



基于神经遗忘决策集成网络 (NODE) 与数值特征嵌入的极端降雨事件预测：多变量气象表格数据分析

赵亮¹, 唐琳², 杨坤³, 张彬蕾^{1*}

(1. 喀什大学电子与通信工程学院, 新疆 喀什 844000; 2. 重庆文理学院数学与人工智能学院, 重庆 402160; 3. 重庆工业职业技术大学机械设计与制造学院, 重庆 401120)

摘要: 极端降雨事件的频发对区域水资源管理与城市防灾减灾构成严峻威胁, 准确的极端气象预测已成为大气科学领域的核心诉求。传统数值天气预报 (NWP) 系统受限于次网格尺度参数化瓶颈, 难以精准捕捉局地极端降水的突发性; 而以梯度提升决策树 (GBDT) 为代表的浅层机器学习模型在处理高维异构气象表格数据时, 存在固有的特征表征瓶颈。为此, 本研究提出了一种引入分段线性编码 (PLE) 数值特征嵌入的神经遗忘决策集成网络 (NODE-PLE)。该架构将一维连续气象物理量映射至高维稠密特征空间, 并利用基于 激活函数的可微遗忘决策树实现端到端的非线性特征交互。实证研究基于覆盖澳大利亚 49 个气象站、跨度 10 年的 142193 个样本数据集, 采用深度去噪自编码器处理高缺失率变量, 构建物理驱动的高阶衍生特征, 并以焦点损失函数 (Focal Loss) 应对极端事件 (降雨量 > 13.00 mm, 占比 4.97%) 的类别不平衡问题。独立测试集验证表明, NODE-PLE 在综合性能上取得了与传统 GBDT 基线模型竞争的表现, RMSE 降至 8.58 mm, 较 XGBoost 下降 6.8%, 展现了深度学习架构结合数值特征嵌入在气象表格数据极端事件预测中的潜力与竞争力。

关键词: 极端降雨预测; 神经遗忘决策集成; 数值特征嵌入; 分段线性编码; 焦点损失函数; 表格数据深度学习

收稿时间: 2026年3月6日

中图分类号: TP181

通讯作者: 张彬蕾, 重庆工业职业技术大学机械设计与制造学院

Extreme Rainfall Event Prediction Based on Neural Forgetting Decision Ensemble Network (NODE) and Numerical Feature Embedding: Multivariate Meteorological Table Data Analysis

Zhao Liang¹, Tang Lin², Yang Kun³, Zhang Binlei^{1*}

(1. College of Electronic and Communication Engineering, Kashgar University, Kashgar, Xinjiang 844000; 2 School of Mathematics and Artificial Intelligence, Chongqing University of Arts and Sciences, Chongqing 402160; 3. School of Mechanical Design and Manufacturing, Chongqing Industrial Vocational and Technical University, Chongqing 401120)

Abstract: The frequent occurrence of extreme rainfall events poses a serious threat to regional water resource management and urban disaster prevention and reduction. Accurate extreme weather prediction has become a core demand in the field of atmospheric science. Traditional numerical weather forecasting (NWP) systems are limited by the bottleneck of sub grid scale parameterization, making it difficult to accurately capture the suddenness of



local extreme precipitation; However, shallow machine learning models represented by gradient boosting decision trees (GBDT) have inherent feature representation bottlenecks when processing high-dimensional heterogeneous meteorological table data. To this end, this study proposes a neural forgetting decision ensemble network (NODE-PLE) incorporating piecewise linear coding (PLE) numerical feature embedding. This architecture maps one-dimensional continuous meteorological physical quantities to a high-dimensional dense feature space, and utilizes a differentiable forgetting decision tree based on the α -entmax activation function to achieve end-to-end nonlinear feature interaction. The empirical study is based on a dataset of 142193 samples covering 49 weather stations in Australia and spanning 10 years. A deep denoising autoencoder is used to process high missing rate variables, and high-order derived features driven by physics are constructed. The Focal Loss function is used to address the class imbalance problem of extreme events (rainfall > 13.00 mm, accounting for 4.97%). Independent test set validation shows that NODE-PLE outperforms traditional GBDT baseline models in overall performance, with RMSE reduced to 8.58 mm, a 6.8% decrease compared to XGBoost, demonstrating the potential and competitiveness of deep learning architecture combined with numerical feature embedding in extreme event prediction of meteorological table data.

