

Application of Ensemble Machine Learning Methods to Improve Effort Prediction Accuracy

Bhaskar Marapelli

Abstract: Estimation of Effort in the Software Development Life Cycle is crucial for project planning where resources are decided. Since a good estimate of resources leads to a good estimate of the cost of the project and helps in the success of the project, effort estimation becomes one of the research topics. Much research using machine learning techniques has been done for effort estimation and no single technique is superior. Machine learning techniques use statistical analysis on datasets and help us in the prediction and estimation of effort. This study proposes effort estimation methods based on ensemble techniques to improve the performance accuracy of individual machine learning methods. Three ensemble techniques Voting, Bagging and Boosting were studied to show the performance accuracy of individual machine learning techniques could be improved using ensemble methods, which combine several machine learning methods into one predictive model. It has been observed that Boosting adds weak predictors to an ensemble and each corrects its predecessor can give more accurate performance than single machine learning methods and also other ensembles.

Index Terms: Accuracy, Bagging, Boosting, Effort Estimation, Ensemble, Linear Regression, K-Nearest Neighbors, Machine Learning Support Vector machines.

1 INTRODUCTION

Machine learning techniques have been studied for a long time and many models were proposed using machine learning techniques [15, 18, 19, 20, 21, 22, and 23] but each model has its own advantages and disadvantages and the accuracy of the models still need to be improved. The learning methods when observed individually give better performance under some configuration and worse with some other configuration, solo techniques are not stable and inconsistent on different data sets [7]. So there is a need to solve the problem for that we propose ensemble machine learning methods that combine several individual methods into an ensemble or a combined predictive model. It is a good idea to get predictions from the combined predictive model instead of individual models which can improve the performance and accuracy of the prediction. In this research, we studied four individual machine learning models created from best configurations of 1. Linear Regression, 2. K-Nearest Neighbors, 3. Support Vector Machines and 4. Decision Tree algorithms. Also, we have studied three ensemble models created from Voting, Bagging and Boosting algorithms. The results show that compared with individual models the ensemble models give the best performance in prediction and improved accuracy score. This paper is organized as follows 2. Related work, 3. Machine learning algorithms, 4. The Ensemble Machine Learning Methods, 5. Methodology 6. Results and Discussion 7. Conclusion, 8. References.

2 RELATED WORK

Omar et.al, (2017) Omar studies machine learning techniques for effort estimation of the software projects and suggests the theoretical methods are time-consuming. He studied individual techniques support vector machine and k-nearest neighbor and he says combining these two techniques gives improved results.

When individually studied Omar says SVM given better performance than KNN and with boosting techniques estimation accuracy can be improved. Petrônio et.al, (2007) in this paper author introduced a machine learning model to estimate the effort along with confidence interval. He proposes the bagging or bootstrap aggregation for classification and regression problems which improves prediction and reduces variance in predictions. He studied the regression methods SVR, MLP, M5P and found the bagging method improves prediction. He used confidence intervals for reliability indication of predictions. For evaluating prediction accuracy Petronio used MMRE and the PRED as measures. Przemysław et.al, (2018) this article gives effective machine learning approaches using industry best practices. He built two models for predicting effort and duration each model uses SVM, MLP and GLM algorithms, from the best performing machine learning algorithms SVM, MLP and GLM he built Ensemble models for effort and duration estimation then mean magnitude relative error and percentage relative error derivation metrics were used to compare the results. He concludes that ensemble models give accurate results compared to other approaches for project initial lifecycle. Pichai et.al, (2018) this paper proposes an ensemble estimation model for effort from estimation methods created from three algorithms. This approach uses the correlation-based feature selection and genetic algorithm. Correlation is used to find the individual estimator's estimation accuracy; this study finds methods with high correlation in terms of correlation coefficients which gives high accuracy then the genetic algorithm is used to create combined estimation from selected methods. By using the popular error metrics the results show that the ensemble technique created gives more accurate results compared to the best method from the algorithms. Leandro et.al, (2013) this paper aims at selecting and analyzing machine learning models and to improve the accuracy of these models ensembles of machine learning models is proposed. To improve effort estimation accuracy bagging ensemble of regression trees is done on different data sets. He further studies to determine existing machine learning algorithms individually can give less performance compared to the ensemble learning approach proposed. The evaluation of technique is performed by metrics like MMRE, PRED. By using ensembles the performance of methods can be

