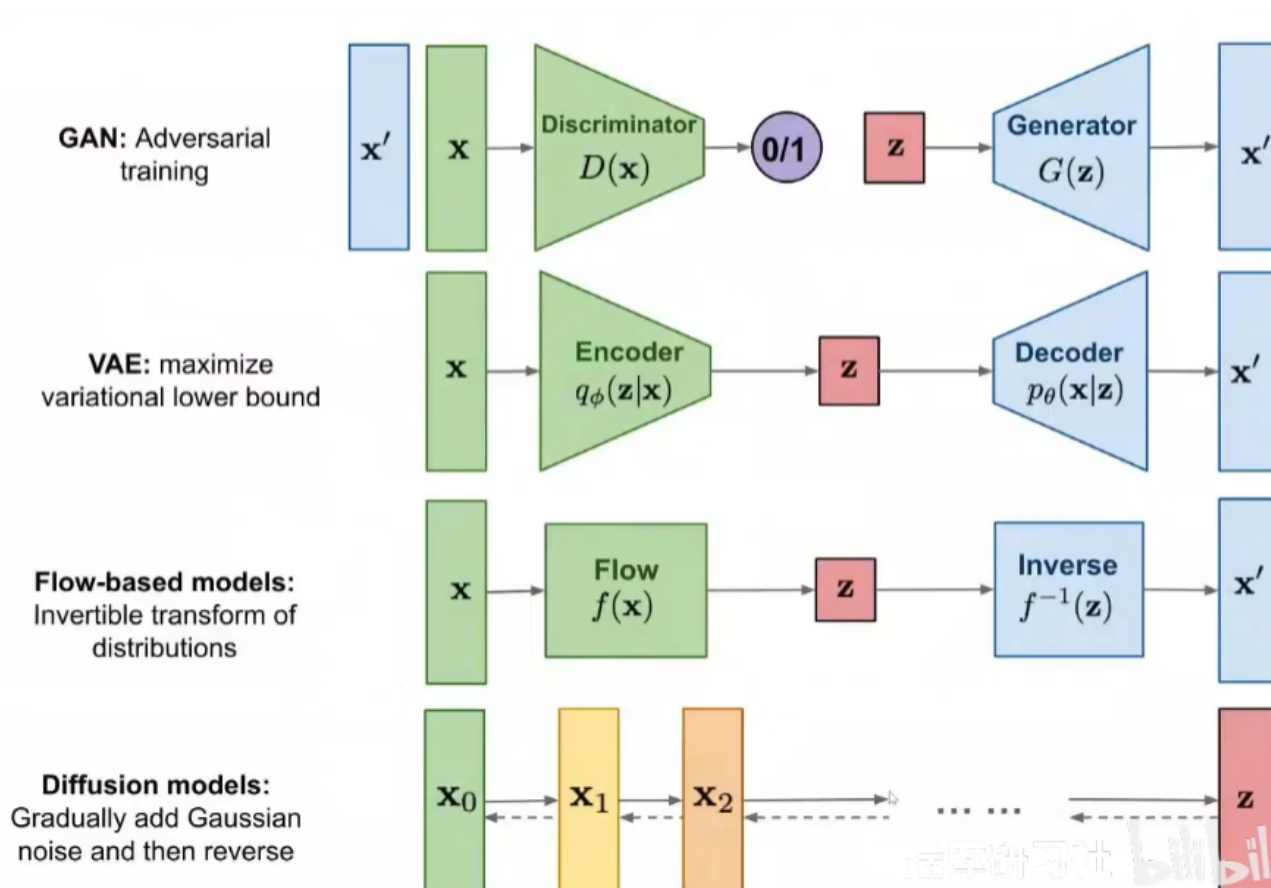


# vae

author: doomx

time: 2023.9.23

## 生成模型的研究背景(fig1)



- GAN model由于对抗性训练的性质，有潜在的不稳定的训练和较低多样性。
- VAE 依赖于代用损失，需要 maximize ELBO (variational lower bound)，不是直接去做 $p(x)$ ，而是用一个东西去代替。
- 流模型(Flow-based model) 必须使用专门的架构来构建逆变换

