

面向开放特征空间的概念演化检测方法

苏 睿¹, 郭虎升^{1,2}, 王 婧¹, 王文剑^{2*}

(1. 山西大学计算机与信息技术学院, 山西太原 030006;

2. 山西大学计算智能与中文信息处理教育部重点实验室, 山西太原 030006)

摘 要: 在众多现实场景中数据以流的形式持续产生, 由于流数据具有动态变化的特点, 在生成过程中可能产生新的类别, 也被称为概念演化. 概念演化是流数据挖掘模型预测性能衰退甚至预测失效的主要原因. 因此, 能及时发现类空间变化并提醒模型做出适应性调节的概念演化检测方法受到广泛关注. 然而, 目前多数概念演化检测方法基于特征空间静态不变的假设构建算法. 在现实场景中, 特征空间同样具有动态性, 属于开放的空间. 具体来说, 随时间推移可能出现部分特征消失和新特征产生的现象, 从而破坏上述假设并导致已有算法失效. 针对这一问题, 本文提出一种面向开放特征空间的概念演化检测方法 (Concept evolution Detection method for Open Feature space, CD_OF). 该方法通过构建微簇集成模型对新进实例分类, 对于开放特征空间中的旧特征消失问题, 通过利用转移矩阵将旧特征中包含的信息转换到共享特征中; 对于新出现的特征, 拓展共享特征空间并重构集成模型. 在此基础上, 根据样本的共享邻域信息定义样本间相似度以检测概念演化, 并建立动态衰减模型, 以解决开放特征空间下的类消失和类循环问题. 实验结果表明, 本文所提出的方法能够对开放特征空间中特征的变化作出及时响应, 增强概念演化检测的能力, 在特征空间变化的真实流数据中与现有方法相比, 错误率降低了 1.7%~11.4%.

关键词: 概念演化; 开放特征空间; 特征相似性度量; 共享邻域; 动态衰减模型; 在线学习

基金项目: 国家自然科学基金 (No.U21A20513, No.62276157, No.62476157); 山西省重点研发计划 (No.202202020101003)

中图分类号: TP391 **文献标识码:** A **文章编号:** 0372-2112(2025)10-3718-12

电子学报 URL: <http://www.ejournal.org.cn> **DOI:** 10.12263/DZXB.20250416

Concept Evolution Detection Method for Open Feature Space

SU Rui¹, GUO Hu-sheng^{1,2}, WANG Jing¹, WANG Wen-jian^{2*}

(1. School of Computer and Information Technology, Shanxi University, Taiyuan, Shanxi 030006, China;

2. Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, Shanxi University, Taiyuan, Shanxi 030006, China)

Abstract: In many real-world scenarios, data is continuously generated in the form of streams. Due to the dynamic nature of streaming data, new categories may emerge during the generation process, which is known as concept evolution. Concept evolution is one of the primary challenges leading to the degradation or even failure of predictive performance in stream mining models. Therefore, concept evolution detection methods capable of promptly identifying changes in the class space and alerting models to perform adaptive adjustments have attracted widespread attention. However, most of the current concept evolution detection methods construct algorithms based on the assumption that the feature space is static and unchanging. In real scenarios, the feature space is also dynamic and belongs to the open space. Specifically, over time, some features may disappear and new features may emerge, thus violating the above assumption and causing existing algorithms to fail. To address this problem, this paper proposes a concept evolution detection method for open feature space (CD_OF). The method constructs a micro-cluster ensemble model to classify incoming instances. For the problem of disappearing old features in the open feature space, the information contained in the old features is converted to the shared features through the transfer matrix; for the newly emerged features, the shared feature space is expanded and the integration model is reconstructed. On this basis, the inter-sample similarity is defined based on the shared neighborhood information of the samples to detect concept evolution, and the dynamic decay model is established to solve the class vanishing and classifications cycling problems under the open feature space. The experimental results show that the method

proposed in this paper is able to respond to the changes of features in the open feature space in a timely manner and enhance the ability of concept evolution detection. The error rate is reduced by 1.7% to 11.4% compared to existing methods on real streaming data with feature space variations.

Key words: concept evolution; open feature space; feature similarity measure; shared neighborhood; dynamic decay model; online learning

Foundation Item(s): National Natural Science Foundation of China (No.U21A20513, No.62276157, No.62476157); Key Research and Development Program of Shanxi Province (No.202202020101003)

1 引言

在许多现实世界应用领域中,例如交通控制、社交媒体、市场分析和传感器网络,数据以流的形式不断产生^[1,2],这些数据也被称为流数据。流数据是连续高速到达的无限有序的数据项序列,具有动态性、时序性、无限性和不可再现性等特点,给数据的收集、存储、分析和处理带来了严峻的挑战^[3,4]。由于流数据存在上述特性,为提升泛化能力而设计的离线学习算法难以适应流数据的变化。相比于离线学习算法,流数据挖掘算法还需要进一步考虑算法的实时推理能力和适应能力^[5]。其中,实时推理能力要求算法在优化过程中不能影响模型的推理过程,因此相关算法的设计不仅需要考虑计算效率与存储效率,并且还需要选择高效的在线优化策略以保证学习效果。适应能力要求算法适应流数据随时间推移在特征空间与标签空间可能产生的动态变化。近年来,随着机器学习与人工智能技术快速发展,大量的在线学习算法被部署在实际场景中,然而这些算法常常会受到环境动态变化的干扰导致泛化能力逐步衰退^[6,7]。

流数据的动态变化主要会引起概念漂移与概念演化这2种问题。其中,概念漂移是指由流数据特征空间和标签空间联合分布的变化导致的模型预测表现衰退问题^[8-11]。概念演化是指流数据可能动态生成新的类别,而这些类别往往是已部署模型所未知的^[12,13],从而导致模型预测错误,如在网络安全领域,每当发生网络攻击入侵时,计算机系统的安全就会受到威胁,而新的计算机病毒可能以不同方式对计算机系统进行攻击,从而在对网络入侵进行防护时,检测到入侵方式的类别标签会增加,这些新的类标签应及时处理,否则分类器将把与新类相关联的所有数据实例错误分类为现有类,严重影响已有分类表现。这两种问题是动态流数据挖掘算法的主要挑战。

目前,虽然有多种概念演化和概念漂移检测适应算法,但是这些工作大多基于特征空间静态不变的假设,即假设特征空间维度不会发生变化。然而在许多实际应用中,数据流的特征空间会随时间动态变化,如在环境监测中需要增加新的传感器并结合部分旧传感器数据以完成新物种的监测任务;在垃圾邮件监测任务

中新型垃圾邮件往往具有主题标签、短URL(Uniform Resource Locator)等特征;在循环流化床工艺中合适的补偿机制需要考虑日益变化的新燃料供给和不均匀性数据等^[2,14]。因此,静态特征空间假设难以满足不断快速发展变化的真实场景。这些包含新特征的增加、旧特征的消失等特征变化场景的空间也被称为开放特征空间。

针对上述问题,本文提出了一种面向开放特征空间的概念演化检测方法(Concept evolution Detection method for Open Feature space, CD_OF),旨在提升当流数据局部旧特征消失或局部新特征出现场景下对新类别的检测表现。CD_OF通过寻找旧特征空间到共享特征空间的转移矩阵,利用转移矩阵实时将当前特征映射到共享特征空间,以缓解旧特征消失的影响。对于新增特征,CD_OF基于共享特征空间进行拓展,通过学习新特征空间知识重构集成模型以适应。在对齐特征空间之后,CD_OF结合共享近邻与样本相似度2项指标检测概念演化是否发生,并通过构建动态衰减模型,以适应动态流数据空间存在的类消失和类循环现象。

