

# 控制与决策

Control and Decision

## 基于数据分布特性的代价敏感宽度学习系统

徐鹏飞, 王敏, 刘金平, 唐朝晖, 马天雨

引用本文:

徐鹏飞, 王敏, 刘金平, 等. 基于数据分布特性的代价敏感宽度学习系统[J]. *控制与决策*, 2021, 36(7): 1686–1692.

在线阅读 View online: <https://doi.org/10.13195/j.kzyjc.2019.1484>

---

## 您可能感兴趣的其他文章

Articles you may be interested in

### 基于深度时序特征迁移的轴承剩余寿命预测方法

Remaining useful life prediction of bearing based on deep temporal feature transfer

*控制与决策*. 2021, 36(7): 1699–1706 <https://doi.org/10.13195/j.kzyjc.2019.1809>

### 嵌入重采样技术的C4.5决策树集成分类算法的临床医学预测

Clinical prediction of C4.5 decision tree classification algorithm with embedded resampling technique

*控制与决策*. 2021, 36(6): 1342–1350 <https://doi.org/10.13195/j.kzyjc.2019.1247>

### 面向复杂网络的异常检测研究进展

Research progress of anomaly detection for complex networks

*控制与决策*. 2021, 36(6): 1293–1310 <https://doi.org/10.13195/j.kzyjc.2020.0055>

### 基于HI-DD-AdaBoost.RT的锂离子动力电池SOH预测

Prediction of Li-ion battery SOH based on HI-DD-AdaBoost.RT

*控制与决策*. 2021, 36(3): 686–692 <https://doi.org/10.13195/j.kzyjc.2019.0764>

### 基于广义罚函数可行性准则的DE算法对不确定数据的处理

Application of improved DE algorithm based on generalized penalty function feasibility criteria in uncertain data processing

*控制与决策*. 2021, 36(2): 498–504 <https://doi.org/10.13195/j.kzyjc.2019.0728>

# 基于数据分布特性的代价敏感宽度学习系统

徐鹏飞<sup>1</sup>, 王敏<sup>1</sup>, 刘金平<sup>1†</sup>, 唐朝晖<sup>2</sup>, 马天雨<sup>1</sup>

(1. 湖南师范大学 信息科学与工程学院, 长沙 410081; 2. 中南大学 自动化学院, 长沙 410083)

**摘要:** 宽度学习系统(broad learning system, BLS)作为深度神经网络的替代框架,具有快速自适应模型结构选择和在线增量学习能力,被认为是知识发现和工程数据领域中一种极具前途的技术. 传统的BLS主要应用于数据分布均衡且误分类代价相同的模式分类任务,但大多数实际应用的数据是非均衡分布的,如网络入侵监测、医疗诊断、信用卡欺诈检测等. 基于此,提出一种基于数据分布特性的代价敏感BLS (data distribution-based cost-sensitive-BLS, DDbCs-BLS),解决数据分布不均、误分代价不同的模式分类任务. DDbCs-BLS在充分考虑数据统计分布特性的基础上寻找代价敏感型BLS分类器的最佳分类边界,保证少数类样本信息不被丢失,从而提高BLS在各类数据集上的模式分类性能. 在多种公共数据集(包括均衡和不均衡数据集)上进行大量的验证性和对比性实验,结果表明DDbCs-BLS能有效确定分类边界线的最佳位置,无论是在均衡数据集还是在非均衡数据集上均能获得更好的分类性能.

**关键词:** 非均衡数据; 代价敏感; 宽度学习系统; 自适应模型结构选择; 增量学习

中图分类号: TP273

文献标志码: A

DOI: 10.13195/j.kzyjc.2019.1484

开放科学(资源服务)标识码(OSID):



**引用格式:** 徐鹏飞,王敏,刘金平,等. 基于数据分布特性的代价敏感宽度学习系统[J]. 控制与决策, 2021, 36(7): 1686-1692.

## Data distribution-based cost-sensitive broad learning system

XU Peng-fei<sup>1</sup>, WANG Min<sup>1</sup>, LIU Jin-ping<sup>1†</sup>, TANG Zhao-hui<sup>2</sup>, MA Tian-yu<sup>1</sup>

(1. College of Information Science and Engineering, Hunan Normal University, Changsha 410081, China; 2. College of Automation, Central South University, Changsha 410083, China)

**Abstract:** Broad learning system (BLS) provides a flexible modeling framework, which is a potential substitute of deep neural network models. Due to its fast adaptive ability of automatic model structure selection and online incremental learning strategies, BLS is referred to as a promising technology in the field of knowledge discovery and data engineering. However, traditional BLS model are mainly aimed at pattern classification tasks with approximately even-distributed data and equal misclassification cost. In real applications, most of pattern recognition tasks are unevenly-distributed, such as credit card fraud detection, network intrusion detection, medical diagnosis, etc. In this paper, a data distribution-based cost-sensitive-BLS (DDbCs-BLS) is proposed for solving the problem of pattern classification tasks with imbalance data and varying misclassification costs on different classes. The DDbCs-BLS can achieve the best classification boundary by adopting the cost sensitive BLS learners, and ensure the lossless of the information of sparse classes, so as to ensure the classification performance of the BLS classifier in various data sets. The DDbCs-BLS is validated on multiple public data sets (including balanced and imbalanced data sets). Extensive validation and comparative results show that the DDbCs-BLS can effectively determine the best location of the classification boundary line, consequently, it can achieve better classification performance on both balanced and imbalanced data sets.

