**STUDENTS' PERFORMANCE PREDICTION USING CLASSSIFICATION ALGORITHMS**

**ABSTRACT**

It is imperative to analyze educational data especially as it relates to students' performance. Educational institutions need to have a fairly accurate prior admitted students' knowledge to predict their future academic performance. This helps to identify the good students and also provides an opportunity to pay attention to and improve those who would possibly not perform too well. As a solution, this paper proposed a system which can predict the performance of students from their previous academic record using concepts of data mining techniques under Classification. The dataset containing information about students, such as gender, age, SSCE grade, UTME score, post UTME score and grade in students first year. ID3 (Iterative Dichotomiser 3) and C4.5 classification algorithms was applied on the data to predict the academic performance of students in future examinations.

## 1. INTRODUCTION

Classification method is the most frequent technique which is used to classify data set. Presently classification is used in several fields such as education, industrial, medical and other many places [1]. Classification is basically a data mining technique in which some input pattern is applied to get desired output by using any classification algorithm. The task of developing effective academic prediction system is a critical issue for educators [2].On yearly basis, higher institutions admit students from different locations and educational background with varying scores in entrance examinations into various departments. Previous studies have revealed that various factors are responsible for students' failure which includes low socio-economic background, student's intellectual capacity, school and home environment, or the support given by parents and other family members [3]. Methodologically, analysis of the previous academic performance of students admitted can be used to better predict their future performance using the concept of machine learning. In this regard, the data of students enrolled in 2008/2009 academic session of Joseph Ayo Babalola University was obtained and used in this study. This data includes attributes such as gender, age, SSCE grade, UTME score, post UTME score and grade in student's first year, category and admission type. Two decision tree algorithms (ID3 and C4.5

algorithms) were used to predict the futureperformance of the student using the dataset. The results of the two algorithms were then compared to determine the most effective algorithm for the prediction.

## 2. LITERATURE REVIEW

A review of relevant literatures was carried out.Abeer and Elaraby [4] analysed previously enrolled students' data in a specific course program across 6 years (2005–2010), with multiple features collected from the university database. The work predicted the students' final grades in the particular course program. Pandey and Pal [5] presented a data mining approach to classify students' according to performers or underperformers class using Naïve Bayes algorithm, classify. Bhardwaj and Pal [6] did a comparative study to test multiple decision tree algorithms on an academic dataset in order to classify the student's academic performance. The work primarily concentrates on choosing the best decision tree algorithm from among commonly used decision tree algorithms, and then provides a standard for them individually. It was discovered that the CART decision tree technique performed reasonably better on the dataset used for testing, that was obtained based on the accuracy and precision produced at the validation stage. Livieris, *et al.* [7] developed an Artificial Neural Network (ANN) classifier to predict the performance of students in Mathematics. From their experiments they discovered that the modified spectral Perry trained artificial neural network performs better classification compared to other classifiers.  Kotsiantis, *et al* [8] explored machine learning techniques for dropout prediction of students in distance learning. This study contributed in that it carved the path for educational data mining and one of the first works to implemented machine learning methods in an academic environment. Their algorithm was fed on demographic data and several project assignment rather than class performance data to make prediction of students. Moucary, *et al.* [9] applied a hybrid technique on K-Means Clustering and Artificial Neural Network for students who are pursuing higher education while adopting a new foreign language as a means of instruction and communication. Firstly, Neural Network was used to predict the student's performance and then fitting them in a particular cluster which was form using the K-Means algorithm. This clustering helped in serving a powerful tool to the instructors to identify students capabilities during their early stages of academics. Hongsuk, *et al*. [10] develop a Deep Neural Network supervised model to estimate link based flow of traffic conditions. A Traffic Performance Index was used for logistic regression to distinguish between a congested traffic condition and a non-congested traffic condition. The 3 layer model was able to estimate the congestion with a 99% of accuracy.  Yadavto estimate the congestion with 99% accuracy. Yadav and Pal [11] proposed a

prediction model for students' performance based on data mining methods with some few features called student's behavioral features. The model was evaluated using three different classifiers; Naïve Bayesian, Artificial Neural Network and Decision Tree. Random Forest, Bagging and Boosting were used as ensemble methods to improve the classifier's performance. The model achieved up to 22.1% more in accuracy compared when behavioral features were removed. It increased up to 25.8% accuracy after using the ensemble methods.

