

## 基于Weighted Majority的集成在线学习算法

**摘要:** 数据流挖掘问题广泛存在于物流配送、采购管理、在线广告、量化投资、众包数据分析等实际应用场景。深度神经网络、支持向量机等传统基于批处理模式的学习模型在处理此类问题时存在更新效率低、实时性差等诸多不足。在线学习算法具有计算复杂度低、模型更新快的优点,是处理数据流挖掘问题的理想工具,已逐步成为大规模数据流挖掘的新范式。为有效解决数据流挖掘问题中普遍存在的人力成本控制和算法选优问题,本文提出采用影响函数法筛选少量数据通过咨询领域专家获取可靠样本标签,以有效控制人力成本;基于可靠学习样本,选取常用在线学习算法作为基学习算法,采用 Weighted Majority 算法框架对基学习算法进行集成,根据在线学习算法的性能表现动态调整对其信任度,并依赖信任度给出预测;最后,通过基准数据集以及垃圾邮件识别问题上的数值实验验证了所提算法的效率和有效性。

**关键词:** 数据流挖掘; 在线学习; 影响函数法; Weighted Majority 算法

中图分类号: TP381

文献标志码: A

DOI: 10.13195/j.kzyj.2019.0000

开放科学(资源服务)标识码(OSID):



## Ensemble Online Learning Algorithm based on Weighted Majority

**Abstract:** Data stream mining widely exists in practical applications such as online advertising, quantitative investment and crowdsourcing data analysis. For these tasks, the existing methods such as deep neural networks and support vector machines, which are based on batch learning scheme in general, have non-satisfactory performances in terms of efficiency and timeliness. As a new learning paradigm, online learning has become an ideal vehicle for large-scale data stream mining due to low computational complexity. To address the ubiquitous issues of labor burden as well as algorithm selection in data stream mining tasks, in this paper, we propose Ensemble Online Learning algorithm with the following characteristics: (1) This algorithm only needs to fetch the labels of partial observations that are determined by influence function; (2) Utilizing Weighted Majority method, this algorithm aggregates base learning algorithms effectively through the trust degree that is adjusted dynamically. Extensive experiments are conducted on benchmark datasets and real-world dataset to verify the effectiveness and efficiency of the proposed algorithm.

**Keywords:** data stream mining; online learning; influence function method; Weighted Majority algorithm

### 0 引言

伴随大数据时代<sup>[1]</sup>的来临,文本<sup>[2]</sup>、图像<sup>[3]</sup>、语音<sup>[4]</sup>等非结构化类型数据层出不穷,数据呈现出容量越来越大、类型越来越多的特征,这给传统数据分析方式带来了巨大挑战<sup>[5]</sup>。统计机器学习<sup>[6]</sup>集良好的统计理论性质和高效求解算法于一身,已成为现代数据分析不可或缺的重要工具。然而,供应链、电子商务、社交媒体、移动互联网、物联网、复杂工业生产等实际应用场景不断产生携带系统动态信息的数据,这些数据具有明显的流式特征<sup>[7]</sup>。物流配送、社交媒体、工业过程等复杂系统具有动态、时变、概念漂移的本质特性,而采样数据流中隐含着复杂系统的动态变化等重要信息,从数据流中快速挖掘动态信息对于决策制定至关重要<sup>[8]</sup>。基于批处理方式的机器学习算法在处理数据流问题时存在计算复杂度高、模型更新效率低等不

足。基于在线凸优化理论 Zinkevich 建立了在线学习框架<sup>[9]</sup>,在线学习基于逐点学习策略,通过动态更新模型逐点学习样本数据,模型更新一次的计算复杂度仅为 $O(1)$ ,具有计算复杂度低、模型更新效率高、实时性强等优势,是处理数据流分析问题的天然工具。近年来,逐渐成为机器学习领域的研究热点,是数据流挖掘的主流方法<sup>[10]</sup>。

近年来,国内外学者提出了一系列在线学习算法包括: Passive-Aggressive<sup>[11]</sup>、Kernel Perceptron<sup>[12]</sup>、Pegasos<sup>[13]</sup>、Online SVM<sup>[14]</sup>、NORMA<sup>[15]</sup>等,这些在线学习算法各有优劣,需要根据实际应用问题对应数据的分布规律等先验信息进行优选。然而,由于实际应用问题的复杂性,其数据分布规律难以提取,故存在算法选择的难题。此外,高质量学习样本是实现高效数据流挖掘目标的核心。然而,获取高质量学习样本标签

收稿日期: 2018-xx-xx; 修回日期: 2018-xx-xx.

责任编辑: 编委1.

