

《神经网络与深度学习》



深度生成模型 Deep Generative Models

<https://nndl.github.io/>

参考

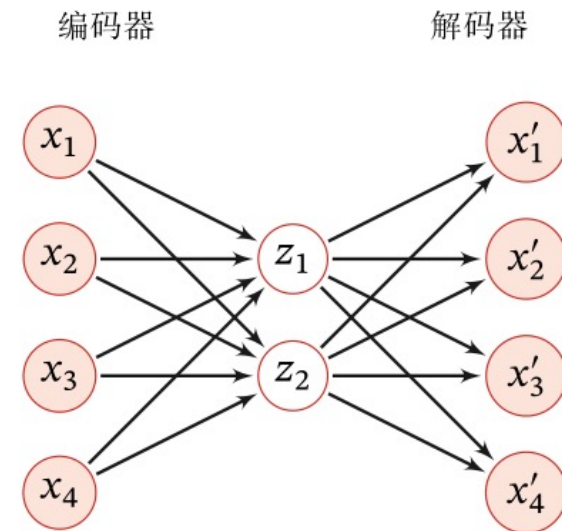
- ▶ 《神经网络与深度学习》 - 深度生成模型
 - ▶ <https://nndl.github.io/>
- ▶ 一些例子来自于李宏毅 《Introduction of Generative Adversarial Network (GAN)》

深度生成模型

- ▶ 深度生成模型就是利用神经网络构建生成模型。
 - ▶ 变分自编码器 (Variational Autoencoder, VAE)
 - ▶ [Kingma and Welling, 2013, Rezende et al., 2014]
 - ▶ 生成对抗网络 (Generative Adversarial Network, GAN)
 - ▶ [Goodfellow et al., 2014]

自编码器 (Auto-Encoder)

- ▶ 编码器 (Encoder) $f : \mathbb{R}^D \rightarrow \mathbb{R}^M$
- ▶ 解码器 (Decoder) $g : \mathbb{R}^M \rightarrow \mathbb{R}^D$

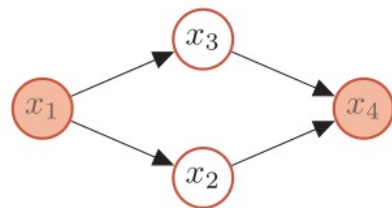




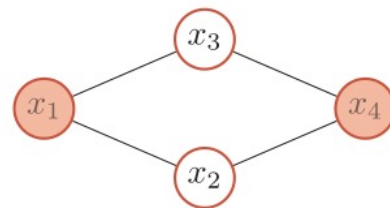
概率图的参数估计

概率图模型

- ▶ 概率图模型是指一种用图结构来描述多元随机变量之间条件独立关系的概率模型。



(a) 有向图：贝叶斯网络



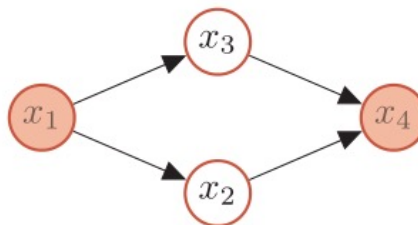
(b) 无向图：马尔可夫随机场

- ▶ 每个节点都对应一个随机变量，可以是观察变量，隐变量或是未知参数等；
- ▶ 每个连接表示两个随机变量之间具有依赖关系。

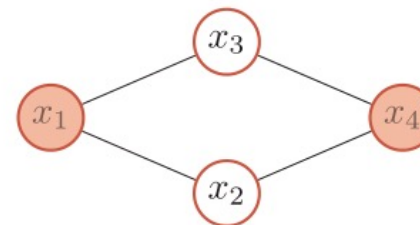
概率图模型

▶ 模型表示 (图结构)

- ▶ 有向图
- ▶ 无向图



(a) 有向图：贝叶斯网络



(b) 无向图：马尔可夫随机场

▶ 推断 (Inference)

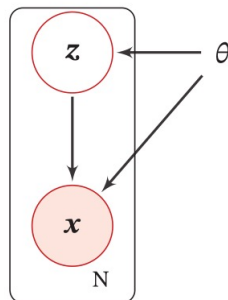
- ▶ 给定部分变量，推断另一部分变量的后验概率。

▶ 学习 (Learning)

- ▶ 参数学习：给定一组训练样本，求解模型参数

含隐变量的参数学习

▶ 隐变量即变量是不可观测的



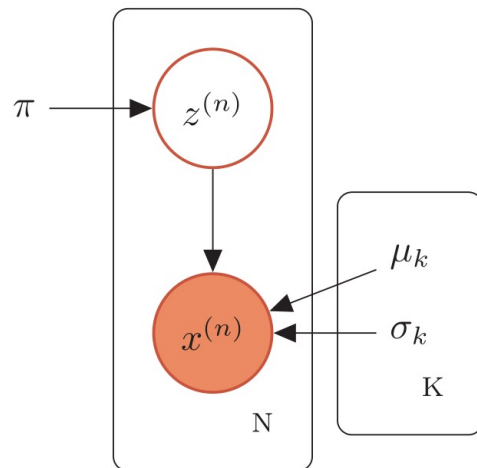
▶ 边际似然函数 (Marginal Likelihood)

$$p(\mathbf{x}; \theta) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \theta)$$

▶ 需要用EM算法进行参数估计

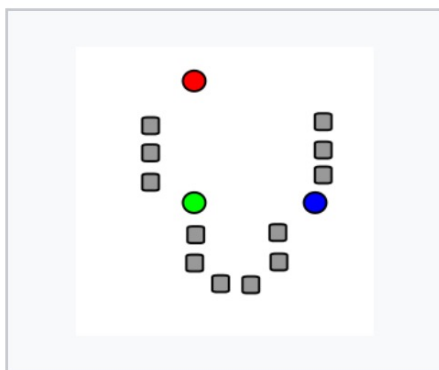
高斯混合模型

图模型表示

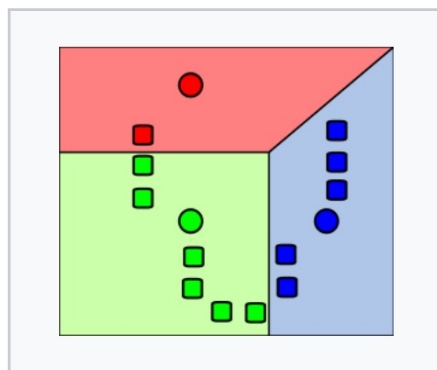


$$\mathcal{N}(x|\mu_k, \sigma_k) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)$$

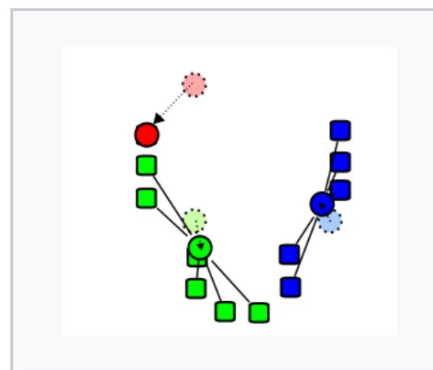
一个简单的解法：K-means



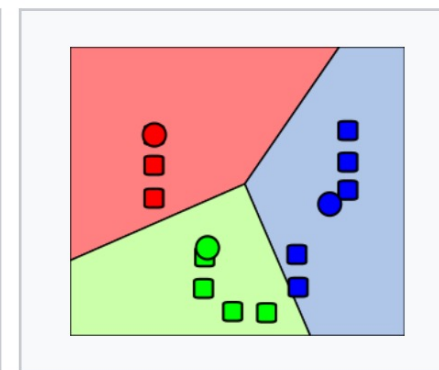
1. k initial "means" (in this case $k=3$) are randomly generated within the data domain (shown in color).



2. k clusters are created by associating every observation with the nearest mean. The partitions here represent the **Voronoi diagram** generated by the means.



3. The **centroid** of each of the k clusters becomes the new mean.



4. Steps 2 and 3 are repeated until convergence has been reached.

K-means算法

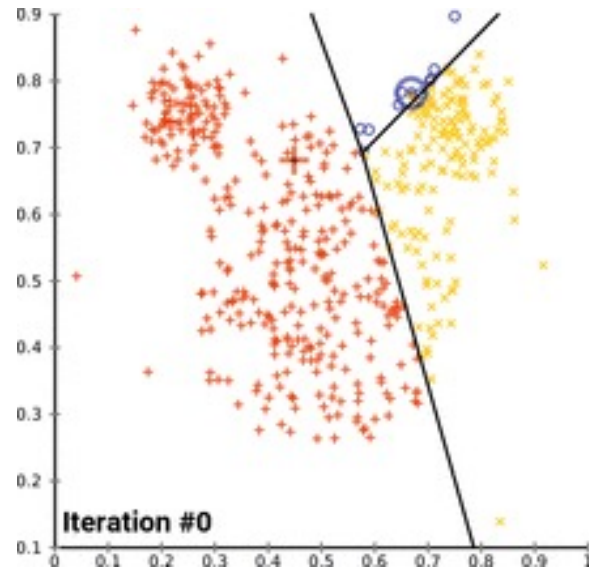
- ▶ 初始化中心点 $m_1^{(1)}, \dots, m_k^{(1)}$
- ▶ 迭代执行下面两步
 - ▶ 分配步 (Assignment step) :

$$S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\|^2 \leq \|x_p - m_j^{(t)}\|^2 \forall j, 1 \leq j \leq k\},$$

- ▶ 更新步 (Update step)

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

K-means算法



期望最大化 (Expectation-Maximum , EM) 算法

▶ 假设有一组变量，有部分变量是不可观测的，如何进行参数估计呢？

D.2.7.1 Jensen 不等式

如果 X 是随机变量, g 是凸函数, 则

$$g(E[X]) \leq E[g(X)]. \quad (D.41)$$

等式当且仅当 X 是一个常数或 g 是线性时成立, 这个性质称为 Jensen 不等式.

特别地, 对于凸函数 g 定义域上的任意两点 x_1, x_2 和一个标量 $\lambda \in [0, 1]$, 有

$$g(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda g(x_1) + (1 - \lambda)g(x_2), \quad (D.42)$$

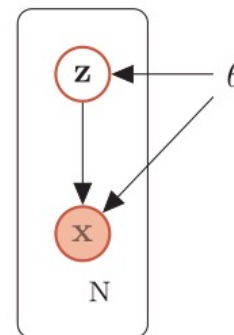
即凸函数 g 上的任意两点的连线位于这两点之间函数曲线的上方.

