

潜在空间自回归用于新颖性检测

达维德·阿巴蒂 安杰洛·波雷洛 西蒙内·卡尔德拉拉 丽塔·库基亚拉

摩德纳和雷焦艾米利亚大学

{名.姓}@unimore.it

摘要

异常检测通常指对不符合已学习常态模型的观测数据进行区分的任务。尽管在不同应用场景中具有重要意义，但由于异常现象的不可预测性及其在训练过程中的不可触及性——这些因素揭示了该问题的无监督本质——设计异常检测器变得极为复杂。在我们的方案中，我们设计了一个通用框架：通过自回归过程为深度自编码器配备参数化密度估计器，从而学习其潜在表征所遵循的概率分布。我们证明，在正常样本重构过程中联合优化的最大似然目标函数，可通过最小化潜在向量所构成分布的微分熵，有效充当当前任务的正则化器。除了提供高度通用的 formulation，我们在公开数据集上的大量实验表明：与单类别异常检测和视频异常检测领域的先进方法相比，本模型实现了相当或更优的性能。与先前研究不同，我们的方案不对异常性质做任何假设，使其能够直接适用于多样化场景。

1. 引言

新颖性检测的定义是，识别那些相对于基于正常样本集合建立的规律性基础模型表现出显著不同特征的样本。自主系统识别未知事件的能力使其在多个领域得到应用，范围从视频监控 [7, 11]，到缺陷检测 [22] 再到医学成像 [38]。此外，由未见事件引发的惊奇感正逐渐成为强化学习环境中的关键要素，作为驱动好奇心探索的赋能因素 [34]。

然而，在此设定下，无法对新型示例进行定义和标记。因此，文献

在通过建模正常样本的内在特性来逼近区分正常样本与新颖样本的理想边界形状方面达成共识。因此，先前的研究遵循无监督学习范式 [9, 37, 11, 26, 30] 衍生的原则来解决这一问题。由于缺乏监督信号，特征提取过程及其正常性评估规则只能通过代理目标进行引导，并假设后者能为当前应用定义合适的边界。

根据认知心理学 [4]，新奇性既可以通过记忆事件的能力来表征，也可通过观测事件时引发的惊异程度 [42] 来衡量。后者在数学上通过低概率事件发生于预期模型之下的概率来建模，或通过降低变分自由能 [16] 来实现。在此框架中，先验模型采用参数化 [49] 或非参数化 [14] 密度估计器。与之不同，记忆事件意味着采用某种记忆表征：或通过正常原型字典（如稀疏编码方法 [9]），或通过输入空间的低维表征（如自组织映射 [20]，或近年兴起的深度自编码器）。因此在新奇性检测中，特定样本的记忆能力可通过测量重构误差 [11, 26] 或执行判别性分布内测试 [37] 来评估。

我们的提案通过将记忆与惊奇度方面融合到一个独特框架中，为该领域做出贡献：我们设计了一种生成式无监督模型（即自编码器，如图 1i 所示），该模型利用端到端训练来最大化正常样本的记忆效率，同时最小化其潜在表示的惊奇度。后一点是通过自回归密度估计器最大化潜在表示的可能性实现的，该过程与重构误差最小化同步进行。我们表明，通过联合优化这两个项，模型会隐式地寻求保持其记忆 / 重构能力的最小熵表示。虽然熵最小化方法已在深度神经压缩中得到应用 [3]，但据我们所知，这是首个将 ...

专为新颖性检测而设计。在内存方面，我们的程序类似于使用尽可能少的模板来构建正常性原型的概念。此外，通过评估估计器的输出，能够衡量给定样本所引起的惊讶程度。

2. 相关工作

基于重构的方法。一方面，许多研究倾向于学习正常数据的参数化投影与重构，假定异常值会产生更高残差。传统稀疏编码算法 [48, 9, 27] 遵循此框架，在假设新样本会在学习子空间中呈现非稀疏表示的前提下，将正常模式表示为若干基组分的线性组合。近期研究中，投影步骤通常源于深度自编码器 [11]。[30] 作者通过对学习表征施加稀疏正则化来恢复稀疏编码原理，同时利用循环神经网络增强其在时间维度上的平滑性。[37]，则另辟蹊径，采用对抗框架——将判别器网络作为实际的新颖性检测器，通过执行离散分布内测试来识别异常。与之相反，未来帧预测 [26] 通过利用历史帧知识来最大化下一帧的期望值；测试时，观测内容与预测内容之间的显著偏差即提示异常。与上述研究不同，我们的方案依赖于对潜在表征先验分布的建模。这一选择与密度估计领域的最新研究 [41, 6] 相契合。但据我们所知，本研究是首个论证该设计选择对新奇检测重要性的工作。

概率化方法。另一条互补的研究路线致力于探索近似正常外观与运动特征密度函数的不同策略。该领域的核心问题在于如何在高维复杂特征空间中估算此类密度函数。对此，早期研究采用光流或轨迹分析等手工设计特征，并在此基础上结合了非参数 [1] 与参数

[5, 31, 25] 估计器，以及图模型 [17, 23]。现代方法则依托深度表征（如自编码器提取的特征），例如高斯分类器 [36] 和高斯混合模型 [49] 的应用。[14] 研究者引入核密度估计器（KDE）对辅助物体检测网络的激活值进行建模。近期研究趋势关注在正常样本上训练生成对抗网络（GANs），但因此类模型逼近的是隐式密度函数，可通过查询生成新样本。

但不适用于似然值。因此，基于 GAN 的模型采用不同的启发式方法来评估新颖性。例如，在 [38] 中利用引导潜空间搜索进行推断，而 [35] 则直接向判别器查询正态性得分。

3. 提出的模型

最大化潜在表示的概率类似于降低模型对正常配置的惊奇度，其定义为潜在变量实例 [42] 的负对数密度。反之，记忆能力可以通过给定样本在其潜在表示下的重构精度来评估。

我们在潜变量模型设置中对上述方面进行建模，其中训练样本 $p(\mathbf{x})$ 的概率密度函数通过辅助随机变量 \mathbf{z} 进行建模，该变量描述了所有观测值背后的因果因子集合。通过因子分解

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}, \quad (1)$$

其中 $p(\mathbf{x}|\mathbf{z})$ 是给定潜在表示 \mathbf{z} 及先验分布 $p(\mathbf{z})$ 的观测条件似然函数，我们可以显式分解记忆与信息惊奇对新奇性的贡献。通过采用负责识别潜在空间向量的推理模型，我们近似边缘化处理，该模型旨在找到使 $p(\mathbf{x}|\mathbf{z})$ 贡献最大的潜向量。形式上，我们采用深度自编码器，在假设 $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\tilde{\mathbf{x}}, I)$ 的前提下，其重构误差扮演 $p(\mathbf{x}|\mathbf{z})$ 负对数的角色（此处 $\tilde{\mathbf{x}}$ 表示输出重构）。此外，通过为自编码器配备辅助深度参数估计器来学习潜在向量的先验分布 $p(\mathbf{z})$ ，并采用最大似然估计进行训练，将信息惊奇注入流程。因此我们的架构包含三个核心模块（图 1i）：编码器 $f(\mathbf{x}; \theta_f)$ 、解码器 $g(\mathbf{z}; \theta_g)$ 以及概率模型 $h(\mathbf{z}; \theta_h)$ ：

$$\begin{aligned} f(\mathbf{x}; \theta_f) : \mathbb{R}^m &\rightarrow \mathbb{R}^d, & g(\mathbf{z}; \theta_g) : \mathbb{R}^d &\rightarrow \mathbb{R}^m, \\ h(\mathbf{z}; \theta_h) : \mathbb{R}^d &\rightarrow [0, 1]. \end{aligned} \quad (2)$$

编码器处理输入 \mathbf{x} 并将其映射为压缩表示 $\mathbf{z} = f(\mathbf{x}; \theta_f)$ ，而解码器则提供输入的重构版本 $\tilde{\mathbf{x}} = g(\mathbf{z}; \theta_g)$ 。概率模型 $h(\mathbf{z}; \theta_h)$ 通过自回归过程估计 \mathbf{z} 中的密度，从而避免采用可能对当前任务无益的特定分布族（例如高斯分布）。关于后一点，请参阅补充材料以与变分自编码器 [19] 进行比较。

