

13.2.5 再参数化

再参数化 (Reparameterization) 是将一个函数 $f(\theta)$ 的参数 θ 用另外一组参数表示 $\theta = g(\vartheta)$, 这样函数 $f(\theta)$ 就转换成参数为 ϑ 的函数 $\hat{f}(\vartheta) = f(g(\vartheta))$. 再参数化通常用来将原始参数转换为另外一组具有特殊属性的参数. 比如当 θ 为一个很大的矩阵时, 可以使用两个低秩矩阵的乘积来再参数化, 从而减少参数量.

再参数化的另一个例子是逐层归一化, 参见第7.5节.

在公式(13.21)中, 期望 $\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}|\mathbf{z}; \theta)]$ 依赖于分布 q 的参数 ϕ . 但是, 由于随机变量 \mathbf{z} 采样自后验分布 $q(\mathbf{z}|\mathbf{x}; \phi)$, 它们之间不是确定性关系, 因此无法直接求解 \mathbf{z} 关于参数 ϕ 的导数. 这时, 我们可以通过再参数化方法来将 \mathbf{z} 和 ϕ 之间随机性的采样关系转变为确定性函数关系.

我们引入一个分布为 $p(\epsilon)$ 的随机变量 ϵ , 期望 $\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}|\mathbf{z}; \theta)]$ 可以重写为

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}|\mathbf{z}; \theta)] = \mathbb{E}_{\epsilon \sim p(\epsilon)} [\log p(\mathbf{x}|g(\phi, \epsilon); \theta)], \quad (13.25)$$

其中 $\mathbf{z} \triangleq g(\phi, \epsilon)$ 为一个确定性函数.

假设 $q(\mathbf{z}|\mathbf{x}; \phi)$ 为正态分布 $N(\mu_I, \sigma_I^2 I)$, 其中 $\{\mu_I, \sigma_I\}$ 是推断网络 $f_I(\mathbf{x}; \phi)$ 的输出, 依赖于参数 ϕ , 我们可以通过下面方式来再参数化:

$$\mathbf{z} = \mu_I + \sigma_I \odot \epsilon, \quad (13.26) \quad \text{参见习题13-1.}$$

其中 $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$. 这样 \mathbf{z} 和参数 ϕ 的关系从采样关系变为确定性关系, 使得 $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \phi)$ 的随机性独立于参数 ϕ , 从而可以求 \mathbf{z} 关于 ϕ 的导数.

13.2.6 训练

通过再参数化, 变分自编码器可以通过梯度下降法来学习参数, 从而提高变分自编码器的训练效率.

给定一个数据集 $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$, 对于每个样本 $\mathbf{x}^{(n)}$, 随机采样 M 个变量 $\epsilon^{(n,m)}, 1 \leq m \leq M$, 并通过公式(13.26)计算 $\mathbf{z}^{(n,m)}$. 变分自编码器的目标函数近似为

$$\mathcal{J}(\phi, \theta|\mathcal{D}) = \sum_{n=1}^N \left(\frac{1}{M} \sum_{m=1}^M \log p(\mathbf{x}^{(n)}|\mathbf{z}^{(n,m)}; \theta) - \text{KL}(q(\mathbf{z}|\mathbf{x}^{(n)}; \phi), \mathcal{N}(\mathbf{z}; \mathbf{0}, I)) \right). \quad (13.27)$$

如果采用随机梯度方法, 每次从数据集中采集一个样本 \mathbf{x} 和一个对应的随机变量 ϵ , 并进一步假设 $p(\mathbf{x}|\mathbf{z}; \theta)$ 服从高斯分布 $\mathcal{N}(\mathbf{x}|\mu_G, \lambda I)$, 其中 $\mu_G = f_G(\mathbf{z}; \theta)$ 是生成网络的输出, λ 为控制方差的超参数, 则目标函数可以简化为

$$\mathcal{J}(\phi, \theta|\mathbf{x}) = -\frac{1}{2} \|\mathbf{x} - \mu_G\|^2 - \lambda \text{KL}(\mathcal{N}(\mu_I, \sigma_I), \mathcal{N}(\mathbf{0}, I)), \quad (13.28) \quad \text{参见习题13-2.}$$

其中第一项可以近似看作输入 \mathbf{x} 的重构正确性，第二项可以看作正则化项， λ 可以看作正则化系数。这和自编码器在形式上非常类似，但它们的内在机理是完全不同的。 参习题13-3.

