

Intelligent Bug Triaging System Using Machine Learning and Random Forest Classification

NELAPATI ANU

PG Scholar, Department of MCA, DNR College, Bhimavaram, Andhra Pradesh

A. Durga Devi

(Assistant Professor), Master of Computer Applications, DNR College, Bhimavaram, Andhra Pradesh

ABSTRACT

Software development projects often involve handling a large number of bug reports submitted by users and testers. Efficiently assigning these bugs to the appropriate developers is a critical task known as bug triaging. Traditional bug triaging methods rely heavily on manual intervention, which can be time-consuming, error-prone, and inefficient, especially in large-scale software systems. This project proposes an Intelligent Bug Triaging System using Machine Learning to automate the process of assigning bug reports to the most suitable developers. The system utilizes Natural Language Processing (NLP) techniques and machine learning algorithms to analyze bug descriptions and predict the developer responsible for resolving the issue. The implementation is carried out using Python with a graphical user interface (GUI) built using Tkinter, making it user-friendly and interactive. The dataset used in this system consists of bug descriptions and their corresponding assigned developers. The textual data is preprocessed and transformed into numerical feature vectors using the Term Frequency-Inverse Document Frequency (TF-IDF) technique. This method effectively captures the importance of words in bug descriptions while reducing the impact of common but less informative terms. A Random Forest Classifier is employed for classification due to its robustness, high accuracy, and ability to handle large datasets with complex relationships. The model is trained on a portion of the dataset and tested on unseen data to evaluate its performance. The system calculates accuracy as a performance metric and visualizes it using a bar graph. The GUI allows users to load datasets, train the model, input new bug descriptions, and obtain predictions in real-time. Once a bug description is entered, the system processes the text, converts it into feature vectors, and predicts the most suitable developer. This approach significantly reduces manual effort, improves efficiency, and enhances accuracy in bug assignment. It is scalable and can be extended with advanced models such as deep learning or integrated into issue-tracking platforms. The proposed system demonstrates the potential of machine learning in automating software engineering tasks and improving workflow efficiency in modern development environments.

Keywords:

Bug Triaging, Machine Learning, Random Forest, TF-IDF, Natural Language Processing, Software Maintenance, Text Classification, Automation

I. INTRODUCTION

In modern software development, managing and resolving bugs efficiently is essential for maintaining software quality and user satisfaction. As software systems grow in complexity, the number of bug reports increases significantly, making manual bug triaging a challenging task. Bug triaging involves analyzing bug reports and assigning them to the appropriate developers based on their expertise. Traditional bug triaging methods rely on project managers or team leads to manually review each bug report and assign it to a developer. This process is time-consuming and prone to human errors, especially when dealing with large-scale projects involving thousands of bug reports. Incorrect assignments can lead to delays in bug resolution and increased development costs. With the advancement of machine learning and natural language processing, automated bug triaging systems have gained significant attention. These systems aim to analyze textual bug descriptions and predict the most suitable developer using classification algorithms. The proposed system leverages TF-IDF for feature extraction and Random Forest for classification. TF-IDF converts textual data into numerical form by measuring the importance of words relative to documents. This helps in capturing meaningful patterns in bug descriptions. Random Forest is an ensemble learning method that combines multiple decision trees to improve prediction accuracy and reduce overfitting. It is particularly effective in handling high-dimensional data, making it suitable for text classification tasks. The system is implemented using Python and provides a graphical user interface using Tkinter. This allows users to easily load datasets, train models, and predict developers for new bug reports. The system also visualizes model performance using accuracy metrics. Overall, this project aims to automate the bug triaging process, reduce manual effort, and improve efficiency in software maintenance. It demonstrates how machine learning can be applied to real-world problems in software engineering.

II. LITERATURE SURVEY (WITH EXISTING METHODS)

Bug triaging has been widely studied in the field of software engineering, with various approaches proposed to improve efficiency and accuracy. Early methods relied on manual assignment based on developer expertise and experience. However, these methods were not scalable and often resulted in incorrect assignments. With the introduction of machine learning, automated bug triaging systems began to emerge. Researchers have explored various text classification techniques, including Naïve Bayes, Support Vector Machines (SVM), and K-Nearest Neighbors (KNN), to predict bug assignments. These methods showed promising results but faced challenges in handling large datasets and complex relationships. Natural Language Processing (NLP) techniques play a crucial role in bug triaging. Feature extraction methods such as Bag-of-Words and TF-IDF are commonly used to convert textual data into numerical form. TF-IDF, in particular, has been widely adopted due to its ability to highlight important words while reducing noise. Ensemble learning methods, such as Random Forest and Gradient Boosting, have gained popularity due to their high accuracy and robustness. Random Forest combines multiple decision trees to improve prediction performance and reduce overfitting. Studies have shown that Random Forest outperforms traditional classifiers in

