

# Generative Deep Learning

Sudong Lee

 [sudonglee@ulsan.ac.kr](mailto:sudonglee@ulsan.ac.kr)

 <https://dais.ulsan.ac.kr/>

# Contents

- Introduction to Generative Deep Learning
- Variational Autoencoders (VAEs)
- Generative Adversarial Networks (GANs)
- Diffusion Models

# Introduction to Generative Deep Learning

# A Probabilistic View of Deep Learning: Discriminative vs. Generative Models

“Discriminative and generative models are two fundamental approaches in deep learning, each with distinct characteristics and applications.”

## ▪ Discriminative models

- Learn the *decision boundary* (e.g., classification and regression)
- Focus on the inference of the output,  $Y$ , given the input,  $X$ .

⇒ Model the *conditional probability* of the output given the input,  $P(Y|X)$ .

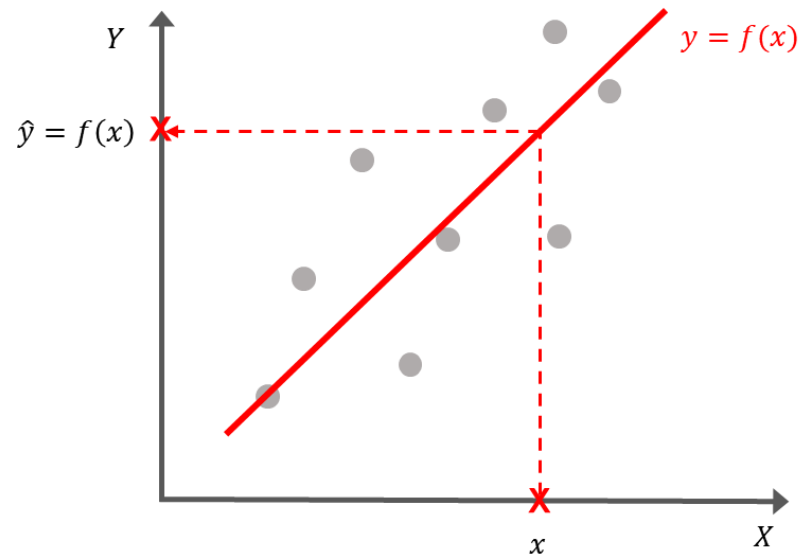
## ▪ Generative models

- Learn the underlying distribution of data.
- Focus on the generation of new data instances.

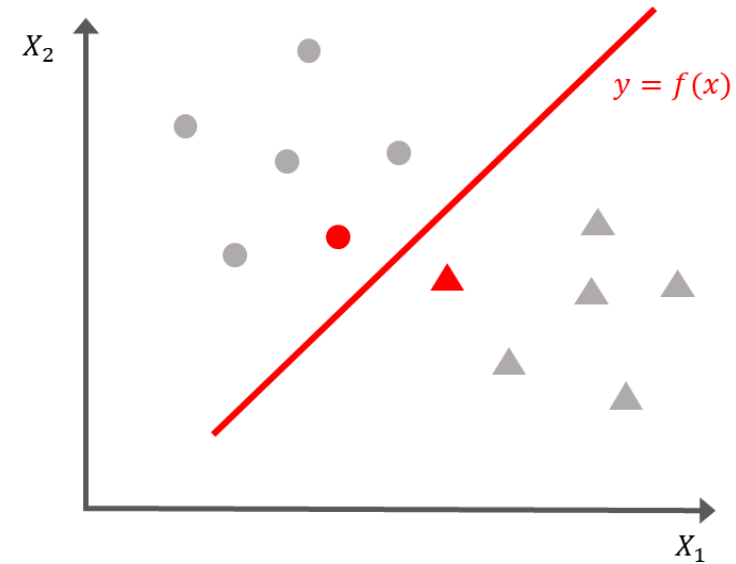
⇒ Model the *joint probability*  $P(X, Y)$ . Depending on the problem,  $P(X|Y)$  or  $P(X)$  can be of interest.

# Discriminative Models

“Discriminative models learn the **decision boundary** to infer the output given the input.”



Regression



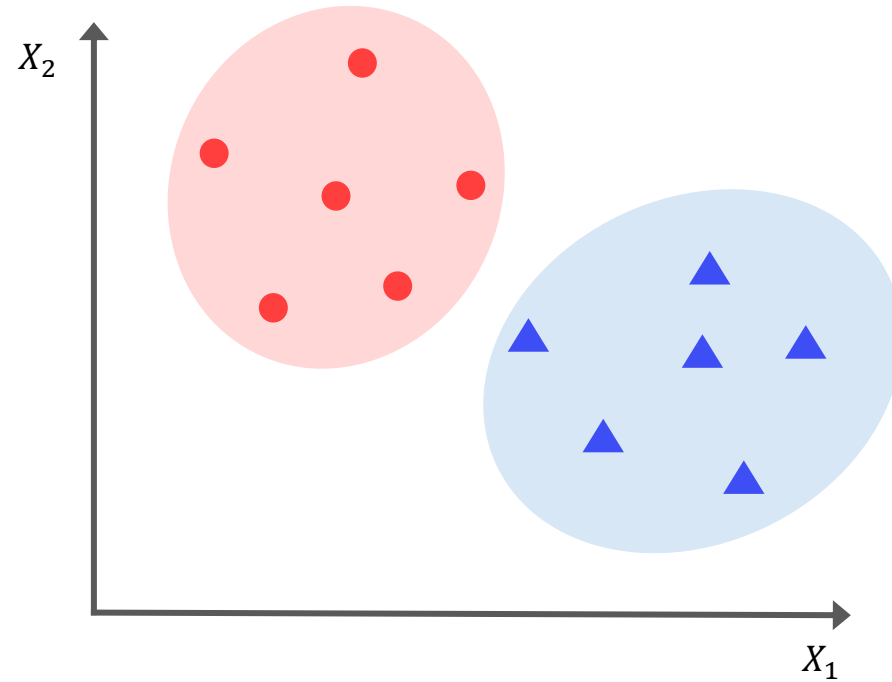
Classification

# Applications of Discriminative Models

- ANY classification and regression
- Image classification, object detection, and segmentation
- Natural Language Processing tasks (e.g., sentiment analysis, text classification, ...)
- Anomaly detection
- Medical diagnosis and disease prediction
- ...

# Generative Models

“Generative models learn the underlying distribution of data enabling them to generate new data instances.”



# Applications of Generative Models

- **Data generation:**

sample unobserved but realistic new data points (e.g., image/text generation)

- **Density estimation:**

predict the rareness of events (e.g., anomaly detection)

- **Latent variable extraction:**

extract low-dimensional latent space for dimensionality reduction (e.g., feature embedding)

- **Simulation and forecasting**

in various fields (e.g., digital twins, weather prediction, financial modeling, etc.)

- ...



# A Different Perspective on Generative Models

“While you learn to generate, you can learn the patterns in your data.”



Richard Phillips Feynman

The esteemed scientist *Richard Feynman* once said,  
“What I cannot create, I do not understand.”

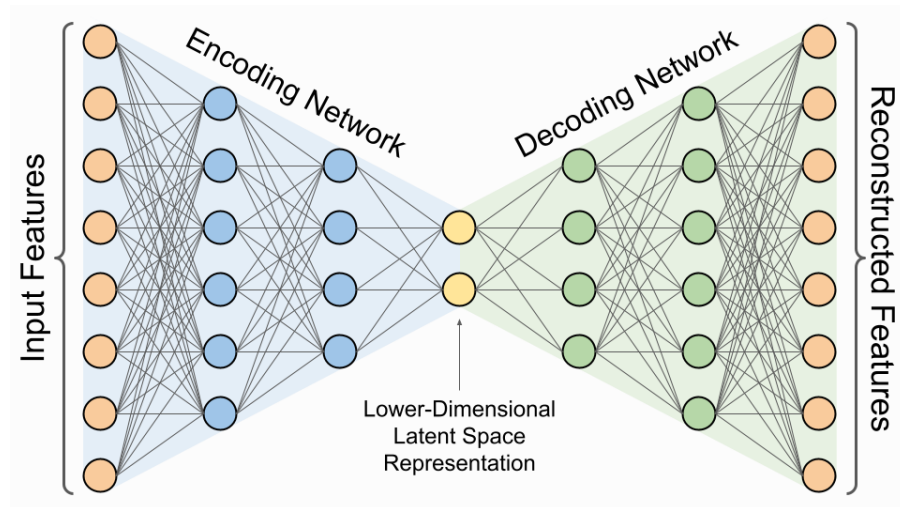
# Examples of Deep Generative Models

- Variational Autoencoders (VAEs)
- Generative Adversarial Networks (GANs)
- Autoregressive models (e.g., Transformers like GPT)
- Flow-based models
- Diffusion models
- ...