本文的主要贡献如下:(1)提出了一种基于共享特征空间的新特征空间对齐方法,旨在缓解特征空间变化对概念演化检测与适应的影响;(2)提出了一种融合共享近邻与样本相似度的概念演化检测方法以识别新类,构建了动态衰减模型识别消失类,提高新类检测性能的同时降低了在循环类样本的误报率。

2 相关工作

2.1 概念演化检测

对于流数据中存在的概念演化问题,现有的检测方法依据分类器构造方式,可以分为2种类型:基于聚类分析的概念演化检测方法和基于模型的概念演化检测方法。

基于聚类分析的概念演化检测方法^[9,12,13,15-19]旨在利用传统的聚类算法对已知类进行建模,并据此定义决策边界。通常假设同一类别的样本在特征空间中彼此之间更接近,而不同类别的样本则相对较远。通过检测传入的样本实例是否位于决策边界内,来有效地识别和区分样本中的异常点。对于异常值,Masud等人^[14]

提出 ECSMiner (Enhanced Classifier for data Streams with novel class Miner) 方法, 通过分析其在特征空间中与已知类样本和其他异常点之间的距离关系, 计算异常点的 q -邻域轮廓系数 (q -Neighborhood Silhouette Coefficient, q -NSC), 来确定该样本是否属于新类. 在这种方法中, 每个分类器的决策边界保持固定, 这可能会把概念漂移导致的落入决策边界外的实例错误识别为异常值. Al-khateeb 等人^[20]提出的 CLAM (Class based Micro classifier ensemble) 方法在每一类初始数据上进行聚类, 对于落在聚类边界外的测试实例, 通过分析是否有足够的彼此接近的异常值来寻找新类. CLAM 方法能够很好地区分循环类和新类, 然而该方法不能区分循环类和已知类, 通常将循环类标识为已知类.

基于模型的概念演化检测方法旨在找到一个分类模型来识别新的类, 包括基于图的概念演化检测方法^[21,22]、基于决策树的概念演化检测方法^[23-25]和基于草图的概念演化检测方法^[26]. 其中, 基于图的检测方法有: SACCOS^[21]方法使用相互图聚类技术, 识别特征空间中局部区域内相似的数据实例; GNG^[22]方法是一种基于图的自适应方法, 在不断发展的数据流中进行增量学习. 以样本实例作为节点, 不同节点之间的连接作为边构建图. 当新样本实例到达时, 通过创建新节点和节点之间的连接更新图模型. 基于决策树的检测方法有: AhtNODE^[24]方法使用 Hoeffding 树学习器和 ADWIN 变化检测器, 检测数据流中的概念漂移和概念演化. 根据类的内聚性和分离性, 识别落在 Hoeffding 树中未使用空间内的具有强内聚性的样本实例, 声明新类. 基于草图的概念演化检测方法有: Du 等人^[26]为整个训练数据创建全局草图, 为数据中的每个类创建局部草图, 全局草图用于检测数据流中的新类实例, 局部草图用于对与已知类关联的数据实例进行分类.

2.2 开放特征空间学习

数据流学习环境的动态可变性, 通常会导致特征

空间随着时间的推移在特征维度以及属性上发生改变. 最早与开放特征空间相关的研究是增量属性学习^[27]. 当数据中出现新特征时, 一方面在旧特征上继续训练已有的模型, 另一方面在新特征上训练新的模型, 最后将得到的新旧模型结合到一起. 这种学习算法没有解决旧特征消失的问题. Hou 等人^[28]提出的 FESL (Feature Evolvable Streaming Learning) 方法采用最小化二乘损失的方法来学习新出现特征与旧特征之间的映射关系, 用新的特征来恢复旧特征, 并进一步引入集成学习思想, 通过加权结合新旧模型对样本进行预测. Liu 等人^[29]提出面向特征继承性增减的在线分类算法, 将被动-主动方法与结构风险最小化原则相结合, 在原始特征空间使用被动方法训练模型, 在新特征空间上使用主动方法训练模型. 同时, 采取矩阵补全算法, 补全恢复数据流的缺失部分. He 等人^[30]提出了 OCDS (Online learning from Capricious Data Streams) 算法, 根据特征之间的相关性, 将不同阶段的数据特征整合到一个通用特征空间中, 同时在该空间上训练一个学习器, 使得在旧特征空间中学习的模型可以应用于新特征空间.

在真实场景中, 特征空间变化与概念演化往往同时存在, 然而现有方法仅能适应单一场景, 在真实的复杂动态变化场景下难以保持方法的泛化性.

3 面向开放特征空间的概念演化检测方法

针对上述问题, 本文提出 CD_OF. 该方法首先构建微簇集成模型对新进实例分类, 对于开放特征空间中的旧特征消失问题, 通过度量旧特征和共享特征之间的相似性, 将旧特征中包含的信息转换到共享特征中; 反之, 对于新特征出现的问题, 添加相应的新特征维度拓展特征空间, 在拓展后的特征空间上重构集成模型. 在此基础上, 结合共享近邻与样本相似度检测概念演化, 并建立动态衰减模型, 以适应类消失和类循环. 图 1 为该方法的整体框架示意图.

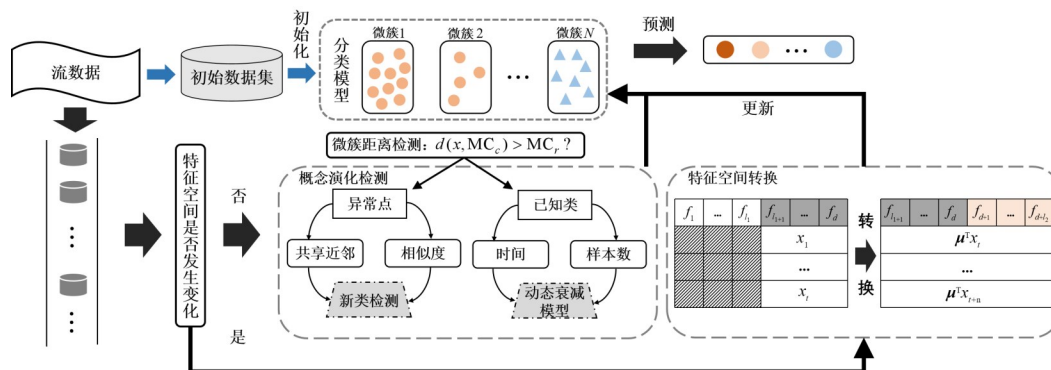


图 1 CD_OF 方法整体框架图

3.1 问题定义

流数据是一组大量、快速、持续到达的数据序列,一般情况下,流数据可被视为一个随时间延续而无限增长的动态数据集. 可以将其表示如式(1)所示:

$$D = \{(x_i, y_i)\}_{i=1}^t \quad (1)$$

其中, (x_i, y_i) 是 t 时刻到达的样本实例, x_i 是特征向量, $x \in \mathbf{R}^d$, y_i 是该样本实例对应的标签 $y \in \mathbf{R}^L$.

若在流数据分类任务中 t 时刻发生了概念演化,可表示为标签空间的概率分布发生变化,如式(2)所示:

$$P_{t-1}(y) \neq P_t(y) \quad (2)$$

这往往是流数据中类别标签增加或者消失而导致的.

在开放特征空间下,随着时间的推移,流数据特征空间可能发生变化. 这种变化可表示为从旧特征空间到新特征空间的转换. 假设初始阶段的旧特征空间为 $S_t = \{f_1, f_2, \dots, f_d\}$, 表示从最初时刻到时刻 t 的样本 x_t 特征空间, 即 $x_t \in \mathbf{R}^d$. 在 $t+n$ 时刻, 特征空间发生变化, 空间中有 l_1 个特征消失且有 l_2 个特征出现, 并且还有 $d-l_1$ 个特征没有发生变化, 因此新的特征空间表示为 $S_{t+n} = \{f_{d-l_1+1}, \dots, f_d, f_{d+1}, \dots, f_{d+l_2}\}$, 此时 $x_{t+1} \in \mathbf{R}^{d-l_1+l_2}$. 特征空间的变化过程如图 2 所示.

