

01 GAN

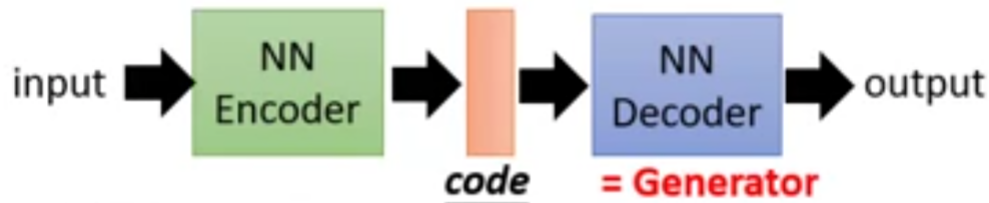
- 效果不佳
- 萌生想法
- 遇到问题
- 可行性证明
- 算法实现流程
- 针对GAN的认识

效果不佳

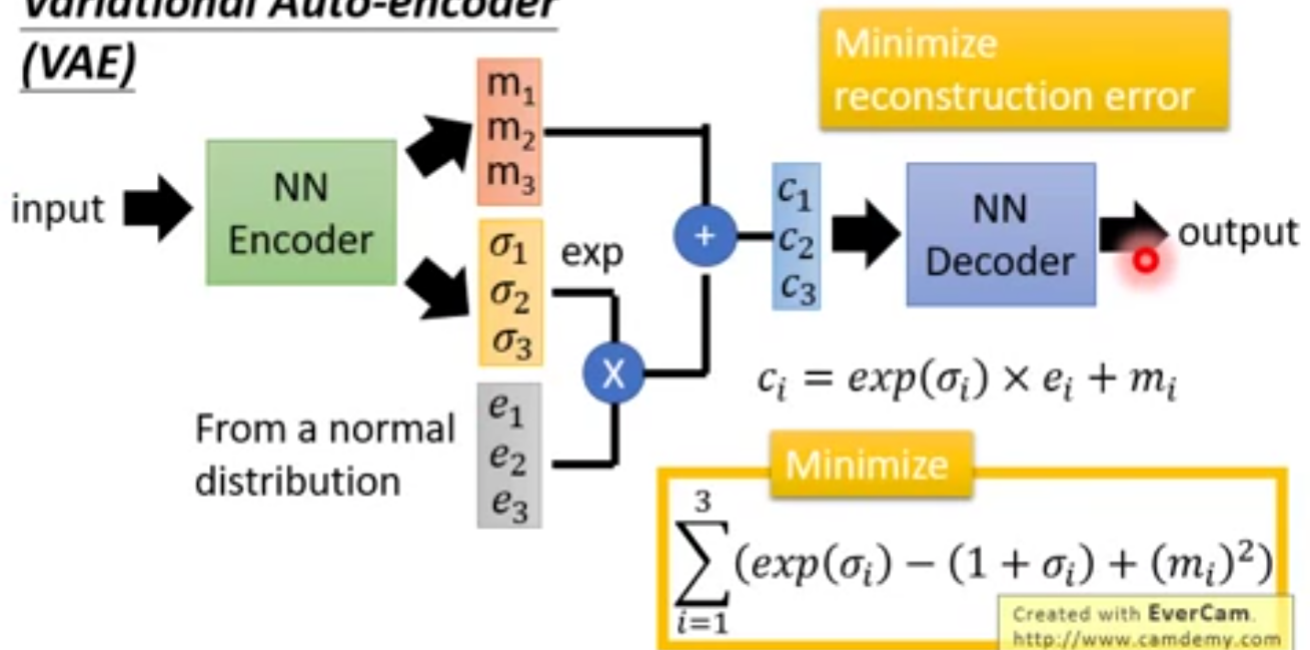
- 普通自编码器通过训练数据学习到的是某个确定的函数；
- 变分自编码器（VAE）通过训练数据学习到的是参数的概率分布。[万字长文带你了解变分自编码器 VAEs](#)



Auto-encoder



Variational Auto-encoder (VAE)



