

无监督深度学习 ——生成模型 (1)

主讲：王皓 副研究员| 硕士导师

邮箱：wanghao3@ustc.edu.cn

主页：http://staff.ustc.edu.cn/~wanghao3



为什么无监督学习

- 有监督学习依赖于大量的**标记数据**, 但获取数据标签开销巨大
- 若没有足够的有标签数据, 应该怎么办呢?
- 可以利用**无监督**深度学习方法, 利用**大量未标记数据**



本章内容

➤ 自编码器 (AutoEncoder)

- PCA
- 降噪自编码器 (Denoising AutoEncoder)
- 稀疏自编码器 (Sparse AutoEncoder)

➤ 变分自编码器 (VAE)

➤ 生成对抗网络 (GAN)

$$S = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top = \frac{1}{n-1} \tilde{X} \tilde{X}^\top,$$

$$\begin{aligned} \max_{G \in \mathbb{R}^{d \times K}} \quad & \frac{1}{n-1} \sum_{i=1}^n \|GG^\top \mathbf{x}_i - GG^\top \bar{\mathbf{x}}\|^2, \\ \text{s.t.} \quad & G^\top G = I. \end{aligned}$$

Notice that

$$\begin{aligned} \frac{1}{n-1} \sum_{i=1}^n \|GG^\top \mathbf{x}_i - GG^\top \bar{\mathbf{x}}\|^2 &= \frac{1}{n-1} \sum_{i=1}^n \langle GG^\top (\mathbf{x}_i - \bar{\mathbf{x}}), GG^\top (\mathbf{x}_i - \bar{\mathbf{x}}) \rangle \\ &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^\top GG^\top GG^\top (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^\top GG^\top (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \frac{1}{n-1} \sum_{i=1}^n \text{tr} \left((\mathbf{x}_i - \bar{\mathbf{x}})^\top GG^\top (\mathbf{x}_i - \bar{\mathbf{x}}) \right) \\ &= \frac{1}{n-1} \sum_{i=1}^n \text{tr} \left(G^\top (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top G \right) \\ &= \text{tr} \left(G^\top \left(\frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top \right) G \right) \\ &= \text{tr} \left(G^\top S G \right). \end{aligned}$$

Denote

$$Q = U^\top G.$$

We can see that

$$Q^\top Q = I.$$

Thus, the problem in (11) reduces to

$$\begin{aligned} \max_{Q \in \mathbb{R}^{d \times K}} \quad & \text{tr}(Q^\top \Sigma_d^2 Q), \\ \text{s.t.} \quad & Q^\top Q = I. \end{aligned}$$

AutoEncoder

➤在MNIST中，用 $28 \times 28 = 784$ 维矩阵表示每个数字

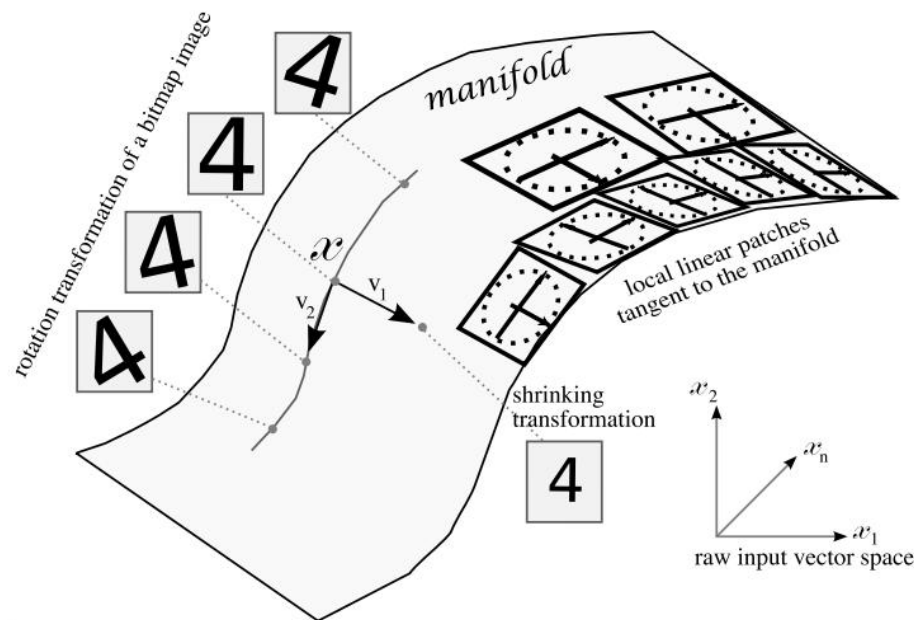


➤但不是每个784维的矩阵都表示数字

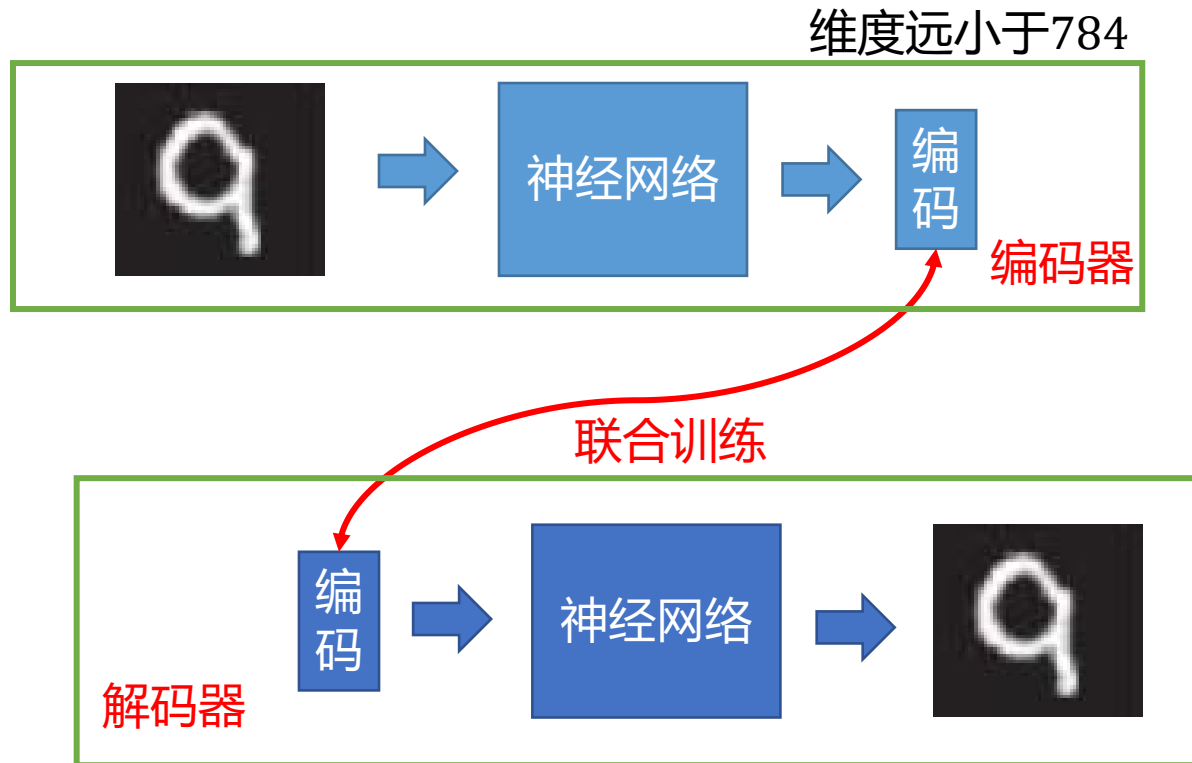


这些数字图像分布在784维空间内嵌的低维流形空间 (manifold) 上

因此可以用更紧凑的向量来表示这些数字图像



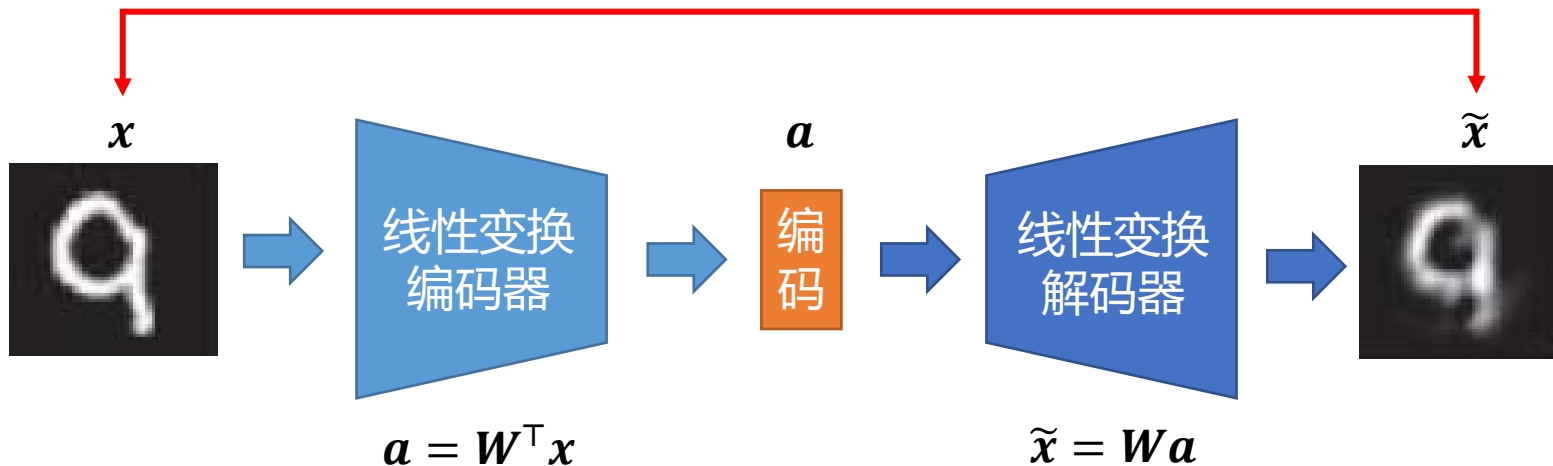
AutoEncoder



对比PCA&AutoEncoder

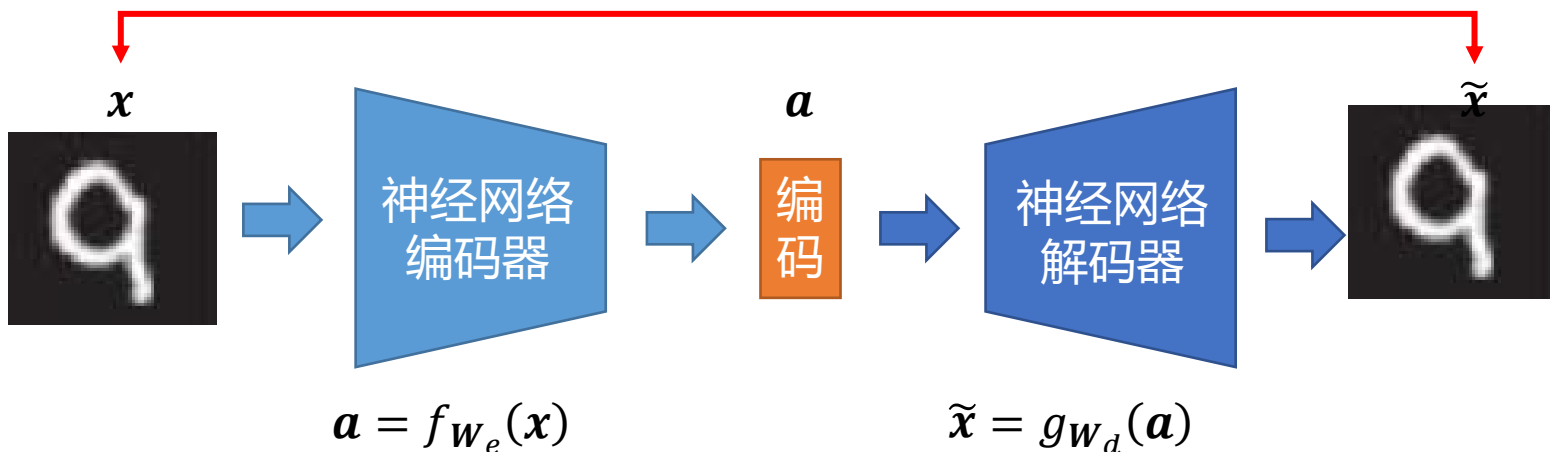
PCA

最小化输入和输出之间的差异 $(x - \tilde{x})^2$



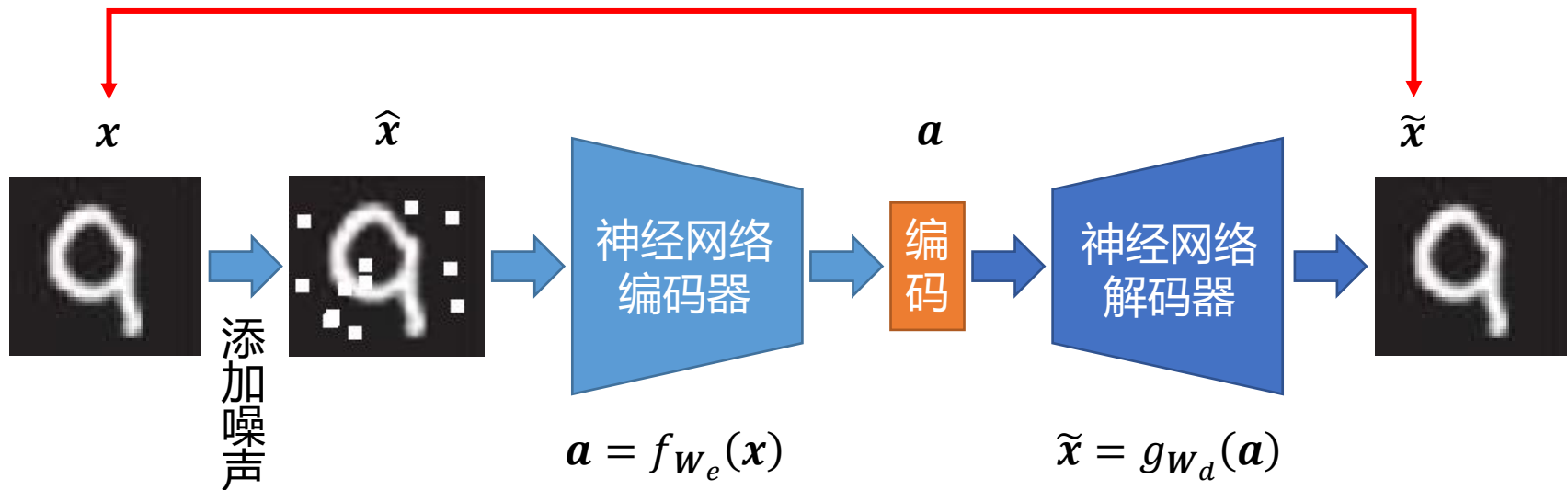
AutoEncoder

最小化输入和输出之间的差异 $(x - \tilde{x})^2$



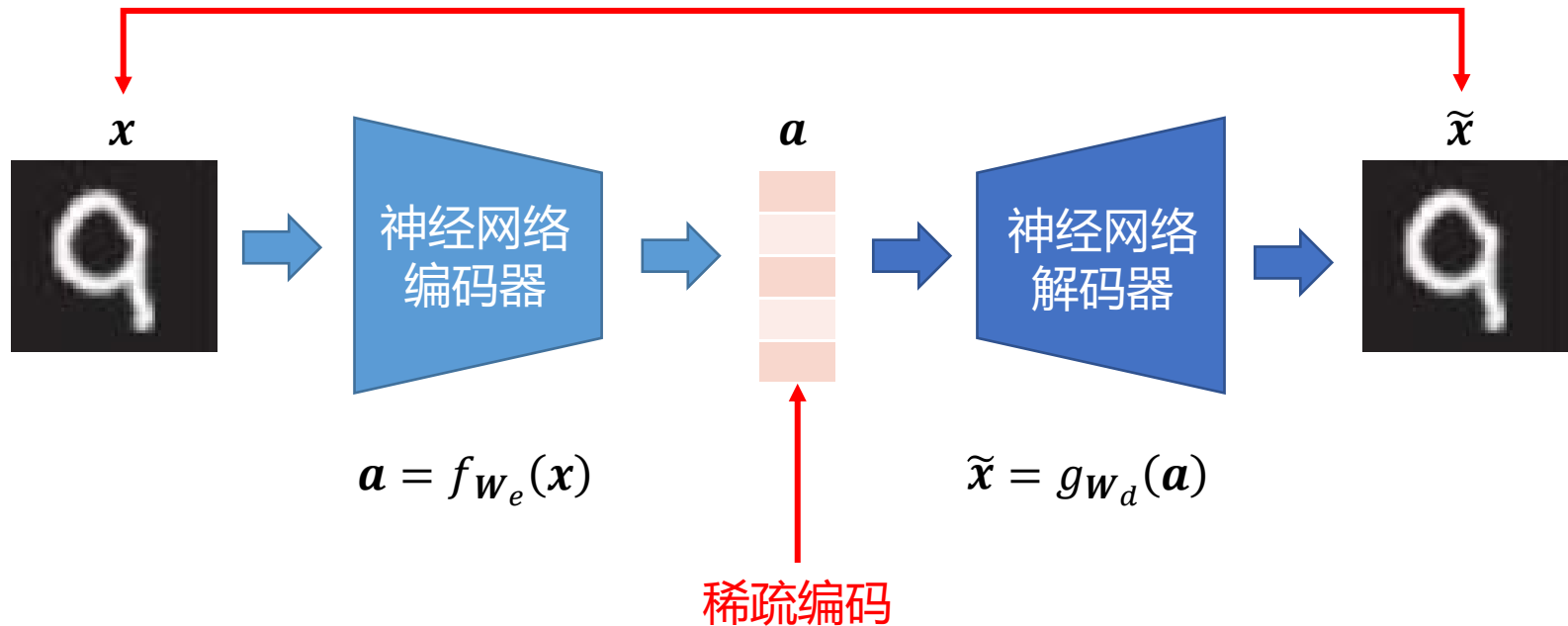
Denoising AutoEncoder

最小化输入和输出之间的差异 $(x - \tilde{x})^2$



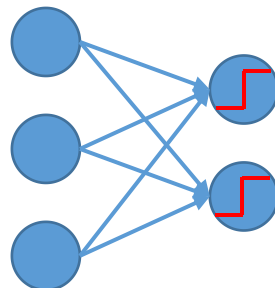
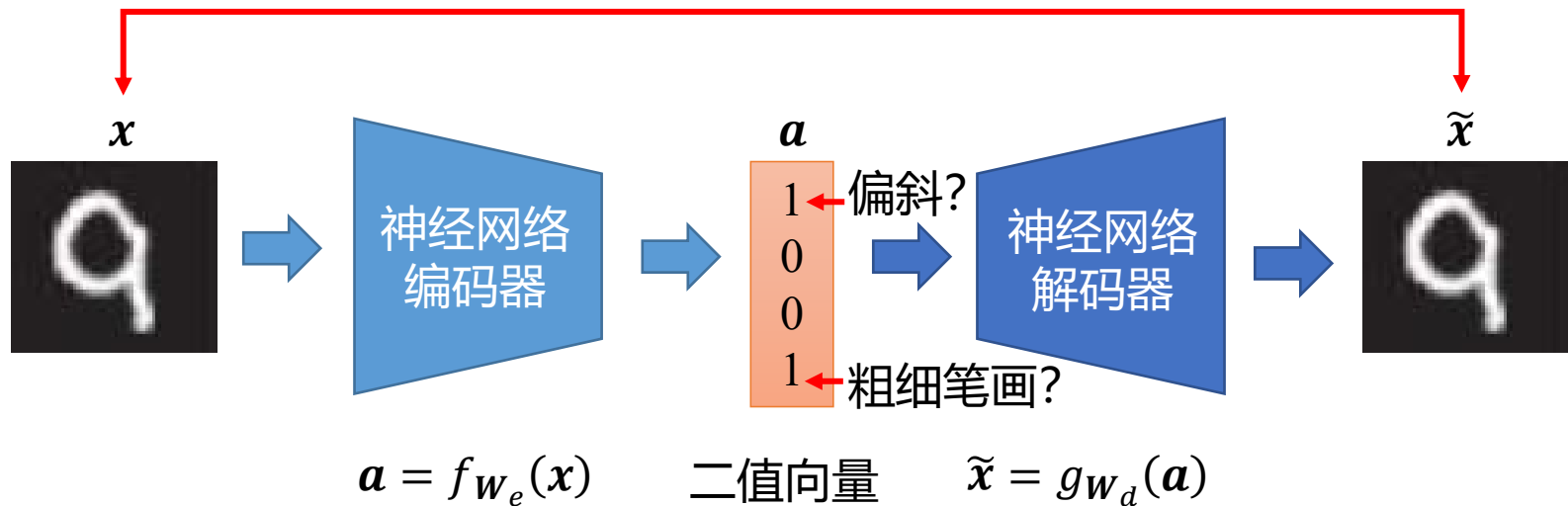
Sparse AutoEncoder

