



CfExplainer: Explainable just-in-time defect prediction based on counterfactuals

Fengyu Yang^{a,b}, Guangdong Zeng^{a,b,*}, Fa Zhong^{a,b}, Peng Xiao^{a,b}, Wei Zheng^{a,b}, Fuxing Qiu^{a,b}

^a School of software, Nanchang Hangkong University, Nanchang 330063, PR China

^b Software Testing and Evaluation Center, Nanchang Hangkong University, Nanchang 330063, PR China



ARTICLE INFO

Keywords:
Explainable just-in-time defect prediction
Counterfactual
Weighted class association rule
Explanation effectiveness
Explanation reliability

ABSTRACT

Just-in-time (JIT) defect prediction helps rationally allocate testing resources and reduce testing costs. However, most JIT defect prediction models lack explainability, which significantly affects their credibility. Recently, the local interpretable model-agnostic explanations (LIME) method has been used widely in model-explainable research, and many improved LIME-based methods have been proposed. However, problems with respect to explanation effectiveness and reliability remain, which seriously affects the practical use of LIME. To address this problem, CfExplainer, a local rule-based model-agnostic approach, is proposed. The approach first applies counterfactuals to generate synthetic instances. It then mines weighted class association rules based on synthetic instances, and it optimizes the process of generating, ranking, pruning, and predicting the class association rules. Next, it employs the rules with the highest priority to explain the prediction results of the model. Experiments were conducted using the public datasets employed in related studies. Compared to other state-of-the-art methods, in terms of explanation effectiveness, CfExplainer's instance similarity improves by 26.5%–31.2%, and local model fitness improves by 2.0%–3.5%, 2.3%–3%, and 0.7%–7.5% on the AUC, F1-score, and Popt metrics, respectively. In terms of the reliability of the explanation, explanations that are 2.6%–4.7% more unique and 2.5%–5.9% more consistent with the actual characteristics of defect-introducing commits than other state-of-the-art methods. Thus, the explanations of the proposed approach can enhance the model credibility and help guide developers in fixing defects and reducing the risk of introducing them.

1. Introduction

In existing high-complexity software systems, it is difficult for limited software quality assurance resources to cover all test cases. Thus, software testing is not capable of finding all defects. To address this issue, software defect prediction (SDP) techniques have been developed to search for more efficient software testing methods. SDP techniques can help testers more effectively allocate test resources and improve test efficiency (Giray et al., 2023).

Software defect prediction is classified into three categories: within-project (WP) (Uddin et al., 2022), cross-project (CP) (Khatri and Singh, 2022), and just-in-time (JIT) defect prediction (Zhuang et al., 2022). Of these, JIT defect prediction predominates. It utilizes change data as the basis for making accurate predictions. However, most JIT defect prediction models are black-box models that lack explainability, which hinders the practical use of JIT defect prediction models (Zheng et al.,

2021). Therefore, researchers have proposed an explainable SDP technique to help developers understand the causes of incurred defects and provide guidance for fixing them. Recently, the local interpretable model-agnostic explanations (LIME) approach has been proposed (Ribeiro et al., 2016). It is a method for explaining model predictions. When given an instance to be explained, LIME first generates synthetic instances similar to the instance to be explained using the random perturbation method. It then predicts synthetic instances using the black-box model and applies the prediction results as class labels of synthetic instances. Finally, it simulates the predictive behaviours of the black-box model through local models to generate specific explanations of instances to be explained. The above process is shown in Fig. 1.

It should be noted, however, that LIME-based explanations face the following problems. (1) LIME-generated explanations are heavily dependent on synthetic instances. Moreover, the use of randomized perturbation methods may generate synthetic instances that are not

* Corresponding author.

E-mail address: 2117421644@qq.com (G. Zeng).

similar to the instance to be explained, thereby failing to generate explanations that are specific to the instance to be explained. (2) The low fitness of local models leads to explanations that cannot effectively explain the decision-making mechanisms of black-box models, which limits the effectiveness of the explanations. (3) Explanations generated by LIME without considering associations between features, and generating only single-feature rules—i.e. only considering the degree of contribution of single features to prediction results—do not well explain the prediction results and lack operability. The above problems seriously affect the effectiveness of LIME explanations. To address the lack of operability of LIME, (Rajapaksha et al., 2021) proposed SQAPlanner, which uses cross-mutation to generate synthetic instances, explains the predictions of black-box models based on rules and generates four types of guidance and risk thresholds, experimental results show that the rule explanations generated by SQAPlanner are effective, stable and operational. However, SQAPlanner constructs the local model by treating the features as equally important and ignoring the degree of contribution of the features to the defects, which will seriously affect the quality of the mined rules, which will in turn affect the fitness of the local model, and ultimately lead to the lack of effectiveness of the explanations generated by the SQAPlanner, as well as the explanations generated by the SQA-Planner fail to provide the degree of importance of the rules, i.e., they do not know the contribution of the features to the defects, which will make the generated explanations difficult to be applied in software development practice. In addition, to address the effectiveness of LIME explanations and the limitations of SQAPlanner, (Pornprasit et al., 2021) proposed PyExplainer, which adopts the method of cross-mutation (Ibrahim et al., 2021) for generating synthetic instances, and it constructs the local model of RuleFit. Experimental results showed that the explanations generated by PyExplainer outperformed those of other baseline methods. Although LIME and improved methods based on LIME (e.g. SQAPlanner, PyExplainer methods) have achieved better results, random perturbation or cross-mutation methods for generating synthetic instances ignore feature constraints, which could seriously affect the truthfulness of the synthetic instances. For example, when generating synthetic instances, instances may be generated that defy real-world laws (e.g. developers=1.5), synthetic instances lack truthfulness, which will seriously affect the reliability of the explanation, the generated explanation will not provide developers with reliable guidance for correcting defects. In addition, when explaining the same instance, the synthetic instances generated by the improved methods of LIME and LIME using random perturbation or cross mutation have some randomness, which is due to the fact that both of them have random factors in the perturbation of the features, which makes the generated synthetic instances with randomness, and the generated synthetic instances stock in different iteration processes may have large differences, which results in lacking of uniqueness in the generated explanations, i.e. The generated explanation is not specific to the instance to be explained, which will seriously affect the reliability of the explanation. In summary, current explainable software defect prediction faces the following two problems: 1) lack of truthfulness of synthetic instances generated by random perturbation or cross mutation; 2) lack of uniqueness of generated explanations.

To investigate a model's explainability, researchers have proposed

the counterfactual explanation (CE) (Wachter et al., 2017). The CE obtains the prediction results desired by users with minimal input changes by using counterfactual instances, thereby elucidating the decision-making behaviour of the models. To address the lack of explainability of the models, (Cito et al., 2022) investigated the counterfactual explanations of source code models by making minimal changes to the source code for changing the model predictions. Further, they explored the factors affecting the generation of realistic and trustworthy counterfactual explanations, hence verifying the validity of the approach on three different source code models. Meanwhile, (Cheng et al., 2020) designed an interactive visualisation system called the decision explorer with counterfactual explanations (DECE) to provide decision support to users, to help users understand the decision-making mechanisms of the model, provide decision support to users, and also introduce a series of interactive functions that enable users to customise the generation of counterfactual explanations to find actionable explanations that better meet their needs.