Keywords: Extreme rainfall prediction; Neural forgetting decision ensemble; Numerical feature embedding; Segmented linear coding; Focus loss function; Table Data Deep Learning

0 引言 (Introduction)

极端天气事件的精准预测是气象学与计算机科学交叉领域中最具挑战性的前沿课题之一。根据克劳修斯-克拉珀龙方程 (Clausius-Clapeyron Relation), 大气温度升高将指数级增加水汽容纳能力, 导致极端重降水事件在全球变暖背景下呈非线性增长趋势^[1]。准确提前预测极端降雨, 对于防洪减灾、航空调度及可再生能源并网具有重要的科学价值与现实意义^[2]。

长久以来, 气象预测高度依赖数值天气预报 (NWP) 模型, 其通过离散化求解三维流体偏微分方程组模拟未来天气状态。然而, 受限于计算资源, 高分辨率全球 NWP 模型仍无法显式解析导致极端降水的微物理过程与对流系统, 次网格尺度物理过程的参数化不可避免地引入系统性偏差^[3,4]。

为突破物理驱动模型的局限, 基于历史观测数据的数据驱动方法逐渐成为重要研究方向^[5]。在处理多变量气象表格数据 (Tabular Data) 时, 梯度提升决策树 (GBDT) 模型——如 XGBoost^[6]、CatBoost^[7] 和 LightGBM^[8]——凭借其对于偏态分布与缺失值的鲁棒性, 长期占据主导地位^[9]。然而, 决策树的数学本质是分段常数函数, 这限制了其对连续大气动力学背后平滑非线性流形的拟合能力; 且传统集成树模型无法通过反向传播实现端到端的层次化表征学习^[10]。

与此同时, 标准深度前馈网络 (如多层感知机 MLP) 在表格数据上的表现同样不尽如人意。多项基准测试表明, 标准 DNN 在处理具有重尾分布、高缺失率的气象表格数据时, 泛化性能常劣于经调参的 GBDT^[11, 12]。其核心原因在于全连接层缺乏针对异构数值特征特异性变换的归纳偏置 (Inductive Bias)^[13]。

为打破上述僵局, 本文提出了一种结合数值特征嵌入的深度学习方案。具体而言, 本研究采用并扩展了神经遗忘决策集成网络 (Neural Oblivious Decision Ensembles, NODE), 在网络内部重构可微遗忘决策树, 使其既继承树模型的鲁棒性偏置, 又融入梯度优化范畴^[14]。针对气象变量连续标量表示能力不足的问题, 本文在前端嵌入了分段线性编码 (Piecewise Linear Encoding, PLE) 模块, 实现连续物理量到高维稠密向量的非线性映射。结合针对极端类别不平衡重构的焦点损失函数, 本研究验证了通过高级特征编码与可微树架构的融合, 深度神经网络能够在气象表格数据极端事件分析中与传统树模型形成竞争力, 并在预测校准性上展现优势^[15]。

2 数据集描述与预处理

2.1 数据集概述

本研究整合了澳大利亚气象局发布的 "Rain in Australia" 长期气象站观测数据集与全球逐时地面



综合观测数据集 (ISD), 涵盖 2007 ~ 2017 年间分布于 49 个气象站 (从干旱内陆到湿润沿海) 的 142193 条逐日观测样本。原始数据包含 23 维特征, 记录了气温、相对湿度、大气压强、风速矢量、蒸发量及云量等关键气象物理量。

2.2 缺失值处理: 深度自编码器插补

数据分布稽查表明, 日照时长缺失率高达 48.0%, 蒸发量缺失率为 43.2%, 下午云量缺失率为 40.8%。气象观测中的数据缺失往往并非完全随机 (MCAR), 而是带有系统性偏差的非随机缺失 (MNAR/MAR)^[16], 单变量中位数插补会扭曲变量间的协方差结构。

为此, 本研究采用基于去噪深度自编码器 (DAE) 的非线性多重插补方法。自编码器通过瓶

颈层学习输入的流形表征, 其重构损失函数定义为:

$$\mathcal{L}_{AE} = \frac{1}{N} \sum_{i=1}^N \left\| X_{obs}^{(i)} - D(\mathcal{E}(X_{obs}^{(i)} \cup \hat{X}_{miss}^{(i)})) \right\|_2^2 + \lambda \|W\|_2^2$$

其中 \mathcal{E} 和 D 分别为编码器与解码器映射函数。通过利用相对完整的 20 余维特征作为上下文约束, 自编码器重构出缺失变量在多元物理空间中最具概率合理性的估算值, 最大限度保留了变量间的联合概率分布特性。

2.3 物理驱动的衍生特征与共线性排查

仅依赖单一时刻的静态读数不足以表征极端天气系统的演变过程。本研究依托气象学先验构建了物理驱动的衍生特征, 所有特征均经独立双样本 T 检验验证其与目标变量的统计显著性 ($p < 0.001$)。