最小化输入和输出之间的差异 $(x - \tilde{x})^2$ + 稀疏约束 $\rho(a)$

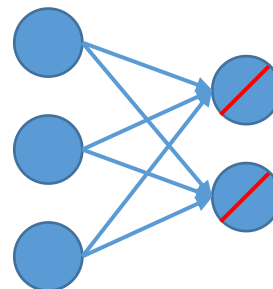


Discrete AutoEncoder

最小化输入和输出之间的差异 $(x - \tilde{x})^2$

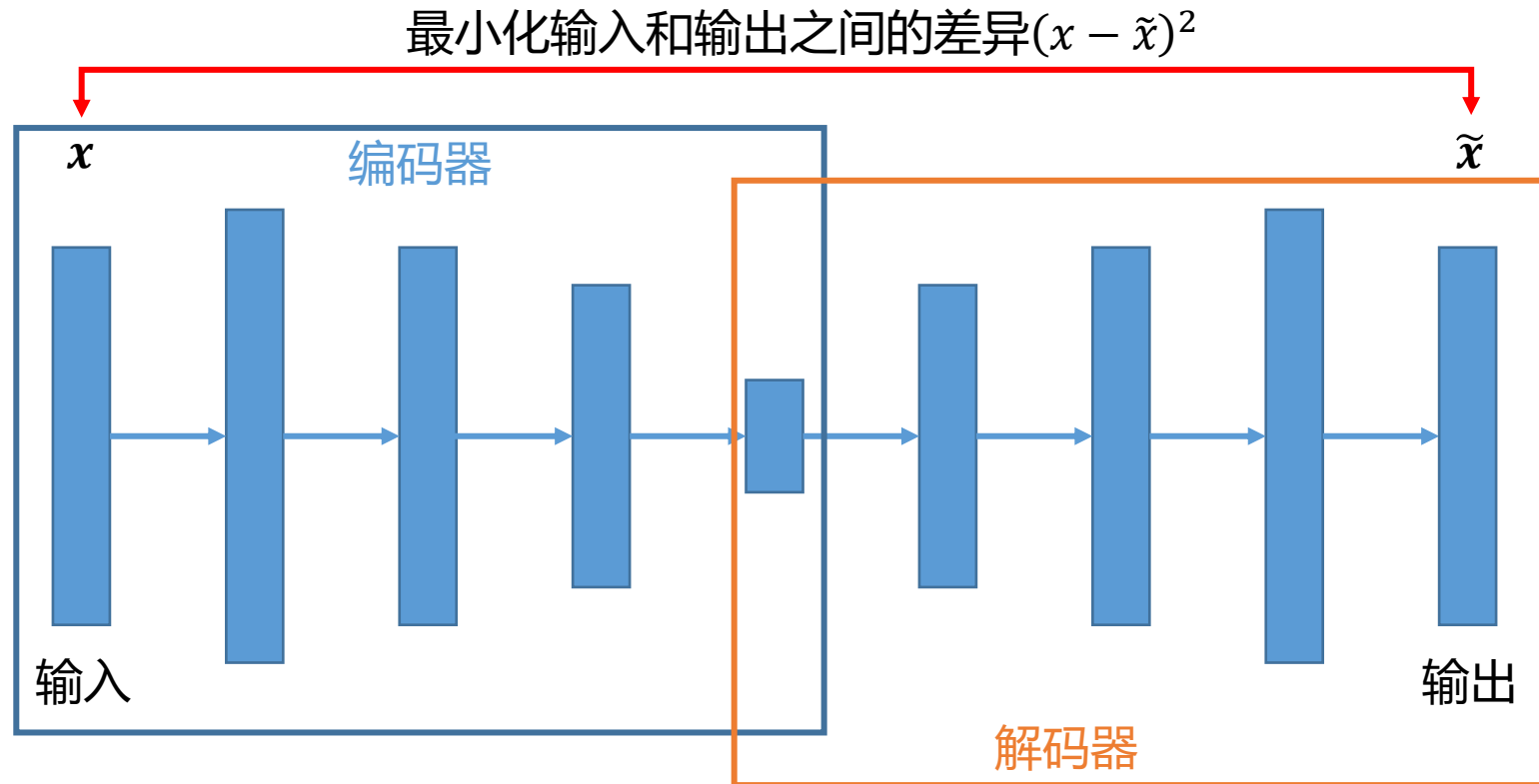


前向

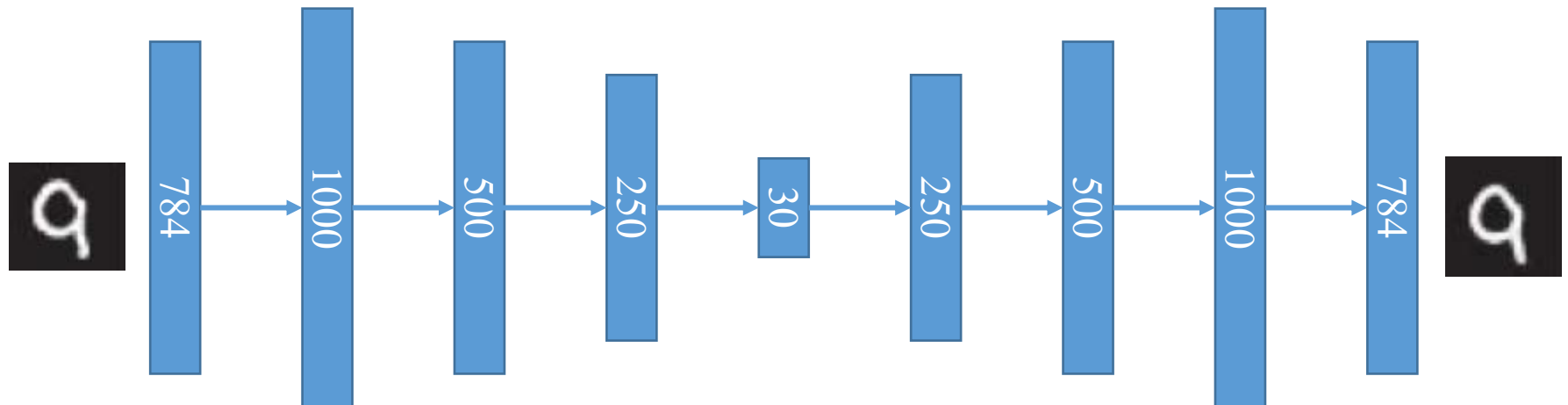


后向

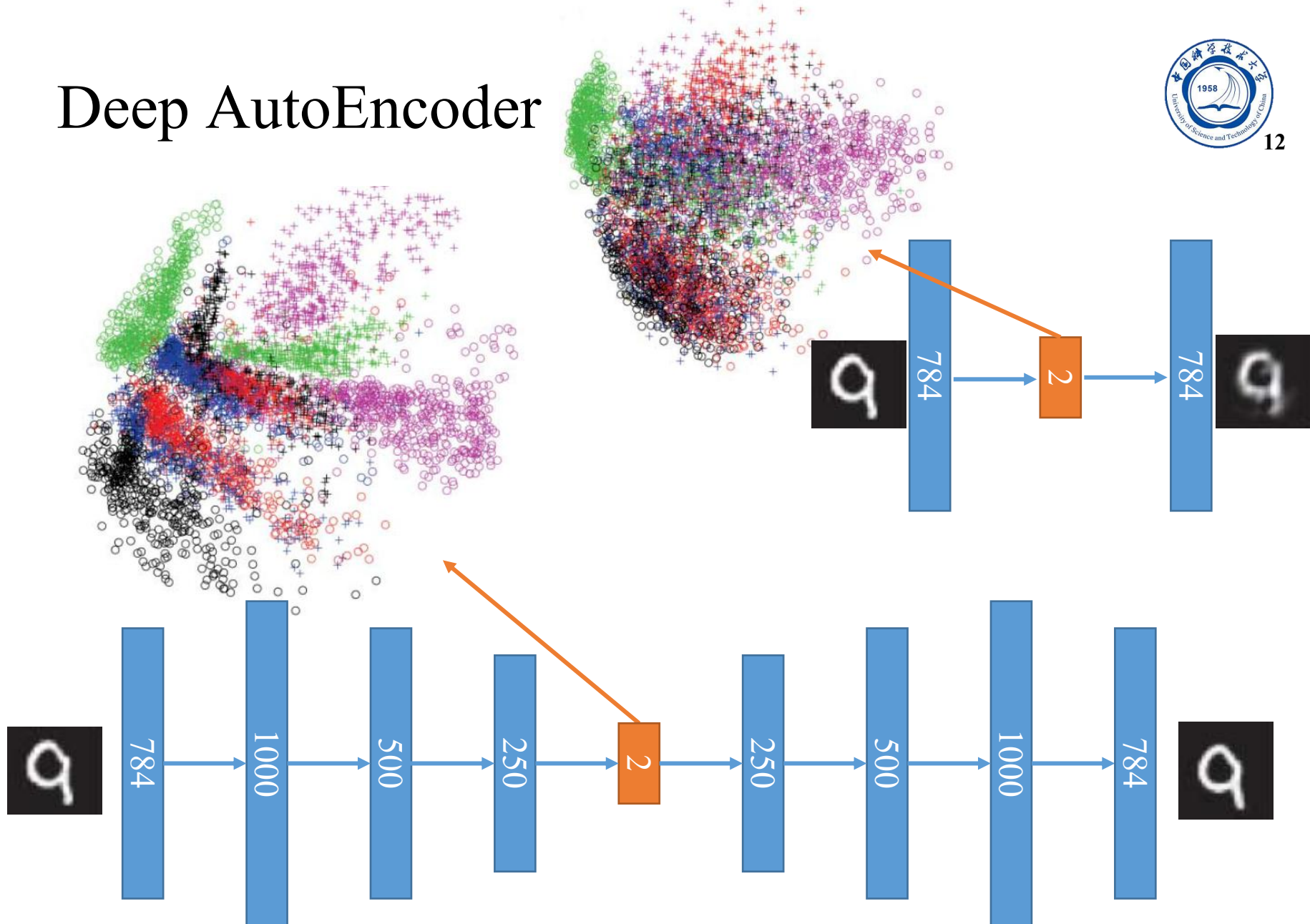
Deep AutoEncoder



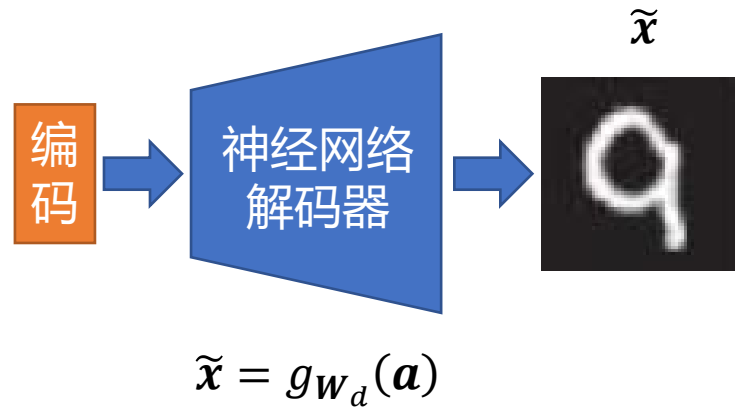
Deep AutoEncoder



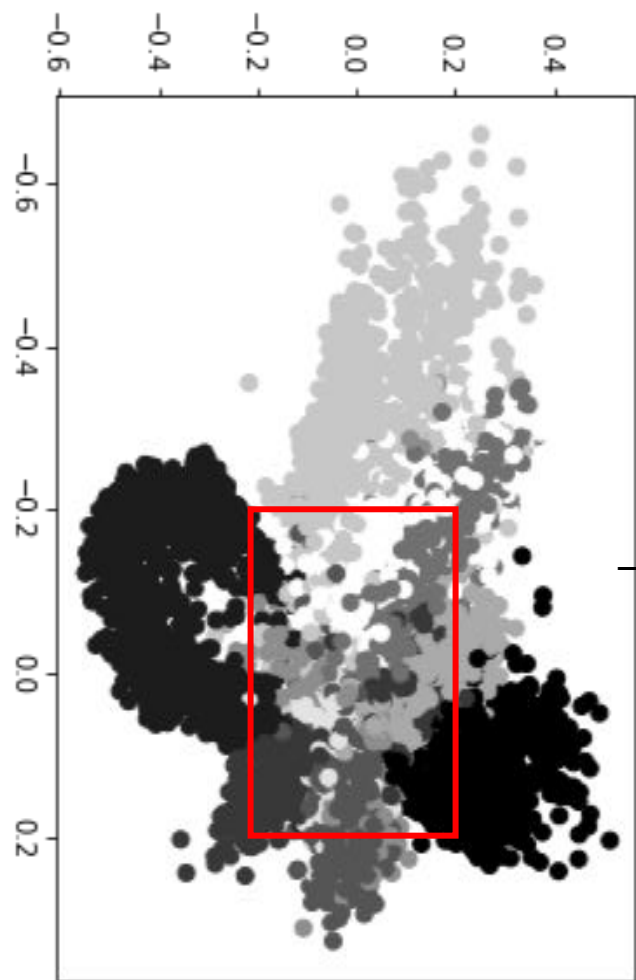
Deep AutoEncoder



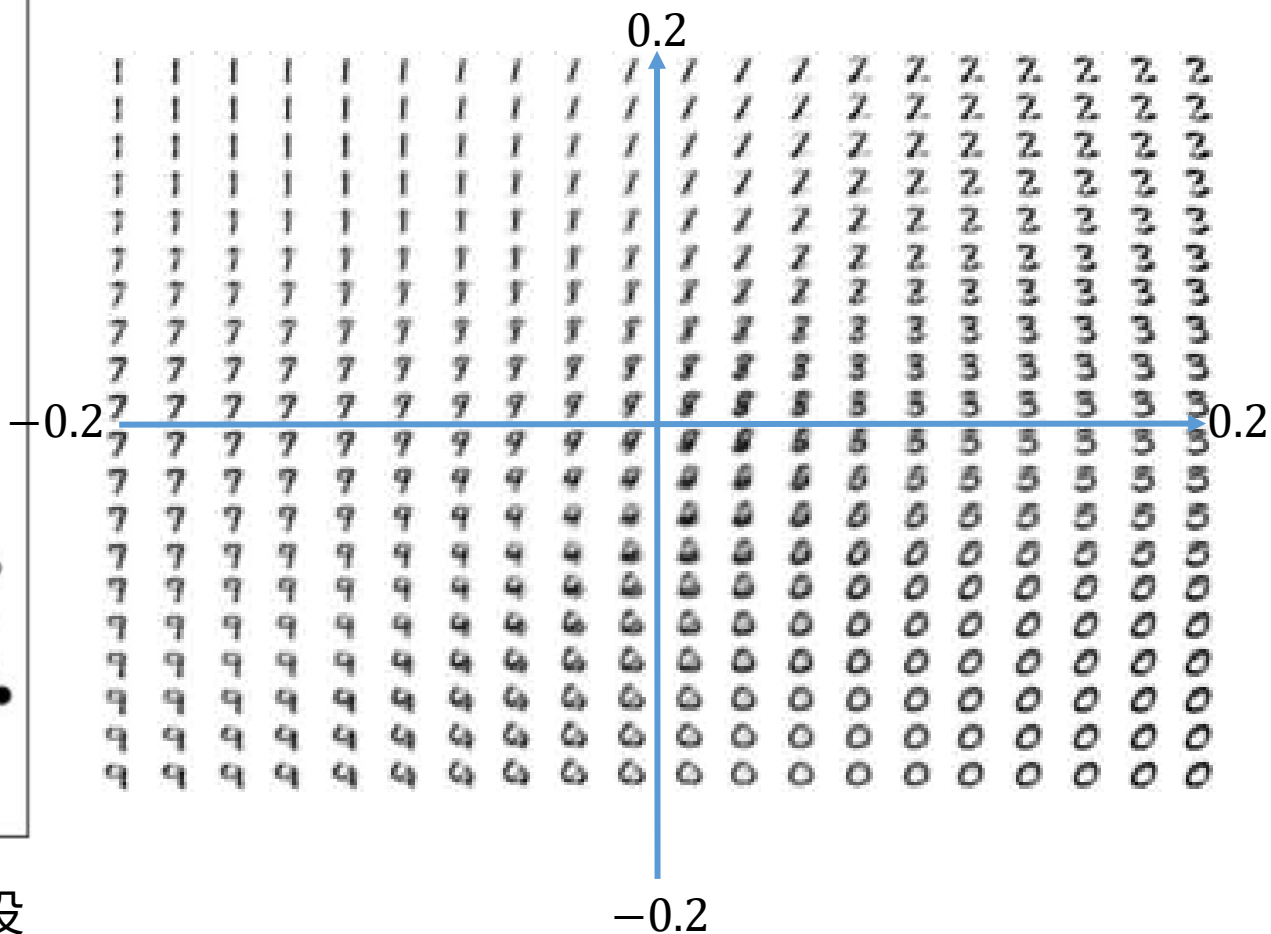
用AutoEncoder生成图片



用AutoEncoder生成图片

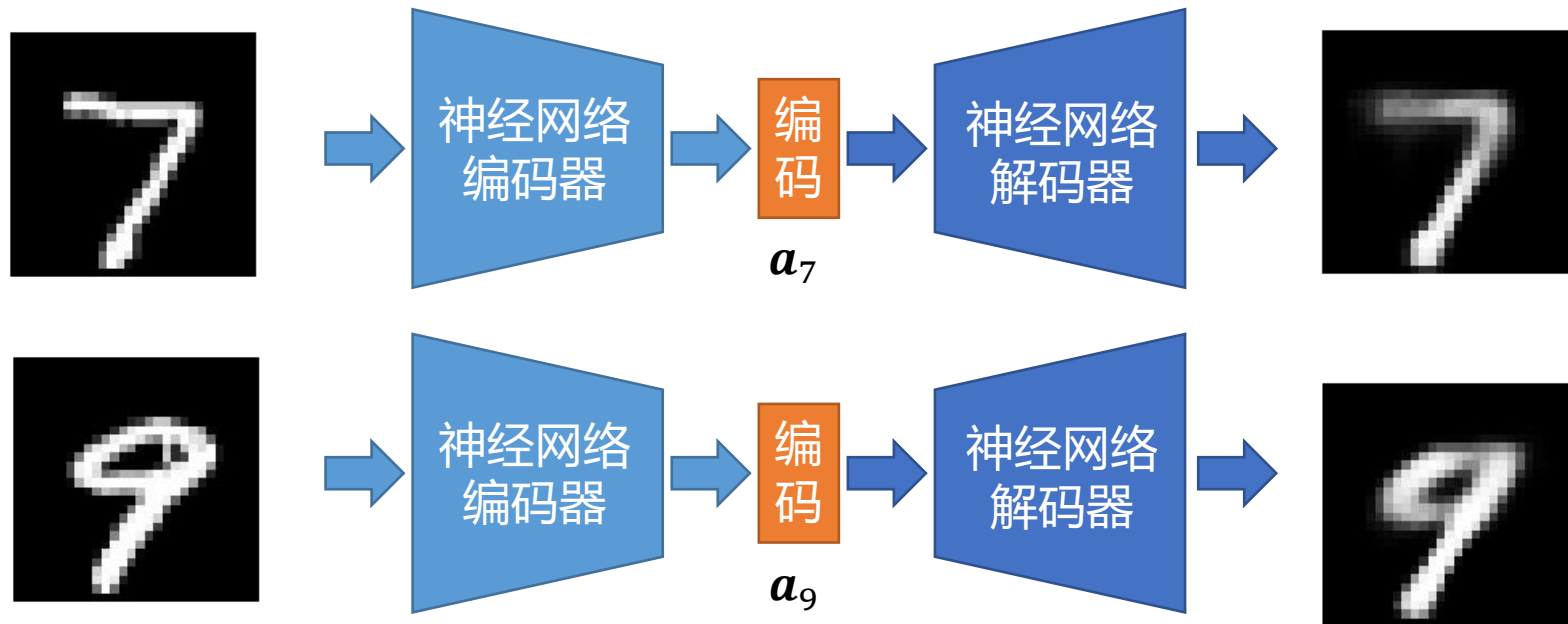


在红色区域内均匀采样400个坐标，作为解码器输入

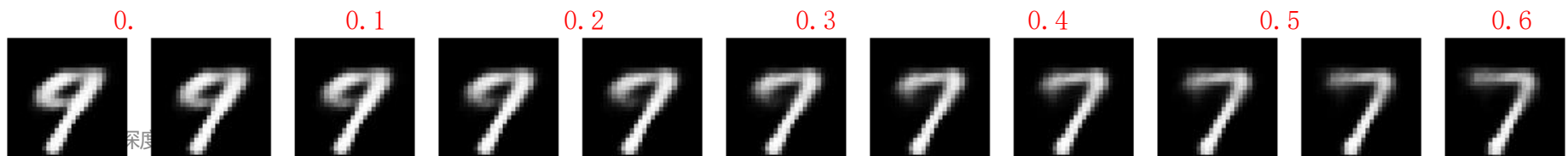
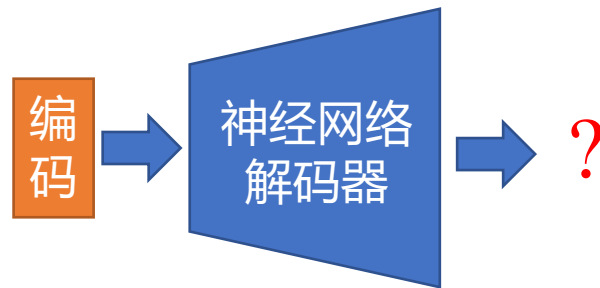