变分

对数边际似然函数

$$\begin{aligned} \log p(\mathbf{x}; \theta) &= \log \sum_{\mathbf{z}} q(\mathbf{z}) \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q(\mathbf{z})} \\ &\geq \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q(\mathbf{z})} \\ &\triangleq ELBO(q, \mathbf{x}; \theta), \end{aligned}$$

证据下界



期望最大化 (Expectation-Maximum , EM) 算法

▶ 假设有一组变量，有部分变量是不可观测的，如何进行参数估计呢？

对数边际似然函数

$$\begin{aligned}\log p(x; \theta) &= \log E_{z \sim q(z)} \left(\frac{p(z, x; \theta)}{q(z)} \right) \\ &\geq E_{z \sim q(z)} \log \left(\frac{p(z, x; \theta)}{q(z)} \right) \\ &= \text{ELBO}\end{aligned}$$

另外一种推导

$$\begin{aligned}\log p(\mathbf{x}; \theta) &= \sum_{\mathbf{z}} q(\mathbf{z}) \log p(\mathbf{x}; \theta) \\ &= \sum_{\mathbf{z}} q(\mathbf{z}) \left(\log p(\mathbf{x}, \mathbf{z}; \theta) - \log p(\mathbf{z}|\mathbf{x}; \theta) \right) \\ &= \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q(\mathbf{z})} - \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}; \theta)}{q(\mathbf{z})} \\ &= ELBO(q, \mathbf{x}; \theta) + \text{KL}(q(\mathbf{z}) \| p(\mathbf{z}|\mathbf{x}; \theta)),\end{aligned}$$

EM算法

$$ELBO(q, \mathbf{x}; \theta) + \text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x}; \theta))$$

▶ E步

$$q_{t+1}(\mathbf{z}) = \arg \max_q ELBO(q, \mathbf{x}; \theta) + \text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x}; \theta))$$

▶ M步

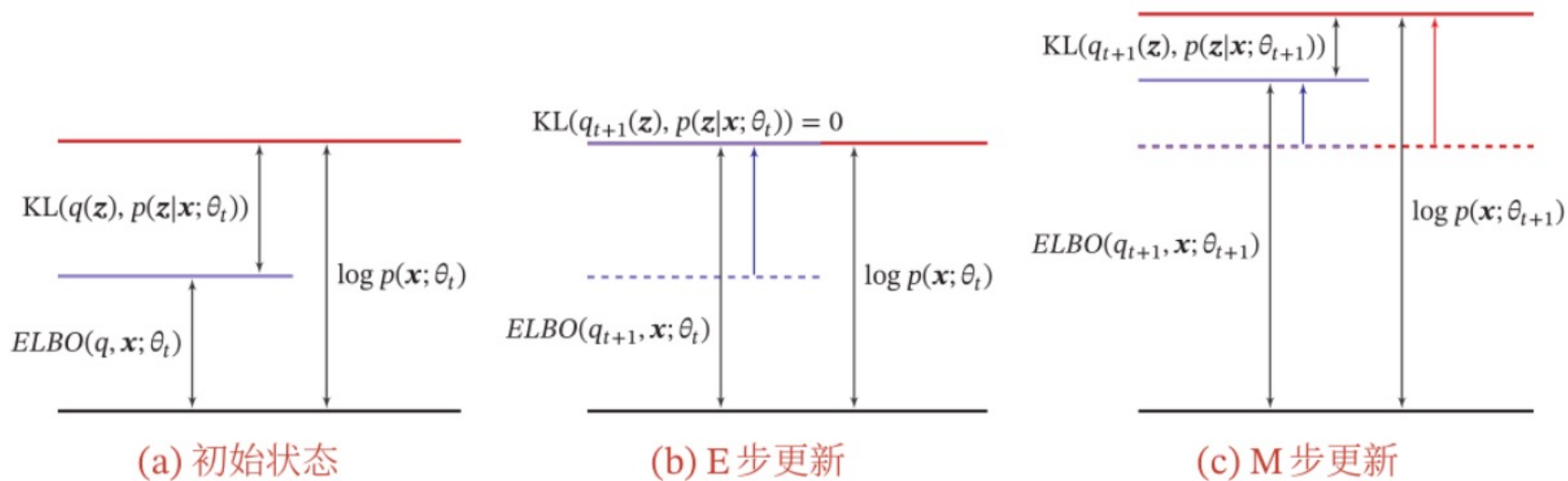
$$\theta_{t+1} = \arg \max_{\theta} ELBO(q, \mathbf{x}; \theta) + \text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x}; \theta))$$

收敛性

$$\log p(\mathbf{x}; \theta_{t+1}) \geq \boxed{ELBO(q_{t+1}, \mathbf{x}; \theta_{t+1})} \geq \boxed{ELBO(q_{t+1}, \mathbf{x}; \theta_t) = \log p(\mathbf{x}; \theta_t)}$$

M 步

E 步

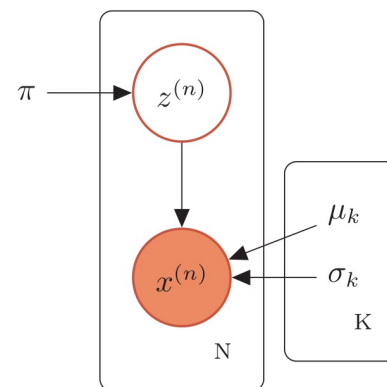


GMM Revisit

E步 先固定参数 μ, σ , 计算后验分布 $p(z^{(n)}|x^{(n)})$

$$\begin{aligned}\gamma_{nk} &\triangleq p(z^{(n)} = k|x^{(n)}) \\ &= \frac{p(z^{(n)})p(x^{(n)}|z^{(n)})}{p(x^{(n)})} \\ &= \frac{\pi_k \mathcal{N}(x^{(n)}|\mu_k, \sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(x^{(n)}|\mu_k, \sigma_k)},\end{aligned}$$

其中 γ_{nk} 定义了样本 $x^{(n)}$ 属于第 k 个高斯分布的后验概率。



$$\mathcal{N}(x|\mu_k, \sigma_k) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)$$

GMM Revisit

M步 令 $q(z = k) = \gamma_{nk}$, 训练集 \mathcal{D} 的证据下界为

$$\begin{aligned} ELBO(\gamma, \mathcal{D} | \pi, \mu, \sigma) &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \log \frac{p(x^{(n)}, z^{(n)} = k)}{\gamma_{nk}} \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(\log \mathcal{N}(x^{(n)} | \mu_k, \sigma_k) + \log \frac{\pi_k}{\gamma_{nk}} \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(-\frac{(x - \mu_k)^2}{2\sigma_k^2} - \log \sigma_k + \log \pi_k \right) + C, \end{aligned}$$

其中 C 为和参数无关的常数。

$$\text{s.t. } \sum_{k=1}^K \pi_k = 1.$$

$$\mathcal{N}(x | \mu_k, \sigma_k) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)$$

$$p(x^n, z^n = k) = \pi_k \mathcal{N}(x^n | \mu_k, \sigma_k)$$

$$\log\left(\frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)\right) = \frac{-(x - \mu_k)^2}{2\sigma_k^2} - \log \sigma_k - \log \sqrt{2\pi}$$

$$C = \sum_{n=1}^N \sum_{k=1}^K (-\log \gamma_{nk} - \log \sqrt{2\pi}) * \gamma_{nk}$$

GMM Revisit

M步 令 $q(z = k) = \gamma_{nk}$, 训练集 \mathcal{D} 的证据下界为

$$\begin{aligned} ELBO(\gamma, \mathcal{D} | \pi, \mu, \sigma) &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \log \frac{p(x^{(n)}, z^{(n)} = k)}{\gamma_{nk}} \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(\log \mathcal{N}(x^{(n)} | \mu_k, \sigma_k) + \log \frac{\pi_k}{\gamma_{nk}} \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(\frac{-(x - \mu_k)^2}{2\sigma_k^2} - \log \sigma_k + \log \pi_k \right) + C, \end{aligned}$$

其中 C 为和参数无关的常数。

$$\text{s.t. } \sum_{k=1}^K \pi_k = 1.$$

$$ELBO(\gamma, \mathcal{D} | \pi, \mu, \sigma) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(\frac{-(x - \mu_k)^2}{2\sigma_k^2} - \log \sigma_k + \log \pi_k \right) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) + C$$

$$\frac{\partial L}{\partial \pi_k} = \sum_{n=1}^N \gamma_{nk} \frac{1}{\pi_k} + \lambda = 0$$

$$\frac{\partial L}{\partial \mu_k} = \sum_{n=1}^N \gamma_{nk} \frac{(x^n - \mu_k)}{\sigma_k^2} = 0$$

$$\frac{\partial L}{\partial \sigma_k} = \sum_{n=1}^N \gamma_{nk} \left(\frac{(x - \mu_k)^2}{\sigma_k^3} - \frac{1}{\sigma_k} \right) = 0$$

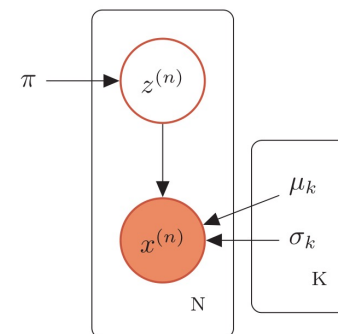
GMM Revisit

$$N_k = \sum_{n=1}^N \gamma_{nk}$$

$$\pi_k = \frac{N_k}{N}$$

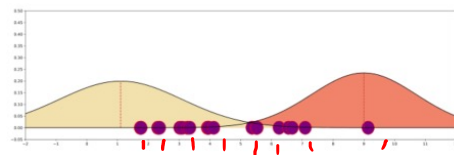
$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} x^{(n)}$$

$$\sigma_k^2 = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (x^{(n)} - \mu_k)^2$$

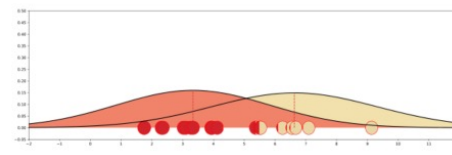


$$\mathcal{N}(x|\mu_k, \sigma_k) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)$$

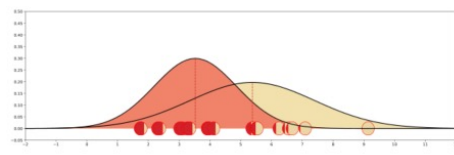
GMM的参数学习



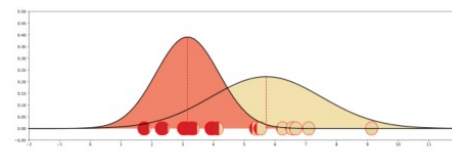
(a) 初始化



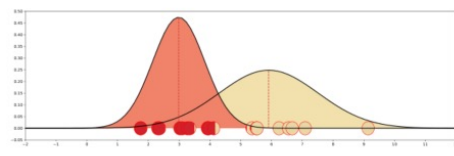
(b) 第1次迭代



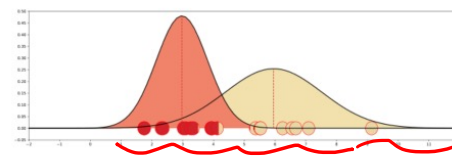
(c) 第4次迭代



(d) 第8次迭代

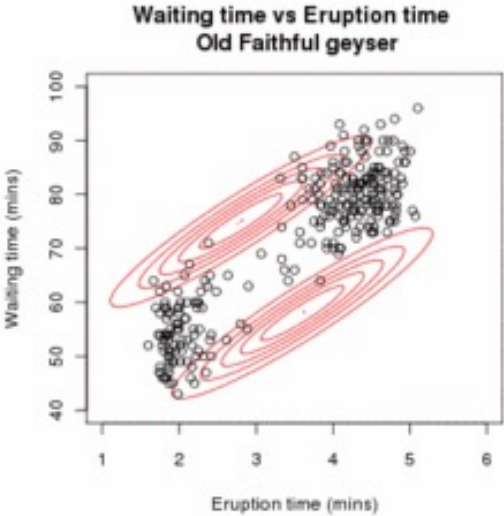
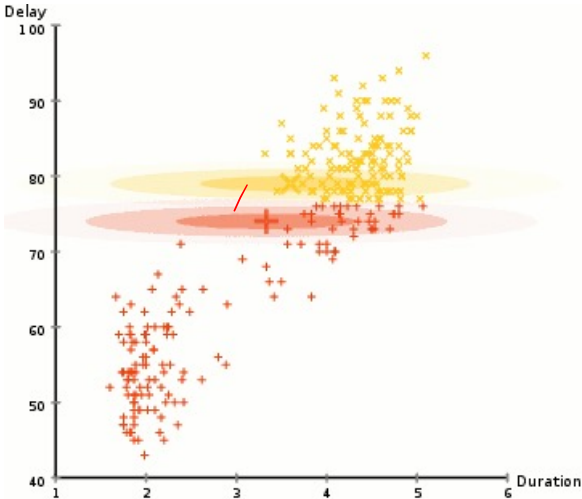


