

An Introduction to Variational Autoencoders

Notes to get started

SummerSemester 2024

01.08.2024

University of Chinese Academy of Sciences

ZTS

plote5024@mail.com

Beijing

Inhaltsverzeichnis

1. Introduction	1
1.1. Motivation	1
1.2. Aim	2
1.3. Probabilistic Models and Variational Inference	2
1.3.1. Conditional Models	2
1.4. Directed Graphical Models and Neural Networks	2
1.5. Dataset	2
1.6. Maximum Likelihood and Minibatch SGD	3
1.6.1. 小批量数据的梯度期望值等于整个数据集的梯度	3
1.6.1.1. Coefficient proof	4
1.7. Learning and Inference in Deep Latent Variable Models	5
1.7.1. Latent Variables	5
1.7.2. Deep Latent Variable Models	6
1.8. Intractabilities	6
2. Variational Autoencoders	6
2.1. Encoder or Approximate Posterior	6

1. Introduction

1.1. Motivation

机器学习的一个主要分支是生成模型和判别模型。在判别建模中，目标是学习给定观察值的预测器，而在生成建模中，目标是解决学习所有变量 **All variables** 的联合分布。

生成模型模拟数据在现实世界中的生成方式。在几乎每一门科学中，“建模”都被理解为通过假设理论和通过观察测试这些理论来揭示这个生成过程。例如，当气象学家为天气建模时，他们使用高度复杂的偏微分方程来表达天气的潜在物理特性。生物学家、化学家、经济学家等等也是如此。科学中的建模实际上几乎全是生成模型。

试图理解数据生成过程的另一个原因是，它自然地表达了世界的因果关系。因果关系的优点，就是它们比单纯的相关性更能概括新情况。

为了将生成模型转化为 Discriminator，我们需要使用贝叶斯规则。即我们可以将生成模型的输出转换为分类任务中需要的条件概率，从而实现判别功能。

在判别方法中，我们直接学习一个与未来预测方向一致的映射。这与生成模型的方向相反。例如，可以认为一张图像在现实世界中是通过先识别物体，然后生成三维物体，再将其投影到像素网格上来生成的。

而判别模型则直接将像素值作为输入，并将其映射到标签上。虽然生成模型能够有效地从数据中学习，但它们往往比纯粹的判别模型对数据做出更强的假设，这通常会在模型出错时导致更高的渐近偏差 (**Asymptotic bias**)。

因此，如果模型出错（事实上几乎总是有一定程度的误差），如果我们只关注于学习如何区分，并且我们有足够多的数据，那么纯粹的判别模型在判别任务中通常会导致更少的错误。然而，取决于数据量的多少，研究数据生成过程可能有助于指导判别器（如分类器）的训练。例如，在半监督学习的情况下，我们可能只有少量的标记样本，但有更多的未标记样本。在这种情况下，可以利用数据的生成模型来改进分类。

这种方法可以帮助我们构建有用的世界抽象，这些抽象可以用于多个后续的预测任务。这种追求数据中解缠、语义上有意义、统计上独立和因果关系的变化因素的过程，通常被称为无监督表示学习，而变分自编码器 (VAE) 已经广泛应用于此目的。

或者，可以将其视为一种隐式正则化形式：通过强制表示对数据生成过程有意义，我们将从输入到表示的映射过程约束在某种特定的模式中。通过预测世界这一辅助任务，我们可以在抽象层面上更好地理解世界，从而更好地进行后续的预测。

变分自编码器 (VAE) 可以看作是两个耦合但独立参数化的模型：编码器或识别模型 (**Encoder or Recognition Model**)，以及解码器或生成模型 (**Decoder or Generative Model**)。这两个模型相互支持。识别模型向生成模型提供其后验分布的近似值，后者需要这些近似值在“期望最大化”学习的迭代过程中更新其参数。反过来，生成模型为识别模型提供了一个框架，使其能够学习数据的有意义表示，包括可能的类别标签。根据贝叶斯规则，识别模型是生成模型的近似逆。

