

Software Design Error Prediction in Software Using SVM Learning Algorithm

Rennila, Sajalingam

²Department of CSE, , BMS College of Engineering, India

Abstract— Software Engineering is a branch of computer science that enables tight communication between system software and training it as per the requirement of the user. We have selected seven distinct algorithms from machine learning techniques and are going to test them using the data sets acquired for NASA public promise repositories. The results of our project enable the users of this software to bag up the defects are selecting the most efficient of given algorithms in doing their further respective tasks, resulting in effective results.

Corresponding Author:

Rennila, Sajalingam

Keywords: *Software quality metrics, Software defect prediction, Software fault prediction, Machine learning algorithms*

Copyright © 2025: Rennila, Sajalingam, This is an open access distribution, and reproduction in any medium, provided Access article distributed under the Creative Commons Attribution License the original work is properly cited License, which permits unrestricted use.

Citation: Rennila, Sajalingam, “Software Design Error Prediction in Software Using SVM Learning Algorithm”, Journal of Science, Computing and Engineering Research, Volume-9, Issue-1, January 2026.

I. INTRODUCTION

A) Problem Statement:

Now-a-days developing software system is a difficult process which involves planning, analyzing, designing, implementing, test, integrate and maintenance. A software engineer work is developing a system in time with limited budget which is done in planning phase. While doing the development process we can have few defects like not proper design, where the logic is poor, data handling is improper, etc. and these defects cause errors which lead to re-do the work, increasing in development and cost of maintenance. This all are responsible for the decrease in customer satisfaction. In this point of view, faults are grouped on the basis of sternness, corrective and advance Actions are taken as per the sternness defined. The selected machine learning algorithms for comparison are:

- a. Multilayer Perceptron (MLP)
- b. KNN classifier
- c. Guassian Naïve Bayes
- d. Decision tree
- e. Support Vector Classifier (SVC).
- f. Ensemble method

B) Objective:

The Objective of this project is to estimate the defect of software using machine learning algorithms. On training the various ML Algorithms we need to get good accuracy percentage so that the particular algorithm fits the best in order to estimate the defects Support Vector Classifier (SVC) supports both classification as well as regression. It is productive and straight-lined method which is used in classification. For classification it divides two groups by making boundaries between the group of data.

C) Proposed System:

The proposed system includes SVC, Multilayer perceptron, Naive Baye’s algorithm, Decision Tree, KNN Classifier, Ensemble method, Functions to solve the class misbalancing problem which causes in the decreasing performance of defect prediction. The dataset has been trained and spitted according to the constraints and using the accuracies has been defined in order to measure the defect estimation capability of various algorithms proposed.

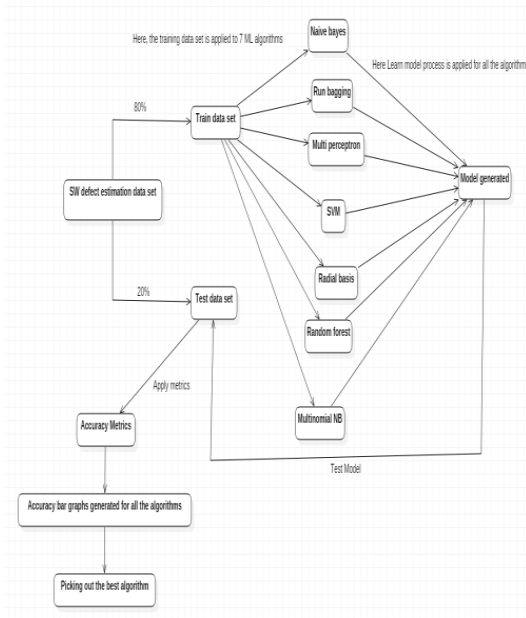
D) Advantages of proposed system:

1. Predicted model is used for evaluating the performance measures.
2. We can apply various datasets in this project. But we are using NASA datasets in our project.
3. Software defects are classified to the extent.

4. Advance measures can be taken on selection of algorithm
5. Provides Better results.
6. Identify defects in the early stage of the project which in turn results in Customer loyalty.

II. SYSTEM DESIGN

A) System Architecture:



III. RESEARCH METHODOLOGY

The proposed methodology for software failure prediction employs a systematic approach, encompassing critical stages essential for comprehensive data analysis and model performance assessment.

