

Data Balancing Improves Self-Admitted Technical Debt Detection

Murali Sridharan¹, Mika Mantyla², Leevi Rantala³, Maelick Claes⁴

M3S, ITEE

University of Oulu

Oulu, Finland

{murali.sridharan¹,mika.mantyla²,leevi.rantala³,maelick.claes⁴}@oulu.fi

Abstract—A high imbalance exists between technical debt and non-technical debt source code comments. Such imbalance affects Self Admitted Technical Debt (SATD) detection performance, and existing literature lacks empirical evidence on the choice of balancing technique. In this work, we evaluate the impact of multiple balancing techniques, including Data level, Classifier level, and Hybrid, for SATD detection in Within-Project and Cross-Project setup. Our results show that the Data level balancing technique SMOTE or Classifier level Ensemble approaches with Random Forest or XGBoost are reasonable choices depending on whether the goal is to maximize Precision, Recall, F1, or AUC-ROC. We compared our best-performing model with the previous SATD detection benchmark (cost-sensitive Convolution Neural Network). Interestingly the top-performing XGBoost with SMOTE sampling improved the Within-project F1 score by 10% but fell short in Cross-Project set up by 9%. This supports the higher generalization capability of deep learning in Cross-project SATD detection, yet while working within individual projects, classical machine learning algorithms can deliver better performance. We also evaluate and quantify the impact of duplicate source code comments in SATD detection performance. Finally, we employ SHAP and discuss the interpreted SATD features. We have included the replication package¹ and shared a web-based SATD prediction tool² with the balancing techniques in this study.

Index Terms—Self-Admitted Technical Debt, data imbalance, classification, data sampling techniques, cost-sensitive technique, ensemble techniques

I. INTRODUCTION

Software development is often hampered by time-pressure. The quick-fix mentality that focuses merely on the immediate goal and short-term benefit has been the norm [1]. The quick fixes often turn out to be sub-optimal, as they lack a holistic approach to software maintenance making the source code more rigid without room for future enhancements. They incur a substantial cost in terms of time and man-effort to refactor the code at a later stage. Such accumulated debt by choosing quick delivery over quality delivery is known as *technical debt* in Software Engineering.

The early detection of such technical debt would be instrumental in reducing the increased software maintenance cost. The established approach among the software practitioners have been to utilize static code analysis for improving

code quality. Often, the developer acknowledged or developer induced hacky patch/workaround goes unnoticed. Such patch/workaround are expressed by the software developers through source code comments. These are termed as Self-Admitted Technical Debt comments by Potdar et al. [2]. Such SATD comments from the source code has vital information about the source code segments that needs refactoring.

The detection of technical debt from SATD source code comments has gathered significant interest in the recent past. A crucial challenge associated with technical debt detection from source code comments is the imbalanced distribution among SATD and non-SATD data instances. The supervised machine learning approaches expect a uniform distribution among the data samples for optimal prediction performance but quite often end up with imbalanced data which will affect the prediction capability. Longadge et al. [3] state that imbalanced/skewed data increase the False Negatives (FN), which will decrease the Recall of the minor class (SATD comments in our context). Thabtah et al. [4] highlight the varying effect of data imbalance for each evaluation metric particularly Precision, Recall and ROC-AUC. The challenge of generating more accurate results while accounting for the class imbalance in the training data is paramount for reliable inference. In reality, this data imbalance problem is very common among multiple tasks in software engineering domain and other domains as well. Previous work on technical debt detection discusses different machine learning approaches but very few have employed balancing techniques for addressing the class imbalance. For example Pecorelli et al. [5] used data sampling on code metrics for detecting code smells and Ren et al. [6] used classifier level balancing technique COST to improve SATD detection from source code comments.

To the best of our knowledge, we have performed the first extensive empirical study on evaluating multiple balancing techniques for technical debt detection from source code comments. These include data-level balancing techniques (SMOTE, ADASYN, BorderLine SMOTE and SVMSMOTE), classifier-level balancing techniques (COST and Ensemble) and hybrid balancing techniques that combine either data-level and classifier-level balancing techniques or employ customised combined algorithms, for example COST based ensemble algorithm or sampling based ensemble, for handling data imbalance scenarios.

¹<https://figshare.com/s/87a4b5002c7488822e60>

²<https://balancing-technical-debt.herokuapp.com/>

We empirically evaluate three categories of balancing techniques and their impact on technical debt detection from source code comments using regression, bagging and boosting classifiers. More specifically, we compare the performance, in terms of precision, recall, F1 score and area under the ROC curve (ROC-AUC), against the BASELINE approach which employs the same classifier without COST or data sampling balancing techniques. Our results enable software engineering researchers and practitioners to choose relevant balancing technique depending on the use case and the evaluation metric in focus.

In this paper, we answer the following research questions:

- **Balancing the Imbalanced Data: Main RQ** Which balancing technique contribute to better Technical Debt Detection in highly imbalanced source code comments data?
 - **RQ 1 Precise detection of SATD comments:** Which balancing technique consistently contribute to better precision in classifying SATD comments?
 - **RQ 2 Extensive detection of SATD comments:** Which balancing technique consistently contribute to better recall while classifying SATD comments?
 - **RQ 3 Improving distinction capability of classifier for SATD detection:** Which balancing technique consistently improve the distinction capability of the machine learning model (ROC-AUC) to distinguish between SATD and non-SATD comments?
 - **RQ 4 Balancing Precision and Recall for SATD Detection:** Which balancing technique consistently improve the overall classification performance (F1) for detecting SATD comments?

Our main contributions in this paper include:

- Evaluation of multiple balancing schemes for SATD detection from highly imbalanced source code comments data.
- Recommendations for choosing a balancing technique for SATD detection through source code comments, depending on the evaluation metric in focus.
- Impact analysis of data level, classifier level and hybrid balancing techniques on evaluation metrics such as Precision, Recall, F1 and ROC-AUC scores.
- Web-based SATD detection tool in batch and online modes based on the data from 10 open source projects.
- Replication package of our experiments for verification and further extension.

The rest of this paper is organized as follows. First, in Section II, we discuss the evolution of TD and the past works associated with the detection of technical debt. In Section III, we discuss the classifiers, data characteristics, and the different balancing techniques studied in this work. We tabulated and discussed our experimental results in Section IV. In Section V, we discuss the technique used for feature interpretation and the impact of balancing techniques on features. We discuss the implications associated with choosing an appropriate balancing technique in Section VI. The potential factors that could

mine the validity of our results are elaborated in Section VII. We conclude with our main findings from this study in Section VIII followed by acknowledgements.

II. BACKGROUND

A. Technical Debt

The Technical Debt (TD) metaphor originates from the early 90's, when Cunningham said that shipping immature code resembles like taking a debt [7]. TD as a category does not represent a uniform type of debt. Previous study by Alves et al. [8] lists 13 categories for TD including architectural debt, code debt, design debt and people debt. The most recent definition for TD can be found from Avgeriou et al.'s. [9] work, where it is defined as:

In software-intensive systems, technical debt is a collection of design or implementation constructs that are expedient in the short term, but set up a technical context that can make future changes more costly or impossible. Technical debt presents an actual or contingent liability whose impact is limited to internal system qualities, primarily maintainability and evolvability.

In 2014, Potdar et al. [2] coined the term self-admitted technical debt (SATD). SATD is a subset of TD, where the developers have left a note admitting that they intentionally incurred TD. The software practitioners leave comments in the source code which is crucial for detecting SATD in the source code. These have two significant benefits, firstly, it is easier and more efficient to detect SATD from source code, and it also does not need to rely on predefined metrics which are difficult to determine as suggested by Maldonado et al. [10].

B. Machine Learning in SATD Detection

Zampetti et al. [11] studied how five machine learning techniques were able to recommend the software practitioners when to admit design debt using structural and readability metrics from source code and warnings from static analysis tools. They achieved an average precision 0.5 and recall 0.52 for SATD detection within project. For cross-project SATD detection they achieved 0.67 precision and 0.55 recall. Unlike them, we do not study different code metrics, but rather evaluate the impact of multiple balancing schemes using different machine learning techniques when detecting SATD from source code comments.

Maldonado et al. [10], developed a maximum entropy classifier for detecting design and requirement SATD from code comments. The classifier builds a maximum entropy model, which is described to be equivalent to multi-class regression model. Their data set consist of 10 projects, which represent different application domains. They used leave-one-out cross-project validation, which means that the classifier was trained on 9 projects and tested against the one which was left out. They achieved an average F1 of 0.620 (min. 0.470, max. 0.814) for design SATD, and F1 of 0.403 (min. 0.154, max. 0.804) for requirement SATD.

In another related work by Huang et al. [12], they developed a machine learning model, which relied on several sub-classifiers. Each of the sub-classifiers used Naive Bayes (NB) multinomial with feature selection technique, and the classification of whether a code comment has SATD or not was based on majority voting done by these classifiers. The data set consisted of 8 different open source projects. Each sub-classifier was trained with leave-one-out cross-project validation, where they used 7 projects for training, and left one project out for testing. Each sub-classifier left out a different project. They achieved an average F1 of 0.737 (min. 0.518, max. 0.841).

Ren et al. [6], trained a Convolution Neural Network (CNN) for SATD detection from source code comments using data from 10 open source Java projects. They evaluated their approach in Within and Cross-project setup. For Within-project they obtained an average F1 of 0.752 (min. 0.445, max. 0.932), and for Cross-Project the average F1 was 0.766 (min. 0.599, max. 0.878).

C. Imbalanced Data

Technical Debt detection through source code comments suffers from imbalanced data problem. Besides SATD detection [6], it has also been investigated on code smell detection [5], [13].

Ren et al. [6], in their work experiment with normal and weighted cross-entropy loss functions. For addressing imbalance in data, they employed COST (assigning class weight). For weight calculation, n represent the total number of SATD comments, and m the total amount of non-SATD comments

- SATD Comment Weight = $\frac{n}{n+m}$
- Non-SATD Comment Weight = $\frac{m}{n+m}$

The CNN with weighted cross-entropy loss function produced better results than the CNN using normal cross-entropy loss function. The authors note that the effectiveness of the weighted cross-entropy loss function increased based on the imbalance in the data [6].

Huang et al. [12] employed feature selection technique to extract features from source code comments during training. They extracted features from 8 open source projects and combined multiple classifiers into a composite classifier for SATD comments detection task. They claimed a performance improvement of 409% over the work by Potdar et al. [2].

In [13], the authors investigated five different data balancing methods when detecting code smells with source code metrics. They examined 125 releases from 13 open source systems from a dataset introduced in previous study [5]. The dataset contained over 8,500 manually validated code smells. Even with such high sounding number, the code smells formed a very small minority in the data. One example is the instance of God Class smell, which at its highest count in a release formed only around 1% of all the classes present in that release [5]. To overcome this problem, the authors in [13] tried five different data balancing methods, which were Class Balancer, resampling with balancing the dataset, SMOTE, Cost Sensitive Classifier and One Class Classifier. When compared to the

BASELINE of no data balancing, the authors report highest performance overall with SMOTE, although there were differences between different smells. Our work is the first extensive empirical study to address the data imbalance for NLP-based SATD detection. We extend Pecorelli et al., [13] by evaluating multiple variants of SMOTE and include COST, Bagging, Boosting (classifier level ensemble techniques), and hybrid techniques for SATD detection from source code comments. We focus on the impact of balancing technique on evaluation metrics and provide a recommendation for each metric for both Within-Project and Cross-Project SATD detection from empirical results.

III. METHODOLOGY

In this section we describe the dataset used, discussed the various balancing techniques employed in our empirical study along with the experiment set up and evaluation criteria.

A. Machine learning classifiers

Cruz et al. [14] empirically evaluated seven different machine learning algorithms for bad smell detection and found Random Forest and XGBoost achieved better overall performance. They have acknowledged the imbalanced data distribution in their data but have not employed explicit data balancing technique. This created undue advantage for the machine learning ensemble algorithms Random Forest and XGBoost which employ implicit balancing and are categorized under algorithm level (classifier level) balancing techniques, refer in III-B2. In our study, we include three machine learning classifiers: Regression (Logistic Regression), Bagging (Random Forest) and Boosting (XGBoost). We choose Logistic regression for its superior performance in binary classification and its widespread usage in the past for NLP related tasks in software engineering [14]–[17]. Random Forest [18], is an ensemble algorithm (classifier level balancing technique explained in III-B2), based on the Bagging concept proposed by Breiman et al. [19]. It is aimed at reducing the variance associated with a model to increase the prediction capability. Random Forest has shown significant performance improvement in many software engineering tasks [11], [14], [20]–[22]. XGBoost, proposed by [23] is a scalable, tree-based ensemble algorithm based on the Boosting concept proposed by Freund and Shapire [24].

B. Balancing Techniques

The balancing techniques could be broadly categorized into three groups which include Data Level techniques, Classifier Level techniques and Hybrid techniques, see Figure 1.

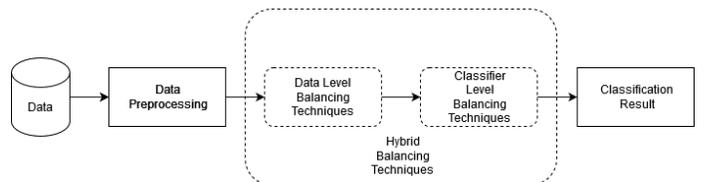


Fig. 1. Balancing Techniques

1) *Data Level Techniques*: These techniques primarily focus on different data sampling schemes to reduce the skewness ratio in a dataset. Data Sampling is a host of techniques that transforms the distribution of the dataset with the objective of reducing the imbalance ratio in the dataset. Such schemes could be grouped into two categories, Oversampling and Undersampling. Oversampling is creating new samples from minority class either by mere duplication or by synthesizing samples. Undersampling is reducing the data samples from majority class by removing or selecting a subset of data samples. Multiple research including [25]–[27] establish the significance of oversampling over undersampling techniques. We include multiple variants of oversampling technique in our study.

a) *SMOTE*: One of the widely used oversampling technique, Synthetic Minority Oversampling Technique (SMOTE) [28]. In this, the minority class is over-sampled and new samples are formed based on randomly chosen samples that are close to each other. At first, random neighbours that are close to each other in the minority class are identified, then new samples are synthesized between any two neighbours and the same is replicated with other chosen neighbours. SMOTE has received a lot of success in data mining community. This success has been attributed to four factors, which are its simplicity, adaptability, generalizability and performance improvements [29].

b) *Borderline SMOTE*: Another variant of SMOTE, wherein data samples are not synthesized from the minority class randomly but with a specific target scheme. Borderline SMOTE (BLINE) was developed [30] with the focus to reduce the misclassification rate. The focus is on those data samples that are incorrectly classified and are lying near the ambiguous decision boundary of any two class X and Y. The target scheme is such that the data samples are synthesized near the decision boundary which segregates a class X from another class Y.

c) *ADASYN*: Adaptive Synthetic (ADASYN) sampling is a SMOTE-based over-sampling technique. The target scheme for this sampling technique is such that the new samples are synthesized in those areas where the density of the data samples is sparse, or where they are very few instances. The data samples in the high density regions are left undisturbed. This technique [31] has claimed to account for correct classification of difficult instances as it adaptively shifts the decision boundary.

d) *SVMSMOTE*: SVMSMOTE(SVMSMT) is another SMOTE based over-sampling technique which employs SVM algorithm [32] to generate synthetic samples. The target scheme for synthesizing new samples of SVMSMOTE algorithm is focused around the hyperplane that separates each class.

2) *Classifier Level Techniques*: Typically, the balancing schemes implemented inherently in the classification algorithms are referred to as classifier-level balancing techniques. Such techniques include, COST (assigning class weight) and a group of homogeneous classifiers (Bagging and Boosting) collectively working on the data samples for classification or

prediction task. The ensemble algorithms Random Forest and XGBoost are explored as balancing techniques in [33]–[36]. Another classifier level technique is COST (or class weighted) in which higher weights are assigned to the class samples which are less in the training set so that it would penalize the model more during training for every incorrect prediction of the less represented class. This technique is simpler and efficient as it does not involve additional step of sample synthesis as in data sampling techniques.

3) *Hybrid Techniques*: The hybrid balancing technique combines multiple balancing techniques to effectively handle the imbalance scenario. The ensemble technique which can handle the imbalance scenario is further augmented with either sampling or cost balancing techniques. In our study, we explore, Random Forest and XGBoost ensemble classifiers augmented with sampling and cost techniques.

C. Dataset

For our experiments, we use the dataset created with 10 open source Java projects from Maldonado et al. [10]. In their work, they extract and filter the source code comments from the 10 open source projects using heuristics including removal of license comments, task annotations, Javadoc comments and automatically-generated comments. The filtered comments were manually labelled and grouped into the following categories: Requirement Debt, Design Debt, Implementation Debt, Test Debt, Documentation Debt and Non-Technical Debt. We consolidated all the types of TD into a single TD class to facilitate binary classification. We have two classes: SATD and non-SATD class for simplicity and easier evaluation across several balancing techniques. Table I lists the statistics of the source code comments data.

TABLE I
DATASET STATISTICS

| Project | Release | # of Cmnts [6] | # of Cmnts* | SATD | % of SATD | % of Non-SATD |
|----------------|---------|----------------|--------------|------------|-------------|---------------|
| Apache Ant | 1.7.0 | 4,137 | 2,992 | 123 | 4.11 | 95.89 |
| ArgoUML | 0.34 | 9,548 | 4,725 | 1,137 | 24.06 | 75.94 |
| Columba | 1.4 | 6,478 | 3,842 | 158 | 4.11 | 95.89 |
| EMF | 2.4.1 | 4,401 | 2,382 | 81 | 3.40 | 96.60 |
| Hibernate | 3.3.2 | 2,968 | 2,458 | 423 | 17.21 | 82.79 |
| JEdit | 4.2 | 10,322 | 4,506 | 228 | 5.06 | 94.94 |
| JFreeChart | 1.0.19 | 4,423 | 2,280 | 106 | 4.65 | 95.35 |
| JMeter | 2.10 | 8,162 | 4,029 | 314 | 7.79 | 92.21 |
| JRuby | 1.4.0 | 4,897 | 2,991 | 458 | 15.31 | 84.69 |
| Squirrel | 3.0.3 | 7,230 | 4,219 | 216 | 5.12 | 94.88 |
| Average | | 6,257 | 3,442 | 324 | 9.08 | 90.91 |

*After Duplicate Removal and Pre-processing in our study

D. Processing Code Comments

We processed the code comments using following steps:

- **Duplicate removal**: Duplicate comments are those that appear more than once in the training data. Example: '// Need to calculate this... just fudging here for now' source code comment appear three times in the training data. We identified such duplicates and removed 23,511 (37.75%) duplicate source code comments from the training data

TABLE II
CONSOLIDATED RESULTS: WITHIN-PROJECT AND CROSS-PROJECT
PRECISION SCORES

| SCHEME | | | | | |
|----------------|-----------------|---------------------|--------------------------|---------------------|-------|
| Within-Project | | HYBRID | | | |
| | METHOD | LR | RF | XGB | Avg. |
| | BASELINE | <u>0.710</u> | 0.898 | 0.839 ^w | 0.816 |
| COST SENSITIVE | WEIGHTED | 0.743 ^b | 0.904 | 0.788 | 0.812 |
| DATA SAMPLING | SMOTE | 0.779 | 0.904 | 0.851 ^w | 0.845 |
| | BLINE | 0.825 ^{bw} | 0.902 | 0.816 | 0.848 |
| | ADASYN | 0.776 | 0.905 | 0.804 | 0.828 |
| | SVMSMT | 0.850 ^{bw} | 0.895 | 0.835 ^w | 0.860 |
| | Avg. | 0.780 | 0.901 | 0.822 | |
| Cross-Project | | | | | |
| | BASELINE | 0.885 ^{ws} | 0.893^s | 0.863 ^{ws} | 0.880 |
| COST SENSITIVE | WEIGHTED | <u>0.623</u> | 0.881 ^s | 0.713 | 0.739 |
| DATA SAMPLING | SMOTE | 0.690 | 0.788 | 0.764 | 0.747 |
| | BLINE | 0.724 ^w | 0.825 | 0.810 | 0.786 |
| | ADASYN | 0.683 | 0.744 | 0.793 | 0.740 |
| | SVMSMT | 0.719 ^w | 0.854 | 0.811 | 0.795 |
| | Avg. | 0.721 | 0.831 | 0.788 | |

Highest score is bolded and lowest score is underlined.
b,w,s - Statistically significant change from baseline, weighted, and sampling respectively based on Wilcoxon-Signed Rank test with 95% confidence level

In Within-Project set up, the data sampling balancing technique SVMSMT/BLINE has improved the BASELINE and COST overall. However, the BASELINE ensemble classifier Random Forest has significantly improved precision scores of COST and Sampling techniques of other classifiers. For cross-project SATD detection, the BASELINE has achieved the highest precision score while the COST and the data sampling techniques decreases the precision. Although the BASELINE scores appear almost equivalent, the ensemble classifier Random Forest has the lowest number of FP per TP with 0.11, followed by Logistic Regression with 0.13. Overall, neither COST nor Sampling consistently improve the precision scores over BASELINE across classifiers. The BASELINE ensemble classifier Random Forest contributes more to the precise detection of SATD comments in both Within and Cross-Project setup. However, between COST and sampling technique, the sampling techniques SVMSMT/BLINE has significantly improved the precision scores across all the classifiers over COST, in both Within and Cross-Project setup.

Recommendation: For consistent precise detection of SATD comments from imbalanced data, use classifier level Bagging ensemble technique (Random Forest).

B. RQ 2 Extensive detection of SATD comments

Our motivation is to determine the balancing technique that consistently contribute to extensive detection of SATD comments in highly skewed data. Tables III lists the Recall scores for both Within and Cross-Project set up respectively.

The Logistic Regression with COST balancing technique achieved the highest Recall scores in both Within and Cross-Project setup with 0.735 and 0.754 respectively. The sampling

TABLE III
CONSOLIDATED RESULTS: WITHIN-PROJECT AND CROSS-PROJECT
RECALL SCORES

| SCHEME | | | | | |
|----------------|-----------------|--------------------------|---------------------|---------------------|-------|
| Within-Project | | HYBRID | | | |
| | METHOD | LR | RF | XGB | Avg. |
| | BASELINE | <u>0.277</u> | 0.580 | 0.611 | 0.489 |
| COST SENSITIVE | WEIGHTED | 0.735^b | 0.586 | 0.730 ^{bs} | 0.684 |
| DATA SAMPLING | SMOTE | 0.726 ^b | 0.653 ^{bw} | 0.688 ^b | 0.689 |
| | BLINE | 0.704 | 0.619 | 0.641 | 0.655 |
| | ADASYN | 0.731 ^b | 0.609 ^b | 0.678 ^b | 0.673 |
| | SVMSMT | 0.658 | 0.644 | 0.644 | 0.649 |
| | Avg. | 0.638 | 0.615 | 0.665 | |
| Cross-Project | | | | | |
| | BASELINE | <u>0.502</u> | 0.617 | 0.657 | 0.592 |
| COST SENSITIVE | WEIGHTED | 0.754^b | 0.590 | 0.723 ^b | 0.689 |
| DATA SAMPLING | SMOTE | 0.727 ^b | 0.674 ^{bw} | 0.702 ^b | 0.701 |
| | BLINE | 0.697 | 0.662 | 0.683 | 0.681 |
| | ADASYN | 0.730 ^b | 0.662 ^{bw} | 0.697 ^b | 0.696 |
| | SVMSMT | 0.676 | 0.640 | 0.683 | 0.666 |
| | Avg. | 0.681 | 0.641 | 0.691 | |

Highest score is bolded and lowest score is underlined.
b,w - Statistically significant change from baseline and weighted respectively based on Wilcoxon-Signed Rank test at 95% confidence level

techniques SMOTE and ADASYN has consistently improved the Recall scores for all three classifiers over the BASELINE in Within and Cross-Project setup. Although the COST technique has recorded the highest recall scores in both the Within and Cross-Project setup, the improvement is not statistically significant over SMOTE and ADASYN sampling techniques. Further, SMOTE and ADASYN consistently improved the BASELINE with a statistically significant change while the COST technique is not consistent across all three classifiers.

Recommendation: For consistent extensive detection of SATD comments from imbalanced data, use SMOTE/ADASYN sampling technique for both Within and Cross-Project setup.

C. RQ 3 Improving distinction capability of classifier for SATD detection

Here our objective is to determine the balancing technique that contributes more to ROC-AUC score. ROC-AUC characterizes the distinctive capability of the machine learning model between technical debt and non-technical debt classes. Tables IV lists the Recall scores for both Within and Cross-Project set up respectively.

The COST balancing technique with Logistic Regression/XGBoost recorded the highest ROC-AUC scores in both Within and Cross-Project setup. The sampling techniques SMOTE/ADASYN consistently improved the BASELINE ROC-AUC scores in both Within and Cross-Project SATD detection. The COST, on the other hand, is not consistent in improving the BASELINE across all the classifiers. Further the top score achieved with COST is not statistically significant over sampling techniques SMOTE/ADASYN. It is evident

TABLE IV

CONSOLIDATED RESULTS: WITHIN-PROJECT AND CROSS-PROJECT AUC SCORES

| SCHEME | | HYBRID | | | | |
|----------------|--|-----------------|--------------------------|---------------------|--------------------------|-------|
| Within-Project | | METHOD | LR | RF | XGB | Avg. |
| | | <u>BASELINE</u> | <u>0.636</u> | 0.786 | 0.799 | 0.740 |
| COST SENSITIVE | | WEIGHTED | 0.855 ^b | 0.790 | 0.857^b | 0.834 |
| DATA SAMPLING | | SMOTE | 0.854 ^b | 0.822 ^{bw} | 0.838 ^b | 0.838 |
| | | BLINE | 0.844 | 0.806 | 0.814 | 0.821 |
| | | ADASYN | 0.856 ^b | 0.818 ^b | 0.832 ^b | 0.835 |
| | | SVMSMT | 0.822 | 0.786 | 0.816 | 0.808 |
| | | Avg. | 0.811 | 0.804 | 0.826 | |
| Cross-Project | | <u>BASELINE</u> | <u>0.747</u> | 0.804 | 0.824 | 0.792 |
| COST SENSITIVE | | WEIGHTED | 0.857^b | 0.790 | 0.850 ^b | 0.832 |
| DATA SAMPLING | | SMOTE | 0.849 ^b | 0.829 ^{bw} | 0.840 | 0.839 |
| | | BLINE | 0.837 | 0.825 | 0.835 | 0.832 |
| | | ADASYN | 0.849 ^b | 0.823 ^w | 0.840 | 0.837 |
| | | SVMSMT | 0.827 | 0.814 | 0.835 | 0.825 |
| | | Avg. | 0.828 | 0.814 | 0.838 | |

Highest score is bolded and lowest score is underlined.
b,w - Statistically significant from baseline, weighted respectively based on Wilcoxon-Signed Rank test with 95% confidence level

from Table III and Table IV that, the COST decreases the recall and ROC-AUC scores for Random Forest, while the sampling techniques SMOTE/ADASYN improves them.

Recommendation: For consistent better distinctive capability of the classifier (AUC-ROC) from imbalanced data, use SMOTE in both Within and Cross-Project setup.

D. RQ 4 Balancing Precision and Recall for SATD Detection

Here, our motivation is to determine the balancing technique that has consistently contributed more towards higher F1 score. Tables V lists the Recall scores for both Within and Cross-Project set up respectively.

XGBoost with COST and BASELINE XGBoost have recorded top F1 scores in Within and Cross-Project setup with 0.755 and 0.729 respectively. In Within-Project setup, the SMOTE balancing technique has consistently improved the F1 scores over BASELINE for all the classifiers, while COST is not consistent across classifiers. Further, the top F1 score with COST is not statistically significant over the F1 scores with BLINE/SMOTE. In Cross-Project setup, although the Boosting based ensemble classifier BASELINE XGBoost achieved the top Cross-Project F1 score of 0.729, it is not statistically significant over sampling techniques or over COST. The sampling techniques BLINE/SMOTE have a statistically significant improvement over BASELINE for Logistic Regression. For the ensemble classifiers Random Forest and XGBoost, the sampling technique have improved over COST and achieved an equivalent performance as that of the BASELINE. Overall, the results show that the sampling

TABLE V

CONSOLIDATED RESULTS: WITHIN-PROJECT AND CROSS-PROJECT F1 SCORES

| SCHEME | | HYBRID | | | | |
|----------------|--|-----------------|---------------------|---------------------|--------------------------|-------|
| Within-Project | | METHOD | LR | RF | XGB | Avg. |
| | | <u>BASELINE</u> | <u>0.361</u> | 0.670 | 0.694 | 0.575 |
| COST SENSITIVE | | WEIGHTED | 0.730 ^b | 0.680 | 0.755^b | 0.722 |
| DATA SAMPLING | | SMOTE | 0.748 ^b | 0.738 ^{bw} | 0.753 ^b | 0.746 |
| | | BLINE | 0.749 | 0.706 | 0.710 | 0.722 |
| | | ADASYN | 0.750 ^b | 0.703 | 0.730 | 0.728 |
| | | SVMSMT | 0.730 | 0.660 | 0.717 | 0.702 |
| | | Avg. | 0.678 | 0.704 | 0.726 | |
| Cross-Project | | <u>BASELINE</u> | <u>0.612</u> | 0.705 ^w | 0.729 | 0.682 |
| COST SENSITIVE | | WEIGHTED | 0.674 ^b | 0.679 | 0.712 | 0.688 |
| DATA SAMPLING | | SMOTE | 0.700 ^{bw} | 0.713 ^w | 0.723 | 0.712 |
| | | BLINE | 0.700 ^{bw} | 0.718 ^w | 0.726 | 0.715 |
| | | ADASYN | 0.697 | 0.708 | 0.721 | 0.709 |
| | | SVMSMT | 0.684 | 0.711 | 0.728 | 0.708 |
| | | Avg. | 0.678 | 0.706 | 0.723 | |

Highest score is bolded and lowest score is underlined.
b,w - Statistically significant from baseline, weighted respectively based on Wilcoxon-Signed Rank test with 95% confidence level

techniques BLINE/SMOTE are more consistent in improving the F1 scores for all the classifiers in Cross-Project setup.

Recommendation: For consistent higher balance of Precision and Recall (F1-score), use SMOTE for Within-Project and BLINE/SMOTE for Cross-project SATD detection.

E. Comparison to prior work

Here we compare the F1 scores of our consistent top-performing model in both Within and Cross-project setup with previous works. First, we compare against the state of the art results from cost-sensitive CNN by Ren et al. [6]. Then, we evaluate the performance improvement over Huang et al. [12] that employs feature selection as balancing technique for ensemble of Naive Bayes (NB) Multinomial classifiers. For identical comparison with CNN that had duplicate source code comments in the training data, we retained the same data as in [6] containing duplicate source code comments unlike [12] in which the duplicate source code comments were removed during the data preprocessing. Table VI lists the F1 comparison scores with both "Duplicates" and "No Duplicates" scenarios.

Our classical machine learning classifier XGBoost with SMOTE achieved an overall average of 0.828 for ten projects exceeding the overall detection performance by 10.10% against cost-sensitive CNN [6]. XGBoost with SMOTE sampling improved the within-project F1 scores of all the projects except ArgoUML, Hibernate and JMeter. It appears that the deep learning model is more accurate for projects that have higher percentage of SATD instances. For the other projects, the XGBoost with SMOTE sampling has better F1 scores indicating improved performance of machine learning algorithm over deep learning in Within-Project setup. While for

TABLE VI
F1 COMPARISON

| Project | WITHIN-PROJECT | | | | CROSS-PROJECT | | | |
|----------------|-----------------------------|------------------|---------------|------------------|-----------------------------|------------------|---------------|------------------|
| | Duplicates [6] [*] | | No Duplicates | | Duplicates [6] [*] | | No Duplicates | |
| | CNN [6] | Our ² | NB [12] | Our ² | CNN [6] | Our ³ | NB [12] | Our ³ |
| Ant | 0.445 | 0.690 | - | 0.476 | 0.660 | 0.627 | - | 0.554 |
| ArgoUML | 0.932 | 0.907 | 0.705 | 0.846 | 0.878 | 0.861 | 0.828 | 0.891 |
| Columba | 0.741 | 0.936 | 0.732 | 0.815 | 0.852 | 0.823 | 0.801 | 0.865 |
| EMF | 0.532 | 0.778 | - | 0.696 | 0.679 | 0.483 | - | 0.522 |
| Hibernate | 0.887 | 0.851 | 0.752 | 0.937 | 0.826 | 0.825 | 0.788 | 0.822 |
| JEdit | 0.622 | 0.623 | 0.619 | 0.457 | 0.599 | 0.423 | 0.518 | 0.433 |
| JFreeChart | 0.795 | 0.919 | 0.581 | 0.815 | 0.739 | 0.574 | 0.687 | 0.724 |
| JMeter | 0.867 | 0.857 | 0.751 | 0.882 | 0.828 | 0.790 | 0.781 | 0.847 |
| JRuby | 0.881 | 0.885 | 0.782 | 0.921 | 0.863 | 0.908 | 0.841 | 0.891 |
| Squirrel | 0.813 | 0.838 | 0.628 | 0.681 | 0.739 | 0.692 | 0.651 | 0.737 |
| Avg. (10 proj) | 0.752 | 0.828 | - | 0.753 | 0.766 | 0.701 | - | 0.726 |
| Avg. (8 proj) | - | - | 0.693 | 0.794 | - | - | 0.736 | 0.780 |
| % Improv. | - | 10.10 | - | 14.57 | - | -9.27 | - | 5.45 |

* Same data set up as in [6] (with duplicate source code comments)

² consistent top performer in Within-Project setup: XGBoost with SMOTE

³ consistent top performer in Cross-Project setup: XGBoost with BLINE

cross-project SATD detection the deep learning CNN model outperformed our consistent top performer in Cross-Project setup XGBoost with BLINE model by 9.27% characterizing the improved feature generalization of source code comments across projects by deep learning model. This hints that for Within-Project SATD detection, classical machine learning with data balancing is better approach than deep learning.

Our consistent top performer outperforms the prior work’s [12] ensemble of Naive Bayes (NB) classifiers with feature selection balancing technique in both Within and Cross-Project set up by 14.57% and 5.45% respectively. All the improvements in the F1 scores are statistically significant based on Wilcoxon-Signed rank test.

F. Impact of Duplicates

Here we highlight the impact of duplicate source code comments in the training data. From the Table VI, it appears that the inclusion of duplicates can either improve or reduce our models’ performance. We can see that the inclusion of duplicates increases the F1 score from 0.753 to 0.828 in the Within-Project setup. At the same time, the inclusion of duplicates decreases the Cross-Project average F1 to 0.726 to 0.701. It is reasonable to think that within a project duplicates can be useful as similar commenting style or copy-paste commenting code can result in exact duplicates. On the other hand for cross-project learning, it appears that learning from duplicates of some other projects reduces performance as such comments are unlikely to be present in the target project.

More detailed project level comparison shows that all the F1 scores under “No Duplicates” have dropped except for the projects Hibernate, JMeter, and JRuby that have a higher number of duplicate SATD comments, with an average of 10% of duplicate SATD comments. The other remaining projects have less than 6% of duplicate SATD comments. This somewhat surprising finding suggests that fewer duplicates below a particular threshold increase the detection performance. While

on the other hand, removal of higher duplicates increases the prediction performance as observed in Hibernate, JMeter and JRuby projects. This indicate higher duplicates have a negative impact on the prediction score. Perhaps the feature learning over the same SATD features of duplicate source code comments in the training data does not enable the machine learning classifier to detect new SATD features when subjected to new source code comments in the test data.

In Cross-Project setup, each target project is trained with the data of the 9 other projects, and thus has a higher average of 778 duplicate SATD source code comments. This in contrast with the Within-Project setup’s per project average of 87 duplicate SATD source code comments, is much higher. The removal of SATD duplicate comments has improved the overall F1 performance from 0.701 to 0.726 asserting our previous observation that higher duplicates have a negative impact on the prediction performance.

V. SATD FEATURE INTERPRETABILITY AND ANALYSIS

In this section, we interpret features from the trained machine learning classifiers and study the impact of the balancing techniques on SATD feature selection. There have been significant recent advances in explaining black-box ML models [40]–[42]. We employ the explainability framework SHAP (SHapley Additive exPlanations) [42], due to its feature interpretation capability. SHAP is based on cooperative game theory, where the goal is to predict which groups (coalitions) will emerge from pool of player to collect payoffs. In other words, we use SHAP to see features as players that form coalitions to detect SATD. SHAP uses Shapley value [43], and the Shapley value of a feature represents the average contribution of a feature with all possible coalitions. With Shapley values, we identified the number of features that contributed to the prediction, features that have a high impact (features that have higher mean of absolute Shapley values) on the prediction, and the features that do not contribute to the prediction.

We choose the Logistic Regression classifier for JMeter project in Within-Project setup for a brief insight on the impact of balancing techniques for SATD detection. We chose JMeter as an example project as it has 7.79% of SATD instances, close to the overall SATD average of 9.08%. We can see in Table VII that in the JMeter example the performance is in line with previous results. BASELINE has the highest precision but suffers from poor recall while balancing techniques sacrifice little precision to improve recall and therefore offer better overall performance in F1 and ROC-AUC.

TABLE VII
LOGISTIC REGRESSION’S SATD DETECTION PERFORMANCE FOR JMETER (WITHIN-PROJECT)

| Technique | Precision | Recall | F1 | ROC-AUC |
|-----------------|--------------|--------------|--------------|--------------|
| BASELINE | 0.875 | 0.389 | 0.538 | 0.692 |
| COST | 0.784 | 0.806 | 0.795 | 0.892 |
| SMOTE | 0.848 | 0.778 | 0.812 | 0.882 |
| ADASYN | 0.829 | 0.806 | 0.817 | 0.895 |
| BLINE | 0.844 | 0.750 | 0.794 | 0.868 |
| SVMSMT | 0.871 | 0.750 | 0.806 | 0.870 |

Next, we look if balancing techniques differ in terms of features. Table VIII shows the number of features that are contributing or not contributing in each balancing technique, e.g., BASELINE has 875 contributing and 2116 non-contributing features. The table shows also how many new features each balancing technique introduces in comparison to BASELINE. COST echoes its inherent nature (of adjusting instance weight) and introduces 1.4% of new features, while the sampling technique introduces new features between 12-15%, reflecting its new sample synthesis ability.

TABLE VIII
SHAP FEATURE CONTRIBUTION STATISTICS (WITHIN-PROJECT JMETER)

| Technique | BASE-LINE | COST | SMOTE | ADASYN | SVMSMT | BLINE |
|--|-----------|------|-------|--------|--------|-------|
| Contributing SHAP value \neq 0 | 875 | 870 | 916 | 905 | 914 | 919 |
| Non-Contributing SHAP value = 0 | 2116 | 2121 | 2075 | 2086 | 2077 | 2072 |
| New Feature count | - | 12 | 124 | 108 | 112 | 140 |
| New Feature % | - | 1.4% | 13.5% | 11.9% | 12.2% | 15.2% |

Feature count alone cannot tell us if the balancing techniques really are meaningfully different from BASELINE so we also look in to the top contributing features of each. Table IX contains the top 10 features of SATD detection by each of the balancing technique.

TABLE IX
TOP 10 CONTRIBUTING FEATURES (JMETER WITHIN-PROJECT)

| BASE-LINE | COST | SMOTE | ADASYN | SVMSMT | BLINE |
|------------|------------|-----------|---------|--------------|------------|
| todo | todo | todo | todo | todo | todo |
| hack | hack | hack | hack | perhaps | perhaps |
| file | later | appear | bug | appear | hack |
| nonjavadoc | file | used | used | doe | used |
| fix | helper | doe | need | cleaning | could |
| later | fix | nonnls | number | fix | exception |
| helper | one | exception | appear | used | improve |
| yet | yet | improve | improve | disconnected | wrapped |
| encoding | nonjavadoc | perhaps | allow | nonnls | always |
| found | currently | later | reason | note | nonjavadoc |

The feature **'todo'** is the most contributing SATD feature across all the techniques. **'hack'** is the second most contributing feature across sampling techniques but does not list in the top 10 of SVMSMT. Similarity to BASELINE is the highest in COST with 8 shared features. Data sampling techniques have lower number of shared features with the BASELINE. SMOTE has 3 (todo, hack, later), ADASYN has 2 (todo, hack), SVMSMT (todo, fix), BLINE (todo, hack, nonjavadoc). So COST is the most similar to BASELINE also in terms of the top 10 contributing features while data sampling techniques have larger differences to BASELINE. When investigating feature importance, in a single project one would suspect that project specific words would have more weight. Indeed this is a difficulty others have experienced when using regression

coefficients as measures of feature importance [44]. It appears that SHAP avoids this problem as only one project specific word **'nonnls'** is seen in the table.

Finally, we explore the newly introduced features in comparison to the BASELINE to understand whether the new features make sense in SATD detection. Table X contains the top 3 features which are new contributing features, i.e., shifted from non-contributing to contributing features by the application of balancing techniques. Feature **'decision'** appears in two

TABLE X
TOP 3 INTRODUCED FEATURES (JMETER WITHIN-PROJECT)

| COST | SMOTE | ADASYN | SVMSMT | BLINE |
|----------------|-----------|-----------|----------|---------|
| validation | decision | although | decision | stopped |
| initialization | important | generally | sensible | nothing |
| fetchd | anything | notused | extra | extra |

techniques and it could very well be used when discussing decisions made while coding. Also feature **'extra'** is in two techniques and it could be used referring to code having extra logic. COST technique appears to have words that are the most technical (**'validation'**, **'initialization'**, **'fetchd'**), while data sampling techniques introduce new features that appear more appropriate in SATD detection like **'important'**, **'anything'**, **'nothing'**, **'sensible'**, **'notused'**, **'extra'**, **'decision'**.

VI. DISCUSSION

In this section, we discuss the factors affecting the choice of a balancing technique and the implications to practitioners and researchers. Our results show that no single balancing technique can provide consistent higher performance across multiple metrics. Even though the sampling technique SMOTE has consistently improved Recall, ROC-AUC, and F1 scores of BASELINE across multiple classifiers, the absolute average difference between COST and SMOTE for the same metrics is less than 4%. This suggests performance tradeoff as the crucial parameter for choosing between COST and SMOTE in which the SMOTE incurs an additional data synthesis step while the former does not.

Ensemble classifiers Random Forest and XGBoost handle class imbalance with their algorithmic design. Random Forest creates a random subset of samples to generate multiple decision trees while XGBoost creates more decision trees based on the misclassified samples in the training data. XGBoost's approach results in higher F1, ROC-AUC, and Recall while Random Forest results in higher Precision. However, the maintenance of Random Forest is more expensive with model training time being 2.5 times than that of the XGBoost.

The real-world application of AI/NLP techniques for SATD detection through source code comments depends heavily on time and resource availability. Limited time and resources could not allow a higher number of incorrectly classified SATD instances, warranting the precise detection of SATD comments. On the contrary, if enough time and resources are available to validate SATD comment classification manually, then extensive detection (recall) of SATD comments should

be the ideal choice of metric. Our study would enable practitioners and researchers to choose an appropriate balancing technique for NLP-based SATD detection in Within-Project (limited data) and Cross-Project (ample data) setup.

Another important implication for researchers is the removal of duplicate source code comments in the training data. Even though the duplicate source code comments accompanied with duplicate source code are to be treated as separate technical debt candidates, such duplicate comments should be avoided while training machine learning models as they affect the prediction scores.

VII. THREATS TO VALIDITY

The superior performance of classical machine learning algorithm XGBoost in Within-Project set up over Deep Learning algorithm Convolution Neural Network might be due to the limited training data size with an average of 3,442 source code comments. Even though the source code comments within a project will always be lesser than consolidated source code comments from multiple projects, the Within-Project setup's training data size might play a vital role in the classical machine learning algorithm's improved performance. The data sampling rate of all the evaluated sampling techniques have been evaluated at 1.0 (maximum), in our study, for realizing an even class distribution between technical debt and non-technical debt classes. The SATD detection performance of each evaluation metric is highly susceptible to change if the sampling rate changes.

The balancing technique recommendations from this study are based on the source code comments data with an average class imbalance ratio of 9.08% between SATD and non-SATD classes. All the balancing techniques employed in this study have been evaluated for the class imbalance ratio of 1:10 approximately (for every one technical debt source code comment, there are ten non-technical debt source code comments). A higher or lower class imbalance ratio before applying a balancing technique might have a different impact on the evaluation metrics. More research should identify the threshold for the ideal class imbalance ratio for improved SATD detection performance.

Another crucial factor that might impact the generalizability of our results is the choice of machine learning models. The applicability of our recommendations might vary for different machine learning classifiers.

VIII. CONCLUSION

In this study, we investigated the effectiveness of multiple balancing techniques to detect SATD comments from imbalanced data. Our study is based on a manually labeled dataset from 10 open source projects. We make four contributions.

First, we studied the performance of balancing techniques. We found that SMOTE provides the most consistent improvement for Recall, ROC-AUC and F1 in Within and Cross-Project SATD detection for all the classifiers included in this study. However, if one does not want the extra data synthesis step in SMOTE sampling, the classifier level COST

(class weight) balancing technique is a potential alternate. The ensemble (classifier level balancing technique) algorithm (Random Forest and XGBoost) achieved the best Cross-Project Precision and F1 scores in imbalanced data without explicit balancing technique such as COST or SMOTE. This shows that, for Cross-Project SATD detection in imbalanced data focusing on high Precision or high F1, the Bagging or Boosting ensemble techniques serves as a potential choice apart from sampling technique.

Second, we compare our results to previous works in SATD detection that include feature selection based Naive Bayes Multinomial [12] and COST based CNN [6]. Our top-performing XGBoost with SMOTE/BLINE sampling technique outperforms Huang et al.'s [12] approach by 14.57% and 5.45% in both Within and Cross-Project setup respectively. Then, in comparison with the state of the art CNN, we find that our consistent top performer in terms of F1 score (XGBoost with SMOTE) beats CNN in Within-Project SATD detection, offering a 10% improvement. However, in Cross-Project SATD detection, our top performer loses to CNN by 9%. For SATD detection, it suggests that in a Within-Project setup, classical ML is a better choice. However, the deep learning approach shows improved generalization capability in cross-project setup.

Third, our feature analysis shows that all techniques use a similar set of features with at least 85% of features being the same in all techniques. Sampling techniques have a larger difference to BASELINE than to COST technique. Highly shared top features were the words "todo" and "hack". Sampling techniques also appear to be able to introduce new feature words that are sensible in the SATD context which further support their consistent improved performance over COST.

Fourth, we developed a SATD classification tool based on this study that works in two modes, online and batch. Our web tool's batch mode would help reduce the initial efforts for labeling a huge volume of source code comments data.

In the future, we intend to apply and evaluate deep learning architectures for SATD detection from source code comments. In particular we are keen to explore the capability of BERT (a transfer learning scheme, that is already pre-trained on a huge corpus) for SATD detection in imbalanced data.

ACKNOWLEDGEMENT

The authors acknowledge the financial support by the Academy of Finland (grant ID 328058) and computational infrastructure by CSC Finland.

REFERENCES

- [1] M. Kuutila, M. Mäntylä, U. Farooq, and M. Claes, "Time pressure in software engineering: A systematic review," *Information and Software Technology*, vol. 121, p. 106257, 2020.
- [2] A. Potdar and E. Shihab, "An exploratory study on self-admitted technical debt," in *2014 IEEE International Conference on Software Maintenance and Evolution*. IEEE, 2014, pp. 91–100.
- [3] R. Longadge and S. Dongre, "Class imbalance problem in data mining review," *arXiv preprint arXiv:1305.1707*, 2013.
- [4] F. Thabtah, S. Hammoud, F. Kamalov, and A. Gonsalves, "Data imbalance in classification: Experimental evaluation," *Information Sciences*, vol. 513, pp. 429–441, 2020.

- [5] F. Pecorelli, F. Palomba, D. Di Nucci, and A. De Lucia, "Comparing heuristic and machine learning approaches for metric-based code smell detection," in *2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)*. IEEE, 2019, pp. 93–104.
- [6] X. Ren, Z. Xing, X. Xia, D. Lo, X. Wang, and J. Grundy, "Neural network-based detection of self-admitted technical debt: from performance to explainability," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 28, no. 3, pp. 1–45, 2019.
- [7] W. Cunningham, "The wycash portfolio management system," *ACM SIGPLAN OOPS Messenger*, vol. 4, no. 2, pp. 29–30, 1992.
- [8] N. S. Alves, L. F. Ribeiro, V. Caires, T. S. Mendes, and R. O. Spínola, "Towards an ontology of terms on technical debt," in *2014 Sixth International Workshop on Managing Technical Debt*. IEEE, 2014, pp. 1–7.
- [9] P. Aygeriou, P. Kruchten, I. Ozkaya, and C. Seaman, "Managing technical debt in software engineering (dagstuhl seminar 16162)," in *Dagstuhl Reports*, vol. 6, no. 4. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2016.
- [10] E. da Silva Maldonado, E. Shihab, and N. Tsantalis, "Using natural language processing to automatically detect self-admitted technical debt," *IEEE Transactions on Software Engineering*, vol. 43, no. 11, pp. 1044–1062, 2017.
- [11] F. Zampetti, C. Noiseux, G. Antoniol, F. Khomh, and M. Di Penta, "Recommending when design technical debt should be self-admitted," in *2017 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 2017, pp. 216–226.
- [12] Q. Huang, E. Shihab, X. Xia, D. Lo, and S. Li, "Identifying self-admitted technical debt in open source projects using text mining," *Empirical Software Engineering*, vol. 23, no. 1, pp. 418–451, 2018.
- [13] F. Pecorelli, D. Di Nucci, C. De Roover, and A. De Lucia, "On the role of data balancing for machine learning-based code smell detection," in *Proceedings of the 3rd ACM SIGSOFT international workshop on machine learning techniques for software quality evaluation*, 2019, pp. 19–24.
- [14] D. Cruz, A. Santana, and E. Figueiredo, "Detecting bad smells with machine learning algorithms: an empirical study," in *Proceedings of the 3rd International Conference on Technical Debt*, 2020, pp. 31–40.
- [15] A. Genkin, D. D. Lewis, and D. Madigan, "Large-scale bayesian logistic regression for text categorization," *Technometrics*, vol. 49, no. 3, pp. 291–304, 2007.
- [16] T. Hall, S. Beecham, D. Bowes, D. Gray, and S. Counsell, "A systematic literature review on fault prediction performance in software engineering," *IEEE Transactions on Software Engineering*, vol. 38, no. 6, pp. 1276–1304, 2011.
- [17] M. V. Mäntylä, F. Calefato, and M. Claes, "Natural language or not (nlon) a package for software engineering text analysis pipeline," in *Proceedings of the 15th International Conference on Mining Software Repositories*, 2018, pp. 387–391.
- [18] A. Liaw, M. Wiener *et al.*, "Classification and regression by randomforest," *R news*, vol. 2, no. 3, pp. 18–22, 2002.
- [19] L. Breiman, "Bagging predictors," *Machine learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [20] H. Mustapha, N. Abdelwahed *et al.*, "Investigating the use of random forest in software effort estimation," *Procedia computer science*, vol. 148, pp. 343–352, 2019.
- [21] M. A. Hossen, M. S. Islam, N. A. T. Yusof, M. S. Rahman, F. Siddika, M. Rahman, S. Khatun, M. S. A. Karim, and S. H. Mahmud, "Hybrid sampling and random forest based machine learning approach for software defect prediction," in *InECCE2019*. Springer, 2020, pp. 541–553.
- [22] F. A. Fontana, M. V. Mäntylä, M. Zanoni, and A. Marino, "Comparing and experimenting machine learning techniques for code smell detection," *Empirical Software Engineering*, vol. 21, no. 3, pp. 1143–1191, 2016.
- [23] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, 2016, pp. 785–794.
- [24] Y. Freund and R. Shapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *J Comput Syst Sci*, vol. 55, pp. 119–139, 1997.
- [25] E. Fitkov-Norris and S. O. Folorunso, "Impact of sampling on neural network classification performance in the context of repeat movie viewing," in *International Conference on Engineering Applications of Neural Networks*. Springer, 2013, pp. 213–222.
- [26] P. Kaur and A. Gosain, "Comparing the behavior of oversampling and undersampling approach of class imbalance learning by combining class imbalance problem with noise," in *ICT Based Innovations*. Springer, 2018, pp. 23–30.
- [27] E. Biswas, K. Vijay-Shanker, and L. Pollock, "Exploring word embedding techniques to improve sentiment analysis of software engineering texts," in *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*. IEEE, 2019, pp. 68–78.
- [28] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [29] S. García, J. Luengo, and F. Herrera, "Tutorial on practical tips of the most influential data preprocessing algorithms in data mining," *Knowledge-Based Systems*, vol. 98, pp. 1–29, 2016.
- [30] H. Han, W.-Y. Wang, and B.-H. Mao, "Borderline-smote: a new over-sampling method in imbalanced data sets learning," in *International conference on intelligent computing*. Springer, 2005, pp. 878–887.
- [31] H. He, Y. Bai, E. A. Garcia, and S. Li, "Adasyn: Adaptive synthetic sampling approach for imbalanced learning," in *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*. IEEE, 2008, pp. 1322–1328.
- [32] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Proceedings of the fifth annual workshop on Computational learning theory*, 1992, pp. 144–152.
- [33] A. Ali, S. M. Shamsuddin, A. L. Ralescu *et al.*, "Classification with class imbalance problem: a review," *Int. J. Advance Soft Compu. Appl*, vol. 7, no. 3, pp. 176–204, 2015.
- [34] M. Galar, A. Fernandez, E. Barrenechea, H. Bustince, and F. Herrera, "A review on ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based approaches," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 42, no. 4, pp. 463–484, 2011.
- [35] T. M. Khoshgoftaar, A. Fazelpour, D. J. Dittman, and A. Napolitano, "Ensemble vs. data sampling: Which option is best suited to improve classification performance of imbalanced bioinformatics data?" in *2015 IEEE 27th International Conference on Tools with Artificial Intelligence (ICTAI)*. IEEE, 2015, pp. 705–712.
- [36] G. Haixiang, L. Yijing, J. Shang, G. Mingyun, H. Yuan Yue, and G. Bing, "Learning from class-imbalanced data: Review of methods and applications," *Expert Systems with Applications*, vol. 73, pp. 220–239, 2017.
- [37] E. Loper and S. Bird, "NLtk: the natural language toolkit," *arXiv preprint cs/0205028*, 2002.
- [38] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, "Scikit-learn: Machine learning in python," *the Journal of machine Learning research*, vol. 12, pp. 2825–2830, 2011.
- [39] G. Lemaître, F. Nogueira, and C. K. Aridas, "Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning," *Journal of Machine Learning Research*, vol. 18, no. 17, pp. 1–5, 2017. [Online]. Available: <http://jmlr.org/papers/v18/16-365.html>
- [40] E. Strumbelj and I. Kononenko, "An efficient explanation of individual classifications using game theory," *The Journal of Machine Learning Research*, vol. 11, pp. 1–18, 2010.
- [41] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should i trust you?" explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.
- [42] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 4765–4774. [Online]. Available: <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- [43] L. S. Shapley, "A value for n-person games," *Contributions to the Theory of Games*, vol. 2, no. 28, pp. 307–317, 1953.
- [44] L. Rantala, M. Mäntylä, and D. Lo, "Prevalence, contents and automatic detection of kl-satd," in *2020 46th Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*. IEEE, 2020, pp. 385–388.