

# Evolutionary Computing based an Efficient and Cost Effective Software Defect Prediction System

Racharla Suresh Kumar<sup>1</sup> and Prof. Bachala Sathyanarayana<sup>2</sup>

<sup>1</sup> Sri Krishnadevaraya University

*Received: 13 April 2015 Accepted: 3 May 2015 Published: 15 May 2015*

## Abstract

The earlier defect prediction and fault removal can play a vital role in ensuring software reliability and quality of service. In this paper Hybrid Evolutionary computing based Neural Network (HENN) based software defect prediction model has been developed. For HENN an adaptive genetic algorithm (A-GA) has been developed that alleviates the key existing limitations like local minima and convergence. Furthermore, the implementation of A-GA enables adaptive crossover and mutation probability selection that strengthens computational efficiency of our proposed system. The proposed HENN algorithm has been used for adaptive weight estimation and learning optimization in ANN for defect prediction. In addition, a novel defect prediction and fault removal cost estimation model has been derived to evaluate the cost effectiveness of the proposed system. The simulation results obtained for PROMISE and NASA MDP datasets exhibit the proposed model outperforms Levenberg Marquardt based ANN system (LM-ANN) and other systems as well. And also cost analysis exhibits that the proposed HENN model is approximate 21.66

**Index terms**— software defect prediction, artificial neural network, adaptive genetic algorithm, levenberg marquardt, object oriented software metrics, cost estimat

## 1 I. Introduction

ith the increase in information technologies and associated software applications, the inevitable requirement of software reliability has alarmed scientific societies, industries as well as academicians to develop certain optimal paradigm to ensure defect free software applications for long run reliability.

Furthermore, the cost factor for software products and services also suggests the defect free software solutions, so as to eliminate probability of faults in future and iterative maintenance. In order to accomplish these objectives, the efficient software defect prediction (SDP) systems are of great significance. In order to ensure optimal software reliability, the defect prediction has become an inevitable part of software development life cycle (SDLC) that intends to eliminate the probability of software failure in run time. The earlier defect prediction can enable software professional to identify fault-prone modules and thus can debug the defects to ensure quality of service provisioning. In recent years the application of open source software has increased tremendously and professional prefer to customize software modules and implement as per need. Still, these modules are prone to defect in real time scenarios, thus demanding for fault prediction and verification [1, ??,3,4] before introducing product to the users. The SDP might be functional on the basis of certain software metrics [3,4,5] like changes in source code, earlier defect or fault etc. Typically, software metrics do represent certain quantitative factor that characterizes the properties of software source code, which can be employed to predict fault proneness of software during function. On the other hand, in recent years majority of software applications are being developed using Object-Oriented (OO) paradigm. The object oriented paradigm enables certain metrics that that can be employed to examine the quality of software application and associated fault proneness. Some of the predominantly proposed software metrics are MOOD [6], QMOOD [7], Bieman and Kang [8], Briand et al. [9],

Etzkorn et al. [10], Halstead [11], Henderson-sellers [12], L and H metrics suite [13], McCabe [14], Tegarden et al. [15], Lorenz and Kidd [16] and CK metric suite [17]. The implementation of object oriented metrics enables software practitioners to examine quality of software in terms of precision, accuracy, fault-resilience, reliable functionality, adaptability, supportability, usability, portability, and cost effectiveness etc. In fact, it makes testing enhanced for large scale software applications. This is the matter of fact that a number of researches have been made for defect prediction. Some of the predominantly employed SDP techniques are based on machine learning and artificial neural network [18,19,20, 21,22], clustering techniques, statistical method, data mining based fault identification, random forest [23,24,25] approaches etc. However, the emerging software complexities, critical software applications, reliable service assurance, quality oriented service provisioning, and cost effective or economical solutions etc., motivate researchers to develop certain cost effective defect prediction solution. In recent years, primarily, support vector machine (SVM) and artificial neural network (ANN) approaches are being explored for SDP utilities. The emergence of artificial intelligence based applications have motivated researchers to explore ANN based defect prediction that works based on the human brain functions, while encompassing multiple neurons and directed edges possessing certain weights values between input and output layers. In fact, ANN is a complex non-linear mapping process that employs output as the input for learning certain complex non-linear input-output relationship between input and output layers. In function ANN encompasses data sets to optimize key factors such as weight parameters, risk minimization mechanism for stopping training once the learning error enters in expected margin level. Although, ANN has established itself as a potential candidate for prediction and classification applications, still its limitations in terms of slow learning ability, local minima and convergence can't be ignored. In order to enhance the performance of ANN based defect prediction some researchers [26,27] have suggested evolutionary computing paradigm that could enable optimal classification and prediction without introducing any computational complexity and premature convergence.

Considering efficiency of evolutionary computing techniques such as Genetic Algorithm (GA) in this paper a robust Adaptive genetic algorithm based ANN learning algorithm has been developed, which has been used for software defect prediction. In addition, to enhance the performance of GA for huge data elements and efficient performance, the genetic parameters (crossover and mutation probability) have been selected dynamically that makes overall system much robust as compared to conventional approaches. In order to examine the performance of the proposed HENN system, a Levenberg Marquardt based ANN (LM-ANN) algorithm has been developed and the comparative performance analysis with the object oriented software metrics, CK metrics [17] has revealed that the proposed HENN algorithm provides better fault detection as compared to LM-ANN. Furthermore, the fault removal cost analysis for both the algorithms has stated that the proposed system is cost effective and can be used for real time defect prediction utilities.

The remaining sections discusses, related work in Section II, the research contributions and problem definitions for the proposed software defect prediction model are presented in III, which has been followed by proposed HENN and LM-ANN based SDP model discussion and implementation in Section IV. Section V presents the results and analysis and conclusion has been discussed in Section VI. The references used in this paper are given at the last of the manuscript.

