

CERTIFICATE

This is to certify that the work contained in the thesis entitled “Observations on Software Defect Prediction”, submitted by Sushant Kumar Pandey (Roll. No.:16071501) for the award of the degree of doctor of philosophy to the Indian Institute of Technology, BHU, Varanasi, is a record of bonafide research works carried out by him under my direct supervision and guidance. It is further certified that the student has fulfilled all the requirements of Comprehensive Examination, Candidacy and SOTA for the award of Ph.D. Degree.

Signature of Supervisor

Prof. Anil Kumar Tripathi

Department of Computer Science and Engineering
Indian Institute of Technology (Banaras Hindu University)

Varanasi-221005, India

DECLARATION

I, **SUSHANT KUMAR PANDEY**, certify that the work embodied in this thesis is my own bona fide work and carried out by me under the supervision of **ANIL KUMAR TRIPATHI** from **January 2017** to **June 2021**, at the Computer Science and Engineering department, Indian Institute of Technology (BHU), Varanasi. The matter embodied in this thesis has not been submitted for the award of any other degree/diploma. I declare that I have faithfully acknowledged and given credits to the research workers wherever their works have been cited in my work in this thesis. I further declare that I have not willfully copied any other's work, paragraphs, text, data, results, etc., reported in journals, books, magazines, reports dissertations, theses, etc., or available at websites and have not included them in this thesis and have not cited as my own work.

Date:

Signature of Student

Place

(Sushant Kumar Pandey)

CERTIFICATE BY THE SUPERVISOR

It is certified that the above statement made by the student is correct to the best of my/our knowledge.

Signature of Supervisor

(Prof. Anil Kumar Tripathi)

Signature of Head of Department/Coordinator of School

COPYRIGHT TRANSFER CERTIFICATE

Title of the Thesis: **Observations on Software Defect Prediction**

Name of Student: **Sushant Kumar Pandey**

COPYRIGHT TRANSFER

The undersigned hereby assigns to the Institute of Technology (Banaras Hindu University) Varanasi, all rights under copyright that may exist in and for the above thesis submitted for the award of the Doctor of Philosophy.

Date:

Signature of Student

Place

(Sushant Kumar Pandey)

Note: However, the author may reproduce or authorize others to reproduce material extracted verbatim from the thesis or derivative of the thesis for author's personal use provided that the source and the Institutes's copyright notice are indicated.

Acknowledgments

“Knowledge is in the end based on acknowledgment.” -Ludwig Wittgenstein

Firstly, I would like to express my sincere gratitude to my advisor Prof. Anil Kumar Tripathi, for the continuous support of my Ph.D. study and related research, for his patience, motivation, and immense knowledge. His guidance helped me in all the time of research, conducting experiments, and writing of this thesis. I could not have imagined having a better advisor and mentor for my Ph.D. study. In addition, I thank retired Prof. Ravi Bhushan Mishra for his support and guidance in the coursework and initial state of my research work.

Besides my advisor, I would like to thank to all the members of my Research Program Evaluation Committee, Prof. Sanjay Kumar Pandey, and Dr.Sukomal Pal, for their insightful comments and encouragement and the hard question which incited me to widen my research from various perspectives. Their guidance and directions help me throughout my journey of PhD.

I would like to express my sincere thanks to Prof. K.K Shukla and Prof. Sanjay Kumar Singh, for their kindness and valuable support in carrying out the research work. I would like to convey my sincere gratitude to the other faculty members of the Department of Computer Science and Engineering, Prof. Rajeev Shrivastava, Dr. R.S Singh, Dr. Bhaskar Biswas, Dr. H.P Gupta, Dr. R.N.Chaudhary, Dr. A.K Singh, Dr. Tanima Dutta and Dr. Vinayak Shrivastava for their extensive and inspiring guidance throughout this tenure.

