

Performance Analysis of Machine Learning Algorithms in Prediction of Student Academic Performance

ABSTRACT

The advancement in technology has contributed largely to the application of data mining in education in recent times. However, selecting appropriate algorithm(s) to “mine” knowledge about educational data presents a difficult challenge to researchers and analyst. This paper contributes to the use of classification algorithms in academic performance prediction. The predictive ability of four popular algorithms; C4.5 Decision tree (CDT), Multilayer Perceptron (MLP), Naïve Bayes (NB) and Random Forest (RF) algorithms were compared. The models were built using student dataset from selected private senior high schools in Ghana. The comparative analysis of the algorithms was made based on their Accuracy, Recall, Specificity, F-Measure and Running time. On all the training and test ratios; 80:20, 70:30 and 10-fold cross validation, the results indicated that all the algorithms performed well in the classification. However, the Naïve Bayes algorithm performed significantly better than the MLP and CDT on some ratios. The running time of the NB, CDT and RF were the quickest while MLP took the longest time.

Keywords: Data Mining, Algorithms, Machine Learning, Classification, Prediction, Student Performance, Multilayer Perceptron Algorithm, Naïve Bayes Algorithm, C4.5 Decision Tree

1. INTRODUCTION

The academic performance of students to some extent is vital to their success in society. Students who are successful academically are more likely to have higher chances of getting good employment opportunities. These students are often well prepared to meet the demands of the constantly changing world than poor performing students. The various academic institutions also benefit from the success of their students as they can attract more students. As the number of students continue to increase, the data available has also increased exponentially. The use of these data by educational institutions has often been to make simple searches and prepare reports. However, more can be done with this data to help with decision making. Institutions can use knowledge extracted from data to determine or predict students' performance, attitude towards learning, possible school dropouts, etc. These possible benefits have contributed to the rise of data mining in the field of education. There are several data mining approaches that can be used on a dataset depending on the task that requires knowledge. Choosing a suitable modeling approach for a given task is important when building models for prediction. One of such methods is Classification which is a supervised learning method. This approach of data mining maps training instances to pre-determined classes. A training set of pre-determined classes is studied by a learning algorithm in the first phase of classification to develop a model and a test set is utilized in the second phases to measure the model's accuracy [9].

Many data mining algorithms such as Decision trees (DT) and Naïve Bayes (NB) have been used to solve various classification problems with varying results which allows for improvements by combining multiple techniques or altering the ratios of training and testing. Much work is still needed in classifying students' data due to its high dimensionality. This study seeks to assess the ability of four commonly used algorithms in prediction of students' academic success on a binary classification task.

We put forward the following contributions:

1. Pre-processing techniques employed on data from selected private senior high schools in Ghana.
2. Models built for prediction using C4.5 decision tree (CDT), Multilayer Perceptron (MLP), Naïve Bayes and Random Forest (RF)
3. Evaluation of the performance of the algorithms based on Accuracy, Specificity, Recall, F-Measure and Running time.

The findings of this study provide researchers and school administrators' direction in discovering and selecting algorithm(s) to classify student datasets in order to identify good and poor students. It also provides insight into the application of classification algorithms on real-world sets and bridges the gap between theoretically projected outcomes against those that were observed.

2. RELATED WORKS

Educational Data Mining makes it possible to find answers to fundamental issues about education using data relating to specific educational contexts [16]. Studies have been done in this area because of the relevance of this field to stakeholders in discovering answers to questions about education using data related to educational settings and making future decisions.

Nguyen et al. [12] assessed the performance of Bayesian networks and Decision trees in prediction of undergraduate and postgraduate student academic achievement using data from two educational institutions. The decision tree algorithm achieved the best accuracy in

predicting the performance of students. It was found that decision trees were suitable for finding students who excelled academically.

In a study by Prabha and Shanavas [15], the performance of Naïve Bayes, Multilayer Perceptron, ZeroR, C4.5 decision tree and Random Forest were compared using dataset collected from “Maths Tutor”, an e-learning tool to find the best algorithm in predicting students’ performance. The records of 120 students from the sixth standard were used to evaluate the performance of the algorithms using five (5) different criteria; number of correctly classified instances, mean absolute error and RMSE rates, the time taken to build the model and the ROC area. The 10-fold cross validation was used for classification. From the result the Multilayer Perceptron and C4.5 decision tree achieved the highest accuracy and ROC curve weight average of 1, however, the error rates and time taken by the C4.5 decision tree were significantly lower than the MLP [15].

