

Deep Generative Models

Chapter 4: Generative Adversarial Networks

Ali Bereyhi

`ali.bereyhi@utoronto.ca`

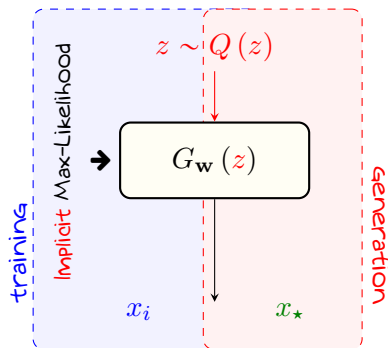
Department of Electrical and Computer Engineering
University of Toronto

Summer 2025

Modern Data Generation

As we mentioned at the end of Chapter 2

modern generative models mainly learn *how to sample!*



Low-dimensional Latent

From our discussions on **latent space**: we expect that the **latent** be of **much lower dimension** $m \ll d$

! We could **not** work with such latent in **flow-based models**

↳ The **flow** needs to be **invertible**

↳ **Invertibility** requires the **latent** and **data space** to be of **the same dimension!**

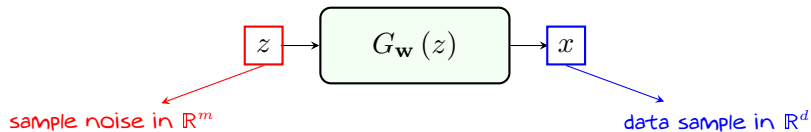
💡 We however **intuitively** expect that **a much smaller latent could do the job**

- + What if we make a model that maps a **low-dimensional latent** into the **data space**?
- Well! We can try; however, we get into **trouble training** it!
- + Can't we use simply **MLE**?!
- **Not explicitly!** Let's take a look

Generator: *Non-invertible Flow*

Generator Model

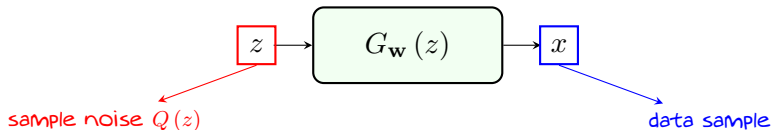
Generator model $G_{\mathbf{w}} : \mathbb{R}^m \mapsto \mathbb{R}^d$ is a mapping that maps a *latent sample* $z \sim Q(z)$, typically Gaussian noise, into a data sample



- + Isn't that what we called *flow*?!
- Not exactly! *Flow* is *invertible*

Sampling

We can sample from a **generator** exactly as in *flow-based models*



```
Sample_Generator():
```

- 1: Sample latent as $z \sim Q(z)$ #Gaussian noise
- 2: Pass through generator $x \leftarrow G_{\mathbf{w}}(z)$
- 3: return x

MLE Learning: Generator Distribution

Given $z \sim Q(z)$: *generated sample* $x = G_w(z)$ is distributed by some $P_w(x)$

$$\text{CDF}_w(x) = \Pr \{G_w(z) \leq x\}$$

The *generator distribution* $P_w(x)$ is then given by

$$P_w(x) = \frac{\partial^d}{\partial x_1 \dots \partial x_d} \text{CDF}_w(x)$$

Since G_w is *not invertible*, we should define

$$\mathbb{Z}_w(x) = \{z \in \mathbb{R}^m : G_w(z) \leq x\}$$

and compute the *generator distribution* $P_w(x)$ as

$$\text{CDF}_w(x) = \int_{\mathbb{Z}_w(x)} P(z) dz$$

MLE Learning: Generator Distribution

The *generator distribution* is hence given by

$$P_{\mathbf{w}}(x) = \frac{\partial^d}{\partial x_1 \dots \partial x_d} \int_{\mathbb{Z}_{\mathbf{w}}(x)} P(z) dz$$

which is a *challenging* computation!

- + Isn't this *simply* computed by *chain rule*?!
 - *Not* that *simple*! As $G_{\mathbf{w}}$ is *not* invertible, the set $\mathbb{Z}_{\mathbf{w}}(x)$ *cannot* be *simply* characterized! Let's see an example

Recall

When $G_{\mathbf{w}}$ is *invertible*, we can *readily* specify $\mathbb{Z}_{\mathbf{w}}(x)$ using $G_{\mathbf{w}}^{-1}$

$$\mathbb{Z}_{\mathbf{w}}(x) = G_{\mathbf{w}}^{-1} \left(\left\{ u \in \mathbb{R}^d : u \leq x \right\} \right)$$