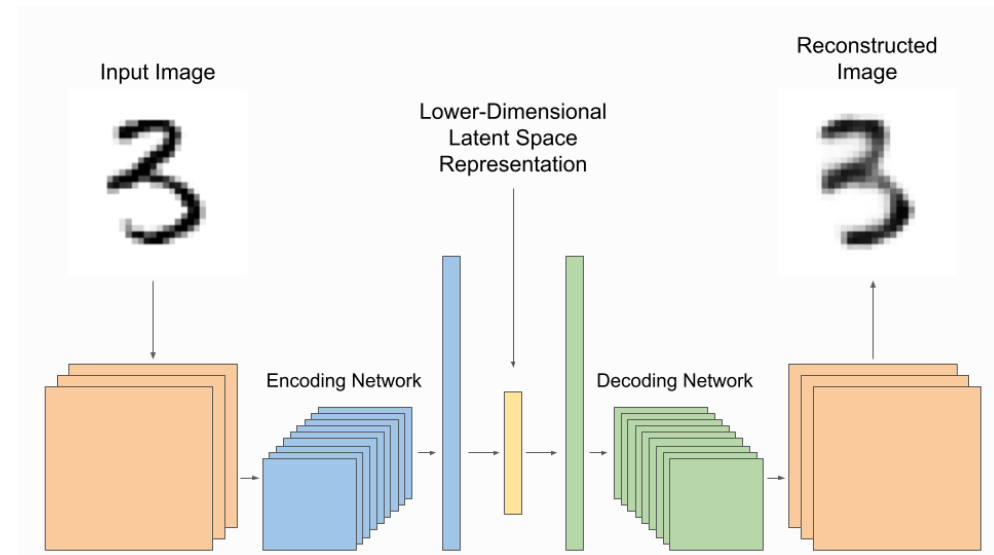
# Variational Autoencoders (VAEs)

# Autoencoder

“An autoencoder compresses the features into a lower dimension by **reconstruction**.”



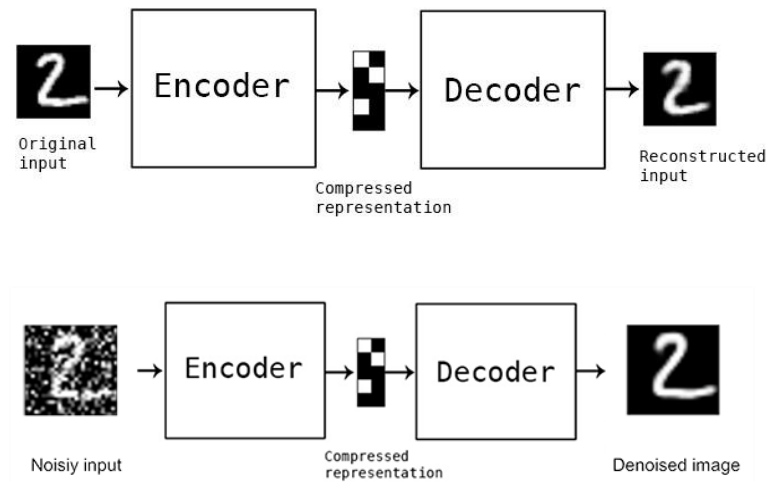
Network architecture for an ordinary autoencoder



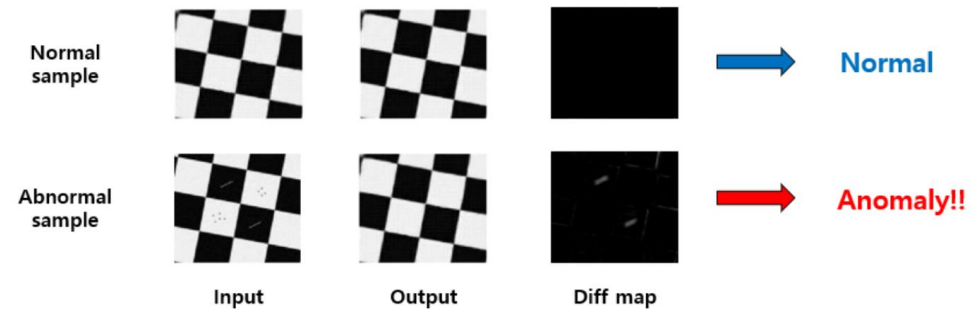
Network architecture for an ordinary convolutional autoencoder

# Autoencoder

“We can use autoencoders for various applications such as denoising and anomaly detection.”



Denoising autoencoder



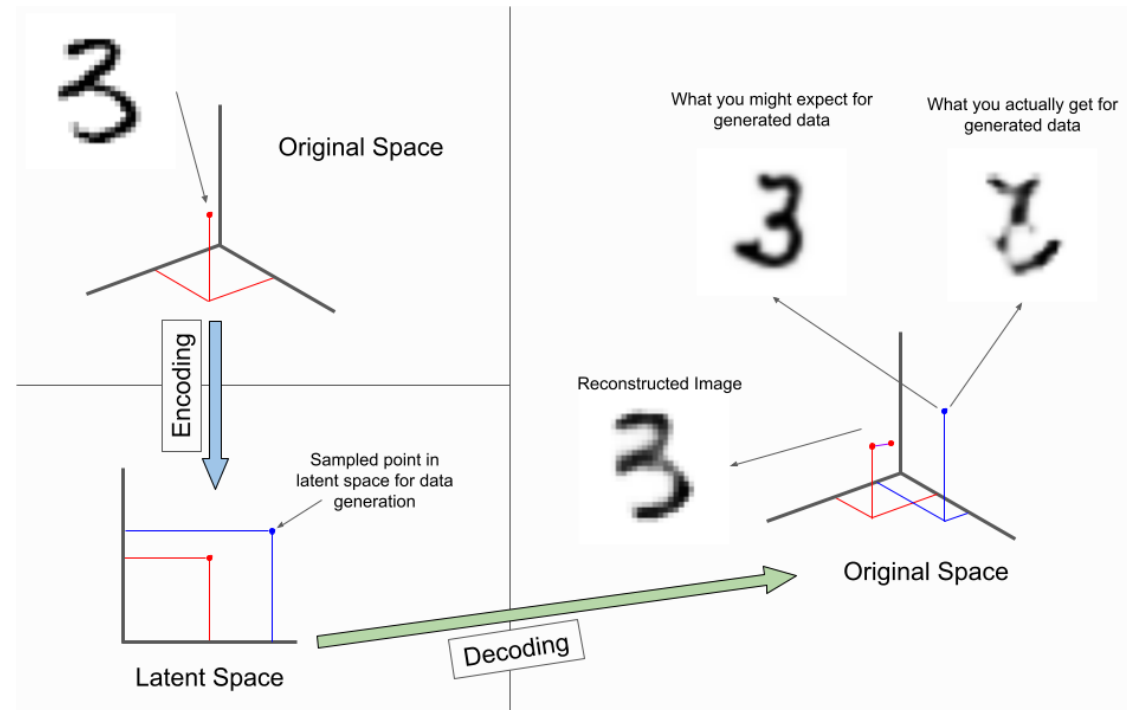
Autoencoder for anomaly detection

# Autoencoders for Data Generation

“Can autoencoders be used for data generation?”

- “Nope. They optimize for *faithful reconstructions*.”

Ordinary autoencoders learn to use the latent space as an embedding space to create **optimal compressions**.



# Manifold Learning

“Manifold learning aims to learn the mapping of high-dimensional data, potentially existing non-linear manifolds,  
onto lower-dimensional latent space.”

