

## RESEARCH ARTICLE

WILEY

# On the relative value of imbalanced learning for code smell detection

Fuyang Li<sup>1</sup> | Kuan Zou<sup>1,2</sup> | Jacky Wai Keung<sup>3</sup> | Xiao Yu<sup>1,4,5</sup>  | Shuo Feng<sup>6</sup> | Yan Xiao<sup>7</sup>

<sup>1</sup>School of Computer Science and Artificial Intelligence, Wuhan University of Technology, Wuhan, China

<sup>2</sup>School of Computer, Electronics and Information, Guangxi University, Nanning, China

<sup>3</sup>Department of Computer Science, City University of Hong Kong, Hong Kong, China

<sup>4</sup>Sanya Science and Education Innovation Park of Wuhan University of Technology, Sanya, China

<sup>5</sup>Wuhan University of Technology Chongqing Research Institute, Chongqing, China

<sup>6</sup>School of Computer and Artificial Intelligence, Zhengzhou University, Zhengzhou, China

<sup>7</sup>School of Cyber Science and Technology, Shenzhen Campus, Sun Yat-sen University, Shenzhen, China

## Correspondence

Xiao Yu, Sanya Science and Education Innovation Park of Wuhan University of Technology, Sanya, China.  
Email: [xiaoyu@whut.edu.cn](mailto:xiaoyu@whut.edu.cn)

## Funding information

Project of Sanya Yazhou Bay Science and Technology City, Grant/Award Number: SCKJ-JYRC-2022-17; Natural Science Foundation of China, Grant/Award Number: 62272356; Youth Fund Project of Hainan Natural Science Foundation, Grant/Award Number: 622QN344; Natural Science Foundation of Chongqing, Grant/Award Number: cstc2021jcyj-msxmX1115; Start-up Grant from Wuhan University of Technology, Grant/Award Number: 104-40120693

## Summary

Machine learning-based code smell detection (CSD) has been demonstrated to be a valuable approach for improving software quality and enabling developers to identify problematic patterns in code. However, previous researches have shown that the code smell datasets commonly used to train these models are heavily imbalanced. While some recent studies have explored the use of imbalanced learning techniques for CSD, they have only evaluated a limited number of techniques and thus their conclusions about the most effective methods may be biased and inconclusive. To thoroughly evaluate the effect of imbalanced learning techniques for machine learning-based CSD, we examine 31 imbalanced learning techniques with seven classifiers to build CSD models on four code smell data sets. We employ four evaluation metrics to assess the detection performance with the Wilcoxon signed-rank test and Cliff's  $\delta$ . The results show that (1) Not all imbalanced learning techniques significantly improve detection performance, but deep forest significantly outperforms the other techniques on all code smell data sets. (2) SMOTE (Synthetic Minority Over-sampling TEchnique) is not the most effective technique for resampling code smell data sets. (3) The best-performing imbalanced learning techniques and the top-3 data resampling techniques have little time cost for code smell detection. Therefore, we provide some practical guidelines. First, researchers and practitioners should select the appropriate imbalanced learning techniques (e.g., deep forest) to ameliorate the class imbalance problem. In contrast, the blind application of imbalanced learning techniques could be harmful. Then, better data resampling techniques than SMOTE should be selected to preprocess the code smell data sets.

## KEYWORDS

code smell detection, empirical software engineering, imbalanced learning, machine learning

## 1 | INTRODUCTION

As the scale of software systems becomes increasingly large and complex, software quality attracts the attention of researchers in the software engineering domain.<sup>1-8</sup> One important aspect of maintaining software quality is the detection of poor programming practices or design weaknesses, known as “Code Smells”.<sup>9,10</sup> Detecting code smells is crucial in supporting software developers to improve software design and quality<sup>11-13</sup> and the existence of code smells will negatively impact key software quality attributes such as reliability and changeability.<sup>14,15</sup> With the growing complexity in software systems, there is a need for innovative approaches and techniques to accurately detect the smelly code blocks. Recently, several approaches to code smell detection (CSD) have been proposed, which include two categories, that is, heuristic-based and machine learning-based. The former<sup>16-18</sup> can be broadly classified into two classes: metrics-based approaches and rule-based approaches. In metrics-based approaches, quality metrics are defined and threshold values are established for each metric. However, selecting appropriate threshold values can be challenging in this approach.<sup>19</sup> In contrast, rule-based approaches involve defining rules to identify code smells. These rules may be manually generated by domain experts.<sup>20</sup> However, both of these approaches can be time-consuming and cognitively demanding for software engineers,<sup>21</sup> leading to a shift towards the use of machine learning approaches.

To address the limitations, researchers have proposed to utilize machine learning techniques for CSD. The machine learning-based approaches first employ some binary classification algorithms to train a detection model based on extracted code smell features and then utilize the built model to predict whether a code block is smelly or not. It means that the existing studies<sup>13,21-23</sup> in code smell detection typically approach the problem as a binary classification task. These approaches<sup>24-26</sup> do not require experts to define heuristic rules and decide thresholds since they perform these tasks internally.

The recent studies<sup>20,22,27,28</sup> have raised the concern that machine learning-based detection models that are trained on imbalanced class data sets (i.e., the existence of more non-smelly instances than smelly instances) may lead to inaccurate prediction results. Therefore, some research works have studied the impact of some imbalanced learning techniques for CSD. For example, Alkharabsheh et al.<sup>13</sup> and Pecorelli et al.<sup>23,29</sup> found that synthetic minority over-sampling technique (SMOTE) could not significantly improve detection performance. Alazba et al.<sup>21</sup> and Aljamaan<sup>30</sup> showed that the stacking ensemble and voting ensemble are better than individual classifiers.

However, they only compare a limited number of imbalanced learning techniques, and many other imbalanced learning techniques have not been investigated. Therefore, their conclusions about the best-performing imbalanced learning techniques are biased and still do not have enough conviction.

Considering the aforementioned issues, we conduct a comprehensive investigation to study the practical impact of 31 imbalanced learning techniques on the detection performance of CSD with seven machine learning classifiers to build detection models. The 31 imbalanced learning techniques can be divided into four families, that is, (1) data resampling techniques, (2) ensemble learning techniques, (3) cost-sensitive learning techniques, and (4) imbalanced ensemble learning techniques. The seven machine learning classifiers divided into the five groups, that is, statistic-based (logistic regression and naive Bayes), support vector machine-based (support vector machine), decision tree-based (decision tree and random forest), nearest neighbor-based (K-nearest neighbour), and neural networks-based (multi-layer perceptron). Similar to Nucci et al.,<sup>22</sup> we merge the same level code smell data sets to construct new data sets containing more than one type of code smells. But we further remove redundant and conflicting instances in our data sets. We use the four performance metrics (i.e., Precision, Recall, F1, and Matthews correlation coefficient [MCC]) to comprehensively evaluate the effects of the above-mentioned imbalanced learning techniques with each classifier on our data sets. In addition, we apply both the Wilcoxon signed-rank test<sup>31</sup> and the Cliff's  $\delta$ <sup>32</sup> to examine the difference of detection performance between the imbalanced learning techniques and None (without imbalanced learning). Our experimental results show that:

- (1) On the original code smell data sets, the classifiers except naive Bayes have great performance for CSD and score at least 0.8 on Precision, Recall, F1, and MCC. However, the Precision, Recall, F1, and MCC values of the best-performing classifiers are reduced by 9.31%–15.66%, 10.04%–18.95%, 10.45%–21.21%, and 28.36%–56.34% on our data sets, respectively. Different classifiers achieve the best detection performance on different code smells across our data sets, for example, naive Bayes on Data Class and God Class, and logistic regression on Feature Envy and Long Method.
- (2) Compared with None (without any imbalanced learning), there are significant performance differences among the imbalanced learning techniques. Not all techniques can enhance the detection performance, namely, 32.26%,

29.03%, 30.65%, and 34.68% techniques show a significant positive effect for CSD on our data sets in terms of Precision, Recall, F1, and MCC, respectively. Deep forest always achieves the best performance on code smell data sets and shows statistically significant improvement of 9.99%–26.18% and 31.31%–105.81% than None in terms of F1 and MCC.

- (3) We are not able to achieve better CSD performance using synthetic minority over-sampling technique (SMOTE) over other data resampling techniques. The best data resampling techniques obtain the detection improvements than SMOTE ranging from 2.63% to 17.73% in terms of MCC.
- (4) The best-performing imbalanced learning method, that is, deep forest, has a relatively low computational time of approximately 2s. The top-3 data resampling techniques also have low computational times, with a maximum of less than 0.3s. However, the computational times for these top-3 data resampling techniques are generally higher for Feature Envy and Long Method compared to Data Class and God Class.

The contributions of our work can be concluded as following:

1. We first conduct such comprehensive empirical research to explore the practical effect of 31 imbalanced learning techniques on the detection performance of machine learning-based CSD.
2. We identify and make an analysis of a set of 18 machine learning-based CSD researches with imbalanced learning techniques from different perspectives, including code smell types, classifiers, data set metrics, and the used imbalanced learning techniques. Researchers are able to utilize the set as a starting point to conduct subsequent investigations of CSD with imbalanced learning.

The remainders of the paper are organized as follow: Section 2 introduces the related work on machine learning-based CSD and imbalanced learning techniques for CSD. Section 3 gives a brief description of the investigated imbalanced learning techniques on CSD. Sections 4 and 5 show the details of our experiment setup and provide the experimental results and research summaries in detail. Section 6 shows the differences between our results and previous results, discusses the threats to the validity, and gives some implications from the experimental results. Section 7 presents the conclusion of our research.

## 2 | LITERATURE REVIEW

### 2.1 | Machine learning-based CSD

Kreimer<sup>33</sup> proposed to use the decision tree (DT) to detect the two code smells (i.e., Blob and Long Method) in two small systems, and the detection model achieved high accuracy. Khomh et al.<sup>34,35</sup> and Vaucher et al.<sup>36</sup> investigated the feasibility of Bayesian belief networks (BN) to detect different code smells. Bryton et al.<sup>37</sup> proposed a Binary Logistic Regression model to detect the Long Method smell. Maiga et al.<sup>38</sup> employed the support vector machine (SVM) method to detect the Blob smell incrementally. Amorim et al.<sup>39</sup> confirmed Kreimer's findings<sup>33</sup> on four medium-scale projects. Fontana et al.<sup>20,40</sup> conducted empirical research on identifying code smells. In their work, the six machine learning classification models (i.e., DT, naive Bayes (NB), random forest (RF), sequential minimal optimization (SMO), LibSVM, and JRip) were trained to detect the four code smells (i.e., Data Class, God Class, Feature Envy, and Long Method). Their results showed that the RF classifier obtained better performance in detecting almost all code smells. Then, they<sup>27</sup> extended their work by considering the different severity of code smells. Nucci et al.<sup>22</sup> conducted a replicated study of Fontana et al.<sup>20</sup> They used the same experimental setup but merged the data sets of Fontana et al. to make the data sets more realistic. Their results suggested that the conclusions of previous studies<sup>20,40</sup> cannot be generalized in real-world scenarios, and the effect of machine learning techniques on CSD needs to be further explored. Kim<sup>41</sup> and Liu et al.<sup>42</sup> used neural networks for CSD with good results. Pecorelli et al.<sup>43</sup> conducted an empirical investigation to compare the detection performance of machine learning-based (i.e., NB) and heuristic-based methods. Their results presented that the accuracy and validity of these two methods for detecting code smell still need further investigation. Some studies<sup>44–48</sup> tried to apply deep learning techniques for CSD, and their conclusions examined that some deep learning models accept a better performance on CSD, that is, convolutional neural networks (CNN-1),<sup>44,48</sup> recurrent neural network,<sup>44</sup> long short-term memory,<sup>45,47</sup> residual network,<sup>46</sup> and attention.<sup>48</sup>

## 2.2 | Imbalanced learning for CSD

A frequently encountered problem is that the code smell data consists of only a few smelly instances and many non-smelly instances. The imbalanced data distribution makes the built detection model prone to predict the new instances to be non-smelly and consequently perform poorly in finding smelly instances. Therefore, some researchers have investigated whether imbalanced learning methods can alleviate the problem and improve the performance of CSD models. To have knowledge of the research progress of machine learning-based CSD with imbalanced learning techniques, we perform a literature investigation to obtain most related articles published between 2005 to May 2023\*. To ensure that our literature review is comprehensive, we have established the following inclusion criteria: (1) the article is about CSD with imbalanced learning techniques and is written in English. (2) the full text is available. As far as we know, the first article that adopted machine learning techniques to detect code smells is published by Kreimer et al.<sup>33</sup> in *Electronic Notes in Theoretical Computer Science* 2005 (abbreviated as Kreimer 2005 article in the rest of this paper). Thus, the starting year of the literature investigation is set as 2005. We retrieve related papers written in English using Google Scholar and DBLP. We follow Zhou et al.<sup>62</sup> forward snowballing search method to recursively retrieve the papers that cited the Kreimer 2005 article. More precisely, we first retrieve the related papers that cited the Kreimer 2005 article and employ machine learning algorithms for CSD, and repeat the procedure on other related papers. Then, we set up “imbalanced learning” + “code smell detection” and “class rebalancing” + “code smell detection” as the search terms to search. Finally, after applying our inclusion criteria, we find the relevant papers that cited in these approaches. By the above steps, we totally identify 18 relevant papers about imbalanced learning techniques on machine learning-based CSD from the two sources. Table 1 lists the summary of the literature review, where the first column presents the references of studies, the second column presents the types of code smells, the third column presents the used machine learning algorithms, and the fourth column presents the code smell features, and the last one presents the imbalanced learning methods used in the study.