- *Bhaskar Marapelli is currently pursuing PhD in Computer Science and Engineering Department, Shri Jagdishprasad Jhabarmal Tibrewala University*
- *Rajasthan, India, bhaskar.marapelli@gmail.com*

improved on smaller data sets. The locality approach studied in this paper suggests that performances of methods on higher data sets could be improved. Abdelali et.al, (2018) In this paper tree-based machine learning techniques random forest were studied. The random forest model was studied by varying all its parameters, and the accuracy of the model with different parameters is analyzed. The best performing parameters were selected and the RF model performance is compared with the regression tree. He used the performance metrics Pred, MMRE, and MdMRE to evaluate the models and in his study he says that the Random forest model gives better performance than the regression tree. He also suggests that random forest with selected parameters performs better than the optimized MLP, M5p, Analogy and SVR based models. He concludes telling random forest is the preferred technique for effort estimation. Tuğçe et.al, (2017) this paper proposes a machine learning model based on neural networks, the network model has been examined using data taken from the largest bank in Turkey. He studied that the neural network model for software estimation which is built with the features using data from the bank is able to overcome uncertainties and complexities effectively. He proposes that ANNs are the best models to calculate software effort because these models are learning-based models and these use previous project data so they can predict software effort accurately. Jørgensen et.al, (2002) this paper gives a review of studies software development effort expert estimation. He contributes by giving a review of expert estimation and examines the validity of expert estimations with human judgment. He gives 12 best practices for expert estimation and he evaluates those practices. He concludes that the expert estimation gives more accurate results when the experts have more domain knowledge and when simple strategies give more accurate results.

3 MACHINE LEARNING ALGORITHMS

In machine learning problems are solved by predicting unknown values using sample data as input. The numbers of entries in the input sample are called the features. While evaluating machine learning algorithms the sample input data is divided into training and testing data. Machine learning problems can be two types supervised and unsupervised. In Supervised learning, the feature needs to be predicted is known and it can be divided into classification and regression. If the values need to be predicted consists of one or more classes the problem is called Classification; if the values need to be predicted contains continuous values then the problem is Regression. If the sample data consists only input and we are not aware of the output or target vales then the problem is unsupervised learning. This kind of problem group similar example data into one type of output this is called clustering. Regression is a method that predicts target values by expressing the data as a statistical equation. Regression predicts a numerical outcome (dependent variable (eg. Effort) from a set of inputs independent variables (eg., effort multipliers).

3.1 Linear Regression

Linear Regression (LR) model predicts the linear relationship between numerical variables. In simple Linear Regression, there will be one dependent variable and one independent variable in multiple regression more than one independent

variable describes the value of the dependent variable. In our model, it tries to predict up to what extent the dependent variable (Effort) is described by independent variables (Effort multipliers).

We could describe Linear Regression with an equation as follows

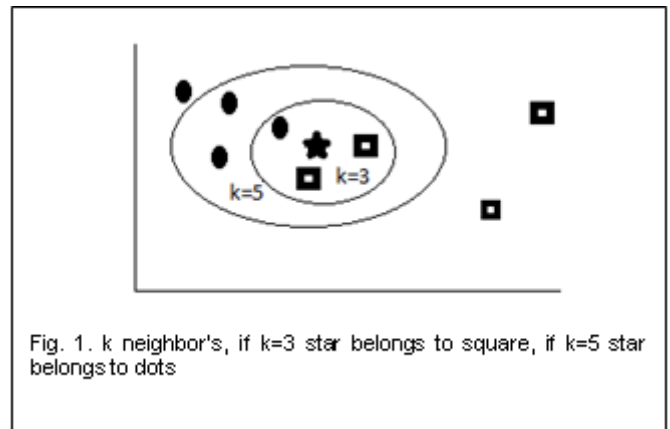
$$y = \beta + \beta_1 X_1 + \beta_2 X_2 + \dots \quad (1)$$

y (dependent variable) is linearly related to each x (independent variables)