通过L2正则编码器将图片投影到2维空间

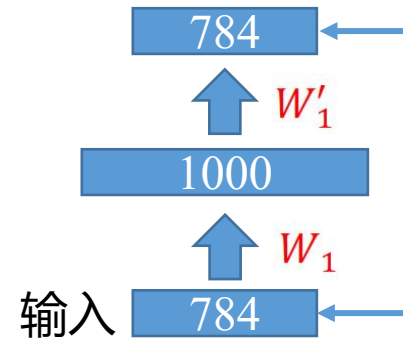
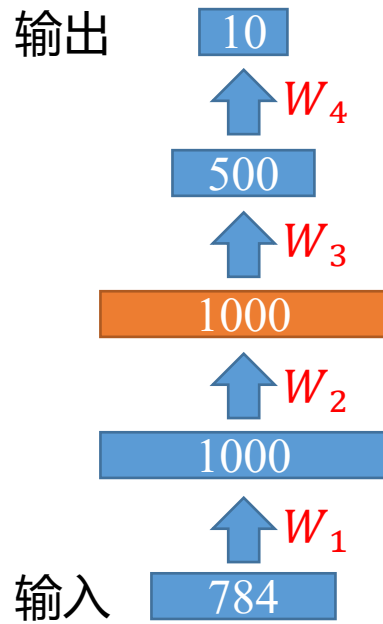
用AutoEncoder生成图片



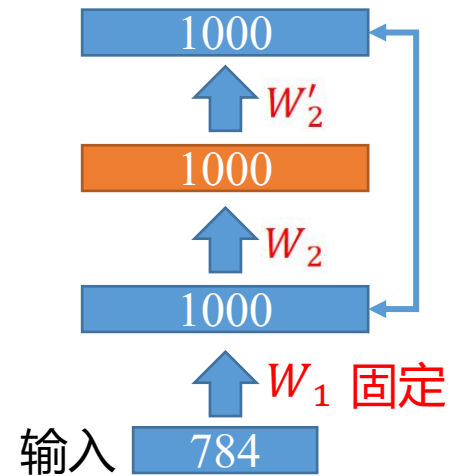
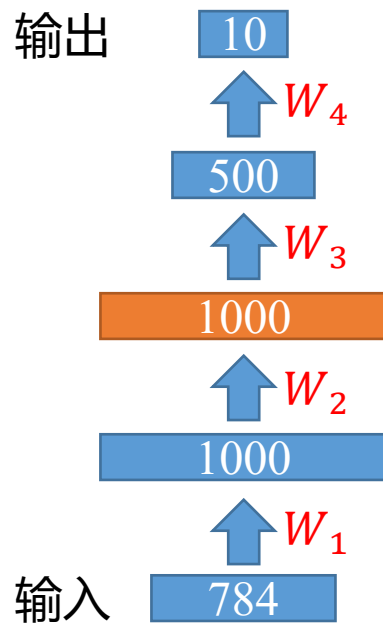
$$a = \alpha a_7 + (1 - \alpha) a_9$$



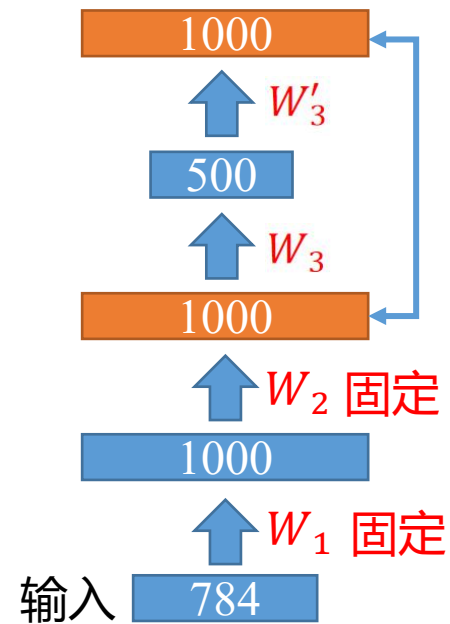
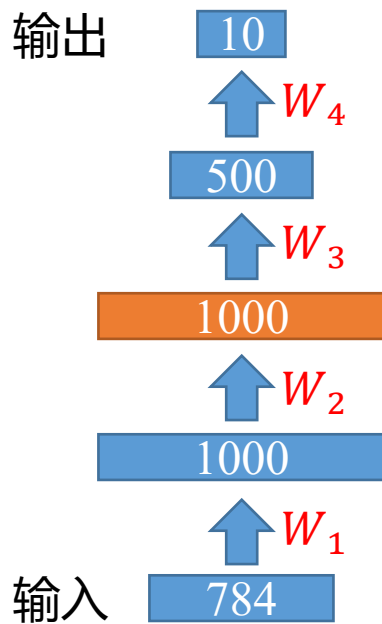
AutoEncoder用于预训练



AutoEncoder用于预训练

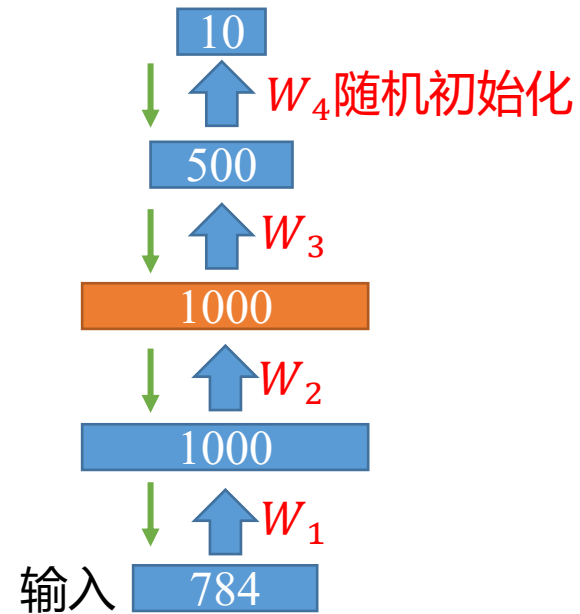
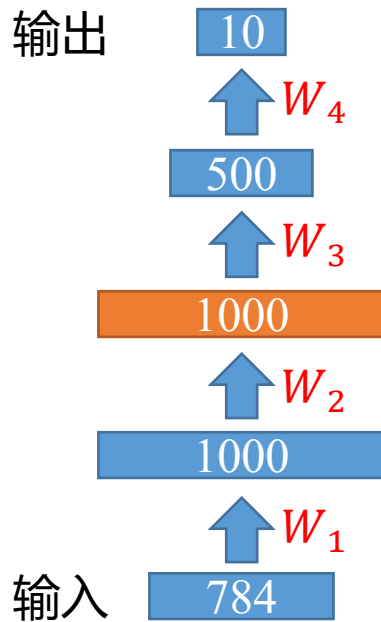


AutoEncoder用于预训练



AutoEncoder用于预训练

通过后向传播来做微调



AutoEncoder用于相似图片检索

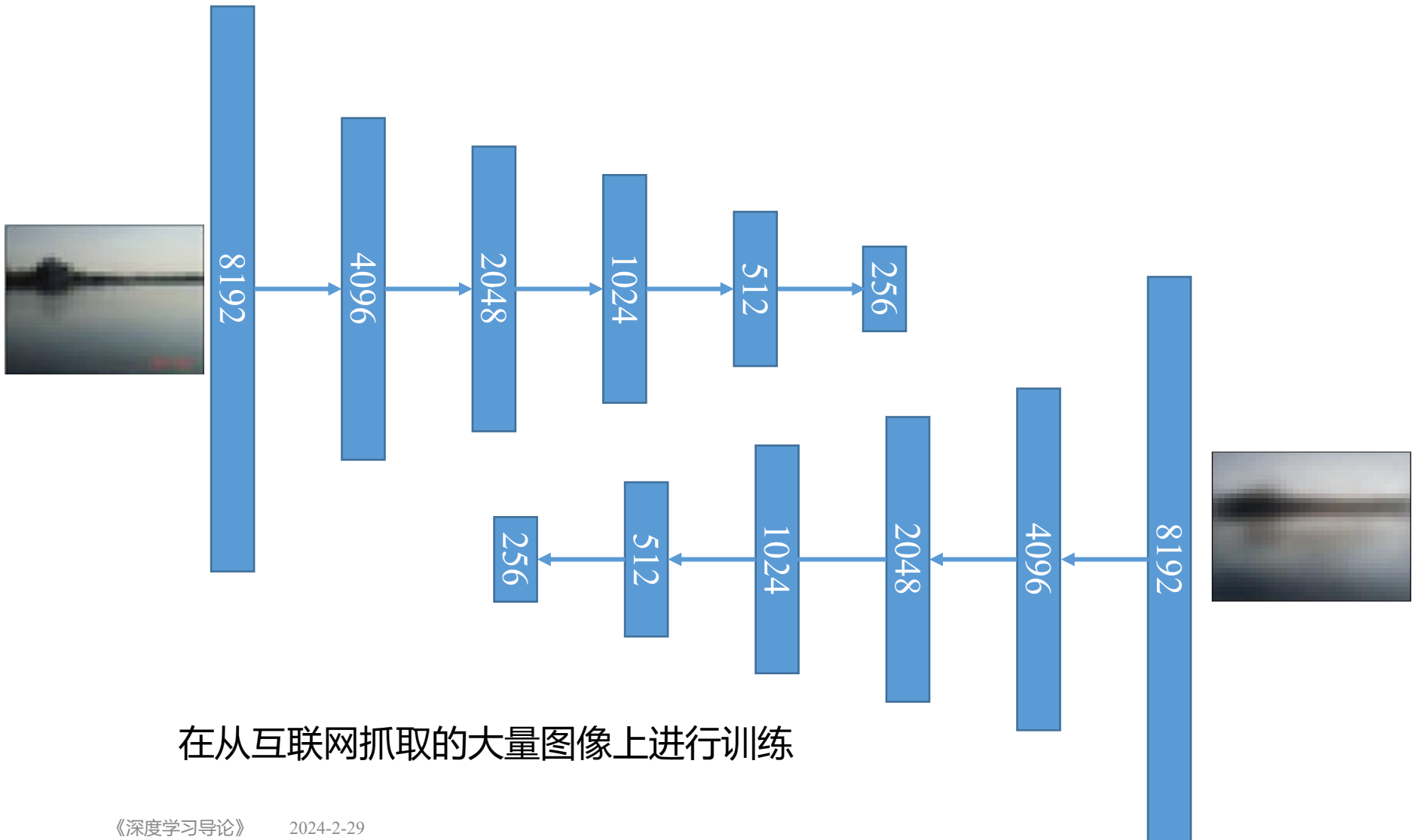


用迈克杰克逊的头
像进行检索

按照图片像素相似
度，最相似的图片



AutoEncoder用于相似图片检索



AutoEncoder用于相似图片检索

dist: 0.0



用迈克杰克逊的头
像进行检索

dist: 61



dist: 64



dist: 65



dist: 66



dist: 67



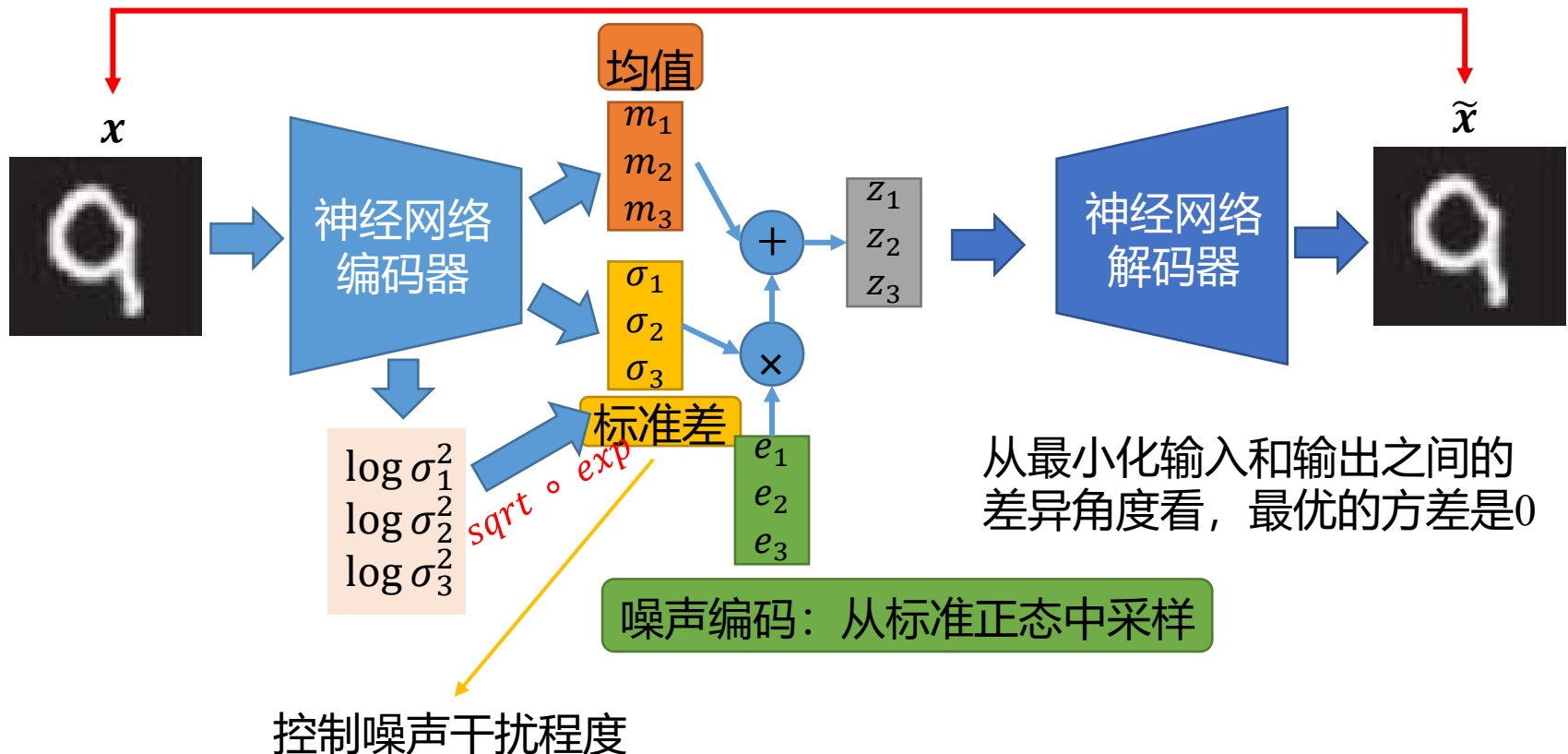
dist: 67



按照图片的256维
的编码向量，最相
似的图片

Variational AutoEncoder (VAE)

最小化输入和输出之间的差异 $(x - \tilde{x})^2$

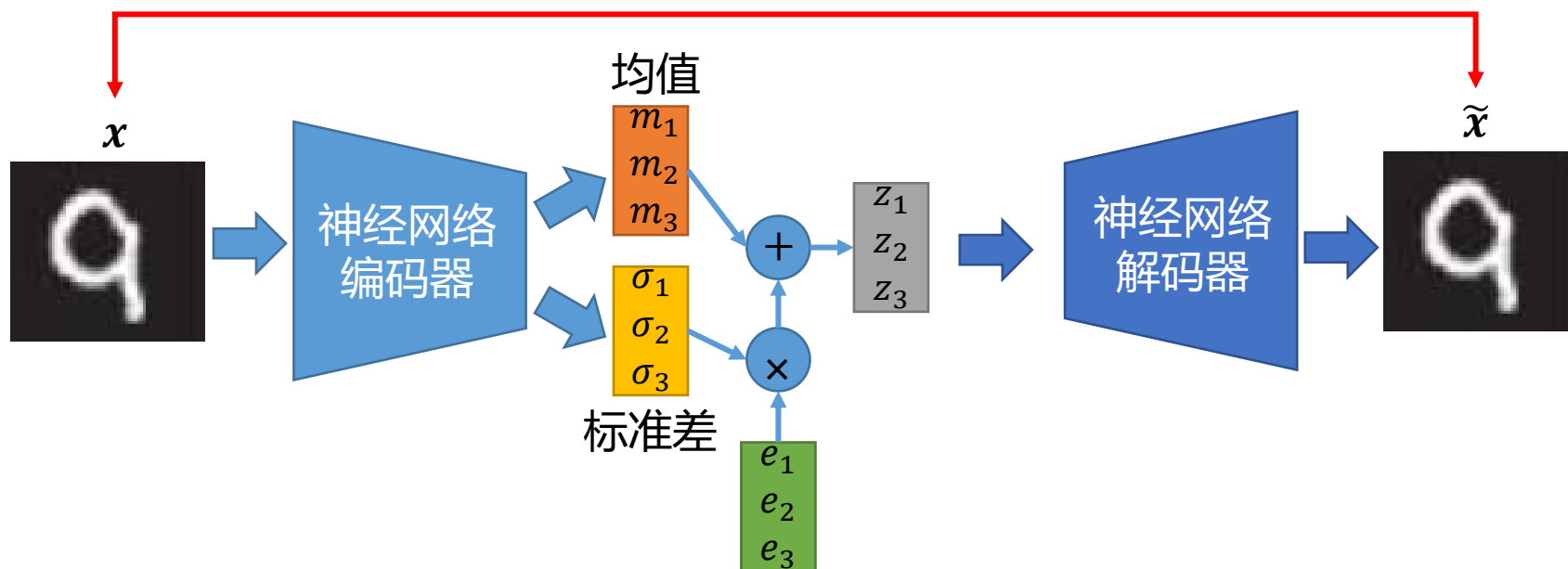


Variational AutoEncoder (VAE)