Generator Distribution: *Example*

Consider a dummy example: $z \in \mathbb{R}$ and $x \in \mathbb{R}^3$ with *generator*

$$G(z) = [-z^2, \exp\{-z\}, z]$$

Say, want to find *the CDF* at $x = [x_1, x_2, x_3]$: we need to find

$$\begin{aligned} \mathbb{Z}(x) &= \{z \in \mathbb{R} : G(z) \leq x\} \\ &= \{z \in \mathbb{R} : -z^2 \leq x_1 \text{ and } \exp\{-z\} \leq x_2 \text{ and } z \leq x_3\} \end{aligned}$$

① *First constraint requires*

$$-z^2 \leq x_1 \rightsquigarrow z^2 \geq -x_1 \rightsquigarrow \begin{cases} z \geq \sqrt{-x_1} \\ z \leq -\sqrt{-x_1} \end{cases}$$

Generator Distribution: *Example*

We need to find

$$\mathbb{Z}(x) = \{z \in \mathbb{R} : -z^2 \leq x_1 \text{ and } \exp\{-z\} \leq x_2 \text{ and } z \leq x_3\}$$

② *Second constraint requires*

$$\exp\{-z\} \leq x_2 \rightsquigarrow -z \leq \log x_2 \rightsquigarrow z \geq -\log x_2$$

③ *And, finally the **third** one restricts z to satisfy*

$$z \leq x_3$$

Generator Distribution: *Example*

So, at $x = [x_1, x_2, x_3]$ we have

$$\mathbb{Z}(x) = ([-\infty, -\sqrt{-x_1}] \cup [\sqrt{-x_1}, +\infty]) \cap [-\log x_2, \infty] \cap [-\infty, x_3]$$

► At $x = [-1, \exp\{2\}, 2]$, we have

$$\begin{aligned}\mathbb{Z}(x) &= ([-\infty, -1] \cup [1, \infty]) \cap [-2, \infty] \cap [-\infty, 2] \\ &= [-2, -1] \cup [1, 2]\end{aligned}$$

► At $x = [-1, 1, 2]$, we have

$$\mathbb{Z}(x) = [1, 2]$$

► At $x = [a, 1, 2]$ with $a > 0$, we have

$$\mathbb{Z}(x) = \emptyset$$

Generator Distribution: *Example*

To find *generator distribution*, we need $\text{CDF}(x)$ at all $x \in \mathbb{R}^3$; however,

computation of $\text{CDF}(x)$ changes as x changes

For instance, in our earlier sample points, we have

$$\text{CDF}(x) = \begin{cases} \int_{-2}^{-1} P(z) dz + \int_1^2 P(z) dz & x = [-1, \exp\{2\}, 2] \\ \int_1^2 P(z) dz & x = [-1, 1, 2] \\ 0 & x = [a, 1, 2] \text{ and } a > 0 \\ \vdots & \vdots \end{cases}$$

Expressing $P(x)$ for all $x \in \mathbb{R}^3$ is quite *cumbersome!*

MLE Challenge: Non-Computability of Likelihood


Computability of Generator Distribution

With a **non-invertible** generator $G_{\mathbf{w}}$, it is challenging to characterize

$$\mathbb{Z}_{\mathbf{w}}(x) = \{z \in \mathbb{R}^m : G_{\mathbf{w}}(z) \leq x\}$$

and thus the output distribution $P_{\mathbf{w}}(x)$ is hard to be computed for all $x \in \mathbb{R}^d$

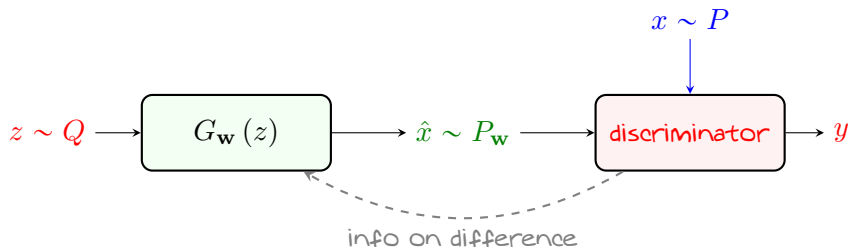
To train the **generator** by MLE \rightsquigarrow we need **likelihood** at **all samples** x^j

$$\hat{R}(\mathbf{w}) = -\frac{1}{n} \sum_j \log P_{\mathbf{w}}(x^j)$$


Non-Computability of Likelihood

With **non-invertible** generator, **likelihood** is **not** directly computable

Alternative Learning Approach: Adversarial Architecture



Discriminator

Discriminator is a binary classifier that classifies sample x as **real** or **fake**

To train with discriminator: we can **easily** make a **labeled** dataset

- **Data samples** are **labeled** as **real**
- **Generated samples** are **labeled** as **fake**

Dataset for Adversarial Architecture

From **discriminator** viewpoint: we train a **classifier** on $2n$ **labeled** samples

$$\mathbb{D} = \left\{ \boxed{(x^j, 1)}, \boxed{(\hat{x}^j, 0)} : j = 1, \dots, n \right\}$$

