

Introduction to Machine Learning CS182

Lu Sun

School of Information Science and Technology
ShanghaiTech University

December 28, 2023

Today:

- Deep Generative Networks (DGN)

Readings:

- Deep Learning (DL), Chapters 14&20

Today's Agenda

深度生成网络：元监督

Deep Generative Models (DGM)

- Overview
- Representation Learning with Autoencoder 自编码器
- Generative Adversarial Network (GAN) 生成对抗网络
- Applications of GANs

Acknowledgement: Hugo Larochelle's, Mehryar Mohri@NYU's, Yingyu Liang@Princeton's, Bhiksha Raj@CMU's & Feifei Li@Stanford's course notes

Deep Generative Networks-DGM

- **Overview**
- Representation Learning with Autoencoder
- Generative Adversarial Network (GAN)
- Applications of GANs

Unsupervised Learning

■ Task formulation

Unsupervised Learning

Data: x

Just data, no labels! 

Goal: Learn some underlying
hidden *structure* of the data

Examples: Clustering,
dimensionality reduction, feature
learning, density estimation, etc.



Unsupervised Learning

■ Task formulation

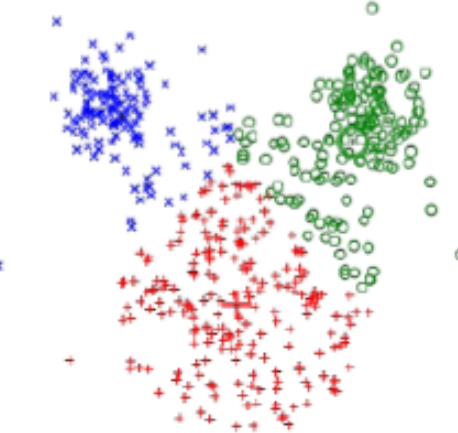
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



K-means clustering



Unsupervised Learning

■ Task formulation

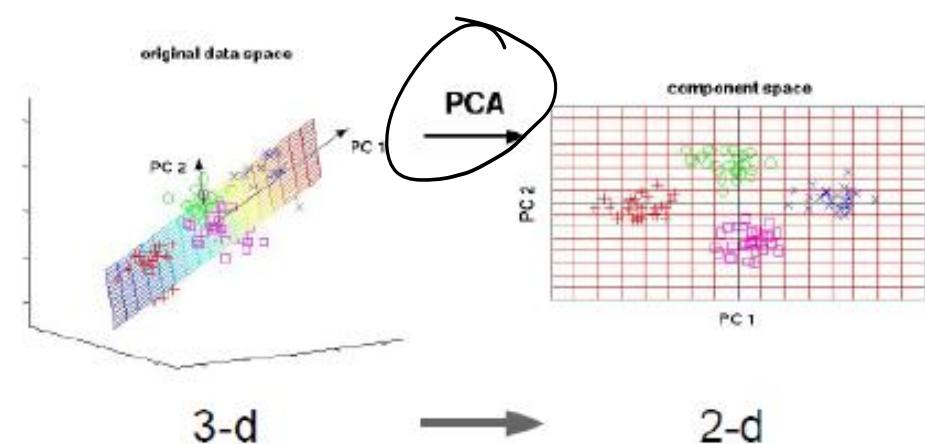
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Principal Component Analysis
(Dimensionality reduction)

Unsupervised Learning

■ Task formulation

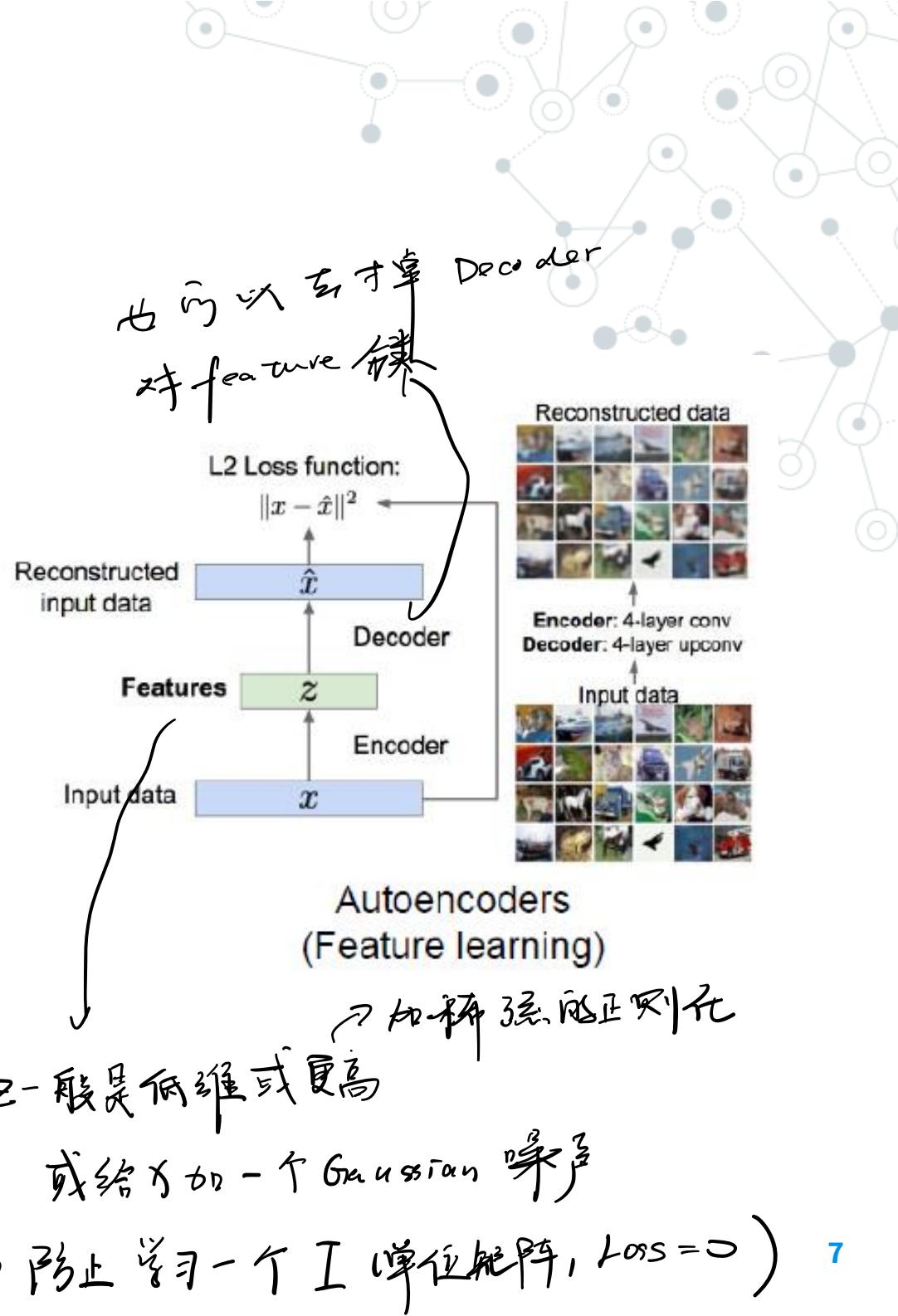
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Unsupervised Learning

■ Task formulation

Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.

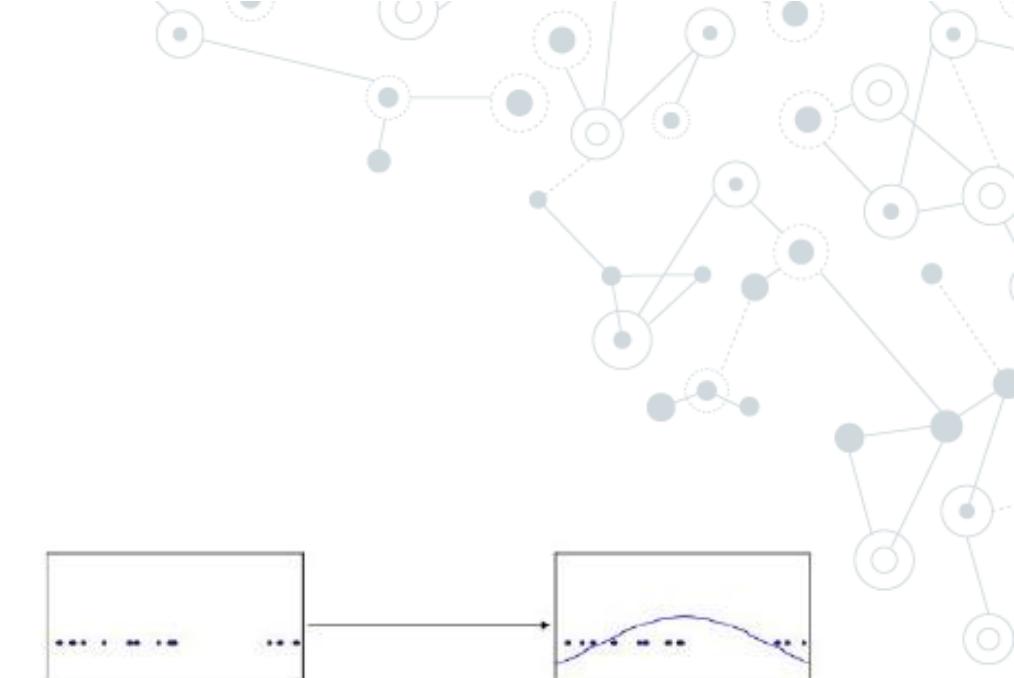
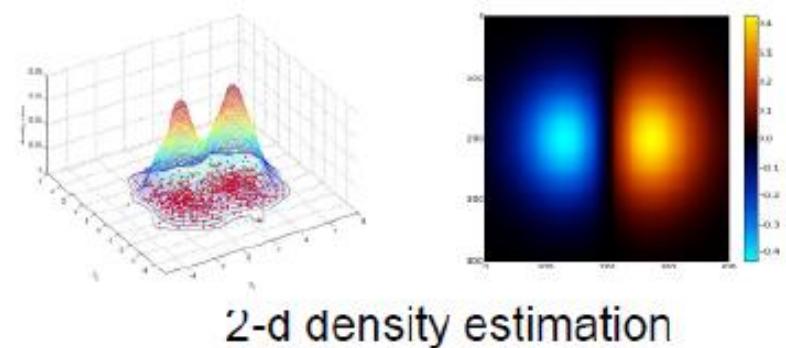


Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



2-d density estimation



Deep Generative Networks-DGM

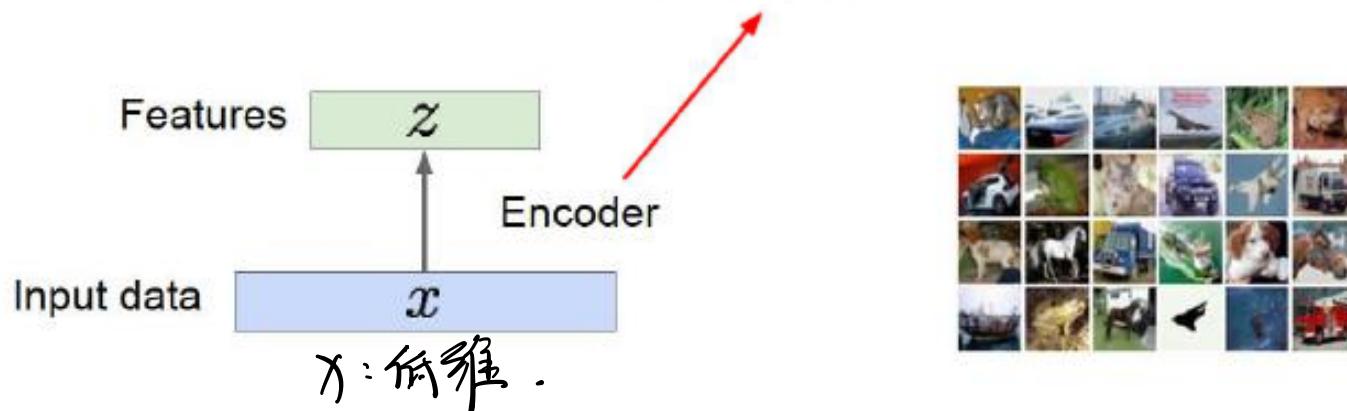
- Overview
- **Representation Learning with Autoencoder**
- **Generative Adversarial Network (GAN)**
- **Applications of GANs**

Autoencoder

■ Feature representation learning

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

Originally: Linear +
nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN

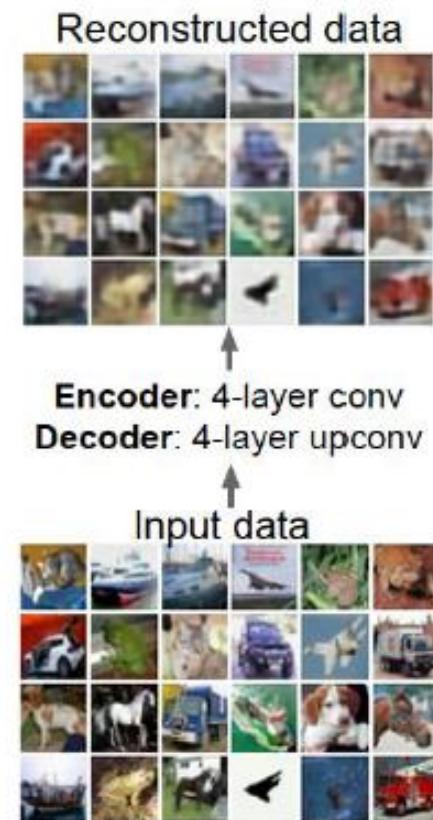
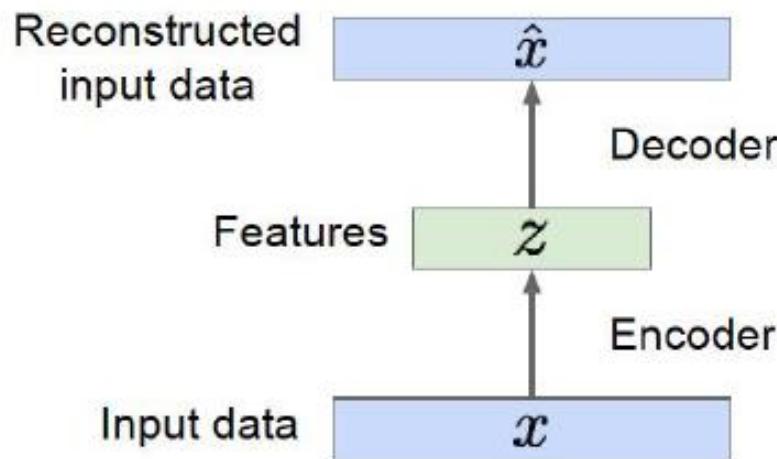


Autoencoder

■ Feature representation learning

How to learn this feature representation?