## 2 II. Related Work

Software reliability is of course an inevitable need for quality service provisioning. The reliability oriented software defect prediction (SDP) has motivated researchers to develop optimal system for cost efficient defect prediction. Researchers examined the relationship between object oriented software metrics and associated faults [28,29,30,31,32,33] by means of machine learning algorithms and detected fault proneness of software. To achieve better prediction some other approaches such as decision trees, naïve Bayes, and 1-rule [34] based fault detection scheme were developed, where the standard datasets such as NASA MDP was used to examine classification accuracy of the SDP approaches. Chug et al [35] demonstrated fault identification using data mining and employed conventional J48, Random Forest, and Naive Bayesian Classifier (NBC) schemes for performance comparison but still couldn't employ the benefits of advanced classification approaches. To optimize conventional random forest based defect prediction Pushphavathi et al [36] incorporated a hybrid random forest (RF) and Fuzzy C Means (FCM) clustering model for software defect prediction. Unfortunately, these approaches could not address the issue of unbalanced datasets, which motivated researchers to come up with Adaboost. Nc [37] which implemented a number of class imbalance approaches, re-sampling, threshold variations, and ensemble algorithms. Exploring insight, this approach can be found to be complicate and not a cost effective solution for large scale dynamic data. Researchers used SVM based defect prediction scheme [38,39] and a dynamic SVM model was proposed that intended to detect faults in source code by means of error data and faulty code execution. In [40,41]an ANN based defect prediction model was developed. A defect severity model using conventional back-propagation learning based ANN was developed in [42]. Similarly in [43] a Radial Basis ANN was used for SDP. ANN based SDP for Halstead data metrics has been done in [44]. In [45] the Bayesian Regularization (BR) technique based ANN model was developed for software fault detection. Almost all ANN based defect prediction model employs conventional learning and weight estimation techniques that confines applicability with huge datasets with dynamic functional environment. The conventional learning and weight estimation approaches can't eliminate the key issues of local minima and convergence issue that limit the performance of generic ANN. The enhancement of learning scheme and further optimization through certain evolutionary computing approaches can make ANN robust for SDP

---

applications. In fact, cost feasibility is one of the key factors that decide employability of certain SDP model, but till no any research work has addressed the issue of cost estimation of the defect prediction model. This paper has considered these limitations as motivation and has developed an evolutionary computing A-GA based SDP model which has been compared with Levenberg Marquardt based ANN and respective fault removal cost estimation has been done.

### III.

## 4 Our Contribution

In SDLC the fault resilience and reliability is of great significance. The implementation of efficient SDP strengthen early fault detection and thus it enables software practitioner to remove faults to ensure reliability and QoS of the software solution. The predominant question in this paper is whether the implementation of Adaptive Genetic Algorithm can enable efficient and cost effective SDP solutions? In this paper, object oriented software metrics [17] has been considered for defect prediction and using proposed SDP models, the fault proneness of metrics data has been retrieved, whether the data is faulty or non-faulty. In order to perform classification of faulty and non-faulty data, initially the conventional ANN learning scheme with Leven berg Marquardt (LM) algorithm [45] has been developed and respective performance towards software defect prediction with NASA defect datasets has been done. This is the matter of fact that LM based ANN performs better as compared to other approaches such as back-propagation or feed-forward learning based NN, still it suffers due to prime limitations of ANN, such as local minima and weight update issues. Thus, considering higher employability of artificial intelligence techniques and respective limitations for critical software applications, in this paper an evolutionary computing based optimization scheme called Genetic Algorithm has been used for weight estimation during ANN learning. Further to ensure optimal performance of GA, in this paper a novelty has been introduced in terms of adaptive GA parameter (Crossover and Mutation probability) selection. The proposed Adaptive Genetic Algorithm (A-GA) performs adaptive weight estimation and learning optimization so as to ensure optimal fault classification and accuracy. The A-GA optimization scheme alleviates the issue of premature convergence and local minima. Such enhancement has lead better classification and accuracy for fault detection in huge datasets.

In order to examine the performance of the proposed SDP model, the object oriented software metrics (here, CK metrics [17]) has been considered. The implemented metrics characterizes various software features. In this paper, six predominant software metrics have been considered in fault identification. The considered metrics are WMC, NOC, DIT, CBO, RFC, and LCOM. The individual metrics has been feed as the input of the ANN and performing learning with the proposed HENN model the classification for faults has been done. The discussion of the proposed A-GA based ANN (HENN) has been discussed in the next section of the presented manuscript. In this paper, in order to examine the cost effectiveness of the developed SDP models, certain cost efficiency model can be used [46, 47, and 48] and with certain standard threshold the applicability of the proposed SDP model for large scale software data can be examined. The performance analysis of the proposed model has been done in terms of accuracy, precision, recall, F-Measures and fault removal cost efficiency. The discussion of the proposed SDP models and its implementation is discussed in the following sections.

### IV.

## 6 System Model

In this section, the proposed Levenberg Marquardt learning based ANN and our proposed HENN based software defect prediction schemes and its algorithmic implementation have been discussed.

### a) Artificial Neural Network based Software Defect

Prediction This is the matter of fact that the Artificial Neural networks (NN) have seen an explosion of interest over the years, and it has been implemented across a range of problem domains, specifically classification and prediction. In fact, the major problems dealing with prediction and classification, ANN is considered to be the dominating solution. For SDP scenario, ANN can be used with different learning schemes like Gradient Descent (GD), Gauss Newton, and Levenberg Marquardt (LM) etc. Unfortunately majority of existing learning paradigm are ineffective to alleviate the key limitations of ANN such as local minima and convergence issue. Even though, researches have revealed that Levenberg Marquardt (LM) can be a potential candidate for ANN learning due to its stable nature and flexible implementation. In this paper, in addition to LM-ANN algorithm, an evolutionary computing technique called Adaptive Genetic Algorithm (A-GA) has been used for dynamic weight estimation for prediction enhancement. In the proposed ANN model and ultimately intended SDP system, it has been intended to find relation between object oriented software metrics and fault prone classes of the six CK metrics; WMC, NOC, DIT, RFC, CBO, LCOM, which has been considered as independent variable. The fault data has been taken as the dependent data. Figure -1 illustrates the architecture of our proposed ANN model comprising three layers i.e., input layer, hidden layer and output layer. Here, 6 input nodes have been defined that takes six CK matrix [17] having multiple classes as individual input. Since, in the proposed ANN model, the expected outputs are either FAULTY or NO-FAULTY, therefore only one output node is needed. Here, we have

## 8 B) LEVENBERG MARQUARDT (LM) LEARNING BASED ANN FOR SOFTWARE DEFECT PREDICTION

considered 8 hidden layers so as to avoid unwanted computational complexity. Thus in the defined Generally, the ANN model is defined in terms of a function  $y = \delta(w^T x)$  where  $y$  states for the output vector and  $w$  and  $x$  represent the weight vector and the input vector respectively. In learning process, the weight factor  $w$  is updated iteratively so as to minimize the Root Mean Square Error (RMSE), which can be estimated by:  $RMSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  (3)

Where  $y$  depicts the actual output and  $\hat{y}$  represents the expected output.

