

1 算例 H: 《放弃重构, 极致的压缩与预测——变分信息瓶颈 (VIB)》

在标准变分自编码器 (VAE) 中, 我们通过摊销推断网络 ϕ 提取隐变量 Z , 并用它来重构输入 X 。然而, 在许多判别性任务 (分类、回归) 中, 人们不关心能否复原 X 的所有细节。例如, 当输入 X 是一张农作物高光谱图像, 人们关心从中提取的特征 Z 能不能准确预测该植株的抗病标签 Y ?

Naftali Tishby 于 1999 年提出的信息瓶颈 (Information Bottleneck, IB) 理论为这一目标提供了完美的数学框架。而变分信息瓶颈 (Variational Information Bottleneck, VIB) 则利用与 VAE 想法类似的变分下界 (Variational Lower Bound, VLB) 推导技巧, 将这个抽象的信息论概念转化为深度学习中可计算、可端到端优化的损失函数。读完此篇后, 请问, 这里的 VLB, 和 ELBO 是相同的概念吗?

1.1 变量定义与逻辑对齐 (对比 VAE)

为了保持严谨的逻辑连贯性, 我们将 VIB 的变量与 VAE 的体系进行一一对齐。VIB 的核心区别在于显式地引入了监督信号 Y 。

符号	描述	在 VAE 中的逻辑对应
X	观测变量 (输入特征, 如高光谱图像、基因表达谱)	对应 VAE 的输入 X
Y	目标变量 (监督标签, 如抗病类别、肿瘤状态)	VIB 新增的观测变量
Z	低维隐变量 (Latent Representation)	对应 VAE 的隐变量 Z
ϕ	推断网络参数 (Encoder, 提取特征)	对应 VAE 的 Encoder ϕ
θ	预测网络参数 (Decoder/Classifier)	对应 VAE 的 Decoder θ

1.2 概率模型设定与马尔可夫假设

1.2.1 全概率分布与马尔可夫链

在 VIB 框架中, 数据生成与特征提取过程必须遵循一个严格的马尔可夫链 (Markov Chain) 假设:

$$Y \leftrightarrow X \leftrightarrow Z \quad (1-1)$$

双向箭头符号被称为 Markov Chain Shorthand (马尔可夫链简写)。含义为:

$$Y \perp\!\!\!\perp Z \mid X$$

即: 给定 X 时, Y 与 Z 条件独立:

$$P(Y, Z|X) = P(Y|X)P(Z|X)$$

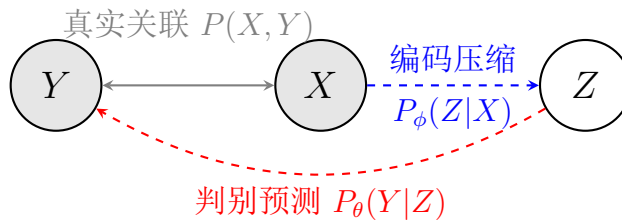


图 1: VIB 的信息流向图。隐变量 Z 只能通过观测数据 X 提取 (蓝色虚线) 为预测目标 Y (红色虚线), Z 必须且只能从 X 中“抽取”与 Y 相关的核心信息, 同时必须像“瓶颈”一样过滤 X 中的冗余背景噪声。

这意味着: 隐变量 Z 只能通过输入 X 来获取关于目标 Y 的信息。一旦给定 X , Z 与 Y 在数学上就是条件独立的 (即 $P(Z|X, Y) = P(Z|X)$)。

因此, 在给定客观真实数据分布 $P(X, Y)$ 的前提下, 包含隐变量的联合分布可以精确分解为:

$$P_{true}(X, Y, Z) = P(X, Y)P(Z|X) \quad (1-2)$$

其中, $P(X, Y)$ 是由客观世界决定的、我们无法改变的真实数据分布; 而 $P(Z|X)$ 则是我们将要使用神经网络 (参数为 ϕ) 去参数化的特征提取器 (即编码器 $q_\phi(Z|X)$)。

1.2.2 变分分布族

由于真实的条件概率往往难以计算, 我们引入参数化的变分分布:

1. **推断网络 (Encoder, 对应 $P(Z|X)$):** 我们用神经网络 ϕ 来参数化推断过程, 假设隐空间服从高斯分布:

$$q_\phi(Z|X) = \mathcal{N}(Z|\mu_\phi(X), \sigma_\phi^2(X))$$

2. **预测网络 (Decoder/Classifier, 近似 $P(Y|Z)$):** 我们引入一个参数为 θ 的网络, 放弃对 X 的重构, 转而基于隐变量 Z 专注预测 Y :

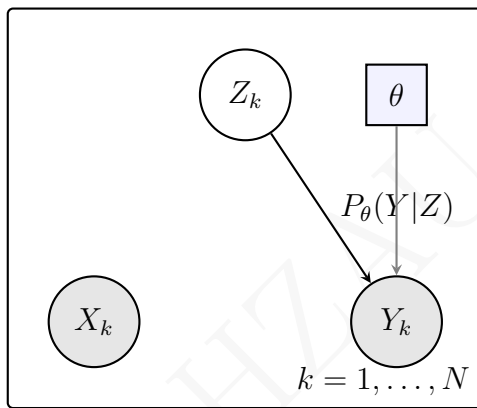
$$P_\theta(Y|Z) = \text{Categorical}(f_\theta(Z)) \quad (\text{若为分类任务})$$

3. **隐变量先验 (Prior):** 假设隐空间服从标准正态分布, 这是限制信息流量、实现特征“瓶颈”压缩的关键:

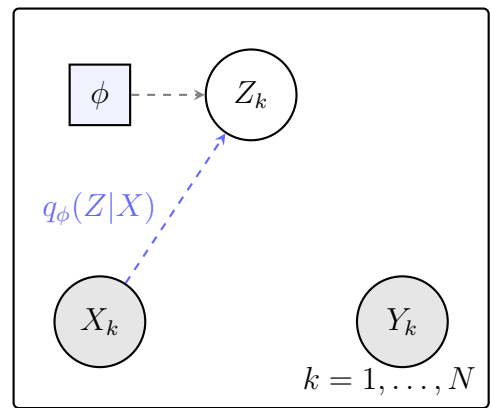
$$P(Z) = \mathcal{N}(Z|\mathbf{0}, \mathbf{I})$$

1.2.3 概率图模型 (PGM): 分离生成与推断, 以及瓶颈的作用

为了更清晰地理解 VIB 的信息流动机制, 我们将联合 PGM 拆解为“预测 (对应 VAE 生成阶段)”和“推断”两个独立的视图。



(a) **预测阶段:** 隐变量 Z 仅生成 Y 。图中不存在 $Z \rightarrow X$ 的连线——即所谓“放弃重构”。



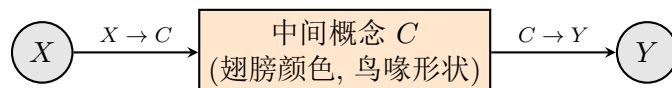
(b) **推断阶段:** 图中不存在 $Y \rightarrow Z$ 的连线, 编码器在提特征时对标签一无所知。即所谓“条件独立” $Y \perp\!\!\!\perp Z | X$ 。

图 2: 变分信息瓶颈 (VIB) 的解耦 PGM。将网络切分为独立的预测网络 θ 和推断网络 ϕ 。两个阶段共同作用, 使得信息流只能从 X 榨取, 穿过 Z 这个压缩“瓶颈”, 最终抵达 Y 。

如何理解模型设计中的“瓶颈”意义呢? 我们可以谈谈 CBM。¹ 无论是 VIB 中的 Z , 还是 CBM 中的 C , 它们都在 PGM 中扮演着马尔可夫链 $Y \leftarrow \text{Bottleneck} \leftarrow X$ 的核心阻断者。VIB (隐式瓶颈) 通过数学优

¹在可解释 AI (XAI) 领域, 斯坦福大学的研究者们 (Koh et al., 2020) 提出了另一种极具启发性的架构: 概念瓶颈模型 (Concept Bottleneck Models, CBM)。