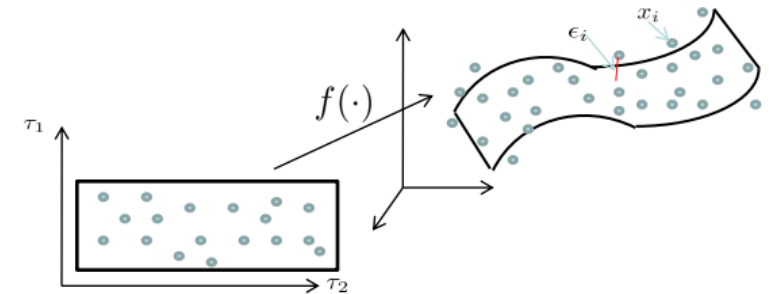
- Definition. *Manifold Learning*

- A  $d$  dimensional manifold  $\mathcal{M}$  is embedded in an  $m$  dimensional space, and there is an explicit mapping  $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$  where  $d \leq m$ .
- Given samples  $x_i \in \mathbb{R}^m$ , we assume that

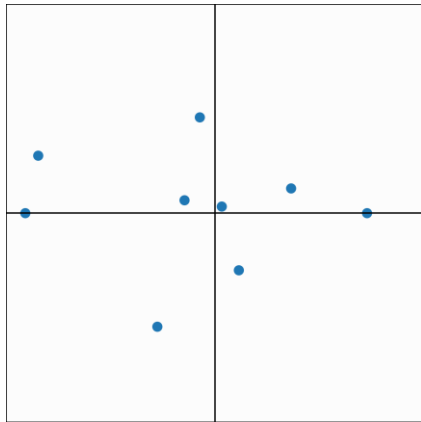
$$x_i = f(\tau_i) + \epsilon_i$$

where  $f$  is called mapping function,  $\tau_i \in \mathbb{R}^d$  is the embedded feature of  $x_i$ , and  $\epsilon_i$  is the random noise.

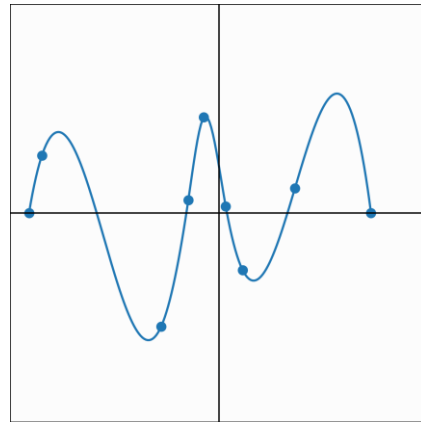
- We assume  $p(\tau)$  is smooth, is distributed uniformly, and noise is small.



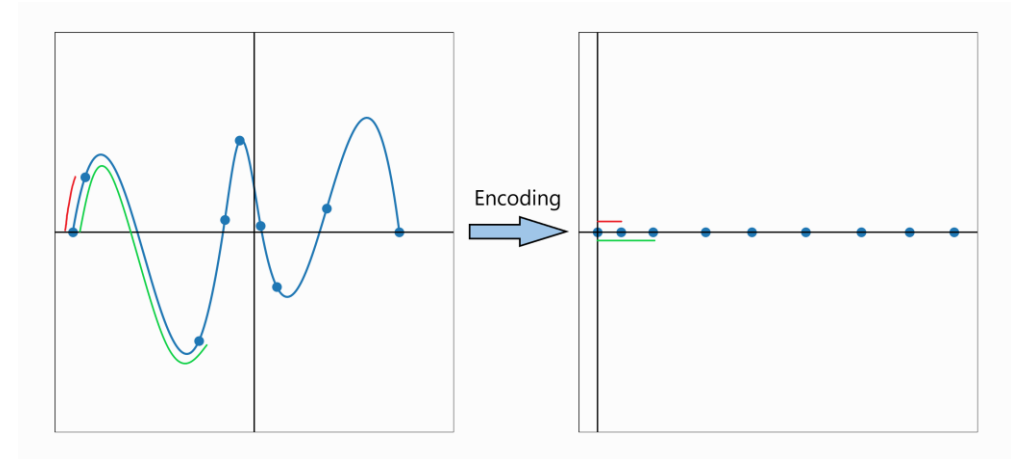
# Manifold Learning by Autoencoders as a Compression



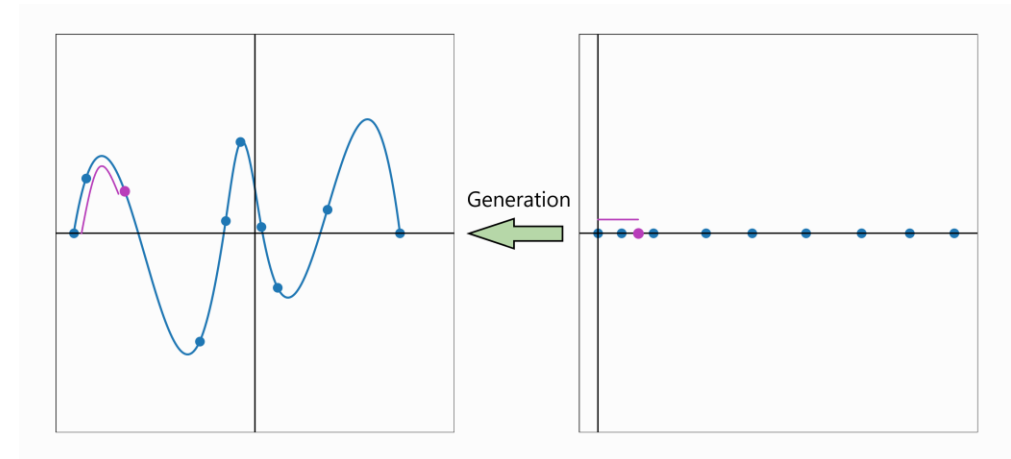
Observations



Hypothesized Manifold



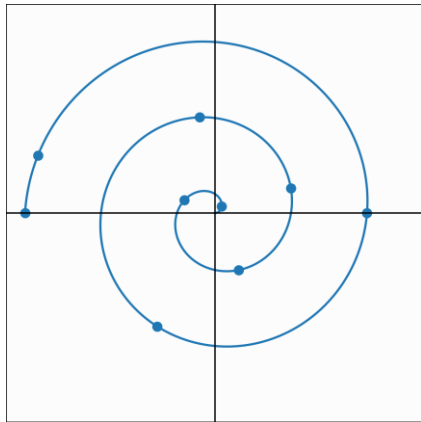
Encoder



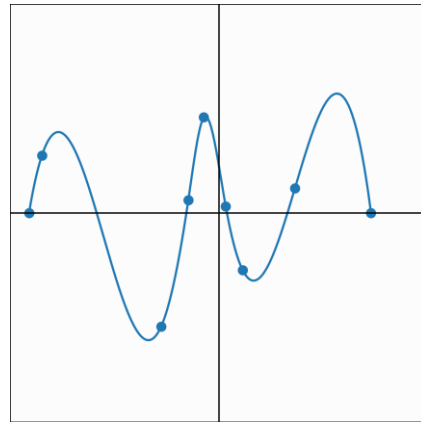
Decoder



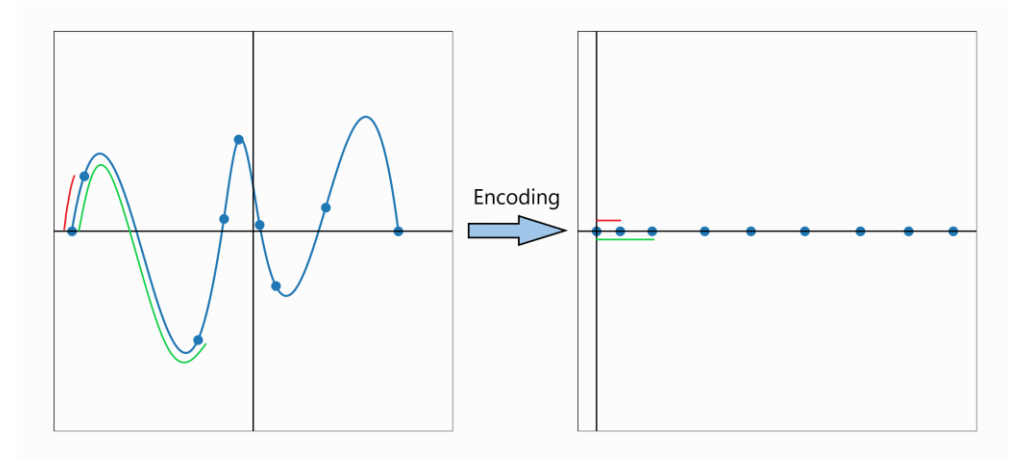
# Manifold Learning by Autoencoders as a Compression