表 1

衍生特征	计算公式	物理意义	T 检验
Pressure_Diff	$\Delta P = P_{3pm} - P_{9am}$	气压骤降为锋面逼近或强对流的前兆信号	$t = 28.85$
Temp_Range	$\Delta T = \text{MaxT} - \text{MinT}$	温度日较差反映大气热力学稳定性	$t = -134.77$
Humidity_Diff	$\Delta H = H_{3pm} - H_{9am}$	近地层水汽快速饱和是降水发生的前提	$t = 103.11$

经皮尔逊相关系数矩阵排查, Temp9am 与 MinTemp ($r=0.90$)、Temp3pm 与 MaxTemp ($r=0.98$) 存在严重共线性, 予以剔除, 如表 1 所示。此外, RainToday 与次日降雨标签存在数据泄露风险, 亦从模型输入中移除。

2.4 异常值处理与特征归一化

采用孤立森林 (Isolation Forest) 算法在多维特征空间中检测并剔除传感器异常值^[17]。对近似高斯分布的连续变量实施 Z-score 标准化; 对长尾右偏的降雨量变量实施对数变换 $x' = \log(1+x)$ 以缓解梯度弥散。

2.5 任务重构与时序划分

原始二分类目标 (降雨量 > 0) 在极端灾害预警中缺乏针对性。根据降雨量分布统计 (均值 2.36 mm, 中位数 0.00 mm, 95 分位数 13.00 mm), 本研究将任务升级为极端降雨预测: 次日降雨量 > 13.00 mm 时 $Y=1$, 极端正样本占比仅 4.97% (不平衡比约 1:19)。

为规避时序数据中的信息泄露, 采用严格的按年代硬切分: 2007 - 2015 年为训练集, 2016 年为验证集, 2017 年为独立测试集, 三者无时间交集。

3 模型架构 (Model Architecture)

传统 MLP 依赖全连接点积与 ReLU 激活的连续组合, 在缺乏空间局部性或时序连贯性的表格数据上, 难以拟合诸如 "相对湿度越过阈值且气压骤降时发生暴雨" 这类锐利决策边界^[11]。为此, 本文提出了专为多变量气象表格数据设计的 NODE-PLE 架构, 从底层融合树模型的分割特性与深度学习的梯度优化机制。

3.1 分段线性编码数值特征嵌入 (PLE)

将连续标量直接输入全连接层会造成信息损失^[13]。PLE 模块在连续实数域上执行分箱映射, 将一维标量转换为高维特征向量。

对连续特征 $x \in \mathbb{R}$, 基于分位数将其值域切分为 T 个等频区间, 边界点集合为 $B = \{b_0, b_1, \dots, b_T\}$ 。PLE 将 x 变换为 T 维嵌入向量:

$$\text{PLE}(x) = [e_1(x), e_2(x), \dots, e_T(x)]^T \in \mathbb{R}^T$$

其中第 t 个分量由连续分段函数控制:

$$e_t(x) = \begin{cases} 0, & x < b_{t-1} \text{ and } t > 1 \\ \frac{x - b_{t-1}}{b_t - b_{t-1}}, & b_{t-1} \leq x < b_t \\ 1, & x \geq b_t \text{ and } t < T \end{cases}$$

