

IMPROVING SOFTWARE FAULT PREDICTION THROUGH CROSS-PROJECT ANALYSIS A FOCUS ON IMBALANCED DATA AND GENERALIZATION

^{#1}Subbalakshamma Thaadi, M.Tech- Student Dept. of CSE-SE

^{#2}P.Viswanatha Reddy, Associate Professor Department of CSE

^{#3}Dr. V. Hemasree, Professor & HoD Department of AI&DS

Viswam Engineering College Madanapalle, AP

Abstract: This paper delves into the challenges of generalizing models and dealing with contradicting evidence. Additionally, it delves into the potential for enhancing software failure prediction by integrating several research endeavors. Conventional methods of failure prediction could not be highly task-specific due to the fact that not all tasks had access to the same data. To overcome these challenges and achieve better prediction accuracy, you can employ feature selection techniques, data resampling tactics, and machine learning procedures. The project involves exploring the usage of various datasets and enhancing model training to expedite problem detection and ensure that solutions are compatible with different software configurations. Software quality assurance methods can be improved and made more adaptable as a direct consequence of the findings.

Keywords: Software Fault Prediction, Cross-Project Analysis, Imbalanced Data, Machine Learning, Generalization, Feature Selection, Data Resampling, Software Quality Assurance.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author(s) and the source are properly cited.

1. Introduction

Software fault prediction (SFP) is an essential component of software quality assurance since it identifies problematic modules at an early stage of development. The fact that traditional methods of defect prediction need data that is unique to each project makes them useless in this context. One potential answer to this issue is cross-project software fault prediction (CP-SFP), which reduces the requirement for project-specific data while improving the model's flexibility by integrating data from several projects. An issue with CP-SFP and a potential source of inaccurate calculations are data imbalances, such as a much lower number of faulty modules compared to non-defective ones.

If you want to get better at predicting software failures across projects using analytics, you have to fix data discrepancies. When asked to predict units that aren't defective, machine learning systems actually do worse because fault datasets aren't varied enough. Some solutions to this issue include cost-sensitive learning, hybrid techniques, and resampling. Projects can become more interdependent with the use of feature

selection and transfer learning approaches, which increase model generalizability and provide accurate predictions across various software environments.

The intricacy of projects, development processes, and coding standards are all susceptible to change, which impacts the accuracy of predictions, making generalization in CP-SFP a challenging effort. Use of deep learning, ensemble methods, and domain adaptation techniques might lead to a more dependable CP-SFP model. The incorporation of explainable AI algorithms has given software developers more confidence when using prediction models. It is well-known that CP-SFP can improve generalization and fix unequal data, which in turn reduces software defects and increases program dependability.

2. Literature Review

Chen, H., Yang, L., & Wang, A. (2024): The primary objective of this endeavor is to enhance CPDP (cross-project software defect prediction) by the utilization of federated meta-learning. Issues with privacy and

inconsistent initiatives plague traditional CPDP methods. Federated learning allows many businesses to train models jointly without exposing any source data. The research employs meta-learning, which allows for rapid project adaptation. In order to reduce processing costs, increase the accuracy of failure prediction, and protect data privacy, the proposed method integrates multiple methodologies.

Saeed, M. S. (2024): This page gives a thorough introduction to the various ensemble learning mechanisms employed by CPDP. To get the best possible outcomes, a technique called ensemble learning combines multiple prediction models. The paper compares and contrasts several methods to demonstrate how they can improve generalization across projects. These methods include bagging, boosting, stacking, and random forests. The author delves into topics including feature variability, data imbalance, and domain adaptability. In addition, they suggest ways that CPDP can be enhanced in the future by incorporating ensemble techniques.

Zhang, Y., & Yang, Y. (2024): The purpose of this research is to improve CPDP transfer learning by strengthening methods of information transfer. Using transfer learning, models built for one project can be used to another, even though there's a chance that performance would suffer because of domain variances. The research found that projects can be linked using domain adaption methodologies such as adversarial learning and feature transformation. The testing findings demonstrate that the models' increased flexibility results in higher rates of defect detection.

