

Objective: Develop a Mathematical Foundation for Gen AI.

Family of DGM's 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_θ .
 P_θ : Represented using Deep NNs.
(Model)
- ii) Define & Compute a divergence metric between the P_θ & P_x . → We dunno P_x tho!?
- iii) Solve an optimization problem over the parameters of P_θ , 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)?

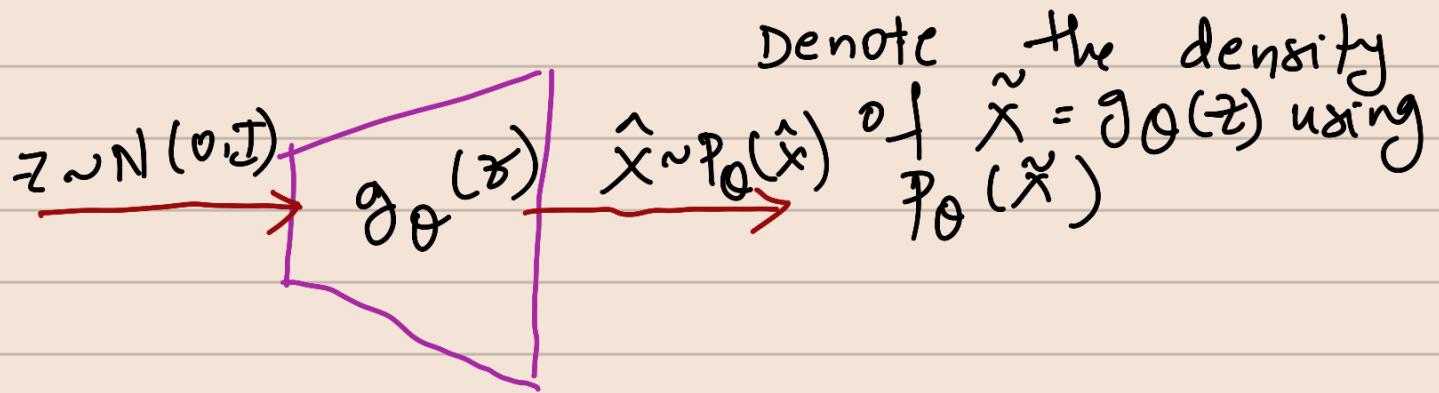
Example:

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



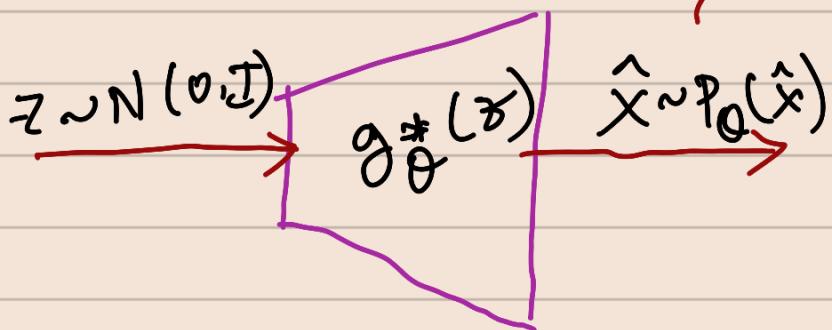
Suppose $D(P_x || P_\theta)$ denote a divergence measure between P_x & P_θ ,

Next step: solve the following optimization problem

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

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_θ & p_x .
- ② What should be the choice of the divergence metric \mathcal{D} ?

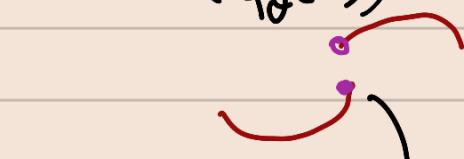
- ③ How to choose the $g_\theta(z)$, in turn P_θ ?
- ④ How to solve the optimization problem of minimizing the divergence metric?

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.