(e) 第12次迭代



(f) 第16次迭代

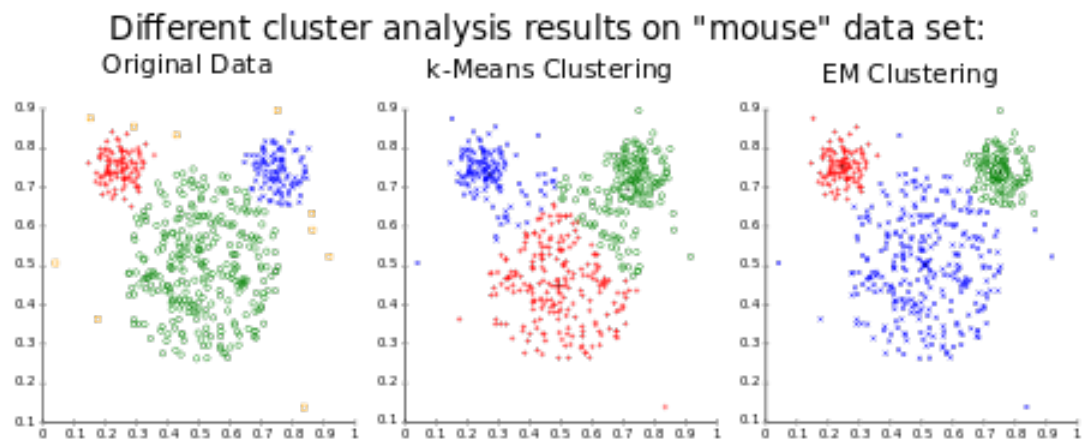
GMM的参数学习



k-means clustering vs. EM clustering

- 随机确定k中心点
 - 计算每个样本点到中心点的距离
 - 把每个样本点分配到离它最近的中心点
 - 对每一个中心点，根据分配到的样本点重新计算中心点位置
 - 迭代直到收敛
- 随机初始化k个高斯分布的均值和方差
 - 计算每个样本点属于每个高斯分布的后验概率
 - 用后验概率重新估计每个高斯分布的均值和方差
 - 迭代直到收敛

k-means clustering vs. EM clustering

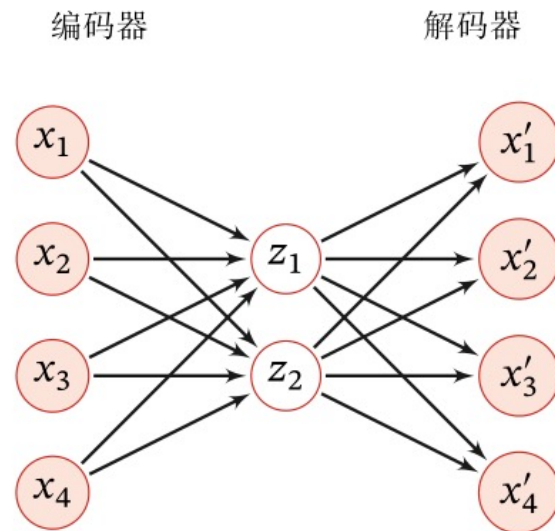


图模型与神经网络的关系

1. 图模型的节点是随机变量，图结构主要描述变量之间的依赖关系；而神经网络中的节点是神经元，是一个计算节点。
2. 图模型中每个变量一般有着明确的解释，变量之间依赖关系一般是人工来定义；而神经网络中单个神经元没有直观的解释。
3. 神经网络是判别式模型，直接用来分类；而图模型既可以是判别式模型，也可以是生成式模型。
4. 神经网络和概率图的结合越来越紧密。

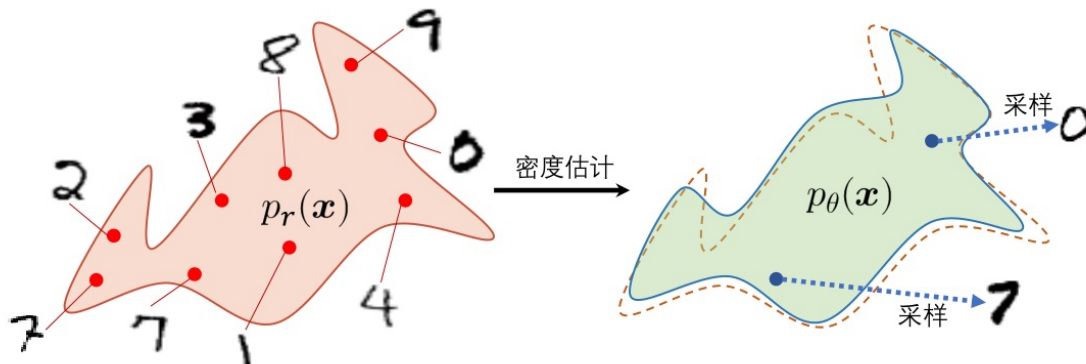
自编码器 (Auto-Encoder)

- ▶ 编码器 (Encoder) $f : \mathbb{R}^D \rightarrow \mathbb{R}^M$
- ▶ 解码器 (Decoder) $g : \mathbb{R}^M \rightarrow \mathbb{R}^D$

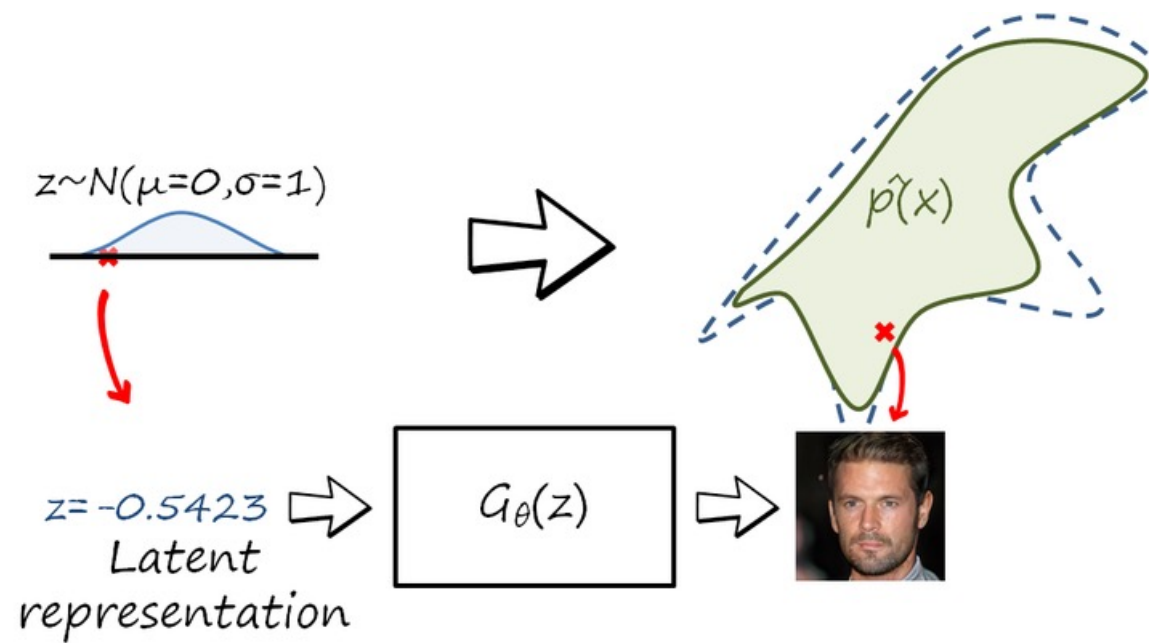


生成模型：一系列用于随机生成可观测数据的模型

- ▶ 生成模型包含两个步骤：
 - ▶ 密度估计
 - ▶ 采样



生成数据的另一种思路

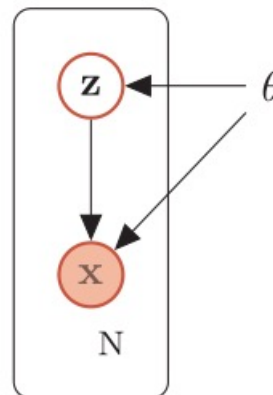




变分自编码器
Variational Autoencoder , VAE

概率生成模型

- ▶ 含隐变量的概率图模型
- ▶ 生成数据 x 的过程可以分为两步进行：
 - ▶ 根据隐变量的先验分布 $p(z;\theta)$ 采样得到样本 z ;
 - ▶ 根据条件分布 $p(x|z;\theta)$ 采样得到 x 。



EM算法回顾

▶ 给定一个样本 \mathbf{x} ，其对数边际似然 $\log p(\mathbf{x}|\theta)$ 可以分解为

$$\log p(\mathbf{x}|\theta) = \underbrace{ELBO(q, \mathbf{x}|\theta, \phi)}_{\text{M step}} + \underbrace{D_{\text{KL}}(q(\mathbf{z}|\phi) \| p(\mathbf{z}|\mathbf{x}, \theta))}_{\text{E step}}$$

$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\phi)} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\phi)} \right]$

变分自编码器(VAE)

▶ 变分自编码器的模型结构可以分为两个部分：

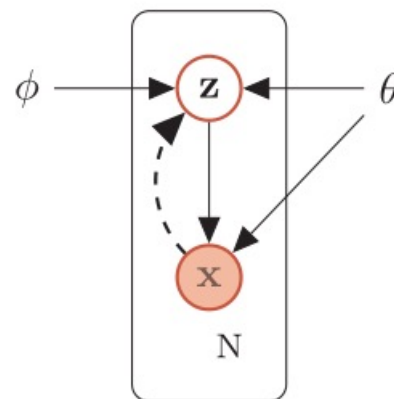
▶ 寻找后验分布 $p(z|x;\theta)$ 的变分近似

$q(z|x;\phi^*)$;

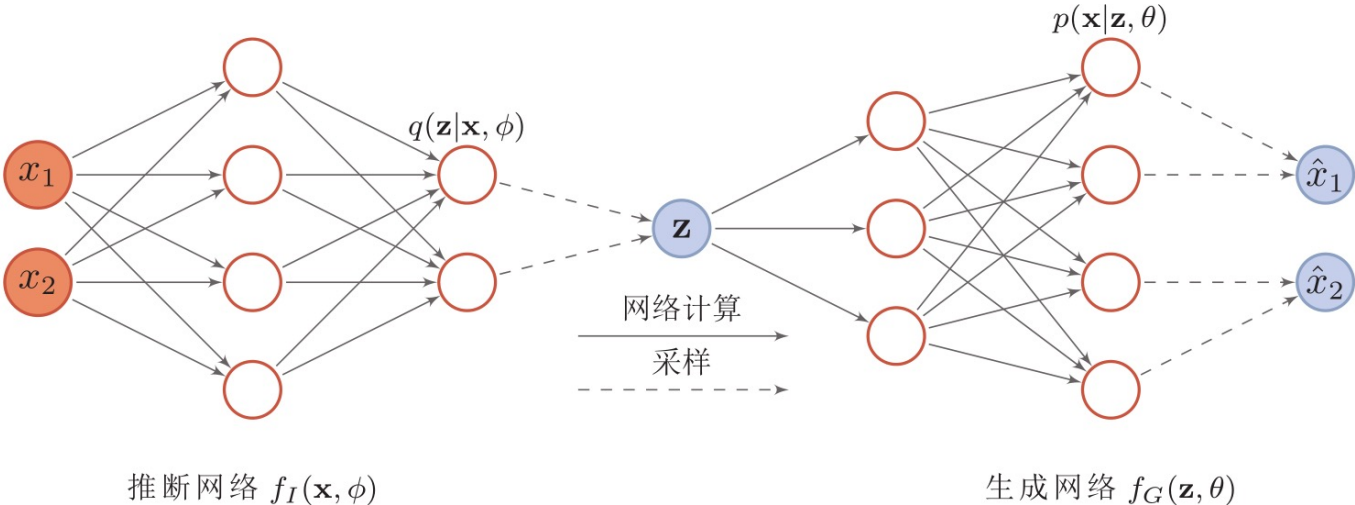
▶ 变分推断：用简单的分布 q 去近似复杂的分布 $p(z|x;\theta)$