与普通的变分推断 (VI) 相比，VAE 框架的一个优势在于识别模型（也称为推断模型）现在是输入变量的一个（随机）函数。这与 VI 不同，后者对每个数据实例都有一个独立的变分分布，这对于大数据集来说效率低下。识别模型使用一组参数来建模输入与潜变量之间的关系，因此被称为“摊销推断”。这种识别模型可以是任意复杂的，但由于其构造方式，只需一次从输入到潜变量的前馈传递即可完成，因此仍然相对快速。然而，我们付出的代价是，这种采样会在学习所需的梯度中引入采样噪声。

或许 VAE 框架最大的贡献是认识到我们可以使用现在被称为“重参数化技巧”的简单方法来重新组织我们的梯度计算，从而减少梯度中的方差。

1.2. Aim

provides a principled method for jointly learning *deep latent-variable models* and *corresponding inference models* using *stochastic gradient descent*.

该框架在生成建模、半监督学习和表示学习等方面有广泛的应用。

1.3. Probabilistic Models and Variational Inference

由于概率模型包含未知数且数据很少能完整地描述这些未知参数，我们通常需要对模型的某些方面假设一定程度的不确定性。这种不确定性的程度和性质通过（条件）概率分布来描述。在某种意义上，最完整的概率模型形式通过这些变量的联合概率分布来指定模型中所有变量之间的相关性和高阶依赖关系。

我们用 \mathbf{x} 表示所有观测变量的向量，其联合分布是我们希望建模的。（使用小写粗体来表示底层的观测随机变量集，即将其展平和拼接，使该集合表示为一个单一的向量。）

假设观测变量 \mathbf{x} 是来自未知底层过程的随机样本，其真实的概率分布 $p^*(\mathbf{x})$ 是未知的。我们尝试用一个选定的模型 $p_\theta \mathbf{x}$ 来近似这个底层过程，其中参数为 θ ：

$$\mathbf{x} \sim p_\theta(\mathbf{x})$$

学习最常见的是一个搜索参数 θ 的过程，使得由模型给出的概率分布函数 $p_\theta(\mathbf{x})$ 近似于数据的真实分布 $p^*(\mathbf{x})$ ，即对于任何观测到的 \mathbf{x} ，两者尽可能接近：

$$p_\theta(\mathbf{x}) \approx p^*(\mathbf{x})$$

1.3.1. Conditional Models

在分类或者回归问题上，我们不关心无条件模型 $p_\theta(\mathbf{x})$ ，更倾向于条件模型 $p_\theta(y|\mathbf{x})$ ，它近似于底层的条件分布 $p^*(y|\mathbf{x})$ ：即在观测变量 \mathbf{x} 的值上，对变量 y 的值进行分布。在这种情况下， $\text{bold}(\mathbf{x})$ 通常被称为模型的输入。与无条件情况类似，我们选择并优化一个模型 $p_\theta(y|\mathbf{x})$ ，使其接近未知的底层分布，即对于任何 \mathbf{x} 和 y ：

$$p_\theta(y|\mathbf{x}) \approx p^*(y|\mathbf{x})$$

1.4. Directed Graphical Models and Neural Networks

我们使用有向概率图模型(Directed probability graph model)，或贝叶斯网络。有向图模型是一种概率模型，其中所有的变量被拓扑组织成一个有向无环图。这些模型的变量的联合分布被分解为先验分布和条件分布的乘积：

$$p_\theta(\mathbf{x}_1, \dots, \mathbf{x}_M) = \prod_{j=1}^M p_\theta(\mathbf{x}_j | P_a(\mathbf{x}_j))$$

$P_a(\mathbf{x}_j)$ 是有向图中节点 j 的父向量的集合。

1.5. Dataset

我们收集 $N \geq 1$ 个数据点组成的数据集 \mathcal{D} ：

$$\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\} \equiv \{\mathbf{x}^{(i)}\}_{i=1}^N \equiv \mathbf{x}^{(1:N)}$$

数据集被认为由同一（不变）系统的独立测量值组成。在这种情况下，观测数据 $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ 被称为独立同分布(i.i.d.)。在独立同分布的假设下，给定参数的情况下，数据点的概率因子化为各个数据点概率的乘积。因此，模型分配给数据的对数概率表示为：

$$\log p_{\theta}(\mathcal{D}) = \sum_{x \in \mathcal{D}} \log p_{\theta}(x)$$

1.6. Maximum Likelihood and Minibatch SGD

概率模型最常见的标准是最大对数似然(ML). **Maximization of the log-likelihood criterion is equivalent to minimization of a Kullback Leibler divergence between the data and model distributions.**