**Keywords:** imbalanced dataset; cost sensitive; broad learning system; adaptive model structure selection; incremental learning

## 0 引言

随着多层神经网络结构及其学习方法的不断发

展,深度神经网络在诸如计算机视觉<sup>[1-2]</sup>、工业故障检测<sup>[3-4]</sup>、音视频识别<sup>[5-7]</sup>等多个领域得到广泛应用. 例

收稿日期: 2019-10-23; 修回日期: 2020-01-16.

基金项目: 国家自然科学基金项目(61971188); 湖南省自然科学基金项目(2018JJ3349); 湖南省教育厅优秀青年项目(19B364); 湖南省知识产权战略推进专项项目(2019F012K); 湖南省研究生科研创新项目(CX20190415).

责任编辑: 阳春华.

<sup>†</sup>通讯作者. E-mail: ljp@hunnu.edu.cn.

如,生成对抗网络(GAN)可被用于图像生成的超分辨率任务<sup>[8]</sup>,利用深度置信网络(DBN)进行深度特征提取<sup>[9]</sup>,基于递归神经网络(RNN)对视频图像目标进行监控<sup>[10]</sup>等. 尽管深层结构如此强大,但仍面临着一些问题:复杂结构和大量的超参数,使得多数网络受到训练过程耗时的影响. 为解决这一问题,Chen等<sup>[11]</sup>提出了宽度学习系统(broad learning system, BLS),通过扩展神经网络的宽度而不是深度来解决高维数据问题,不仅克服了深度学习训练过程耗时的缺点,还可以快速增量构建网络模型.

Chen等<sup>[12]</sup>已经证明BLS具有通用的逼近能力,提出一种基于正则化稳健的BLS来学习不确定的数据模型<sup>[13]</sup>. Xu等<sup>[14]</sup>基于BLS提出一种时间序列预测模型(recurrent BLS)用于多变量预测时间序列. BLS具有通用的逼近能力、良好的泛化,且拥有时间记忆和灵活的重塑过程,被认为是一种极具前途的深度学习替代方法. 然而,传统的BLS在分类问题上默认数据样本是类别分布均衡或者错分类的代价相同,但实际应用中的大部分模式分类问题面临的都是一些非均衡数据,而且少数类往往占据更重要的地位,如网络入侵检测、医疗诊断、网络欺诈监测等<sup>[15-17]</sup>. 为此,本文提出一种基于数据分布特性的代价敏感宽度学习系统,旨在提供一种简单、高效的模式分类方法,在快速收敛、成本节约的基础上针对非均衡数据获得较高的准确率,主要贡献如下:

1) 针对数据非均衡分布问题,将通用的代价敏感学习框架引入BLS,提出代价敏感宽度学习系统(cost-sensitive-BLS, Cs-BLS),将传统的BLS扩展到不均衡数据集的模式分类任务.

2) 在Cs-BLS中考虑样本数据的分布特性,提出基于数据分布特性的代价敏感宽度学习系统(data distribution-based cost-sensitive BLS, DDbCs-BLS),根据不同类别数据的分布特性确定最佳分类边界,使少数类样本不被边缘化,有效应对各种分布特性数据集的模式分类任务.

## 1 BLS基本结构

宽度学习系统是以随机向量函数链接神经网络(random vector functional link neural network, RVFL)为载体,通过神经节点的增量实现网络横向扩展的一种单层神经网络学习系统,模型结构如图1所示.

BLS使用 $n$ 组特征映射窗口和 $m$ 组增强节点取代RVFL中输入与输出之间的直接连接,每组特征映射窗口包含 $N$ 个神经节点,每组增强节点包含 $M$ 个神经节点.

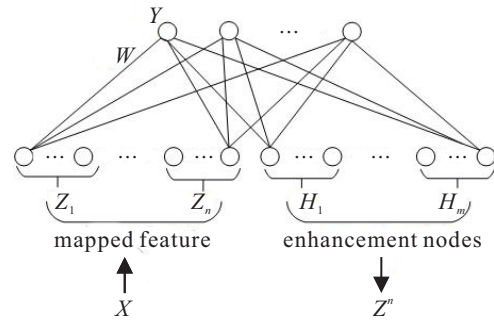


图1 BLS基本结构

首先,样本矩阵  $\mathbf{X} = [x_1, x_2, \dots, x_L]^T_{L \times D}$  经过如下公式的线性变换形成特征映射窗口  $\mathbf{Z}_i$ :

$$\mathbf{Z}_i = \phi(\mathbf{X}\mathbf{W}_{zi} + \beta_{zi})_{L \times N}, i = 1, 2, \dots, n. \quad (1)$$

其中: $L$ 为样本总数, $D$ 为样本 $x_l \in \mathbf{R}^{1 \times D}$ 的维度;权重矩阵 $\mathbf{W}_{zi} \in \mathbf{R}^{D \times N}$ 和偏置矩阵 $\beta_{zi} \in \mathbf{R}^{L \times N}$ 随机产生,各特征映射窗口的线性函数可以不同;所有特征映射窗口 $\mathbf{Z}_i$ 拼接组成特征映射层 $\mathbf{Z}^n$ 为

$$\mathbf{Z}^n = [\mathbf{Z}_1, \dots, \mathbf{Z}_n]_{L \times Nn}. \quad (2)$$

然后,特征映射层 $\mathbf{Z}^n$ 经过非线性变换形成增强节点,第 $j$ 组增强节点的输出如下所示:

$$\mathbf{H}_j = \delta(\mathbf{Z}^n\mathbf{W}_{hj} + \beta_{hj})_{L \times M}, j = 1, 2, \dots, m. \quad (3)$$