对于许多模态，我们可以将我们观察到的数据视为由**相关的**看不见的**隐变量**表示或生成，我们用**随机变量** $z$ 表示。表达这个想法的最佳直觉是通过柏拉图的洞穴寓言。在这个寓言中，一群人一生都被锁在一个山洞里，只能看到投射在他们面前的墙上的二维阴影，这是由看不见的三维物体在火前经过而产生的。对于这样的人来说，他们所观察到的一切，实际上都是由他们永远看不到的更高维度的抽象概念决定的。

## ELBO-Evidence Lower Bound

将隐变量 $z$ 和我们观察到的数据 $x$ ，用联合分布 $p(x, z)$ 建模。用最大似然“likelihood-based”的方法，学习一个模型以最大化所有观察到的 $x$ 的likelihood  $p(x)$ 。接下来，由两种方法去重构数据

$p(x)$ 。

显式地边缘化隐变量 $z$ ，我们不关心 $z$ ，我们用一个积分把它积掉，得到我们所关心的 $p(x)$ ：  
eq.1

$$p(x) = \int p(x, z) dz$$

或者，我们可以用概率链法则：  
eq.2

$$p(x) = \frac{p(x, z)}{p(z | x)} = \frac{p(x | z) p(z)}{p(z | x)}$$

上面两个公式，我们想去直接算 $p(x)$ 是做不到的，首先 $z$ 的分布我们是很难去处理得到的，所以接下来的思路是用一个代用的损失去替代，也就是提出一个东西，叫ELBO(Evidence Lower Bound)，顾名思义，他是**证据**的下界。

这种情况下，evidence被量化为 $x$ 的log likelihood，这样，maximize ELBO就成为了优化隐变量模型的代理任务。在最好的情况下，当ELBO被参数化且完美优化时，它与证据完全等价。

我们定义一个变分分布近似， $q_\theta(z | x)$ 是一个**灵活的**变分分布近似， $\theta$ 是我们寻求优化的参数，我们通过调整 $\theta$ 来增加下限以最大化ELBO。通过这种方式，**用一个变分分布近似去替换隐变量 $z$ 的分布**，从而获得可用于对真实数据分布 $p(x)$ 进行建模并从中采样的模型，从而学习 $\theta$ 。  
eq.3

$$\log p(x) = \int p(x, z) dz \quad (\text{Apply eq.1})$$

eq.4

$$= \log \int \frac{p(x, z) q_\theta(z | x)}{q_\theta(z | x)} dz \quad (\text{Multiply by } 1 = \frac{q_\theta(z | x)}{q_\theta(z | x)})$$

eq.5

$$= \log E_{q_\theta(z | x)} \left[ \frac{p(x, z)}{q_\theta(z | x)} \right] \quad (\text{Definition of Expectation})$$

eq.6

$$\geq E_{q_\theta(z | x)} \left[ \log \frac{p(x, z)}{q_\theta(z | x)} \right] \quad (\text{Apply Jensen's Inequality})$$

Jensen's Inequality(fig2&3)

## 凸函数

凸函数是一个定义在某个向量空间的凸子集  $C$  (区间) 上的实值函数  $f$ , 如果在其定义域  $C$  上的任意两点  $x_1, x_2$ ,  $0 \leq t \leq 1$ , 有

$$tf(x_1) + (1 - t)f(x_2) \geq f(tx_1 + (1 - t)x_2) \quad (1)$$

也就是说凸函数任意两点的割线位于函数图形上方, **这也是Jensen不等式的两点形式。**

## Jensen不等式

若对于任意点集  $\{x_i\}$ , 若  $\lambda_i \geq 0$  且  $\sum_i \lambda_i = 1$ , 使用**数学归纳法**, 可以证明凸函数  $f(x)$  满足:

$$f(\sum_{i=1}^M \lambda_i x_i) \leq \sum_{i=1}^M \lambda_i f(x_i) \quad (2)$$

公式(2)被称为 Jensen 不等式, **它是式(1)的泛化形式。**

**在概率论中**, 如果把  $\lambda_i$  看成取值为  $x_i$  的离散变量  $x$  的概率分布, 那么公式(2)就可以写成

$$f(E[x]) \leq E[f(x)]$$

其中,  $E[\cdot]$  表示期望。

对于连续变量, Jensen不等式给出了积分的凸函数值和凸函数的积分值间的关系:

$$f(\int xp(x)dx) \leq \int f(x)p(x)dx$$

用了Jensen's Inequality, 从等式变到了不等式, 需要找到缺失的东西。

现在我们尝试用eq.2去推导ELBO:

eq.7

$$\log p(x) = \log p(x) \int q_\theta(z | x) dz \quad (\text{Multiply by } 1 = \int q_\theta(z | x) dz)$$

eq.8

$$= \int q_\theta(z | x)(\log p(x)) dz \quad (\text{Bring evidence into integral})$$

eq.9

$$= E_{q_{\theta}(z|x)} [\log p(x)] \quad (\text{Definition of Expectation})$$

eq.10

$$= E_{q_{\theta}(z|x)} \left[ \log \frac{p(x, z)}{p(z|x)} \right] \quad (\text{Apply equation 2})$$

eq.11

$$= E_{q_{\theta}(z|x)} \left[ \log \frac{p(x, z) q_{\theta}(z|x)}{p(z|x) q_{\theta}(z|x)} \right] \quad (\text{Multiply by } 1 = \frac{q_{\theta}(z|x)}{q_{\theta}(z|x)})$$

eq.12

$$= E_{q_{\theta}(z|x)} \left[ \log \frac{p(x, z)}{q_{\theta}(z|x)} \right] + E_{q_{\theta}(z|x)} \left[ \log \frac{q_{\theta}(z|x)}{p(z|x)} \right] \quad (\text{Split the Expectation})$$

eq.13

$$= E_{q_{\theta}(z|x)} \left[ \log \frac{p(x, z)}{q_{\theta}(z|x)} \right] + D_{KL}(q_{\theta}(z|x) || p(z|x)) \quad (\text{Definition of KL Divergence})$$

eq.14:  $\log p(x) \geq \text{ELBO}$

$$\geq E_{q_{\theta}(z|x)} \left[ \log \frac{p(x, z)}{q_{\theta}(z|x)} \right] \quad (\text{KL Divergence always } \geq 0)$$

我们这时候也可以得知，用了Jensen's Inequation缺失的就是KL Divergence。

这时可以得知，ELBO项为下界的原因是，evidence与ELBO之间的差异是一个严格非负的KL Divergence，因此ELBO的值永远不会超过evidence。

我们想要优化我们的变分后验（变分分布近似） $q_{\theta}(z|x)$ 的参数以逼近真实后验分布 $p(z|x)$ 。 $\log p(x)$ 为常数，因此ELBO和KL Divergence的和为常数，使ELBO项变大等于使KL Divergence项变小，ELBO可以作为近似学习后验分布的代理任务。

我们需要优化一系列由  $\theta$  参数化的后验分布（变分近似分布） $q_{\theta}(z|x)$ ，他被称为自编码器（auto-encoder）。接下来我们进一步解析ELBO：

eq.15

$$E_{q_{\theta}(z|x)} \left[ \log \frac{p(x, z)}{q_{\theta}(z|x)} \right] = E_{q_{\theta}(z|x)} \left[ \log \frac{p_{\phi}(x|z) p(z)}{q_{\theta}(z|x)} \right] \quad (\text{Chain Rule of Probability})$$

eq.16

$$= E_{q_{\theta}(z|x)} [\log p_{\phi}(x|z)] + E_{q_{\theta}(z|x)} \left[ \log \frac{p(z)}{q_{\theta}(z|x)} \right] \quad (\text{Spilt the Expectation})$$