Inspired by CEs that generate counterfactual instances by minimising changes, we propose CfExplainer, a model-agnostic approach based on local rules, to explain the prediction results of JIT defect prediction models. CfExplainer draws on the thinking of counterfactual explanation to generate counterfactual instances by minimising changes, while the generated counterfactual instances are stable during different iterations, using counterfactuals to generate synthetic instances, and mining class association rules based on synthetic instances, finally explaining the prediction results of the black-box model based on class association rules. Our scripts are available online(<https://github.com/software-defect-prediction-research/CfExplainer>). The main contributions of this study are as follows:

- (1) CfExplainer's instance similarity is higher than that of other methods and provides more accurate explanations for the model's prediction results.
- (2) CfExplainer has a better local model of fitness than the other methods, generating explanations that are closer to the instances to be explained.
- (3) CfExplainer explanations are reliable and outperform other methods on %Unique and %Consistency metrics.

The remainder of this paper is organised as follows. Section II describes related work, Section III presents CfExplainer, Section IV outlines the experimental design, Section V presents the analysis of the experimental results, Section VI provides discussion, Section VII explains the validity threat analysis, and Section VIII presents conclusions.

2. Related work

Research surveys have shown that developers often doubt the results of defect prediction (Zheng et al., 2022), a sentiment that stems from the unexplainability of the model. Therefore, the explainability of models is as critical as their predictive performance. This section focuses on LIME-based explainable SDP and their counterfactual.

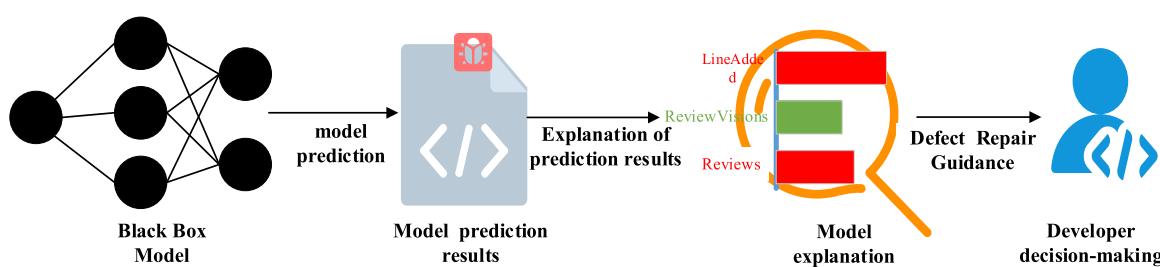


Fig. 1. Explanation of the proposed model.

2.1. Explainable software defect prediction based on lime

LIME explains the prediction of a model by constructing an explainable local model, which is a model-level explanation. It explains the prediction results of the model using the following steps: It (1) selects the instances to be explained and generates synthetic instances by random perturbation; (2) predicts synthetic instances; and (3) constructs explainable local models based on synthetic instances and explains the prediction results of black box models by fitting the predictions of black box models through local models.

Although LIME has been widely used in explainable studies of models, it still has some limitations. For example, the synthetic instances generated by LIME through random perturbations may not be similar to the instances to be explained, and the explanations generated by LIME are heavily reliant on the synthetic instances (Jia et al., 2019), which will seriously affect the effectiveness of the explanations. To address the above limitations, (Jiarpakdee et al., 2020) empirically evaluated three model-agnostic methods, (Gosiewska and IBreakDown, 2019), LIME, and hpo-LIME, where hpo-LIME is a hyper-parameter optimisation of LIME. The experimental results showed that the instance explanations generated by hpo-LIME had a high overlap with the global explanations of defect prediction models, and the generated explanations were reliable. Furthermore, Zafar et al. (Zafar and Khan, 2021) proposed DLIME, which first groups the training data using hierarchical aggregation and then selects the relevant clusters of the instances to be explained using K-nearest neighbours instead of generating synthetic instances using the random perturbation method. Finally, they construct explainable local models based on the relevant clustered instances of the instances to be explained to generate explanations. The experimental results showed that DLIME can generate more reliable explanations. To solve the problem of the effectiveness of the LIME explanation, Pornprasit et al. (Pornprasit et al., 2021) proposed PyExplainer, which uses the cross-mutation method to generate instances and construct local rule models. Their experimental results showed that the effectiveness of the PyExplainer explanation outperforms that of other methods.

However, most current LIME-based-model explainability studies generate synthetic instances mainly by using random perturbation or cross-mutation, ignoring the constraints of the instances. That is, the generated instances may not conform to real-world laws, which leads to the lack of reliability of the generated explanations. Therefore, the generation of effective and reliable explanations based on LIME requires further research.

2.2. Counterfactual

Counterfactual explanation is a method of explaining the predicted results of models. It generates counterfactual instances by constructing hypothetical conditions that change the model's outputs contrary to the actual conditions to explain the model's decision-making mechanism, i.e., the predicted results of the model are explained by the generated counterfactual instances, which is an instance-level explanation method. The main steps in generating counterfactual explanations include: (1) selecting an instance to be explained, generating counterfactual instances based on the instance to be explained, and modifying the feature values of the instance to be explained by minimizing perturbations; (2) inputting counterfactual instances into the model for prediction results that differ from the original prediction. To this end, it compares the differences between the original and counterfactual instances to explain the decision-making mechanisms for model prediction.

To address the issue of model explainability, (Mothilal et al., 2020) proposed a counterfactual explanation that generates diversified counterfactuals that are well approximated to the decision boundaries. It generates counterfactual instances of a quality superior to that of previous methods. Hanmer and Mendiratta (2022) utilized counterfactual explanation in the field of SDP. Their prediction results from a model based on the counterfactual explanation improved the model

transparency, thereby providing guidance to developers for defect correction. Furthermore, Temraz and Keane (2022) proposed the Counterfactual Augmentation (CFA) method to solve the data imbalance problem, CFA generates counterfactual instances in most classes of data to reduce the data imbalance rate compared with other data sampling techniques (e.g., SMOTE (Joloudari et al., 2023)). However, constructing effective counterfactual instances is a key challenge in counterfactual explanation. Maragno et al. (2022) proposed an optimal learning framework for generating counterfactual explanations with constraints that make the generated counterfactual explanations more reliable.

In summary, counterfactual data augmentation is achieved by generating counterfactual instances to reduce the class imbalance rate, while counterfactual explanation is achieved by generating counterfactual instances to explain the prediction results of the model, both via counterfactual instances. Although counterfactual explanations can be used to explain the model's prediction results, counterfactual explanations are designed for defect prediction for continuous releases of software, which requires multiple versions of the software to be trained and evaluated, whereas just-in-time defect prediction is based on developer-submitted data rather than being trained and evaluated based on multiple training versions. Therefore, Pornprasit et al. pointed out in their discussion that counterfactual explanation is not applicable to JIT defect prediction (Pornprasit et al., 2021).

Currently, LIME and LIME-improved explainable defect prediction face a lack of truthfulness of synthetic instances and lack of uniqueness of generated explanations. Inspired by counterfactuals and the fact that there is no research on generating synthetic instances using counterfactuals, these issues motivate us to propose a method for generating synthetic instances based on counterfactuals. The core idea of the counterfactuals approach for generating synthetic instances is to improve the reliability of the explanations by minimally changing the eigenvalues and artificially constraining the range of perturbations.

