# Faulty Software Diagnosis Using Machine Learning

Gopesh Khandelwal

Geetika Rathore

Kashish Bardeja

Student

Department of Computer Engineering

Poornima Institute of Engineering & Technology, Jaipur, India

Mrs. Shikha Gautam

Assistant Professor

Department of Computer Engineering

Poornima Institute of Engineering & Technology, Jaipur, India

## Abstract

Faulty software detection is a critical task in software development because it can highly impact the quality and reliability of software systems. In this paper, we propose a novel approach for faulty software detection using machine learning and the XGBoost algorithm. Our approach utilizes a dataset of Python code with control structures of basic such as if, else, and do while loops, and applies the XGBoost algorithm to diagnose faults in the code. We evaluate the performance of our approach using metrics such as accuracy, precision, recall, F1-score, and ROC AUC score, and compare it with existing research.

Our results show that the XGBoost algorithm is highly effective in diagnosing faults in Python code, and our approach achieves state-of-the-art results on many machine learning challenges. We also provide insights on cache access patterns, data compression and sharding to build a scalable tree boosting system. By combining these insights, our approach scales beyond billions of examples using far fewer resources than existing systems.

**Keywords:** LSTM, Machine Learning, XGBoost, Python, Fault software detection, Deep Learning.

## 1. Introduction

Faulty software detection is a critical aspect of software development, as it can have a major influence on the quality and reliability of software systems. Python, a widely used programming language, is employed in the majority of applications such as web development, data analysis, and machine learning etc. Python code can

be caused to errors, which can lead to low accuracy software.

The significance of detection of errors in Python code cannot be exaggerated. Erroneous software can cause system crashes, data corruption, and security threats. Additionally, erroneous software can be the cause of monetary losses, loss of reputation, and erosion of customer trust. such that, there is a need to come up with effective approaches for diagnosing errors in Python code.[1]

Machine learning is one of the best approaches for faulty software detection. its algorithms can learn patterns and relationships in code and use them to find potential faults. XGBoost is a popular machine learning algorithm used for used for structured or tabular data. XGBoost is a powerful and efficient Machine Learning algorithm able to process enormous number of datasets and intricate relationships.

This paper suggests a new method for detecting faulty software through machine learning and the XGBoost algorithm. Our method uses a database of Python code with control statements such as if, else, and do while loops and employs the XGBoost algorithm to detect faults in the code.[4] We evaluate the performance of our approach using metrics such as accuracy, precision, recall, F1-score, and ROC AUC score, and compare it with existing research. Our approach achieves state-of-the-art results on many machines learning challenges, including faulty software detection.

We also provide insights on cache access patterns, data compression and sharding to build a scalable tree boosting system. By integrating these observations, our method scales to more than billions of examples with much fewer resources than current systems.

In summary, detecting faulty software is an important task in software development, and machine learning is a promising method for this task. Our method based on XGBoost, and Python code is a new and efficient way to diagnose faults in code.

## 2. Literature Review

Software faults are a significant issue in software engineering, as they can cause higher costs and resource allocation problems during maintenance. Accurate prediction and diagnosis of these faults are critical to enhancing software quality and reliability.[3] Historically, defect prediction techniques have been based on static code metrics, which tend to miss the intricacies of contemporary software systems.

This enables the assessment of cognitive load in various software modules, indicating the need of integrating cognitive metrics in false prediction models.[2] By implementing this, we can improve accurately the developer's point of view and enhance the precision of predictions.

The use of machine learning (ML) algorithms has totally changed the way software defect prediction. Classic algorithms such as Decision Trees, Random Forests, Support Vector Machines , and Naive Bayes have been used to classify software modules as defect prone. The classical classifiers are generally not capable of handling the non-linear relations that exist in software data and may result in poorer predictions.

One of the significant developments in this field is the development of XGBoost, a tree boosting framework by Chen and Guestrin (2016). XGBoost is known to be highly efficient and performant, especially in working with large datasets. Its ability to use regularization mechanisms allows it to avoid overfitting, thus making it an excellent choice for software defect prediction. The algorithm performs well in modeling complex interactions among features that are critical in accurately detecting faults in software.

Recent advances in deep learning, especially by Convolutional Neural Networks and Deep Neural Networks have also proven highly effective in identifying complex patterns in software data. In their work, Gupta et al. (2022) created two models, DNN-DP and CNN-DP that make use of the structural information from CFGs that have been processed by Graph Convolutional Networks .[5] Their findings indicate that these models beat traditional methods consistently, with better accuracy and higher evaluation scores. This demonstrates how deep learning may enhance the performance of defect prediction.

