# Investigating Adversarial Attacks in Software Analytics via Machine Learning Explainability

MD ABDUL AWAL\*, MRIGANK ROCHAN, and CHANCHAL K. ROY, Department of Computer Science, University of Saskatchewan, Canada

With the recent advancements in machine learning (ML), numerous ML-based approaches have been extensively applied in software analytics tasks to streamline software development and maintenance processes. Nevertheless, studies indicate that despite their potential usefulness, ML models are vulnerable to adversarial attacks, which may result in significant monetary losses in these processes. As a result, the ML models' robustness against adversarial attacks must be assessed before they are deployed in software analytics tasks. Despite several techniques being available for adversarial attacks in software analytics tasks, exploring adversarial attacks using ML explainability is largely unexplored. Therefore, this study aims to investigate the relationship between ML explainability and adversarial attacks to measure the robustness of ML models in software analytics tasks. In addition, unlike most existing attacks that directly perturb input-space, our attack approach focuses on perturbing feature-space. Our extensive experiments, involving six datasets, three ML explainability techniques, and seven ML models, demonstrate that ML explainability can be used to conduct successful adversarial attacks on ML models in software analytics tasks. This is achieved by modifying only the top 1-3 important features identified by ML explainability techniques. Consequently, the ML models under attack fail to accurately predict up to 86.6% of instances that were correctly predicted before adversarial attacks, indicating the models' low robustness against such attacks. Finally, our proposed technique demonstrates promising results compared to four state-of-the-art adversarial attack techniques targeting tabular data.

CCS Concepts: • Software and its engineering → Machine learning explainability; • Just-in-Time (JIT) defect prediction, Code clone detection.;

Additional Key Words and Phrases: Software analytics, Machine learning, Adversarial attacks, Explainability, Robustness

#### **ACM Reference Format:**

#### 1 INTRODUCTION

In recent years, machine learning (ML) models have been extensively applied in software analytics tasks such as code clone detection [21, 30, 56, 94], API recommendation [54, 57], automatic code summarization [49, 109], malware detection [22, 68], code completion [18, 81], Just-in-Time (JIT) defect prediction [14, 37, 98], and source code authorship attribution [2, 89]. However, research shows that despite the effectiveness of the state-of-the-art ML models, they are vulnerable to adversarial attacks [9, 16, 34, 44, 64, 70, 75, 82, 97, 101, 107]. For example, carefully crafted data from large open-source software repositories in an API recommender system can lead to malicious API calls by developers [58]. As a result, researchers have proposed various techniques for generating

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, , © 2024 Association for Computing Machinery. ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

adversarial examples to evaluate the robustness of the ML models in software analytics tasks [9, 19, 34, 44, 45, 58, 64, 75, 77, 83, 97, 100, 101, 107].

Adversarial attacks generate perturbations to deceive and manipulate ML model predictions, potentially leading to system failures in real-life applications. For example, failing to predict JIT defect introducing commit using ML models under adversarial attacks can significantly impact software development and maintenance processes and optimize the allocation of limited software quality assurance (SQA) resources [36, 62]. It can also affect the quality maintenance and execution of exhaustive code review activities for all incoming commits with limited SQA resources in modern code review [10, 51]. Consequently, the overall software development and maintenance process may experience significant delays, resulting in financial losses and wasted time, as the underlying ML models struggle to make accurate decisions in the face of adversarial attacks [14, 37]. Furthermore, attacks on ML models used in code review comment classification tasks can affect automatic support for improving non-useful comments in code review [65]. Mondal et al. [53] empirically showed that code clones are directly related to bugs, and the intensity of bug propagation through code cloning is significant. In addition, as software systems evolve and the demand for code grows, the occurrence of code clones may pose a significant threat to vulnerability propagation within these systems [21]. Thus, failure to detect code clones accurately under adversarial attacks can lead to the development of buggy software, resulting in technical debt and a huge loss (\$3.61 per line of code) in monetary value [25, 50]. Therefore, having a thorough understanding of the performance and resilience of ML models in various software analytics tasks regarding input manipulation is crucial for developing robust ML models.

Most research on adversarial attacks in software analytics tasks focuses on deep learning (DL) models trained on the source code [9, 19, 34, 44, 45, 58, 64, 75, 83, 97, 100, 101, 107]. Furthermore, the attack approach targets the input-space (e.g., source code) when generating adversarial examples. However, in other software analytics tasks such as JIT defect prediction, clone detection, and code review comment classification, many classical <sup>1</sup> ML models (not DL) have been extensively used. For example, Feng et al. [21], Hu et al. [30], and Wu et al. [94] demonstrated that ML models (e.g., Random Forest) outperformed several state-of-the-art DL-based approaches (e.g., ASTNN [102], SCDetector [95], DeeSIM [105], FCCA [31], CDLH [91], and TBCNN [55]) in software clone detection on the BigCloneBench [80] dataset. Moreover, in these tasks, features from the input data are extracted using a pre-trained machine learning model or some other metrics first [14, 37, 56] and then represented in tabular format. These tabular data are later used to tackle the tasks of interest. Ravid et al. [72] demonstrated that ML models trained on tabular data outperformed DL models trained on the same data in various tasks. Therefore, it is essential to assess the robustness of ML models trained on tabular data through feature-space manipulation in software analytics tasks. Adversarial attacks in feature-space have been studied in computer vision [33, 73, 96]. However, to the best of our knowledge, this area remains largely unexplored in software analytics tasks. Therefore, our study focuses on bypassing the input-space and instead examines attacks on the feature-space for ML models in software analytics tasks.

Despite significant efforts to advance adversarial attacks targeting DL models trained on source code [19, 100], to the best of our knowledge, adversarial attacks on classical ML models trained on tabular data for various software analytics tasks have not been explored. Moreover, due to fundamental differences in the characteristics of the training datasets and the model's architecture, the methods proposed for adversarial attacks on DL models cannot be directly applied to ML models trained on tabular data. Therefore, there is an urgent need to propose novel techniques to perform adversarial attacks on ML models trained specifically on tabular data in software analytics tasks.

<sup>&</sup>lt;sup>1</sup>In this study, unless otherwise stated, ML models refer to classical machine learning models, not deep learning models.

Furthermore, investigating adversarial attacks using machine learning explainability in software analytics tasks remains largely unexplored. Recently, Severi et al. [70] proposed explanation-guided backdoor poisoning attacks to assess the robustness of malware classifiers. The primary distinction between Severi et al. [70] and our work is that they focused on influencing the ML training process, while our approach is specifically designed to attack the ML inference process. Similarly, Amich et al. [3, 4] employ ML explainability to boost and diagnose evasion attacks on ML models, particularly those trained on image classifier datasets. However, our approach diverges from previous studies [3, 4, 70] in terms of how we select important features, modify them, and attack feature-space while generating adversarial examples, making our work *novel*.

This study aims to assess the robustness of ML models in software analytics tasks by generating adversarial examples using ML explainability techniques such as SHapley Additive exPlanation (SHAP) [46], Local Interpretable Model-agnostic Explanations (LIME) [66], and PyExplainer [62]. The key contributions of our study are listed below:

- We apply existing ML explainability techniques (e.g., SHAP, LIME, and PyExplainer) to figure
  out important features of an instance that influence ML models toward decision-making. We
  also determine the impacts of each feature on the model's prediction for individual instances.
- Our explanation-guided adversarial example generation creates transformed instances, significantly reducing the accuracy of the ML models in software analytics tasks.
- We comprehensively evaluate the proposed adversarial example generation technique using six datasets and seven ML models. Additionally, we compare our approach with four state-of-the-art adversarial attack techniques (e.g., Zoo [17], Boundary attack [11], PermuteAttack [27], and HopSkipJump [15]) designed for ML models in a model-agnostic way. We also compare our approach with two makeshift tools. Our proposed approach demonstrates promising results compared to the baselines and makeshift tools, showcasing its effectiveness in generating adversarial examples to evaluate the robustness of ML models.
- Our code and the corresponding dataset are publicly available to enhance further research<sup>2</sup>.

#### 2 BACKGROUND

This section briefly discusses some basic terminologies related to this study.

Adversarial Attack: When a subtle modification to an original input results in a transformed version resembling the original, we refer to the changed version as an adversarial example. If the modified sample alters prediction made by ML model, it is categorized as an adversarial attack [82].

**Adversarial Robustness**: If a ML model can perform well in adversarial attacks and maintain accuracy on adversarial examples, we consider it robust against such attacks.

**Risk Score**: The Risk Score represents the likelihood of an instance being classified into a predefined class. For instance, let's consider a ML model trained on the JIT defect prediction dataset, such as cross-project mobile apps. In Figure 2(c), the selected instance is predicted as a *clean* commit with a Risk Score of 33%. This indicates a higher likelihood of the instance being classified as a *clean* commit, as the Risk Score is below 50%. Conversely, a higher Risk Score implies a greater probability of the instance being classified as a *buggy* commit.

**Feature Importance Rank**: Feature importance rank refers to the assessment of how much each input feature contributes to the prediction made by a ML model for any given input instance. The rank helps to identify the most significant features that affect the model's output. For example, Figure 2(a) demonstrates the name of each feature in the *y*-axis and the contribution of each feature to the magnitude of the model output in the *x*-axis. This figure also shows that the most important feature is *ndev* (number of developers).

<sup>&</sup>lt;sup>2</sup>Replication-package

**Machine Learning Explainability**: ML explainability refers to the ability to understand and interpret the decisions made by a ML model [66]. In other words, it allows us to explain why a model makes a particular prediction or decision.

#### 3 RESEARCH METHODOLOGY

In this section, we discuss the overall methodology of our explanation-guided adversarial attack approach. Our objective is to assess the robustness of ML models through an extensive study by generating adversarial examples based on altering the top-k important features identified in ML explainability. Figure 1 depicts the workflow of our research methodology, which is divided into five distinct steps described below:

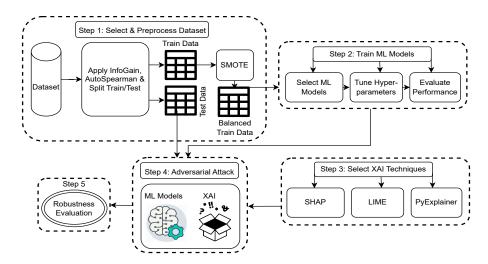


Fig. 1. Our proposed workflow for evaluating the robustness of ML models under adversarial attacks. Step 1 includes the dataset selection and preprocessing methods. Step 2 involves training ML models, while Step 3 involves selecting ML explainability techniques. Step 4 then describes how we perform adversarial attacks using ML explainability, and Step 5 investigates the robustness of ML models under adversarial attacks.

# 3.1 Dataset Selection and Preprocessing

This study aims to evaluate the robustness of ML models by generating adversarial examples through feature-space manipulation. Therefore, dataset selection and preprocessing are vital for enhancing the validity and generalizability of the ML results, ensuring that the models perform optimally before undergoing adversarial attacks. We meticulously select and process the datasets as follows:

3.1.1 Dataset Selection. This study investigates the robustness of ML models against adversarial attacks across various software analytics tasks, including defect prediction, clone detection, and code review comment classification. Additionally, our study shifts the focus of attack from the input-space to the feature-space when generating adversarial examples. Thus, we aim to select software analytics tasks where features extracted from the input (e.g., source code) are utilized to perform the tasks of interest. For example, in cross-language clone detection, various features such as the 'Number of Variables Declared,' 'Total Number of Operators,' 'Number of Arguments,' 'Number of Expressions,' 'Total Number of Operands,' 'Number of Loops (for, while),' 'Number of Exceptions

Thrown,' 'Number of Exceptions Referenced,' and 'McCabe's Cyclomatic Complexity' were extracted first, and ML models were then trained based on these extracted features [56]. Consequently, we select six different datasets: three for defect prediction, two for clone detection, and one for code review comment classification, as the targets for adversarial attacks, where features are extracted from the input (e.g., source code) and stored in tabular format.

For JIT defect prediction, we utilize a dataset for cross-project mobile apps [14] and ten defect prediction datasets from various Java projects [98]. Additionally, we select the Postgres dataset [37] from six subject systems for the Desktop application. The code review comment dataset [65] comprises 15 features and approximately 1,100 data points for detecting useful comments. Finally, we extend our experimental analysis by incorporating the cross-language clone dataset used by Nafi et al. [56] and the BigCloneBench (BCB) dataset [80]. The selected cross-project mobile apps dataset contains 14 features and approximately 30,000 data points, while different Java projects have datasets with 65 features and approximately 25,000 data points. The CLCDSA dataset consists of 18 features and approximately 30,000 data points. Feng et al. [21] extracted six features from the BCB dataset based on some similarity metrics. In our study, we also use these six features to train the ML models. The preprocessed dataset ended up with six features and approximately 546K samples. Please refer to these articles [14, 37, 56, 65, 80] for details about the features of each dataset.

3.1.2 Feature Engineering. The rationale for selecting six datasets in this paper is described in Section 3.1.1. Therefore, to optimize the performance of ML models in our experiment, we adhere to the same feature engineering techniques previously applied to these datasets by other studies. In the study conducted by Catolino et al. [14], the InfoGain method was employed to select a subset of features. From the original 14 features, this technique identified 6 as particularly relevant, which we incorporated into our model. Similarly, Nafi et al. [56] conducted feature selection for their code fragments, retaining 9 features from an initial pool of 24. This process resulted in the derivation of 18 feature values, each associated with a single class label for every clone pair. We utilized these filtered features to train ML models in our study. When dealing with the Java project dataset, we encounter the issue of collinearity among features. To address this challenge, we apply the AutoSpearman technique, which has been utilized by Yatish et al. [98] on the same dataset. Applying this technique allows us to identify and retain 27 pertinent features from an initial set of 65. In our investigations on the Postgres dataset, we also encounter the challenge of the multicollinearity problem. To address this issue effectively, we implement the same feature selection and normalization methodologies outlined in Kamei et al. [37] on this dataset.

Feng et al. [21] proposed a code clone detection method, namely *Toma*, using the BCB dataset and classical machine learning models. They converted the source code into type sequences and applied six similarity metrics, such as the *Jaccard similarity coefficient*, *Dice similarity coefficient*, *Levenshtein distance*, *Levenshtein ratio*, *Jaro similarity*, and *Jaro-Winkler similarity*, to extract six features from the type sequences. Thus, the BCB dataset has six features and is used in our study to train the ML models.

3.1.3 Dataset Balancing. The JIT defect prediction task faces the challenge of imbalanced binary classification, with one class (e.g., clean commit) having significantly more examples than the other (e.g., buggy commit). To tackle this issue, we employ the Synthetic Minority Over-sampling Technique (SMOTE), a technique previously utilized in previous studies [37, 62, 67] on the defect prediction datasets selected for our study. Additionally, we apply SMOTE to rectify imbalanced datasets within the cross-language clone dataset. We must emphasize that we exclusively use this technique for the training data, ensuring that the test dataset remains unaffected by the oversampling process. Finally, the BCB dataset has more clone pairs than non-clone pairs. Since the BCB dataset contains

only 278,838 non-clone pairs, we randomly select approximately 270,000 clone pairs from eight million clones, following the study by Feng et al. [21], ensuring that the BCB dataset used in our study is balanced.

## 3.2 Machine Learning Model Selection and Training

Selecting and training ML models are crucial for the success of any ML project. This process entails choosing the appropriate algorithm or model architecture, training it on the preprocessed dataset, and tuning hyperparameter values to enhance generalization while avoiding overfitting.

- Machine Learning Models. Ravid et al. [72] demonstrated that classical ML models trained on tabular data outperformed DL models trained on the same data in various tasks. In addition, Feng et al. [21], Hu et al. [30], and Wu et al. [94] demonstrated that ML models (e.g., Random Forest) outperformed several state-of-the-art deep learning models ((e.g., ASTNN [102], SCDetector [95], DeeSIM [105], FCCA [31], CDLH [91], and TBCNN [55])) in software clone detection on the BCB dataset. Since the datasets selected for this study contain tabular data, we aim to evaluate the robustness of ML models trained on features extracted from the original input (e.g., source code) presented in a tabular format. Consequently, we focus on ML models commonly used for tabular data in various software analytics tasks, including defect prediction, clone detection, and classification of code review comments. For example, Logistic Regression and Random Forest are commonly used classification techniques in software defect prediction. Thus, we select widely used ML models for our experiments, encompassing three standard ML models (Logistic Regression (LR), Multi-Layered Perceptron (MLP), and Decision Tree (DT)), as well as four ensemble methods (Gradient Boosting Classifier (GBC), Random Forest (RF), AdaBoost (ADA), and Bagging (BAG)). It is worth noting that all of these models have been successfully employed in previous software analytics studies [6, 14, 21, 30, 37, 62, 65, 67, 94].
- 3.2.2 Hyperparameters Tuning. In the realm of model optimization, the tuning of hyperparameters plays a pivotal role in facilitating effective generalization and mitigating the risk of overfitting. To maximize model performance, we explore various hyperparameter combinations, following established standard settings for different ML models [14, 21, 37, 62, 65, 67, 84, 85]. For example, Feng et al. [21] observed that setting the depth parameters to 32, 32, 64, and 16 for RF, DT, Adaboost, and GBDT models, respectively, achieves the highest F1-scores for clone detection on the BCB dataset. In addition, our approach leverages various optimization techniques, including grid search [42], random search [8], and Bayesian optimization [93], to meticulously fine-tune hyperparameters, ultimately selecting the optimal combinations for our specific objectives. This rigorous optimization process is a cornerstone of our ML model training, ensuring that our models are finely calibrated to yield the best possible performance. It is important to note that we only tune hyperparameters while leaving the model's parameters unaltered.
- 3.2.3 Evaluate Model Performance. Model validation is a critical step in ML, ensuring the model's ability to perform well on new, unseen data. Since PyExplainer often requires excessive time to generate an explanation for a single instance, we adopt a dataset-splitting approach, dividing it into training (90%) and testing (10%) subsets to expedite the overall process. To validate the ML models using the 90% training data, we employ a 10-fold cross-validation technique, which was applied in other software analytics studies [6, 21, 30, 37, 67, 94, 98]. Finally, we conduct a thorough evaluation of model performance, utilizing key metrics, including accuracy, F1-score, and the Area Under the Receiver Operating Characteristic Curve (AUC), as recommended by recent studies [14, 67, 98].

The AUC assesses the discriminatory power of models by measuring the true negative rate (coverage of the negative class) on the *x*-axis and the true positive rate (coverage of the positive

class) on the *y*-axis. AUC values range from 0 (indicating the worst performance) to 0.5 (no better than random guessing) up to 1 (representing the best performance) [26]. Therefore, AUC-ROC is a comprehensive tool for evaluating the effectiveness of binary classification models in distinguishing between positive instances (e.g., *buggy* commit) and negative instances (e.g., *clean* commit). Finally, since we aim to perform explanation-guided adversarial attacks on ML models, a test-AUC value of 0.75 is the minimum threshold for ensuring the reliability of explanations [47, 67]. Thus, we evaluate the performance of ML models trained on the selected datasets for this study.

# 3.3 ML Explainability Technique Selection

This study aims to investigate adversarial attacks in software analytics tasks using the ML explainability technique. Our goal is to modify the fewest features possible when generating adversarial examples [48]. Therefore, ML explainability techniques can serve as an effective solution for identifying the important features on which ML models base their decisions. We then modify the top-k features from the important feature list to generate adversarial examples, where k ranges between 1 and half of the total features. The details of adversarial attacks using ML explainability techniques are described in Section 3.4.

While numerous ML explainability techniques have been proposed in recent years to shed light on the decisions made by ML models [46, 66, 86], our focus is solely on model-agnostic techniques, which are applied across various software analytics tasks. Jiarpakdee et al. [36] empirically demonstrated the effectiveness of model-agnostic ML explainability techniques (e.g., LIME [66]) in elucidating the outcomes of defect prediction models. Recently, Roy et al. [67] demonstrated that model-agnostic techniques such as SHAP and LIME could be successfully used for post-hoc analysis of black-box ML models in software analytics tasks. Additionally, Amich demonstrated that LIME [66] and SHAP [46] yield more accurate results when compared to other model-agnostic explanation methods, such as DeepLIFT [71] and LEMNA [24]. Therefore, we select two highly cited model-agnostic ML explainability techniques, SHAP and LIME. Furthermore, Pornprasit et al. [62] introduced the model-agnostic technique PyExplainer, which builds upon the concept of LIME for JIT defect prediction models and supports all ML models available in the *scikit-learn* library. Hence, we also include PyExplainer in our study due to its model-agnostic nature.

#### 3.4 Adversarial Attacks

In computer vision [82], a carefully crafted pixel-level perturbation added to an input image yields an adversarial example that cannot be differentiated from the original input upon visual inspection. The main property of adversarial examples is that they are inputs intentionally modified to cause incorrect predictions while remaining imperceptible to the human eye. However, this is not true for any software analytics task. In software analytics, any transformations to feature values in tabular format or tokens of source code lead to perceptible changes to the human eye. Therefore, adversarial example generation in software analytics tasks differs fundamentally from that in computer vision. Ballet et al. [7], and Cartella et al. [13] introduced a distinct concept of imperceptibility concerning adversarial examples within financial domain data. They altered non-important features based on human judgment while generating these examples. For instance, according to expert analysis, only a subset of features, such as income and age, might be crucial for specific predictions (e.g., fraudulent transaction detection). The attacker should not aim to change these important features; instead, they should target the non-important features so that the generated adversarial examples remain imperceptible. However, these imperceptibility concepts are domain-specific and are not directly applicable to tabular data in software analytics tasks. Therefore, to ensure imperceptibility while attacking, we consider the criterion introduced by Mathov et al. [48], which states that "an attacker should aim to minimize the number of modified features (e.g., a minimal  $\ell_0$  perturbation) when dealing with tabular data." The  $\ell_0$  perturbation is defined as follows:

$$\|\mathbf{x} - \mathbf{x}'\|_{0} = \sum_{i=1}^{n} \mathbf{1} \quad (x_{i} \neq x_{i}' \& f(\mathbf{x}) \neq f(\mathbf{x}'))$$
 (1)

Where **x** is the original input, **x**' is the generated adversarial input obtained by modifying the smallest number of top-k feature values, f(.) is the model's prediction and n is the number of features. The goal of our adversarial attack is to minimize the  $\ell_0$  perturbation as much as possible.

To conduct the adversarial attack while maintaining imperceptibility, as defined in equation 1, this study aims to address two key aspects: (1) which features in a given instance have the most significant influence on guiding ML models to make specific predictions? and (2) how can these influential features be strategically manipulated to create transformed examples effective in adversarial attacks? To address the first aspect, we present a use case using PyExplainer to generate an explanation for an instance (e.g., a data point of the test dataset), as shown in Table 1. This instance includes six distinct features: nd (number of modified directories), nf (number of modified files), la (lines of code added), la (lines of code deleted), ndev (number of developers), and nuc (number of unique changes to the modified files).

Table 1. An instance that we use to generate an explanation in Fig. 2(c) and calculation of one standard deviation (STD) value of each feature using the whole training dataset.

Feature name	nd	nf	la	ld	ndev	nuc
Feature value	1	4	25	0	4	1
One standard deviation	14	72	8377	5592	10	229

Figures 2(c) and 2(d) illustrate how the features *nd* and *ld* influence the LR model's decision-making process. A comparison between Figures 2(c) and 2(d) reveals an increase in the **Risk Score** from 33.0% to 54.0% after changing the feature values in the *red* zone, effectively flipping the ML model's prediction. These two figures collectively demonstrate that increasing the feature values in the *red* zone will elevate the probability of an instance being predicted as a positive class (e.g., *buggy* commit). Conversely, raising the feature values in the *green* zone will reduce the probability of an instance being categorized as a positive class, implying that the given instance will be predicted as a negative class (e.g., *clean* commit). Our manual investigation typically found that PyExplainer identified the number of important features for the selected datasets to be within the range of 1–3.

While PyExplainer directly illustrates which features are responsible for guiding ML models toward a decision, SHAP and LIME provide us with a feature importance rank. In adversarial example generation, the goal is to identify minimal perturbations or changes (e.g., a minimal  $\ell_0$  perturbation) to the input that will cause the model to misclassify the transformed instance. For example, in the case of ML models trained on tabular data, our objective is to modify the minimum number (e.g., top-k) of feature values from the feature importance rank to generate adversarial examples. However, a fundamental question arises: how do we determine the optimal value of k for the top-k features?

To determine the value of k for the top-k important features, we leveraged the concept of the *Elbow* method in the K-means clustering algorithm [39], which we refer to as the *Reverse Elbow Method*. Initially, we modify the top-1 feature, proceed to the top-2 features, and so forth to calculate the ASR (definition of ASR is provided in Section 3.5.) metric value until we reach half the total number of features in the feature importance rank. For example, in Figure 2(a), the feature importance rank of the cross-project mobile apps dataset is displayed from top to bottom.

Investigating Adversarial Attacks in Software Analytics via Machine Learning Explainability Conference acronym 'XX,,

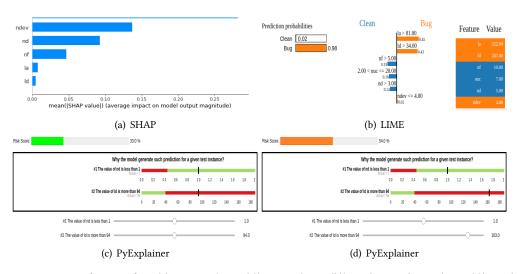


Fig. 2. Important features found by SHAP (Fig. 2(a)), LIME (Fig. 2(b)), and PyExplainer (Fig. 2(c)) on the cross-project mobile apps dataset. Figure 2(c) depicts the explanation and Risk Score when PyExplainer works on the original input. Figure 2(d) shows the explanation and Risk Score when PyExplainer works on the modified features in the guided direction.

As a result, when considering half of the features (e.g., 3 out of 6 features), we obtain three feature combinations (Line 7, Algorithm 1): {ndev}, {ndev}, nd}, and {ndev}, nd, nf}, as shown in Table 2.

Table 2. Determining optimal top-k features for adversarial attacks using ASR metric values.

	Top-1	Top-2	Top-3
ASR	{ndev}	{ndev, nd}	{ndev, nd, nf}
	49.9	58.4	62.4

In each iteration, we select a combination to modify feature values when generating adversarial examples, allowing us to calculate ASR metric (Lines 7–13, as outlined in Algorithm 1). If, in consecutive iterations, the difference in ASR values does not increase, we can halt the process and set the value of k equal to the length of the previous feature combination (Lines 15–34, Algorithm 1). Additionally, if the ASR value decreases after the first iteration, we set the value of k to 1 (Lines 21–24, Algorithm 1). Let us consider a concrete example to illustrate this process. Table 2 presents the ASR values for different feature combinations in each iteration, considering the DT model trained on the cross-project mobile apps dataset. The ASR difference between the  $2^{nd}$  and  $1^{st}$  iterations is 58.4 - 49.9 = 8.5. In contrast, the ASR difference between the  $3^{rd}$  and  $2^{nd}$  iterations is 62.4 - 58.4 = 4. It is evident that the ASR value does not significantly increase after the  $2^{nd}$  iteration. Therefore, we can set the value of k = 2 to determine the top-k features.

We leverage the technique introduced in the previous study by Pornprasit et al. [62] to address the second aspect and perform the adversarial attacks (Lines 8–12, Algorithm 1), and this is one of the many ways to change the feature values [27]. PyExplainer explicitly identifies the features that influence ML models' decision-making processes. It also offers threshold values for each feature, enabling us to make additions or subtractions of one STD with the threshold values in alignment with the guided directions for generating adversarial examples. We calculate STD from the training set as mentioned in the previous study [62]. As SHAP and LIME do not provide specific threshold values with their explanations, we perform one STD addition and subtraction operation relative

## **Algorithm 1:** Explanation-guided Adversarial Attacks

```
Input: ML Model (M), ML Explainability Techniques (XAI) (SHAP, LIME, PyExplainer), Train (X<sub>train</sub>),
           and Test (X_{test}) Data.
   Output: Adversarial Examples (X'_{test}), and ASR Metric Value.
1 ASR \leftarrow 0
                                                                             // ASR for all test samples
                                                                    // ASR for changing top-k features
2 ASR_{topk} \leftarrow 0
                                                              // Store k as top-k important features
3 Top_k \leftarrow 0
4 FC_{ASR} \leftarrow []
                                                           // Store ASR for changing top-k features
5 \ FC_{comb} \leftarrow getFeatureCombination(getFeatureImportanceRank(XAI, X'_{test}))
  // Calculate the ASR and top-k values for different feature combinations by
       optimizing the multi-objective optimization function
6 foreach FC \in FC_{comb} do
       // X'_{test} contains the correctly predicted test instances
       foreach x \in X'_{test} do
            x' \leftarrow getTransformedInstance(x, FC)
8
            y_{pred} = getPrediction(M, x')
            if y_{pred} \neq y_{test} then
10
                ASR_{topk} \leftarrow ASR_{topk} + 1
11
            end
12
       end
13
       FC_{ASR}[FC] \leftarrow ASR_{topk}
14
       for i in range(len(FC<sub>ASR</sub>)) do
15
                                      // Initialize and store ASR difference between iteration
            ASR_{diff} \leftarrow -\infty
16
            if i == 0 then
                ASR_{diff} \leftarrow FC_{ASR}[i]
18
                ASR \leftarrow FC_{ASR}[0]
            end
20
            else if FC_{ASR}[0] > FC_{ASR}[1] then
21
                Top_k = 1
22
                ASR \leftarrow FC_{ASR}[0]
23
                break
24
            end
25
            else if FC_{ASR}[i+1] - FC_{ASR}[i] > ASR_{diff} then
26
                ASR_{diff} \leftarrow FC_{ASR}[i+1] - FC_{ASR}[i]
                ASR \leftarrow FC_{ASR}[i]
28
29
            end
            else
30
                Top_k \leftarrow len(FC_{ASR}[i])
31
                break
32
            end
33
       end
35 end
```

to the original feature values. It is crucial to note that we discard changes that would result in a negative value after subtracting one STD from the current value. This precaution ensures that the transformed instances remain valid, do not violate the fundamental properties of the features, and

do not generate out-of-distribution samples. For example, the feature *number of modified directories* (nd) can't have a negative value after its alteration (e.g., 1 - 14 = -13).

Our objective is to generate adversarial examples such that the feature value changes satisfy the  $\ell_0$  perturbations, causing the ML models under attack to fail to accurately predict the maximum number of these adversarial examples, which were predicted correctly before the adversarial attack. To achieve this, we use a multi-objective optimization function as an aggregation of (1) Finding the minimum number of modified features (e.g.,  $\ell_0$  perturbations) using the Reverse Elbow Method; and (2) Maximizing the ASR metric values. Thus, the multi-objective optimization function can be defined as follows:

$$\min \|\mathbf{x} - \mathbf{x}'\|_0 (f(\mathbf{x}) \neq f(\mathbf{x}')) - \lambda \cdot ASR$$

In the context of the above multi-objective optimization function, the minus sign (–) indicates that we are simultaneously minimizing the number of modified features and maximizing the ASR metric value. Additionally,  $\lambda$  is a weighting factor that balances the two objectives by denoting the number of modified features. Thus, we generate transformed instances to conduct adversarial attacks on the ML models.

#### 3.5 Robustness Evaluation

We evaluate the robustness of the selected ML models using the generated adversarial examples. To assess the quality of these adversarial examples and evaluate the robustness of the ML models, we employ the Attack Success Rate (ASR) metric, as previously used by Yang et al. [97]. The ASR quantifies the percentage of instances that were correctly predicted but are no longer predicted correctly by the ML models following adversarial attacks. The ASR is defined as follows:

$$ASR = \frac{|\{x | x \in X \land M(x') \neq M(x)\}|}{|X|}$$
 (2)

where X is a dataset,  $x \in X$  represents an instance, x' denotes the generated adversarial example, and M denotes the ML model under adversarial attack. A higher ASR value indicates that the generated examples possess sufficient quality to challenge the robustness of the ML models. Conversely, a lower ASR value demonstrates the models' resistance to adversarial attacks.

#### 4 EXPERIMENTAL RESULTS AND ANALYSIS

We implement and evaluate our experiment using the scikit-learn [12], SHAP [46], LIME [66], PyExplainer [62], and Adversarial Robustness Toolbox (ART) [59] libraries. Our study involves six datasets, three ML explainability techniques, and seven classical ML models, resulting in 126 ( $6 \times 3 \times 7 = 126$ ) experimental combinations. Similarly to prior studies [7, 27, 97, 101], we focus exclusively on the proportion of correctly predicted instances in the test dataset. Thus, we start by evaluating the performance of our chosen ML models, and then, we assess the effectiveness of our proposed approach through extensive experiments by addressing the following two research questions:

**RQ1:** Do changes in important feature values found from ML explainability techniques affect ML model's prediction probability?

**RQ2:** Can we use ML explainability techniques to generate adversarial examples to assess the robustness of machine learning models in software analytics tasks?

The study proceeds with a comprehensive analysis of the performance of selected ML models. Subsequently, we investigate how altering the top-k important features influences the prediction probabilities of these models (**RQ1**). We then delve into examining the impact of adversarial attacks

on ML model accuracy, utilizing the *ASR* metric (**RQ2-a**). Finally, we assess the effectiveness of our approach in comparison to baselines in terms of imperceptibility and the ASR metric (**RQ2-b**).

**Baseline Attacks:** To the best of our knowledge, adversarial attacks targeting classical ML models trained on tabular data for different software analytics tasks, such as JIT defect prediction, code review comment classification, and clone detection, have not been explored. Hence, we compare our approach against four state-of-the-art attack techniques: *Zoo* [17], *Boundary attack* [11], *PermuteAttack* [27], and *HopSkipJump* [15], specifically applied for ML models trained on tabular data. We deliberately selected these as baselines due to the black-box nature of adversarial attacks and their publicly available implementations. We briefly describe each of the baseline attack approaches below:

- Zeroth Order Optimization (ZOO) Attack: The ZOO attack is a black-box adversarial attack method used to generate adversarial examples for machine learning models. Unlike gradient-based attacks, ZOO does not require direct access to the model's gradients. Instead, it estimates the gradients by querying the model and using zeroth-order optimization techniques. This makes ZOO particularly useful for attacking black-box models where only input-output pairs are accessible. The key steps involve: (1) Initialization: Start with a clean input sample, initialize the perturbation vector and define parameters such as learning rate, number of iterations, and batch size for gradient estimation; (2) Gradient Estimation: Estimate the gradient of the loss function with respect to the input using finite difference methods. This involves querying the model with slightly perturbed versions of the input and observing the change in the output; (3) Gradient Perturbation: Use the estimated gradients to update the perturbation vector. Apply gradient descent or a similar optimization technique to minimize the loss function with respect to the input; (4) Projection: Project the perturbed input back to the valid input-space to ensure it remains within allowable bounds (e.g., valid pixel values for images); (5) Repeat: Iterate over steps 2, 3, and 4 for a predefined number of iterations or until the adversarial example successfully fools the model.
- Boundary Attack: The Boundary Attack is an iterative, gradient-free method for finding adversarial examples. It starts with an initial adversarial example or a sample near the decision boundary and refines the perturbation through multiple iterations to find one that successfully fools the model. The steps involve: (1) Initialization: Start with a set of valid inputs close to the decision boundary of the model; (2) Boundary Perturbation: Perturb the input along the direction of the decision boundary, making minor adjustments; (3) Optimization: Adjust the perturbation iteratively to find the minimal perturbation that causes the model to misclassify the input. This method is computationally efficient and can be applied to a wide range of models.
- HopSkipJump Attack: The HopSkipJumpAttack is a hyperparameter-free and query-efficient adversarial attack technique designed to generate adversarial examples by perturbing the input data to maximize the model's misclassification while remaining within a specified distance from the original input. It works in three steps: (1) Hop: Perturb the input data using small, incremental changes; (2) Skip: Skip intermediate steps that do not show significant progress towards generating a successful adversarial example; (3) Jump: Apply more significant perturbations if needed to achieve the desired adversarial example. HopSkipJump does not rely on gradients to generate adversarial examples, making it suitable for models where gradients are unavailable. Additionally, HopSkipJump is computationally efficient, as it skips unnecessary steps and focuses on promising perturbations.
- **PermuteAttack:** It is a black-box adversarial attack technique that generates counterfactual examples to evaluate the robustness of machine learning models trained on tabular data,

including discrete and categorical variables. PermuteAttack uses gradient-free optimization based on a genetic algorithm to generate adversarial examples. The goal of PermuteAttack is to find the  $x_{perm}$  where the number of altered features is minimized to  $\delta_{0,max}$ , and the change in feature values is minimized within an  $\ell_2$ -ball of  $\delta_{2,max}$ . The permuted sample that meets these two conditions is considered the counterfactual  $x_{cnt}$ . PermuteAttack solves equation 3 by randomly selecting features and changing their values based on the fitness function defined in the equation.

$$\arg\max_{c \in C} f(x_{\rm cnt}) = t \text{ such that}$$
 
$$\|x_{orig} - x_{cnt}\|_0 \le \delta_{0,max} \text{ and } \|x_{orig} - x_{cnt}\|_2 \le \delta_{2,max}$$
 (3)

ComputeFitness(X) = 
$$||f(x)_t - f(x_{orig})_t||_2 - \rho_0 ||x_{orig} - x||_0 - \rho_1 ||x_{orig} - x||_2$$
 (4)

Where  $f(x)_t$  is the outcome for the target class t, samples x' with higher fitness values are selected for the next iteration of the genetic algorithm. The set  $\rho = \{\rho_0, \rho_1\}$  denotes the two conditions mentioned in equation 3.

#### 4.1 Analyse Model Performance

The hyperparameter tuning process for ML models is described in Section 3.2.2. The hyperparameters and their optimized values for the ML models trained on the cross-project mobile app dataset are shown in Table 3. Please also refer to our replication package for the hyperparameter settings of other datasets. After determining the optimized hyperparameters, we trained the ML models.

ML Model	Hyperparameters
LR	C: 1.0, dual: False, max_iter: 140, penalty: 'l2', solver: 'lbfgs'
DT	criterion: 'entropy', max_depth: 15, min_samples_leaf: 20, min_samples_split: 8
RF	bootstrap: False, max_depth: None, max_features: 'sqrt', min_samples_leaf: 1, min_samples_split: 2, n_estimators: 50
MLP	solver: 'adam', learning_rate: 'adaptive', hidden_layer_sizes: (10, 30, 10), alpha: 0.05, activation: 'relu'
ADA	algorithm: 'SAMME.R', learning_rate: 1.02, n_estimators: 20
BAG	max_features: 13, max_samples: 100, n_estimators: 800
GBC	learning_rate: 1, max_depth: 7, min_samples_leaf: 0.1, min_samples_split: 0.1, n_estimators: 200

Table 3. Optimized hyperparameters for different ML models on the CLCDSA dataset after tuning.

Table 4 presents the accuracy, F1-score, and AUC values for our ML models trained on various datasets without adversarial attacks. Notably, when considering the cross-project mobile apps dataset, the LR model displays the lowest AUC value (0.73), while the BAG model achieves the highest AUC value (0.85). On the other hand, ML models trained on the CLCDSA dataset consistently exhibit high AUC values, ranging from 0.82 to 0.96. In contrast, for the code review dataset, Table 4 reveals relatively lower AUC values; however, the AUC values range from 0.63 to 0.96 for all the ML models. A closer look at Table 4 emphasizes that all the ML models maintain acceptable accuracy, F1-scores, and AUC values, indicating their high accuracy and non-overfitting characteristics. Furthermore, all test-AUCs (with the exception of two cases and the code review dataset) exceed 0.75, a benchmark recommended by previous studies as the minimum required to guarantee the reliability of explanations [47, 67].

# 4.2 Present Experimental Results

We present the findings of our experimental study regarding our two research questions. The following sections provide a detailed analysis and the corresponding results, with each research question addressed separately.