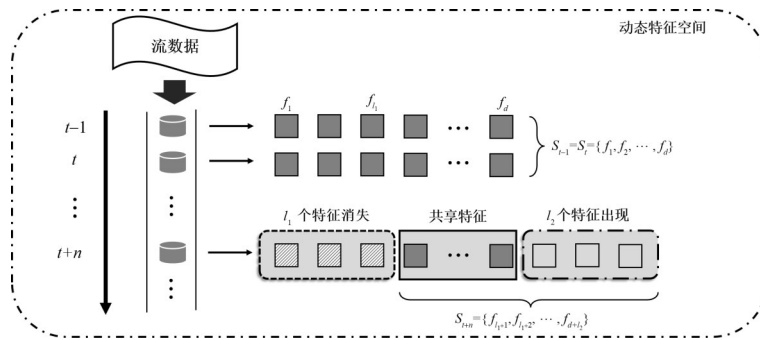


图 2 开放特征空间示例图

消失特征指只在旧特征空间中存在的特征, 共享特征指当前特征空间和上一阶段特征空间共同存在的特征, 而新特征指只在新特征空间中存在的特征.

3.2 初始集成模型构建

假设在流数据的整个过程中, 只有初始训练数据的一个子集的真实标签是已知的, 而其他数据的标签是未知的. 在此训练期间, 本文首先将初始有标签数据集按照数据标签类别, 划分为标签不相交的集合, 使得每个集合仅包含一个类的实例. 其次在每个集合上分别应用聚类算法将其划分为 n 个簇, 利用得到的聚类对训练集中的无标签样本进行分类, 直到所有无标签数据都被标记完. 最后利用标记后的数据集更新聚类簇.

由于流数据具有数据量大的特征, 存储原始数据会消耗内存空间, 因此计算每个聚类簇的特征并将其存储在微簇 (Micro-Cluster, MC)^[31] 中来代替存储原始数据. 具体地, 若初始训练数据有 L 个类, 则通过 K -means 聚类方法生成 $L \times n$ 个微簇. 微簇定义如式(3)所示:

$$MC = (Ls, Ss, N, W) \quad (3)$$

其中, $Ls = \sum_{i=1}^N x_i$ 是微簇中样本特征的线性和, $Ss = \sum_{i=1}^N x_i^2$

表示微簇中样本特征的平方和, N 是微簇中样本的数量, W 是微簇随时间推移的重要性, 初始值设置为 1. 根据微簇的定义, 可计算出每个微簇的中心 MC_c 和半径 MC_r , 计算公式见式(4)、式(5):

$$MC_c = \frac{Ls}{N} \quad (4)$$

$$MC_r = \left(\frac{Ss}{N} - \frac{Ls^2}{N^2} \right)^{\frac{1}{2}} \quad (5)$$

3.3 特征相似性度量与映射

新特征主要来源于真实场景中由环境因素的变化而引发的特征的变化, 如在环境监测中, 传感器的更换或新增, 均会产生新特征. 在旧特征空间到新特征空间的变化过程中, 存在部分旧特征消失和新特征产生的现象, 但有一些特征始终没有发生变化, 因此旧特征分布与新特征分布不完全一致. 新增特征与原有特征之间并没有确定的联系. 为了从旧特征空间转换到新特征空间, 利用不变的共享特征进行辅助, 以降低适应难度. 因此, 假设在旧特征和共享特征之间存在确定的映射关系 $\psi: \mathbf{R}^d \rightarrow \mathbf{R}^{d-l_1}$. 假设 x_t 是初始阶段缓冲窗口中的一个样本, x'_t 是由 x_t 样本降维到 $d-l_1$ 维度后的部分特征再进行补 0 得到, x'_t 去掉尾部 l_1 个特征即为初始共享特征. 为获取当前特征分布与共享特征空间之间的最佳映射参数矩阵 μ^* , 本文采用局部加权线性回归方法, 具体表示如式(6)所示:

$$\min_{\mu \in \mathbf{R}^{d-l_1}} \sum_{i=1}^{T_w} v_i (\mu^T x_i - x'_i)^2 \quad (6)$$

其中, μ 表示旧特征空间到共享特征空间的转移矩阵, T_w 即缓冲窗口, 高斯核权重 v 可使用 x'_i 与 x_i 计算得到,

如式(7)所示:

$$v_t = e^{-\frac{(x_t - x'_t)^2}{2\theta^2}} \quad (7)$$

对应时刻的最优转移矩阵求解如式(8)所示:

$$\mu^* = \left(\sum_{t=1}^{T_v} x_t v_t x_t^T \right)^{-1} \left(\sum_{t=1}^{T_v} x_t v_t x'_t{}^T \right) \quad (8)$$

在得到转换矩阵 $\hat{\mu}$ 后,便可以将旧空间特征转换到共享特征空间,得到转换后的样本 x'_{t+1} .当特征空间发生变化后,利用转换矩阵 $\hat{\mu}$ 将新特征空间下的数据样本转换为基于旧特征重构的新样本再重新进行聚类,根据聚类结果更新微簇.其中,对于没有对应数据的微簇,其重要性系数将不断衰减,如果重要性系数降为0,该微簇将被移除;对于具有对应数据的微簇,其聚类特征将得到更新并且重要性系数将得到提升;

此外,还有数据难以被聚类到合适的微簇中,那么将新增微簇,但该微簇所属类别还需要概念演化检测器识别.通过上述转换,特征空间与微簇将统一被转换到新特征空间中,实现特征空间对齐,既避免了特征空间变化对模型预测表现的影响,同时使模型更关注于新特征学习,有助于更好地识别新特征空间下样本类别.

综上所述,当特征空间发生变化,特征维度转换如图3所示.假设 t 时刻特征空间发生了变化,有 l_1 个特征消失,因此先学习特征空间 R^d 和共享特征空间 R^{d-l_1} 之间的线性映射系数矩阵 $\hat{\mu}$,其次通过转移矩阵 $\hat{\mu}$ 将旧微簇转换到共享特征空间.同时,对于新样本,使得 $t+1$ 时刻样本特征经转移矩阵 $\hat{\mu}$ 映射,包含消失的特征中所携带的旧特征空间信息,然后利用维度扩充得到新特征空间下的样本.

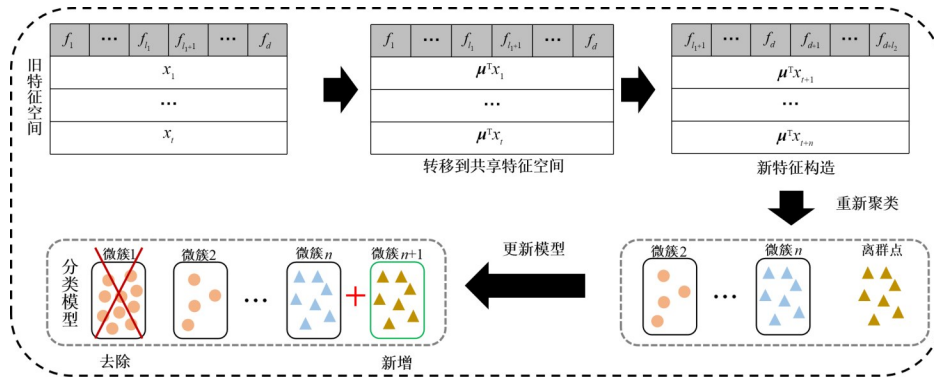


图3 空间转化过程图

3.4 特征相似性度量与映射

当新样本 x_t 到达并通过特征空间检测与处理后,新样本和微簇都将被转换到新的特征空间.在转换后的样本空间中,CD_OF将对当前批次样本进行概念演化检测,以观察是否存在新类别样本.本文利用微簇集成模型判断新样本是否为异常点,主要通过比较样本到每个微簇中心的距离与微簇半径的大小关系,对样本进行区分.具体地,若新样本到每个微簇中心的距离大于微簇半径,则认为该样本是异常点,存入临时缓冲区 B 中;反之,则认为该样本属于已知类,对其进行分类,标签获取如式(9)所示:

$$y = \arg \max_{i \in \{1, 2, \dots, k\}} \frac{W_i N_i}{D(x_t, MC_{c_i})} \quad (9)$$

其中, W_i 是第 i 个微簇的重要性, N_i 是第 i 个微簇中所含样本数量, $D(x_t, MC_{c_i})$ 是样本到第 i 个微簇中心的距离.对于已知类样本,选择微簇集成中距离该样本最近的 k 个微簇,根据这 k 个微簇所含样本数量、微簇重要性以及样本到微簇中心的距离这三个指标综合判断样本标签,并更新相应微簇.