借助此类模块，在测试时我们可以评估两种新颖性来源：观测效果较差的元素

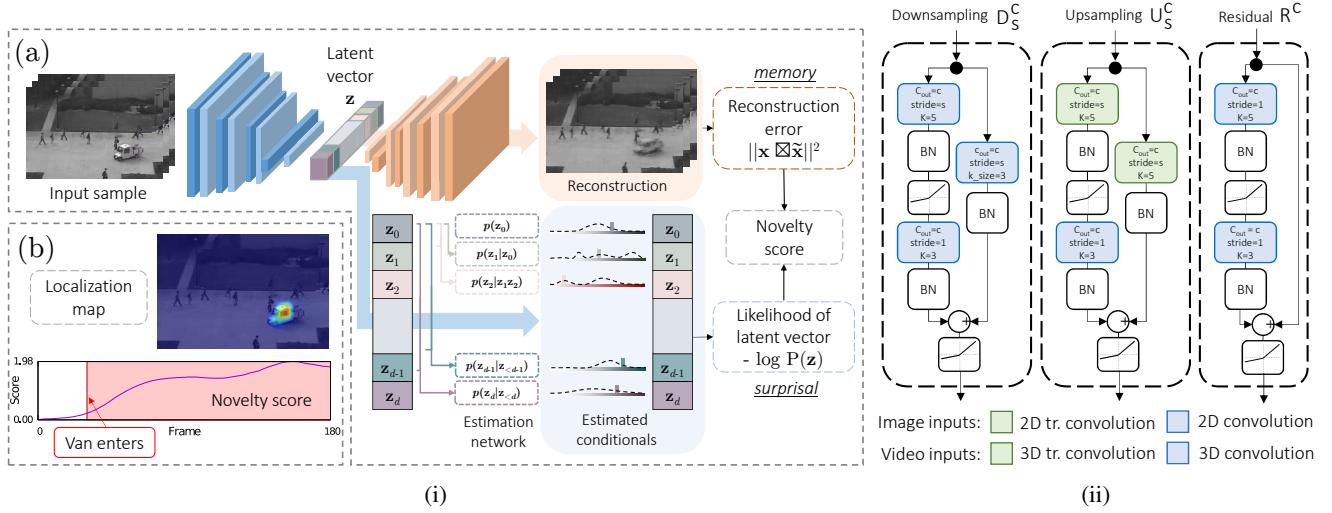


图 1: (i) 提出的新颖性检测框架。整体架构如 (a) 所示, 包含一个深度自编码器及其潜在空间上运行的自回归估计网络。通过联合最小化各自目标函数, 可获得新颖性度量 ——(b)— 该度量通过评估模型观察新样本时的记忆强度, 并结合因果因素引发的惊异值来实现。 (ii) 自编码器架构中采用的构建模块。

由正常样本引入的因果因素 (即高重构误差) 所解释; 在学习的先验条件下, 表现出良好重构效果, 同时显示出令人惊讶的潜在表征的元素。

自回归密度估计。自回归模型为涉及序列预测的任务提供了通用框架, 其中每个输出都依赖于先前的观测结果 [28, 32]。我们采用这种技术来分解联合分布, 从而避免先验地定义其分布形态 [24, 43]。形式上, $p(\mathbf{z})$ 被分解为

$$p(\mathbf{z}) = \prod_{i=1}^d p(z_i | \mathbf{z}_{<i}), \quad (3)$$

因此, 估计 $p(\mathbf{z})$ 可简化为对每个条件概率密度 (CPD) $p(z_i | \mathbf{z}_{<i})$ 的单独估计, 其中符号 $<$ 表示随机变量间的顺序关系。部分先验模型遵循人工设计的顺序 [46, 45], , 而其他模型则采用顺序无关的训练方法 [44, 10]。然而, 如何确定给定变量集的最优顺序仍不明确。在我们的模型中, 该问题通过优化过程直接解决。由于我们对学到的潜在表示进行自回归操作, 最大似然估计目标会驱动自编码器在这些表示上施加预定义的因果结构。补充材料中提供了这一现象的实验证据。

从技术角度来看, 估计器 $h(\mathbf{z}; \theta_h)$ 会输出 d 分布 $p(z_i | \mathbf{z}_{<i})$ 的参数。在我们的实现中, 每个条件概率分布被建模为 $B=100$ 量化箱上的多项式分布。为确保对每个变量的条件估计

在底层密度的基础上, 我们设计了合适的层, 确保每个符号 z_i 的条件概率分布仅从输入 $\{z_1, \dots, z_{i-1}\}$ 计算得出。

目标及其与微分熵的关联。三个组件 f 、 g 和 h 通过联合训练最小化 $\mathcal{L} \equiv \mathcal{L}(\theta_f, \theta_g, \theta_h)$, 具体实现方式如下:

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{\text{REC}}(\theta_f, \theta_g) + \lambda \mathcal{L}_{\text{LLK}}(\theta_f, \theta_h) \\ &= \mathbb{E}_{\mathbf{x}} \left[\underbrace{\|\mathbf{x} - \tilde{\mathbf{x}}\|^2}_{\text{reconstruction term}} - \lambda \underbrace{\log(h(\mathbf{z}; \theta_h))}_{\text{log-likelihood term}} \right], \end{aligned} \quad (4)$$

其中 λ 是控制 \mathcal{L}_{LLK} 项权重的超参数。值得注意的是, 可以将对数似然项表示为

$$\begin{aligned} &\mathbb{E}_{\mathbf{z} \sim p^*(\mathbf{z}; \theta_f)} [-\log h(\mathbf{z}; \theta_h)] \\ &= \mathbb{E}_{\mathbf{z} \sim p^*(\mathbf{z}; \theta_f)} [-\log h(\mathbf{z}; \theta_h) + \log p^*(\mathbf{z}; \theta_f) - \log p^*(\mathbf{z}; \theta_f)] \\ &= D_{\text{KL}}(p^*(\mathbf{z}; \theta_f) \| h(\mathbf{z}; \theta_h)) + \mathbb{H}[p^*(\mathbf{z}; \theta_f)], \end{aligned} \quad (5)$$

其中 $p^*(\mathbf{z}; \theta_f)$ 表示编码器生成代码的真实分布, 因此由 θ_f 参数化。这种对 MLE 目标的重构为优化中涉及的实体提供了有意义的见解。一方面, Kullback-Leibler 散度确保参数模型 h 与真实分布 p^* 之间的信息差距很小。另一方面, 该框架导致编码器 f 生成的代码所基于分布的微分熵最小化。在学习正态性时, 这种约束构成了一个关键点。直观地说, 如果我们考虑

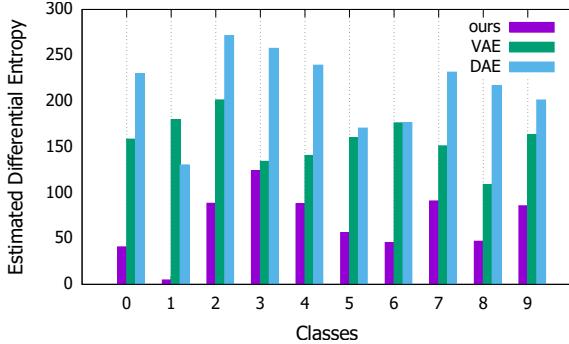


图 2：在不同正则化策略下，各 MNIST 类别估计得到的微分熵：我们的方法、高斯先验散度（VAE）和输入扰动（DAE）。针对每个类别，该估计值基于训练样本的隐层表示计算得出，其分布通过三维空间中的高斯核密度估计进行拟合。在模型结构相同的情况下，我们的方法在所有类别上均展现出更低的熵值。

关于编码器作为发射符号（即潜在表示）的源，其在建模数据中的正常方面时，期望行为应收敛于一种以内在低熵为特征的“乏味”过程——因为训练阶段不太可能出现令人惊讶的新颖事件。相应地，在所有可能的隐表示设置中，该目标要求编码器展现出低微分熵，从而提取易于预测的特征，这些特征因此成为训练集中常见且重复出现的模式。这类特征确实最能有效区分新颖样本与正常样本，使得我们的提案成为异常检测场景中合适的正则化器。

我们在图 2 中报告了微分熵递减的经验证据，该图比较了同一模型在不同正则化策略下的行为。

3.1. 架构组件

自编码器模块。 编码器和解码器分别由图 1ii 所示的下采样和上采样残差块构成。编码器末端为全连接层。处理视频输入时，我们在编码器中采用因果三维卷积 [2]（即仅访问先前时间步的信息）。此外，在编码器末端采用时序共享全连接层（TFC，即在输入特征图上沿时间轴共享参数的线性投影），生成特征向量的时间序列。这种编码方式不会跨时间步混淆信息，从而确保时序顺序。

自回归层。 为保证每个输出条件概率分布的自回归特性，我们需要确保

估计器 h 每一层的连接模式。此外，由于潜在表征根据输入性质（图像或视频）会呈现不同形态，我们提出了两种不同的解决方案。

处理图像时，编码器会生成维度为 d 的特征向量。自回归估计器通过堆叠多个掩码全连接层（MFC，图 3-(a)）构成。其形式化计算可表示为：给定输入特征图 $\mathbf{h} \in \mathbb{R}^{d \times ci}$ （假定输入层维度为 $ci = 1$ ），输出特征图 $\mathbf{o} \in \mathbb{R}^{d \times co}$ （其中 co 表示输出通道数）。位于位置 i 、通道 k 的输入元素 \mathbf{h}_i^k 与输出元素 \mathbf{o}_j^l 之间的连接关系由以下参数定义：

$$\begin{cases} w_{i,j}^{k,l} & \text{if } i < j \\ w_{i,j}^{k,l} & \text{if type = B} \\ 0 & \text{if type = A} \\ 0 & \text{if } i > j. \end{cases} \quad (6)$$