▶ 在已知 $q(z|x;\phi^*)$ 的情况下，估计更好的分布 $p(x|z;\theta)$ 。

用神经网络来替代



变分自编码器



变分自编码器 vs. 自编码器

- ▶ 变分自编码器的名字来自于其整个网络结构和自编码器比较类似。
 - ▶ 推断网络可看作“编码器”，将可观测变量映射为隐变量
 - ▶ 生成网络看作“解码器”，将隐变量映射为可观测变量
- ▶ 变分自编码器的原理和自编码器原理不同
 - ▶ 变分自编码器中的编码器和解码器的输出为分布（或分布的参数），而不是确定的编码。

推断网络

▶ 推断网络

$$q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_I, \boldsymbol{\sigma}_I^2 I)$$

$$\begin{bmatrix} \boldsymbol{\mu}_I \\ \boldsymbol{\sigma}_I \end{bmatrix} = f_I(\mathbf{x}, \phi)$$

$$h = \text{ReLU}(W^1 x + b^1)$$

$$\boldsymbol{\mu}_I = W^2 x + b^2$$

$$\boldsymbol{\sigma}_I = W^3 x + b^3$$

▶ E步目标

$$\phi^* = \arg \min_{\phi} D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}, \phi) || p(\mathbf{z}|\mathbf{x}, \theta))$$

$$= \arg \min_{\phi} \log p(\mathbf{x}|\theta) - ELBO(q, \mathbf{x}|\theta, \phi)$$

$$= \arg \max_{\phi} ELBO(q, \mathbf{x}|\theta, \phi).$$

生成网络

▶ 先验分布 $p(\mathbf{z}|\theta)$

- ▶ 一般假设隐变量 \mathbf{z} 的先验分布为各向同性的标准高斯分布 $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$

▶ 条件概率分布 $p(\mathbf{x}|\mathbf{z}, \theta)$

- ▶ 假设 $p(\mathbf{x}|\mathbf{z}, \theta)$ 服从对角化协方差的高斯分布

$$p(\mathbf{x}|\mathbf{z}, \theta) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_G, \boldsymbol{\sigma}_G^2 \mathbf{I})$$

▶ M步目标

$$\theta^* = \arg \max_{\theta} ELBO(q, \mathbf{x}|\theta, \phi)$$

模型汇总

$$\begin{aligned}\max_{\theta, \phi} ELBO(q, x|\theta, \phi) &= \max_{\theta, \phi} E_{z \sim q(z|x; \phi)} \left[\log \frac{p(x, z|\theta)}{q(z|x; \phi)} \right] = \max_{\theta, \phi} E_{z \sim q(z|x; \phi)} \left[\log \frac{p(x|z; \theta)p(z|\theta)}{q(z|x; \phi)} \right] \\ &= \int q(z|x; \phi) \left[\log \frac{p(x|z; \theta)p(z|\theta)}{q(z|x; \phi)} \right] dz \\ &= \int q(z|x; \phi) [\log p(x|z; \theta)] dz + \int q(z|x; \phi) \left[\log \frac{p(z|\theta)}{q(z|x; \phi)} \right] dz \\ \\ \max_{\theta, \phi} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}, \phi)} \left[\log p(\mathbf{x}|\mathbf{z}, \theta) \right] &- D_{\text{KL}} \left(q(\mathbf{z}|\mathbf{x}, \phi) || p(\mathbf{z}|\theta) \right)\end{aligned}$$

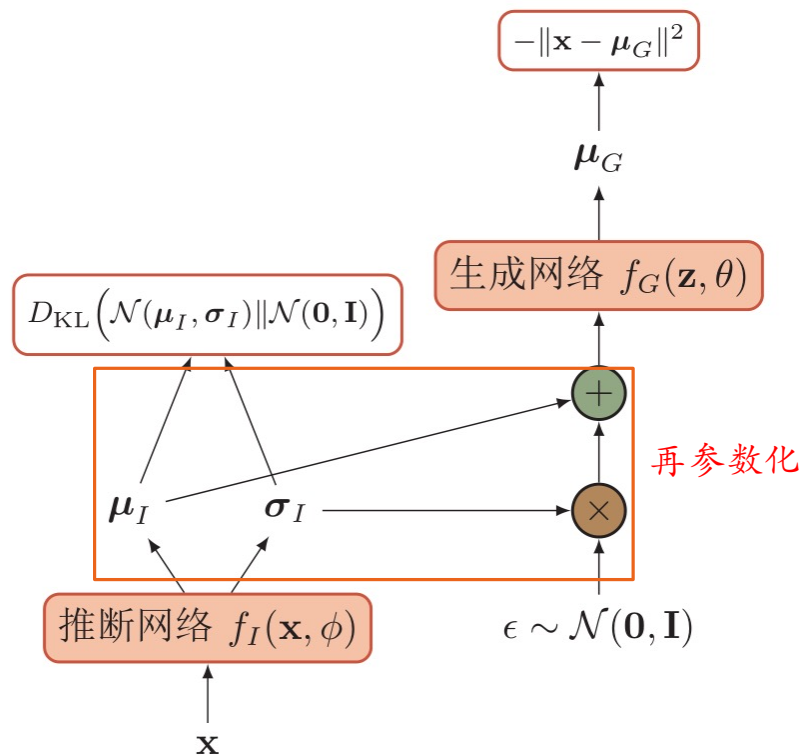
模型汇总

$$\max_{\theta, \phi} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}, \phi)} \left[\log p(\mathbf{x}|\mathbf{z}, \theta) \right] - D_{\text{KL}} \left(q(\mathbf{z}|\mathbf{x}, \phi) \parallel p(\mathbf{z}|\theta) \right)$$

最大化图像重构概率

最小化后验概率分布 $q(\mathbf{z}|\mathbf{x}, \phi)$ 和高斯分布 $p(\mathbf{z}|\theta)$ 的差异

变分自编码器的训练过程



$D_{kl}(N(\mu_I, \sigma_I) | N(0, 1))$ 推导

$$N(\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$N(0, 1) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

$$D_{kl}(N(\mu, \sigma) | N(0, 1)) = \int \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} \log \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}} / \sqrt{2\pi}\sigma}{e^{-\frac{x^2}{2}} / \sqrt{2\pi}} dx$$

$D_{kl}(N(\mu_I, \sigma_I) | N(0, I))$ 推导

$$\begin{aligned} D_{kl}(N(\mu, \sigma^2) | N(0, 1)) &= \int \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} \log \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}/\sqrt{2\pi}\sigma}{e^{-x^2/2}/\sqrt{2\pi}} dx \\ &= \int \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} \log \frac{1}{\sigma} e^{\frac{1}{2}[x^2 - \frac{(x-\mu)^2}{\sigma^2}]} dx \\ &= \int \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} [\log(\sigma^2)^{-\frac{1}{2}} + \frac{1}{2}(x^2 - \frac{(x-\mu)^2}{\sigma^2})] dx \\ &= \frac{1}{2} * \int \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} [-\log(\sigma^2) + x^2 - \frac{(x-\mu)^2}{\sigma^2}] dx \\ &= \frac{1}{2} * (-\log \sigma^2 + \mu^2 + \sigma^2 - 1) \end{aligned}$$

μ 表示 x 在高斯分布下的平均值。相似的，二阶矩为 (the second order moment) :

$$E[x^2] = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2)x^2 dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = E[x^2] - E[x]^2 = \sigma^2$$

证明上式: 令 $z = x - \mu$, 则

$$\begin{aligned} E[x^2] &= \int_{-\infty}^{\infty} N(x|\mu, \sigma^2)x^2 dx \\ &= \int_{-\infty}^{\infty} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}z^2\right\} (z + \mu)^2 dz \\ &= \int_{-\infty}^{\infty} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}z^2\right\} z^2 dz + 2\mu \int_{-\infty}^{\infty} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}z^2\right\} z dz \\ &\quad + \mu^2 \int_{-\infty}^{\infty} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}z^2\right\} dz \end{aligned}$$

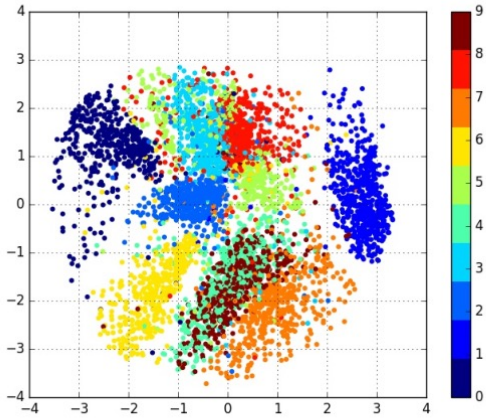
式中第二项为奇函数, 且在 $(-\infty, \infty)$ 内积分, 所以为0; 第三项积分部分为1, 所以第三项为 μ^2 , 因此现在只要求出第一项即可。在第一项中 $z = x - \mu$, 所以第一项等于 $E[z^2] = E[(x - \mu)^2]$:

$$E[(x - \mu)^2] = \left\{ \sum_{i=1}^N (x_i - \mu)^2 + \dots + (x_N - \mu)^2 \right\} / N = \sigma^2, \quad N \rightarrow \infty$$

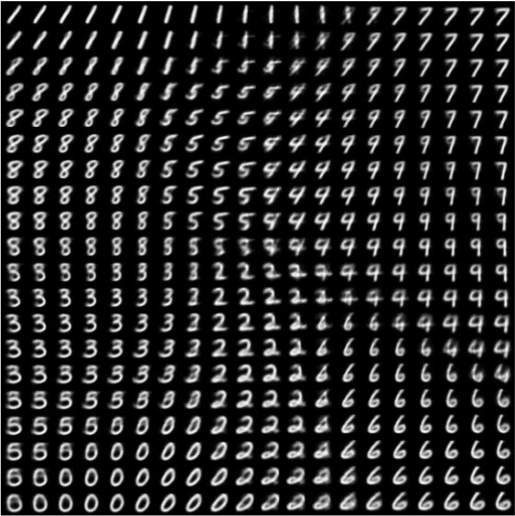
所以 $E[(x - \mu)^2] = \sigma^2$, 即

$$E[x^2] = \sigma^2 + 0 + \mu^2 = \sigma^2 + \mu^2$$

变分自编码器学习到的隐变量流形



(a) 训练集上所有样本在隐空间上的投影。



(b) 隐变量 z 在图像空间的投影。



生成对抗网络

Generative Adversarial Network (GAN)

显式密度模型和隐式密度模型

▶ 显式密度模型

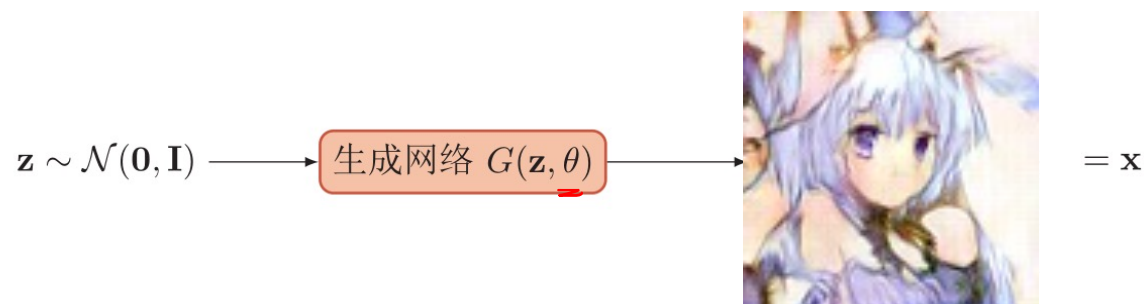
- ▶ 显式地构建出样本的密度函数 $p(x|\theta)$ ，并通过最大似然估计来求解参数；
- ▶ 变分自编码器、深度信念网络

▶ 隐式密度模型

- ▶ 不显式地估计出数据分布的密度函数
- ▶ 但能生成符合数据分布 $p_{\text{data}}(x)$ 的样本
- ▶ 无法用最大似然估计

生成网络

- ▶ 生成网络从隐空间 (latent space) 中随机采样作为输入，其输出结果需要尽量模仿训练集中的真实样本。



如何学习生成网络?