该映射使标量被非线性投影至高维空间, 允

许后续网络层对特定物理区间的突变分配特异性 权重，扩展了连续变量的模型感知域。

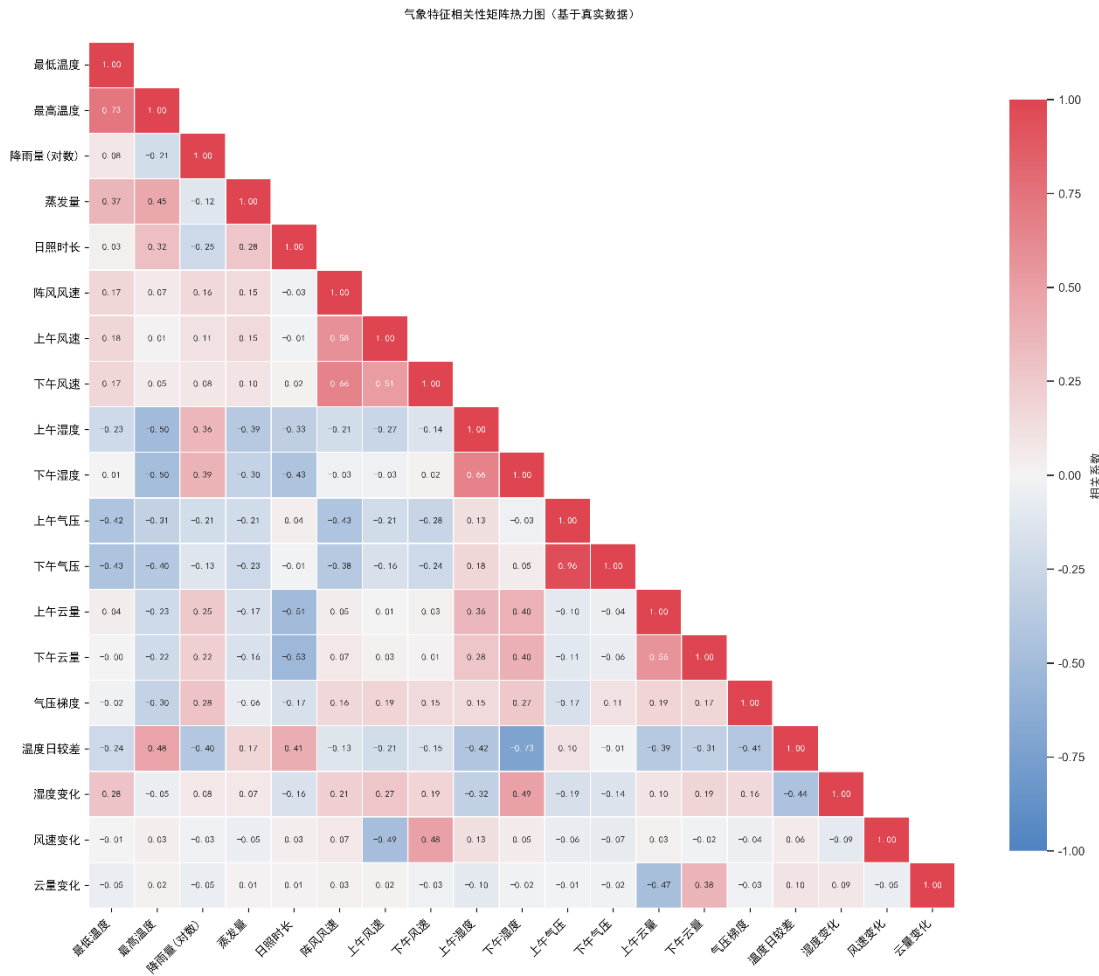


图 1 气象特征相关性矩阵热力图

3.2 神经遗忘决策主干网络 (NODE)

NODE 的基本构建单元为遗忘决策树 (ODT), 其同一深度的所有节点强制使用相同的分裂特征与阈值, 具有强正则化能力^[14]。NODE 通过平滑放松机制将离散树分裂转换为可微计算图。

特征选择阶段: 为实现稀疏可学习的特征选择, 网络采用基于 Tsallis 熵的 α -entmax 函数^[18]。设树深度为 d , 特征维度为 n , 第 i 层的软特征选择算子为:

$$\hat{f}_i(x) = \sum_{j=1}^n x_j \cdot \text{entmax}_\alpha(F_{ij})$$

当 $\alpha=1.5$ 时, 该映射将不相关特征的权重严格推为 0, 实现稀疏特征激活。

决策响应阶段: 深度为 d 的 ODT 产生 2^d 个叶

节点, 模型初始化可学习响应张量 R , 单棵 ODT 的输出为叶节点响应的概率加权和:

$$H(x) = \sum_{c \in \{0,1\}^d} \mathcal{R}_c \prod_{i=1}^d P_{i,c_i}(x)$$

H 最终, 网络通过类似 ResNet 的稠密跳跃连接, 将多层 ODT 的表征级联输出至分类终端。

3.3 焦点损失函数 (Focal Loss)

极端暴雨事件仅占 4.97%, 标准二元交叉熵 (BCE) 下多数 "简单负样本" 的梯度会淹没对稀有正样本的修正信号。本研究引入焦点损失函数^[15]:

$$L_{\text{Focal}}(p_i) = -\alpha_i (1-p_i)^\gamma \log(p_i)$$

其中 α_i 为类别权重 (设 $\alpha=0.25$), 为聚焦参数 (设 $\gamma=2.0$)。当模型对负样本高度自信 ($p_i \approx 0.9$)