									Dat	aset									
ML	Cro	ss Pro	ject	Jav	va Pro	ject	I	Postgres			CLCDSA			Code Review			BigCloneBench		
Model	Acc	F-1	AUC	Acc	F-1	AUC	Acc	F-1	AUC	Acc	F-1	AUC	Acc	F-1	AUC	Acc	F-1	AUC	
LR	0.68	0.49	0.73	0.75	0.37	0.76	0.75	0.57	0.76	0.76	0.81	0.82	0.62	0.69	0.63	0.75	0.74	0.76	
DT	0.79	0.60	0.83	0.80	0.49	0.80	0.78	0.59	0.79	0.86	0.88	0.93	0.67	0.73	0.68	0.84	0.84	0.85	
RF	0.81	0.58	0.84	0.88	0.55	0.85	0.80	0.59	0.82	0.91	0.93	0.96	0.71	0.76	0.72	0.88	0.87	0.88	
MLP	0.75	0.57	0.83	0.80	0.48	0.81	0.71	0.51	0.73	0.83	0.86	0.92	0.60	0.68	0.63	0.83	0.83	0.84	
ADA	0.73	0.56	0.81	0.79	0.47	0.81	0.78	0.58	0.79	0.81	0.77	0.85	0.63	0.65	0.67	0.83	0.82	0.83	
BAG	0.77	0.59	0.85	0.78	0.50	0.83	0.78	0.61	0.82	0.81	0.84	0.89	0.65	0.73	0.71	0.9	0.89	0.9	
GBC	0.81	0.59	0.84	0.81	0.48	0.80	0.78	0.57	0.80	0.85	0.89	0.94	0.63	0.69	0.65	0.87	0.86	0.87	

Table 4. Performance of the ML models on different datasets without adversarial attacks considering different evaluation metrics such as Accuracy (Acc), F1-score (F1), and AUC values.

Finding the Top-k Important Features and Their Impact on Model's Predictions when

**Changed:** Our initial step to answering RQ1 involves identifying the top-k important features from the feature importance rank using the *Reverse Elbow Method* and considering the test dataset. Figure 2 illustrates the features that SHAP, LIME, and PyExplainer identified for the LR model trained on the cross-project mobile apps dataset. SHAP provides a feature importance rank, allowing us to select the top-k important features. Figure 2(a) presents the feature importance rank identified by SHAP from top to bottom, with each feature's relative importance score on the x-axis. For example, the figure highlights that the most important feature is ndev with an approximate relative importance score of 0.13, followed by nd with an approximate relative importance score of 0.096.

Similar to SHAP, LIME also reveals the relative importance of each feature in the model's prediction process. For instance, the explanation in Figure 2(b) indicates that the ML model predicts the given instance as a buggy commit with a 98% probability. Furthermore, we observe that the top-3 important features are la, ld, and nf based on their contribution scores in influencing the model's prediction. In contrast to SHAP and LIME, PyExplainer explicitly identifies the most important features in the explanation that influence the model's prediction. For example, Figure 2(c) displays the important features nd and ld identified by PyExplainer.

Figure 3 illustrates variations in ASR values resulting from alterations in the values of the top-k feature combinations in SHAP and LIME approaches. We present experimental results for ML models trained on the cross-project mobile apps, Java project, CLCDSA, and BCB datasets. Regarding the SHAP and cross-project mobile apps dataset, the BAG models exhibit a sharp increase in ASR value when changing from top-1 to top-2 feature values, followed by a flat increase. The GBC model shows a linear increase in the ASR metric value. The LR model displays a decrease in ASR metric value after changing the top-1 feature. The ASR metric value remains relatively flat for the DT model after changing the top-2 feature values. A similar trend is observed across all ML models trained on the cross-project mobile apps dataset when we employ the LIME explainability to identify the top-k important features.

Figure 3(e) illustrates that when considering *SHAP*, there is a notable increase in the ASR metric value after changing the top-1 or top-2 features for all ML models trained on the CLCDSA dataset except the LR model. We observe a sharp decrease in ASR metric value for the LR model after changing top-2 features. We observe a very similar result for all the ML models trained on the BCB dataset except the LR model. In the case of the LR model, we observe a sharp decrease after changing the top-2 features. In the case of LIME, a sharp decrease in the ASR metric value is observed after altering the top-1 feature for all ML models trained on the Java dataset, except for LR. Overall, Figure 3 visualizes that significant impacts on the ASR metric value occur after changing the top-2 or top-3 feature values. Thus, we can select the top-*k* feature values from the feature importance rank using *Reverse Elbow Method*.

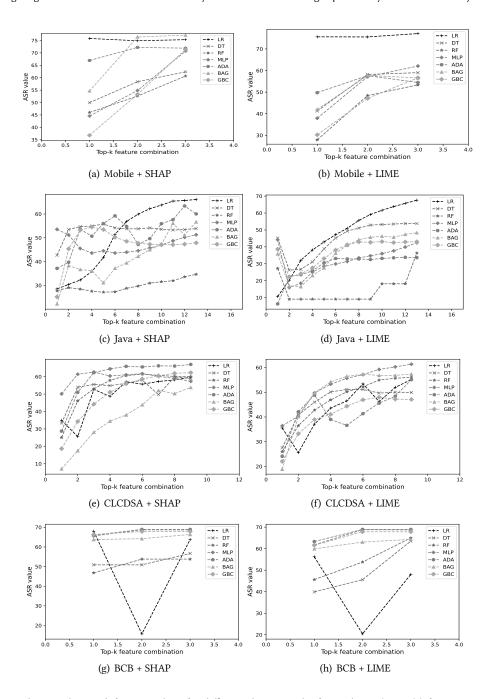


Fig. 3. selecting the top-k feature values for different datasets. The first column (e.g., 3(a) for cross-project mobile apps, 3(c) for Java projects, 3(e) for CLCDSA, and 3(g) for BCB) shows how the ASR metric value changes when different combinations of features are selected from the feature importance rank using SHAP explainability. The second column (e.g., 3(b) for cross-project mobile apps, 3(d) for Java projects, 3(f) for CLCDSA, and 3(h) for BCB) shows how the ASR metric value changes when different combinations of features are selected from the feature importance rank using LIME explainability.

Our *Reverse Elbow Method* effectively identifies the least amount of altered feature values required to impact the model's prediction accuracy to a great extent.

Figures 2(c) and 2(d) illustrate how the features nd and ld influence the LR model's decision-making process. A comparison between these figures reveals that changing the feature values in the guided direction increases the **Risk Score** from 33.0% to 54.0%. The lower Risk Score initially denotes that the given instance is predicted as a clean commit. However, changing the values of these important features increases the Risk Score, indicating that the same instance is more likely to be classified as a buggy commit. Therefore, our objective is to investigate whether there are differences, and to what extent, in prediction probabilities between the original and transformed instances across all test instances (e.g., considering the entire test dataset) and among all the ML models.

Figures 4 and 5 illustrate the distribution of prediction probability differences between the original and transformed instances for the ML models trained on the cross-project mobile apps and BCB datasets. The both figures clearly indicate substantial differences in prediction probabilities between the original and transformed instances, except for the ADA model and SHAP (Figures 4(e) and 5(e)). For example, in the case of the LR model, the mean difference is approximately 0.61 when changing the top-k important feature values identified by SHAP and LIME. The mean probability difference is approximately 0.47 when modifying the top-k important feature values identified by PyExplainer (PyExp). The mean difference consistently exceeds 0.22 for all other ML models and explainability techniques. This denotes a significant prediction probability difference exists between the original and the transformed instances when we modify the original instance based on the top-k important features identified by ML explainability techniques. The only exception is in the case of the ADA model, where SHAP shows a comparatively lower prediction probability difference. In contrast, LIME and PyExplainer maintain consistency compared to SHAP for all the ML models.

The experimental results for the ML models trained on the other four datasets are available in our replication package. However, our experimental findings demonstrate similar results across all other datasets. Therefore, based on our experimental results, we conclude that changing the top-k important feature values identified by ML explainability techniques can be employed to generate adversarial examples to assess the robustness of ML models in software analytics tasks.

## Result RQ1

Changing important feature values identified by ML explainability techniques significantly affects the prediction probability of the ML models in software analytics tasks.

Conducting Adversarial Attacks to assess the Robustness of ML Models: To address this research question, we hypothesize that altering the values of the top-k important features identified through ML explainability techniques should result in a change in the predictions made by the ML models. If there is concrete experimental evidence supporting this hypothesis, it suggests that ML explainability techniques effectively generate adversarial examples capable of assessing the robustness of ML models.

Figures 2(c) and 2(d) illustrate how changing the important features identified through ML explainability affects the prediction probability of the ML model. The research question (RQ1) findings further demonstrate a significant difference in prediction probabilities between the original

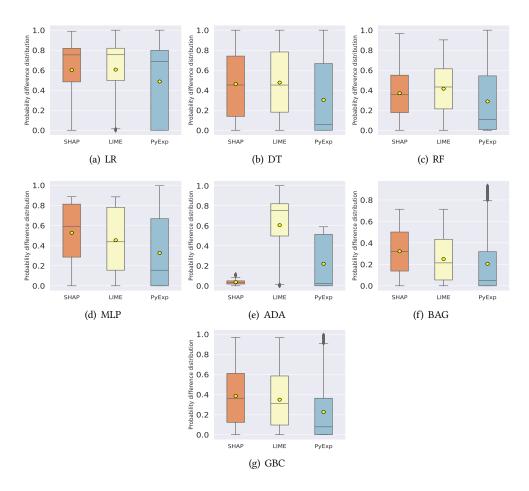


Fig. 4. Distribution of prediction probability differences between the original and transformed instances for various ML models trained on the cross-project mobile apps dataset are shown in Figures 4(a), 4(b), 4(c), 4(d), 4(e), 4(f), and 4(g). These figures represent the distribution for LR, DT, RF, MLP, ADA, BAG, and GBC models, respectively.

and transformed instances when changing the top-k important feature values. Thus, a follow-up question arises: Is the prediction probability difference sufficient to reverse the decision of the ML model? We promptly find the answer to this question in Figures 2(c) and 2(d). These two figures collectively reveal that altering the top-k important features identified by ML explainability flips the prediction of the ML model (changing the colour from green to orange). Now, we investigate whether this scenario applies to all the test datasets and ML models.

Table 5 presents the ASR metric and top-k values both for the explanation-guided adversarial attacks and the selected state-of-the-art adversarial attacks on tabular data. The ASR metric quantifies the degree to which the models' accuracy is compromised when specific features are altered to generate adversarial examples. For instance, in Table 5, it is evident that the BAG model trained on the cross-project mobile apps dataset fails to accurately predict 76.4% of instances that it correctly predicted before adversarial attacks when considering SHAP. Similarly, the LR model fails to predict 75.6% and 56.3% instances accurately after adversarial attacks when considering

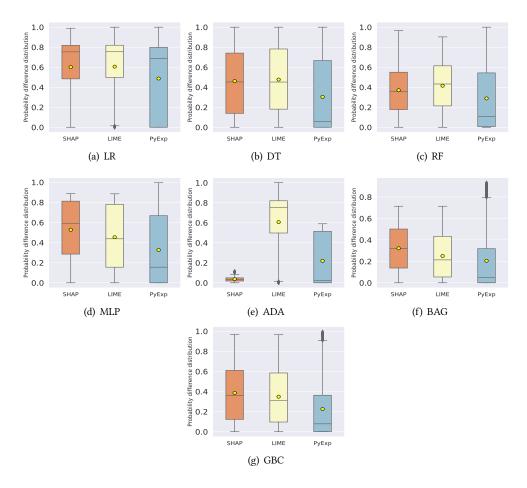


Fig. 5. Distribution of prediction probability differences between the original and transformed instances for various ML models trained on the BigCloneBench dataset are shown in Figures 5(a), 5(b), 5(c), 5(d), 5(e), 5(f), and 5(g). These figures represent the distribution for LR, DT, RF, MLP, ADA, BAG, and GBC models, respectively.

LIME and PyExplainer, respectively. The ASR metric value for the cross-project mobile apps dataset ranges between 23.3% and 76.4%. Similarly, the range is 14.8%–53.6% for the Java project dataset, 11.5%–86.6% for the Postgres dataset, 13.4%–61.4% for the CLCDSA dataset, 17.87%–69.0% for the BCB dataset, and 26.9%–71.6% for the code review dataset, respectively. In summary, the ASR metric value ranges between 11.5%–86.6% for all the ML models, considering all the datasets.

From Table 5, it is evident that ML models fail to accurately predict a large number of instances that they correctly predicted before adversarial attacks, even when altering just up to *three* feature values. The only exception is the LR model trained on the Java project dataset, where we modified *six* feature values to generate adversarial examples. In most cases, we only needed to alter the values of the top-1 or top-2 important features. Interestingly, the LR model trained on the Postgres dataset exhibited the highest ASR metric value (e.g., 86.6%), even when changing only *one* feature. Therefore, in terms of imperceptibility defined in Section 3.4, our proposed explanation-guided adversarial attack shows promising results. These experimental results underscore the effectiveness

Investigating Adversarial Attacks in Software Analytics via Machine Learning Explainability Conference acronym 'XX,,

Table 5. Attack Success Rate (ASR) metric values and the Top-k modified features (e.g., a minimal  $\ell_0$  perturbation) when we apply explanation-guided adversarial attacks and selected baseline attacks on the ML models for various datasets.