<sup>†</sup>通讯作者. E-mail: zuozheyi@163.com.

需要消耗大量人力资源。例如在垃圾邮件分类问题中,获取邮件的真实标签往往需要对邮件进行人工检查。在有效控制人力成本前提下,尽可能多的获取典型学习样本成为数据流挖掘领域亟待解决的又一关键性难题。鉴于此,本文提出了基于Weighted Majority(WM)算法框架的集成在线学习算法,主要贡献包括:

1) 利用逻辑回归模型的影响函数,高效筛选部分样本,通过咨询领域专家获取可信标签,进而获取典型学习样本,实现人力成本的有效控制;

2) 采用WM算法框架,实现对常用在线学习算法的有效集成,并分析所提集成在线学习算法的理论性质,解决了算法优选的难题;

3) 通过在MNIST等基准数据集和垃圾邮件分类数据上的应用验证,说明所提算法的高效性和有效性。

## 1 相关工作

本文的研究工作主要涉及在线学习、集成学习、影响函数三个方面,本节将着重回顾国内外在上述三个方面的相关研究。

### 1.1 在线学习算法

近年来,国内外学者围绕在线学习算法与理论开展了大量研究工作,提出了一系列实用算法并广泛应用于机器学习和数据挖掘等诸多领域,包括:Dekel等研究了预算核感知机模型并提出Forgetron算法<sup>[16]</sup>; Crammer等基于随机梯度下降法研究了预算SVM模型<sup>[14]</sup>; Zhao等基于在线梯度下降算法提出了可尺度化在线核学习模型<sup>[17]</sup>; Hazan等人基于损失函数的二阶信息,提出了在线牛顿算法,证明在一定条件下该算法的regret上界仅为 $\mathcal{O}(\log n)$ <sup>[18]</sup>; Jian等基于核技巧提出了软间隔在线二分类算法,并分析了算法的渐近收敛性<sup>[19]</sup>。

### 1.2 集成学习算法

集成学习的目标是将基学习算法整合为一个元算法,以期有效解决单一算法泛化能力差的不足。Freund和Schapire<sup>[20]</sup>提出的AdaBoost(Adaptive Boosting)是机器学习与数据挖掘领域最为常用的集成学习算法框架。AdaBoost通过对弱分类算法进行集成以提高其性能表现。理论上,只要基学习算法的性能表现优于随机猜测,可以证明最终的集成学习算法可收敛为一个强学习算法<sup>[20]</sup>。AdaBoost具有自适应性,可以对弱分类算法进行调整,以支持被先前分类算法错分的样本<sup>[20]</sup>。但AdaBoost算法框架是基于批处理模式提出的,在处理数据流问题

时,存在模型计算复杂度过高的缺点。

为了解决基于专家意见的预测问题,基于在线凸优化理论Littlestone等提出了Weighted Majority算法<sup>[21]</sup>。该算法在在线学习的基本框架下根据专家在一系列任务上的性能表现动态优化对其信任度并依赖于信任度选择最受信赖的专家意见作为算法预测,理论上可以证明该算法是low regret(低反悔)算法,其性能表现可以收敛到最优专家的性能表现。

鉴于Weighted Majority算法良好的理论性质,为实现数据流挖掘场景下的集成学习问题,本文将采用Weighted Majority算法框架对基学习算法进行集成,并根据基学习算法的性能动态调整对其信任度。

### 1.3 影响函数法

影响函数方法是来自稳健统计的经典技术<sup>[22]</sup>,该函数刻画的是增加训练点权重对学习模型参数的影响程度。尽管在统计学领域已被广泛使用,但影响函数尚未在机器学习和数据挖掘领域得到广泛应用。比较有影响力的研究主要有Wojnowicz等<sup>[23]</sup>以及Percy Liang等<sup>[24]</sup>的工作。Wojnowicz等提出了一种在广义线性模型中近似影响量的方法,该方法需通过二阶导数来计算影响函数,要求模型具有可微性和凸性,而深度学习等流行机器学习模型通常不具有可微性、凸性,导致该方法在机器学习与数据挖掘领域难以广泛应用。为此,Percy Liang等借助二阶优化技术有效近似影响函数以解决上述挑战,并成功将影响函数法推广应用到深度学习等机器学习模型中,取得了良好应用效果。

## 2 基于WM的集成在线学习算法

本节将首先引入影响函数法,结合经典监督学习方法,基于成本控制原则实现可信样本标签的获取。其次,选择常用在线学习算法作为基在线学习算法,并基于WM算法框架对基在线学习进行动态集成。最后,完成所提算法的理论分析,包括regret界分析和计算复杂度分析。

### 2.1 构建逻辑回归模型的影响函数

本节选择在分类问题中广泛使用的经典监督学习方法逻辑回归模型作为基础模型。逻辑回归与线性回归联系紧密,线性回归模型 $t = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n + \varepsilon$ 中 $x_i$ 表示模型的输入变量, $\theta_i$ 是模型参数。为应用线性回归的思想处理分类问题,逻辑回归模型引入Sigmoid函数 $\sigma(t) = \frac{1}{1+e^{-t}}$ ,将线性模型的输出值域从 $(-\infty, +\infty)$ 映射到 $(0,1)$ 。

以 $-1/1$ 二分类问题为例, 模型首先计算出 $\sigma(t)$ , 然后以 $0.5$  为分界线, 如果 $\sigma(t)$ 大于 $0.5$ , 则预测分类结果为 $1$ , 反之 $-1$ , 由此实现了连续型标签到离散二值标签的转换。

假设训练数据集为 $S = \{z_1, z_2, \dots, z_n\}$ (其中 $z_i = (\mathbf{x}_i, y_i) \in R^d \times Y, Y = \{-1, 1\}$ ), 在 $S$ 上学习逻辑回归模型 $\mathcal{L}(S)$ 的过程为: 记逻辑回归模型的判别函数为 $\sigma(t) = h_\theta(\mathbf{x}) = \frac{1}{1+e^{-y\theta^T \mathbf{x}}}$ , 其中 $t = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ 可向量化表示为 $t = \theta^T \mathbf{x}, \theta^T = [\theta_0, \theta_1, \dots, \theta_n]^T$ 为模型参数,  $\mathbf{x} = [x_0, x_1, \dots, x_n]^T$  ( $x_0 = 1$ )为输入变量。通常采用极大似然估计方法对逻辑回归模型的最优模型参数 $\hat{\theta}$ 进行估计, 将极大似然函数做对数处理可得逻辑回归模型的损失函数为:  $\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \theta^T \mathbf{x}_i))$ 。简便起见, 记损失函数为 $L(z, \theta) = \log(1 + \exp(-y\theta^T \mathbf{x}))$ , 则损失函数 $L(z, \theta)$ 关于参数 $\theta$ 二阶连续可微, 可求一阶、二阶导数分别为

$$\nabla_\theta L(z, \theta) = -\sigma(-y\theta^T \mathbf{x}) y \mathbf{x} \quad (1)$$

$$H_\theta = \frac{1}{n} \sum_{i=1}^n \sigma(\theta^T \mathbf{x}_i) \sigma(-\theta^T \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T \quad (2)$$

记根据经验风险最小化准则确定的最优模型参数为:  $\hat{\theta} \stackrel{\text{def}}{=} \underset{(\theta \in \Theta)}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n L(z_i, \theta)$ , 在训练集中删除样本点 $z$ 时模型参数将发生改变, 记新的模型参数为 $\hat{\theta}_{-z} \stackrel{\text{def}}{=} \underset{(\theta \in \Theta)}{\operatorname{argmin}} \frac{1}{n} \sum_{z_i \neq z} L(z_i, \theta)$ 。记在样本点 $z$ 的损失上施加一个小的权重 $\varepsilon$ 时, 所得新的模型参数为 $\hat{\theta}_{\varepsilon, z} \stackrel{\text{def}}{=} \underset{(\theta \in \Theta)}{\operatorname{argmin}} \frac{1}{n} \sum_{z_i \neq z} L(z_i, \theta) + \varepsilon L(z, \theta)$ , 则模型参数受到学习样本 $z$ 的影响, 可用下述影响函数来刻画 $I_{up, \text{params}}(z) \stackrel{\text{def}}{=} \left. \frac{d\hat{\theta}_{\varepsilon, z}}{d\varepsilon} \right|_{\varepsilon=0} = -H_{\hat{\theta}}^{-1} \nabla_\theta L(z, \hat{\theta})$ 。进一步的, 训练样本点的改变亦会对测试样本点的预测结果产生影响, 对应的影响函数的表示形式为 $I_{up, \text{loss}}(z, z_{\text{test}}) \stackrel{\text{def}}{=} -\nabla_\theta L(z_{\text{test}}, \hat{\theta})^T H_{\hat{\theta}}^{-1} \nabla_\theta L(z, \hat{\theta})$ , 结合公式(1)和公式(2)可得逻辑回归模型的影响函数计算公式为

$$I_{up, \text{loss}}(z, z_{\text{test}}) = -y_{\text{test}} y \cdot \sigma(-y_{\text{test}} \theta^T \mathbf{x}_{\text{test}}) \cdot \sigma(-y \theta^T \mathbf{x}) \cdot \mathbf{x}_{\text{test}}^T H_{\hat{\theta}}^{-1} \mathbf{x} \quad (3)$$

## 2.2 构建基于WM算法框架的集成在线学习算法

假设选取 $k$ 个基在线学习算法作为集成分类器, 对于第 $t$ 个任务, 记 $\mathbf{x}_t$ 为该任务的输入特征,  $y_t$ 为其真实标签,  $\tilde{y}_t = (\tilde{y}_t[1], \tilde{y}_t[2], \dots, \tilde{y}_t[k])$ 是基分类器预测标签,  $p_t \in \{1, 2, \dots, k\}$ 为集成学习算法在处

理第 $t$ 个任务时所选出的基学习算法序号,  $\tilde{z}_t[i] = \text{loss}(y_t, \tilde{y}_t[i])$ 为 $p_t = i$ 时的损失,  $w_t$ 为模型权重参数。在这个预测问题中, 可行集是离散的, 因而优化模型损失的问题本质上是凸问题。非凸优化问题的算法设计和理论分析均是难题。为此, 我们通过随机化预测 $p_t$ 的方式将该优化问题转化为凸优化问题。令 $\mathcal{P} = \{w \in R^k : w \geq 0 \wedge \|w\|_1 = 1\}$ 为可行集, 显然该可行集是个单纯形, 故是凸集合。在处理第 $t$ 个任务时, 集成学习算法基于多项式概率分布 $w_t \in \mathcal{P}$ , 以概率 $\mathbb{P}[p_t = i] = w_t[i]$ 随机选取一个基在线学习算法, 并采用该基在线学习算法的预测标签对当前输入特征进行预测。这里 $w_t[i]$ 刻画的是算法在完成第 $t$ 个任务时对第 $[i]$ 个基在线学习算法的信任度。因此, 集成在线学习算法在完成第 $t$ 个任务时的期望损失为:  $\mathbb{E}[\tilde{z}_t[p_t]] = \sum_{i=1}^k \mathbb{P}[p_t = i] \tilde{z}_t[i] = \langle w_t, \tilde{z}_t[i] \rangle$ , 为保证基础学习算法在迭代过程中保持稳定, 采用Follow-the-Regularied-Leader策略动态优化累计期望损失, 由此可得优化问题为

$$w_{t+1} = \arg \min_w \sum_{j=1}^t \langle w, \tilde{z}_j \rangle + R(w) \quad (4)$$

其中,  $R(w)$ 表示算法选择的正则化算子。由于可行集 $\mathcal{P}$ 为凸集, 且损失函数 $\ell_t(w) = \langle w, \tilde{z}_t \rangle$ 为线性函数(凸损失函数), 正则化算子 $R(w)$ 为凸函数, 故优化模型期望损失的问题公式(4)是一个凸优化问题。WM算法是处理该类型问题的有效算法, 具有坚实的理论保障。故本项目拟采用WM算法框架, 结合影响函数法, 开发集成在线学习算法如下:

- 1: 初始化: 可信样本集 $S$ , 学习率参数 $\eta$ ,  $w_1 = [\frac{1}{d}, \dots, \frac{1}{d}]^T$ , 校正阈值参数 $\delta$
- 2: **for**  $t = 1, 2, \dots$  **do**
- 3:   收到第 $t$ 个任务 $\mathbf{x}_t \in \mathbb{R}^d$ ;
- 4:   第 $i$ 个基学习器作出预测 $\tilde{y}_t[i], i = 1, \dots, d$ ;
- 5:   根据信任度 $w_t$ 产生随机指标 $p_t$ ;
- 6:   集成学习器预测标签为 $\tilde{y}_t[p_t]$ ;
- 7:   记 $z_t \triangleq (\mathbf{x}_t, \tilde{y}_t[p_t])$ ;
- 8:   **if**  $-I_{up, \text{loss}}(z_t, z_t) > \delta$  **then**
- 9:     咨询领域专家获得可信标签 $y_t$ ;
- 10:    更新可信样本集 $S := S \cup \{(\mathbf{x}_t, y_t)\}$ ;
- 11:    更新基础模型 $\mathcal{L}(S)$ 的参数 $\theta$ ;
- 12:    更新基学习器损失 $\tilde{z}_t[i] = \text{loss}(y_t, \tilde{y}_t[i])$ ;
- 13:    更新信任度 $w_{t+1}[i] = \frac{w_t[i] e^{-\eta \tilde{z}_t[i]}}{\sum_i w_t[i] e^{-\eta \tilde{z}_t[i]}}$ ,  $i = 1, \dots, d$ ;
- 14:   **else**
- 15:     模型参数 $w_{t+1} = w_t$

16: **end if**

17: **end for**

图1 A: 基于WM算法框架的集成在线学习算法

在执行指定数据流挖掘任务之前, 咨询领域专家获得少量真实样本标签, 这部分样本构成可信训练集 $S$ , 选择经典监督学习方法逻辑回归模型作为基础模型, 在 $S$ 上学习基础模型 $\mathcal{L}(S)$ , 得到基础模型参数记为 $\theta$ 。构造影响函数用来衡量采样数据对基础模型的影响程度, 若某样本点的影响函数值超出给定的阈值 $\delta$ , 则通过咨询专家的方式对此样本点进行人工数据校正, 并使用校正后的真实标签 $y_t$ 对模型参数 $\theta$ 进行调整。这里阈值参数 $\delta$ 的设置有待深入探讨。基于WM的集成在线学习算法主要包括在线预测与数据矫正两大部分, 工作流程如图1所示。

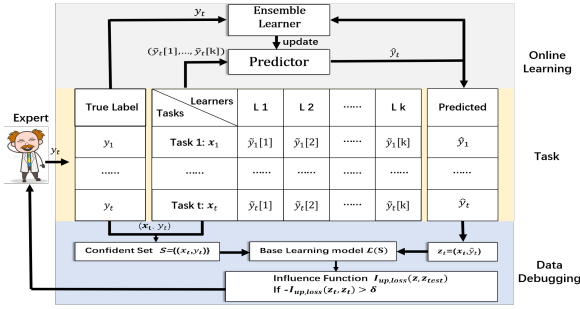


图2 基于WM的集成在线学习算法示意图

### 2.3 算法的理论性质

**引理1** 假设损失函数 $\ell(w)$ 为线性损失, 运行WM算法得到一系列预测 $w_1, \dots, w_T \in \mathcal{P}$ ,  $\mathcal{P} = \{w \in R^k : w \geq 0 \wedge \|w\|_1 = 1\}$ , 若WM算法的学习率参数 $\eta$ 满足条件 $\eta \tilde{z}_t[i] \geq -1$ , 则有

$$\sum_{t=1}^T \langle w_t, \tilde{z}_t \rangle - \sum_{t=1}^T \langle u, \tilde{z}_t \rangle \leq \frac{\log k}{\eta} + \eta \sum_{t=1}^T \sum_{i=1}^k w_t[i] \tilde{z}_t[i]^2$$

其中,  $u$ 为可行集 $\mathcal{P}$ 中任意变量。

证明: 详细证明过程可参考文献[25]定理2.22。

**定理1** 基于WM算法框架的集成在线学习算法的期望损失反悔界为 $\frac{\log k}{\eta} + \eta T$ , 即

$$\mathbb{E} \sum_{t=1}^T [\tilde{z}_t[p_t]] - \min_{1 \leq i \leq k} \sum_{t=1}^T \tilde{z}_t[i] = \sum_{t=1}^T \langle w_t, \tilde{z}_t \rangle - \min_{1 \leq i \leq k} \sum_{t=1}^T \tilde{z}_t[i] \leq \frac{\log k}{\eta} + \eta T$$

在 $\tilde{z}_t$ 有界的条件下( $|\tilde{z}_t[i]| \leq L, i = 1, \dots, k$ ), 令 $\eta = \frac{1}{L} \sqrt{\log k / T}$ 可得所提算法A的反悔界(关于最优在线学习算法)为

$$\text{regret}_T(\mathcal{A}) = \sum_{t=1}^T \langle w_t, \tilde{z}_t \rangle - \min_{1 \leq i \leq k} \sum_{t=1}^T \tilde{z}_t[i] \leq$$

$$2L \sqrt{(\log k)T}$$

证明: 在本文所提基于WM的在线学习算法中, 第 $i$ 个在线学习算法在完成第 $t$ 轮任务时的损失为 $\tilde{z}_t[i] = 0$ 或 $1, i = 1, \dots, k$ , 而学习率参数 $\eta > 0$ , 因此引理1中条件 $\eta \tilde{z}_t[i] \geq -1$ 显然满足。本文所提算法A在第 $t$ 轮的期望损失为 $\mathbb{E} [\tilde{z}_t[p_t]] = \langle w_t, \tilde{z}_t \rangle$ , 其中 $w_t \in \mathcal{P}$ ,  $\mathcal{P} = \{w \in R^k : w \geq 0 \wedge \|w\|_1 = 1\}$ 表示多项式分布, 对于可行集 $\mathcal{P}$ 中的任意变量 $u$ , 由引理1可得

$$\sum_{t=1}^T \langle w_t, \tilde{z}_t \rangle - \sum_{t=1}^T \langle u, \tilde{z}_t \rangle \leq \frac{\log k}{\eta} + \eta \sum_{t=1}^T \sum_{i=1}^k w_t[i] \tilde{z}_t[i]^2 \quad (5)$$

又公式(5)对 $\forall u \in \mathcal{P}$ 均成立, 特别地对于单纯形 $\mathcal{P}$ 的顶点即 $u = (0, \dots, 1, \dots, 0)$ 亦成立。由此可得

$$\sum_{t=1}^T \langle w_t, \tilde{z}_t \rangle - \min_{1 \leq i \leq k} \sum_{t=1}^T \tilde{z}_t[i] \leq \frac{\log k}{\eta} + \eta \sum_{t=1}^T \sum_{i=1}^k w_t[i] \tilde{z}_t[i]^2$$

由于本文采用经典0-1损失, 故 $L = 1$ , 进而可得 $|\tilde{z}_t[i]| \leq 1, i = 1, \dots, k$ 则有

$$\sum_{t=1}^T \langle w_t, \tilde{z}_t \rangle - \min_{1 \leq i \leq k} \sum_{t=1}^T \tilde{z}_t[i] \leq \frac{\log k}{\eta} + \eta T$$

在设置学习率参数 $\eta = \sqrt{\log k / T}$ 的条件下, 可得所提算法A的反悔界为 $2\sqrt{(\log k)T}$ 。

### 2.4 计算复杂度

本文所提算法A在执行在线预测任务的过程中, 若在执行第 $t$ 个任务时计算所得影响函数值小于给定阈值即 $-I_{up, loss}(z_t, z_t) \leq \delta$ 成立, 无需咨询专家意见亦无需对基在线学习算法和基础模型进行更新; 相反, 当第 $t$ 个样本的影响函数值大于给定阈值即 $-I_{up, loss}(z_t, z_t) \geq \delta$ 时, 需要咨询专家意见得到可信标签, 并对基在线学习算法和基础模型进行更新, 由于基在线学习算法在执行模型更新时的计算量均为 $\mathcal{O}(1)$ , 此时算法A的计算量主要集中在计算基础模型的海瑟矩阵, 该环节的计算量为 $\mathcal{O}((\gamma nd)^2)$ , 其中 $n$ 表示学习样本总量,  $\gamma$ 表示执行专家咨询的样本点占比,  $d$ 表示输入样本的特征维数。由此可见, 本文所提出的基于影响函数法的集成在线学习算法的计算量主要由专家咨询次数 $\gamma n$ 以及样本输入的特征维度 $d$ 决定。

### 3 数值实验

本节通过数值实验对所提基于WM的集成在线学习算法的有效性和效率进行验证。实验所使用

的软硬件环境具体配置为处理器Intel(R) Core(TM) i5-5250U CPU @ 1.60GHz内存为8GB, 操作系统为Windows 10, 编程语言为Python3.6.5, 集成开发环境为Spyder3.2。

### 3.1 数据集介绍

首先, 本节在3个基准数据集和1个垃圾邮件分类数据集上验证所提算法的有效性。3个基准数据集包括2个常见疾病数据集和1个手写数字数据集, 均通过scikit-learn库导入。

表1列出了3个基准数据集和垃圾邮件数据集的详细信息, 为了模拟数据流挖掘的真实场景, 我们设置3个基准数据集学习的轮次均为29次, 表1中P/N列描述了数据集中正负两类样本点的比例。

表1 基准数据集数据特征

Dataset	特征维度	样本量	轮次	P/N
Breast cancer	30	569	29	1.684
digits3,5	64	365	29	0.995
diabetes	10	442	29	1.415
spambase	57	4601	1	0.650

### 3.2 基在线学习算法

本节选取了5个常用的在线线性分类器包括: Perceptron、Passive Aggressive算法、在线线性SVM、在线逻辑回归、在线BernoulliNB算法, 各在线学习算法的参数设置采用scikit-learn中的默认设置。

### 3.3 算法评价指标

本节选择两个在分类问题中广泛使用的评价指标即错误率和F1-score对所提算法进行评价。F1-score指标是precision和recall的调和平均, 计算公式为 $F1 = \frac{2TP}{2TP+FP+FN}$ 。

根据定理1, 为达到最优学习效果, 理论上的最优学习率参数 $\eta$ 应为 $\sqrt{\log k/T}$ 但在处理具体的实际应用问题时, 通常情况下任务总量 $T$ 无法事先预知, 通常采用Doubling trick 或者预估的方式确定学习率参数 $\eta$ 。本实验中我们设置学习率参数为 $\eta = 0.05$ 。