同时最小化 $c = \sum_{i=1}^3 \underbrace{(\sigma_i^2 - (1 + \log \sigma_i^2))}_{\text{自动学习方差}} + \underbrace{(m_i)^2}_{\text{正则化, 使其不离0远}}$

记 $\sigma_i^2 = \exp(\epsilon_i)$

最小化输入和输出之间的差异 $(x - \tilde{x})^2$



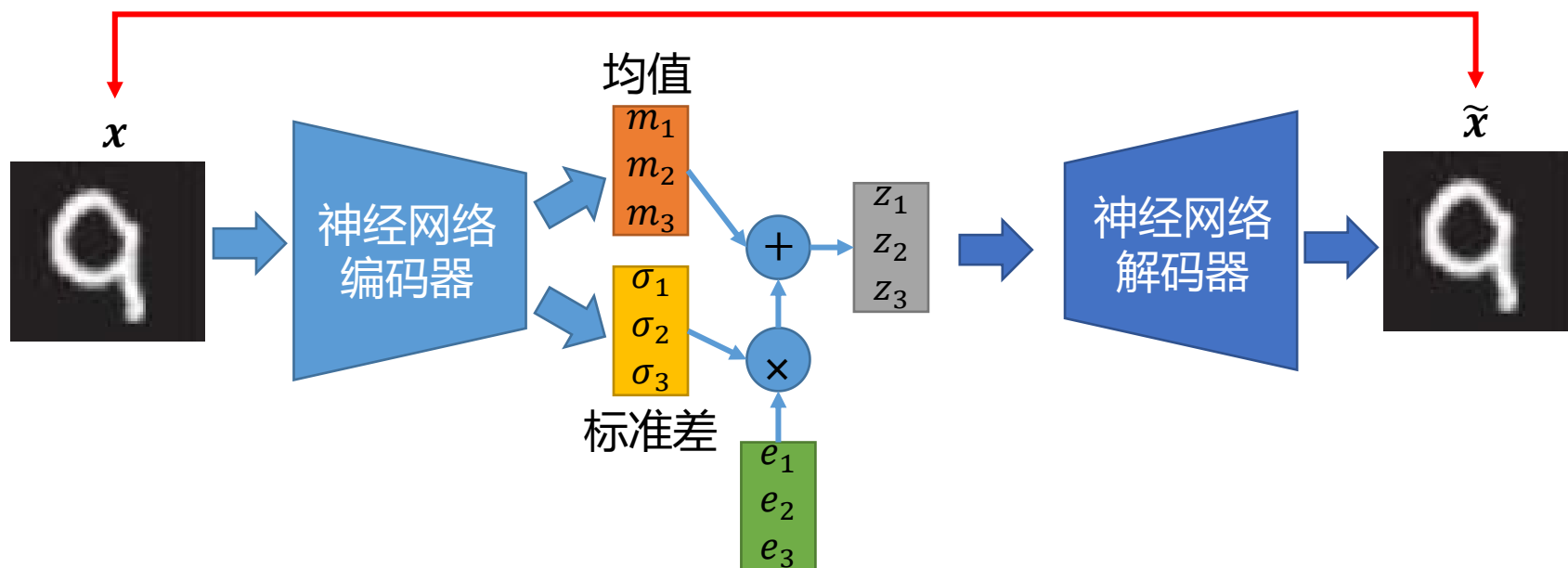
从标准正态中采样

Variational AutoEncoder (VAE)

记 $\sigma_i^2 = \exp(\epsilon_i)$

同时最小化 $c = \sum_{i=1}^3 \underbrace{(\exp(\epsilon_i) - (1 + \epsilon_i))}_{\text{自动学习方差}} + \underbrace{(m_i)^2}_{\text{正则化, 使其不离0远}}$

最小化输入和输出之间的差异 $(x - \tilde{x})^2$



Variational AutoEncoder (VAE)

记 $\sigma_i^2 = \exp(\epsilon_i)$

同时最小化 $c = \sum_{i=1}^3 \underbrace{(\exp(\epsilon_i) - (1 + \epsilon_i))}_{\text{自动学习方差}} + \underbrace{(m_i)^2}_{\text{正则化, 使其不离0远}}$

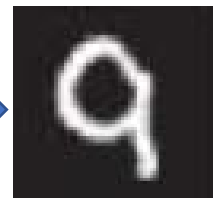
自动学习方差

正则化, 使其
不离0远

的差异 $(x - \tilde{x})^2$

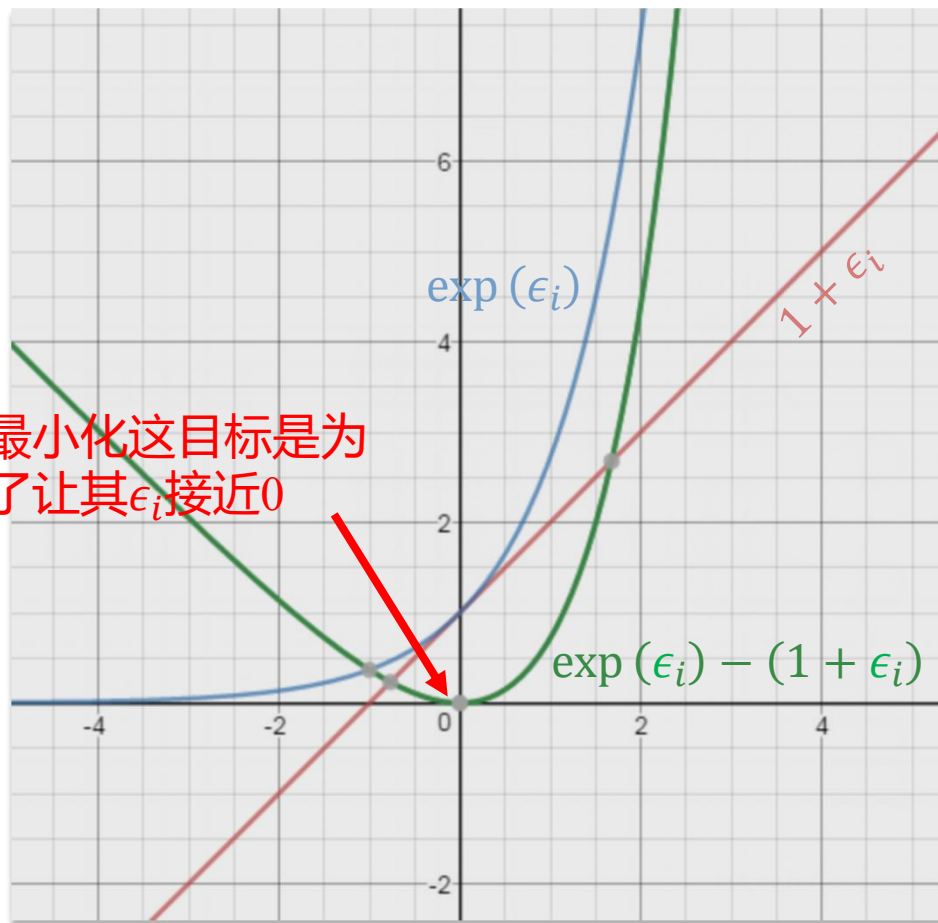
\tilde{x}

神经网络
解码器



z_1
 z_2
 z_3

中采样



一个例子理解VAE

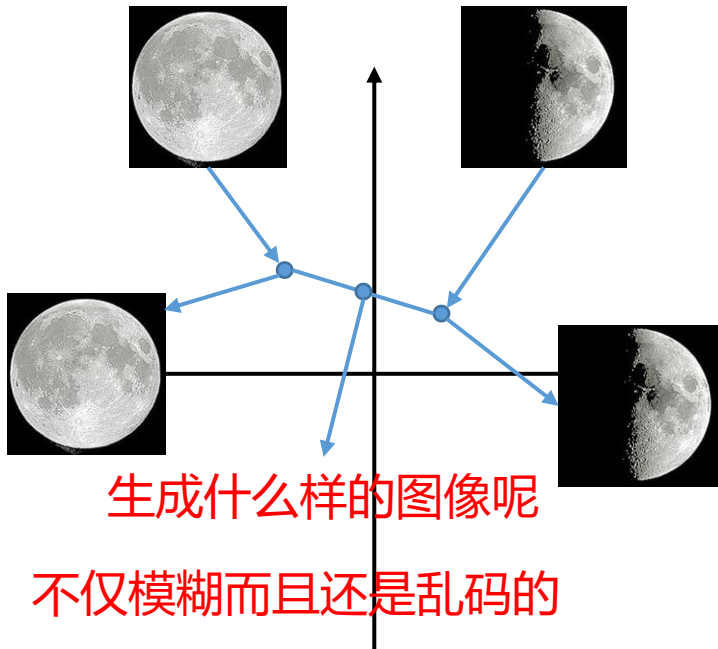
➤为什么引入噪声和重构误差损失？

➤把VAE正在做的这件事情比作是在参加高考：

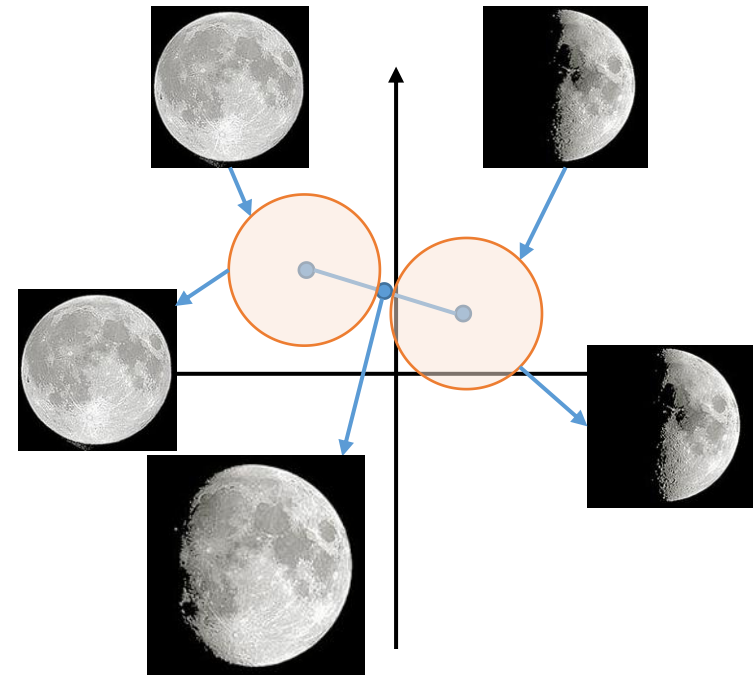
- 一般为了能够真正在考场取得好成绩（模型预测），
- 学生在平常的学习生活中需要做各式各样的测试（模型训练）。
- 这些测试题的考试难度（方差编码 σ ）应该由老师来定，因为只有这样才能客观的检验学生的学习能力。
- 假如没有老师监督（辅助loss），而是让学生（模型）来决定考题的难度（分配给噪声的权重），
- 那么学生肯定是偏向让测试题难度降到最低（使噪声影响最小），最好一点难度（噪声）都没有，从而能够考满分（使最终的重构误差为0）。
- 因此，为了能够真正在高考上取得好成绩，而不是以这种投机取巧的方式，所以必须引入老师（辅助的loss）这个中间人来监督这个课堂测试（训练过程）的难度。

VAE V.S. AE

AutoEncoder



Variational AutoEncoder



VAE V.S. AE

编码



神经网络
解码器



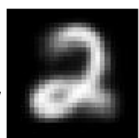
?

$$\mathbf{a} = \alpha \mathbf{a}_{\text{start}} + (1 - \alpha) \mathbf{a}_{\text{end}}$$

\mathbf{a}_{end}



AE



0.

0.1

0.2

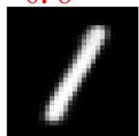
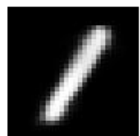
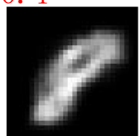
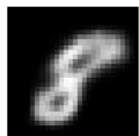
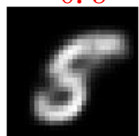
0.3

0.4

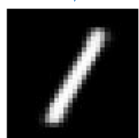
0.5

0.6

VAE



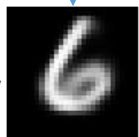
$\mathbf{a}_{\text{start}}$



\mathbf{a}_{end}



AE



0.

0.1

0.2

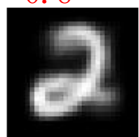
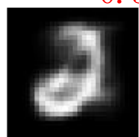
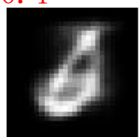
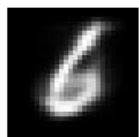
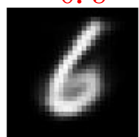
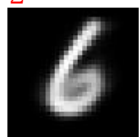
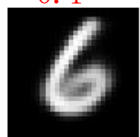
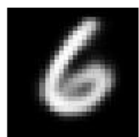
0.3

0.4

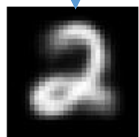
0.5

0.6

VAE



$\mathbf{a}_{\text{start}}$



高斯混合模型

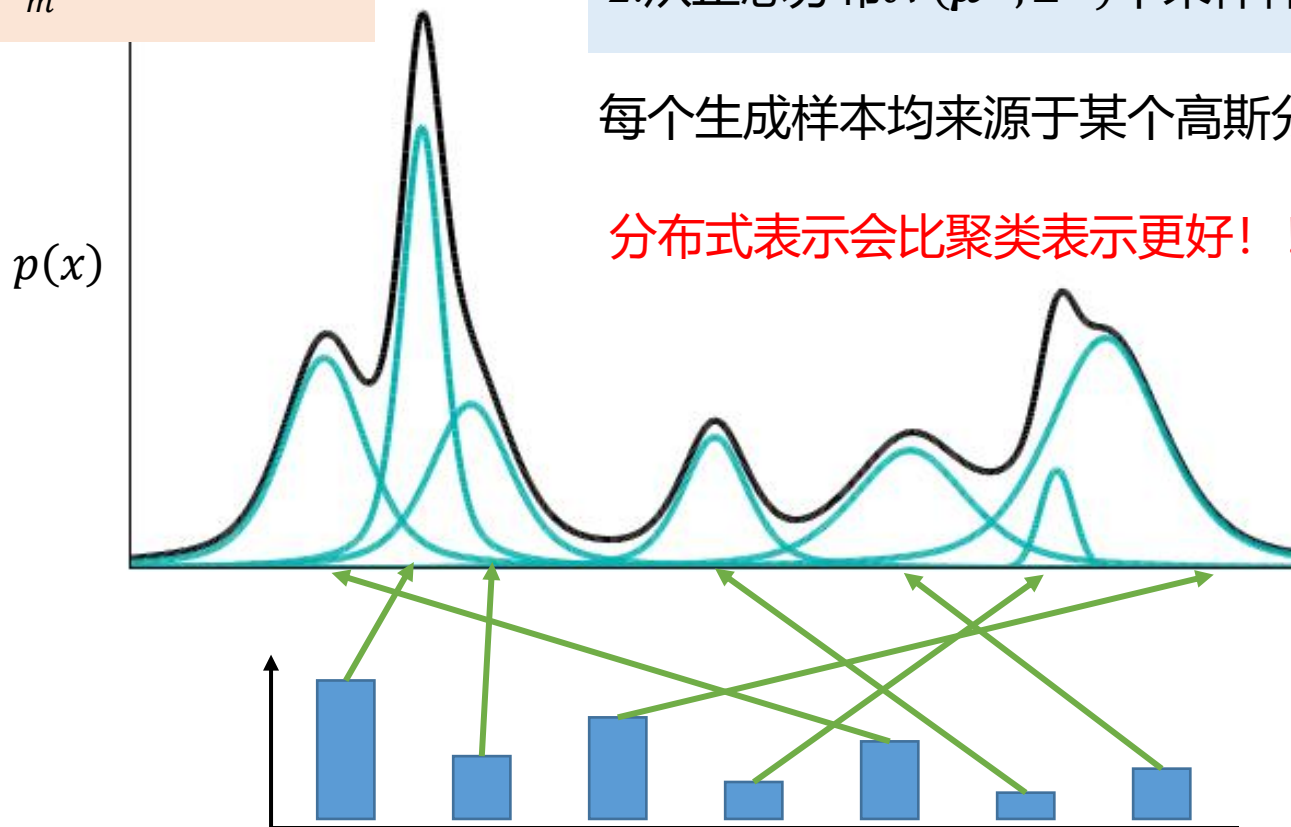
$$p(x) = \sum_m P(m)p(x|m)$$

样本生成过程

1. 先从多项分布中 $P(m)$ 采样一个 m
2. 从正态分布 $\mathcal{N}(\mu^m, \Sigma^m)$ 中采样样本 x

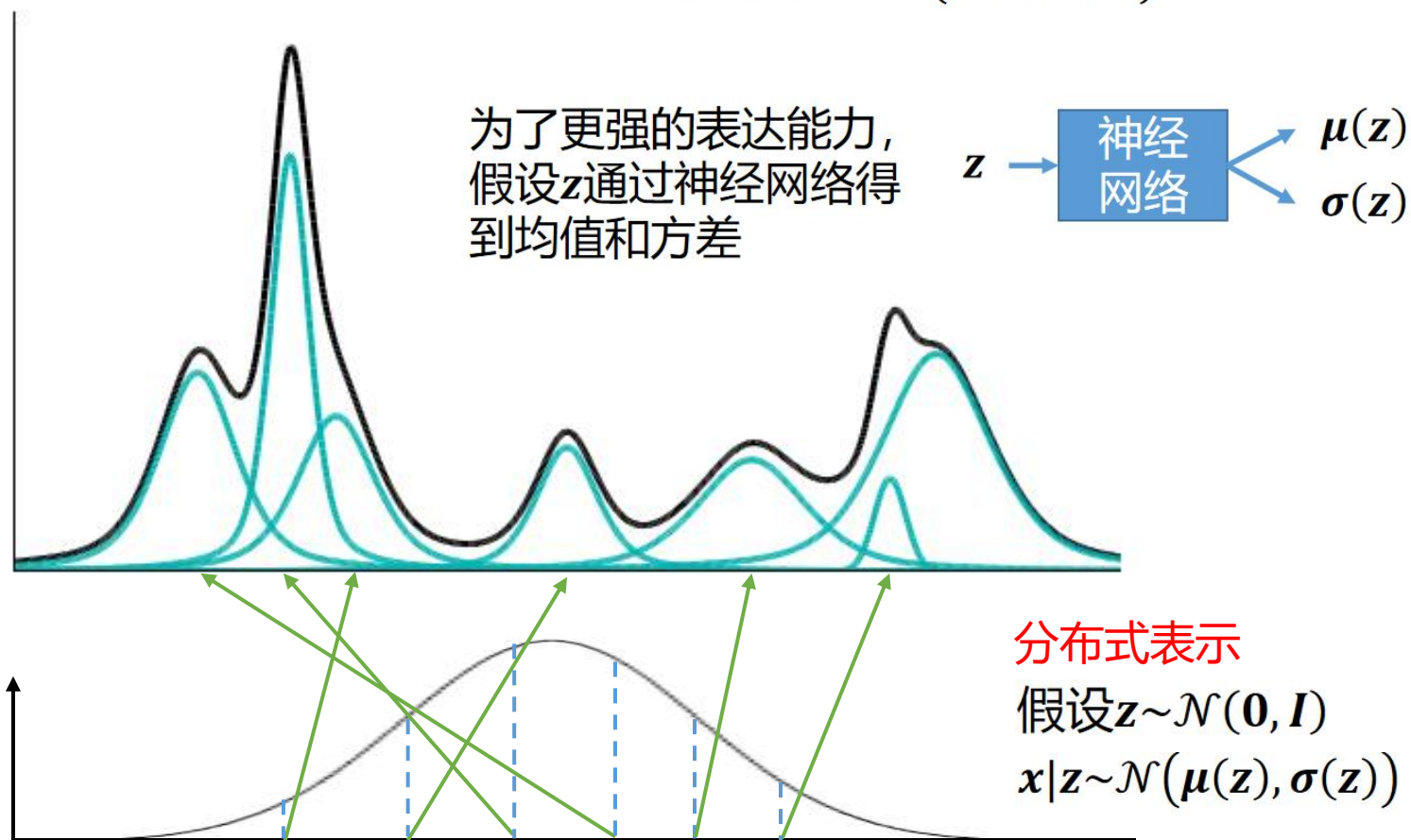
每个生成样本均来源于某个高斯分布

分布式表示会比聚类表示更好!!!