I extend my sincere thanks to other faculty members of the CSE department for their support. I want to thank my friends and colleague Mr. Tribikram Pradhan, Ms. Manisha Singh, Mr. Ashwini Singh, Mr. Shivang Agarwal, Mr. Ashish Shurma, Mr. Nigmendra Yadav, Mr. Ashish Ranjan, Ms. Dipty Tripathi, Mr. Amit Biswas, Mr. Anshul Gupta, and Mr. Ashish Gupta for their motivation and valuable comments. I also extend my thanks to other colleagues members of our department Mr. Ankit Jaiswal, Mr. Nirbhay

Tagore, Ms. Manisha Sighla, Ms. Pratishta Verma, Mr. Chintoo Kumar, Ms. Anita Saroj, Ms. Naina Yadav, and along with other departmental colleagues of computer science and engineering department and different departmental colleagues. This thesis would not have been possible without their invaluable remarks and persistent help. I extend special thanks to the non-teaching and technical staff in the department, particularly, Mr. Ravi Kumar Bharti, Mr. Ritesh Singh, Mr. Shubham Pandey, Mr. Prakhar Kumar, Mr. Akhilesh Kumar Pal, Mr. Manoj Kumar Rai, Mr. Viplav Biswas, and Mr. Dinesh Tiwari.

In the end, I am grateful to my parents, grandpa, siblings, cousins, friends, and acquaintances, who remembered me in their prayers for the ultimate success. I consider myself nothing without them. They gave me enough moral support, encouragement, and motivation to accomplish my personal goals. My two lifelines (parents) and Granpa who have always supported me financially so that I only pay attention to the studies and achieving my objective without any obstacle on the way. Finally thanks to lord Baba Vishwanath for their belling.

Sushant kumar Pandey

Contents

Acknowledgments	iv
Contents	vi
List of Figures	xi
List of Tables	xiii
Abbreviations	xv
Symbols	xvi
PREFACE	xix
1 Introduction	1
1.1 Software Defect Prediction	2
1.1.1 Cross Version Defect Number Prediction	4
1.1.2 Cross Project Defect Prediction	5
1.1.3 Hybrid Regression Analysis	6
1.2 Limitations of Existing Software Defect Prediction Approaches	6
1.2.1 Limitation on Cross Version Defect Number Prediction	7
1.2.2 Cross Project Defect Prediction	7
1.2.3 Limitations on Hybrid Regression Analysis	8
1.3 Motivation	9
1.4 Research Goal	12
1.5 Contributions	13
1.5.1 BPDET: Classification Based Software Defect Prediction Method	14
1.5.2 BCV-Predictor: Cross-Version Defect Count Vector Predictor	15
1.5.3 DNNAttention: Cross-Project Defect Number Prediction Ap- proach	16
1.5.4 Hybrid Regression Analysis: Entire Software Prediction	18

1.6	Literature Review of Various SDP Methods	18
1.6.1	Literature Review of Classification based Model	19
1.6.2	Literature Review over Cross-Version Defect Prediction	20
1.6.3	Literature Review over Cross-Project Fault Prediction and At- tention layer	22
1.6.4	Literature Review Over Hybrid Regression and Deep Learning	23
1.7	Structure of the Thesis	23
2	Preliminary	25
2.1	Software Metrics	25
2.2	Dataset Description	26
2.2.1	Challenges Over Datasets	29
2.3	Methodologies	30
2.3.1	Techniques Involve in Preprocessing	30
2.3.2	Learning Techniques	31
2.4	Performance metrics	36
2.5	Baseline Methods	38
3	BPDET: Classification Based Software Defect Prediction Model	40
3.1	Introduction	40
3.2	Architectural design	41
3.2.1	Proposed model	42
3.2.1.1	Implication of combining SDA and EL techniques	43
3.2.1.2	Data preprocessing	44
3.2.1.3	Deep learning phase	44
3.2.1.4	Ensemble learning phase	46
3.2.2	Experiment setup	53
3.2.2.1	Datasets	54
3.2.2.2	Tools and its parameters	54
3.2.2.3	Software metrics	54
3.2.2.4	Performance evaluation parameters	54
3.3	Result and discussion	55
3.3.1	Effect of corruption rate over performance of BPDET	63
3.3.2	Research queries discussion	65
3.3.2.1	Effectiveness of BPDET in terms of performance met- rics compared with other fault prediction model	65
3.3.2.2	How useful BPDET model compared with traditional methods regarding the class imbalance and over-fitting problem	66
3.3.2.3	How much training time taken by the BPDET model	68
3.3.3	Significant analysis	69
3.3.4	Insightful discussion	71
3.4	Threats to validity	72