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

$$= D_{KL}$$

$\underbrace{D_{KL}(P_x || P_\theta)}_{\text{Forward KL}} \neq \underbrace{D_{KL}(P_\theta || P_x)}_{\text{Reverse KL}}$ 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) $f(u) = \frac{1}{2} [u \log u - (u+1) \log \left(\frac{u+1}{2} \right)]$: JS-Divergence.
 - c) $f(u) = \frac{1}{2} |u-1|$: Total Variational Distance.
- \hookrightarrow Used famously in GANs

Algorithm for F-Divergence Minimization.

Objective: Algorithm - to minimize D_f b/w

$P_x \not\subseteq P_\theta$, without knowing both,
but having samples from both.

Samples from P_x : dataset D

Samples from P_θ : 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
 $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} P_x$

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

Law of Large Numbers:

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

$x_i \sim \text{iid } P_X$

If the f -divergence metric can be expressed in terms of expectation of functions w.r.t $P_X \neq P_0$, then one can compute & optimize them.

Expressing D_f in terms of Expectations over $P_X \neq P_0$.

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

Even though this expression looks really similar to the $\int_x P_0(x)h(x).dx$, 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 prob density 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

$$\{ut - f(u)\}$$

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

Properties of Conjugate:

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

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

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

DETOUR DONE.

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

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

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

$$\begin{aligned}
 D_f(P_X || P_\theta) &= \int_X p_\theta(x) \cdot \sup_t (t u - f^*(t)) \cdot dx \\
 &= \int_X p_\theta(x) \sup_t \left\{ - \left[\frac{p_X(x)}{p_\theta(x)} \right] - f^*(t) \right\} p_\theta(x) \cdot dx \\
 &\quad \text{cannot directly pull out the } \sup.!!
 \end{aligned}$$

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

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

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

Space of functions
 containing solutions
 for the inner
 optimization problem

But,

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

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.

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

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

18/05/25, 02h5

Realization of VDM

Given Data $D = \{x_1, x_2, \dots, x_n\} \rightsquigarrow P_X$

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

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

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

but we are delegated 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 $T(x) \in \mathcal{T}$