Train such that features can be used to reconstruct original data
“Autoencoding” - encoding itself

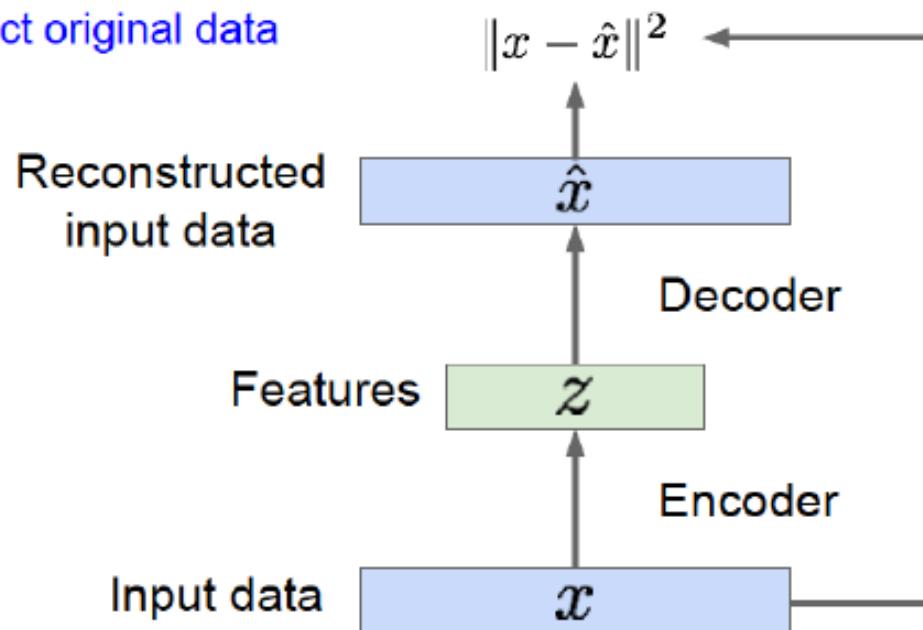


Autoencoder

■ Feature representation learning

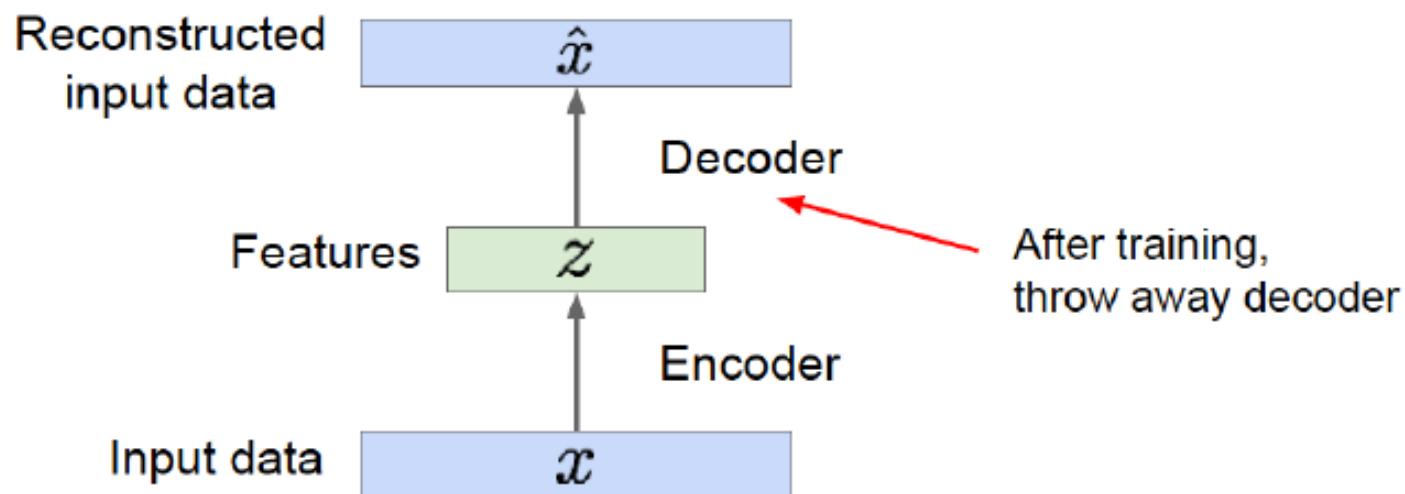
Train such that features
can be used to
reconstruct original data

L2 Loss function:



Autoencoder

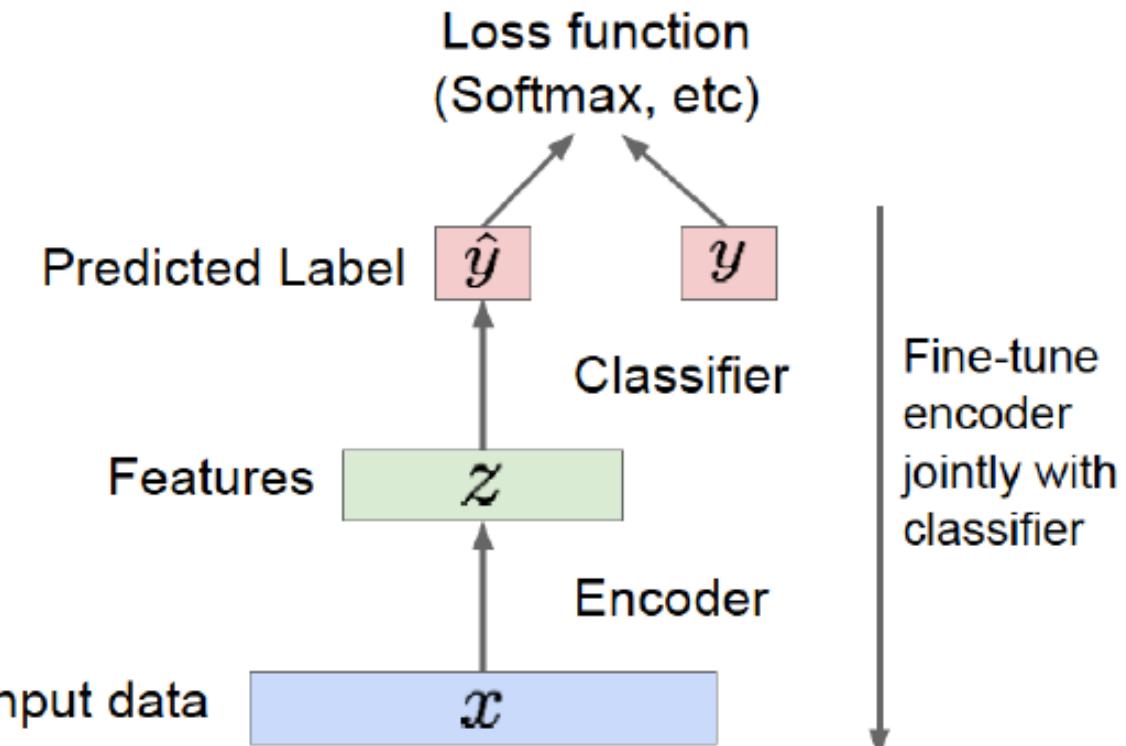
- Feature representation learning



Autoencoder

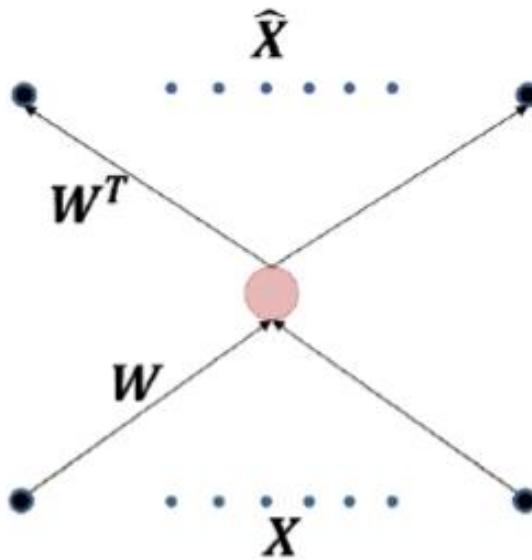
- Feature representation learning

Encoder can be used to initialize a **supervised** model



Autoencoder

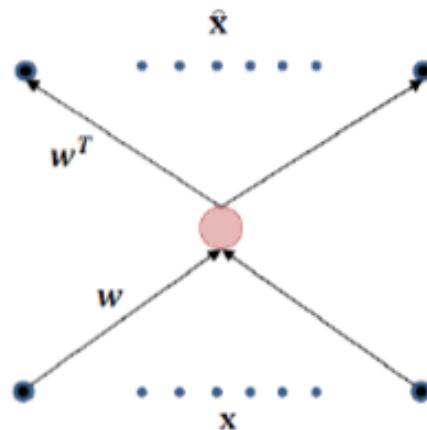
■ Linear hidden layer example



- A single hidden unit
- Hidden unit has *linear* activation
- What will this learn?

Autoencoder

■ Linear hidden layer example



Training: Learning W by minimizing
L2 divergence

$$\hat{x} = \underline{w^T w x} \quad \text{(保证矩阵乘法)} \\ div(\hat{x}, x) = \|x - \hat{x}\|^2 = \|x - w^T w x\|^2$$

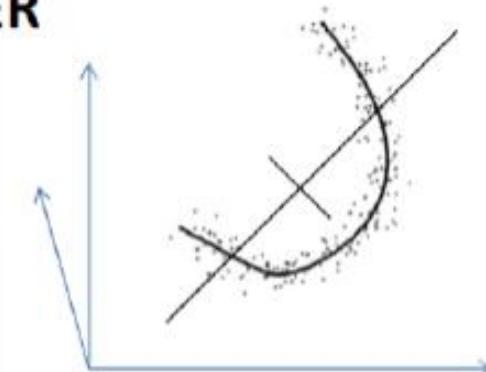
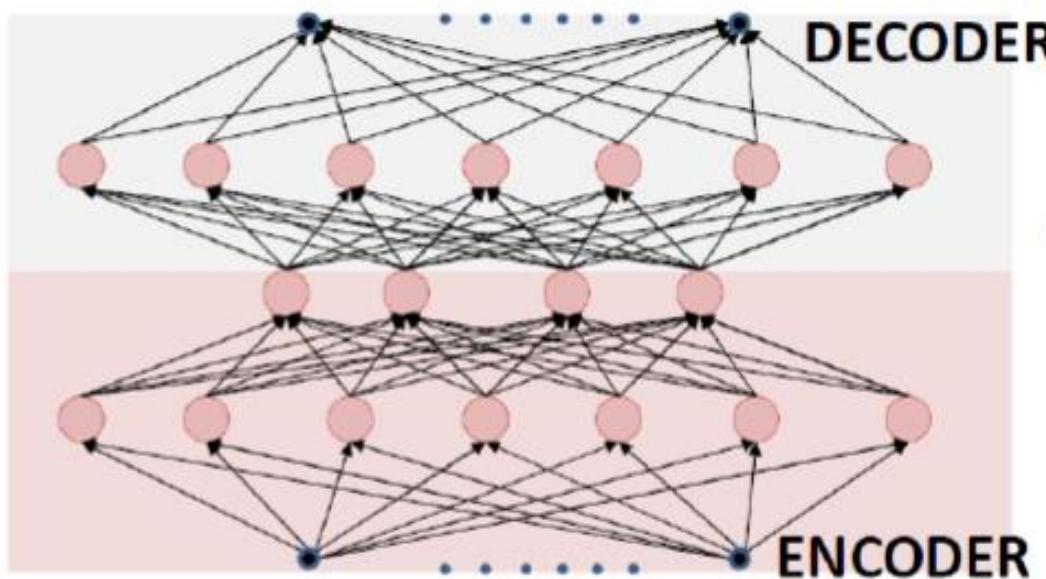
$$\hat{W} = \operatorname{argmin}_W E[div(\hat{x}, x)]$$

$$\hat{W} = \operatorname{argmin}_W E[\|x - w^T w x\|^2]$$

- This is just PCA! \Rightarrow 单隐层且只有一个隐层单元时
就是一个线性PCA。

Autoencoder

■ Nonlinear hidden layer



- With non-linearity
 - “Non linear” PCA
 - Deeper networks can capture more complicated manifolds

多层会导致计算慢并可能陷入局部最优
(初始化困难)

这两层
二> v h \rightarrow 使用 PBM 来初始化
受限玻尔兹曼机

$$P_w(v, h) = \frac{\exp(-v^T W h)}{Z(w)}$$

离散自编码器

$$x = w_0 + \epsilon, \epsilon \sim N(0, \sigma^2) \quad \text{且} \quad P(x|z) = N(w_0 + \epsilon)$$

$$\text{拓展: } P(x|z) = N(\mu_\phi(z), \sigma_\phi^2(z))$$

变成函数

给定 x 去估计 θ . 但目前未知 z

类似 EM 算法: $\begin{cases} \text{估计 } \theta(z) \\ \text{重算 } \log P_\theta(x) \end{cases}$

$$\max_{\theta} \log P_\theta(x) = \log \int P_\theta(x, z) dz = \log \int q(z) \frac{P_\theta(x, z)}{q(z)} dz$$

$$\begin{aligned} &= \log \mathbb{E}_{q(z)} \left[\frac{P_\theta(x, z)}{q(z)} \right] \geq \mathbb{E}_{q(z)} \left[\log \frac{P_\theta(x, z)}{q(z)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{P_\theta(x, z)}{q_\phi(z|x)} \right] \end{aligned}$$

琴生不等式: \log 是 convex

VAE: 去做一个概率的自编码器. 使用 $q_\phi(z|x)$ 拟合 $q(z)$
(后验) $q_\phi(z|x) \sim N(\mu_\phi(x), \sigma_\phi^2(x))$

$$\Rightarrow \theta = \arg \max_{\theta} \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{P_\theta(x, z)}{q_\phi(z|x)} \right] = -KL(q_\phi(z|x) || P_\theta(x, z))$$

$$= \int q_\phi(z|x) \log \frac{P_\theta(x, z)}{q_\phi(z|x)} dz$$

一般是给一个较简单的先验

$P_\theta(x|z)$ $P_\theta(z)$ \downarrow 如: $z \sim N(0, I)$

$$= \int q_\phi(z|x) \left(\log \frac{q(z)}{q_\phi(z|x)} + \log P_\theta(x|z) \right) dz$$

是隐空间的初始采样

$$= \int q_\phi(z|x) \log \frac{P_\theta(z)}{q_\phi(z|x)} dz + \int q_\phi(z|x) \log \frac{P_\theta(x|z)}{q_\phi(z|x)} dz$$

$\sim N(\mu_\theta(z), \sigma_\theta^2(z))$

$$= -KL(q_\phi(z|x) || P_\theta(z)) + \mathbb{E}_{q_\phi(z|x)} [\log P_\theta(x|z)]$$

$$ELBO = \frac{1}{K} \sum_{i=1}^K \log \frac{P_\theta(x, z_i)}{q_\phi(z_i|x)} = \frac{1}{K} \sum_{i=1}^K \log \frac{P_\theta(x, \mu_\phi(x) + \sigma_\phi(x) \epsilon_i)}{q_\phi(\mu_\phi(x) + \sigma_\phi(x) \epsilon_i | x)}$$

$$z_i \sim q_\phi(z|x), z_i = \mu_\phi(x) + \sigma_\phi(x) \epsilon_i, \epsilon_i \sim N(0, I)$$