其中:权重矩阵 $\mathbf{W}_{hj} \in \mathbf{R}^{Nn \times M}$ 和偏置矩阵 $\beta_{hj} \in \mathbf{R}^{L \times M}$ 随机产生,各增强节点的非线性函数可以不同;所有增强节点 $\mathbf{H}_j$ 拼接组成增强层 $\mathbf{H}^m$ 为

$$\mathbf{H}^m = [\mathbf{H}_1, \dots, \mathbf{H}_m]_{L \times Mm}. \quad (4)$$

最后,将特征映射层 $\mathbf{Z}^n$ 和增强层 $\mathbf{H}^m$ 直接连接到BLS的输出端,不妨设 $\mathbf{A} = [\mathbf{Z}^n | \mathbf{H}^m]$ ,则BLS的最终估计输出为

$$\mathbf{Y} = \mathbf{A}\mathbf{W}. \quad (5)$$

其中: $\mathbf{A} \in \mathbf{R}^{L \times (Nn + Mm)}$ 为样本矩阵 $\mathbf{X}$ 变换后的输入矩阵, $\mathbf{Y} \in \mathbf{R}^{L \times c}$ 为样本矩阵 $\mathbf{X}$ 对应的预测值, $c$ 为样本类别数, $\mathbf{W} \in \mathbf{R}^{(Nn + Mm) \times c}$ 为输出权重矩阵.

求解输出权重矩阵 $\mathbf{W}$ 的目标优化函数如下所示:

$$\arg \min_{\mathbf{W}} : \|\mathbf{A}\mathbf{W} - \mathbf{Y}\|_2^2 + \lambda \|\mathbf{W}\|_2^2, \quad (6)$$

其中 $\lambda$ 为正则化参数.

## 2 所提出的方法

本节先简要推导Cs-BLS,并分析数据的分布特性;然后详细介绍所提出的DDbCs-BLS,并对方法进行简单的理论分析.

### 2.1 Cs-BLS

代价敏感学习利用代价矩阵的不同错误分类成本来解决不均衡数据集的模式分类问题<sup>[18]</sup>.

典型的代价敏感矩阵如下所示:

$$C = \begin{bmatrix} 0 & C_{12} & \dots & C_{1L} \\ C_{21} & 0 & \dots & C_{2L} \\ \vdots & \vdots & \ddots & \vdots \\ C_{L1} & C_{L2} & \dots & 0 \end{bmatrix}_{L \times L},$$

其中  $C_{ij}$  表示将样本  $x_i$  误判为样本  $x_j$  的代价. 通常, 少数类误判为多数类的代价高于多数类误分为少数类的代价.

在式(6)的BLS优化目标函数中, Cs-BLS为误差函数引入一个代价敏感矩阵, 如下所示:

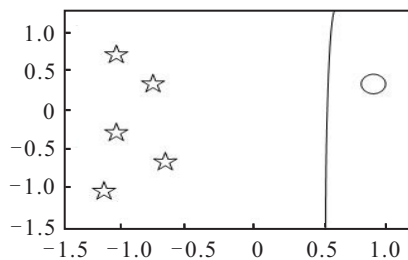
$$\arg \min_{\mathbf{W}} : \lambda \|\mathbf{W}\|_2^2 + \mathbf{C} \|\mathbf{A}\mathbf{W} - \mathbf{Y}\|_2^2, \quad (7)$$

使少数类获得一个较大权值, 提高分类器对少数类的敏感度, 解决非均衡样本的有效分类问题, 其中代价敏感矩阵  $\mathbf{C}$  随机产生.

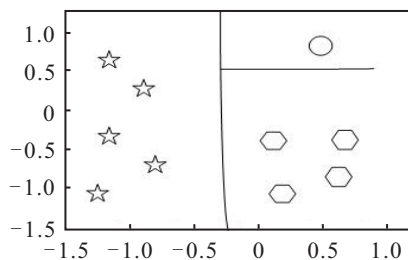
### 2.2 DDbCs-BLS

在使用Cs-BLS进行模式分类时, 若样本数据分布不均衡, 则多数类会凭借自身数量优势积累训练误差, 忽略少数类积累训练误差, 导致多数类边界距离最大化, 即分类边界线偏向少数类.

例如, 在图2所示的Cs-BLS分类示意图中, 发现分类边界线偏向于少数类, 这种少数类边界的模糊性会加大分类难度, 降低少数类的泛化性能.



(a) 二分类



(b) 三分类

○ ☆ maj class    ○ min class

图2 Cs-BLS分类示意

显然, 适当将分类边界向中间移动是一种较为理想的结果. 为此, DDbCs-BLS根据文献[19]中加权ELM的思想, 为少数类样本分配一个较大的权值, 以提高少数类积累训练误差的影响. 为多数类分配一个较小的权值, 以降低多数类积累训练误差的影响,

即再次为每个样本  $x_i$  赋予一个权值  $T_{ii}$ , 具体计算过程如下所示:

$$T_{ii} = \begin{cases} \frac{\varepsilon}{\text{Class}(x_i)}, & \text{Class}(x_i) > \text{AVG}(\text{class}); \\ 1, & \text{otherwise.} \end{cases} \quad (8)$$

其中:  $0 < \varepsilon \leq 1$ ,  $\text{Class}(x_i)$  表示与样本  $x_i$  同类别的样本总数,  $\text{AVG}(\text{class})$  表示样本类别的平均数.

从本质上, 权值  $T_{ii}$  决定了用户寻求的重新均衡程度, 以及边界向多数阶级推进的程度, 但其值并不是越大越好. 例如, 在图3(a)和(d)中,  $\varepsilon = 0.2$ , 出现与式(7)相似的情况, 关注点放在多数类, 少数类的泛化性能下降; 在图3(c)和(f)中,  $\varepsilon = 1$ , 更多的关注误差最小化, 边界是弯向少数类, 如果数据非均衡度提高, 边界线会更加偏向少数类, 少数类性能提高的并不明显; 在图3(b)和(e)中,  $\varepsilon = 0.6$ , 分类边界线偏向多数类, 边界线偏向于平滑, 没有出现明显弯曲现象, 具有更好的分类特性.

将所有样本的权值  $T_{ii}$  加入式(7)后, 可得DDbCs-BLS目标优化函数

$$\arg \min_{\mathbf{W}} : \lambda \|\mathbf{W}\|_2^2 + \mathbf{Q} \|\mathbf{A}\mathbf{W} - \mathbf{Y}\|_2^2. \quad (9)$$

其中:  $\mathbf{Q} = \mathbf{T}\mathbf{C}$ ,  $\mathbf{T} = \text{diag}(T_{ii})$ ,  $i = 1, 2, \dots, L$ .

求解式(9)的目标优化函数, 可得DDbCs-BLS的最终输出权重矩阵  $\mathbf{W}$  为

$$\mathbf{W} = (\lambda \mathbf{I} + \mathbf{A}^T \mathbf{Q} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Q} \mathbf{Y}. \quad (10)$$

### 2.3 DDbCs-BLS增量学习

假设现有样本为  $\mathbf{X}_{L_0}$ , 依据式(10), 现有样本的输出权重矩阵为

$$\mathbf{W}_0 = (\lambda \mathbf{I} + \mathbf{A}_0^T \mathbf{Q}_0 \mathbf{A}_0)^{-1} \mathbf{A}_0^T \mathbf{Q}_0 \mathbf{Y}_0. \quad (11)$$

设工作方阵  $\mathbf{P}_0 = (\lambda \mathbf{I} + \mathbf{A}_0^T \mathbf{Q}_0 \mathbf{A}_0)^{-1}$ , 则有

$$\mathbf{P}_0^{-1} = \lambda \mathbf{I} + \mathbf{A}_0^T \mathbf{Q}_0 \mathbf{A}_0, \quad (12)$$

$$\mathbf{W}_0 = \mathbf{P}_0 \mathbf{A}_0^T \mathbf{Q}_0 \mathbf{Y}_0. \quad (13)$$

当一批新样本  $\mathbf{X}_{L_1}$  加入在线学习时, 设  $\mathbf{A}_1$  为新样本的变换输入矩阵, 依据式(10), 现有样本合并新样本后的输出权重矩阵为

$$\mathbf{W}_1 =$$

$$\left( \lambda \mathbf{I} + \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}^T \begin{bmatrix} \mathbf{Q}_0 & 0 \\ 0 & \mathbf{Q}_1 \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix} \right)^{-1} \times$$

$$\begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}^T \begin{bmatrix} \mathbf{Q}_0 & 0 \\ 0 & \mathbf{Q}_1 \end{bmatrix} \begin{bmatrix} \mathbf{Y}_0 \\ \mathbf{Y}_1 \end{bmatrix} =$$

$$(\lambda \mathbf{I} + \mathbf{A}_0^T \mathbf{Q}_0 \mathbf{A}_0 + \mathbf{A}_1^T \mathbf{Q}_1 \mathbf{A}_1)^{-1} (\mathbf{A}_0^T \mathbf{Q}_0 \mathbf{Y}_0 + \mathbf{A}_1^T \mathbf{Q}_1 \mathbf{Y}_1) =$$

$$(\mathbf{P}_0^{-1} + \mathbf{A}_1^T \mathbf{Q}_1 \mathbf{A}_1)^{-1} (\mathbf{A}_0^T \mathbf{Q}_0 \mathbf{Y}_0 + \mathbf{A}_1^T \mathbf{Q}_1 \mathbf{Y}_1). \quad (14)$$

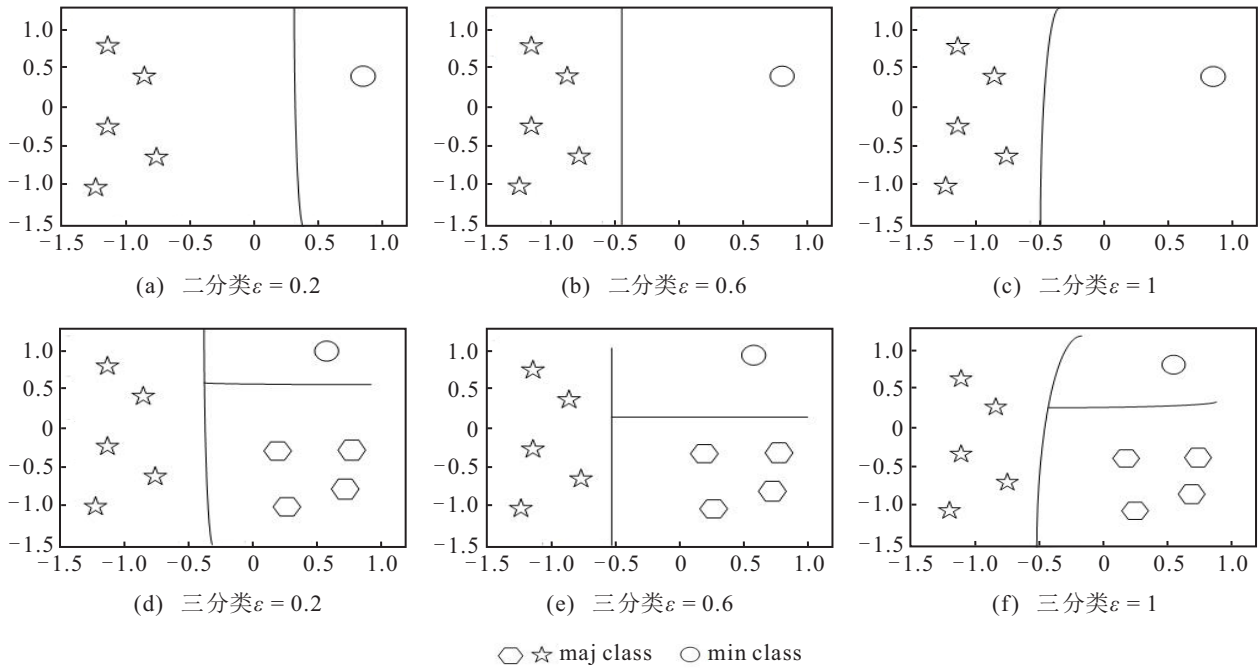


图3 DDbCs-BLS分类示意

设工作方阵  $P_1 = (P_0^{-1} + A_1^T Q_1 A_1)^{-1}$ , 则有

$$P_0^{-1} = P_1^{-1} - A_1^T Q_1 A_1. \quad (15)$$

从而输出权重矩阵  $W_1$  又可以改写为

$$\begin{aligned} W_1 &= \\ P_1 P_0^{-1} P_0 A_0^T Q_0 Y_0 + P_1 A_1^T Q_1 Y_1 &= \\ P_1 (P_1^{-1} - A_1^T Q_1 A_1) W_0 + P_1 A_1^T Q_1 Y_1 &= \\ W_0 + P_1 A_1^T Q_1 (Y_1 - A_1 W_0). \end{aligned} \quad (16)$$

因此, 当第  $k$  批样本  $X_{L_k}$  加入后, 可以递推得到

$$W_k = W_{k-1} + P_k A_k^T Q_k (Y_k - A_k W_{k-1}), \quad (17)$$

其中工作方阵  $P_k = (P_{k-1}^{-1} + A_k^T Q_k A_k)^{-1}$ .

### 2.4 算法流程

综合上述, DDbCS-BLS 算法流程描述如下.

#### 初始化阶段:

step 1: 随机初始化权重矩阵  $W_{zi}$  和偏置矩阵  $\beta_{zi}$ , 根据式(1)计算特征映射窗口  $Z_i$ , 将所有  $Z_i$  组成特征映射层  $Z^n$ ;

step 2: 随机初始化权重矩阵  $W_{hj}$  和偏置矩阵  $\beta_{hj}$ , 根据式(3)计算增强节点  $H_j$ , 将所有  $H_j$  组成增强层  $H^m$ ;

step 3: 将特征映射层  $Z^n$  和增强层  $H^m$  拼接组成初始样本的变换输入矩阵  $A_0$ ;

step 4: 随机初始化代价敏感矩阵  $C$ , 根据式(8)计算对角矩阵  $T$ , 计算初始样本的权值矩阵  $Q_0 = TC$ ;

step 5: 计算初始样本的工作方阵  $P_0$ , 根据式(13)

计算初始样本的输出权重矩阵  $W_0$ .

#### 增量学习阶段:

step 6: 当新一批样本  $X_{L_k}$  加入时, 重复 step 1 ~ step 4, 计算新样本的变换输入矩阵  $A_k$  和权值矩阵  $Q_k$ ;

step 7: 依据工作方阵  $P_{k-1}$ , 递推更新加入新样本后的工作方阵  $P_k$ ;

step 8: 根据式(17), 递推更新加入新样本后的输出权重矩阵  $W_k$ ;

step 9: return  $P_k$  和  $W_k$ .

### 2.5 理论分析

BLS 样本矩阵  $X$  为连续值, 所以式(6)在该定义域是可导的, 故式(6)有唯一最小解  $W$ , 根据式(9)推出的  $W$  也是唯一解. 由于加入新样本  $X_{L_k}$  后的样本矩阵仍然为连续值, 增量学习模型在该定义域仍是可导的, 故增量学习模型的输出权重矩阵  $W_k$  有唯一最小解.

假设特征映射层和增强层分别有  $N_n$  和  $M_m$  个神经节点, 现有样本的总数为  $L$ , 样本的维数为  $D$ , 则计算一个特征向量的时间是  $O(N_n M_m D)$ , 所有样本的时间复杂度约为  $O(L N_n M_m D)$ . 设新加入样本  $X_{L_k}$  的数量为  $L_k$ , 递推更新  $W_k$  的主要计算是式(17)的第2项, 时间复杂度约为  $O(L_k N_n M_m D)$ ; 随着新样本的不断加入, 将有  $L_k \ll L$ , 递推更新时间复杂度远小于非递推所需时间复杂度  $O(L N_n M_m D + L_k N_n M_m D)$ .