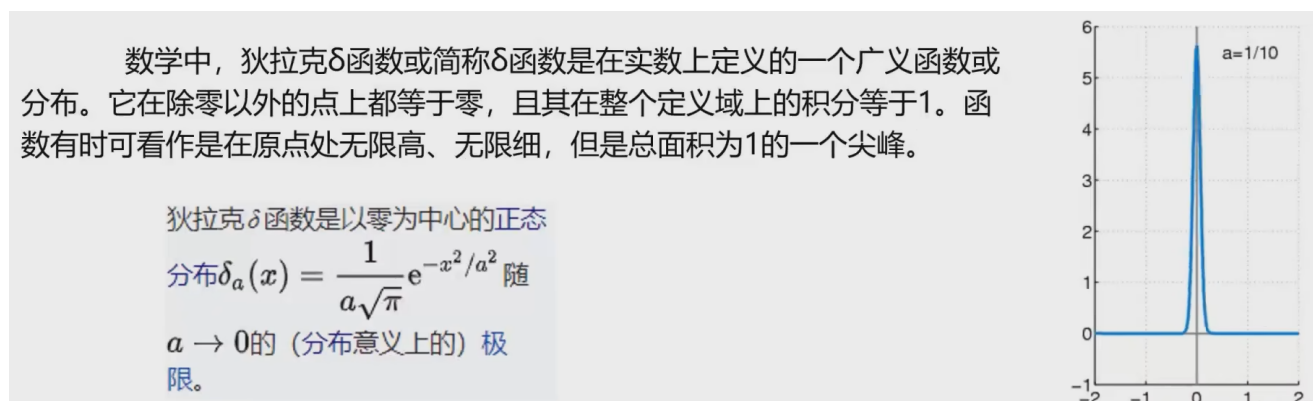
eq.17

$$= E_{q_{\theta}(z|x)} [\log p_{\phi}(x|z)] - D_{KL}(q_{\theta}(z|x) || p(z)) \quad (\text{Definition of KL Divergence})$$

在这种情况下，我们学习了一个变分近似分布  $q_{\theta}(z|x)$ ，它被视为**编码器**，可以将输入转换为

可能的隐变量的分布。同时，我们学习了一个**确定性函数**  $p_\phi(x | z)$ ，将给定的隐变量 $z$ 转换为观测值 $x$ ，可以将其当作**解码器**。

狄拉克  $\delta$  函数(fig4)



eq17两个项的直观描述：第一项衡量**解码器**从变分分布中重构的概率，第二项衡量学到的变分分布与隐变量的先验分布的相似程度。因此，最大化ELBO等同于最大化其第一项并最小化其第二项。最小化第二项可以鼓励编码器学习一个优质的高斯分布，而不是崩溃成一个狄拉克 $\delta$ 函数。所以maximize ELBO作为代理任务既可以更好地学习编码器  $q_\theta(z | x)$ ，也可以更好地学习解码器  $p_\phi(x | z)$ 。

## VAE：如何在 $\delta$ 和 $\theta$ 上联合优化 ELBO

VAE的编码器通常选择具有对角协方差的多元高斯进行建模，而先验通常被选择为标准的多元高斯：

eq.18

$$q_\theta(z | x) = \mathcal{N}(z; \mu_\theta(x), \sigma_\theta^2(x)I)$$

eq.19

$$p(z) = \mathcal{N}(z; 0, I)$$

最后，ELBP的KL Divergence可以用解析计算，而重建项可以用蒙特卡洛估计来近似。目标函数可以改写为(约等)：

eq.20

$$\arg \max \sum_{l=1}^L \log p_\phi(x | z^{(l)}) - D_{KL}(q_\theta(z | x) || p(z))$$

重参数化(fig5)

对于数据集中的数据 $x$ ，其中隐变量  $\{z^{(l)}\}_{l=1}^L$  是从  $q_\phi(z|x)$  中采样的。然而，在这个默认设置中出现了一个问题：计算损失的每个  $z^{(l)}$  都是随机抽样生成的，该过程通常是不可微的。幸运的是，当  $q_\phi(z|x)$  被设计为对某些分布（例如多元高斯分布）进行建模时，可以通过重新参数化技巧 (*reparameterization trick*) 来解决这个问题。

重参数化技巧将随机变量重写为噪音变量的确定函数；这允许通过梯度下降优化非随机项。例如，从任意均值 $\mu$ 和方差  $\sigma^2$  的正态分布  $x \sim \mathcal{N}(x; \mu, \sigma^2)$  采样得到的样本可以重写为：

$$x = \mu + \sigma\epsilon \quad \text{with } \epsilon \sim \mathcal{N}(\epsilon; 0, I)$$