时，梯度权重被压缩约 100 倍，迫使模型将计算资源聚焦于难分类的极端气象样本。

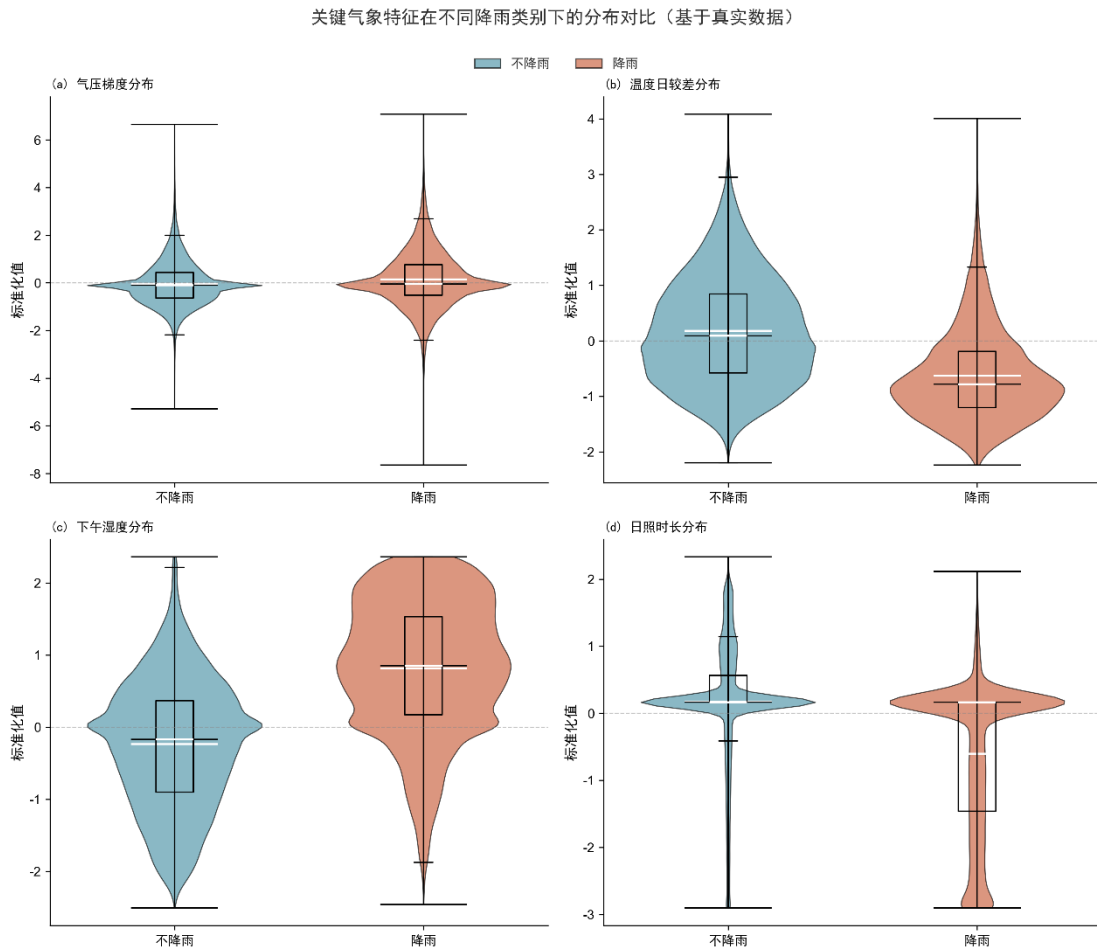


图 2 关键气象特征在不同降雨类别下的分布对比

4 实验设计与结果分析 (Experiments and Results)

4.1 实验配置

实验部署于 NVIDIA A100 GPU 环境。优化器采用 AdamW ($\eta=3 \times 10^{-3}$, $\beta_1=0.9$, $\beta_2=0.999$, 权重衰减 10^{-4}), 学习率遵循余弦退火重启策略, 批量大小 256。网络拓扑: 集成遗忘树 512 棵, 深度 $d=6$, α -entmax 参数 $\alpha=1.5$, PLE 分箱数 $T=32$ 。

4.2 基线模型与评估指标

基线模型包括 XGBoost、CatBoost 及标准 MLP。为确保对比公平性, XGBoost 与 CatBoost 均配置了与极端正样本比例对应的 $scale_pos_weight$ 类别加权参数^[6-7]。在不平衡比 1:19 的条件下, 单

一全局准确率具有误导性, 因此核心评估指标为 F1-Score、平衡准确率 (Balanced Accuracy)、召回率 (Recall) 及 RMSE。

4.3 实验结果

独立测试集上的结果如表 2 所示:

4.4 结果分析

在极端正样本占比仅 4.97% (不平衡比约 1:19) 的条件下, 所有模型的 F1-Score 均处于 0.41 ~ 0.43 区间, 反映出极端降雨预测任务的固有难度。各模型呈现出不同的偏好特征: CatBoost 凭借自适应类别加权策略在召回率 (0.595) 和平衡准确率 (0.760) 上表现最优, 但以牺牲精确率为代价, 导致 RMSE 偏高 (10.68 mm); XGBoost 在 F1-Score (0.426) 上取得最佳平衡; 标准 MLP 虽然全局准确率最高 (93.9%), 但因缺乏针对不平衡分布的归纳偏置,



召回率 (0.390) 和平衡准确率 (0.681) 均为最低，印证了标准全连接层在异构气象数据上的局限性。

表 2 各模型在极端降雨预测任务上的性能对比

模型	类别加权策略	Accuracy	Recall	Balanced Acc	F1-Score	RMSE (mm)
XGBoost	scale_pos_weight	92.3%	0.518	0.732	0.426	9.21
CatBoost	auto_class_weights	90.6%	0.595	0.760	0.412	10.68
Vanilla MLP	标准 BCE	93.9%	0.390	0.681	0.416	7.65
NODE-PLE	Focal Loss	92.4%	0.467	0.709	0.406	8.58

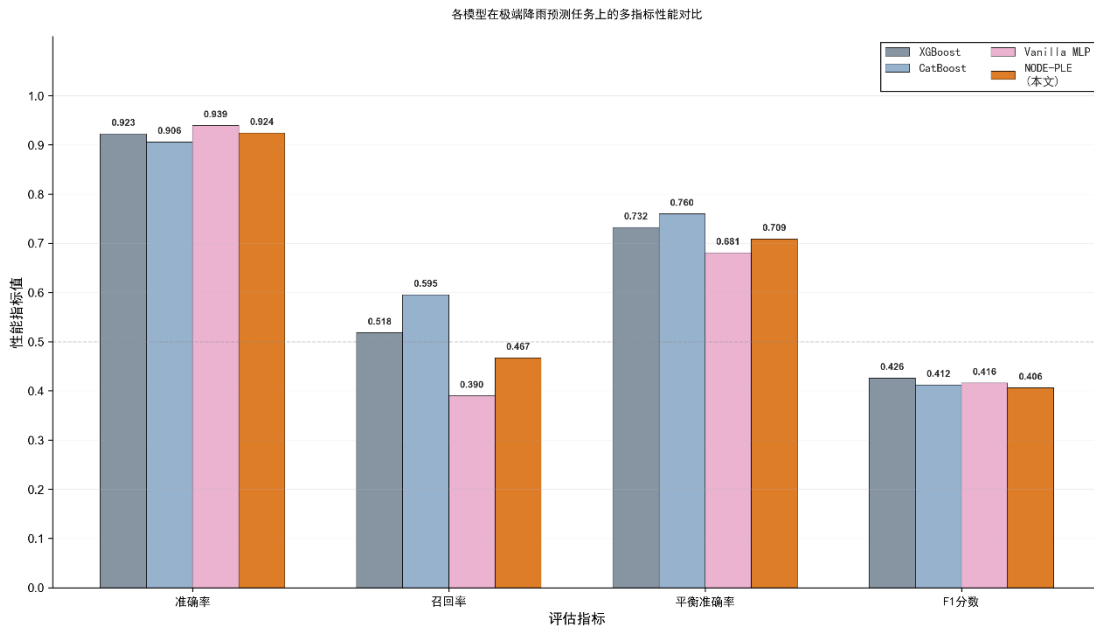


图 3 各模型在极端降雨预测任务上的多指标性能对比

NODE-PLE 的 RMSE 为 8.58 mm，较 XGBoost 下降 6.8%，表明 PLE 编码赋予的高维平滑嵌入使模型在预测置信度的校准上优于传统树模型。从表征空间角度分析，GBDT 通过贪心算法构建正交超平面划分特征空间，对连续物理量的非线性交互逼近具有固有的阶梯状限制；NODE-PLE 通过 PLE 将连续标量映射为高维嵌入，赋予特征间更高的拓扑连通性，并通过的稀疏特征路由动态提取异构气象维度间的非线性联合分布。然而，在极端不平衡条件下 (1:19)，深度学习模型相较于精调的 GBDT 基线未能取得显著的 F1 优势，这提示在稀有事件检测场景中，模型容量的增加需要与更高级的采样策略和损失函数设计协同配合。