### 3 实验验证

为验证所提方法的有效性,在多种公共数据集进行验证性和对比性实验.实验环境如下:处理器型号为Intel (R) Core (TM) i5-4210 U CPU @1.70 GHz,运行内存为4 GB;操作系统为64位Windows;实现语言为Python 3.7.

BLS和DDbCs-BLS设置相同实验参数:特征映射层和增强层的神经节点个数分别为150和200,正则化系数 $\lambda = 2^{-30}$ , $\varepsilon = 0.618$ ,非线性函数 $\delta = \tanh(x)$ 表示为

$$\tanh(x) = \frac{2}{1 + \exp(-2x)} - 1. \quad (18)$$

#### 3.1 验证性实验

##### 3.1.1 均衡数据集

MNIST数据集由黑白手写0~9数字图像组成,每个数字尺寸为 $28 \times 28$ 的图像,其中训练集有60 000张图像,测试集有10 000张图像.

NORB数据集是以不同照明及摆放方式摄制50种不同的3D玩具模型图像,其中训练集有24 300张图像,测试集有24 300张图像.

表1 MNIST和NORB数据集accuracy %

	BLS	DDbCs-BLS
MNIST	98.74	99.07
NORB	89.24	89.31

从表1的结果可以看出,DDbCs-BLS分别以99.07%和89.31%的准确度展示了更好的性能,比传统BLS分别高0.33%和0.04%.

##### 3.1.2 非均衡数据集

本组实验选用UCI的8组非均衡数据集,如表2所示,其中IR为不均衡率. Adult与CMC为普通的分类数据集,其余6组为医学应用数据集.

表2 UCI的8组非均衡数据集基本信息

名称	样本数	特征数	类别	训练集	测试集	IR
CMC	1 473	10	3	1 031	442	3.42
Adult	15 003	14	2	10 503	4 500	3.08
Blood	671	5	2	470	201	3.21
Pima	768	8	2	538	230	1.87
Breast1	683	9	2	479	204	1.86
Breast2	569	30	2	399	170	1.61
Haberman	306	3	2	214	92	2.78
ILPD	583	10	2	408	175	2.49

针对非均衡数据集,分类准确度accuracy已不合适,为此本组实验使用G-mean值来衡量整体的分类性能,其定义如下:

$$G\text{-mean} = \sqrt{\frac{TP}{TP + FN} \times \frac{TN}{TN + FP}}, \quad (19)$$

其中各参数如表3所示.通常,只有在多数类和少数类的分类精度同时都高时,G-mean值才会更大.

表3 混淆矩阵

	预测正类	预测负类
真正类	TP	FN
真负类	FP	TN

两种方法的G-mean值如表4所示,DDbCs-BLS明显优于传统BLS,仅在Breast2数据集中的训练准确率比BLS低0.1%,在Adult数据集上训练准确率要比BLS高5.7%,在Pima数据集训练和测试的准确率分别比BLS高6.5%和7.9%.

表4 两种方法在不同数据集上的G-mean

dataset	方法	G-mean / %	
		train	test
CMC	BLS	80.1	79.2
	DDbCs-BLS	83.4	82.3
Adult	BLS	73.4	73.2
	DDbCs-BLS	79.1	77.6
Blood	BLS	70.4	72.3
	DDbCs-BLS	73.4	74.2
Pima	BLS	71.2	68.3
	DDbCs-BLS	77.7	76.2
Breast1	BLS	78.4	78.9
	DDbCs-BLS	80.1	78.9
Breast2	BLS	84.3	83.8
	DDbCs-BLS	84.2	84.0
Haberman	BLS	76.6	75.3
	DDbCs-BLS	79.9	81.0
ILPD	BLS	75.9	74.4
	DDbCs-BLS	77.3	76.9

#### 3.1.3 实验结果与分析

通过上述两组对比实验可以发现:在均衡数据集上,DDbCs-BLS不会降低传统BLS的分类精度;而在非均衡数据集上,DDbCs-BLS较传统BLS获得了更为理想的分类性能.其主要原因是:DDbCs-BLS利用代价敏感学习使数据得到均衡,并充分考虑数据集的分布特性,进而得到分类模型的最佳边界线,使每个类别都能最大限度地得到正确分类.

#### 3.2 对比性实验

##### 3.2.1 数据集与评价标准

本组实验选用KEEL非均衡数据库的8组数据集,如表5所示,其中训练集与测试集按照7:3的比例随机划分.计算F-value、G-mean、AUC三个评价指标值,并与CWsRF<sup>[20]</sup>、WOS-ELM<sup>[21]</sup>、CS-CNN<sup>[22]</sup>等非

均衡方法进行比较.