生成网络示例



Each dimension of input vector represents some characteristics.



Longer hair



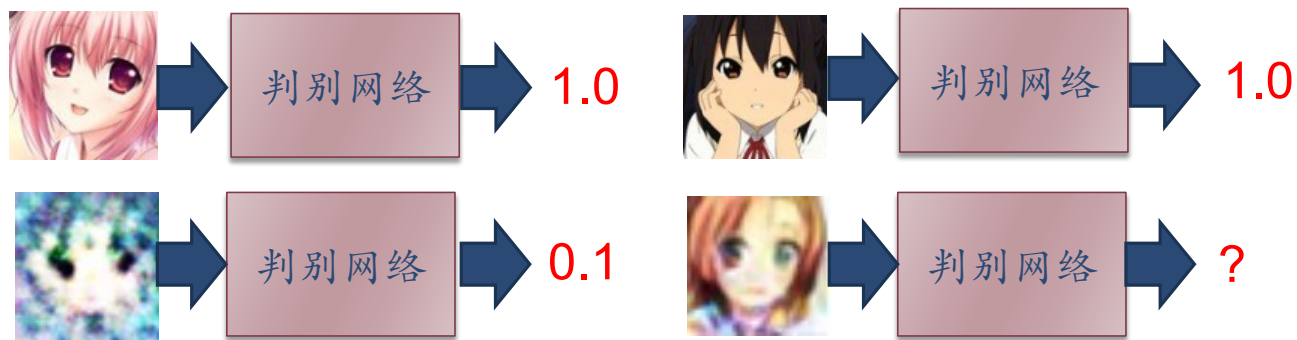
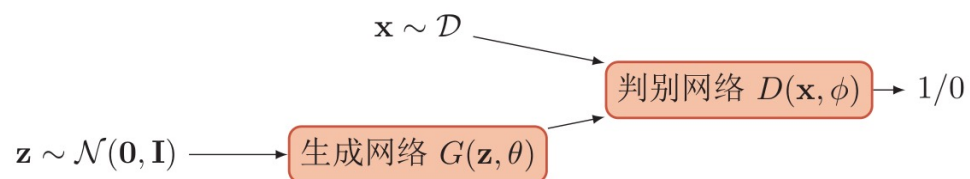
blue hair



Open mouth

判别网络

- ▶ 判别网络的输入则为真实样本或生成网络的输出，其目的是将生成网络的输出从真实样本中尽可能分辨出来。



MinMax Game

▶ 对抗训练

- ▶ 生成网络要尽可能地欺骗判别网络。
 - ▶ 判别网络将生成网络生成的样本与真实样本中尽可能区分出来。
- ▶ 两个网络相互对抗、不断调整参数，最终目的是使判别网络无法判断生成网络的输出结果是否真实。

对抗过程



MinMax Game

▶ 判别网络

$$\max_{\phi} \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[\log D(\mathbf{x}; \phi) \right] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[\log(1 - D(G(\mathbf{z}; \theta); \phi)) \right]$$

▶ 生成网络

$$\max_{\theta} \left(\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[\log D(G(\mathbf{z}; \theta); \phi) \right] \right)$$

▶ Minimax Game

$$\min_{\theta} \max_{\phi} \left(\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[\log D(\mathbf{x}; \phi) \right] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[\log(1 - D(G(\mathbf{z}; \theta); \phi)) \right] \right)$$

训练过程

算法 13.1 生成对抗网络的训练过程

输入: 训练集 \mathcal{D} , 对抗训练迭代次数 T , 每次判别网络的训练迭代次数 K , 小批量样本数量 M

1 随机初始化 θ, ϕ ;

2 **for** $t \leftarrow 1$ **to** T **do**

 // 训练判别网络 $D(\mathbf{x}; \phi)$

3 **for** $k \leftarrow 1$ **to** K **do**

 // 采集小批量训练样本

4 从训练集 \mathcal{D} 中采集 M 个样本 $\{\mathbf{x}^{(m)}\}, 1 \leq m \leq M$;

5 从分布 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ 中采集 M 个样本 $\{\mathbf{z}^{(m)}\}, 1 \leq m \leq M$;

6 使用随机梯度上升更新 ϕ , 梯度为

$$\frac{\partial}{\partial \phi} \left[\frac{1}{M} \sum_{m=1}^M \left(\log D(\mathbf{x}^{(m)}; \phi) + \log(1 - D(G(\mathbf{z}^{(m)}; \theta); \phi)) \right) \right];$$

7 **end**

 // 训练生成网络 $G(\mathbf{z}; \theta)$

8 从分布 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ 中采集 M 个样本 $\{\mathbf{z}^{(m)}\}, 1 \leq m \leq M$;

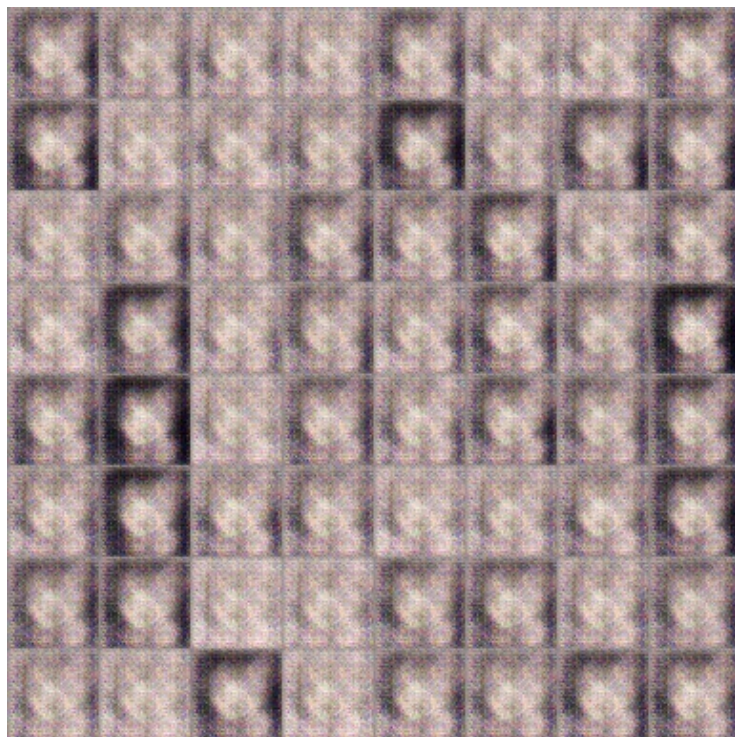
9 使用随机梯度上升更新 θ , 梯度为

$$\frac{\partial}{\partial \theta} \left[\frac{1}{M} \sum_{m=1}^M D(G(\mathbf{z}^{(m)}; \theta), \phi) \right];$$

10 **end**

输出: 生成网络 $G(\mathbf{z}; \theta)$

Anime Face Generation



100 updates



1000 updates

Anime Face Generation



2000 updates

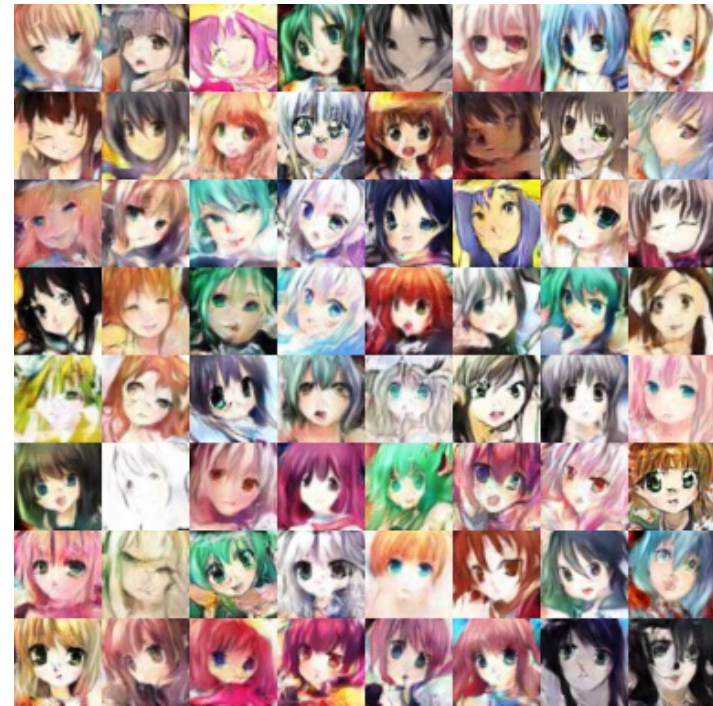


5000 updates

Anime Face Generation



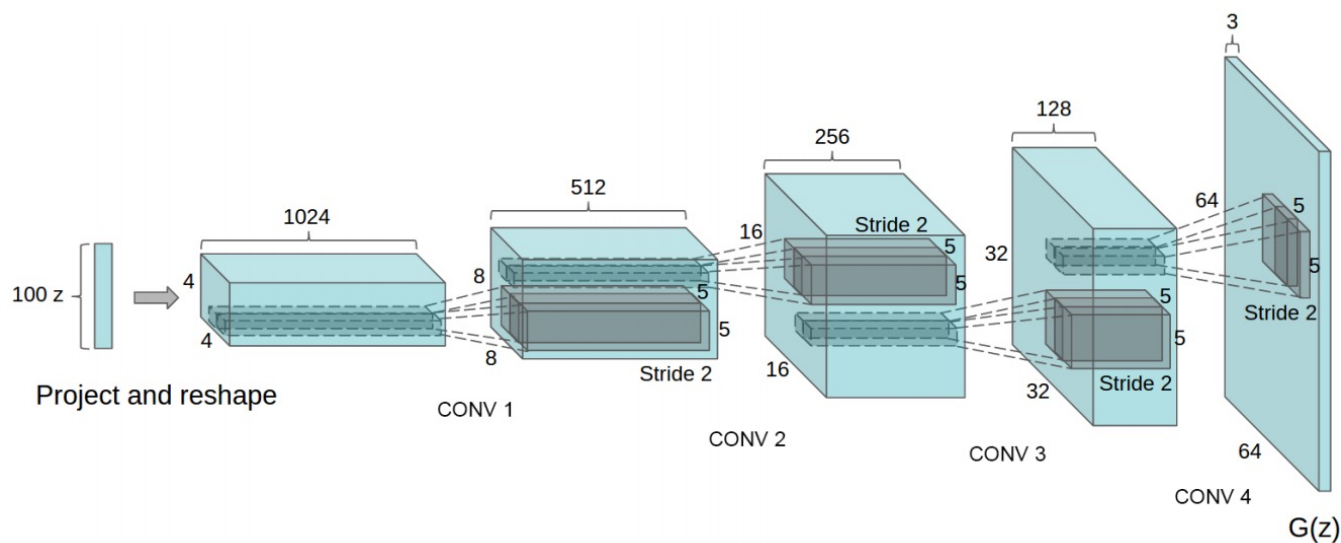
10,000
updates



50,000
updates

一个具体的模型：DCGANs

- ▶ 判别网络是一个传统的深度卷积网络，但使用了带步长的卷积来实现下采样操作，不用最大汇聚（pooling）操作。
- ▶ 生成网络使用一个特殊的深度卷积网络来实现使用微步卷积来生成 64×64 大小的图像。