Shi, H., Ai, J., Liu, J., & Xu, J. (2023): The issue of CPDP software failure datasets that are both chaotic and unequal is the primary focus of this research. Missing or incorrect data, along with an uneven distribution of defective and non-defective examples across several datasets, can hinder predictions. Noise filtering, cost-sensitive learning, and class balancing methods like SMOTE (Synthetic Minority Over-sampling Technique) are recommended for data preprocessing. Defect detection keeps its greater accuracy even when dealing with widely variable datasets, as this research shows.

Zheng, Y., & Zhang, H. (2023): The primary emphasis of this inquiry is on the management of different data sets in CPDP. The dispersed nature of fault records across projects makes it impossible to generalize using existing models. In their work, the authors introduce a multi-source learning method that enables simultaneous examination of numerous

distributions. Their method improves prediction accuracy on various software datasets by implementing adjustments to dynamic feature representation that are specific to each project.

Khleel, M., & Nehéz, K. (2023): Using oversampling approaches and bidirectional Long Short-Term Memory (BiLSTM) networks, this research aims to improve CPDP. Improving defect prediction, BiLSTM networks—a more advanced form of recurrent neural network (RNN)—find long-range correlations in defect data. The research used oversampling approaches, such as SMOTE, to make sure the models were more equitable because the sample distribution was uneven across the classes. In comparison to more conventional machine learning methods, the proposed method achieves better results when applied to sequential defect data.

Kumar, K. B., Kanchanapally, S. P., & Thota, M. K. (2023): This research introduces a centroid-based project selection method to enhance CPDP's precision. One of the most difficult aspects of CPDP is deciding which efforts to use as training data. The authors detail a clustering method that uses similarity between projects to form training sets. This enhances the accuracy of predictions in a wide variety of software applications, increases data representativeness, and more evenly distributes classes.

Nevendra, M., & Singh, P. (2022): The primary objectives of this research are to determine the data balance for CPDP and to recommend software metrics. When noise is added into some software metrics, they become unsuitable for defect prediction. The paper proposes a feature selection technique to keep only the most significant metrics while reducing processing expenses. The impartiality of failure predictions is guaranteed by use of data balancing procedures. The results demonstrate a considerable enhancement in the precision of the predictions.

Majd, A., & Salimi, M. (2022): This research uses a hybrid approach, integrating ensemble learning with Generative Adversarial Networks (GANs), to predict imbalanced defects. A GAN's ability to diversify training samples is enhanced by its use of fabricated error data. Researchers found that CPDP performed better when they combined GAN-generated data with popular ensemble methods like boosting and random forests. The results demonstrate that this mixed-method strategy improves responsiveness to new initiatives while decreasing class imbalance.

Goel, L., Sharma, M., Khatri, S. K., & Damodaran, D. (2021): This empirical research investigates multiple

data collection methods with the goal of resolving the CPDP class gap.

Several resampling procedures are discussed, with SMOTE serving as an example. These tactics include hybrid approaches, under-sampling, and over-sampling. Examining the effects of fair training datasets on machine learning models makes it much easier to detect mistakes in real-world applications.

Wu, J., Wu, Y., Niu, N., & Zhou, M. (2021): MHCPDP, which stands for "multi-source heterogeneous CPDP," is one of the systems employed in the research. To extract useful details regarding failures across different projects, multi-source transfer learning and autoencoders are employed. The approach enables the detection of software system issues by integrating data from many sources while preserving feature integrity.

Sun, Z., Li, J., Sun, H., & He, L. (2021): Instead of CPDP, the authors suggest CFPS, an acronym for "Collaborative Filtering-based Project Selection." Collaborative filtering, a technique often employed in recommendation systems, is employed to ascertain the most effective training programs. The use of only the most relevant training data increases the probability of accurate predictions using this selection strategy.

Jiang, K., Wang, A., Zhang, Y., & Wu, H. (2020): Resolving notable class inequalities through the application of transfer learning in CPDP is the primary objective of this work. Traditional methods fail miserably when faced with highly skewed defect data while trying to identify important trends. Finding a happy medium between the defect and non-defect classes is the goal of the research, which is why it suggests a domain adaptation method to improve defect prediction.