对于缓冲区中的异常点,当缓冲区 B 满时,通过对

缓冲区内样本检测以判断是否出现新类,即是否存在概念演化.本文基于新类实例是一组具有较强内聚性的相似数据假设,根据样本的共享邻域信息定义样本间的相似度,通过与已知类样本的相似度进行比较,从而判断样本是否属于新的类别.其中,样本 x_i 与 x_j 的共享邻域^[32]是指 x_i 与 x_j 的 k 近邻交集,定义如式(10)所示:

$$\text{SNN}(x_i, x_j) = \text{KNN}(x_i) \cap \text{KNN}(x_j) \quad (10)$$

其中, $\text{KNN}(x_i)$ 是 x_i 的一组最近邻, $\text{KNN}(x_j)$ 是 x_j 的一组最近邻.

根据共享邻域,计算样本 x_i 与 x_j 的相似度,计算如式(11)所示:

$$\text{Sim}(x_i, x_j) = \begin{cases} e^{-\sum_{q \in \text{SNN}(x_i, x_j)} (d_{iq} + d_{jq})} \times e^{\frac{|\text{SNN}(x_i, x_j)|}{k}} & \text{if } i, j \in \text{SNN}(x_i, x_j) \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

其中, d_{iq} 表示样本 x_i 与共享近邻 x_q 之间的欧氏距离.当 x_i 与 x_j 是彼此的 k 近邻时才会计算相似度,样本点之间的距离越近,不同样本所属空间的相似度越高.

$|\text{SNN}(x_i, x_j)|$ 是样本 x_i 与 x_j 共享邻域中样本点的个数,共

享邻域中样本点个数越多,样本之间的相似度越高.

若缓冲区中的一个样本与其在缓冲区中的共享近邻相似度高,但是在已知类样本中的共享近邻相似度高,则认为该点是新类实例.具体计算如式(12)所示:

$$\gamma(x_i) = \frac{\sum_{x_j \in \text{KNN}(x_i)} \text{Sim}(x_i, x_j) - \sum_{x_p \in \text{KNN}(x_i)} \text{Sim}(x_i, x_p)}{\max\left(\sum_{x_j \in \text{KNN}(x_i)} \text{Sim}(x_i, x_j), \sum_{x_p \in \text{KNN}(x_i)} \text{Sim}(x_i, x_p)\right)} \quad (12)$$

其中, x_j 是 x_i 在缓冲区中的最近邻, x_p 表示 x_i 在微簇中的最近邻. $\gamma(x_i)$ 的取值范围是 $(-1, 1)$, 当 $\gamma(x_i) > 0$ 时, 该样本为潜在的新类实例. 此时, CD_OF 将得到的新类实例移入新类实例缓冲区中, 过滤掉缓冲区中的所有新类实例点, 剩余的异常点属于已知类.

若对于某个微簇, 较长时间内没有新进样本落入该微簇决策边界内, 则认为该微簇的重要性随着时间的推移而降低, 即不能够捕获数据流中的最新概念. 当微簇的重要性近似达到 0 时, 集成模型将消除这些过时的微簇. 具体地, 使用应用指数衰减函数来计算微簇的重要性, 定义如式(13)所示:

$$W = W \times e^{-\lambda \frac{\sum_{j=1}^{j+1} (t_{j+1} - t_j)}{N}} \quad (13)$$

其中, W 的初始值设置为 1, t_{j+1} 和 t_j 分别是微簇中第 $j+1$ 和第 j 个样本出现的时刻, λ 是重要性衰减参数, N 表示该微簇中的样本个数. 随着该微簇中数据出现的平均时间间隔越来越大, 该微簇的重要性将逐渐降低. 当 W 趋向于 0 时, 从模型中消除该微簇.

3.5 CD_OF 算法

CD_OF 的算法伪代码如算法 1 所示.

3.6 时间复杂度分析

本文中时间复杂度主要由 2 个部分构成, 分别是局部线性回归和簇重新分配. 其中, 在局部线性回归中, 由于矩阵直接求逆具有较高的时间复杂度, 因此本文先通过 LU 分解法 (Lower-Upper Decomposition) 将矩阵分解为上三角矩阵和下三角矩阵, 然后利用三角矩阵求解逆矩阵. 虽然时间复杂度仍受特征空间维度 d 影响为 $O(d^3)$, 但在实际在线学习过程中, 算法的计算效率可以得到良好的控制. 首先, 矩阵分解过程可以在特征发生变化前做好准备, 从而避免对在线学习过程的影响. 具体来说, 本文中在特征变化前使用缓冲窗口 T_w 存储未变特征, 然后使用窗口中特征求解分解结果, 求解后的结果将被存储在内存当中, 在特征空间发生变化后使用内存中的求解结果直接完成求逆运算, 具有更高的计算效率. 其次, 在在线学习过程中数据的特征维度 d 较小, 而对于高维数据可通过特征表示方法控制

算法 1 面向开放特征空间的概念演化检测算法(CD_OF)

初始化: 流数据 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_t, y_t), \dots\}$; 弱监督集成模型 $MC = \{MC_1, MC_2, \dots, MC_{L \times n}\}$

1. while 流数据未结束, 第 t 个时刻样本 x_t 进入
2. while 对于新进入的样本 x_t
3. 检测新特征空间是否发生改变, S_{t-1} 与 S_t 是否相同
4. if $S_{t-1} \neq S_t$
5. 根据式(8)求解转移矩阵 μ
6. 将微簇与当前样本特征转换到新特征空间
7. end if
8. end while
9. if $d(x_t, MC_c) < MC_r$
10. 根据式(9)对已知类进行分类
11. else
12. $B = B \cup x_t$
13. if $|B| > T_0$
14. 根据式(11)计算相似度 $\text{Sim}(x_t, x_j)$
15. 根据式(12)计算 x_t 是新类的可能性 $\gamma(x_t)$
16. if $\gamma(x_t) > 0$
17. 判定 x_t 属于新类
18. end if
19. end if
20. 根据式(13)计算指数衰减函数 W
21. if $W \rightarrow 0$
22. 消除微簇 MC_i
23. end if
24. end if
25. end while

特征维度大小, 从而避免对实时预测过程的影响. 簇重新分配即簇更新过程, 由于需要使用 K -means 方法重新聚类, 因此时间负载度与批次大小 b 和聚类数 $L \times n$ 相关, 而其余统计量的计算时间复杂度均为线性, 仅与批次大小相关, 因此时间复杂度为 $O(b + Lnb)$.

4 实验与性能分析

为验证 CD_OF 的处理性能, 本文在不同的含概念演化的标准数据集和真实数据集上进行实验验证, 并与基于集成学习的 CLAM^[20] 方法、基于霍夫丁树的 Aht-NODE^[24] 方法、基于完全随机树的 SENCForest^[23] 方法、基于半监督学习的 OSSL^[33] 方法、基于聚类的 EMC^[17] 方法和 ECSSMiner^[14] 方法进行对比.