假设我们要从一张鸟类图像 X 预测鸟的种类 Y 。普通的黑盒网络会直接 $X \rightarrow Y$ 。但在 CBM 中, 人们在中间强制插入一个“概念层 C ”:



网络必须先根据图像 X 预测出一组人类预定义的、可理解的中概念 C (如: 翅膀是红色的吗? 鸟喙是尖的吗?)。然后, 分类

化，让网络自己学出最浓缩的不可解释特征。CBM (显式瓶颈) 通过人类先验，让网络对齐到最干净的可解释概念。

理解了“瓶颈 = 过滤噪音的安检机”这一核心物理直觉，我们接下来将看到，VIB 是如何用极其优美的变分数学手段，在深度网络中全自动实现这种极致的信息挤压的。

1.3 理论起点：信息瓶颈 (IB) 目标函数

在拆解信息瓶颈理论之前，我们需要首先回顾互信息，其本质上是在衡量 X 和 Y 的联合分布 $P(X, Y)$ ，与假设它们完全独立时的分布 $P(X)P(Y)$ 之间的“距离”。两者越不独立，包含的互信息就越大。²

信息瓶颈理论的直觉是：一个好的特征 Z 应该像一个极其挑剔的“筛子”或“瓶颈”。它必须保留足够多关于目标 Y 的信息，同时尽最大可能漏掉（遗忘）关于输入 X 的冗余细节。

用香农互信息 $I(\cdot; \cdot)$ 表达，我们寻求最大化以下目标函数：

$$\max \mathcal{J}_{IB} = I(Z; Y) - \beta I(Z; X) \quad (1-3)$$

- $I(Z; Y)$: 充分性 (Sufficiency)。我们需要最大化这一项。这意味着 Z 需要尽可能多地包含 Y 的信息，使得 $H(Y|Z)$ 降到最低（即给定特征 Z 后，预测标签 Y 不再有悬念），要求 Z 能极其准确地预测 Y 。
- $I(Z; X)$: 最小性/压缩 (Minimality/Compression)。由于互信息前带负号，我们需要最小化这一项。虽然 Z 是从 X 中提取的，但我们强迫 Z 尽可能“忘记” X 。它要求 Z 剥离掉 X 中一切与 Y 无关的背景噪音。
- $\beta > 0$: 控制压缩强度的超参数（拉格朗日乘子）。调节 β 本质上是在“预测准确度”与“特征的极致压缩”之间进行零和博弈寻找平衡点。

1.4 VIB 的变分下界 VLB

直接计算高维空间中的互信息是不可解的 (Intractable)。我们必须利用变分推断（引入近似分布）分别对这两项求取变分下界 (Variational Lower Bound)。

器只能仅凭这些概念得分去预测最终的种类 Y 。这里的 C 就是一个语义瓶颈。如果背景里有一片蓝天，由于人类没有定义“背景是蓝天”这个中间概念，这部分冗余信息在概念层 C 被拦截，不会渗透到最后的 Y 中。

²信息熵 (Entropy) $H(X)$ 度量了一个随机变量 X 内部的不确定度（或所含的固有信息量）。对于连续变量，其信息熵定义为：

$$H(X) = - \int P(x) \log P(x) dx = \mathbb{E}_{P(X)}[-\log P(X)]$$

数据的分布越平缓、越难以预测，其包含的不确定度（熵值）就越大。

条件熵 (Conditional Entropy) $H(X|Y)$ 度量了在已经观测到变量 Y 的取值后，变量 X 仍然保留的残余不确定度：

$$H(X|Y) = - \iint P(x, y) \log P(x|y) dx dy = \mathbb{E}_{P(X, Y)}[-\log P(X|Y)]$$

互信息 (Mutual Information) $I(X; Y)$ 是信息瓶颈理论的绝对主角。它度量了两个变量之间共享的信息量，或者说：因为知道了 Y ，使得我们对 X 的不确定度减少了多少。在数学上，它被极其优美地定义为边缘熵与条件熵之差：