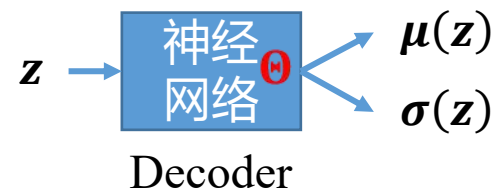
无限高斯混合模型

$$P(\mathbf{x}) = \int \underbrace{P(\mathbf{z})}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} \underbrace{P(\mathbf{x}|\mathbf{z})}_{\mathcal{N}(\boldsymbol{\mu}(\mathbf{z}), \boldsymbol{\sigma}(\mathbf{z}))} d\mathbf{z}$$



VAE最大化似然概率

$$P(\mathbf{x}) = \int \underbrace{P(\mathbf{z})}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} \underbrace{P(\mathbf{x}|\mathbf{z})}_{\mathcal{N}(\boldsymbol{\mu}(\mathbf{z}), \boldsymbol{\sigma}(\mathbf{z}))} d\mathbf{z}$$



优化目标: $\log P(\mathbf{x}) = \log \int P(\mathbf{z})P(\mathbf{x}|\mathbf{z})d\mathbf{z}$

优化参数: 神经网络中的参数 Θ

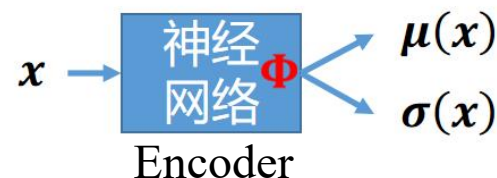
$$\ln P(\mathbf{x}|\Theta) = \ln \int Q(\mathbf{z}) \frac{P(\mathbf{x}, \mathbf{z}|\Theta)}{Q(\mathbf{z})}$$

当 $Q(\mathbf{z}) = P(\mathbf{z}|\mathbf{x}, \Theta)$ 时等式成立

$$\geq \underbrace{\int Q(\mathbf{z}) \ln \frac{P(\mathbf{x}, \mathbf{z}|\Theta)}{Q(\mathbf{z})}}_{\mathcal{L}(Q, \Theta)}$$

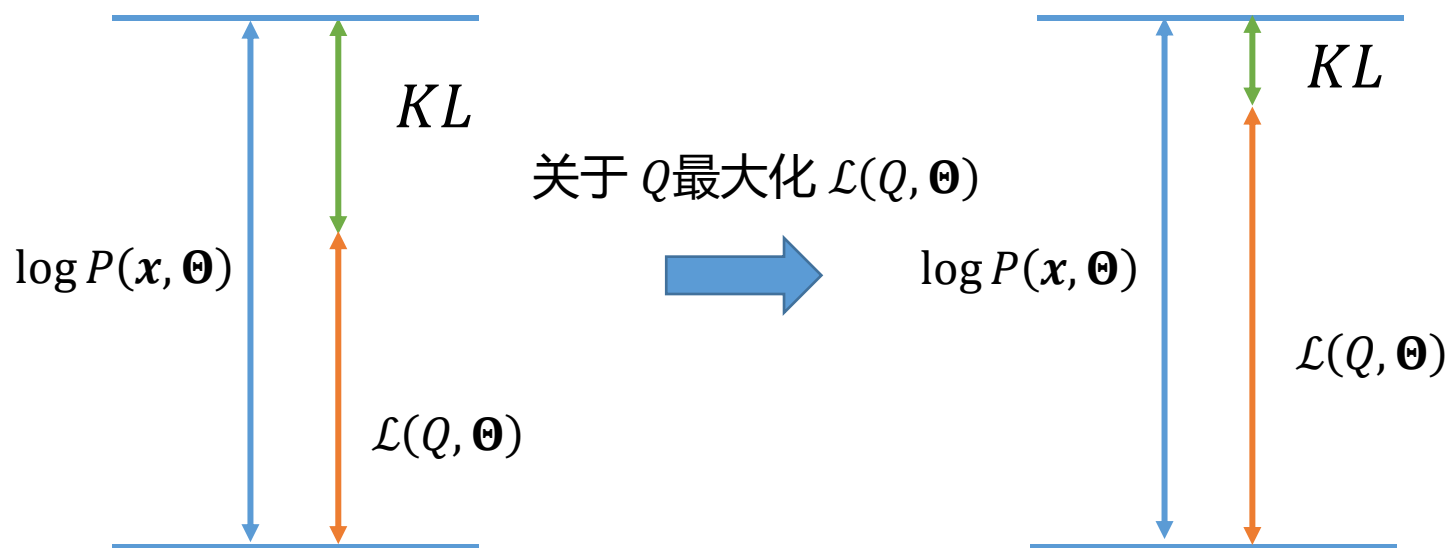
由于 $P(\mathbf{z}|\mathbf{x}, \Theta)$ 通常难以计算, 假设其为正态分布, 记为 $Q(\mathbf{z}|\mathbf{x}, \Phi)$

用神经网络近似其参数



VAE最大化似然概率

$$\mathcal{L}(Q, \Theta) = \int Q(z|x, \Phi) \ln \frac{P(x, z|\Theta)}{Q(z|x, \Phi)} = \ln P(x|\Theta) - KL(Q(z|x, \Phi) || P(z|x, \Theta))$$



关于 Q 最大化 $\mathcal{L}(Q, \Theta)$ 会使得 $Q(z|x, \Phi)$ 与 $P(z|x, \Theta)$ 差异减小
即 $Q(z|x, \Phi)$ 是 $P(z|x, \Theta)$ 的一个好的近似

VAE最大化似然概率

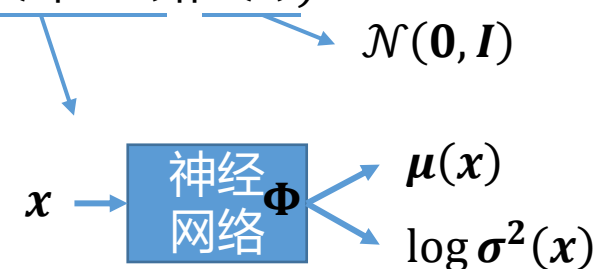
$$\begin{aligned}\ln P(\mathbf{x}|\Theta) &= \ln \int Q(\mathbf{z}|\mathbf{x}, \Phi) \frac{P(\mathbf{x}, \mathbf{z}|\Theta)}{Q(\mathbf{z}|\mathbf{x}, \Phi)} \\&\geq \underbrace{\int Q(\mathbf{z}|\mathbf{x}, \Phi) \ln \frac{P(\mathbf{x}, \mathbf{z}|\Theta)}{Q(\mathbf{z}|\mathbf{x}, \Phi)}}_{\mathcal{L}(Q, \Theta)} = \int Q(\mathbf{z}|\mathbf{x}, \Phi) \ln \frac{P(\mathbf{x}, \mathbf{z}|\Theta)}{Q(\mathbf{z}|\mathbf{x}, \Phi)} \\&= \int Q(\mathbf{z}|\mathbf{x}, \Phi) \ln \frac{P(\mathbf{x}|\mathbf{z}, \Theta)P(\mathbf{z})}{Q(\mathbf{z}|\mathbf{x}, \Phi)} \\&= \int Q(\mathbf{z}|\mathbf{x}, \Phi) \ln P(\mathbf{x}|\mathbf{z}, \Theta) + \underbrace{\int Q(\mathbf{z}|\mathbf{x}, \Phi) \ln \frac{P(\mathbf{z})}{Q(\mathbf{z}|\mathbf{x}, \Phi)}}_{-KL(Q(\mathbf{z}|\mathbf{x}, \Phi)||P(\mathbf{z}))}\end{aligned}$$

VAE最大化似然概率

$$\mathcal{L}(Q, \Theta) = \int Q(\mathbf{z}|\mathbf{x}, \Phi) \ln P(\mathbf{x}|\mathbf{z}, \Theta) - KL(Q(\mathbf{z}|\mathbf{x}, \Phi) || P(\mathbf{z}))$$

$$\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\exp(-1/2(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}))}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}}$$

用神经网络近似其参数



$$\begin{aligned} & KL(\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2) || \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})) \\ &= \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \frac{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2)}{\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})} \quad \mathbf{z} \text{ 为 } D \text{ 维向量} \\ &= \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2) \left(-\frac{1}{2} \log |\text{diag}(\boldsymbol{\sigma}^2)| - \frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^\top \text{diag}(\boldsymbol{\sigma}^2)^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \frac{1}{2} \mathbf{z}^\top \mathbf{z} \right) \\ &= -\frac{1}{2} \log |\text{diag}(\boldsymbol{\sigma}^2)| - \frac{D}{2} + \frac{1}{2} \text{tr}(\boldsymbol{\mu} \boldsymbol{\mu}^\top + \text{diag}(\boldsymbol{\sigma}^2)) \\ &= -\frac{1}{2} \sum_i \log \sigma_i^2 - \frac{D}{2} + \frac{1}{2} (\sum_i \mu_i^2 + \sum_i \sigma_i^2) \\ &= \frac{1}{2} \sum_i (-1 - \log \sigma_i^2 + \mu_i^2 + \sigma_i^2) \end{aligned}$$

$$\begin{aligned} & \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})] \\ &= \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]) \\ &= \text{tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}) \\ &= D \end{aligned}$$

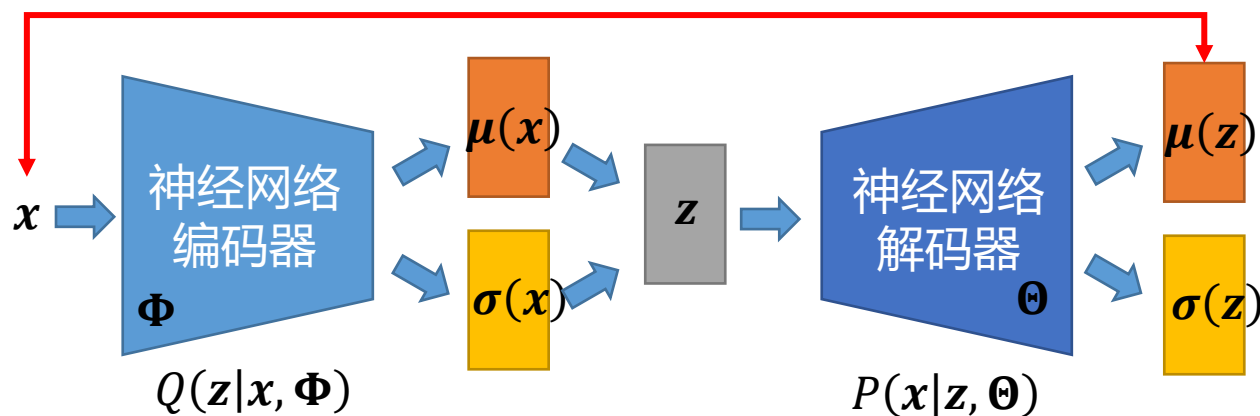
VAE最大化似然概率

$$\text{最大化 } \mathcal{L}(Q, \Theta) = \underbrace{\int Q(z|x, \Phi) \ln P(x|z, \Theta)}_{\text{orange line}} - \underbrace{KL(Q(z|x, \Phi) || P(z))}_{\text{blue line}}$$

$$\begin{aligned} & \mathbb{E}_{Q(z|x, \Phi)} [\ln P(x|z, \Theta)] \\ & \approx \frac{1}{L} \sum_l \ln P(x|z_l, \Theta) \quad z_l \sim Q(z|x, \Phi) \\ & \quad \quad \quad \sim \mathcal{N}(\mu(x), \sigma(x)) \end{aligned}$$

$$\frac{1}{2} \sum_i (-1 - \log \sigma_i^2 + \mu_i^2 + \sigma_i^2)$$

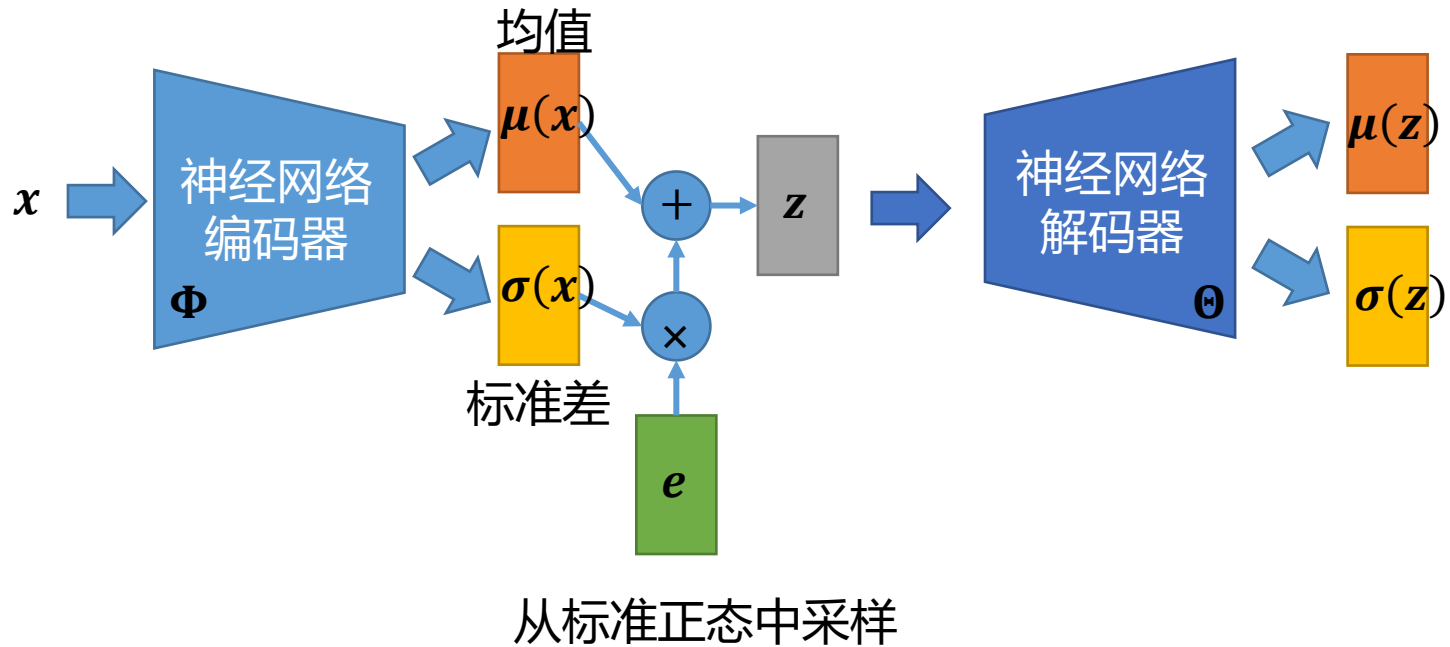
最小化输入和输出之间的差异 $(x - \tilde{x})^2$ 同时最小化



直接从 $\mathcal{N}(\mu(x), \sigma(x))$ 采样编码会导致梯度无法回传

VAE的重采样技术

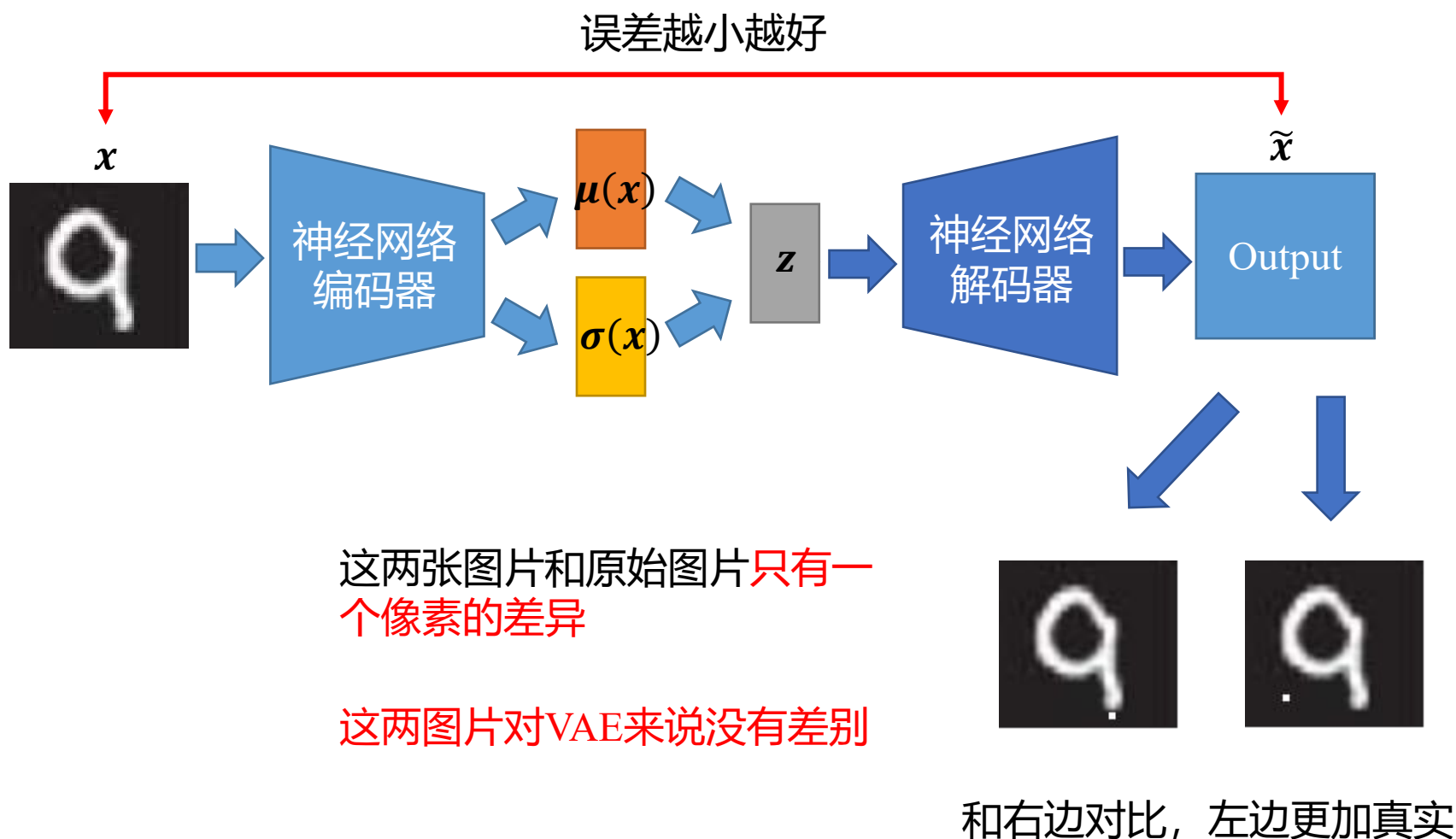
直接从 $\mathcal{N}(\mu(x), \sigma(x))$ 采样 z 等价于 从 $\mathcal{N}(\mathbf{0}, I)$ 采样 e 后, 运用如下变换
 $z = e \odot \sigma(x) + \mu(x)$