2.3. Motivating example

We illustrate these problems through an example. As shown in Fig. 2, the LIME random perturbation approach may generate instances that are not similar to those to be explained and which defy the laws of real-world data. Furthermore, the low fitness of the local model cannot effectively simulate the decision-making mechanism of the black box model, which affects the effectiveness of the explanation.

Explanations 1 and 2 in Fig. 2 are the explanations generated by LIME in the first and second iterations, respectively. On the left is the explanation generated by LIME, where the horizontal axis indicates the degree of the feature contribution to the model prediction result. On the right is the actual feature value of the Java file. Orange indicates that the feature supports the prediction of defects; blue indicates that the feature supports the prediction of non-defects. The contribution of "ns" to the prediction of defects is 0.054 since the actual feature value of "ns" is 2, according to the explanation generated by LIME. When "ns>1", "ns" is considered to have a supporting effect on predicting defects, and its contribution of 0.054 is the weight coefficient of the local model. Under different iterations, there are large differences in explanations generated by both, which indicates the lack of reliability of the LIME-generated explanations.

As mentioned previously, explainable software defect prediction based on LIME and improved methods of LIME faces the following problems: 1) lack of truthfulness of generated synthetic instances, i.e., synthetic instances do not have truthfulness; and 2) lack of uniqueness of generated explanations. To address the above problems, we propose CfExplainer.

3. Methodology

This section describes CfExplainer in detail. Its general framework is shown in Fig. 3.

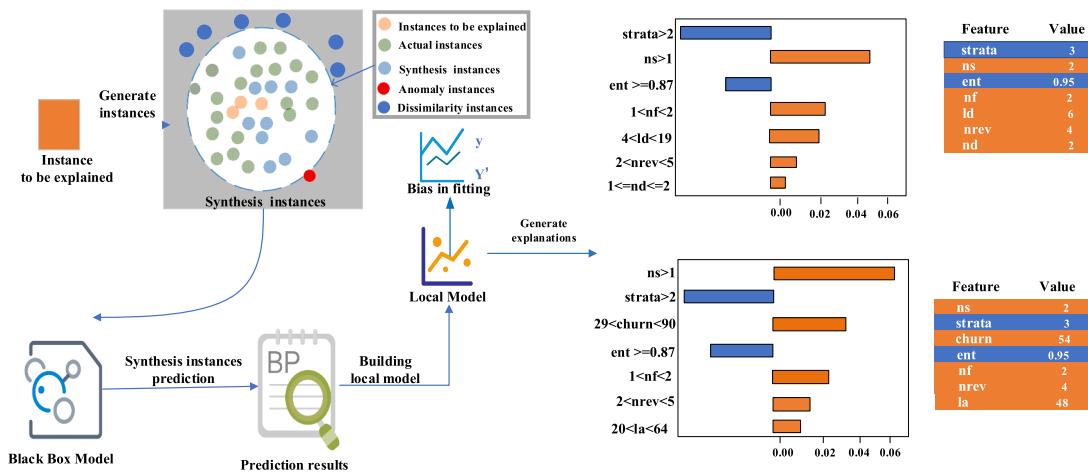


Fig. 2. Explanation of LIME generation.

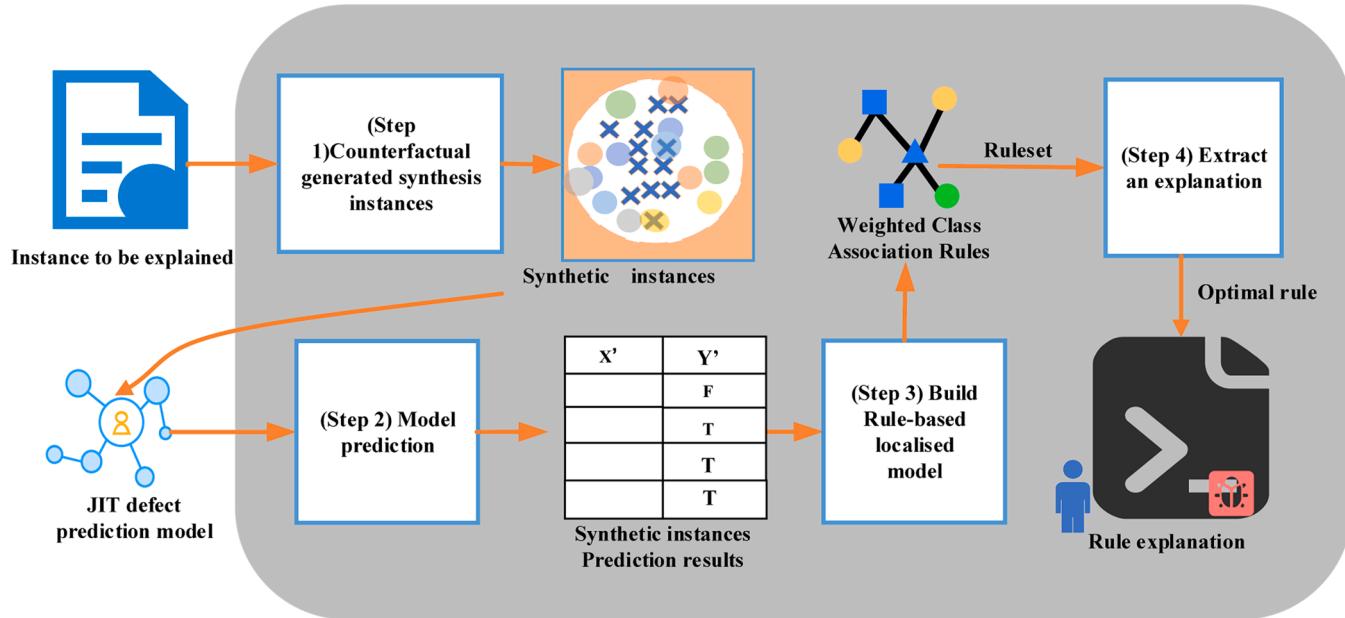


Fig. 3. CfExplainer framework diagram.

CfExplainer generates explanations through the following three steps: (1) generating synthetic instances, (2) constructing local models, and (3) generating explanations.

3.1. Generating synthetic instances

To ensure instance similarity, we selected the top N instances that were the most similar to the instance to be explained from the training data, as used by Rajapaksha et al. (2021), and we set N to 40. However, there may have been an insufficient number of instances, resulting in inadequate mining of the characteristics of the defects introduced by the instances to be explained, which would have affected the reliability of the explanation. Therefore, counterfactuals were introduced to generate synthetic instances.

Similar Instance Selection: select N instances with high similarity from the neighbourhood of the instance to be explained, N instances selection consists of the following three steps.

1. Feature normalisation: To eliminate the effects of different feature scales on instance selection, we performed Z-score normalisation on the training sets.

2. Calculate the instance similarity score: We first calculated the Euclidean distance between instances in each category in the training set and the instance to be explained, i.e., the Euclidean distance between each instance in the defective and non-defective classes and the instances to be explained, and then sorted them in descending order according to the Euclidean distance, and selected the top N instances in each class based on the similarity scores, and the similarity scores are calculated as shown in Eq. (1).

$$K(i_k, i_e) = \exp\left(-\frac{dist(i_k, i_e)^2}{2w^2}\right) \quad (1)$$

where $dist(i_k, i_e)$ is the Euclidean distance between instances i_k and i_e , and w is the kernel width advocated by Ribeiro et al. (2016), that is, the product of 0.75 and the number of instance features.

3. Select top-N similar instances: based on the calculated instance similarity scores, we sort the similarity scores of the instances of each class in descending order, and then we select the top-N instances of each class from the sorted instances.

Counterfactual generation of synthetic instances: Counterfactual

Table 1

Algorithm for rule ranking.

Algorithm 1: Pseudo-code for rule ranking

```

Inputs: DR, NDR;
Outputs: rules
1: DR_time= DR.index
2: NDR_time= NDR.index
3: DR'=DR.sort_values(by=["support", "length","DR_time"], ascending=False)
4: NDR'= NDR.sort_values(by=["support", "length","NDR_time"], ascending=False)
5: return rules ={DR',NDR'};
```

Table 2

Rule pruning algorithm.

Algorithm 2: Pseudo-code for rule pruning.

```

Inputs : DR, NDR;
Outputs : rules
1 : if  $\forall r_1 : X_i \rightarrow C(Wsupp = s1)$  and  $\forall r_2 : X_i \rightarrow C(Wsupp = s2)$  then
2 : if ( $s2 > s1$ ) then
3 : DR' = DR' - {r1};
4 : NDR' = NDR' - {r1};
5 : else
6 : DR' = DR' - {r2};
7 : NDR' = NDR' - {r2};
8 : if  $\forall r' : X_i \rightarrow \text{Defect}$  and  $\forall r' : X_i \rightarrow \text{No Defect}$  then
9 : DR' = DR' - {r'};
10: ND = NDR' - {r'};
11: return rules = {DR',NDR'};
```

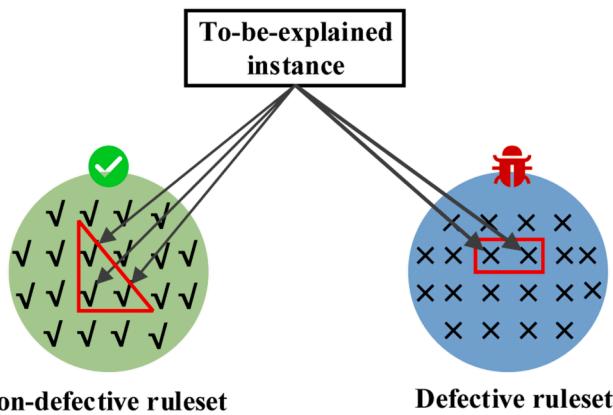


Fig. 4. Multi-rule prediction.

defects; (3) relationships between features and the model's predicted results; i.e. whether features play a supporting or opposing role in model prediction results; and (4) risk probability, or the risk probability of prediction of defects, that is, weighted support.

4. Experimental design

Our experiments were conducted on a computer with a GPU (NVIDIA GTX 3060), CPU (i7-11,700 K, 16 cores), Windows 10 operating system, and Python 3.9 programming language. To reduce the randomness of the conclusions, we repeated the experiment five times, taking the median of the results. This section presents the research dataset and the details of the experimental design, the experimental design is shown in Fig. 5.

4.1. Research questions

To verify the efficacy of CfExplainer, we proposed three research questions from the perspectives of both effectiveness and reliability of CfExplainer explanation, in which RQ1 and RQ2 are to verify the

effectiveness of CfExplainer explanation, and RQ3 is to verify the reliability of CfExplainer explanation, and the research questions are specified as follows.

RQ1: How similar are the synthetic instances generated by CfExplainer to the instance to be explained?

Motivation: The validity of an explanation depends heavily on the quality of the synthetic instances (Pornprasit et al., 2021). Therefore, if the synthetic instances are not similar to the to-be-explained instances, the local model may not provide accurate insights about the logical reasoning of the global model. Therefore, the generated synthetic instances should be similar to the to-be-explained instances, whereas LIME using a random perturbation approach may generate synthetic instances that are not similar to the to-be-explained instances, which can seriously affect the effectiveness of the explanation. For this reason, we need to evaluate the effectiveness of CfExplainer's explanations, and the fitness of CfExplainer's local models will be studied.

RQ2: Does our CfExplainer have higher local model fitness than other baseline methods?

Motivation: the key to model agnostic methods is to build local models with good fitness to simulate the predictive behaviour of the global model. Goodness of fit is often used as a measure of how well a local model approximates the predictions of the global model (Ribeiro et al., 2016), with a higher fit indicating that the local model is able to better simulate the predictive behaviour of the global model. Therefore, local models should have good fitness, however, LIME's local models often fail to fit the global model effectively, i.e., they have low fittedness, which can seriously affect the effectiveness of the explanation. Therefore, to make the CfExplainer explanation effective, it is necessary to construct a high fit local model, i.e., a local model that can accurately fit the predictions of the global model. For this reason, we need to evaluate the effectiveness of the CfExplainer explanation by studying the fitness of the local model of CfExplainer.

RQ3: How reliable is the CfExplainer explanation?

Motivation: different instance synthesis methods generate different instances, which in turn leads to different explanations generated by model-agnostic methods. Currently, model-agnostic methods (LIME, PyExplainer) often use random perturbation or cross-mutation methods to generate synthetic instances, which have a certain degree of

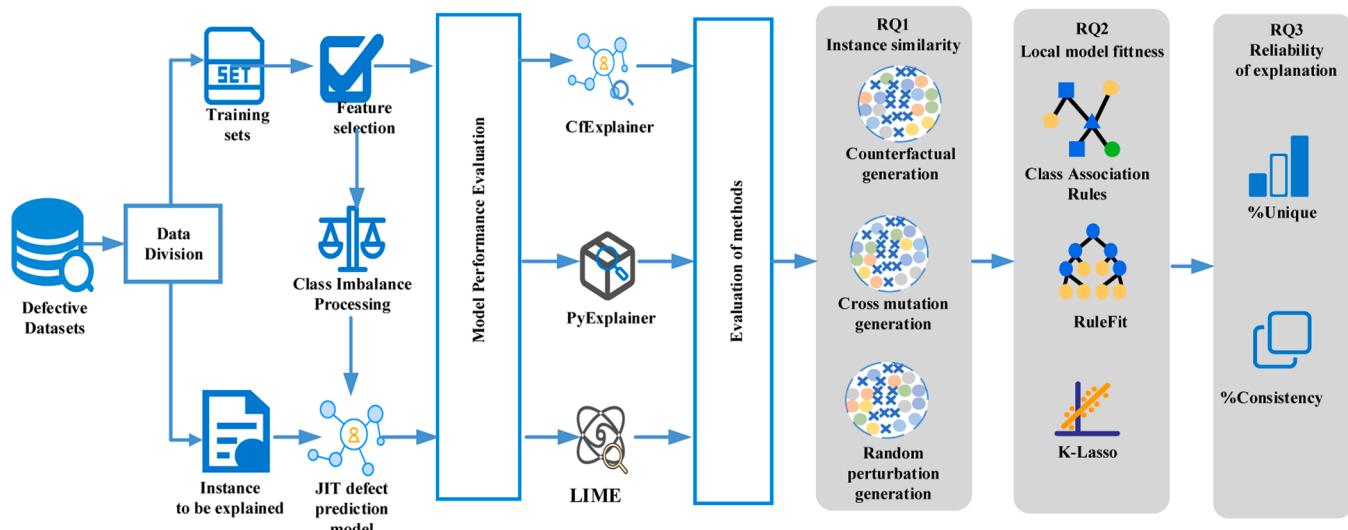


Fig. 5. Experimental design diagram.

randomness, i.e., the difference of instances generated by multiple iterations will produce different explanations of the rules, which in turn affects the reliability of the explanations. In addition, random perturbation or cross mutation methods generate synthetic instances that do not satisfy real-world laws, and the generated synthetic instances lack truthfulness, which can seriously affect the reliability of the explanation. Therefore, we will study the reliability of explanations and evaluate the reliability of explanations generated by the CfExplainer method through metrics.