Some researchers<sup>50-53,55-61</sup> applied SMOTE as a data preprocessing technique to alleviate the class imbalance problems, and utilized or proposed some more advanced algorithms for CSD. For example, Akhter et al.<sup>50</sup> used the four machine learning classifiers (i.e., NB, RF, DT, and SVM) to investigate the effect of machine learning techniques on CSD. Alkharabsheh et al.<sup>51</sup> extended two software metrics project domain and size category to detect code smells by using the eight machine learning classifiers (i.e., LibSVM, IBK, DT, JRip, SMO, NB, RF, and random committee (RC)). Gupta et al.<sup>52</sup> employed eight deep learning models for CSD. Jain et al.<sup>53</sup> used hybrid feature selection and ensemble learning techniques to improve detection performance by adopting the 11 classifiers (i.e., SVM, KNN, NB, DT, linear discriminant analysis (LDA), logistic regression (LR), Bagging, AdaBoost, XGBoost, gradient boost (GB), and Stacking). Stefano et al.<sup>55</sup> proposed a cross-project method, which used the four machine learning classifiers (i.e., LR, DT, NB, and RF) to train detection models and predicted the smelliness of within-project instances. Khleel et al.<sup>56</sup> employed CNN-1 with SMOTE for CSD. Kovavcevic et al.<sup>57</sup> used code embeddings (i.e., code2vec, code2sep, and CuBERT) and over-sampling strategies (i.e., SMOTE and SMOTEENN) to detect code smells. Nanda et al.<sup>58</sup> proposed SSHM method with SMOTE and Stacking for severity CSD of four code smells. Yedida et al.<sup>59</sup> studied the effect of the fuzzy oversampling technique with SMOTE on CSD. Shen et al.<sup>49</sup> used RUS to solve the imbalanced problem across six machine learning classifiers (i.e., DT, NB, RF, SVM, K-nearest neighbour (KNN), and Rule ensemble (RULE)) when they analyzed the influence of hyper-parameter optimization on CSD. Patnaik et al.<sup>54</sup> used RUS and ROS to handle the imbalanced problem in their data sets and developed a model for CSD by using the seven classifiers (i.e., DT, RF, SVM, NB, Ridge, least absolute shrinkage and selection operator (LASSO), and least angle regression (LAR)). However, the above-mentioned studies do not compare the performance difference of CSD models trained on the original data sets and the balanced data sets processed by SMOTE, RUS, and ROS. Therefore, no conclusive empirical evidence from their experimental results show that using SMOTE, RUS, and ROS can significantly positively impact machine learning-based CSD models.

Alkharabsheh et al.<sup>13</sup> used the machine learning classifiers (i.e., LDA, quadratic discriminant analysis (QDA), NB, Multi-Layer Perceptron (MLP), SVM, DT, GB, CatBoost, light gradient boosting machine (LGBM), XGBoost, XGBoost with random forest (XGBRF), AdaBoost, Bagging, RF, extra trees (ET), KNN, nearest centroid (NC), Gaussian process (GP), Ridge, LR, Perceptron, passive aggressive (PA), and stochastic gradient descent (SGD)) to compare whether using SMOTE would improve the detection performance on God Class detection. Their results showed that SMOTE could not improve the God Class detection performance.

Alazba et al.<sup>21</sup> studied the effect of the stacking ensemble on six code smells, and their conclusions examined that the performance of Stacking with LR and SVM is better than all individual classifiers. Aljamaan<sup>30</sup> investigated the

\*The literature investigation was conducted in May 2023.

TABLE 1 The summary of related works.

Study	Code smell type	Classifiers	Data set metrics	Imbalanced learning methods
29	Class-level: 4 Method-level: 1	NB	9 Object-Oriented Metrics	ClassBalancer, Resample, SMOTE, Cost-Sensitive Classifier, One-Class Classifier
23	Class-level: 5 Method-level: 6	NB	17 Object-Oriented Metrics	Over-Sampling, Under-Sampling, SMOTE, Cost-Sensitive Classifier, One-Class Classifier
49	Class-level: 1 Method-level: 1	DT, KNN, NB, RF, RULE, SVM	11 Object-Oriented Metrics	RUS
50	Class-level: 3 Method-level: 1	SVM, NB, RF, DT	50 Object-Oriented Metrics	SMOTE
21	Class-level: 2 Method-level: 4	DT, SVM, NB, LR, MLP, SGD, GP, KNN, LDA	55-82 Object-Oriented Metrics	Stacking Ensemble
51	Class-level: 1	LibSVM, IBK, DT, JRip, SMO, NB, RC, RF	19 Object-Oriented Metrics	SMOTE
30	Class-level: 2 Method-level: 4	DT, LR, SVM, MLP, SGD	61-82 Object-Oriented Metrics	Voting Ensemble
52	Class-level: 4 Method-level: 4	8 Deep Learning Models	Dimensional Metrics Complexity Metrics Object-Oriented Metrics Android-Oriented Metrics	SMOTE, ADASYN
53	Class-level: 2 Method-level: 4	SVM, KNN, NB, DT, LDA, LR, Bagging, AdaBoost, XGBoost, GB	61-82 Object-Oriented Metrics	SMOTE
54	Class-level: 3 Method-level: 2	LASSO, Ridge, LAR, NB, SVM, RF, DT	45 Object-Oriented Metrics	Random Sampling
55	Class-level: 2 Method-level: 1	LR, DT, NB, RF	6 Object-Oriented Metrics	SMOTE
13	Class-level: 1	LDA, QDA, NB, MLP, SVM, DT, GB, CatBoost, LGBM, XGB, XGBRF, AdaBoost, Bagging, RF, ET, KNN, NC, GP, Ridge, LR, Perceptron, PA, SGD	16 Object-Oriented Metrics	SMOTE
56	Class-level: 2 Method-level: 2	CNN-1	43 Object-Oriented Metrics	SMOTE
57	Class-level: 1 Method-level: 1	code2vec, code2sep, CuBERT	10 Object-Oriented Metrics	SMOTE, SMOTEENN
58	Class-level: 2 Method-level: 2	DT, RF, SVM, JRip, NB	63-84 Object-Oriented Metrics	SMOTE
59	Class-level: 2 Method-level: 2	Deep Neural Network	Structural Information Lexical Information	SMOTE
60	Class-level: 9	CatBoost, RF, DT, ET, LR, KNN, GBM, XGB, LDA, AdaBoost, LBGM, NB, Dummy, SVM, QDA	65 Object-Oriented Metrics	SMOTE
61	Class-level: 3 Method-level: 2	NB, RF, J48	Object-Oriented Metrics	SMOTE



performance of the voting ensemble on CSD. Their results showed that the voting ensemble has a better detection performance on all code smells.

Pecorelli et al.<sup>29</sup> conducted a investigation to examine the role of imbalanced learning techniques (ClassBalancer, Resample, SMOTE, Cost-Sensitive Classifier, and One-Class Classifier) on machine learning-based CSD. Then, they<sup>23</sup> combined five imbalanced learning techniques (Over-Sampling, Under-Sampling, SMOTE, Cost-Sensitive Classifier, and One-Class Classifier) with a NB classifier to further investigate this issue. Their results examined that the performance of CSD models has not been significantly improved using imbalanced learning techniques, even though the performance has a slight improvement when using SMOTE.

However, the limitations of the past research are that only a few imbalanced learning techniques be used, and the data sets used for experiments need to be improved. To overcome these limitations of past research and their findings, in our study, we conduct a larger empirical research to examine the practical effect of the detection performance of CSD on machine learning algorithms using imbalanced learning techniques.

### 3 | PRELIMINARIES

We briefly outline the studied imbalanced learning techniques. We also regard the None method as a baseline method (without any imbalanced learning). We apply 31 imbalanced learning techniques in this empirical study, as shown in Table 2. The 31 techniques are widely adopted in previous CSD studies<sup>13,21,23,63</sup> and contain the four families, that is, nine data resampling techniques, eight ensemble learning techniques, five cost-sensitive learning techniques, and nine imbalanced ensemble learning techniques. The techniques we selected cover the algorithms investigated in these two researches.<sup>13,29</sup>

#### 3.1 | Data resampling

Data resampling techniques can mainly divide into two categories, that is, over-sampling and under-sampling. The former produces a superset of the original code smell data sets by duplicating existing smelly instances or creating new smelly instances from existing smelly ones, while the latter produces a subset of the original code smell data sets by eliminating non-smelly instances. To maintain consistency with previous studies<sup>13,23</sup> and common practices, we set the default smelly ratio to 0.5, resulting in an equal number of smelly and non-smelly instances in the balanced data sets.

(1) Random over-sampling (ROS) randomly replicates  $n$  smelly instances from the original code smell data sets ( $n$  is calculated according to the ratio value of the smelly instances) and then adds them into the original data sets to obtain balanced data sets.

(2) Synthetic minority over-sampling technique (SMOTE) first randomly selects  $n$  smelly instances from the original code smell data sets. For each smelly instance  $M_A$ , we randomly choose one of its  $k$  nearest neighbors using the K-nearest neighbors (KNN) algorithm, that is,  $M_B$ . The feature of the synthetic instance  $M_{synthetic}$  is

$$\mathbf{x}_{synthetic} = \mathbf{x}_A + \text{Random}(0, 1) \times (\mathbf{x}_B - \mathbf{x}_A), \quad (1)$$

where  $\mathbf{x}_{synthetic}$ ,  $\mathbf{x}_A$ , and  $\mathbf{x}_B$  are the feature of  $M_{synthetic}$ ,  $M_A$ , and  $M_B$ , respectively. Finally, we add the  $n$  synthetic smelly instances into the original data sets to obtain balanced data sets.

(3) Adaptive synthetic sampling approach (ADASYN) is an extension of SMOTE, but it creates new synthetic instances near the boundary between two classes rather than within the smelly instances.

(4) Borderline-SMOTE (BSMOTE) is an improved oversampling algorithm based on SMOTE, which only uses the smelly instances on the boundary to synthesize new instances.

(5) Random under-sampling (RUS) randomly eliminates  $m$  non-smelly instances to balance the class distribution ( $m$  is calculated according to the ratio value of the smelly instances).

(6) Near miss (NM) calculates the distance between the smelly and non-smelly instances and randomly deletes the non-smelly instances according to the distance, mainly to alleviate the information loss problem in RUS.

(7) Condensed nearest neighbor (CNN) first finds a subset of non-smelly instances that leads to no loss in model performance and then deletes the instances from the original code smell data sets.

(8) Tomek links (TL) is a modification from CNN to obtain balanced data sets by finding all the non-smelly instances closest to the smelly instances and then removing them.

TABLE 2 The overview of the 31 imbalanced learning methods.

Family	Methods	Abbreviation
Data Resampling Techniques	Random Over-Sampling	ROS
	Synthetic Minority Over-Sampling Technique	SMOTE
	Adaptive Synthetic Sampling Approach	ADASYN
	Borderline-SMOTE	BSMOTE
	Random Under-Sampling	RUS
	Near Miss	NM
	Condensed Nearest Neighbor	CNN
	Tomek Links	TL
	Edited Nearest Neighbors	ENN
Ensemble Learning Techniques	Bootstrap aggregating	Bagging
	Adaptive Boosting	AdaBoost
	CatBoost	CatBoost
	eXtreme Gradient Boosting	XGBoost
	Deep Forest	DF
	Stacking Logistic Regression	StackingLR
	Stacking Decision Tree	StackingDT
	Stacking Support Vector Machine	StackingSVM
Cost-Sensitive Learning Techniques	AdaCost	AdaCost
	Asymmetric Boosting	AsymBoost
	AdaUBoost	AdaUBoost
	Cost-Sensitive Support Vector Machine	CSSVM
	Cost-Sensitive Decision Tree	CSDT
Imbalanced Ensemble Learning Techniques	SMOTEBoost	SMOTEBoost
	ROSBoost	ROSBoost
	SMOTEBagging	SMOTEBagging
	ROSBagging	ROSBagging
	RUSBoost	RUSBoost
	RUSBagging	RUSBagging
	EasyEnsemble Classifier	EasyEnsemble
	BalanceCascade Classifier	BalanceCascade
	Balanced Random Forest	BRF

(9) Edited nearest neighbors (ENN) selects a non-smelly instance and uses the KNN algorithm to get the  $k$  neighbors of the instance. If more than half of the  $k$  neighbors are not the non-smelly instances, the instance will be deleted.

### 3.2 | Ensemble learning

Ensemble learning is regarded as a meta-algorithm, which is not a single machine learning technique but finishes the learning task by building and combining multiple machine learners to improve the performance of CSD.

(1) Bootstrap aggregating (Bagging) trains multi classification models based on multi bootstrap samples and form a final stronger classifier by voting their individual predictions.

(2) Adaptive Boosting (AdaBoost) integrates multiple base classifiers and assigns new weights to the samples misclassified by the previous base classifier to reduce the error rate.

(3) CatBoost is a gradient boosting decision tree framework with few parameters and high accuracy based on an oblivious tree algorithm, which could efficiently and reasonably process categorical features and has strong generalization ability.

(4) eXtreme Gradient Boosting (XGBoost) is a gradient boosting decision tree framework that provides a novel tree learning algorithm for processing sparse data and enables faster learning of machine learning classifiers through parallel and distributed computing, thus enabling faster model exploration.

(5) Deep forest (DF) integrates different kinds of forests in width and depth, which improves the classification ability of machine learning models by using multi-grained scanning and cascade forests.

(6) Stacking Ensemble uses a meta-learning algorithm to best integrate the predictors from two or more base well-performing base classification models. The three based classifiers commonly used in stacking ensembles are logistic regression (LR), decision tree (DT), and support vector machine (SVM). We call the stacking ensemble integrated with the three classifiers StackingLR, StackingDT, and StackingSVM.

In our experiment, the base classifier utilized in the AdaBoost, Bagging, XGBoost, and CatBoost models is decision tree, while the base classifier used in deep forest is random forest. On the other hand, the Stacking model is based on a combination of seven base classifiers (k-nearest neighbors, support vector machine, logistic regression, decision tree, naive Bayes, random forest, and multi-layer perceptron), which are then combined with three meta-classifiers (support vector machine, logistic regression, and decision tree).

### 3.3 | Cost-sensitive learning

In actuality, non-smelly instances are more common than smelly instances.<sup>21</sup> As a result, when utilizing machine learning classifiers for classification, instances are more likely to be misclassified as non-smelly due to the disproportionate representation of non-smelly instances. Cost-sensitive learning techniques aim at building code smells detection models with minimal misclassification costs by specifying a cost-sensitive matrix.

(1) AdaCost applies the misclassification cost to redistribute the training data sets on successive boosting rounds. The main idea of AdaCost is to comprise the cost and produce more advanced machine learning classifiers, which can reduce the cost of misclassifications better than AdaBoost.

(2) Asymmetric Boosting (AsymBoost) combines AdaBoost with cost-sensitive learning techniques and uses asymmetric misclassification costs to update the data distribution during classifier training to reduce the likelihood of misclassification.

(3) AdaUBoost is a variant of AdaBoost and is designed to optimize an unequal loss on imbalanced training data sets by preprocessing and manipulating the data distribution during classifier training to reduce the cumulative misclassification cost.

(4) Cost-sensitive support vector machine (CSSVM) weighs the margin by integrating the misclassification cost during the training of SVM.

(5) Cost-sensitive decision tree (CSDT) incorporates the misclassification cost into the process of separating software instances to two groups during the training of the decision tree.

### 3.4 | Imbalanced ensemble learning

Imbalanced ensemble learning techniques resample or reweight the training data sets in the training process of ensemble learning, which could improve the correct classification ability of the ensemble learning classifier for imbalanced data sets.

(1) SMOTEBoost iteratively builds  $t$  weak classification models similar to AdaBoost. In each iteration, the weak classification model is trained based on the balanced data sets generated by SMOTE.

(2) ROSBoost is similar to SMOTEBoost but uses ROS instead of SMOTE.

(3) RUSBoost is similar to SMOTEBoost but uses RUS instead of SMOTE to achieve data balance.

(4) SMOTEBagging involves SMOTE in the process of Bagging to generate balanced bootstrap samples and trains multi classification models based on the balanced data sets.



- (5) ROSBagging is similar to SMOTEBagging but uses ROS instead of SMOTE.
- (6) RUSBagging is similar to SMOTEBagging but uses RUS instead of SMOTE to achieve data balance.
- (7) EasyEnsemble Classifier builds an ensemble of AdaBoost predictors trained on multi-balanced bootstrap samples generated by RUS.
- (8) BalanceCascade Classifier iteratively drops non-smelly instances that were already well-classified by the current ensemble. After that, it performs RUS on the remaining non-smelly instances and trains a new base predictor.
- (9) Balanced random forest (BRF) trains the random forest predictor based on multi-balanced bootstrap samples generated by RUS.

## 4 | EXPERIMENTAL SETUP

### 4.1 | Data sets

In our study, we use the same experimental data sets as Nucci et al.,<sup>22</sup> Aljamaan,<sup>30</sup> Jain et al.,<sup>53</sup> and Nanda et al.,<sup>58</sup> which are built by Fontana et al.<sup>20</sup> and contain four code smells collected from 74 software systems. Each type of code smell data sets includes 140 smelly instances and 280 non-smelly ones (420 instances in total). The definitions of the four code smell data sets utilized in our study, namely Data Class, God Class, Feature Envy, and Long Method, are outlined as follow:

Data Class (class-level) is characterized by the existence of a class that primarily stores data with limited functionality. This class represents a deviation from the principles of object-oriented design, where classes are expected to embody both data and behavior, and can be accessed by other classes through simple getter and setter methods and is notable for its large number of properties and access points. A Data Class smell can lead to code redundancy, as it may result in multiple classes that have the same or similar functionality.

God Class (class-level), also referred to as Blob, refers to a class with centralized system behavior represented by multiple classes. This class is characterized by a large number of codes, methods, and properties and has a tendency to rely on data from other classes rather than its own. A God Class smell can make the system more difficult to modify or maintain as it tends to have many interdependent features and a complex control flow.

Feature Envy (method-level) pertains to methods displaying a greater attraction towards other classes as opposed to the class they belong to. This code smell is characterized by methods that extensively utilize the properties of other classes, often resulting in code that is difficult to maintain, understand, and modify. A Feature Envy smell can be caused by a design flaw in the software system, where the responsibilities of a class are not adequately distributed.

Long Method (method-level) is characterized by its large size and centralization of class operations. It encompasses a significant amount of code in order to perform multiple functions and is often considered complex and difficult to comprehend. A Long Method smell can start to become a problem when it grows to be excessively large and complex, making it difficult to understand and modify.

But the original data sets are questioned by Nucci et al.<sup>22</sup> Their replicated investigation found that a trained machine learning classifier could easily identify smelly and non-smelly instances since their feature distribution is very different in the original data sets. In addition, each data set only contains one type of code smells, which does not correspond to real-world scenarios (Nucci et al.<sup>22</sup> pointed that a software instance usually has two or more types of smells instead of one). Therefore, they merged the same level (i.e., class-level and method-level) of code smells to construct new data sets. For example, they merged the God Class and Data Class data sets, then set the label of all instances in the Data Class data set as no-smelly, and finally returned a new God Class data set.

Alazba et al.<sup>21</sup> raised the question about Nucci's behavior<sup>22</sup> of merging the same level code smell data sets. They showed that this merging strategy would result in data sets with more than 30% redundant instances and more than 15% conflicting instances, which would degrade the performance of machine learning classifiers. Therefore, they used the original data sets in their study.

**Example 1.** Suppose a God Class data set contains 2 smelly instances (i.e., *a* and *b*) and 3 non-smelly instances (i.e., *c*, *d*, and *e*). A Data Class data set contains 2 smelly instances (i.e., *b* and *f*) and 3 non-smelly instances (i.e., *e*, *g*, and *h*). If we adopt the Nucci et al.<sup>22</sup> merging strategy, we will return a new God Class data set, which contains 2 smelly instances (i.e., *a* and *b*) and 8 non-smelly instances (i.e., *c*, *d*, *e*, *b*, *f*, *e*, *g*, and *h*). From Alazba et al.<sup>21</sup> point of view, *e* is the redundant instance and *b* is the conflicting instance in the new God Class data set.

TABLE 3 The details of the our data sets.

Code smells	Level	Number of instances		
		Smelly	Non-smelly	Total
God class	Class	140	559	699
Data class		138	561	699
Feature envy	Method	140	378	518
Long method		140	378	518

TABLE 4 The confusion matrix.

	Actual smelly	Actual non-smelly
Predicted smelly	TP	FP
Predicted non-smelly	FN	TN

To avoid these issues, we use the original data sets provided by Fontana et al.,<sup>20</sup> and the merging strategy of Nucci et al.<sup>22</sup> is also used on the original data sets. But we remove the same instances so the redundant and conflicting instances would not appear in the new data sets. For example, we delete *b* and *e* from Data Class, then merge both the code smell data sets to create a new God Class data set, which contains 2 smelly instances (i.e., *a* and *b*) and 6 non-smelly instances (i.e., *c*, *d*, *e*, *f*, *g*, *h*). Therefore, we can obtain the experimental data sets that contain more than one type of code smells. Table 3 shows the detailed information of our experimental data sets.

## 4.2 | Performance measures

In our empirical study, we use the three threshold-dependent evaluation metrics (Precision, Recall, and F-measure (F1)) and one threshold-independent evaluation metric (Matthews correlation coefficient, MCC) to evaluate the performance of CSD models. The metrics are widely used in both software engineering studies<sup>64-71</sup> and artificial intelligence researches.<sup>72-75</sup> In the binary classification problem, these four evaluation metrics can be calculated according to a confusion matrix, as shown in Table 4.

**Precision** is the percentage of the actual smelly instances to all the predicted smelly instances in the data sets. A higher Precision means that the CSD model could help the software testing team more exactly to find smelly instances.

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

**Recall** is the percentage of the correctly predicted smelly instances to all the actual smell ones in the data sets. A higher Recall means that the software testing team could capture more smelly instances.

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

**F1** is viewed as the harmonic mean between Precision and Recall and its value ranges from 0 to 1, where 0 implies the worst prediction, and 1 implies the best prediction. In other words, a CSD model achieves better performance if it has a higher F1 value.

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

**MCC** provides an overview by taking into account all the terms of the confusion matrix, which could be used in an imbalanced environment to reduce bias. MCC has a value ranging from  $-1$  to  $1$ , where  $-1$  means that the predicted and actual results are completely inconsistent, and  $1$  means an opposite aspect to  $-1$ .  $MCC = 0$  means that the CSD model applies a random prediction.

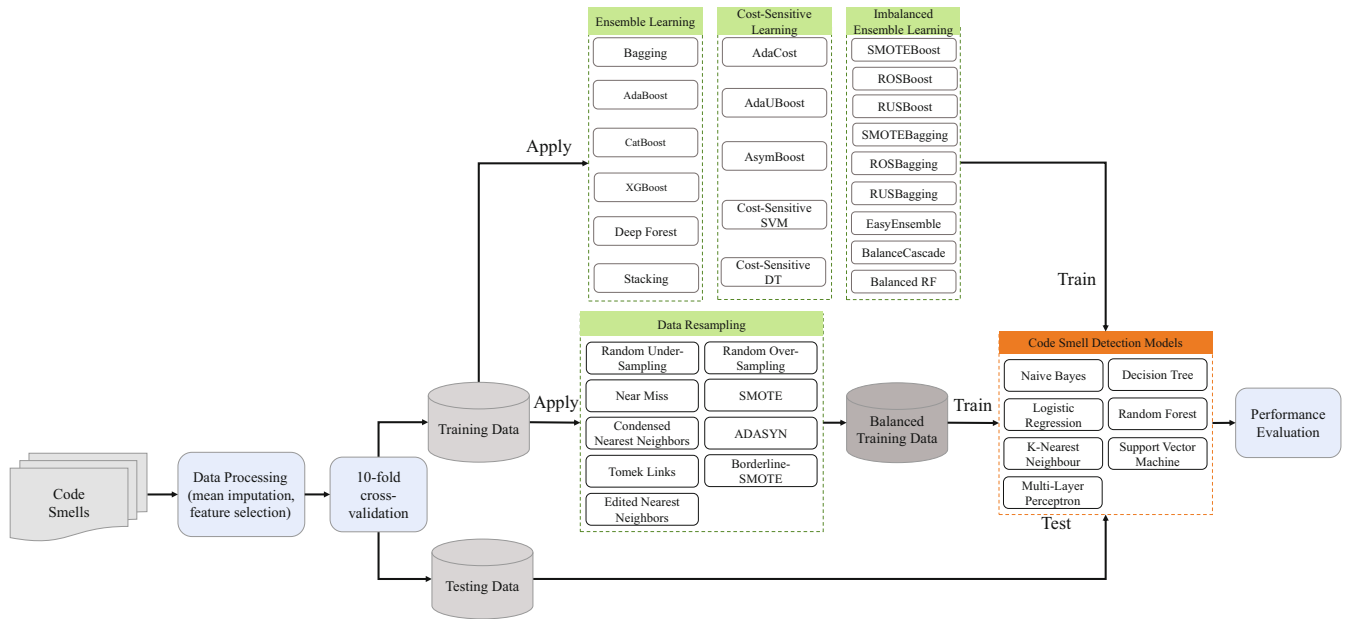


FIGURE 1 Framework of our study.

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (5)$$

### 4.3 | Experimental process

The original data sets constructed by Fontana et al.<sup>20</sup> contain missing values, which may impair the performance of classification models.<sup>76</sup> There are many ways to deal with missing data,<sup>77</sup> such as mean imputation and median imputation. The same as Alazba et al.<sup>21</sup> investigation, we use the mean imputation method to handle missing data for each code smell data set. Mean imputation is a commonly used approach for managing missing values, which produces a good classification result in supervised classification problems.<sup>78</sup> These missing values will be replaced with the mean values belonging to the current column properties.

Irrelevant software features in the original data sets may degrade the detection performance of machine learning classifiers.<sup>79</sup> The same as Nucci et al. study,<sup>22</sup> we remove irrelevant software features by the feature selection method information gain ratio (IGR), whose value is between 0 and 1. We calculate IGR values of all software features in each code smell data set and remove the software features whose IGR values are less than 0.1.<sup>80</sup>

The same as previous CSD studies,<sup>13,21-23</sup> we use 10-fold cross-validation method to verify the detection performance of classification models. Each data set is randomly fall into 10 folds of same size. We select nine folds from all ten folds as training data sets and apply imbalanced learning techniques to them. The last fold is referred as the testing data sets. Next, we iterate 10 times to ensure all ten folds are used as both training and testing data sets. Finally, we obtain 100 results of each detection model in each data set and take the median value of the 100 results as the final result of data sets. The framework of our study as shown in Figure 1.

### 4.4 | Classifiers

Many machine learning classifiers have been used for CSD.<sup>13</sup> To study the influence of the imbalanced learning techniques with more classifiers for the detection performance of CSD, seven classifiers are used to construct the training model, which are the most commonly utilized for code smell detection.<sup>13,21,50,53</sup>

(1) **Naive Bayes (NB)** is based on the Bayes theory and assumes the software features are independent. It needs to estimate a few parameters and is not sensitive to missing data, and the implementation of the algorithm is relatively simple.

**TABLE 5** The optimized hyper-parameters of the classifiers.

Classifier	Hyper-parameters	Tuning range	Description
DT	min_samples_split	[2, 3, 4, 5, 6]	The minimum number of samples required to split an internal node
KNN	n_neighbors	[1, 5, 9, 13, 17]	Number of neighbors to use by default for k neighbors queries
LR	tol	[0.1, 0.01, 0.001, <b>0.0001</b> , 0.00001]	Tolerance for stopping criteria
MLP	alpha	[0, 0.1, 0.01, 0.001, <b>0.0001</b> ]	L2 penalty (regularization term) parameter
	hidden_layer_sizes	[4, 8, 16, 32, 64, <b>100</b> ]	The number of neurons in the hidden layers
RF	n_estimators	[10, 20, 30, 40, 50, <b>100</b> ]	The number of trees in the forest
SVM	<i>C</i>	[0.25, 0.5, <b>1.0</b> , 2.0, 4.0]	The strength of the regularization is inversely proportional to C

Note: The default parameter value is in bold font.

(2) **Support vector machine (SVM)** uses a kernel function to transform the original linearly inseparable data to a feature space, where the data is linearly separable and then finds the optimal hyperplane to partition the feature space to separate instances with different class labels.

(3) **Logistic regression (LR)** is used to classify elements of a set into two groups (binary classification) by calculating the probability of each element of the set, which makes it able to classify code smells into discrete outcomes.

(4) **Decision tree (DT)** expresses a tree structure, which depends on the attributes of the object type by its branch sort.

(5) **K-nearest neighbour (KNN)** is one of the simplest machine learning algorithms. A sample also belongs to the category, if most of the *K* most similar (closest distance) samples in the feature space belong to a specific category.

(6) **Random forest (RF)** generates multi decision trees based on bootstrap samples and makes the final prediction by voting.

(7) **Multi-layer perceptron (MLP)** is a feedforward artificial neural network model that maps multiple inputs to a single output. It contains the input layer, the hidden layers, and the output layer, each of which is fully connected to the previous and subsequent layers.

These classifiers belong to five groups, including statistic-based (i.e., NB and LR), support vector machine-based (i.e., SVM), decision tree-based (i.e., DT and RF), nearest neighbor based (i.e., KNN), and neural networks-based (i.e., MLP). For the machine learning classifiers except NB used in our study, some hyper-parameters are needed to configure to achieve the best detection performance before training models. Optimizing the performance of machine learning models often involves identifying the most effective configuration of hyper-parameters, which collectively influence the learning process and generalization capabilities. In pursuit of this objective, researchers and practitioners face the challenge of exploring a multi-dimensional space encompassing various hyper-parameter combinations. The same as investigations of Shen et al.<sup>49</sup> and Stefano et al.,<sup>55</sup> the grid search algorithm<sup>81</sup> is employed to optimize those hyper-parameters, which involves combining multiple hyper-parameters of a classifier and searching for the optimal values in a grid to obtain the best detection model.<sup>82</sup> The process of hyper-parameter optimization involves discretizing the hyper-parameter space into a grid-like structure. Subsequently, exhaustive exploration is performed by systematically evaluating each combination of hyper-parameter values within this grid. During this evaluation, performance metrics are computed using cross-validation techniques, which provide robust estimates of the model's effectiveness. Table 5 shows the details of the hyper-parameters of the classifiers.

## 4.5 | Statistic test

We employ the Wilcoxon signed-rank test<sup>31</sup> with a Benjamini-Hochberg correction and Cliff's  $\delta$ <sup>32</sup> to analyze the practical significance of the prediction performance between the two models. The Benjamini-Hochberg method involves the following steps. Firstly, we rank the p-values in ascending order from the smallest to the largest. Secondly, we calculate the critical value (q-value) for each p-value. Thirdly, we identify the largest *k* such that the *k*th q-value is less than or equal to *k*/multiplier, where the multiplier is typically set to 1. Then, we reject all null hypotheses corresponding to p-values less than or equal to the critical threshold determined in Step 3. Finally, we report the adjusted p-values (q-values) for each

significant result. The Wilcoxon signed-rank test is a non-parametric statistical hypothesis test used to compare pairs of models. The Cliff's  $\delta$  is a non-parametric effect size that measures the magnitude of difference between two models. In our empirical study, we provide the 100 results of the 10-fold cross-validation of the imbalanced learning techniques and None to the Wilcoxon signed-rank test and calculate Cliff's  $\delta$  value. When the adjusted p-value is less than 0.05 and  $|\delta|$  is larger than 0.147, the corresponding imbalanced learning method and None have a significant difference.<sup>31,32</sup>

## 5 | EXPERIMENTAL RESULTS

This section discusses the research questions with their motivations and results.

### 5.1 | What is the performance of the classifiers for CSD on our data sets?

**Motivations:** Fontana et al.<sup>20</sup> constructed the original code smell data sets, and their results showed that machine learning classifiers achieved high performance on CSD. However, Nucci et al. conclusions<sup>22</sup> implied that the original data sets are not suitable for real-world scenarios where an instance should contain more than one type of code smells. Therefore, we wonder how is the practical effect of machine learning classifiers on CSD.

**Approaches:** We use the baseline method None to train the seven classifiers on the original data sets and our data sets to compare the detection performance on different data sets. We list the median Precision, Recall, F1, and MCC values of the seven classifiers in Tables 6 and 7. In binary classification problems, Precision, Recall and F1 are the commonly used evaluation metrics. However, it may not be appropriate for datasets with imbalanced positive-negative ratios. In the context of imbalanced datasets, machine learning classifiers often exhibit a bias towards learning the characteristics of the majority class. Consequently, when evaluating the classifier's performance on an imbalanced test set, it is conceivable for the classifier to predict all instances as belonging to the majority class and still achieve a relatively high accuracy rate.

TABLE 6 The performance of the seven classifiers on the original code smell data sets.

Classifier	Data class				God class			
	Precision	Recall	F1	MCC	Precision	Recall	F1	MCC
DT	<b>0.978</b>	<b>0.976</b>	<b>0.976</b>	<b>0.949</b>	0.955	0.952	0.952	0.901
KNN	0.927	0.917	0.917	0.833	0.936	0.929	0.928	0.901
LR	0.928	0.929	0.928	0.837	0.954	0.952	0.952	0.892
MLP	0.942	0.929	0.930	0.853	<b>0.977</b>	<b>0.977</b>	<b>0.976</b>	0.934
NB	0.873	0.833	0.835	0.695	0.948	0.940	0.942	0.860
RF	<b>0.978</b>	<b>0.976</b>	<b>0.976</b>	0.947	<b>0.977</b>	0.976	<b>0.976</b>	<b>0.946</b>
SVM	0.977	<b>0.976</b>	<b>0.976</b>	0.944	<b>0.977</b>	0.976	<b>0.976</b>	0.939
Classifier	Feature envy				Long method			
	Precision	Recall	F1	MCC	Precision	Recall	F1	MCC
DT	0.956	0.952	0.951	0.894	<b>0.978</b>	0.976	<b>0.976</b>	0.949
KNN	0.916	0.917	0.915	0.801	0.952	0.952	0.952	0.889
LR	0.933	0.929	0.926	0.834	0.977	<b>0.977</b>	<b>0.976</b>	<b>0.950</b>
MLP	0.930	0.929	0.929	0.847	<b>0.978</b>	0.976	<b>0.976</b>	0.949
NB	0.905	0.905	0.905	0.777	0.959	0.952	0.953	0.896
RF	<b>0.977</b>	<b>0.976</b>	<b>0.976</b>	<b>0.945</b>	<b>0.978</b>	0.976	0.976	0.949
SVM	0.954	0.952	0.951	0.886	0.977	0.976	<b>0.976</b>	<b>0.950</b>

Note: The bold font means the corresponding technique obtains the best performance.



TABLE 7 The performance of the seven classifiers on our data sets.

Classifier	Data class				God class			
	Precision	Recall	F1	MCC	Precision	Recall	F1	MCC
DT	0.622	0.632	0.625	−0.195	0.621	0.631	0.623	−0.189
KNN	0.749	0.743	0.742	0.191	0.724	0.731	0.720	0.113
LR	0.761	<b>0.791</b>	0.757	0.192	0.760	0.788	0.762	0.195
MLP	0.765	0.790	<b>0.769</b>	0.241	0.771	0.789	0.773	0.252
NB	<b>0.872</b>	0.737	0.764	<b>0.501</b>	<b>0.824</b>	<b>0.794</b>	<b>0.804</b>	<b>0.413</b>
RF	0.637	0.645	0.639	−0.148	0.641	0.640	0.637	−0.126
SVM	0.671	0.787	0.721	0.058	0.640	0.785	0.704	−0.029
Classifier	Feature envy				Long method			
	Precision	Recall	F1	MCC	Precision	Recall	F1	MCC
DT	0.822	0.821	0.819	0.532	0.761	0.750	0.749	0.350
KNN	0.860	0.859	0.854	0.623	0.826	0.817	0.814	0.515
LR	<b>0.886</b>	<b>0.878</b>	0.871	<b>0.677</b>	<b>0.834</b>	<b>0.828</b>	<b>0.825</b>	<b>0.557</b>
MLP	0.880	0.876	<b>0.874</b>	0.675	0.829	0.819	0.820	0.545
NB	0.837	0.828	0.830	0.570	0.830	0.819	0.820	0.543
RF	0.841	0.836	0.836	0.579	0.779	0.765	0.767	0.408
SVM	0.870	0.865	0.856	0.632	0.816	0.817	0.809	0.512

Note: The bold font means the corresponding technique obtains the best performance.

Pecorelli et al.<sup>23</sup> suggested the use of MCC as a more reliable metric, as it considers all four categories of the confusion matrix to generate high scores. We highlight the best classifiers on each evaluation metric for each data set in bold.

**Results:** (1) On the original data sets, the best classifiers achieve the following performance: DT scores 0.976 of F1 and 0.949 of MCC on Data Class, RF scores 0.976 of F1 and 0.946 of MCC on God Class, RF scores 0.976 of F1 and 0.945 of MCC on Feature Envy, LR and SVM score 0.976 of F1 and 0.950 of MCC on Long Method. The classifiers except NB score at least 0.8 on Precision, Recall, F1, and MCC, which means that the classifiers perform pretty well for CSD on the original data sets. The main reason of high performance is that the feature distribution between smelly instances and non-smelly ones is significantly different in most cases according to Nucci et al. analysis.<sup>22</sup>

(2) On our data sets, the best classifiers achieve the following performance: MLP and NB score 0.769 of F1 and 0.501 of MCC on Data Class, NB scores 0.804 of F1 and 0.413 of MCC on God Class, LR and MLP score 0.874 of F1 and 0.677 of MCC on Feature Envy, and LR scores 0.825 of F1 and 0.557 of MCC on Long Method.

(3) The MCC values of the best-performing classifiers on our data sets are 28.36%–56.34% lower than that on the original data sets. The Precision, Recall, and F1 values are reduced by 9.31%–15.66%, 10.04%–18.95%, and 10.45%–21.21%, respectively. As Nucci et al.<sup>22</sup> pointed out that if there are substantial differences in the metrics distributions between smelly and non-smelly instances, machine learning techniques are easily able to differentiate between the two instances. However, this does not accurately reflect real-world situations where the boundary between the characteristics of smelly and non-smelly instances is not always clearly defined. Additionally, Palomba et al.<sup>28</sup> has found that smelly instances make up just a small portion of the all instances in a software system, with approximately one third of the all instances (including non-smelly and smelly instances) in the original dataset being smelly instances. The results imply that the classifiers cannot properly detect the smelliness of software instances in a more realistic case where software instances contain more than one type of smell. Therefore, there is still plenty of room for performance improvement of the CSD models built with the seven classifiers.

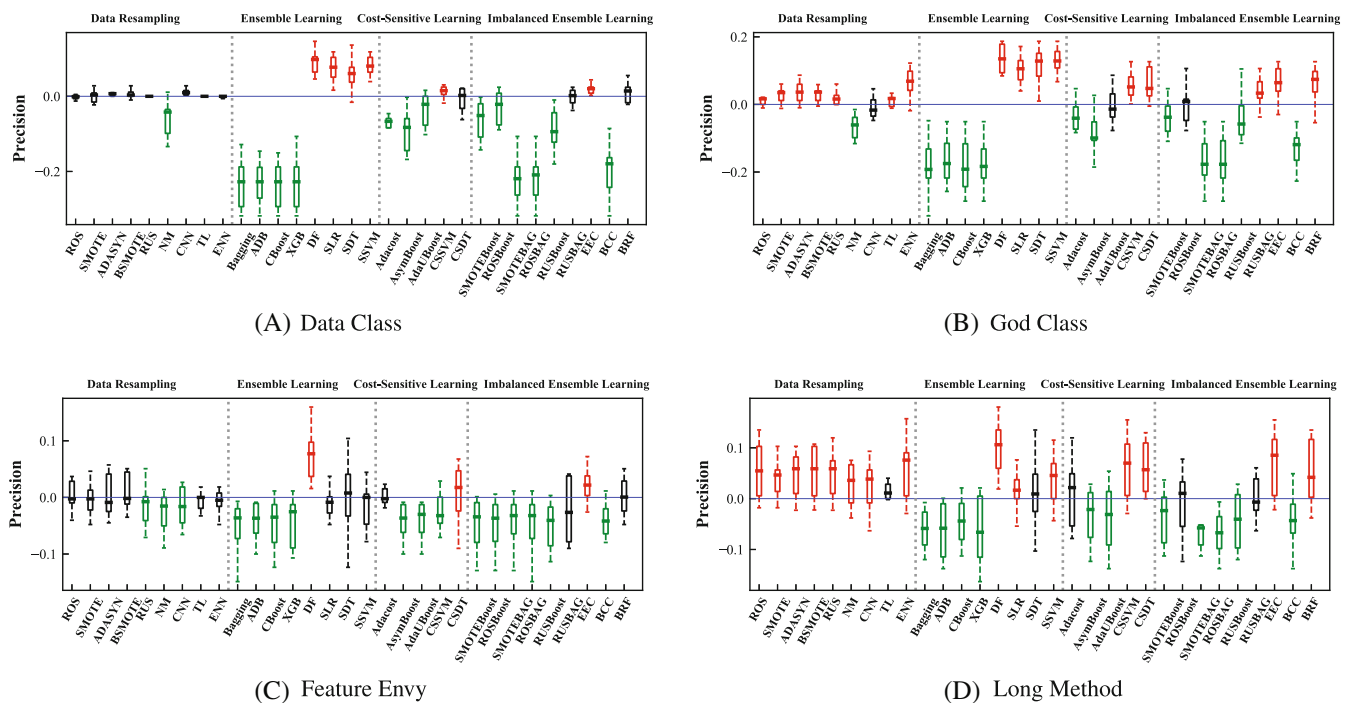
**Summary 1:** The F1 and MCC values of the best-performing classifiers are 10.45%–21.21% and 28.36%–56.34% lower than those on the original data sets. The best performance on different code smells is achieved by different classifiers, that is, NB on Data Class and God Class, LR on Feature Envy and Long Method.