3.5	Summary	72
4	BCV-Predictor: Cross-Version Defect Count Vector Predictor	74
4.1	Introduction	74
4.2	Problem statement	76
4.2.1	Dataset description	77
4.2.2	Software metrics	77
4.2.3	Meta-dataset collection process	77
4.3	Proposed approach	78
4.3.1	Preprocessing	80
4.3.2	Deep learning architecture	81
4.3.3	Proposed model	81
4.3.4	Experimental arrangements	85
4.3.5	Performance metrics	85
4.3.6	Baseline methods	85
4.3.7	Results and explanations	87
4.3.7.1	Justification of RQ-1	87
4.3.7.2	Justification of RQ-2	92
4.3.8	Insightful discussion	92
4.3.8.1	Justification of RQ-3	96
4.3.9	Non-Parametric test	96
4.4	Summary	98
5	DNNAttention: Cross-Project Defect Number Prediction Approach	100
5.1	Introduction	100
5.1.1	Problem statement	102
5.2	Premises	103
5.2.1	Cross project defect number prediction	103
5.2.2	Software projects	104
5.2.3	Evaluation metrics	104
5.3	Proposed work	104
5.3.1	Data synthesis	105
5.3.2	Deep neural network	107
5.3.2.1	Long short term memory	107
5.3.2.2	Attention layer	108
5.3.3	Preprocessing	108
5.3.4	Proposed algorithm	109
5.3.5	Transfer learning phase	111
5.3.6	Experimental setup	112
5.4	Result and discussion	113
5.4.1	What is the accuracy that DNNAttention reached while predicting defect numbers over the target projects?	114

5.4.2	What losses the DNNAttention sufferers while predicting defect numbers over a new project?	116
5.4.3	Compare the performance of the proposed model over exiting baseline methods?	117
5.4.4	How much the DNNAttention subjugate the CI and overfitting problem? How much is the proposed model stable?	119
5.4.5	Some insightful discussion	119
5.5	Threats to validity	120
5.6	Summary	122
6	Predicting the next version of the Software System: A Hybrid Regression Analysis	129
6.1	Introduction	129
6.2	Problem Identification	131
6.2.1	Problem definition	131
6.3	Premises	133
6.3.1	Deep Learning Architecture	133
6.3.2	Software systems	133
6.3.3	Performance Measures	133
6.4	Proposed Approach	135
6.4.1	Proposed algorithm	136
6.4.2	Experimental setup	136
6.4.3	Optimization and Hyperparameters Tuning	137
6.5	Results and Discussion	138
6.5.1	How is the proposed approach effective in predicting the bug count of every module in the next version of a software system?	139
6.5.2	How is the proposed approach effective in predicting the each software metric values of every module in the next version of a software system?	140
6.5.3	How much the proposed model is significantly effective over existing learning methods?	140
6.5.4	Insightful discussion	141
6.6	Threats to validity	146
6.7	Summary	146
7	Summary	150
8	Conclusion and Future Direction	153
8.1	Conclusion	153
8.2	Future Direction	156
A	Publications	158
A.1	Journal Papers	158