### 3.4 基准数据集上的表现

为了清晰表现出信任度 $w_t$ 的收敛性变化趋势, 本实验采用了机器学习领域广泛使用的多轮学习策略对数据集进行多次重复的方式将各数据集的样本总量扩大到10000以上。图2给出了学习率参数 $\eta = 0.01$ 情况下信任度 $w_t$ 的趋势曲线。由图2可见breast cancer数据集上Passive Aggressive算法的性能最佳, 而集成学习算法对该算法的信任度具有明显收敛到1的趋势, 对其余学习算法的信任度均收敛到0;

在digits 3 vs. 5上除Bernoulli之外其余4种在线学习算法的性能相当, 集成学习算法对Bernoulli算法的信任度收敛到0; 在diabetes数据集上, Perceptron 算法的性能最佳, 集成学习算法对该算法的信任度收敛到1。以上观察意味着所提集成在线学习算法能收敛到最优基在线学习算法, 进而保证了集成在线学习算法的有效性。

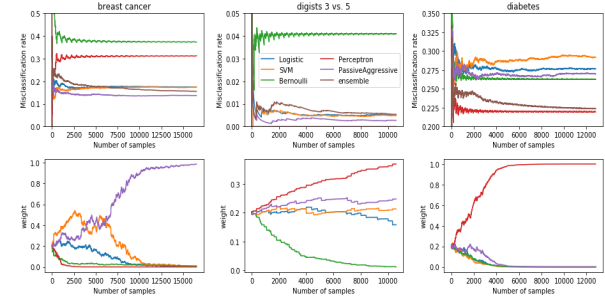


图3 权重动态变化与平均在线分类表现图

图2绘制了在线平均错误率关于轮次 $t$ 的函数曲线, 刻画了所提集成学习算法与基在线学习算法的平均在线性能表现, 由图可见所提算法A通过动态调整对不同基在线学习算法的信任度, 其性能表现能逼近最优基在线学习算法, 表明在缺少先验信息(处理具体应用问题的最佳学习算法未知)的情况下所提算法的有效性。

在本实验中影响函数阈值参数设置为 $\delta = 0.01$ , 在breast cancer、digits 3 vs. 5、diabetes 3个数据集上咨询专家的样本比例 $\gamma$ 分别为0.27、0.20、0.39, 所提集成在线学习算法的总体运行时间分别为794s、683s、253s。实际运行时间与2.4节计算复杂度分析契合。未来, 我们将引入数值优化技术进一步降低更新样本信息时计算基础模型海瑟矩阵的计算量, 以期有效降低所提算法的计算复杂度, 提高解决大规模实际应用问题的实时性。

### 3.5 spambase数据集上的表现

本节主要考察所提集成在线学习算法在垃圾邮件分类问题上的表现。由图3中的权重子图清晰看见Bernoulli算法的信任度在大约第700个样本点附近已近几乎收敛到1。上述观察说明所提算法A通过动态调整对不同基在线学习算法信任度, 其性能表现能快速逼近最优基在线学习算法, 是处理实际数据流挖掘问题可靠算法。件分类数据集上的性能表现, 选取F1-score作为模型性能评价指标。图3绘制了在线平均F1-score关于轮次 $t$ 的函数曲线, 刻画了所提集成学习算法与基在线学习算法的平均在线性能表现, 由图3中的性能表现子



图可见Bernoulli算法的表现最优，而所提集成学习算法可以较好的逼近Bernoulli算法；由图3中的权重子图清晰看见Bernoulli算法的信任度在大约第700个样本点附近已近几乎收敛到1。上述观察说明所提算法A通过动态调整对不同基在线学习算法信任度，其性能表现能快速逼近最优基在线学习算法，是处理实际数据流挖掘问题可靠算法。

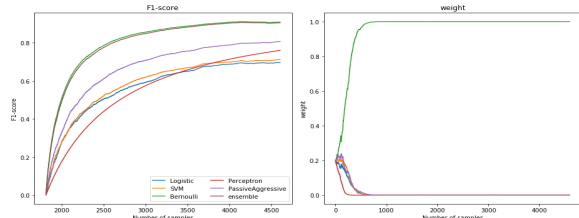


图4 spambase数据集上的权重动态变化与平均在线分类表现图

## 4 结论

本文综合运用影响函数法、Weighted Majority算法，基于成本控制原则提出了一种集成在线学习算法，根据各基学习算法在处理系列任务时的动态表现不断调整对其信任度，并使用信任度作为概率分布，基于概率分布随机选择基学习算法进行标签预测。在基准数据集和垃圾邮件分类问题上的数值实验显示，所提算法的平均在线精度可逼近最优基学习算法，上述实验结果表明了所提算法的有效性，可为数据流挖掘问题的有效解决提供方法支撑。