## 5.2 | Do imbalanced learning techniques improve the performance of code smell detection models?

**Motivations:** In our study, we have considered the review work of Azeem et al.<sup>25</sup> regarding the application of machine learning techniques in CSD. However, we have specifically focused our literature review on the utilization of imbalanced learning technique in CSD. Therefore, our comparison of detection performance is centered on the differences between various imbalanced learning techniques, rather than a deeper exploration of the impact of imbalanced learning techniques on the detection performance of machine learning classifiers. The imbalanced ratio of our data sets ranges from 20.03% to 27.03%. The previous studies<sup>13,23</sup> show that CSD models trained on imbalanced data sets might lead to inaccurate prediction results. Therefore, we aim to investigate whether imbalanced learning techniques can alleviate the class imbalance problem to improve the performance of CSD models.

**Approaches:** We apply the nine data resampling techniques to our data sets and then employ the best-performing classifier on each data set to build CSD models (i.e., NB on Data Class and God Class, LR on Feature Envy and Long Method). We apply the eight ensemble learning techniques, the five cost-sensitive learning techniques, and the nine imbalanced ensemble learning techniques to directly build CSD models. Figures 2, 3, 4, and 5 show the distribution of performance difference between the imbalanced learning techniques and None in terms of Precision, Recall, F1, and MCC, respectively. We employ both statistic tests (i.e., Wilcoxon signed-rank test and Cliff's  $\delta$ ) and plot boxplots to examine the significant difference, where the red box means the corresponding technique significantly outperforms None, the green box means the corresponding technique significantly performs worse than None, and the black box means no significant difference between the corresponding technique and None. In order to make the boxplots more clear, we abbreviate AdaBoost, CatBoost, XGBoost, Deep Forest, StackingLR, StackingDT, StackingSVM, SMOTEBagging, ROS-Bagging, RUSBagging, EasyEnsemble Classifier, and BalanceCascade Classifier as ADB, CBoost, XGB, DF, SLR, SDT, SSVM, SMOTEBAG, ROSBAG, RUSBAG, EEC, and BCC. Table 8 shows the top-3 data resampling techniques with the best performance and the improvement ratio compared with None.

**Results:** (1) There is great variability in the performance of the imbalanced learning techniques compared with None. Not all imbalanced learning techniques can enhance the performance, and an only average of 32.26%, 29.03%, 30.65%, and 34.68% of the imbalanced learning techniques have a significant positive effect on CSD across our data sets in terms of Precision, Recall, F1, and MCC, respectively.



**FIGURE 2** The performance difference between the imbalanced learning techniques and none in terms of precision.

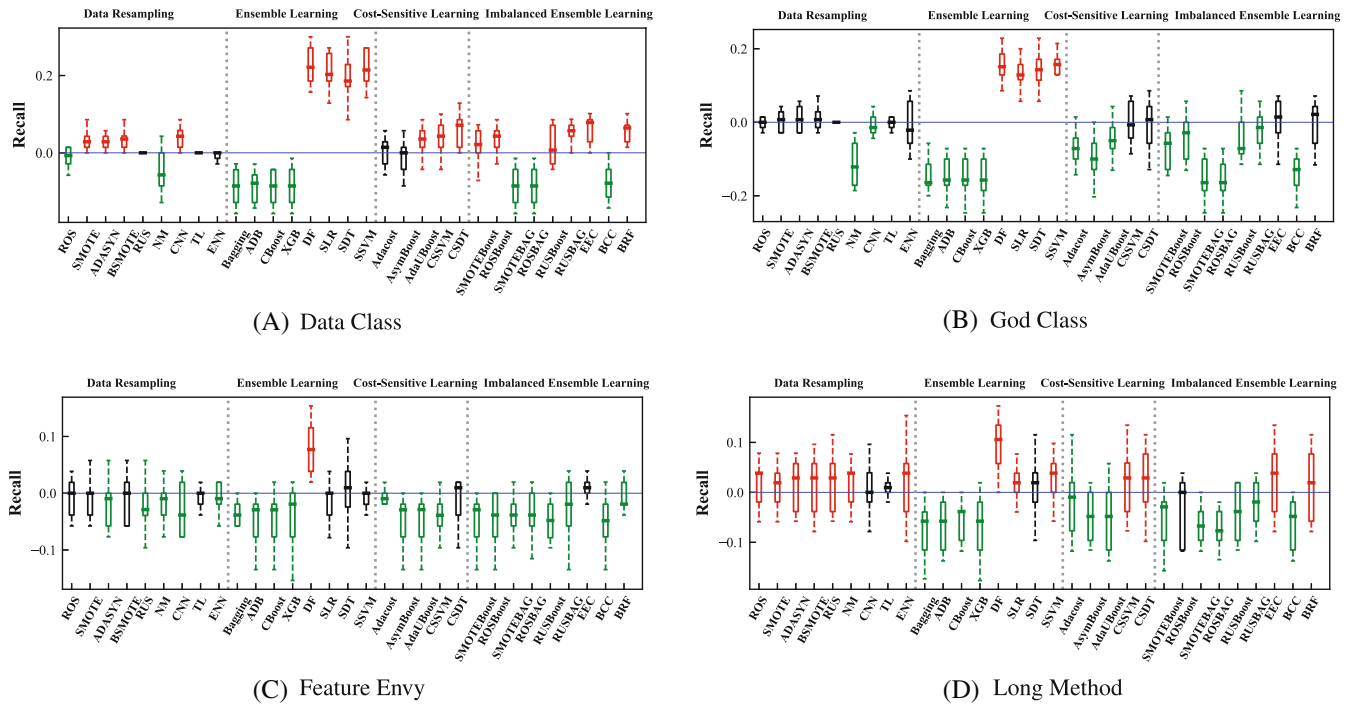


FIGURE 3 The performance difference between the imbalanced learning techniques and none in terms of recall.

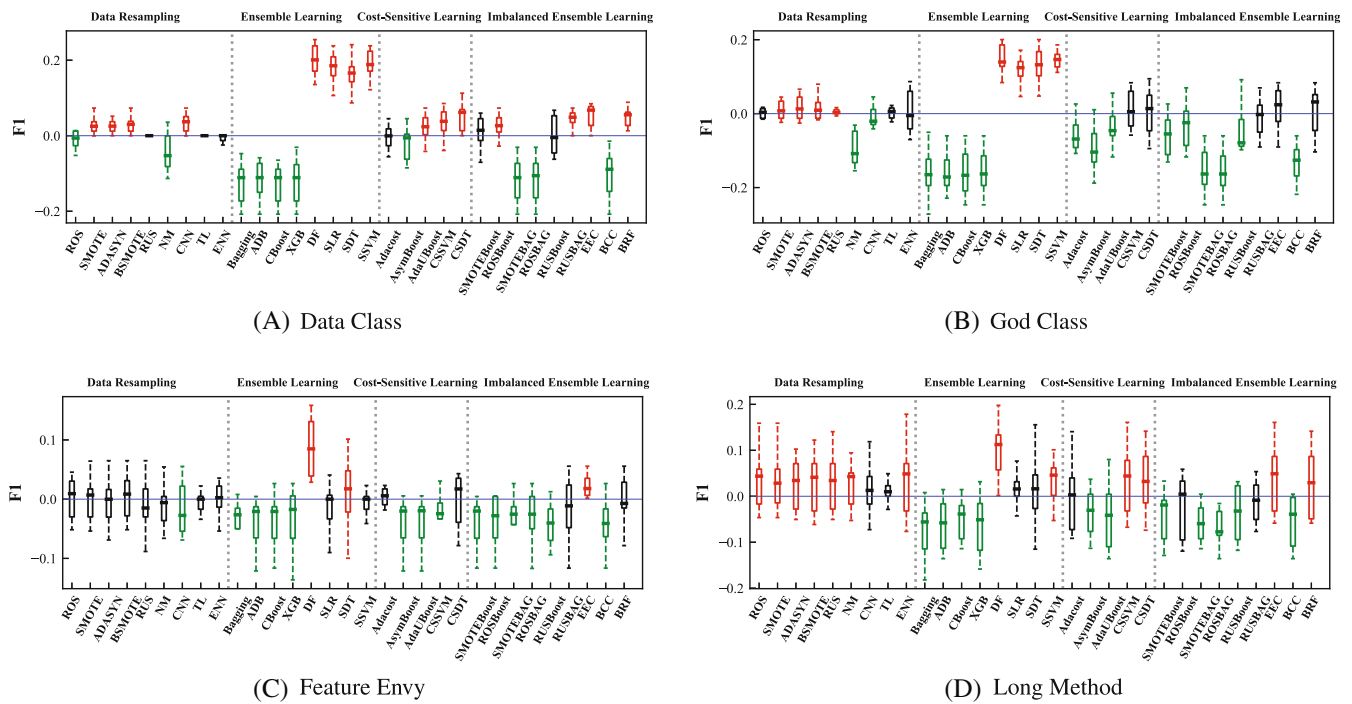
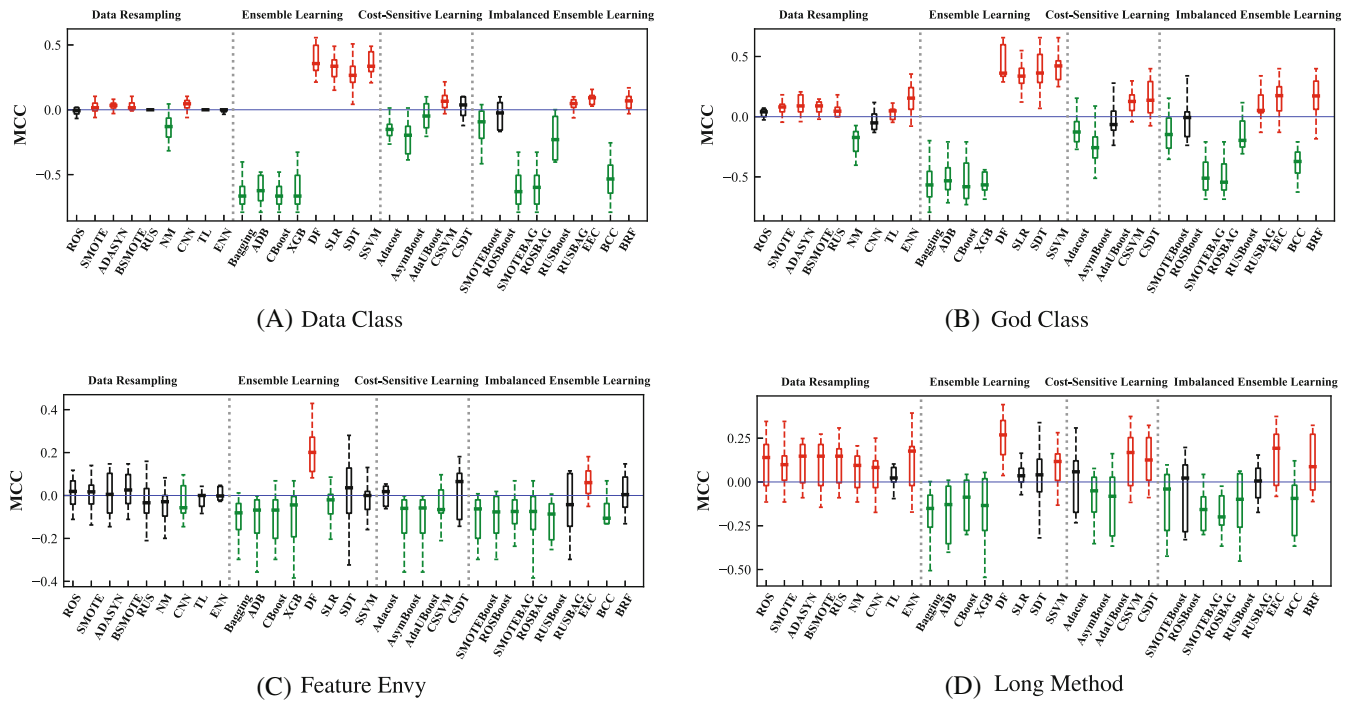


FIGURE 4 The performance difference between the imbalanced learning techniques and none in terms of F1.