4.2. Datasets

Mcintosh and Kamei (2018) proposed three important criteria to be satisfied for just-in-time defect dataset selection: traceability, fast iteration, and code review policy, and whether or not a just-in-time defect dataset satisfies the criteria of traceability, fast iteration, and code review policy could seriously affect the effectiveness of the method. Therefore, we selected Openstack and Qt from the open source software projects provided by Mcintosh and Kamei, which both datasets satisfy the above three criteria. In addition, we also chose the Openstack and Qt datasets because these two datasets (1) are often used as benchmarks for defect prediction studies (Pornprasit and JITLine, 2021; Gesi et al., 2021; Zhang et al., 2023); and (2) are able to manually validate the effectiveness of the SZZ algorithm (Śliwski et al., 2005) in order to reduce the number of false positives and false negatives. Where Openstack is an open source software for cloud infrastructure services and Qt is a cross-platform application development framework. Table 3 lists the details of the datasets, where each dataset is categorized into five types of commit-level features, i.e., size (e.g. lines of code added, lines of code deleted), diffusion (e.g. #modified files), history (#developers), experience, and code review activities. The features are show in Table 4.

Table 4
Description of Characteristics.

	Property	Description
Size	Lines added Lines deleted	The number of lines added by a change. The number of lines deleted by a change.
Diffusion	Subsystems Directories Files	The number of modified subsystems. The number of modified directories. The number of modified files.
History	Entropy Unique changes Developers	The spread of modified lines across file. The number of prior changes to the modified files. The number of developers who have changed the modified files in the past.
	Age	The time interval between the last and current changes.
Author/Rev. Experience	Prior changes Recent changes	The number of prior changes that an actor has participated in. The number of prior changes that an actor has participated in weighted by the age of the changes (older changes are given less weight than recent ones).
	Subsystem changes Awareness	The number of prior changes to the modified subsystem(s) that an actor has participated in. The proportion of the prior changes to the modified subsystem(s) that an actor has participated in.
Review	Iterations Reviewers Comments Review window	Number of times that a change was revised prior to integration. Number of reviewers who have voted on whether a change should be integrated or abandoned. The number of non-automated, non-owner comments posted during the review of a change. The length of time between the creation of a review request and its final approval for integration.

Table 3
Statistics on Dataset Information.

Project	Training Sets				Testing Sets			
	Start Date	End Date	# Commits	# Defective Commits	Start Date	End Date	# Commits	# Defective Commits
OpenStack	11/30/2011	08/13/2013	9246	980 (11 %)	08/13/2013	0228/2014	3.963	646(16 %)
Qt	06/18/2011	05/08/2013	19,312	1577 (8 %)	05/08/2013	03/18/2014	8.277	476(6 %)

Eqs. (13) and 14.

$$I_{crossover} = i_1 + (i_2 - i_1) * \alpha \quad (13)$$

$$I_{mutation} = i_1 + (i_2 - i_3) * u \quad (14)$$

where i_1 , i_2 , i_3 are randomly selected instances, α is a random number between 0 and 1, and u is a random number between 0.5 and 1.

Although PyExplainer is able to generate more similar synthetic instances compared to the LIME method, the selection of initial instances can seriously affect the results of synthetic instance generation. In addition, the synthetic instances generated by the random numbers α and u during the process of cross-mutation have a certain degree of randomness, and the values of the features may be changed significantly during the process of cross-mutation, which can all seriously affect the

similarity of the instances. However, CfExplainer generates synthetic instances by minimally perturbing the features through counterfactuals, which ensures the similarity with the original instances. Therefore, the similarity between the synthetic instances generated by our CfExplainer and the instances to be explained is higher than PyExplainer, LIME. In conclusion, the similarity between the synthetic instances generated based on counterfactuals and the instances to be explained is higher than that of random perturbations and cross mutations, and the generation of synthetic instances based on counterfactuals can provide more accurate explanations of the model's prediction results. Therefore, the explanations generated by CfExplainer are closer to the instances to be explained than those generated by LIME and PyExplainer.