Kabakchieva [10] evaluated algorithms from four different families in predicting the performance of students at a Bulgarian university based on their personal and pre-university characteristics. The Naïve Bayes, Bayes net, OneR, JRip, C4.5 decision tree and K-Nearest Neighbor algorithms were applied in classifying students into variety of classes. To evaluate the models, the True Positive Rate and Precision were used as criteria. The results indicated that C4.5 decision tree was the best algorithm as it achieved the highest overall accuracy. In another study [5], Cortez and Silva concluded that Decision Tree, Random Forest, and Neural Networks had a high predictive accuracy in classifying students’ performance into two and five classes. They compared the predictive accuracy of these four algorithms in predicting students’ failure in Mathematics and Portuguese [5].

3. METHODOLOGY

3.1 Study Area and Data Collection

The study was carried out in the Kwahu West Municipal District of Ghana. The municipal lies between latitudes 6°30’ North, and 7° North and longitudes 0° 30’ West and 1° West. It has an approximate total land area of 414km². There are 17 communities in the municipal with Nkawkaw as its capital. Five (5) private senior high schools were selected randomly from the municipal. Data was obtained from 456 students of the senior high schools. The dataset consisted of students’ socio economic and prior enrollment variables. Table 1 shows the detailed description of the dataset.

Table 1. Description of Variables

Variable	Description	Possible Values	Data Type
SSG	Junior High School Grade in Social Studies	1, 2, 3, 4, 5, 6, 7, 8, 9	Numeric
SG	Junior High School Grade in Science	1, 2, 3, 4, 5, 6, 7, 8, 9	Numeric
EG	Junior High School Grade in English Language	1, 2, 3, 4, 5, 6, 7, 8, 9	Numeric
MG	Junior High School Grade in Mathematics	1, 2, 3, 4, 5, 6, 7, 8, 9	Numeric

fInc	Income level of family	High, Medium, Low	Nominal
mumEdu	Education level of mother	No formal education, Primary, Junior High School, Secondary, Tertiary	Nominal
fatherEdu	Education level of father	No formal education, Primary, Junior High School, Secondary, Tertiary	Nominal
mStatus	Parents' marital status	Single, Married	Nominal
fSize	Size of family	<=3, >3	Nominal
Sex	Sex of student	Male, Female	Nominal
eStatus	Employment status of parent	Employed, Unemployed	Nominal

3.2 Classification Framework

In order to "mine" information from data, several important processes must be completed. A data mining process may be based on a specific framework. Figure 1 presents the study's framework. This consist of handling data with missing values and outliers, feature selection, implementation of algorithms, comparison, and results.

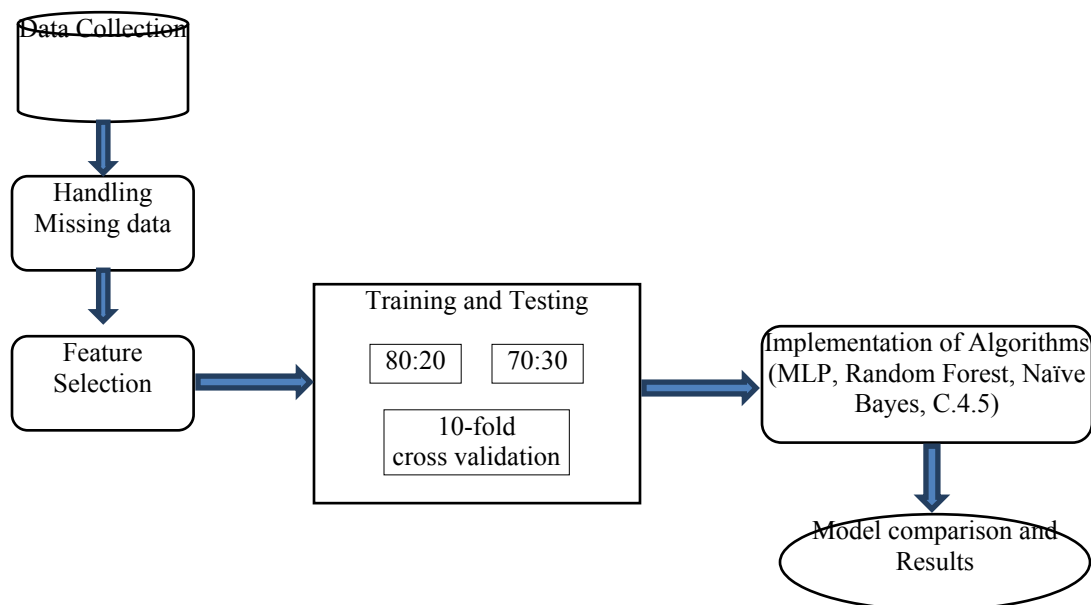


Fig. 1. Classification Framework

3.3 Data Pre-processing

In order to "mine" information from data, several important processes must be completed. A data mining process may be based Each figure should have a caption. The caption should be concise and typed separately, not on the figure area. Figures should be self-explanatory. Information presented in the figure should not be repeated in the table. All symbols and

abbreviations used in the illustrations should be defined clearly. Figure legends should be given below the figures.