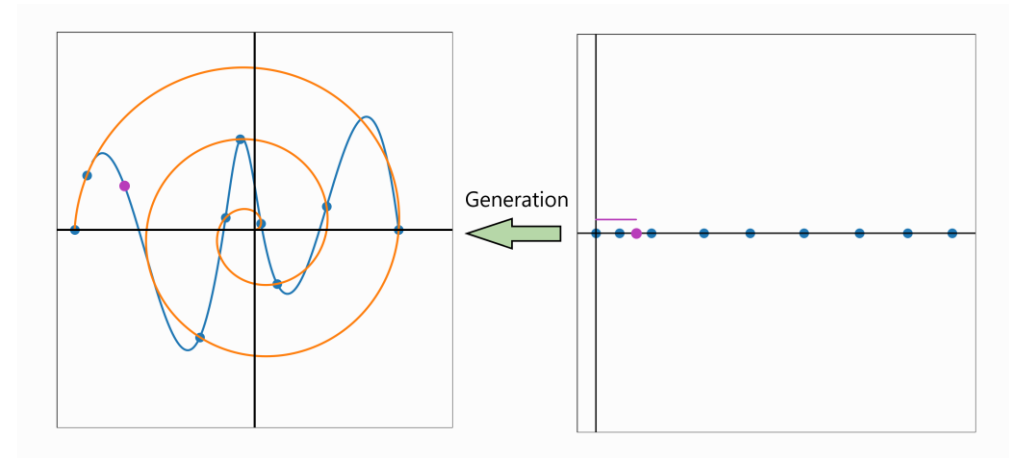
Actual Manifold



Hypothesized Manifold



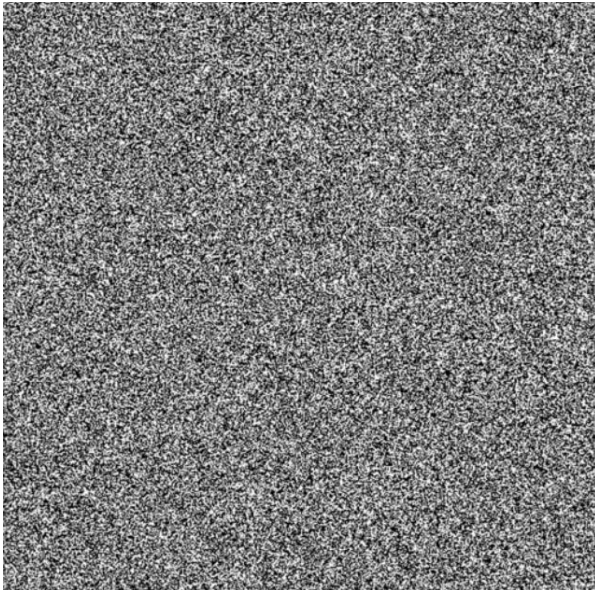
Encoder



Decoder

# Data Generation from Noise

Random Noise Sampling

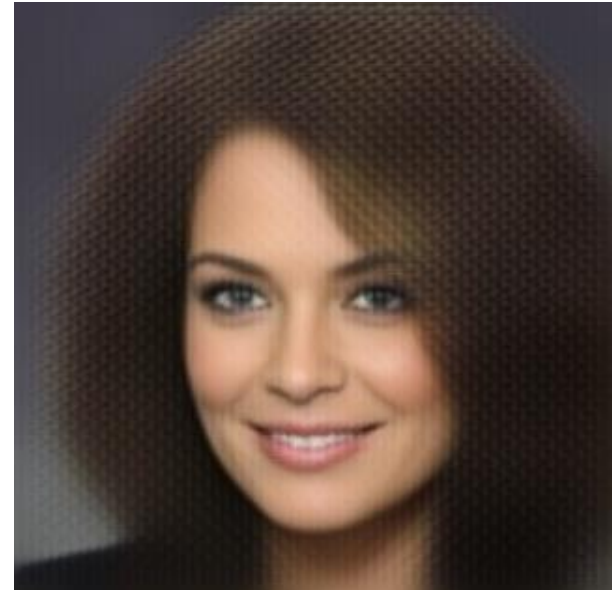


$Z$

$P(X|Z)$



Image Generation



$X$

# Data Generation from Noise

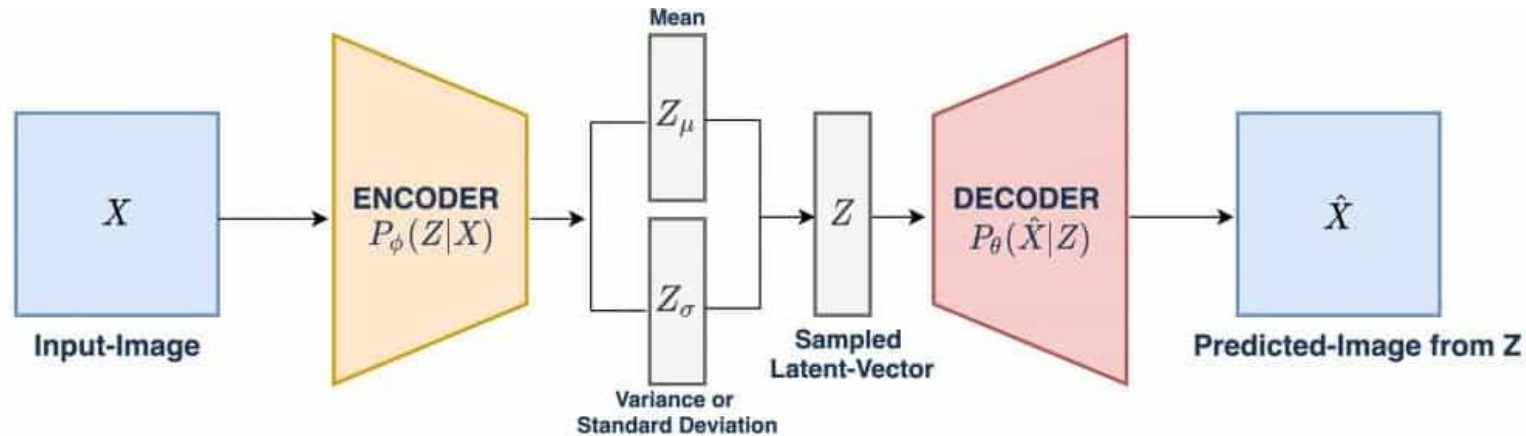
“How can we effectively encode the *variations* of generated data?”



	<i><b>Z</b></i>	<i><b>Z</b></i>	<i><b>Z</b></i>
Smile	-0.21	0.03	0.56
Skin tone	0.72	0.72	0.72
Gender	-0.63	-0.63	-0.63
Glasses	-0.43	-0.43	-0.43
Hair color	0.44	0.44	0.44

# Variational Autoencoder (VAE)

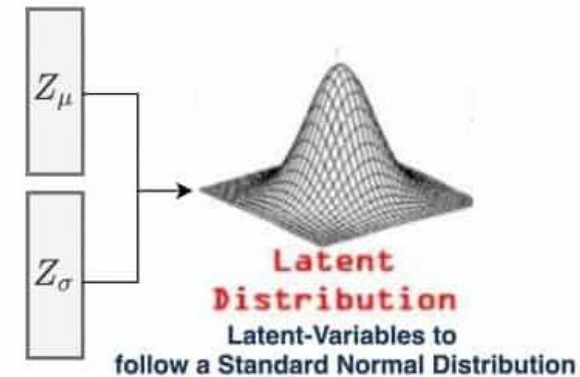
“VAEs encode inputs as **distributions** and **sample** from them to produce input **variations** for data generation.”