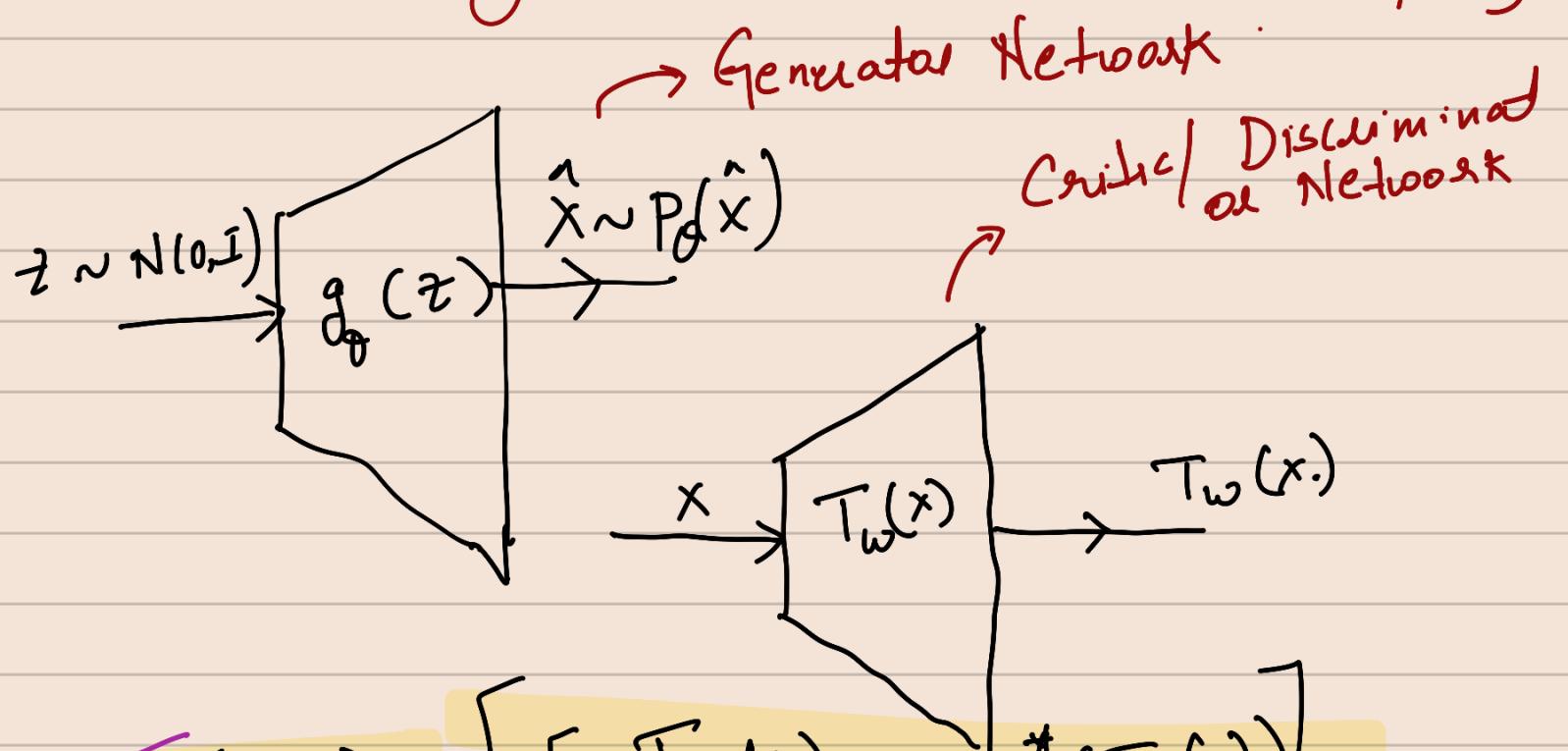
This cannot be done analytically?
Optimizing over space
of non-trivial functions.

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:

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

Implementing VDM for Generative Sampling.



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

$$\theta^*, \omega^* = \arg \min_{\theta} \max_{\omega} J(\theta, \omega)$$

↳ Saddle-point optimization

Problem where we have alternate minimization & maximization.

Blueprint for adversarial learning

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

19 May 2025, 10:21:06:

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

In practice, we do this by:

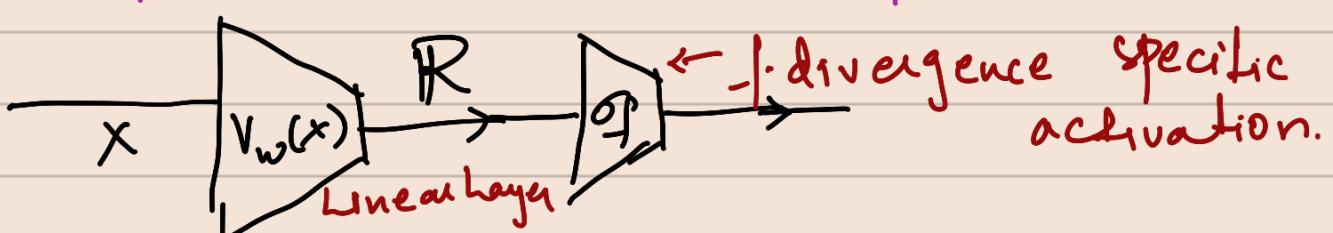
$T_w(x)$ is represented as: $\sigma_f(V_w(x))$, so that the range of T network corresponds to domain of f^* .