3.4 Feature Selection

Reducing total attributes in a dataset by removing unimportant variables can contribute to improved performance of data mining algorithms [8]. In addition to the full dataset (FSet), a reduced dataset (SetA) containing five highest ranked attributes were generated and used for the classification. The value of Information gain (IG) of attributes were used to rank them. To calculate Information gain, the posterior entropy was subtracted from the Class entropy. Entropy measures the degree of "impurity". When it's close to 0, it suggests the dataset contains less impurity. The entropy is reduced the most by an attribute that contains the most information. Information gain (IG(Y)) and Entropy were calculated using the formula:

Where E_j is the Entropy of the class and $E_{j|X}$ is the Posterior Entropy.

Where X is the data sample and p_j is the proportion of X in respect to class j.

3.5 Training, Testing and Implementation of Algorithms

To carry out the experiments, models were built from two datasets; FSet and SetA, each containing different combinations of variables using Decision tree, Multilayer Perceptron, Naïve Bayes and Random Forest algorithms. The two distinct datasets for the study are shown in Table 2. The models were built by the algorithms using training and test ratios of 80:20, 70:30 and 10-fold cross validation. In building models using the 80:20 and 70:30 ratios, the dataset was split into two (training: testing) for training and testing purposes.

The 10-fold cross validation divides dataset into 10 mutually exclusive subsets of approximately equal size. The training and testing of the algorithm are done 10 times, with one subset serving as the test set and the remaining nine serving as the training set each time. To calculate for the classifier's accuracy, the total number of accurate cases is divided by the total classification. Weka was used for the implementation of the algorithms. The details of each algorithm are presented in this session.

Table 2. Datasets for Classification

Set	Number of Variables	Variables
FSet	11	EG, SSG, SG, MG, eStatus, mumEdu, fatherEdu, mStatus, flnc, fSize, Sex
SetA	5	EG, SSG, SG, MG, mumEdu

3.5.1 C4.5 Decision Tree

The C4.5 algorithm learns and generates trees from a training dataset that is used to classify new datasets. The C4.5 algorithm comes with added features and methods which make it easier to create simpler trees without compromising accuracy. It also handles both continuous and discrete data and uses pessimistic pruning to improve classification accuracy by removing non-essential branches from the decision tree. The root node of the

tree, which is the Class's best splitting attribute was expanded and used to partition samples into many parts.

The information gain was computed for each variable and the one with the highest value was chosen as the root node. When selecting the next split attribute, the same procedure was utilized, without considering previously used split attributes. This was repeated until all the attributes had been used or the classification was completed. Overfitting was avoided by pruning or removing non-essential nodes. Pruning was done after generating the trees by substituting a child node for a parent node if the error rate validation does not reduce.

3.5.2 Naïve Bayes

The Bayes theorem provides the foundation for the Naive Bayes algorithm which assumes a high level of independency. It assumes that attributes are conditionally independent. In supervised settings, NB performs simple calculations and learns quickly. The NB algorithm predicts the class of unknown dataset using the formula:

Where $P(s)$ = the posterior probability of s given s , $P(s)$ = the likelihood of s occurring given, $P(s)$ = the likelihood of s occurring and $P(s)$ = the likelihood of s occurring.

3.5.3 Multilayer Perceptron

Multilayer Perceptron algorithm consists of layers that cycle through weights of supplied data and maps them to correct outputs. It is from the family of Artificial Neural Networks. The MLP is a widely used algorithm for classification that locates weight that capture hidden Input / Output mapping in a set of instances. The number of weights to be updated by the algorithm was set at 0.3. This enables the training of the network while the weight amount is controlled. The momentum rate was set at 0.2 to prevent the network from reaching a local minimum too soon and conversely avoid overshooting the function's global minimum.

3.5.4 Random Forest

Random Forest is an ensemble machine learning algorithm that uses several classification trees, thus earning it the name, Random Forest. RF algorithm is considered to be a highly rigorous learning algorithm in recent times, in spite of its biasness towards classification of high-level categorical attributes of datasets consisting of varied levels of categorical data [6]. The algorithm functions by first gathering from the given dataset, random samples. Next, each sample is used to build a decision tree. The predicted results are then voted on to select the classification with the most votes. In training and testing of the two datasets, the maximum depth of trees configuration of the algorithm was set to 16.