类型 A 强制对先前元素进行严格依赖（仅作为第一个估计器层使用），而类型 B 仅遮蔽后续元素。假设每个条件概率分布被建模为多项式分布，最后一个自回归层（位于 $\mathbb{R}^{d \times B}$ 中）的输出为构成空间量化的 B 区间提供概率估计。

另一方面，视频片段的压缩表示维度为 $t \times d$ ，其中 t 表示时间步数， d 表示代码长度。相应地，估计网络被设计用于捕捉代码观测元素中的二维模式。然而，直接使用二维卷积层会假设输入映射的两个轴都具有平移不变性，而由于压缩表示的构建方式，这种假设仅沿时间轴成立。为解决此问题，我们沿代码轴应用 d 个不同的卷积核，使其既能观测前一时刻的完整特征向量，也能捕捉当前时刻的部分特征。每次卷积可沿时间轴自由滑动以捕获时序模式。在这种称为掩码堆叠卷积（MSC，图 3-(b)）的操作中，第 i 个卷积配备的卷积核 $\mathbf{w}^{(i)} \in \mathbb{R}^{3 \times d}$ 会与二进制掩码 $\mathbf{M}^{(i)}$ 相乘，该掩码定义为

$$m_{j,k}^{(i)} \in \mathbf{M}^{(i)} = \begin{cases} 1 & \text{if } j = 0 \\ 1 & \text{if } j = 1 \text{ and } k < i \text{ and type=A} \\ 1 & \text{if } j = 1 \text{ and } k \leq i \text{ and type=B} \\ 0 & \text{otherwise,} \end{cases} \quad (7)$$

其中 j 索引时间轴， k 索引代码轴。
每次卷积都会产生一个列向量，这是其在时间维度上滑动步长的结果。由这些卷积操作产生的一系列列向量

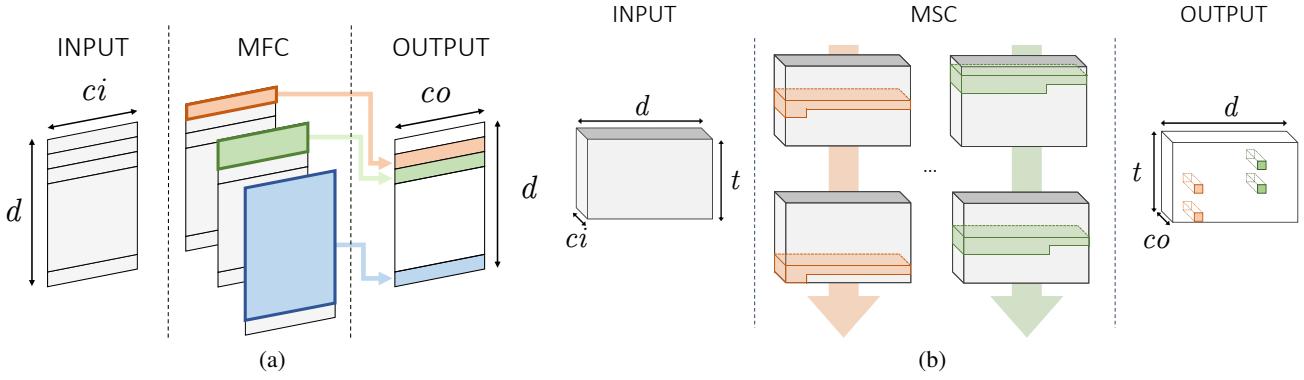


图3：提出的自回归层，即掩码全连接层（a，公式6）与掩码堆叠卷积层（b，公式7）。两种层结构均采用A型架构，不同颜色的卷积核代表不同的参数化形式。

将 d 卷积应用于输入张量 $\mathbf{h} \in \mathbb{R}^{t \times d \times ci}$ 后，通过水平堆叠构建输出张量 $\mathbf{o} \in \mathbb{R}^{t \times d \times co}$ ，具体如下：

$$\mathbf{o} = \left\| \left[(\mathbf{M}^{(i)} \odot \mathbf{w}^{(i)}) * \mathbf{h} \right] \right\|_{i=1}^d, \quad (8)$$

其中 $\|$ 表示水平拼接操作。

4. 实验¹

我们在三种不同设置下测试我们的解决方案：图像、视频和认知数据。所有实验中，对第 i 个样本的新颖性评估均通过将重构项（ REC_i ）和似然对数项（ LLK_i ）在公式4中合并为单一新颖性分数 NS_i 来实现：

$$NS_i = norm_S(REC_i) + norm_S(LLK_i). \quad (9)$$

个体分数通过参考示例集 S 进行归一化（每个实验有所不同），

$$norm_S(L_i) = \frac{L_i - \max_{j \in S} L_j}{\max_{j \in S} L_j - \min_{j \in S} L_j}. \quad (10)$$

进一步的实现细节和架构超参数详见补充材料。

4.1. 图像单类新颖性检测

为了评估模型在单类设置中的性能，我们分别在MNIST或CIFAR-10的每个类别上单独训练模型。在测试阶段，我们使用由所有类别共10000个样本组成的对应测试集，并期望模型能为与训练样本标签相同的图像分配较低的新颖性分数。我们采用标准的训练/测试分割，并隔离10%的训练样本用于

出于验证目的，我们将其作为归一化集（公式9中的 S ）用于新颖性得分的计算。关于基线方法，我们考虑以下方案：

- 采用PCA白化提取特征的标准方法，如OC-SVM [39] 和核密度估计器（KDE）；
- 使用与我们提案相同架构的去噪自编码器（DAE），但缺少密度估计模块。其重构误差被用作正常性与新颖性的衡量指标；
- 变分自编码器（VAE）[19]，同样采用与我们模型相同的架构，其中使用证据下界（ELBO）作为评分标准；
- 在图像空间直接应用自回归进行密度建模的Pix-CNN [45]；
- 基于GAN的方法如文献[38]所述。

我们在表1中报告了比较结果，其中性能通过接收者操作特征曲线下面积（AUROC）进行衡量，这是该任务的标准评估指标。如表所示，我们的方法在两种设置下均优于所有基线模型。

考虑到MNIST数据集，大多数方法都表现良好。值得注意的是，Pix-CNN在建模除一个数字外的所有数字分布时均告失败，这可能是由于直接在像素空间建模密度并遵循固定自回归顺序的复杂性。尽管在训练过程中我们观察到了高质量的样本，但测试性能仍然很差：事实上，样本质量与模型测试对数似然之间的弱相关性已在{v1}中得到解释。令人惊讶的是，在这种设置下OC-SVM的表现优于大多数基于深度学习的模型。相反，CIFAR10则代表了更为严峻的挑战，大多数模型的低性能便是明证，这可能是由于图像分辨率低以及类别间视觉杂乱造成的。具体而言，我们观察到

¹本节结果重现代码发布于 <https://github.com/aimagelab/novelty-detection>。

	MNIST						CIFAR10							
	OC SVM	KDE	DAE	VAE	Pix CNN	GAN	ours	OC SVM	KDE	DAE	VAE	Pix CNN	GAN	ours
0	0.988	0.885	0.991	0.998	0.531	0.926	0.993	0.630	0.658	0.718	0.688	0.788	0.708	0.735
1	0.999	0.996	0.999	0.999	0.995	0.995	0.999	0.440	0.520	0.401	0.403	0.428	0.458	0.580
2	0.902	0.710	0.891	0.962	0.476	0.805	0.959	0.649	0.657	0.685	0.679	0.617	0.664	0.690
3	0.950	0.693	0.935	0.947	0.517	0.818	0.966	0.487	0.497	0.556	0.528	0.574	0.510	0.542
4	0.955	0.844	0.921	0.965	0.739	0.823	0.956	0.735	0.727	0.740	0.748	0.511	0.722	0.761
5	0.968	0.776	0.937	0.963	0.542	0.803	0.964	0.500	0.496	0.547	0.519	0.571	0.505	0.546
6	0.978	0.861	0.981	0.995	0.592	0.890	0.994	0.725	0.758	0.642	0.695	0.422	0.707	0.751
7	0.965	0.884	0.964	0.974	0.789	0.898	0.980	0.533	0.564	0.497	0.500	0.454	0.471	0.535
8	0.853	0.669	0.841	0.905	0.340	0.817	0.953	0.649	0.680	0.724	0.700	0.715	0.713	0.717
9	0.955	0.825	0.960	0.978	0.662	0.887	0.981	0.508	0.540	0.389	0.398	0.426	0.458	0.548
avg	0.951	0.814	0.942	0.969	0.618	0.866	0.975	0.586	0.610	0.590	0.586	0.551	0.592	0.641