A.2 Conference Papers	159
A.3 Communicated Papers	159
Bibliography	160

List of Figures

1.1	Basic architecture of software defect prediction technique.	3
1.2	Thesis Structure.	24
2.1	Percentage of studies for different dataset that are used in SFP.	28
2.2	Underlying architecture of LSTM.	32
2.3	Attention layer architecture.	32
2.4	The architecture of Sequence to Sequence model using LSTM cell.	33
3.1	Working architecture of BPDET.	41
3.2	Design of Denoising auto encoder, the initial input is a and it is corrupted to \hat{a} , \hat{a} is mapped with b, then finally c tries to regenerate to a. The regeneration error is represented by $M_H(a,c)$. Taken from [254].	42
3.3	Staked denoising auto-encoder (SDA), the input layer nodes represented by the solid circle(green) at the bottom, the middle circles are hidden layer nodes. The class label c_i of deep representation is the same as the input.	43
3.4	Comparison of MCC, ROC, PRC, F-measure between best fault prediction techniques over all dataset.	58
3.5	Comparison of MCC, ROC, PRC, F-measure of BPDET over all 12 dataset.	64
3.6	Boxplot of positive parameters λ and δ for corruption rate(σ_d).	65
3.7	Variation of MCC, ROC, PRC and F-measure of BPDET with respect to corruption rate for every dataset.	66
3.8	Boxplot representation of MCC, ROC, PRC and F-measure for every dataset when the corruption rate(σ_d) varies.	67
4.1	Basic architecture of software bug count prediction model.	75
4.2	Meta-data collection process.	78
4.3	Graphical abstract of Proposed work.	79
4.4	Basic architecture of LSTM network, and LSTM cell.	80
4.5	Accuracy of all seven metadata.	90
4.6	Loss of all seven metadata.	93
4.7	Bar-graph of MAE, MSE, & accuracy over all meta-datasets.	94

4.8	Boxplot of MAE, MSE, & accuracy over all meta-datasets after 300 epoch.	95
5.1	Underlying framework of cross project defect number prediction.	103
5.2	Graphical abstract.	106
5.3	Underlying framework of Data Synthesis.	106
5.4	Bar-graph of MSE, MAE, & Accuracy over every project.	114
5.5	Accuracy of datasets.	123
5.6	Accuracy of datasets.	124
5.7	Accuracy of datasets.	125
5.8	Loss of datasets.	126
5.9	Loss of datasets.	127
5.10	Loss of datasets.	128
6.1	Structure of next version prediction of a software system.	129
6.2	The graphical abstract of proposed approach.	135
6.3	Comparison of average MSE, MAE, and Accuracy value(%) of all learning model at different iteration.	142
6.4	Cost effectiveness of Proposed model compared with regular LSTM in terms of MSE, MAE, and Accuracy.	144
6.5	Accuracy and loss versus epoch of all eight metadataset in data augmentation phase.	147
6.6	Accuracy and loss versus epoch of all eight meta dataset in next version prediction phase.	148

List of Tables

1.1	Method level metrics description.	2
2.1	Method level metrics description.	26
2.2	Object oriented level metrics description.	27
2.3	Dataset description before and after preprocessing.	28
2.4	PROMISE Dataset description.	29
3.1	Comparison of MCC over baseline and EL methods with BPDET. Note** “-” symbol implies model cannot give output.	49
3.2	Comparison of MCC with various classification methods with BPDET. Note** “-” symbol implies model cannot give output.	50
3.3	Comparison of ROC over baseline and EL methods with BPDET. . .	51
3.4	Comparison of ROC with various classification methods against BPDET.	52
3.5	Comparison of PRC with various EL method.	53
3.6	Comparison of PRC with various classification method.	57
3.7	Comparison of F-measure over baseline and EL methods with BPDET. Note** “-” symbol implies model cannot produce output.	59
3.8	Comparison of F-measure with various classification method. Note** “-” symbol implies model cannot give output.	60
3.9	Comparison of total training time(seconds) taken by the various en- semble learning methods and some base classifier with respect to BPDET.	61
3.10	Comparison of MCC, ROC, PRC, f-score on k-fold cross validation over BPDET	62
3.11	Adjustment of parameters(λ, δ) with defective rate to achieve thresh- old $T(\sigma_d)$	69
3.12	Two different sample of BPDET and other best models. The full description of the table in the section 3.3.3.	69
3.13	Ranking of values from samples. The full description of the table in the section 3.3.3.	70
4.1	Meta-datasets description.	77
4.2	Performance of BCV-Predictor.	87
4.3	MAE comparison of BCV-Predictor with other techniques.	91
4.4	MSE comparison of BCV-Predictor with other techniques.	91