Sample a point from  $G(Z_\mu, Z_\sigma)$

$$Z = \mu + \sigma \odot \epsilon$$
$$\epsilon \sim \mathcal{N}(0, 1)$$

**Reparameterization trick**  
(sampling + differentiability)



# Variational Autoencoder (VAE)

- Objective function

$$\text{maximize } \log p_{\theta}(\mathbf{x})$$

- Explicit density modeling:

$$p_{\theta}(\mathbf{x}) = \frac{\overset{\text{Prior}}{p_{\theta}(\mathbf{z})} \overset{\text{Likelihood}}{p_{\theta}(\mathbf{x}|\mathbf{z})}}{\underset{\text{Marginal likelihood}}{p_{\theta}(\mathbf{z}|\mathbf{x})}} \quad \text{Posterior} \rightarrow \text{Intractable!!}$$

- Variational Inference

- Assume that  $q(\mathbf{z}|\mathbf{x})$  follows a Gaussian distribution.
- Variational inference approximates  $p_{\theta}(\mathbf{z}|\mathbf{x})$  using  $q_{\phi}(\mathbf{z}|\mathbf{x})$ .  $\rightarrow \text{minimize } D(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))$
- $q^{*}(\mathbf{z}|\mathbf{x}) = \arg \min_{q \in Q} (D_{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})))$  where the Kullback-Leibler divergence  $D_{KL}(q||p) = \int p(x) \log \frac{p(x)}{q(x)} dx$ .

# Bayes' Theorem

“Bayes' theorem gives a mathematical rule for inverting conditional probabilities.”

- Bayes' rule

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- Bayesian inference

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \propto P(D|\theta)P(\theta)$$

- $P(\theta|D)$ : posterior
- $P(D|\theta)$ : likelihood
- $P(\theta)$ : prior

# Variational Autoencoder (VAE)

$$\text{maximize } \log p_{\theta}(\mathbf{x}) = \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{z}$$

$$= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{Z})p(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{Z})p(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \frac{q_{\theta}(\mathbf{z}|\mathbf{x})}{q_{\theta}(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z})) + D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \leftarrow \text{Intractable}$$

$$\geq \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z})) \leftarrow \text{Evidence lower bound (ELBO)}$$

Reconstruction error

Regularization

# Generative Adversarial Networks (GANs)



# Generative Adversarial Network (GAN)

“GAN produces data that is hard to distinguish from the real thing, but not real.”

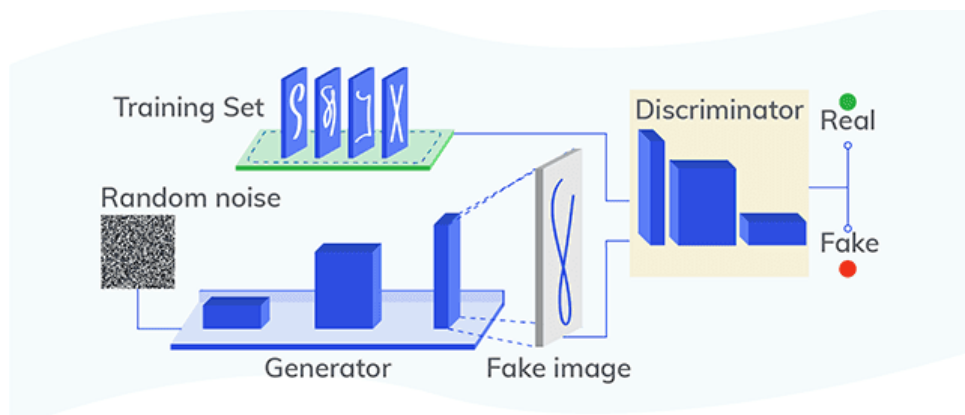


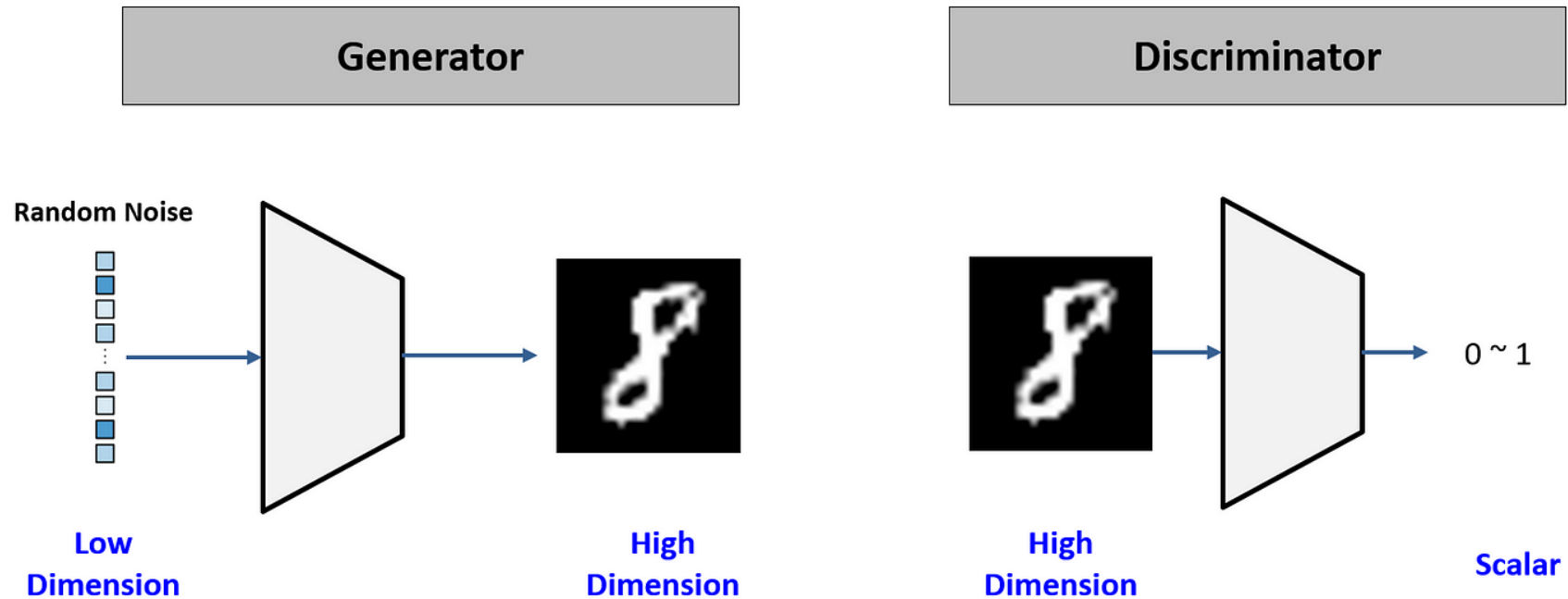
Image Source | xenonstack.com



The standard GAN loss function, also known as the min-max loss, was first described in a 2014 paper by Ian Goodfellow et al., titled “Generative Adversarial Networks”:

$$E_x[\log(D(x))] + E_z[\log(1 - D(G(z)))]$$

# Generative Adversarial Network (GAN)

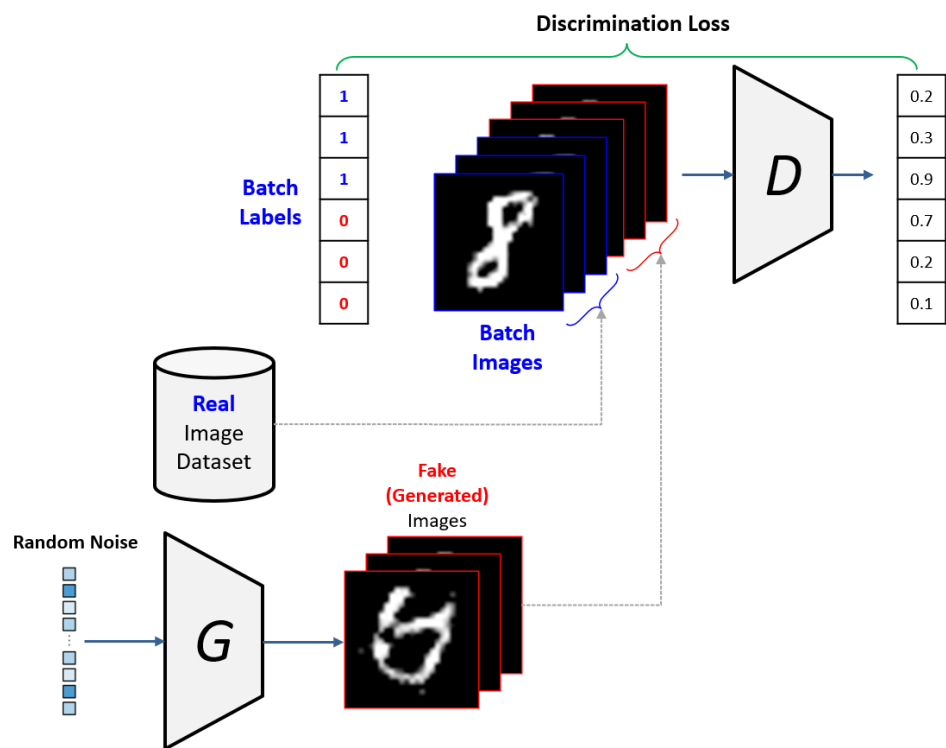


- It works like a 'decoder' in a VAE.
- GOAL: Generate "**REALISTIC**" data.