## 3. METHODOLOGY

The methodology of this study is composed of: identification of the required variables for students performance prediction, the collection and preparation of data, formulation of the predictive models using the supervised machine learning algorithm (Decision Tree), simulation of the predictive models using the WEKA simulation environment and the performance evaluation metrics applied during model validation for the predictive models performance evaluation.

### a. Data Collection and Preparation

The dataset used in this study was obtained from the academic record office, Joseph Ayo Babalola University. The data was anonymously obtained without any bias. Personal and academic record of students admitted in 2008/2009 into the university from Six (6) major departments namely: Computer Science (CSC), Accounting (ACC), Political Science (POL), Microbiology (MCB), Economics (ECO), Business Administration (BUS) was used. The size of the dataset is 100 records.

### b. Model Formulation

In this study, decision tree algorithm was used in formulating the model for prediction because the pattern explaining the link between the attributes identified (input attributes) and the student's performance (the target attribute) was needed. The pattern identified was then converted into a set of rules that can assist in making informed decisions regarding the performance of students. In formulating a predictive model using supervised machine learning algorithm, a mapping function is used to easily state the general expression. The dataset $S$ which consists of the records of students containing fields representing the set of classification factors (i number of input variables for j students), $X_{ij}$ alongside the respective target variable (student's performance) denoted by the variable $Y_j$ – the student's performance for the j[th] individual in the j records of data for the study. The mapping function that defines the link between the classification features and the target attribute – classification of student's performance is given in equation (1).

$$\varphi : X \rightarrow Y \tag{1}$$

$$defined\ as : \varphi(X) = Y$$

93    The equation shows the relationship between the set of classification factors represented by a vector, $X$

94    consisting of the values of i variables and the label $Y$ which defines the student's performance – First,

95    Two-1, Two-2 and Third of each student as expressed in equation (2). Assuming the values of the set of

96    variables for a student is represented as $X = \{X_1, X_2, X_3, \ldots\ldots, X_i\}$ where $X_i$ is the value of each

97    variablei = 1 to j; then the mapping $\varphi$ which represents the predictive model for student's performance

98    maps the variables of each one to their corresponding student's performance according to equation (2).

$$\varphi(X) = \begin{cases} First \\ Two-1 \\ Two-2 \\ Third \end{cases} \tag{2}$$

99    The decision trees developed for the performance of students was used to propose a set of rules that can

100    be used to determine the student's performance directly just by observing the value of the variables

101    identified by the model and the succession of events. Also, the set of attributes identified in the final

102    decision trees model for student's performance are the variables which have the most relevant

103    importance to the determination of each student's performance. It was proposed to be given much

104    consideration during performance assessment of students.

105    For the training dataset, $S$ is a set containing $S_1, S_2, \ldots, S_j$ of samples that have been classified already

106    of the students' records which consist the values of their variables, $X = \{X_1, X_2, \ldots, X_i\}$ together with

107    the classification of student's performance, $Y = \{First, Two-1, Two-2, Third\}$ such that, $S =$

108    $(X, Y)$ for all students from 1 to j.