many text classification tasks. Recent research has also explored deep learning approaches, such as Recurrent Neural Networks (RNNs) and Transformers, for bug triaging. These models can capture contextual information and semantic relationships in text, leading to improved accuracy. However, they require large datasets and significant computational resources. Another important aspect of bug triaging is dataset quality. Public datasets such as Bugzilla and GitHub issue trackers are commonly used for research. Data preprocessing techniques, including stop-word removal, stemming, and normalization, are essential for improving model performance. The proposed system builds upon these existing methods by combining TF-IDF with Random Forest and providing a user-friendly interface. This approach offers a balance between accuracy and computational efficiency, making it suitable for practical applications.

III. EXISTING SYSTEM

The existing bug triaging systems primarily rely on manual processes or basic automated tools. In manual systems, project managers or team leads review bug reports and assign them to developers based on their knowledge and experience. This approach is time-consuming and prone to errors, especially in large-scale projects. Some automated systems use simple keyword matching or rule-based approaches to assign bugs. While these methods reduce manual effort, they lack accuracy and fail to capture the semantic meaning of bug descriptions. Traditional machine learning approaches, such as Naïve Bayes and Support Vector Machines, have been used for bug triaging. However, these models may struggle with high-dimensional text data and may not perform well when the dataset is large or imbalanced. Another limitation of existing systems is the lack of user-friendly interfaces. Many systems are command-line based, making them difficult to use for non-technical users. Additionally, most systems do not provide visualization tools for performance evaluation. Scalability is also a major concern. As the number of bug reports increases, manual and basic automated systems become inefficient. Moreover, incorrect assignments can lead to delays in bug resolution and increased workload for developers. The proposed system addresses these limitations by using advanced machine learning techniques, providing a graphical user interface, and offering real-time predictions with improved accuracy and efficiency.

IV. PROPOSED METHOD

The proposed system is an **Intelligent Bug Triaging System** that automates the process of assigning bug reports to the most suitable developers using machine learning and natural language processing techniques. The system is designed to reduce manual effort, improve accuracy, and enhance the efficiency of software maintenance processes. The system takes bug descriptions as input and predicts the appropriate developer based on historical data. It uses a dataset containing bug descriptions and corresponding developer assignments. The textual data is processed using Term Frequency–Inverse Document Frequency (TF-IDF), which converts unstructured text into numerical feature vectors. A Random Forest classifier is used to train the model. This algorithm is chosen due to its

robustness, ability to handle high-dimensional data, and resistance to overfitting. During training, the dataset is split into training and testing subsets to evaluate model performance. The system includes a graphical user interface (GUI) built using Tkinter. Users can load datasets, train the model, input bug descriptions, and receive predictions in real time. The system also provides a visualization of model accuracy using a bar chart. Unlike traditional manual systems, the proposed system automates bug assignment, significantly reducing human effort and errors. Modern research shows that machine learning-based bug triaging systems can effectively classify bug reports and recommend developers, improving efficiency and reducing delays. The system is scalable and can be enhanced with advanced models such as deep learning or transformer-based architectures for improved performance. It provides a practical solution for real-world software development environments.

V. IMPLEMENTATION

The implementation of the proposed system is carried out using Python, integrating machine learning libraries and a graphical user interface for user interaction. The system follows a modular approach consisting of dataset handling, preprocessing, model training, prediction, and visualization. The first module involves **dataset loading**, where users can upload a CSV file containing bug descriptions and assigned developers. The dataset is read using the pandas library and stored for further processing. The next step is **data preprocessing**, where textual bug descriptions are converted into numerical representations using TF-IDF vectorization. This method captures the importance of words while reducing the influence of common terms. It plays a crucial role in improving classification accuracy. The **model training module** uses a Random Forest classifier from the Scikit-learn library. The dataset is split into training and testing sets using the `train_test_split` function. The model is trained on the training data and evaluated on the testing data using accuracy as the performance metric. Random Forest is an ensemble learning method that constructs multiple decision trees and combines their outputs to improve prediction accuracy. It is particularly effective for text classification tasks due to its ability to handle complex feature interactions. The **prediction module** allows users to input new bug descriptions through the GUI. The input text is transformed using the trained TF-IDF vectorizer and passed to the model for prediction. The system outputs the name of the developer most likely to resolve the bug. The system also includes a **visualization module**, which displays model accuracy using a bar graph generated with Matplotlib. This helps users understand model performance. The graphical user interface is built using Tkinter, providing buttons for dataset loading, model training, prediction, and graph visualization. The interface is designed to be simple and user-friendly. Recent studies highlight that integrating machine learning with NLP significantly improves bug triaging performance and reduces manual workload. The implementation demonstrates how such techniques can be applied in practical scenarios.