4.5	Accuracy comparison of BCV-Predictor with other techniques.	91
4.6	Training time (second) of various models compared with BCV-Predictor.	97
4.7	Actual and predicted values of bugs in different samples.	98
5.1	Software projects description.	105
5.2	Performance of DNNAttention.	113
5.3	MSE comparison of DNNAttention over baseline methods.	115
5.4	MAE comparison of DNNAttention over baseline methods.	116
5.5	Accuracy comparison of DNNAttention over baseline methods.	118
5.6	Training time (sec) comparison over cross-heap data of DNNAttention over baseline methods.	118
5.7	Comparison of accuracy produced by the DNNAttention over most optimal baseline model.	121
6.1	Dataset description.	132
6.2	Performance at data augmentation phase.	139
6.3	Performance at next version prediction phase.	142
6.4	Comparison of MSE (%) of next version prediction phase with baseline learning methods. The results are in the form of mean \pm standard deviation.	143
6.5	Comparison of MAE (%) of next version prediction phase with baseline learning methods.	143
6.6	Comparison of accuracy (%) of next version prediction phase with baseline learning methods.	143
6.7	P-Value and Cliff's Delta (δ) of proposed model (PM) compared with the existing approaches in terms of MSE.	145
6.8	P-Value and Cliff's Delta (δ) of proposed model (PM) compared with the existing approaches in terms of MAE.	145
6.9	P-Value and Cliff's Delta (δ) of proposed model (PM) compared with the existing approaches in terms of Accuracy.	145
6.10	Effectiveness level [43] of cliff's delta(δ).	146
6.11	Number of projects in which PM is statistically significantly improves over existing approaches (+); performs more or less or approximately equally well (=); and performs loses (-) compared to the baseline techniques in form of MSE/MAE/Accuracy.	146

Abbreviations

S[D, B, F]P	Software [D efect, B ug, F ault] P rediction
S[D, B, F]NP	Software [D efect, B ug, F ault] N umber P rediction
S[D, B, F]CP	Software [D efect, B ug, F ault] C ount P rediction
CPDP	C ross P roject D efect P rediction
CP[D, B, F]DNP	C ross P roject[D efect, B ug, F ault] N umber P rediction
CPDCVP	C ross P roject D efect C ount V ector P rediction
CV[D, B, F]P	C ross V ersion [D efect, B ug, F ault] P rediction
CV[D, B, F]CP	C ross V ersion [D efect, B ug, F ault] C ount P rediction
BCV	B ug C ount V ector
SDA	S taked D enoing A utoencoder
MLP	M ulti L ayer P erceptron
SVM	S upport V ector M achine
EL	E nsemble L earning
NB	N aive B ayes
LR	L ogistic R egression

Symbols

x_t	input sequence at t time step
y_t	output sequence at t time step
w_c	weight matrix for candidate cell
w_i	weight matrix for input gate
w_f	weight matrix for forget gate
w_o	weight matrix for output gate
b_c	bias value for candidate cell
b_i	bias value for input cell
b_f	bias value for forget cell
b_o	bias value for output cell
σ	sigma activation function
f^t	forget gate
a^t	output sequence of t time step
C^t	cell state value at t time step
C'^t	new cell state value at t time step
i^t	update state at t time step
o^t	output state at t time step
o^t	output state at t time step
T_x	Fixed length input
T_y	Fixed length output
s_i	hidden state
Y_t	amount of attention