优化方式比较好用的方法是 *stochastic gradient descent*

考虑一个数据集 \mathcal{D} 包含 N 个数据点, 我们无法每次都使用整个数据集来计算梯度, 因为这样计算代价太高。因此, 我们随机抽取一个小批量数据 \mathcal{M} (例如, 大小为 $N_{\mathcal{M}}$), 并计算这个小批量数据上的梯度。我们希望这个小批量数据上的梯度能够在期望值上等于整个数据集上的梯度。

在随机梯度下降中, 这意味着 (公式推导放在 1.6.1, 不占用篇幅继续写结论):

$$\mathbb{E}[\nabla_{\theta} \log p_{\theta}(\mathcal{M})] = \nabla_{\theta} \log p_{\theta}(\mathcal{D})$$

\mathcal{M} 是从 \mathcal{D} 中随机抽出来的小批量数据, $\log p_{\theta}(\mathcal{M})$ 是这个小额批量的对数似然, $\log p_{\theta}(\mathcal{D})$ 是整个数据集的对数似然。

对数似然无偏估计:

对于整个数据集 \mathcal{D} , 其对数似然是:

$$\log p_{\theta}(\mathcal{D}) = \sum_{x \in \mathcal{D}} \log p_{\theta}(x)$$

对于小批量数据 \mathcal{M} , 其对数似然是:

$$\log p_{\theta}(\mathcal{M}) = \sum_{x \in \mathcal{M}} \log p_{\theta}(x)$$

在期望上, 小批量数据的对数似然可以近似整个数据集的对数似然:

$$\frac{1}{N_{\mathcal{D}}} \log p_{\theta}(\mathcal{D}) \approx \frac{1}{N_{\mathcal{M}}} \log p_{\theta}(\mathcal{M})$$

通过这样的迷你批次 *minibatches*, 我们可以形成最大似然准则的无偏估计:

$$\begin{aligned} \frac{1}{N_{\mathcal{D}}} \log p_{\theta}(\mathcal{D}) &\approx \frac{1}{N_{\mathcal{M}}} \log p_{\theta}(\mathcal{M}) = \frac{1}{N_{\mathcal{M}}} \sum_{x \in \mathcal{M}} \log p_{\theta}(x) \\ \frac{1}{N_{\mathcal{D}}} \nabla_{\theta} \log p_{\theta}(\mathcal{D}) &\approx \frac{1}{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(\mathcal{M}) = \frac{1}{N_{\mathcal{M}}} \sum_{x \in \mathcal{M}} \nabla_{\theta} \log p_{\theta}(x) \end{aligned}$$

1.6.1. 小批量数据的梯度期望值等于整个数据集的梯度

假设数据集 \mathcal{D} 有 N 个数据点, 小批量 \mathcal{M} 之中有 $N_{\mathcal{M}}$ 个数据点, 并且小批量数据是从数据集中随机抽取的。

整个数据集的对数似然梯度:

$$\nabla_{\theta} \log p_{\theta}(\mathcal{D}) = \nabla_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i) = \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(x_i)$$

小批量数据的对数似然梯度:

$$\nabla_{\theta} \log p_{\theta}(\mathcal{M}) = \frac{N}{N_{\mathcal{M}}} \sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)$$

这里乘以 $\frac{N}{N_{\mathcal{M}}}$ 证明写在 1.6.1.1 了, 这里就不占篇幅继续向下推导了。

假设我们抽取了 $N_{\mathcal{M}}$ 个样本, 样本是独立同分布的, 我们可以用数学期望来表示这种抽样过程:

$$\mathbb{E}[\nabla_{\theta} \log p_{\theta}(\mathcal{M})] = \mathbb{E}\left[\frac{N}{N_{\mathcal{M}}} \sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)\right]$$

由于 \mathcal{M} 中的样本是独立同分布的, 所以:

$$\mathbb{E}\left[\sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)\right] = N_{\mathcal{M}} \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)]$$