3.6 Model Performance Evaluation

The algorithms were evaluated using the Accuracy, Recall, Specificity, F-measure and Running time. All these metrics, except for Running time were calculated using the classification results in the Confusion Matrix (see Table 3), where TP is True Positives, TN is True Negatives, FP is False Positives and FN is False Negatives.

Table 3. Confusion Matrix

Classified as			
Positive	Negative		
TP	FN	Positive	← Actual Class
FP	TN	Negative	

3.6.1 Accuracy

Accuracy is the proportion of correctly classified instances of a model. It was calculated as:

3.6.2 Recall

Recall is the percentage of correctly classified positive class samples. It was calculated as:

3.6.3 Specificity

Specificity is the percentage of correctly classified negative instances. Specificity was computed as:

3.6.4 F-Measure

F-Measure is the harmonic mean of the Precision and Recall. The Precision and Recall values were used in calculating F-Measure. The formulas are as follow:

4. RESULTS AND DISCUSSION

The Weka Experimenter was used to run the full datasets as well as the reduced set using the four algorithms to see if a base algorithm was considerably better or worse than the others. The Weka Experimenter examines algorithms outputs based on a set criterion. The Naïve Bayes algorithm was selected as the base algorithm against which the C4.5 decision tree, Random Forest and MLP algorithms were compared. A two-tailed t-test [11] with significant level of 0.5 was used for the experiment. A classifier which was significantly better than the base classifier was tagged with (v) and if worse, tagged with asterisk (*). No tag simply means the classifier was neither better nor worse as compared to the base classifier.

4.1 Classification Algorithms on FSet

Fig. 2a, 2b and 2c visualize the performance of the algorithms on FSet in terms of Accuracy, Recall, Specificity and F-Measure using training and test ratios of 80:20, 70:30 and 10-fold cross validation respectively. The results show that all four algorithms performed well on the dataset. The highest accuracy of the NB and RF were 86.76% and 80.95% respectively at 10-fold cross validation. The MLP and CDT's highest accuracy was 81.11% and 85.08%,

which was achieved at 80:20 and 70:30 ratios respectively. The classification accuracy of the algorithms on all the three ratios were significantly the same except for MLP at 10-fold cross validation, which was significantly worse than the base algorithm, NB.

In terms of the True positives or Recall, neither algorithms performed worse nor better than the base algorithm, NB. The algorithms showed very high results on all the different training and testing ratios. From the results, the NB and RF achieved highest Recall of 89% and 87% respectively at 10-fold cross validation. The MLP and C4.5 performed best in Recall at 70:30 ratio with 83% and 91% respectively.

The performance of the NB, MLP, CDT and RF on all the training and testing ratio were consistent, with the algorithms achieving a very high rate of 80% and above. However, the MLP performed significantly worse than NB at 10-fold cross validation. It achieved an F-Measure of 81% as compared to the 89% of the NB algorithm. The results also shows that Specificity for NB, CDT, MLP and RF were not significantly different, except at 80:20 ratio where the CDT was significantly lower than the NB. The CDT achieved a Specificity of 71% as compared to the 86% by the base algorithm.

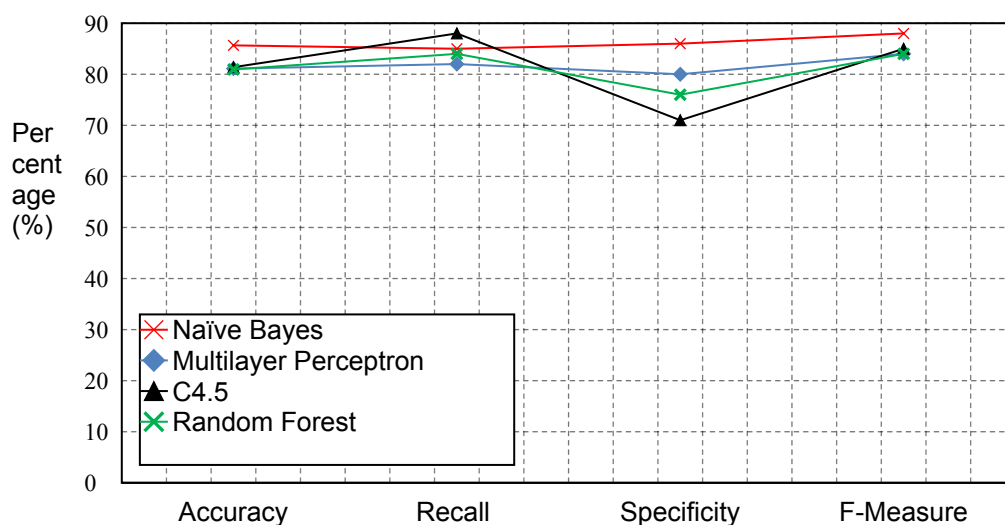


Fig 2a: Performance of Algorithms on 80:20 ratio (FSet)

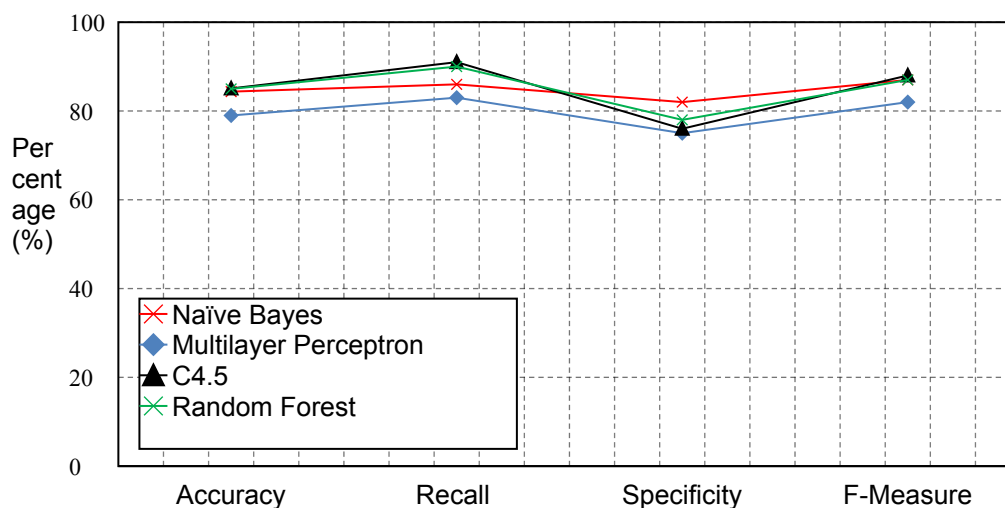


Fig 2b: Performance of Algorithms on 70:30 ratio (FSet)

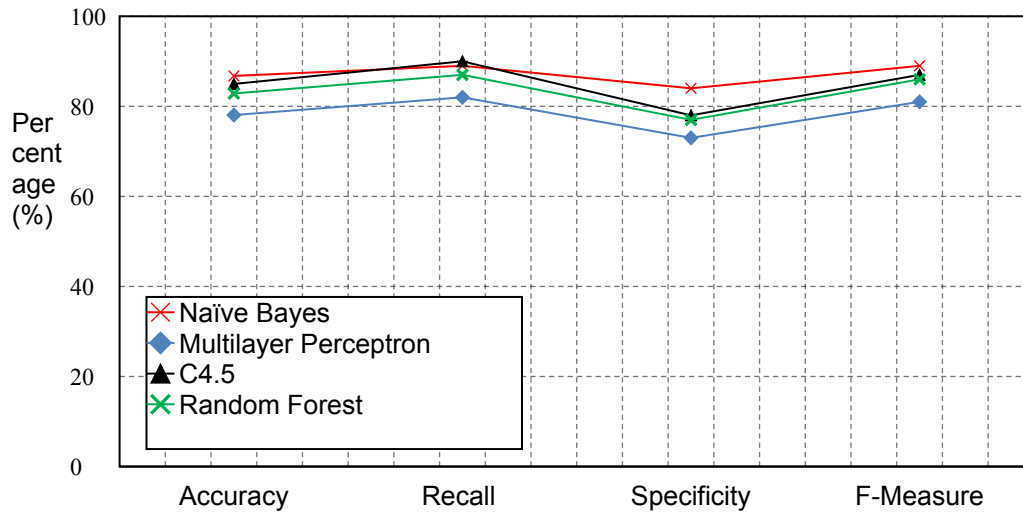


Fig 2c: Performance of Algorithms on 10-fold cross-validation (FSet)

The running time of the algorithms as presented in Table 4 were also significantly the same for NB, CDT and RF. However, the MLP took a longer running time than all the algorithms on all the ratios.

Table 4. FSet Running Time Results (in seconds)

Model	80:20	70:30	10-fold cross- validation
Naïve Bayes	0.01	0.01	0.01
Multilayer Perceptron	0.15*	0.14*	0.17*
C4.5	0.01	0.1	0.01
Random Forest	0.02	0.01	0.02

4.2 Classification Algorithms on SetA

The results for SetA are presented in Fig 3a, 3b and 3c. The results show that the classification accuracy of the algorithms were significantly the same. Each of the algorithms performed very well in the binary classification. At 70:30 ratio, the results of the NB and RF were 88.76% and 83.38% respectively which were the classifiers' best accuracy. The MLP performed best in Accuracy (82.63%) at 80:20 whereas CDT's was 83.94% at 10-fold cross validation.