In order to make computation efficient and to process multidimensional data with ANN, it is inevitable to perform the normalization. In the proposed ANN based SDP models; the data normalization has been done using Min-Max approach, which is discussed as follows:

i. Data normalization In this paper, normalization has been performed on the defect datasets that strengthens the proposed ANN based software defect prediction systems for better readability and classification. In the proposed SDP model, the data normalization has been done over the range of  $[0, 1]$  so as to adjust the defined range of input feature value and avoid the saturation of neurons. There a number of normalization approaches such as Min-Max normalization, Z-Score normalization and decimal scaling etc. We have normalized the defect data using Min-Max normalization scheme that performs a linear transformation on the original data and then maps individual data  $x$  of attribute  $i$  to the normalized value  $x_{norm}$  in the range of  $[0, 1]$ . The normalization using Min-Max approach has been done using following equation:  $x_{norm} = \frac{x - \min(x)}{\max(x) - \min(x)}$  (4)

where  $\max(x)$  and  $\min(x)$  are the maximum and minimum values of the attribute  $i$  respectively. In the proposed SDP model, performing data normalization the ANN model has been implemented for fault classification.

In ANN based systems, the efficient weight estimation and learning approach is of great significance. Till a number of approaches have been developed for learning optimization in ANN based artificial intelligence applications. Some of the predominant approaches are: Gauss Newton, Gradient descent, Levenberg Marquardt (LM) etc. Interestingly LM can work as both gradient descent as well as gauss Newton. Some researchers also have advocated that LM can outperform other existing learning schemes in ANN. Thus considering significance of LM for effective learning for SDP, in this paper initially LM based ANN (LMANN) has been developed for SDP model. The discussion of the proposed LMANN model for SDP application is given as follows:

### 8 b) Levenberg Marquardt (LM) Learning based ANN for Software Defect Prediction

The prime scope for ANN optimization is the enhancement of its weight estimation and respective learning optimization. Therefore, considering these factors, a number of algorithms have been proposed for weight update in ANN learning (Table ?? Levenberg Marquardt (M) algorithm performs localization of the bare minimum value of multivariate function in a repetitive manner, which is expressed as the sum of squares of non-linear real-valued functions. Similar to GD algorithm, in HENN, LM algorithm updates the weights during NN learning process. Considering the performance novelty, the proposed LM algorithm comprises the functional ability of Steepest Descent and Gauss Newton method. The proposed LM algorithm can update the weight vector by following expression:  $w_{k+1} = w_k + (J^T J + \mu I)^{-1} J^T e$  (1)

Where  $w_{k+1}$  is the updated weights,  $w_k$  is the current weights,  $I$  represents the identity or unit matrix,  $\mu$  is the Jacobian matrix and  $\mu$ , the combination coefficient is always positive. With  $\mu$  as very small, it functions as Gauss Newton method while making  $\mu$  as very large makes it functional as Gradient descent method. The Jacobian matrix derived as given as:  $J = \begin{bmatrix} \frac{\partial y_1}{\partial w_1} & \frac{\partial y_1}{\partial w_2} & \dots & \frac{\partial y_1}{\partial w_n} \\ \frac{\partial y_2}{\partial w_1} & \frac{\partial y_2}{\partial w_2} & \dots & \frac{\partial y_2}{\partial w_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial w_1} & \frac{\partial y_m}{\partial w_2} & \dots & \frac{\partial y_m}{\partial w_n} \end{bmatrix}$  (2)

Where  $w$  refers the weight counts and the input patterns are  $P$ . The output patterns are indicated by  $y$ . In the proposed SDP model, in the initial phase the LM algorithm has been used to estimate the weights for the learning scheme. Figure ?? represents the adaptive weight estimation approach using LM algorithm. The weights are updated dynamically so as to reduce RMSE and satisfying the stopping criteria, the classification has been done for fault prediction. On the basis of fault classification, the confusion matrix has been obtained which has been employed further to examine performance of the proposed SDP model. This is the matter of fact that LM-ANN has been employed for varied classification utilities but considering the specific requirements of fault prediction and robust function with huge data sets in real time software utilities, the local minima problem and convergence issues of ANN can't be ignored. Thus, considering these limitations, in this paper, the evolutionary algorithm Adaptive-Genetic Algorithm (A-GA) has been used for parameter optimization that can strengthen the function of the proposed system to yield more precise, accurate and efficient outputs. The implementation of A-GA for ANN based SDP utility has been discussed in the following section



criteria is not achieved and the 95% chromosomes in gene pool achieves unique fitness value, as beyond it the fitness level of chromosomes get saturated.

Step 6 Fault Classification: Considering step-3, and stopping criteria, with the optimal RMSE, the final output at output layer of ANN is obtained that more than 0.5 signifies towards FAULTY class otherwise NON-FAULTY.

Step 7 Confusion Matrix: On the basis of FAULTY and NON-FAULTY label of comprising classes, a Confusion Matrix is derived that is used for performance evaluation. Thus, implementing the above mentioned approaches, the proposed HENN model performs Software Defect Prediction. This is the matter of fact that a number of SDP systems have been developed but only prediction accuracy and precision can't be the justification for a system to be employable in real time scenarios. Industries demands for certain cost effective and efficient system for defect prediction. A system with higher computational efficiency with minimal cost of fault detection and removal can be of great significance and can be suggested to be used in real time SDP applications.

Thus, considering the need of a novel cost analysis mechanism, in this paper a novel cost estimation approach has been developed which has been used to assess the computational (Fault detection and removal) cost analysis for both our proposed HENN based SDP as well as reference, LM-ANN based SDP model. The discussion of the proposed cost estimation model is given as follows:

### 11 d) Software Fault Estimation and Removal Cost analysis

In this paper, a novel cost estimation approach has been developed that estimates the cost of fault detection and removal, as the efficiency to be considered as a criterion that decides whether the system should be used or not in real time applications. The proposed cost estimation model has been derived from [46]. In the developed cost estimation approach, certain constraints have been assumed such as, varied testing phases might take different cost for certain fault removal as different softwares are developed in varied software platform and with varied development standards, and it is impractical to perform comprising unit testing on all the associated modules [47]. In the proposed cost estimation model, the identification efficiency model proposed in [48] has been incorporated that suggests following efficiencies to be used for cost estimation model. In this paper, the following notations have been used to formulate mathematical model for fault estimation and removal cost. Cost Norm =  $\frac{\text{Cost}_{\text{SDP}}}{\text{Cost}_{\text{WSDP}}}$   $\frac{\text{Cost}_{\text{SDP}}}{\text{Cost}_{\text{WSDP}}} = ? < 1$ ,  $\frac{\text{Cost}_{\text{SDP}}}{\text{Cost}_{\text{WSDP}}} > 1$  Not Suitable(11)

Here, Cost Estm SDP represents the estimated fault removal cost for software with fault prediction scheme, Cost Estm \_WSDP is the fault removal cost without using any SDP system. The variable Cost Norm refers the normalized cost with the SDP models. As illustrated in above expression, the minimal normalized cost signifies better employability of a defect prediction system. In this paper, the cost analysis for both the proposed HENN as well as Levenberg Marquardt based ANN (LMANN) has been done. The results obtained are given in Table 7.