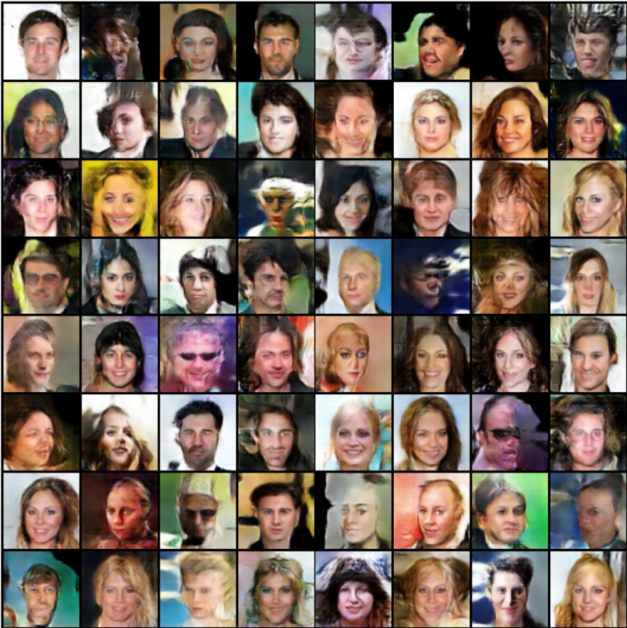


DCGANs

Real Images



Fake Images



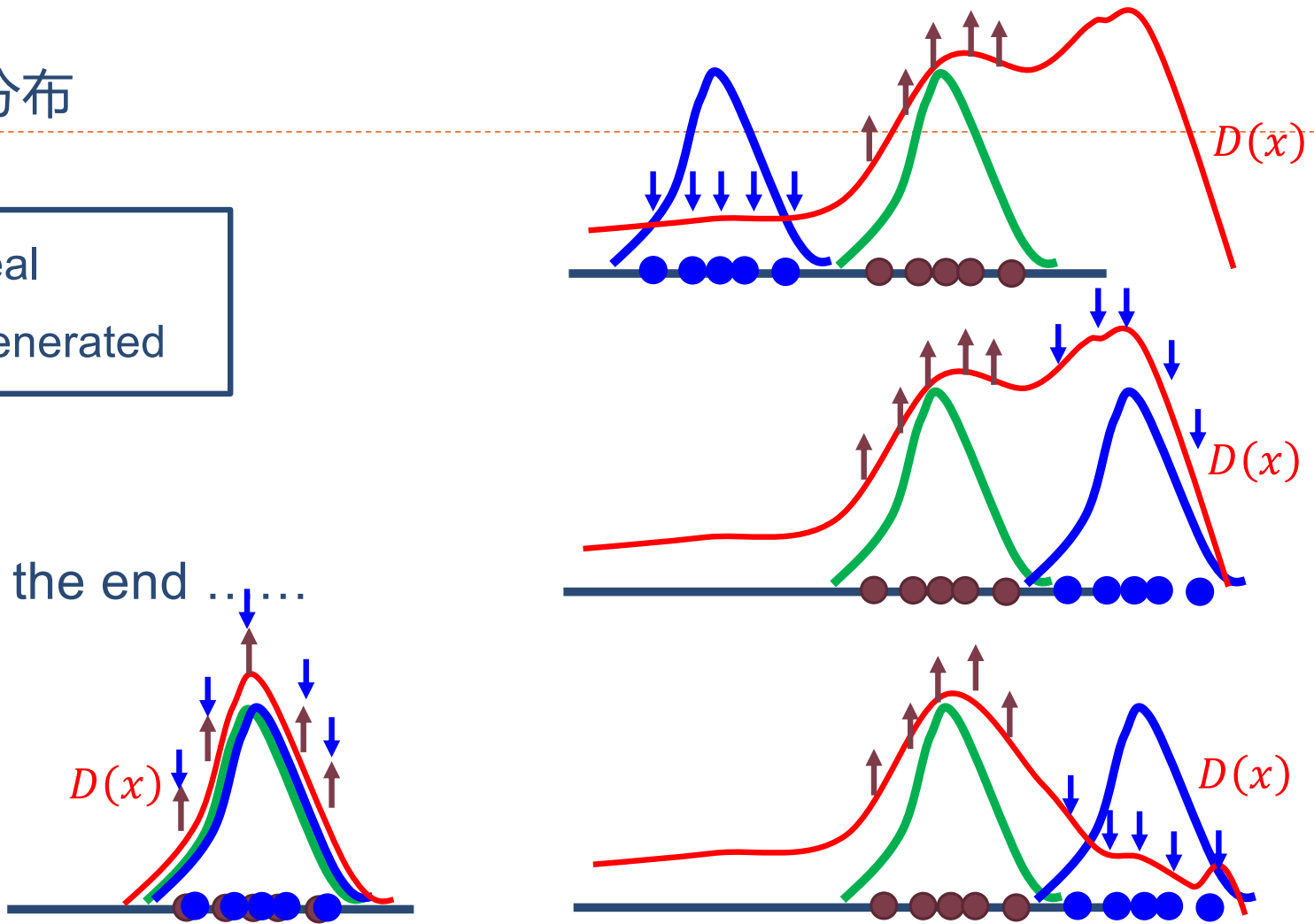


模型分析

数据分布

● real
● generated

In the end



模型分析

▶ 假设 $p_r(x)$ 和 $p_\theta(x)$ 已知，则最优的判别器为

$$\min_{\theta} \max_{\phi} \left(\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[\log D(\mathbf{x}; \phi) \right] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[\log(1 - D(G(\mathbf{z}; \theta); \phi)) \right] \right)$$

$$f(y) = a \log y + b \log(1 - y)$$

$$f'(y) = 0 \Rightarrow \frac{a}{y} - \frac{b}{1-y} = 0 \Rightarrow y = \frac{a}{a+b}$$

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_\theta(\mathbf{x})}$$

通过在其驻点进行二阶导求解可得：

$$f''\left(\frac{a}{a+b}\right) = -\frac{a}{\left(\frac{a}{a+b}\right)^2} - \frac{b}{1 - \left(\frac{a}{a+b}\right)^2} < 0$$

模型分析

▶ 假设 $p_r(x)$ 和 $p_\theta(x)$ 已知，则最优的判别器为

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_\theta(\mathbf{x})}$$

▶ 目标函数变为

$$\begin{aligned}\mathcal{L}(G|D^*) &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[\log D^*(\mathbf{x}) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log(1 - D^*(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[\log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_\theta(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x})}{p_r(\mathbf{x}) + p_\theta(\mathbf{x})} \right] \\ &= D_{\text{KL}}(p_r \| p_a) + D_{\text{KL}}(p_\theta \| p_a) - 2 \log 2 \\ &= 2D_{\text{JS}}(p_r \| p_\theta) - 2 \log 2,\end{aligned}$$

$$D_{\text{JS}}(p \| q) = \frac{1}{2} D_{\text{KL}}(p \| m) + \frac{1}{2} D_{\text{KL}}(q \| m)$$

$$m = \frac{1}{2}(p + q)$$

不稳定性：生成网络的梯度消失

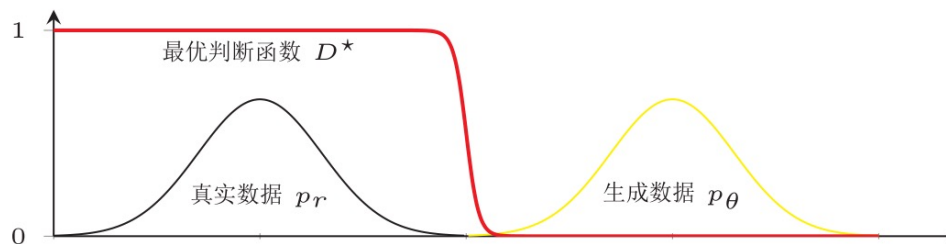
$$\min_{\theta} \max_{\phi} \left(\mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} \left[\log D(\mathbf{x}, \phi) \right] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[\log(1 - D(G(\mathbf{z}, \theta), \phi)) \right] \right)$$

在生成对抗网络中，当判断网络为最优时，生成网络的优化目标是 minimize 真实分布 $p_r(x)$ 和模型分布 $p_{\theta}(x)$ 之间的 JS 散度。

当两个分布相同时，JS 散度为 0，最优生成网络对应的损失为 $-2\log 2$ 。

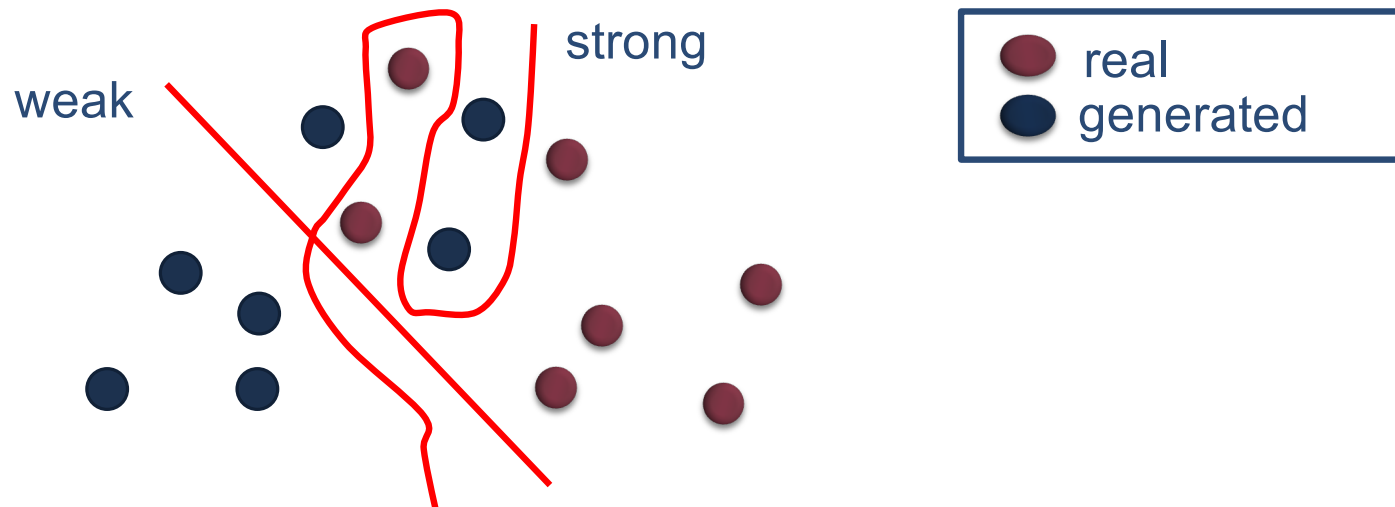
使用 JS 散度来训练生成对抗网络的一个问题是当两个分布没有重叠时，它们之间的 JS 散度恒等于常数 $\log 2$ 。对生成网络来说，目标函数关于参数的梯度为 0。