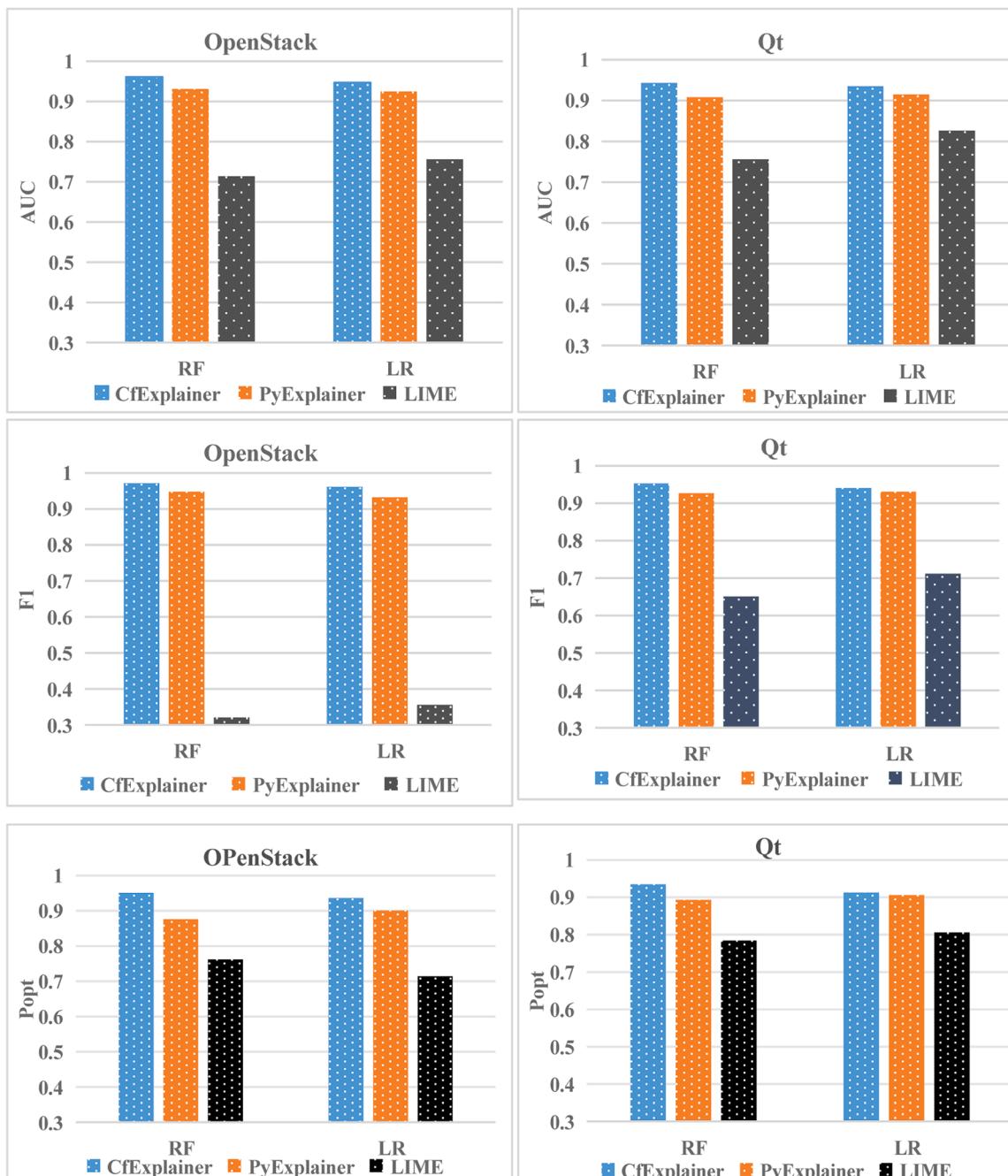


Fig. 7. Plot of the fitness of local model.

5.2. RQ2: does our CfExplainer have higher local model fitness than other baseline methods?

To evaluate whether the local model fit of CfExplainer outperforms other methods, we compared the fitness of the local models of CfExplainer (Association Rule), PyExplainer (RuleFit), and LIME (K-Lasso) using AUC, F1-score, and Popt metrics.

Analysis of local model fitness results: Compared with PyExplainer and LIME, CfExplainer improves 2.0 %–3.5 % and 10.9 %–24.7 % in AUC metrics, 2.3 %–3.0 % and 56.2 %–66.3 % in F1-score metrics, and 0.7 %–7.5 % and 10.7 %–22.2 % in Popt metrics. Meanwhile, CfExplainer has significant improvement in local model fitness in AUC, F1-score, and Popt metrics compared to PyExplainer and LIME.

On the AUC metric, CfExplainer was 0.961 and 0.943 for the OpenStack and Qt datasets, respectively. On the F1-score metric, CfExplainer achieved 0.971 and 0.953 on the OpenStack and Qt datasets, respectively. On the Popt metric, CfExplainer was 0.951 and 0.935 on OpenStack and Qt datasets, respectively. The experimental results are shown in Fig. 7. Furthermore, on the AUC metric, the p-values of CfExplainer versus PyExplainer and LIME are 0.125 and 0.005, respectively; on the F1-score metric, the p-values of CfExplainer versus PyExplainer and LIME are 0.109 and 0.008, respectively; on the Popt metric on the Popt metric, the p-value of CfExplainer versus PyExplainer and LIME is 0.125, 0.006, respectively, the p-value values are limited by the number of datasets, but the Cliff's $|\delta|$ values of CfExplainer versus PyExplainer are 1.0, 0.875, and 1, respectively. The Cliff's $|\delta|$ values of CfExplainer and LIME are 1, this suggests that the CfExplainer approach has large effect sizes and practical significance.

The experimental results show that our local model of CfExplainer outperforms PyExplainer and LIME, which shows that the local model of CfExplainer is able to learn the prediction behaviour of the global model better. The local model fitness is an important factor affecting the explanation effectiveness of CfExplainer. From the experimental results of RQ1, it can be seen that the similarity between the synthetic instances generated by CfExplainer and the to-be-explained instances is better than that of PyExplainer and LIME, and the higher instance similarity of CfExplainer means that the generated synthetic instances are more similar to the distribution of the to-be-explained instances' data in the test set, which makes CfExplainer's local model able to learn the characteristics of the to-be-explained instances better. In addition, CfExplainer's local model-like association rules have better prediction performance in SDP scenarios, and the contribution of features to defect prediction is taken into account, and the process of generating, sorting, pruning, and predicting rules is optimised by the new rule metric weighted support, which improves the fitness of the local model. Therefore, the local model fitness of CfExplainer is better than that of PyExplainer and LIME.

5.3. RQ3: how reliable is the CfExplainer explanation?

Model-agnostic methods often use random perturbation or cross mutation methods to generate synthetic instances, and the generated synthetic instances are random, and different synthetic instances will be generated during different iterations, leading to the generation of different rule explanations for the local model, which seriously affects the uniqueness of the explanations, i.e., the generated explanations are not able to be specific to the to-be-explained instances. Therefore, we investigate the uniqueness of the explanation by measuring whether the generated explanation can be specific to the to-be-explained instance through the $\%Unique$ metric.

Analysis of $\%Unique$ results: The explanations generated by CfExplainer are improved by 2.6 %–4.7 % and 28.2 %–35.8 % in $\%Unique$ metrics compared to PyExplainer and LIME. Meanwhile,

On the OpenStack dataset, the median $\%Unique$ of CfExplainer is 0.971, which increases by 2.7 %–4.7 % and 33 %–34.9 % compared with PyExplainer and LIME methods, respectively. On Qt dataset, the

median $\%Unique$ of CfExplainer method is 0.953, which increases by 2.6 %–4.0 % and 28.2 %–35.8 % compared with PyExplainer and LIME methods, respectively. The experimental results are shown in Fig. 8. On OpenStack and Qt datasets, The p-values of CfExplainer versus PyExplainer and LIME on the $\%Unique$ metric are 0.062 and 0.005, respectively, and Cliff's $|\delta|$ values are all 1, this suggests that the CfExplainer approach has large effect sizes and practical significance

The experimental results show that the uniqueness of our CfExplainer explanations is better than PyExplainer and LIME, which suggests that CfExplainer is able to generate explanations with higher uniqueness. PyExplainer and LIME generate synthetic instances by using cross-mutation and random perturbation, respectively. There is a certain degree of randomness in the synthetic instances, and the synthetic instances generated by both of them in different iteration processes will have large differences, which seriously affects the uniqueness of the explanation, while CfExplainer generates synthetic instances based on the counterfactuals in different iteration processes, which are made by minimally perturbing the features so as to make the global model prediction results change, and by performing conditional constraints on the perturbation of the feature values, which make the synthetic instances generated in different iterative processes are stable and not random. Therefore, CfExplainer generates synthetic instances that are stable and generate more unique explanations than LIME and PyExplainer.

Model-agnostic methods use random perturbation or cross-mutation methods to generate synthetic instances, and when features are perturbed, instances that do not conform to real-world laws (anomalous instances) are generated, and these anomalous instances can seriously affect the consistency of the explanation. CfExplainer mines the weighted class association rules through the synthetic instances, and when the instances to be explained are predicted to be defective, the rules with the highest weighted support in the set of defective rules are selected to explain the prediction results of the global model, and conversely, the rule with the highest weighted support in the set of non-defective rules is selected to explain the prediction of the global model. Therefore, we studied the consistency of the explanations, measured by the $\%Consistency$ metric, whether CfExplainer can efficiently generate explanations that are consistent with the introduced defects.