表 1: MNIST 和 CIFAR10 上新奇检测的 AUROC 结果。每行代表基线模型与我们的模型所训练的不同类别。

我们的提案是唯一优于简单 KDE 基线的模型；但这一发现需结合非参数估计器的特性进行辩证看待。事实上，非参数模型被允许访问整个训练集以评估每个样本。因此，尽管它们在大样本集的密度建模方面具有优势，但随着数据集规模增大，会导致推理过程难以实现。

关于与 DAE 性能差异的可能原因有两点。首先，DAE 仅能基于重构误差识别新样本，因此依赖其记忆能力；而我们的方案还考虑了这些表征在所学先验下的似然度，从而同时利用了惊奇度。其次，通过最小化潜在分布的微分熵，我们的方案增强了重构的判别能力。直观而言，最后这一论断的合理性在于：新颖样本被迫位于潜在空间的高概率区域，而该区域仅能捕捉训练集中出现的非惊奇变异因素。另一方面，与 VAE 的差距表明，对于当前任务，应采用更灵活的自回归先验——

与各向同性的多元高斯分布相比，我们的方法更具优势。关键区别在于：VAE 追求的是平均信息量收敛于固定期望值的表征（即先验分布的微分熵），而我们的解决方案则在其最大似然估计目标中最小化该数值。这种灵活性能够调节潜在表征的丰富度与模型重建能力之间的平衡。相反在 VAE 中，固定先验如同盲目正则化器，可能导致过度平滑的表征；这一特点在补充材料所示的模型采样过程中也得到印证。

图 4 展示了针对公式 9 中损失函数聚合方式的消融研究。该图呈现了三种不同新颖性评分下的 ROC 曲线：i) 对数似然项，ii) 重构项，以及 iii) 同时考虑两者的 proposed 方案。如图所示，在每项数据集中同时考虑记忆性与意外性两方面因素均能带来性能提升。更多佐证请参阅补充材料。

4.2. 视频异常检测

在视频监控场景中，异常通常指异常的人类行为。为此，我们将提出的方法与最先进的异常检测模型进行比较评估。我们采用文献中两个标准基准测试集——UCSD Ped2 [8] 与 ShanghaiTech[30]。尽管这两个数据集在视频数量和分辨率上存在差异，但都包含了监控场景中典型的异常行为（例如机动车驶入人行道、扒窃、斗殴）。针对 UCSD Ped 数据集，我们对 16 帧输入片段进行预处理以提取更小的图像块（详见补充材料），并使用 $\sigma = 0.025$ 的高斯噪声对这些输入进行随机扰动。每个输入片段的新颖性得分通过计算所有图像块的平均新颖值得出。对于 ShanghaiTech 数据集，我们消除了

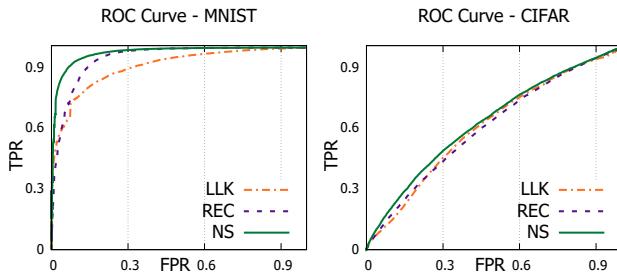


图 4: 不同评分策略在 MNIST 和 CIFAR-10 测试集上呈现的 ROC 曲线。每条曲线均为对十个类别的插值结果。

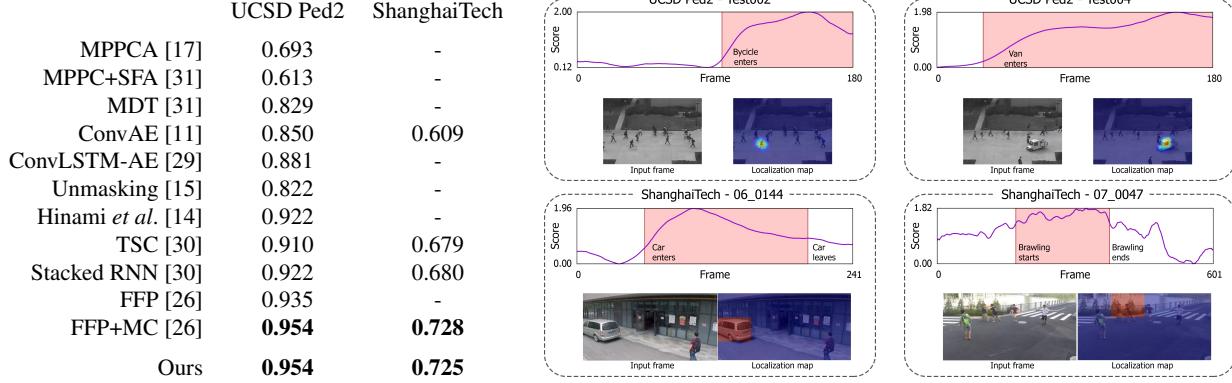


图 5：左侧为本模型相较于顶尖竞争者的 AUROC 性能表现；右侧为从 UCSD Ped2 和 ShanghaiTech 数据集中抽取样本的新颖性评分与定位图谱。每个示例均展示评估得分的走势，并使用不同颜色高亮异常目标进入场景的时间区间。

我们通过使用基于标准 MOG 的方法估计视频片段每一帧的前景并移除背景来构建场景。模型输入为 16 帧片段，但真实异常标注是在帧级别进行的。为恢复每帧的新颖性评分，我们计算该帧所在所有片段的平均得分。随后按照公式 9 所示的相同策略合并损失函数的两项，但遵循异常检测文献的标准方法，以逐序列方式计算归一化系数。各序列的评分最终被拼接起来计算模型的整体 AUROC。此外，我们为两个数据集设计了定位策略：对于 UCSD 数据集，将帧中具有最高新颖性评分的图像块标记为异常；而

在 ShanghaiTech 数据集中，我们采用滑动窗口方法 [47]——当用矩形遮挡块覆盖异常源时，新颖性评分会如预期般显著下降。

图 5 展示了与现有研究的对比结果，同时提供了关于新颖性评分和定位能力的定性评估。尽管采用了更通用的公式化方法，我们的方案在性能上与当前专门为视频应用设计的先进解决方案相当——这些方案利用了光流估计和运动约束。实际上，在不依赖这些假设的情况下（图 5 中 FFP 条目所示），我们的方法在 UCSD Ped2 数据集上超越了未来帧预测的表现。

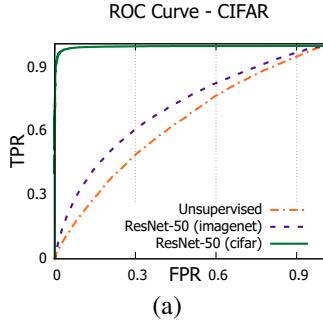
4.3. 模型分析

带有语义特征的 CIFAR-10 数据集。我们研究了在不同假设条件下模型的行为，这些假设涉及新样本的预期性质。我们预期，随着这些假设正确性的提高，novelty 检测性能也会相应提升。这一特性对于事先对新样本存在认知的应用场景尤为重要。

可以设想。为此，我们利用第 4.1 节中描述的 CIFAR-10 基准，并改变输入信息的类型。具体而言，我们不再输入原始图像，而是向模型提供通过 ResNet-50 [12] 提取的语义表示——该网络或在 Imagenet 上预训练（即假设存在语义新颖性），或在 CIFAR-10 本身上训练（即假设存在数据特定新颖性）。两个模型在各自测试集上分别达到了 79.26% 和 95.4% 的 top-1 分类准确率。尽管这种方法在新颖性检测中应被视为不公平，但它可作为验证基准，展示我们的模型在应用于更优质特征时能达到的性能上限。为处理密集输入，我们在估计网络中采用了全连接自动编码器和 MFC 层。

图 6-(a) 展示了最终的 ROC 曲线，其中语义描述符相较于原始图像输入（“无监督”条目）提升了 AUROC 指标。这些结果表明，我们的模型有效利用了正常与异常输入表征之间的分离性并进行相应扩展，甚至在该任务中达到接近最优的性能。但值得注意的是，不同监督程度会带来显著差异的性能表现。正如预期，数据集特定监督将 AUROC 从 0.64 提升至 0.99（满分）。令人惊讶的是，基于 Imagenet（包含所有 CIFAR 类别）训练的语义特征向量带来的提升较为有限，仅产生 0.72 的 AUROC。这一现象表明，即使在某些罕见情况下可以预知新颖性的语义信息，其对正常性建模的贡献仍然有限，这主要是因为新颖性可能依赖于其他线索（例如低层统计特征）。