← real data sample → fake generated sample

We can compactly write the **dataset** as

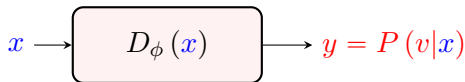
$$\mathbb{D} = \left\{ (x^j, v^j) : j = 1, \dots, 2n \right\}$$

where the **true label** v is defined as

$$v = \begin{cases} 0 & x \text{ if } \text{fake} \equiv x = G_{\mathbf{w}}(z) \\ 1 & x \text{ if } \text{real} \equiv x \sim P(x) \end{cases}$$

Training Computational Discriminator

We now consider a **computational discriminator**: a model $D_\phi : \mathbb{R}^d \mapsto [0, 1]$ which classifies sample x



Recall

This is a **discriminative** model! This explains the **appellation** 😊

This is a classical **classification problem** \rightsquigarrow we use **cross-entropy loss**

$$\begin{aligned} \text{CE}(y, v) &= -v \log y - (1 - v) \log (1 - y) \\ &= \begin{cases} -\log y & \text{if } x \text{ is real} \equiv v = 1 \\ -\log (1 - y) & \text{if } x \text{ is fake} \equiv v = 0 \end{cases} \end{aligned}$$

Empirical Risk of Adversarial Architecture

To train on this **labeled dataset** \rightsquigarrow we minimize the **empirical risk**, i.e.,

$$\begin{aligned}
 \hat{R}(\mathbf{w}, \phi) &= \frac{1}{2n} \sum_j \text{CE}(y^j, v^j) \\
 &= \frac{1}{2n} \sum_j \text{CE}(D_\phi(x^j), v^j) \\
 &= \frac{1}{n} \sum_j \text{CE}(D_\phi(x^j), 1) + \frac{1}{n} \sum_j \text{CE}(D_\phi(G_{\mathbf{w}}(z^j)), 0) \\
 &= \hat{\mathbb{E}}_{x \sim P} \{ \text{CE}(D_\phi(x), 1) \} + \hat{\mathbb{E}}_{z \sim Q} \{ \text{CE}(D_\phi(G_{\mathbf{w}}(z)), 0) \} \\
 &= \hat{\mathbb{E}}_{x \sim P} \{ -\log D_\phi(x) \} + \hat{\mathbb{E}}_{z \sim Q} \{ -\log(1 - D_\phi(G_{\mathbf{w}}(z))) \} \\
 &= - \left[\hat{\mathbb{E}}_{x \sim P} \{ \log D_\phi(x) \} + \hat{\mathbb{E}}_{z \sim Q} \{ \log(1 - D_\phi(G_{\mathbf{w}}(z))) \} \right]
 \end{aligned}$$

Empirical risk depends on both **generator** and **discriminator**


Training by Min-Max Game: *Discriminator's Role*

To train the **adversarial architecture** by the **empirical risk**

$$\hat{R}(\mathbf{w}, \phi) = - \left[\hat{\mathbb{E}}_{x \sim P} \{ \log D_{\phi}(x^j) \} + \hat{\mathbb{E}}_{z \sim Q} \{ \log (1 - D_{\phi}(G_{\mathbf{w}}(z^j))) \} \right]$$

we let the **discriminator** and **generator** play a game

Discriminator tries its best to **correctly classify fake** sample from **true ones**, i.e.,

$$\begin{aligned} \min_{\phi} \hat{R}(\mathbf{w}, \phi) &\equiv \max_{\phi} \left[\hat{\mathbb{E}}_{x \sim P} \{ \log D_{\phi}(x^j) \} + \hat{\mathbb{E}}_{z \sim Q} \{ \log (1 - D_{\phi}(G_{\mathbf{w}}(z^j))) \} \right] \\ &\equiv \max_{\phi} \mathcal{L}(\mathbf{w}, \phi) \end{aligned}$$


After maximization, **discriminator's best try** gets the objective¹

$$\mathcal{L}^*(\mathbf{w}) = \max_{\phi} \mathcal{L}(\mathbf{w}, \phi)$$

¹Recall that $\mathcal{L}(\mathbf{w}, \phi)$ is indeed the **classifier likelihood**!