萌生想法

对抗网络希望通过改变参数 θ ，使生成模型概率分布 $p_{model}(x; \theta)$ 能够逼近真实数据概率分布 $p_{data}(x)$ 。

即通过 $\theta^* = \arg \max_{\theta} L = \arg \max_{\theta} \prod_{i=1}^m p_{model}(x^{(i)}; \theta)$ 使得 $p_{model}(x; \theta)$ 逼近 $p_{data}(x)$ 。

$$\begin{aligned}
\theta^* &= \arg \max_{\theta} \prod_{i=1}^m p_{model}(x^{(i)}; \theta) = \arg \max_{\theta} \log \prod_{i=1}^m p_{model}(x^{(i)}; \theta) = \arg \max_{\theta} \sum_{i=1}^m \log p_{model}(x^{(i)}; \theta) \\
&= \arg \max_{\theta} E_{x \sim p_{data}} \log p_{model}(x; \theta) = \arg \max_{\theta} \int p_{data}(x) \log p_{model}(x; \theta) dx \\
&= \arg \max_{\theta} \left(\int p_{data}(x) \log p_{model}(x; \theta) dx - \int p_{data}(x) \log p_{data}(x) dx \right) \\
&= \arg \max_{\theta} \int p_{data}(x) (\log p_{model}(x; \theta) - \log p_{data}(x)) dx \\
&= \arg \max_{\theta} \int p_{data}(x) \log \frac{p_{model}(x; \theta)}{p_{data}(x)} dx = \arg \max_{\theta} \left(- \int p_{data}(x) \log \frac{p_{data}(x)}{p_{model}(x; \theta)} dx \right) \\
&= \arg \min_{\theta} \left(\int p_{data}(x) \log \frac{p_{data}(x)}{p_{model}(x; \theta)} dx \right) \\
&= \arg \min_{\theta} KL(p_{data}(x) || p_{model}(x))
\end{aligned}$$

(1)

即 $\theta^* = \arg \min_{\theta} KL(p_{data}(x) || p_{model}(x))$

知识点tip:

KL散度：一种计算概率分布之间相似程度的计算方法。

在假定为连续随机变量的前提下，且两个概率分布分别为 P 和 Q ，则

$$KL(P || Q) = \int p(x) \log \frac{p(x)}{q(x)} dx。$$

注意：非负性，对称性。

遇到问题

问题:

生成模型 $x = G(z)$

输入数据 z 的概率分布函数 已知

但生成函数 G 是一种神经网络的形式 $\Rightarrow P_{model}(x) \Rightarrow L$

解决:

生成对抗网络中的代价函数

判别器D: $J^{(D)}(D, G) = -\frac{1}{2}E_{x \sim p_{data}} [\log D(x)] - \frac{1}{2}E_{z \sim p_z} [\log(1 - D(G(z)))]$

生成器G: $J^{(G)}(D, G) = \frac{1}{2}E_{x \sim p_{data}} [\log D(x)] + \frac{1}{2}E_{z \sim p_z} [\log(1 - D(G(z)))]$

价值函数 $V(D, G)$ 表示 $J^{(G)}$ 和 $J^{(D)}$ 。

$$V(D, G) = E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p_z} [\log(1 - D(G(z)))]$$

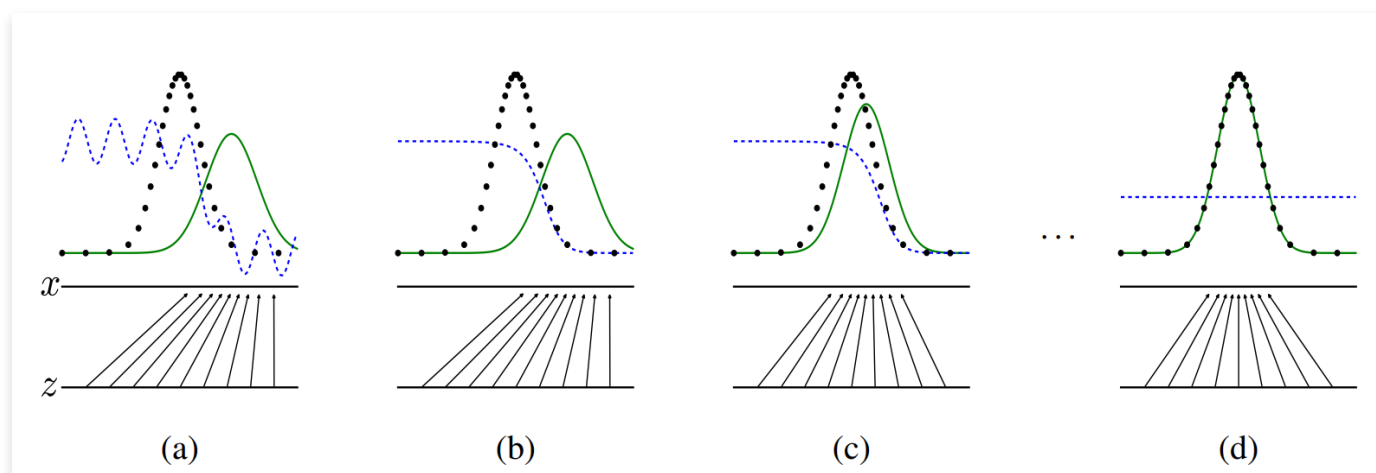
$$J^{(D)} = -\frac{1}{2}V(D, G) \quad J^{(G)} = \frac{1}{2}V(D, G)$$