V.

## 12 Result and Analysis

This section discusses the experimental setup, benchmark fault data, results and performance analysis.

In this paper, the overall algorithms for artificial neural network, Levenberg Marquardt based ANN, Adaptive Genetic Algorithm and its implementation with ANN for defect prediction, etc have been developed using MATLAB2012b software model. In addition, the toolboxes of machine learning and artificial neural network have been considered to perform simulation. In order to examine the performance of the proposed HENN model, object oriented software metrics suite, CK Metrics [17] has been considered, which has been derived from the fault data taken from PROMISE [49] and NASA MDP [50] fault data repository. The software metrics from the fault datasets (JEdit, Ant, Camel and IVY) have been derived using Chidamber and Kemerer Java Metrics tool (CKJM) tool that extracts software metrics by executing byte code of compiled Java cases and assigns a definite weight of the comprising classes having feature vectors. In this paper, six predominant CK metrics have been considered as depicted in the Table-4. A set of approaches that can be executed in response to a message received by an object of that class LCOM Dissimilarity measurement of varied methods in a class using instanced attributes/variables In our work, the six software metrics have been considered as the independent data while the fault data has been taken as dependent variable.

The considered data JEdit, Ant, Camel and IVY comprise static code measures along with varied modules sizes, defective modules and defect rates. In the proposed SDP models the respective extracted weights and features of the data classes have been taken as input to the ANN as illustrated in Figure -1. On the basis of final outcome of the both SDP models, LM-ANN as well as HENN for individual datasets, the confusion matrix has been obtained. A confusion matrix comprises two rows and columns representing true positive (TP), false negatives (FN), false positive (FP) and true Negative variables. The variables in confusion matrix represent the faulty and non-faulty data and its severity. As depicted in Table-5, TP depicts modules which are classified as FAULTY, FN represents the modules which are FAULTY but are classified incorrectly as NON-FAULTY. Similarly, FP represents the modules which are non-faulty but are classified as faulty.

---

## 13 a) Result Analysis

The following section represents the results obtained from the proposed HENN based SDP model and a reference model based on Leven berg Marquardt based ANN. Here, from the results obtained it can be found that the proposed HENN based SDP model performs better than Leven berg Marquardt algorithm based ANN (LMANN). Here, it can be found that the average fault prediction accuracy of the proposed HENN model is 87.23%, on contrary, the LM-ANN based SDP models delivers 75.48% and hence the proposed system outperforms the existing and till most efficient ANN model, LMANN. In addition, the analysis results states that the proposed system provides 98.2% precision, 92.74% F-measure, 88.55% of recall, which is 87.7% 85.7%, and 85.4% for LMANN based SDP system, respectively. The following figures (Figure ??-8) represent the average performance of the proposed system with four benchmark datasets (JEdit, Ant, Camel and IVY). The performance results for the developed SDP models with individual datasets are given in Table-7. Considering cost effectiveness of HENN and LMANN based SDP models, Figure ?? depicts that the proposed HENN based system is most cost efficient as compared to LMAMM, and hence it can be implemented for real time applications intending software defect prediction and removal. 7 depicts that the proposed defect prediction approach is highly robust and efficient as compared to Levenberg-Marquardt based ANN system, which is supposed to be the most effective ANN system till. The proposed HENN model has exhibited better cost effectiveness for the fault detection and removal than LMANN. Further to explore effectiveness of the proposed HENN model as compared to other existing systems, a comparison has been done (Table -8) and results revealed that the proposed system can be the best optimal solution for defect prediction for object oriented software applications. [57] 94.2 –Symbolic Regression [57] 89.50 –RBP-NN [57] 80.0 –LP [52] 86.6 86.6 87.4 Naive Based [52] 85.6 83.1 83.9 CPSO [53] 69.2 67.6 -T-SVM [54] 75.8 84.1 80.9 GANN [53] 73.4 81.6 -AdaBoost [53] 79.1 82.3 -Random Forest [58] 91.4 -k-NN [56] 91.8 –C4.5 [56] 88. J 48 [56] 90.9 Levenberg-Marquardt-NN [56] 88.0 –NNEP-Evolutionary [53] 88.8 81.2 -PSO [55] 78.7 –PSO-NN [57] 97.7 –HENN SDP 97.9 1 98.9

VI.

## 14 Conclusion

In order to ensure optimal software reliability and quality of service the earlier prediction of faults and its removal is of great significance. In addition, the cost effective solution for defect prediction and fault removal has motivated industries as well as academicians to develop a novel SDP solution that could ensure cost effective and optimal defect prediction solutions. In this paper, an object oriented software matrix based defect prediction model has been developed.

Considering the limitations of artificial intelligence techniques such as artificial neural network, in this paper an evolutionary computing technique named Adaptive Genetic Algorithm (A-GA) has been developed for ANN dynamic weight estimation and learning optimization. The proposed Hybrid Evolutionary computing based Neural Network (HENN) based system has been employed for SDP system. Furthermore, Levenberg Marquardt algorithm based ANN algorithm (LMANN) has been developed for defect prediction. Considering cost effectiveness of the defect prediction systems, a novel mathematical model has been derived and the cost analysis results confirms that the proposed HENN model is cost effective as well as performs better as compared to other existing systems. The simulation results obtained with PROMISE and NASA MDP datasets exhibits that the proposed model performs on average 87.23% accuracy and the best classification accuracy obtained is 97.99% with 100% precision. The proposed model delivers 98.97% of Fmeasure. The cost analysis exhibits that the proposed HENN model is approximate 21.66% cost effective as compared to LMANN. The comparative analysis in this paper reveals that the proposed HENN model performs better as compared to other existing techniques. This paper could perform cost analysis of only HENN and LMANN, hence in future other defect prediction models can also be examined for their cost effectiveness for real time applications.