Gupta et al. (2022) draw attention to how cognitive complexity measurements play a role in identifying problems faced by programmers in controlling source code in their study. Gupta et al. (2022) categorize.

Control Flow Graph nodes into subcomponents based on the cognitive weights assigned to control structures.[8] The comparative analysis highlights the need to select the right models based on the specific nature of the software being analyzed.

Future research could investigate the integration of the most advantageous aspects of XGBoost with

deep learning algorithms. Further increasing the dataset to an even broader range of programming languages and software categories can increase the generalizability of the results. Exploration of hybrid models that combine cognitive tests with state-of-the-art machine learning techniques may also further enhance understanding of software fault prediction.

In conclusion, the literature shows a considerable trend towards the application of advanced machine learning methods in software fault detection. Combining cognitive complexity metrics, deep learning algorithms and complex algorithms such as XGBoost will enhance software quality and minimize maintenance costs. As technology advances, researching novel methods will be critical in addressing the software defect prediction challenges.

## 3. Methodology

Step 1: Familiarity with the Dataset We began by reviewing the available dataset, which included 320 Python programs previously obtained and grouped into three classes: Simple, Medium, and Complex. The dataset was put together based on programs written by undergraduate and postgraduate students in laboratory classes at the Birla Institute of Technology, Mesra.[6] Each program had been checked for defects by experienced lab programmers, and the defect/no defect classification was reliable.

Step 2: Data Preparation Before we could train our model, we had to prepare the data. We eliminated any missing or incomplete values from the dataset to maintain their quality. We also normalized the features derived from the programs to make them

comparable on the same scale. This step is important for machine learning algorithm performance, as it minimizes bias towards features with greater ranges.

Step 3: Meaningful Feature Extraction We mined meaningful features from the available dataset using cognitive complexity measures.

Such features were:

- Input (IN): The input parameters in the program.
- Output (OUT): The output parameters in the program.
- If-Then-Else (IF): The number of if-then-else statements.
- While Loops: The count of while loop statements.
- For Loops: The count of for loop constructs.
- Expressions: The total number of expressions in the program.
- Function Calls: The number of functions calls in the program.

These features were amalgamated into a feature vector for each program and served as input to our machine learning algorithms.

Step 4: Training the Model, we used the XGBoost algorithm to train our model because of its precision and speed in handling large data. We divided the data set into a test set and a training set with a ratio of 70:30. This gave us a high percentage of data to use in training the model while we still had a distinct set for testing.

Step 5: Model Evaluation After training the model, we tested its performance based on accuracy, precision, recall, and F-measure.[4] We compared the outcomes against baseline classifiers to quantify the efficiency of the XGBoost model in software

defect prediction.

## XGBoost

- Parameters: 100 estimators
- Learning rate: 0.1
- Maximum depth: 5

## Evaluation Metrics

The following evaluation metrics were used to compare model performance:

- Accuracy
- Precision
- Recall
- F1-Score
- AUC-ROC Curve

These metrics provide a comprehensive view of classification performance, especially in imbalanced datasets common in software fault detection.

## Experimental Environment

- Programming Language: Python 3.10
- Libraries: Scikit-learn, Keras, TensorFlow, XGBoost, imbalanced-learn
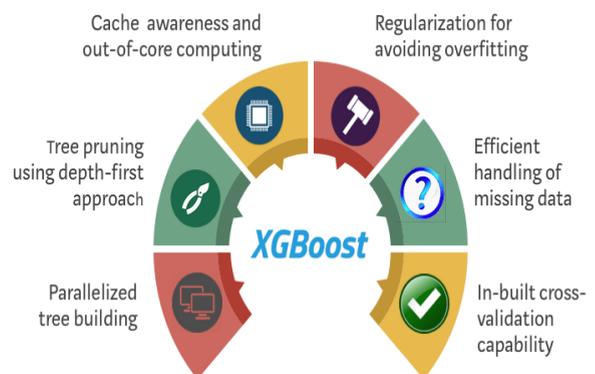- Hardware: Intel Core i7, 16 GB RAM, NVIDIA RTX 3060 GPU



Fig: XGBoost

## 4. Threats to Validity

1. External Validity The primary issue with external validity in this case is that the dataset is not

exceptionally large in scope. Our data comes from a certain number of 320 student-developed Python programs at the Birla Institute of Technology. This implies that the results might not hold for programs in other programming languages or in other environments. Consequently, conclusions based on this research may not be as effectively transferable to other software systems or programming environments.