CD_OF 的微簇初始化过程中, 每个类别下的聚类数量 n 为 25, KNN 近邻参数 k 初始化为 3, 缓冲区大小 $|B|$ 初始化为 100, 训练数据集数量为 1 000, 标签总数为 2. 因此在在线学习之前, 模型的已知类别仅为 2, 在线学习过程中新类别样本将逐渐按序到达.

4.1 数据集

本文使用 10 个数据集进行评估,每个数据集属性如表 1 所示.

表 1 实验数据集

数据集	属性维数	类别数	样本数	特征空间是否变化
IoT botnet	115	11	663 795	否
Forest cover	54	7	581 012	否
Kddcup99	41	23	494 021	否
Shuttle	9	7	58 000	是
Electricity	8	2	45 312	是
Poker Hand	10	10	11 250	是
Avila	10	12	20 867	是
HAR	561	6	10 299	否
Mnist	784	10	10 000	否
Page Blocks	10	5	5 473	是

IoT botnet 数据集收集自真实交通数据,包含 9 个商业 IoT 设备的收集结果. 每种设备结果包含由 2 种僵尸网络干扰的 10 种攻击类型流数据和常规场景中的流数据,拥有 115 种属性和 11 种类别. 本文选用了 663 795 个样本.

Forest cover 数据集收集自 7 种 30 m × 30 m 网格中的森林覆盖类型数据,包含 54 种属性和 7 种类别,共 581 012 个样本.

Kddcup99 数据集收集自 MIT Lincoln 实验室的 LAN 网络流量中的 TCP 连接,包含 23 种类别和 41 种属性,本文选择该数据集 10% 版本,共 494 021 个样本.

Shuttle 数据集包含 9 种属性和 7 种类别,共 58 000 个样本. 其中约有 80% 的数据属于类别 1.

Electricity 数据集收集自电力市场真实数据,包含 8 种属性和 2 种类别,共 45 312 个样本.

Poker Hand 数据集收集自由 52 张标准扑克牌中随机抽取 5 张组成的手牌,每张手牌包含 2 种属性花色与点数,包含 10 个属性和 10 种类别. 本文选择 11 250 个样本进行测试.

Avila^[34]数据集收集自 800 张图片并已经过标准化处理,包含 10 种属性和 12 种类别,共 20 867 个样本.

HAR 数据集收集自 30 名实验人员日常活动数据,包含 561 种属性和 6 种类别,共 10 299 个样本.

Mnist 数据集收集自 250 名实验人员手写数字图片,图片包含 28 × 32 共 784 种特征和 10 种类别,本文选择 10 000 个样本进行测试.

Page Blocks 数据集收集自 54 篇不同的文档,包含 10 种属性和 5 种类别,共 5 473 个样本.

本文所使用的数据集均来自真实数据. 为了评估 CD_OF 在动态特征空间场景下的表现,本文选择 Shuttle、Electricity、Poker Hand、Avila 和 Page Blocks 数据

集进行模拟. 为了模拟开放特征空间中特征随时间的出现与消失,本文对数据集进行了预处理. 首先将数据集平均划分为多个块,然后采用随机选择方法对特征进行处理. 通过随机选择方法确保每个块中特征组合方式不同,因此部分特征可能在后续块中消失或重新出现,新块中可能包含模型从未见过的新特征. 具体来说,块 A 和块 $A+1$ 中的特征均从原始特征空间中随机选择,在块 A 中被选中特征而在块 $A+1$ 中未被选中特征的情景模拟了特征的消失,反之,在块 A 中未被选中特征而在块 $A+1$ 中被选中特征的情景模拟了新特征的出现. 随机选择方法能够有效模拟真实环境中传感器特征失效、数据采集通道变化或信息通路重构等现象,具有更强的现实相关性和动态复杂性. 非动态特征空间下使用的数据集特征维度不发生变化,旨在单独测试 CD_OF 中概念演化检测方法和模型适应表现.

4.2 评价指标

为验证所提 CD_OF 算法的性能,本文采用了 3 种指标从新类检测情况、模型的测试精度 2 个方面对模型进行评估,具体指标如下所述.

M_{new} : 新类检测指标,新类实例被错误分类为已知类的百分比^[35]. 具体形式如式(14)所示:

$$M_{\text{new}} = \frac{F_{\text{exi}}}{T_{\text{nov}} + F_{\text{exi}}} \quad (14)$$

其中, F_{exi} 表示模型错误地将新类识别为已知类样本的数目, T_{nov} 表示模型正确识别新类样本的数目.

F_{new} : 新类检测指标,表示已知类实例被错误分类为新类的百分比^[35]. 具体形式如式(15)所示:

$$F_{\text{new}} = \frac{F_{\text{nov}}}{F_{\text{nov}} + T_{\text{exi}}} \quad (15)$$

其中, F_{nov} 表示模型错误地将已知类识别为新类样本的数目, F_{exi} 表示模型正确识别已知类样本的数目.

Err: 模型分类指标,通过计算模型错误识别的样本数目占总样本数目的比例表示累积错误率. 具体形式如式(16)所示:

$$\text{Err} = \frac{N - T_{\text{nov}} - T_{\text{exi}}}{N} \quad (16)$$

其中, N 表示总样本数目.

4.3 实验结果与分析

本文分别从特征空间固定和特征空间变化 2 种环境进行分析研究. 本文所有实验采用的计算机 CPU 为 Intel(R) Core(TM) i7-14700, 内存为 64 GB, 操作系统为 Windows 11, 编程工具为 Matlab.

4.3.1 特征空间不变场景

图 4 和图 5 分别为对比方法在特征空间固定不变数据集上得到的 M_{new} 和 F_{new} 测试结果. 通过图 4 对比实验结果,我们可以看到在 M_{new} 指标上,本文方法与其他对比

方法相比具有较低的 M_{new} 值,大多数情况下小于 20%,这说明 CD_OF 在检测概念演化的过程中,能够以较高的准确性识别新类别的数据,避免将其错误归类为已知类别,从而有助于模型及时适应新的数据分布.对于 F_{new} 指标,本文方法在这一指标上的结果都小于

5%,这表明本文方法能够有效地区分已知类的实例,避免将它们错误地归类为新类,对已知类具有很好的分类效果.上述结果表明,本文方法不仅能够有效地检测和识别新类别样本,还能够保持对已知类别样本的准确分类.