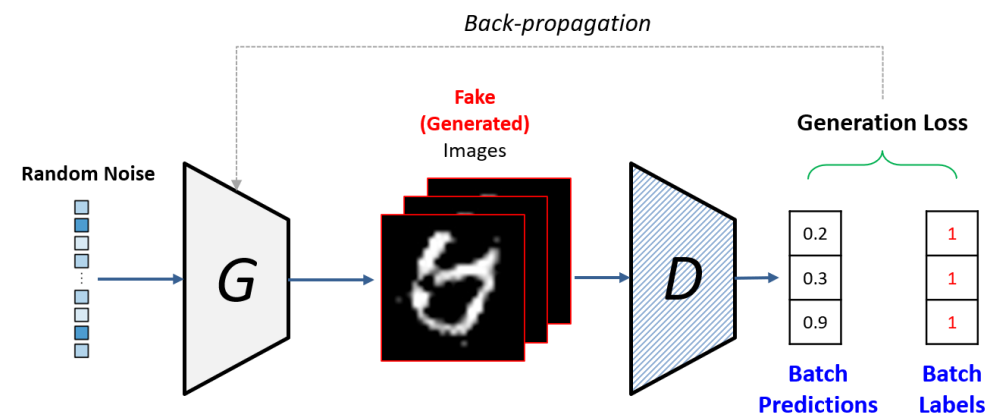
- It is a typical 'image classification'.
- GOAL: Discriminate between **real** and **fake** data.

# Generative Adversarial Network (GAN)

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



Training of discriminator

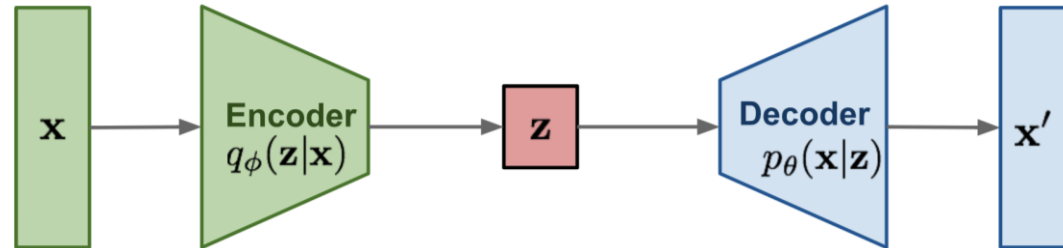


Training of Generator

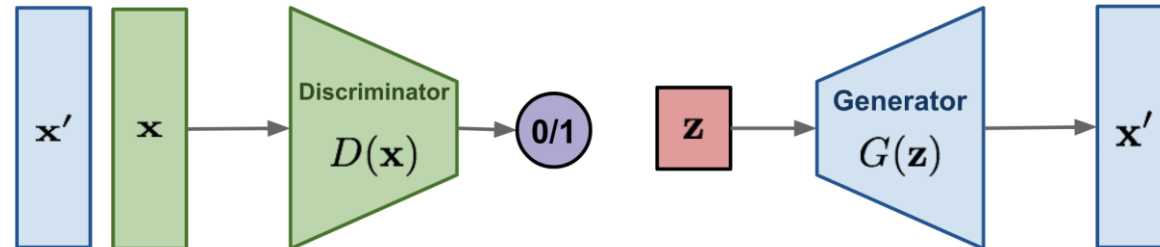
# Flow-based Models

# Comparison of Different Generative Model Types

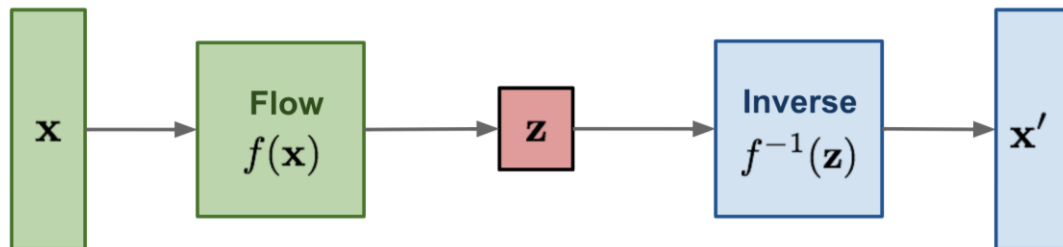
**VAE:** maximize ELBO.



**GAN:** minimax the classification error loss.

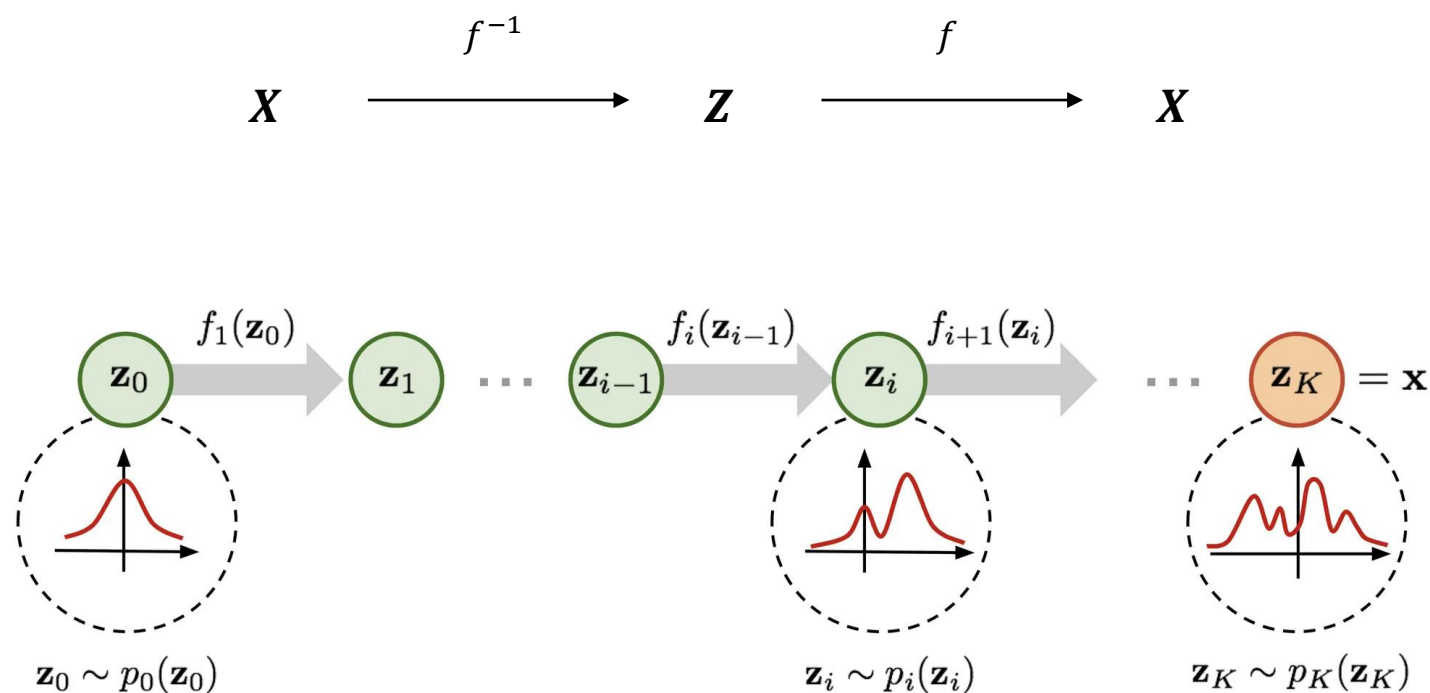


**Flow-based generative models:**  
minimize the negative log-likelihood

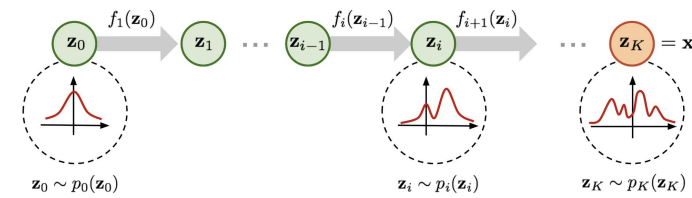


# Normalizing Flows

“A normalizing flow transforms a simple distribution into a complex one by applying a sequence of *invertible* transformation functions.”