通过循环层进行自回归。为衡量第 3 节所述 MFC 与 MSC 层的贡献，我们在 CIFAR-10 和 UCSD 数据集上进行测试



CIFAR-10	
LSTM _[100]	0.623
LSTM _[32,32,32,32,100]	0.622
MFC _[100]	0.625
MFC _[32,32,32,32,100]	0.641
UCSD Ped2	
LSTM _[100]	0.849
LSTM _[4,4,4,4,100]	0.845
MSC _[100]	0.849
MSC _[4,4,4,4,100]	0.954

(b)

图 6: (a) 使用语义输入向量的 CIFAR-10 ROC 曲线。每条曲线均为十个类别间的插值结果。(b) 特征空间中自回归密度估计的不同架构对比。LSTM_[F₁, F₂, ..., F_N] 表示估计器各 N 层的输出形状 (MFC 与 MSC 标注方式相同)。实验结果以测试 AUROC 值呈现。

Ped2，自回归密度估计器的替代解决方案。具体而言，我们研究了循环网络，因为它们代表了最具自回归特性的自然替代方案。我们将提出的构建模块与由 LSTM 层组成的估计器进行基准测试，该估计器被设计为顺序观察潜在符号 $\mathbf{z}_{<i}$ 并输出 z_i 的条件概率分布作为最后一层的隐藏状态。我们在单层和多层设置中测试了 MFC、MSC 和 LSTM，并在图 6-(b) 中报告了所有结果。

结果表明，尽管我们的解决方案在浅层设置中表现与循环基线相似，但在连续堆叠多层时能显著发挥其深度优势。MFC 和 MSC 确实为每个输出 CPD 采用解缠结的参数化方法，这一特性等同于为每个 z_i 配置专用估计器网络，从而提升对其指定 CPD 密度建模的专业能力。相比之下，LSTM 网络虽将所有历史信息（即已观测符号）嵌入记忆单元，但通过相同的权重矩阵处理序列中的每个输入。在这种机制下，循环模块需要学习符号间共享的参数，从而丧失 specialization 并削弱其建模能力。

4.4. 认知时间过程中的新颖性

作为我们方案的一个潜在应用，我们研究了其在模拟人类注意力行为方面的能力。为此，我们采用 DR(eye)VE 数据集 [33]，——该数据集专为驾驶场景中的注意力焦点预测而开发。该数据集包含 74 段驾驶视频，其中逐帧提供了注视点分布图，用以标示驾驶员关注的场景区域。为了捕捉注意力模式的动态特性，我们特意舍弃了视觉内容中的

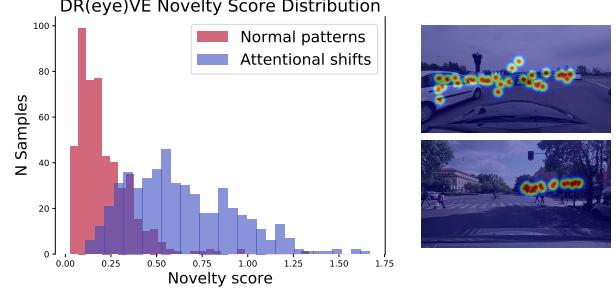


图 7: 左图展示了在 DR(eye)VE 数据集中标注的注意力转移情况下，正常模式所获得的新颖性分数分布。右图呈现了 DR(eye)VE 数据集中产生最高新颖性分数的视频片段（即注意力模式偏离预期行为的片段）。值得注意的是，这些片段描绘了某些特殊场景，例如等待交通灯或接近环岛时的情形。

场景并基于从训练集中随机提取的注视点地图片段优化我们的模型。训练完成后，我们依据每个片段的新颖度评分作为注意力模式非普遍性的代理指标。此外，由于数据集包含特殊及罕见模式（如分心、录制错误）的标注，我们可以衡量所捕获新颖度与这些标注之间的相关性。在 AUROC 指标上，我们的模型达到 0.926 分，表明新颖度可能源于驾驶员的意外行为，例如分心或其他注意力转移。图 7 展示了普通事件与特殊事件在新颖度评分上的差异分布。

5. 结论

我们提出了一个用于新颖性检测的综合框架。我们将模型形式化以捕捉新颖性的双重本质：其既涉及对未见数据记忆的不可行性，也涉及观察其潜在表征所引发的惊奇感。从技术角度来看，这两个方面均通过深度生成自编码器进行建模，并辅以一个额外的自回归密度估计器，该估计器通过最大似然原则学习潜在向量的分布。为此，我们引入了两种适用于图像和视频数据的掩码层。我们证明，在潜在空间中引入此类辅助模块可导致编码器微分熵的最小化，这被证明是适用于当前任务的正则化器。实验结果显示，在单类别和异常检测设置中均实现了最先进的性能，从而证明了我们的框架在不同任务中的灵活性，且无需做出任何与数据相关的假设。

致谢。我们衷心感谢 Facebook 人工智能研究院和松下硅谷实验室为本研究捐赠 GPU 设备。

参考文献

- [1] A. Adam, E. Rivlin, I. Shimshoni 和 D. Reinitz。使用多固定位置监控器的鲁棒实时异常事件检测。《IEEE 模式分析与机器智能汇刊》，30(3):555–560，2008 年。2
- [2] S. Bai, J. Z. Kolter 和 V. Koltun。通用卷积与循环网络在序列建模中的实证评估。《arXiv:1803.01271》，2018 年。4
- [3] J. Ball'e, V. Laparra 和 E. P. Simoncelli。端到端优化的图像压缩。《国际学习表征大会》，2017 年。1[4] A. Barto, M. Mirolli 和 G. Baldassarre。新颖性或惊奇度？《心理学前沿》，4:907，2013 年。1[5] A. Bharat, A. Gritai 和 M. Shah。通过学习物体运动模式实现异常检测与改进的物体检测。见《IEEE 国际计算机视觉与模式识别会议》，第 1–8 页。IEEE，2008 年。2[6] M. Bauer 和 A. Mnih。变分自编码器的重采样先验。《国际人工智能与统计大会》，2019 年。2[7] S. Calderara, U. Heinemann, A. Prati, R. Cucchiara 和 N. Tishby。使用谱图分析检测人员轨迹中的异常。《计算机视觉与图像理解》，115(8):1099–1111，2011 年。1[8] A. Chan 和 N. Vasconcelos。UCSD 行人数据库。《IEEE 模式分析与机器智能汇刊》，2008 年。6[9] Y. Cong, J. Yuan 和 J. Liu。基于稀疏重构成本的异常事件检测。见《IEEE 国际计算机视觉与模式识别会议》，第 3449–3456 页。IEEE，2011 年。1, 2[10] M. Germain, K. Gregor, I. Murray 和 H. Larochelle。MADE：用于分布估计的掩码自编码器。见《国际机器学习大会》，第 881–889 页，2015 年。3[11] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury 和 L. S. Davis。学习视频序列中的时间规律性。见《IEEE 国际计算机视觉与模式识别会议》，第 733–742 页。IEEE，2016 年。1, 2, 7[12] K. He, X. Zhang, S. Ren 和 J. Sun。用于图像识别的深度残差学习。见《IEEE 国际计算机视觉与模式识别会议》，第 770–778 页，2016 年。7[13] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler 和 S. Hochreiter。通过双时间尺度更新规则训练的 GAN 收敛至局部纳什均衡。见《神经信息处理系统》，第 6626–6637 页，2017 年。12[14] R. Hinami, T. Mei 和 S. Satoh。通过学习深度通用知识实现异常事件的联合检测与叙述。见《IEEE 国际计算机视觉大会》，第 3639–3647 页，2017 年。1, 2, 7[15] R. T. Ionescu, S. Smeureanu, B. Alexe 和 M. Popescu。视频中异常事件的解蔽。《IEEE 国际计算机视觉大会》，2017 年。7[16] L. Itti 和 P. Baldi。贝叶斯惊奇度吸引人类注意力。《视觉研究》，49(10):1295–1306，2009 年。1
- [17] J. Kim 与 K. Grauman。《局部观察，全局推断：支持增量更新的异常活动检测时空马尔可夫随机场》。发表于 IEEE 国际计算机视觉与模式识别会议，第 2921–2928 页。IEEE 出版社，2009 年。2, 7[18] D. P. Kingma 与 J. Ba。《Adam：一种随机优化方法》。发表于国际学习表征会议，2015 年。11[19] D. P. Kingma 与 M. Welling。《变分贝叶斯自编码器》。发表于国际学习表征会议，2014 年。2, 5, 11[20] T. Kohonen。《自组织与关联记忆》第 8 卷。Springer 科学与商业媒体，2012 年。1[21] D. Koller 与 N. Friedman。《概率图模型：原理与技术 - 自适应计算与机器学习》。MIT 出版社，2009 年。13[22] A. Kumar。《基于计算机视觉的织物缺陷检测综述》。载于 IEEE 工业电子学汇刊，第 55 卷第 1 期，第 348–363 页，2008 年。1
- [23] J. Kwon 与 K. M. Lee。《多视角事件摘要与稀有事件检测的统一框架》。载于 IEEE 模式分析与机器智能汇刊，第 37 卷第 9 期，第 1737–1750 页，2015 年。2[24] H. Larochelle 与 I. Murray。《神经自回归分布估计器》。发表于第十四届人工智能与统计国际会议论文集，第 29–37 页，2011 年。3[25] W. Li、V. Mahadevan 与 N. Vasconcelos。《拥挤场景中的异常检测与定位》。载于 IEEE 模式分析与机器智能汇刊，第 36 卷第 1 期，第 18–32 页，2014 年。2[26] W. Liu、W. Luo、D. Lian 与 S. Gao。《基于未来帧预测的异常检测——新基线》。发表于 IEEE 国际计算机视觉与模式识别会议，2018 年。1, 2, 7
- [27] C. Lu、J. Shi 与 J. Jia。《在 MATLAB 中以 150 帧 / 秒检测异常事件》。发表于 IEEE 国际计算机视觉会议，第 2720–2727 页。IEEE 出版社，2013 年。2[28] P. Luc、N. Neverova、C. Couprie、J. Verbeek 与 Y. Le-Cun。《语义分割的未来深度预测》。发表于 IEEE 国际计算机视觉会议，第 648–657 页，2017 年。3[29] W. Luo、W. Liu 与 S. Gao。《利用卷积 LSTM 记忆历史实现异常检测》。发表于 IEEE 国际多媒体与博览会会议，第 439–444 页。IEEE 出版社，2017 年。7[30] W. Luo、W. Liu 与 S. Gao。《基于稀疏编码的堆叠 RNN 异常检测框架再探索》。发表于 IEEE 国际计算机视觉会议，2017 年。1, 2, 6, 7[31] V. Mahadevan、W. Li、V. Bhalodia 与 N. Vasconcelos。《拥挤场景中的异常检测》。发表于 IEEE 国际计算机视觉与模式识别会议，第 1975–1981 页。IEEE 出版社，2010 年。2, 7[32] A. v. d. Oord、S. Dieleman、H. Zen、K. Simonyan、O. Vinyals、A. Graves、N. Kalchbrenner、A. Senior 与 K. Kavukcuoglu。《WaveNet：原始音频的生成模型》。arXiv 预印本 arXiv:1609.03499，2016 年。3