Analysis of $\%Consistency$ Results: The explanations generated by CfExplainer are improved by 2.5 %–5.9 % and 9.3 %–42.6 % in $\%Consistency$ metrics compared to PyExplainer and LIME.

On the OpenStack dataset, the median $\%Consistency$ of CfExplainer was 0.806. Compared to PyExplainer and LIME, CfExplainer yielded improvement by 3.6 %–3.9 % and 24.5 %–24.9 %, respectively. On the Qt dataset, the median $\%Consistency$ of CfExplainer was 0.784. Compared to PyExplainer and LIME, CfExplainer yielded improvement by 2.5 %–5.9 % and 9.3 %–42.6 %, respectively. The experimental results are shown in Fig. 9. Furthermore, the p-values of CfExplainer versus PyExplainer and LIME on the $\%Consistency$ metric are 0.109, 0.008, respectively, and Cliff's $|\delta|$ values are all 1, this suggests that the CfExplainer approach has large effect sizes and practical significance

The experimental results show that the consistency of our CfExplainer explanations is better than PyExplainer and LIME, which suggests that CfExplainer is able to generate higher consistency explanations. PyExplainer and LIME use cross-mutation and random perturbation to generate synthetic instances, respectively, which have randomness, and both may generate synthetic instances (anomalous instances) that deviate from the real-world laws, i.e., the synthetic instances lack truthfulness, and the local model is constructed based on synthetic instances, which in turn leads to the knowledge learnt by the local model is not in line with the real world, i.e., the mined rules lack reliability. However, CfExplainer generates synthetic instances based on counterfactuals, sets artificial constraints when perturbing features, and generates synthetic instances with truthfulness, which enables CfExplainer's local model to learn knowledge that conforms to the real-world laws, i.e., it can correctly explain the real reasons for introducing defects. In addition, CfExplainer constructs local models by mining

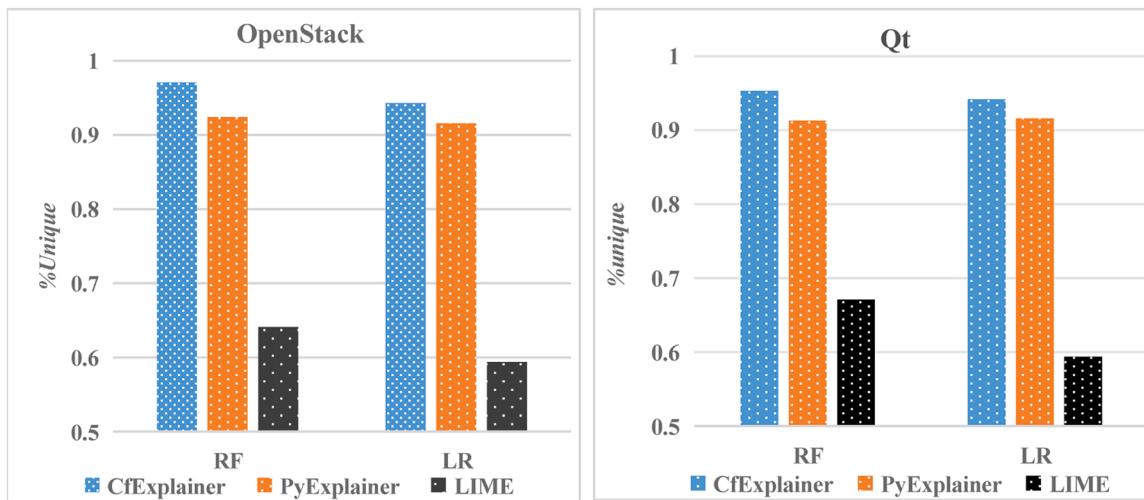


Fig. 8. Explanation of %Unique.

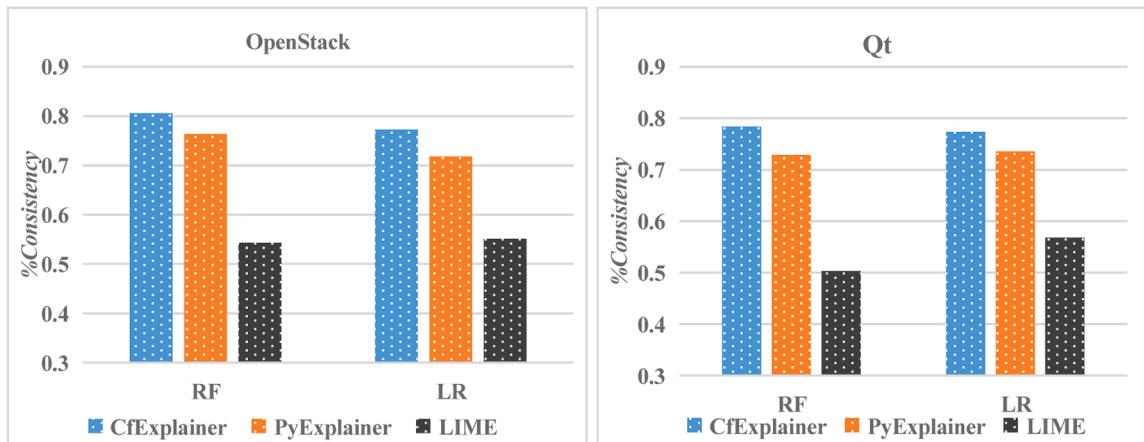


Fig. 9. Explanation of %Consistency.

weighted class association rules, considers the contribution of features to defects, uses weighted support measure to measure the importance of rules, and optimises the process of generating, ranking, pruning, and prediction of rules, and does filtering of redundant and conflicting rules, which improves the quality of the rule set, and thus improves the reliability of explanations generated by the local models. Therefore, the explanations generated by CfExplainer have higher consistency than LIME and PyExplainer.

5.4. Parameter setting

To reduce the influence of the parameters on the experiments, we explored the influence of the minimum support (minsupport) threshold on the fitness of the weighted class association rules. Minsupport thresholds for class association rules on different datasets tend to be different, while different minsupport thresholds affect the number of rules generated; the smaller the minsupport, the greater the number of rules generated, and vice versa. This in turn leads to different rule ranking results, which ultimately affects the quality of rule mining and the accuracy of prediction. Therefore, we studied the effect of different minsupport values on the fitness of the weighted class association rules. To set the appropriate minsupport threshold, we determined the best minsupport based on a grid search.

As shown in Fig. 10, when minsupport is 0.05–0.15, the fitness of the weighted class association rules increases as minsupport increases.

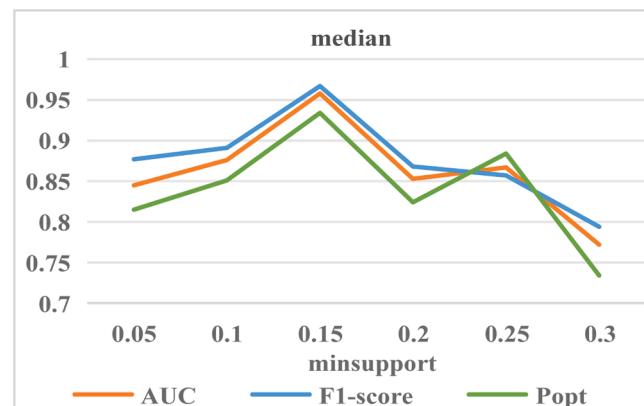


Fig. 10. Effect of minsupport on goodness of fit.