## 15 Global Journal of Computer Science and Technology

Volume XV Issue II Version I Year ( )

1 2 3

---

<sup>1</sup>© 2015 Global Journals Inc. (US)

<sup>2</sup>© 2015 Global Journals Inc. (US) 1

<sup>3</sup>Evolutionary Computing based an Efficient and Cost Effective Software Defect Prediction System



Figure 1: Figure 1 :

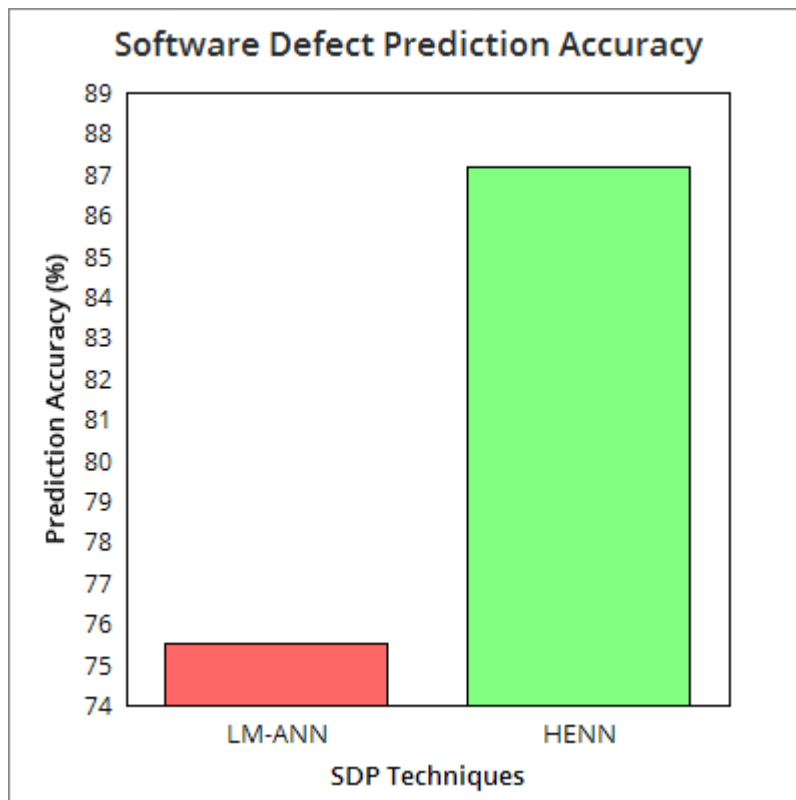


Figure 2:



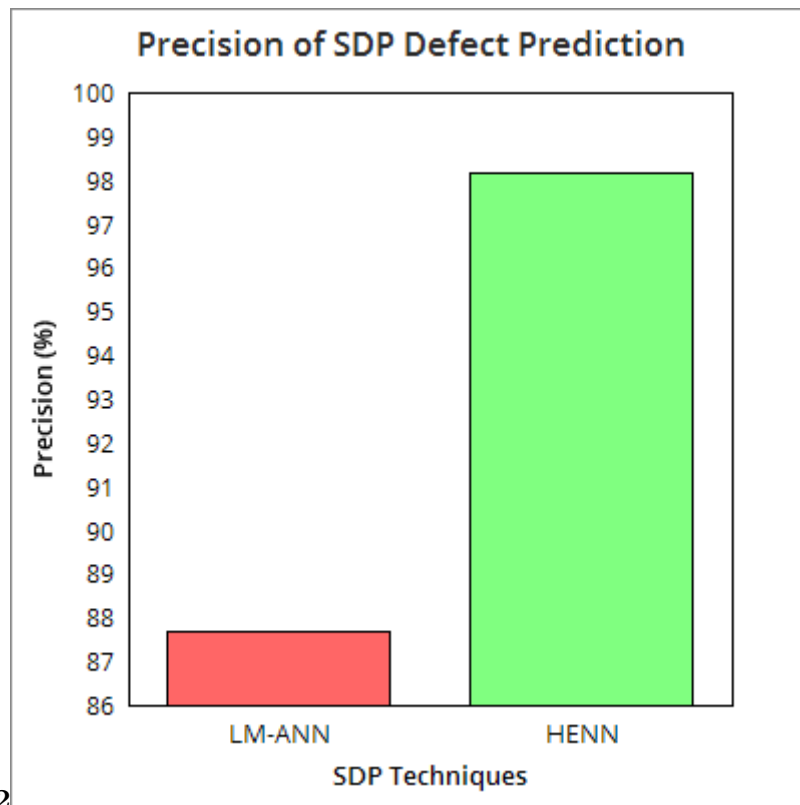


Figure 3: Figure 2 :

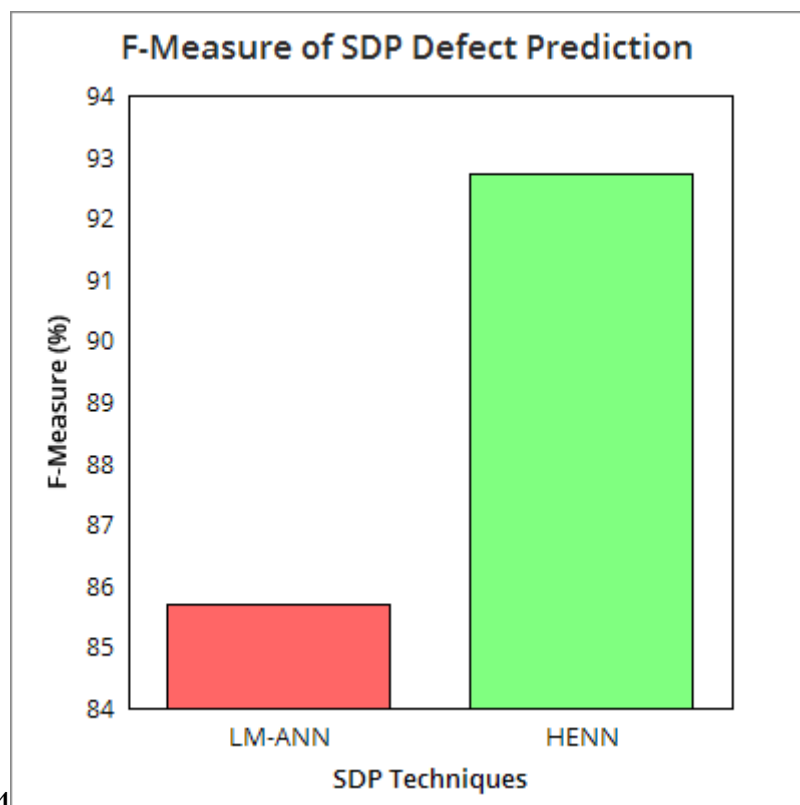


Figure 4: Figure 4 :