# Normalizing Flows



“A normalizing flow transforms a simple distribution into a complex one by applying a sequence of *invertible* transformation functions.”

## ▪ Change of variable theorem

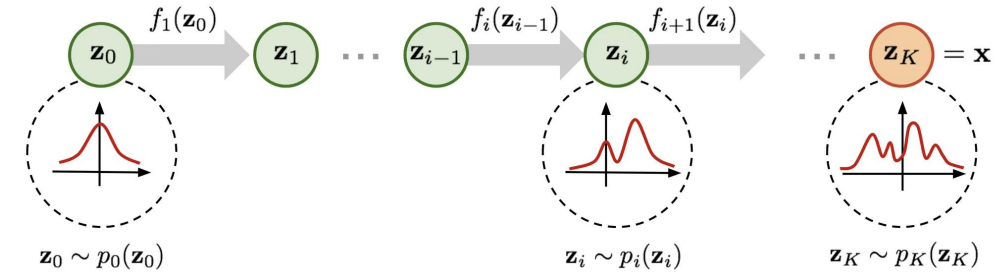
- Assume that  $z \sim \pi(z)$ ,  $x = f(z)$ ,  $z = f^{-1}(x)$  for random variables  $x$  and  $z$ .
- $\int p(x)dx = \int \pi(z)dz = 1$  (Definition of probability distribution)
- $p(x) = \pi(z) \left| \frac{dz}{dx} \right| = \pi(f^{-1}(x)) \left| \frac{df^{-1}}{dx} \right| \Rightarrow p(\mathbf{x}) = \pi(\mathbf{z}) \left| \det \frac{d\mathbf{z}}{d\mathbf{x}} \right| = \pi(f^{-1}(\mathbf{x})) \left| \frac{df^{-1}}{d\mathbf{x}} \right|$

## ▪ Inverse function theorem

- If  $y = f(x)$  and  $x = f^{-1}(y)$ ,

$$\frac{df^{-1}(y)}{dy} = \frac{dx}{dy} = \left( \frac{dy}{dx} \right)^{-1} = \left( \frac{df(x)}{dx} \right)^{-1}$$

# Normalizing Flows



By the definition of the normalizing flow model,

$$z_{i-1} \sim p_{i-1}(z_{i-1})$$

$$z_i = f_i(z_{i-1}), \text{ thus } z_{i-1} = f_i^{-1}(z_i)$$

$$p_i(z_i) = p_{i-1}(f_i^{-1}(z_i)) \left| \det \frac{df_i^{-1}}{dz_i} \right| ; \text{ according to the change of variable theorem.}$$

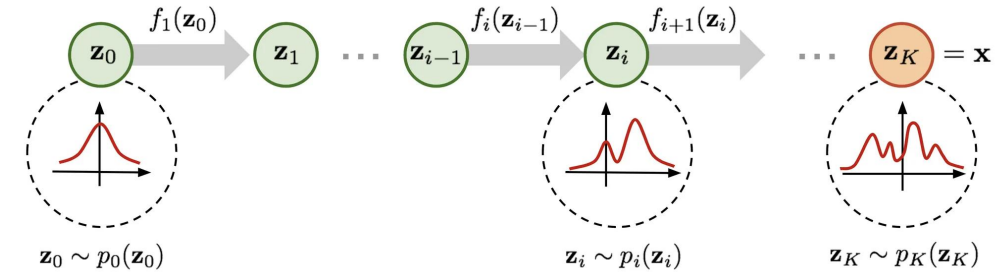
$$= p_{i-1}(z_{i-1}) \left| \det \left( \frac{df_i}{dz_{i-1}} \right)^{-1} \right| ; \text{ according to the invers function theorem.}$$

$$= p_{i-1}(z_{i-1}) \left| \det \frac{df_i}{dz_{i-1}} \right|^{-1} ; \text{ according to a property of Jacobians of invertible function.}$$

$$\Rightarrow \log p_i(z_i) = \log p_{i-1}(z_{i-1}) - \log \left| \det \frac{df_i}{dz_{i-1}} \right|.$$



# Normalizing Flows



$$x = z_K = f_K \circ f_{K-1} \circ \dots \circ f_1(z_0)$$

$$\log p(x) = \log \pi_K(z_K)$$

$$= \log \pi_{K-1}(z_{K-1}) - \log \left| \det \frac{df_K}{dz_{K-1}} \right|$$

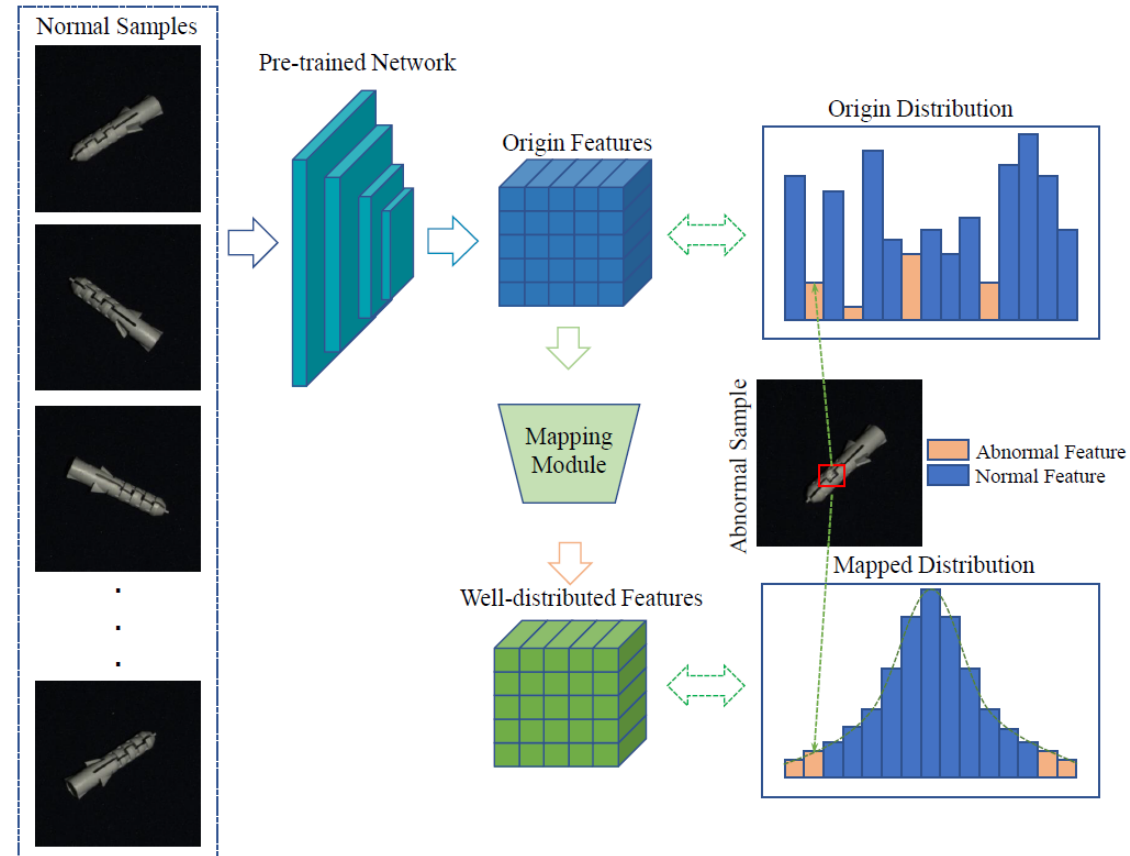
$$= \log \pi_{K-2}(z_{K-2}) - \log \left| \det \frac{df_{K-1}}{dz_{K-2}} \right| - \log \left| \det \frac{df_K}{dz_{K-1}} \right|$$

$$= \dots$$

$$= \log \pi_0(z_0) - \sum_{i=1}^K \log \left| \det \frac{df_i}{dz_{i-1}} \right|$$

# Anomaly Detection Using Normalizing Flows

“Normalizing flows can learn the distribution of normal data to detect anomalies lie out of the distribution.”