Li, K., Xiang, Z., Chen, T., & Tan, K. C. (2020); A two-level programming method for the automatic identification of CPDP models is introduced in this research as BiLO-CPDP. Selecting features, hyperparameters, and learning algorithms becomes much easier using the framework's two-step method. The research found that across different datasets, CPDP performed substantially better with automatic model adjustment.

Saifudin, A., Hendric, S. W. H. L., Soewito, B., Gaol, F. L., Abdurachman, E., & Heryadi, Y. (2019): To resolve the issue of class incompatibility in CPDP, this research employs Ensemble SMOTE. The authors utilize several resampling strategies to guarantee a training set that is more representative of the population at large. In addition to drastically

decreasing prediction bias, their approach makes the models more suitable for future applications.

3. Existing System

Software defect prediction is an essential component of software quality assurance because it identifies potential issues prior to their implementation. The traditional approach to training machine learning models typically employs data that is already included in the project. The initial stage is to train these models using the fault data from the project. These data administration strategies may be advantageous for startups and smaller businesses, as they have fewer fault records. Typically, supervised learning techniques such as decision trees, support vector machines, and neural networks are employed to classify modules with a high defect risk. Nevertheless, the efficacy of these methods is contingent upon the quantity and quality of the classified data.

One approach to managing data concerns is to anticipate flaws across initiatives (CPDP). This method is predicated on statistical data that illustrates errors in comparable assignments. Domain adaptation and transfer learning techniques facilitate the alignment of various project datasets to enhance the accuracy of predictions. However, the potential for project sizes, defect numbers, and coding patterns to fluctuate raises concerns. The biased models that result from the imbalanced data make it more challenging to generate precise predictions when there are fewer defective modules than non-defective ones. These procedures prioritize all classes that comprise the preponderance.

Generalization in CPDP remains challenging due to the fact that models trained on one project frequently perform poorly on another. Setting up models with extensive applicability is a challenging task in the dynamic realm of software development. Ensemble learning, feature mapping, and instance selection are among the techniques employed to enhance generalizability. Nevertheless, the issue of achieving high-quality predictive performance across various initiatives persists.

- **Data Imbalance Issues** – A limited number of defective modules in software projects presents a greater challenge for machine learning algorithms to identify significant trends. Erroneous predictions may result from falsified or over- or under-sampled data.
- **Domain Mismatch** – The utilization of transfer learning in new projects is not as prevalent due to variations in defect types, project designs, and

- coding methodologies. This is due to the models' inability to generalize sufficiently.
- **Feature Discrepancies** – It is challenging to standardize data across projects due to the variability of software metrics and feature distributions. Predictions are rendered inaccurate due to the model's decreased consistency.
- **Computational Complexity** – Complex approaches, such as deep learning, ensemble learning, and domain adaptation, necessitate substantial computing capacity and are incapable of accurately predicting the failure of large projects in real time.
- **Evaluation and Validation Challenges** – It is difficult to compare CPDP models due to the absence of defined datasets or evaluation standards. Performance comparisons may not be as reliable as they could be due to discrepancies in dataset labeling and preparation procedures.

4. Proposed System

The suggested method—a cross-project analytical framework—improves software failure prediction by addressing data imbalance and enhancing generalization. Our solution ensures that feature regions are consistent across all projects by employing advanced transfer learning techniques, rather than relying solely on defect data from a single project. The strategy implements domain adaptation techniques to account for variations in software measurements and coding practices. As a result, models are capable of learning from a wide range of software development duties. In order to enhance the accuracy of error detection, hybrid machine learning models combine traditional classifiers with deep learning.

To address the issue of uneven data, the proposed method integrates sophisticated resampling techniques, such as SMOTE, with cost-sensitive learning. Class bias is mitigated by these methods, which train on a set that includes both functional and malfunctioning modules. An ensemble learning technique is employed to enhance the stability of the model and ensure that the defect forecast performance is superior across a variety of projects. This technique involves the integration of multiple classifiers. The system prioritizes software metrics by employing feature selection techniques. Consequently, forecasts are more precise and calculations are less costly.