其中 x 是从 \mathcal{D} 中随机抽取的一个样本。因此:

$$\begin{aligned} \mathbb{E}[\nabla_{\theta} \log p_{\theta}(\mathcal{M})] &= \mathbb{E}\left[\frac{N}{N_{\mathcal{M}}} \sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)\right] \\ &= \frac{N}{N_{\mathcal{M}}} N_{\mathcal{M}} \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] \\ &= N \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] \end{aligned}$$

由于数据点 x 是从整个数据集中抽取的, 所以 $\mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)]$ 实际上就是总体梯度的平均值:

$$N \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] = \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(x_i) = \nabla_{\theta} \log p_{\theta}(\mathcal{D})$$

1.6.1.1. Coefficient proof

是为了补偿小批量数据与整个数据集的大小差异, 使得梯度估计在期望值上保持一致。这里乘系数的原因是:

1. 比例调整: 由于小批量数据 \mathcal{M} 只是整个数据集 \mathcal{D} 的一个子集, 其样本数量 $N_{\mathcal{M}}$ 小于 N 。乘以比例系数可以调整小批量数据的梯度估计, 使其在期望上与整个数据集的梯度一致。
2. 无偏估计: 这种调整确保了小批量数据梯度估计的期望值等于整体数据集的梯度, 从而在期望上保持无偏。

定义整个数据集的对数似然梯度:

$$\nabla_{\theta} \log p_{\theta}(\mathcal{D}) = \nabla_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i) = \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(x_i)$$

定义小批量数据的对数似然梯度：

$$\nabla_{\theta} \log p_{\theta}(\mathcal{M}) = \sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)$$

现在假设我们从数据集中随机抽取小批量数据 \mathcal{M} ，由于 \mathcal{M} 中的数据点是独立同分布的，我们有：

$$\mathbb{E}[\nabla_{\theta} \log p_{\theta}(\mathcal{M})] = \mathbb{E}\left[\sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)\right] = N_{\mathcal{M}} \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)]$$

为了使小批量数据的梯度在期望上等于整个数据集的梯度，我们需要将小批量数据的梯度进行比例调整。我们定义调整后的梯度估计为：

$$\tilde{\nabla}_{\theta} \log p_{\theta}(\mathcal{M}) = \frac{N}{N_{\mathcal{M}}} \sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)$$

我们现在计算这个调整后的梯度的期望值：

$$\begin{aligned} \mathbb{E}[\tilde{\nabla}_{\theta} \log p_{\theta}(\mathcal{M})] &= \mathbb{E}\left[\frac{N}{N_{\mathcal{M}}} \sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)\right] \\ &= \frac{N}{N_{\mathcal{M}}} \mathbb{E}\left[\sum_{j=1}^{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(x_j)\right] \\ &= \frac{N}{N_{\mathcal{M}}} N_{\mathcal{M}} \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] \\ &= N \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] \end{aligned}$$

由于 x 是从整个数据集中随机抽取的样本，我们有：

$$N \mathbb{E}[\nabla_{\theta} \log p_{\theta}(x)] = \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(x_i) = \nabla_{\theta} \log p_{\theta}(\mathcal{D})$$

1.7. Learning and Inference in Deep Latent Variable Models

1.7.1. Latent Variables

我们可以将前一节讨论的完全观测有向模型扩展到包含潜变量的有向模型。潜变量是模型的一部分，但我们并不观测到它们，因此它们不属于数据集的一部分。我们通常用 \mathbf{z} 来表示这样的潜变量。在无条件建模观测变量 \mathbf{x} 的情况下，有向图模型将表示一个观测变量 \mathbf{x} 和潜变量 \mathbf{z} 的联合分布 $p_{\theta}(\mathbf{x}|\mathbf{z})$ 。观测变量 \mathbf{x} 的边缘分布由以下公式给出：

$$p_{\theta} = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

这也被称为（单个数据点的）边缘似然或模型证据，当它作为 θ 的函数时。

这种关于 \mathbf{x} 的隐式分布可能是非常灵活的。如果 \mathbf{z} 是离散的，并且 $p_{\theta}(\mathbf{x}|\mathbf{z})$ 是一个高斯分布，那么 $p_{\theta}(\mathbf{x})$ 就是一个高斯混合分布。对于连续的 \mathbf{z} , $p_{\theta}(\mathbf{x})$ 可以看作是一个无限混合，这种混合潜在地比离散混合更强大。这样的边缘分布也称为复合概率分布。