Each x contributes additively to y

3.2 K-Nearest Neighbors

The K-Nearest Neighbors (KNN) method works by choosing the closest neighbor points to the predefined training samples. If a point is nearer to the maximum number of some points (which are of the same group) in training set then this point belongs to that class. The distance can be calculated using standard metric Euclidean distance or Manhattan distance. This method remembers all its train data, for this reason, it is also known as non-generalizing machine learning methods [27]. In the below diagram for k=3 (when neighbors are 3) then star belongs to the square class and when k=5 then star belongs to circle class. KNN can be used in solving classification and also regression problems, here our case is Effort prediction the outcome is a continuous value so we used the KNN Regression algorithm.



Euclidian distance between two data points can be calculated using following formula [7].

$$d(x, y) = \sum_{i=1}^N \sqrt{(x_i - y_i)^2} \quad (2)$$

3.3 Support Vector Machine (SVM)

Support Vector Machine (SVM) algorithms can be used for classification and also for regression. Support Vector Regressor (SVR) implements Regression with SVM [6]. SVM can give highly accurate results and it can be used to model complex problems [12]. SVM uses a hyper plane to separate instances of two classes; the hyper plane maximizes the margin. The points that are near to or lies on border are classed as support vectors. The complexity is decided by the number of support vectors [1]. SVM is also used for applications like pattern classification and regression such as handwritten character recognition [20]. Since in this paper we

are predicting effort which is a continuous value we use the Support Vector Regressor (SVR) algorithm.

3.4 Decision Tree

Decision Tree (DT) algorithms can be used for both classification and regression which are supervised learning algorithm. The DT works by creating models with learning decision rules created from training data and tries to predict the class or value of the target values. DT uses IF-THEN rules to make decisions; the best fit model will have a deeper tree with decision rules, the decision making tree complexity will affect the accuracy of the results [26]. DT can handle both numerical and categorical data. DT is applied for regression problems like effort estimation by using Decision Tree Regressor [27].

4 THE ENSEMBLE MACHINE LEARNING METHODS

There is no evidence that anyone single machine learning method is best in estimating the effort accurately so we use ensemble learning methods that combine individual estimator results into one estimated result. When single methods not performing better, then an ensemble of individual methods can give the best performance [4]. For performance improvement of single machine learning methods many researches has been done on ensemble models [16, 8, 1, 4]. In our study, we tried to see the ensemble models based on voting (Voting Regressor), Bagging (Random Forest Regressor) and Boosting (Gradient Boost Regressor) algorithms. When we observed the results certainly the ensemble models improve the performance of single machine learning models. In the three ensemble methods when we compare the accuracy given by Gradient boost Regressor is quite impressive on all three data sets the GBR gives more than 98% accuracy.

4.1 Voting.

Voting is the simplest ensemble algorithm that can be used for classification and also for Regression. Voting creates two or more sub-models they predict by taking the mean or the mode of the predictions and voting algorithms allow sub-model to vote to the best outcome. The creation of sub-models helps in selecting models with different predictions. In our case to predict effort, we used the Voting Regressor (VR) algorithm and we used individual models as LR, KNR, SVR, DTR from these four models voting is done. VR with Deshrains and China data sets it shown improved accuracy.

4.2 Bagging.

Bagging is also called Bootstrap Aggregation which is an ensemble algorithm that can be used for Classification and also for Regression. This is a statistical estimation technique where the mean of random samples of the data (with replacement) is estimated. This technique is useful when the data is limited. In this technique, multiple machine learning models are trained with random samples taken from training data with replacement and each model predicted results are averaged to give the best prediction. Decision Tree is the most used algorithm for bagging. Since the random sample is drawn with replacement default bag size is 100% to make possible the training sample is chosen again and again. In our case, we have used Random Forest Regressor (RFR) as a Bagging algorithm with a number of bags or estimators as 500 and we

got good performance compared with individual Decision tree and other algorithms. From the Results, we found RFR with Deshrains and China data sets shown improved accuracy.

4.3 Boosting.

Boosting makes weak learners modified to become better predictors. The first boosting application is Adaptive Boosting which is a good success and is also called as AdaBoost. In AdaBoost, predictions will depend on the number of votes given by weak learner's, weighted by their individual accuracy. The AdaBoost algorithm was successfully applied on binary classification problems [25]. The statistical framework created by Breiman, is improved by Friedman and is called Gradient Boosting Machines or gradient tree boosting [28]. Gradient boosting algorithm is a greedy algorithm and it can overfit a training dataset quickly [25].

Gradient boosting Regressor (GBR) is a type of ensemble decision tree model; for building predictive models GBR is the best chosen techniques. In our approach for effort prediction, we used GBR for improving prediction accuracy. The results show on all three data sets Deshrains, China and COCOMONASA2 the performance accuracy of GBR is more than 98% which is improved from individual learners and also the other ensemble methods.

5 METHODOLOGY

This study aims to find the improved accuracy in effort estimation. For this experiment, we used three data sets Deshrains, China and COCOMONASA2. The data sets are partitioned into the train and test sets and all models are trained with a train set and tested with the test set. For each data set individual models using LR, KNR, SVR, and DTR are created and accuracy score is calculated then for each data set ensemble voting, bagging and boosting algorithms are applied and scores are calculated. Finally, the scores obtained by individual models and ensemble models are compared.

5.1 Data Sets

In this paper, we analyzed different data sets (Deshrains, COCOMONASA2) from PredictOr Models In Software Engineering (PROMISE) Repository and China data set. The COCOMONASA2 contains 15 cost drivers as input features; the actual effort in person-months is the dependent variable. Deshrains' data set contains 12 features in them; the effort is the dependent variable. The missing values in Deshrain's data sets are filled with zeros.

TABLE I.

Dataset	# of Projects	# attribute s	Train and Test Split	
			# Train samples	#Test samples
DESHARNAIS	81	12	64	17
COCOMONASA2	93	9	74	19

PROMISE Repository data sets description

The china dataset contains 19 features in them project id is removed and n_effort which is repeated also removed. So we used 17 features and in them effort will be dependent variable.

TABLE II.

Dataset	#of Projects	# attributes	Train and Test Split	
			# Train samples	# Test samples
CHINA	499	15	399	100

China data set Description

5.2 Evaluation Criteria

To evaluate the models created we have used Regression metrics implemented by sklearn.metrics module like Mean Squared Error (MSE), Mean Absolute Error (MAE), Explained Variance Score and R2 score.

5.2.1 Mean Squared Error (MSE).

The MSE gives average of differences between predicted values and actual values..

$$MSE(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} (y_i - \hat{y}_i)^2 \quad (3)$$

\hat{y}_i is calculated value of i -th sample

y_i is actual value

5.2.2 Mean Absolute Error (MAE).

The MAE gives how far the estimated values are from actual values. In this measure all individual difference are weighted equally.

$$MAE(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} |y_i - \hat{y}_i| \quad (4)$$

\hat{y}_i is calculated value of i -th sample

y_i is actual value

5.2.3 Explain Variance Score.

The Explained Variance Score computes the explained variance regression score, the best possible score is 1.0, lower values are worse.

$$\text{explained_variance}(y, \hat{y}) = 1 - \frac{\text{var}(y - \hat{y})}{\text{var}(y)} \quad (5)$$

\hat{y} is the calculated target output

y the actual output

var is variance

5.2.4 R2 Score.

R2 score computes the coefficient of determination. R2 score gives the measure of goodness of fit which indicates how well unseen samples are likely to be predicted by the model. The best possible score could be 1.0 and for some worst models it can be negative also.

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (6)$$

Where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

And $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \epsilon_i^2$

\hat{y}_i is calculated value of i -th sample

y_i is actual value

6 RESULTS AND DISCUSSION

We presented the results of our experiments that we have been performed using machine learning techniques in this section. We have shown the experimental results of individual machine learning methods on three data sets and also we presented the improved scores with Ensemble machine learning methods. All the experiments were conducted by using Python's Scikit-learn machine learning package. Python's Scikit-learn is an open-source machine learning package built on NumPy, SciPy, and matplotlib which is an efficient tool for data mining and data analysis.

TABLE III.

Sn	Machine Learning Model	Metric	Data sets		
			DESHA RNAIS	CHINA	COCOMONA SA2
1	Linear Regression model	R2 Score	-0.0563	0.7096	0.2149
2	K-Nearest Neighbors Regressor model		0.1998	0.6161	0.0920
3	Support Vector Machines Regressor model		-0.0679	0.6200	0.0490
4	Decision Tree Regressor		0.6565	0.7381	0.2181

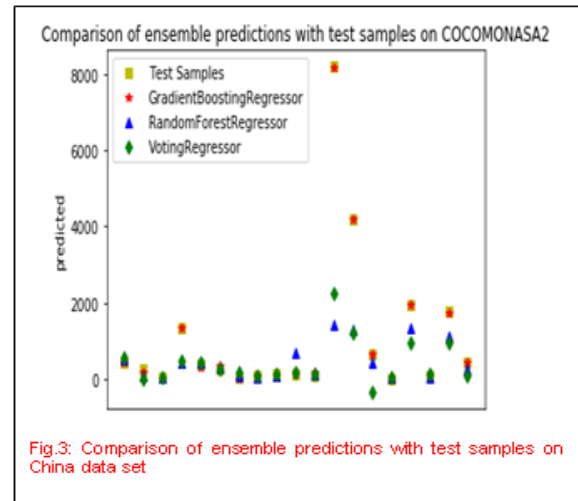
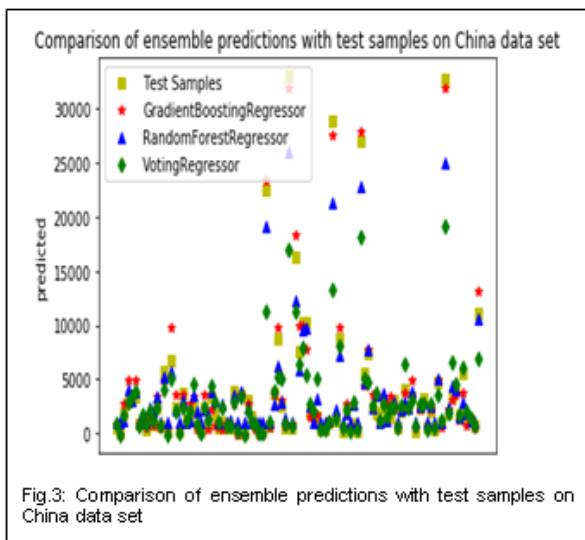
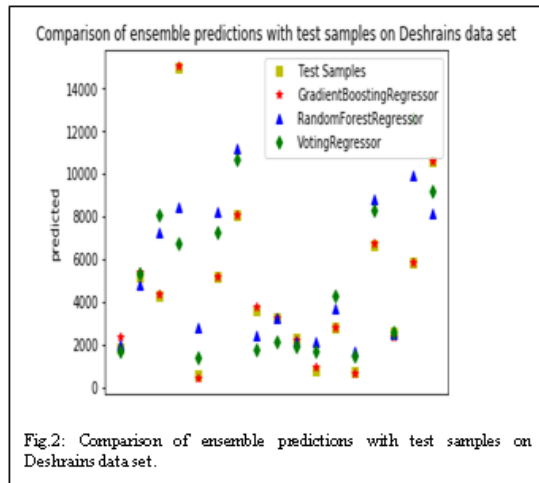
Accuracy scores of individual machine learning models

TABLE IV.

S N O	Machine Learning Model	Metric	Data sets		
			DESHA RNAIS	CHINA	COCOMON ASA2
1	Voeting Regressor	R2 Score	0.4542	0.7375	0.1555
2	Random forest Regressor		0.5033	0.9261	0.2213
3	Gradient Boosting Regressor		0.9970	0.9821	0.9994

Accuracy scores of Ensemble Machine learning methods

From the above tables information it is clear that the scores by individual machine learning methods are improved by ensemble machine learning methods. In Table III we have shown the R2 scores of individual machine learning methods and in Table IV we presented ensemble methods Regression scores. From Table IV we can observe and identify that Gradient boosting Regressor out performs all other individual methods and also other ensemble methods.



