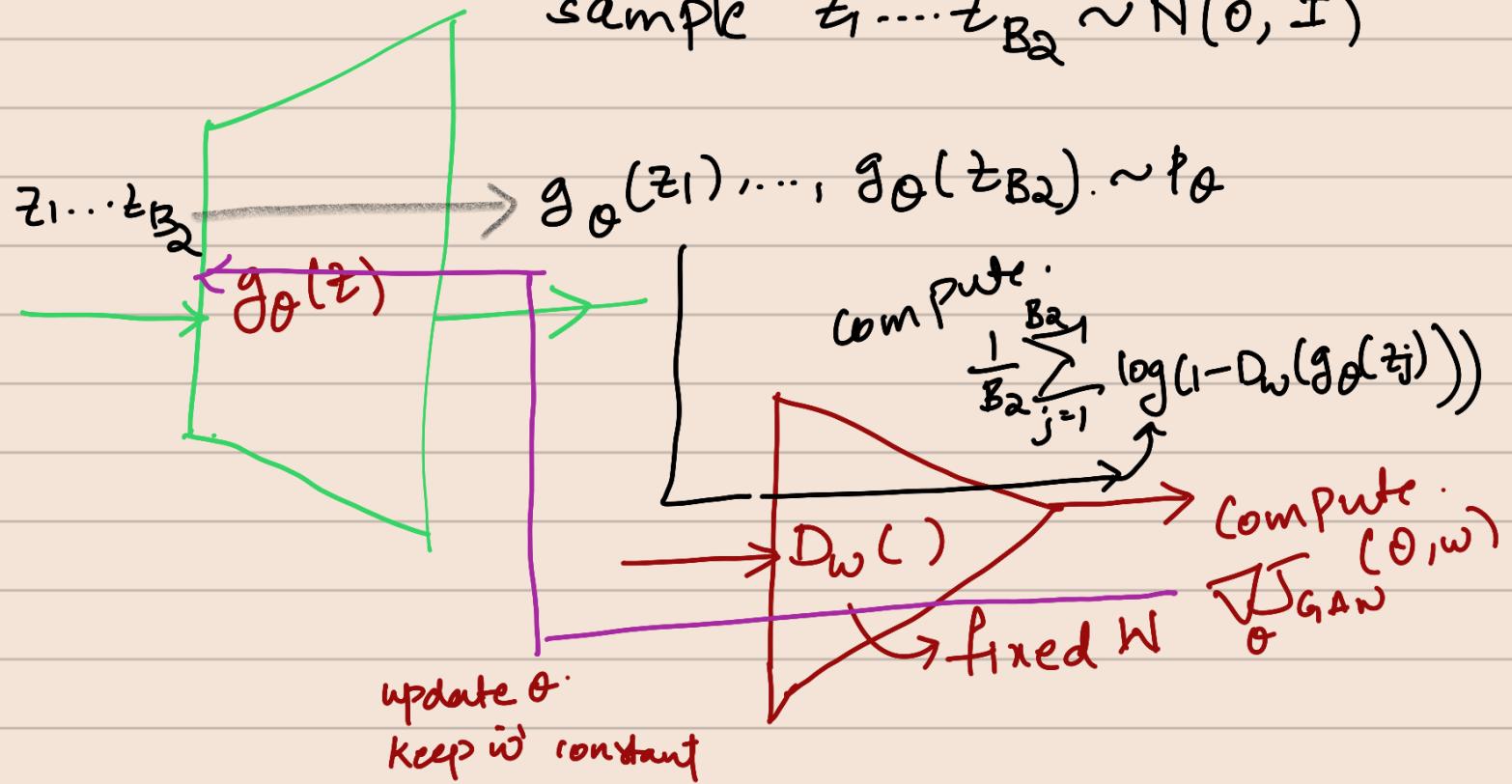
$$\frac{d \underbrace{f(D(x))}_{d D(x)}}{d D(x)} = \frac{P_x(x)}{D(x)} + \frac{P_\theta(x) (-1)}{1-D(x)} = 0$$

$$\frac{P_x(x)}{D(x)} = \frac{P_\theta(x)}{1-D(x)} \quad \text{if} \quad 1-D(x)=D(x)$$
$$2D(x)=1$$
$$D(x)=\sqrt{2}$$

Think more deeply about why optimization over  $D(x)$  is treated independently for each  $x$ .

## TO TRAIN THE GENERATOR:

sample  $z_1, \dots, z_{B_2} \sim N(0, I)$



$$\left[ \frac{1}{B_2} \sum_{j=1}^{B_2} \log(1 - D_w(g_\theta(z_j))) \right] = \bar{J}_{GAN}$$

To update the params of generator, we do not need the 1<sup>st</sup> term (we do not need the samples from the data).

$$\theta^{t+1} \leftarrow \theta^t - \alpha_2 \nabla_{\theta} \bar{J}_{GAN}(\theta, w)$$

Typically,

Alternate between 1 step of Generator and 1 step of discriminator whilst training

## Stopping Criterion 2)

When the quality metrics of the generators are reached.