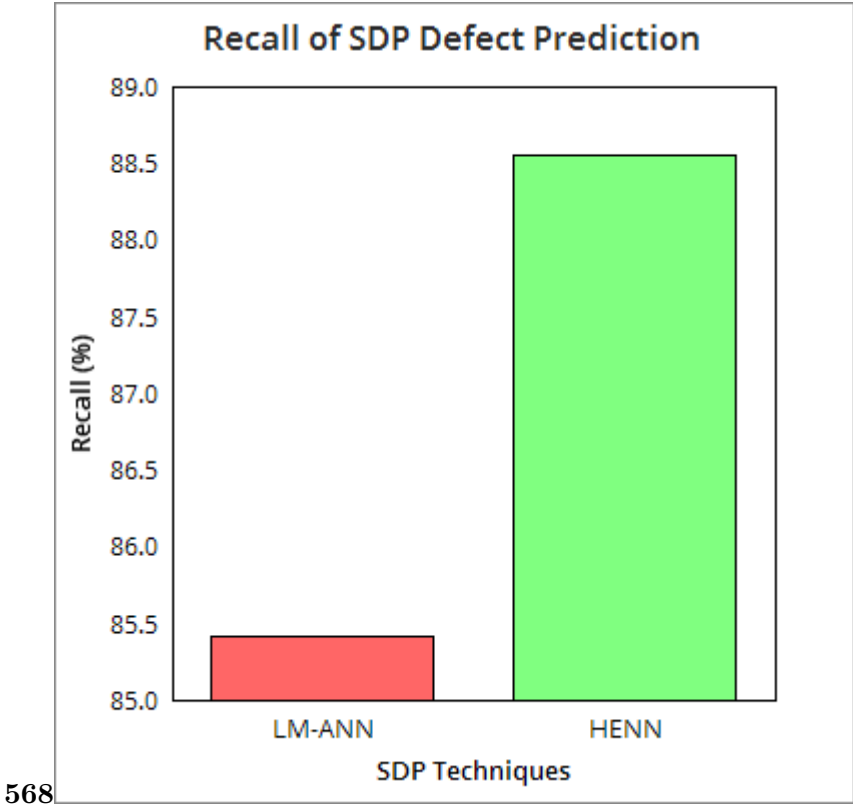


Figure 5: Figure 5 :Figure- 6 :Figure- 8 :

1

1). In this paper,

[Note: G]

Figure 6: Table 1 :

---

Algorithm for Fitness Estimation

Input:?? ?

?? ?? = ?      ??ð ??”ð ??” 0 ? ?? ????? +1 < 5 ?? ????? +2 \* 10 ???2 + ?? ?????  
?      ?  
?      ?  
?      ?  
?      ?  
?      ?  
?      ?  
?      ?  
+      ?

Phase-4: Calculate RMSE of chromosome?? ??

?? ?? = ?      ?  
??      ??  
=??      ??  
??      =1  
??      ??

Where ?? is the number of training data

Phase-5: Calculate the fitness value for chromosome?? ??

[Note: ?? = (?? 1?? ,)]

Figure 7:

Generate Random Population  
of

'n=50' genes or Chrosomes

Extract the Weight of each  
chromosomes

Perform Crossover

Replace the Minimum Fitness  
value Chromosome by Maxi-  
mum fitness value Chromosome

Fed the Weight values for training  
in HENN model

Estimate the Fitness value for  
each chrosomes

NoIs threshold meet? (If Stop Cri-  
teria is accomplished. Implement  
the Model for Software Defect  
Prediction Yes

Year

Volume XV Is-  
sue II Version I

( ) G  
Global Journal  
of C omp uter  
S cience and T  
echnology

[Note: fuctioning till 95% of chromosomes are having similar fitness value. Once the stopping criterion is achieved the A-GA terminates and the final output at output layerO o is obtained.If the final estimated output is more than 0.5, it signifies class as FAULTY otherwise NON-FAULTY. On the basis of retrieved FAULTY and NON-FAULTY data, a confusion matrix is obtained, which is further used for performance assessment.Figure-4 represents the flow diagram of the proposed HENN based SDP model.]

Figure 8:

**2**

Testing	Min	Max	Median
Unit	1.5	6	2.5
System	2.82	8.37	6.2
Field	3.9	27.24	27

Figure 9: Table 2 :

**3**

Cost Estm __SDP	Estimated fault re- moval cost of the software when fault prediction is performed	Cost Estm __WSDP	Estimated fault removal cost of the software
-----------------	---	------------------	--

TP	Number of true pos- itive
TN	Number of true neg- ative
TC	Total number of classes
FC	Total number of faulty classes
? u	Fault identification efficiency of unit testing
? s	Fault identification efficiency of system testing

Figure 10: Table 3 :

4

---

WMC	Overall complexities of the methods in comprising classes
NOC	Number of sub-classes subordinate to a class in the class hierarchy
DIT	Maximum height of the class hierarchy
CBO	Number of other classes to which it is allied with
RFC	

Figure 11: Table 4 :

5

	Predicted Defective	Predicted Defect Free
FAULTY	True Positive	False Negative
NON-FAULTY	False Positive	True Negative

In this paper, the performance of the proposed HENN as well as LM-ANN SDP models has been examined in terms of fault prediction accuracy, precision, F-measure, recall, specification and fault detection and removal cost. The mathematical expression for considered performance parameters are given in Table-6.

Figure 12: Table 5 :

6

Construct	Mathematical Expression
Recall	$TP / (TP + FN)$
Precision	$TP / (TP + FP)$
Specification	$TP / (TP + FP)$
F-measure	$2 * TP / (2 * TP + FP + FN)$
Accuracy	$(TP + TN) / (TP + FP + FN + TN) * 2$

Figure 13: Table 6 :

7

Figure 14: Table 7 :

8

SDP Techniques	Accuracy (%)	Precision (%)	F-Measure (%)
LLE-SVM[51]	81.1	82.5	80.4
SVM [51]	69.4	68.1	69.7
SVM [52]	55.3	88.0	83.2
Natural Gas			

Figure 15: Table 8 :

Figure -9 : Data Modules Tech.				Accuracy	Precision	F-Measure	Recall	Specification	Norm. Fault Removal Cost (Norm.)
JEDIT	492	HENN	0.9799		1	0.9897	1	0.9756	0.2406
		LMANN	0.8394		0.8503	0.9119	0.9832	0.0526	0.2927
ANT	744	HENN	0.8145		0.9343	0.8867	0.8438	0.6346	0.9149
		LMANN	0.7675		0.9879	0.8684	0.7748	0	0.9763
IVY	352	HENN	0.8835		0.9936	0.9380	0.8883	0.3333	0.7115
		LMANN	0.6278		0.6955	0.7681	0.8577	0.0404	0.8936
CAMEL	965	HENN	0.8114		1	0.8952	0.8102	1	0.8771
		LMANN	0.7845		0.9743	0.8792	0.8011	0	1.3401
					3	-	-		

Figure 16:

---

383 [Colonnade Road Suite] , *Colonnade Road Suite* 204.