A) DATA INFORMATION

The JM1 dataset, a component of the PROMISE repository, is designed for software defect prediction and has been made publicly available by NASA and the NASA Metrics Data Program. Naive Baye’s, derived from JM1, has exhibited promising performance, out performing J48 for defect detection, and the dataset has highlighted the nuanced relationship between accuracy and the effectiveness of defect detectors. The JM1 dataset comprises 10,885 instances, each characterized by 22 attributes.

Notably, there are no missing attributes in the dataset. The class distribution indicates that 19.35% of instances are labeled

as "false" (modules without reported defects), while 80.65% are labeled as "true" (modules with one or more reported defects).

id	v1(E)	v2(E)	v3(E)	v4	v5	v6	v7	v8	v9	v10	v11	v12	v13	v14	v15	v16	v17	v18	v19	v20	v21	v22	loc_code	loc_comment	loc_lines	loc_colno	loc_rowno
0	1.1	1.4	1.4	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3	1.3
1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
2	72.0	7.0	1.0	6.0	190.0	1134.13	0.95	20.31	55.85	23029.10	...	51	10	8	1	1	1	1	1	1	1	1	1	1	1	1	1
...
10882	42.0	4.0	1.0	2.0	103.0	616.57	0.04	36.40	19.60	13716.72	...	20	1	10	0	0	0	0	0	0	0	0	0	0	0	0	0
10883	10.0	1.0	1.0	1.0	30.0	142.15	0.12	0.44	17.44	1241.57	...	6	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0
10884	19.0	3.0	1.0	1.0	50.0	272.63	0.09	11.57	23.36	3164.67	...	13	0	2	1	1	1	1	1	1	1	1	1	1	1	1	1

Figure: data frame

B) DATA PRE-PROCESSING

In the preprocessing pipeline for software failure prediction using machine learning, several critical steps are undertaken to refine textual data and optimize it for analysis, the standard scaler technique is applied to normalize and scale the data, enhancing its numerical stability. Subsequently, the text is transformed to lowercase to ensure consistency and mitigate case-related variations. To facilitate readability and semantic analysis, punctuation signals are then systematically removed from the text. The final preprocessing step involves lemmatization, wherein words are transformed into their root forms. This not only improves consistency but also aids in further reducing dimensionality while retaining the essential semantic information. The overarching goal of this preprocessing approach is to enhance the quality of the textual data, minimizing noise and ensuring its readiness for subsequent analysis or utilization in machine learning algorithms tailored for software failure prediction

C) EDA

Exploratory Data Analysis (EDA) stands as a pivotal phase in the continuum of data analysis, marked by its systematic exploration, visualization, and comprehension of a dataset. It functions as a fundamental tool, empowering data analysts, researchers, and scientists to unearth pertinent insights, identify patterns, detect anomalies, and formulate hypotheses. EDA serves as a crucial precursor to in-depth studies and informed decision-making. Its significance lies in its dual capability: first, to unveil latent information inherent in the data set, and second, to provide a framework for conducting thorough and comprehensive research.

By cultivating a more comprehensive insight into the data, EDA plays a pivotal role in guiding subsequent stages of analysis, contributing to the formulation of hypotheses, model selection, and the overall refinement of analytical methodologies.

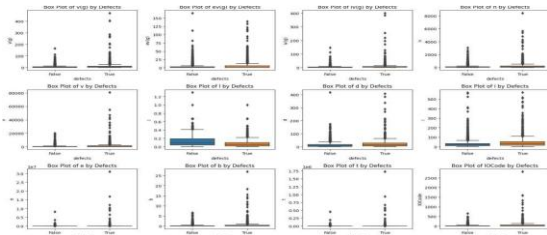


Figure: a box plot depicts the distribution of software defect metrics

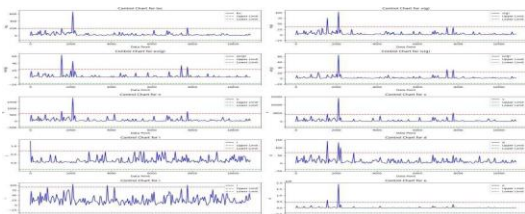


Figure: Informative control charts for the specified software Metrics

D) DATA SPLITTING

In the domain of machine learning for software failure prediction, a widely adopted methodology is the "80:20 data partitioning" strategy, which involves dividing a dataset into training and testing subsets, allocating 80% for training and reserving the remaining 20% for testing. This approach serves as a standard practice due to its effectiveness. By maintaining a clear distinction between training and testing data, this methodology establishes a framework for assessing the credibility of predictive and analytical outcomes. The evaluation process involves contrasting model-generated results with benchmarks derived from the reserved testing subset, ensuring a thorough assessment of the model's capabilities.

1) Accuracy

In the domain of classification, accuracy is quantitatively defined as the ratio of correctly classified instances to the total size of the dataset, expressed as a percentage. This fundamental metric serves as a rigorous measure of a classification model's effectiveness in accurately assigning data examples to their respective categories.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

2) Loss :

In instances where anticipated outcomes diverge from actual observations, the consequential outcome is often disappointment. This process aims to optimize the model's performance by mitigating discrepancies between predicted

and actual outcomes, thereby enhancing its utility and reliability.

$$Loss = -1 \sum_{i=1}^m f_i \log(f_i)$$

3) Precision

Precision is a metric that assesses the model's performance in making accurate predictions, specifically measuring how often the model correctly anticipates a favorable outcome. This statistical measure addresses the question of how many times the model accurately predicts positive instances. Mathematically, precision is expressed as the ratio of true positives to the sum of true positives and false positives

$$Precision = \frac{TP}{TP+FP}$$

4) Recall

The evaluation of model performance hinges on its ability to recall and correctly identify all pertinent data points. Specifically, when posed with the question, "Among all the genuine positive instances, how many did the model accurately predict as positive?" recall furnishes the relevant answer. In essence, while a high recall indicates the model's proficiency in capturing positive instances, a balanced consideration of other metrics is imperative for a comprehensive assessment of its overall performance.

$$Recall = \frac{TP}{TP+FN}$$

5) F1 - Score

The F1-score serves as a consolidated metric that integrates a classifier's recall and precision through the calculation of their harmonic mean. This metric is specifically designed for the comparative evaluation of two classifiers.

$$F1-score = \frac{2}{\frac{1}{precision} + \frac{1}{recall}}$$

Performance Evaluation of Machine Learning Model

Models	Accuracy	Precision	Recall	F1-score
--------	----------	-----------	--------	----------

MLP model	0.93	0.70	0.92	0.79
Gaussian NB	0.98	0.92	0.96	0.94
Decision tree	0.99	1.0	1.0	1.0
SVC	0.97	0.97	0.80	0.87
KNN	0.99	0.98	0.96	0.93
Ensemble	0.99	1.0	1.0	1.0

The table provides a thorough assessment of machine learning models for software defect prediction, using important performance measures. The rows in the table represent individual models, while the columns provide information on several metrics, including accuracy, precision, recall, and F1-score.

The MLP model demonstrates a commendable accuracy of 0.93, suggesting that it consistently makes correct predictions. Gaussian Naive Bayes (NB) demonstrates exceptional performance in defect prediction, with a high accuracy of 0.98 and well-balanced precision, recall, and F1-score of 0.92, 0.96, and 0.94 respectively.

The Decision Tree algorithm demonstrates outstanding prediction capabilities, achieving perfect scores (1.0) in accuracy, precision, recall, and F1-score, across all measures.

The Support Vector Classifier (SVC) demonstrates a commendable accuracy of 0.97.

KNN has exceptional performance with a high level of accuracy (0.99) and well-balanced precision (0.98), recall (0.96), and F1-score (0.93).

The Ensemble method, which mirrors the Decision Tree, receives perfect scores (1.0) in all measures, confirming its effectiveness in predicting software defects.

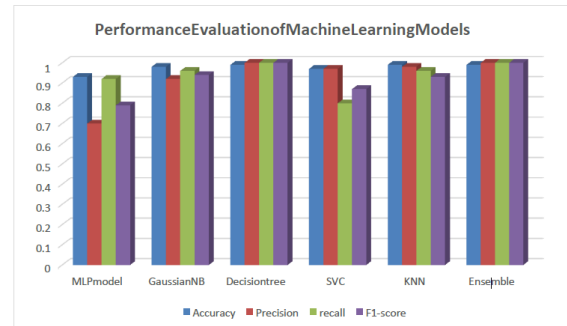


Figure: Performance evaluation graph of machine learning models

The Decision Tree and Ensemble models demonstrate exceptional predictive capabilities with perfect scores across all metrics, making them suitable for tasks where precision, recall, and overall accuracy are crucial. Gaussian NB excels with a balanced approach, showcasing high accuracy and well-maintained precision, recall, and F1-score, making it a reliable choice for defect prediction. KNN, with its high accuracy and balanced precision and recall, offers a robust solution for identifying software defects.