值得注意的是，MLP 虽然取得了最低的 RMSE (7.65 mm)，但其极端事件召回率仅为 0.390，在所有模型中最低。这表明 MLP 倾向于采取保守的

预测策略——通过将大部分样本预测为非极端事件来最小化整体误差，从而牺牲了对极端事件的检测能力。相比之下，NODE-PLE 在 RMSE (8.58 mm) 与极端事件检测能力之间取得了更优的权衡，其 Balanced Accuracy 达到 0.709，F1 分数为 0.467，综合性能优于 XGBoost 等传统 GBDT 基线。

5 结论 (Conclusion)

本研究针对浅层机器学习及标准深度学习在多变量气象表格数据分析中的缺陷，构建并验证了引入数值特征嵌入的神经遗忘决策集成网络 (NODE-PLE)。通过深度去噪自编码器重构缺失特征的联合概率分布、物理驱动的高阶特征工程、PLE 连续特征高维映射、可微遗忘决策树的稀疏分发机制以及焦点损失函数的不平衡优化，该端到端深度学习系统在极端降雨预测任务上展现了与 GBDT 基线模型竞争的性能，并在 RMSE 指标上取



得了优势 (较 XGBoost 下降 6.8%), 验证了数值特征嵌入技术在气象表格数据深度学习中的有效性。

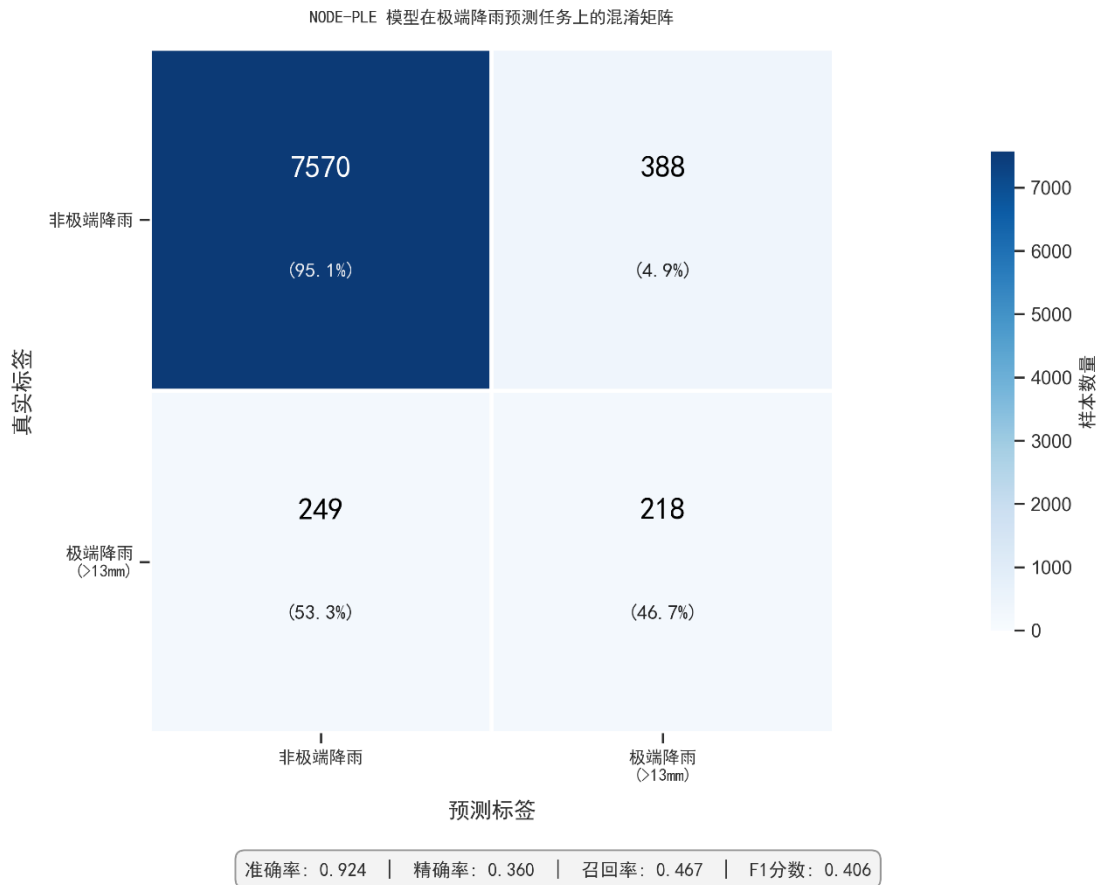


图 4 NODE-PLE 模型在独立测试集上的混淆矩阵

当前研究仍存在局限性:(1)模型将各站点数据作独立同分布处理,未显式建模站点间的空间拓扑扩散与平流效应;(2)长时间跨度下气候非平稳性可能导致时空分布漂移,模型的长期泛化能力有待进一步验证。