Training by Min-Max Game: Generator's Role

Generator tries its best to *confuse discriminator*

💡 This way *discriminator cannot* distinguish between *fake* and *true samples*

↳ This means that *fake samples* look the same as the *true ones*

💡 The *generator* can *only* play with its *parameters w*

Generator tries to make *discriminator's best try as bad as possible*, i.e.,

$$\min_w \mathcal{L}^*(w) = \min_w \max_{\phi} \mathcal{L}(w, \phi)$$

Statistical Perspective

Generator minimizes *discriminator's maximal likelihood* \rightsquigarrow *discriminator* remains *confused* between *generated* and *true* samples even if it does *its best*

GAN: Generative Adversarial Network

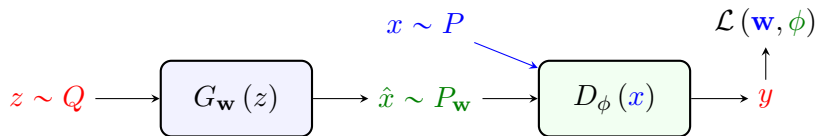
Generative Adversarial Network

GAN consists of a **generator** $G_{\mathbf{w}}$ and **discriminator** D_{ϕ} that are jointly learned through the **min-max game**

$$\min_{\mathbf{w}} \max_{\phi} \mathcal{L}(\mathbf{w}, \phi)$$

with objective function

$$\mathcal{L}(\mathbf{w}, \phi) = \hat{\mathbb{E}}_{x \sim P} \{\log D_{\phi}(x)\} + \hat{\mathbb{E}}_{z \sim Q} \{\log (1 - D_{\phi}(G_{\mathbf{w}}(z)))\}$$



GAN: Training Loop

Train_vanillaGAN(\mathbb{D} :dataset):

```

1: Initiate the generator  $G_{\mathbf{w}}$  and discriminator  $D_{\phi}$  with some  $\mathbf{w}$  and  $\phi$ 
2: for multiple epochs do
3:   Sample a batch of data samples  $\{x^j : j = 1, \dots, n\}$  from  $\mathbb{D}$ 
4:   Sample latents  $\{z^j : j = 1, \dots, n\}$  using  $Q$  and compute  $\hat{x}^j \leftarrow G_{\mathbf{w}}(z^j)$ 
5:   for  $\xi = 1, \dots, \Xi$  do
6:     for  $j = 1, \dots, n$  do
7:       Compute  $\mathcal{L}^j = \log D_{\phi}(x^j) + \log(1 - D_{\phi}(\hat{x}^j))$ 
8:       Backpropagate over discriminator to compute  $\nabla_{\phi} \mathcal{L}^j$ 
9:       if  $\xi = \Xi$  then Backpropagate over generator to compute  $\nabla_{\mathbf{w}} \mathcal{L}^j$ 
10:    end for
11:    Update  $\phi$  using  $\text{Opt\_avg} \{ \nabla_{\phi} \mathcal{L}^j \}$ 
12:  end for
13:  Update  $\mathbf{w}$  using  $\text{Opt\_avg} \{ \nabla_{\mathbf{w}} \mathcal{L}^j \}$ 
14: end for
15: return trained generator  $G_{\mathbf{w}}$ 

```

GAN Training: Few Notes

There are a few tricks to make training work

- 1 We update **discriminator** for Ξ steps **before** updating **generator once**

↳ This allows the **inner maximization** to slightly converge

$$\min_{\mathbf{w}} \max_{\phi} \mathcal{L}(\mathbf{w}, \phi)$$

↳ Typical choices are $5 \leq \Xi \leq 10$

- 2 To update **generator**, we only need to compute derivative of second term

$$\nabla_{\mathbf{w}} \mathcal{L}^j = \underbrace{\nabla_{\mathbf{w}} \log D_{\phi}(x^j)}_0 + \nabla_{\mathbf{w}} \log(1 - D_{\phi}(G_{\mathbf{w}}(z^j)))$$

↳ This gives **very small gradients first** and **exploding later**

↳ **First** $D_{\phi}(G_{\mathbf{w}}(z^j)) \approx 0$ and later when it's **fooled** $D_{\phi}(G_{\mathbf{w}}(z^j)) \approx 1$

↳ This is **opposite** of what we want!

↳ Heuristic remedy is to use $-\nabla_{\mathbf{w}} \log D_{\phi}(G_{\mathbf{w}}(z^j))$ instead

↳ Some people call $\mathcal{L}_G^j = -\log D_{\phi}(G_{\mathbf{w}}(z^j))$ **generator loss**

GAN: Sampling

Sampling is *exactly* as in *flow-based models*

Sample_GAN():

- 1: Sample latent as $z \sim Q(z)$ #Gaussian noise
- 2: Pass through generator $x \leftarrow G_{\mathbf{w}}(z)$
- 3: return generated sample x

Vanilla GAN

- + Why you called the training “*vanilla*”?!
 - Because we used a *basic loss* for our *adversarial network*
- + Can we do *better*?
 - Yes! That's what we do in *Wasserstein GAN*!

Though working with *tricks*, the *vanilla* GAN training is *generally unstable*

Wasserstein GAN addresses this issue by a very *smart trick*

To understand the *trick* used by *Wasserstein GAN*, we need to first learn the *statistical interpretation* of *vanilla* GAN training

Next Stop: Interpreting GAN training as *divergence minimization*