VAE生成出来的图像



VAE的问题



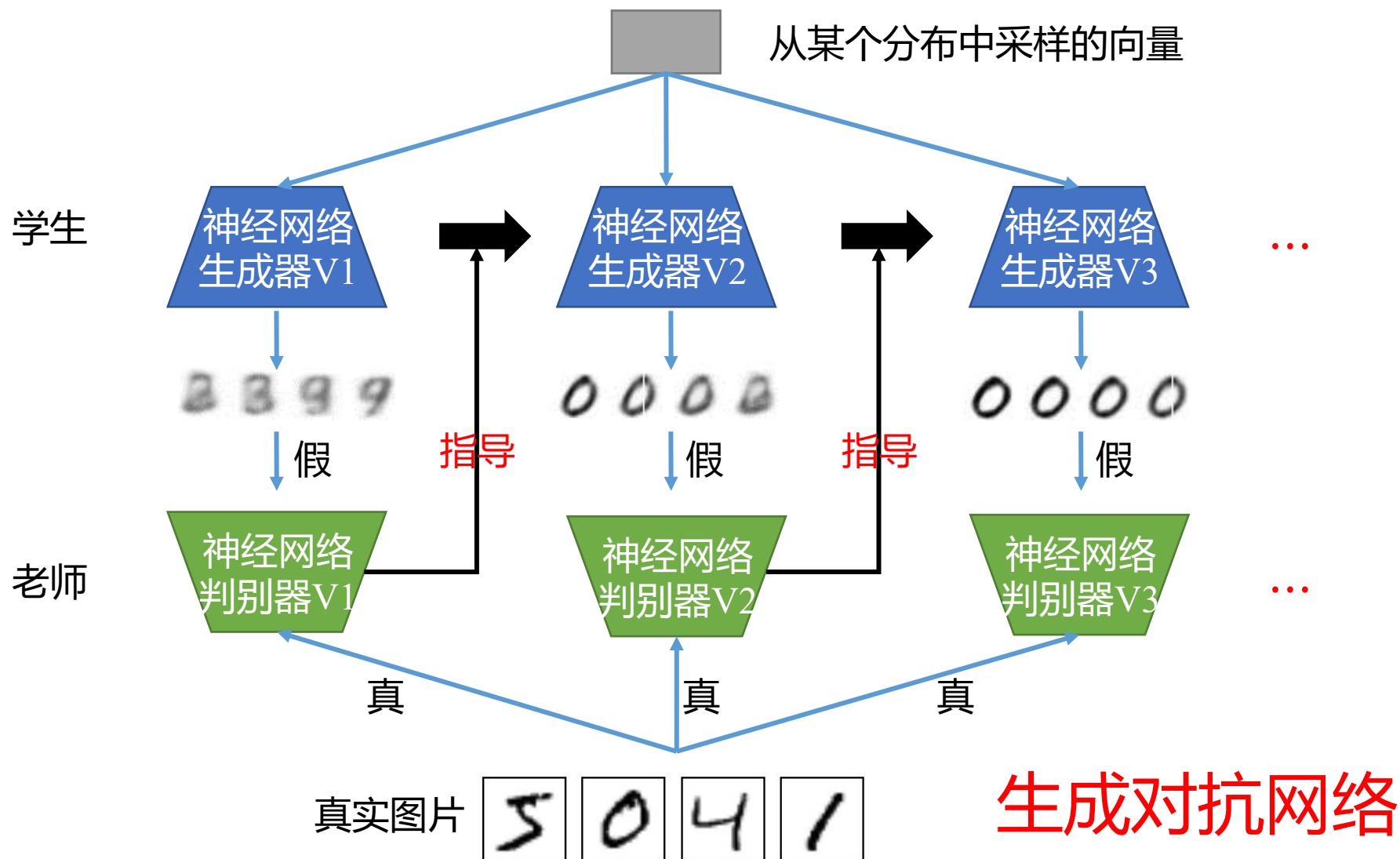
VAE只是记住了数据库里面的图片，并没有学习到如何生成新图片

VAE生成的图片



VAE模型往往产生不现实的、模糊的样本。这是由于它是通过直接计算生成图片和原始图片之间的均方误差来保证生成图片的真实性，所以得到的是——张“平均图像”。

解码器训练的新思路—引入判别器

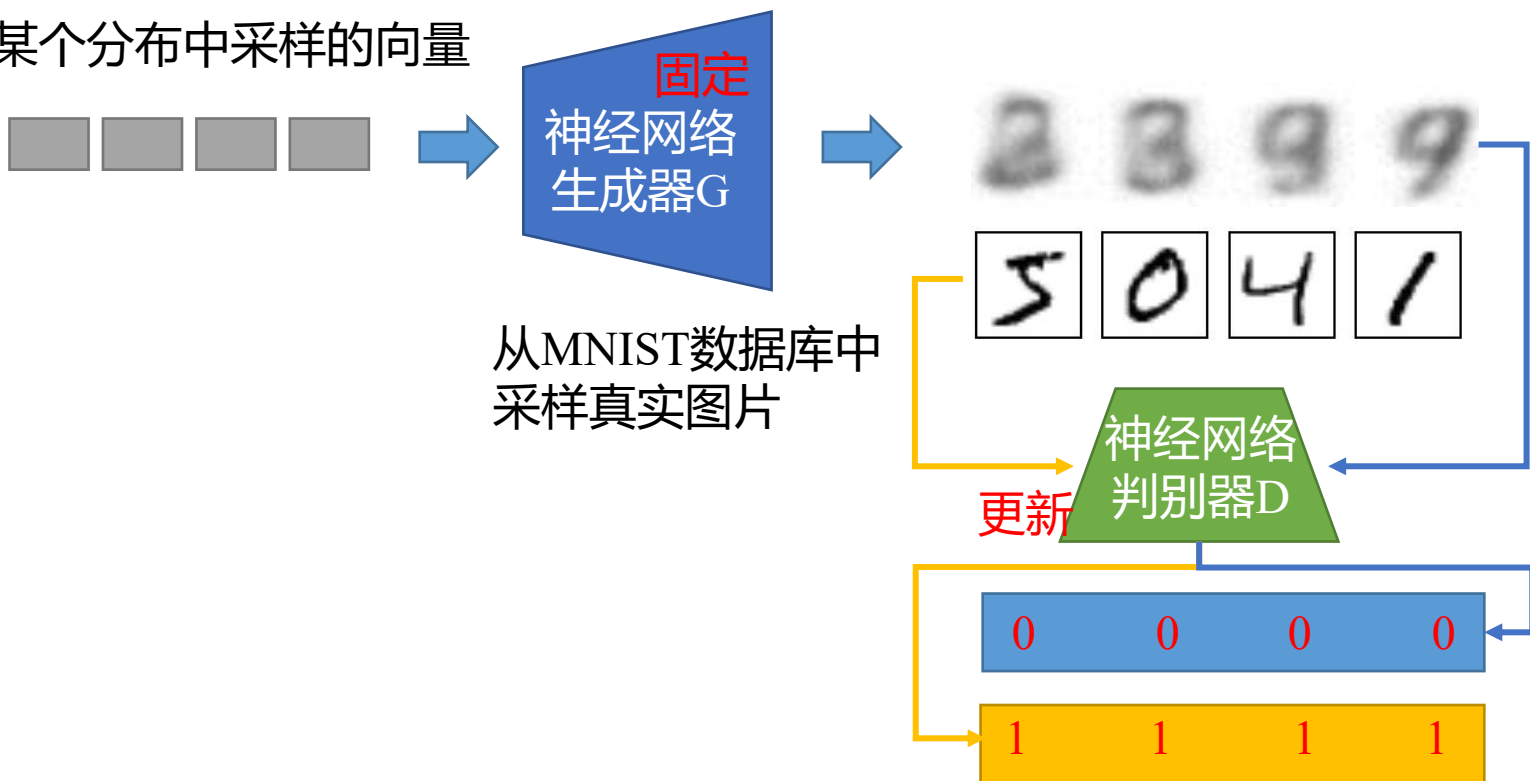


生成对抗网络的算法

- 初始化生成器和判别器
- 在每个训练步
- 1) 固定生成器, 更新判别器D

学会给真实图片高分, 给生成图片低分

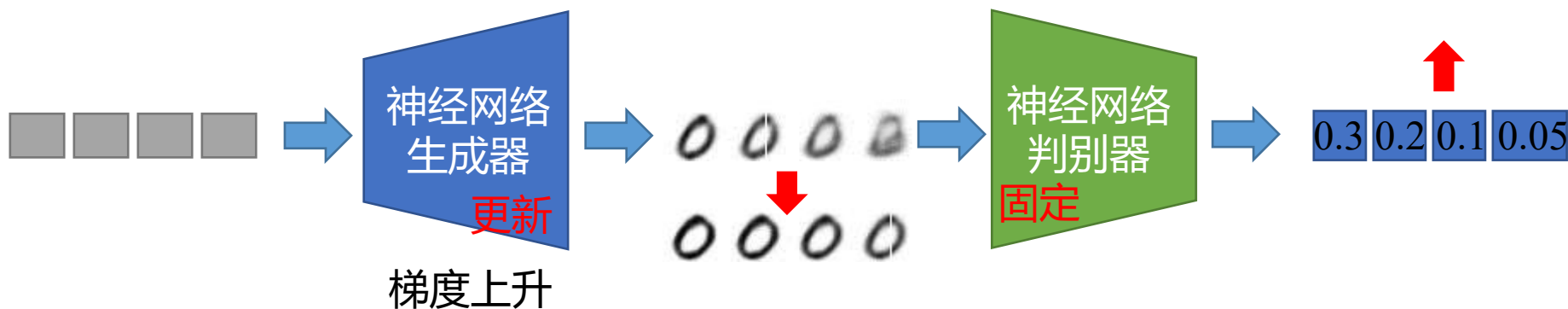
从某个分布中采样的向量



生成对抗网络的算法

- 初始化生成器和判别器
- 在每个训练步
- 2) 固定判别器，更新生成器G

生成得分更高的图片
让判别器更易混淆的图片



生成对抗网络的算法

- 随机初始化判别器参数 θ_d 和生成器参数 θ_g
- While not convergent:

更新D

- 从数据库中采样 m 样本 $\{x^1, x^2, \dots, x^m\}$
- 从某个分布中采样 m 个噪声样本 $\{z^1, z^2, \dots, z^m\}$
- 用生成器生成样本 $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- 更新判别器参数 θ_d , 通过最大化如下目标

$$\tilde{\mathcal{L}} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(\tilde{x}^i))$$

$$\theta_d \leftarrow \theta_d + \eta \nabla \tilde{\mathcal{L}}(\theta_d)$$

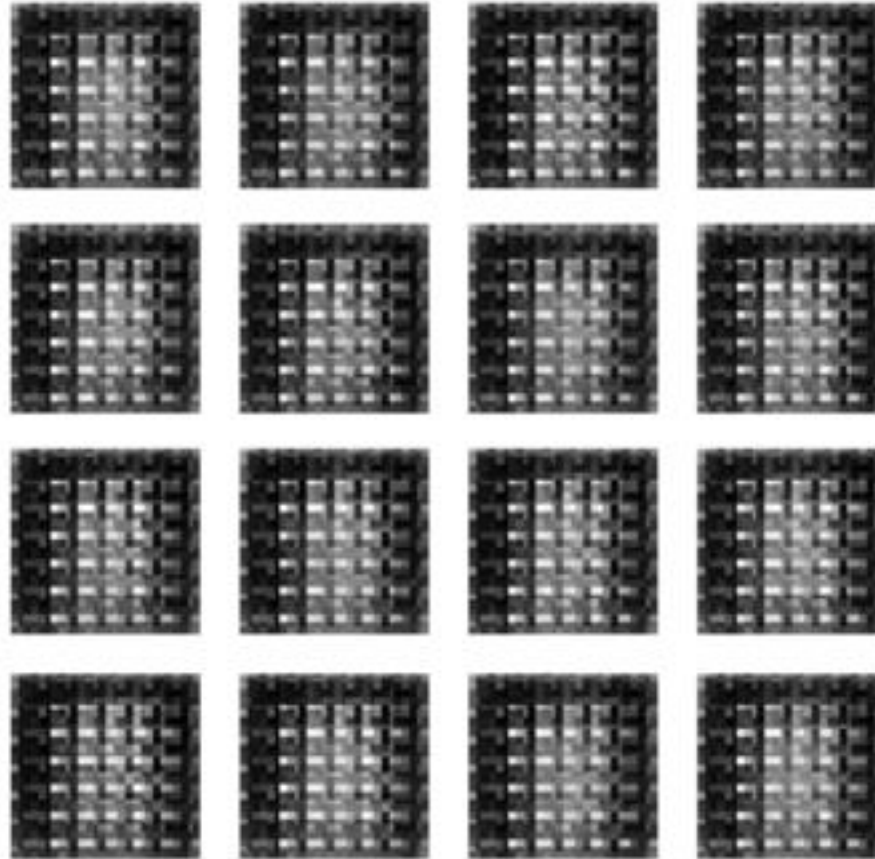
更新G

- 从某个分布中采样 m 个噪声样本 $\{z^1, z^2, \dots, z^m\}$
- 更新生成器参数 θ_g , 通过最大化如下目标

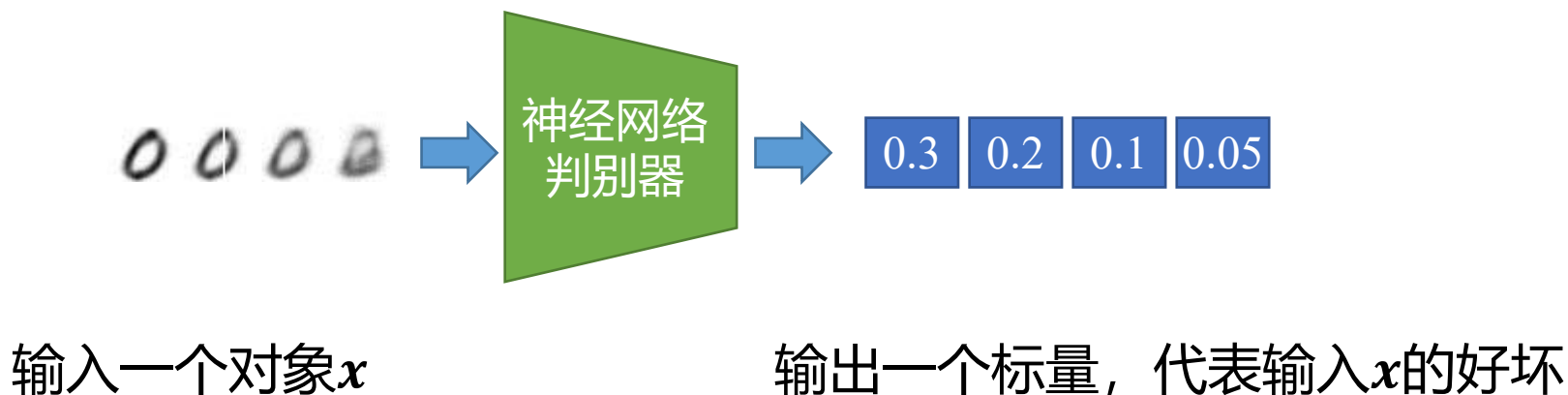
$$\tilde{\mathcal{L}} = \frac{1}{m} \sum_{i=1}^m \log (D(G(z^i)))$$

$$\theta_g \leftarrow \theta_g + \eta \nabla \tilde{\mathcal{L}}(\theta_g)$$

生成对抗网络用于MNIST数据集



判别器D的更多理解——假设只有判别器



能否用判别器生成对象呢？ 答案是可以的！！

判别器D的更多理解—假设只有判别器

- 假设 $D(x)$ 已经训练的很好，那么可以通过如下方式生成图像

$$x^* = \arg \max_{x \in \mathcal{X}} D(x)$$

- 即找到 判别器 打分最高（最真实） 的图片
 - 可以枚举所有的 $x \in \mathcal{X}$ 找最高分的 x^*
 - 也可以通过梯度上升法 $x^t \leftarrow x^t + \eta D(x^t)$

