



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 10    **Issue:** V    **Month of publication:** May 2022

**DOI:** <https://doi.org/10.22214/ijraset.2022.42306>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# AI Based False Positive Analysis of Software Vulnerabilities

Priya Patil

Assistant Professor Sharnbasva University BCA, Gulbarga

**Abstract:** *Programming measurements and shortcoming information having a place with a past programming variant are utilized to assemble the product issue expectation model for the following arrival of the product. Notwithstanding, there are sure situations when past issue information are absent. As such foreseeing the shortcoming inclination of program modules when the issue marks for modules are inaccessible is a difficult assignment oftentimes arised in the product business There is need to foster a few strategies to assemble the product issue forecast model in light of unaided realizing which can assist with anticipating the shortcoming inclination of a program modules when shortcoming names for modules are absent. One of the strategies is utilization of grouping methods. Solo methods like grouping might be utilized for issue expectation in programming modules, all the more so in those situations where shortcoming names are not accessible. In this review, we propose a Machine Learning grouping based programming shortcoming forecast approach for this difficult issue.*

## I. INTRODUCTION

Programming quality models are helpful devices toward accomplishing the goals of a product quality confirmation drive. A product quality model can be utilized to recognize program modules that are probably going to be deficient. Therefore, the restricted assets assigned for programming quality review and improvement can be designated toward just those program modules, accomplishing financially savvy asset usage. A product quality assessment model permits the product improvement group to follow and recognize potential programming absconds generally from the beginning during advancement, which is basic to some high-confirmation frameworks.

A product quality model is normally prepared utilizing programming estimation and imperfection (quality) information of a formerly evolved discharge or comparative venture. The prepared model is then applied to modules of the ongoing venture to appraise their quality. Such a managed learning approach expects that the improvement association has insight with frameworks like the ongoing undertaking and that deformity information are accessible for all program modules in the preparation information. In programming improvement practice, be that as it may, different pragmatic issues limit the accessibility of deformity information for modules in the preparation information. For instance, an association might have not recorded or gathered programming deformity information from past deliveries or comparable ventures. Likewise, since the association might not have experience fostering a comparable framework, the utilization of programming estimation and imperfection information of past undertakings for it is unseemly to demonstrate purposes. In present situations, where globalization of innovation has picked up speed, disseminated programming advancement is entirely expected. Under such circumstances, programming deformity information may not be gathered by all advancement locales relying upon the authoritative construction and assets of individual destinations. The shortfall of deformity information or quality-based class names from the preparation information forestalls following the usually utilize directed learning way to deal with programming quality demonstrating. Subsequently, the errand of programming quality assessment or marking program modules as shortcoming inclined (fp) or not issue inclined (nfp) falls on the programming master. The method involved with marking each program module each in turn is a difficult, costly, and tedious exertion. We propose a semi directed grouping plan to help the master in the marking system. The proposed conspire depends on imperative based grouping involving k-implies as the fundamental calculation. During the k-implies bunching process, the imperative keeps up with enrollment of modules (occurrences) to groups that are as of now marked as either fp or nfp. The proposed approach is to examine the nature of unlabelled S/W modules utilizing two phase draws near. One is sing measurements limits and another is utilizing FuzzyC Means grouping and afterward contrasting both as far as time and mean squared mistake esteem. The bunching techniques make gatherings of endlessly protests in a single bunch are comparative while objects in Clustering strategies can be utilized to bunch the modules having comparable measurements by utilizing comparability measures or divergence measures (distances). Subsequent to bunching stage, a specialist or a robotized approach can actually take a look at the agent modules of each group and afterward, choose to mark the bunch as shortcoming inclined or not issue inclined. In this review, we demonstrate the way that product measurements limits

can be utilized to mark the bunches as opposed to involving a specialist for this time-it are unlike consume stage different group. Grouping is essentially Partitional Clustering: Given a data set of  $n$  protests, a partitional bunching calculation builds  $k$  segments of the information, where each bunch improves a grouping standard, like the minimization of the amount of squared separation from the mean inside each bunch. Kinds of Partitional bunching:  $K$ -implies Clustering and Machine Learning Clustering.

### A. Machine Learning

In machine learning, data plays an important role, and the machine learning is used catch on and learn properties from the data. The learning and prediction performance will effect on the quantity and quality of data set.

### B. Types Of Machine Learning

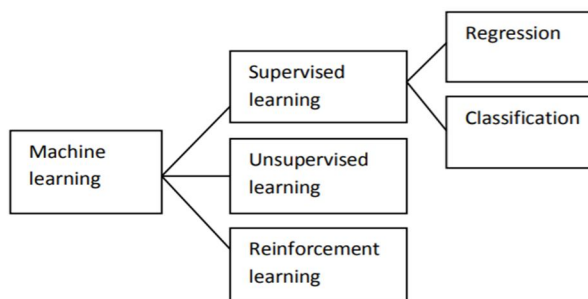


Figure 1 Types of Machine Learning

- 1) *Supervised Learning*: Directed learning is prepared a named information. Directed advancing commonly utilized in applications where authentic information foresee like future use. This strategy further isolated into two classes as relapse and characterization. In relapse the mark is ceaseless amount. Then again, in grouping the mark is discrete
- 2) *Unsupervised Learning*: Unsupervised learning used unlabeled data that has no historical labels.
- 3) *Reinforcement Learning*: It is used for gaming, navigation and robotics. This type of learning has three types of primary components: the learner, the environment and actions.
- 4) *Machine Learning Techniques*: A machine learning algorithms are developed to build machine learning models and important machine learning process. In this paper we discuss three classifiers like decision tree, naïve bayes and support vector machine.
- a) *Decision Tree*: decision tree based on supervised learning algorithm. It is one of the predictive model approaches used in machine learning, data mining and statistics. In decision analysis, a decision tree can used to usually represent decision making and decisions. In data mining, a decision tree characterize data but not decision s, to some extent the resulting allocation tree can be an input for decision making.

There are many decision tree algorithms:

C4.5 (successor of ID3) CART (classification and regression tree)

MARS: extends decision tree to better handle numerical data.

- b) *Naïve bayes*: Naïve bayes classifier is widely studied probabilistic learning method. Naïve Bayesian classifier conclude that there are no addiction among attributes. This presumption is called conditional independence.

Advantages of naïve bayes:

- It is used a very perceptive technique.
- It widely studied expectation learning method.
- Naïve bayes classifier is computational fast when executing decision.

- c) *Support vector machine (SVM)*: “Support Vector Machine” (SVM) is a supervised machine learning algorithm which can be used for both classification and regression challenges. However it is mostly used in classification problems. SVM used for classification of both linear and non-linear data. In this algorithm, we plot each data item as a point in  $n$ -dimensional space (where  $n$  is number of features you have) with the value of each feature being the value of a particular coordinate.

Support vector machine some steps:

- Set up the training data
- Set up of SVM parameters
- Types of SVM kernel
- Train the SVM
- Regions classified by SVM
- Support vectors

## II. RELATED WORK

- 1) **Deep Learning based Vulnerability Detection: Are They have There Yet [1]**- In this paper, they have inquire, "how they have do the cutting edge DL-based procedures act in a certifiable weakness expectation situation". Shockingly, they have observe that their presentation drops by over half. An orderly examination of what causes such steep execution drop uncovers that current DL-based weakness forecast approaches experience the ill effects of difficulties with the preparation information (e.g., information duplication, ridiculous conveyance of weak classes, and so forth) and with the model decisions (e.g., straightforward token-based models). Subsequently, these methodologies frequently don't learn highlights connected with the real reason for the weaknesses. All things considered, they gain inconsequential relics from the dataset (e.g., explicit variable/work names, and so on.). Utilizing these exact discoveries, they have exhibit how a more principled way to deal with information assortment and model plan, in view of reasonable settings of weakness forecast, can prompt improved arrangements. The subsequent instruments perform essentially better compared to the concentrated on standard up to 33.57% lift in accuracy and 128.38% lift in review contrasted with the best performing model in the writing. In general, this paper clarifies existing DL-based weakness expectation frameworks' possible issues and draws a guide for future DL-based weakness forecast research.
- 2) **VIVA: Binary Level Vulnerability Identification via Partial Signature[2]**- In this paper they have propose a VIVA, a twofold level weakness and fix semantic rundown and matching device for precise repeating weakness identification. It utilizes novel parallel program cutting methods with the guide of pseudo-code follow refinement to create fractional weakness and fix marks, which catch the semantics. It coordinates the marks with pre-separating to recognize 1-day and repeating weaknesses productively. The exploratory outcomes show that VIVA beats other source code and twofold coordinating devices with an accuracy of 100 percent for 1-day weaknesses and 87.6% for repeating weaknesses and great execution (28.58s per signature search in 4M capacities). It identifies 92 new weaknesses in various series and various variants of genuine ventures, with 11 exist without fixing in the most recent rendition
- 3) **D2A: A Dataset Built for AI-Based Vulnerability Detection Methods Using Differential Analysis [3]**- They have propose D2A, a differential examination based way to deal with mark issues detailed by static investigation apparatuses. The D2A dataset is worked by dissecting rendition matches from different open source projects. From each undertaking, they have select bug fixing commits and they have run static examination on the forms when such commits. Assuming a few issues identified in a preceding commit rendition vanish in the relating after-commit variant, they are probably going to be genuine bugs that sorted out by the commit. They have use D2A to produce a huge marked dataset to prepare models for weakness ID. They have show that the dataset can be utilized to assemble a classifier to recognize conceivable misleading problems among the issues revealed by static examination, subsequently helping engineers focus on and explore potential genuine up-sides first.
- 4) **A Practical Approach for Ranking Software Warnings from Multiple Static Code Analysis Reports [4]**- Static examination instruments inspect source code to search for programming imperfections and likely weaknesses. It is a typical practice to utilize numerous apparatuses so they have don't ignore code which really has an issue. Hothey havever, the issue of utilizing various bugs observing apparatuses is they identify comparable programming absconds as well as create new advance notice messages. The exorbitant alerts make code examination tedious and costly. In this paper, they have depict our strategies to combine programming advance notice classifications from various bugs finding devices from two famous programming dialects like Java and C++, and focus on the documents solidified cautioning messages by building an insightful model utilizing head part investigation. Results have shown that documents genuine programming deserts involved first spot on the list, and bogus up-sides involved the base openings.
- 5) **Techniques and Tools for Advanced Software Vulnerability Detection[5]**- This paper aims to study the combination of different techniques to improve the effectiveness of vulnerability detection (increasing the detection rate and decreasing the number of false-positives). Static Code Analysis (SCA) has a good detection rate and is the central technique of this work.

Hothey havever, as SCA reports many false-positives, they have will study the combination of various SCA tools and the integration with other detection approaches (e.g., software metrics) to improve vulnerability detection capabilities. They have will also study the use of such combination to prioritize the reported vulnerabilities and thus guide the development efforts and fixes in resource-constrained projects.

### III. EXISTING SYSTEM

In Existing, Quad Tree-based K-Means calculation has been applied for anticipating issues in program modules. To start with, Quad Trees are applied for viewing the underlying bunch communities as contribution to the K-Means Algorithm. An information edge boundary oversees the quantity of introductory group habitats and by fluctuating the client can produce wanted beginning bunch communities. The idea of bunching gain has been utilized to decide the nature of groups for assessment of the Quad Tree-based introduction calculation when contrasted with other instatement methods. The groups acquired by Quad Tree-based calculation were found to have most extreme increase values. Second, the Quad Treebased calculation is applied for anticipating issues in program modules.

Disadvantage in Existing System

- 1) The K-Means algorithm is very sensitive to noise.
- 2) The Quad Tree-based method assigns the appropriate initial cluster centers.
- 3) Using both Quad Tree and K-Means the Overall processing time was increased.
- 4) In software fault prediction, the value of threshold is not clearly described.

### IV. PROPOSED SYSTEM

In the proposed framework we propose a Machine Learning bunching method for programming weaknesses expectation. First we select the info dataset and apply the trait choice strategy. Characteristic determination strategy is utilized to choose the significant quality and lessen the quantity of traits. Characteristic Evaluation approach is utilized to assess the property and return the heaviness of the trait. Then, at that point, apply the positioning strategy for get the most weighted ascribes. After the Attribute determination we get the decreased number of characteristics. This decreased trait is accustomed to bunching approach. Given an information base of n protests, a partitional grouping calculation develops k segments of the information, where each bunch enhances a bunching measure, like the minimization of the amount of squared separation from the mean inside each bunch. Fluffy C means is a one kind of partitional bunching approach. To group the data of interest, we are utilizing Fuzzy-C means bunching. After the bunching approach each group keep up with the centroid esteem. Metric limit is approach used to characterize the edge esteem. Look at this metric edge and centroid information values. Assuming any measurement worth of the centroid data of interest of a bunch was more prominent than the edge, that group was named as VULNERABILITIES y and in any case it was marked as nonVULNERABILITIES.

Advantage in proposed system

- 1) Propose a single technique
- 2) FUZZY C means is used for Clustering.
- 3) Processing time is reduced
- 4) Clear description of Metric Threshold

### V. SYSTEM ARCHITECTURE

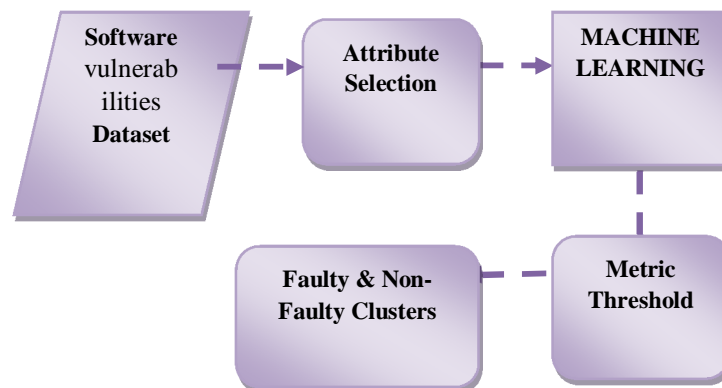


Figure 2 Proposed System

**A. Vulnerabilities Dataset**

Collect Software Dataset (AR3, AR4 and AR5). Read the data of the Dataset and stored into DB. Get the Attribute names for Attributes Selection.

- 1) **Attribute Selection:** Get all Attribute in our dataset, The dataset contains totally 30 numbers of attributes. Attribute Selection is the technique of selecting a subset of relevant features for building robust learning models. We take all attributes into our process it takes so much time for processing and increase the work burden. So, we reduce the total number of attributes and consider high relevance attributes only. Calculate the relevance of an attribute using Attribute Evaluation. We use Weka tool for attribute selection and ranker.
- 2) **Machine Learning Clustering:** A cluster is a collection of objects which are “similar” between them and are “dissimilar” to the objects belonging to other clusters. Machine Learning (FCM) is a method of clustering which allows one piece of data to belong to two or more clusters. The data’s are initially clustered and then performing Machine Learning algorithm. After the clustering process we separate clustered data’s and cluster centroid
- 3) **Metric Threshold:** Determine the acceptable metrics thresholds using some parameters. The Parameters are,
  - a) Lines of Code (LoC),
  - b) Cyclomatic Complexity (CC),
  - c) Unique Operator (UOp),
  - d) Unique Operand (UOpnd),
  - e) Total Operator (TOp),
  - f) Total Operand (TOpnd).

Finally, we have the threshold vector [LoC, CC, UOp, UOpnd, TOp, TOpnd]

**B. Detect Vulnerabilities**

After the clustering process each cluster maintain the data point. Get the metric threshold for each cluster. If any metric value of the centroid data point of a cluster was greater than the threshold, that cluster was labeled as faulty and otherwise it was labeled as non faulty

Just-in-time defect prediction aims to predict if a particular file involved in a commit (i.e., a change) is buggy or not. Traditional just-in-time defect prediction techniques typically follow the following steps:

- 1) **Training Data Extraction:** For each change, label it as buggy or clean by mining a project’s revision history and issue tracking system. Buggy change means the change contains bugs (one or more), while clean change means the change has no bug.
- 2) **Feature Extraction:** Extract the values of various features from each change. Many different features have been used in past change classification studies.
- 3) **Model Learning:** Build a model by using a classification algorithm based on the labeled changes and their corresponding features.
- 4) **Model Application:** For a new change, extract the values of various features. Input these values to the learned model to predict whether the change is buggy or clean.

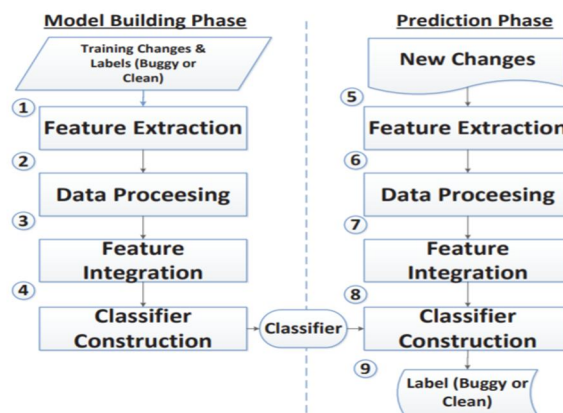


Figure 3 Flow of Proposed System

The framework mainly contains two phases: a model building phase and a prediction phase. In the model building phase, our goal is to build a classifier (i.e., a statistical model) by leveraging deep learning and machine learning techniques from historical changes with known labels (i.e., buggy or clean). In the prediction phase, this classifier would be used to predict if an unknown change would be buggy or clean.

C. Performance Measures

Machine learning checks the prediction performance with the help of various performance measures:

- 1) Precision
- 2) Recall
- 3) Accuracy
- 4) F measure
- 5) ROC (receiver operating characteristics)

There were calculated using the prediction classification confusion matrix table:

		PREDICTED	
		P	N
ACTUAL	P	TRUE POSITIVE (TP)	FALSE NEGATIVE (FN)
	N	FALSE POSITIVE (FP)	TRUE NEGATIVE (TN)

Table 1: Confusion Matrix

TP (true positive): Number of correct predictions that an instance is positive.

TN (true negative): Number of correct predictions that an instance is negative.

FN (false negative): Number of incorrect predictions that an instance is positive.

FP (false positive): Number of incorrect predictions that instance is negative.

- Accuracy: The total number of predictions that were correct

$$\text{Accuracy (\%)} = (TP+TN)/(TP+FP+FN+TN)$$

- Precision: The predicted true pages those were correct:

$$\text{Precision (\%)} = TP/(TP+EP)$$

- Recall: The predicted true pages that were correctly identify.

$$\text{Recall (\%)} = TP/(FN+TP)$$

- F-Measure: Derives from precision and recall values:

$$\text{FMeasure (\%)} = (2 \times \text{recall} \times \text{precision}) / (\text{recall} + \text{precision})$$

## VI. RESULTS AND ANALYSIS

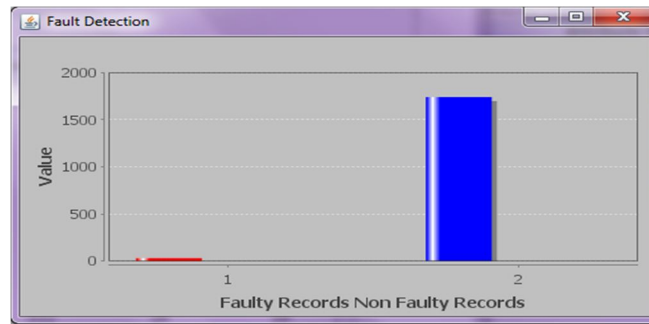


Figure 4 Faulty Detection Rate



Figure 5 Error Rate

## VII. CONCLUSION

We proposed a Machine Learning clustering and Metrics threshold-based software VULNERABILITIES prediction approaches for the cases where there is no priori VULNERABILITIES data. Attribute selection method is used to select a weighted attribute and returns a most weighted attribute only. After the clustering method the data points are clustered and each cluster have the own centroid value. Metric Threshold was used to define the threshold and it's compared to each cluster centroid. Finally, VULNERABILITIES y and non-VULNERABILITIES y data points was displayed.

## REFERENCES

- [1] Saikat Chakraborty;Rahul Krishna;Yangruibo Ding;Baishakhi Ray "Deep Learning based Vulnerability Detection: Are They have There Yet" IEEE Transactions on Software Engineering Year: 2021 | Early Access Article | Publisher: IEEE
- [2] Yang Xiao;Zhengzi Xu;They haveithey havei Zhang;Chendong Yu;Longquan Liu;They havei Zou;Zimu Yuan;Yang Liu;Aihua Piao;They havei Huo "VIVA: Binary Level Vulnerability Identification via Partial Signature" 2021 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER) Year: 2021 | Conference Paper | Publisher: IEEE
- [3] Yunhui Zheng;Saurabh Pujar;Burn Lewis;Luca Buratti;Edward Epstein;Bo Yang;Jim Laredo;Alessandro Morari;Zhong Su "D2A: A Dataset Built for AI-Based Vulnerability Detection Methods Using Differential Analysis" 2021 IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP) Year: 2021 | Conference Paper | Publisher: IEEE
- [4] Binh Hy Dang "A Practical Approach for Ranking Software Warnings from Multiple Static Code Analysis Reports" 2020 SoutheastCon Year: 2020 | Volume: 2 | Conference Paper | Publisher: IEEE
- [5] José D' Abruzzo Pereira "Techniques and Tools for Advanced Software Vulnerability Detection" 2020 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW) Year: 2020 | Conference Paper | Publisher: IEEE
- [6] "Software VULNERABILITIES Prediction and Defect Estimation Using Machine Learning and KMedoids Algorithm" – V. Bhattacharjee and P.S. Bishnu, 2011
- [7] "Application of K-Medoids with kd-Tree for Software VULNERABILITIES Prediction" – P.S. Bishnu and V. Bhattacharjee, 2011
- [8] "Outlier Detection Technique Using Quad Tree" – P.S. Bishnu and V. Bhattacharjee, 2009
- [9] "Initializing K-Means Batch Clustering: A Critical Evaluation of Several Techniques" – D. Steinley and M.J. Brusco, 2007
- [10] "Unsupervised Learning Approach to VULNERABILITIES Prediction in Software Module" – V. Bhattacharjee and P.S. Bishnu, 2010
- [11] "Complexity Metrics for Analogy Based Effort Estimation" – V. Bhattacharjee, P.K. Mohanti, and S. Kumar, 2009
- [12] "Analyzing Software Measurement Data with Clustering Techniques" – S. Zhong, T.M. Khoshgoftaar, and N. Seliya, 2004
- [13] "Clustering and Metrics Threshold Based Software VULNERABILITIES Prediction of Unlabeled Program Modules" - C. Catal, U. Sevim, and B. Diri, 2009





10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)