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

where σ_f is a \mathbb{f} -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 \mathbb{f} '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: $\text{wlogu} - (\text{uH})\log(\text{uH})$

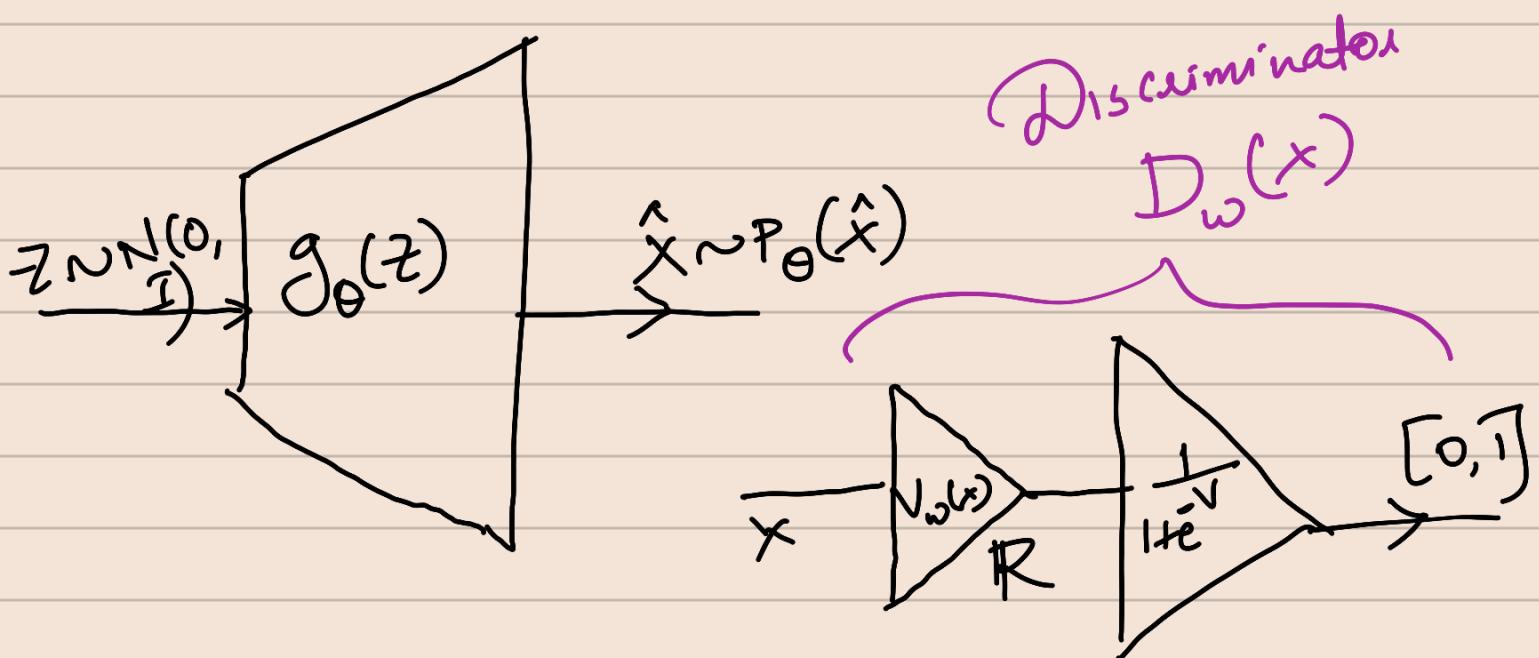
$$f^*(t) = -\log(1-e^t), \quad \text{dom } f^* = \mathbb{R}$$