**FIGURE 5** The performance difference between the imbalanced learning techniques and none in terms of MCC.

**TABLE 8** The top-3 imbalanced learning techniques with the best performance of each data set.

Data set	Techniques	Precision	Recall	F1	MCC
Data class	DF	<b>0.966(10.78%)</b>	<b>0.964(30.80%)</b>	<b>0.964(26.18%)</b>	<b>0.892(78.04%)</b>
	SSVM	<b>0.954(9.40%)</b>	<b>0.953(29.31%)</b>	<b>0.953(24.74%)</b>	<b>0.854(70.46%)</b>
	SLR	<b>0.950(8.94%)</b>	<b>0.950(28.90%)</b>	<b>0.949(24.21%)</b>	<b>0.839(67.47%)</b>
	None	0.872	0.737	0.764	0.501
God class	DF	<b>0.956(16.02%)</b>	<b>0.953(20.03%)</b>	<b>0.953(18.53%)</b>	<b>0.850(105.81%)</b>
	SSVM	<b>0.953(15.66%)</b>	<b>0.952(19.90%)</b>	<b>0.951(18.28%)</b>	<b>0.838(102.91%)</b>
	SDT	<b>0.941(14.20%)</b>	<b>0.937(18.01%)</b>	<b>0.937(16.54%)</b>	<b>0.797(92.98%)</b>
	None	0.824	0.794	0.804	0.413
Feature envy	DF	<b>0.959(8.24%)</b>	<b>0.958(9.11%)</b>	<b>0.958(9.99%)</b>	<b>0.889(31.31%)</b>
	EEC	<b>0.910(2.71%)</b>	0.882(0.46%)	<b>0.888(1.95%)</b>	<b>0.744(9.90%)</b>
	CSDT	<b>0.894(0.90%)</b>	0.867(−1.25%)	0.872(0.11%)	0.705(4.14%)
	None	0.886	0.878	0.871	0.677
Long method	DF	<b>0.934(11.99%)</b>	<b>0.921(11.23%)</b>	<b>0.923(11.88%)</b>	<b>0.808(45.06%)</b>
	EEC	<b>0.903(8.27%)</b>	<b>0.855(3.26%)</b>	<b>0.864(4.73%)</b>	<b>0.703(26.21%)</b>
	CSSVM	<b>0.898(7.67%)</b>	<b>0.851(2.78%)</b>	<b>0.851(3.15%)</b>	<b>0.691(24.06%)</b>
	None	0.834	0.828	0.825	0.557

Note: The bold font means the corresponding technique significantly outperforms None.

(2) In terms of Precision and Recall, no data resampling techniques have a significant positive effect on Data Class and Feature Envy. However, 77.78% and 88.89% data resampling techniques significantly enhance the performance on Precision, and 44.44% and 77.78% data resampling techniques significantly improve the Recall value on God Class and Long Method. 12.50%–50.00% ensemble learning techniques, 20.00%–40.00% cost-sensitive learning techniques, and 11.11%–33.33% imbalanced ensemble learning techniques significantly improve the detection performance in terms of Precision on the four data sets. 12.50%–50.00% ensemble learning techniques significantly improve the detection performance in terms of Recall on the four data sets. 60.00% and 40.00% cost-sensitive learning techniques and 66.67% and 22.22% imbalanced ensemble learning techniques significantly increase the Precision value on Data Class and Long Method. But no cost-sensitive learning and imbalanced ensemble learning techniques have a significant positive effect on God Class and Feature Envy.

(3) In terms of F1, no data resampling techniques have a significant positive effect on Feature Envy. However, 44.44%, 44.44%, and 77.78% data resampling techniques significantly enhance the performance on Data Class, God Class, and Long Method, respectively. 50.00%, 50.00%, 25.00%, and 25.00% ensemble learning techniques significantly improve the performance on Data Class, God Class, Feature Envy, and Long Method, respectively. 60.00% and 40.00% cost-sensitive learning techniques significantly improve the performance on Data Class and Long Method, respectively. But no cost-sensitive learning techniques have a significant positive effect on God Class and Feature Envy. 44.44%, 11.11%, and 22.22% imbalanced ensemble learning techniques significantly improve the performance on Data Class, Feature Envy, and Long Method, respectively. However, no imbalanced learning techniques have a significant positive effect on God Class.

(4) In terms of MCC, no data resampling techniques have a significant positive effect on Feature Envy. However, 44.44%, 66.67%, and 88.89% data resampling techniques significantly enhance the performance on Data Class, God Class, and Long Method. 50.00%, 50.00%, 12.50%, and 25.00% ensemble learning techniques significantly improve the performance on Data Class, God Class, Feature Envy, and Long Method, respectively. 20.00%, 40.00%, and 40.00% cost-sensitive learning techniques significantly improve the performance on Data Class, God Class, and Long Method, respectively. But no cost-sensitive learning techniques have a significant positive effect on God Class and Feature Envy. 33.33%, 33.33%, 11.11%, and 22.22% imbalanced ensemble learning techniques significantly improve the performance on Data Class, God Class, Feature Envy, and Long Method, respectively.

(5) StackingSVM is the second best performing technique on Data Class and God Class, which increases the F1 and MCC values by 18.28%–24.74% and 70.46%–102.91%. EasyEnsemble achieves the second best performance on Feature Envy and Long Method, which increases the F1 and MCC values by 1.95%–4.73% and 9.90%–26.21%. Deep forest is stable and always obtains the best performance in terms of all evaluation metrics on all four data sets. It increases the Precision, Recall, F1, and MCC values by 8.24%–16.02%, 9.11%–30.80%, 9.99%–26.18%, and 31.31%–105.81% across four data sets. The main reason for high performance is that deep forest is adept in handling code smell features by applying multi-grained scanning, and its cascade forest structure is an ensemble of decision tree and random forest, which produces more accurate CSD results under layer-by-layer forests training.

**Summary 2:** Not all imbalanced learning techniques can enhance the performance of CSD models, but deep forest always shows statistically significant improvement of 9.99%–26.18% and 31.31%–105.81% in terms of F1 and MCC.

### 5.3 | What is the best data resampling technique for CSD?

**Motivations:** Some researchers might only use the dataset-level imbalanced learning technique (i.e., data resampling) as a preprocessing method to alleviate the imbalance problems and then employ more advanced machine learning algorithms (e.g., deep neural networks) to build CSD models. The main reason is that the machine learning algorithms are difficult to embed the algorithm-level imbalanced learning techniques (i.e., ensemble learning, cost-sensitive learning, and imbalanced ensemble learning). In addition, the previous studies<sup>13,23</sup> showed that using SMOTE can not significant improve the detection performance on code smell data sets. Therefore, we discuss the reason for non-improvement and conduct a investigation to explore the practical effect of data resampling techniques on our data sets using an individual research question.

**Approaches:** We use the nine data resampling techniques to balance our data sets and compare the performance with SMOTE. We show the top-3 data resampling techniques with the best performance and the improvement ratio compared with SMOTE in Table 9.

**Results:** (1) SMOTE does not achieve the best performance on all data sets. The top-1 data resampling techniques achieve the following MCC performance: CNN scores 0.541 on Data Class, ENN scores 0.571 on God Class, BSMOTE



**TABLE 9** The top-3 data resampling techniques with the best performance of each data set.

Data set	Techniques	Precision	Recall	F1	MCC
Data class	CNN	<b>0.877(0.57%)</b>	0.778(0.91%)	0.801(0.88%)	<b>0.541(3.05%)</b>
	BSMOTE	0.876(0.46%)	0.771(0.00%)	0.794(0.00%)	0.535(1.90%)
	ADASYN	<b>0.876(0.46%)</b>	0.767(−0.52%)	0.791(−0.38%)	0.531(1.14%)
	SMOTE	0.872	0.771	0.794	0.525
God class	ENN	<b>0.890(4.58%)</b>	0.788(−1.13%)	0.811(−0.25%)	<b>0.571(17.73%)</b>
	ADASYN	<b>0.860(1.06%)</b>	0.803(0.75%)	0.819(0.74%)	<b>0.510(5.15%)</b>
	BSMOTE	0.858(0.82%)	0.803(0.75%)	0.819(0.74%)	<b>0.508(4.12%)</b>
	SMOTE	0.851	0.797	0.813	0.485
Feature envy	BSMOTE	<b>0.893(1.02%)</b>	0.869(0.00%)	0.874(0.11%)	<b>0.702(2.63%)</b>
	ROS	0.885(0.11%)	0.871(0.23%)	0.874(0.11%)	<b>0.691(1.02%)</b>
	SMOTE	0.884	0.869	0.873	0.684
	ADASYN	0.887(0.34%)	0.857(−1.38%)	0.863(−1.15%)	0.681(−0.44%)
Long method	ROS	0.889(1.02%)	0.855(0.59%)	0.862(0.70%)	0.677(3.68%)
	ENN	0.892(1.36%)	0.849(−0.12%)	0.858(0.23%)	0.675(3.37%)
	RUS	0.887(0.80%)	0.848(−0.24%)	0.856(0.00%)	0.670(2.60%)
	SMOTE	0.880	0.850	0.856	0.653

Note: The bold font means the corresponding technique significantly outperforms SMOTE.

scores 0.702 on Feature Envy, and ROS scores 0.677 on Long Method. The MCC values of the top-1 data resampling techniques are 2.63%–17.73% higher than SMOTE on our data sets. In addition, CNN significantly outperforms SMOTE on Data Class; ENN, ADASYN, and BSMOTE significantly outperform SMOTE on God Class; BSMOTE and ROS significantly outperform SMOTE on Feature Envy. The Precision, Recall, and F1 improvements of the top-1 data resampling techniques range from 0.57%–4.58%, −1.13%–0.91%, and −0.25%–0.88%, respectively.

(2) The best classifiers are DT, RF, RF, and LR on the original four data sets, respectively. On the preprocessed data sets by SMOTE, the F1 and MCC values of DT reduce by 1.23% and 3.37% on Data Class; The MCC value of RF increases by 0.53% on God Class, but the improvement is not significant; The F1 and MCC values of RF and LR do not change on Feature Envy and Long Method, respectively. In summary, using SMOTE to balance the original data sets cannot significantly improve the performance of CSD models. The main reasons are as follows. Since the feature distribution between smelly instances and non-smelly instances in the original data sets is quite different, the best classifiers on each data set have achieved relatively high performance in terms of F1 and MCC. Using SMOTE cannot further improve the feature distribution difference, so there is little room for improvement on each data set by SMOTE.

**Summary 3:** The best data resampling techniques are CNN, ENN, BSMOTE, and ROS on the four data sets, respectively.

## 5.4 | How efficient are the best-performing imbalanced learning techniques?

**Motivations:** Since the top-3 imbalanced learning techniques and the top-3 data resampling techniques need to build multi base classifiers or increase the size of training data, the training time increases accordingly. Hence, we aim to explore the efficiency of the imbalanced learning techniques.

**Approaches:** We conduct experiments on a personal computer with an Intel i7-8750H CPU and 16GB RAM. We record the time cost of ten rounds of 10-fold cross-validation on all data sets. We select the top-3 imbalanced learning techniques and the top-3 data resampling techniques to show their time cost in Table 10.

**Results:** The stacking ensemble learning techniques consume a long time about 20s. They use the seven machine learning classifiers as the base classifiers to generate the first detection predictions, then select a classifier as the

**TABLE 10** The efficiency of top-3 imbalanced learning techniques and the top-3 data resampling technique.

Data set	Techniques	Time cost (s)
Data class	DF	2.32
	SSVM	21.97
	SLR	22.11
	CNN	0.01
	BSMOTE	0.01
	ADASYN	0.01
God class	DF	2.37
	SSVM	20.86
	SDT	20.56
	BSMOTE	0.01
	ROS	0.01
	SMOTE	0.01
Feature envy	DF	2.03
	EEC	1.90
	CSDT	4.25
	ENN	0.20
	ADASYN	0.24
	BSMOTE	0.25
Long method	DF	2.12
	EEC	1.90
	CSSVM	0.02
	ROS	0.24
	ENN	0.20
	RUS	0.19

meta-classifier to take the outputs of the first predictions as inputs and produce the final detection results. Therefore, the stacking ensemble takes the longest training time. Among the top-3 imbalanced learning techniques, the best-performing technique, deep forest, can achieve a low time cost about of 2s. The time costs of the top-3 data resampling techniques are less than 0.3s. The main reason is that we use only a classifier to build the detection model rather than train multi classifiers. We employ NB as the classifier on Data Class and God Class and use the default parameter. We employ LR as the classifier on Feature Envy and Long Method and use the grid search algorithm to find the optimal parameter of LR. Therefore, the time cost of the top-3 data resampling techniques on Feature Envy and Long Method is higher than that of Data Class and God Class.

**Summary 4:** The best-performing imbalanced learning technique deep forest runs fast, and the time cost of the top-3 data resampling techniques is less than 0.3s.

## 6 | DISCUSSION

### 6.1 | The differences between our findings and previous findings

As outlined in our review of the literature in Section 2, we will now discuss the results obtained by recent researches.<sup>13,20-23</sup> Fontana et al.<sup>20</sup> conducted an investigation using 16 different machine learning algorithms to detect four specific code

smells. They found that all of the algorithms achieved high performance and concluded that the use of machine learning for the detection of these code smells can result in high accuracy ( $> 96\%$ ). However, they also noted that the imbalanced nature of the data (low prevalence of code smells) resulted in inconsistent performance across the dataset, and this issue should be addressed in future research. Nucci et al.<sup>22</sup> acknowledged that the work of Fontana et al. presented a new perspective on CSD, but pointed out that it only considers instances affected by a single type of code smell in the data sets used for training and testing the machine learning models. In order to address this limitation, Nucci et al. replicated the investigation using a data set containing instances of multiple types of code smells. The results showed that when the machine learning-based CSD model was tested on the Nucci et al. data sets, the F1-measure value was 90% lower than previously reported. As previously mentioned by Fontana et al., the highly imbalanced nature of CSD data sets can make machine learning techniques unsuitable. Pecorelli et al.<sup>23</sup> examined the impact of five imbalanced learning methods for CSD in object-oriented systems. They found that the methods can not significantly improve the performance. Alazba et al.<sup>21</sup> explored the use of stacking ensemble models for CSD. They applied information gain feature selection to select relevant features and evaluated the performance of 14 individual classifiers on the Fontana et al. data sets. They then constructed three stacking ensembles using these individual classifiers as base models and three different meta-classifiers (LR, SVM, and DT) and compared the detection performance of the stacking ensembles to the individual models. The results showed that the stacking ensembles using LR and SVM meta-classifiers consistently achieved high detection performance for both class-level and method-level code smells. Alkharabsheh et al.<sup>13</sup> investigated the effectiveness of machine learning for CSD and conducted a comparative investigation to assess the impact of data imbalance on accuracy and behavior during CSD. They carried out experiments using 28 machine learning classifiers, and found that most classifiers achieved high performance, with CatBoost showing particularly good performance. In their experiments, the imbalanced learning technique SMOTE did not significantly improve the detection of God Class code smell.