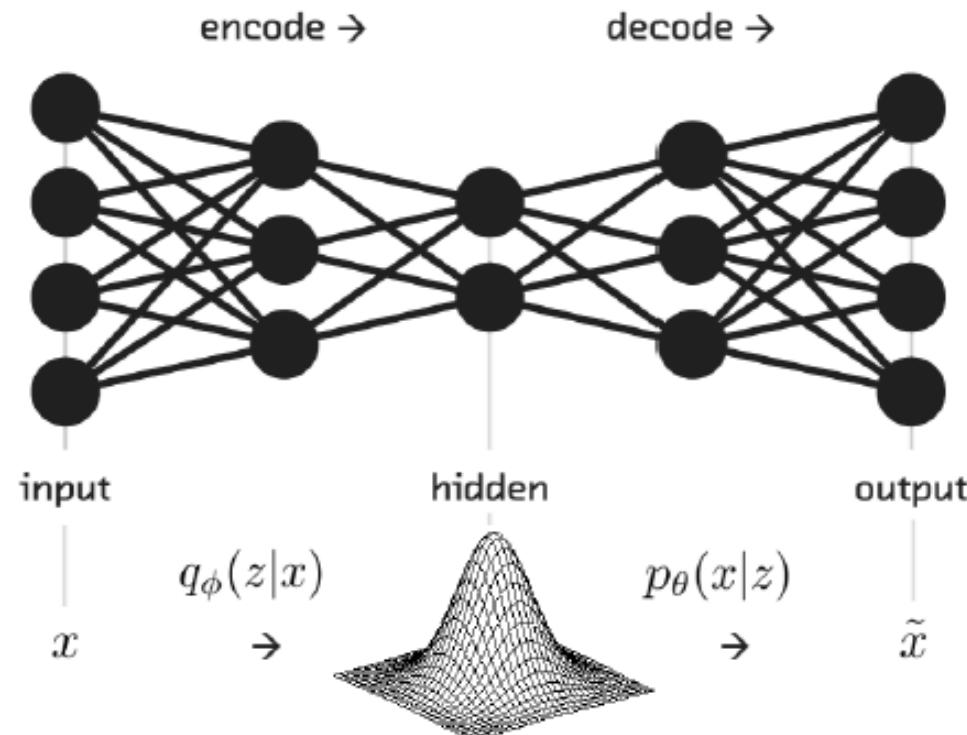
通过梯度下降计算: $\theta \leftarrow \theta - \gamma \frac{\partial L}{\partial \theta}, \phi \leftarrow \phi - \gamma \frac{\partial L}{\partial \phi}$

Variational Autoencoder (VAE)

- Objective $\mathcal{L}(x, \phi, \theta) = -D_{KL}(q_\phi(z|x)||p_\theta(z)) + E_{q_\phi(z|x)}[\log p_\theta(x|z)]$

Regularization term

Reconstruction term

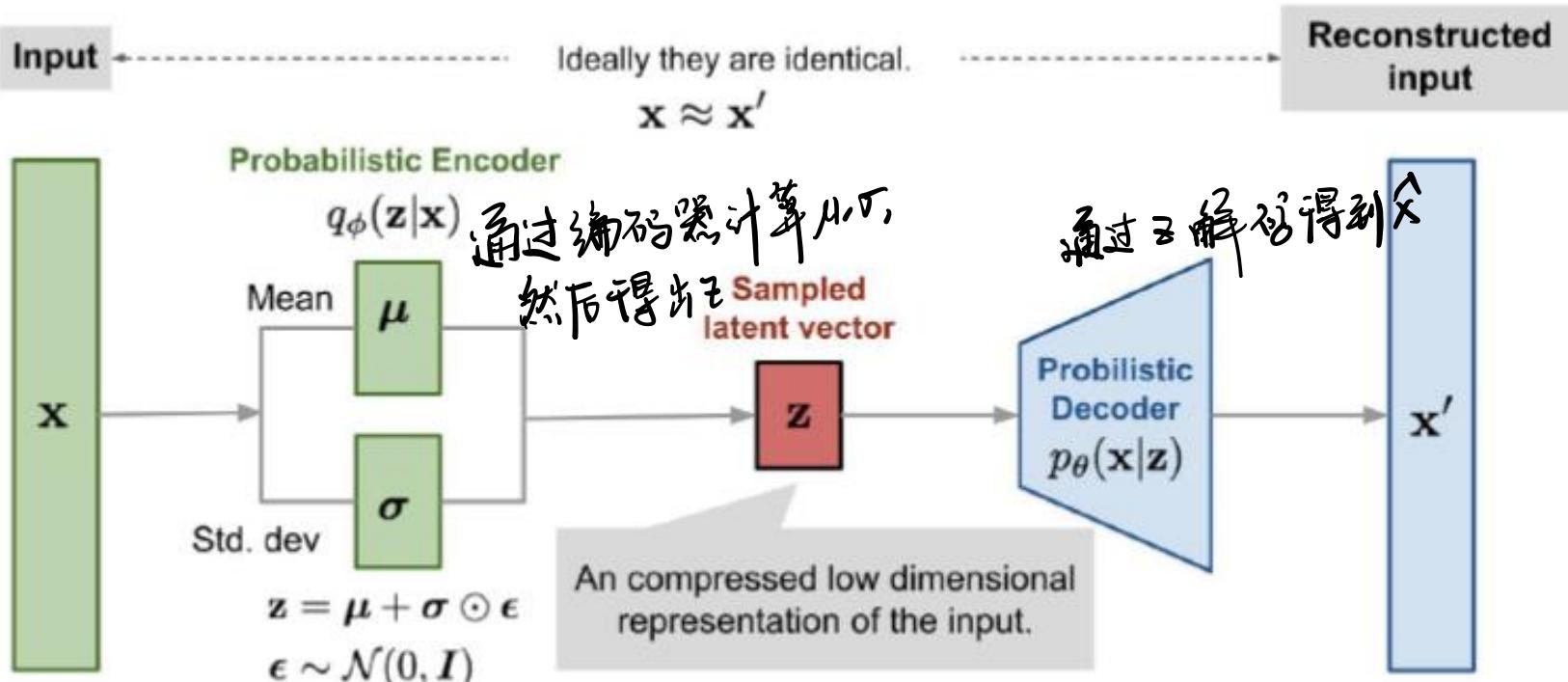


Variational Autoencoder (VAE)

■ Objective $\mathcal{L}(x, \phi, \theta) = -D_{KL}(q_\phi(z|x)||p_\theta(z)) + E_{q_\phi(z|x)}[\log p_\theta(x|z)]$

Regularization term

Reconstruction term



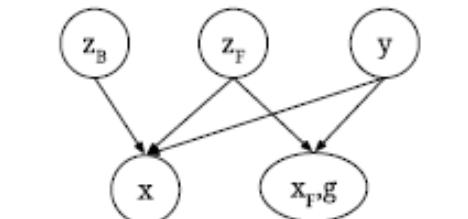
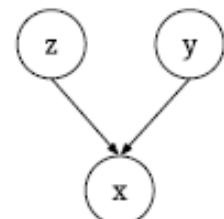
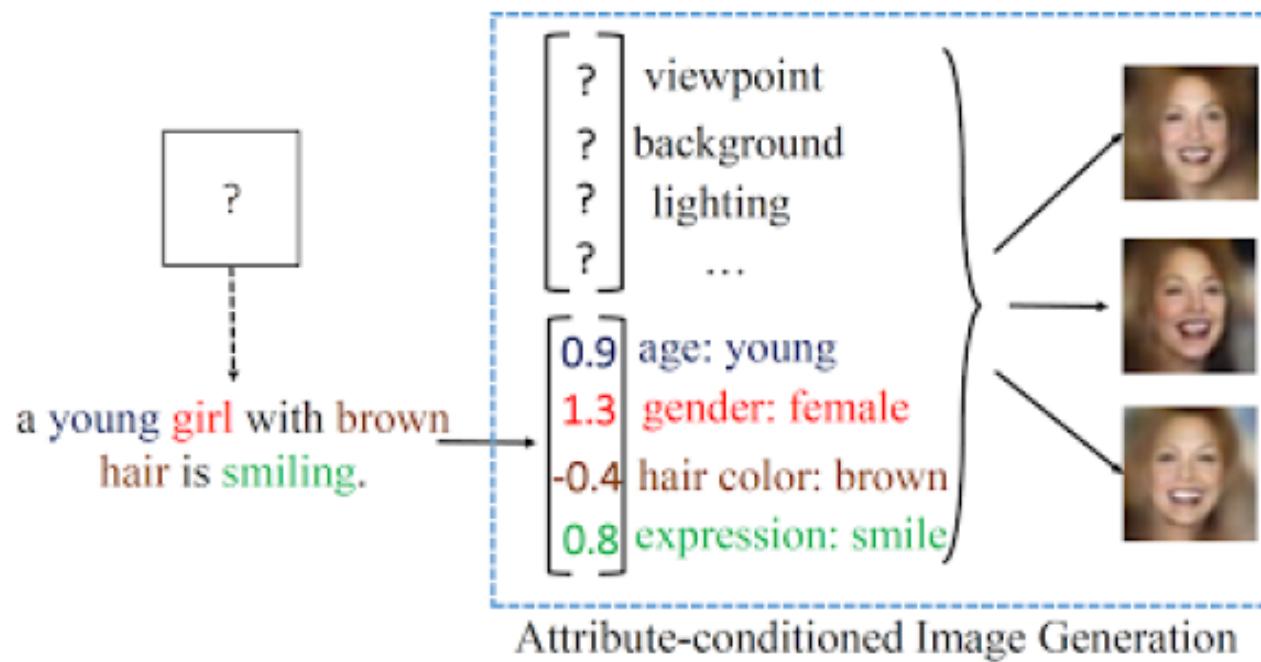
Interpreting the Latent Space