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

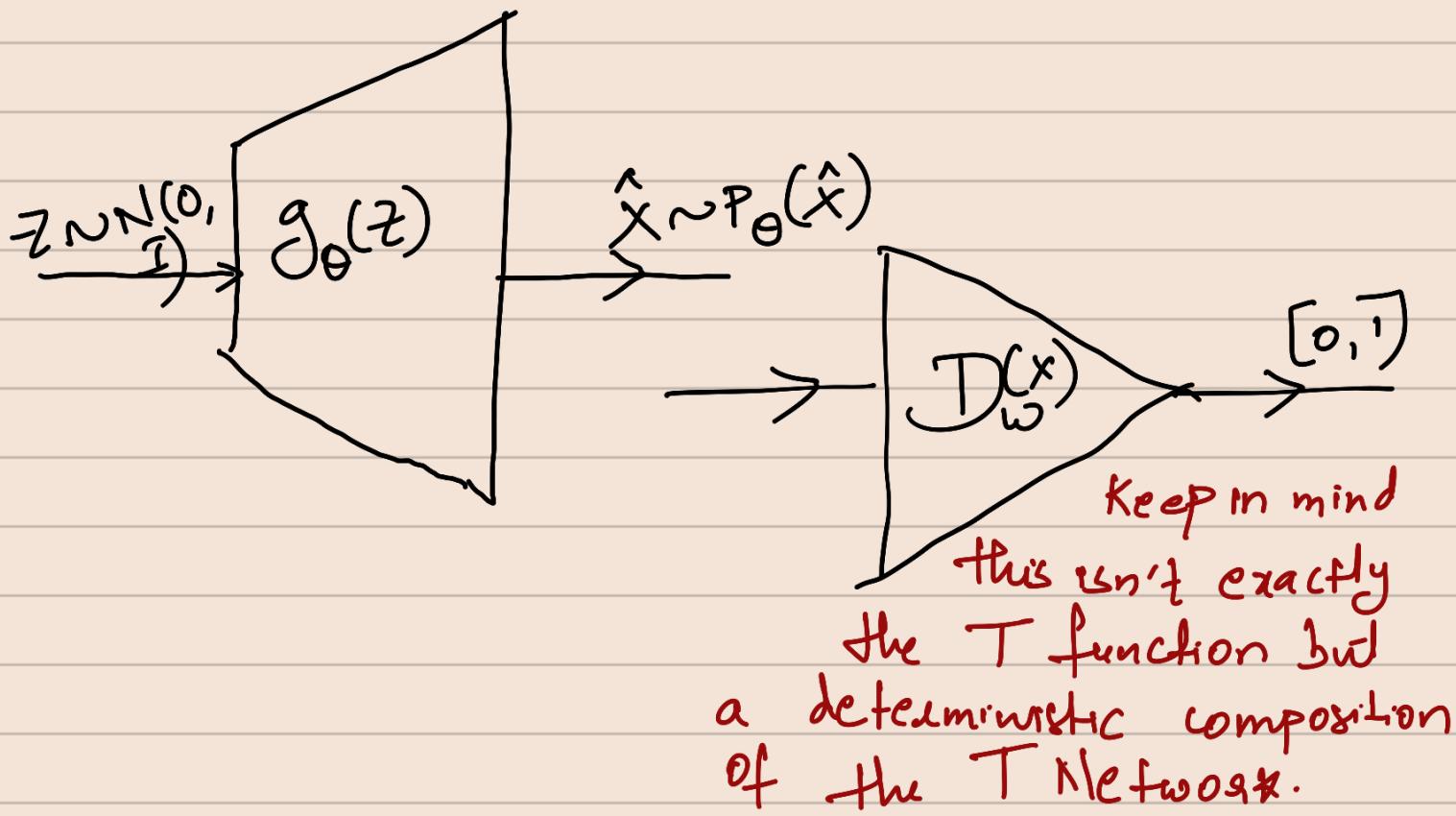
Anfully similar
to JS Div.

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

$$\text{where } D_\omega(x) := \frac{1}{1+e^{-V_\omega(x)}}$$



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

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

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

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

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

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

Independent of θ .