ML					Cr	oject Mobil	ject Mobile Apps							
		IAP		ME	PyExplainer		uteAttack		Zoo		ındary		kipJump	
Model	ASR	Top-k	ASR	Top-k		ASR	Top-k	ASR	Top-k	ASR	Top-k	ASR	Top-k	
LR	75.8	Top-1	75.6	Top-1	56.3	88.1	Top-4	78.2	Top-6	56.8	Top-6	55.7	Top-6	
DT	58.4	Top-2	58.2	Top-2	23.3	99.6	Top-3	25.2	Top-2	85.7	Top-6	84.9	Top-6	
RF	60.8	Top-3	48.2	Top-2	32.1	98.1	Top-4	40.8	Top-3	99.7	Top-6	99.9	Top-6	
MLP	70.7	Top-3	54.1	Top-2	38.7	92.1	Top-4	71.5	Top-6	99.9	Top-6	99.9	Top-6	
ADA	72.2	Top-2	57.7	Top-2	48.0	11.4	Top-2	5.1	Top-2	23.3	Top-6	23.1	Top-6	
BAG	76.4	Top-2	57.7	Top-2	28.6	74.1	Top-4	31.1	Top-6	22.2	Top-6	22.2	Top-6	
GBC	71.4	Top-3	47.1	Top-2	28.4	99.3	Top-4	42.8	Top-2	84.4	Top-6	84.9	Top-6	
	Java Project  I.D.   51.5   Top. 6   22.0   Top. 2   22.1													
LR	51.5	Top-6	32.9	Top-2	32.1	99.4	Top-4	35.5	Top-13	11.4	Top-27	11.8	Top-27	
DT	53.5	Top-2	45.2	Top-1	18.4	98.4	Top-5	0.0	Top-0	87.9	Top-27	87.9	Top-27	
RF	29.1	Top-2	27.3	Top-1	14.8	98.6	Top-7	0.0	Top-0	99.9	Top-27	100	Top-27	
MLP	53.5	Top-1	44.3	Top-1	26.3	99.6	Top-6	26.1	Top-13	14.2	Top-27	14.3	Top-27	
ADA	53.6	Top-3	22.4	Top-2	32.9	0.0	Top-0	0.0	Top-0	12.4	Top-27	12.4	Top-27	
BAG	38.3	Top-2	35.8	Top-1	21.7	73.1	Top-7	0.25	Top-3	99.7	Top-27	99.8	Top-27	
GBC	45.4	Top-2	38.6	Top-1	22.2	99.6	Top-6	0.0	Top-0	11.5	Top-27	11.4	Top-27	
		_		_	T		tgres	1					_	
LR	86.6	Top-1	83.4	Top-2	32.0	99.4	Top-5	99.5	Top-9	19.9	Top-12	19.9	Top-12	
DT	80.3	Top-2	82.8	Top-2	25.0	100	Top-4	88.7	Top-3	35.9	Top-12	33.2	Top-12	
RF	70.9	Top-2	71.6	Top-3	20.7	99.6	Top-5	91.6	Top-4	45.5	Top-12	45.3	Top-12	
MLP	67.2	Top-1	63.9	Top-3	36.6	93.5	Top-6	91.1	Top-12	31.9	Top-12	30.8	Top-12	
ADA	77.3	Top-1	64.4	Top-1	21.7	5.6	Top-2	66.8	Top-3	19.9	Top-12	19.9	Top-12	
BAG	66.7	Top-2	72.3	Top-2	11.5	100	Top-4	87.2	Top-7	29.4	Top-12	29.4	Top-12	
GBC	57.3	Top-2	63.7	Top-2	30.7	100	Top-5	88.2	Top-4	88.2	Top-12	90.1	Top-12	
		m .		<b>m</b> .	T		DSA		-		m			
LR	34.5	Top-1	35.5	Top-1	55.2	99.0	Top-6	9.7	Top-13	100	Top-18	100	Top-18	
DT	53.9	Top-2	40.9	Top-2	13.4	95.5	Top-5	0.0	Top-0	74.8	Top-18	74.5	Top-18	
RF	46.1	Top-2	36.5	Top-2	15.3	96.2	Top-8	0.0	Top-0	64.1	Top-18	63.5	Top-18	
MLP	61.4	Top-2	49.8	Top-3	36.8	96.3	Top-6	15.3	Top-13	100	Top-18	100	Top-18	
ADA	51.1	Top-2	42.1	Top-2	17.3	0.0	Top-0	0.0	Top-0	96.2	Top-18	96.1	Top-18	
BAG	28.1	Top-3	40.2	Top-2	16.1	82.9	Top-5	0.63	Top-7	84.2	Top-18	85.2	Top-18	
GBC	34.3	Top-2	33.3	Top-2	14.9	98.5	Top-6	0.0	Top-0	88.3	Top-18	87.8	Top-18	
ID	F2.0	T 1	40.0	Т 1	41.0	_	Review	100	T 10	100	T 15	100	Т 15	
LR DT	53.9	Top-1	48.9	Top-1	41.9	99.3	Top-3	100	Top-10	100	Top-15	100	Top-15	
RF	38.0	Top-2	35.3 32.9	Top-3	26.9 21.1	99.6 99.3	Top-3 Top-4	50.0 86.7	Top-2	68.0	Top-15	66.0	Top-15	
MLP	71.6	Top-3 Top-1	57.5	Top-2	28.3	79.8	Top-4 Top-5	89.5	Top-4 Top-13	100	Top-15 Top-15	100	Top-15 Top-15	
ADA	62.9	Top-1	47.9	Top-2	38.1	79.8 40.0	Top-3	66.4	Top-13	57.8	Top-15	52.1	Top-15	
BAG	27.6	Top-1	23.5	Top-2	26.9	79.3	Top-3	80.7	Top-8	100	Top-15	100	Top-15	
GBC	44.4	Top-2	38.1	Top-2	26.9	100	Top-4	70.4		92.9	_	89.4		
GBC	44.4	10p-2	30.1	10h-2			neBench	70.4	Top-4	74.7	Top-15	07.4	Top-15	
LR	68.0	Top-1	56.2	Top-1	52.4	100	Тор-3	100	Top-5	99.9	Top-6	99.9	Top-6	
DT	50.9	Top-1	63.5	Top-1	17.8	100	Top-3	87.8	Top-3	51.2	Top-6	51.2	Top-6	
RF	53.8	Top-1	64.9	Top-3	20.5	99.9	Top-5	98.6	Top-4	52.6	Top-6	52.6	Тор-6	
MLP	69.5	Top-2	68.6	Top-2	44.4	100	Top-3	100	Top-6	99.4	Top-6	99.0	Тор-6	
ADA	69.0	Top-1	68.9	Top-2	37.3	92.6	Top-4	99.6	Top-6	53.3	Top-6	<b>53.3</b>	Top-6	
BAG	64.7	Top-1	63.0	Top-2	17.8	100	Top-5	97.2	Top-5	53.0	Top-6	53.0	Тор-6	
GBC	67.6	Top-1	67.8	Top-2	18.7	99.8	Top-5	95.5	Top-5	52.2	Top-6	52.2	Тор-6	
GDC	07.0	10b-1	07.0	10p-2	10./	77.0	10p-3	93.3	rop-5	32.2	10h-0	34.4	10h-0	

of the important features identified by ML explainability techniques in generating adversarial examples. Furthermore, our findings emphasize that ML models struggled to accurately predict up to 86.6% of correctly predicted instances after undergoing adversarial attacks.

Following existing studies, we apply sampling techniques (e.g., SMOTE) only to the training data, leaving the test data unaffected. As a result, the test datasets remain highly imbalanced, which impacts the F1-scores of some models. For instance, after the train-test split, the Postgres test dataset includes 1,524 samples of the negative class (e.g., *clean* commits) and 520 samples of the positive class (e.g., *buggy* commits) when testing the LR model. Due to the high imbalance in the test dataset, we perform additional experiments using the undersampling technique to balance the test dataset. Table 6 shows that after balancing the test dataset, we observe improved F1-scores along with nearly similar accuracy and AUC values.

Table 6. Performance of the ML models on different datasets without adversarial attacks considering different evaluation metrics such as Accuracy (Acc), F1-score (F1), and AUC values before balancing (BB) and after balancing (AB) the test data.

-		Dataset														
ML		Java Project						Postgres								
Model	A	сс	F.	-1	Al	JC	A	Acc		Acc		Acc		-1	AUC	
	BB	AB	BB	AB	BB	AB	BB	AB	BB	AB	BB	AB				
LR	0.75	0.69	0.37	0.67	0.76	0.76	0.75	0.72	0.57	0.68	0.76	0.77				
DT	0.80	0.74	0.49	0.72	0.80	0.79	0.78	0.72	0.59	0.67	0.79	0.79				
RF	0.88	0.76	0.55	0.70	0.85	0.87	0.80	0.73	0.59	0.67	0.82	0.82				
MLP	0.80	0.73	0.48	0.70	0.81	0.80	0.71	0.66	0.51	0.63	0.73	0.74				
ADA	0.79	0.76	0.47	0.75	0.81	0.81	0.78	0.72	0.58	0.68	0.79	0.79				
BAG	0.78	0.75	0.50	0.75	0.83	0.82	0.78	0.75	0.61	0.73	0.82	0.82				
GBC	0.81	0.72	0.48	0.68	0.80	0.79	0.78	0.72	0.57	0.67	0.80	0.81				

Table 7 presents the ASR metric and top-k values for both the explanation-guided and selected state-of-the-art adversarial attacks using the balanced test data. From this table, it is evident that by changing only up to three feature values, the accuracy of the ML models drops by up to 85.89%. Interestingly, when changing only one feature value, the LR model fails to predict 85.89% of instances after undergoing adversarial attacks. Therefore, we still get a similar conclusion after balancing the test dataset.

# Result RQ2-a

When altering the top-k feature values to generate adversarial examples and testing ML models on such examples, the accuracy of the ML models under attack can be compromised by up to 86.6%.

Comparison with the baselines: We compare our approach with four state-of-the-art techniques: Zoo [17], Boundary attack [11], PermuteAttack [27], and HopSkipJump [15], in terms of ASR metric values and imperceptibility (e.g., minimal  $\ell_0$  perturbation) across various ML models. It is important to note that we consider the best result (highlighted in bold text) from the three explanation-guided attacks for comparison with the baselines. Table 5 demonstrates that our technique outperforms HopSkipJump and Boundary attacks significantly in terms of imperceptibility. In many cases (highlighted in bold text), we observe that our approach outperforms HopSkipJump and Boundary attacks both in ASR metric values and imperceptibility. Additionally, HopSkipJump and Boundary attacks modify all features of test instances, thereby contradicting the imperceptibility property of the adversarial examples defined in Section 3.4. Shifting to the Zoo attack, in the

Table 7. Attack Success Rate (ASR) metric values and the Top-k modified features (e.g., a minimal  $\ell_0$  perturbation) when we apply explanation-guided adversarial attacks and selected baseline attacks on the ML models for the balanced Java project and Postgres datasets.

ML	Java Project												
Model	SHAP		LIME		PyExplainer	PermuteAttack		Zoo		Boundary		HopSkipJump	
Model	ASR	Top-k	ASR	Top-k		ASR	Top-k	ASR	Top-k	ASR	Top-k	ASR	Top-k
LR	40.98	Top-3	48.78	Top-3	25.06	97.98	Top-7	23.38	Top-13	45.43	Top-27	44.98	Top-27
DT	60.71	Top-2	47.49	Top-1	20.59	97.95	Top-5	0.0	Top-0	55.12	Top-27	55.12	Top-27
RF	44.58	Top-2	39.58	Top-3	9.07	98.75	Top-6	0.0	Top-0	100	Top-27	100	Top-27
MLP	50.09	Top-2	38.4	Top-1	22.26	99.27	Top-7	16.67	Top-12	45.82	Top-27	45.08	Top-27
ADA	46.62	Top-3	41.59	Top-3	19.11	0.0	Top-0	0.0	Top-0	46.61	Top-27	46.42	Top-27
BAG	51.96	Top-2	44.21	Top-2	17.86	56.19	Top-6	0.02	Top-1	99.63	Top-27	99.63	Top-27
GBC	55.97	Top-3	46.27	Top-1	14.37	99.25	Top-6	0.0	Top-0	41.6	Top-27	41.6	Top-27
						Post	gres						
LR	85.89	Top-1	83.06	Top-3	34.54	99.46	Top-6	99.32	Top-8	41.8	Top-12	41.8	Top-12
DT	85.63	Top-2	82.92	Top-2	6.9	85.87	Top-4	85.87	Top-12	66.54	Top-12	68.55	Top-12
RF	65.45	Top-2	72.91	Top-2	21.36	100	Top-4	91.99	Top-4	86.08	Top-12	84.77	Top-12
MLP	56.99	Top-1	59.45	Top-2	22.4	87.94	Top-6	88.08	Top-11	48.9	Top-12	48.76	Top-12
ADA	84.09	Top-3	56.98	Top-1	19.61	8.34	Top-3	56.97	Top-4	40.93	Top-12	40.93	Top-12
BAG	83.59	Top-2	79.56	Top-3	36.98	100	Top-4	86.71	Top-7	50.13	Top-12	49.47	Top-12
GBC	70.99	Top-2	73.54	Top-2	19.75	100	Top-4	85.69	Top-4	93.24	Top-12	92.05	Top-12

majority of cases (highlighted in bold text), our approach surpasses the *Zoo* attack in terms of ASR metric values and imperceptibility. Furthermore, the *Zoo* attack fails to generate a single adversarial example for a few models (highlighted in grey cells). Moreover, regarding imperceptibility, our approach exhibits superior performance over the *Zoo* attack.

Moving to the *PermuteAttack*, our approach shows promising results over *PermuteAttack* in terms of imperceptibility. Except for a few cases (highlighted in bold text), *PermuteAttack* exhibits superior performance regarding ASR metric values. However, for a few ML models, *PermuteAttack* fails to generate a single adversarial example. Moreover, our manual investigation found that, since *PermuteAttack* employs the *Genetic* algorithm and randomly selected features and their values for generating adversarial examples, it produces inconsistent outcomes for the same instance across multiple executions in a row. Therefore, our approach outperforms *PermuteAttack* regarding consistency and imperceptibility.

Table 7 presents the ASR metric and top-k values for both the explanation-guided and selected state-of-the-art adversarial attacks using the balanced test data. From Table 7, it is evident that our explanation-guided adversarial attack approach outperforms the ZOO, HopSkipJump, and Boundary attack techniques in terms of the ASR metric and imperceptibility. Although the PermuteAttack performs well in terms of the ASR metric, our approach achieves a better balance between the ASR metric and imperceptibility. Additionally, the PermuteAttack fails to generate a single adversarial example for the ADA model trained on the Java project dataset and performs poorly on the same model trained on the Postgres dataset. Therefore, we reach a similar conclusion even after balancing the test dataset.