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

互信息还有一个非常关键的等价定义，它直接揭示了互信息与变分推断中常用的 KL 散度的关系：

$$I(X; Y) = \iint P(x, y) \log \frac{P(x, y)}{P(x)P(y)} dx dy = \mathcal{D}_{KL}(P(X, Y) \parallel P(X)P(Y))$$

互信息的本质，就是在衡量“联合分布”与“独立分布”之间的 KL 距离。距离越大，独立性越明显，说明 X 和 Y 绑得越死，共享的信息就越多。

1.4.1 第一步：用变分下界近似 $I(Z; Y)$ (充分性项)

互信息 $I(Z; Y)$ 定义为：³

$$I(Z; Y) = \int P(Z, Y) \log \frac{P(Y|Z)}{P(Y)} dZ dY$$

这里 $P(Y|Z)$ 是真实的后验，算不出来。所以我们需要给其寻找一个可以计算的下界。

我们引入带有参数 θ 的变分预测器 $P_\theta(Y|Z)$ ，利用 KL 散度的非负性定理 ($\mathcal{D}_{KL}(P(Y|Z) \parallel P_\theta(Y|Z)) \geq 0$)，可以推导出：

$$\int P(Y|Z) [\log P(Y|Z) - \log P_\theta(Y|Z)] dY \geq 0.$$

因为概率密度函数必然是非负的 (即 $P(Z) \geq 0$)，所以在不等式两边同时乘以 $P(Z)$ 并对 Z 积分，不等号方向不会改变。

$$\int P(Z, Y) \log P(Y|Z) dZ dY \geq \int P(Z, Y) \log P_\theta(Y|Z) dZ dY$$

将其代入互信息公式：

$$\begin{aligned} I(Z; Y) &\geq \int P(Z, Y) \log P_\theta(Y|Z) dZ dY - \int P(Y) \log P(Y) dY \\ &= \mathbb{E}_{P(X, Y)} [\mathbb{E}_{q_\phi(Z|X)} [\log P_\theta(Y|Z)]] + H(Y) \end{aligned}$$

由于熵 $H(Y)$ 是一个常数，最大化互信息 $I(Z; Y)$ 等价于最大化对数预测似然的期望 $\mathbb{E}_{q_\phi(Z|X)} [\log P_\theta(Y|Z)]$ 。

1.4.2 第二步：用变分上界近似 $I(Z; X)$ (压缩项)

互信息 $I(Z; X)$ 定义为：

$$I(Z; X) = \int P(Z, X) \log \frac{q_\phi(Z|X)}{P_{true}(Z)} dZ dX = \iint P(X) q_\phi(Z|X) \log \frac{q_\phi(Z|X)}{P_{true}(Z)} dZ dX$$

真实的边缘分布 $P_{true}(Z) = \int q_\phi(Z|X) P(X) dX$ 极其复杂，算不出来。我们引入一个简单的先验分布 $P(Z) = \mathcal{N}(0, \mathbf{I})$ 来近似它。同样利用 KL 散度非负性，推导出上界：

$$\begin{aligned} I(Z; X) &= \iint P(X) q_\phi(Z|X) \log \left(\frac{q_\phi(Z|X)}{P(Z)} \cdot \frac{P(Z)}{P_{true}(Z)} \right) dZ dX \\ &= \underbrace{\iint P(X) q_\phi(Z|X) \log \frac{q_\phi(Z|X)}{P(Z)} dZ dX}_{\text{目标变分项}} + \underbrace{\iint P(X) q_\phi(Z|X) \log \frac{P(Z)}{P_{true}(Z)} dZ dX}_{\text{误差补偿项}} \\ &= \text{目标变分项} + \int \left[\int P(X) q_\phi(Z|X) dX \right] \log \frac{P(Z)}{P_{true}(Z)} dZ \\ &= \text{目标变分项} + \int P_{true}(Z) \log \frac{P(Z)}{P_{true}(Z)} dZ \\ &= \text{目标变分项} - \mathcal{D}_{KL}(P_{true}(Z) \parallel P(Z)) \\ &\leq \text{目标变分项} \\ &= \int P(X) \left[\int q_\phi(Z|X) \log \frac{q_\phi(Z|X)}{P(Z)} dZ \right] dX \\ &= \mathbb{E}_{P(X)} [\mathcal{D}_{KL}(q_\phi(Z|X) \parallel P(Z))] \end{aligned}$$