- [33] A. Palazzi, D. Abati, S. Calderara, F. Solera 与 R. Cucchiara。预测驾驶员注意力焦点: dr(eye)ve 项目。IEEE 模式分析与机器智能汇刊, 2018 年。8[34] D. Pathak, P. Agrawal, A. A. Efros 与 T. Darrell。通过自监督预测实现好奇心驱动探索。发表于国际机器学习大会, 2017 年第 1 卷。1[35] M. Ravankhah, E. Sangineto, M. Nabi 与 N. Sebe。训练对抗判别器用于跨通道人群异常事件检测。arXiv 预印本 arXiv:1706.07680, 2017 年。2[36] M. Sabokrou, M. Fayyaz, M. Fathy, Z. Moayed 与 R. Klette。Deep-anomaly: 面向拥挤场景快速异常检测的全卷积神经网络。计算机视觉与图像理解, 2018 年。2[37] M. Sabokrou, M. Khalooei, M. Fathy 与 E. Adeli。基于对抗学习的单类分类器用于新颖性检测。发表于 IEEE 国际计算机视觉与模式识别会议, 第 3379-3388 页, 2018 年。1, 2[38] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth 与 G. Langs。利用生成对抗网络进行无监督异常检测以指导标记发现。发表于医学影像信息处理国际会议, 第 146-157 页。Springer, 2017 年。1, 2, 5[39] B. Scholkopf, R. C. Williamson, A. J. Smola, J. Shawe-Taylor 与 J. C. Platt。新颖性检测的支持向量方法。发表于神经信息处理系统大会, 2000 年。5[40] L. Theis, A. v. d. Oord 与 M. Bethge。关于生成模型评估的说明。国际学习表征大会, 2016 年。5[41] J. M. Tomczak 与 M. Welling。采用 VAMP 先验的变分自编码器。国际人工智能与统计会议, 2018 年。2[42] M. Tribus。热静力学与热力学: 能量、信息与物态导论及工程应用。van Nostrand 出版社, CS7 系列, 1961 年。1, 2[43] B. Uria, I. Murray 与 H. Larochelle。RNNADE: 实值神经自回归密度估计器。发表于神经信息处理系统进展, 第 2175-2183 页, 2013 年。3[44] B. Uria, I. Murray 与 H. Larochelle。一种深度可处理的密度估计器。发表于国际机器学习大会, 第 467-475 页, 2014 年。3[45] A. van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals 与 A. Graves。基于 PixelCNN 解码器的条件图像生成。发表于神经信息处理系统大会, 2016 年。3, 5[46] A. van den Oord, N. Kalchbrenner 与 K. Kavukcuoglu。像素循环神经网络。国际机器学习大会, 2016 年。3, 11[47] M. D. Zeiler 与 R. Fergus。卷积网络的可视化与理解。发表于欧洲计算机视觉大会, 第 818-833 页。Springer, 2014 年。7[48] B. Zhao, L. Fei-Fei 与 E. P. Xing。通过动态稀疏编码实现视频异常事件在线检测。发表于 IEEE 国际计算机视觉与模式识别会议, 第 3313-3320 页。IEEE, 2011 年。2[49] B. Zong、Q. Song、M. R. Min、W. Cheng、C. Lumezanu、D. Cho 和 H. Chen。基于深度自编码高斯混合模型的无监督异常检测。见国际学习表征大会, 2018 年。1, 2

补充材料

6. 关于实现细节

各实验采用的架构与超参数详见表 2，包括模块类型、自回归层数、小批量规模、学习率以及对数似然目标的权重。所有中间层均采用 Leaky ReLU 激活函数。使用 Adam [18] 优化器对目标函数进行优化，所有超参数通过在保留验证集上最小化原始目标函数（公式 4 结合 $\lambda = 1$ ）进行调优。

7. 关于对数似然目标函数

在本节中，我们将详细说明如何计算和优化对数似然项（主论文中的公式 4）。重要的是，如主论文所述，我们通过多项分布对每个条件概率分布进行建模。为此，我们首先需要

	MNIST	CIFAR-10	UCSD Ped2	ShanghaiTech	DR(eye)VE
Input Shape	1,28,28	3,32,32	1,8,32,32*	3,16,256,512	1,16,160,256
Encoder Network	2D Conv _{3x3} ³²	D _{1,2,2} ⁸	D _{1,2,2} ⁸	D _{1,2,2} ⁸	D _{1,2,2} ⁸
	R ³²	D _{2,1,1} ⁸	D _{1,2,2} ⁸	D _{2,2,2} ⁸	D _{2,2,2} ⁸
	D _{2,2} ³²	D _{2,2,2} ⁸	D _{2,2,2} ⁸	D _{2,2,2} ⁸	D _{2,2,2} ⁸
	D _{2,2} ³²	D _{1,2,2} ⁸	D _{1,2,2} ⁸	D _{1,2,2} ⁸	D _{1,2,2} ⁸
	FC ⁶⁴	D _{2,1,1} ⁸	D _{1,2,2} ⁸	D _{1,2,2} ⁸	D _{1,2,2} ⁸
	FC ⁶⁴	D _{1,2,2} ⁸	TFC ₃₁₂	TFC ₃₁₂	TFC ₆₄
Decoder Network	FC ²⁵⁶	TFC ⁶⁴	TFC ⁶⁴	TFC ⁶⁴	TFC ⁶⁴
	FC ⁶⁴	TFC ⁶⁴	TFC ₅₁₂	TFC ₅₁₂	TFC ₅₁₂
	FC ⁶⁴	U _{2,2} ¹²⁸	U _{2,2,2} ⁶⁴	U _{2,2,2} ⁶⁴	U _{2,2,2} ⁶⁴
	U _{2,2} ³²	U _{2,2} ⁶⁴	U _{2,2,2} ³²	U _{2,2,2} ³²	U _{2,2,2} ³²
	U _{2,2} ¹⁶	U _{2,2} ³²	U _{2,2,2} ¹⁶	U _{2,2,2} ¹⁶	U _{2,2,2} ¹⁶
	2D Conv _{1x1} ¹	R ³²	U _{2,1,1} ⁸	U _{2,2,2} ⁸	U _{1,2,2} ⁸
Estimator Network	2D Conv _{1x1} ³	2D Conv _{1x1} ³	3D Conv _{1x1} ³	3D Conv _{1x1} ³	3D Conv _{1x1} ³
	MFC ³²	MFC ³²	MSC ⁴	MSC ⁴	MSC ⁴
	MFC ³²	MFC ³²	MSC ⁴	MSC ⁴	MSC ⁴
	MFC ³²	MFC ³²	MSC ⁴	MSC ⁴	MSC ⁴
	MFC ³²	MFC ³²	MSC ⁴	MSC ¹⁰⁰	MSC ⁴
	MFC ¹⁰⁰	MFC ¹⁰⁰	MSC ¹⁰⁰	MSC ¹⁰⁰	MSC ¹⁰⁰
Mini Batch	256	256	2760	8	16
Learning Rate	10^{-4}	10^{-3}	10^{-3}	10^{-3}	10^{-3}
λ	1	0.1	0.1	1	1

*Patches extracted from input clips having shape 1,16,256,384.