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

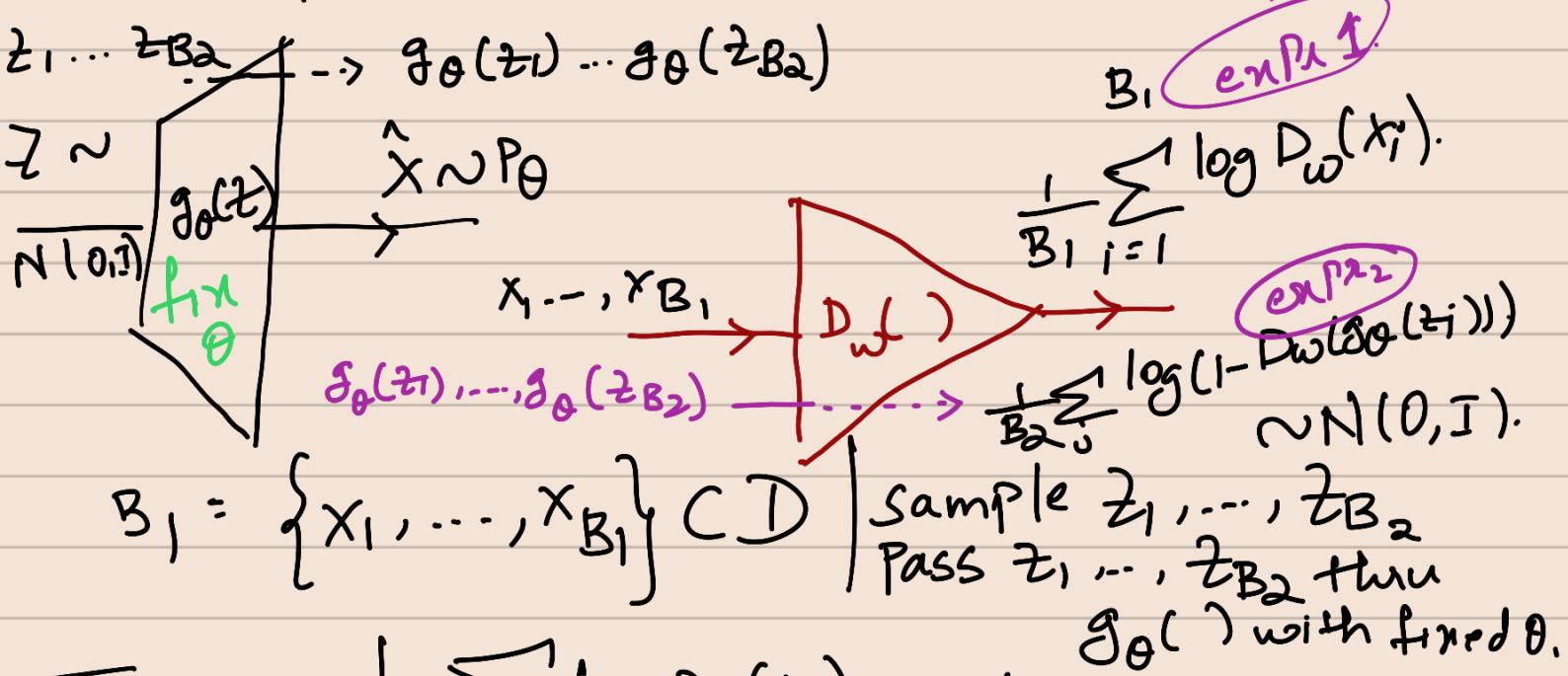
$$\approx \operatorname{argmin}_{\theta} \left[\frac{1}{B_2} \sum_{j=1}^{B_2} \log (1 - D_{\omega}(g_{\theta}(z_j))) \right]$$