但是如何学习判别器呢？

判别器D的更多理解——假设只有判别器

- 在没有生成器的情况下，如何学习判别器呢

- 我们只有真实的图像，如果只用真实数据训练，

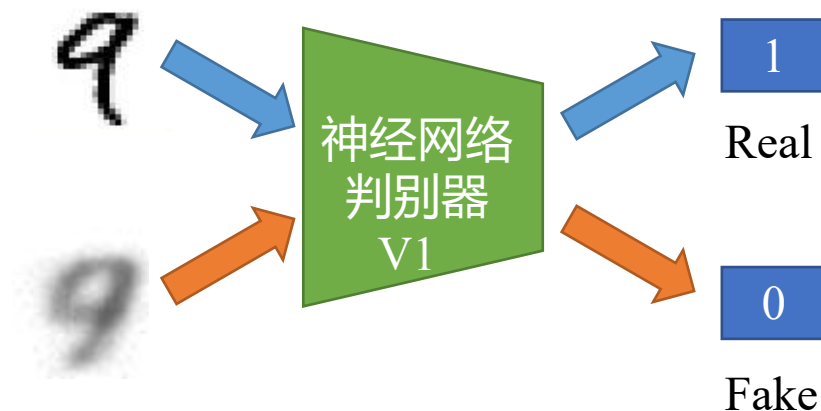
那么判别器会学习到总是输出1



- 因此，判别器学习需要负样本，即不真实样本

判别器D的更多理解——假设只有判别器

- 要训练好判别器，负样本非常关键



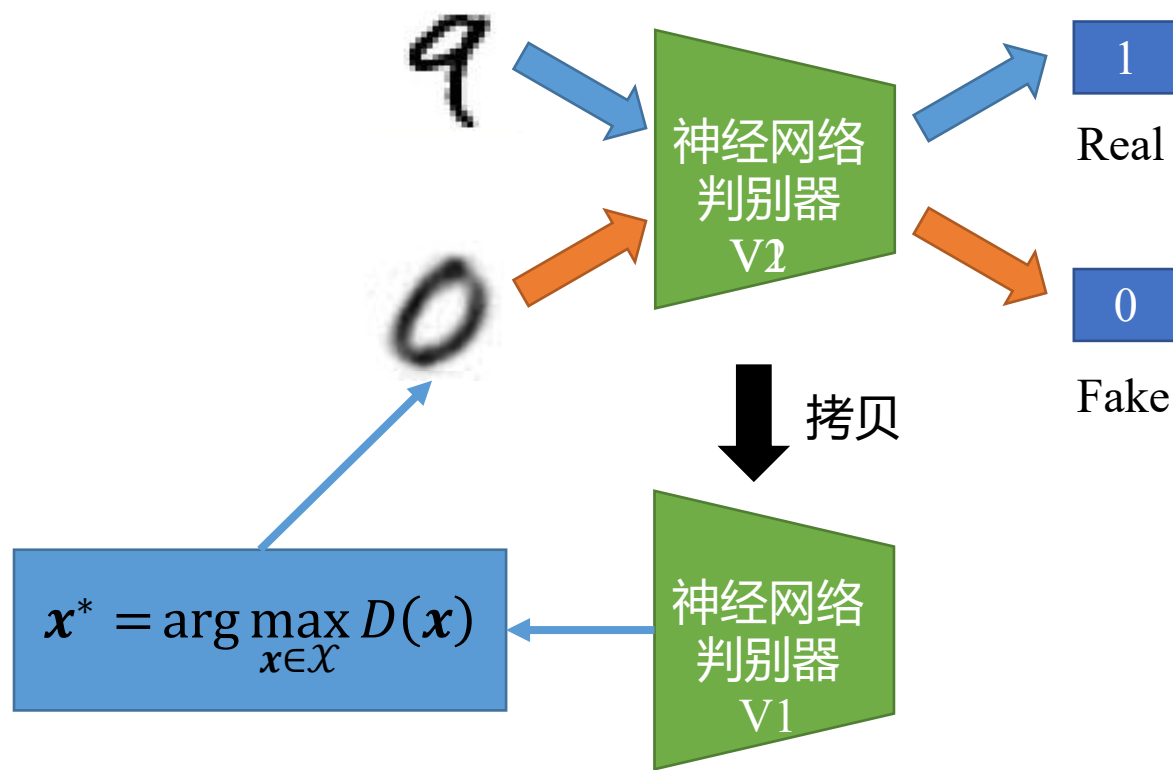
在训练完后，输入



输出这么高的分
是因为训练时负样本太烂了，
使得判别器会误认为已经非常好了

判别器D的更多理解—假设只有判别器

- 要训练好判别器，负样本非常关键



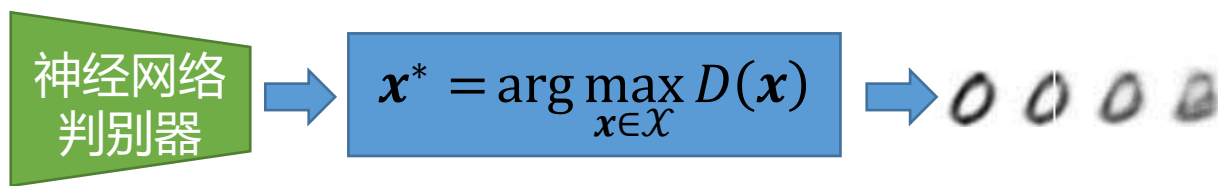
判别器D的更多理解—假设只有判别器

在没有生成器的情况下判别器D的训练算法

- 给定真实样本集，随机生成负样本集
- While not convergent
 - 学习判别器D，使得给正样本高分，负样本低分



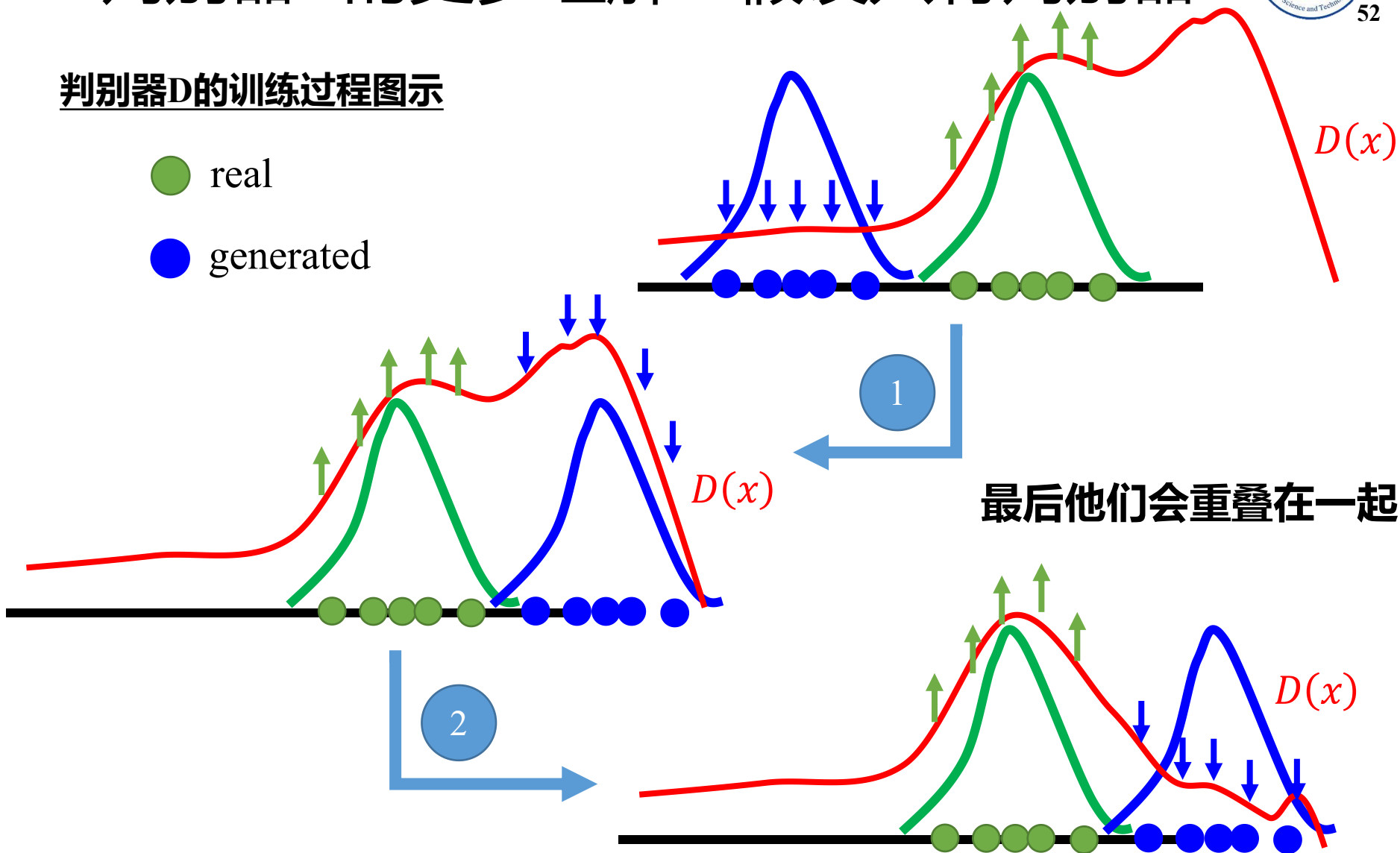
- 用判别器D 生成负样本集



判别器D的更多理解—假设只有判别器

判别器D的训练过程图示

- real
- generated



判别器+生成器

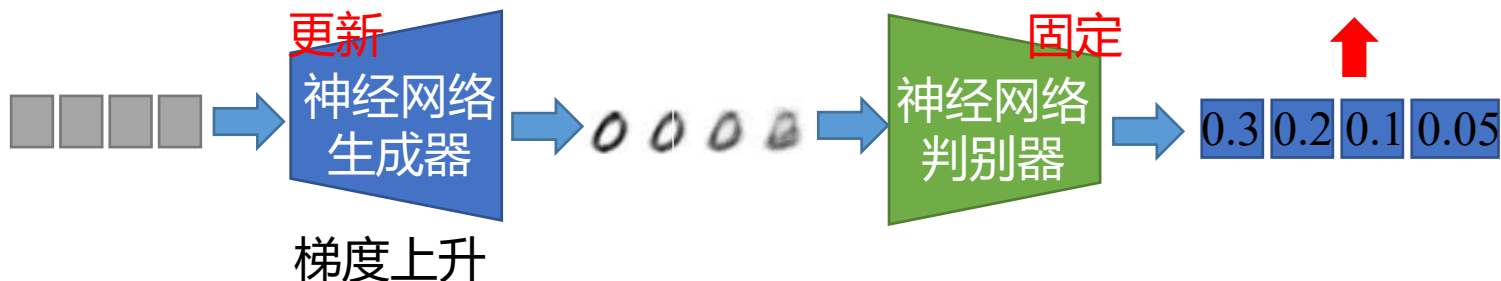
GAN的训练算法

- 给定真实样本集，随机生成负样本集
- While not convergent
 - 学习判别器D，使得给正样本高分，负样本低分



- 用判别器D 生成负样本集

$$x^* = \arg \max_{x \in \mathcal{X}} D(x)$$



GAN的优势

生成器

- 优点
 - 容易生成样本
- 缺点
 - 缺乏全局观

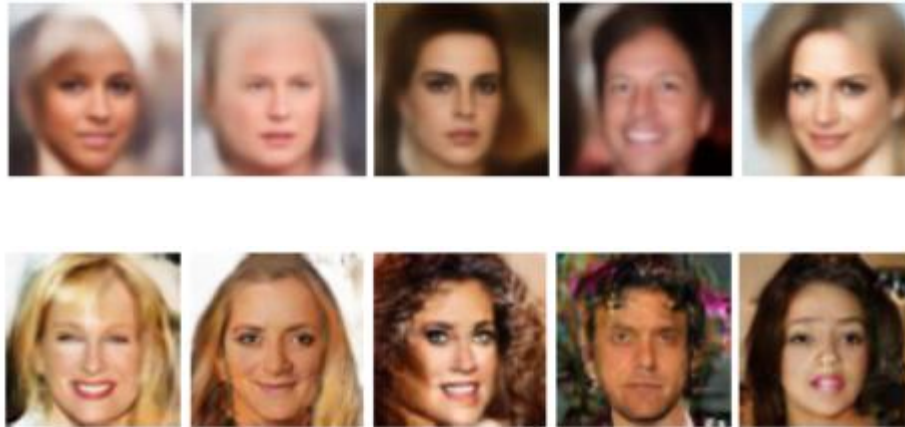
判别器

- 优点
 - 有全局意识
- 缺点
 - 难生成

综合了生成器和判别器的优点，弥补了各自的缺点

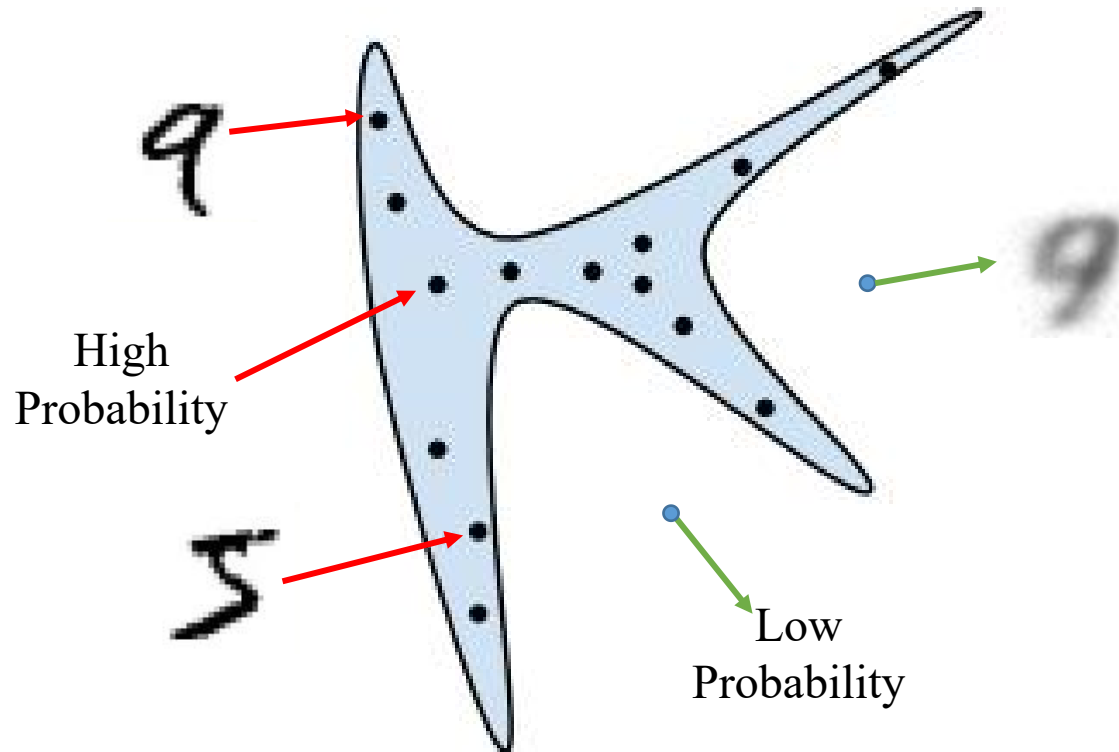
通过生成器快速生成样本，通过判别器提升生成器的全局观

GAN和VAE的生成质量对比



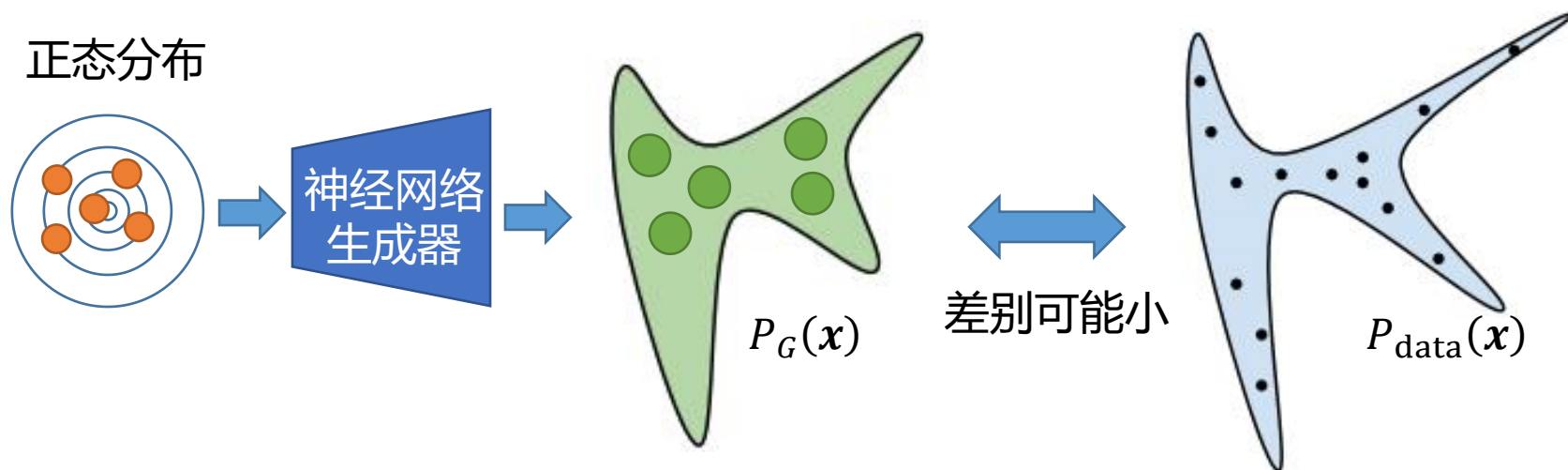
GAN背后的数学

- GAN的目标是找到真实的数据分布



GAN背后的数学

- 生成器G实际上隐式地定义了一个分布 $P_G(\mathbf{x})$ ，因为从生成器中生成样本相当于从分布 $P_G(\mathbf{x})$ 采样



$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

P_G 和 P_{data} 间的距离，但是应该如何计算呢？

GAN背后的数学

- 虽然不知道 $P_G(\mathbf{x})$ 和 $P_{\text{data}}(\mathbf{x})$ 的具体形式，但是都很容易从中采样
- 判别器的训练就是将给正样本（从 $P_{\text{data}}(\mathbf{x})$ 采样的样本）高分，给负样本（从 $P_G(\mathbf{x})$ 中采样的）低分，其目标函数如下

$$V(G, D) = E_{\mathbf{x} \sim P_{\text{data}}} [\log D(\mathbf{x})] + E_{\mathbf{x} \sim P_G} [\log (1 - D(\mathbf{x}))]$$

- 在给定G的情况下，最优的判别器 $D^* = \arg \max_D V(G, D)$

其中 $\max_D V(G, D)$ 与 P_G 和 P_{data} 间的散度（JS散度）是相关的

GAN背后的数学

★ 从 $P_{\text{data}}(x)$ 采样
★ 从 $P_G(x)$ 采样

• 直观上理解

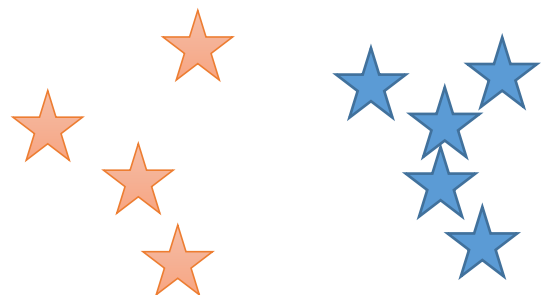
$$D^* = \arg \max_D V(G, D)$$



$P_G(x)$ 和 $P_{\text{data}}(x)$ 差异小



样本更难区分, 所以
 $\max_D V(G, D)$ 更小



$P_G(x)$ 和 $P_{\text{data}}(x)$ 差异大



样本很容易区分
 $\max_D V(G, D)$ 更大

GAN背后的数学

$$\begin{aligned} V(G, D) &= E_{\mathbf{x} \sim P_{data}} [\log D(\mathbf{x})] + E_{\mathbf{x} \sim P_G} [\log (1 - D(\mathbf{x}))] \\ &= \int P_{data}(\mathbf{x}) \log D(\mathbf{x}) d\mathbf{x} + \int P_G(\mathbf{x}) \log (1 - D(\mathbf{x})) d\mathbf{x} \\ &= \int [P_{data}(\mathbf{x}) \log D(\mathbf{x}) + P_G(\mathbf{x}) \log (1 - D(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

$$\frac{\partial V(G, D)}{\partial D(\mathbf{x})} = P_{data}(\mathbf{x}) \frac{1}{D(\mathbf{x})} + P_G(\mathbf{x}) \frac{1}{1 - D(\mathbf{x})} (-1) = 0$$

$$\Rightarrow \frac{1 - D(\mathbf{x})}{D(\mathbf{x})} = \frac{P_G(\mathbf{x})}{P_{data}(\mathbf{x})} \Rightarrow D^*(\mathbf{x}) = \frac{P_{data}(\mathbf{x})}{P_{data}(\mathbf{x}) + P_G(\mathbf{x})}$$

GAN背后的数学

$$M = \frac{1}{2}(P + Q)$$

$$\text{JSD}(P \parallel Q) = \frac{1}{2}D(P \parallel M) + \frac{1}{2}D(Q \parallel M)$$

$$V(G, D^*) = E_{\mathbf{x} \sim P_{data}}[\log D^*(\mathbf{x})] + E_{\mathbf{x} \sim P_G}[\log(1 - D^*(\mathbf{x}))]$$

$$= \int P_{data}(\mathbf{x}) \log \frac{P_{data}(\mathbf{x})}{P_{data}(\mathbf{x}) + P_G(\mathbf{x})} d\mathbf{x} + \int P_G(\mathbf{x}) \log \left(\frac{P_G(\mathbf{x})}{P_{data}(\mathbf{x}) + P_G(\mathbf{x})} \right) d\mathbf{x}$$

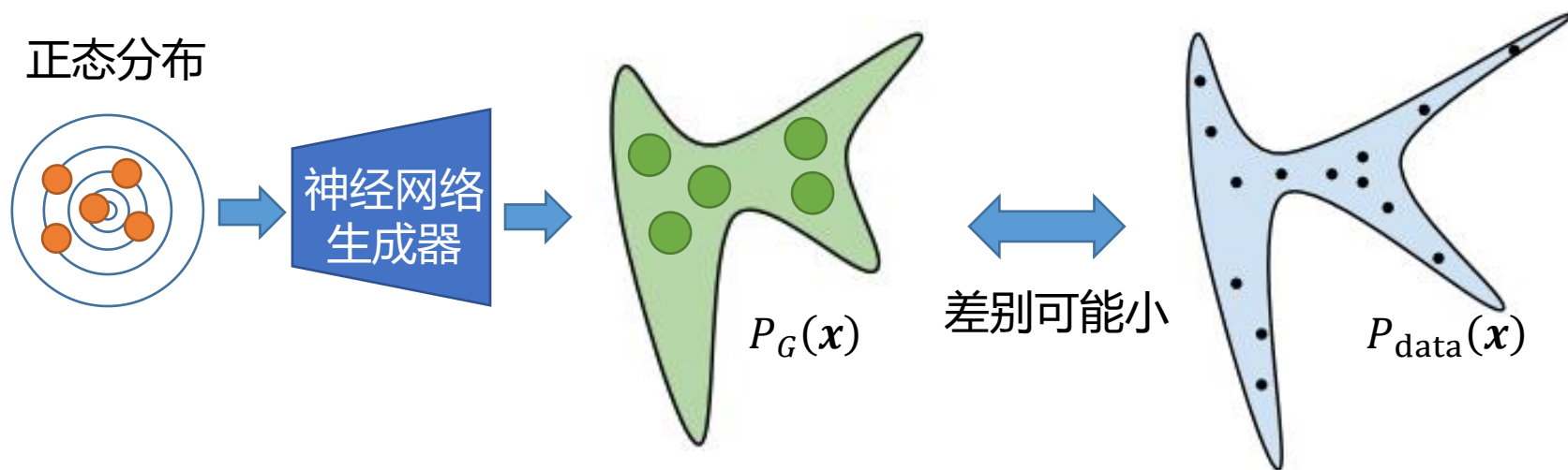
$$= \int P_{data}(\mathbf{x}) \log \frac{P_{data}(\mathbf{x})}{1/2(P_{data}(\mathbf{x}) + P_G(\mathbf{x}))} d\mathbf{x} \\ + \int P_G(\mathbf{x}) \log \left(\frac{P_G(\mathbf{x})}{1/2(P_{data}(\mathbf{x}) + P_G(\mathbf{x}))} \right) d\mathbf{x} - 2\log 2$$

$$= KL \left(P_{data}(\mathbf{x}) \parallel \frac{P_{data}(\mathbf{x}) + P_G(\mathbf{x})}{2} \right) + KL \left(P_G(\mathbf{x}) \parallel \frac{(P_{data}(\mathbf{x}) + P_G(\mathbf{x}))}{2} \right) - 2\log 2$$

$$= 2\text{JS}(P_{data}(\mathbf{x}) \parallel P_G(\mathbf{x})) - 2\log 2$$

GAN背后的数学

- 生成器G实际上隐式地定义了一个分布 $P_G(\mathbf{x})$ ，因为从生成器中生成样本相当于从分布 $P_G(\mathbf{x})$ 采样



$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$



$$G^* = \arg \min_G \max_D V(G, D)$$

P_G 和 P_{data} 间的距离，但是应该如何计算呢？

GAN背后的数学

$$V(G, D^*) = E_{x \sim P_{data}} [\log D^*(x)] + E_{x \sim P_G} [\log (1 - D^*(x))]$$