变分自编码器的训练过程如图13.5所示，其中空心矩形表示“目标函数”。

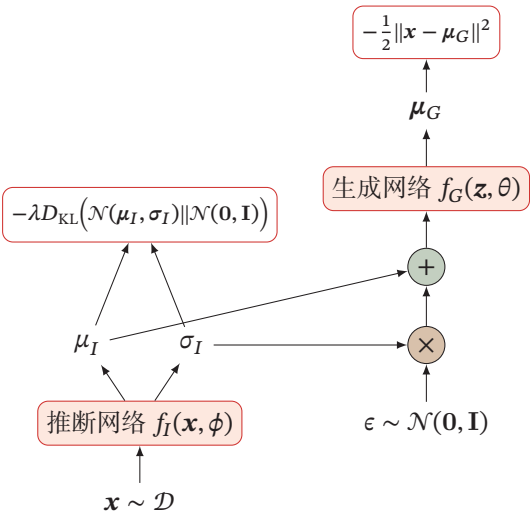


图 13.5 变分自编码器的训练过程

图13.6给出了在 MNIST 数据集上变分自编码器学习到的隐变量流形的可视化示例。图13.6a是将训练集上每个样本 \mathbf{x} 通过推断网络映射到 2 维的隐变量空间，图中的每个点表示 $\mathbb{E}[\mathbf{z}|\mathbf{x}]$ ，不同颜色表示不同的数字。图13.6b是对 2 维的标准高斯分布上进行均匀采样得到不同的隐变量 \mathbf{z} ，然后通过生成网络产生 $\mathbb{E}[\mathbf{x}|\mathbf{z}]$ 。

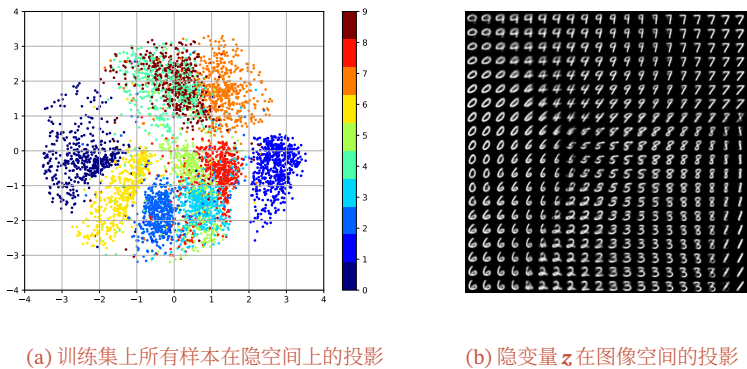


图 13.6 在 MNIST 数据集上变分自编码器学习到的隐变量流形的可视化示例

13.3 生成对抗网络

13.3.1 显式密度模型和隐式密度模型

之前介绍的深度生成模型，比如变分自编码器、深度信念网络等，都是显式地构建出样本的密度函数 $p(\mathbf{x}; \theta)$ ，并通过最大似然估计来求解参数，称为**显式密度模型**（Explicit Density Model）。比如，变分自编码器的密度函数为 $p(\mathbf{x}, \mathbf{z}; \theta) = p(\mathbf{x}|\mathbf{z}; \theta)p(\mathbf{z}; \theta)$ 。虽然使用了神经网络来估计 $p(\mathbf{x}|\mathbf{z}; \theta)$ ，但是我们依然假设 $p(\mathbf{x}|\mathbf{z}; \theta)$ 为一个参数分布族，而神经网络只是用来预测这个参数分布族的参数。这在某种程度上限制了神经网络的能力。

如果只是希望有一个模型能生成符合数据分布 $p_r(\mathbf{x})$ 的样本，那么可以不显式地估计出数据分布的密度函数。假设在低维空间 \mathcal{Z} 中有一个简单易采样的分布 $p(\mathbf{z})$ ， $p(\mathbf{z})$ 通常为标准多元正态分布 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ 。我们用神经网络构建一个映射函数 $G : \mathcal{Z} \rightarrow \mathcal{X}$ ，称为**生成网络**。利用神经网络强大的拟合能力，使得 $G(\mathbf{z})$ 服从数据分布 $p_r(\mathbf{x})$ 。这种模型就称为**隐式密度模型**（Implicit Density Model）。所谓隐式模型就是指并不显式地建模 $p_r(\mathbf{x})$ ，而是建模生成过程。图13.7给出了隐式模型生成样本的过程。

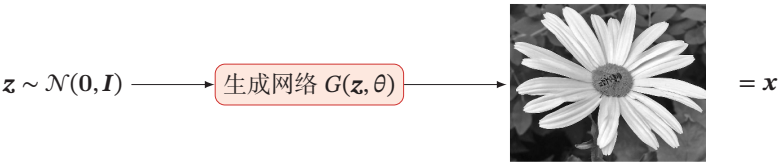


图 13.7 隐式模型生成样本的过程

13.3.2 网络分解

隐式密度模型的一个关键是如何确保生成网络产生的样本一定是服从真实的数据分布。既然我们不构建显式密度函数，就无法通过最大似然估计等方法来训练。**生成对抗网络**（Generative Adversarial Networks, GAN）[Goodfellow et al., 2014] 是通过对抗训练的方式来使得生成网络产生的样本服从真实数据分布。在生成对抗网络中，有两个网络进行对抗训练。一个是**判别网络**，目标是尽量准确地判断一个样本是来自于真实数据还是由生成网络产生；另一个是**生成网络**，目标是尽量生成判别网络无法区分来源的样本。这两个目标相反的网络不断地进行交替训练。当最后收敛时，如果判别网络再也无法判断出一个样本的来源，那么也就等价于生成网络可以生成符合真实数据分布的样本。生成对抗网络的流程图如图13.8所示。