The Recall and F-Measure achieved by the four algorithms were neither significantly worse nor better than the other. The NB, MLP, CDT and RF performed very well in classifying the dataset. In terms of F-Measure, all the algorithms performed better at 10-fold cross validation. The results of the NB, MLP, CDT and RF were 90%, 85%, 87% and 85%

respectively. The results also show that the highest Recall by all the algorithms were achieved at 10-fold cross validation. The Recall for NB, MLP, CDT and RF were 89%, 86%, 89% and 86% in that order.

The Specificity of the algorithms on all the ratios were neither significantly worse nor better, with the exception of MLP which performed slightly lower (71%) than the base algorithm NB (93%) at 80:20. All the other algorithms did very well classifying True negatives. The MLP, in terms of running time performed significantly worse than the three algorithms as shown in Table 5.

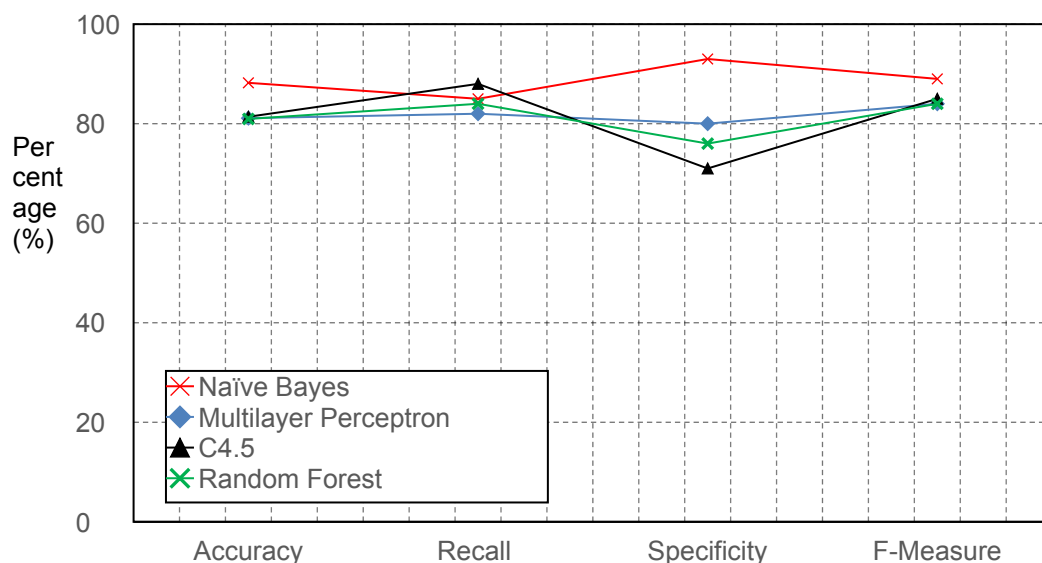


Fig 3a: Performance of Algorithms on 80:20 ratio (SetA)

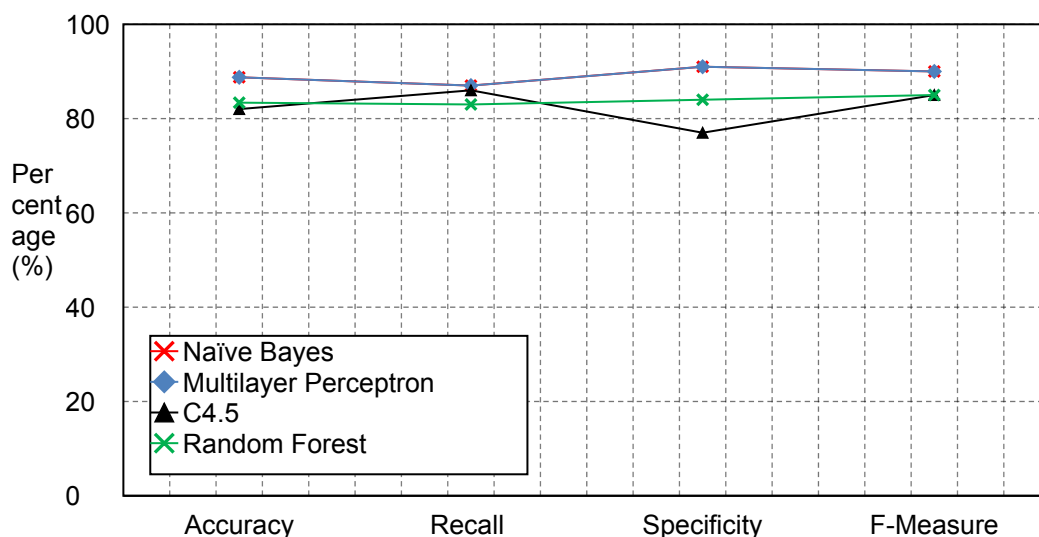


Fig 3b: Performance of Algorithms on 70:30 ratio (SetA)

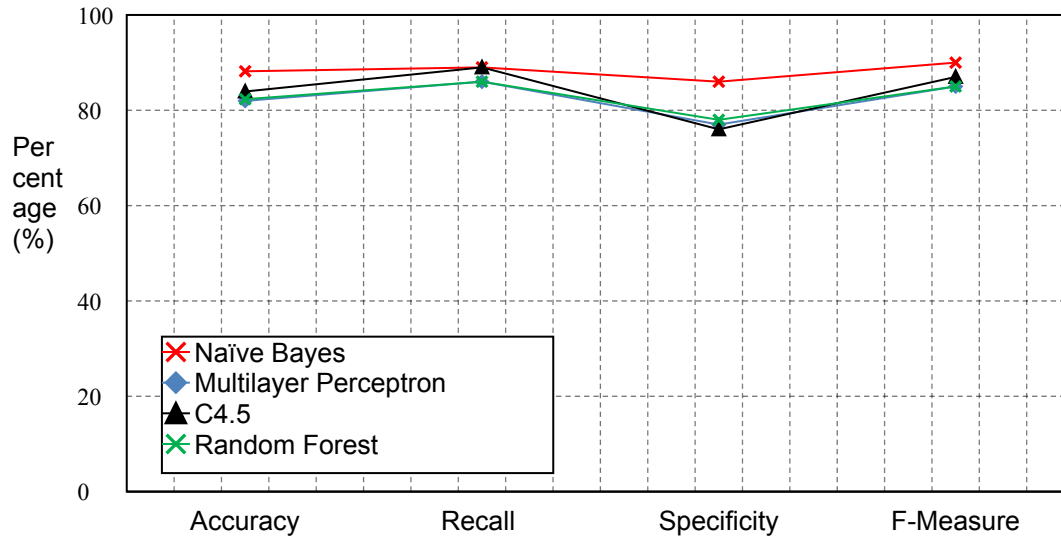


Fig 3c: Performance of Algorithms on 10-fold cross validation (SetA)

Table 5. SetA Running Time Results (in seconds)

Model	80:20	70:30	10-fold cross- validation
Naïve Bayes	0.01	0.01	0.01
Multilayer Perceptron	0.05*	0.04*	0.05*
C4.5	0.01	0.1	0.01
Random Forest	0.01	0.01	0.01

4.3 Critical Analysis of Classification Algorithms

The machine learning classifiers utilized in this work performed well on the two distinct datasets (FSet and SetA) using the different training and testing ratios. The NB, particularly, performed better on the reduced dataset (SetA) in classifying the Accuracy, Recall, Specificity and F-Measure on all the ratios. This is consistent with the results of [2, 3, 4, 13]. Several reasons may have contributed to the performance of the NB algorithm. Even though most probability estimates are poor, it is possible for the right class to have the highest estimate, resulting in the accurate classification. The NB also has the capability to handle categorical data well, especially in small datasets than most algorithms [1].

From the result, there was no significant difference in the running time of the NB, RF and CDT. This is consistent with the results of [3] where decision tree was found to be very fast in terms of running time in building trees. The speed of NB could be because it has a linear

learning time [14]. The results from this study suggest that NB, RF and CDT are computationally and resource-wise very efficient. Across all training sets, the MLP algorithm took much longer than the others. This could be because setting up MLP's network and learning from a training set requires more memory and runtime.

5. SUMMARY AND CONCLUSION

The study was aimed at evaluating the predictive ability of four popular algorithms; C4.5 Decision Tree, Random Forest, Naïve Bayes, and Multilayer Perceptron algorithms on a binary classification task. Assessment metrics which include Accuracy, Recall, F-Measure, Specificity and Running time were used in comparing the performance of the algorithms in predicting students' academic performance. Data was obtained from 456 students from senior high school in the eastern region of Ghana.

To ensure effective prediction of the four algorithms, the dataset was thoroughly analyzed and pre-processed before the classification stage. A new dataset (SetA) which consists of the 5 highest attributes with Information Gain was generated from the full dataset (FSet) and together used for the classification. The training and testing ratios used for building models by the algorithms were 80:20, 70:30 and 10-fold cross validation.