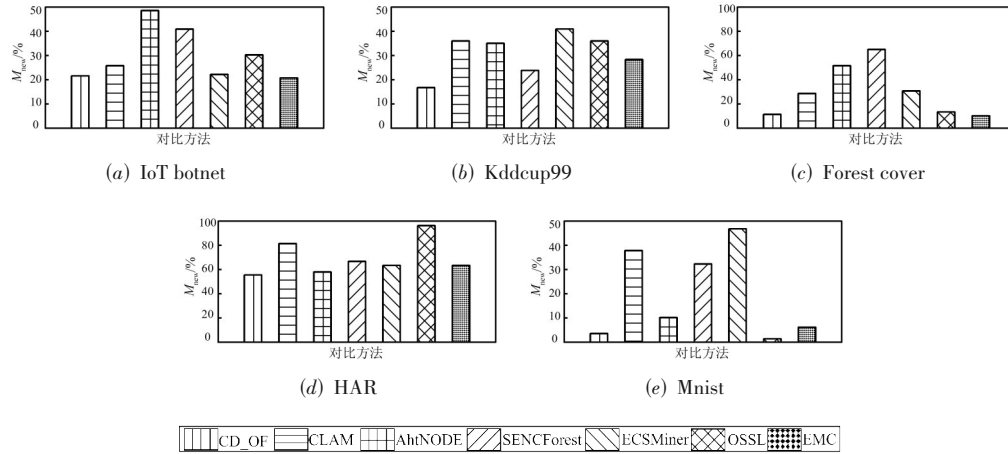


图 4 在特征空间固定数据集上的 M_{new} 指标实验结果

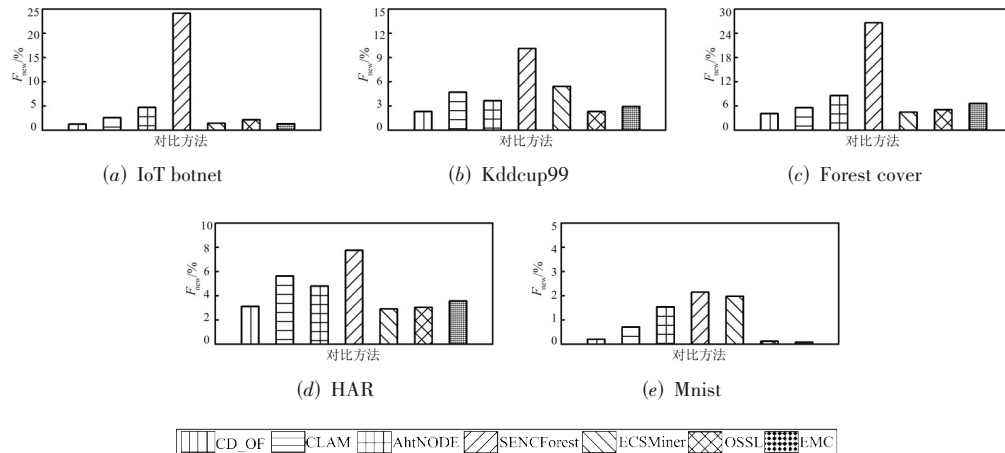


图 5 在特征空间固定数据集上的 F_{new} 指标实验结果

图 6 为对比方法在特征空间固定数据集上得到的累积错误率 Err 值.图 6 中的实验结果表明本文方法的累积错误率明显低于大部分对比方法,且在与其他对比方法的比较中也表现出微弱的优势,略低于其他对比方法.这表明本文方法能维持较为稳定的性能并保持较为准确的分类效果.

4.3.2 特征空间变化场景

图 7 和图 8 分别为对比方法在特征空间变化数据集上得到的 M_{new} 和 F_{new} 测试结果.图 7 和图 8 的实验结果展示了与对比方法相比,CD_OF 的 M_{new} 值和 F_{new} 值都低于对比方法.实验结果说明即使在特征空间发生变化时,本文方法将新类样本错误归类为已知类的概率

相对较低,且能够较好地对待已知类样本进行分类,不会将其错误地检测为新类别.上述结果表明本文方法能够较好地处理开放特征空间下的流数据概念演化检测问题,在检测新类实例方面表现出色.

图 9 为对比方法在特征空间变化数据集上得到的累积错误率 Err 测试结果.从图 9 中可以发现,CD_OF 的累积错误率明显低于其他对比方法,主要是因为 CD_OF 能够及时检测到特征空间的变化,并有效地学习旧特征与共享特征之间的映射关系,将旧特征中携带的信息迁移到共享特征中,保证了模型性能的稳定性.对比方法缺乏对特征空间变化的快速响应能力,无法迁移旧特征中的信息到共享特征空间.因此,当面对

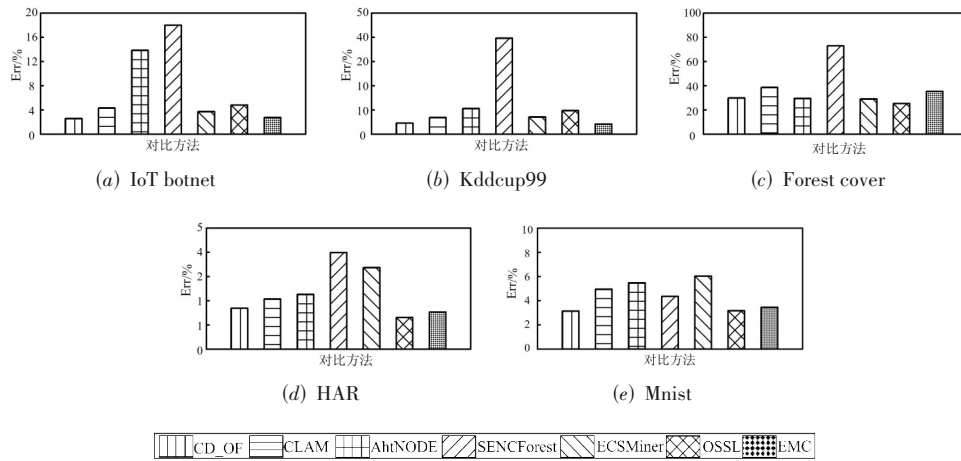


图6 在特征空间固定数据集上的Err指标实验结果

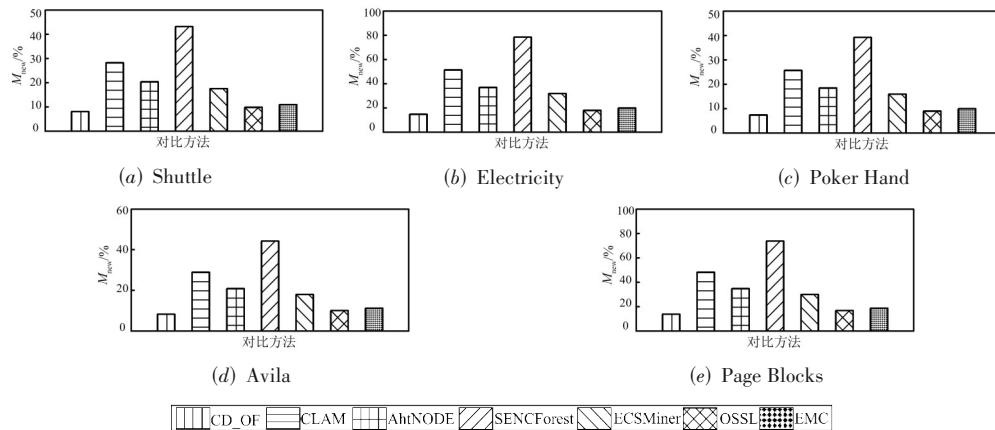


图7 不同方法在各数据集上的M_{new}比较

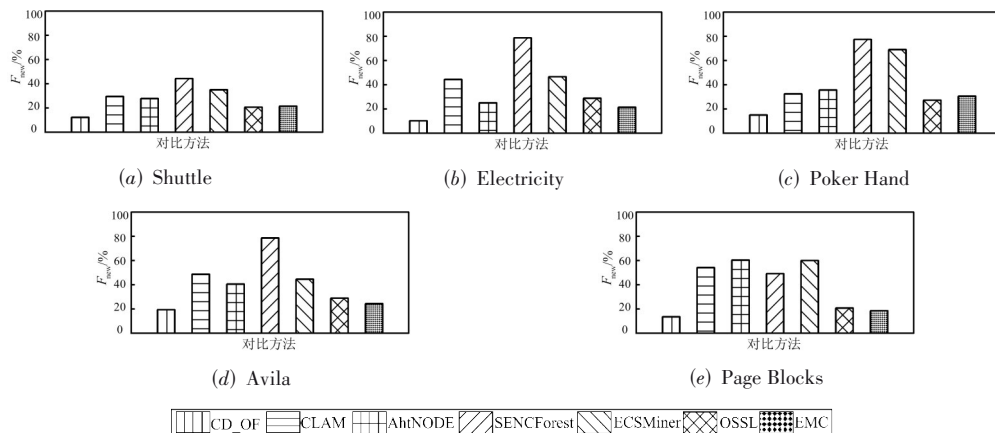


图8 不同方法在各数据集上的F_{new}比较

特征空间变化的数据流时,对比方法的累积错误率会显著增加.本文提出的方法确保了模型能够在开放特征空间中维持高水平的性能表现,提高了模型对开放特征空间的适应能力.

