

24 生成扩散模型漫谈（十九）：作为扩散ODE的GAN

Jun By 苏剑林 | 2023-06-24 | 31058位读者 引用

在文章《生成扩散模型漫谈（十六）：W距离 \leq 得分匹配》中，我们推导了Wasserstein距离与扩散模型得分匹配损失之间的一个不等式，表明扩散模型的优化目标与WGAN的优化目标在某种程度上具有相似性。而在本文，我们将探讨《MonoFlow: Rethinking Divergence GANs via the Perspective of Wasserstein Gradient Flows》中的研究成果，它进一步展示了GAN与扩散模型之间的联系：**GAN实际上可以被视为在另一个时间维度上的扩散ODE！**

这些发现表明，尽管GAN和扩散模型表面上是两种截然不同的生成式模型，但它们实际上存在许多相似之处，并在许多方面可以相互借鉴和参考。

思路简介

我们知道，GAN所训练的生成器是从噪声 z 到真实样本的一个直接的确定性变换 $g_\theta(z)$ ，而扩散模型的显著特点是“渐进式生成”，它的生成过程对应于从一系列渐变的分布 $p_0(x_0), p_1(x_1), \dots, p_T(x_T)$ 中采样（注：在前面十几篇文章中， x_T 是噪声， x_0 是目标样本，采样过程是 $x_T \rightarrow x_0$ ，但为了便于下面的表述，这里反过来改为 $x_0 \rightarrow x_T$ ）。看上去确实找不到多少相同之处，那怎么才能将两者联系起来呢？

很明显，如果想要从扩散模型的视角理解GAN，那么就要想办法构造出一系列渐变的分布出来。生成器 $g_\theta(z)$ 本身就是一个一步到位的变换，不存在渐变，然而我们知道模型的优化是渐变的，可否用参数 θ 的历史轨迹 θ_t 来构建这一系列渐变分布呢？具体来说，假设生成器初始化为 θ_0 ，经过 T 步对抗训练后得到最优参数 θ_T ，训练过程的中间参数为 $\theta_1, \theta_2, \dots, \theta_{T-1}$ ，那么我们定义 $x_t = g_{\theta_t}(z)$ ，不就定义了一系列渐变的 x_0, x_1, \dots, x_T ，从而也就定义了渐变的分布 $p_0(x_0), p_1(x_1), \dots, p_T(x_T)$ 了？

如果这个思路可行的话，那么GAN就可以诠释为梯度下降的（虚拟）时间维度上的扩散模型！下面我们就沿着这个思路进行探索。

梯度之流

首先，我们需要重温上一篇文章《梯度流：探索通向最小值之路》关于Wasserstein梯度流的结果：它指出方程

$$\frac{\partial q_t(\mathbf{x})}{\partial t} = -\nabla_{\mathbf{x}} \cdot (q_t(\mathbf{x}) \nabla_{\mathbf{x}} \log r_t(\mathbf{x})) \quad (1)$$

在最小化 $p(\mathbf{x})$ 和 $q_t(\mathbf{x})$ 的KL散度，即 $\lim_{t \rightarrow \infty} q_t(\mathbf{x}) = p(\mathbf{x})$ ，这里 $r_t(\mathbf{x}) = \frac{p(\mathbf{x})}{q_t(\mathbf{x})}$ 。如果 $p(\mathbf{x})$ 代表真实样本的分布，那么如果能实现从 $q_t(\mathbf{x})$ 采样的话，那么逐渐推到 $t \rightarrow \infty$ 时，就可以实现从 $p(\mathbf{x})$ 采样了。根据《测试函数法推导连续性方程和Fokker-Planck方程》，从 $q_t(\mathbf{x})$ 采样可以通过下述ODE实现：

$$\frac{d\mathbf{x}}{dt} = \nabla_{\mathbf{x}} \log r_t(\mathbf{x}) \quad (2)$$

然而，上式中的 $r_t(\mathbf{x})$ 是未知的，所以我们还无法通过上式进行采样，需要先想办法估算 $r_t(\mathbf{x})$ 。

判别估计

这时候登场的是GAN的判别器。以最早的Vanilla GAN为例，它的训练目标是

$$\max_D \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\log \sigma(D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\log(1 - \sigma(D(\mathbf{x})))] \quad (3)$$

这里的 D 是判别器， $\sigma(t) = 1/(1 + e^{-t})$ 是sigmoid函数， $p(\mathbf{x})$ 是真样本的分布， $q(\mathbf{x})$ 是假样本的分布。可以证明（不清楚的读者可以参考《RSGAN：对抗模型中的“图灵测试”思想》中的“补充证明”一节），上式中判别器 D 的理论最优解是

$$D(\mathbf{x}) = \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad (4)$$

更一般化的f-GAN（参考《f-GAN简介：GAN模型的生产车间》、《Designing GANs：又一个GAN生产车间》）结果会稍有不同，但可以证明的是它们判别器的理论最优解都