Meta-learning and adaptive feature transformation enable enhanced generalizability and real-time model adaptability to novel circumstances. The system

dynamically updates its feature representation and decision-making capabilities by leveraging new project data, rather than relying on static model training. It consistently achieves exceptional accuracy across various software setups by significantly relying on benchmark datasets in its validation approach. The emergence of AI that is explicable has facilitated the availability of a solution that was previously incomprehensible. To identify errors prior to their development into significant problems, software engineers can implement these strategies to notify developers of potential problems.

- **Better Handling of Imbalanced Data** – Advanced resampling techniques and cost-sensitive learning are employed to enhance the accuracy of fault predictions and accurately represent the minority class, which is composed of defective modules.
- **Enhanced Generalization Across Projects** – Transfer learning and domain adaptation can enhance the model's ability to learn from a variety of sources. This guarantees compatibility with an extensive array of applications.
- **Improved Fault Detection Accuracy** – The hybrid method enhances prediction accuracy and decreases the false positive/negative rate by integrating traditional classifiers with deep learning and ensemble techniques.
- **Reduced Feature Discrepancies** – Adaptive feature transformation and feature selection facilitate the standardization of raw data, the elimination of discrepancies between projects, and the assurance that the model operates as intended indefinitely.
- **Scalability and Continuous Learning** – The system is more easily able to adapt to shifting software production processes and make improvements when a consistent stream of new data is maintained.

5. Implementation

Modules:

Service Provider

The service provider must possess an active, password-protected account in order to utilize this feature. Upon completion of the check-in procedure, he will have the ability to select between training and evaluation activities and access data sets. You were anticipating that these files would be located in the near future. A bar chart that illustrates the efficacy of various testing and learning methodologies and the rates of various software fault forecasts should be

incorporated into your forthcoming visit. The overall accuracy of the learned and evaluated data, the types and quantities of software fault predictions, all user profiles, and much more can be accessed remotely.

View and Authorize Users

The manager can now access a comprehensive list of all users as a result of this module. Using a user's name, email address, and physical address, administrators can assign or revoke rights.

Remote User

Before proceeding, verify that the application form has been completed in its entirety. The database will contain the individual's information after they have completed the registration process. Upon completion of the registration procedure, he will be prompted to provide his login and password. After the user's identity has been verified, they will be able to view their photo, evaluate the likelihood of a software error, and proceed with the registration and authentication process.