<https://arxiv.org/pdf/1610.00291.pdf>

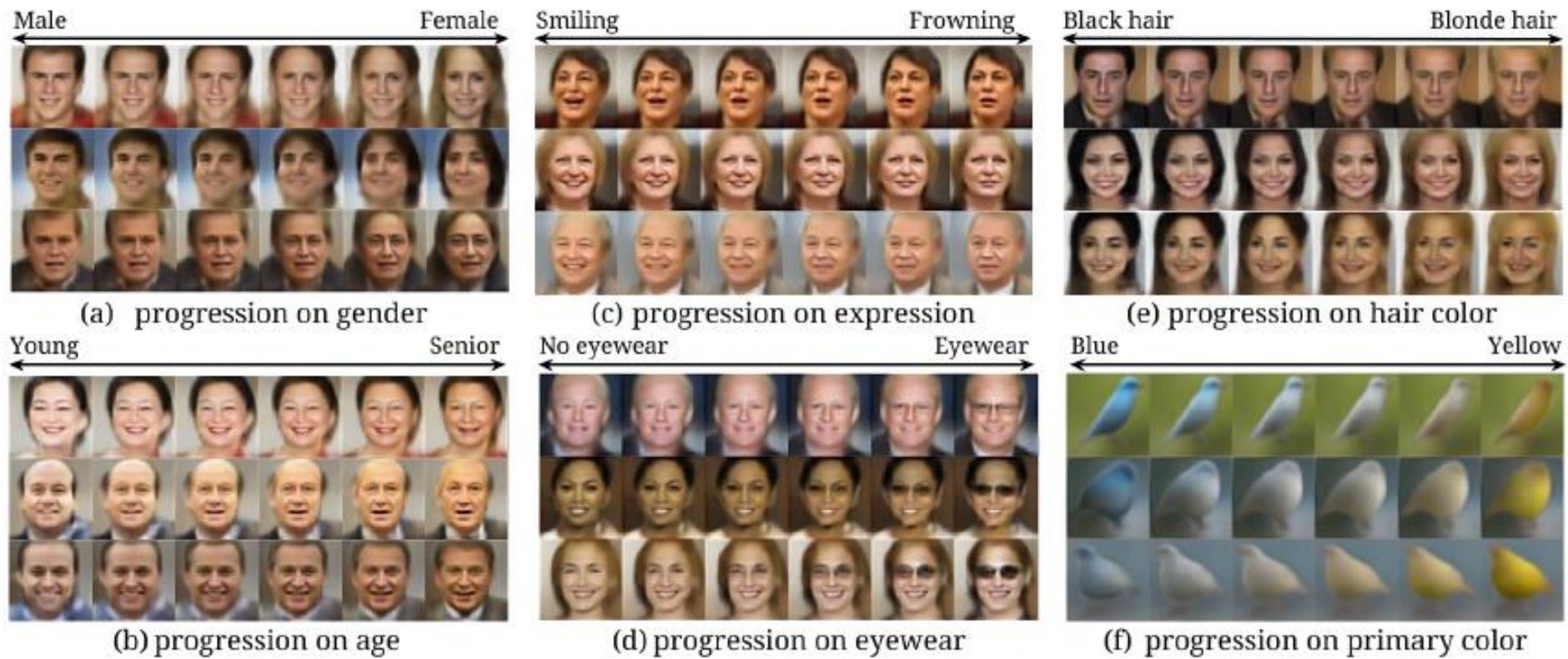
Example: Attribute2Image

■ Main ideas



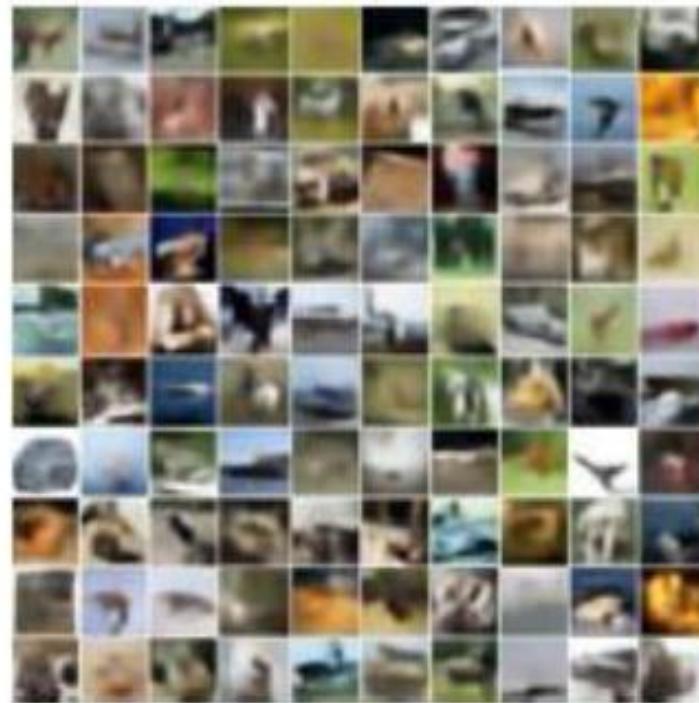
Example: Attribute2Image

■ Results



Problem of VAE

- Blurry images

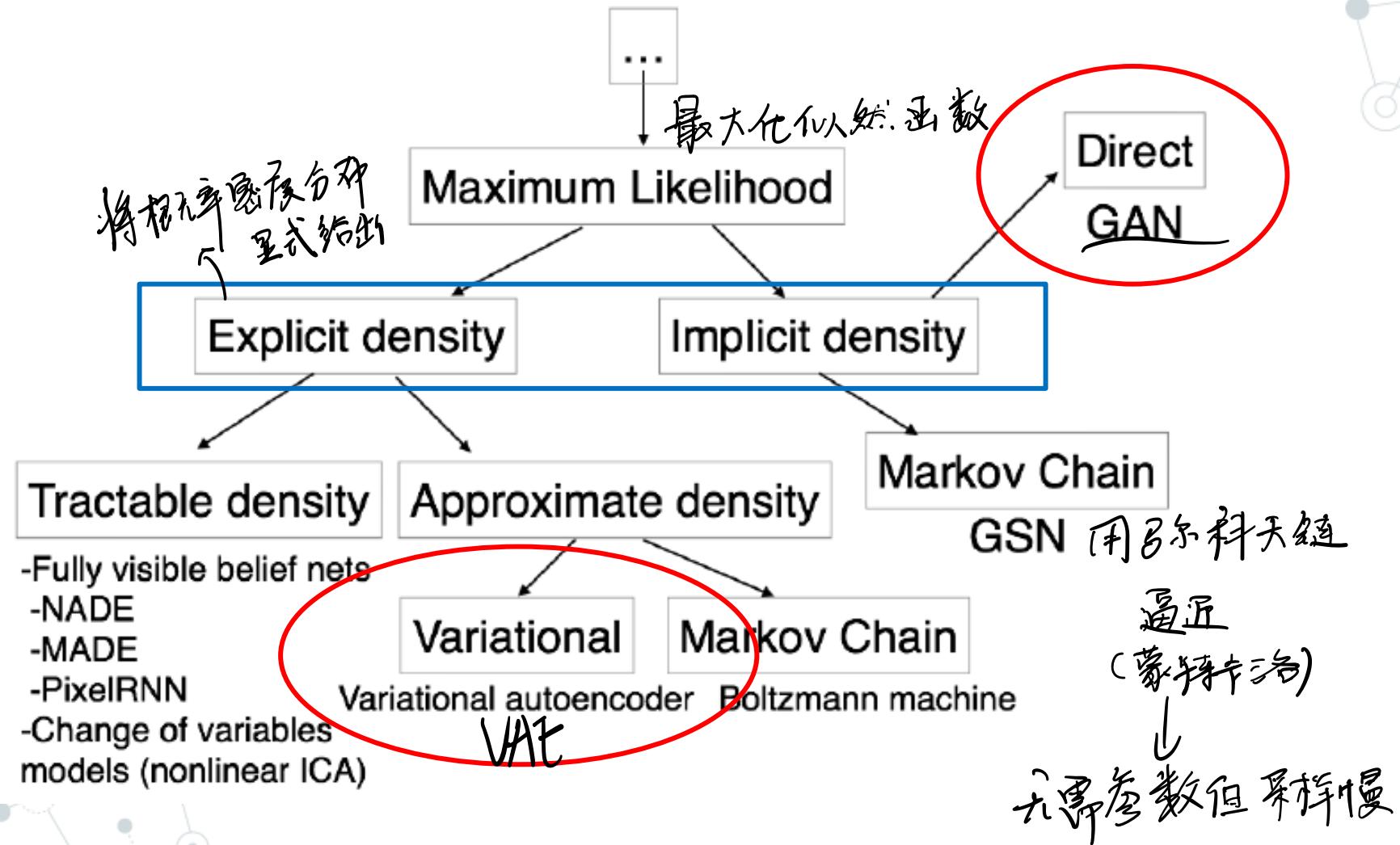


<https://blog.openai.com/generative-models/>

Deep Generative Networks-DGM

- Overview
- Representation Learning with Autoencoder
- **Generative Adversarial Network (GAN)**
- Applications of GANs
生成对抗网络

Taxonomy of Generative Models



Implicit Generative Models

更加复杂

- Working with explicit model $p(x)$ could be expensive
 - Variational Autoencoder (variational inference)
 - Boltzmann Machines (MCMC)
- Representation learning may not require $p(x)$
 - Sometimes we are more interested in taking samples from $p(x)$ instead of p itself

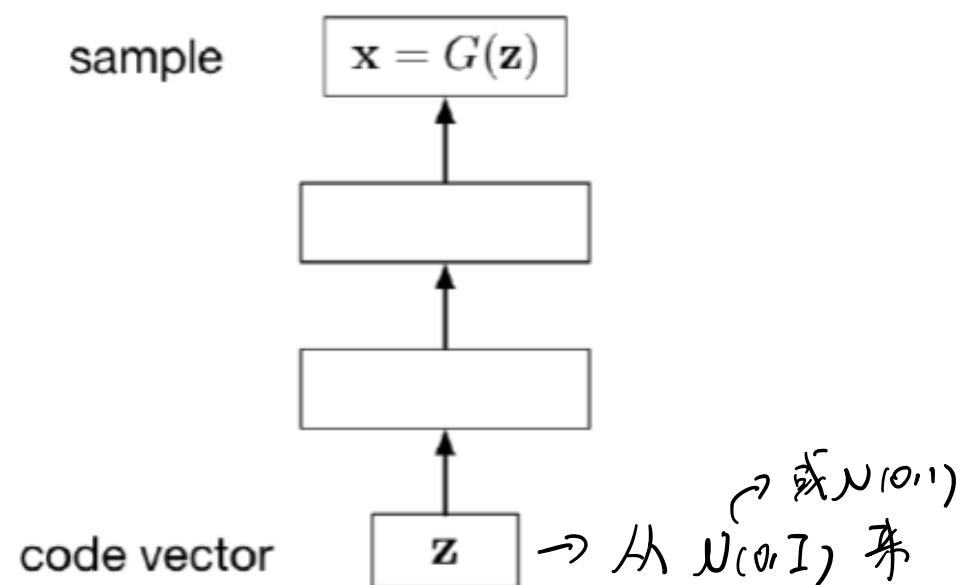
不需要知道样本点的概率分布

只用知道样本本身

Implicit Generative Models

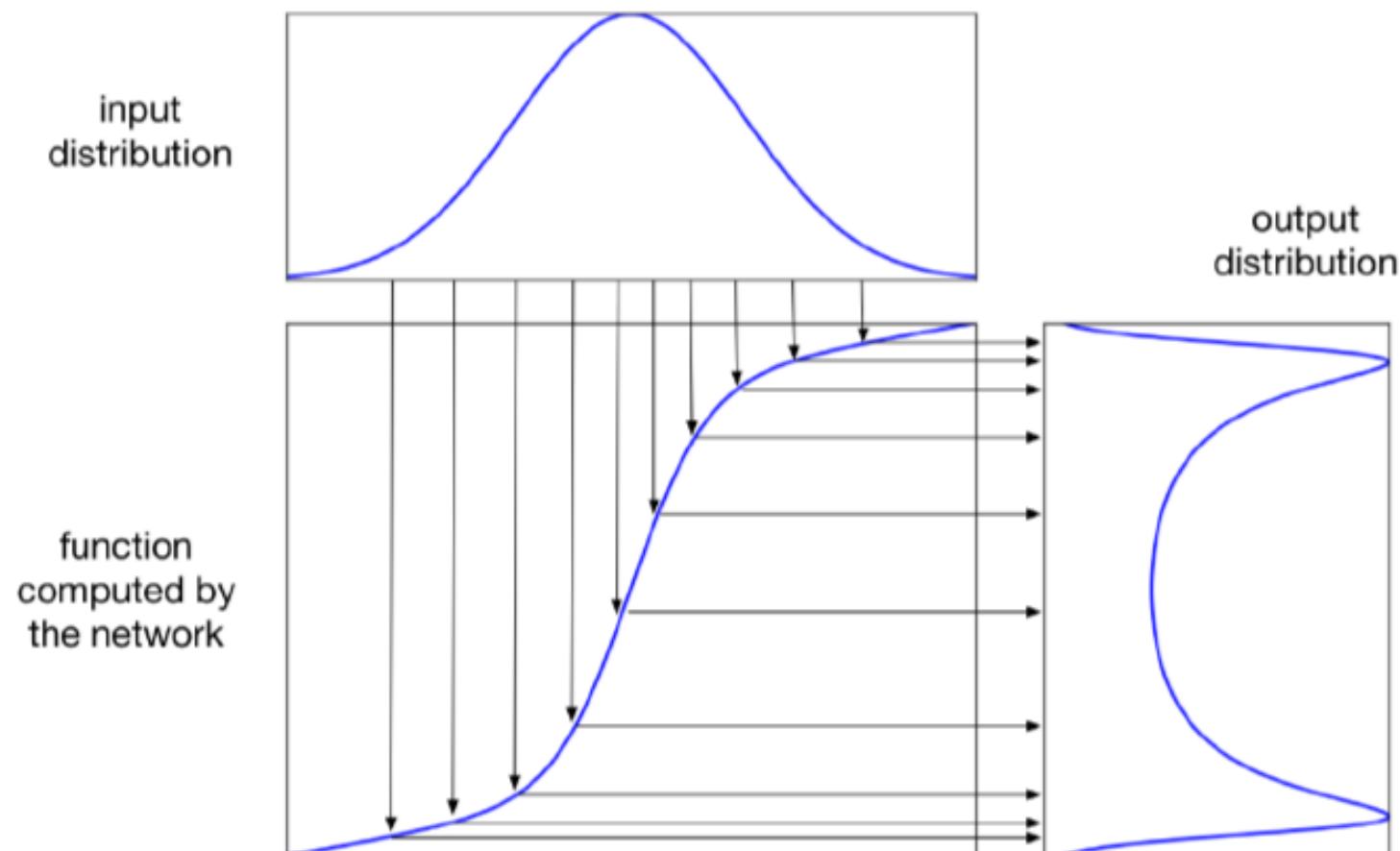
- Implicitly define a probability distribution
- Start by sampling the code vector z from a fixed, simple distribution
- A generator network computes a differentiable function G mapping z to an x in data space

生成器
Generator



Implicit Generative Models

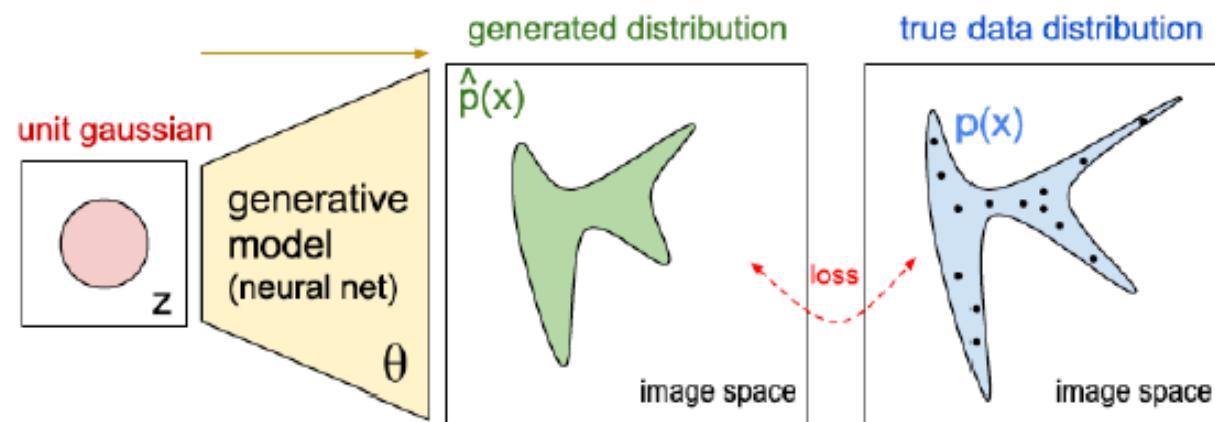
- Intuition: 1D example



所以用 NN 模拟一个分布
足够复杂则可以生成图像。

Implicit Generative Models

■ Intuition

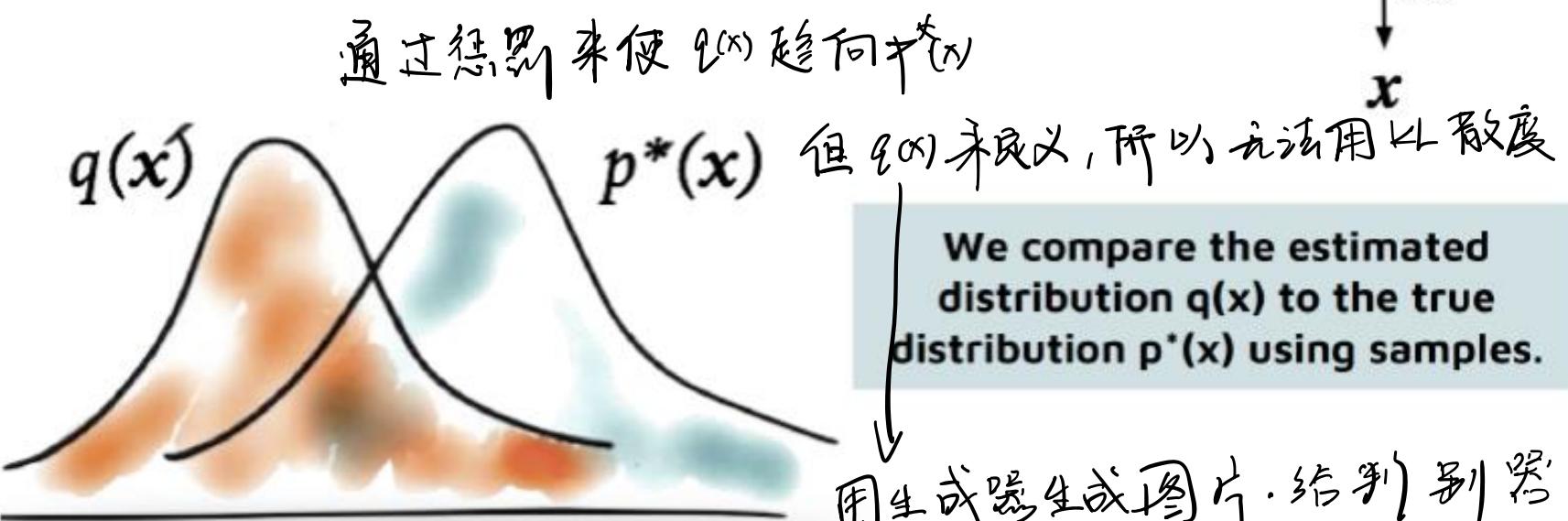


advocate/penalize samples within the blue/white region.

Learning by Comparison

■ Basic idea

For some models, we only have access to an unnormalised probability, partial knowledge of the distribution, or a simulator of data.