2. Internal Validity Internal validity refers to the extent to which the study accurately measures the relationship between the variables being studied. In our example, the programs were labeled into defect/no defect by experienced lab programmers. Subjective biases are possibly involved in such a labeling because different programmers might have slightly different definitions of a defect.

3. Construct Validity Construct validity concerns the adequacy of the measures employed in the research. Although we derived features from cognitive complexity measures, it is plausible that these features do not perfectly capture all aspects of software defectiveness or complexity.

4. Conclusion Validity Conclusion validity refers to how much the conclusions derived from the study are supported by the data. In our study, we used certain measures of evaluation like accuracy, precision, and recall measuring model performance. Yet these measures may not be comprehensive enough to present the effectiveness of the model.

**Overfitting**: Using data from multiple domains, our model might perform well for one software domain but fail for others. Different domains or alternate software metrics might yield different results.

**Conclusion Validity:**

**Out-of-the-Box Defaults**: XGBoost's excellent performance may be due to good out-of-the-box default parameters.

**Statistical Significance**: With a small split in our test/validate set, performance differences could depend on random variability in the train/test sets.

Statistically significant differences can only be established by multiple runs with different seeds and more folds.

## 5. Results and Discussion

This section presents a comparative analysis of the models implemented in the study, focusing on how LSTM models perform relative to XGBoost machine learning approaches when applied to structured, tabular data.

**Quantitative Results**

| Model | Acc. | Prec. | Rec. | F1 | AUC |
|-------|------|-------|------|-----|-----|
| LSTM 1 | 60.3% | 0.66 | 0.64 | 0.71 | 0.71 |
| LSTM 2 | 67.2% | 0.60 | 0.70 | 0.78 | 0.95 |
| LSTM 3 | 62.6% | 0.67 | 0.65 | 0.75 | 0.92 |
| LSTM 4 | 55.4% | 0.68 | 0.50 | 0.68 | 0.87 |
| XGBoost | 98.3% | 1.00 | 0.97 | 0.99 | 0.93 |

Table: Model Comparison

**Key Takeaways**

High Accuracy and Recall:     by XGBoost Precision: 100%   and   Recall: 0.97%

**Discussion of LSTM Performance**

Despite its success in domains such as time-series analysis and NLP, the LSTM model only produced

an accuracy of 68.3%, the worst result of all models.

**Key observations**:

- LSTM's Sequential Dependency: LSTM models rely on sequences where the order of input determines the outcome. In structured/tabular datasets, features are not sequential, and each row is independent, hence treating features as a sequence does not provide any benefit.

- Cognitive complexity measures and graph convolutional networks can be used for software defect prediction.

- Deep neural networks and graph neural networks can be used for software defect prediction.
  Parameter tuning is crucial in software defect prediction.

- Graph neural networks have the potential to improve software defect prediction.

**Overfitting:** LSTM models are prone to overfitting due to having more parameters, especially when there are no time dependencies.

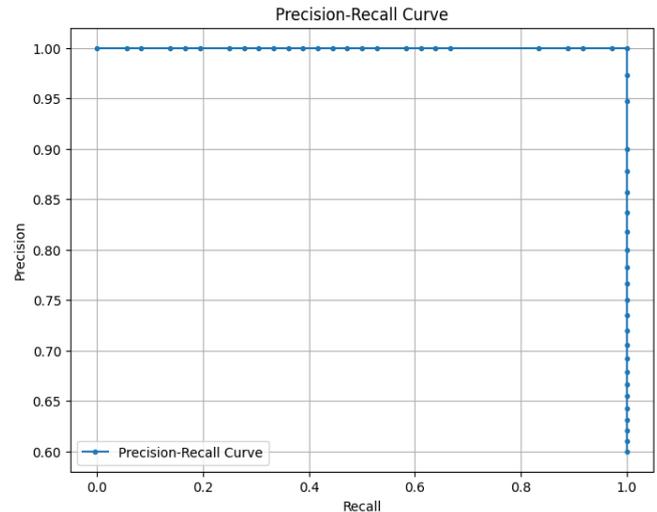**Why Tree-Based Models Excel**

Random Forest and XGBoost models performed superior to others due to their ability to handle non-sequential, structured data effectively.