Finally, we extend our experiments by (1) randomly changing the feature values and (2) altering the least important features ranked at the bottom of the list. We refer to these methods as **makeshift tools**: (1) **BL**, where we change the number of features equal to the top-k from the bottom of the feature importance rank, and (2) **BR**, where we randomly select and change the number of features equal to the top-k. It is important to note that for random changes, we exclude the top-k features from the feature importance rank from consideration. This is our design choice, and we deliberately do this because there is a possibility that the most important feature could be randomly selected to

conduct adversarial attacks. Thus, the ASR metric value for **BR** might be high. Finally, we apply the same feature value transformation strategy as mentioned in Section 3.4 for these makeshift tools.

Figure 6 demonstrates that our technique outperforms makeshift tools regarding adversarial attacks on ML models, as measured by the ASR metric. For example, in our adversarial attacks on the MLP model trained on the CLCDSA dataset, changing the top-k important feature values identified by SHAP results in a maximum ASR value of 61%. In contrast, makeshift tools BL and BR lead to ASR values of only 11% and 24%, respectively. We observe a similar result for the MLP trained on the CLCDSA dataset for LIME explainability. For the LR model trained on the CLCDSA dataset, altering the important feature values identified by PyExplainer results in a maximum ASR value of 55%. In comparison, makeshift tools BL and BR lead to ASR values of only 14% and 4%, respectively.

A closer examination of Figure 6 reveals a substantial difference between our approach and the makeshift tools regarding the ASR metric for adversarial attacks when using SHAP explainability. Similar results are observed for LIME and PyExplainer, with only a few exceptions. For instance, the difference between our approach and BL is only 4 for the RF model trained on the CLCDSA dataset with LIME, and it's 0 for the DT model trained on the CLCDSA dataset with PyExplainer. Similar results are observed for the code review datasets. In contrast, a significant difference in the ASR metric is observed for the Postgres and BCB datasets between our approach and the makeshift tools.

Similar to RQ2(a), we conducted additional experiments to compare our approach with makeshift tools regarding the ASR metric using the balanced test datasets. Figure 7 clearly shows that our explanation-guided adversarial attack approach outperforms the makeshift tools by a large margin. Therefore, our approach surpasses the baselines and makeshift tools in terms of the ASR metric and imperceptibility when performing adversarial attacks on the ML models in software analytics tasks.

#### Result RQ2-b

Our approach demonstrates promising results compared to the baselines and makeshift tools in terms of ASR metric and imperceptibility when conducting adversarial attacks on machine learning models in software analytics tasks.

A significant advantage of explanation-guided adversarial attacks is that we do not need to access the model's parameters to generate adversarial examples. Based on our experimental results, the key observations are:

- (1) Modifying just one or two features can significantly affect the accuracy of ML models in software analytics tasks. Therefore, researchers should prioritize the development of more robust ML models and consider implementing countermeasures to mitigate such attacks.
- (2) We find that changing the top-k important features identified by SHAP performs better compared to LIME and PyExplainer in generating adversarial examples. The second best performer is LIME, while PyExplainer's performance is not as good as SHAP and LIME. Thus, our explanation-guided adversarial attacks could be a valuable tool to assess the effectiveness of the explanations offered by different explainability techniques.

#### 5 THREATS TO VALIDITY

This section briefly describes the internal, construct and external threats related to our study.

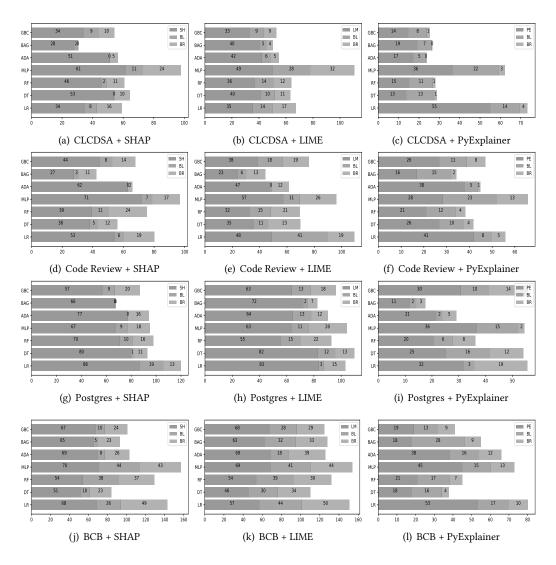


Fig. 6. Comparison of our approach with the makeshift tools regarding ASR metric. The first row compares our approach with the makeshift tools when we change the top-k important features identified by SHAP (SH, Fig. 6(a)), LIME (LM. Fig. 6(b)), and PyExplainer (PE, Fig. 6(c)) for the CLCDSA dataset. Similarly, the second, third, and fourth rows compare the code review, Postgres, and BCB datasets, respectively. Note: ASR metric values are rounded for better visualization.

#### 5.1 Internal Threat

The first internal threat is the accuracy of ML models. However, we countered this threat by choosing the hyperparameter settings described in previous studies [14, 62]. We used SMOTE and Autospearman and applied grid search, random search, and Bayesian optimization to find the best hyperparameter combinations for each ML model. Thus, our selected ML models have achieved good accuracy. We addressed the potential threat concerning the quality of important features identified by ML explainability techniques by choosing two widely recognized model-agnostic

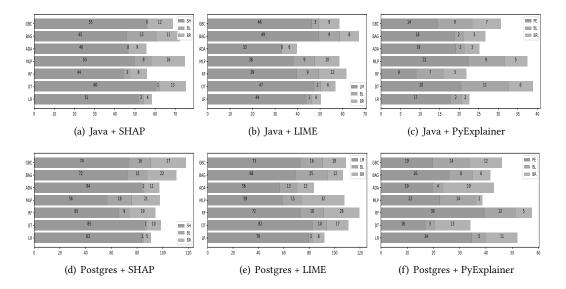


Fig. 7. Comparison of our approach with the makeshift tools regarding ASR metric for the balanced test data. The first and second rows compare our approach with the makeshift tools when we change the top-k important features identified by SHAP, LIME, and PyExplainer for the Java project and Postgres datasets. Note: ASR metric values are rounded for better visualization.

methods, SHAP and LIME. Additionally, we used PyExplainer [62], explicitly designed to explain JIT defect models. Therefore, the important features found by these three techniques are good enough to generate adversarial examples. Another internal threat could be selecting the top-k features to generate adversarial examples. However, Section 3.4 describes how we mitigated this threat. Our experiments show that LIME and PyExplainer sometimes generate different feature importance rankings across multiple executions of the same instance. This variability in experimental outcomes is not a limitation of our study.

#### 5.2 Construct Threat

One potential threat could be how changes to top-k features might affect source code. However, we addressed this concern by considering the characteristics of black-box adversarial attacks performed on extracted features stored in tabular format. Similar to studies in the financial domain [13, 19, 48] and in the software analytics [103], we assume access to both training and test datasets for black-box attacks. We constructed surrogate ML models to carry out adversarial attacks, following the approaches outlined in [13, 48]. In addition, similar to studies in computer vision [33, 73, 96], we attack the feature-space (e.g., metrics or properties of source code in tabular format) while generating adversarial examples. Since we only have access to the extracted features in tabular format (e.g., training and testing data) rather than the source code itself, and considering that these extracted features serve as the direct input to ML models during training, we explicitly state that observing the impact of altered features on the source code falls outside the scope of this study.

Another potential concern is that the preprocessing of the dataset and feature engineering, along with the validation of the model performance and the accuracy of the ML model, could affect the effectiveness of our approach in generating adversarial examples. However, we addressed this concern by comparing our approach with four state-of-the-art adversarial attack techniques on the

same trained ML models and datasets. Our approach demonstrates promising results compared to baselines and makeshift tools discussed in Section 4. Therefore, as long as we compare our approach with the baselines using the same trained ML models and datasets, this concern does not affect the effectiveness of our approach.

#### 5.3 External Threat

The generalizability of our experimental results might threaten the external validity. However, we resolved this threat by choosing seven classical ML models extensively studied in different software analytics tasks, three ML explainability techniques and six distinct datasets used in the previous studies [14, 16, 21, 30, 62, 67, 94]. Thus, we considered  $6 \times 3 \times 7 = 126$  experimental combinations for the analysis. Therefore, our extensive experimental findings are possibly sufficiently robust to generalize our study.

#### 6 RELATED WORK

There has been a growing interest in adversarial learning, which encompasses adversarial training, attacks, countermeasures, and generating adversarial examples. This section briefly discusses existing works focusing on adversarial attacks on ML models in software analytics tasks.

Szegedy et al. [82] first introduced the concept of adversarial attack for image classifiers. To the best of our knowledge, we have not found any research on adversarial attacks targeting classical ML models trained on extracted features stored in tabular format within software analytics tasks, such as JIT defect prediction, clone detection, and classification of useful code review comments. Several studies have been done on adversarial attacks targeting ML models trained on tabular data in other domains [1, 7, 20, 23, 27, 29, 32, 38, 48, 52, 60, 61, 74]. Ballet et al. [7] were the first to introduce systematic adversarial attacks on fully connected neural networks trained on tabular financial data. They altered feature values and modified less significant features based on human judgment to craft imperceptible adversarial examples. Similar approaches have been proposed by Grosse et al. [23], Hashemi et al. [27], Levy et al. [48], and Papernot et al. [60]. Additionally, Cartella et al. [13] adapted state-of-the-art adversarial attack techniques such as the Zoo attack [17], boundary attack [11], and HopSkipJump attack [15] from computer vision to the financial domain using tabular data. A different study by Levy et al. [48] introduced an approach for transforming a surrogate model while preserving the characteristics of the original model, enabling the application of established generation techniques. They demonstrated that a slight human imperceptible perturbation of the input image could change the prediction of the image classifiers. All of the aforementioned studies considered various constraints when generating adversarial examples, none of which are applicable to tabular data in software analytics tasks.

A plethora of research has been conducted on adversarial attacks targeting various software analytics tasks. Liu et al. [44], and Erwin et al. [64] proposed a practical black-box attack on source code authorship identification classifiers based on a set of semantically equivalent program transformations. Nguyen et al. [58] showed that state-of-the-art API recommender systems are vulnerable to adversarial attacks if attackers corrupt the training corpus by injecting malicious data. Chen et al. [16] developed the EvenAttack model considering the contributions of the features of instances to the malware detection problem. Liu et al. [45] proposed a novel method (ATMPA) considering the gradient descent and L-norm optimization method. An attacker can use their technique to generate adversarial examples by introducing tiny perturbations to the input data. Grosse et al. and Suciu et al. [77] demonstrated how adversarial examples fool the malware detection classifiers with slightly modified input data. However, none of the methods proposed for adversarial attacks on ML models in software analytics consider the importance of the features identified by ML explainability techniques. Our work aims to manipulate the important features identified by

ML explainability techniques to generate adversarial examples. Moreover, adversarial attacks in the feature-space have been studied in computer vision [33, 73, 96]. However, this area remains largely unexplored in software analytics tasks. Therefore, our work aims to advance research in this direction.

Severi et al. [70] applied ML explainability to select important features and values for malware classifier adversarial attacks. Unlike their work, our approach is designed specifically for ML inference, not training. Furthermore, our methodology diverges in selecting and modifying the top-k important features while generating adversarial examples. Amich et al. [3, 4] employ ML explainability to enhance and diagnose evasion attacks on ML models. Their study focuses on ML models trained on image classifier datasets like MNIST and CIFAR-10, with a different approach to perturbation and adversarial example generation than ours. Zhang et al. [104] proposed a novel ensemble-based adversarial attack approach that focuses on balancing two key aspects: 1) transferability and 2) imperceptibility, based on model interpretability for image data. Sun et al. [78] proposed an explainability-guided, model-agnostic testing framework to assess the robustness of malware detectors through feature-space manipulation. The framework uses Accured Malicious Magnitude (AMM) to identify the fragile features for manipulation while generating adversarial examples. The primary distinction between our work and that of Sun et al. [78] is that we introduced the concept of the Reverse Elbow Method for identifying fragile features (e.g., top-k important features) for feature-space manipulation. Another important distinction is that their adversarial attack technique is a white-box method, assuming the attacker has complete knowledge of the model architecture and training dataset. In contrast, our explanation-guided adversarial attack technique does not have access to any information about the model's architecture. Additionally, we applied a multi-objective optimization function, as defined in Section 3.4, to generate adversarial examples.

Numerous studies have investigated adversarial attacks targeting deep learning models and large language models (LLMs) trained on source code for various software analytics tasks such as clone detection, source code authorship attribution, vulnerability detection, code comment generation, and code summarization [9, 28, 35, 40, 41, 43, 75, 76, 79, 87, 88, 90, 92, 99, 101, 103, 106–108]. Yefet et al. [99] proposed a white-box attack technique, DAMP, that changes the identifier's name in the code snippet based on the gradient information of the victim model. Srikant et al. [76] generated adversarial examples for models of code using optimized obfuscation. Zhang et al. [101] proposed a black-box attack technique MHM based on Metropolis-Hasting sampling-based variable renaming approach. Tian et al. [87] introduced QMDP (Q-Learning-based Markov decision process) to generate adversarial examples for models of code based on the semantic equivalent program transformations. Their experimental results demonstrated that QMDP generates adversarial examples and enhances the robustness of the source code classification models by over 44%.

Pour et al. [63] developed a search-based adversarial test generation framework to measure the robustness of the neural source code embedding methods (i.e., Code2vec, Code2seq, and CodeBERT). Their experimental results show that the generated adversarial examples can, on average, decrease the performance of these embedding methods from 5.41% to 9.58%. Applis et al. [5] proposed a testing framework: LAMPION, to assess the robustness of the ML-based program analysis tools adopting metamorphic program transformations. Zhou et al. [107] studied the robustness of the code comment generation tasks in adversarial settings. They proposed ACCENT to generate adversarial examples by substituting identifiers in code snippets to generate syntactically correct and semantically similar ones to original code snippets. Recently, several works have focused on assessing the robustness of pre-trained code models under adversarial attacks [19, 34, 97, 100]. Yang et al. [97] performed a natural attack on the pre-trained code models. They transformed the code snippet, preserving the original inputs' operational and natural semantics. Bielik et al. [9],

and Springer et al. [75] both evaluated the robustness of models of code in adversarial settings. Schuster et al. [69] demonstrated that carefully designed adversarial files attached to the training corpus make the code completion model vulnerable.