最小化互信息 $I(Z; X)$ 等价于最小化推断分布与先验之间的 KL 散度 $\mathcal{D}_{KL}(q_\phi(Z|X) \parallel P(Z))$ 。

1.4.3 第三步：组合得到最终的 VIB-变分下界

将第一步 (下界) 和第二步 (上界) 代入原始的 \mathcal{J}_{IB} 中，并忽略常数项 $H(Y)$ ，我们得到了 VIB 可优化的目标函数：

³注意，在信息流中，我们已经讲过，是 X 生成 Z ， Z 再生成 Y 。这里就决定了互信息定义中的条件概率。

$$\mathcal{L}_{VIB}(\phi, \theta) = \underbrace{\mathbb{E}_{q_\phi(Z|X)}[\log P_\theta(Y|Z)]}_{1. \text{ 预测似然 (例如: 交叉熵)}} - \beta \cdot \underbrace{\mathcal{D}_{KL}(q_\phi(Z|X) \parallel P(Z))}_{2. \text{ 正则化/压缩项}} \quad (1-4)$$

如果你将上述公式与算例 C (VAE) 进行对比, 会发现其数学形式如出一辙! 但在物理意义上有着根本的分歧:

- **重构 vs. 预测:** VAE 的第一项是 $\log P_\theta(X|Z)$, 要求 Z 必须包含 X 的所有像素/基因细节; VIB 的第一项是 $\log P_\theta(Y|Z)$, 要求 Z 只需包含区分 Y 的关键特征。
- **KL 散度的不同使命:** 在 VAE 中, KL 散度是为了让隐空间变得平滑, 方便采样生成。在 VIB 中, KL 散度是无情的“信息粉碎机”——它强迫 Z 向 $\mathcal{N}(0, \mathbf{I})$ 坍缩, 只有那些对预测 Y 极其重要 (能够降低大量交叉熵损失) 的特征, 才能抵抗住这股拉力被保留下来。