109    In this work, C4.5 and ID3 classification algorithms were used for the predictive model formulation. The

110    two conditions used by the C4.5 decision trees in developing its decision trees are stated in equations (3)

111    and (4) defined as the information gain and the split criteria respectively. Equation (3) is used in

112    determining which attribute is used to split the dataset at every iteration while equation (4) is used to

113    determine which of the selected attribute split is most effective in splitting the dataset after attribute

114    selection by equation (3).

$$IG(X_i) = H(X_i) - \sum_{t \epsilon T} \frac{|t|}{|X_{ij}|} \cdot H(X_i) \tag{3}$$

115    where:

$$H(X_i) = -\sum_{t\epsilon T} \frac{|t, X_i|}{|X_{ij}|} \cdot \log_2 \frac{|t, X_i|}{|X_{ij}|}$$

$$Split(T) = -\sum_{t\epsilon T} \frac{|t|}{|X_{ij}|} \cdot \log_2 \frac{|t|}{|X_{ij}|} \hspace{2cm} (4)$$

116    *T is the set of values for a given attribute $X_i$.*

117    The simulation of the predictive model was performed in WEKA environment.

118    **c. 10-fold Cross Validation (Model Validation)**

119    Cross-validation procedure was used in this work. This entails splitting the entire datasets into some

120    folds (or partitions). Each fold was selected for testing, with the remaining k – 1 fold; the subsequent

121    fold was used for testing with the remaining fold (together with the first fold used) used for training,

122    pending when all k partitions had been selected for testing. The error rate recorded from each process

123    was added up with the mean the mean error-rate recorded

124    **d. Performance Evaluation of Model Validation Process**

125    In the course of evaluating the predictive model, the models' performance was quantified using some

126    metrics. Basically, four (4) parameters must be known from the model testing of predictions made by the

127    classifier during model testing. These parameters are: true positive (TP), true negative (TN), false

128    positive (FP) and false negative (FP). TPs refers to the accurate prediction of positive cases, TNs refers

129    to the accurate prediction of negative cases, and FPs indicates the negative cases predicted as positives

130    while FNs indicates the positive cases predicted as negatives. The results were then obtainable on

131    confusion matrix which is a 4 x 4 matrix table owing to the four (4) labels of the output class (see Figure

132    1). Correct classifications were plotted along the diagonal from the north-west position for first

133    predicted as first (A), 2-1 predicted as 2-1 (F), followed by 2-2 predicted as 2-2 (K) and third predicted

134    as third (P) on the south-east corner (also called true positives and negatives). The incorrect

135    classifications were plotted in the remaining cells of the confusion matrix (also called false positives).

136    Also, the actual first cases are A+B+C+D, actual 2-1 cases are E+F+G+H, actual 2-2 cases are I+J+K+L

137    while actual third are M+N+O+P and the predicted first are A+E+I+M, 2-1 are B+F+J+N, predicted 2-2

138    are C+G+K+O and predicted third are D+H+L+P.

139    The developed model was validated with a number of performance metrics based on the values

140    of $A – P$ in the confusion matrix for each predictive model. They are presented as follows.

141      a. Accuracy: the total number of correct classification.

$$Accuracy = \frac{A + F + K + P}{total\_cases} \tag{5}$$

142      b. TP rate (recall/sensitivity): the amount of actual cases accurately classified.

$$TP_{first} = \frac{A}{A + B + C + D} \tag{6}$$

143

144

| FIRST | 2-1 | 2-2 | THIRD | |
|---|---|---|---|---|
| **A** | **B** | **C** | **D** | FIRST |
| **E** | **F** | **G** | **H** | 2-1 |
| **I** | **J** | **K** | **L** | 2-2 |
| **M** | **N** | **O** | **P** | THIRD |

145

146

147      **Figure 1: Confusion Matrix**

148

$$TP_{2-1} = \frac{F}{E + F + G + H} \tag{7}$$

$$TP_{2-2} = \frac{K}{I + J + K + L} \tag{8}$$

$$TP_{third} = \frac{P}{M + N + O + P} \tag{9}$$

149      c. FP (false alarm/1-specificity): the amount of negative cases inaccurately classified as positive.

$$FP_{first} = \frac{E + I + M}{actual_{2-1} + actual_{2-2} + actual_{third}} \tag{10}$$

$$FP_{2-1} = \frac{B + J + N}{actual_{first} + actual_{2-2} + actual_{third}} \tag{11}$$

$$FP_{2-2} = \frac{C + G + O}{actual_{first} + actual_{2-1} + actual_{third}} \tag{12}$$

$$FP_{third} = \frac{D + H + L}{actual_{first} + actual_{2-1} + actual_{2-2}} \tag{13}$$

150      d. Precision: the proportion of predictions that are correct.

$$Precision_{first} = \frac{A}{A+E+I+M} \qquad (14)$$

$$Precision_{2-1} = \frac{F}{B+F+J+N} (15)$$

$$Precision_{2-2} = \frac{K}{C+G+K+O} (16)$$

$$Precision_{third} = \frac{P}{D+H+L+P} \qquad (17)$$

151  Using the aforementioned performance metrics, the performance of the predictive model for the
152  classification of student's performance was evaluated by validation, using a dataset. The TP rate and
153  precision lie within the interval [0, 1], accuracy within the interval of [0, 100] % while the FP rate lies
154  within an interval of [0, 1]. The closer the accuracy is to 100% the better the model, the closer the value
155  of the TP rate and precision is to 1 the better. While the closer the value of FP rate is to 0, the better.
156  Therefore, the evaluation of an effective model has a high TP/Precision rates and a low FP rates.

157  **4. RESULTS AND DISCUSSION**
158  Table 1 shows a description of the variables that were discretized and the nominal variables to which
159  they were converted to for clarity of model complexity. Afterwards, the pre-processed dataset was saved
160  in the acceptable format (attribute relation file format (.arff)) for the machine learning simulation
161  environment.

162  **Table 1: Student's performance data showing the discretized numeric variables**

| Name of Variable | Raw Label | Interval | Discretized Value |
|---|---|---|---|
| UTME Score | Numeric (0 – 400) | 1 to 100<br>101 to 200<br>201 to 300<br>301 to 400 | 1 – 100<br>101 – 200<br>201 – 300<br>301 – 400 |
| Age | Numeric (in years) | Less than 18 years<br>18 years and above | Below-18<br>18-above |
| SSCE Score | Numeric (0 – 30) | 1 to 10<br>1 to 20<br>21 to 30 | 1 – 10<br>11 – 20<br>21 – 30 |
| 100 Level Grade | Numeric (0.0 – 5.0) | Below 2.00<br>2.00 - 2.50 | Pass<br>Third |

| | | 2.50 – 2.49 | Two-2 |
| | | 3.50 – 4.49 | Two-1 |
| | | 4.50 – 5.00 | First |

Table 2 gives the narrative of the number of students with their individual classification of student's performance from the record of 47 student chosen for formulating and validating the model which were saved in the Student-Train-Data.arff file. The table shows that of the 100 students used; 2% had first class, 22% had second class upper, 61% had second class lower while 15% had third class degree by the time of graduation. The results showed that majority of the students had second class lower degrees amounting for about 70% of the student population selected for this study.

**Table 2:  Distribution of student performance among historical dataset**

| Student's Performance | Frequency | Percentage (%) |
|---|---|---|
| First | 2 | 2.0 |
| Two – 1 | 22 | 22.0 |
| Two – 2 | 61 | 61.0 |
| Third | 15 | 15.0 |
| Total | 100 | 100.0 |

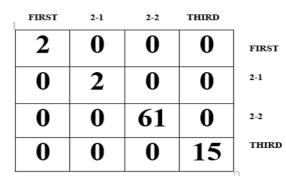**Results of Model Formulation and Simulation**

Following theidentification of the factors that are associated with student performance,  the  next  phase is  model  formulation  using  the aforementioned  decision trees  algorithms  available  in  the  WEKA environment.   The 10-fold cross validation  technique  was  used  in  evaluating  the performance  of the  developed  predictive  model  for  student performance  using  the  historical dataset used for training the model.  This process was performed  for  both  decision trees algorithm used with their respective performance compared for the most effective.

**a.   Results of model formulation and simulation using the ID3 algorithm**

The results of the formulation of the predictive model using the ID3 decision trees algorithm showed that a limited number of variables were the most important classification factors.  Identified variables in the order of their significance are:

185    • 100 level grade;

186    • Subject grades in core subjects such as physics, mathematics, and English;

187    • UTME score;

188    • Age at admission; and

189    • Student's gender.

190    The predictive model was formulated based on ID3 identified variables, using the results of the

191    simulation with the C4.5 algorithm in WEKA simulation environment. The ID3 was used to formulate a

192    tree that was adopted in deducing the set of rules used for the classification of student's performance.

193    Following the simulation of the predictive using the ID3 and C4.5 decision trees algorithm, after 10-fold

194    cross validation, the result of the performance evaluation of the model was recorded. The confusion

195    matrix used to interpret TP and TN alongside the FP and FN of the validation result is shown in Figure 2

196    and Figure 3 respectively. The results showed that out of the figure 2 actual first classes, all were

197    correctly classified, out of the 22 actual two-1, all were correct classified, out of the 61 two-2, all were

198    correctly classified and out of the 15 third class cases, all were correctly classified.  Hence, all 100

199    instances in the dataset were correctly classifier by the ID3 decision trees classifier meaning 100%
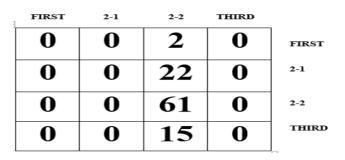
200    accuracy.

| FIRST | 2-1 | 2-2 | THIRD | |
|---|---|---|---|---|
| 2 | 0 | 0 | 0 | FIRST |
| 0 | 2 | 0 | 0 | 2-1 |
| 0 | 0 | 61 | 0 | 2-2 |
| 0 | 0 | 0 | 15 | THIRD |

201

202    **Figure 2:   Confusion matrix of performance evaluation using ID3**

203    **b.  Model Formulation and Simulation in C4.5 algorithm**

204    The C4.5 algorithm was also used to implement predictive model in the simulation environment.  From

205    the result, the algorithm could not identify the variables that were the most important factors of student's

206    performance.

207    The confusion matrix in figure 3 was used to evaluate the performance of the predictive model for

208    classification of student's performance.  The results further showed that using the C4.5 decision trees

209 algorithm to formulate the model for the classification of student's performance, all 61 two-2 cases were
210 correctly classified while all 2 first class cases, 22 two-1 cases and 15 third class cases were
211 misclassified as two-2 cases. Therefore, 61 out of the 100 students' instances were correctly classified
212 by the C4.5 decision trees classifier for the model development owing for an accuracy of 61%.

| FIRST | 2-1 | 2-2 | THIRD | |
|-------|-----|-----|-------|---|
| 0 | 0 | 2 | 0 | FIRST |
| 0 | 0 | 22 | 0 | 2-1 |
| 0 | 0 | 61 | 0 | 2-2 |
| 0 | 0 | 15 | 0 | THIRD |

213
214 **Figure 3: Confusion matrix of performance evaluation using C4.5**

215 **c.      Discussion of results**

216 The result of the performance evaluation of the machine learning algorithms are presented in Table 3
217 which presents the average values of each performance evaluation metrics considered for this study. For
218 the ID3 decision trees algorithm based on the results presented in the confusion matrix presented in
219 figure 3. The results showed that the TP rate which gave a description of the proportion of actual cases
220 that was correctly predicted was 1 which implied that 100% of the actual cases were correctly predicted;
221 the FP rate was 0 which implied that 0% of actual cases were not accurately classified while the
222 precision was 1 which implied that 100% of the predictions made by the classifier were correct.