4.3.3 参数敏感性分析

在参数 n 敏感性实验中,设置了步长为 5,从 10~50

的对比簇数.实验结果如图 10(a)所示.当 $n=25$ 时,在所有数据集上的平均错误率最低.当 n 采用较小的数值时,难以表示类内样本之间的差异;而当 n 采用较大数值时,模型容易产生过拟合问题,难以应用于后续任务中.

在参数 k 敏感性实验中,设置了步长为 1,从 1~5 的近邻数.实验结果如图 10(b)所示.当 $k=3$ 时,在所有数

数据集上的平均错误率最低. 合适的 k 值对分类过程中选择微簇、微簇参数量更新和新类检测过程中共享近邻判断紧密相关. 通过上述实验,当 k 选择3时效果最佳,过大或过小的 k 取值均会对模型预测性能造成较大的影响.

参数 λ 是微簇重要性衰减因子. 在敏感性实验中,

我们选择了0.1,0.01,0.001,0.000 1进行对比. λ 的大小表示未获取样本对微簇重要性的影响,与重要性系数的衰退成正比. λ 的取值可以反映对历史数据与当前数据对模型的重要程度. 实验结果如图10(c)所示. 与其余超参数相比, λ 取值的变化对模型的影响较小,在 $\lambda=0.01$ 时,模型的表现最佳.

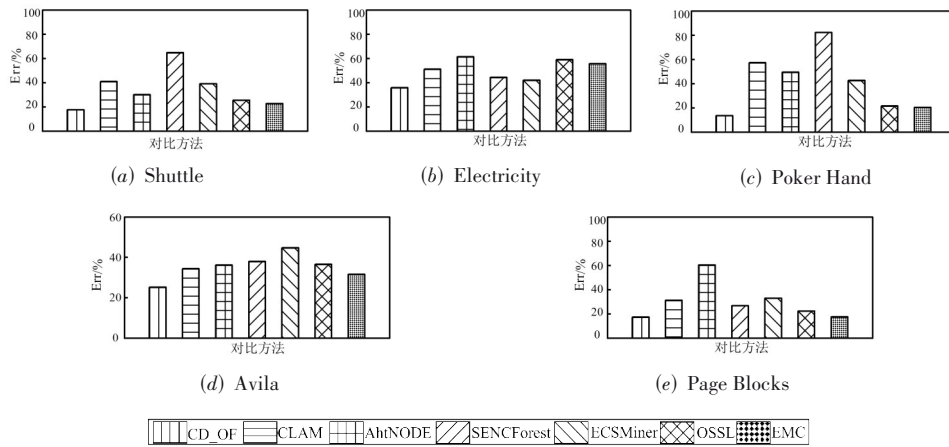


图9 不同方法在各数据集上的Err比较

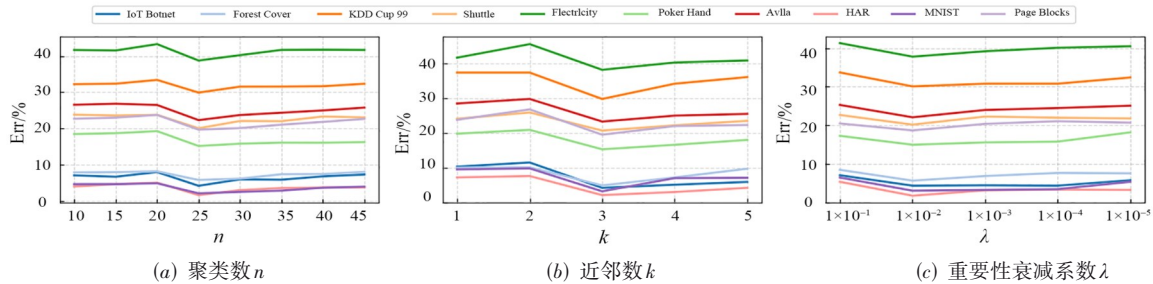


图10 参数敏感性实验

4.3.4 数据内聚性分析

新类数据的轮廓系数和 Davies-Bouldin 指数进行计算对数据内聚性表示进行定量分析,实验结果如表2所示. 虽然一些数据集的内聚性指标低于其余数据,如实验结果在 Avila 数据集的表现低于 Electricity 数据集,但是2个数据集的错误率并没有明显差别. 这主要是因为 CD_OF 在每个类别中进一步聚类得到微簇能够更精确地表示数据,降低同类样本间内聚性不够强的影响.

5 结论

针对传统基于静态非开放特征空间下的概念演化检测难以满足现实场景中特征维度变化情况下的概念演化检测与适应需求,本文提出一种面向开放特征空间的概念演化检测方法 CD_OF. 该方法首先通过将变化后的特征空间转换到共享特征空间以对齐新旧数据分布,然后通过补齐维度从而使模型快速学习补齐后的特征,有效缓解开放特征空间动态变化的严重影响.

表2 数据集内聚性测试

数据集	轮廓系数	Davies-Bouldin 指数
IoT botnet	-0.46	5.56
Forest cover	-0.09	9.69
Kddcup99	0.04	5.29
Shuttle	0.39	1.05
Electricity	0.61	0.58
Poker Hand	0.08	3.30
Avila	0.01	4.93
HAR	0.06	3.48
Mnist	0.05	3.76
Page Block	-0.23	1.89

在此基础上,通过度量特征相似性和共享近邻覆盖范围2种信息共同判断样本点是否属于新样本,提高了模型对有较强内聚性的新类样本的检测能力. 然后,通过动态衰减模型,实现了模型对类消失和类循环数据的适应能力,增强了开放特征空间下概念演化检测的泛

化能力.

本文方法对解决开放特征空间下的概念演化检测问题进行了积极的探索. 然而, 在开放特征空间中数据及特征存在更多样式的变化类型, 未来我们将对这些变化情况进行深入分析从而提升已有方法的适应能力.