1.5 算法实现的若干问题

1.5.1 神经网络的结构

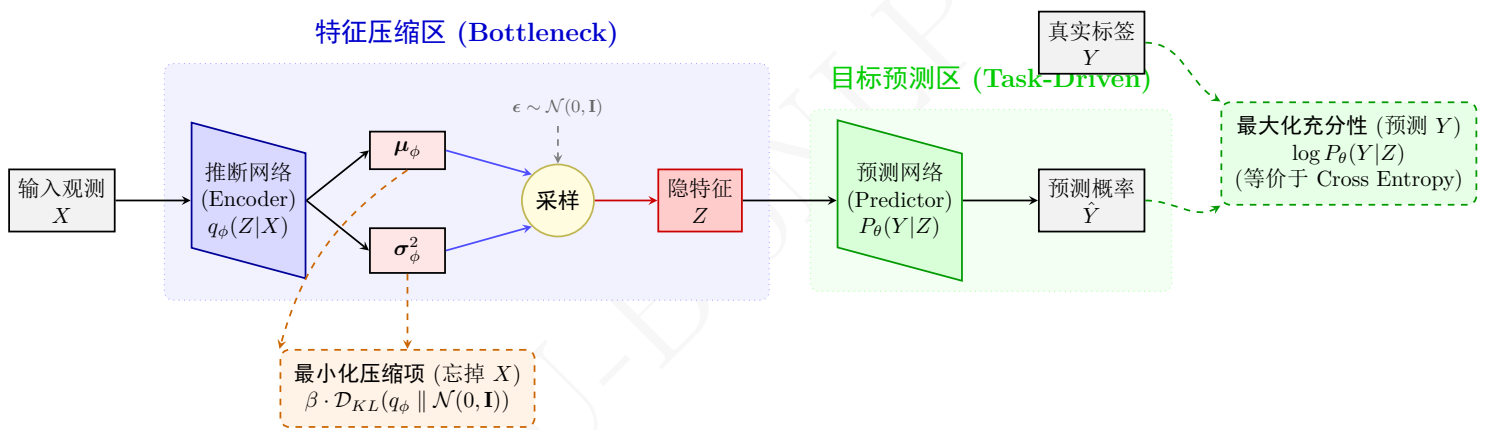


图 3: 变分信息瓶颈 (VIB) 神经网络实现架构图。网络在结构上犹如一个漏斗, 利用 KL 散度在 Z 处实施极致压缩 (遗忘 X 的噪音), 同时利用交叉熵迫使 Z 拼命保留能预测 Y 的核心语义信息。

- Encoder (ϕ): 负责将高维 X 映射为分布参数 μ 和 σ^2 。它受 KL 散度的向内挤压, 被迫丢弃噪音。
- Bottleneck (Z): 通过重参数化生成的随机变量, 是信息的唯一出口。
- Predictor (θ): 负责从 Z 中恢复 Y 。它受交叉熵的向外拉扯, 迫使编码器保留最关键的预测语义。
- β 超参数: 它是“水龙头的开关”。 β 越大, 瓶颈越窄, 特征越压缩, 可解释性可能更强, 但预测准确度可能下降 (欠拟合); β 越小, 模型越接近传统的深度分类器, 泛化能力可能下降 (过拟合)。

1.5.2 从变分下界到 loss 函数

1. 充分性项 \rightarrow 预测损失 (Classification Loss)

这一项 $\mathbb{E}_{q_\phi(Z|X)}[\log P_\theta(Y|Z)]$ 要求从压缩后的特征 Z 中准确预测标签 Y 。对于离散分类任务, $P_\theta(Y|Z)$ 建模为类别分布 (Categorical Distribution)。

最大化对数似然 $\log P_\theta(Y|Z)$ 在代码实现上完全等价于最小化交叉熵 (Cross Entropy)。

工程实现: 通过编码器采样得到 z 。将 z 送入预测网络 (Predictor/Decoder) 得到 logits。计算预测值与真实标签 Y 的距离。

PyTorch 代码

```
# y_pred 是 Predictor 输出的 Logits, y_true 是真实标签
loss_classification = F.cross_entropy(y_pred, y_true)
```

2. 最小性项 → 压缩损失 (KL Divergence Loss)

这一项 $\mathcal{D}_{KL}(q_\phi(Z|X) \parallel P(Z))$ 负责“瓶颈”挤压。它约束编码器输出的分布不要偏离先验分布（通常为标准正态分布 $\mathcal{N}(0, \mathbf{I})$ ）太远。

神经网络实现中，设编码器输出为对角高斯分布 $q_\phi(Z|X) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$ ，先验为 $P(Z) = \mathcal{N}(0, \mathbf{I})$ 。可以计算出解析解。利用高斯分布间 KL 散度的闭式解，该项转化为：

$$\mathcal{D}_{KL} = \frac{1}{2} \sum_{j=1}^d (\mu_j^2 + \sigma_j^2 - \log \sigma_j^2 - 1)$$

此处， μ^2 惩罚均值偏离原点， $\sigma^2 - \log \sigma^2$ 惩罚方差偏离 1。

PyTorch 代码

```
# mu 和 logvar 是 Encoder 输出的变分参数
loss_kl = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())
```