The findings from the experiments showed that all the algorithms exhibited a high level of performance in the classification. On the full dataset, all the algorithms did well on the 3 training and testing ratios performed. However, the Multilayer Perceptron achieved lowest performance. On the reduced dataset (SetA), the performance of the algorithms improved further which is suggestive that using optimal attributes of a dataset can help increase the performance of the algorithm. Generally, the Naïve Bayes algorithm was significantly better than the Multilayer Perceptron, C4.5 and Random Forest algorithms which is suggestive that it is a good candidate for predicting students' academic performance. More investigation needs to be done with a hybrid model at various training and testing ratios and other classification algorithms for predicting students' academic performance.

6. RECOMMENDATIONS FOR FUTURE STUDIES

From the conclusions drawn in the study, the following are recommended for further study:

1. Other variables that affect students' performance should be included to better evaluate the performance of the classifiers. This study considered only selected students' pre-enrollment and socio-economic variables. The performance of the classification algorithms using factors such as students affect, motivation and other non-cognitive variables may be considered.
2. In future studies, different feature selection methods other than the one tested may be employed in ranking and selecting attributes to improve prediction. For example, wrapper methods of selecting attributes may be studied.
3. Studies can be carried out on implementing an algorithm whose functionality can be switched to carry out multiple specific outcomes. Moreover, the use of multiple algorithms in an ensemble method may be used for increased performance in terms of prediction.

REFERENCES

- [1] Al-Barrak M, Al-Razgan M. 2016. Predicting students' performance through classification: A case study, *Journal of theoretical and Applied Information Technology*. 75(2): 167-175.

- [2] Anuradha, C, Velmurugan T. 2015. A Comparative Analysis on the Evaluation of Classification Algorithms in the Prediction of Students Performance, *Indian Journal of Science and Technology*, 8(15):1-12
- [3] Ashari A, Paryudi I, Tjoa A. M. 2013. Performance Comparison between Naïve Bayes, Decision Tree and k-Nearest Neighbor in Searching Alternative Design in an Energy Simulation Tool, *International Journal of Advanced Computer Science and Applications*, 4(11): 33-39.
- [4] Bharti T. 2014. Data Mining with Big Data Using C4.5 and Bayesian Classifier, *International Journal of Advanced Research in Computer Science and Software Engineering*, 4(8).
- [5] Cortez P, Silva A. M. G. 2008. Using Data Mining to Predict Secondary School Student Performance. In A. Brito, & J. Teixeira (Eds.), *Proceedings of 5th Annual Future Business Technology Conference*, Porto, 5-12.
- [6] Cutler A, Cutler D. R, Stevens J. R. "Random forests," in *Ensemble Machine Learning*. New York, NY, USA: Springer, 2012, pp. 157 175
- [7] Fayyad U. M, Piatetsky-Shapiro G, Smyth P. 1996. From data mining to knowledge discovery: An overview. In U.M. Fayyad, et al. (eds.), *Advances in Knowledge Discovery and Data Mining*, 1-35. AAAI/MIT Press.
- [8] Guyon I, Elisseeff A. 2003. An introduction to variable and feature selection. *The Journal of Machine Learning Research*, 3, pp. 1157–1182.
- [9] Han J, Kamber M. 2006. *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers, San Francisco
- [10] Kabakchieva D. 2013. Predicting student performance by using data mining methods for classification, *Cybernetics, and Information Technologies*, 13(1): 61-72.
- [11] Nadeau C, Bengio Y. 2000. Inference for the generalization error. *Advances in Neural Information Processing Systems 12*, MIT Press.
- [12] Nguyen T. N, Janecek P, Haddawy P. 2007. A comparative analysis of techniques for predicting academic performance. *Frontiers in Education Conference - Global Engineering: Knowledge without Borders, Opportunities without Passports, 2007. FIE '07. 37th Annual*, pp.T2G-7, T2G-12.
- [13] Osmanbegovi E, Suljic M. 2012. Data mining approach for predicting student performance, *Journal of Economic and Business*, X (1): 3-12.
- [14] Pazzani M, Billsus D. 1997. "Learning and Revising User Profiles: The Identification of Interesting Web Sites," *Machine Learning*, Vol. 27, 313 – 331.
- [15] Prabha S. L, Shanavas A. R. M. 2015. Performance of Classification Algorithms on Students' Data – A Comparative Study, *International Journal of Computer Science and Mobile Applications*, 3 (9), pp. 1-8.
- [16] Romero C, Ventura S. 2010. Educational data mining: a review of the state of the art. *IEEE Trans. On Syst. Man and Cybern. Part C Appl. Rev.* 40(6), 601–618.