是 $\frac{p(\mathbf{x})}{q(\mathbf{x})}$ 的函数。也就是说，只要我们可以实现从 $p(\mathbf{x})$ 和 $q_t(\mathbf{x})$ 中采样，那么通过GAN的判别器训练(3)就可以估算出 $r_t(\mathbf{x}) = \frac{p(\mathbf{x})}{q_t(\mathbf{x})}$ 出来。

向前一步

这时候可能有读者疑惑：这不就进入“鸡生蛋、蛋生鸡”的循环论证了吗？我们估算 $r_t(\mathbf{x})$ 不就是为了利用式(2)实现从 $q_t(\mathbf{x})$ 中采样吗？现在你又假设能从 $q_t(\mathbf{x})$ 采样才来估算 $r_t(\mathbf{x})$ ？不着急，经典的一笔就要来了。

假设我们有生成器 $g_{\theta_t}(\mathbf{z})$ ，它的采样生成结果就等于从 $q_t(\mathbf{x})$ 采样的结果，即

$$\{g_{\theta_t}(\mathbf{z}) \mid \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})\} = \{\mathbf{x}_t \mid \mathbf{x}_t \sim q_t(\mathbf{x})\} \quad (5)$$

那么现在我们就可以利用它和式(3)来估算 $r_t(\mathbf{x})$ 。注意这只是 t 时刻的 $r_t(\mathbf{x})$ ，其他时刻的 $r_t(\mathbf{x})$ 我们并不知道，所以无法直接通过式(2)完成最终的采样过程，但是我们可以往前推一小步：

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \epsilon \nabla_{\mathbf{x}_t} \log r_t(\mathbf{x}_t) = \mathbf{x}_t + \epsilon \nabla_{\mathbf{x}_t} D(\mathbf{x}_t) \quad (6)$$

这里的 ϵ 是一个很小的正数，代表步长。那么，现在我们就有了下一步采样的结果，我们希望它继续能等价于下一步的生成器的采样结果，即

$$\begin{aligned} \{g_{\theta_{t+1}}(\mathbf{z}) \mid \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})\} &= \{\mathbf{x}_{t+1} \mid \mathbf{x}_{t+1} \sim q_{t+1}(\mathbf{x})\} \\ &= \{\mathbf{x}_{t+1} \mid \mathbf{x}_t + \epsilon \nabla_{\mathbf{x}_t} D(\mathbf{x}_t), \mathbf{x}_t \sim q_t(\mathbf{x})\} \end{aligned} \quad (7)$$

换句话说，我们想要**将扩散模型中样本的运动转化为生成器参数的运动**！为了达到这个目标，我们通过如下损失去求 θ_{t+1} ：

$$\theta_{t+1} = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\left\| g_{\theta}(\mathbf{z}) - g_{\theta_t}(\mathbf{z}) - \epsilon \nabla_{\mathbf{g}} D(g_{\theta_t}(\mathbf{z})) \right\|^2 \right] \quad (8)$$

也就是说，拿 $\mathbf{x}_t = g_{\theta_t}(\mathbf{z})$ 往前迭代一步得到 \mathbf{x}_{t+1} ，然后希望新的 $g_{\theta_{t+1}}(\mathbf{z})$ 能尽量等于 \mathbf{x}_{t+1} 。完成这一轮后，再用 θ_{t+1} 替代原本的 θ_t 开始新一轮的迭代，也就是式(3)和式(8)交替执行，是不是就有GAN的味道了？

点睛之笔

如果这还不够，我们还可以继续完善一下，将它变得跟GAN更加一致。注意到式(8)的被期望函数的梯度是：

$$\begin{aligned} & \nabla_{\theta} \|g_{\theta}(z) - g_{\theta_t}(z) - \epsilon \nabla_g D(g_{\theta_t}(z))\|^2 \\ &= 2 \langle g_{\theta}(z) - g_{\theta_t}(z) - \epsilon \nabla_g D(g_{\theta_t}(z)), \nabla_{\theta} g_{\theta}(z) \rangle \end{aligned} \quad (9)$$

代入当前值 $\theta = \theta_t$ ，那么结果是

$$-2\epsilon \langle \nabla_g D(g_{\theta_t}(z)), \nabla_{\theta_t} g_{\theta_t}(z) \rangle = -2\epsilon \nabla_{\theta_t} D(g_{\theta_t}(z)) \quad (10)$$

也就是说，如果用基于梯度的优化器只优化一步的话，那么以式(8)为损失函数，跟以下式为损失函数，结果是等价的（因为梯度只差一个常数倍）：

$$\theta_{t+1} = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{z \sim \mathcal{N}(\mathbf{0}, I)} [-D(g_{\theta}(z))] \quad (11)$$