384 [Boehm ()] , B W Boehm . *Software Engineering Economics* 1981. Prentice-Hall.

385 [Henderson-Sellers and Metrics ()] , B Henderson-Sellers , *Software Metrics* . 1996. UK: Prentice-Hall.

386 [6th International Conference (2013)] *6th International Conference*, Nov. 2013. 2 p. .

387 [McCabe (1976)] ‘A complexity measure’. T J McCabe . *IEEE Transactions on Software Engineering* December  
388 1976. 2 p. .

389 [Fenton et al. ()] ‘A Critique of Software Defect Prediction Models’. N E Fenton , M Neil , I Bellini , P Bruno ,  
390 D Nesi , Rogai . *IEEE Trans. Softw. Engineering* 1999. 25 (5) p. . University of Florence

391 [Zuse ()] *A Framework of Software Measurement*, H Zuse . 1998. Walter de Gruyter Publish.

392 [Bansiya and Davis (2002)] ‘A hierarchical model for Object-Oriented design quality assessment’. J Bansiya , C  
393 G Davis . *ACM Transactions on Programming Languages and Systems* August 2002. 128 p. .

394 [Shrivastava and Shrivastava (2014)] ‘A Hybrid Model of Soft Computing Technique for Software Fault Prediction’.  
395 A Shrivastava , V Shrivastava . *International Journal of Current Engineering and Tech* Aug 2014. 4 (4)  
396 .

397 [W ()] ‘A literature survey of the quality economics of defect-detection techniques’. W . *Proceedings of  
398 the ACM/IEEE International Symposium on Empirical Software Engineering (ISESE)*, (the ACM/IEEE  
399 International Symposium on Empirical Software Engineering (ISESE)) 2006. p. .

400 [Kutlubay and Bener ()] *A Machine Learning Based Model for Software Defect Prediction*, O Kutlubay , A Bener  
401 . 2005. Boaziçi University, Computer Engineering Department (working paper)

402 [Chidamber and Kemerer (1994)] ‘A metrics suite for Object-Oriented design’. S R Chidamber , C F Kemerer .  
403 *IEEE Transactions on Software Engineering* June 1994. 20 p. .

404 [Jianhong et al.] ‘A Neural network based approach for modeling of severity of defects in function based software  
405 systems’. Z Jianhong , P S Sandhu , S Rani . *International Conference on Electronics and information  
406 Engineering*, 2 p. .

407 [Xia et al. (2014)] ‘A new metrics selection method for software defect prediction’. Y Xia , G Yan , X Jiang , Y  
408 Yang . *Progress in Informatics and Computing (PIC)*, *International Conference*, May 2014. p. .

409 [Xing et al. ()] ‘A novel method for early software quality prediction based on support vector machine’. F Xing  
410 , P Guo , M R Lyu . *Software Reliability Engineering, International Symposium*, 2005. p. .

411 [Pushphavathi et al. ()] ‘A novel method for software defect prediction: Hybrid of FCM and random forest’. T  
412 P Pushphavathi , V Suma , V Ramaswamy . *Electronics and Communication Systems (ICECS)*, 2014.

413 [Tegarden et al. ()] ‘A software complexity model of Object-Oriented systems’. D P Tegarden , S D Sheetz , D  
414 E Monarchi . *Decision Support Systems* 1995. 13 (3) p. .

415 [Bo et al. (2007)] ‘A study on software reliability prediction based on support vector machines’. Bo , Xiang Yang  
416 , Li . *The Annual IEEE International Conference on Industrial Engineering and Engineering Management*,  
417 2-4 Dec. 2007. p. .

418 [Rojas and Fernandez-Reyes (2005)] ‘Adapting multiple kernel parameters for support vector machines using  
419 genetic algorithms’. S A Rojas , D Fernandez-Reyes . *The 2005 IEEE Congress on Evolutionary Computation*,  
420 September, 2005. 1 p. .

421 [Racharla Suresh Kumar and Satyanarayana ()] ‘Adaptive Genetic Algorithm Based Artificial Neural Network  
422 for Software Defect Prediction’. Bachala Racharla Suresh Kumar , Satyanarayana . *Global Journal of Computer  
423 Science and Technology : D* 2015. 15 (1) p. . (Version 1.0)

424 [Khoshgoftaar et al. (2001)] ‘An Application of Zero-Inflated Poisson Regression for Software Fault Prediction.  
425 Software Reliability Engineering’. T M Khoshgoftaar , K Gao , R M Szabo . *ISSRE 2001. Proceedings of 12th  
426 International Symposium*, 27-30 Nov. 2001. p. .

427 [Briand et al. (2002)] ‘Assessing the Applicability of Fault-Prone Models Across Object-Oriented Software  
428 Projects’. L C Briand , W L Melo , J Wu , St . *IEEE Trans. Software Eng* July 2002. 28 (7) p. .

429 [Cai ()] K Cai . *On the Neural Network Approach in Software Reliability Modeling*, 2001. p. .

430 [Kang and Bieman (1995)] ‘Cohesion and reuse in an Object-Oriented system’. B K Kang , J M Bieman .  
431 *Proceedings of the ACM SIGSOFT Symposium on software reusability*, (the ACM SIGSOFT Symposium  
432 on software reusability Seattle) March 1995. p. .

433 [Bellini ()] ‘Comparing Fault-Prone Estimation Models’. P Bellini . *10th IEEE International Conference on  
434 Engineering of Complex Computer Systems (ICECCS’05)*, 2005. p. .

435 [Lanubile et al. (1995)] ‘Comparing Models for Identifying Fault-Prone Software Components’. F Lanubile ,  
436 A Lonigro , G Visaggio . *Proceedings of Seventh International Conference on Software Engineering  
437 and Knowledge Engineering*, (Seventh International Conference on Software Engineering and Knowledge  
438 Engineering) June 1995. p. .