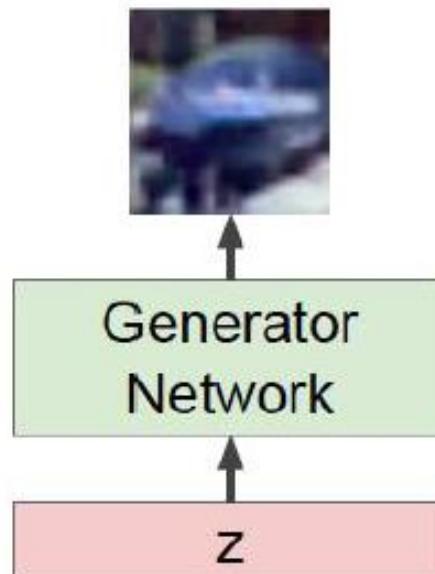
用生成器生成图片·给判别器
判别器也在同时训练
通过判别器的反馈惩罚 $q(x)$

Generative Adversarial Networks (GAN)

- Using a neural network to generate data

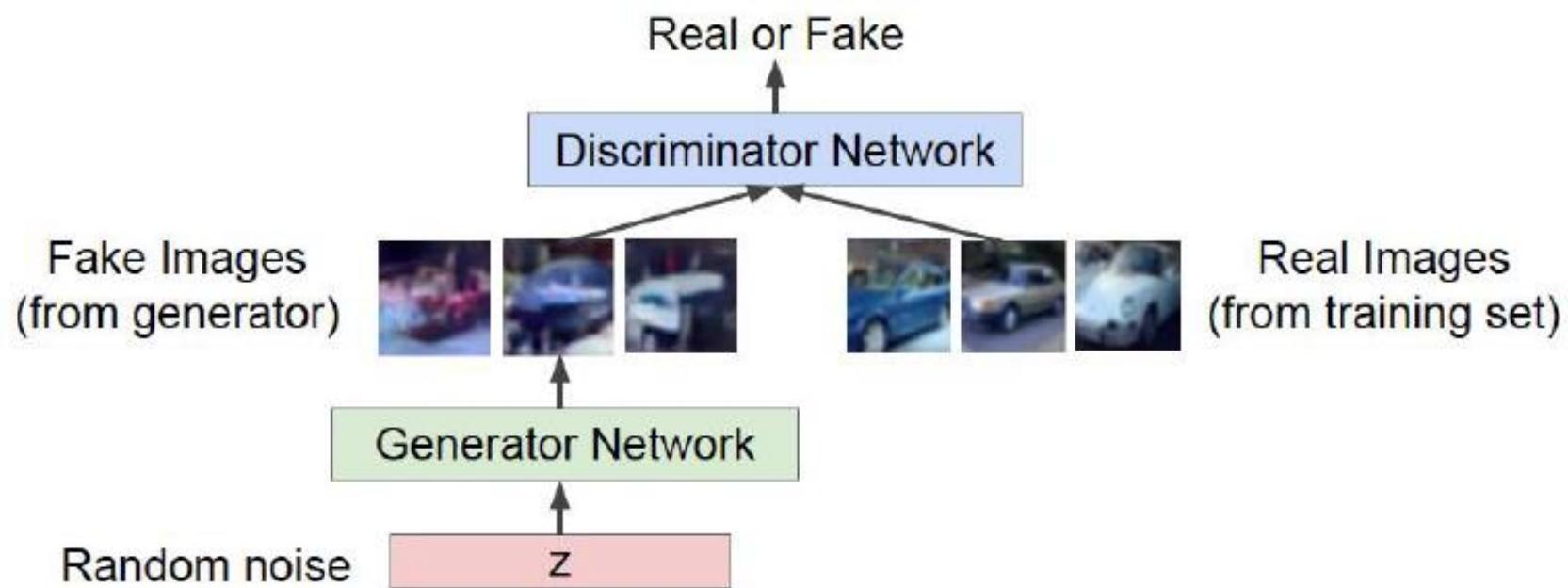
Output: Sample from
training distribution

Input: Random noise



Generative Adversarial Networks (GAN)

- Using another neural network to determine if the data is real or not

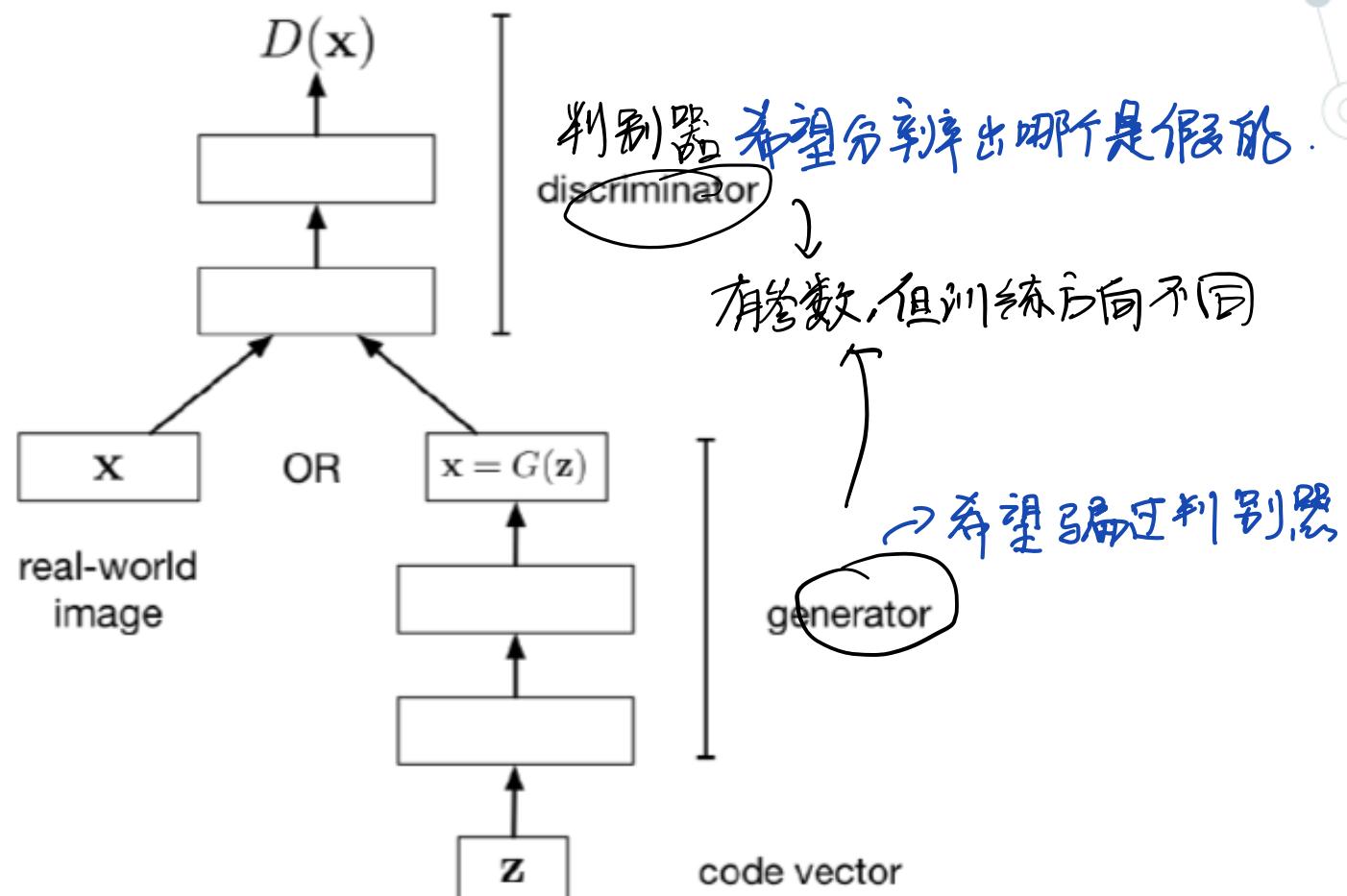


Adversarial Learning

- GAN objective for the generator is some complicated objective function defined by a neural network.
 - This means a new way of thinking about "distance".
 - We are training networks to minimize the "distance" or "divergence" between generated images and real images.
 - Instead of some hand-crafted distance metric like L1 or L2, we can make something completely new.
 - A neural network, with the right architecture, is arguably the definition of perceptual similarity (assuming our visual system is some sort of neural network).

Adversarial Learning

■ Adversarial loss



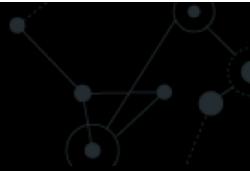
Logistic 分布:	取值	1	0
	概率	$\tau(\beta^T x)$	$1 - \tau(\beta^T x)$

$$P(y_i | x, \beta) = \tau^{y_i} (\beta^T x) (1 - \tau(\beta^T x))^{1-y_i}$$

Crossentropy-loss

$$\Rightarrow \max_{\beta} \sum_{i=1}^n \log P(y_i | x_i, \beta) = \max_{\beta} \sum_{i=1}^n [y_i \log \tau(\beta^T x_i) + (1-y_i) \log (1 - \tau(\beta^T x_i))]$$

Adversarial Learning



- Let D denote the discriminator's predicted probability of being real data
- Discriminator's cost function: cross-entropy loss for task of classifying real vs. fake images

判别器的损失

$$\mathcal{J}_D = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[-\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z}}[-\log(1 - D(G(\mathbf{z})))]$$

- One possible cost function for the generator: the opposite of the discriminator's

$$\mathcal{J}_G = -\mathcal{J}_D$$

$$= \text{const} + \mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$$

因为参数不在 D 中，只在 G 中

Two-Player Game

■ Minimax formulation

- The generator and discriminator are playing a zero-sum game against each other 最大最小化問題

$$\max_G \min_D \mathcal{J}_D$$

- #### Using parametric models

Minimax objective function:

Discriminator outputs likelihood in (0,1) of real image

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log \underbrace{D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \mathbb{E}_{z \sim p(z)} \log (1 - \underbrace{D_{\theta_d}(G_{\theta_g}(z))}_{\text{Discriminator output for generated fake data } G(z)}) \right]$$

注意，有一个负号，所以 max 和 min 需要换一下。

Learning Procedure

■ Minimax objective function

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

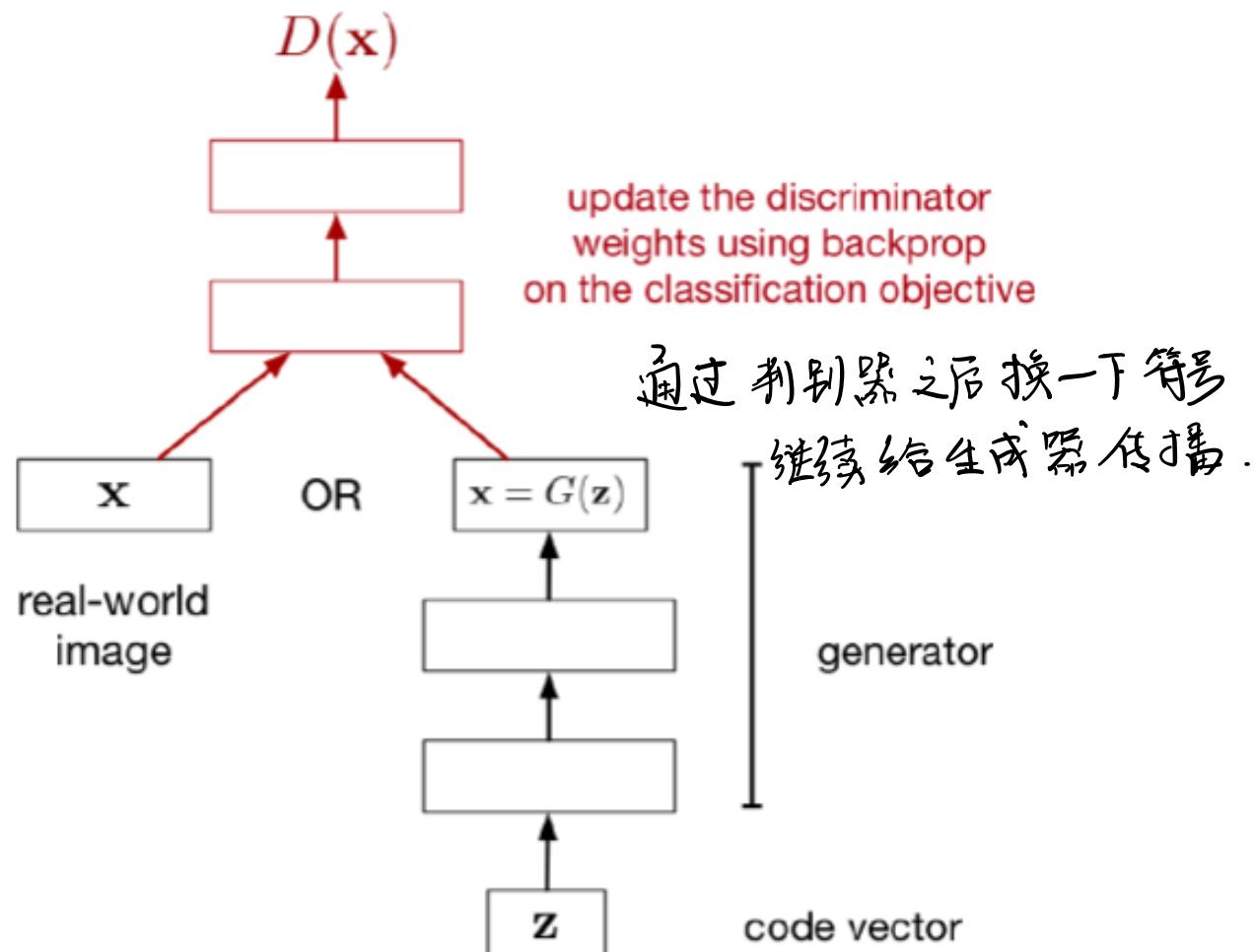
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