6. Results

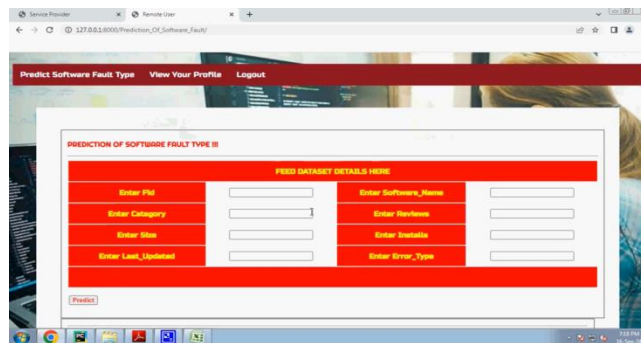


Figure 1: Prediction of Software Fault Type

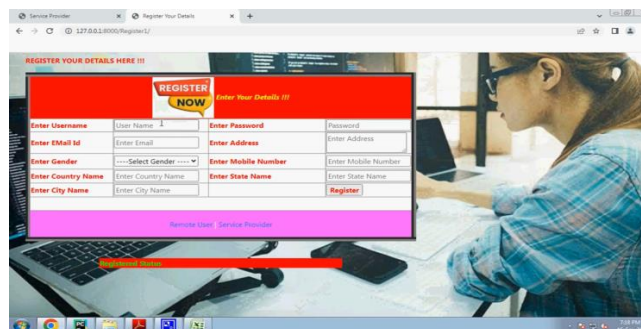


Figure 2: User Registration Page



Figure 3: Software Fault Prediction Type Ratio



PLD	Software Name	Category	Reviews	Size	Installs	Last Updated	Error Type	Prediction
180.163.251.231-10.42.8.211-90-56942-0	TextNow - free text + calls	SOCIAL	441189	Varies with device	10,000,000+	03-Aug-18	No Exception	Fault Not Found
10.42.8.104-64.71.142.124-39935-86-6	The Messenger App	SOCIAL	4919	2.8M	1,000,000+	03-Aug-18	No Exception	Fault Not Found
10.42.8.42-119.146.74.94-40654-89-6	Telegram X	SOCIAL	70616	Varies with device	5,000,000+	27-Jul-18	runtime exception	Fault Found
10.42.8.211-10.42.8.1-36140-58-17	Who Viewed My Facebook Profile - Stalkers Visitors	SOCIAL	271445	0.0M	5,000,000+	24-Jan-18	No Exception	Fault Not Found

Figure 4: Software Fault Prediction Type Details Page

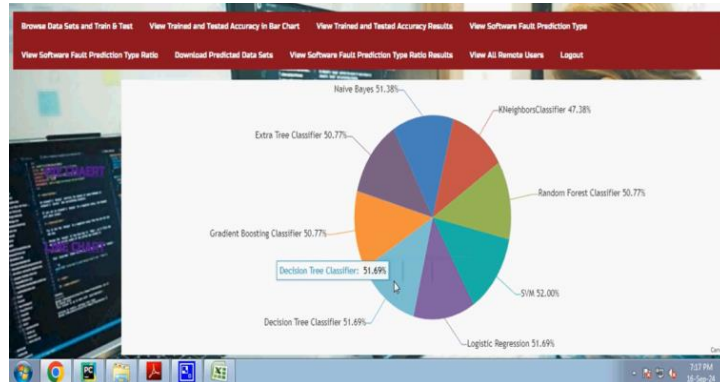


Figure 5: Software Fault Prediction in Pie Chart

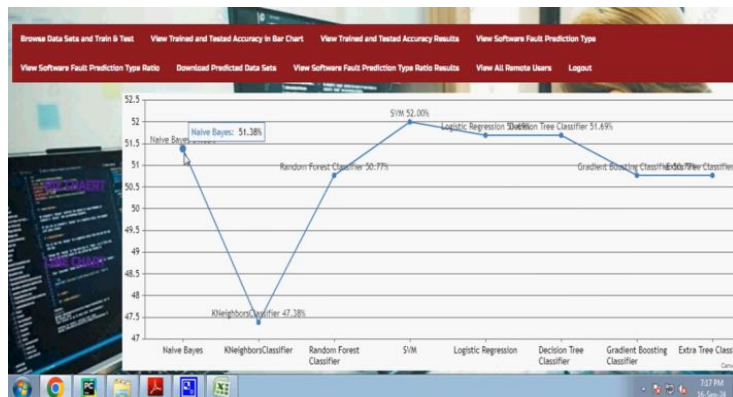


Figure 6: Software Fault Prediction in Line Chart



Figure 7: Software Fault Prediction in Bar Chart



Figure 8: Datasets Trained and Tested Results

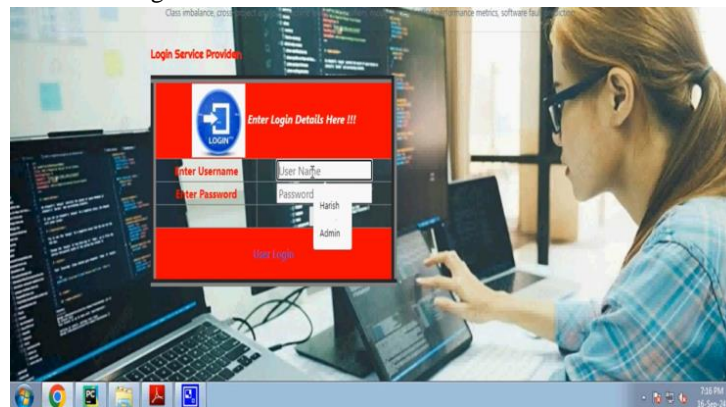


Figure 9: Service Provider Login Page



Figure 10: User Login Page

7. Conclusion

Cross-project analysis may prove advantageous for the prediction of software quality and, as a result, software defects. This is particularly beneficial when confronted with contradictory data or issues that involve generalization. Cutting-edge machine learning techniques, including domain adaptation, transfer learning, and ensemble learning, are employed in the proposed methodology to resolve the challenges associated with existing within-project models. By combining cost-sensitive learning with resampling