When minsupport is 0.15–0.2, the fitness of the weighted class association rules decreases. When minsupport is 0.2–0.25, the fitness slowly increases and then decreases again after 0.25.

The prediction performance of the model has an impact on the effectiveness of the CfExplainer method to a certain extent. Therefore, we studied the impact of just-in-time defect prediction model parameters on the interpretation performance of CfExplainer by selecting

representative parameters of random forests and logistic regression, i.e., n_estimator and the regularization parameter penalty (Contreras et al., 2021; Li et al., 2022).

As shown in Fig. 11, the impact of n_estimators and penalty on method performance is demonstrated. On the %Unique and %Consistency metrics, %Unique and %Consistency increased when n_estimators increased from 50 to 150, but then %Unique and %Consistency gradually decreased as n_estimators continued to increase.

In logistic regression models, the regularization parameter penalty has a significant impact on the performance of the method. L1 regularization, L2 regularization, and Elasticnet are three commonly used regularization methods that prevent overfitting by adding penalties to the loss function associated with the model parameters. When the logistic regression model uses L1 regularization, the %Unique and %Consistency metrics are 0.926 and 0.786, respectively. When the logistic regression model uses L2 regularization, the %Unique and %Consistency metrics are 0.854 and 0.741, respectively. When the logistic regression model uses Elasticnet regularization, the %Unique and %Consistency metrics are 0.883 and 0.762, respectively.

6. Threats to validity

6.1. Internal threats

Minsupport must be set manually when mining weighted association rules, which rely on prior knowledge. Minsupport thresholds for class association rules on different datasets usually differ, and setting min-support usually affects the number of rules generated; the smaller min-support, the more rules are generated, and vice versa. This leads to different results in ranking rules and ultimately affects the rule-mining quality and prediction accuracy. To reduce the randomness of the conclusions, the number of synthetic instances generated by CfExplainer is set to 2000.

6.2. External threats

Because of the limited datasets that meet traceability, rapid iteration, and code review policies, consistent with other approaches in this field, we only selected OpenStack and Qt datasets. In future research, to further validate CfExplainer's generalisability, we will validate it on more standard datasets.

7. Conclusion

We herein proposed CfExplainer, a model-agnostic approach based on local rules. CfExplainer uses counterfactuals to generate synthetic

instances, mines weighted class association rules based on synthetic instances, optimises the processes of generating, ranking, pruning, and predicting class association rules, and uses the highest-priority rules to explain the predictions of the model. CfExplainer explanations outperformed those of other methods in terms of effectiveness and reliability. The explanations generated by CfExplainer improved the transparency of JIT defect prediction models. CfExplainer can thus help developers analyse the causes of defects, provide developers with guidance on defect fixes, and reduce the risks of introducing defects.

Some current model interpretation methods are designed for in-project software defect prediction based on successive versions of software, but this requires multiple successive versions of software for training and evaluation. Therefore, they are not applicable to JIT defect prediction models, while CfExplainer is performed in an just-in-time defect prediction scenario, which is suitable for open scenarios with rapid software iteration, and the generated explanations can improve the transparency of just-in-time defect prediction models, help developers analyze the causes of defects, and provide developers with guidance on defect fixing to reduce the risk of introducing defects. Developers understand the model's prediction results through the explanations generated by CfExplainer. For example, CfExplainer provides rule interpretations for instances to be interpreted (e.g., 29<churn<100 and ns>1⇒Defect), which are predicted to be defective because churn is greater than 29 and <90, and ns is greater than 1. The rule interpretations help testers to focus on the aspects that are most relevant to the risk of defects, rather than focusing on other, less important aspects, and thus improve the efficiency of software fixes. In order to improve the usefulness of CfExplainer in explaining SDP models, more dimensional features will be mined through feature engineering techniques in the future to enhance the model's explaining capability.

The minsupport settings depend on prior knowledge. Therefore, to reduce the potential effects of prior knowledge, we will consider adaptive mminsupport in future studies. Furthermore, the quality of synthetic instances generated based on counterfactuals has significant effects on the effectiveness and reliability of the explanations. Therefore, we will further study optimisation methods for the counterfactual generation of synthetic instances.

Author contributions

All authors are contributed equally.

CRediT authorship contribution statement

Fengyu Yang: Supervision, Funding acquisition. **Guangdong Zeng:** Writing – original draft, Validation, Software, Methodology. **Fa Zhong:**

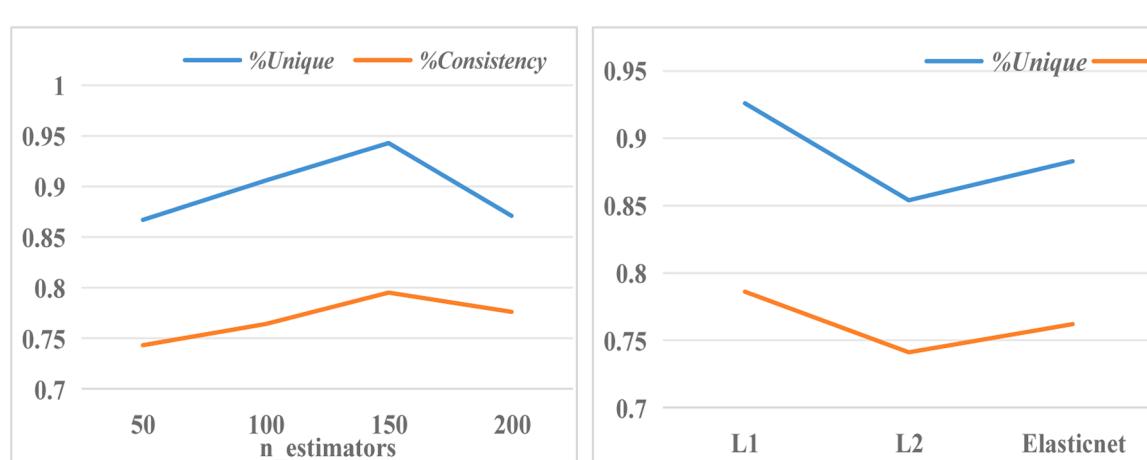


Fig. 11. Effect of model classification parameters on CfExplainer.



Guangdong Zeng is currently pursuing the M.S. degree with the School of Software, Nanchang Hangkong University. His research interests include software defect prediction and association rules.



Peng Xiao received the Ph.D. degree from Beihang University, in 2018. He is currently a Lecturer and a Master Supervisor with the School of Software, Nanchang Hangkong University. His research interest include software testing and software reliability.



Fa Zhong is currently pursuing the M.S. degree with the School of Software, Nanchang Hangkong University. His research interests include software testing and software defect prediction.



Wei Zheng received the Ph.D. degree from Xidian University, in 2010. He is currently a Professor and a Master Supervisor with the School of Software, Nanchang Hangkong University. His research interests include complex networks and software reliability.