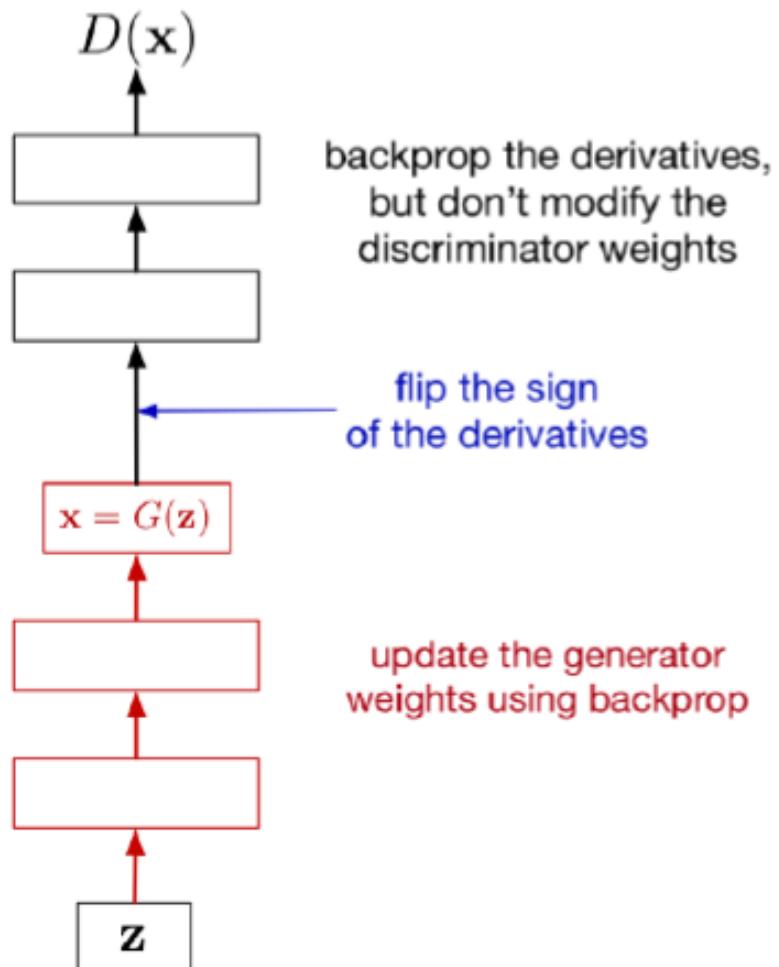
Learning Procedure

■ Updating the discriminator

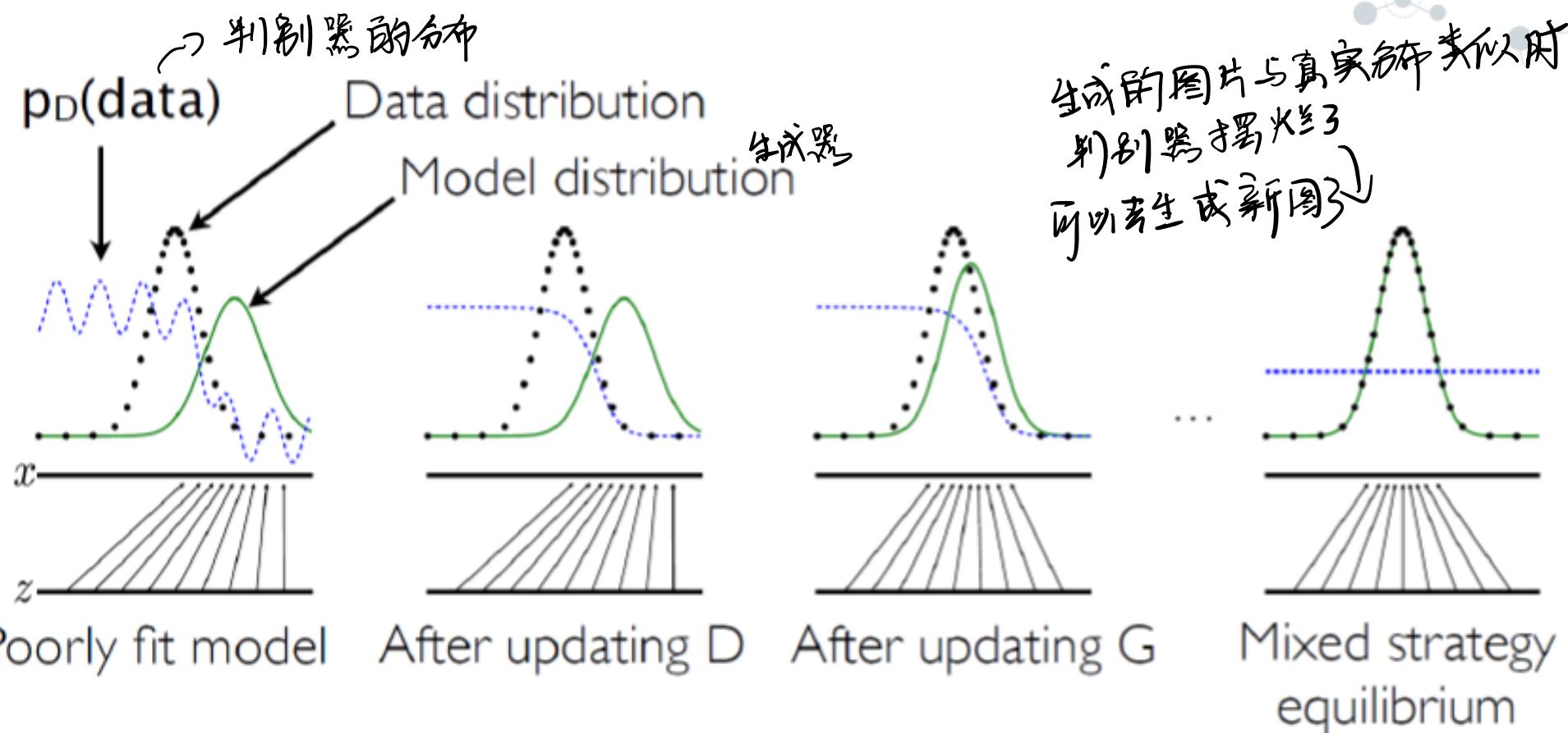


Learning Procedure

■ Updating the generator

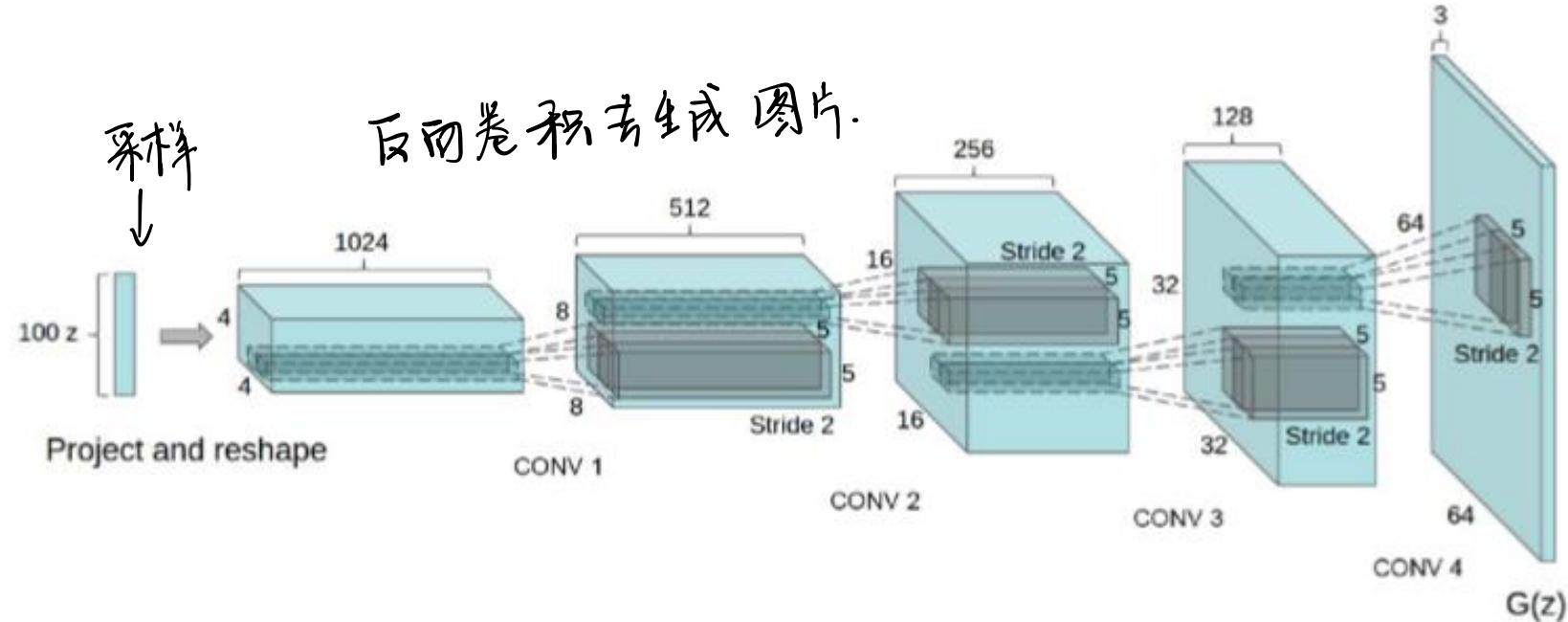


Training GANs



Typical Generator Architecture

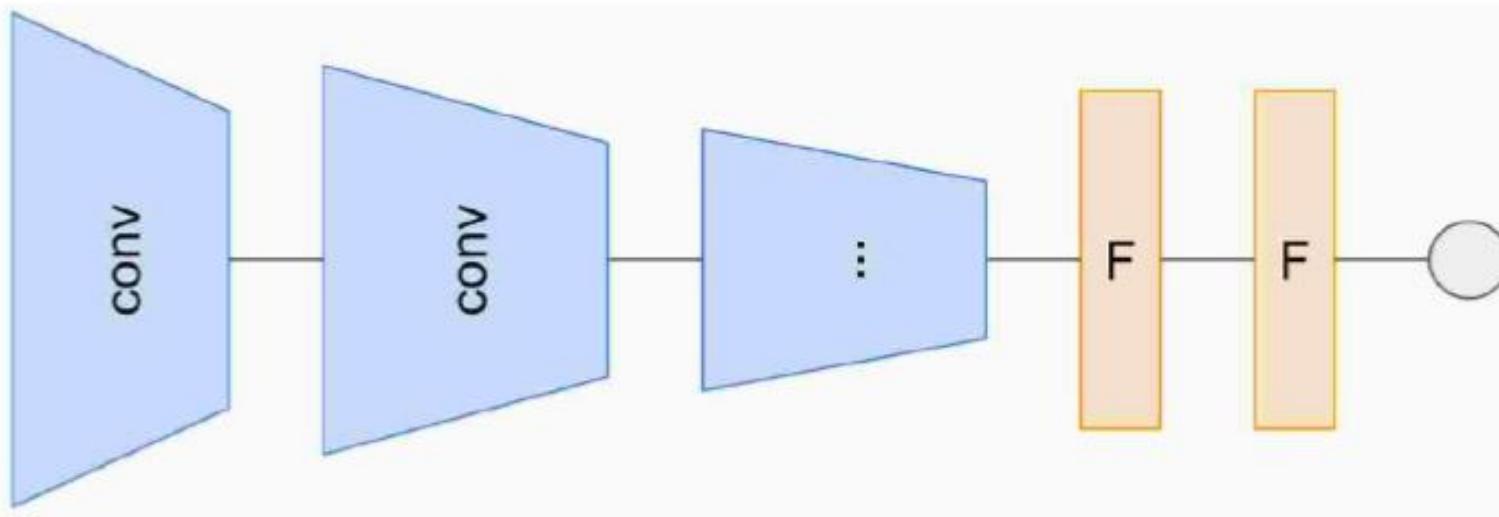
■ For images



- ▶ Unit Gaussian distribution on z , typically 10-100 dim.
- ▶ Up-convolutional deep network (reverse recognition CNN)

Typical Discriminator Architecture

- For images



- ▶ Recognition CNN model
- ▶ Binary classification output: real / synthetic

Training GANs

- Since GANs were introduced in 2014, there have been hundreds of papers introducing various architectures and training methods
- GAN Zoo: <https://github.com/hindupuravinash/the-gan-zoo>
- In general, training a GAN is tricky and **unstable**
 非常不稳定
 一旦有个好模型
 图片质量会极高.
- Many tricks:
 - S. Chintala, How to train a GAN, ICCV 2017 tutorial
 - https://github.com/soumith/talks/blob/master/2017-ICCV_Venice/How_To_Train_a_GAN.pdf

Generative Samples

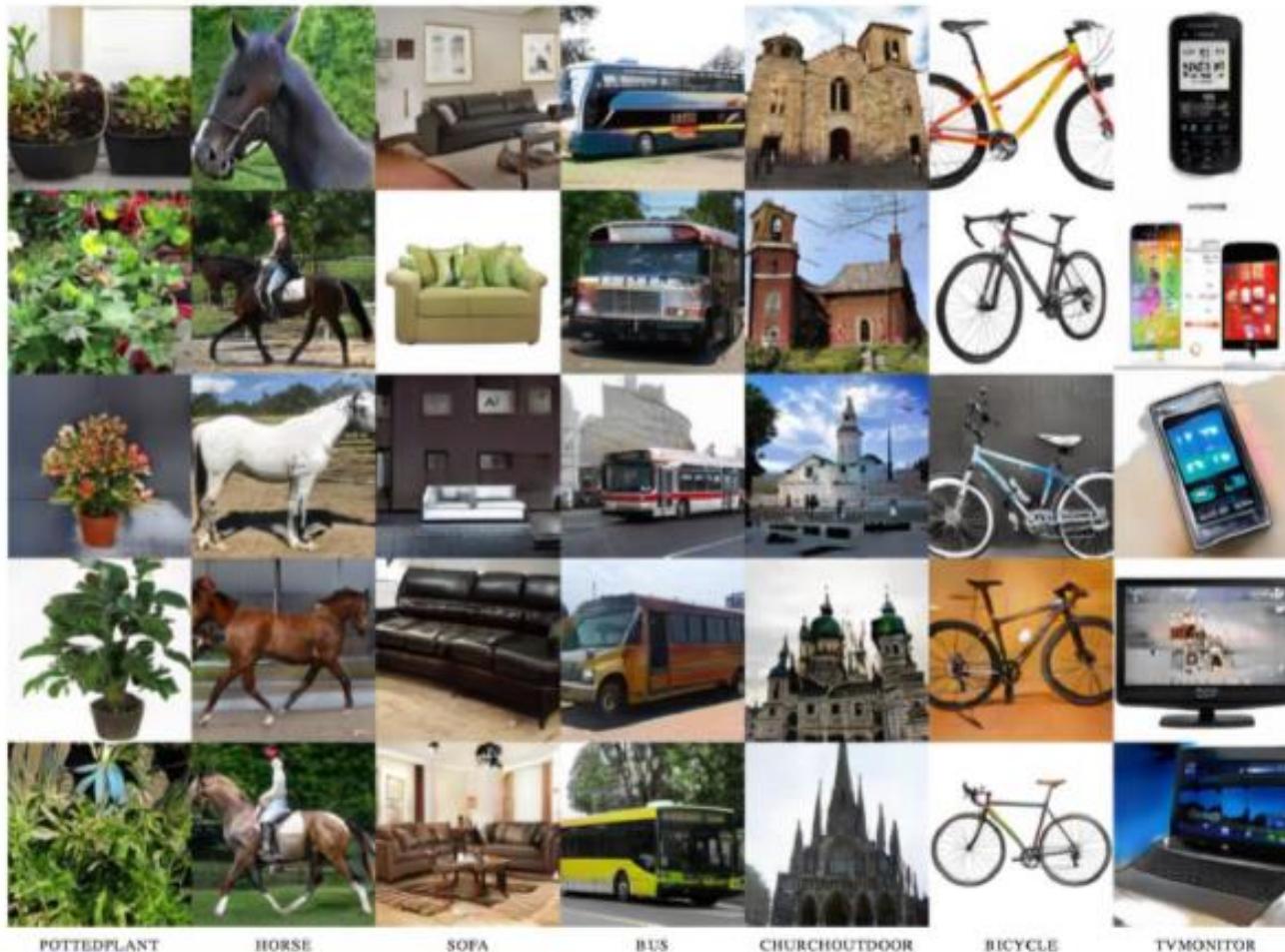
Celebrities:



Karras et al., 2017. Progressive growing of GANs for improved quality, stability, and variation

Generative Samples

Objects:



Walk Around Data Manifold

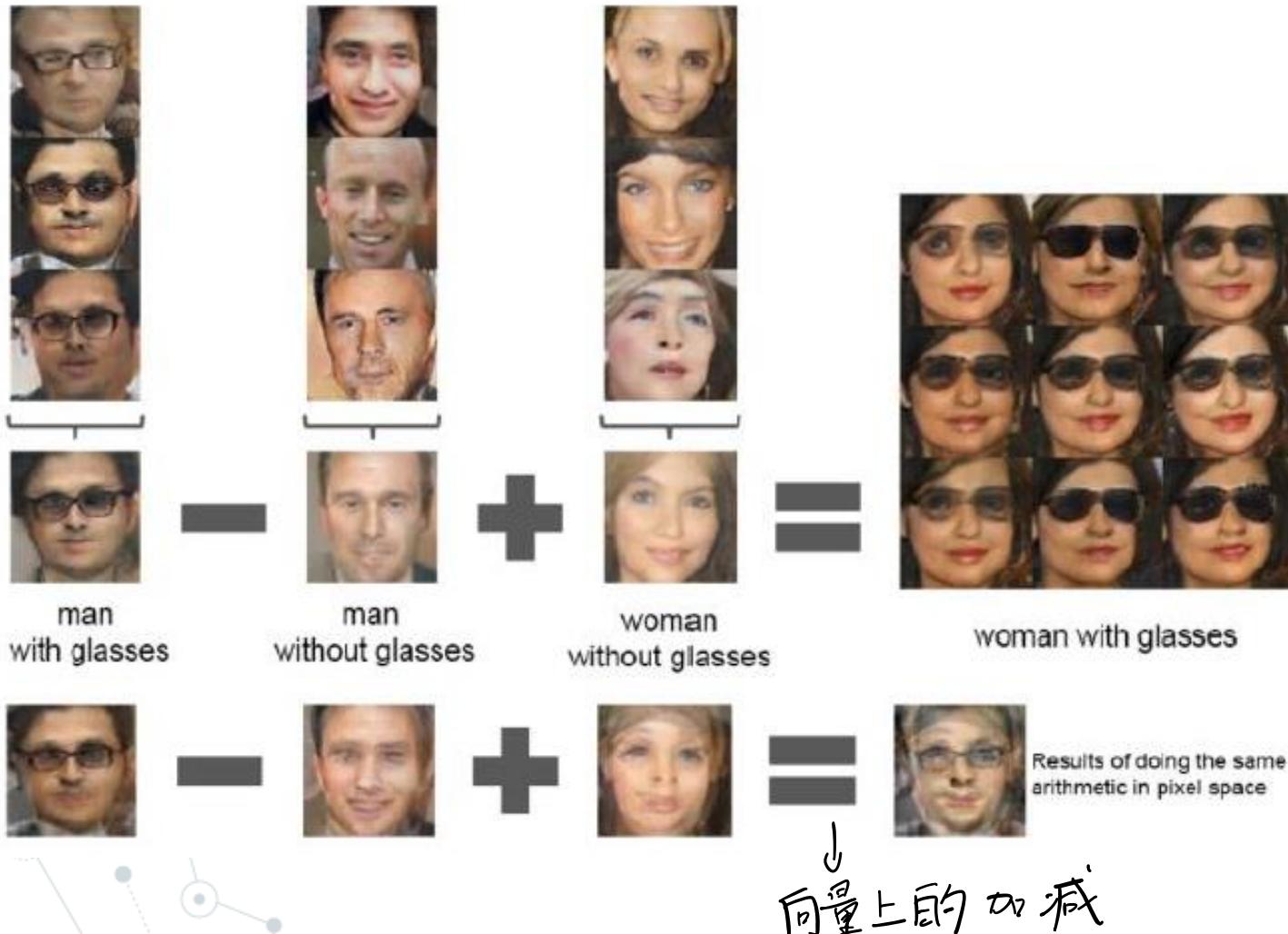
Interpolating
between
random
points in latent
space



Radford et al,
ICLR 2016

Walk Around Data Manifold

■ Vector Arithmetic



Deep Generative Networks-DGM

- Overview
- Representation Learning with Autoencoder
- Generative Adversarial Network (GAN)
- **Applications of GANs**

Conditional GANs

- Conditional GANs include a label and learn $P(X|Y)$
 - Add conditional variable y into G and D
 - Objective function

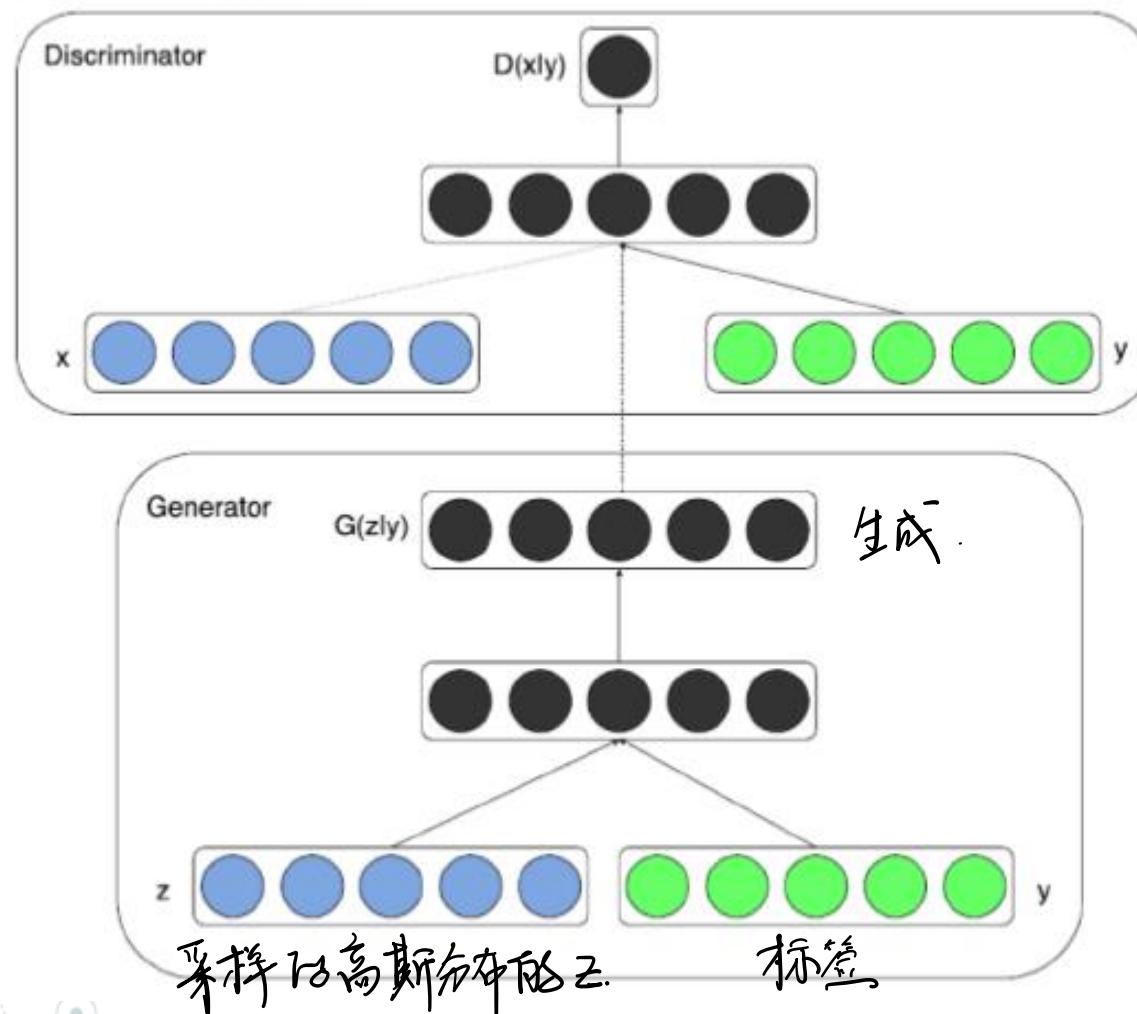
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$



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

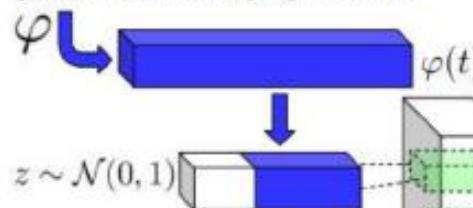
Conditional GANs

■ Model architecture



Conditional GANs

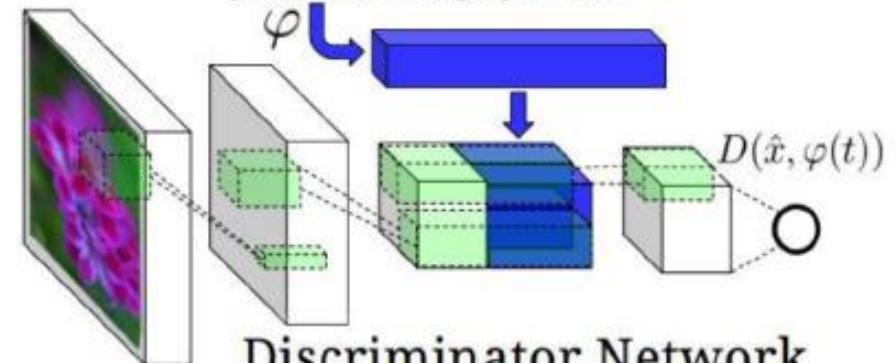
This flower has small, round violet petals with a dark purple center



Generator Network

$$\hat{x} := G(z, \varphi(t))$$

This flower has small, round violet petals with a dark purple center



Discriminator Network

this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



Reed et al 2015

StackGAN

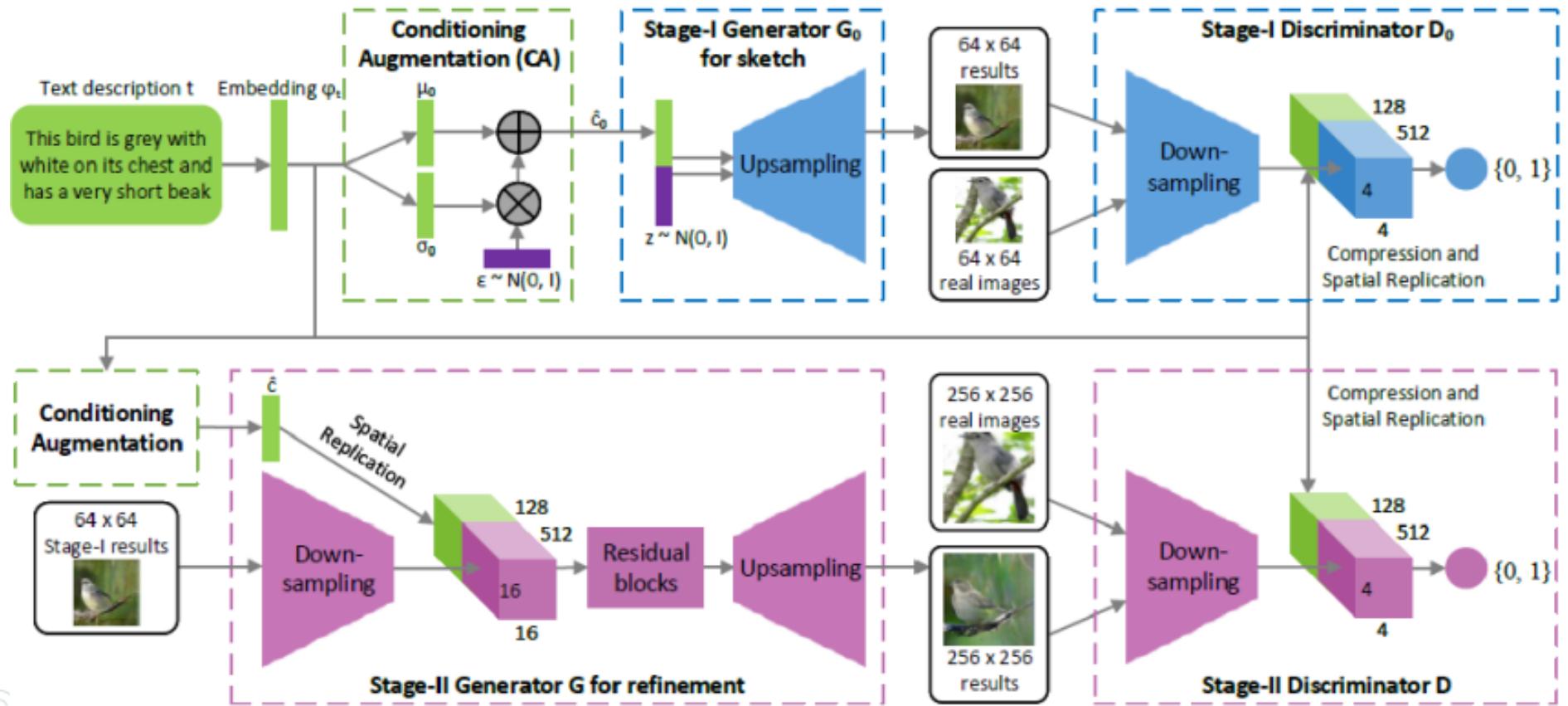
■ A coarse-to-fine manner

Text description	This bird is blue with white and has a very short beak	This bird has wings that are brown and has a yellow belly	A white bird with a black crown and yellow beak	This bird is white, black, and brown in color, with a brown beak	The bird has small beak, with reddish brown crown and gray belly	This is a small, black bird with a white breast and white on the wingbars.	This bird is white black and yellow in color, with a short black beak	
Stage-I images								
Stage-II images								

Zhang et al. 2016

StackGAN

- Use stacked GAN structure



More StackGAN results

Text
description

This flower is pink, white, and yellow in color, and has petals that are striped

This flower has a lot of small purple petals in a dome-like configuration

This flower is white and yellow in color, with petals that are wavy and smooth

This flower has petals that are dark pink with white edges and pink stamen

64x64
GAN-INT-
CLS



256x256
StackGAN



Image-to-Image Translation

- One-to-many or many-to-one mapping [Isola et al., 2016]