techniques, an equitable representation of error-prone modules is ensured. In this manner, bias is diminished and prediction accuracy is enhanced.

A model's generalizability can be improved by employing adaptive feature transformation and meta-learning. It is capable of predicting the software issues of projects with varying fault distributions and code standards. The system can adapt to changing software development environments by analyzing and enhancing its previous errors. To assist software

engineers in making informed decisions in advance, this solution employs explainable AI to provide essential information regarding critical failure predictors.

Our method significantly improves the prediction of software faults by addressing feature discrepancies,

classification biases, and data shortages. Predicting cross-project issues is an indispensable methodology for guaranteeing software quality. Increasing the precision and efficacy of fault identification frequently results in more stable software solutions and reduced maintenance expenses.

References

1. Chen, H., Yang, L., & Wang, A. (2024). Efficient cross-project software defect prediction based on federated meta-learning. *Electronics*, 13(6), 1105.
2. Saeed, M. S. (2024). Cross-project software defect prediction using ensemble learning: A comprehensive review. *International Journal of Computational and Innovative Sciences*, 3(2), 34–42.
3. Zhang, Y., & Yang, Y. (2024). Improving transfer learning for software cross-project defect prediction. *Applied Intelligence*, 54, 1234–1250.
4. Shi, H., Ai, J., Liu, J., & Xu, J. (2023). Improving software defect prediction in noisy imbalanced datasets. *Applied Sciences*, 13(20), 10466.
5. Zheng, Y., & Zhang, H. (2023). Cross-project defect prediction considering multiple data distribution simultaneously. *Symmetry*, 14(2), 401.
6. Khleel, M., & Nehéz, K. (2023). Software defect prediction using a bidirectional LSTM network combined with oversampling techniques. *Cluster Computing*, 26, 1059–1075.
7. Kumar, K. B., Kanchanapally, S. P., & Thota, M. K. (2023). Class imbalance reduction and centroid-based relevant project selection for cross-project defect prediction. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(6s), 293–302.
8. Nevendra, M., & Singh, P. (2022). Cross-project defect prediction with metrics selection and balancing approach. *Applied Computer Systems*, 27(2), 137–148.
9. Majd, A., & Salimi, M. (2022). Leveraging ensemble learning with generative adversarial networks for imbalanced software defects prediction. *Applied Sciences*, 12(24), 13319.
10. Goel, L., Sharma, M., Khatri, S. K., & Damodaran, D. (2021). Cross-project defect prediction using data sampling for class imbalance learning: An empirical research. *International Journal of Parallel, Emergent and Distributed Systems*, 36(2), 130–143.
11. Wu, J., Wu, Y., Niu, N., & Zhou, M. (2021). MHCPDP: Multi-source heterogeneous cross-project defect prediction via multi-source transfer learning and autoencoder. *Software Quality Journal*, 29, 405–430.
12. Sun, Z., Li, J., Sun, H., & He, L. (2021). CFPS: Collaborative filtering based source projects selection for cross-project defect prediction. *Applied Soft Computing*, 99, 106940.
13. Jiang, K., Wang, A., Zhang, Y., & Wu, H. (2020). Heterogeneous defect prediction based on transfer learning to handle extreme imbalance. *Applied Sciences*, 10(1), 396.
14. Li, K., Xiang, Z., Chen, T., & Tan, K. C. (2020). BiLO-CPDP: Bi-level programming for automated model discovery in cross-project defect prediction. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering* (pp. 573–584).