

基于 SMOTE+ENN 的个人信用评估方法

吕颖, 邢进生

(山西师范大学 数学与计算机科学学院, 山西 临汾 041004)

摘要:个人信用评估作为商业银行判定借贷风险的直接依据,在金融领域显得尤为重要。针对传统个人信用评估模型存在数据不平衡、模型结构单一、易受主观因素干扰等问题,提出一种基于 SMOTE+ENN (synthetic minority oversampling technique+edited nearest neighbours)算法与集成学习的个人信用评估方法。首先,该方法在数据预处理的基础上,采用 SMOTE+ENN 算法对样本数据进行数据平衡分布处理,增强了分类算法性能;然后,基于网格搜索优化算法,搜寻适用于多种分类器的最优超参数,进而构造出相应的最优单一评估模型,达到了提高个人信用评估精确度的目的;最后,利用相关的集成学习策略将表现最优的三种分类器结果集成,构造出信用评估的最优预测模型,从而实现更为准确的个人信用评估。实验结果表明,在现有公开数据集 Give Me Some Credit 上,与传统数据不平衡处理方法相比,该方法的预测准确率高达 97%,精确度提升约 2%,验证了算法改进的有效性。

关键词:信用评估;数据不平衡;数据预处理;网格搜索;集成学习

中图分类号:TP391

文献标识码:A

文章编号:1673-629X(2022)06-0045-07

doi:10.3969/j.issn.1673-629X.2022.06.008

Personal Credit Evaluation Method Based on SMOTE+ENN

LYU Ying, XING Jin-sheng

(School of Mathematics and Computer Science, Shanxi Normal University, Linfen 041004, China)

Abstract: Personal credit evaluation as a direct basis for commercial banks to judge loan risk is particularly important in the financial field. Aiming at the problems of the traditional personal credit evaluation model, such as data imbalance, single model structure and being easily interfered by subjective factors, a personal credit evaluation method based on SMOTE+ENN algorithm and ensemble learning is proposed. First of all, SMOTE+ENN algorithm is used to balance and distribute the sample data on the basis of data preprocessing, which enhances the performance of the classification algorithm. Then, based on the grid search optimization algorithm, the optimal super parameters suitable for a variety of classifiers are searched, and the corresponding optimal single evaluation model is constructed to achieve the purpose of improving the accuracy of personal credit evaluation. Finally, the results of the three classifiers with the best performance are integrated with the related ensemble learning strategy to construct the optimal prediction model of credit evaluation, so as to achieve a more accurate personal credit evaluation. Experiment shows that on the existing public dataset Give Me Some Credit, compared with the traditional data imbalance processing method, the proposed method is as high as 97% in prediction accuracy, and the accuracy is improved by about 2%, which verifies the effectiveness of the improved algorithm.

Key words: credit evaluation; data imbalance; data preprocessing; grid search; ensemble learning

0 引言

随着大数据时代的到来,互联网金融得到了快速的发展,个人信用在金融领域越来越重要。银行根据个人信用报告所记录的内容,对个人信用进行评估,不仅有助于帮助客户树立正确的信用观念,而且有利于提高银行的授信效率,扩大消费信贷的发放^[1]。

传统的信用评估方法主要以 BP (back propagation, BP) 神经网络、支持向量机 (support vector

machine, SVM)^[2]、逻辑回归 (logistics regression, LR)^[3]、决策树 (decision tree, DT)^[4] 等理论为基础。Chen Jie 等人^[5] 基于 BP 神经网络,建立绿色供应链合作信用评价模型,从而更好地学习和评估不同层次的绿色供应链信用。姜凤茹^[6] 利用支持向量机建立了个人信用评估模型,并引入遗传算法进行参数优化,实验证明遗传算法-支持向量机 (genetic algorithm-support vector machine, GA-SVM) 模型有效解决了 P2P (peer

收稿日期:2021-06-22

修回日期:2021-10-22

基金项目:山西省软科学基金资助项目(2011041033-03)

作者简介:吕颖(1995-),女,硕士研究生,研究方向为数据挖掘、计算机视觉;邢进生,教授,硕导,研究方向为计算智能、数据挖掘、计算机视觉。

to peer lending, P2P) 网贷平台的个人信用评估问题。李太勇等人^[7]充分利用稀疏贝叶斯学习 (sparse Bayesian learning, SBL) 的优势,使得特征权重尽量稀疏,以此实现了个人信用评估和特征选择。但随着金融数据量的增加和数据不平衡问题的出现,这些传统机器学习算法已不能更好地满足市场模型评估的需求^[8]。

近年来,随着深度学习^[9]相关理论以及信用评估方法^[10]的深入发展,国内外众多学者在数据处理及模型构建方面进行了相应研究,取得了良好的成绩。赵雪峰等人^[11]利用银行个人信贷数据,构建了基于卷积神经网络的信用贷款评估模型,提高了信用贷款模型的鲁棒性。Shen Feng 等人^[12]提出了一种改进的合成少数类过采样技术 (synthetic minority oversampling technique, SMOTE)^[13]用于不平衡数据处理,并将长短时记忆网络 (long short-term memory network, LSTM) 和自适应提升 (adaptive boosting, AdaBoost) 算法整合到集成框架中,解决了不平衡信用风险评估问题。Jie Sun 等人^[14]提出一种基于 SMOTE 和引导聚集 (bootstrap aggregating, Bagging) 集成学习^[15]算法,为不平衡的企业信用评估建立有效的决策树集成模型。Lean Yu 等人^[16]提出一种基于深度信念网络 (deep belief nets, DBN) 的重采样支持向量机集成学习范例,以解决信用分类中的不平衡数据问题。在现有数据不均衡的条件下,分类器集成俨然成为信用评估建模的发展趋势^[17]。

鉴于目前个人信用评估方法的研究现状,提出了基于 SMOTE+ENN 算法与集成学习的个人信用评估模型。该方法在解决样本数据不均衡的基础上,采用网格搜索^[18]算法对多种分类器模型进行超参数优化,并利用集成学习技术把最优模型结果集成,从而构建出最优的个人信用评估模型。为了验证算法的有效性,采用公开数据集 Give Me Some Credit 对该算法进行测试。与传统单一分类器模型预测结果相比,该模型不仅良好地解决了数据不均衡的问题,而且提高了信用评估的精度。

1 相关理论

1.1 SMOTE 算法

SMOTE 算法是一种合成少数类的过采样技术,广泛应用于处理数据不平衡的问题。其基本思想是基于对现有少数样本的分析,人工生成额外少量样本,并添加到原始数据集中,形成样本数量均衡的新数据集,从而有效改善了样本数据不均衡问题。SMOTE 算法具体步骤如下:

(1) 对于每个少数类中样本 x_a , 以欧氏距离为标

准,找到该样本最近的 k 个样本。

(2) 从 k 个近邻样本中随机选择一个样本 x_b 。

(3) 在 x_a 和 x_b 两个样本点之间随机插入一个新的样本 x_c , 合成新样本的数学公式如式(1):

$$x_c = x_a + \text{rand}(0,1) \times |x_a - x_b| \quad (1)$$

式中, $\text{rand}(0,1)$ 表示区间 $(0,1)$ 内的任何一个数。

1.2 网格搜索

网格搜索通过遍历所有超参数组合,寻找一组合适的超参数配置,从而达到优化模型的目的。假设总共有 n 个超参数,第 n 个超参数可以取 m_n 个值,即总共的配置组合数量为 $m_1 \times m_2 \times \dots \times m_n$ 。网格搜索根据这些超参数的不同组合分别训练一个模型,通过循环遍历进行测试,从中选择一组性能最好的参数。网格搜索优化参数的基本原理如下:

(1) 在初始状态下,网格搜索法依据经验先设置好待搜索的超参数区域,为防止关键参数被遗漏,设置较广的待搜索超参数区域范围。

(2) 设置参数搜索步长,从起始点出发,沿着参数不同增长方向以单位步长运动,将所到区域以网格形式表示出来,网格中的交叉点就是所要搜索的参数。

(3) 在搜索区域内选取网格的节点,通过交叉验证,测试该参数下模型性能,选取表现最优的参数。

1.3 集成学习

集成学习作为当前比较流行的机器学习算法之一,主要通过某种策略将多个弱监督模型集成,并利用群体决策来提高决策准确率,得到一个更好更全面的强监督模型。即便其中某个弱分类器预测错误,其他的弱分类器也可以及时纠正,从而获得比个体学习器更好的泛化性能。集成学习首先构建一组个体学习器,然后通过某种结合策略将其集成构建成一个强学习器,如图 1 所示。目前,集成学习常用的结合策略主要有以下几类:

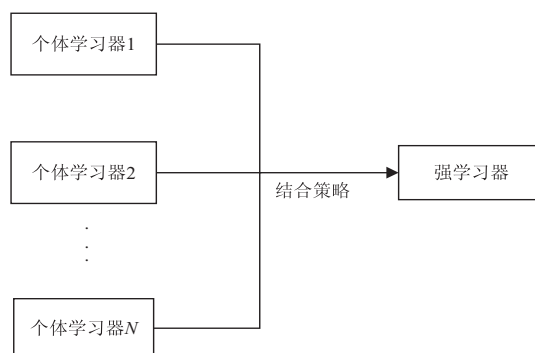


图 1 集成学习示意图

1.3.1 平均法

对于数值输出问题,平均法是最常用的集成学习结合策略。

(1) 简单算术平均的集成模型可由式(2)表达:

$$F(X) = \frac{1}{M} \sum_{m=1}^M f_m(x) \tag{2}$$

式中, $f(x)$ 为个体学习器预测结果, M 为模型数量。

(2) 加权算术平均的集成模型可由式(3)表达:

$$F(X) = \sum_{m=1}^M f_m(x) w_m \tag{3}$$

式中, w_m 为个体学习器 f_m 学习的权重, 通过训练数据获取, $w_m \geq 0$, $\sum_{m=1}^M w_m = 1$ 。当个体分类器性能差异较

大时, 宜采用加权算术平均法; 个体差异相差无几时, 宜采用简单算术平均法。

1.3.2 投票法

对于分类问题, 最常用的集成学习策略就是投票法, 投票法主要分为硬投票和软投票两种。

(1) 硬投票。

根据少数服从多数的原则, 选择所有个体分类器中输出最多的标签作为最终预测结果, 如图 2 所示。

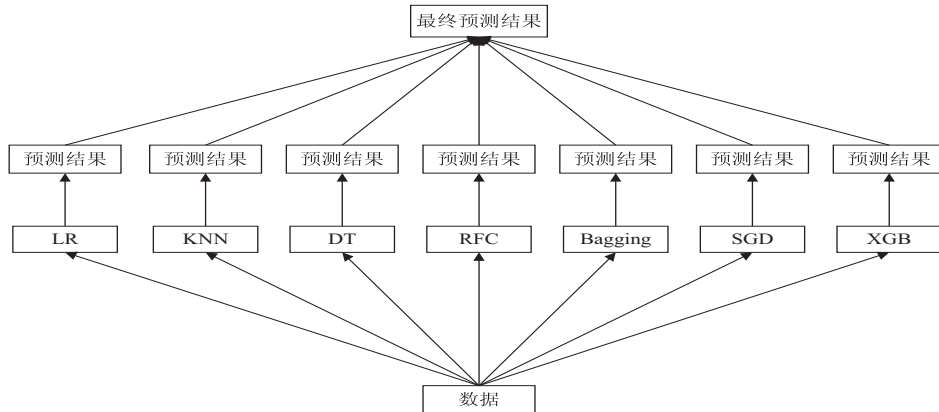


图 2 集成学习硬投票机制

(2) 软投票。

将预测结果为某一类别的所有分类器的预测概率平均, 平均概率最高的类别为最终预测结果。

1.3.3 学习法

学习法的代表方法是 Stacking, 它是一种分层模型集成框架。首先采用原始数据集训练出多个不同个体学习器, 然后再用训练好的个体学习器的输出作为新训练集的输入, 进而训练出一个新模型, 得到最终输出结果。

2 个人信用评估模型的构建

2.1 数据不平衡处理

GiveMe Some Credit 数据集一共包含 15 万条样本数据, 其中违约数据 10 026 条, 占比 6.68%, 无违约数据 139 974 条, 占比 93.32%, 样本数据严重不平衡, 因而需要进行样本数据的平衡处理。但经典 SMOTE 算法随机选取少数类样本合成新样本, 忽略了周边样本情况, 容易导致新合成的少数类样本与周围多数类样本产生重叠, 形成较多噪音。数据清洗技术 ENN 广泛应用于重叠样本数据的处理, 因而可以将 SMOTE 算法与 ENN 结合起来形成一个 pipeline, 即先过采样再进行数据清洗, 以改善 SMOTE 算法处理数据的不足。该文采用 SMOTE+ENN 算法对少数类样本进行合成处理, 以提高分类器优化性能。

SMOTE+ENN 算法主要是在 SMOTE 过采样的基础上, 通过 ENN 算法清洗重叠数据, 达到均衡样本数

据的目的。具体步骤如下: 首先使用 SMOTE 算法生成新的少数类样本, 得到扩充数据集。然后对新数据集的每一个样本使用 KNN (一般取 $K=5$) 方法进行预测, 若预测结果和实际结果差异较大, 则剔除该样本。

2.2 多分类器集成学习

选取给定数据集的 75% 作为训练集, 剩余数据集作为测试集, 并使用训练集分别建立逻辑回归 LR、K 最近邻 (K-nearest neighbor, KNN)、决策树 DT、随机森林分类器 (random forest classifier, RFC)、引导聚集 Bagging、随机梯度下降 (stochastic gradient descent, SGD)、极端梯度提升 (extreme gradient boosting, XGB) 和集成学习 (由 KNN、Bagging、和 XGB 模型集成) 8 种模型。LR 是一种概率模型, 通过使用 Logistic 函数将实数值映射到 $[0, 1]$ 之间, 以逼近真实标记的对数概率; KNN 通过计算样本点与其近邻的 5 个样本数据距离, 从而达到分类的目的; DT 通过递归选择最优特征对训练数据进行分割, 使得各子数据集达到一个最好分类。RFC 是一种基于决策分类树的 Bagging 集成学习方法, 其输出结果由包含的各决策树输出类别的众数决定。Bagging 方法通常是独立并行学习多个弱学习器, 并按照某种确定性的平均过程将它们组合。SGD 每次随机选择一个样本, 不断更新模型参数, 从而使得目标函数达到极小值点。XGB 是一种极限提升树模型, 将许多树模型集成在一起, 从而形成一个很强的分类。

该文采用基于投票法的集成学习结合策略进行预测分类,该方法主要是先并行学习几个弱分类器,根据输出预测结果从中选取 3 个最优个体分类器作为基分类器,并通过硬投票法将基分类器结果集成,输出最终的预测结果。算法流程如图 3 所示。

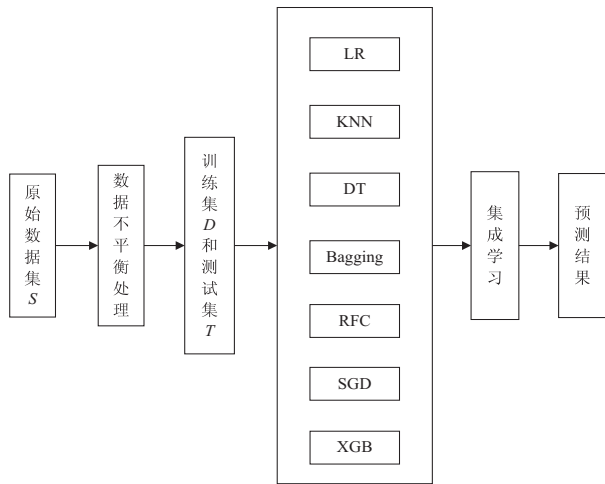


图 3 集成学习方法流程

(1) 将 Give Me Some Credit 数据集 S 划分为训练集 D (75%) 和测试集 T (25%)。

(2) 将训练集 D 按照五折交叉验证的方法随机均等划分为 D_1 、 D_2 、 D_3 、 D_4 和 D_5 五个子集,依次选取各子集 D_i ($i = 1, 2, 3, 4, 5$) 作为测试子集,剩下的 4 份作为训练子集。

(3) 网格搜索算法分别对 LR、KNN、DT、RFC、Bagging、SGD、XGB 模型自动交叉验证、训练,通过调节每一个模型超参数来跟踪评分结果,获取最优超参数。

(4) 利用测试集 T ,对所有最优超参数的个体分类器模型进行预测,将预测结果 AUC 评分可视化,结

果如图 4 所示。

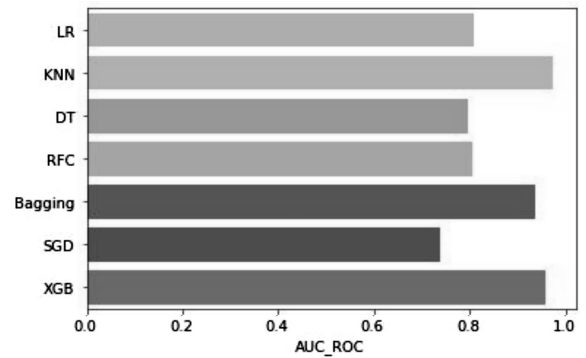


图 4 单一分类器可视化 AUC 评分结果

(5) 根据 AUC 评分结果,选择 KNN、Bagging 和 XGB 为基模型,并通过基于硬投票法的集成学习结合策略将基分类器结果集成新预测分类模型。用训练集 D 重新训练,测试集 T 验证模型效果。

3 实验设计

3.1 实验环境设置

所有实验运行环境均为:Windows 10 操作系统, Intel(R) Xeon(R) W - 2123 CPU@ 3.6 GHz 中央处理器, NVIDIA GeForce GTX1080 显卡,显存大小为 8 GB、Python 3.6 版本、Anaconda 3.6 集成开发环境。

3.2 数据采集与分析

实验数据集来自 Kaggle 数据算法比赛中的 Give Me Some Credit 开源数据集,该数据集包含 12 个变量,15 万条的样本数据。为了方便特征描述,将变量 SeriousDlqin2yrs 设置为 Y , Y 中的数据 1 表示违约,0 表示没有违约。其他特征变量按顺序依次设置为 $X_1 \sim X_{10}$ 。Give Me Some Credit 数据集相关特征变量信息如表 1 所示。

表 1 个人信用特征及其对应解释

变量名	变量描述	特征编号
SeriousDlqin2yrs	超过 90 天或者更糟的逾期拖欠	Y
RevolvingUtilizationOfUnsecuredLines	循环贷款无抵押额度	X_1
age	借款人当时的年龄	X_2
NumberOfTime30-59DaysPastDueNotWorse	35 ~ 59 天逾期但不糟糕次数	X_3
DebtRatio	负债比率	X_4
MonthlyIncome	月收入	X_5
NumberOfOpenCreditLinesAndLoans	开放式信贷和贷款数量,开放式贷款(分期付款如汽车贷款或抵押)	X_6
NumberOfTimes90DaysLate	90 天逾期次数;借款者有 90 天或更高逾期的次数	X_7
NumberRealEstateLoansOrLines	不动产贷款或额度数量;抵押贷款和不动产放款包括房屋净值信贷额度	X_8
NumberOfTime60-89DaysPastDueNotWorse	60 ~ 89 天逾期但不糟糕次数;借款人在在过去两年内有 60 ~ 89 天逾期还款但不糟糕的次数	X_9
NumberOfDependents	家属数量;不包括本人在内的家属数量	X_{10}

3.3 数据预处理

实验所采用的数据集是非结构化数据,因而包含多种噪声,需要对数据进行预处理操作。实验数据预处理主要包括缺失值处理、异常值处理、数据归一化操作。

经统计分析发现,特征变量 X_5 、 X_{10} 存在数据缺失,变量 X_1 、 X_2 、 X_4 、 X_5 、 X_6 、 X_8 存在不同程度的数据异常。本次实验为了处理噪声数据的干扰,采用前一个非缺失值去填充该缺失值,并将异常值视为缺失值,按照缺失值处理方法来进行处理,以保证原数据的均值和标准差不发生大的改动。

同时,为了消除特征变量之间的差异,提高模型预测准确度,需要对数据进行归一化处理。本实验采用最小最大值归一化方法,通过缩放将每一个特征的取值范围压缩到 $[0, 1]$ 或 $[-1, 1]$ 之间,从而减少变量 X_2 、 X_5 与其他特征变量间的差异。假设有 N 个样本 $\{X^{(n)}\}_{n=1}^N$, 对于每一维特征 X , 归一化的特征为:

$$\hat{X}^{(n)} = \frac{X^{(n)} - \min_n(X^{(n)})}{\max_n(X^{(n)}) - \min_n(X^{(n)})} \quad (4)$$

式中, $\min(X)$ 和 $\max(X)$ 分别是特征 X 在所有样本上的最小值和最大值。

3.4 超参数优化

为了确保实验有效进行,需要为每个分类器仔细设置超参数,如表 2 所示。网格搜索法作为一种自动参数寻优算法,根据给定模型自动进行交叉验证,通过调节每一个参数跟踪评分结果,由于其参数寻优性能的稳定可靠,被广泛应用于多个领域。

表 2 不同算法下网格搜索参数设置

算法	参数设置
LR	penalty: ['l1', 'l2'], C: [0.001, 0.01, 0.1, 1, 10]
KNN	n_neighbors: list(range(2, 5, 1)), algorithm: ['auto', 'ball_tree', 'kd_tree', 'brute']
DT	criterion: ['gini', 'entropy'], # 衡量标准 max_depth: list(range(2, 5, 1)), # 树的深度 min_samples_leaf: list(range(3, 7, 1)) # 最小叶子节点数
RFC	n_estimators: [100, 150, 200], # 多少棵树 criterion: ['gini', 'entropy'], # 衡量标准 max_depth: list(range(2, 5, 1)) # 树的深度
BAG	n_estimators: [10, 15, 20]
SGD	penalty: ['l2', 'l1'], max_iter: [1 000, 1 500, 2 000]
XGB	max_depth: [3, 4, 5, 6]

3.5 评价指标

为了衡量一个机器学习模型的好坏,通常采用准

确率 (Accuracy)、精确率 (Precision)、召回率 (Recall) 和 F 值 (F1-Score) 对分类模型预测效果进行评估,各个指标的具体描述如公式(5)~公式(8)所示。

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \quad (5)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (6)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (7)$$

$$\text{F1} = \frac{2\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (8)$$

式中, TP 表示预测正确的正样本数量; FP 表示预测错误的负样本数量; FN 表示预测错误的正样本数量; TN 表示预测正确的负样本数量。

3.6 实验测试与结果分析

3.6.1 SMOTE+ENN 算法数据不平衡处理

SMOTE+ENN 算法通过选择性地产出少数类样本,生成了质量相对较高的新样本集合,增强了分类算法性能。对比实验结果如图 5 所示。

从图 5 可以看出,与经典 SMOTE 算法相比, SMOTE+ENN 算法生成的新样本中,离群点数目明显减少,生成的新样本分布更加均匀。SMOTE+undersampling 作为欠采样与过采样结合的简单算法,其生成的新样本离群点数目较多,样本数据存在欠拟合。与之相比, SMOTE+ENN 算法生成的新样本拟合度更高,样本数据分布更佳。

经 SMOTE、SMOTE + undersampling、SMOTE + ENN 三种不同的数据不平衡算法处理,样本数据结果如表 3 所示。

表 3 不同算法下样本数据的数量分布

	SMOTE	SMOTE+ undersampling	SMOTE+ENN
0	139 974	27 994	110 433
1	139 974	13 997	129 873
总样本	279 948	41 991	240 306

SMOTE 算法只是单纯重复了违约样本,会过分强调已有的违约样本,如果部分样本点标记错误或者是噪音,那么 SMOTE 算法处理后的数据错误也会被成倍放大,造成违约样本过拟合。SMOTE + undersampling 算法抛弃了大部分正常样本数据,从而弱化了中间部分正常样本的影响,但样本数据的大量遗弃,可能会形成偏差很大的模型。SMOTE + ENN 算法与 SMOTE、SMOTE+undersampling 两种算法有着明显的不同,不仅是单纯地重复违约样本,而且在局部区域通过 KNN 生成了新的违约样本数据,降低了数据过拟合的风险。

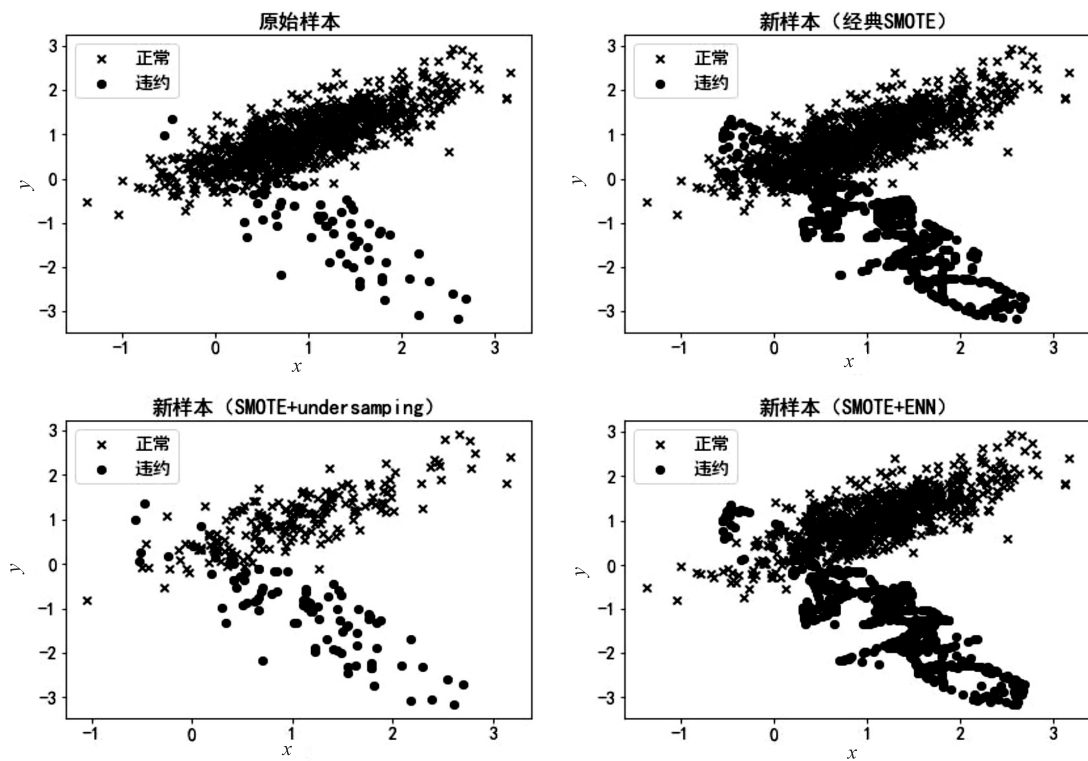


图 5 SMOTE+ENN 算法生成新样本分布对比

3.6.2 实验结果分析

针对 SMOTE、SMOTE + undersampling、SMOTE + ENN 三种不同的数据预处理算法,分别使用简单 5 折交叉验证、网格搜索和集成学习的算法预测模型分类准确率,具体实验结果如表 4 ~ 表 7 所示。

表 4 不同算法下简单交叉验证的准确率

模型	SMOTE	SMOTE+undersampling	SMOTE+ENN
LR	0.73	0.74	0.81
KNN	0.87	0.78	0.94
DT	0.87	0.73	0.89
RFC	0.90	0.81	0.93
BAG	0.91	0.79	0.92
SGD	0.74	0.76	0.80
XGB	0.94	0.82	0.96

表 5 不同算法下网格搜索优化模型的准确率

模型	SMOTE	SMOTE+undersampling	SMOTE+ENN
LR	0.66	0.74	0.80
KNN	0.75	0.77	0.97
DT	0.72	0.78	0.79
RFC	0.71	0.78	0.81
BAG	0.79	0.79	0.93
SGD	0.70	0.78	0.75
XGB	0.87	0.81	0.95

表 6 不同算法下集成学习的准确率

模型	SMOTE	SMOTE+undersampling	SMOTE+ENN
accuracy	0.95	0.81	0.97
Macro avg	0.95	0.80	0.97
Weighted avg	0.95	0.81	0.97

表 7 集成学习评估结果

	Precision	Recall	F1-score	Support
0	0.97	0.97	0.97	27 598
1	0.98	0.97	0.98	32 479
accuracy			0.97	60 077
Macro avg	0.97	0.97	0.97	60 077
Weighted avg	0.97	0.97	0.97	60 077

从表 4、表 5 的实验结果可以看出,在两种不同的信用评估方法下,SMOTE+ENN 算法相较于 SMOTE 和 SMOTE+undersampling 两种算法,个体分类器模型的预测准确率最高。在简单交叉验证的条件下,SMOTE+ENN 算法预测分类准确率最高,但由于简单交叉验证无任何参数优化,且产生的训练集数据分布和原始数据集有所不同,预测结果偏差较大,可信度低;网格搜索通过优化超参数,获取最优个体分类器模型,SMOTE+ENN 算法预测分类准确率相较于简单交叉验证虽有所降低,但可信度增加。

根据表 4 ~ 表 7 的评估结果,可以看出在不同算法下,集成学习模型的准确率远远高于个体分类器模

型。该文提出的 SMOTE+ENN 算法对数据不平衡处理后,相较于 SMOTE 和 SMOTE+undersampling 数据不平衡处理方法,无论是在个体分类器还是集成学习模型上,预测准确率均有大幅度提升。在基于网格搜索和集成学习的情况下,所提出的 SMOTE+ENN 算法能够在一定程度上改善了 SMOTE 算法的不足,预测精度提升了 2%,同时实验结果表明集成学习模型相较于个体分类模型而言,预测效果更好。

4 结束语

首先介绍了常用的信用评估方法,在相关研究的基础上,提出了一种应用于不平衡数据集样本处理的算法,即 SMOTE+ENN 算法。通过设计多种分类器、超参数优化以及集成学习,实现最优评估模型的构建,从而完成对个人信用的评估。为了验证该算法的可行性,针对 SMOTE 和 SMOTE+undersampling 两种算法进行实验对比分析,结果表明了基于 SMOTE+ENN 与集成学习方法的信用评估模型的稳健性。

参考文献:

- [1] DASTILE X, CELIK T, POTSANE M. Statistical and machine learning models in credit scoring: a systematic literature survey[J]. Applied Soft Computing, 2020, 91: 106263.
- [2] WU Yue, XU Yunjie, LI Jiaoyang. Feature construction for fraudulent credit card cash-out detection[J]. Decision Support Systems, 2019, 127: 113155.
- [3] MIHALOVIC M. Performance comparison of multiple discriminant analysis and logit models in bankruptcy prediction[J]. Economics & Sociology, 2016, 9(4): 101-118.
- [4] CHANG YUNG-CHIA, CHANG KUEI-HU, CHU HENG-HSUAN, et al. Establishing decision tree-based short-term default credit risk assessment models[J]. Communications in Statistics - Theory and Methods, 2016, 45(23): 6803-6815.
- [5] CHEN J, HUANG S. Evaluation model of green supply chain cooperation credit based on BP neural network[J]. Neural Comput & Applic. 2021, 33: 1007-1015.
- [6] 姜凤茹. 基于 GA-SVM 的网络借贷个人信用评估模型研究[J]. 控制工程, 2020, 27(6): 1025-1031.
- [7] 李太勇, 王会军, 吴江, 等. 基于稀疏贝叶斯学习的个人信用评估[J]. 计算机应用, 2013, 33(11): 3094-3096.
- [8] MOSCATO V, PICARIELLO A, SPERLÍ G. A benchmark of machine learning approaches for credit score prediction[J]. Expert Systems with Applications, 2021, 165: 113986.
- [9] SEZER O B, GUDELEK U, OZBAYOGLU M. Financial time series forecasting with deep learning: a systematic literature review: 2005-2019[J]. Applied Soft Computing, 2020, 90: 106181.
- [10] CHI G, YU S, ZHOU Y. A novel credit evaluation model based on the maximum discrimination of evaluation results[J]. Emerging Markets Finance and Trade, 2020, 56(11): 2543-2562.
- [11] 赵雪峰, 吴伟伟, 时辉凝. 基于自然语言处理与深度学习的信用贷款评估模型[J]. 系统管理学报, 2020, 29(4): 629-638.
- [12] SHEN Feng, ZHAO Xingchao, KOU Gang, et al. A new deep learning ensemble credit risk evaluation model with an improved synthetic minority oversampling technique[J]. Applied Soft Computing, 2021, 98: 106852.
- [13] CHAWLA N V, BOWYER K W, HALL L O, et al. Smote: synthetic minority over-sampling technique[J]. Journal of Artificial Intelligence Research, 2002, 16: 321-357.
- [14] SUN Jie, LANG Jie, FUJITA H, et al. Imbalanced enterprise credit evaluation with DTE-SBD: decision tree ensemble based on SMOTE and bagging with differentiated sampling rates[J]. Information Sciences, 2018, 425: 76-91.
- [15] ZHOU Zhihua, TANG Wei. Clusterer ensemble[J]. Knowledge-Based Systems, 2005, 19(1): 77-83.
- [16] YU L, ZHOU Rongtian, TANG Ling, et al. A DBN-based resampling SVM ensemble learning paradigm for credit classification with imbalanced data[J]. Applied Soft Computing, 2018, 69: 192-202.
- [17] KRAWCZYK B. Learning from imbalanced data: open challenges and future directions[J]. Progress in Artificial Intelligence, 2016, 5(4): 221-232.
- [18] TORRES-BARRÁN A, ALAÍZ C M, DORRONSORO J R. Faster SVM training via conjugate SMO[J]. Pattern Recognition, 2021, 111: 107644.