All the studies above pertain to adversarial attacks on DL models or LLMs trained solely on source code for software analytics tasks. However, the attack methods developed for DL models or LLMs are not directly applicable to classical ML models due to differences in the characteristics of the datasets used for training. Additionally, many attacks on DL models rely on renaming identifiers to generate successful adversarial examples, which does not affect the extracted features from the source code. For instance, in the CLCDSA dataset, the feature 'Number of Variables Declared' remains unchanged regardless of identifier renaming. Finally, while the above adversarial attack techniques target input-space manipulation for generating adversarial examples, our approach focuses on feature-space manipulation for this purpose.

#### 7 CONCLUSION

In recent years, ML models have achieved widespread success in software analytics tasks such as JIT defect prediction, clone detection, code comment generation, code completion, API recommendation, malware detection, and code authorship attribution. However, ML models are vulnerable to adversarial attacks when the test input is changed with minimal perturbations, which may lead to substantial monetary losses in the software development and maintenance process. Many techniques have been proposed in the literature to assess the robustness of ML models under adversarial attacks. However, in this paper, we investigated how ML explainability leads to generating adversarial examples to assess the robustness of ML models against attacks on the feature-space. Our experimental results demonstrate a positive correlation between ML explainability and adversarial attacks. Modifying up to the top-3 most influential features identified by ML explainability techniques can create adversarial examples that can be used to assess the robustness of ML models. Our experimental results underscore the importance of developing robust ML models and the countermeasures against explanation-guided adversarial attacks on ML models in software analytics tasks.

#### **DATA AVAILABILITY**

Our code and the corresponding dataset are publicly available to enhance further research<sup>3</sup>.

#### **ACKNOWLEDGMENT**

This research is supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC), and by the industry-stream NSERC CREATE in Software Analytics Research (SOAR).

# DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

During the preparation of this work, the author(s) used Grammarly  $^4$  and ChatGPT  $3.5^5$  to find grammatical mistakes and improve sentence clarity/ presentation. After using these tools/services, the author(s) reviewed and edited the content as needed and take(s) full responsibility for the content of the publication.

<sup>&</sup>lt;sup>3</sup>Replication-package

<sup>&</sup>lt;sup>4</sup>https://app.grammarly.com/

<sup>5</sup>https://chat.openai.com/

#### REFERENCES

- Abdallah Alshantti, Damiano Varagnolo, Adil Rasheed, Aria Rahmati, and Frank Westad. 2023. CasTGAN: Cascaded Generative Adversarial Network for Realistic Tabular Data Synthesis. arXiv preprint arXiv:2307.00384 (2023).
- [2] Bander Alsulami, Edwin Dauber, Richard Harang, Spiros Mancoridis, and Rachel Greenstadt. 2017. Source code authorship attribution using long short-term memory based networks. In Computer Security–ESORICS 2017: 22nd European Symposium on Research in Computer Security, Oslo, Norway, September 11-15, 2017, Proceedings, Part I 22. Springer, 65–82.
- [3] Abderrahmen Amich and Birhanu Eshete. 2021. Explanation-guided diagnosis of machine learning evasion attacks. In Security and Privacy in Communication Networks: 17th EAI International Conference, SecureComm 2021, Virtual Event, September 6–9, 2021, Proceedings, Part I 17. Springer, 207–228.
- [4] Abderrahmen Amich and Birhanu Eshete. 2022. EG-Booster: explanation-guided booster of ML evasion attacks. In Proceedings of the Twelfth ACM Conference on Data and Application Security and Privacy. 16–28.
- [5] Leonhard Applis, Annibale Panichella, and Arie van Deursen. 2021. Assessing robustness of ML-based program analysis tools using metamorphic program transformations. In 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 1377–1381.
- [6] Md Abdul Awal and Chanchal K Roy. 2024. EvaluateXAI: A framework to evaluate the reliability and consistency of rule-based XAI techniques for software analytics tasks. *Journal of Systems and Software* (2024), 112159.
- [7] Vincent Ballet, Xavier Renard, Jonathan Aigrain, Thibault Laugel, Pascal Frossard, and Marcin Detyniecki. 2019. Imperceptible adversarial attacks on tabular data. arXiv preprint arXiv:1911.03274 (2019).
- [8] James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *Journal of machine learning research* 13, 2 (2012).
- [9] Pavol Bielik and Martin Vechev. 2020. Adversarial robustness for code. In International Conference on Machine Learning. PMLR, 896–907.
- [10] Amiangshu Bosu, Michaela Greiler, and Christian Bird. 2015. Characteristics of useful code reviews: An empirical study at microsoft. In 2015 IEEE/ACM 12th Working Conference on Mining Software Repositories. IEEE, 146–156.
- [11] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2017. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. arXiv preprint arXiv:1712.04248 (2017).
- [12] Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake VanderPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. 2013. API design for machine learning software: experiences from the scikit-learn project. In ECML PKDD Workshop: Languages for Data Mining and Machine Learning. 108–122.
- [13] Francesco Cartella, Orlando Anunciacao, Yuki Funabiki, Daisuke Yamaguchi, Toru Akishita, and Olivier Elshocht. 2021. Adversarial attacks for tabular data: Application to fraud detection and imbalanced data. arXiv preprint arXiv:2101.08030 (2021).
- [14] Gemma Catolino, Dario Di Nucci, and Filomena Ferrucci. 2019. Cross-project just-in-time bug prediction for mobile apps: An empirical assessment. In 2019 IEEE/ACM 6th International Conference on Mobile Software Engineering and Systems (MOBILESoft). IEEE, 99–110.
- [15] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. 2020. Hopskipjumpattack: A query-efficient decision-based attack. In 2020 ieee symposium on security and privacy (sp). IEEE, 1277–1294.
- [16] Lingwei Chen, Yanfang Ye, and Thirimachos Bourlai. 2017. Adversarial machine learning in malware detection: Arms race between evasion attack and defense. In 2017 European intelligence and security informatics conference (EISIC). IEEE, 99–106.
- [17] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 2017. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*. 15–26.
- [18] Subhasis Das and Chinmayee Shah. 2015. Contextual code completion using machine learning. *Stanford university, Stanford University* (2015).
- [19] Xiaohu Du, Ming Wen, Zichao Wei, Shangwen Wang, and Hai Jin. 2023. An Extensive Study on Adversarial Attack against Pre-trained Models of Code. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 489–501.
- [20] Mohammad Esmaeilpour, Nourhene Chaalia, Adel Abusitta, Francois-Xavier Devailly, Wissem Maazoun, and Patrick Cardinal. 2022. RCC-GAN: Regularized Compound Conditional GAN for Large-Scale Tabular Data Synthesis. arXiv preprint arXiv:2205.11693 (2022).
- [21] Siyue Feng, Wenqi Suo, Yueming Wu, Deqing Zou, Yang Liu, and Hai Jin. 2024. Machine Learning is All You Need: A Simple Token-based Approach for Effective Code Clone Detection. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.

- [22] Dragoş Gavriluţ, Mihai Cimpoeşu, Dan Anton, and Liviu Ciortuz. 2009. Malware detection using machine learning. In 2009 International multiconference on computer science and information technology. IEEE, 735–741.
- [23] Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick McDaniel. 2017. On the (statistical) detection of adversarial examples. arXiv preprint arXiv:1702.06280 (2017).
- [24] Wenbo Guo, Dongliang Mu, Jun Xu, Purui Su, Gang Wang, and Xinyu Xing. 2018. Lemna: Explaining deep learning based security applications. In proceedings of the 2018 ACM SIGSAC conference on computer and communications security. 364–379.
- [25] Yuepu Guo, Rodrigo Oliveira Spínola, and Carolyn Seaman. 2016. Exploring the costs of technical debt management—a case study. *Empirical Software Engineering* 21 (2016), 159–182.
- [26] James A Hanley and Barbara J McNeil. 1982. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 143, 1 (1982), 29–36.
- [27] Masoud Hashemi and Ali Fathi. 2020. Permuteattack: Counterfactual explanation of machine learning credit scorecards. arXiv preprint arXiv:2008.10138 (2020).
- [28] Jingxuan He and Martin Vechev. 2023. Large language models for code: Security hardening and adversarial testing. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security.* 1865–1879.
- [29] Aoting Hu, Renjie Xie, Zhigang Lu, Aiqun Hu, and Minhui Xue. 2021. TableGAN-MCA: Evaluating membership collisions of GAN-synthesized tabular data releasing. In ACM SIGSAC Conference on CCS. 2096–2112.
- [30] Yutao Hu, Deqing Zou, Junru Peng, Yueming Wu, Junjie Shan, and Hai Jin. 2022. TreeCen: Building tree graph for scalable semantic code clone detection. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*. 1–12.
- [31] Wei Hua, Yulei Sui, Yao Wan, Guangzhong Liu, and Guandong Xu. 2020. FCCA: Hybrid code representation for functional clone detection using attention networks. IEEE Transactions on Reliability 70, 1 (2020), 304–318.
- [32] Jihyeon Hyeong, Jayoung Kim, Noseong Park, and Sushil Jajodia. 2022. An empirical study on the membership inference attack against tabular data synthesis models. In Proceedings of the 31st ACM International Conference on Information & Knowledge Management. 4064–4068.
- [33] Nathan Inkawhich, Wei Wen, Hai Helen Li, and Yiran Chen. 2019. Feature space perturbations yield more transferable adversarial examples. In CVPR. 7066–7074.
- [34] Akshita Jha and Chandan K Reddy. 2023. Codeattack: Code-based adversarial attacks for pre-trained programming language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 14892–14900.
- [35] Jinghan Jia, Shashank Srikant, Tamara Mitrovska, Chuang Gan, Shiyu Chang, Sijia Liu, and Una-May O'Reilly. 2023. Clawsat: Towards both robust and accurate code models. In 2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER). IEEE, 212–223.
- [36] Jirayus Jiarpakdee, Chakkrit Kla Tantithamthavorn, Hoa Khanh Dam, and John Grundy. 2020. An empirical study of model-agnostic techniques for defect prediction models. *IEEE Transactions on Software Engineering* 48, 1 (2020), 166–185.
- [37] Yasutaka Kamei, Emad Shihab, Bram Adams, Ahmed E Hassan, Audris Mockus, Anand Sinha, and Naoyasu Ubayashi. 2012. A large-scale empirical study of just-in-time quality assurance. *IEEE Transactions on Software Engineering* 39, 6 (2012), 757–773.
- [38] Klim Kireev, Bogdan Kulynych, and Carmela Troncoso. 2022. Adversarial Robustness for Tabular Data through Cost and Utility Awareness. arXiv preprint arXiv:2208.13058 (2022).
- [39] Trupti M Kodinariya, Prashant R Makwana, et al. 2013. Review on determining number of Cluster in K-Means Clustering. *International Journal* 1, 6 (2013), 90–95.
- [40] Yuanchun Li, Jiayi Hua, Haoyu Wang, Chunyang Chen, and Yunxin Liu. 2021. Deeppayload: Black-box backdoor attack on deep learning models through neural payload injection. In 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE). IEEE, 263–274.
- [41] Zhen Li, Guenevere Chen, Chen Chen, Yayi Zou, and Shouhuai Xu. 2022. Ropgen: Towards robust code authorship attribution via automatic coding style transformation. In *Proceedings of the 44th International Conference on Software Engineering*. 1906–1918.
- [42] Petro Liashchynskyi and Pavlo Liashchynskyi. 2019. Grid search, random search, genetic algorithm: a big comparison for NAS. arXiv preprint arXiv:1912.06059 (2019).
- [43] Dexin Liu and Shikun Zhang. 2024. ALANCA: Active Learning Guided Adversarial Attacks for Code Comprehension on Diverse Pre-trained and Large Language Models. In 2024 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER). IEEE, 602–613.
- [44] Qianjun Liu, Shouling Ji, Changchang Liu, and Chunming Wu. 2021. A practical black-box attack on source code authorship identification classifiers. IEEE Transactions on Information Forensics and Security 16 (2021), 3620–3633.
- [45] Xinbo Liu, Jiliang Zhang, Yaping Lin, and He Li. 2019. ATMPA: attacking machine learning-based malware visualization detection methods via adversarial examples. In *Proceedings of the International Symposium on Quality of Service*. 1–10.