参考文献

- [1] RAMZAN F, AYYAZ M. A comprehensive review on data stream mining techniques for data classification; and future trends[J]. *EPH - International Journal of Science and Engineering*, 2023, 9(3): 1-29.
- [2] LU J, LIU A J, DONG F, et al. Learning under concept drift: A review[J]. *IEEE Transactions on Knowledge and Data Engineering*, 2019, 31(12): 2346-2363.
- [3] LI J P, YU H, ZHANG Z Y, et al. Concept drift adaptation by exploiting drift type[J]. *ACM Transactions on Knowledge Discovery from Data*, 2024, 18(4): 1-22.
- [4] 杜航原, 王文剑, 白亮. 一种基于优化模型的演化数据流聚类方法[J]. *中国科学: 信息科学*, 2017, 47(11): 1464-1482.
DU H Y, WANG W J, BAI L. A novel evolving data stream clustering method based on optimization model[J]. *Scientia Sinica (Informationis)*, 2017, 47(11): 1464-1482. (in Chinese)
- [5] KSIENIEWICZ P, ZYBLEWSKI P. Stream-learn: Open-source Python library for difficult data stream batch analysis[J]. *Neurocomputing*, 2022, 478: 11-21.
- [6] CACCIARELLI D, KULAHCI M. Active learning for data streams: A survey[J]. *Machine Learning*, 2024, 113(1): 185-239.
- [7] 翟婷婷, 高阳, 朱俊武. 面向流数据分类的在线学习综述[J]. *软件学报*, 2020, 31(4): 912-931.
ZHAI T T, GAO Y, ZHU J W. Survey of online learning algorithms for streaming data classification[J]. *Journal of Software*, 2020, 31(4): 912-931. (in Chinese)
- [8] KLIKOWSKI J. Concept drift detector based on centroid distance analysis[C]//2022 International Joint Conference on Neural Networks. Piscataway: IEEE, 2022: 1-8.
- [9] LI X J, ZHOU Y, JIN Z Y, et al. A classification and novel class detection algorithm for concept drift data stream based on the cohesiveness and separation index of mahalanobis distance[J]. *Journal of Electrical and Computer Engineering*, 2020, 2020(1): 4027423.
- [10] 韩光洁, 赵腾飞, 刘立, 等. 基于多元区域集划分的工业数据流概念漂移检测[J]. *电子学报*, 2023, 51(7): 1906-1916.
HAN G J, ZHAO T F, LIU L, et al. Concept drift detection of industrial data flow based on multivariate region set partition[J]. *Acta Electronica Sinica*, 2023, 51(7): 1906-1916. (in Chinese)
- [11] 代劲, 李昊, 王国胤. 基于动态样本选择的概念漂移自适应预测方法[J]. *电子学报*, 2024, 52(9): 3228-3239.
DAI J, LI H, WANG G Y. Concept drift adaptive prediction method based on dynamic sample selection[J]. *Acta Electronica Sinica*, 2024, 52(9): 3228-3239. (in Chinese)
- [12] ZUBAROĞLU A, ATALAY V. Online embedding and clustering of evolving data streams[J]. *Statistical Analysis and Data Mining: the ASA Data Science Journal*, 2023, 16(1): 29-44.
- [13] 王婧. 基于集成学习的概念演化检测方法研究[D]. 太原: 山西大学, 2024.
WANG J. Research on Concept Evolution Detection Method Based on Ensemble Learning[D]. Taiyuan: Shanxi University, 2024. (in Chinese)
- [14] MASUD M, GAO J, KHAN L, et al. Classification and novel class detection in concept-drifting data streams under time constraints[J]. *IEEE Transactions on Knowledge and Data Engineering*, 2011, 23(6): 859-874.
- [15] 王婧, 郭虎升, 王文剑. 基于弱监督集成的概念演化自适应检测方法[J]. *吉林大学学报(信息科学版)*, 2024, 42(3): 406-420.
WANG J, GUO H S, WANG W J. Adaptive detection method for concept evolution based on weakly supervised ensemble[J]. *Journal of Jilin University (Information Science Edition)*, 2024, 42(3): 406-420. (in Chinese)
- [16] GUO H S, XIA H S, LI H, et al. Concept evolution detection based on noise reduction soft boundary[J]. *Information Sciences*, 2023, 628: 391-408.
- [17] DIN S U, SHAO J M. Exploiting evolving micro-clusters for data stream classification with emerging class detection[J]. *Information Sciences*, 2020, 507: 404-420.
- [18] GARCIA K D, DE FARIA E R, DE SÁ C R, et al. Ensemble clustering for novelty detection in data streams[C]// *Discovery Science*. Cham: Springer, 2019: 460-470.
- [19] ZHENG X L, LI P P, HU X G, et al. Semi-supervised classification on data streams with recurring concept drift and concept evolution[J]. *Knowledge-Based Systems*, 2021, 215: 106749.
- [20] AL-KHATEEB T, MASUD M M, AL-NAAMI K M, et al. Recurring and novel class detection using class-based ensemble for evolving data stream[J]. *IEEE Transactions on Knowledge and Data Engineering*, 2016, 28(10): 2752-2764.
- [21] GAO Y, CHANDRA S, LI Y F, et al. SACCOS: A semi-

- supervised framework for emerging class detection and concept drift adaption over data streams[J]. IEEE Transactions on Knowledge and Data Engineering, 2022, 34(3): 1416-1426.
- [22] BOUGUELIA M R, NOWACZYK S, PAYBERAH A H. An adaptive algorithm for anomaly and novelty detection in evolving data streams[J]. Data Mining and Knowledge Discovery, 2018, 32(6): 1597-1633.
- [23] MU X, TING K M, ZHOU Z H. Classification under streaming emerging new classes: A solution using completely-random trees[J]. IEEE Transactions on Knowledge and Data Engineering, 2017, 29(8): 1605-1618.
- [24] GANDHI J, GANDHI V. Novel class detection with concept drift in data stream - AhtNODE[J]. International Journal of Distributed Systems and Technologies, 2020, 11(1): 15-26.
- [25] 赵鹏, 周志华. 基于决策树模型重用的分布变化流数据学习[J]. 中国科学: 信息科学, 2021, 51(1): 1-12.
ZHAO P, ZHOU Z H. Learning from distribution-changing data streams via decision tree model reuse[J]. Scientia Sinica (Informationis), 2021, 51(1): 1-12. (in Chinese)
- [26] ZHANG Z L, LI Y, ZHANG Z W, et al. Adaptive matrix sketching and clustering for semisupervised incremental learning[J]. IEEE Signal Processing Letters, 2018, 25(7): 1069-1073.
- [27] GUAN S U, LI S C. Incremental learning with respect to new incoming input attributes[J]. Neural Processing Letters, 2001, 14(3): 241-260.
- [28] HOU B J, ZHANG L J, ZHOU Z H. Learning with feature evolvable streams[J]. IEEE Transactions on Knowledge and Data Engineering, 2021, 33(6): 2602-2615.
- [29] 刘兆清, 古仕林, 侯臣平. 面向特征继承性增减的在线分类算法[J]. 计算机研究与发展, 2022, 59(8): 1668-1682.
LIU Z Q, GU S L, HOU C P. Online classification algorithm with feature inheritably increasing and decreasing[J]. Journal of Computer Research and Development, 2022, 59(8): 1668-1682. (in Chinese)
- [30] HE Y, WU B J, WU D, et al. Online learning from capricious data streams: A generative approach[C]//Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence. International Joint Conferences on Artificial Intelligence Organization, 2019: 2491-2497.
- [31] LIAO G B, ZHANG P, YIN H P, et al. A novel semi-supervised classification approach for evolving data streams[J]. Expert Systems with Applications, 2023, 215: 119273.
- [32] LIU R, WANG H, YU X M. Shared-nearest-neighbor-based clustering by fast search and find of density peaks[J]. Information Sciences, 2018, 450: 200-226.
- [33] UD DIN S, SHAO J M, KUMAR J, et al. Online reliable semi-supervised learning on evolving data streams[J]. Information Sciences, 2020, 525: 153-171.
- [34] DE STEFANO C, MANIACI M, FONTANELLA F, et al. Reliable writer identification in medieval manuscripts through page layout features: The “Avila” Bible case[J]. Engineering Applications of Artificial Intelligence, 2018, 72: 99-110.
- [35] ZAREMOODI P, KAMALI SIAHROUDI S, BEIGY H. Concept-evolution detection in non-stationary data streams: A fuzzy clustering approach[J]. Knowledge and Information Systems, 2019, 60(3): 1329-1352.

作者简介



苏睿男, 1999年1月出生于山西省忻州市. 现为山西大学计算机与信息技术学院博士研究生. 主要研究方向为流数据挖掘、机器学习、深度学习.

E-mail: 202212407015@email.sxu.edu.cn



王婧女, 1999年8月出生于山西省吕梁市. 现为山西大学计算机与信息技术学院硕士研究生. 主要研究方向为流数据挖掘.

E-mail: 1756645158@qq.com



郭虎升男, 1986年10月出生于山西省晋中市. 现为山西大学教授、博士生导师. 主要研究方向为数据挖掘、机器学习.

E-mail: guohusheng@sxu.edu.cn



王文娟女, 1968年10月出生于山西省太原市. 现为山西大学计算智能与中文信息处理教育部重点实验室副主任、教授、博士生导师. 主要研究方向为机器学习、数据挖掘.

E-mail: wjwang@sxu.edu.cn