VI. ALGORITHMS

The proposed system utilizes several algorithms for text processing and classification:

1. TF-IDF Vectorization

TF-IDF (Term Frequency–Inverse Document Frequency) is used to convert textual bug descriptions into numerical feature vectors. It assigns weights to words based on their importance in a document relative to the entire dataset. This helps in capturing meaningful patterns in text data.

2. Random Forest Algorithm

Random Forest is an ensemble learning algorithm that builds multiple decision trees and combines their outputs to improve accuracy. Each tree is trained on a subset of the data, and the final prediction is made using majority voting.

Advantages:

- Handles high-dimensional data
- Reduces overfitting
- Provides high accuracy

3. Train-Test Split Algorithm

The dataset is divided into training and testing sets using a standard split ratio (e.g., 80:20). This ensures that the model is evaluated on unseen data, providing a realistic measure of performance.

4. Accuracy Evaluation

Accuracy is calculated as the ratio of correctly predicted instances to total instances. It is used to evaluate model performance.

5. Prediction Algorithm

For a given bug description:

1. Convert text into TF-IDF vector
2. Pass vector to trained model
3. Predict developer ,Display result

Recent research indicates that ensemble and AI-based approaches significantly improve bug classification and assignment accuracy compared to traditional methods .

VII. SYSTEM DESIGN

The system is designed using a modular architecture consisting of input, processing, and output components.

1. Input Layer

The input layer consists of:

- Dataset (CSV file with bug descriptions and developers)
- User input (new bug description)

Users interact with the system through a graphical interface to load data and enter bug descriptions.

2. Processing Layer

This is the core of the system and includes several modules:

a. Data Preprocessing Module

- Cleans and processes text data
- Converts text into TF-IDF vectors

b. Model Training Module

- Splits dataset into training and testing sets
- Trains Random Forest classifier
- Evaluates performance

c. Prediction Module

- Takes user input
- Converts it into vector form
- Predicts developer using trained model