表5 KEEL的8组非均衡数据集基本信息

数据集	特征	样本	少数类	多数类	IR
Segment	19	2 308	329	1 979	6.02
Vehicle0	18	846	199	647	3.25
Vehicle1	18	846	217	629	2.90
Vehicle2	18	846	218	628	2.88
Vehicle3	18	846	212	634	2.99
Yeast1	8	1 484	429	1 055	2.46
Yeast3	8	1 484	161	1 321	8.10
Page-block0	10	5 472	559	4 913	8.79

$F$ -value 是查全率和查准率的调和均值, 可以衡量少数类样本的精确度, 其定义如下:

$$F\text{-value} = \frac{(1 + \eta^2)\text{Precision} \times \text{Recall}}{\eta^2 \times (\text{Recall} + \text{Precision})} \quad (20)$$

其中:  $\text{Precision} = \frac{TP}{TP + FP}$ , 表示真正类占所有预测正类的比例;  $\text{Recall} = \frac{TP}{TP + FN}$ , 表示真正类占实际正类的比例; 参数  $\eta$  用来调节 Precision 和 Recall 的相对重要度, 一般取值为 1. 通常,  $F$ -value 数值越大, 少数类样本的分类精度越高.

AUC 是由 ROC 曲线演变而来, 常用曲线下面积 AUC 来代替 ROC 曲线作为评价方法, 提供了评价模型平均性能的另一方法, 它表示预测的正例排在负例前面的概率, AUC 值越大, 模型越好.

### 3.2.2 实验结果与分析

1)  $F$ -value 值, 如表 6 所示. DDbCs-BLS 的  $F$ -value 明显高于 CWsRF 和 CS-CNN, 仅仅在 Vehicle3 数据集上比 WOS-ELM 低 0.19%. 在 Page-block0 和 Vehicle1 数据集上分别比 WOS-ELM 高 4.68% 和 2.91%. 即本文方法能更好地细化出数据的分类边界线, 使少数类样本的分类精度更高.

表6 4种方法在不同数据集上的  $F$ -value %

Datasets	CWsRF	WOS-ELM	CS-CNN	DDbCs-BLS
Segment	98.43	98.70	98.62	99.71
Vehicle0	98.02	96.97	97.68	99.32
Vehicle1	85.21	83.54	86.32	86.45
Vehicle2	99.11	98.73	97.56	99.32
Vehicle3	83.95	85.01	82.77	84.82
Yeast1	76.96	76.95	75.45	78.29
Yeast3	93.22	92.91	92.46	93.56
Page-block0	96.78	93.40	96.06	98.08

2)  $G$ -mean 值, 如表 7 所示. 在 8 组数据集上 DDbCs-BLS 所获得的  $G$ -mean 值都要比 CWsRF 和 WOS-ELM 高. 虽然在 Vehicle1 数据集上略比 CS-CNN 低 1.11%, 但在其他数据集上的  $G$ -mean 值有不同程度的提高, 其中在 Page-block0 与 Segment 数据

集上的  $G$ -mean 值要分别比 CS-CNN 方法高 2.2% 和 0.96%.

表7 4种方法在不同数据集上的  $G$ -mean %

Datasets	CWsRF	WOS-ELM	CS-CNN	DDbCs-BLS
Segment	96.33	96.59	96.66	97.62
Vehicle0	96.59	96.97	97.68	98.02
Vehicle1	84.92	82.42	86.42	85.31
Vehicle2	98.21	95.67	97.31	98.55
Vehicle3	82.75	83.41	82.77	83.82
Yeast1	75.46	75.59	74.54	76.91
Yeast3	92.12	91.90	91.61	93.65
Page-block0	94.97	92.82	95.66	97.68

3) AUC 值, 如表 8 所示. 在 8 组数据集上, 本文方法的 AUC 值基本高于其他 3 种方法, 说明本文方法的平均性能得到显著提升.

表8 4种方法在不同数据集上的 AUC %

Datasets	CWsRF	WOS-ELM	CS-CNN	DDbCs-BLS
Segment	96.31	96.98	96.33	97.61
Vehicle0	96.76	97.07	97.64	97.92
Vehicle1	85.02	83.02	85.23	85.31
Vehicle2	97.93	95.76	97.64	98.65
Vehicle3	82.89	83.56	82.77	83.82
Yeast1	75.35	76.61	75.94	77.00
Yeast3	92.21	92.00	91.96	93.03
Page-block0	95.79	92.97	96.46	97.01

## 4 结论

本文提出了一种基于数据分布特性的代价宽度学习框架 DDbCs-BLS, 以 BLS 为理论基础, 首先采用代价敏感矩阵对非均衡数据进行加权调整比例, 然后通过分析数据的分布特性确定最佳分类边界线, 提高了方法的稳定性. 在多种公共数据集(包括均衡数据集和不均衡数据集)上进行了大量的验证性和对比性实验, 表明 DDbCs-BLS 一定程度上提高了 BLS 对非均衡数据分类的性能, 能够在不降低 BLS 对均衡数据分类精度的同时, 保证分类器对少数类的正确分类. 下一步的研究工作是将 DDbCs-BLS 在复杂工业过程中进行工业验证与应用, 以进一步改进提高所提方法的性能.

### 参考文献(References)

[1] 刘金平, 何捷舟, 唐朝晖. 基于 WCGAN 的矿物浮选泡沫图像光照不变颜色提取[J]. 自动化学报, 2019, DOI: <https://doi.org/10.16383/j.aas.c190330>.  
(Liu J P, He J Z, Tang Z H. WCGAN-based illumination-invariant color measuring of mineral flotation frothImages[J]. Acta Automatica Sinica, 2019, DOI: <https://doi.org/10.16383/j.aas.c190330>.)