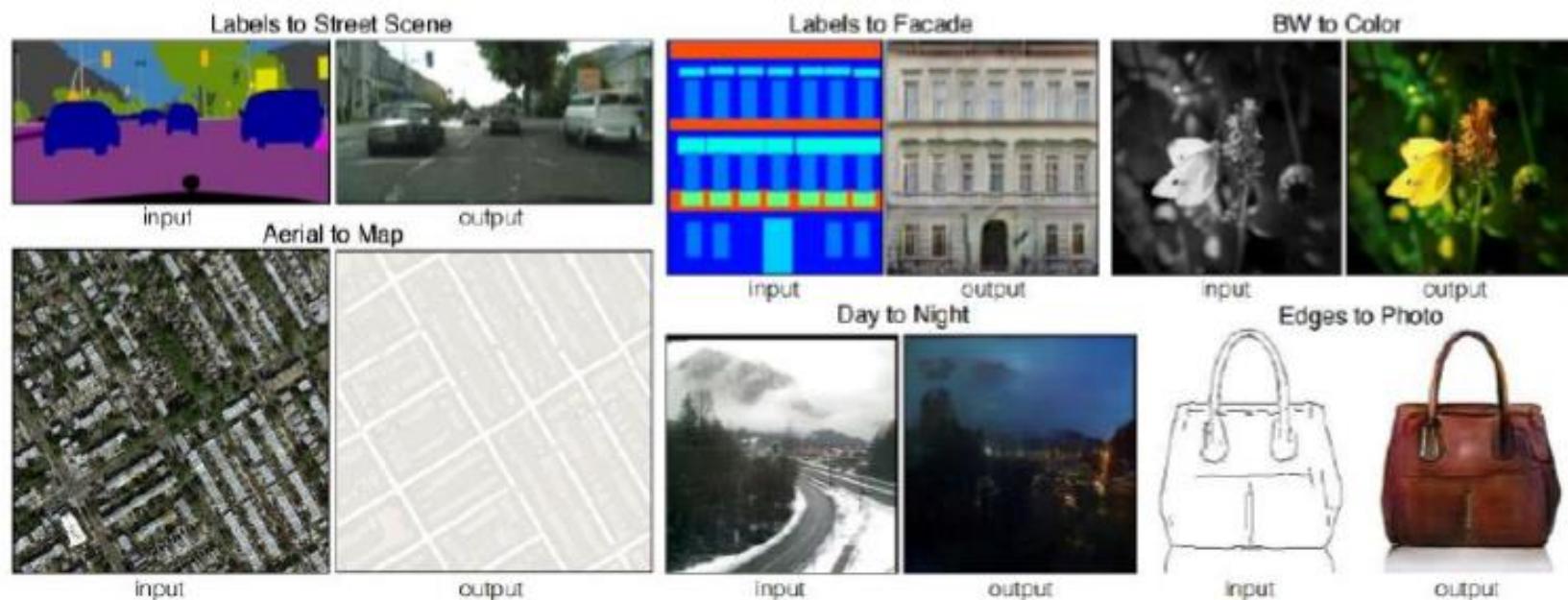


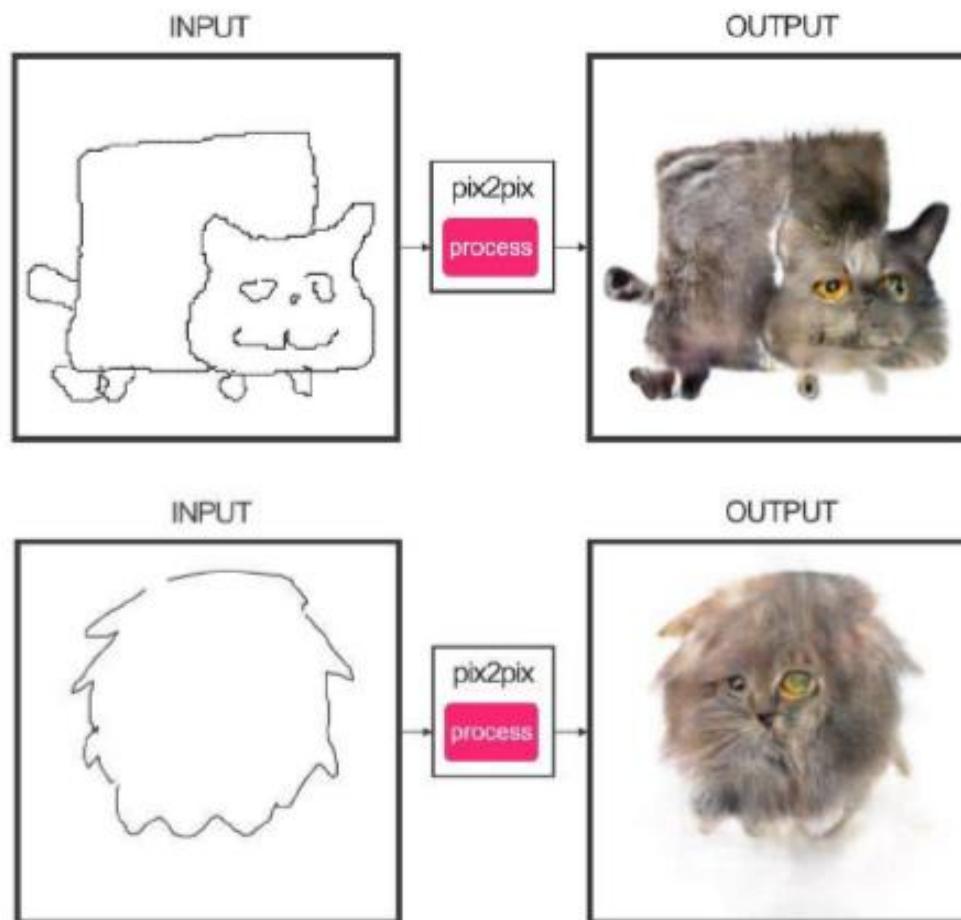
Image-to-Image Translation

■ More results



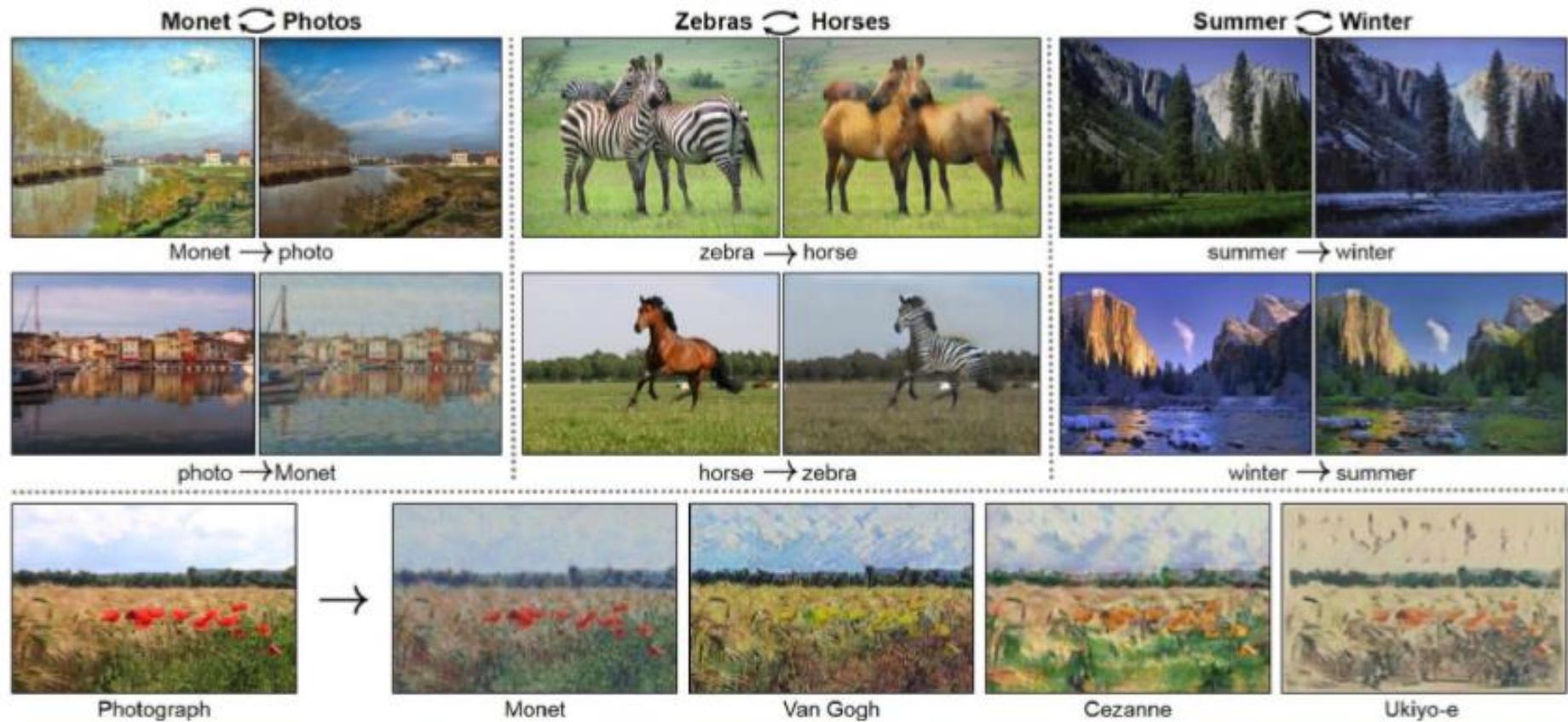
Image-to-Image Translation

■ More results



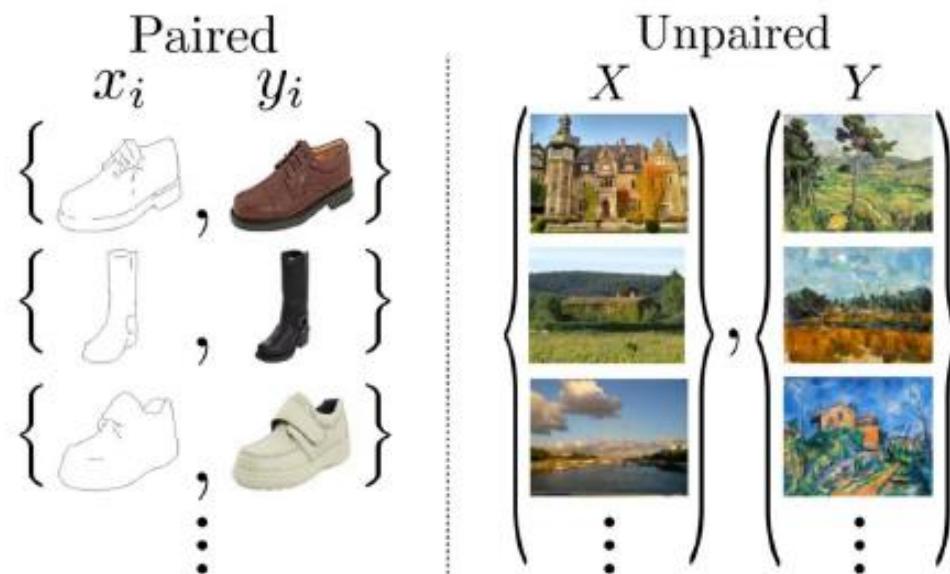
CycleGAN

- Image-to-image translation without paired data



CycleGAN

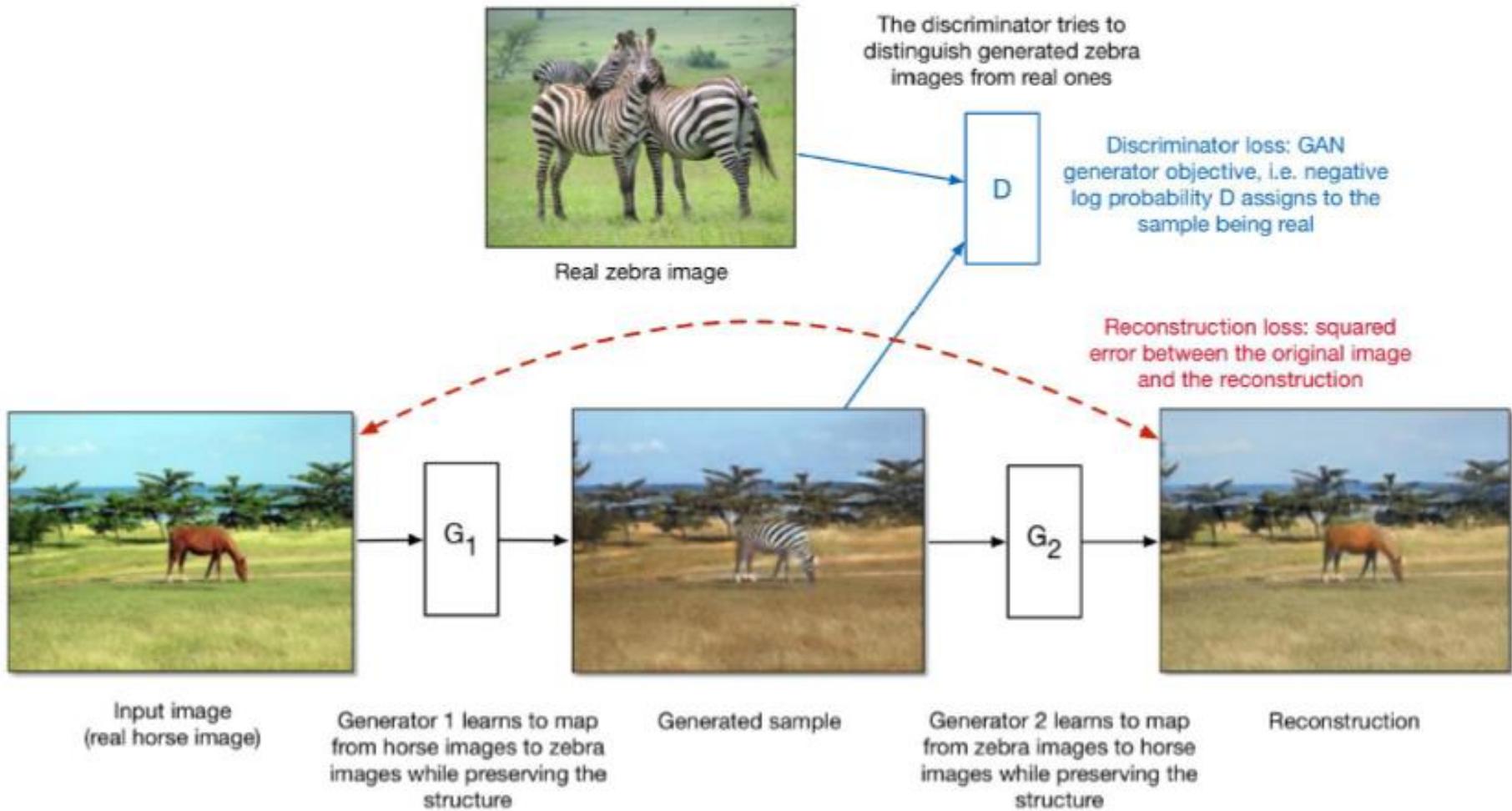
- If we had paired data (same content in both styles), this would be a supervised learning problem. But this is hard to find.



CycleGAN

- If we had paired data (same content in both styles), this would be a supervised learning problem. But this is hard to find.
- The CycleGAN architecture learns to do it from unpaired data.
 - Train two different generator nets to go from style 1 to style 2, and vice versa.
 - Make sure the generated samples of style 2 are indistinguishable from real images by a discriminator net.
 - Make sure the generators are **cycle-consistent**: mapping from style 1 to style 2 and back again should give you almost the original image.

CycleGAN



CycleGAN

- Results



- More details

<https://hardikbansal.github.io/CycleGANBlog/>