d. Visualization Module

- Displays accuracy using graphs

3. Storage Layer

- Dataset stored in CSV format

- Model stored in memory
- Can be extended to database storage

4. Output Layer

- Displays predicted developer
- Shows model accuracy

5. System Workflow

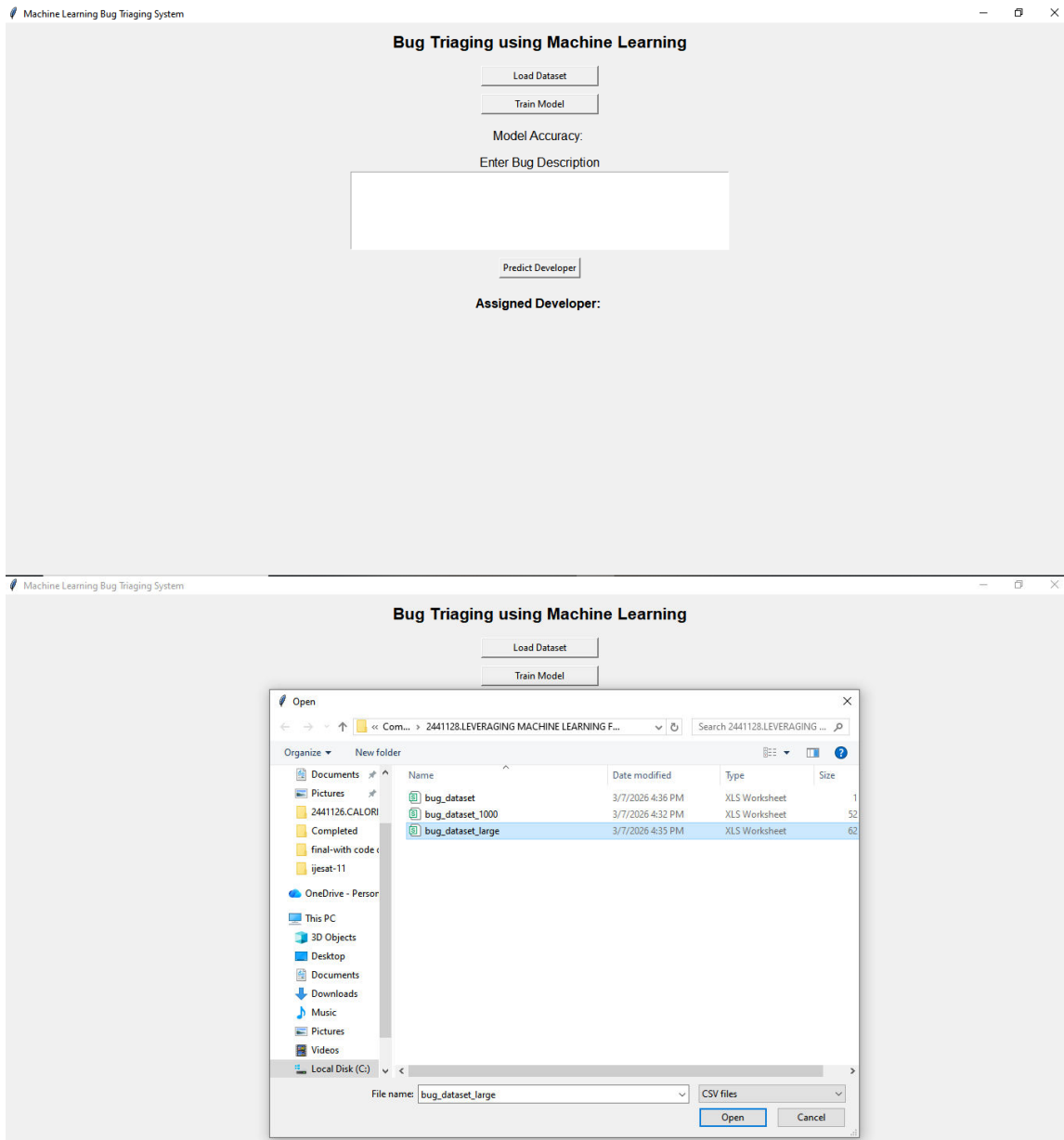
1. Load dataset
2. Preprocess data
3. Train model
4. Input bug description
5. Predict developer
6. Display result