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

ω is kept a constant

TO TRAIN THE DISCRIMINATOR:

keep θ constant



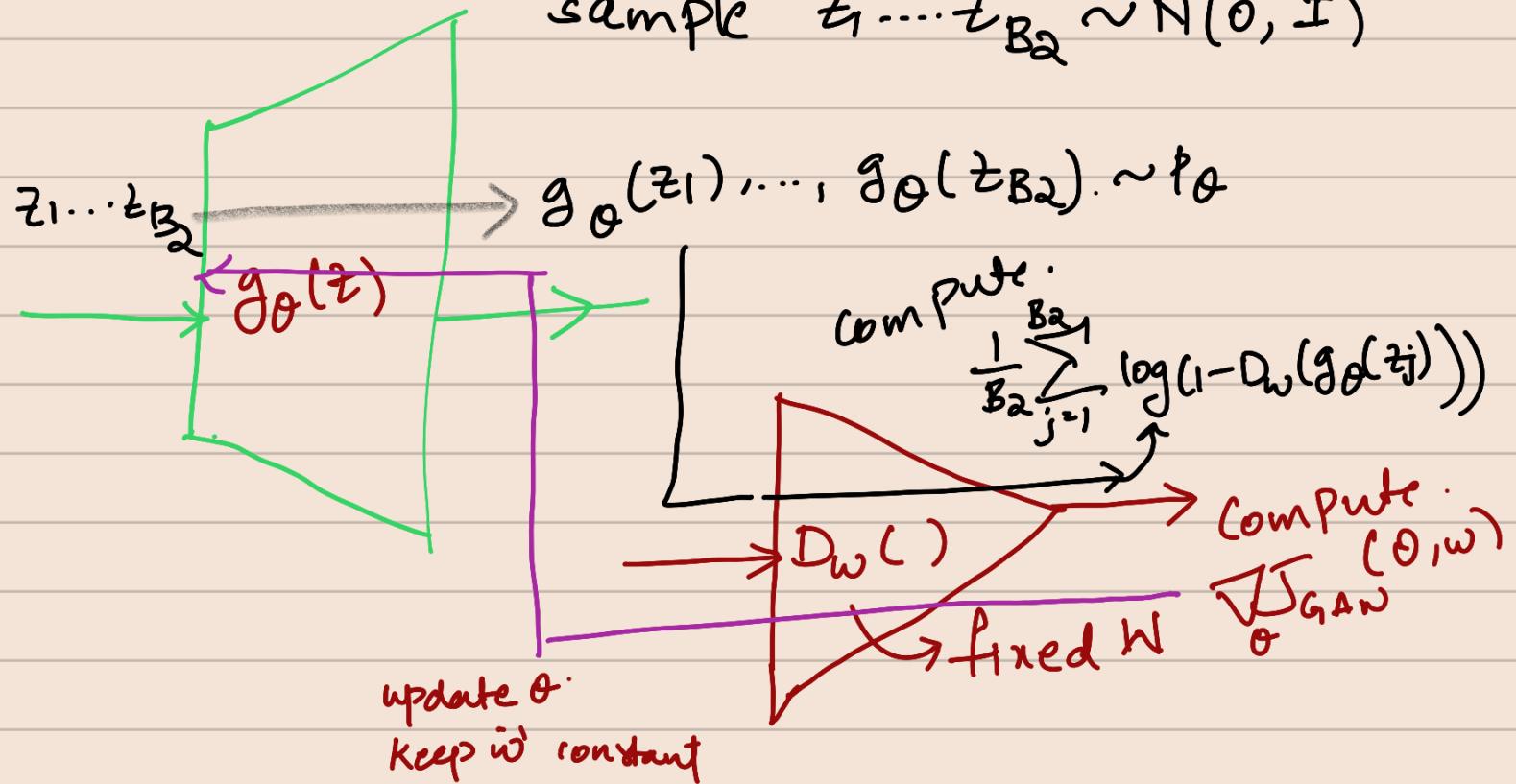
$$J_{GAN} = \frac{1}{B_1} \sum_{j=1}^{B_1} \log D_w(x_i) + \frac{1}{B_2} \sum_{j=1}^{B_2} \log(1 - D_w(\hat{x}_j))$$

$\underbrace{g_\theta(z_j)}$

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

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 1st 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 generator are reached.

Training VDM or GAN's done

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