IV. CONCLUSION

Software defects can have a severe impact on software quality, causing problems for customers and developers. With growing complexities in software designs and technology, manual software detection becomes a challenging and time-consuming task. Thus, automatic software detection has become a hotspot for industrial research in the past couple of years. In this paper, we try to apply machine learning and deep learning to solve this problem. We use datasets provided by the NASA Promise dataset repository and compare the state of the art machine learning algorithms' results. The strengths and weaknesses of each model are highlighted in these measures, helping choose the best model depending on specific goals and trade-offs between precision and recall in software defect prediction tasks. This field still has much scope for improvement. We can think of some novel approaches which use complex deep learning algorithms, and also researchers should focus on more data collection.

REFERENCES

- [1]. P. Nirmala, T. Manimegalai, J. R. Arunkumar, S. Vimala, G. Vinoth Rajkumar, Raja Raju, "A Mechanism for Detecting the Intruder in the Network through a Stacking Dilated CNN Model", Wireless Communications and Mobile Computing,

- vol. 2022, Article ID 1955009, 13 pages, 2022. <https://doi.org/10.1155/2022/1955009>.
- [2]. D. Sathyanarayanan, T. S. Reddy, A. Sathish, P. Geetha, J. R. Arunkumar and S. P. K. Deepak, "American Sign Language Recognition System for Numerical and Alphabets," 2023 International Conference on Research Methodologies in Knowledge Management, Artificial Intelligence and Telecommunication Engineering (RMKMATE), Chennai, India, 2023, pp. 1-6, doi: 10.1109/RMKMATE59243.2023.10369455.
- [3]. J. R. Arunkumar, Tagele berihun Mengist, 2020" Developing Ethiopian Yirgacheffe Coffee Grading Model using a Deep Learning Classifier" International Journal of Innovative Technology and Exploring Engineering (IJITEE) ISSN: 2278-3075, Volume-9 Issue-4, February 2020. DOI: 10.35940/ijitee.D1823.029420.
- [4]. Ashwini, S., Arunkumar, J.R., Prabu, R.T. et al. Diagnosis and multi-classification of lung diseases in CXR images using optimized deep convolutional neural network. *Soft Comput* (2023). <https://doi.org/10.1007/s00500-023-09480-3>
- [5]. J.R.Arunkumar, Dr.E.Muthukumar," A Novel Method to Improve AODV Protocol for WSN" in *Journal of Engineering Sciences* ISSN NO: 0377-9254 Volume 3, Issue 1, Jul 2012.
- [6]. R. K, A. Shameem, P. Biswas, B. T. Geetha, J. R. Arunkumar and P. K. Lakineni, "Supply Chain Management Using Blockchain: Opportunities, Challenges, and Future Directions," 2023 Second International Conference on Informatics (ICI), Noida, India, 2023, pp. 1-6, doi: 10.1109/ICI60088.2023.10421633.
- [7]. Arunkumar, J. R. "Study Analysis of Cloud Security Challenges and Issues in Cloud Computing Technologies." *Journal of Science, Computing and Engineering Research* 6.8 (2023): 06-10.
- [8]. J. R. Arunkumar, R. Raman, S. Sivakumar and R. Pavithra, "Wearable Devices for Patient Monitoring System using IoT," 2023 8th International Conference on Communication and Electronics Systems (ICES), Coimbatore, India, 2023, pp. 381-385, doi: 10.1109/ICES57224.2023.10192741.
- [9]. S. Sugumaran, C. Geetha, S. S, P. C. Bharath Kumar, T. D. Subha and J. R. Arunkumar, "Energy Efficient Routing Algorithm with Mobile Sink Assistance in Wireless Sensor Networks," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ACCAI58221.2023.10201142.
- [10]. R. S. Vignesh, V. Chinnammal, Gururaj.D, A. K. Kumar, K. V. Karthikeyan and J. R. Arunkumar, "Secured Data Access and Control Abilities Management over Cloud Environment using Novel Cryptographic Principles," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-8, doi: 10.1109/ACCAI58221.2023.10199616.
- [11]. Syamala, M., Anusuya, R., Sonkar, S.K. et al. Big data analytics for dynamic network slicing in 5G and beyond with dynamic user preferences. *Opt Quant Electron* 56, 61 (2024). <https://doi.org/10.1007/s11082-023-05663-2>
- [12]. Krishna Veni, S. R., and R. Anusuya. "Design and Study Analysis Automated Recognition system of Fake Currency Notes." *Journal of Science, Computing and Engineering Research* 6.6 (2023): 16-20.
- [13]. V. RamKumar, S. Shanthi, K. S. Kumar, S. Kanageswari, S. Mahalakshmi and R. Anusuya, "Internet of Things Assisted Remote Health and Safety Monitoring Scheme Using Intelligent Sensors," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-8, doi: 10.1109/ACCAI58221.2023.10199766.
- [14]. R. S. Vignesh, R. Sankar, A. Balaji, K. S. Kumar, V. Sharmila Bhargavi and R. Anusuya, "IoT Assisted Drunk and Drive People Identification to Avoid Accidents and Ensure Road Safety Measures," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ACCAI58221.2023.10200809.
- [15]. I. Chandra, G. Sowmiya, G. Charulatha, S. D, S. Gomathi and R. Anusuya, "An efficient Intelligent Systems for Low-Power Consumption Zigbee-Based Wearable Device for Voice Data Transmission," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ICECONF57129.2023.10083856.
- [16]. G. Karthikeyan, D. T. G, R. Anusuya, K. K. G, J. T and R. T. Prabu, "Real-Time Sidewalk Crack Identification and Classification based on Convolutional Neural Network using Thermal Images," 2022 International Conference on Automation, Computing and Renewable Systems (ICACRS), Pudukkottai, India, 2022, pp. 1266-1274, doi: 10.1109/ICACRS55517.2022.10029202.
- [17]. R. Meena, T. Kavitha, A. K. S, D. M. Mathew, R. Anusuya and G. Karthik, "Extracting Behavioral Characteristics of College Students Using Data Mining on Big Data," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ICECONF57129.2023.10084276.
- [18]. S. Bharathi, A. Balaji, D. Irene, J. C. Kalaivanan and R. Anusuya, "An Efficient Liver Disease Prediction based on Deep Convolutional Neural Network using Biopsy Images," 2022 3rd International Conference on Smart Electronics and Communication (ICOSEC), Trichy, India, 2022, pp. 1141-1147, doi: 10.1109/ICOSEC54921.2022.9951870.
- [19]. I. Chandra, G. Sowmiya, G. Charulatha, S. D, S. Gomathi and R. Anusuya, "An efficient Intelligent Systems for Low-Power Consumption Zigbee-Based Wearable Device for Voice Data Transmission," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ICECONF57129.2023.10083856. I. Chandra, K. V. Karthikeyan, R. V, S. K, M. Tamilselvi and J. R. Arunkumar, "A Robust and Efficient Computational Offloading and Task Scheduling Model in Mobile Cloud Computing," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-8, doi: 10.1109/ICECONF57129.2023.10084293.
- [20]. Revathi, S., et al. "Developing an Infant Monitoring System using IoT (INMOS)." *International Scientific Journal of*