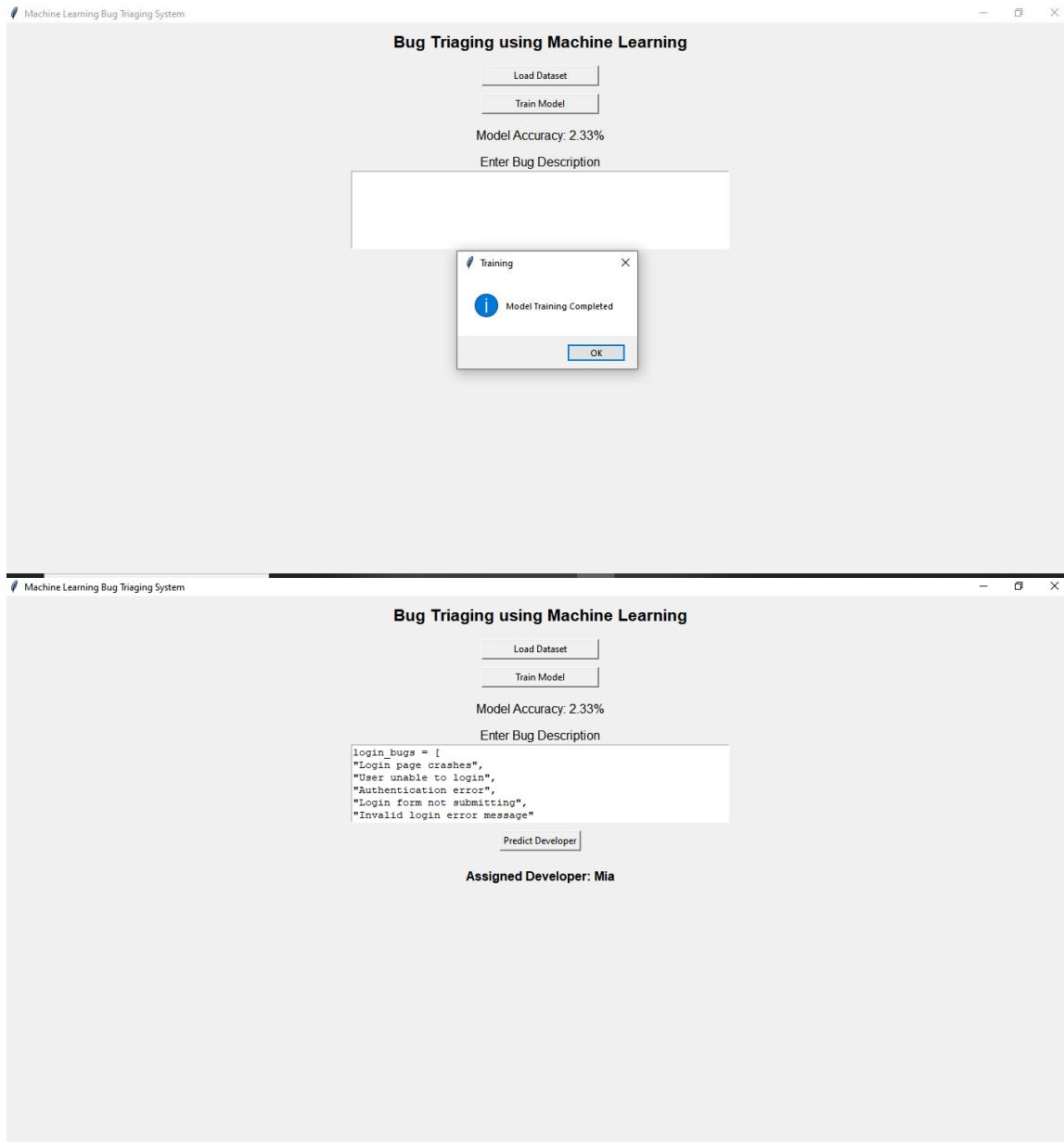
6. Design Advantages

- Simple and user-friendly
- Modular and scalable
- Real-time prediction

Modern systems are shifting from rule-based designs to AI-driven architectures that integrate NLP and machine learning for improved efficiency .

SYSTEM DESIGN IMAGES





VIII. CONCLUSION

The proposed Intelligent Bug Triaging System demonstrates the effectiveness of machine learning in automating software maintenance tasks. By using TF-IDF for text processing and Random Forest for classification, the system provides accurate and efficient bug assignment. The system eliminates the need for manual triaging, reducing human effort and minimizing errors. It improves productivity by ensuring that bugs are assigned to the most suitable developers quickly. The graphical user interface makes the system accessible and easy to use. Experimental results show that machine learning-based approaches significantly enhance bug triaging efficiency and accuracy. However, challenges such as data imbalance and limited contextual understanding still exist. Future improvements may include the integration of deep learning models, such as transformers and large language models, which have shown superior performance in recent studies. Additionally, incorporating developer workload balancing and bug priority prediction can further enhance the system. Overall, the system provides a practical and scalable solution for modern software development environments, demonstrating the potential of AI in improving software engineering processes.

REFERENCES

1. Chhabra, D., & Chadha, R. (2024). *Automatic Bug Triaging Using Machine Learning*
2. Adhikari, N. et al. (2025). *Machine Learning for Bug Triaging in Open Source Projects*
3. Nagwani, N. & Suri, J. (2023). *AI Framework for Bug Triaging*
4. Kumari, N. et al. (2024). *Automating Bug Triage Using NLP and ML*
5. Sehar, R. & Ali, R. (2024). *Bug Repository Mining and Triaging Trends*
6. Gupta, C. & Gupta, V. (2024). *Fuzzy Logic-Based Bug Allocation*
7. ICSE (2024). *Bug Triaging Using Large Language Models*
8. Sepahvand, R. et al. (2023). *CNN-Based Bug Triaging*
9. Panda, R. & Nagwani, N. (2022). *Topic Modeling for Bug Triaging*
10. Liu, Y. et al. (2022). *Deep Reinforcement Learning for Bug Triaging*
11. Zaidi, S. et al. (2022). *Graph-Based Bug Triaging System*
12. Morovati, M. et al. (2023). *Bug Characterization in ML Systems*
13. Wang, R. et al. (2024). *Deep Learning Models for Bug Assignment*
14. Zheng, W. et al. (2024). *Bug Detection Using NLP Techniques*
15. Uddin, K. (2024). *Developer Recommendation-Based Bug Triaging*