表 2：各配置的架构与优化超参数。我们以 D_S^C （下采样）、 U_S^C （上采样）和 R^C （残差）表示所采用构建模块的参数化配置（详见主论文图 1ii）。其中 C 表示输出通道数， S 表示模块中首个卷积层的步长。全连接层 FC^C 与时序共享全连接层 TFC^C 分别表示稠密层和时序共享的全连接层（此时 C 指输出特征数）。最后， MFC^C 和 MSC^C 对应本文提出的自回归层（详见稿件图 3）。关于各类层的完整说明，请参阅主论文第 3.1 节。

编码器充当有界函数。为实现这一目标，我们简单地采用 S 形激活函数，确保潜在表示 $\mathbf{z} = f(\mathbf{x}; \theta_f)$ 位于 $[0, 1]^d$ 范围内。因此，对于每个 z_j （其中 $j = 1, 2, \dots, d$ ），我们对空间 $[0, 1]$ 执行线性量化，将其划分为 B 个区间（其中 B 为超参数）。此步骤为 z_j 生成一个 B 维分类分布 $\phi(z_j)$ ，标示出 z_j 所属的正确区间。对于每个条件概率分布，该分布将作为估计器 $h(\mathbf{z}; \theta_h)$ 的真实标签，该估计器通过 softmax 激活函数一致地预测 d 个分布 $p(z_j | \mathbf{z}_{<j})$ ，覆盖全部 B 个区间。如此，如公式 11 所示， \mathcal{L}_{LLK} 损失转化为有效的似然项，定义为每个估计条件概率分布与其对应分类分布之间的交叉熵损失：

$$\mathcal{L}_{LLK}(\theta_f, \theta_h) = \mathbb{E}_{\mathbf{x} \sim P} \left[- \sum_{j=1}^d \sum_{k=1}^B \phi(z_j)_k \log(p(z_j | \mathbf{z}_{<j}))_k \right]. \quad (11)$$

值得注意的是，多项式只是 CPD 的合理模型之一。实际上，若将其替换为高斯模型，整体框架仍能成立。但正如我们在不同试验中观察到的，这种选择并未带来显著改进，反而会引发数值不稳定问题——正如先前研究 [46] 所述。

8. 关于与变分自编码器的关系

我们的框架与变分自编码器（VAE）[19] 存在一些相似之处。事实上，它们都通过最小化重构误差来近似主论文中公式 1 的积分，同时受到涉及隐向量先验分布的正则化约束。然而，有几个根本区别值得注意。首先，我们的模型没有提供从后验分布采样的显式策略，因此形成了从输入到隐层表示的确定性映射。其次，VAE 为建模先验 $p(\mathbf{z})$ 指定了显式且固定的形式，而在我们的公式中，其分布形态不受任何假设约束，并可直接通过估计器的自回归特性进行学习。这一点使我们的方案具有两个优势：第一，由于 VAE 强制编码分布与先验匹配，其微分熵会收敛至与先验相同。这种行为导致不同设置下的熵值近似恒定（可参见主论文图 2，我们在其中讨论了新颖性检测任务中熵最小化的直观意义）；第二，使用过于简单的先验可能导致过度正则化的表示，而我们的方案不易出现此类风险。相关实证证据

FID	VAE Samples	Our Samples	FID
149.72			72.96
172.02			72.53
181.56			76.27
188.37			67.33
202.06			68.33
207.47			73.92
186.48			62.26
220.79			64.38
164.36			52.53
204.84			67.17

图 8：针对所有 CIFAR-10 类别（按不同行排列），我们展示了从 VAE（左侧）和采用自回归先验的提议自编码器中采样的图像。可见，我们的样本在视觉上呈现出细腻的细节和清晰度，这与 VAE 产生的严重模糊图像形成鲜明对比。最后，通过观察 FID 分数（位于图表两端，数值越低越好）可以证实 VAE 存在的过度正则化问题。

此类行为亦可在图 8 中得到印证——我们分别从 VAE 与本文模型中抽取新样本，二者均在 CIFAR-10 数据集上完成训练。在同等条件下，本模型生成的幻象视觉上远比 VAE 产生的更为真实，后者会导致过度平滑的形态且缺乏细节特征，这一结论进一步得到了两者在 Fr{v1} echet Inception Distance (FID) 分数 {v2} 上的显著差异所证实。

9. 论新颖性的双重本质

在本节中，我们重点强调了同时保留两项要素对于获得高区分度新颖性评分 (NS, 主论文公式 9) 的重要性：即建模记忆能力的重构误差 (REC)，以及从潜在表示中提取意外信息的对数似然项 (LLK)。为强化这一在稿件图 4 中仅简要说明的观点，我们在表 3 中汇报了主论文所述各场景下不同评分策略的表现（以 AUROC 表示）。除 ShanghaiTech 数据集外，我们系统性地观察到兼顾两项要素的评估方式能持续带来性能提升——

此外，对于 MNIST 和 CIFAR-10 数据集，我们发现由我们的重构误差产生的性能差距与去噪自编码器 (DAE) 变体产生的性能差距尤为引人关注（根据主论文表 1 报告，两个数据集上的数值分别为 0.942 和 0.590）。在这方面，我们收集到的新证据表明，惊奇最小化起到了面向新颖性的 {v*} 作用。

	LLK	REC	NS
MNIST	0.926	0.949	0.975
CIFAR-10	0.627	0.603	0.641
UCSD Ped2	0.933	0.909	0.954
ShanghaiTech	0.695	0.726	0.725
DR(eye)VE	0.917	0.863	0.926

表 3：针对每种设置，三种不同新颖性评分下的 AUROC 性能表现：i) 对数似然项 (LLK)，ii) 重构项 (REC)，以及 iii) 同时考虑二者的 proposed 方案 (NS)。

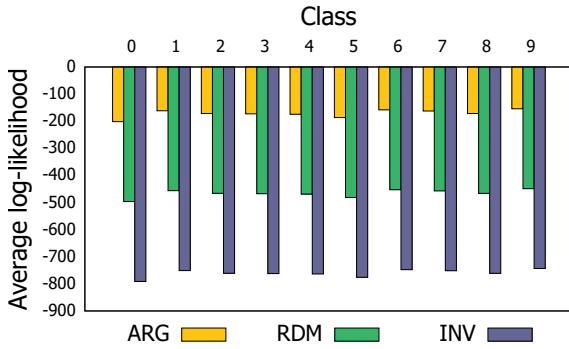


图 9：贝叶斯网络建模我们模型在 MNIST 数字数据集上训练时编码器所生成隐码分布的样本训练对数似然。当 BN 结构接近训练时施加的自回归顺序，可获得显著更高的似然值。该现象在所有类别中保持一致，证实了编码器能够生成遵循预设自回归结构的编码。

整体架构的正则化器，因为它提升了重建的判别能力（正如主论文第 4.1 节中已推测的那样）。

10. 关于表征的因果结构

我们现在研究编码器生成表征的能力，这些表征遵循主论文第 3 节中提到的 LLK 损失所施加的自回归因果结构。为此，我们从在 MNIST 数字上训练的十个模型中提取表征，并使用结构化密度估计器拟合它们的分布。具体来说，我们采用具有不同自回归结构的贝叶斯网络（BNs）。在这方面，每个 BN 都使用线性高斯模型 [21]，进行建模，使得每个条件概率分布 $p(z_i|Pa(z_i))$ 在 $i = 1, 2, \dots, d$ 条件下由以下公式给出：

$$p(z_i|Pa(z_i)) = \mathcal{N}(z_i | w_0^{(i)} + \sum_{z_j \in Pa(z_i)} w_j^{(i)} z_j, \sigma_i^2), \quad (12)$$

其中每个 $w_j^{(i)}$ 、 σ_i^2 都是可学习参数。我们用 $Pa(z_i)$ 表示贝叶斯网络中 z_i 的父变量。上述等式对所有节点都成立，除根节点外——根节点通过高斯分布进行建模。关于贝叶斯网络结构，我们测试了以下方案：

- 自回归顺序：BN 结构遵循训练期间施加的自回归顺序，即 $Pa(z_i) = \{z_j | j = 1, 2, \dots, i-1\}$
- 随机顺序：BN 结构遵循随机的自回归顺序。
- 逆序：BN 结构遵循自回归顺序，该顺序与