- [Etzkorn et al. ()] 'Design and code complexity metrics for Object-Oriented classes'. L Etzkorn , J Bansiya , C Davis . *Object-Oriented Programming* 1999. 12 (10) p. .
- [Mahajan et al. ()] 'Design Of Software Fault Prediction Model Using BR Technique'. R Mahajan , S Gupta , R K Bedi . *Procedia Computer Science* 2015. 46 p. .
- [Hu et al. (2006)] 'Early software reliability prediction with extended ANN model'. Q Hu , Y S Dai , M Xie , S H Ng . *Proceedings of the 30th Annual International Computer Software and Applications Conference (COMPSAC'06)*, (the 30th Annual International Computer Software and Applications Conference (COMPSAC'06)) September 2006. 2 p. .
- [Denaro (2000)] 'Estimating Software Fault-Proneness for Tuning Testing Activities'. Giovanni Denaro . *Proceedings of the 22nd International Conference on Software Engineering*, (the 22nd International Conference on Software Engineering Limerick, Ireland) June 2000.
- [Briand et al. (2000)] 'Exploring the relationships between design measures and software quality in Object-Oriented systems'. L C Briand , J Wust , J W Daly , D V Porter . *The Journal of Systems and Software* May 2000. 51 p. .
- [Yousef ()] 'Extracting software static defect models using data mining'. A H Yousef . *Ain Shams Engineering Journal* 2015. 6 p. .
- [Brun and Michael (2004)] 'Finding Latent Code Errors via Machine Learning over Program Executions'. Y Brun , D E Michael . *Proceedings of the 26th International Conference on Software Engineering*, (the 26th International Conference on Software Engineering) May, 2004.
- [Grosan and Abraham ()] C Grosan , A Abraham . *Hybrid Evolutionary Algorithms: Methodologies, Architectures, and Reviews*, 2011. 75 p. .
- [Halstead ()] M Halstead . *Elements of Software Science*, (New York, USA) 1977. Elsevier Science.
- [Sandhu et al. (2007)] 'Intelligence System for Software Maintenance Severity Prediction'. Parvinder Sandhu , Sunil Singh , Hardeep Kumar , Singh . *Journal of Systems and Software* 2007. Feb. 2008. 3 (5) p. . (Journal of Computer Science)
- [International Conference] *International Conference*, 5 p. .
- [Benlarbi et al. ()] 'Issues in Validating Object-Oriented Metrics for Early Risk Prediction'. Saida Benlarbi , Khaled El Emam , Nishith Geol . *Cistel Technology* 1999. p. 210.
- [Li and Henry ()] 'Maintenance metrics for the Object-Oriented paradigm'. W Li , S Henry . *Proceedings of First International Software Metrics Symposium*, (First International Software Metrics Symposium) 1993. p. .
- [Huitt and Wilde ()] 'Maintenance support for object-oriented programs'. R Huitt , N Wilde . *IEEE Transactions on Software Engineering* 1992. 18 (12) p. .
- [Armah et al.] 'Multilevel data preprocessing for software defect prediction'. G K Armah , Guangchun Luo , Ke Qin . *Information Management, Innovation Management and Industrial Engineering* p. 2013. (ICIII)
- [Abreu and Carapuca ()] 'Object-Oriented software engineering: Measuring and controlling the development process'. F B E Abreu , R Carapuca . *Proceedings of the 4th International Conference on Software Quality*, (the 4th International Conference on Software Quality) 1994. 186.
- [Lorenz and Kidd ()] *Object-Oriented Software Metrics*, M Lorenz , J Kidd . 1994. NJ, Englewood: Prentice-Hall.
- [Malhotra et al. (2014)] 'On the applicability of evolutionary computation for software defect prediction'. R Malhotra , N Pritam , Y Singh . *Advances in Computing, Communications and Informatics (ICACCI, 2014 International Conference*, Sept. 2014. p. .
- [Deodhar (2002)] *Prediction Model and the Size Factor for Fault-proneness of Object Oriented Systems*, Manasi Deodhar . Dec. 2002. Michigan Tech. University (MS Thesis)
- [Rosenberg and Sheppard (1994)] L Rosenberg , S B Sheppard . *Metrics in Software Process Assessment, Quality Assurance and Risk Assessment*, (London) October, 1994. (2nd International Symposium on Software Metrics)
- [Shan et al.] C Shan , B Chen , C Hu , J Xue , N Li . *SOFTWARE DEFECT PREDICTION MODEL BASED ON LLE AND SVM*,
- [Singh and Salaria (2013)] 'Software Defect Prediction Tool based on Neural Network'. M Singh , D S Salaria . *International Journal of Computer Applications* May 2013. 70 p. .
- [Mahdi Askari and Bardsiri ()] 'Software Defect Prediction using a High Performance Neural Network'. Mo-hamad Mahdi Askari , Vahid Khatibi Bardsiri . *International Journal of Software Engineering and Its Applications* 2014. 8 (12) p. .
- [Askari and Bardsiri ()] 'Software Defect Prediction using a High Performance Neural Network'. M M Askari , V K Bardsiri . *International Journal of Software Engineering and Its Applications* 2014. 8 (12) p. .



---

493 [Jindal et al. (2014)] ‘Software defect prediction using neural networks’. R Jindal , R Malhotra , A Jain . *3rd*  
 494 *International Conference on Reliability, Infocom Technologies and Optimization (ICRITO)*, 8-10 Oct, 2014.  
 495 p. .

496 [Chug and Dhall (2013)] ‘Software defect prediction using supervised learning algorithm and unsupervised  
 497 learning algorithm’. A Chug , S Dhall . *Confluence 2013: The Next Generation Information Technology*  
 498 *Summit*, Sept. 2013. p. .

499 [Chug and Dhall (2013)] ‘Software defect prediction using supervised learning algorithm and unsupervised  
 500 learning algorithm’. A Chug , S Dhall . *Confluence 2013: The Next Generation Information Technology*  
 501 *Summit*, Sept. 2013. p. .

502 [Chug and Dhall (2013)] ‘Software defect prediction using supervised learning algorithm and unsupervised  
 503 learning algorithm’. A Chug , S Dhall . *Confluence 2013: The Next Generation Information Technology*  
 504 *Summit*, Sept. 2013. p. .

505 [Verma and Gupta (2012)] ‘Software defect prediction using two level data pre-processing’. R Verma , A Gupta  
 506 . *Recent Advances in Computing and Software Systems (RACSS)*, *International Conference*, April 2012. p. .

507 [ J ()] ‘Software quality in 2010: a survey of the state of the art’. J . *Founder and Chief Scientist Emeritus*, 2010.

508 [Wang and Yao (2013)] ‘Using Class Imbalance Learning for Software Defect Prediction’. S Wang , X Yao .  
 509 *Reliability, IEEE Transactions*, June 2013. 62 p. .

510 [Harman ()] ‘Why the Virtual Nature of Software makes it Ideal for Search Based Optimization’. M Harman .  
 511 *Fundamental Approaches to Software Engineering*, 2010.