From the above figures (Fig 1, Fig 2, Fig3) we can observe on each data set ensemble methods predictions with test samples and these figures shows Gradient Boosting Regressor is the best predictor. As we have seen the R2 scores of GBR more than 0.98 the predictions done by GBR are closer to the test sample values.

7 CONCLUSION

This research study evaluates individual machine learning methods LR, KNR, SVR and DTR, ensemble methods Voting (VR), Bagging (RFR) and Boosting (GBR) on three datasets Deshrains, China and COCOMONASA2. Since there is no evidence that anyone single machine learning method is best in estimating the effort accurately the ensemble learning methods which combine individual estimator results into one estimated result, will give improved performance [4]. The voting ensemble algorithms allow sub-models to vote to the best outcome, the Bagging trains multiple models and each model predicted results are averaged to give the best prediction, the Boosting makes weak learners modified to become better predictors. The results of the study show ensemble machine learning methods improve the accuracy of predictions. The ensemble methods accuracy calculated with R2 score of Scikit-learn metrics is far better than the individual methods. GBR accuracy is more than 98% on all three data sets we have studied. With the evidence of the results, we can conclude by telling Gradient Boost Regressor to predict the effort with improved accuracy on all three data sets.

8 REFERENCES

- [1] Omar Hidmi, and Betul Erdogan Sakar, Software Development Effort Estimation Using Ensemble Machine Learning, Int'l Journal of Computing, Communications & Instrumentation Engg. (IJCCIE) Vol. 4, Issue 1 (2017) ISSN 2349-1469 EISSN 2349-1477.
- [2] Morakot Choetkierikul, Hoa Khanh Dam, Truyen Tran, Trang Pham, Aditya Ghose, and Tim Menzies, A deep learning model for estimating storypoints, JOURNAL OF LATEX CLASS FILES, VOL. 14, NO. 8, AUGUST 2015.
- [3] M. Jorgensen, A review of studies on expert estimation of software development effort, The Journal of Systems and Software 70 (2004) 37–60.