## 参考文献(References)

- [1] Lohr S. The age of big data [N]. New York Times, 2012, 11.
- [2] Halevy A, Norvig P, Pereira F. The Unreasonable Effectiveness of Data [J]. IEEE Intelligent Systems, 2009, 24(2):8-12.
- [3] Deng J, Dong W, Socher R, et al. Imagenet: A large-scale hierarchical image database [C]//2009 IEEE Conference on Computer Vision and Pattern Recognition. 2009:248-255.
- [4] Warden P. Speech commands: A dataset for limited-vocabulary speech recognition [J]. arXiv preprint arXiv:1804.03209, 2018.
- [5] Yu B, Kumbier K. Three principles of data science: predictability, computability, and stability (PCS) [J]. arXiv preprint arXiv:1901.08152, 2019.
- [6] Sugiyama M. Introduction to statistical machine learning [M]. Morgan Kaufmann, 2015.
- [7] Srinivasan A, Bain M. An empirical study of on-line models for relational data streams [J]. Machine Learning, 2017, 106(2):243-276.
- [8] Wu X, Zhu X, Wu G, Ding W. Data mining with big data [J]. IEEE Transactions on Knowledge and Data Engineering, 2014, 26(1):97-107.
- [9] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent [C]// In Proceedings of the Twentieth International Conference on Machine Learning, 2003:928-936.
- [10] Hoi S. C. H, Sahoo D, Lu J, Zhao P. Online learning: a comprehensive survey [J]. arXiv:1802.02871, 2018.
- [11] Crammer K, Dekel O, Keshet J, Shalev-Shwartz S, Singer Y. Online passive-aggressive algorithms [J]. Journal Machine Learning Research, 2006,7(1):551-585.
- [12] Crammer K, Singer Y. Ultraconservative online algorithms for multiclass problems [J]. Journal Machine Learning Research, 2003, 3(1):951-991.
- [13] Shalev-Shwartz S, Singer Y, Srebro N, Cotter A. Pegasos: primal estimated sub-gradient solver for SVM [J]. Mathematical Programming, 2011, 127(1):3-30.
- [14] Z. Wang, K. Crammer, S. Vucetic. Breaking the curse of kernelization: budgeted stochastic gradient descent for large-scale SVM training [J]. Journal of Machine Learning Research, 2012, 13(1):3103-3131.
- [15] Kivinen J, Smola A. J, Williamson R. C. Online learning with kernels [J]. IEEE Transactions on Signal Processing, 2004, 52(8):2165-76.
- [16] Dekel, S. Shalev-Shwartz, Y. Singer. The Forgetron: A kernel based perceptron on a budget [J]. SIAM Journal on Computing, 2006, 37(1):1342-1372.
- [17] P. Zhao, J. Wang, P. Wu, R. Jin, S. C. H. Hoi. Fast bounded online gradient descent algorithms for scalable kernel based online learning [C]//In Proceedings of the 29th International Conference on Machine Learning, 2012:169-176.
- [18] E. Hazan, A. Agarwal, S. Kale. Logarithmic regret algorithms for online convex optimization [J]. Machine Learning, 2007, 69(2-3):169-192.
- [19] L. Jian, J. Li, H. Liu. Toward online node classification on streaming networks [J]. Data Mining and Knowledge Discovery, 2018, 32(1):231-257.
- [20] Freund Y, Schapire R. E. Experiments with a new boosting algorithm [C]//In Proceedings of the 13th International Conference on Machine Learning. 1996:148-156.
- [21] Littlestone N, Warmuth M. K. The weighted majority algorithm [T]. University of California, Santa Cruz, Computer Research Laboratory. pages 256-261.
- [22] Cook R. D, Weisberg S. Characterizations of an empirical influence function for detecting influential cases in regression [J]. Technometrics, 1980, 22(4):495-508.
- [23] Wojnowicz, M., Cruz, B., Zhao, X., Wallace, B., Wolff, M., Luan, J., and Crable, C. "Influence sketching": Finding influential samples in large-scale regressions. arXiv preprint arXiv:1611.05923, 2016.
- [24] Koh, P.W. and Liang, P. Understanding black-box predictions via influence functions [C]// In Proceedings

of the 34th International Conference on Machine Learning, 2017:1885-1894.

optimization [J]. Foundations and Trends in Machine Learning, 2011, 4(2):107-194.

[25] Shai Shalev-Shwartz. Online learning and online convex

(责任编辑: X X)