在训练过程中施加的 g 即 y
 $Pa(z_i) = \{z_j | j = i+1, i+2, \dots, d\}$

值得注意的是，由于这三种结构具有相同数量的边和独立参数，它们拟合能力的差异仅源于变量间施加的因果顺序。

图 9 展示了所有 BN 模型的训练对数似然样本。值得注意的是，自回归排序实现了更好的拟合效果，这印证了编码器网络能够提取具有学习自回归特性的特征。此外，为证明该结果并非由过拟合或其他潜在行为导致，我们在表 4 中同步呈现了训练集、验证集与测试集的对数似然值。

11. 关于熵最小化

为了进一步理解表示熵最小化的作用，我们聚焦于单个 MNIST 数字（类别 7），并在图 10 中展示了从训练集中随机抽取的部分重建样本。这些重建结果是在三种不同正则化机制下学习得到的，通过对数似然目标函数 (λ , 主论文公式 4) 的不同权重来体现。如图 10 所示，更高强度的正则化（即对熵的更严格约束）会生成近乎模式坍塌的重建结果，牺牲了鲜明特征变化以捕捉输入分布中更少的原型。

12. 论自回归层的复杂性

在本节中，我们简要讨论掩码全连接（MFC）和掩码堆叠卷积的复杂度。

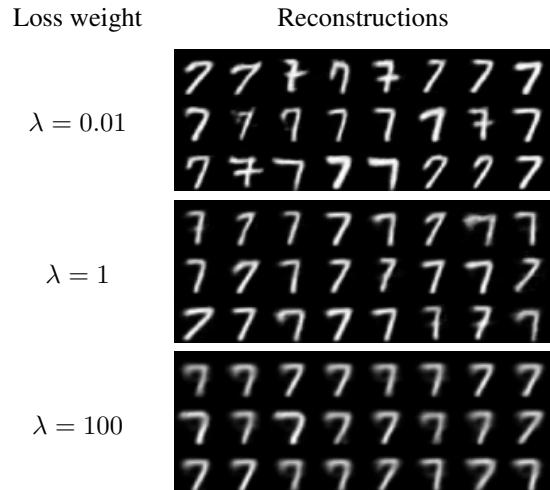


图 10：由不同 λ 值提供的 MNIST 重建结果，后者控制着微分熵最小化的影响程度。

		Classes									
		0	1	2	3	4	5	6	7	8	9
ARG	Train	-201.60	-161.60	-171.43	-172.73	-174.17	-186.48	-158.22	-162.37	-171.65	-154.11
	Val	-200.96	-160.38	-170.10	-172.29	-173.85	-185.25	-157.22	-162.20	-171.42	-154.02
	Test	-200.89	-159.73	-169.64	-170.75	-172.40	-184.27	-157.74	-161.65	-170.10	-152.70
RDM	Train	-496.33	-456.34	-466.16	-467.47	-468.90	-481.21	-452.95	-457.10	-466.39	-448.84
	Val	-495.69	-455.11	-464.83	-467.02	-468.58	-479.98	-451.95	-456.93	-466.15	-448.75
	Test	-495.62	-454.47	-464.37	-465.48	-467.13	-479.00	-452.48	-456.38	-464.83	-447.43
INV	Train	-791.06	-751.07	-760.89	-762.20	-763.63	-775.94	-747.68	-751.83	-761.12	-743.57
	Val	-790.42	-749.84	-759.56	-761.75	-763.31	-774.71	-746.68	-751.66	-760.88	-743.48
	Test	-790.35	-749.20	-759.11	-760.22	-761.86	-773.73	-747.21	-751.12	-759.56	-742.16

表 4: 不同 BN 结构在拟合 MNIST 表征时获得的样本对数似然。每个 BN 均在单类别训练集的潜码上训练，遵循自回归序 (ARG)、随机序 (RDM) 或逆自回归序 (INV)。我们同时报告了验证集和测试集上的对数似然。关于训练 - 验证 - 测试集划分方法，请参阅论文第 4.1 节。本评估仅使用 “正常” 测试样本。

(MSC) 层 (主论文图 3)²: 遵循主论文第 3 节引入的符号表示，MFC 具有 $d^2 + d \cdot ci \cdot co + d \cdot co$ 个可训练参数，计算复杂度为 $\mathcal{O}(d^2 \cdot ci \cdot co)$ 。而 MSC 则具有 $3d^2 + d \cdot ci \cdot co + d \cdot co$ 个自由参数，时间复杂度为 $\mathcal{O}(d^2 \cdot ci \cdot co \cdot t)$ 。

13. 视频异常检测中的定位与新颖性评分

我们在图 11 中展示了模型在视频异常检测场景下行为的其他定性证据，即 UCSD Ped2 和上海科技大学数据集。

²我们将这两层的类型称为 “B”，因为它是类型 “A”的上界

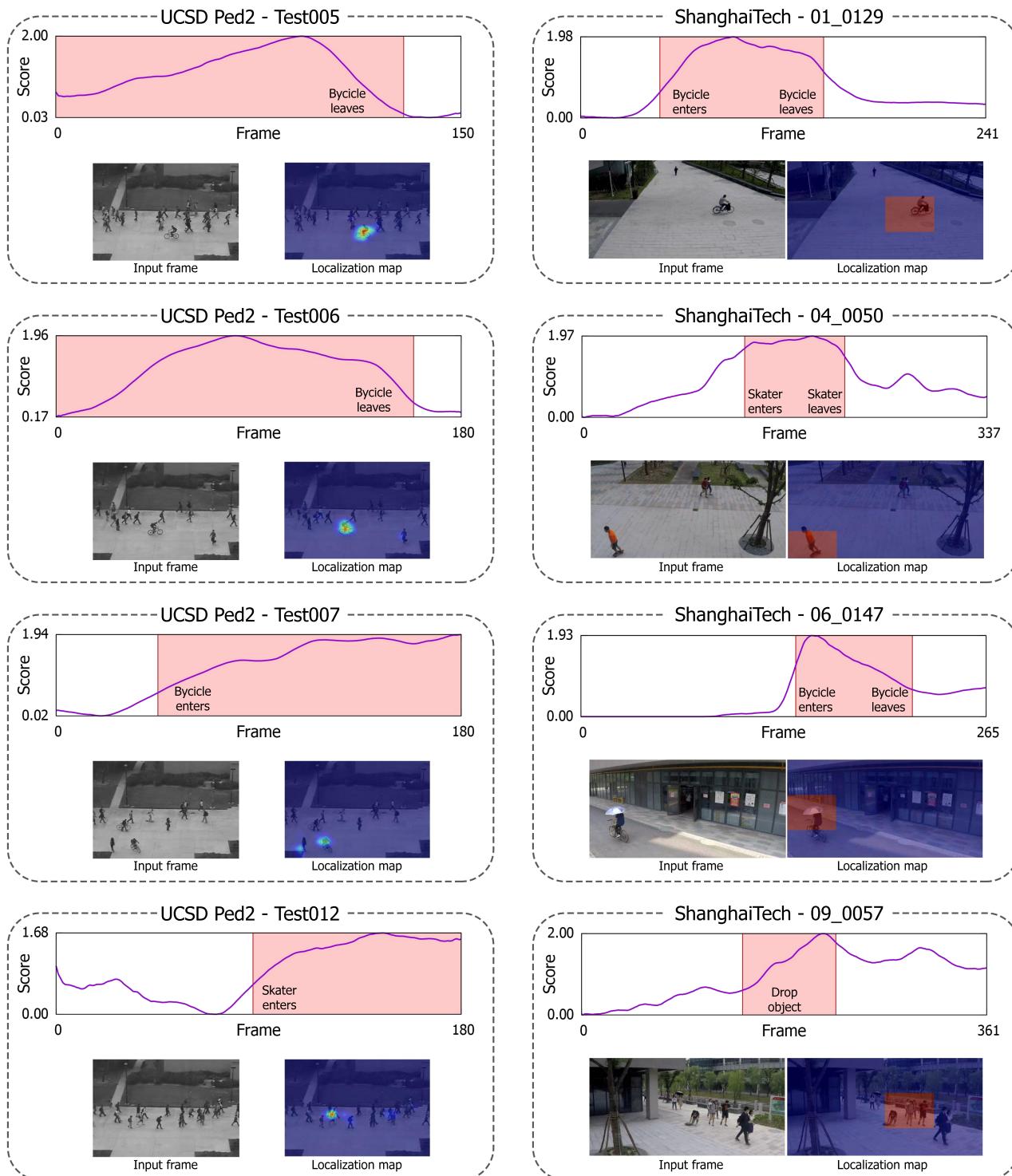


图 11：来自 UCSD Ped2（左图）和 ShanghaiTech（右图）若干测试片段的新颖性评分及定位图谱。