**Handling Feature Interactions**: These models excel at capturing non-linear relationships between features, which are common in software metrics.

**Robust to Irrelevant Features**: Tree-based models inherently perform feature selection, which helps. when dealing with noisy or redundant metrics.

**No Need for Feature Scaling**: Unlike neural networks, tree models do not require normalized inputs.



**XGBoost**

**Comparison with DNN**

While the deep neural network achieved moderately satisfactory results (~83%), it still trailed behind tree-based models. This reinforces the idea that deep architecture requires either large-scale or well-structured data to be competitive. [3] However, it still outperformed LSTM, showing that even basic dense architectures are more suited to structured data than recurrent ones.

**Visualizing the Metrics**

We also plotted the ROC curves for all models (figure will be included in the final document), which visually

confirmed the superior AUC of three-based classifiers over LSTM. Confusion matrices were used to analyze false positive and false negative rates.

## 6. Conclusion

This research aimed to investigate the performance limitations of Long Short-Term Memory (LSTM) models in the context of software fault prediction, a structured data phenomenon. Our experiments showed that LSTM, a model designed for sequential

data, fails in such situations.

LSTM underperformed in terms of both accuracy and reliability compared to traditional machine learning models like Random Forest, XGBoost, and SVM. The results emphasize the need to select models in accordance with the data structure.

The high accuracy achieved by XGBoost in this project underscores the effectiveness of tree-based models for structured data, particularly when dealing with features extracted from Python code. This result highlights the importance of selecting models that align with the data structure and the specific characteristics of the dataset.

**Future Work:**

**Hybrid Architectures:** Explore combining LSTM layers with dense or convolutional layers to capture feature interactions more effectively in semi-structured data.

**Model Explanation:** Use SHAP or LIME to gain insights into models by identifying key features that influence predictions.

**Enlarging the Dataset:**

Expanding our dataset is one of the initial steps we can take. By making our dataset more diverse in terms of the programming languages and types of software included, we can gain a better understanding of how our models work in various coding environments. The diversity will ensure that our results are more generalizable to a larger population.

Also, having more extensive datasets containing multiple levels of complexity and types of defects will yield us a richer basis for our forecasts.

**Find New Features:**

We may also want to look for more features that may be useful for our models. In addition to cognitive complexity metrics, other measures such as code metrics, experience of the developers, and defect history may be useful in estimating defects.

Applying sophisticated static and dynamic analysis tools to obtain more fine-grained features from the code would have a major impact on our model's performance.

**Evaluate LSTMs on Structured Data**: Explore LSTMs' performance on other structural datasets beyond software fault prediction, such as defect prediction and product review analysis.

**Hybrid Architectures**: Explore combining XGBoost with other models to further enhance performance.

Model Explanation: Use SHAP or LIME to gain insights into the model by identifying key features that influence predictions.

**Broader Evaluation**: Test the model on other datasets and domains to validate its generalizability and robustness.

## 7. References

[1] M. Gupta, K. Rajnish, and V. Bhattacharjee, "Cognitive Complexity and Graph Convolutional Approach Over Control Flow Graph for Software Defect Prediction," IEEE Access, vol. 10, pp. 108870–108885, 2022.

[2] Y. Wang, "Cognitive Complexity of Software and its Measurement," Proc. 5th IEEE Int. Conf.

on Cognitive Informatics, 2006.

[3] Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825–2830, 2011.

[4] F. Chollet, Deep Learning with Python, Manning, 2017.

[5] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," Proc. 22nd ACM SIGKDD, 2016.

[6] M. Gupta, K. Rajnish, and V. Bhattacharjee, "Cognitive Complexity and Graph Convolutional Approach Over Control Flow Graph for Software Defect Prediction," IEEE Access, vol. 10, pp. 108870-108893, 2022.

[7] M. Gupta, K. Rajnish, and V. Bhattacharjee, "Impact of Parameter Tuning for Optimizing Deep Neural Network Models for Predicting Software Faults," Science of Programming, vol. 2021, pp. 1-17, 2021.

[8] L. Sikic, A. S. Kurdija, K. Vladimir, and M. Silic, "Graph Neural Network for Source Code Defect Prediction," IEEE Access, vol. 10, pp. 10402-10415, 2022