## Training VDM or GAN's done

Next we'll look at inference of GAN, improvisations [different f-functions' result in different kind of VDM's]

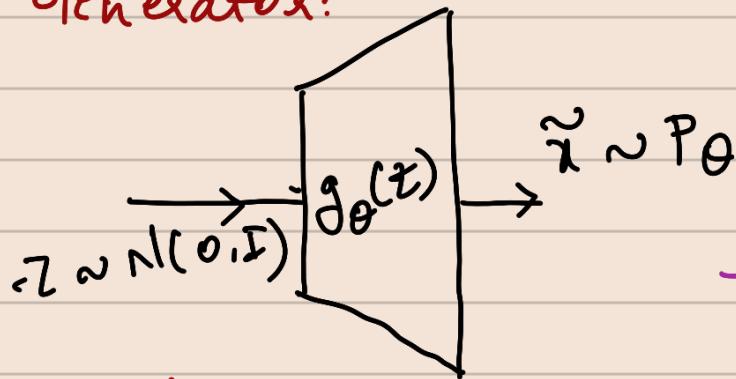
17:53, 25 May, Sun

Vanilla GAN:

Input:  $D = \{x_1, x_2, \dots, x_n\} \stackrel{iid}{\sim} P_x$

Data: MNIST  $28 \times 28$

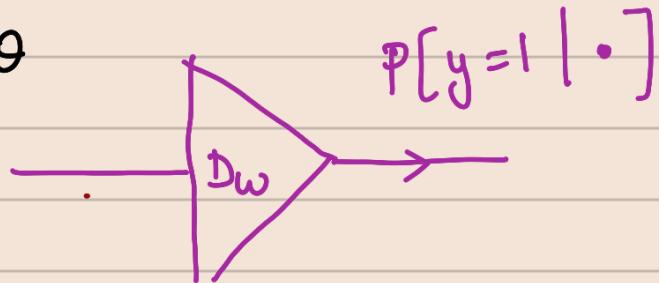
Generator:



Decide on z-dim: 100

Output: 784, then  
reshape to  $28 \times 28$

Discriminator



Input: 784  
Output:  $P(y=1 | \cdot)$

Loss function:

$$J(\theta, \omega) = \underset{P_X}{\mathbb{E}} \left[ \log D_\omega(x) \right] + \underset{P_\theta}{\mathbb{E}} \left[ -\log D_\omega(\hat{x}) \right]$$

(WHEN)

Approximate expectation using sample mean  
and use Batches.

$$J(\theta, \omega) = \frac{1}{B_1} \sum_{j=1}^{B_1} \log D_\omega(x_j) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_\omega(\hat{x}_j))$$

$$x_1, x_2, x_3, \dots, x_{B_1} \stackrel{iid}{\sim} P_X$$

$$\hat{x}_1, \hat{x}_2, \dots, \hat{x}_{B_2} \stackrel{iid}{\sim} P_\theta$$

$$J(\theta, \omega) = \frac{1}{B_1} \sum_{j=1}^{B_1} \log D_\omega(x_j) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_\omega(g_\theta(\hat{x}_j)))$$

$$z_1, z_2, \dots, z_{B_2} \stackrel{iid}{\sim} N(0, I)$$

Discriminator  $\pi/\omega$ :

$$\omega^* = \operatorname{argmax}_\omega \frac{1}{B_1} \sum_{j=1}^{B_1} \log (D_\omega(x_j)) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_\omega(g_\theta(z_j)))$$

$$\omega^{t+1} \leftarrow \omega^t + \alpha_1 \nabla_\omega J(\theta, \omega)$$

While training the discriminator the parameters of the generator network are

Kept constant.

\*-

Sample  $z_1, z_2, \dots, z_{B_2}$   $\stackrel{\text{iid}}{\sim} N(0, I)$   
Pass thru  $g_\theta(z)$

$g_\theta(z_1), g_\theta(z_2), \dots, g_\theta(z_{B_2})$

Pass these through  $D_w(\cdot)$  & obtain  
2<sup>nd</sup> term in loss.

\*

$x_1, x_2, \dots, x_{B_1} \sim p_x$

Pass through  $D_w(\cdot)$  & get the 1<sup>st</sup> term.

Total loss = term 1 + term 2

This is how the gradient ascent is implemented for the discriminator network.

Generator N/w:

$$\begin{aligned}\theta^* &= \arg \min_{\theta} J(\theta, w) \\ &= \arg \min_{\theta} \frac{1}{B_1} \sum_{i=1}^{B_1} \log(D_w(x_i)) \\ &\quad + \frac{1}{B_2} \sum_{j=1}^{B_2} \log(1 - D_w(g_\theta(z_j)))\end{aligned}$$

Independent of  $\theta$ .

$$\theta^{t+1} \leftarrow \theta^t - \alpha_2 \nabla_{\theta} J(\theta, w)$$

whole training the generator keep the params  
of the discriminator constant.