# Diffusion Models

# Diffusion

“A diffusion model gradually converts a well-known distribution like Gaussian noise into a complex data distribution by iteratively removing a bit of noise.”

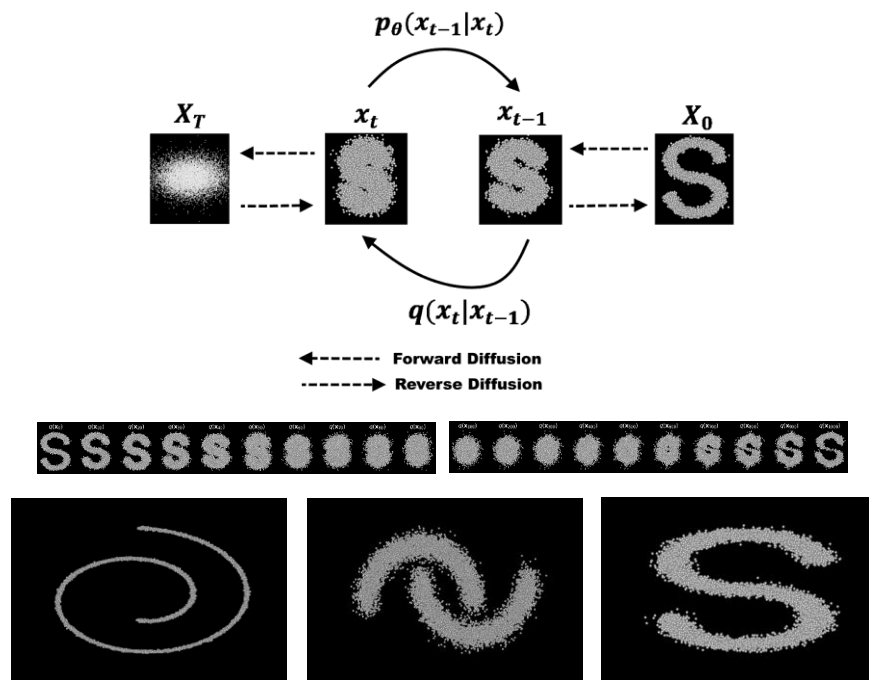


Image Source | Towards Data Science ([link](#))

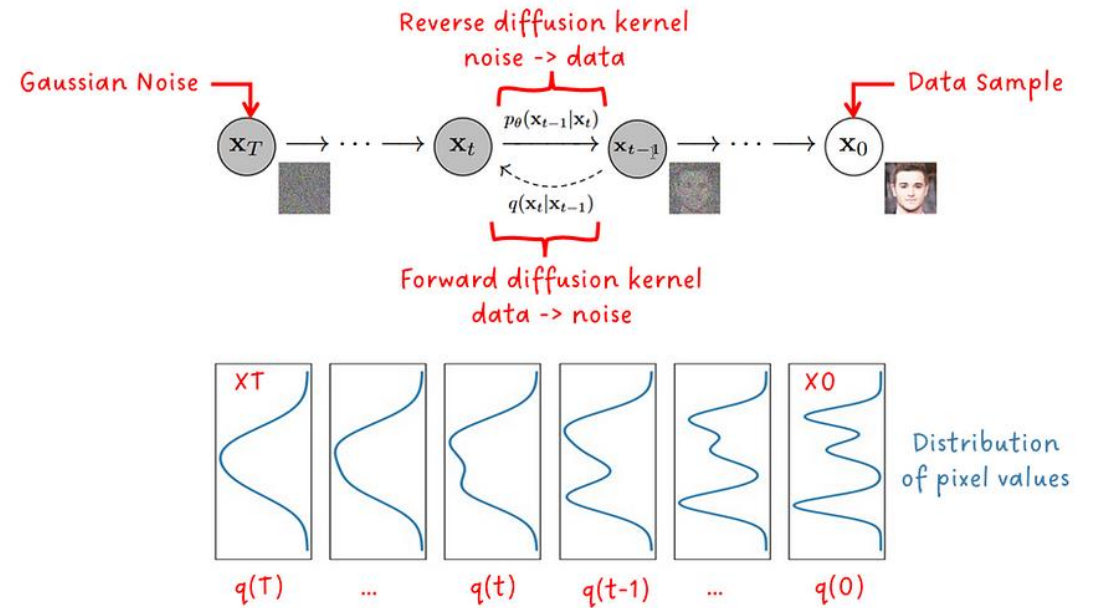


Image Source | Towards Data Science ([link](#))

# DDPM: Denoising Diffusion Probabilistic Model

“DDPM is diffusion probabilistic model inspired by considerations from nonequilibrium thermodynamics.”

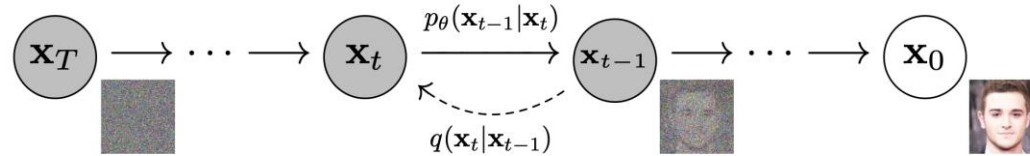


Figure 2: The directed graphical model considered in this work.

## ▪ DDPM at a glance

- [Forward diffusion process] We know  $q(x_t|x_{t-1})$ , but not  $q(x_{t-1}|x_t)$ .  $q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$ ,  $q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$
- [Reverse diffusion process] We approximate  $q(x_{t-1}|x_t)$  by modelling  $p_\theta(x_{t-1}|x_t)$ .
- Finally, we can model  $p_\theta(x_T|x_0)$  through  $p_\theta(x_{t-1}|x_t)$  by the Markov chain based on Gaussian transitions.

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)), \quad p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$

# DDPM: Denoising Diffusion Probabilistic Model

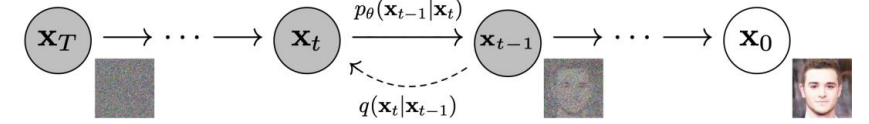


Figure 2: The directed graphical model considered in this work.

## ▪ Objective function

Minimize the negative log-likelihood:

$$\mathbb{E} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] = \mathbb{E}_q \left[ -\log p(\mathbf{x}_T) - \sum_{t \geq 1} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] =: L$$

## ▪ Loss function

$$\mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{L_T} + \sum_{t \geq 1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right]$$

- $L_T$ : the difference in distribution between noise of  $\mathbf{x}_T$  generated by  $p$  and by  $q$  given  $\mathbf{x}_0$ .
- $L_{t-1}$ : the difference between the reverse ( $p$ ) and the forward ( $q$ ) processes.
- $L_0$ : the likelihood of  $\mathbf{x}_0$  given  $\mathbf{x}_1$ .

# Takeaways

# A Probabilistic View of Deep Learning: Discriminative vs. Generative Models

“Discriminative and generative models are two fundamental approaches in deep learning, each with distinct characteristics and applications.”

## ▪ Discriminative models

- Learn the *decision boundary* (e.g., classification and regression)
- Focus on the inference of the output,  $Y$ , given the input,  $X$ .

⇒ Model the *conditional probability* of the output given the input,  $P(Y|X)$ .

## ▪ Generative models

- Learn the underlying distribution of data.
- Focus on the generation of new data instances.

⇒ Model the *joint probability*  $P(X, Y)$ . Depending on the problem,  $P(X|Y)$  or  $P(X)$  can be of interest.



# Examples of Deep Generative Models

- Variational Autoencoders (VAEs)
- Generative Adversarial Networks (GANs)
- Autoregressive models (e.g., Transformers like GPT)
- Flow-based models
- Diffusion models
- ...

Thank you! 😊