Contemporary Research in Engineering Science and Management 6.1 (2021): 111-115.

- [21].R. K. A. Shameem, P. Biswas, B. T. Geetha, J. R. Arunkumar and P. K. Lakineni, "Supply Chain Management Using Blockchain: Opportunities, Challenges, and Future Directions," 2023 Second International Conference on Informatics (ICI), Noida, India, 2023, pp. 1-6, doi: 10.1109/ICI60088.2023.10421633.
- [22].J.R.Arunkumar. "Comprehensice Analysis of Security Issues in Cloud Computing Technologies", Journal of Science, Computing and Engineering Research, 6(5), 06-10, June 2023.
- [23].S. Sugumaran, C. Geetha, S. S, P. C. Bharath Kumar, T. D. Subha and J. R. Arunkumar, "Energy Efficient Routing Algorithm with Mobile Sink Assistance in Wireless Sensor Networks," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ACCAI58221.2023.10201142.
- [24].I. Chandra, K. V. Karthikeyan, R. V, S. K, M. Tamilselvi and J. R. Arunkumar, "A Robust and Efficient Computational Offloading and Task Scheduling Model in Mobile Cloud Computing," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-8, doi: 10.1109/ICECONF57129.2023.10084293.
- [25].R. S. Vignesh, A. Kumar S, T. M. Amirthalakshmi, P. Delphy, J. R. Arunkumar and S. Kamatchi, "An Efficient and Intelligent Systems for Internet of Things Based Health Observance System for Covid 19 Patients," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-8, doi:

