

Software data Reduction Orders with Bug Prediction using Machine Learning

Patnala Lavanya¹, Murali Krishna Vasantha²

¹M.Tech Scholar(CSE) in Department of Computer Science and Engineering,

²Assistant Professor, Department of Computer Science and Engineering, Avanathi Institute of Engineering and Technology, Vizianagaram (Dist), Andhra Pradesh, India.

Abstract

Software Bug Prediction is an esteemed issue that has occurred in software development and within the process of maintenance, which concerns with the success of overall software. This can be because in earlier innovate prediction of the software faults improves the software quality, reliability, efficiency and reduces the software cost. Bug-prediction techniques are aiming towards prediction of the software modules that are faulty in order that it can be beneficial within the upcoming phases of software development. Difference performance criteria are being employed so as to spice up the performance of the already existing ways. However, the most in performing the prediction of the software faults is ignored constantly. Classification is that the most used technique that's getting used for the exclusion of faulty from non-faulty modules. The work which is completed previously under this subject has been applied using different techniques. However, it's challenging task in developing robust bug prediction model and lots of approaches are proposed within the literature. This project represents a software bug prediction model by machine learning (ML) algorithms.

Keywords: Bug Prediction, Data Reduction, Instance Selection, Feature Selection, Data Mining

I. Introduction

Information mining innovation is used as a piece of programming advancement measure can't simply fabricates the exactness and satisfaction of programming improvement yet moreover grows the reliability of the product. A product bugs is an error or fault in a PC program or system that reasons it to make a misguided or unexpected outcome. Most bugs start from oversights and errors made by people in either a program's source code or its blueprint and a couple are initiated by compilers conveying wrong code. Reports as for bugs in a program are generally called as bug reports or accuse reports. The guideline

space of mining bug storage facilities which has destinations to use information mining to oversee programming building issues. Programming vaults have generous scale information bases which are used for taking care of the yield of programming improvement. When in doubt for tremendous scope and complex information in programming storage facilities, programming examination isn't thoroughly suitable. So information mining strategies, mining programming vaults can discover fascinating information in programming documents and deal with real programming issues [1]. A bug chronicle is generally called programming vault which is used taking care of information of bugs. The bug expectation is fundamental steps for settling a bug which are fittingly giving out an engineer to new bug. For open source huge scope programming expands, the amount of ordinary bugs is so enormous which makes the triaging cycle hard and testing. Programming associations pay most of cost in settling bugs. In a bug store, a bug is kept up as a bug report, which records the abstract portrayal of impersonating the bug and updates as demonstrated by the status of bug settling. In bug vault, bug reports are called bug information. There are two driving difficulties in programming improvement identified with bug information that may envision on bug stores that are the goliath scale and the low Superiority. Because of day-today declared bugs, enormous number of new bugs is taken care of in bug storage facilities. Additionally, inferior quality bugs are loud and abundance. Disorderly bugs may futile information that are associated creators and dull bugs infers a comparative quality may have unmistakable name in different data set. They discarded confined season of bug managing [1]. The central target of information diminishing for bug forecast to build up a little scope and high predominance set of bug information by ousting bug reports and words which are dreary or not-accommodating. So case choice and highlight choice methods are used meanwhile to reduce the bug estimation and the word estimation.

The lessened bug information have humble number of bug reports and more unassuming number of words than one of a kind bug information. Furthermore, they also give intently taking after information than one of a kind bug information. The case choice suggests subset of related occasions for example bug report in bug information and the component choice infers subset of related highlights for example words in bug information [1].

II. Related work

Sandusky, Gasser and Ripoche gave the element of the distortion following stores, which showed the evolvement of the bug report sort out which demonstrates formal and easygoing relationship in the bug information and moreover to review the dependence among the bug report [2]. Q. Hong, S. Kim et.al gave theory for seeing tremendous programming advancement and backing demand participation of social affair of engineers for improving turn of events and upkeep quality and diminishment cost[3]. J. Xuan, H. Jiang et.al perceived the engineer prioritization which can perceive architects and helps task in programming upkeep [4]. Thomas Zimmermann and Silvia Breu suggested that various resulting requests are ought to have been acted to the writers of bug. Along these lines, there is high need of capable and effective correspondence with the gatherings in open source adventure. These gave the high consequences of bug settling practices and have updates of the bug. Compromise and dynamic participation of customers in bug following will bring about powerful bug tracking.[11]. The product improvement, bug gives crucial information to architect. Nevertheless, they contrast in quality. Zimmerman and Battenberg perceived the idea of bug information, the sketching out overviews to engineers and customers. Considering this studies, we can depict what impacts a nice bug to report and what classifier are to be set up to recognize the idea of bug and methods to improve them[5]. To recognize the duplicate bug report which cripples the idea of bug, Sun and Jaing portrayed duplicate bug area approach by smoothing out a recuperation deal with different highlights [6]. D. Cubranic and G. C. Murphy handled the issue of consigning content document into at any rate one

orders or classes, we apply the substance plan on this [7]. J. Anvik and G. C. Murphy recommended applying bug forecast which upgrades the product advancement measure. An expectation chooses, on the off chance that the report are huge, these report are then made for blend out of the endeavor's advancement cycle [8]. A. Lamkanfi, S. Demeyer foreseen that the reality of nitty gritty bug is a fundamental factor and picks that soon it should be settled. We can do this by printed depiction using content mining algorithms[9]. Rogati and Yang uncovered a four doubtlessly comprehended plan computations a Naïve Bayesian(NB) approach, a Rocchio-style classifier, a k-Nearest Neighbour(kNN) method and a Support Vector Machine(SVM) system. This gave examination of four gathering computations and planned the execution measures[9]. Above outcome gave that the Naïve Bayes portrayal improves with more diminutive planning set. The informational index was isolated into a test set and get ready set by self-assertively picking a degree of a records from the informational index. The informational collection are to be placed into get ready set and game plan is performed. J.A. OlveraLopez and J.A. Carrasco evaluated with the occurrence determination counts and we have picked Learning Vector Quantization for the case choice [10]. We address the issue of information diminishment for bug expectation. Rather than bug expectation, distortion gauge is an equal plan issue, which means to foresee a product knick-knack contains flaws as demonstrated by removed highlights. In this way, T. M. Khoshagotkar, K. Gao, N. Seliya take a gander at the strategies on component determination to manage imbalanced information.

III. Information Reduction for Bug Prediction

Here, we present a calculation which tells the best way to apply case determination and highlight choice which is information decrease of bug expectation.

- Algorithm for Data Reduction Data decrease dependent on FS->IS : Input: preparing set T with n words and m bug reports, decrease request FS->IS last number nF of words, last number mI of bug reports, Output: diminished informational index T FI for bug forecast

- 1) Calculate target esteems n for all the words by applying FS to n expressions of T;
- 2) Generate a preparation set TF by choosing top nF words;
- 3) Apply IS to mI bug reports of TF;
- 4) When the quantity of bug reports is equivalent to or not as much as mI , end IS and produce the last preparing set TFI.

A. Applying example determination and highlight choice

In bug forecast, a bug informational index is changed over into a book network with two measurements which are known as the bug measurement and the word measurement. In our work, we exploit the blend of example determination and highlight choice procedure which will create a decreased bug informational collection. We at that point supplant the first informational collection with the diminished informational index for bug forecast. These procedures are generally famous and much of the time utilized in information preparing. For a given informational collection, example choice is to give a subset of important cases (i.e., bug reports in bug information) while highlight determination intends to give a subset of significant highlights (i.e., words in bug information). We give the accompanying meaning to recognize the sets of applying occurrence determination and highlight choice. Given an occurrence choice calculation IS and an element choice calculation FS, we use FS->IS to signify the bug information decrease, which first we need to apply FS and afterward IS; then again, IS->FS means initially applying IS and afterward FS. In Algorithm, we quickly present how to lessen the bug information dependent on FS-> IS. Given a bug informational index, the yield of bug information decrease is another and diminished informational index. Two calculations FS and IS are applied successively. That is in Step 2, some of bug reports might be clear during highlight choice since all the words in a bug report are eliminated. Such clear bug reports are additionally taken out in the component choice. Here, FS-> IS a lot - > FS are seen as two sets of bug information decrease. We present these calculations as follows:

B. Occurrence Selection

Occurrence choice is a method to decrease the quantity of cases by eliminating loud and repetitive occasions. Decreased informational index by eliminating non-delegate occurrences. There are four occurrence determination calculations, specifically Iterative Case Filter (ICF), Learning Vectors Quantization (LVQ), Decremental Reduction Optimization Procedure (DROP), and Patterns by Ordered Projections (POP).

C. Highlight Selection

Highlight determination is a preprocessing procedure for choosing a decreased arrangement of highlights for huge scope informational indexes. The diminished set is considered as the agent highlights of the first list of capabilities. We center around the component determination calculations in content information. There are four well calculations in content information and programming information, Information Gain, x_2 measurement, Symmetrical Uncertainty characteristic assessment (SU), and Relief-F Attribute determination (RF). Bug reports are arranged by their component esteems additionally on given number of words with enormous qualities is chosen as delegate highlights dependent on element choice.

D. Advantage of Data Reduction

We have two advantages of the Data Reduction technique specifically, diminishing the information scale and improving the precision.

We spare the work cost of designers by utilizing lessening the size of information. Bug measurement: The point of bug forecast is to dole out designers for bug fixing. Word measurement: We use highlight choice to eliminate loud or copy words in an informational collection. The decreased informational collection can be taken care of all the more effectively via programmed strategies (e.g., bug expectation draws near) than the first informational collection dependent on element determination. Other than bug forecast, the decreased informational index can be further additionally be utilized for other programming undertakings after bug expectation.

2. Improving the Accuracy.

Information decrease eliminates boisterous or copy data in informational indexes. This can help in improving exactness. Bug measurement: Instance choice strategy gives a simplicity in expulsion of uninformative bug reports; despite what might be expected, evacuation of bug reports may brings about diminished precision. Word measurement: Removing uninformative words can help in improving precision of bug forecast. This can recuperate the exactness misfortune by example determination.

IV. Proposed Methodology

In this part, we present the information decrease methods to lessen sizes of bug informational collection. The fundamental objective of our work is to join the occurrence determination and highlight choice in right request to eliminate the boisterous, repetitive and non-useful bug reports.

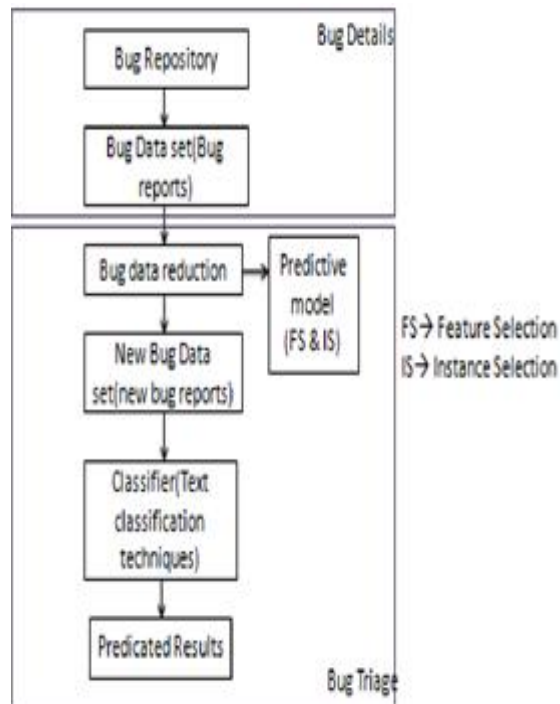


Fig. 1 System Architecture

Bug Details

The bug subtleties comprise of bug vault and bug reports. In a bug storehouse, a bug is continued as a

bug report, which follows the literary delineation to rehash the bug and updates as per the status of bug fixing.

Bug Repository

A bug storehouse is a regular programming vault, for putting away subtleties of bugs, e.g., a mainstream and open source bug archive, Bugzilla [2]. Large programming ventures convey bug stores is likewise called as bug or issue global positioning frameworks, which is utilized to help data assortment and to help designers to deal with bugs. Each bug is kept up as a bug report, which follows the narrative depiction of duplicating the bug and overhauls as indicated by the centrality of bug fixing. The utilization of bug vault can improve the advancement cycle and nature of programming created. It presents an information stage to support numerous types of task on bugs, e.g., imperfection forecast, bug confinement and resumed bug investigation.

Bug Report

A recorded bug is known as a bug report or bug information. It has various things for enumerating the data of duplicating the bug. In a bug report, the diagram and the report are two key things about the data of the bug, which are followed in characteristic dialects. Rundown indicates the overall assertion for distinguishing a bug and depiction gives the subtleties to imitate the bug [2]. The bug report may likewise contain different things additionally, for example, Product, Platform, and Importance.

Bug Prediction

The strategy for assigning a right engineer for redesigning the bug is called bug forecast. When the bug report is framed, a bug prediction apportions the bug to an engineer who can fix this bug and designer is recorded in a thing doled out to with no throwing.

Bug Data Reduction

By utilize the gathering of highlight determination and occasion choice calculations to dispose of undesirable and non-useful bug reports. With the involvement with text classification techniques, a case in bug forecast determines bug reports while a component in bug expectation shows the bug words.

The fundamental objective of our work is to diminish the content grid with two measurements to be specific, bug report measurement and word measurement.

The two-stage mix of case choice and highlight choice calculations are utilized to decrease the bug informational collection on two measurements i.e., bug report measurement and word measurement. To decide the request for information decrease method, the prescient model is applied by the proposed framework. This model assists with anticipating the right request i.e., FS to IS or IS to FS to diminish the work cost and time cost. The decreased bug information contain less bug information than the first bug information and gives related data over the first bug information.

Highlight Selection

Highlight determination is a pre-preparing strategy for picking a reduced arrangement of highlights for colossal scale informational collections [4,7]. The pre-handling strategies are tokenization, stop word expulsion, stemming cycle and vector space model. The tokenization technique is utilized to tokenize the rundown and depiction of the bug reports into word vectors. Nonalphabetic words and uncommon character are eliminated to evade the loud bug words. Stop word evacuation strategy eliminate the stop words in high recurrence and give no accommodating data to bug forecast. Stemming strategy utilizes watchman stemming calculation for lessening curved words their statement stem/root structure. Vector space model/Term vector model is a logarithmic model for speaking to message archive as vector of identifier. The limited set is considered as the delegate highlights of the first component set[5].

The four all around performed calculations are picked in content information [3, 9] and programming information, specifically Information Gain (IG), 2 measurement (CH)[4], Symmetrical Uncertainty trait assessment (SU), and Relief-F Attribute determination (RF).

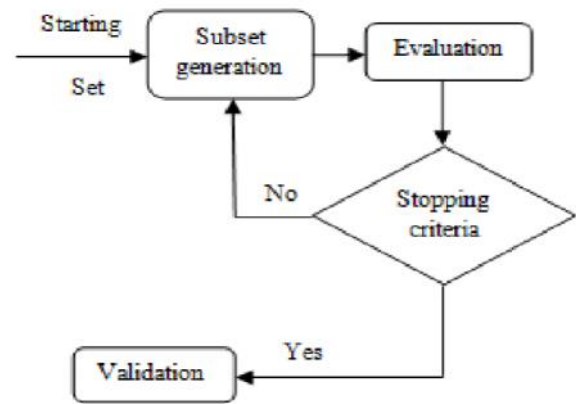


Fig. 2 General Feature Selection Structure

Based on feature selection, words in bug reports are organized according to their feature importance and a given number of words with large values are selected as representative features.

The chi-squared distribution also known as chi-square or χ^2 distribution with k degrees of freedom is the distribution of a sum of the squares of k independent criterion normal random variables. It is a unique case of the gamma distribution and the most widely used probability distributions in inferential statistics. If Z_1, \dots, Z_k are independent, standard normal random variables, then the sum of their squares, Q

$$Q = \sum_{i=1}^k Z_i^2 \quad (1) \quad i=1$$

is distributed according to the chi-squared distribution with k degrees of freedom. This is usually denoted as

$$Q \sim \chi^2(k) \text{ or } Q \sim \chi_k^2 \quad (2)$$

where k is a positive integer that specifies the number of degrees of freedom (i.e. the number of Z_i 's). The Chi-squared attribute evaluation evaluates the worth of a feature by computing the importance of the chi-squared gauge with respect to the class. The initial hypothesis H_0 is the assumption that the two features are dissimilar and it is checked by chisquared formulae:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad (3)$$

where O_{ij} is the observed frequency and E_{ij} is the expected (theoretical) frequency, asserted by the null hypothesis.

V. Results and Discussion

The neural network is implemented using Keras[2] running on jupyter notebook. The Keras Sequential model is set up as shown.

```

In [83]: classifier.add(Dense(output_dim = 6, init = 'uniform', activation = 'relu'))

<jupyter-input-89-fc312ea702>:1: UserWarning: update your 'Dense' call to the Keras 2 API: 'Dense(activation='relu', units=6, kernel_initializer='uniform')
classifier.add(Dense(output_dim = 6, init = 'uniform', activation = 'relu'))

In [90]: # Adding the output layer
classifier.add(Dense(output_dim = 1, init = 'uniform', activation = 'sigmoid'))

<jupyter-input-98-6c7181af53bb>:2: UserWarning: update your 'Dense' call to the Keras 2 API: 'Dense(activation='sigmoid', units=1, kernel_initializer='uniform')
classifier.add(Dense(output_dim = 1, init = 'uniform', activation = 'sigmoid'))

In [95]: classifier.compile(optimizer = 'adam', loss = 'binary_crossentropy', metrics = ['accuracy'])

In [94]: classifier.fit(X_train, y_train, batch_size = 10, epochs = 100)

Epoch 1/100
5792/5792 [#####] - 2s 170ms/step - loss: 0.4204 - accuracy: 0.8282
Epoch 2/100
5792/5792 [#####] - 2s 169ms/step - loss: 0.4279 - accuracy: 0.8284
Epoch 3/100
5792/5792 [#####] - 2s 169ms/step - loss: 0.4288 - accuracy: 0.8191
Epoch 4/100
5792/5792 [#####] - 2s 169ms/step - loss: 0.4283 - accuracy: 0.8196
Epoch 5/100
5792/5792 [#####] - 2s 165ms/step - loss: 0.4288 - accuracy: 0.8283
Epoch 6/100
5792/5792 [#####] - 2s 172ms/step - loss: 0.4288 - accuracy: 0.8283
Epoch 7/100
5792/5792 [#####] - 2s 181ms/step - loss: 0.4283 - accuracy: 0.8195
Epoch 8/100
5792/5792 [#####] - 2s 188ms/step - loss: 0.4278 - accuracy: 0.8212
Epoch 9/100
5792/5792 [#####] - 2s 181ms/step - loss: 0.4279 - accuracy: 0.8196
Epoch 10/100

```

The accuracy obtained on the above fold was 82.36% with a good AUC of 0.8. The ROC Curve is given below

NAÏVE BAYES,DECISION TREE,KNN AND SVM

```

In [85]: # sklearn component analysis
from sklearn.decomposition import PCA
pca = PCA(n_components=4)
X_train_pca = pca.transform(X_train)
X_test_pca = pca.transform(X_test)
explained_variance = pca.explained_variance_ratio_

In [14]: # fitting svm model with 4 support vectors
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier()
classifier.fit(X_train, y_train)

Out[14]: KNeighborsClassifier()

In [15]: # predicting output
y_pred_svm = classifier.predict(X_test)

In [16]: from sklearn import metrics
cm = metrics.confusion_matrix(y_test, y_pred_svm)
accuracy_score_svm = metrics.accuracy_score(y_test, y_pred_svm)
print(cm)
print("accuracy_score", accuracy_score_svm)

[[ 264   72]
 [  58  508]
 accuracy_score 0.96236116496882

```

```

In [80]: import numpy as np
import pandas as pd

In [82]: # Importing the dataset
dataset = pd.read_csv("90.csv")
X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, 2].values

In [81]: from sklearn.preprocessing import LabelEncoder
labelencoder_y = LabelEncoder()
y = labelencoder_y.fit_transform(y)

In [84]: # Splitting the dataset into the training set and test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.1, random_state = 0)

In [86]: # Feature scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

In [ ]: import keras
from keras.models import Sequential
from keras.layers import Dense

using TensorFlow backend.

In [87]: classifier = Sequential()

In [88]: classifier.add(Dense(output_dim = 6, init = 'uniform', activation = 'relu', input_dim = 21))

```

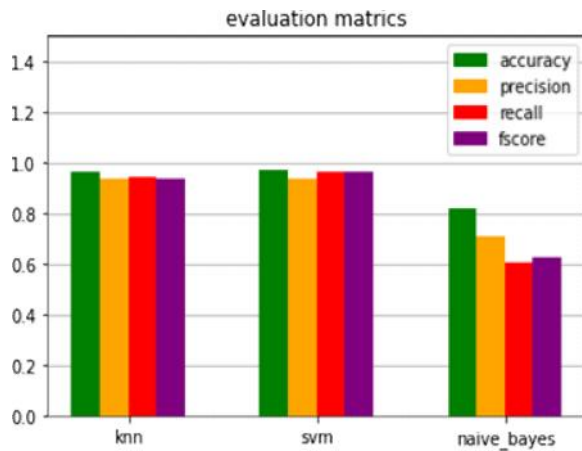


Fig Bar graph of KNN,SVM, Naive Bayes metrics

VI. Conclusion

The outcomes show that the composite proportion of accuracy and review (F-measure) is the most elevated in the proposed model for both the preparation and the test informational collections, in correlation with other existing models. This model, when coordinated with the product advancement conditions, will prompt sensibly precise expectations. The principle target of this postulation is to propose a totally novel methodology for improving the exhibition of programming bug forecast models. The methodology is "AI for programming bug forecast". In particular, we apply ML calculations for improving grouping ability of classifiers like k-closest neighbor. At the start, a logical audit has been led to accomplish an understanding into the writing, and to break down whether our methodology is conceivable. The audit gives subtleties of different kinds of programming measurements, generally utilized AI calculations, information preprocessing strategies and assessment measures inside the field of programming bug forecast. Additionally, the audit presents assortment of remarkable examinations for building programming bug forecast models that may effectively aid programming exertion assessment. A significant number of these examinations utilize programming measurements and distance-based characterization calculations like k-closest neighbor with Mahalanobis distance for mining informational indexes from programming storehouses. Be that as it

may, the normal executions of the Mahalanobis distance work don't utilize the name data from the preparation information. Inspired by this issue, we have chosen to propose the methodology of utilizing AI, which utilizes the mark data, to upgrade the presentation of programming bug indicators. In the wake of presenting general thoughts of AI and giving subtleties of our methodology, we have done the examinations followed by the investigation and contrasted the calculations with confirm the adequacy of the methodology. Right off the bat, a trial system with various learning plans was planned. At that point, the analyses are led on NASA informational index created from PROMISE programming designing store. Results from the examinations were assessed utilizing equilibrium and f-measure. At last, the examination and correlations are allotted to show that the system making bug.

References

- [1] Jifeng Xuan, He Jiang, Yan Hu, Zhilei Ren, Weiqin Zou, Zhongxuan Luo, and Xindong Wu, "Towards Effective Bug Prediction with Software Data Reduction Techniques" IEEE transactions on knowledge and data engineering, vol. 27, no. 1, January 2015.
- [2] R. J. Sandusky, L. Gasser, and G. Ripoché, "Bug report networks: Varieties, strategies, and impacts in an F/OSS development community," in Proc. 1st Intl. Workshop Mining Softw. Repositories, May 2004.
- [3] Q. Hong, S. Kim, S. C. Cheung, and C. Bird, "Understanding a developer social network and its evolution," in Proc. 27th IEEE Int. Conf. Softw. Maintenance, Sep. 2011.
- [4] J. Xuan, H. Jiang, Z. Ren, and W. Zou, "Developer prioritization in bug repositories," in Proc. 34th Int. Conf. Softw. Eng., 2012.
- [5] T. Zimmermann, R. Premraj, N. Bettenburg, S. Just, A. Schröter, and C. Weiss, "What makes a good bug report?" IEEE Trans. Softw. Eng., vol. 36, no. 5, Oct. 2010.
- [6] C. Sun, D. Lo, S. C. Khoo, and J. Jiang, "Towards more accurate retrieval of duplicate bug reports," in

Proc. 26th IEEE/ACM Int. Conf. Automated Softw. Eng., 2011.

[7] D. Cubrani_c and G. C. Murphy, "Automatic bug prediction using text categorization," in Proc. 16th Int. Conf. Softw. Eng. Knowl. Eng., Jun. 2004.

[8] J. Anvik and G. C. Murphy, "Reducing the effort of bug report prediction: Recommenders for development-oriented decisions," ACM Trans. Soft. Eng. Methodol., vol. 20, no. 3, article 10, Aug. 2011.

[9] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals, "Predicting the severity of a reported bug," in Proc. 7th IEEE Working Conf. Mining Softw. Repositories, May 2010.

[10] T. M. Khoshgoftaar, K. Gao, and N. Seliya, "Attribute selection and imbalanced data: Problems in software defect prediction," in Proc. 22nd IEEE Int. Conf. Tools Artif. Intell., Oct. 2010.

Authors



Patnala Lavanya Holds a B.Tech Degree in Computer Science & Engineering from Coastal Institute of Technology And Management, Narapam, Veerabhadrapuram (P), Kothavalasa(Mandal),

Vizianagaram, Andhra Pradesh 535183. She is presently Pursuing M.Tech (CSE) in Department of Computer Science and Engineering from Avanthi Institute of Engineering and Technology Cherukupalli(V), Near Tagarapuvalsa Bridge,Vizianagaram (Dist) - 531162,Andhra Pradesh. Area of interest include Database Programming, Web Technologies, web programming and C programming, Machine learning and Data science.



Mr.Murali Krishna Vasantha Working as Asst. Professor, Department Of Computer Science and Engineering in Avanthi Institute of Engineering and Technology Cherukupalli(V), Near

Tagarapuvalsa Bridge,Vizianagaram (Dist) - 531162,Andhra Pradesh. He has More than 11 Years of Teaching Experience in Various Colleges in Andhra Pradesh. His Area of interests include Machine Learning, Artificial intelligence,Blockchain Technology,Web Programming,Data mining with Outlier Detection and Database Programming.