In contrast to previous studies, our study aims to evaluate the effectiveness of a variety of imbalanced learning techniques for CSD and to identify the best performing technique. To this end, we combine the same level code smell data sets and remove the the redundant and conflicting instances. Then, we compare the performance of 31 imbalanced learning techniques across four evaluation metrics on different code smells, taking into account both detection performance and time cost. Our results show that the difference between our study and previous investigations.

(1) Our study shows that machine learning techniques do not achieve particularly high detection performance on our data sets, with the F1 value and MCC value of the best-performing classifier dropping by 10.45%–21.21% and 28.36%–56.34%, respectively, compared to the results of Fontana et al. However, the performance drop is less than that reported by Nucci et al.<sup>22</sup> on their data sets, possibly due to the removal of redundant and conflicting instances in our data sets.

(2) In our investigation of 31 imbalanced learning techniques, we find that the deep forest consistently improve the detection of the four code smells across four evaluation metrics, with low time cost and good performance on most code smells compared to the other top-3 imbalanced learning techniques. This differs from the results of Pecorelli et al.<sup>23</sup> We also find that the use of stacking ensemble technique can significantly improve the performance of CSD, as previously concluded by Alazba et al.<sup>21</sup> However, their study only examines stacking ensemble technique and does not find that the deep forest could achieve better detection performance. Specifically, the performance of the deep forest is better than that of the stacking ensemble model on our data sets.

(3) Additionally, our empirical results do not support the claim that SMOTE is the best data resampling technique for handling imbalanced code smell data sets, as it does not consistently achieve the best detection performance across all code smells when compared to other data resampling techniques. Contrary to the conclusion of Alkharabsheh et al.,<sup>13</sup> we find that SMOTE has a detection performance improvement for all three code smells (Data Class, God Class, and Long Method) in our data sets. However, even though SMOTE improves detection performance, it is not able to become the best data resampling technique for CSD.

## 6.2 | Threats to validity

**Construct validity.** To ensure the generalizability of our conclusions and avoid bias as far as possible, we have employed a wide range of 31 imbalanced learning techniques and seven classifiers in our study. While we recognize the existence of other imbalanced learning techniques, their inclusion in our study is left for future work. To minimize technical errors, all steps in our empirical research process, including data preprocessing, model building, and result analysis, were

carried out using third-party Python libraries. For instances, Scikit-learn<sup>†</sup> for model construction, Imbalanced-learn<sup>‡</sup> and Imbalanced-ensemble<sup>§</sup> for imbalanced learning technique implement. To minimize potential errors in the experimental process, all experiments are conducted using 10-fold cross-validation.

**Internal validity.** One potential threat to the validity of our results is the choice of evaluation metrics used to assess the performance of CSD models. Using only one evaluation metric can be biased and lead to misleading conclusions. To address this issue, we employ both threshold-dependent (i.e., Precision, Recall, and F1) and threshold-independent (MCC) evaluation metrics, which are commonly used in recent CSD studies.<sup>43,64,83</sup>

**External validity.** Our findings are based on data sets provided by Fontana et al.,<sup>20</sup> which are derived from 74 systems in the Qualitas corpus. Despite the widespread usage of these code smell datasets in recent studies,<sup>21,22,53</sup> it should be noted that the conclusions reached in our study may not necessarily hold true for other code smell datasets. Current researches<sup>13,23,25</sup> in CSD tend to treat code smells as a binary classification problem, meaning that a code block is either classified as having a particular smell or not having that smell. This means that when a code block contains multiple code smells, separate models must be used to predict the presence of each smell. However, it may be more practical to detect all code smells in a single model, as this would eliminate the need for multiple separate predictions. This study represents an area for future research in the field. To facilitate replication of our experiment by future researchers and mitigate the effects of model variability caused by randomly processing data, we employ 10-fold cross-validation, a well-established validation method in machine learning. Despite these efforts, we acknowledge that our conclusions may not be generalizable to other data sets.

**Conclusion validity.** To optimize the performance of our classifiers, we use grid search to identify the optimal hyper-parameters. We set the smelly ratio of the data resampling techniques to 0.5, as this is a common practice and has been used in previous empirical CSD studies.<sup>13,23</sup> For other imbalanced learning techniques, we use the default hyper-parameters provided by the corresponding third-party libraries. Future research may involve further exploration of hyper-parameter tuning.

### 6.3 | Implications

Based on our experimental results, we summarize the following implications for future CSD studies.

(1) Our results indicate that the blind application of imbalanced learning techniques may not always be optimal, and researchers and practitioners should consider using deep forest when building CSD models. While previous studies, such as that by Alazba et al.,<sup>21</sup> have found that the stacking ensemble technique performs well on CSD, our results show that deep forest consistently outperforms the best-performing stacking ensemble technique across all data sets in terms of Precision, Recall, F1, and MCC by 0.36%–1.38%, 0.13%–1.49%, 0.25%–1.44%, and 2.90%–7.58%, respectively. Based on our findings, the deep forest technique is recommended for improving CSD performance. However, not all imbalanced learning techniques are effective in enhancing CSD, and an average of 32.26%, 29.03%, 30.65%, and 34.68% of the techniques we tested have a significant positive impact on CSD across our data sets in terms of Precision, Recall, F1, and MCC, respectively. For example, the Near Miss technique perform similarly or worse than the None technique on Data Class, God Class, and Feature Envy. Therefore, we advise researchers and practitioners to carefully select appropriate imbalanced learning techniques, such as deep forest, to achieve more accurate CSD models.

(2) Researchers and practitioners should consider following the data sets processing methods outlined by Nucci et al.<sup>22</sup> and Alazba et al.<sup>21</sup>. Nucci et al.<sup>22</sup> noted that the instances in original data sets often only contain one type of code smell, which does not accurately reflect real-world scenarios and makes it easier for CSD models to distinguish smelly instances. As a result, experimental results based on these data sets may be misleading. To address this issue, Nucci et al. merged data sets at the same level to create new data sets. However, Alazba et al.<sup>21</sup> pointed out that this merging strategy can lead to redundant and conflicting instances. Therefore, we believe that two different code smell data sets can be merged into a new code smell data set, as suggested by Nucci et al., but it is important to also follow Alazba et al. recommendation to remove redundant and conflicting instances from the resulting data sets. By following these steps, future research can more accurately study the performance of machine learning techniques in CSD and obtain more reliable results.

<sup>†</sup><https://github.com/scikit-learn/scikit-learn>

<sup>‡</sup><https://github.com/scikit-learn-contrib/imbalanced-learn>

<sup>§</sup><https://github.com/ZhiningLiu1998/imbalanced-ensemble>

(3) Researchers and practitioners should consider using a more effective data resampling technique to preprocess code smell data sets instead of relying solely on SMOTE. Previous empirical studies<sup>13,23</sup> have indicated that SMOTE does not consistently improve detection performance on original code smell data sets, likely due to its inability to significantly increase the difference in feature distribution. As shown in Figures 4 and 5, SMOTE does achieve significant improvement over the None technique on Data Class, God Class, and Long Method across our data sets, and obtains non-significant improvement on Feature Envy. Therefore, researchers and practitioners may still consider using SMOTE as a preprocessing method in line with previous studies,<sup>50-53,55-59</sup> but should also consider exploring other techniques that may be more effective. Our results in Section 5.3 demonstrate that SMOTE does not consistently achieve the best performance on all four data sets, and the top-performing data resampling technique outperforms SMOTE by 2.63%–17.73% in terms of MCC. In other words, SMOTE is not the best data resampling technique. Therefore, we recommend that researchers and practitioners consider using alternative techniques to preprocess code smell data sets.

## 7 | CONCLUSION

In our investigation, we evaluate the impact of 31 imbalanced learning techniques on seven machine learning classifiers for CSD. To obtain more effective data sets, we follow the methods of Nucci et al. and Alazba et al. to combine two different code smell data sets and remove redundant and conflicting instances. We use Precision, Recall, F1, and MCC to comprehensively evaluate the performance of the detection models and applied both the Wilcoxon signed-rank test and Cliff's  $\delta$  to analyze the results statistically. Our findings show that the effect of the imbalanced learning techniques varies across different code smells, and deep forest consistently improves performance across all data sets. Therefore, we recommend that researchers and practitioners use deep forest to enhance detection performance. Additionally, we find that certain data resampling techniques, such as CNN, ENN, BSMOTE, and ROS, perform better than SMOTE. Future research should carefully consider the selection of different data resampling techniques to preprocess different code smell data sets.

## ACKNOWLEDGMENTS

This work was in part supported by the Project of Sanya Yazhou Bay Science and Technology City (SCKJ-JYRC-2022-17), the Natural Science Foundation of China (62272356), the Youth Fund Project of Hainan Natural Science Foundation of China (622QN344), the Natural Science Foundation of Chongqing (cstc2021jcyj-msxmX115), and the Start-up Grant from Wuhan University of Technology (104-40120693).

## CONFLICT OF INTEREST STATEMENT

The authors have no conflicts of interest to declare that are relevant to the content of this article.

## AUTHOR CONTRIBUTIONS

**Fuyang Li:** Methodology, Formal analysis, Writing - Original Draft & Review. **Kuan Zou:** Data Curation, Software, Writing - Original Draft. **Jacky Wai Keung:** Project administration, Validation. **Xiao Yu:** Formal analysis, Methodology, Supervision, Writing - review & editing. **Shuo Feng:** Investigation, Software, Visualization. **Yan Xiao:** Data Curation, Writing - Review & Editing.

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are available in IST2021-CodeSmellStackingEnsemble at <https://github.com/hjamaan/IST2021-CodeSmellStackingEnsemble>, reference Alazba and Aljamaan (2021).

## ORCID

Xiao Yu  <https://orcid.org/0000-0002-4473-3068>

## REFERENCES

1. Sabir F, Palma F, Rasool G, Guéhéneuc YG, Moha N. A systematic literature review on the detection of smells and their evolution in object-oriented and service-oriented systems. *Softw: Pract Exp*. 2019;49:3-39.
2. Pornprasit C, Tantithamthavorn C. Deeplinedp: towards a deep learning approach for line-level defect prediction. *IEEE Trans Softw Eng*. 2022;49:84-98.



3. Jia K, Yu X, Zhang C, Hu W, Zhao D, Xiang J. Software aging prediction for cloud services using a gate recurrent unit neural network model based on time series decomposition. *IEEE Trans Emerg Top Comput*. 2023. doi:10.1109/TETC.2023.3258503
4. Yu X, Keung J, Xiao Y, Feng S, Li F, Dai H. Predicting the precise number of software defects: are we there yet? *Inf Softw Technol*. 2022;146:106847.
5. Li Z, Zhang H, Jing XY, Xie J, Guo M, Ren J. Dssdpp: data selection and sampling based domain programming predictor for cross-project defect prediction. *IEEE Trans Softw Eng*. 2022;49:1941–1963.
6. Tong H, Lu W, Xing W, Liu B, Wang S. Shse: a subspace hybrid sampling ensemble method for software defect number prediction. *Inf Softw Technol*. 2022;142:106747.
7. Yu X, Dai H, Li L, et al. Finding the best learning to rank algorithms for effort-aware defect prediction. *Inf Softw Technol*. 2023;157:107165. doi:10.1016/j.infsof.2023.107165
8. Yang Z, Keung J, Yu X, et al. A multi-modal transformer-based code summarization approach for smart contracts. Paper presented at: 2021 IEEE/ACM 29th International Conference on Program Comprehension (ICPC), IEEE. 2021 1–12.
9. Fowler M. Refactoring: improving the design of existing code. Paper presented at: 11th European Conference. Jyväskylä, Finland. 1997.
10. Rahman MM, Riyadh RR, Khaled SM, Satter A, Rahman MR. Mmruc3: a recommendation approach of move method refactoring using coupling, cohesion, and contextual similarity to enhance software design. *Softw: Pract Exp*. 2018;48:1560–1587.
11. Sousa BL, Bigonha MA, Ferreira KA. An exploratory study on cooccurrence of design patterns and bad smells using software metrics. *Softw: Pract Exp*. 2019;49:1079–1113.
12. Rahad K, Badreddin O, Mohsin Reza S. The human in model-driven engineering loop: a case study on integrating handwritten code in model-driven engineering repositories. *Softw: Pract Exp*. 2021;51:1308–1321.
13. Alkharabsheh K, Alawadi S, Kebande VR, Crespo Y, Delgado MF, Taboada JA. A comparison of machine learning algorithms on design smell detection using balanced and imbalanced dataset: a study of god class. *Inf Softw Technol*. 2022;143:106736.
14. Brown WH, Malveau RC, McCormick HWS, Mowbray TJ. *AntiPatterns: Refactoring Software, Architectures, and Projects in Crisis*. John Wiley & Sons, Inc; 1998.
15. García FJP. *Refactoring Planning for Design Smell Correction in Object-Oriented Software*. Universidad de Valladolid; 2011 Ph.D. thesis.
16. Tsantalís N, Chatzigeorgiou A. Identification of move method refactoring opportunities. *IEEE Trans. Softw Eng*. 2009;35:347–367.
17. Moha N, Guéhéneuc Y, Duchien L, Meur AL. DECOR: a method for the specification and detection of code and design smells. *IEEE Trans Softw Eng*. 2010;36:20–36.
18. Palomba F, Bavota G, Penta MD, Oliveto R, Lucia AD, Poshyvanyk D. Detecting bad smells in source code using change history information. In: Denney E, Bultan T, Zeller A, eds. *2013 28th IEEE/ACM International Conference on Automated Software Engineering, ASE 2013, Silicon Valley, CA, USA, November 11–15, 2013*. IEEE; 2013:268–278.
19. Fernandes E, Oliveira J, Vale G, Paiva T, Figueiredo E. A review-based comparative study of bad smell detection tools. *Proceedings of the 20th International Conference on Evaluation and Assessment in Software Engineering*. 2016 1–12.
20. Fontana FA, Mäntylä MV, Zanoni M, Marino A. Comparing and experimenting machine learning techniques for code smell detection. *Empir Softw Eng*. 2016;21:1143–1191.
21. Alazba A, Aljamaan H. Code smell detection using feature selection and stacking ensemble: an empirical investigation. *Inf. Softw. Technol*. 2021;138:106648.
22. Nucci DD, Palomba F, Tamburri DA, Serebrenik A, Lucia AD. Detecting code smells using machine learning techniques: Are we there yet? In: Oliveto R, Penta MD, Shepherd DC, eds. *25th International Conference on Software Analysis, Evolution and Reengineering, SANER 2018, Campobasso, Italy, March 20–23, 2018*. IEEE Computer Society; 2018:612–621.
23. Pecorelli F, Nucci DD, Roover CD, Lucia AD. A large empirical assessment of the role of data balancing in machine-learning-based code smell detection. *J. Syst. Softw*. 2020;169:110693.
24. Alkharabsheh K, Crespo Y, Manso ME, Taboada JA. Software design smell detection: a systematic mapping study. *Softw Qual J*. 2019;27:1069–1148.
25. Azeem MI, Palomba F, Shi L, Wang Q. Machine learning techniques for code smell detection: a systematic literature review and meta-analysis. *Inf. Softw. Technol*. 2019;108:115–138.
26. Al-Shaaby A, Aljamaan H, Alshayeb M. Bad smell detection using machine learning techniques: a systematic literature review. *Arab J Sci Eng*. 2020;45:2341–2369.
27. Fontana FA, Zanoni M. Code smell severity classification using machine learning techniques. *Knowl Based Syst*. 2017;128:43–58.
28. Palomba F, Bavota G, Di Penta M, Fasano F, Oliveto R, De Lucia A. On the diffuseness and the impact on maintainability of code smells: a large scale empirical investigation. *Proceedings of the 40th International Conference on Software Engineering*. 2018 482.
29. Pecorelli F, Di Nucci D, De Roover C, De Lucia A. On the role of data balancing for machine learning-based code smell detection. *Proceedings of the 3rd ACM SIGSOFT International Workshop on Machine Learning Techniques for Software Quality Evaluation*, 2019a; 2018 19–24.
30. Aljamaan H. Voting heterogeneous ensemble for code smell detection. In: Wani MA, Sethi IK, Shi W, Qu G, Raicu DS, Jin R, eds. *20th IEEE International Conference on Machine Learning and Applications, ICMLA 2021, Pasadena, CA, USA, December 13–16, 2021*. IEEE; 2021:897–902.
31. Wilcoxon F. Individual comparisons by ranking methods. *Breakthroughs in Statistics*. Springer; 1992:196–202.
32. Kampenes VB, Dybå T, Hannay JE, Sjøberg DI. A systematic review of effect size in software engineering experiments. *Inf Softw Technol*. 2007;49:1073–1086.
33. Kreimer J. Adaptive detection of design flaws. *Electron Notes Theor Comput Sci*. 2005;141:117–136.

34. Khomh F, Vaucher S, Guéhéneuc Y, Sahraoui HA. A bayesian approach for the detection of code and design smells. In: Choi B, ed. *Proceedings of the Ninth International Conference on Quality Software, QSIC 2009, Jeju, Korea, August 24-25, 2009*. IEEE Computer Society; 2009:305-314.
35. Khomh F, Vaucher S, Guéhéneuc Y, Sahraoui HA. BDTEX: a gqm-based bayesian approach for the detection of antipatterns. *J Syst Softw*. 2011;84:559-572.
36. Vaucher S, Khomh F, Moha N, Guéhéneuc Y. Tracking design smells: lessons from a study of god classes. In: Zaidman A, Antoniol G, Ducasse S, eds. *16th Working Conference on Reverse Engineering, WCRE 2009, 13-16 October 2009*. IEEE Computer Society; 2009:145-154.
37. Bryton S, Abreu FB, Monteiro MP. Reducing subjectivity in code smells detection: experimenting with the long method. In: Abreu FB, Faria JP, Machado RJ, eds. *Quality of Information and Communications Technology, 7th International Conference on the Quality of Information and Communications Technology, QUATIC 2010, Porto, Portugal, 29 September–2 October, 2010, Proceedings*. IEEE Computer Society; 2010:337-342.
38. Maiga A, Ali N, Bhattacharya N, Sabane A, Guéhéneuc Y, Aïmeur E. SMURF: a svm-based incremental anti-pattern detection approach. *19th Working Conference on Reverse Engineering, WCRE 2012, Kingston, ON, Canada, October 15-18*. IEEE Computer Society; 2012:466-475.
39. Amorim L, Costa E, Antunes N, Fonseca B, Ribeiro M. Experience report: evaluating the effectiveness of decision trees for detecting code smells. *26th IEEE International Symposium on Software Reliability Engineering, ISSRE 2015, Gaithersbury, MD, USA, November 2-5, 2015*. IEEE Computer Society; 2015:261-269.
40. Fontana FA, Zaroni M, Marino A, Mäntylä M. Code smell detection: towards a machine learning-based approach. *2013 IEEE International Conference on Software Maintenance, Eindhoven, the Netherlands, September 22-28, 2013*. IEEE Computer Society; 2013:396-399.
41. Kim DK. Finding bad code smells with neural network models. *Int J Electr Comput Eng*. 2017;7:3613.
42. Liu H, Xu Z, Zou Y. Deep learning based feature envy detection. In: Huchard M, Kästner C, Fraser G, eds. *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE 2018, Montpellier, France, September 3-7, 2018*. ACM; 2018:385-396.
43. Pecorelli F, Palomba F, Nucci DD, Lucia AD. Comparing heuristic and machine learning approaches for metric-based code smell detection. In: Guéhéneuc Y, Khomh F, Sarro F, eds. *2019b, Proceedings of the 27th International Conference on Program Comprehension, ICPC 2019, Montreal, QC, Canada, May 25-31, 2019*. IEEE / ACM; 2019:93-104.
44. Sharma T, Efstathiou V, Louridas P, Spinellis D. Code smell detection by deep direct-learning and transfer-learning. *J Syst Softw*. 2021;176:110936.
45. Yu J, Mao C, Ye X. A novel tree-based neural network for android code smells detection. *2021 IEEE 21st International Conference on Software Quality, Reliability and Security (QRS)*. IEEE; 2021:738-748.
46. Zhang Y, Dong C. Mars: detecting brain class/method code smell based on metric-attention mechanism and residual network. *J Softw: Evol Process*. 2021:e2403. doi:10.1002/smr.2403
47. Li Y, Zhang X. Multi-label code smell detection with hybrid model based on deep learning.
48. Zhang Y, Ge C, Hong S, Tian R, Dong C, Liu J. Delesmell: code smell detection based on deep learning and latent semantic analysis. *Knowl-Based Syst*. 2022;255:109737.
49. Shen L, Liu W, Chen X, Gu Q, Liu X. Improving machine learning-based code smell detection via hyper-parameter optimization. *2020 27th Asia-Pacific Software Engineering Conference (APSEC)*. IEEE; 2020:276-285.
50. Akhter N, Rahman S, Taher KA. An anti-pattern detection technique using machine learning to improve code quality. *2021 International Conference on Information and Communication Technology for Sustainable Development (ICICT4SD)*. IEEE; 2021:356-360.
51. Alkharabsheh K, Crespo Y, Delgado MF, Viqueira JRR, Taboada JA. Exploratory study of the impact of project domain and size category on the detection of the god class design smell. *Softw. Qual. J*. 2021;29:197-237.
52. Gupta H, Kulkarni TG, Kumar L, Neti LBM, Krishna A. An empirical study on predictability of software code smell using deep learning models. *International Conference on Advanced Information Networking and Applications*. Springer; 2021:120-132.
53. Jain S, Saha A. Improving performance with hybrid feature selection and ensemble machine learning techniques for code smell detection. *Sci Comput Program*. 2021;212:102713.
54. Patnaik A, Padhy N. Does code complexity affect the quality of real-time projects? Detection of code smell on software projects using machine learning algorithms. *Proceedings of the International Conference on Data Science, Machine Learning and Artificial Intelligence*. ACM; 2021:178-185.
55. Stefano MD, Pecorelli F, Palomba F, Lucia AD. Comparing within- and cross-project machine learning algorithms for code smell detection. In: Ampatzoglou A, Feitosa D, Catolino G, Lenarduzzi V, eds. *MaLTesQuE@ESEC/SIGSOFT FSE 2021: Proceedings of the 5th International Workshop on Machine Learning Techniques for Software Quality Evolution, Athens, Greece, 23 August 2021*. ACM; 2021:1-6.
56. Khleel NAA, Nehéz K. Deep convolutional neural network model for bad code smells detection based on oversampling method. *Indo J Electr Eng Comput Sci*. 2022;26:1725-1735.
57. Kovačević A, Slivka J, Vidaković D, et al. Automatic detection of long method and god class code smells through neural source code embeddings. *Expert Syst Appl*. 2022;204:117607.
58. Nanda J, Chhabra JK. Sshn: smote-stacked hybrid model for improving severity classification of code smell. *Int J Inf Technol*. 2022;14:2701-2707.
59. Yedida R, Menzies T. How to improve deep learning for software analytics (a case study with code smell detection). 2022 arXiv preprint arXiv:2202.01322.
60. Santos G, Santana A, Vale G, Figueiredo E. Yet another model! A study on model's similarities for defect and code smells. *Fundamental Approaches to Software Engineering: 26th International Conference, FASE 2023, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2023, Paris, France, April 22–27, 2023, Proceedings*. Springer; 2023:282-305.

61. Singh S, Jalal M, Kaur S. Bug classification depend upon refactoring area of code. *J Inst Eng (India): Ser B*. 2023;104:1-17.
62. Zhou Y, Yang Y, Lu H, et al. How far we have progressed in the journey? An examination of cross-project defect prediction. *ACM Trans Softw Eng Methodol*. 2018;27:1.
63. Galar M, Fernández A, Tartas EB, Sola HB, Herrera F. A review on ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based approaches. *IEEE Trans Syst Man Cybern Part C*. 2012;42:463-484.
64. Boutaib S, Bechikh S, Palomba F, Elarbi M, Makhoul M, Said LB. Code smell detection and identification in imbalanced environments. *Expert Syst Appl*. 2021;166:114076.
65. Ma X, Keung J, Yang Z, Yu X, Li Y, Zhang H. Casms: combining clustering with attention semantic model for identifying security bug reports. *Inf Softw Technol*. 2022;147:106906.
66. Zhen Y, Keung JW, Xiao Y, Yan X, Zhi J, Zhang J. On the significance of category prediction for code-comment synchronization. *ACM Trans Softw Eng Methodol*. 2022;32:1-41.
67. Chen Y, Dai H, Yu X, Hu W, Xie Z, Tan C. Improving ponzi scheme contract detection using multi-channel textcnn and transformer. *Sensors*. 2021;21:6417.
68. Li F, Lu W, Keung JW, Yu X, Gong L, Li J. The impact of feature selection techniques on effort-aware defect prediction: an empirical study. *IET Softw*. 2023;17:168-193.
69. Yu X, Liu J, Keung JW, et al. Improving ranking-oriented defect prediction using a cost-sensitive ranking svm. *IEEE Trans Reliab*. 2019;69:139-153.
70. Zhu K, Ying S, Ding W, Zhang N, Zhu D. Ivkmp: a robust data-driven heterogeneous defect model based on deep representation optimization learning. *Inform Sci*. 2022;583:332-363.
71. Zhu K, Ying S, Zhang N, Zhu D. Software defect prediction based on enhanced metaheuristic feature selection optimization and a hybrid deep neural network. *J Syst Softw*. 2021;180:111026.
72. Chen Y, Lu X, Wang S. Deep cross-modal image-voice retrieval in remote sensing. *IEEE Trans Geosci Remote Sens*. 2020;58:7049-7061.
73. He C, Wu J, Zhang Q. Characterizing research leadership on geographically weighted collaboration network. *Scientometrics*. 2021;126:4005-4037.
74. Chen Y, Xiong S, Mou L, Zhu XX. Deep quadruple-based hashing for remote sensing image-sound retrieval. *IEEE Trans Geosci Remote Sens*. 2022;60:1-14.
75. He C, Wu J, Zhang Q. Proximity-aware research leadership recommendation in research collaboration via deep neural networks. *J Assoc Inf Sci Technol*. 2022;73:70-89.
76. Acuna E, Rodriguez C. The treatment of missing values and its effect on classifier accuracy. *Classification, Clustering, and Data Mining Applications*. Springer; 2004:639-647.
77. Young W, Weckman G, Holland W. A survey of methodologies for the treatment of missing values within datasets: limitations and benefits. *Theor Issues Ergon Sci*. 2011;12:15-43.
78. Mundfrom DJ, Whitcomb A. Imputing missing values: The effect on the accuracy of classification. 1998.
79. John GH, Kohavi R, Pfleger K. Irrelevant features and the subset selection problem. In: Cohen WW, Hirsh H, eds. *Machine Learning, Proceedings of the Eleventh International Conference, Rutgers University, New Brunswick, NJ, USA, July 10-13, 1994*. Morgan Kaufmann; 1994:121-129.
80. Quinlan JR. Induction of decision trees. *Mach Learn*. 1986;1:81-106.
81. Bergstra J, Bengio Y. Random search for hyper-parameter optimization. *J Mach Learn Res*. 2012;13:281-305.
82. LaValle SM, Branicky MS, Lindemann SR. On the relationship between classical grid search and probabilistic roadmaps. *Int J Robot Res*. 2004;23:673-692.
83. Dewangan S, Rao RS, Mishra A, Gupta M. A novel approach for code smell detection: an empirical study. *IEEE Access*. 2021;9:162869-162883.

**How to cite this article:** Li F, Zou K, Keung JW, Yu X, Feng S, Xiao Y. On the relative value of imbalanced learning for code smell detection. *Softw Pract Exper*. 2023;53(10):1902-1927. doi: 10.1002/spe.3235