h_t	hidden state
Val_{mae}	mean absolute error at validation set
Val_{mse}	mean squared error at validation set
Val_{acc}	accuracy at validation set
Val_{mae}	mean absolute error at validation set
a_{i1}	input feature of a with i^{th} module
$deep_{i1}$	deep feature of a with i^{th} module
Y_{final}	final output of ensemble learning phase
D_{tr}	training data
D_{ts}	testing data
$cPND_j$	combined probability of non-defective of j^{th} module
cPD_j	combined probability of defective of j^{th} module
λ, δ	SDA parameters
M_i	i^{th} module
ζ_l	input vector for layer l
η	initial learning rate
g_t	gradient at time t
f_i	i^{th} feature
v	is exponential average of gradient
κ_q	probability score of each bug q
Seq2Seq	sequence to sequence

Dedicated to my Parents, and Grandpa

PREFACE

Software testing is an essential phase of the software development life cycle (SDLC); it requires approximately 40% of resources in software development. Software defect causes a discrepancy in actual output and expected output, which leads to system failure. A software defect is a condition in a software product that does not meet the software requirement (as stated in the requirement specifications) or end-user expectation (which may not be specified but is reasonable). Software Defect Prediction (SDP) is a method to predict the fault-prone module of a software system. After detecting the fault-prone module, the project leader can optimally allocate testing resources; more tester team members to more fault-prone software modules. It reduces testing effort and leads to a decrease in the cost of software development. Many software practitioners have proposed several SDP methodologies using statistical techniques, machine learning methods, or deep learning architecture. There are many hindrances for detecting faulty modules in software fault prediction systems, such as missing values or samples, data redundancy, irrelevance features, and correlation. Many researchers have built various Software Bug Prediction (SBP) models, which classifies buggy and clean modules associated with software metrics. There are several challenges in the SDP domain, such as skew distribution in the public datasets results in a class imbalance problem, overfitting, quality of a public dataset. This thesis's main objective is to reveal the favorable result by feature selection, machine learning, or deep learning methods to detect defective and non-defective software modules, estimating the number of a module in the target project, and predicting data associated with the next version of a software project. Along with addressing class imbalance and overfitting problems in the SDP domain and produces unbiased results. Further, normally predicting based on one or more similar projects may start showing deficient prediction for newer projects; to counter this possibility of learning from across multiple projects may be useful. In this thesis, we also proposed a cross-project defect prediction method. Additionally, an attempt has been made to include the prediction of certain software metrics apart from classical software bug prediction. The proposed work in the thesis is classified into four sections first that are based on a classification-based approach to identify the faulty or clean module. Then we will discuss how to identify the number of defects in every module in the upcoming version of a software project; it is one step ahead of just identifying a module is faulty or not. Further, the cross-project defect

number prediction and prediction of the entire software project are discussed. These approaches help in optimally allocating development and testing resources and leads to produce high-quality software products.

First, we investigate various existing classification-based defect prediction methods and propose a rudimentary classification-based framework, Bug Prediction using Deep representation and Ensemble learning Technique (BPDET) for SBP. It combinedly applies by ensemble learning and deep representation. The software metrics which are used for SBP are mostly conventional. Staked denoising auto-encoder (SDA) is used for the deep representation of software metrics, which is a robust feature learning method. The proposed model is mainly divided into two stages: the deep learning stage and two layers of the EL stage (TEL). The extraction of features from SDA in the very first step of the model then applied TEL in the second stage. TEL is also dealing with the class imbalance problem. The experiment mainly performed NASA (12) datasets to reveal the efficiency of DR, SDA, and TEL. The performance is analyzed in terms of Mathew co-relation coefficient (MCC), the area under the curve (AUC), precision-recall area (PRC), F-measure, and time. Out of 12 dataset MCC values over 11 datasets, ROC values over 6 datasets, PRC values overall 12 datasets, and F-measure over 8 datasets surpass the existing state-of-the-art bug prediction methods. We have tested BPDET using Wilcoxon rank-sum test, which rejects the null hypothesis at $\alpha = 0.025$. We have also tested the stability of the model over 5, 8, 10, 12, and 15 fold cross-validation and got similar results. Finally, we conclude that BPDET is stable and outperformed on most of the datasets compared with existing EL-based methods and other benchmark techniques.