VIB Loss 函数代码汇总

```
import torch
import torch.nn as nn
import torch.nn.functional as F

def vib_loss_function(y_true, y_logits, mu, log_var, beta=1e-3):
    """
    y_true: 真实的标签 Y
    y_logits: 预测网络 P_theta(Y|Z) 输出的预测 Logits
    mu, log_var: 编码器 q_phi(Z|X) 输出的均值和对数方差
    beta: 信息压缩系数
    """

    # --- 1. 预测似然项 (Prediction Likelihood) ---
    # E_q [log P_theta(Y|Z)]
    # 对于分类任务，这等价于负的交叉熵损失
    # 注意：这里的 y_logits 是通过重参数化后的 Z^{(1)} 计算得出的
    prediction_loss = F.cross_entropy(y_logits, y_true, reduction='mean')

    # --- 2. 压缩项 (Compression via KL Divergence) ---
    # D_KL(q_phi(Z|X) || P(Z))
    # 解析展开与 VAE 完全相同
    kl_div = -0.5 * torch.sum(1 + log_var - mu.pow(2) - log_var.exp(), dim=1)
    kl_loss = torch.mean(kl_div)

    # --- 最终目标 ---
    # 最小化：交叉熵 + beta * KL散度
    total_loss = prediction_loss + beta * kl_loss

    return total_loss, prediction_loss, kl_loss
```

1.5.3 利用重参数化计算梯度

与 VAE 完全一样，由于期望算子内部包含了由 ϕ 参数化的 $q_\phi(Z|X)$ ，我们必须使用重参数化技巧 $Z^{(l)} = \mu_\phi(X) + \sigma_\phi(X) \cdot \epsilon^{(l)}$ 来保证梯度的平滑传导。在代码实现过程中，有 VAE 中使用过的编程技巧，譬

如 `std = torch.exp(0.5 * logvar)`，这是为什么呢？⁴

Pytorch 代码

```
import torch

def reparameterize(mu, logvar):
    """
    重参数化技巧的 PyTorch 实现
    :param mu: 推断网络输出的均值向量 (Mean)
    :param logvar: 推断网络输出的对数方差向量 (Log-Variance)
    :return: 采样得到的隐变量 Z (可导)
    """
    # 1. 恢复标准差 std = exp(0.5 * logvar)
    # 技巧解析：网络直接输出 logvar 而不是方差，是为了让输出值域覆盖整个实数界，
    # 从而自然地规避了方差必须非负的数值约束限制。
    std = torch.exp(0.5 * logvar)

    # 2. 独立采样噪声 epsilon ~ N(0, I)
    # torch.randn_like 会自动生成一个形状和设备(CPU/GPU)与 std 完全一致的张量
    eps = torch.randn_like(std)

    # 3. 核心仿射变换: Z = mu + sigma * epsilon
    # 此处的乘法为张量的逐元素相乘 (element-wise multiplication)。
    # 经过这一步，原本不可导的随机性被转移到了游离的 eps 上，
    # 梯度得以安全地穿过 mu 和 std 传回底层网络。
    z = mu + eps * std

    return z
```

总结：变分信息瓶颈（VIB）是理解现代深度学习中“为什么添加噪声和正则化能提高泛化能力”的一把终极钥匙。它证明了：一个具有强大泛化能力的判别模型，其本质上是在进行一场带有变分近似的信息压缩游戏。

⁴为什么是 `std = torch.exp(0.5 * logvar)`？

在代码中，`logvar` 代表的是方差的自然对数，即 $\log(\sigma^2)$ 。我们需要的是标准差 σ 。根据对数的基本运算性质，指数的提取公式为：

$$\log(\sigma^2) = 2 \log(\sigma)$$

将两边同时除以 2：

$$\log(\sigma) = \frac{1}{2} \log(\sigma^2) = 0.5 \times \logvar$$

最后，两边同时取自然指数（`exp`）以消除对数，就完美得到了代码中的等式：

$$\sigma = \exp(0.5 \times \logvar)$$

如果不这样做，直接输出会遇到什么问题呢？

在概率论中，方差 σ^2 和标准差 σ 必须是严格大于 0 的实数。然而，神经网络的线性层（如 `nn.Linear`）的输出值域是 $(-\infty, +\infty)$ 。如果我们让网络直接输出 σ 或 σ^2 ，在工程上就必须加一个激活函数（比如 `ReLU` 或 `Softplus`）来强制把负数截断或映射到正数区间。如果用 `ReLU`：一旦输出为负，梯度直接变成 0，导致神经元“死亡”，方差再也无法更新。如果用 `Softplus`：在数值极大或极小时容易遇到计算溢出等数值不稳定的问题。