$$\frac{\partial \mathcal{L}(G|D^*)}{\partial \theta} = 0.$$



改进

- ▶ 弱化判别器
- ▶ 使用更好的损失函数



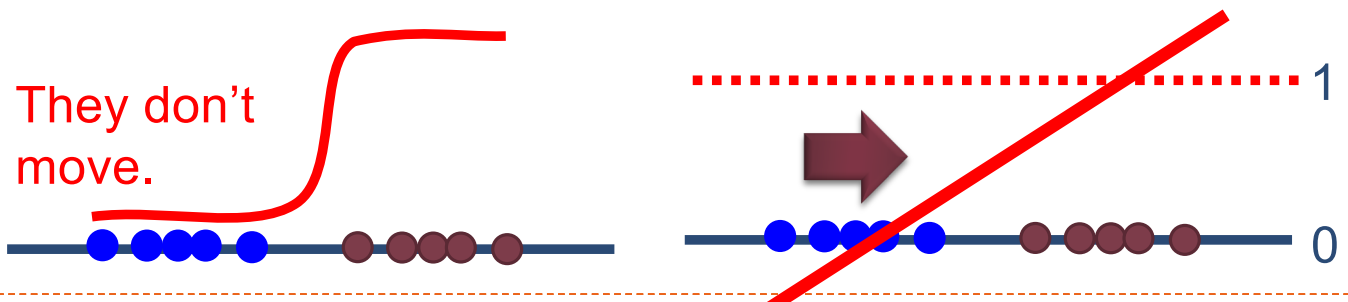
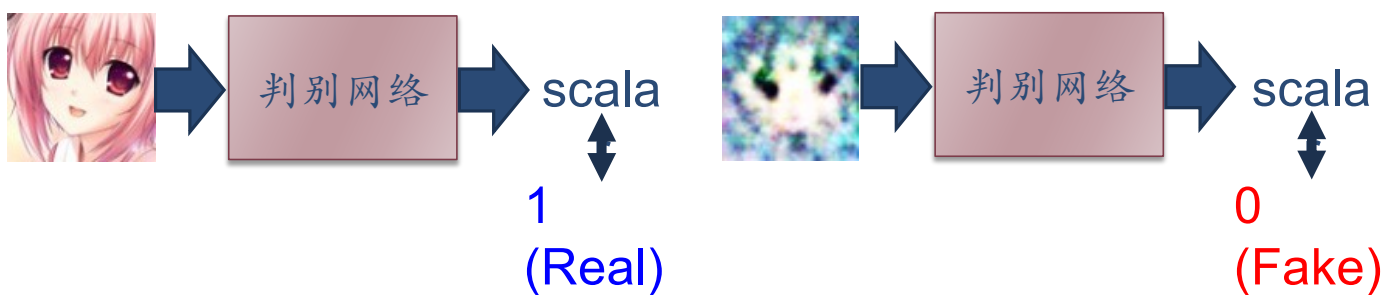
Least Square GAN (LSGAN)

● real
● generated

► Replace sigmoid with linear (replace classification with regression)

$$\min_D V_{LSGAN}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}) - b)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [(D(G(\mathbf{z})) - a)^2]$$

$$\min_G V_{LSGAN}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [(D(G(\mathbf{z})) - c)^2],$$



Wasserstein距离

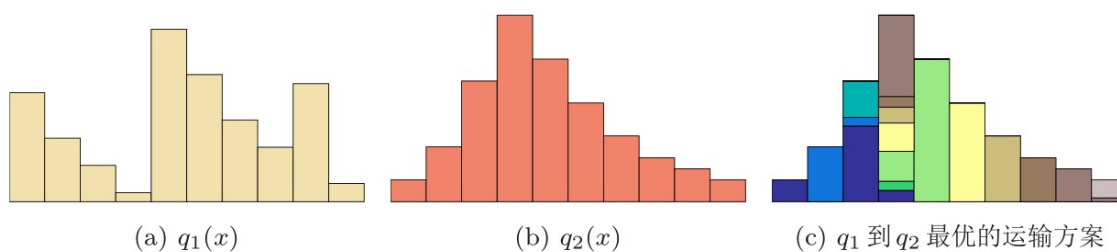
- ▶ Wasserstein距离用于衡量两个分布之间的距离。

$$W_p(q_1, q_2) = \left(\inf_{\gamma(x,y) \in \Gamma(q_1, q_2)} \mathbb{E}_{(x,y) \sim \gamma(x,y)} [d(x,y)^p] \right)^{\frac{1}{p}}$$

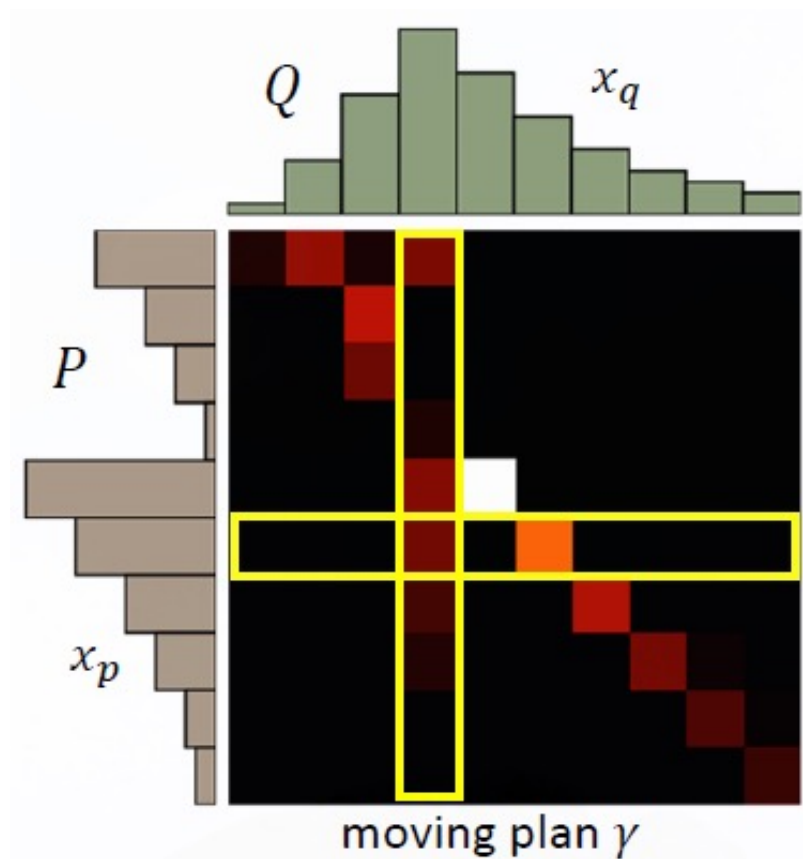
- ▶ 其中 $\Gamma(q_1, q_2)$ 是边际分布为 q_1, q_2 的所有可能的联合分布集合， $d(x,y)$ 为 x 和 y 的距离，比如 l_p 距离等。
- ▶ Wasserstein距离相比JS散度的优势在于：即使两个分布没有重叠或者重叠非常少，Wasserstein距离仍然能反映两个分布的远近。

Wasserstein距离

- ▶ 如果将两个分布看作是两个土堆，联合分布 $\gamma(x,y)$ 看作是从土堆 q_1 的位置 x 到土堆 q_2 的位置 y 的搬运土的数量。Wasserstein距离可以理解为搬运土堆的最小工作量，也称为推土机距离（Earth-Mover's Distance, EMD）。



Wasserstein距离



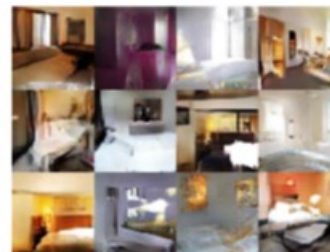
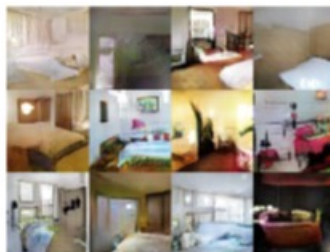
DCGAN

LSGAN

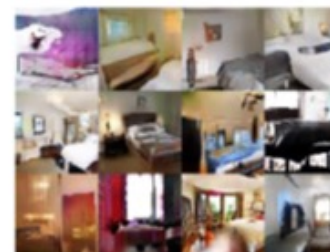
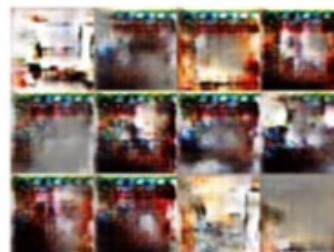
Original
WGAN

Improved
WGAN

G: CNN, D: CNN



G: CNN (no normalization), D: CNN (no normalization)



G: CNN (tanh), D: CNN(tanh)



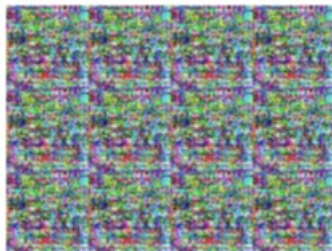
DCGAN

LSGAN

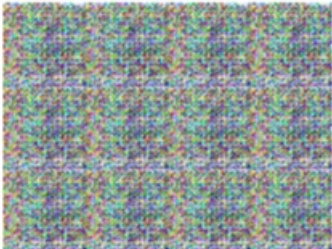
Original
WGAN

Improved
WGAN

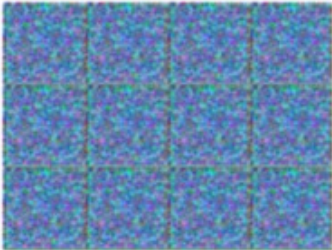
G: MLP, D: CNN



G: CNN (bad structure), D: CNN



G: 101 layer, D: 101 layer





GAN的扩展

条件生成

▶ 根据条件针对性的生成数据



“Girl with red hair and red eyes”

“Girl with yellow ribbon”



条件生成

Caption Generation

Given condition:



“A young girl is dancing”



Chat-bot

Given condition:

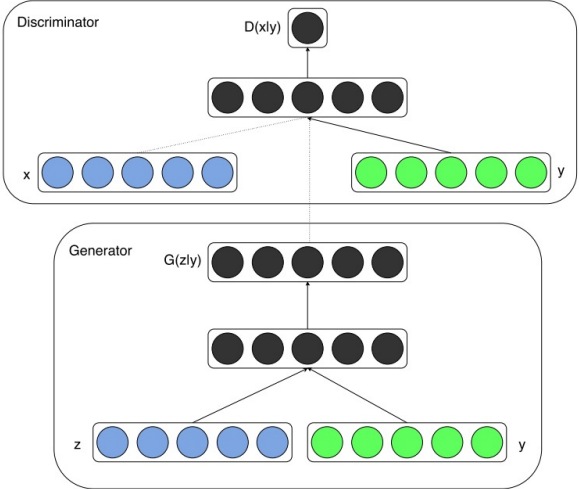
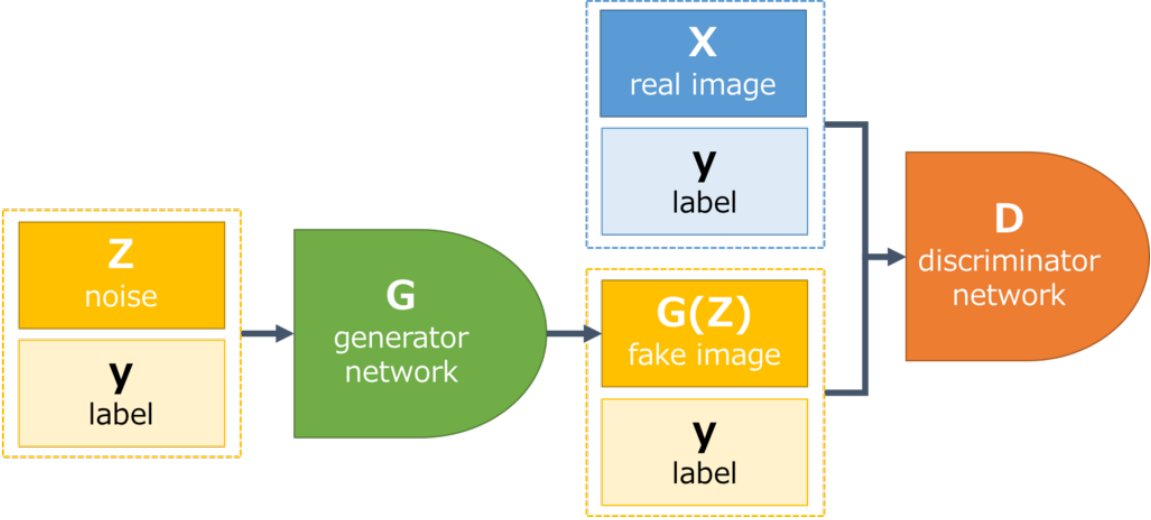


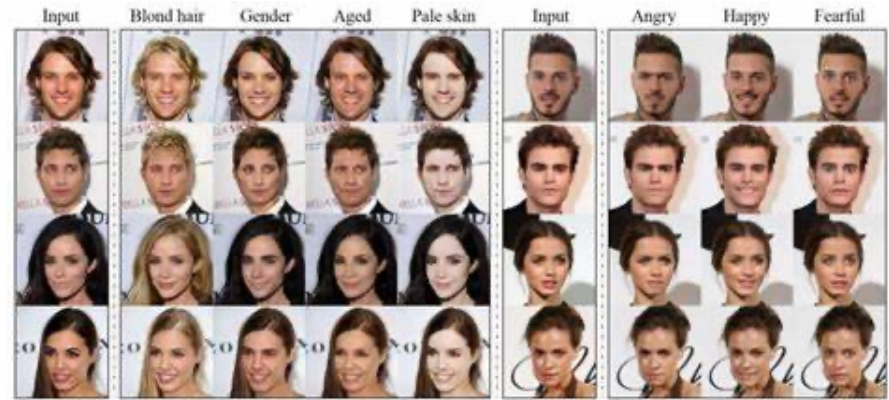
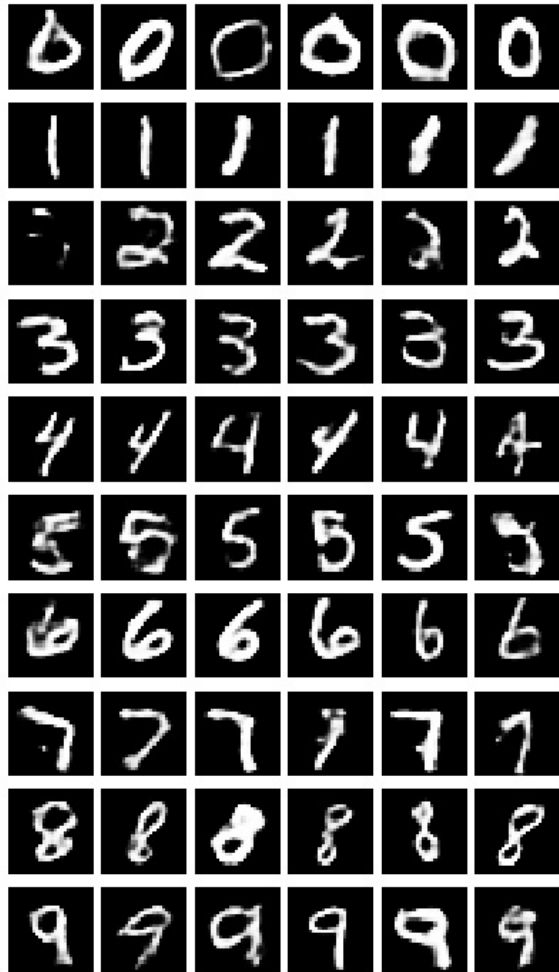
“Hello”

“Hello. Nice to see you.”



Conditional GAN







(a) Rotation

(b) Width

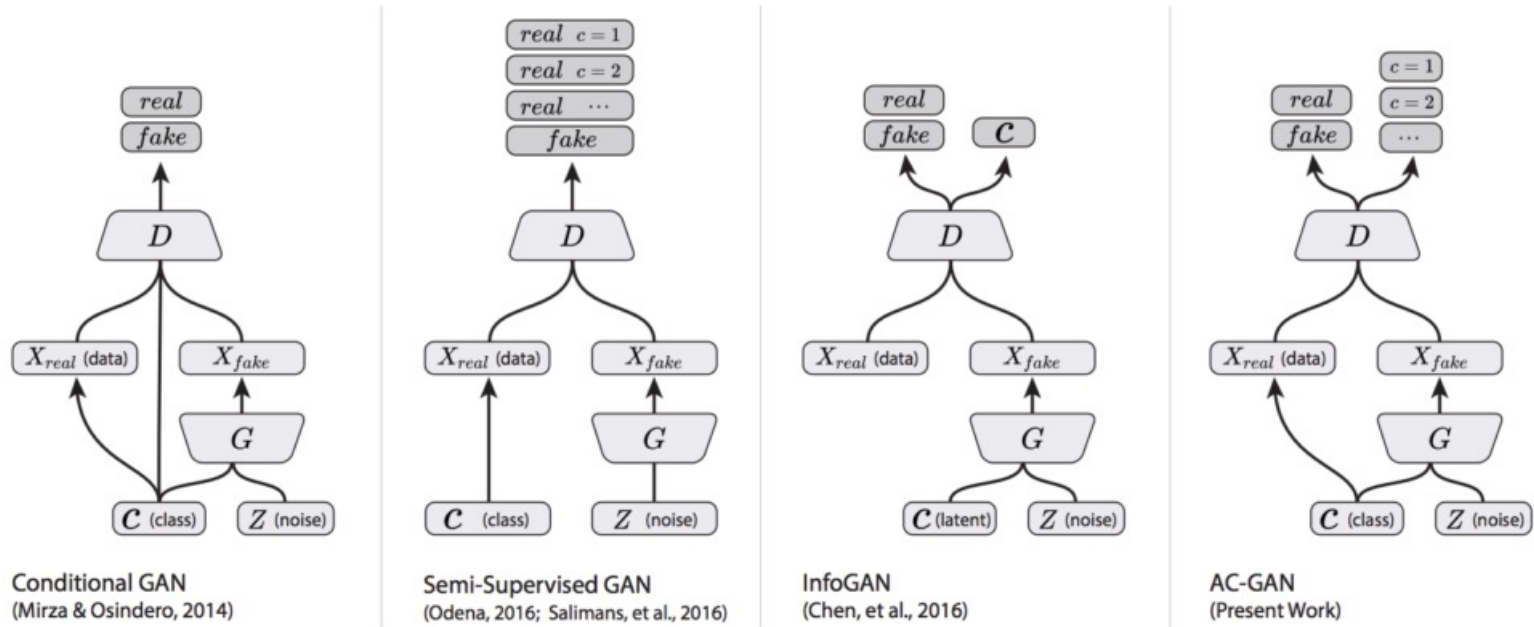


(c) Lighting

(d) Wide or Narrow

<https://arxiv.org/abs/1606.03657>

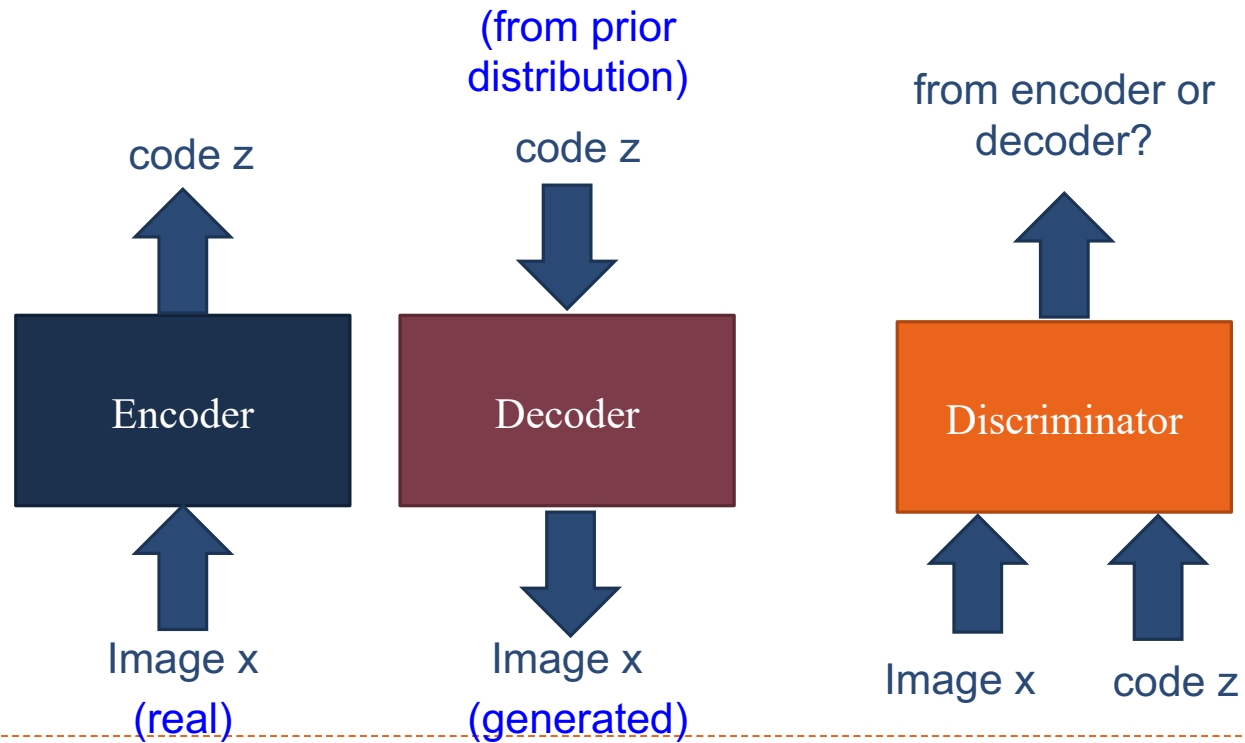
AC-GAN



BiGAN

Jeff Donahue, Philipp Krähenbühl, Trevor Darrell,
“Adversarial Feature Learning”, ICLR, 2017

Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier
Mastropietro, Alex Lamb, Martin Arjovsky, Aaron
Courville, “Adversarially Learned Inference”, ICLR, 2017





谢谢
<https://nndl.github.io/>