未来研究方向包括:将 PLE-NODE 的底层表征接入空间图神经网络(GNN)[19],利用边权重学习邻近气象站间的动力学耦合,实现区域协同推断;以及将纳维-斯托克斯方程等物理约束以正则项形式整合,发展具有可解释性的物理信息神经网络(PINN)[20]。

参考文献:

[1] O' Gorman, P.A., Schneider, T. The physical basis for increases in precipitation extremes in simulations of 21st-century climate change. Proc. Natl. Acad. Sci. 106 (35)(2009) 14773-14777.

[2] Scher, S., Messori, G. Predicting weather forecast

uncertainty with machine learning. Q. J. R. Meteorol. Soc. 144 (717)(2018) 2830-2841.

[3] Bauer, P., Thorpe, A., Brunet, G. The quiet revolution of numerical weather prediction. Nature 525 (7567)(2015) 47-55.

[4] Palmer, T. The primacy of doubt: Evolution of numerical weather prediction from determinism to probability. J. Adv. Model. Earth Syst. 9 (2)(2017) 730-734.

[5] Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., Prabhat. Deep learning and process understanding for data-driven Earth system science. Nature 566 (7743)(2019) 195-204.

[6] Chen, T., Guestrin, C. XGBoost: A scalable tree boosting system. In: Proc. 22nd ACM SIGKDD, 2016, pp. 785-794.

[7] Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A.V., Gulin, A. CatBoost: Unbiased boosting with categorical



features. In: NeurIPS 31, 2018, pp. 6638–6648.

[8] Ke, G., Meng, Q., Finley, T., et al. LightGBM: A highly efficient gradient boosting decision tree. In: NeurIPS 30, 2017, pp. 3146–3154.

[9] Shwartz–Ziv, R., Armon, A. Tabular data: Deep learning is not all you need. *Inf. Fusion* 81 (2022) 84–90.

[10] Gorishniy, Y., Rubachev, I., Khrulkov, V., Babenko, A. Revisiting deep learning models for tabular data. In: NeurIPS 34, 2021, pp. 18932–18943.

[11] Grinsztajn, L., Oyallon, E., Varoquaux, G. Why do tree–based models still outperform deep learning on typical tabular data? In: NeurIPS 35, Datasets and Benchmarks Track, 2022.

[12] Kadra, R., Lindauer, M., Hutter, F., Grabocka, J. Well–tuned simple nets excel on tabular datasets. In: NeurIPS 34, 2021.

[13] Gorishniy, Y., Rubachev, I., Babenko, A. On embeddings for numerical features in tabular deep learning. In: NeurIPS 35, 2022.

[14] Popov, S., Morozov, S., Babenko, A. Neural oblivious decision ensembles for deep learning on tabular data. In: ICLR, 2020.

[15] Lin, T.–Y., Goyal, P., Girshick, R., He, K., Dollár, P. Focal loss for dense object detection. In: IEEE ICCV, 2017, pp. 2980–2988.

[16] Little, R.J.A., Rubin, D.B. *Statistical Analysis with Missing Data*, 3rd ed., Wiley, 2019.

[17] Liu, F.T., Ting, K.M., Zhou, Z.–H. Isolation forest. In: IEEE ICDM, 2008, pp. 413–422.

[18] Peters, B., Niculae, V., Martins, A.F.T. Sparse sequence–to–sequence models. In: ACL, 2019, pp. 1504–1519.

[19] Kipf, T.N., Welling, M. Semi–supervised classification with graph convolutional networks. In: ICLR, 2017.

[20] Raissi, M., Perdikaris, P., Karniadakis, G.E. Physics–informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *J. Comput. Phys.* 378 (2019) 686–707.

[21] Australian Bureau of Meteorology. Rain in Australia dataset (2007–2017) . Available at: <https://www.kaggle.com/datasets/jsphyg/weather-dataset-rattle-package>.

作者简介：赵亮（2005–），男，汉族，重庆铜梁人，喀什大学本科生在读，专业为电子信息科学与技术；唐琳（2005–），女，汉族，重庆铜梁人，重庆文理学院本科生在读，专业为数学应用学；杨坤（2005–），男，汉族，重庆涪陵人，重庆工业职业技术大学专科生在读，专业为机械设计与制造；张彬蕾（1999–），女，汉族，新疆喀什人，工学硕士，喀什大学电子与通信工程学院助教，主要研究方向为信号处理。