这便是常见的生成器损失之一。式(3)和式(11)交替训练，就是一个常见的GAN变体。

特别地，原论文还证明了生成器的损失函数可以一般化为

$$\theta_{t+1} = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{z \sim \mathcal{N}(\mathbf{0}, I)} [-h(D(g_{\theta}(z)))] \quad (12)$$

其中 $h(\cdot)$ 是任意单调递增函数，它也对应于Wasserstein梯度流(1)中的 $\log r_t(\mathbf{x})$ 可以换成 $h(\log r_t(\mathbf{x}))$ ，这应该也是MonoFlow一词的来源（Monotonically increasing function + Wasserstein flow）。这个证明过程就不展开了，大家自行看原论文就好。

意义思考

总的来说，将GAN理解为扩散模型的思路是

$$\cdots \longrightarrow g_{\theta_t}(z) \xrightarrow{\text{式(3)}} r_t(\mathbf{x}) \xrightarrow{\text{式(6)}} \mathbf{x}_{t+1} \xrightarrow{\text{式(8)}} g_{\theta_{t+1}}(z) \longrightarrow \cdots$$

其中，核心的式子是(6)，它源于Wasserstein梯度流的式(1)和式(2)，这部分我们在上一篇文章《梯度流：探索通向最小值之路》讨论过了。

可能有读者想问：这个视角看上去并没有得到比GAN更多的东西，为什么要费这番大功夫去重新理解GAN呢？首先，在笔者看来，从扩散模型角度理解GAN，或者说将扩散模型和GAN统一起来，它本身就是一件很有趣、很好玩的事情，并不一定需要发挥什么实际作用，有趣、好玩就是它最大的意义。

其次，如同作者在原论文所说，已有的GAN的推导过程跟它实际的训练过程是不一致的，而本文所讨论的扩散视角，则是跟训练过程是一致的。也就是说，以训练过程为标准的话，GAN已有的推导过程是错的，本文的扩散视角才是对的。怎么理解这一点呢？以前面提到的GAN为例，判别器和生成器的目标分别是：

$$\max_D \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\log \sigma(D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\log(1 - \sigma(D(\mathbf{x})))] \quad (13)$$

$$\min_q \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [-D(\mathbf{x})] \quad (14)$$

通常的证明方式是，证明 D 的最优解是 $\log \frac{p(\mathbf{x})}{q(\mathbf{x})}$ ，然后代入生成器的损失函数，发现它在最小化 $q(\mathbf{x}), p(\mathbf{x})$ 的KL散度，所以最优解是 $q(\mathbf{x}) = p(\mathbf{x})$ 。但是，这样的证明对应的训练方式应该是先针对任意的 $q(\mathbf{x})$ ，将 \max_D 这一步都求解出来（求出的 D 应该是 $q(\mathbf{x})$ 的函数，或者说应该是生成器参数 θ 的函数），然后再去执行 \min_q 这一步，而不是实际上用的交替训练。而基于扩散模型的理解，它在设计上就是交替的，所以它跟训练过程更加一致。

总的来说，从扩散模型的角度来理解GAN，不单单是理解GAN的一种新视角，而且还是一种更贴近训练过程的视角。比如，我们可以解释为什么GAN的生成器不能训练太多步，因为只有单步优化时，式(11)和式(8)才等价，如果GAN要进行更多步的优化，那么应该使用式(8)为损失函数。事实上，式(8)就相当于笔者之前在《用变分推断统一理解生成模型（VAE、GAN、AAE、ALI）》所提出的 $KL(q(\mathbf{x}) \| q^o(\mathbf{x}))$ 项，它保证了生成器的“传承”而不仅仅是“创新”。

文章小结

本文介绍了MonoFlow，它展示了我们可以将GAN理解为在另一个时间维度上的扩散ODE，从而建立了一种基于扩散模型理解GAN的新视角。特别地，这是一种比GAN的常

规推导更加贴近训练过程的视角。

转载到请包括本文地址：<https://spaces.ac.cn/archives/9662>

更详细的转载事宜请参考：《科学空间FAQ》

如果您需要引用本文，请参考：

苏剑林. (Jun. 24, 2023). 《生成扩散模型漫谈（十九）：作为扩散ODE的GAN》 [Blog post]. Retrieved from <https://spaces.ac.cn/archives/9662>

```
@online{kexuefm-9662,  
  title={生成扩散模型漫谈（十九）：作为扩散ODE的GAN},  
  author={苏剑林},  
  year={2023},  
  month={Jun},  
  url={\url{https://spaces.ac.cn/archives/9662}},  
}
```