- [4] Ekrem Kocaguneli, Tim Menzies, Jacky Keung, On the Value of Ensemble Effort Estimation, JOURNAL OF IEEE TRANSACTIONS ON SOFTWARE ENGINEERING, VOL. X, NO. Y, SOMEMONTH 201Z.
- [5] Mohammad Azzeh, Ali Bou Nassif, Shadi Banitaan, Cuauhtémoc López-Martín, Ensemble of Learning Project Productivity in Software Effort Based on Use Case Points, 2018 17th IEEE International Conference on Machine Learning and Applications, 978-1-5386-6805-4/18.
- [6] Mohamed Hosni and Ali Idri, Ali Bou Nassif, Alain Abran, Heterogeneous Ensembles for Software Development Effort Estimation, 2016 3rd International Conference on Soft Computing & Machine Intelligence, 978-1-5090-3696-7/16.
- [7] Mohamed Hosni, Ali Idri, Alain Abran, Ali Bou Nassif, On the value of parameter tuning in heterogeneous ensembles effort estimation, Springer Nature 2017, Soft Computing, <https://doi.org/10.1007/s00500-017-2945-4>.
- [8] Petrônio L. Braga and Adriano L. I. Oliveira, Software Effort Estimation using Machine Learning Techniques with Robust Confidence Intervals, Seventh International Conference on Hybrid Intelligent Systems, IEEE 2007, 0-7695-2946-1/07.
- [9] Boyce Sigweni, Martin Shepperd, Tommaso Turchi, Realistic Assessment of Software Effort Estimation Models, 2016 ACM. ISBN 978-1-4503-3691-8/16/06.
- [10] Tim Menzies, Zhihao Chen, Jairus Hihn, Karen Lum, Selecting Best Practices for Effort Estimation, IEEE TRANSACTIONS ON SOFTWARE ENGINEERING, VOL. 32, NO. 11, NOVEMBER 2006, 0098-5589/06.
- [11] Mohammad Azzeh, Daniel Neagu, Peter Cowling, Improving Analogy Software Effort Estimation using Fuzzy Feature Subset Selection Algorithm, 2008 ACM, 978-1-60558-036-4/08/05.
- [12] Przemysław Pospieszny, Beata Czarnacka Chrobot and Andrzej Kobylinski, An effective approach for software project effort and duration estimation with machine learning algorithms, The Journal of Systems & Software (2017), doi: 10.1016/j.jss.2017.11.066.
- [13] Leandro L. Minku, Xin Yao, Ensembles and locality: Insight on improving software effort estimation, Information and Software Technology 55 (2013) 1512–1528.
- [14] Pichai Jodpimai*, Peraphon Sophatsathit and Chidchanok Lursinsap, Ensemble effort estimation using selection and genetic algorithms, Int. J. Computer Applications in Technology, Vol. 58, No. 1, 2018.
- [15] Ahmed BaniMustafa, Predicting Software Effort Estimation Using Machine Learning Techniques, 2018 8th International Conference on Computer Science and Information Technology (CSIT), DOI: 10.1109/CSIT.2018.8486222.
- [16] Abdelali Zakrani, Mustapha Hain, Abdelwahed Namir, Software Development Effort Estimation Using Random Forests: An Empirical Study and Evaluation, International Journal of Intelligent Engineering and Systems, Vol.11, No.6, 2018, DOI: 10.22266/ijies2018.1231.30.
- [17] Tuğçe Uğurlu Altuntas., S. Emre Alptekin, Software Development Effort Estimation by Using Neural Networks – A Case Study, International Journal of Computers, Volume 2, 2017, ISSN: 2367-8895.
- [18] Bilge Başkeleş, Burak Turhan, Ayşe Bener, Software Effort Estimation Using Machine Learning Methods, IEEE 2007, 1-4244-1364-8/07.
- [19] Jianfeng Wen, Shixian Li, Zhiyong Lin, Yong Hu, Changqin Huang, Systematic literature review of machine learning based software development effort estimation models, Information and Software Technology, 54 (2012) 41–59.
- [20] Prabhakar and Maitreyee Dutta, Application of machine learning techniques for predicting software effort, Elixir Comp. Sci. & Engg. 56 (2013) 13677-13682.
- [21] Jyoti Shivhare, Santanu Ku. Rath, Software Effort Estimation using Machine Learning Techniques, 2014, ACM 978-1-4503-2776-3/14/02.
- [22] Vlad-Sebastian Ionescu, Babes-Bolyai University, An approach to software development effort estimation using machine learning, 2017 IEEE, 978-1-5386-3368-7/17.
- [23] Younghee Kim, Keumsuk Lee, A Comparison of Techniques for Software Development Effort Estimating, SYSTEM INTEGRATION 2005.
- [24] Prabhakar, Maitreyee Dutta, Prediction of Software Effort Using Artificial Neural Network and Support Vector Machine, International Journal of Advanced Research in computer Science and Software Engineering, Volume 3, Issue 3, March 2013, ISSN: 2277 128X.
- [25] Yoav Freund and Robert E. Schapire, A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting, journal of computer and system sciences 55, 119139 (1997), article no. SS971504, 0022-0000/97.
- [26] Wan M.U. Noormanshah, Puteri N.E. Nohuddin, Zuraini Zainol, Document Categorization Using Decision Tree: Preliminary Study, International Journal of Engineering & Technology, 7 (4.34) (2018) 437-440.
- [27] Scikit-learn, <https://scikit-learn.org/stable/>
- [28] Leo Breiman, PREDICTION GAMES AND ARCING ALGORITHMS, Technical Report 504, December 19, 1997.
- [29] Data sets, <http://promise.site.uottawa.ca/SERepository/datasets-page.html>.
- [30] Data sets (china data set), <https://zenodo.org/record/268446#.XZ9vB1UzbIU>