1.7.2. Deep Latent Variable Models

deep latent variable model (DLVM)

表示那些分布由神经网络参数化的潜变量模型 $p_{\theta}(\mathbf{x}, \mathbf{z})$ 。这样的模型可以基于某些上下文条件，如 $p_{\theta}(\mathbf{x}, \mathbf{z} | \mathbf{y})$ 。DLVM 的一个重要优势是，即使有向模型中的每个因子（先验或条件分布）相对简单（例如条件高斯分布），其边缘分布 $p_{\theta}(\mathbf{x})$ 也可以非常复杂，即包含几乎任意的依赖关系。这种表达能力使得深度潜变量模型在近似复杂的底层分布 $p^*(\mathbf{x})$ 时非常有吸引力。

或许最简单也是最常见的 DLVM 是通过以下结构指定的分解模型：

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$$

其中, $p_{\theta}(\mathbf{z})$ 和 $p_{\theta}(\mathbf{x}|\mathbf{z})$ 是已知的。分布 $p(\mathbf{z})$ 通常被称为潜变量 \mathbf{z} 的先验分布，因为它不以任何观测数据为条件。

1.8. Intractabilities

dlvm 中最大似然学习的主要困难是 **the marginal probability of data under the model** 通常难以处理。

由于计算边际似然(或模型证据)的方程中的积分, $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$, **not having an analytic solution or efficient estimator**。由于这种不可解性，我们不能对其参数进行微分并优化它，就像我们对完全观测模型所做的那样。计算边缘似然 **Need to integrate high-dimensional space**，这在实践中是非常困难的。所以在 MLE 中，我们需要最大化数据在模型下的概率。对于 DLVMs，这涉及到计算边缘似然 $p_{\theta}(\mathbf{x})$ ，但由于其不可解性，直接进行 MLE 是不现实的。由于直接计算边缘似然是不可行的，我们通常使用变分推断 (**Variational Inference**) 或马尔可夫链蒙特卡罗 (**MCMC**) 方法来近似后验分布。这些方法通过引入可解的近似来处理不可解的积分问题。同样，神经网络参数化的(有向模型) $p(\theta|\mathcal{D})$ 的后验通常难以精确计算，需要近似推理技术。

2. Variational Autoencoders

2.1. Encoder or Approximate Posterior

dlvm 估计这种模型中的对数似然分布和后验分布的问题。变分自编码器(VAEs)框架提供了一种计算效率高的方法来优化 dlvm，并结合相应的推理模型使用 SGD 进行优化。为了将 DLVM 的后验推理和学习问题转化为可处理的问题，我们引入了一个参数推理模型 $q_{\Phi}(\mathbf{z}|\mathbf{x})$ 。这个模型也被称为编码器或识别模型。用 Φ 表示该推理模型的参数，也称为变分参数。

Selbstständigkeitserklärung

Hiermit versichere ich, dass ich die vorliegende schriftliche Hausarbeit (Seminararbeit, Belegarbeit) selbstständig verfasst und keine anderen als die von mir angegebenen Quellen und Hilfsmittel benutzt habe. Die Stellen der Arbeit, die anderen Werken wörtlich oder sinngemäß entnommen sind, wurden in jedem Fall unter Angabe der Quellen (einschließlich des World Wide Web und anderer elektronischer Text- und Datensammlungen) kenntlich gemacht. Dies gilt auch für beigegebene Zeichnungen, bildliche Darstellungen, Skizzen und dergleichen. Ich versichere weiter, dass die Arbeit in gleicher oder ähnlicher Fassung noch nicht Bestandteil einer Prüfungsleistung oder einer schriftlichen Hausarbeit (Seminararbeit, Belegarbeit) war. Mir ist bewusst, dass jedes Zuwiderhandeln als Täuschungsversuch zu gelten hat, aufgrund dessen das Seminar oder die Übung als nicht bestanden bewertet und die Anerkennung der Hausarbeit als Leistungsnachweis/Modulprüfung (Scheinvergabe) ausgeschlossen wird. Ich bin mir weiter darüber im Klaren, dass das zuständige Lehrerprüfungsamt/Studienbüro über den Betrugsversuch informiert werden kann und Plagiate rechtlich als Straftatbestand gewertet werden.

Ort: **Unterschrift:**

Datum: