



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 10 Issue: V Month of publication: May 2022

DOI: <https://doi.org/10.22214/ijraset.2022.42562>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Bug Detection and Report: A Case Study on Dataset for Software Management using Security Bug Report

A Vinothini¹, Hariharan R²

¹Assistant Professor Department of Computer Science and Engineering, Veltech Multitech Dr. Rangarajan Dr. Sakunthala Engineering College, Avadi, India

²Veltech Multitech Dr. Rangarajan Dr. Sakunthala Engineering College, Avadi, Chennai.

Abstract: *In order to build a prediction model we must label a large amount of data in order to mine software repositories. The accuracy of the labels has a significant impact on a model's performance. However, there have been few research that have looked into the influence on a prediction model, there are occurrences that have been mislabeled. To close the gap, we conduct a research project on how to report a security bug (SBR) prediction in this paper. Furthermore, it has the potential to mislead SBR prediction research. We first enhance the label validity of these 5 datasets by personally evaluating each and every bug complaint in this study, and we discover 749 SBRs that were previously Non-SBRs have been mislabeled (NSBRs). We then examine the performance of the classification models both on messy (before the alteration) and cleaner (after our reconfiguration) datasets, impact of dataset label correctness. The results suggest that cleaning the datasets improves the performance of classification models.*

Index Terms: *Prediction of security bug reports, data quality, software detection and report*

I. INTRODUCTION

Mining software repositories (MSR) has evolved into a popular research topic for locating interesting and useful data about software systems and projects. MSR is based on big data in software engineering, and developing a prediction model requires a significant amount of labelled data. For appropriately analysing the success of a prediction model, the label correctness of data is crucial. Furthermore, poorly categorised data can lead to inaccurate results, which can lead to a misunderstanding of the target problem's study direction.

Currently, there are few research that focus on the influence of inaccurate labels, and while building a predictive model, attention to label correctness is necessary. Because finding SBRs culled from such a humongous bug database is crucial for lowering a software product's security concerns, this article studies the effects of mislabeled instances on SBR prediction. In recent years, a slew of deep learning-based SBR data resolutions have surfaced. However, the results of the two recent studies proposed by Peters et al. and Shu et al. (e.g., F1-score) are not optimal. When it comes to F1-score, for example, Peters et al. only managed 0.37 in their best scenarios. Shu et al. Straight-forward applied the data-sets garbled by Peters et al. and applied a hyperparameter optimization strategy to increase the performance of prediction models. Their experiment findings indicated that in the best situation, their technique could achieve pd (recall) and pf (false positive rate) of 0.86 and 0.25, respectively, which implies that the number of false positive items in Chromium achieves 5,388 as that the NSBRs with in Chromium testing set are 20,855. It is still undesirable to put the instrument into operation with such a high proportion of false alarms.

In this study, we investigate the factors that contribute to their research' poor performance, and we discover that one factor that cannot be overlooked is the standard of labels provided to the datasets. Peters et al. and Shu et al. employed 5 publicly available datasets: Chromium, Ambari, Camel, Derby, and Wicket. The first dataset, which comprises 40,940 bug complaints, was made available at MSR 2011, while the following 4 datasets were manually tagged by Ohira et al. In these 5 datasets, SBRs have indeed been misrepresented as NSBRs, according to Peters et al. Memory spillage and void pointer flaws have been reported as bugs, for example, are labelled as NSBRs in Chromium. They do, however, fall within the category of traditional vulnerability types, since they are regularly attacked by hackers and are among the top twenty five most hazardous CWE (Common Weakness Enumeration) kinds. 3 examples of mislabeled cases in Chromium are shown in Figure 1. The text indicating that the bug-report in question seems to be an SBR is marked in red.

Furthermore, Peters et al. and Shu et al techniques'swere labour-intensive and also sophisticated. Shu et al hyperparameter's tuning for the learner, for example, took roughly 5 hours to optimise the learner RF parameter on dataset Chromium. It is a recommended practise, according to Fu et al. [23] and Liu et al. [24], to study basic yet successful approaches. As a result, we look into the efficacy of basic text classification in SBR prediction.

The following 2 research questions are the focus of our research:

1) *RQ1*: To what degree does the accuracy of dataset labels affect the performance of classification models?

Both the cleaner and noisier data-sets, we perform an experiment to assess the efficacy of categorization models. We initially use manual annotation to fix the five datasets' mislabeled entries, resulting in 5 clean data-sets. As a consequence, we discovered 749 SBRs that were incorrectly categorised.

We separated each of the 5 clean datasets into two equal portions, which we called the both a training and a testing set based on Peters et al evaluation's [11]. Then, using the clean and noisy datasets, we use the techniques given by Peters et al. and Shu et al. to determine their performance values. Finally, we compare baseline techniques on clean datasets with those on noisy datasets to assess the influence of label accuracy. On cleaner datasets, the performance of three baseline techniques consistently beats that on noisy ones. Recall, Accuracy, F1-score, and G-measures are all up 191 percent, 46 percent, 147 percent, and 136 percent, respectively, on average.

2) *RQ2*: How does basic text categorization perform for SBR prediction on clean and noisy datasets?

We elevate RQ2 because the Description provided using text nature language is the major content of problem complaints for SBR prediction [25], [26].

Furthermore, because text classification is simple and straightforward, much earlier work on bug report processing has employed text classification in conjunction with machine learning. [13], [14], [15], [16], [17], [27].

On both the cleaner and noisy data-sets, we execute all five classification techniques (Random-Forest, Naive-Bayes, K-Nearest-Neighbour, Multi-layer Perception, and Logistical-Regression) applied by Peters et al. and Shu et al. using basic text categorization. The values for overall performance (i.e., Recall, Accuracy, F1-score, and G-measures) on clean datasets are substantially greater than in datasets with noise, according to our findings. In addition, with clean data-sets, the efficiency numbers are higher. (i.e., recall, precision, F1-score, and G-measure) of Random Forest employing text categorization outperform the baseline techniques based on security keywords matrices.

Finally, we look at the reasons behind the Farsec's enactment flaws, the hyperparameter's time cost tweaking methods, and the consequences of our research. We show that data label correctness has a significant impact on classification model performance. When data label accuracy is enhanced, simple text segmentation is much more beneficial for SBR forecasting than properly designed matrix-based baseline approaches.

The followings are the endowment made by this paper:

- a) For five publicly available datasets, we manually annotate the label validity (i.e., Chromium, Ambari, Camel, Derby, and Wicket). We discover none NSBRs that have been mislabeled as SBRs, but we do locate 749 SBRs that have been mislabeled as NSBRs.
- b) I believe that I'm the first one to experimentally examine the impact of data label correctness on SBR prediction.
- c) We discover that (1) performance on cleaner datasets is considerably better than performance on noisy datasets for the same classifier; (2) Text categorization that is straightforward that beats the 3 baseline techniques when training on clean datasets; and (3) simple-text- stratification surpasses the 3 base-guideline approaches when training on clean datasets. For SBR prediction researchers and practitioners, these data provide research hints and recommendations.
- d) The remainder of this thesis is organised as follows: The report's related stuff is described in Section 2. The background is presented in Section 3. Our data annotation strategy and outcomes are described in Section 4. Section 5 describes our experimental settings. Each research question's experimental outcomes are detailed in Section 6. Section 7 examines If Shu et al's tuning technique for basic text segmentation works, the patterns of improperly labelled SBRs in noisier data-sets, accuracy performance levels, the repercussions of our inquiry, and objections to the study's validity. The paper is summarised in Section 8.

II. RELATED WORK

A. Prediction of SBR

Peters et al. [11] and Shu et al. [12] have proposed 2 recent studies that focus on SBR prediction. Both of these research employ the same datasets and classification methods. Shu et al. build on Peters et al.'s work by using the SMOTE oversampling and hyperparameter tweaking technique to the learner (classifier).

Wijayasekara et al. [28] looked at bug detection reports from 2 open-source projects that are well-known, MySQL and the Linux kernel. They discovered that THIRTY TWO percent and SIXTY TWO percent of Linux kernel and MySQL vulnerabilities reported in bug tracking systems between 2006 and 2011 were concealed vulnerabilities that were not identified as SBRs until they were publicly revealed [29]. Their research also revealed an upward tendency in the prevalence of concealed vulnerabilities over the next 2 years. Furthermore, the researchers or the authors developed a way of discovering flaws that aren't readily apparent by generating features from bug reports' lengthy explanation (i.e., the field narration) and short elucidation (i.e., the field salutation) before using a classification strategy.

Michael et al. [13] proposed an earlier work in 2010 that focused on SBR prediction. The authors created an automated method for training a statistical model that is manually updated labelled bug detection reports that leverages descriptions of bug reports are mined for text.. This framework is then used to detect SBRs in the bug tracking system that have been manually mislabeled as NSBRs. On a Cisco software system, their technique is tested.

| | | | |
|------------|--|---|--|
| Title & ID | The root process has a memory leak. (Issue-1643) | On shutdown, the Alsa audio output leaks. (Issue 16036) | Check:CrashForBitmapAllocationFailure () (Issue 3795) |
| Definition | The root process acquired 1.6GiB of RAM after many days of use. And there was no way to restart or selectively destroy this process. The root process must not lose memory, and it must be easy to restart, taking up current page processes. Page navigation, for example, might be handled by a distinct child process. If it dies, it will respawn. A total of 1.6 GB has been set aside. | The shutdown function for Audio Output Stream does not appropriately destroy the output stream object once the thread is terminated. There is a resource leak as a result of this. Because the data source has been changed We should now remove the every-single-stream output task and write audio to avoid blocking. data using a single message loop that spans all output streams. | Other than checking for GDI leaks, I'm really unaware what we'd do. Despite the dearth of evidence thus far, this is a common crash in 154.6.. If we run out of GDI entities and if another problem arises, we won't receive a / bitmap here.. Since the data pointer is NULL, this will lead us to crash later. To ensure that individuals responsible for the crashes may be identified... |

Fig.1 In Chromium, below are some examples of bug reports that have been mislabeled. The contents of the red-colored words assist us in determining whether or not a bug-reports are definitely security-related (i.e., SBR).\

B. The Importance of Accurate Data Labels

Specialists of software development have highlighted concerns about the dataset utilised by machine learning-based categorization systems in recent years. [4], [8], [30], [31], [32], [33], [34].

Tantithamthavorn et al. [4] Investigate if the emergence of mislabeling is a result of chance. The authors looked at 3,931 person sending the message on Apache Jackrabbit and Lucene platforms and came to certain findings. Mislabeled of problem reports, for example, is not arbitrary; mislabeled issue reports have a minor influence on accuracy. Their research found that models are trained on noisy data obtain 56 percent to 68 percent of Using clean data, the recall of training images.

Kim et al. [8] conduct an experiment to see how noise affects inception fault analysis. They do their research on both a file and change level, inserting false negatives into their databases at random. Their findings demonstrate that when 25 percent to 35 percent of the datasets are mislabeled, the model performance in defect prediction suffers dramatically.

Mislabeled alterations' impact on the profitability and perception of only defect prediction models is investigated by Fan et al. [31].

The scientists looked at four SZZ [35] variations and then used the labelled data from these four variants to create prediction models. AG-mislabeled SZZ's modifications produce a considerable performance drop, according to their experimental evaluation of 126,526 alterations from 10 Apache-projects. The mis-labeled modifications of B-SZZ and AG-SZZ result in NINE percent to TEN and ONE percent to FIVE more lost inspection work, respectively, when the developers' inspection work is taken into account.

Tangled modifications affect the number of related defects for 16.6% of all source codes, according to Herzig et al. [33]. Because there are so many files affected by the number of flaws, the noise introduced by tangled modifications has a large influence on the models that forecast a lot of errors in origin source (root) codes.

Kochhar et al. [32] pointed out concern-reports (issues) mislabeling, i.e. complaints classified as bugs but referring to non-bug concerns. The authors use feature values extracted from bug reports to determine whether such a bug report should be classed. Bug-detection-reports from the projects HTTPClient, Jackrabbit, Lucene-Java, Rhino, and Tomcat5 were used to evaluate their methodology. The Precision, Recall, and F1-score performance of their technique spans from 0.58 to 0.71, 0.61 to 0.72, and 0.57 to 0.71, respectively.

Our research contradicts the findings of the previous investigations. First, we use 2 recent SBR prediction works put forwarded by Peters et al. and Shu et al. to enhance the label accuracy of five publicly accessible datasets. Second, we employ Peters et aland .s Shu et almethodologies .s as the study's baselines. Third, we compare baselines' productivity on noisy datasets to baselines' performance on clean datasets to assess the influence of data label accuracy. Using the three baseline methodologies, We conduct an experiment to determine the impact of data label correctness. Finally, we demonstrate that cataloguing is a simple yet effective approach with reliable data labels by using the commonly used text classification approach to both cleaner and noisier data sets.

III. BACKDROP ENVIRONMENT

In this part, we review 2 most recently proceeded SBR prediction studies that influenced our research.

Frame-works for Farsec. The 5 publicly accessible SBR-prediction data-sets, Chromium, Ambari, Camel, Derby, and Wicket, have a mislabeled issue, according to Peters et al. [11]. They create Farsec, a system for improving SBR analysis by screening out noisier data from NSBRs and integrating text mining with a matrix of security key words. There are three primary phases in the Farsec process: (1) Making data matri- ces and identifying security keywords. They start by tokenizing a dataset's SBRs to terms. The tf-idf value for each phrase is then computed, and the top 100 phrases with the largest tf-idf values are retained as security keywords. The prevalence of each security phrase in the Characterization of each bug report is then calculated for both the training and testing sets, For such training and testing sets, security-keywords matrices are created. (2) Filtering the NSBRs for terms linked to security. Filtration in Farsec is used to eliminate NSBRs that include security-related key phrases. They used seven distinct filters to do this: farsec, farsecsq, farsectwo, clni, clnifarsec, clnifarsecssq, and clnifarsectwo. A basic overview of these filters may be found in Table 1. Each one of these filters is applied to the training dataset of each data-set. As a result, seven additional training data-set models are created, each of which is used to fit into the model separately. (3) Bug-reports are ranked. Based on ensemble learning, a list of sorted bug reports was constructed after the SBRs were predicted. In the prediction results, the real SBRs show towards the top of the list.

TABLE 1

Peters et al. [11] employed seven filters.

| Filters | Descriptions |
|---------------|--|
| farsec | No support function should be used. |
| farsecsq | Apply Jalali et al[36] .s support functionalities to frequencies of terms present in SBRs. |
| farsectwo | Multiply the frequency by two in the Graham version [37]. |
| clni | Apply the CLNI (Closed List Noise Iden- tification) noise filter. |
| clnifarsec | Applying CLNI to farsec data which was filtered. |
| clnifarsecsq | ApplyingCLNIfilter to farsecs data which was filtered. |
| clnifarsectwo | ApplyingCLNIfiltertofarsectwo data which was filtered. |

Tuning hyperparameters. Model criterion are used in machine learning to describe the qualities of training data. The process of finding the most ideal settings for the model's hyperparameters is known as hyperparameter optimization [38]. Shu et al. [12] use a hyperparameter optimization strategy to enhance the productivity of SBR prediction. They base their research on Peters et al. findings. The controlled variables of a learner (i.e., classification algorithm) and the data-pre-processing technique SMOTE [39] are optimised using a differential evolution algorithm. They enhance Random Forest, for example, by adjusting the number of trees (i.e., the framework-parameter n estimators); the possible values is ten, with a tuning range of 10 to 150. In SMOTE, they adjust the number of neighbours (i.e., the framework-parameter k), with a default values of 5 and a tuning range of 1 to 20. In their work, Table 2 lists hyperparameters for adjusting classifier Random Forest settings, as well as the oversampling technique SMOTE.

IV. METHODS AND TECHNOLOGY FOR SBR DATA-ANNOTATION

Data elucidation takes a long time to complete [24], [40], [41]. Because of the high demand for competent knowledge, appropriately labelling SBRs is never easy [42], [43], [44].

Chromium, Ambari, Camel, Derby, and Wicket are the 5 data-sets used by Peters et al. and Shu et al. We'll go over the annotation procedure for both noisy and clean datasets in this section.

A. Noisy Datasets Annotation

We'll go over how the SBRs throughout the five datasets we used for this project worked, are identified in this section. In addition, we investigate why the underlying datasets contain mislabeled cases.

- 1) *Chromium*: Chromium1 is a web browser that is both free and open-source.. The Chromium dataset is supplied by the MSR conference's 2011 mining challenge [18], which comprises 40,940 bug-reports. The reports mark the SBR labels at first, and many of these are linked to CVE (Common risks and Subjections) entries [45]. Records with the identifiers Issue 34495 and Issue 34498, for example, are linked to CVE-2010-0048 and CVE-2010-0052, respectively. As a result, the SBR labels in the data-set are trustworthy. Furthermore, as Peters et al. [11] pointed out, There are still a lot of SBRs that have already been misidentified as NSBRs in this dataset. Because many bug reports are submitted by end-terminals, whom were seldom software security specialists, this is common in open-source projects.
- 2) *Four Apache Datasets*: Ohira et al. [19] obtained the datasets for the 4 Apache projects Ambari2, Camel 3, Derby 4, and Wicket 5 - JIRA, This is a well-known approach for trailing problems [46]. In JIRA, there may be several sorts of issues, including Bug, Improvements, Documentation, and Task. For each project, Ohira et al. chose one thousand problems with the categories Bug or Improvement at random. These 4,000 problems were personally evaluated and categorised by graduates and academics. They used the following procedures to identify six high-impact categories of problems (blocking, security, scalability, and breakage flaws):
 - a) For each project, one pupil and one research professor separately label the bug reports.
 - b) They then address their differences in order to find a consensus.

Because there were no set criteria for privacy, performance, and breakage categories, certain concerns were appraised differently, as Ohira et al. [19] pointed out. This might explain why certain SBRs in such data-sets have been mislabeled as NSBRs. Another factor that may contribute to the mislabeling is annotators' lack of security software awareness. If a null-pointer reference causes software to produce an exception, A bug reporting with insufficient security knowledge, for example, would submit the issue as an NSBR. The null-pointer, on the other hand, might be infected with malware, and the null-pointer de-reference is among the CWE top 25, this is a collection of the most common and significant flaws in software that can lead to major vulnerabilities [22].

B. Clean Datasets Annotation

From the five datasets, we will look for SBRs that have been mislabeled as NSBRs (i.e., Chromium, Ambari, Camel, Derby, and Wicket). An SBR model is said to be a defect-reports that describes 1 or more software system susceptibilities [13]. A security flaw is a software flaw that a hacker can exploit to obtain ingress to the structure or web-network [47]. The emergence of a vulnerability might have potentially dangerous effects, such as data leaking and unauthorised privilege escalation.

We don't only recruit highly qualified software-security professionals, but we also create a unique physical analysis procedure to ensure quality of our hand annotation results. As a starting point, we employ CWE [48]'s software vulnerability classes to create a codebook that may be used to determine if the bug-report is a SBR.

1) Observers

We do have a total of 6 annotators. 2 of them seem to be Ph.D. Graduates, while another 4 are Huawei workers. They have had at least 4 years of expertise in software securities and were acquainted with the initiatives that resulted in the problem reports.

Table 3 summarises the backgrounds of the 6 observers across 4 acreages: (1) their organisational roles (column Role); (2) their expertise with the practical software quality-checking (Testing) (column Software-development); (3) expertise with the initiatives or analogous projects (column Projects); and (4) encounter with security-bug-detection-reports and persual (source code exposed assessment, persual of security bug reports) (i.e., column Security-bug). We utilise year (amount of experience) to objectively quantify in order to be objective, the annotators' experience within every acreage. In addition, the table's last pole (Note) contains important annotator supplemental information.

TABLE2

For Random-Forest and SMOTE techniques, the hyperparameter tweaking is set. (Differential evolution is abbreviated as DE.)

| Targets | Parameters of the given Targets | | | Parameters of the DE | | | |
|-----------------|---------------------------------|----------------|--------------|----------------------|-----|-----|------|
| | Parameters | Default values | Tuning Range | NP | F | CR | ITER |
| Random-Forest | <i>nestimators</i> | 10 | [10,150] | 60 | 0.8 | 0.9 | 3,10 |
| | <i>minsamplesleaf</i> | 1 | [1,20] | | | | |
| | <i>minsamplessplit</i> | 2 | [2,20] | | | | |
| | <i>maxleafnodes</i> | None | [2,50] | | | | |
| | <i>maxfeatures</i> | auto | [0.01,1] | | | | |
| | <i>maxdepth</i> | None | [1,10] | | | | |
| SMOTE Technique | <i>k</i> | 5 | [1,20] | 30 | 0.8 | 0.9 | 10 |
| | <i>m</i> | 50% | [50,400] | | | | |
| | <i>r</i> | 2 | [1,6] | | | | |

2) Categories of the Software Vulnerabilities

Taxonomies of vulnerability categories have been developed by software security experts. Security flaws in software systems, according to Landwehr et al. [49], varying from implementation flaws on a local level (e.g., In C/C++, the gets() method call is used), to Errors in the interprocedural interface (for example, a race among a access-connect check as well as a file operation), to far higher level design blunders.

CWE is a society collection of common software flaws that, if left unchecked, might leave systems vulnerable to attack [48]. It includes flaws from major functional system vendors, commercial data security product suppliers, academia, government organisations, and research institutes, among others. The CWE principle states that , endangers that have a similar features are classified together into a single category, with 40 top-level CWE groups and 417 subtypes for software products [50]. CWE-1228, for example, is a top-level category that describes API/function problems. It has something to do with using third-party or built-in functionality APIs. CWE-242 (Use of intrinsically harmful function), CWE-477 (Use of outdated functions), and CWE-479 (Use of obsolete function) are among the seven subcategories (A potentially hazardous approach or function has been revealed.) [51].

Annotators were required to tie each SBR to certain CWE categories in order to illustrate the features of SBRs. (or sub-categories) CWE requires a comprehensive and in-depth knowledge basis for the annotation process. However, the domain region and development languages involved in a particular software system are often stable, reducing the number of CWE associated with the research and therefore minimising the effort of physically correlating CWE categories from SBRs.

3) Generation of Codebook

Before beginning the human review, we create a coding scheme to assist the annotation process, similar to Viviani et al. [52]. The 351 originally labelled SBRs from the 5 noisy datasets were used to create this codebook. Each one of the six authors individually evaluates the SBRs and generates a codebook.

There are two parts to the codebook:

- a) Why is the bug report an SBR? The Annotators must write out why bug-reports is indeed an SBR based on their experience and skills. They're indicated for describing the potential security dangers or implications of SBR after it happens.
- b) Empirical words or phrases: rather than duplicating the whole bug report summary, annotators pick and record the words and expressions that best demonstrate that bug-report is a SBR.

The six code books were then integrated and duplicates were deleted by the same annotators. Figure 2 shows an example of a codebook excerpt.

4) *Card Sorting and Manual Annotation*

We undertake the manual review using the card categorizing [53], [54] technique, in which each data-set is examined and classified separately by two corporate workers and one Ph.D. student based on CWE vulnerability analysis categories and the codebook. Annotators A1, A4, and A6 operate on the 4 tiny datasets (Ambari, Camel, Derby, and Wicket), whereas A2, A3, and A5 made activities on Chromium, according to their experience. As a consequence, each bug report is represented by three cards, for a sum of 137,820 cards across all five datasets. To quantify the agreement among the 3 annotators of each dataset, we use the Fleiss's Kappa Coefficient values [52], [55].

TABLE3

The 6 annotators' backgrounds. (Note: 4 of them are Huawei software engineers, as well as the other 2 are Ph.D. graduates specialising in software security- analysis.)

| Annotators | Roles | Experience(Years) | | | Notes |
|------------|------------------|--------------------|----------------|---------------|---|
| | | Software developed | Projects count | Security bugs | |
| A1 | Senior-Developer | 16 | 4 | 4 | Working on Hadoop-related products as a product developer Apache's products include Ambari, Camel, Derby, and Wicket. |
| A2 | Senior-Developer | 9 | 5 | 2 | is experienced with web crawler development; is a Project chromium volunteer |
| A3 | Test-Manager | 12 | 6 | 7 | For the past 5 years, I've performed as a security tester, analysing bug reports and recording and analysing CWE and CVE data. |
| A4 | Security-Tester | 5 | 4 | 5 | carries out security testing, such as fuzzing and penetration testing. |
| A5 | Ph.D.student | 4 | 3 | 5 | For the past three years, I've worked as an application security tester; my current research focuses on safety&security flaw report persual. |
| A6 | Ph.D-student | 1 | 3 | 5 | specialises in vulnerability mining and has won over \$30,000 in cash wins from vulnerability scanning competitions in the previous three years.. |

| Reasonsofbeingan SBR | Evidenceorphrases | Bug reportID |
|--|---|-----------------|
| disclosure of sensitive data | • • The services does not follow all of OWASP's recommendations for making it much harder for an attacker to crack encrypted passwords. It does not use an arbitrary salt and only performs one hash operation. | Derby-5539 |
| | • • A issue point with the display pass-code; password leaks | Chromium-1758 |
| | • • You might want to include the superior pass-code. | Chromium-1785 |
| | •• • | •• • |
| may result in a system crash or a DOS attack | • • In Google Chrome crawler, there is a memory corruption issue; investigate the underlying cause since this appears to be exploitable. | Chromium-11308 |
| | • • When dragging a file to a new tab, memory corruption occurs. | Chromium-12027 |
| | • • None storage harm is caused by the incorrect cache. | Chromium-27509 |
| | •• • | •• • |
| •• • | •• • | •• • |

Fig. 2 A. A codebook list snippet after getting merged

C. The Data Annotation results

The five datasets have a total of 118 conflicting label outcomes based on our hand annotation results. The mean agreement between many annotators is calculated using the Fleiss kappa statistic. The Fleiss kappa agreement level for the 5 datasets is significant or Near clear (the values of Fleiss kappa vary from 0.73 to 0.87 across the 5 datasets) 6, indicating excellent annotator agreement [52].

We utilise the labels of authors as final outcomes for bug-detection-reports that have received 3 harmonous labels. The annotators, on the other hand, plan a face-to-face review conference for bug reports with conflicting label findings to debate until a consistent outcome is obtained for each item. As a consequence, our hand annotation has certified 749 NSBRs of noise data-sets as SBRs. There are around 616, 27, 42, 91, and 37 newly detected SBRs for Chromium, Ambari, Camel, Derby, and Wicket, respectively, in each dataset. It's worth mentioning how no SBRs in noise data-sets have indeed been verified as NSBRs, implying that SBRs in noise datasets are quite reliable. One explanation might be that experienced users were more likely to report SBRs in bug detection and report methods.

6. Relationship between agreement level and Fleiss kappa value (Kp). Poor: Kp 0; Slight: 0.01 Kp 0.20; Fair: 0.20 Kp 0.40; Moderate: 0.40 Kp 0.60; Significant: 0.60 Kp 0.80; Near-perfect: 0.80 Kp 1.

Security testers who were specialised in system security or were particularly worried about it.

The vast majority of SBRs uncovered are obvious security issues that can be attributed to a single CWE. Figure III shows four SBRs that are associated with well-known susceptible categories (inappropriate permission, null pointer, wrong memory management, privacy leaks, etc.). The Camel-286, for example, is an SBR that is connected to null points. It is indeed component of CWE-465 (Pointer Errors), When an application removes a pointer that it assumes is legitimate but is NULL, it crashes or exits. So because Because void pointer decrementing is only 1 step back from crashing a programme (worst-case scenario) or giving an attacker unrestricted access to the system [58], identifying the null pointer issue is critical to system security. [56], [57].

In compared to the original datasets, our labelled datasets are clearer. As a result, we refer to our labelled data-sets as clean data-sets, whereas the original data-sets are referred to as noisy datasets.

The dispersion of the noisy data-sets, and also our clean datasets, is shown in Table 4. In the 5 datasets for such noisy data-sets, there's only a tiny number of bug-reports identified as SBRs. SBRs make up 0.45 percent to 8.80 percent of the population. For clean datasets, the proportion of SBRs ranges from 1.93 percent to 17.90 percent. We detected 616 SBRs in the huge dataset Chromium, which were categorised into NSBRs in the noisier sample, raising the number of SBRs by – 320% percent.

TABLE4

The five datasets were distributed both in noisy and clean versions.

| Datasets | #BR | Percentage analysis | | Percentage analysis | |
|----------|--------|---------------------|-------------|---------------------|----------------|
| | | #Noisy(%) | #Clean(%) | #Noisy(%) | #Clean(%) |
| Chromium | 41,940 | 193(0.47%) | 809(1.94%) | 41,749(99.55%) | 41,133(98.08%) |
| Ambari | 1,000 | 30 (3.00%) | 57 (5.70%) | 972(97.12%) | 945(94.50%) |
| Camel | 1,000 | 33 (3.30%) | 75 (7.50%) | 969(96.90%) | 927(92.70%) |
| Derby | 1,000 | 8 (8.90%) | 180(18.00%) | 913(91.30%) | 822(82.20%) |
| Wicket | 1,000 | 11 (1.10%) | 48 (4.80%) | 991(99.10%) | 954(95.40%) |

| ID | Title&Description of the Snippets given |
|-------------------|---|
| Camel 286 | ,When there is a connection between both the router and the service, it is called an endpoint. Endpoints of the CXF. NullPointerException occurs in CXF routes. Whenever an endpoint is introduced between such a cxf router and another cxf router... |
| Ambari- 3315 | On a big cluster (100 nodes), the no. of Out of Memory occurs when the number of ExecutionCommandEntity instances grows. Use this script to recreate the issue... |
| Chromium 8388 | Concerns about sites' capacity to trace the user's history outweigh concerns about privacy leaks from utilising non-incognito history databases to... |
| Chromium 29824 | Under Linux, the heap stack and others have permissions to read and write executables... Memory corruption is far easier to exploit with these permissions, especially when using a rwX HEAP. To be defeated. To perform a heap, an attacker might just use javascript. |

Fig.3. The following are some examples of evident SBRs found in our data analysis. The contents of the red-colored words assist us in assessing if a bug report represents a safety concern (i.e., SBR).

V. SETTING UP FOR EXPERIMENTAL USE

We introduce the basic methodologies proposed by prior studies in this section (i.e., Peters et al. and Shu et al.). Following that, we'll show you how we used a simple-text classification setup to explore RQ2. Finally, the classifiers and performance assessment metrics used in our investigation are explained.

A. Baseline Strategies

As a starting point for determining the impact of data label correctness, we adopt the methodologies suggested by Peters et al. and Shu et al., which we've described in Section 3. By adopting the nomenclature of their work, we refer to Peters et al.'s work as Farsec. We acquired two baseline methods from Shu et al. [12]'s work, as they modify the learner and SMOTE hyperparameters, respectively, in their investigation. Because they fully relate the trained and check set matrix data supplied by Farsec throughout his study, FarsecTuned is the name given to the process for tuning the learner's characteristics, as well as the strategy for adjusting the features of SMOTE.

The process of creating these 3 baselines is referred to as follows:

- 1) *Farsec*: A substructure for minimising the prevalence of security-related keywords in bug reports by filtering and rating them. Farsec builds security before aligning the prediction models, several filters are used to filter out non-security bug reports containing safety-related terms. (shown in Table 1). In our research, we used the farsectwo filter since it accounts for 80% of the data. In Peters et al.'s investigation, farsectwo produced the best outcomes across all five datasets. Farsec is a method of rating. It creates a keywords matrix for each dataset with above X security-related keywords. We ran trials (the inner-project experiments of Peters et al.'s study) with X = 50, 100, and 200 to find the best value for X. The average F1-score was the same (0.15) for all three settings, however the G-measure acquired with option 100 was indeed the best (0.29). As a result, we used the top 100 in our analysis, which matches the setup of Peters et al.'s [11]. Peters et al.'s [11] was conducted in this environment.
- 2) *FarsecLearner*: Using the multi - objective evolutionary approach to tune crucial factors of the learners using the computed passcode lock matrix information as inputs (i.e., the train and test set) [59].
- 3) *FarsecSmote*: Modifying the key dimensions of SMOTE with dynamic optimization, utilising the computed security phrases matrix records as inputs (i.e., the training dataset and testing set). To be impartial, we take Shu et al.'s sourcecode directly and leave all of the parameter settings (i.e. For each key parameter, there are crucial key components, attribute values, and tuning ranges.) exactly because they are engrossed in their task. The three baselines employ the same categorization approaches in their research (Random-Forest, Naïve-Bayes, K-Nearest-Neighbour, Multi - layer-perception, and Logistical-Regression) as well as data-sets (Chromium, Ambar, Camel, Derby, and Wicket).

B. Text Categorization Made Easy

The 100 best Matrix of security keywords is used to build models of classification in the three baseline techniques. This is not the same as standard text categorization. We also utilise simple text categorization as another strategy for our experimental assessment in this study since it is common and used by most bug report analyses [13], [14], [16], [17], [27]. To pre-process the content of the bug report using Count-Vectorizer and Select- FromMode from the machine-learning package for text standardisation and dimension reduction scikit-learn [60]. Description.

The CountVectorizer programme turns the content of a bug-report description into a token count matrix. Stop - word like a, the, and are eliminated to prevent being used as a signaling for prediction since they are assumed to be unspecific in expressing the substance of a document.

The token counts are then represented in a sparse representation matrix. A feature selection method is SelectFromModel. It weighs the relevance of characteristics and picks them based on their weights. It is a morpho that may be used in conjunction with other meta-transformers. any estimator with a coef or feature 's significance property, which specifies the feature weights. If the responsive coef or features importance values were less than the set threshold-value, the aspects are taken insignificant and eliminated. There are built-in algorithms for finding a barrier using a text parameter, in addition to defining the threshold numerically. Float, Mean, Median duplicates of these, such as 0.1*mean value, are available heuristics. For the CountVectorizer and SelectFromModel arguments, we utilise the default values.

C. Performance Metrics and Classifiers

We apply the performance indicators and classifiers used by Peters et al. and Shu et al. in this work to make the assessment findings objective.

Classifiers. We employ the five classifiers used by Peters et al. and Shu et al. in this investigation to be more purposeful. Random-Forest (RF), Naïve-Bayes (NB), K-Nearest-Neighbor (KNN), Multilayer-Perceptron (MLP), and Logistics Support Regression are all examples of machine learning algorithms (LR). These algorithms are also commonly used in the mining of software data repositories [61], [62]. To answer RQ1, however, we solely employ classifier RF in combination with baseline techniques to make a comparison clearer and easier. We chose RF [63] as the major classifier because (1) It is among the most popular clustering methodologies and needs to perform well for software development text classification; (2) this is one of to play a dominant role in baseline Farsec test research results; (3) it was the most frequently used classification methods and performs well enough for software development text classification [52], [54].

Metrics for measuring performance. We use all of the performance evaluations to eliminate bias and provide a comprehensive review. Retention (i.e., pd in their job), pf (probability of false-alarm), Accuracy, F1-score, and G-measures are some of these measures. The first 4 are widely used standards in empirical software development [61], [64], whereas G-measures is a prominent statistic in Peters et al study 's [11]. Both the F1-score and the G-measures are harmonic means, with the G-measures taking into account the both majority and minority classes' recalls [11].

There are 3 possible outputs of the forecast result for each bug report:

- Classified as +VE(positive) (TP): an SBR is projected to be SBR;
- False -VE(negative) (FN): an SBR is anticipated to be NSBR;
- True -VE(negative) (TN): an NSBR is expected to be NSBR;

The performance measures Recall, pf, Accuracy, F1-score, and G-measures may be derived based on these results:

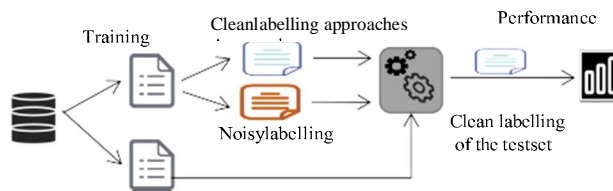


Fig. 4 Setting up data and the basic procedure for conducting studies.

$$\text{Recall} = \text{pd} = \text{TP} / \text{FN} + \text{TP} \quad (1)$$

$$\text{Pt} = \text{FPR} = \text{FP} / \text{TN} + \text{FP} \quad (2)$$

$$\text{Calculating Accuracy} = \text{TP} / \text{FP} + \text{TP} \quad (3)$$

$$\text{F1-measure} = 2 * \text{recall} * \text{precision} / \text{Precision} + \text{Recall} \quad (4)$$

$$\text{G-measures} = 2 * \text{recall} * (1 - \text{pf}) / (1 - \text{pf}) + \text{recall} \quad (5)$$

The higher the recall, accuracy, F1-measure, and G-measures are among these 5 performance indicators, while the lower the pf, the better.

D. Setting of Data

To guarantee that the comparison is fair, we divide the training and testing sets according to the baseline Farsec method. To put it another way, each dataset is categorized sequentially before being separated into 2 equal halves (i.e., 50 percent and 50 percent). The first section acts as a warm-up exercise., The second acts as a practise set., both of which are representative of the actual usage scenario in production. Shu et al two 's baselines (FarsecTuned and FarsecTuned) As they use similar matrix data processed, they use the same data-set partition technique by Peters et al.

In this experimental setup is depicted in Figure 4. On either noisy or clean data, we use several categorization algorithms (i.e., Farsec, FarsecTuned, FarsecTuned, Text, TextTuned, and TextTuned). However, we evaluate each technique using the test set's clean labels, along with the model trained using noisy data. We do it because the most important thing is if the prediction accuracy were right in stating of the reality, i.e., clean statistics.

VI. RESULTS OF THE EXPERIMENT

In this section, we give the outcomes of our experiment by responding to the two questions.

A. RQ1's Response

RQ1: To what degree does the accuracy of dataset labels affect classification performance algorithm?

To address RQ1, we are using the classifier RF to execute the 3 baseline proposition (i.e., Farsec, FarsecTuned, and FarsecTuned) on both the noisier and cleaner datasets.

The results obtained of the 3 baselines learned both with noisy and clean datasets are shown in Table 5.

TABLE 5

Both on Noisy and Clean datasets, performance outcomes of the 3 baseline techniques with classifier RF. All of the labels in the test dataset are from the Cleaner datasets. If the value of a Cleaner dataset is greater than that of a Noisy dataset, it is bold-faced (Note: For Recall, the greater the number, the better. , Accuracy, F1-score, and G-measures; the lower the number, the better for pf).

| Data names | Approaches | Recall | | pf | | Accuracy | | F1-score | | G-measures | |
|------------|--|--------|-------|-------|-------|----------|-------|----------|-------|------------|-------|
| | | Noisy | Clean | Noisy | Clean | Noisy | Clean | Noisy | Clean | Noisy | Clean |
| Chromium | Farsec | 0.01 | 0.66 | 0.00 | 0.00 | 0.83 | 0.95 | 0.02 | 0.78 | 0.02 | 0.79 |
| | <i>Farsec^{Tuned} Learn er</i> | 0.03 | 0.67 | 0.00 | 0.00 | 0.39 | 0.94 | 0.05 | 0.78 | 0.05 | 0.80 |
| | <i>Farsec^{Tuned} Smote</i> | 0.30 | 0.80 | 0.17 | 0.07 | 0.04 | 0.22 | 0.07 | 0.35 | 0.45 | 0.86 |
| Ambari | Farsec | 0.44 | 0.50 | 0.02 | 0.07 | 0.39 | 0.19 | 0.41 | 0.28 | 0.60 | 0.65 |
| | <i>Farsec^{Tuned} Learn er</i> | 0.38 | 0.56 | 0.02 | 0.06 | 0.35 | 0.23 | 0.36 | 0.33 | 0.54 | 0.70 |
| | <i>Farsec^{Tuned} Smote</i> | 0.44 | 0.44 | 0.05 | 0.06 | 0.22 | 0.19 | 0.29 | 0.26 | 0.60 | 0.60 |
| Camel | Farsec | 0.07 | 0.33 | 0.02 | 0.12 | 0.27 | 0.21 | 0.11 | 0.26 | 0.12 | 0.48 |
| | <i>Farsec^{Tuned} Learn er</i> | 0.11 | 0.24 | 0.04 | 0.05 | 0.21 | 0.33 | 0.14 | 0.28 | 0.20 | 0.38 |
| | <i>Farsec^{Tuned} Smote</i> | 0.07 | 0.33 | 0.06 | 0.13 | 0.11 | 0.20 | 0.08 | 0.25 | 0.12 | 0.47 |
| Derby | Farsec | 0.24 | 0.84 | 0.18 | 0.54 | 0.25 | 0.27 | 0.25 | 0.41 | 0.38 | 0.60 |
| | <i>Farsec^{Tuned} Learn er</i> | 0.38 | 0.69 | 0.22 | 0.44 | 0.30 | 0.27 | 0.33 | 0.39 | 0.51 | 0.62 |
| | <i>Farsec^{Tuned} Smote</i> | 0.33 | 0.60 | 0.28 | 0.25 | 0.22 | 0.37 | 0.27 | 0.45 | 0.45 | 0.67 |
| Wicket | Farsec | 0.00 | 0.52 | 0.00 | 0.05 | 0.00 | 0.34 | 0.00 | 0.41 | 0.00 | 0.67 |
| | <i>Farsec^{Tuned} Learn er</i> | 0.00 | 0.52 | 0.00 | 0.08 | 0.00 | 0.24 | 0.00 | 0.33 | 0.00 | 0.67 |
| | <i>Farsec^{Tuned} Smote</i> | 0.00 | 0.39 | 0.00 | 0.05 | 0.00 | 0.26 | 0.00 | 0.32 | 0.00 | 0.55 |
| Average | Farsec | 0.15 | 0.57 | 0.04 | 0.16 | 0.35 | 0.40 | 0.16 | 0.43 | 0.22 | 0.64 |
| | <i>Farsec^{Tuned} Learn er</i> | 0.18 | 0.54 | 0.06 | 0.13 | 0.25 | 0.40 | 0.18 | 0.42 | 0.26 | 0.63 |
| | <i>Farsec^{Tuned} Smote</i> | 0.23 | 0.51 | 0.11 | 0.11 | 0.12 | 0.25 | 0.14 | 0.33 | 0.32 | 0.63 |
| | <i>All</i> | 0.19 | 0.54 | 0.07 | 0.13 | 0.24 | 0.35 | 0.16 | 0.39 | 0.27 | 0.63 |

The average Recall, Accuracy, F1-score, and G-measures scores across the 5 datasets after training with the noisy dataset are 0.15, 0.35, 0.16, and 0.22 in Farsec; 0.18, 0.25, 0.18, and 0.26 in FarsecTuned; and in FarsecTuned, 0.23, 0.12, 0.14, and 0.32. The mean prices of Recall, Accuracy, F1-score, and G-measures across the 5 datasets are 0.57, 0.40, 0.43, and 0.64 in Farsec; 0.54, 0.40, 0.42, and 0.63 in FarsecTuned; and 0.51, 0.25, 0.33, and 0.63 in FarsecTuned. The following findings may be obtained when evaluating performance post training with cleaner and noisier datasets:

- (1) On cleaner data-sets, the maximal of Recall, Accuracy, F1-score, and G-measures of the baseline methods are greater than on noisy datasets.
- (2) Recall, Accuracy, F1-score, and G-measures average values have grown by 191 percent, 46 percent, 147 percent, and 136 percent, respectively.

Finding1

On clean datasets, the same classifier performs considerably better than that on noisier datasets.

B. RQ2's Response

RQ2: How does the basic text classification perform for SBR prediction on cleaner datasets?

To answer RQ2, we use basic text classification to run the five classification algorithms (RF, NB, MLP, LR, and KNN). Each one of the 5 classification techniques is initially trained on a noisy dataset and its clean counterpart. The model is then used to anticipate the project's testing set. The precise data labels are then used to fabricate the model's performance measures.

The outcomes derived of the 3 baseline methods learned both with noisy and clean datasets are shown in Table 6. If the clean dataset's performance is better than the noisy dataset's, the clean dataset's performance is emphasised in boldface. It's worth noting that the clean edition of the five datasets has substantially higher total values of Recall, Accuracy, F1-score, and G-measures than the noisy version. Recall, Accuracy, F1-score, and G-measures have average values of 0.08, 0.29, 0.11, and 0.14 for models, respectively.

bred on noisier dataset, and 0.37, 0.54, 0.42, and 0.51 on cleaner data-sets, which are 362 percent, 86 percent, 281 percent, and 264 percent higher than noisy datasets. The gain on the huge dataset Chromium, in particular, is significant. The average F1-score for the 5 suggest that educators on the clean Chromium dataset is 654 percentage points greater than for the messy Chromium data-set.

Finding2

On clean datasets, simple text categorization performs better than on noisy ones.

In terms of Precision, our results for basic text categorization contradict Tantithamthavorn et al conclusion's that mislabeling has no influence on prediction model Precision [4]. We dug deeper into the causes and discovered that Tantithamthavorn To combat class imbalance, we re-balanced the data to train., which has an influence on prediction model precision [65], [66]. In fact, their conclusions are quite alike to the outcomes of our baseline techniques (Farsec, FarsecTuned), and FarsecTuned), and these baselines reduce class imbalance by removing noise from the dominant group or applying SMOTE.

Table 6

Both on Noisy and Clean datasets, the five classifiers performed well with basic text classification. The test dataset's labels were from the Cleaner datasets. If the values of Cleaner datasets is greater than the value of Noisier datasets, it is bold-faced (Note: For Recall, the greater the number, the better., Accuracy, F1-score, and G-measures; For pf, the lower the better.).

| Data names | Learner | Recall | | pf | | Accuracy | | F1-score | | G-measures | |
|------------|---------|--------|-------|-------|-------|----------|-------|----------|-------|------------|-------|
| | | Noisy | Clean | Noisy | Clean | Noisy | Clean | Noisy | Clean | Noisy | Clean |
| Chromium | RF | 0.01 | 0.72 | 0.00 | 0.00 | 0.99 | 0.92 | 0.02 | 0.81 | 0.02 | 0.84 |
| | NB | 0.15 | 0.66 | 0.02 | 0.04 | 0.17 | 0.29 | 0.16 | 0.40 | 0.26 | 0.79 |
| | MLP | 0.06 | 0.49 | 0.00 | 0.00 | 0.33 | 0.71 | 0.10 | 0.58 | 0.11 | 0.66 |
| | LR | 0.04 | 0.58 | 0.00 | 0.00 | 0.37 | 0.84 | 0.07 | 0.69 | 0.08 | 0.73 |

| | | | | | | | | | | | |
|---------|---------|------|------|------|------|------|------|------|------|------|------|
| | KN N | 0.04 | 0.58 | 0.00 | 0.00 | 0.37 | 0.84 | 0.07 | 0.69 | 0.08 | 0.73 |
| Ambari | RF | 0.13 | 0.25 | 0.02 | 0.02 | 0.20 | 0.31 | 0.15 | 0.28 | 0.22 | 0.40 |
| | NB | 0.13 | 0.44 | 0.02 | 0.02 | 0.15 | 0.39 | 0.14 | 0.41 | 0.22 | 0.60 |
| | ML P | 0.31 | 0.38 | 0.03 | 0.02 | 0.28 | 0.38 | 0.29 | 0.38 | 0.47 | 0.54 |
| | LR | 0.19 | 0.25 | 0.01 | 0.01 | 0.38 | 0.40 | 0.25 | 0.31 | 0.32 | 0.40 |
| | KN N | 0.19 | 0.25 | 0.01 | 0.01 | 0.38 | 0.40 | 0.25 | 0.31 | 0.32 | 0.40 |
| Camel | RF | 0.00 | 0.37 | 0.00 | 0.00 | 0.00 | 0.89 | 0.00 | 0.52 | 0.00 | 0.54 |
| | NB | 0.04 | 0.30 | 0.05 | 0.05 | 0.09 | 0.37 | 0.06 | 0.33 | 0.08 | 0.46 |
| | ML P | 0.00 | 0.20 | 0.02 | 0.03 | 0.00 | 0.43 | 0.00 | 0.27 | 0.00 | 0.33 |
| | LR | 0.02 | 0.15 | 0.00 | 0.02 | 0.33 | 0.39 | 0.04 | 0.22 | 0.04 | 0.26 |
| | KN N | 0.02 | 0.15 | 0.00 | 0.02 | 0.33 | 0.39 | 0.04 | 0.22 | 0.04 | 0.26 |
| Derby | RF | 0.10 | 0.46 | 0.02 | 0.01 | 0.59 | 0.90 | 0.18 | 0.61 | 0.19 | 0.63 |
| | NB | 0.16 | 0.53 | 0.06 | 0.13 | 0.39 | 0.50 | 0.23 | 0.51 | 0.28 | 0.66 |
| | ML P | 0.14 | 0.39 | 0.03 | 0.07 | 0.50 | 0.56 | 0.22 | 0.46 | 0.25 | 0.55 |
| | LR | 0.12 | 0.39 | 0.01 | 0.06 | 0.67 | 0.60 | 0.21 | 0.48 | 0.22 | 0.55 |
| | KN N | 0.12 | 0.39 | 0.01 | 0.06 | 0.67 | 0.60 | 0.21 | 0.48 | 0.22 | 0.55 |
| Wicket | RF | 0.00 | 0.48 | 0.00 | 0.00 | 0.00 | 0.92 | 0.00 | 0.63 | 0.00 | 0.65 |
| | NB | 0.00 | 0.30 | 0.00 | 0.03 | 0.00 | 0.30 | 0.00 | 0.30 | 0.00 | 0.46 |
| | ML P | 0.00 | 0.22 | 0.00 | 0.02 | 0.00 | 0.36 | 0.00 | 0.27 | 0.00 | 0.36 |
| | LR | 0.00 | 0.13 | 0.00 | 0.01 | 0.00 | 0.38 | 0.00 | 0.19 | 0.00 | 0.23 |
| | KN N | 0.00 | 0.13 | 0.00 | 0.01 | 0.00 | 0.38 | 0.00 | 0.19 | 0.00 | 0.23 |
| Average | RF | 0.05 | 0.46 | 0.01 | 0.01 | 0.36 | 0.79 | 0.07 | 0.57 | 0.09 | 0.61 |
| | NB | 0.10 | 0.45 | 0.03 | 0.05 | 0.16 | 0.37 | 0.12 | 0.39 | 0.17 | 0.59 |
| | ML | 0.10 | 0.33 | 0.02 | 0.03 | 0.22 | 0.49 | 0.12 | 0.39 | 0.17 | 0.49 |
| | PL | 0.07 | 0.30 | 0.01 | 0.02 | 0.35 | 0.52 | 0.11 | 0.38 | 0.13 | 0.44 |
| | RK | 0.07 | 0.30 | 0.01 | 0.02 | 0.35 | 0.52 | 0.11 | 0.38 | 0.13 | 0.44 |
| | NN | 0.08 | 0.37 | 0.01 | 0.03 | 0.29 | 0.54 | 0.11 | 0.42 | 0.14 | 0.51 |
| | All | | | | | | | | | | |

On the cleaner and noisier datasets, Figure 5 displays the Recall boxplots, Accuracy, F1-score, and G-measures for the techniques Text, Farsec, tuned learner, and Farsec tuned SMOTE. The cleaner datasets' scores are presented in red boxes, The scores of the noisier data-sets are displayed in blue color boxes. The Precision of the 3 standardising on pure and ambiguous information is comparable, as seen by the Precision charts., which is consistent with Tantithamthavorn et al result .s [4].

When measuring the results of basic text categorization to the three baselines trained on clean datasets (i.e., Farsec, FarsecTuned, and FarsecTuned), the Precision has increased dramatically, with an average increase of 125 percent. In terms of Recall, however, the 3 baselines perform better. The proportion of real SBRs recognised from the overall SBRs of the " proposed is measured by recall, whereas the percentage of authentic SBRs from the identified SBRs is measured by accuracy.. For SBR prediction, recall and accuracy were both significant criteria. To eliminate bias, we employ the F1-score as one of the most significant assessment statistic.

The F1-score is a mean of Accuracy and Recall that determines whether an improvement in Accuracy (Recall) balances a decrease in Recall (Precision) [67]. To eliminate bias in this research, we employ the F1-score as one of the most essential assessment statistic. On the cleaner dataset (red), the box of basic text categorization (txt) is way greater than the 3 benchmarks, as seen in the F1-score boxplots (fsc, ftl, and fts). In terms of F1-score, Simple text classification outperforms the other three baselines by 46 percentage points.

Finding3

Simple text classification outperforms the three baseline approaches in clean datasets.

VII.DISCOURSE

Here we discuss whether Shu et al hyperparameter's tuning method works when incorporating simple-text-categorization, ways of misidentifying in the actual data-sets, current achievement of various authors (Professionals vs Doctoral students), the repercussions of our research, In this part, we'll discuss the threats to the credibility of our research.

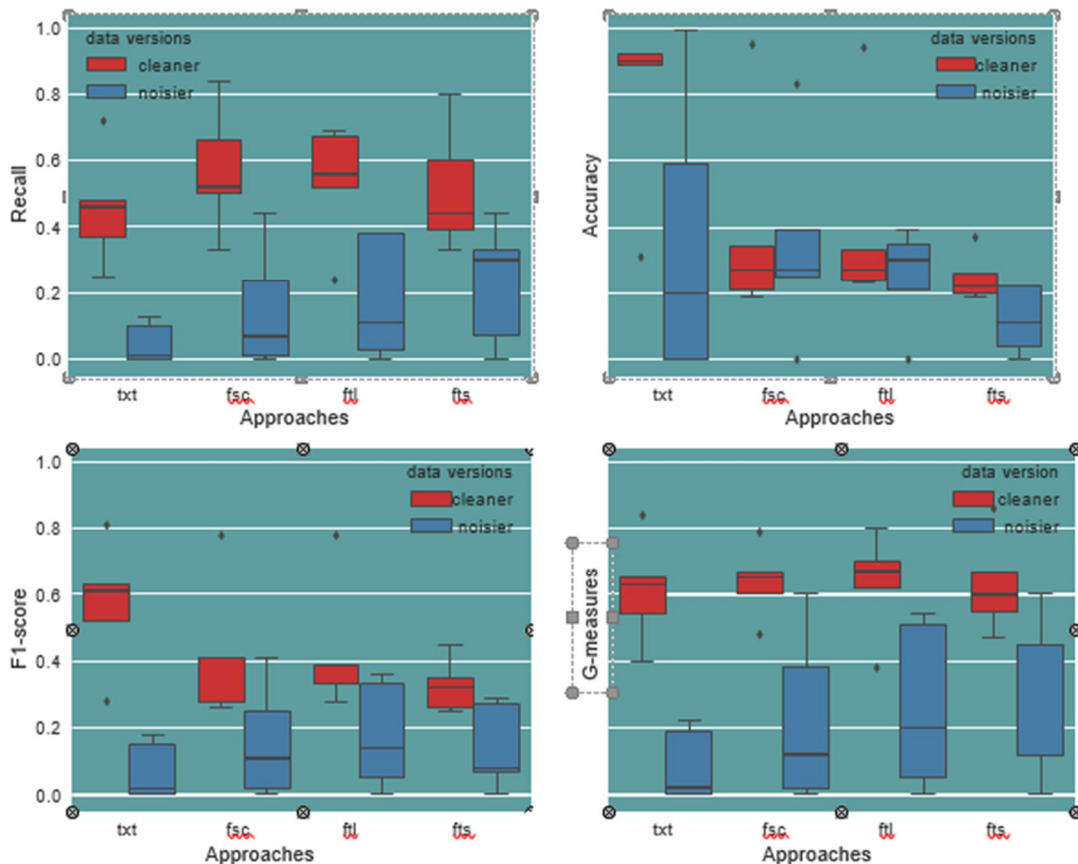


Figure 5: shows the execution of the cleaner (red) and noisy (blue) data-sets for basic text-classifications, as well as the 3 Baselines (For example, Farsec, FarsecTuned), as well as FarsecTuned). (Noting the letters txt, fsc, ftl, and fts stand for Text, Farsec, and FarsecTuned, respectively).

A. Do Shu et Al hyperparameter's Tweaking Techniques work for Simple Text Classification?

We test whether Shu et al hyperparameter's tuning approach works in conjunction with the ordinary text categorization approaches presented in RQ2. Here, we use RF to test Shu et al two hyperparameter tuning techniques. We use FarsecTuned and TextTuned to imply the mixture of learner improvements for text classification (tweaking the important guidelines of learner RF) and the mixture of text-categorization with Smote intonation (able to tune the important parameters of Smote), respectively, following the mentioning rule of our baseline approaches.

The results of integrating text-classification with hyperparameter tweaking are shown in Table 7. (i.e., tune and SMOTE respectively).

Across the five datasets, the best values of Recall, Accuracy, F1-score, and G-measures are 0.76, 0.94, 0.84, and 0.89, respectively; while the mean of Recall, Accuracy, F1-score, and G-measure is 0.56, 0.58, 0.52, and 0.68. FarsecTuned (combo of basic text classification and learner hyperparameter tweaking) increases Precision significantly, while FarsecTuned (combo of simple text categorization and SMOTE hyperparameter adjustment) improves Recall significantly. FarsecTuned beats FarsecTuned on average, given the relevance of the harmonic metric F1-score.

Figure 6 shows boxplots of the performance rating of the six techniques over the five clean datasets to compare the efficacy of FarsecTuned and FarsecTuned with the 4 main ways used to answer RQ1 and RQ2. Each group is represented by a paired colour, with the blue coloredpair representing 2 basic techniques (Farsec and Text), the greener coloredpair representing appeal with learner hyperparameter tuning (FarsecTuned and Tuned Learner), and the red pair representing appeal with SMOTE hyperparameter tuning (FarsecTuned and TextTuned).

TABLE7

The effectiveness of Shu et al hyperparameter .'s optimization methodologies paired with basic text categorization.

| Data-sets | Approaches | Recall | pf | Accuracy | F1-score | G-measures |
|-----------|----------------------|--------|------|----------|----------|------------|
| Chromium | TextTuned Learner | 0.76 | 0.00 | 0.94 | 0.84 | 0.87 |
| | TextTuned Smote | 0.86 | 0.08 | 0.17 | 0.28 | 0.89 |
| Ambari | TextTuned Learner | 0.25 | 0.02 | 0.31 | 0.28 | 0.40 |
| | TextTuned Smote | 0.50 | 0.02 | 0.50 | 0.50 | 0.66 |
| Camel | TextTuned Learner | 0.39 | 0.00 | 0.95 | 0.55 | 0.56 |
| | TextTuned Smote | 0.41 | 0.09 | 0.32 | 0.36 | 0.57 |
| Derby | TextTuned Learner | 0.52 | 0.01 | 0.93 | 0.66 | 0.68 |
| | TextTuned Smote | 0.68 | 0.25 | 0.40 | 0.50 | 0.71 |
| Wicket | TextTuned Learner | 0.52 | 0.01 | 0.71 | 0.60 | 0.68 |
| | TextTuned Smote | 0.70 | 0.03 | 0.57 | 0.63 | 0.81 |
| Average | TextTuned Learner | 0.49 | 0.01 | 0.77 | 0.59 | 0.64 |
| | TextTuned Smote | 0.63 | 0.09 | 0.39 | 0.45 | 0.73 |
| | TextTuned Smote | 0.55 | 0.04 | 0.59 | 0.53 | 0.69 |
| | All | | | | | |

In terms of Accuracy and F1-score, we can see that FarsecTuned (green) is significantly better than FarsecTuned (light green), which is consistent with evaluating Text (blue) with FarsecTuned (green) (light blue). For all four performance parameters, though, TextTuned (red) consistently surpasses FarsecTuned (light red) (Recall, Accuracy, F1-score, and G-measures). TextTuned outperforms Text in terms of the crucial performance metric F1-score, making it the strongest of the 6 techniques.

B. Mislabeling Patterns in the Source Datasets

Data reflection by hand is a time-consuming process. This section evaluates the manually identified SBRs in order to detect patterns of misrepresentation in order to give information on subsequent advanced data cleansing and data labelling efforts.

1) Pervasiveness

The manual assessment found a total of 749 SBRs from the misinterpretation of noisy datasets. Each record has one or more CWE tags applied to it. We were able to comprehend the distribution of mislabeled objects thanks to these tags.

The 749 SBRs also include than 50 CWE categories. However, we grouped these SBRs as from 40 top-level CWE classes of software evolution to derive the most common properties. Finally, the top 3 categories with the most SBRs are CWE-1218 (Cache memory error), CWE-199 (Information administration mistake), and CWE-465 (Pointer Issues). They account for over 75% of all SBRs. Chromium is responsible for roughly 90% of the entries in CWE-1218 and CWE-199, whereas the four tiny datasets are responsible for 95% of the records in CWE-465.

2) Description Patterns

With plain language, a bug description is a brief summary of the observed behaviour. (OB) and/or the expected behaviour (EB). The following are the most prevalent OB and EB patterns, according to Chaparro et al. [26]:

- a) *OB*: ([subject]) ([negative aux. verb]) ([subject]) ([subject]) ([subject]) ([sub [adverb] [adverb] [adverb] [adverb] [Here, [deleterious aux. verb] means "not," "can't," "doesn't," "didn't," "didn't," "didn't," "didn't," "didn't," "didn't," "didn't," "didn't," " For instance, the shutdown code for [Audio Output Stream] [It] [doesn't] [doesn't [After the thread has been shut down, appropriately destroy the output stream object...]) (Image courtesy of Chromium Issue 16036).
- b) *EB*: [topic] should/shall/must/must/must/must/must/must/must/must/must/must/must/must for instance, [It] shouldn't [mind asking back to the Main thread because it'll be pointless at current point.] (Adapted from Chromium 41547). When something comes to SBRs, there are a few things to consider, the [object] and [verb]/[complement] of the SBRs' description contain certain security-related keywords. For such top three CWE classes, Table 8 summarises several high-frequency phrases. (CWE-1218, CWE-199, and CWE-464) based on the 749 SBRs detected. SBRs Always include a mix of terms both from [subject] & [verb]/[complement] classes, or words with comparable roots, in each category.

C. Engineers vs. Ph.D. Students in Terms of Annotation Performance

6 annotators work on our clean datasets annotation, including four Huawei software engineers and 2 Ph.D. Graduates. This section discusses the quality standards of comments from business personnel and university student-personnels.

The accuracy rate of company annotators is superior than those of Ph.D. students on the basis about each annotating result and the clean datasets' final labels. Ph.D. students A5, who has 3 years of test automation expertise and is presently concentrating on bug report analysis, surpasses company accuracy A3, who has 8 years of software design expertise and is acquainted with the project Chromium. – i.e., authors who have done hands-on development and testing are more apt to provide correct labels.; However, knowledge in software security-related operations aids in correct labelling even more.

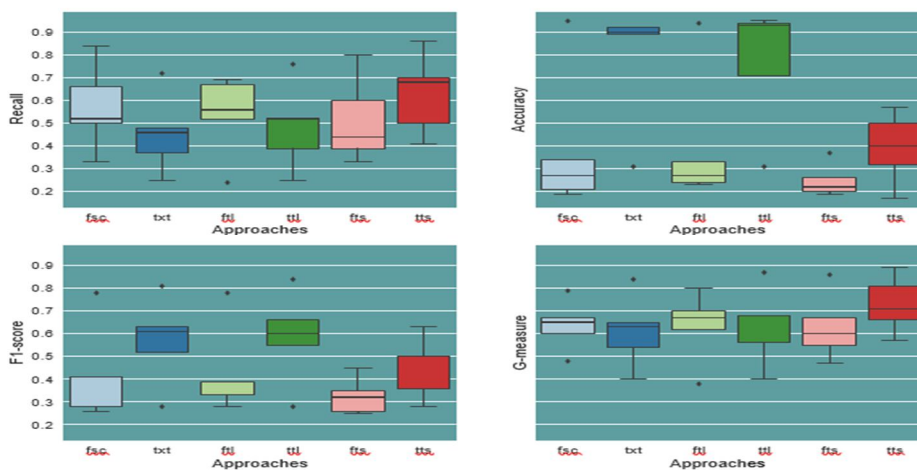


Figure 6 shows boxplots with On the cleaner datasets, paired colour for the effectiveness of the six included classification techniques. The blue pair represents two fundamental techniques (Farsec and Text), the green-colored pair represents learner hyperparameter tuning (Farsec Tuned Learner), and the red pair represents SMOTE hyperparameter tuning (Farsec Tuned SMOTE) (FarsecTuned and TextTuned). (Note that fsc, txt, ftl, ttl, fts, and tts Smote are abbreviations for Farsec, Text, FarsecTuned, FarsecTuned, and TextTuned, respectively.)

D. Indications

We have distilled several general insights from this work that go beyond the specific objective and methodology. The quality of your data is crucial. The most significant challenge for Microsoft's data scientists is data comprehensiveness, based on the findings of Kim et al. [68]. Data quality difficulties, according to researchers, make it tough for data analysts to have faith in the accuracy of their work. The potency of SBR estimation techniques is strongly correlated with data label correctness, as evidenced by our SBR prediction experiment results. A model fit with clean data outperforms a model fit with filthy information with skewed data by a significant margin. In the world of software engineering, comparable verification has indeed been done in the fields of static code defect prediction. [31], [68], legitimate bug prediction [32], and clone detection [69]. Kim et al. [68], for example, looked at the effect of data-labels on fault prediction. The prediction model's performance has definitely worsened, with mislabels exceeding 25%, according to their evaluation results. According to Tantithamthavorn et al. [4], predictive prototype are trained on noisier dataset attain 56 percent to 68 percent of the Remembering of trained models on clean data.

It will need a team effort. A high-quality dataset is required for appropriate experimental assessment of SBR prediction techniques. Because the existing automated technique cannot handle SBR prediction dataset labelling adequately, manual work is required. Despite the fact that manual data labelling is laborious and costly, it may be overcome with a collective effort. There are several examples of collaborative initiatives yielding high-quality datasets. Svajlenko et al. [69], [70], for example, By mining and manually inspecting clones of 10 common characteristics, we were able to construct a huge clone detection benchmark. Six million authentic affirmative mutations of various clone types are available. (Type-1, Type-2, Type-3, and Type-4) make up the benchmark. Over the course of 216 hours of mechanical attestation efforts, three judges discovered these clones. Researchers [71], [72], and [73] often utilise these datasets.

TABLE8
Mislabeling SBRs in the top three categories resulted in the following keywords.

| Categories | | |
|------------|--|--|
| | [Subjects] | [verbs]/[complements] |
| CWE-1218 | heaps, stacks, caches, buffers, pools, cpu, loops, length, ranges, indices, arrays, files, directory, data, exception, bound consumptions. | out, consume, uncontrolled, dereference, use-after-free, out-of-bounds |
| CWE-119 | user, user-name, sessions, pro-file, security, licences, hosts, certificates | leak, exposure, mask, stockpiling, transfers, logs, delicate, clean, hashing, authorize, allow, invalid, malicious are all terms that may be used to describe a situation. |
| CWE-465 | Pointer indication. | dereference, release, wrong, improper, null, outside, invalid, exception, uninitialized, initialise, out-of-range, expired, handle |

A collaborative effort is required to properly identify data and develop high-quality openly accessible data-sets in the realm of SBR prediction. One of the stages toward accomplishing that aim is the job we do. Our datasets and manually annotations files have been made openly available for others to check and extend.

E. Validity Challenges

- 1) *Internal Validity*: It is a term used to describe the validity of a Possible mislabels in the clean data-sets pose a danger to the study's internal validity. We take a number of steps to lessen the hazard. However, ensuring that the clean datasets are free of false positives and negatives remains a challenge. Since there is no specific definition for SBR, as Ohira et al. [19] pointed out, determining whether such a bug (error) reported is either an SBR or an NSBR is difficult. To determine if a bug report is an SBR, we apply the criteria of CWE, the most certified institution in the field of vulnerability management. Furthermore, each of the manual review's six annotators has extensive practical expertise and a thorough understanding of various vulnerability kinds and features. Furthermore, Viviani et al. [52] describe a solid data strategy; we use their method to construct a code-book by evaluating 351 known SBRs from the five studies. The arguments why is it that a bug-report is tagged as SBR, as well as the particular terms that justify this judgement, are recorded in this codebook. We utilise the codebook as a guide for the subsequent evaluation processes, The cards sorting process is also used to confirm that the label findings are accurate. The influence of optimizing for the outcomes is another danger to internal validity. In order to counteract this threat, we tested a variety of approaches, including selecting optimum value for Farsec (i.e., using the farsectwo and the top 100 access control keywords), tuning important parameters of classifier (i.e., FarsecTuned and TextTuned Learner), and tuning metrics of the over-sampling approach SMOTE. The results of the experiments reveal that various tuning procedures have minimal effect on the results. Other tuning techniques, such as facet pick and under investigation, may be studied in the later [74].
- 2) *Validity from the Outside*: The scope of our findings exposes us to external validity challenges. In this work, we compare the outcomes of similar SBR predictions on noisier data-sets to those of clean datasets to determine the impact of data label correctness. The three baseline approaches from Peters et al. and Shu et al. 's SBR prediction work, as well as basic text categorization models, are used. Furthermore, the evaluation is based on a series of performance metrics (most of the performance measures that they used in their study), including Recall, pf, Accuracy, F1-score, and G-measures. The generalisation of the findings is another external hazard. The influence of data integrity on SBR prediction is investigated in this study. Our findings are not universally applicable to all technology analytics activities. Data annotation, on the other hand, is used in many MSR activities. Source-codes defect predictions [4], [8], [31], high-impact-bug-reporting predictions [27], [32], [34], and software productivities and quality analysis [68] are only a few examples. One of the most important aspects impacting a model's efficacy for prediction model-related activities is data quality [64]. The label accuracy of training data has been shown to affect the performance of the classification models in most prior research [31], [33], [68], which really is consistent with our findings. We believe that the insights of this article can be used to MSR scenarios that have comparable features to SBR prediction. These features encompass, but aren't restricted to: (1) a large class imbalance in the scenario; (2) a text mining-based problem-solving technique. The most likely scenario is the forecast of a high-impact defect report. Prediction of performance bug reports [75] and config bug reports [76], for example, are significant jobs that might affect software quality.

VIII. OUTLINE

We enhanced the label accuracy of five publicly accessible SBR predictions data-sets and conducted an extensive experimental examination of the impact of data label consistency on classification methods in this study.

The findings demonstrate that (1) when using the same forecasting models (e.g., Peters et al.'s Shu et al. methods), 's performance on cleaner data-sets is considerably preferable than on noisier datasets. (2)

On clean datasets, uncomplicated text categorization performs substantially better than on noisy ones. (3) With clean datasets, basic text classification surpasses Peters et al. and Shu et al. baseline's techniques.

A. Appreciations

For their contributions to the annotation of our datasets, I'd want to express my gratitude to SensenGuo, Zhong Liu, Weiqiang Fu, Fengyu Liu, and Manqing Zhang. We'd like to express our gratitude to the researchers of finer Security-Bug-Report Classifications via Predictive Algorithms and Note Filter and Rankings for Security-Bug-Reports Predictions for making their data-sets and programmes available to the public with us.

With the help of Veltech University and the Division of Computer Science and Engineering, this study was completed.

TESTIMONIALS OR REFERENCES

- [1] "Mining software repository," 2020, <http://www.msrfconf.org/>.
- [2] M. Shepperd, Q. Song, Z. Sun, and C. Mair, "Data quality: Some comments on the nasa software defect datasets," *IEEE Transactions on Software Engineering*, vol. 39, no. 9, pp. 1208–1215, 2013.
- [3] Y. Tian, N. Ali, D. Lo, and A. E. Hassan, "On the unreliability of bug severity data," *Empirical Software Engineering*, vol. 21, no. 6, pp. 2298–2323, 2016.
- [4] C. Tantithamthavorn, S. McIntosh, A. E. Hassan, A. Ihara, and
- [5] K. Matsumoto, "The impact of mislabelling on the performance and interpretation of defect prediction models," in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, IEEE, 2015, pp. 812–823.
- [6] M. Jimenez, R. Rwemalika, M. Papadakis, F. Sarro, Y. Le Traon, and M. Harman, "The importance of accounting for real-world labelling when predicting software vulnerabilities," in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2019, pp. 695–705.
- [7] B. Ghotra, S. McIntosh, and A. E. Hassan, "Revisiting the impact of classification techniques on the performance of defect prediction models," in *In Proceedings of the 37th International Conference on Software Engineering - Volume 1*, IEEE, 2015, pp. 789–800.
- [8] J. Debattista, S. Auer, and C. Lange, "Luzzu—a methodology and framework for linked data quality assessment," *Journal of Data and Information Quality (JDIQ)*, vol. 8, no. 1, pp. 1–32, 2016.
- [9] S. Kim, H. Zhang, R. Wu, and L. Gong, "Dealing with noise in defect prediction," in *2011 33rd International Conference on Software Engineering (ICSE)*, IEEE, 2011, pp. 481–490.
- [10] O. Alonso, "Challenges with label quality for supervised learning," *Journal of Data and Information Quality (JDIQ)*, vol. 6, no. 1, pp. 1–3, 2015.
- [11] Z. Liu, X. Xia, C. Treude, D. Lo, and S. Li, "Automatic generation of pull request descriptions," in *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, IEEE, 2019, pp. 176–188.
- [12] F. Peters, T. Tun, Y. Yu, and B. Nuseibeh, "Text filtering and ranking for security bug report prediction," *IEEE Transactions on Software Engineering*, vol. PP, no. 99, pp. 1–1, 2017.
- [13] R. Shu, T. Xia, L. Williams, and T. Menzies, "Better Security Bug Report Classification via Hyperparameter Optimization," *in preprint arXiv:1905.06872*, 2019.
- [14] M. Gegick, P. Rotella, and T. Xie, "Identifying security bug reports via text mining: An industrial case study," in *Proceedings of the 7th International Working Conference on Mining Software Repositories, Cape Town, South Africa, 2–3 May 2010*, IEEE, 2010, pp. 11–20.
- [15] D. Behl, S. Handa, and A. Arora, "A bug mining tool to identify and analyze security bugs using naive bayes and tf-idf," in *Proceedings of International Conference on Optimization, Reliability, and Information Technology, Faridabad, India, 6–8 Feb. 2014*, IEEE, 2014.
- [16] S. Zaman, B. Adams, and A. E. Hassan, "Security versus performance bugs: a case study on firefox," in *Proceedings of the 8th working conference on mining software repositories*, 2011, pp. 93–102.
- [17] P. Kamongiang and K. Kavi, "VULCAN: Vulnerability Assessment Framework for Cloud Computing," in *Proceedings of the 7th International Conference on Software Security and Reliability, Gaithersburg, MD, USA, 18–20 June 2013*, IEEE, 2013.
- [18] Z. Thomas, N. Nachiappan, and W. Laurie A., "Searching for needles in a haystack: Predicting security vulnerabilities for windows vista," in *Proceedings of the 3rd International Conference on Software Testing, Verification and Validation, Paris, France, April 7–9, 2010*, IEEE, 2010, pp. 421–428.
- [19] "Mining Software Repository Challenge," 2011, <http://2011.msrfconf.org/msr-challenge.html>.
- [20] M. Ohira, Y. Kashiwa, Y. Yamatani, H. Yoshiyuki, Y. Maeda,
- [21] N. Limsettho, K. Fujino, H. Hata, A. Ihara, and K. Matsumoto, "A dataset of high impact bugs: Manually-classified issue reports," in *Mining Software Repositories*, 2015.
- [22] Q. Gao, Y. Xiong, Y. Mi, L. Zhang, W. Yang, Z. Zhou, B. Xie, and H. Mei, "Safe memory-leak fixing for c programs," in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, IEEE, 2015, pp. 459–470.
- [23] Y. Kim, M. Papamichael, O. Mutlu, and M. Harchol-Balter, "Thread cluster memory scheduling: Exploiting differences in memory access behavior," in *2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture*, IEEE, 2010, pp. 65–76.
- [24] "2019 TOP25 CWEs," 2019, <https://cwe.mitre.org/top25/>.
- [25] W. Fu and T. Menzies, "Easy over hard: A case study on deep learning," in *Proceedings of the 2017 11th joint meeting on foundations of software engineering*, 2017, pp. 49–60.
- [26] Z. Liu, X. Xia, A. E. Hassan, D. Lo, Z. Xing, and X. Wang, "Neural-machine-translation-based commit message generation: how far are we?" in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, 2018, pp. 373–384.
- [27] O. Chaparro, C. Bernal-Cárdenas, J. Lu, K. Moran, A. Marcus,
- [28] M. DiPenta, D. Poshvanyk, and V. Ng, "Assessing the quality of the steps to reproduce in bug reports," in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2019, pp. 86–96.
- [29] O. Chaparro, J. Lu, F. Zampetti, L. Moreno, M. DiPenta, A. Marcus, G. Bavota, and V. Ng, "Detecting missing information in bug descriptions," in *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*, 2017, pp. 396–407.
- [30] X. L. Yang, D. Lo, X. Xia, Q. Huang, and J.-L. Sun, "High-impact bug report identification with imbalanced learning strategies," *Journal of Computer Science and Technology*, vol. 32, no. 1, pp. 181–198, 2017.
- [31] W. J. Wijayasekara, D. Manic, M., "Mining Bug Databases for Unidentified Software Vulnerabilities," in *Proceedings of the 5th International Conference on Human System Interactions, Perth, WA, Australia, 6–8 June 2012*, 2012, pp. 89–96.
- [32] J. Arnold, A. Tim, D. Waseem, P. Gregory, E. Nelson, T. Geoffrey, and A. Kaseorg, "Security Impact Ratings Considered Harmful," in *Proceedings of the 12th Conference on Hot Topics in Operating Systems*, Usenix, 2009.
- [33] S. Kim, T. Zimmermann, K. Pan, E. James, J. Retal, "Automatic identification of bug-introducing changes," in *21st IEEE/ACM international conference on automated software engineering (ASE'06)*, IEEE, 2006, pp. 81–90.

- [34] Y.Fan,X.Xia,D.A. da Costa, D. Lo, A. E. Hassan, and S. Li,“The impact of changes mislabeled by szz on just-in-time defectprediction,”IEEETransactionsonSoftwareEngineering,2019.
- [35] P. S. Kochhar, F. Thung, and D. Lo, “Automatic fine-grained issuereportreclassification,”in201419thInternationalConferenceonEngineering of Complex Computer Systems.IEEE, 2014, pp. 126–135.
- [36] S.J.HerzigKimand A. Zeller, “It’s not a bug, it’s a feature:how misclassification impacts bug prediction,” in 2013 35th In-ternational Conference on Software Engineering (ICSE).IEEE, 2013,pp.392–401.
- [37] X. Xia, D. Lo, W. Qiu, X. Wang, and B. Zhou, “Automated configu-rationbugreportpredictionusingtextmining,”in2014IEEE38thAnnual Computer Software and Applications Conference.IEEE, 2014,pp.107–116.
- [38] J.Śliwerski,T.Zimmermann,andA.Zeller,“Whendochanges induce fixes?” ACM sigsoft software engineering notes, vol. 30, no. 4,pp.1–5,2005.
- [39] O.Jalali,T.Menzies,andM.Feather,“Optimizingrequirementsdecisionswithkeys,”inProceedingsofthe4thinternationalworkshoponPredictormodelsinsoftwar eengineering,2008,pp.79–86.
- [41] Graham and Paul, Hackers and painters : big ideas from the computerage.O’Reilly,2004.
- [42] M. Feurer and F. Hutter, “Hyperparameter optimization,” Auto-matedMachineLearning,pp.3–33,2019.
- [43] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer,“Smote: synthetic minority over-sampling technique,” Journal ofartificialintelligenceresearch,vol.16,pp.321–357,2002.
- [44] X. Zhou, X. Peng, T. Xie, J. Sun, C. Ji, W. Li, and D. Ding, “Faultanalysis and debugging of microservice systems: Industrial sur-vey, benchmark system, and empirical study,” IEEE TransactionsonSoftwareEngineering,2018.
- [45] Y.Gan, Y. Zhang, D. Cheng, A. Shetty, P. Rath, N. Katarki,
- [46] A.Bruno,J.Hu,B.Ritchken,B.Jacksonetal.,“Anopen-sourcebenchmarksuiteformicroservicesandtheirhardware- softwareimplicationsforcloud&edgesystems,”inProceedingsoftheTwenty-Fourth International Conference on Architectural Support forProgrammingLanguagesandOperatingSystems,2019,pp.3–18.
- [47] B. Liu, G. Meng, W. Zou, Q. Gong, F. Li, M. Lin, D. Sun, W. Huo,and C. Zhang, “A large-scale empirical study on vulnerabilitydistributionwithprojectsandthelessonslearned,”in202042ndInternationalConferenceonSoftwareEngineering(ICSE).IEEE,2011.
- [48] K.Ku,T.E.Hart,M.Chechik,andD.Lie,“Abufferoverflowbenchmarkforsoftwaremodelcheckers,”inProceedingsofthetwenty-second IEEE/ACM international conference on Automated soft-wareengineering,2007,pp.389–392.
- [49] I.Pashchenko,S.Dashevskyi,andF.Massacci,“Delta-bench:differential benchmark for static analysis security testing tools,”in2017 ACM/IEEE International Symposium on Empirical SoftwareEngineeringandMeasurement(ESEM).IEEE,2017,pp.163–168.
- [50] “CVE,”2020,<https://www.cvedetails.com/>.
- [51] “JIRA,”2018,<https://www.atlassian.com/software/jira>.
- [52] B. Potter and G. McGraw, “Software security testing,” IEEE Secu-rity&Privacy,vol.2,no.5,pp.81–85,2004.
- [53] “CWE,”2020,<https://cwe.mitre.org/>.
- [54] C.E.Landwehr,B.Alan R., M. John P., and C. William S.,“Ataxonomyofcomputerprogramsecurity,”ACMComputingSurveys,vol.26,no.3,pp.211–254,1994.
- [55] “SoftwareCWEs,”2020,<https://cwe.mitre.org/data/definitions/699.html>.in2016IEEE27Thinternationalsymposiumonsoftwarereliability.
- [56] “CWE1228,”2020, <https://cwe.mitre.org/data/definitions/1228.html/>.
- [57] G. Viviani, M. Famelis, X. Xia, C. Janik-Jones, and G. C. Murphy,“Locating latent design information in developer discussions: Astudy on pull requests,” IEEE Transactions on Software Engineering,2019.
- [58] D. Spencer, “Card sorting: Designing usable categories,” RosenfeldMedia,2009.
- [59] Z.Wan,X.Xia,A.E.Hassan,D.Lo,J.Yin,andX.Yang,“Per-ceptions, expectations, and challenges in defect prediction,” IEEETransactionsonSoftwareEngineering,2018.
- [60] J.L.Fleiss,“Measuringnominalscaleagreementamongmany raters.”PsychologicalBulletin,vol.76,no.5,pp.378–382,1971.
- [61] B. Jack, “Vector rewrite attack: Exploitable null pointer vulnerabili-ties on arm and xscale architectures,” White paper, Juniper Networks,2007.
- [62] T.Mandt,“Kernelattacksthroughuser-modecallback,”BlackHatUSA,Aug,2011.
- [63] A.Kogtenkov,“Voidsafety,”Ph.D.dissertation,ETHZurich,2017.
- [64] S.DasandP.N.Suganthan,“Differential evolution:Asurveyofthe state-of-the-art,” IEEE transactions on evolutionary computation,vol.15,no.1,pp.4–31,2010.
- [65] “Scikitlearn,”2020,<https://scikit-learn.org/stable/index.html>.
- [66] X.Ren,Z.Xing,X.Xia,D.Lo,X.Wang,andJ.Grundy,“Neuralnetwork-based detection of self-admitted technical debt: Fromperformance to explainability,” ACM Transactions on Software Engi-neeringandMethodology(TOSEM),vol.28,no.3,pp.1–45,2019.
- [67] C. Liu, D. Yang, X. Xia, M. Yan, and X. Zhang, “A two-phase trans-fer learning model for cross-project defect prediction,” InformationandSoftwareTechnology,vol.107,pp.125–136,2019.
- [68] L. Breiman, “Random forests,” Machine learning, vol. 45, no. 1, pp.5–32,2001.
- [69] Y. Yang, X. Xia, D. Lo, T. Bi, J. Grundy, and X. Yang, “Predictivemodelsinsoftwareengineering:Challengesandopportunities,”inpreprintarXiv:2008.03656,2020.
- [70] P. Hulot, D. Aloise, and S. D. Jena, “Towards station-level demandpredictionforeffectiverebalancinginbike-sharing systems,”in Proceedingsofthe24thACMSIGKDDInternationalConferenceonKnowledgeDiscovery&DataMining,2018,pp.378–386.
- [71] C.Tantithamthavorn,A.E.Hassan,andK.Matsumoto,“TheImpact of Class Rebalancing Techniques on the Performance andInterpretationof Defect Prediction Models,” IEEE Transactions onSoftwareEngineering,2018.
- [72] X.Xia,D.Lo, E.Shihab, X.Wang, and B. Zhou, “Automatic,high accuracy prediction of reopened bugs,” Automated SoftwareEngineering,vol.22,no.1,pp.75–109,2015.
- [73] M. Kim, T. Zimmermann, R. DeLine, and A. Begel, “Data scientistsin software teams: State of the art and challenges,” IEEE Transac-



- tiononSoftwareEngineering,vol.44, no.11,pp.1024–1038,2017.
- [74] J. Svajlenko, J. F. Islam, I. Keivanloo, C. K. Roy, and M. M. Mia, “Towards big data curated benchmark of inter-project code clones,” in 2014 IEEE International Conference on Software Maintenance and Evolution. IEEE, 2014, pp. 476–480.
- [75] J. Svajlenko and C. K. Roy, “Evaluating clone detection tools with bigclonebench,” in 2015 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 2015, pp. 131–140.
- [76] H. Sajjani, V. Saini, J. Svajlenko, C. K. Roy, and C. V. Lopes, “Sourcererc: Scaling code clone detection to big-code,” in Proceedings of the 38th International Conference on Software Engineering, 2016, pp. 1157–1168.
- [77] J. Svajlenko and C. K. Roy, “Bigcloneeval: A clone detection tool evaluation framework with bigclonebench,” in 2016 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 2016, pp. 596–600.
- [78] H. Yu, W. Lam, L. Chen, G. Li, T. Xie, and Q. Wang, “Neural detection of semantic code clones via tree-based convolution,” in 2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC). IEEE, 2019, pp. 70–80.
- [79] Y. Fan, X. Xin, D. Lo, and A. E. Hassan, “Chaff from the wheat: Characterizing and determining valid bug reports,” IEEE Transactions on Software Engineering, vol. PP, no. 99, pp. 1–1, 2018.
- [80] X. Han, T. Yu, and D. Lo, “Perflearner: learning from bug reports to understand and generate performance test frames,” in Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering. ACM, 2018, pp. 17–28.
- [81] W. Wen, T. Yu, and J. H. Hayes, “Colua: Automatically predicting configuration bug reports and extracting configuration options,” 2016 IEEE 27th International Symposium on Software Reliability Engineering (ISSRE). IEEE, 2016, pp. 150–161



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)