- [46] Scott M Lundberg and Su-In Lee. 2017. 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., 4765–4774. http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf
- [47] Yingzhe Lyu, Gopi Krishnan Rajbahadur, Dayi Lin, Boyuan Chen, and Zhen Ming Jiang. 2021. Towards a consistent interpretation of aiops models. ACM Transactions on Software Engineering and Methodology (TOSEM) 31, 1 (2021), 1–38
- [48] Yael Mathov, Eden Levy, Ziv Katzir, Asaf Shabtai, and Yuval Elovici. 2020. Not all datasets are born equal: On heterogeneous data and adversarial examples. arXiv preprint arXiv:2010.03180 (2020).
- [49] Paul W McBurney and Collin McMillan. 2015. Automatic source code summarization of context for java methods. *IEEE Transactions on Software Engineering* 42, 2 (2015), 103–119.
- [50] New York McGraw-Hill Book Co. 2012. Cast worldwide application software quality study: summary of key findings. Cast report Charette RN (1989) Software engineering, risk analysis and management Intertext publications.
- [51] Shane McIntosh, Yasutaka Kamei, Bram Adams, and Ahmed E Hassan. 2014. The impact of code review coverage and code review participation on software quality: A case study of the qt, vtk, and itk projects. In *Proceedings of the 11th* working conference on mining software repositories. 192–201.
- [52] Tiago Leon Melo, João Bravo, Marco OP Sampaio, Paolo Romano, Hugo Ferreira, João Tiago Ascensão, and Pedro Bizarro. 2023. Adversarial training for tabular data with attack propagation. arXiv preprint arXiv:2307.15677 (2023).
- [53] Manishankar Mondal, Banani Roy, Chanchal K Roy, and Kevin A Schneider. 2019. An empirical study on bug propagation through code cloning. Journal of Systems and Software 158 (2019), 110407.
- [54] Laura Moreno, Gabriele Bavota, Massimiliano Di Penta, Rocco Oliveto, and Andrian Marcus. 2015. How can I use this method?. In 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, Vol. 1. IEEE, 880–890.
- [55] Lili Mou, Ge Li, Zhi Jin, Lu Zhang, and Tao Wang. 2014. TBCNN: A tree-based convolutional neural network for programming language processing. arXiv preprint arXiv:1409.5718 (2014).
- [56] Kawser Wazed Nafi, Tonny Shekha Kar, Banani Roy, Chanchal K Roy, and Kevin A Schneider. 2019. Clcdsa: cross language code clone detection using syntactical features and api documentation. In 2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 1026–1037.
- [57] Phuong T Nguyen, Juri Di Rocco, Davide Di Ruscio, Lina Ochoa, Thomas Degueule, and Massimiliano Di Penta. 2019. Focus: A recommender system for mining api function calls and usage patterns. In 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE). IEEE, 1050–1060.
- [58] Phuong T Nguyen, Claudio Di Sipio, Juri Di Rocco, Massimiliano Di Penta, and Davide Di Ruscio. 2021. Adversarial attacks to api recommender systems: Time to wake up and smell the coffee? In 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 253–265.
- [59] Maria-Irina Nicolae, Mathieu Sinn, Minh Ngoc Tran, Beat Buesser, Ambrish Rawat, Martin Wistuba, Valentina Zantedeschi, Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, Ian Molloy, and Ben Edwards. 2018. Adversarial Robustness Toolbox v1.2.0. CoRR 1807.01069 (2018). https://arxiv.org/pdf/1807.01069
- [60] Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. 2016. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. arXiv preprint arXiv:1605.07277 (2016).
- [61] Bart Pleiter, Behrad Tajalli, Stefanos Koffas, Gorka Abad, Jing Xu, Martha Larson, and Stjepan Picek. 2023. Tabdoor: Backdoor Vulnerabilities in Transformer-based Neural Networks for Tabular Data. arXiv preprint arXiv:2311.07550 (2023).
- [62] Chanathip Pornprasit, Chakkrit Tantithamthavorn, Jirayus Jiarpakdee, Micheal Fu, and Patanamon Thongtanunam. 2021. PyExplainer: Explaining the Predictions of Just-In-Time Defect Models. In Proceedings of th International Conference on Automated Software Engineering (ASE). 12 pages.
- [63] Maryam Vahdat Pour, Zhuo Li, Lei Ma, and Hadi Hemmati. 2021. A search-based testing framework for deep neural networks of source code embedding. In 2021 14th IEEE Conference on Software Testing, Verification and Validation (ICST). IEEE, 36–46.
- [64] Erwin Quiring, Alwin Maier, and Konrad Rieck. 2019. Misleading authorship attribution of source code using adversarial learning. In 28th USENIX Security Symposium (USENIX Security 19). 479–496.
- [65] Mohammad Masudur Rahman, Chanchal K Roy, and Raula G Kula. 2017. Predicting usefulness of code review comments using textual features and developer experience. In 2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR). IEEE, 215–226.
- [66] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "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, San Francisco, CA, USA, August 13-17, 2016. 1135-1144.
- [67] Saumendu Roy, Gabriel Laberge, Banani Roy, Foutse Khomh, Amin Nikanjam, and Saikat Mondal. 2022. Why Don't XAI Techniques Agree? Characterizing the Disagreements Between Post-hoc Explanations of Defect Predictions. In

- 2022 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 444-448.
- [68] Justin Sahs and Latifur Khan. 2012. A machine learning approach to android malware detection. In 2012 European intelligence and security informatics conference. IEEE, 141–147.
- [69] Roei Schuster, Congzheng Song, Eran Tromer, and Vitaly Shmatikov. 2021. You autocomplete me: Poisoning vulnerabilities in neural code completion. In 30th USENIX Security Symposium (USENIX Security 21). 1559–1575.
- [70] Giorgio Severi, Jim Meyer, Scott Coull, and Alina Oprea. 2021. {Explanation-Guided} backdoor poisoning attacks against malware classifiers. In 30th USENIX security symposium (USENIX security 21). 1487–1504.
- [71] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *International conference on machine learning*. PMLR, 3145–3153.
- [72] Ravid Shwartz-Ziv and Amitai Armon. 2022. Tabular data: Deep learning is not all you need. *Information Fusion* 81 (2022), 84–90.
- [73] Thibault Simonetto, Salijona Dyrmishi, Salah Ghamizi, Maxime Cordy, and Yves Le Traon. 2021. A unified framework for adversarial attack and defense in constrained feature space. arXiv preprint arXiv:2112.01156 (2021).
- [74] Thibault Simonetto, Salah Ghamizi, Antoine Desjardins, Maxime Cordy, and Yves Le Traon. 2023. Constrained Adaptive Attacks: Realistic Evaluation of Adversarial Examples and Robust Training of Deep Neural Networks for Tabular Data. arXiv preprint arXiv:2311.04503 (2023).
- [75] Jacob M Springer, Bryn Marie Reinstadler, and Una-May O'Reilly. 2020. STRATA: simple, gradient-free attacks for models of code. arXiv preprint arXiv:2009.13562 (2020).
- [76] Shashank Srikant, Sijia Liu, Tamara Mitrovska, Shiyu Chang, Quanfu Fan, Gaoyuan Zhang, and Una-May O'Reilly. 2021. Generating adversarial computer programs using optimized obfuscations. arXiv preprint arXiv:2103.11882 (2021).
- [77] Octavian Suciu, Scott E Coull, and Jeffrey Johns. 2019. Exploring adversarial examples in malware detection. In 2019 IEEE Security and Privacy Workshops (SPW). IEEE, 8–14.
- [78] Ruoxi Sun, Minhui Xue, Gareth Tyson, Tian Dong, Shaofeng Li, Shuo Wang, Haojin Zhu, Seyit Camtepe, and Surya Nepal. 2023. Mate! Are you really aware? An explainability-guided testing framework for robustness of malware detectors. In Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 1573–1585.
- [79] Zhensu Sun, Xiaoning Du, Fu Song, and Li Li. 2023. Codemark: Imperceptible watermarking for code datasets against neural code completion models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1561–1572.
- [80] Jeffrey Svajlenko, Judith F Islam, Iman Keivanloo, Chanchal K Roy, and Mohammad Mamun Mia. 2014. Towards a big data curated benchmark of inter-project code clones. In 2014 IEEE International Conference on Software Maintenance and Evolution. IEEE, 476–480.
- [81] Alexey Svyatkovskiy, Sebastian Lee, Anna Hadjitofi, Maik Riechert, Juliana Vicente Franco, and Miltiadis Allamanis. 2021. Fast and memory-efficient neural code completion. In 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR). IEEE, 329–340.
- [82] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013).
- [83] Phuong T. Nguyen, Davide Di Ruscio, Juri Di Rocco, Claudio Di Sipio, and Massimiliano Di Penta. 2021. Adversarial machine learning: On the resilience of third-party library recommender systems. In *Evaluation and Assessment in Software Engineering*. 247–253.
- [84] Chakkrit Tantithamthavorn, Shane McIntosh, Ahmed E Hassan, and Kenichi Matsumoto. 2016. Automated parameter optimization of classification techniques for defect prediction models. In Proceedings of the 38th international conference on software engineering. 321–332.
- [85] Chakkrit Tantithamthavorn, Shane McIntosh, Ahmed E Hassan, and Kenichi Matsumoto. 2018. The impact of automated parameter optimization on defect prediction models. *IEEE Transactions on Software Engineering* 45, 7 (2018), 683–711.
- [86] Chakkrit Kla Tantithamthavorn and Jirayus Jiarpakdee. 2021. Explainable ai for software engineering. In 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 1–2.
- [87] Junfeng Tian, Chenxin Wang, Zhen Li, and Yu Wen. 2021. Generating Adversarial Examples of Source Code Classification Models via Q-Learning-Based Markov Decision Process. In 2021 IEEE 21st International Conference on Software Quality, Reliability and Security (QRS). IEEE, 807–818.
- [88] Zhao Tian, Junjie Chen, and Zhi Jin. 2023. Code difference guided adversarial example generation for deep code models. In 2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 850–862.
- [89] Farhan Ullah, Junfeng Wang, Sohail Jabbar, Fadi Al-Turjman, and Mamoun Alazab. 2019. Source code authorship attribution using hybrid approach of program dependence graph and deep learning model. IEEE Access 7 (2019), 141987–141999.

- [90] Yao Wan, Shijie Zhang, Hongyu Zhang, Yulei Sui, Guandong Xu, Dezhong Yao, Hai Jin, and Lichao Sun. 2022. You see what i want you to see: poisoning vulnerabilities in neural code search. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1233–1245.
- [91] Huihui Wei and Ming Li. 2017. Supervised deep features for software functional clone detection by exploiting lexical and syntactical information in source code.. In IJCAI. 3034–3040.
- [92] Wai Kin Wong, Huaijin Wang, Pingchuan Ma, Shuai Wang, Mingyue Jiang, Tsong Yueh Chen, Qiyi Tang, Sen Nie, and Shi Wu. 2022. Deceiving Deep Neural Networks-Based Binary Code Matching with Adversarial Programs. In 2022 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 117–128.
- [93] Jia Wu, Xiu-Yun Chen, Hao Zhang, Li-Dong Xiong, Hang Lei, and Si-Hao Deng. 2019. Hyperparameter optimization for machine learning models based on Bayesian optimization. *Journal of Electronic Science and Technology* 17, 1 (2019), 26–40.
- [94] Yueming Wu, Siyue Feng, Deqing Zou, and Hai Jin. 2022. Detecting semantic code clones by building AST-based Markov chains model. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*. 1–13.
- [95] Yueming Wu, Deqing Zou, Shihan Dou, Siru Yang, Wei Yang, Feng Cheng, Hong Liang, and Hai Jin. 2020. SCDetector: Software functional clone detection based on semantic tokens analysis. In Proceedings of the 35th IEEE/ACM international conference on automated software engineering. 821–833.
- [96] Qiuling Xu, Guanhong Tao, Siyuan Cheng, and Xiangyu Zhang. 2020. Towards feature space adversarial attack. arXiv preprint arXiv:2004.12385 (2020).
- [97] Zhou Yang, Jieke Shi, Junda He, and David Lo. 2022. Natural attack for pre-trained models of code. In *Proceedings of the 44th International Conference on Software Engineering*. 1482–1493.
- [98] Suraj Yatish, Jirayus Jiarpakdee, Patanamon Thongtanunam, and Chakkrit Tantithamthavorn. 2019. Mining software defects: Should we consider affected releases?. In 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE). IEEE, 654–665.
- [99] Noam Yefet, Uri Alon, and Eran Yahav. 2020. Adversarial examples for models of code. *Proceedings of the ACM on Programming Languages* 4, OOPSLA (2020), 1–30.
- [100] Zhengran Zeng, Hanzhuo Tan, Haotian Zhang, Jing Li, Yuqun Zhang, and Lingming Zhang. 2022. An extensive study on pre-trained models for program understanding and generation. In 31st ACM SIGSOFT ISSTA. 39–51.
- [101] Huangzhao Zhang, Zhuo Li, Ge Li, Lei Ma, Yang Liu, and Zhi Jin. 2020. Generating adversarial examples for holding robustness of source code processing models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 1169–1176.
- [102] Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, Kaixuan Wang, and Xudong Liu. 2019. A novel neural source code representation based on abstract syntax tree. In 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE). IEEE, 783-794.
- [103] Quan Zhang, Yifeng Ding, Yongqiang Tian, Jianmin Guo, Min Yuan, and Yu Jiang. 2021. Advdoor: adversarial backdoor attack of deep learning system. In Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis. 127–138.
- [104] Rui Zhang, Hui Xia, Zi Kang, Zhengheng Li, Yu Du, and Mingyang Gao. 2024. Harmonizing Transferability and Imperceptibility: A Novel Ensemble Adversarial Attack. *IEEE Internet of Things Journal* (2024).
- [105] Gang Zhao and Jeff Huang. 2018. Deepsim: deep learning code functional similarity. In Proceedings of the 2018 26th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering. 141–151.
- [106] Shasha Zhou, Mingyu Huang, Yanan Sun, and Ke Li. 2024. Evolutionary Multi-objective Optimization for Contextual Adversarial Example Generation. Proceedings of the ACM on Software Engineering 1, FSE (2024), 2285–2308.
- [107] Yu Zhou, Xiaoqing Zhang, Juanjuan Shen, Tingting Han, Taolue Chen, and Harald Gall. 2022. Adversarial robustness of deep code comment generation. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 31, 4 (2022), 1–30.
- [108] Rui Zhu and Cunming Zhang. 2023. How Robust Is a Large Pre-trained Language Model for Code Generation A Case on Attacking GPT2. In 2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER). IEEE, 708–712.
- [109] Yuxiang Zhu and Minxue Pan. 2019. Automatic code summarization: A systematic literature review. *arXiv preprint arXiv:1909.04352* (2019).