In the next work, we investigate various existing and baseline works over software number fault prediction. Predicting the number of bugs in a software system intensifies the software quality and turns down the testing effort required in software development. It reduces the overall cost of software development. The evolution of hardware, platform, and user requirements leads to develop the next version of a software system. In this chapter, we formulate a problem statement and its novel solution, i.e., we are considering the prediction of the bug count vector of a successive version of a software system. After predicting the bug count vector in the next version of a software, the developer team leader can adequately allocate the developers in respective fault dense modules; in a more faulty dense module, more developers are required. We have conducted our experiment over seven PROMISE repository

datasets of different versions. We build metadata using a concatenation of various versions of the same software system for conducting experiments. We proposed a novel architecture using deep learning called BCV-Predictor. BCV-Predictor predicts the bug count vector of the next version software system; it is trained using metadata (collection of different dataset). To the best of our knowledge, no such work has been done in these aspects. We also address overfitting and class imbalance problems using the random oversampling method and dropout regularization technique. We conclude that BCV-Predictor is conducive to predicting the bug count vector of the next version of the software project. We found five out of seven meta datasets reach more than 80% accuracy. In all seven meta datasets, Mean Squared Error (MSE) lies from 0.71 to 4.715, Mean Absolute Error (MAE) lies from 0.22 to 1.679, MSE and MAE over validation set lie between 0.84 to 4.865, and 0.22 to 1.709 respectively. We also compared the performance of the BCV-Predictor with eleven baseline techniques and found the proposed approach statistically outperforms on most of the dataset.

In the next work, we revisited the software fault number estimation and along with existing work of cross-project number prediction. Both techniques help in allocating resources before the testing phase more optimally. Due to a lack of an adequate dataset, defects can be predicted by employing data from different projects to train the classifier called cross-project defect prediction (CPDP). Cross-project defect number prediction (CPDNP) is one step ahead of CPDP, in which we can also estimate the number of defects in each module of a software system; we contemplate it as a regression problem. This work dealt with the CPDNP mechanism and suggested a CPDNP architecture by employing a deep neural network and attention layer called DNNAttention. We synthesis substantial data named cross-heap by utilizing an amalgamation of 44 projects from the PROMISE data repository. We feed the cross-heap into DNNAttention to train and evaluate the performance over 44 datasets by applying transfer learning. We have also address class imbalance (CI) and overfitting problems by employing multi-label random oversampling and dropout regularization, respectively. We compared the performance of DNNAttention using MSE, MAE, and accuracy over eight baseline methods. We found out of 44 projects, 19 and 20 have minimum MSE and MAE, respectively, and in 19 projects, accuracy yields by the proposed model surpasses exiting techniques. Moreover, we found the improvement of DNNAttention over other baseline methods

in terms of MAE, MSE, and accuracy by inspecting 20% lines of code is substantial. In most situations, the improvements are significant, and it has a large effect size across all 44 projects.

The overall cost of software development is detrimental to determining the utility of any software. However, the uncertain nature of this cost increases the risk associated with the software. Hybrid regression analysis investigates the various regression-based approaches and estimates all software features, including a number of bugs in every module of the software system. Predicting the next version of the software system helps reduce this cost, allowing a better allocation of developers and testers for software development. In this chapter, we develop a novel approach for the prediction of the next version of a software system. Our method predicts the bug count and the metric values for each module in the software using a computational framework consisting of two phases, the data augmentation phase and the next version prediction phase. We perform our experimentation on 8 public datasets from the PROMISE repository using all existing versions. To achieve unbiased results, we conducted each experiment 50 times and took a mean value of it. The proposed methodology has an accuracy of over 60% on 5 out of eight datasets, while MSE and MAE lie between 45.34% to 185.5% and 30.59% to 131.21%, respectively. We have also compared the proposed model's performance with eight baseline learning methods and found the proposed model significantly outperforms these techniques.