[2] He X D, Deng L. Deep learning for image-to-text

- generation: A technical overview[J]. IEEE Signal Processing Magazine, 2017, 34(6): 109-116.
- [3] 陈桥, 丁宝苍, 王雅楠, 等. 基于目标跟踪的双层结构工业预测控制[J]. 控制与决策, 2017, 32(5): 797-803. (Chen Q, Ding B C, Wang Y N, et al. Double-layered industrial predictive control based on target tracking[J]. Control and Decision, 2017, 32(5): 797-803.)
- [4] 丁进良, 杨翠娥, 陈远东, 等. 复杂工业过程智能优化决策系统的现状与展望[J]. 自动化学报, 2018, 44(11): 1931-1943. (Ding J L, Yang C E, Chen Y D, et al. Research progress and prospects of intelligent optimization decision making in complex industrial process[J]. Acta Automatica Sinica, 2018, 44(11): 1931-1943.)
- [5] Tassi A, Khirallah C, Vukobratovic D. Resource allocation strategies for network-coded video broadcasting services over LTE-advanced[J]. IEEE Transactions on Vehicular Technology, 2015, 64(5): 2186-2192.
- [6] Rao Y M, Lu J W, Zhou J. Learning discriminative aggregation network for video-based face recognition and person re-identification[J]. International Journal of Computer Vision, 2019, 127(6/7): 701-718.
- [7] 黄雅婷, 石晶, 许家铭, 等. 鸡尾酒会问题与相关听觉模型的研究现状与展望[J]. 自动化学报, 2019, 45(2): 234-251. (Huang Y T, Shi J, Xu J M, et al. Research advances and perspective on the cocktail party problem and related auditory models[J]. Acta Automatica Sinica, 2019, 45(2): 234-251.)
- [8] Song D Y, Chandolu A, Stojanovic N. Effect of impurity incorporation on emission wavelength in cathodoluminescence spectrum image study of GaN pyramids grown by selective area epitaxy[J]. Journal of Applied Physics, 2008, 104(6): 064309.
- [9] Li J J, Xi B B, Li Y S, et al. Hyperspectral classification based on texture feature enhancement and deep belief networks[J]. Remote Sensing, 2018, 10(3): 396.
- [10] Liu F, Chen Z G, Wang J. Video image target monitoring based on RNN-LSTM[J]. Multimedia Tools Applications, 2019, 78(4): 4527-4544.
- [11] Chen C L P, Liu Z L. Broad learning system: An effective and efficient incremental learning system without the need for deep architecture[J]. IEEE Transactions on Neural Networks and Learning Systems, 2018, 29(1): 10-24.
- [12] Chen C L P, Liu Z L, Feng S. Universal approximation capability of broad learning system and its structural variations[J]. IEEE Transactions on Neural Networks and Learning Systems, 2019, 30(4): 1191-1204.
- [13] Jin J W, Philip Chen C L. Regularized robust broad learning system for uncertain data modeling[J]. Neurocomputing, 2018, 322: 58-69.
- [14] Xu M, Han M, Chen C L P. Recurrent broad learning systems for time series prediction[J]. IEEE Transactions on Cybernetics, 2018, 9(10): 1-13.
- [15] Ogunleye A, Wang Q G, Marwala T. Integrated learning via randomized forests and localized regression with application to medical diagnosis[J]. IEEE Access, 2019, 7: 18727-18733.
- [16] Papamartzivanos D, Gomez Marmol F, Kambourakis G. Introducing deep learning self-adaptive misuse network intrusion detection systems[J]. IEEE Access, 2019, 7: 13546-13560.
- [17] Vinayakumar R, Alazab M, Soman K P, et al. Deep learning approach for intelligent intrusion detection system[J]. IEEE Access, 2019, 7(4): 41525-41550.
- [18] He H B, Garcia E A. Learning from imbalanced data[J]. IEEE Transactions on Knowledge and Data Engineering, 2009, 21(9): 1263-1284.
- [19] Zong W, Huang G B, Chen Y Q. Weighted extreme learning machine for imbalance learning[J]. Neurocomputing, 2013, 101: 229-242.
- [20] Zhu M, Xia J, Jin X Q, et al. Class weights random forest algorithm for processing class imbalanced medical data[J]. IEEE Access, 2018, 6: 4641-4652.
- [21] Mirza B, Lin Z P, Tou K A. Weighted online sequential extreme learning machine for class imbalance learning[J]. Neural Processing Letters, 2013, 38(3): 465-486.
- [22] Khan S H, Hayat M, Bennamoun M. Cost-sensitive learning of deep feature representations from imbalanced data[J]. IEEE Transactions on Neural Networks and Learning Systems, 2016, 29(8): 3573-3587.

### 作者简介

徐鹏飞(1977—), 男, 副教授, 博士, 从事智能信息处理等研究, E-mail: xupf@hunnu.edu.cn;

王敏(1995—), 女, 硕士生, 从事智能算法的研究, E-mail: 201870291208@smail.hunnu.edu.cn;

刘金平(1983—), 男, 副教授, 博士, 从事复杂工业过程自动化监控等研究, E-mail: ljp@hunnu.edu.cn;

唐朝晖(1965—), 男, 教授, 博士, 从事复杂工业过程建模与故障诊断等研究, E-mail: zhtang@csu.edu.cn;

马天雨(1978—), 男, 讲师, 博士, 从事复杂工业过程建模及优化控制等研究, E-mail: mty@hunnu.edu.cn.

(责任编辑: 齐 霖)