GAN的训练算法

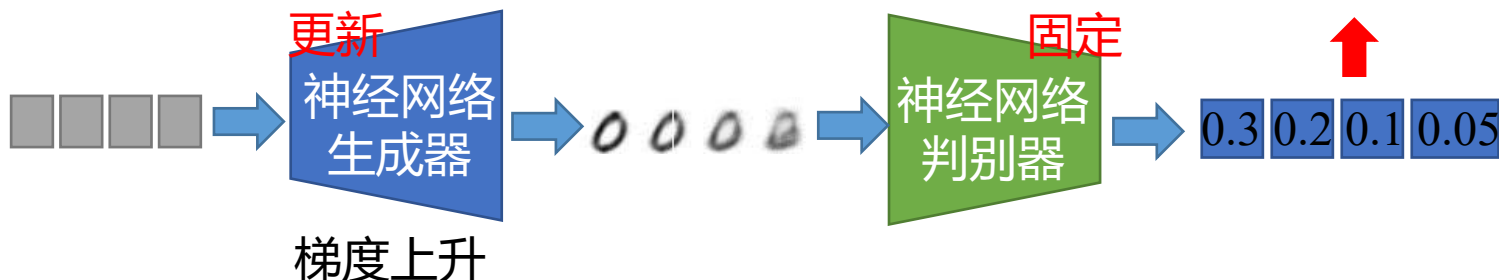
- 给定真实样本集，随机生成负样本集
- While not convergent
 - 学习判别器D，使得给正样本高分，负样本低分

$$D^* = \arg \max_D V(G, D)$$



- 用判别器D 生成负样本集

$$G^* = \arg \min_G \max_D V(G, D)$$



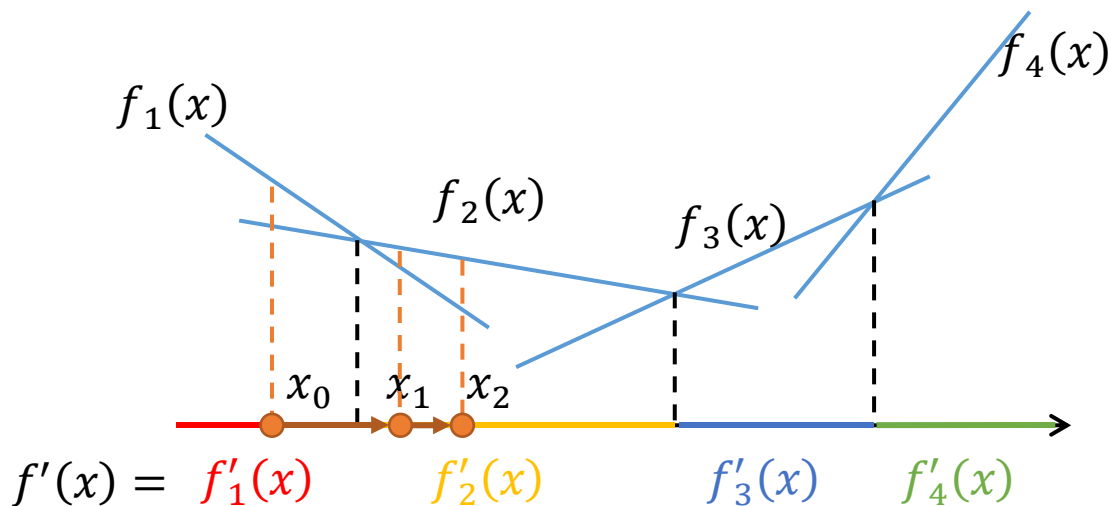
GAN背后的数学

$$G^* = \arg \min_G \max_D V(G, D)$$

- 为了找到最优的 G^* ，可以通过梯度下降法进行优化

$$\theta_G \leftarrow \theta_G - \eta \frac{\partial \left(\max_D V(G, D) \right)}{\partial \theta_G}$$

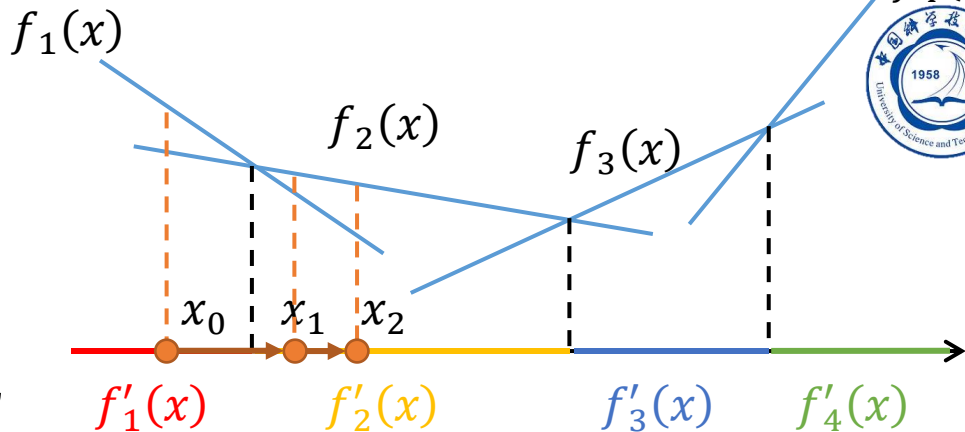
$$f(x) = \max \{f_1(x), f_2(x), f_3(x), f_4(x)\} \quad f'(x) = ?$$



GAN背后的数学

GAN的训练算法

- 初始化 G_0
- 找到 D_0^* 使得 $V(G_0, D)$ 最大化



$V(G_0, D_0^*)$ 是分布 $P_{G_0}(x)$ 和 $P_{data}(x)$ 的JS散度

- $\theta_{G_1} \leftarrow \theta_{G_0} - \eta \frac{\partial V(G_0, D_0^*)}{\partial \theta_G} \rightarrow G_1$ 减小了 $P_{G_0}(x)$ 和 $P_{data}(x)$ 的JS散度
- 找到 D_1^* 使得 $V(G_1, D)$ 最大化

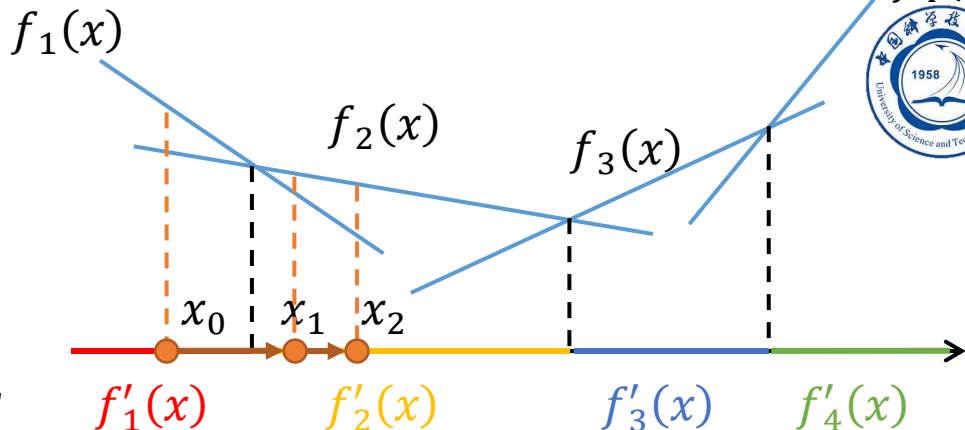
$V(G_1, D_1^*)$ 是分布 $P_{G_1}(x)$ 和 $P_{data}(x)$ 的JS散度

- $\theta_{G_2} \leftarrow \theta_{G_1} - \eta \frac{\partial V(G_1, D_1^*)}{\partial \theta_G} \rightarrow G_2$ 减小了 $P_{G_1}(x)$ 和 $P_{data}(x)$ 的JS散度
- ...

GAN背后的数学

GAN的训练算法

- 初始化 G_0
- 找到 D_0^* 使得 $V(G_0, D)$ 最大化



$V(G_0, D_0^*)$ 是分布 $P_{G_0}(x)$ 和 $P_{data}(x)$ 的JS散度

$$\theta_{G_1} \leftarrow \theta_{G_0} - \eta \frac{\partial V(G_0, D_0^*)}{\partial \theta_G}$$



G_1

减小了 $P_{G_0}(x)$ 和 $P_{data}(x)$ 的JS散度

$$\theta_{G_2} \leftarrow \theta_{G_1} - \eta \frac{\partial V(G_1, D_0^*)}{\partial \theta_G}$$



G_2

每次 D^* 的求解需要不断迭代直到收敛，而 G 只更新一次， G 的更新效率太低

$$\theta_{G_3} \leftarrow \theta_{G_2} - \eta \frac{\partial V(G_2, D_0^*)}{\partial \theta_G}$$



G_3

...

?

只能学习率较小且更新次数不多

否则 $V(G_i, D_0^*)$ 无法准确衡量 $P_{G_i}(x)$ 和 $P_{data}(x)$ 的JS散度

生成对抗网络的算法实现技巧

- 随机初始化判别器参数 θ_d 和生成器参数 θ_g
- While not convergent:

出于效率考虑, 对 D^* 的求解只迭代了 K 步, 因而只得到 $\max_D V(G, D)$ 的下界

更新D

重复
K次

- 从数据库中采样 m 样本 $\{x^1, x^2, \dots, x^m\}$
- 从某个分布中采样 m 个噪声样本 $\{z^1, z^2, \dots, z^m\}$
- 用生成器生成样本 $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- 更新判别器参数 θ_d , 通过最大化如下目标

$$\tilde{\mathcal{L}} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(\tilde{x}^i))$$

$$\theta_d \leftarrow \theta_d + \eta \nabla \tilde{\mathcal{L}}(\theta_d)$$

更新G

只做
一次

- 从某个分布中采样 m 个噪声样本 $\{z^1, z^2, \dots, z^m\}$
- 更新生成器参数 θ_g , 通过最小化如下目标

$$\tilde{\mathcal{L}} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^i)))$$

$$\theta_g \leftarrow \theta_g - \eta \nabla \tilde{\mathcal{L}}(\theta_g)$$

生成对抗网络的算法实现技巧

在更新G的时候，最小化如下目标

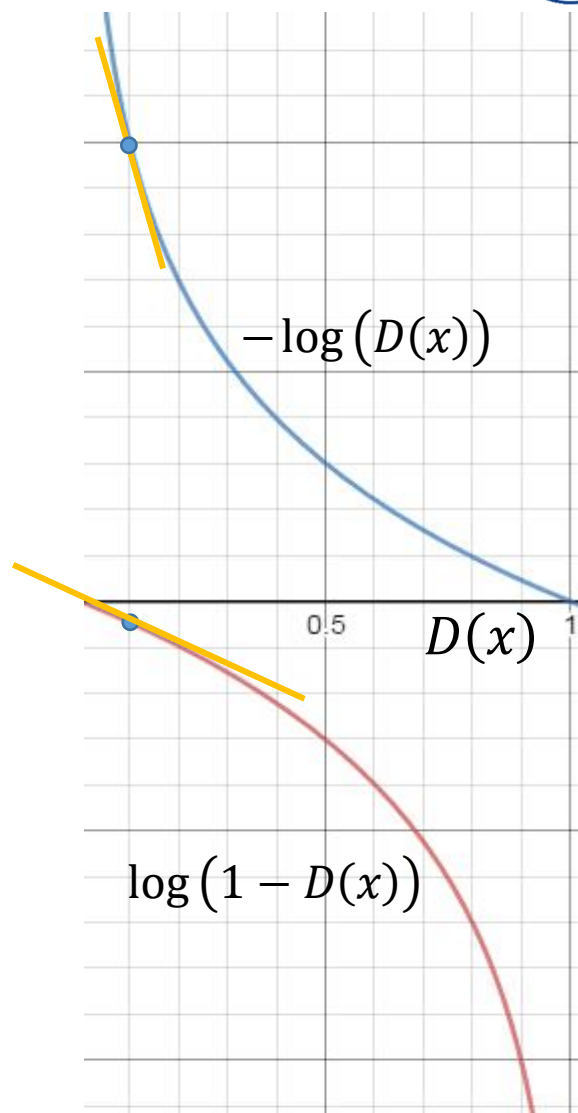
$$\tilde{\mathcal{L}} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^i)))$$

但是在初始时候会很慢，因为G比较差，生成的样本很容易被D分辨为假的，也就是说 $D(G(z^i))$ 是趋于0的， $\log(1 - D(G(z^i)))$ 的梯度很小

在更新G的时候，最小化如下目标

$$\tilde{\mathcal{L}} = -\frac{1}{m} \sum_{i=1}^m \log (D(G(z^i)))$$

Log D技巧



生成对抗网络的算法实现技巧

- 随机初始化判别器参数 θ_d 和生成器参数 θ_g
- While not convergent:

出于效率考虑, 对 D^* 的求解只迭代了K步, 因而只得到 $\max_D V(G, D)$ 的下界

更新D

重复
K次

- 从数据库中采样 m 样本 $\{x^1, x^2, \dots, x^m\}$
- 从某个分布中采样 m 个噪声样本 $\{z^1, z^2, \dots, z^m\}$
- 用生成器生成样本 $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- 更新判别器参数 θ_d , 通过最大化如下目标

$$\tilde{\mathcal{L}} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(\tilde{x}^i))$$

$$\theta_d \leftarrow \theta_d + \eta \nabla \tilde{\mathcal{L}}(\theta_d)$$

更新G

只做
一次

- 从某个分布中采样 m 个噪声样本 $\{z^1, z^2, \dots, z^m\}$
- 更新生成器参数 θ_g , 通过最小化如下目标

$$\tilde{\mathcal{L}} = -\frac{1}{m} \sum_{i=1}^m \log (D(G(z^i)))$$

$$\theta_g \leftarrow \theta_g - \eta \nabla \tilde{\mathcal{L}}(\theta_g)$$

GAN背后的数学

- 不用Log D技巧，生成网络的目标函数为

$$V(G, D^*) = 2JS(P_{data}(\mathbf{x}) || P_G(x)) - 2\log 2$$




- 但用Log D技巧时，生成网络的目标函数为

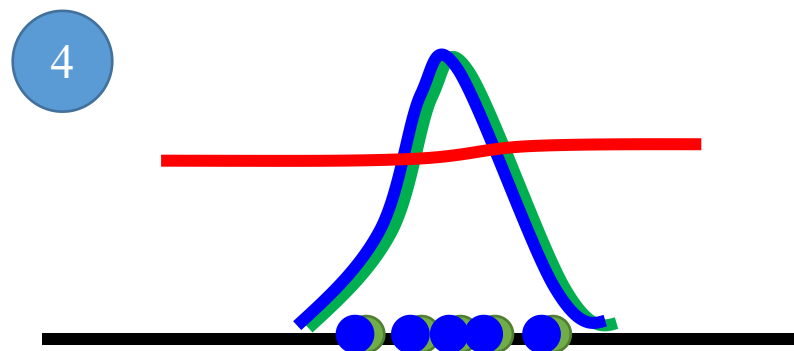
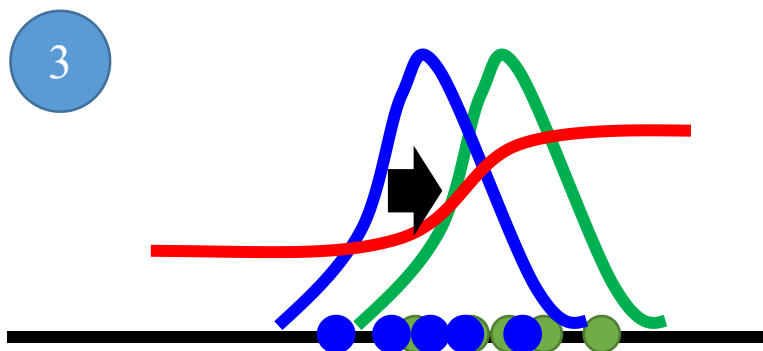
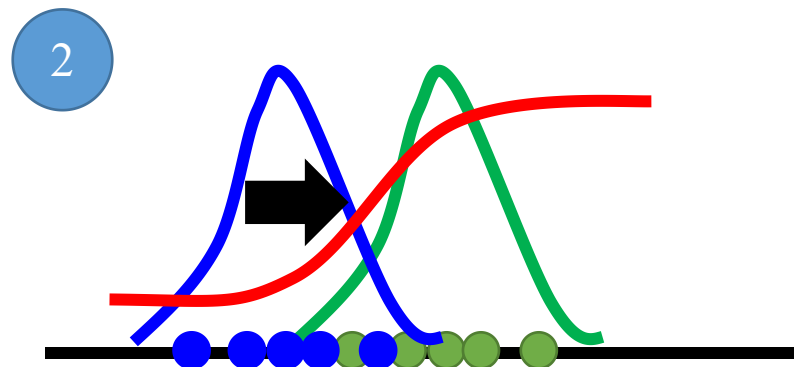
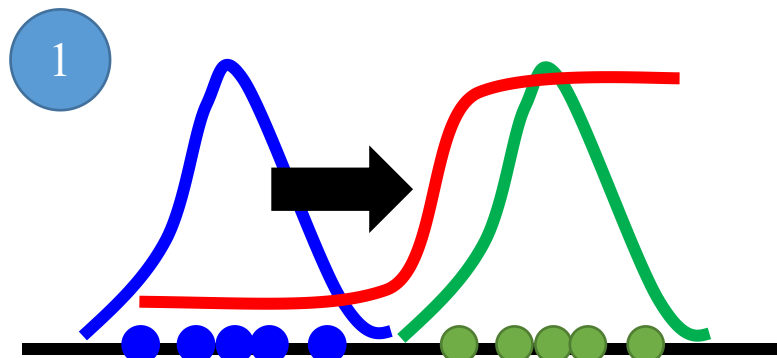
$$V(G, D^*) \propto KL(P_G(\mathbf{x}) || P_{data}(\mathbf{x})) - 2JS(P_{data}(\mathbf{x}) || P_G(x))$$

同时最小化生成分布与真实分布的KL散度，并最大化两者的JS散度，

这和直觉相冲突，在数值上则会导致梯度不稳定

GAN训练的直观理解

 P_{data}

 P_{D^c}





谢谢