在生成对抗网络中，我们要计算的纳什平衡点正是要寻找一个生成器G和判别器D，使得各自的代价函数最小，即希望找到了一个 $V(D, G)$ 对于生成器来说最小而对判别器来说最大。公式：

$\arg \min_G \max_D V(D, G)$ 。例：鞍点。

可行性证明

当生成数据的分布 $p_g(x)$ 趋近于真实数据分布 $p_{data}(x)$ 时，D网络输出的概率 $D^*(x)$ 会趋近于 $\frac{1}{2}$ 。这也是最终希望达到的训练结果，这时候G和D网络也就达到了一个平衡状态。



z 是随机噪声， x 是训练数据中的真实数据data。

分别求出理想的判别器 D^* 和生成器 G^* 。

$$D^* = \arg \max_D V(D, G^*) = \arg \min_G \max_D V(D, G) = \arg \min_G V(D^*, G) = G^* = \frac{1}{2}$$

• 最优判别器

假设生成器G是固定的，令 $G(z) = x$ ，

$$\begin{aligned} V(G, D) &= E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p_z} [\log(1 - D(G(z)))] \\ &= E_{x \sim p_{data}} [\log D(x)] + E_{x \sim p_g} [\log(1 - D(x))] \end{aligned}$$

$$\begin{aligned} \text{则} \quad &= \int p_{data}(x) \log D(x) dx + \int p_g(x) \log(1 - D(x)) dx \quad (2) \\ &= \int p_{data}(x) \log D(x) + p_g(x) \log(1 - D(x)) dx \end{aligned}$$

希望对 $f(x) = p_{data}(x) \log D(x) + p_g(x) \log(1 - D(x))$ ，无论 x 取何值都能最大。

$$\frac{df(x)}{dD(x)} = \frac{p_{data}(x)}{D(x)} - \frac{p_g(x)}{1 - D(x)} = 0 \Rightarrow D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \in (0, 1)$$

• 最优生成器

代入 $D^*(x)$ ，即判别器D是固定的，

$$\begin{aligned} C(G) &= \max_D V(G, D) = V(G, D^*) \\ &= E_{x \sim p_{data}} [\log D^*(x)] + E_{z \sim p_z} [\log(1 - D^*(G(z)))] \\ &= E_{x \sim p_{data}} [\log D^*(x)] + E_{x \sim p_g} [\log(1 - D^*(x))] \\ &= E_{x \sim p_{data}} [\log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}] + E_{x \sim p_g} [\log \frac{p_g(x)}{p_{data}(x) + p_g(x)}] \\ &= E_{x \sim p_{data}} [\log \frac{p_{data}(x)/2}{(p_{data}(x) + p_g(x))/2}] + E_{x \sim p_g} [\log \frac{p_g(x)/2}{(p_{data}(x) + p_g(x))/2}] \\ &= E_{x \sim p_{data}} [\log \frac{p_{data}(x)}{(p_{data}(x) + p_g(x))/2}] + E_{x \sim p_{data}} [-\log 2] + E_{x \sim p_g} [\log \frac{p_g(x)}{(p_{data}(x) + p_g(x))/2}] + E_{x \sim p_g} [-\log 2] \\ &= \int p_{data}(x) \log \frac{p_{data}(x)}{(p_{data}(x) + p_g(x))/2} dx + \int p_g(x) \log \frac{p_g(x)}{(p_{data}(x) + p_g(x))/2} dx - \log(4) \\ &= -\log(4) + KL\left(p_{data}(x) \parallel \frac{p_{data}(x) + p_g(x)}{2}\right) + KL\left(p_g(x) \parallel \frac{p_{data}(x) + p_g(x)}{2}\right) \\ &= -\log(4) + 2 \cdot JSD(p_{data} \parallel p_g) \end{aligned}$$

$$\text{即 } C(G) = -\log(4) + 2 \cdot JSD(p_{data} \parallel p_g)$$

当且仅当 $p_{data} = p_g$ 的时候，上式可以取得全局最小值 $C^* = -\log(4)$ 。

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

算法实现流程

Algorithm Initialize θ_d for D and θ_g for G

- In each training iteration:

Learning
D

- Sample m examples $\{x^1, x^2, \dots, x^m\}$ from database
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Obtaining generated data $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}, \tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to maximize
 - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(\tilde{x}^i))$
 - $\theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$

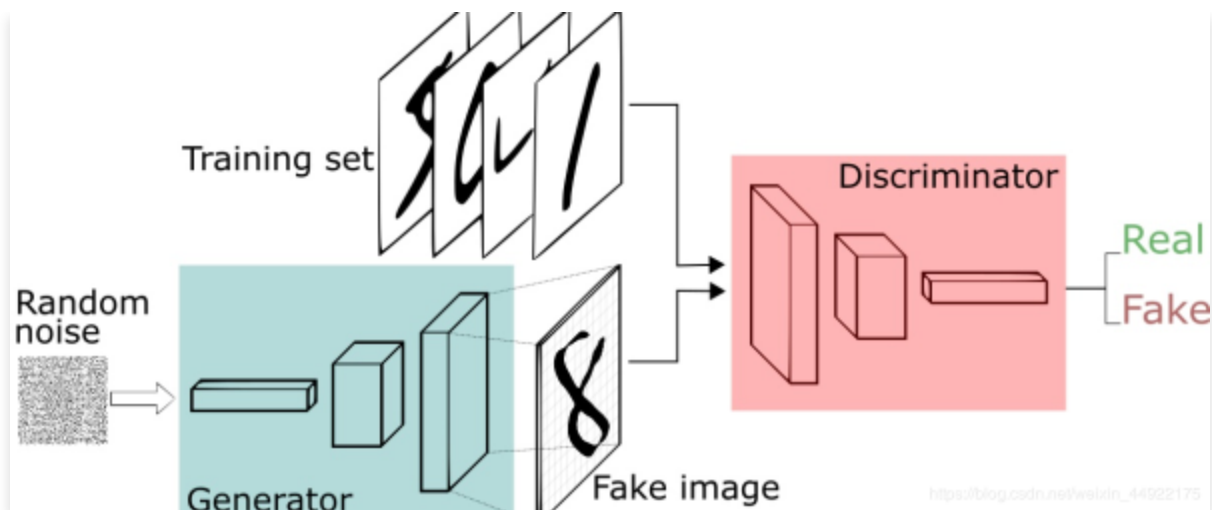
Learning
G

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Update generator parameters θ_g to maximize
 - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log (D(G(z^i)))$
 - $\theta_g \leftarrow \theta_g + \eta \nabla \tilde{V}(\theta_g)$

Created with EverCam.

针对GAN的认识

1. 由判别模型和生成模型两部分组成。
2. 生成模型负责生成假样本，欺骗判别器。
3. 判别模型来分辨真样本和假样本。
4. 两个模型交替训练，类似对抗博弈的过程。
5. 训练判别器，固定生成器，使用反向传播。
6. 固定判别器，训练生成器，使用反向传播。



论文给出的算法实现：

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

训练技巧和方法：

1. 首先G和D是同步训练，但是两者训练次数不一样，通常是D网络训练k次后，G训练一次。主要原因是GAN刚开始训练时候会很不稳定；
2. D的训练是同时输入真实数据和生成数据来计算loss，而不是采用交叉熵（cross entropy）分开计算。不采用cross entropy的原因是这会让 $D(G(z))$ 变为0，导致没有梯度提供给G更新，而现在GAN的做法是会收敛到 $\frac{1}{2}$ ；
3. 实际训练的时候，作者是采用 $-\log(D(G(z)))$ 来代替 $\log(1 - D(G(z)))$ ，这是希望在训练初始就可以加大梯度信息，这是因为初始阶段D的分类能力会远大于G生成足够真实数据的能力，但这种修改也将让整个GAN不再是一个完美的零和博弈。