223 For the C4.5 decision trees algorithm based on the results presented in the confusion matrix presented in
224 figure 3. The results showed that the TP rate was 1 for two-2 but 0 for first/two-1/third which implied
225 that 100% and 0% of the actual two-2 cases and first/two-1/third cases respectively were correctly
226 predicted; the FP rate was 1 for two-2 but 0 for first/two-1/third which implied that 100% and 0% of
227 actual cases were misclassified while the precision which gave a description of the proportion of
228 predictions that were correctly classified was 0.61 for two-2 but 0 for first/two-1/third which implied
229 that 61% and 0% of the predictions made by the classifier were correct.

230 From the study, it was discovered that ID3 decision trees algorithm was able to classify the performance
231 of students by graduation better than the C4.5 decision trees algorithm. The ID3 decision trees
232 algorithm was able to accurately classify all cases of students with a value of 100% showing that it had
233 the capacity to identify the complex patterns that existed within the dataset than the C4.5 decision trees

234 algorithm. The variables identified by the ID3 decision trees algorithm can also be given very close
235 attention and observed in order to better understand the students' performance and proper monitoring.

236

237 **Table 3:  Performance Evaluation Result Summary for the machine learning algorithms selected**

| Algorithm Used | Correct Classification | Accuracy (%) | TP Rate | FP Rate | Precision |
|---|---|---|---|---|---|
| ID3 Decision Tree Algorithm | 61 | 61.0 | 1.000 | 0.000 | 1.000 |
| C4.5Decisio Tree Algorithm | 100 | 100.0 | 0.333 | 0.333 | 0.203 |

238

239 **5. CONCLUSION**

240 The study proposed a predictive model for student performance using relevant classification factors
241 selected from a predefined set of factors of student performance.  The ID3 decision trees algorithms
242 identified few factors which were more related in determining the performance of students.   The
243 predictive model was formulated using the variables identified by ID3 decision trees for this study and
244 the performance evaluation of both models showed that the model developed using the ID3 decision
245 trees algorithm was a better model.  Unlike the C4.5 decision trees algorithm which could not clearly
246 state the relevant attributes, ID3 was able to identify the important variables and used them in
247 developing the decision trees for students' performance classification. The results of the study revealed
248 the variables that were identified by the ID3 decision trees algorithm as relevant for identifying the
249 classification of student's performance.  The ID3 algorithm was observed to show a better accuracy
250 compared to that of the C4.5 algorithm using the training dataset presented in the study.

251

252

253 **6. REFERENCES**

254 [1] B. Baradwaj and S. Pal, 2011. Mining Educational Data to Analyze Students' Performance.
255 (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 2, No. 6, 2011.
256 [2] K. Adhatrao, A. Gaykar A, A. Dhawan, R. Jha. and V. Honrao.  Predicting Students' Performance
257 Using ID3 and C4.5 Classification Algorithms. International Journal of Data Mining & Knowledge
258 Management Process (IJDKP) Vol.3, No.5, September 2013
259 [3] Z. Xiaoliang,., W. Jian, Y. Hongcan, and W. Shangzhuo, Research and Application of the improved
260 Algorithm C4.5 on Decision Tree. International Conference on Test and

Measurement (ICTM), Vol. 2, 2009.

[4] A. Ahmed,. and I. Elaraby, 2014. Data Mining: A prediction for Student's Performance Using Classification Method. World Journal of Computer Application and Technology, 2(2), pp.43-47.

[5] U. Pandeyand S. Pal. Data Mining: A Prediction of performer or underperformer using Classification. (IJCSIT) International Journal of Computer Science and Information Technology, Vol. 2(2), 2011

[6.] B. Bhardwaj and S. Pal. Data Mining: A prediction for performance improvement using classification. (IJCSIS) International Journal of Computer Science and Information Security, Vol. 9, No. 4, 2012.

[7] I. Livieris and P. Pintela. An improved spectral conjugate gradient neural network training algorithm. *International Journal on Artificial Intelligence Tools*, 21(1), 2012.

[8] S. Kotsiantis, C. Pierrakeas and P. Pintelas. Predicting students' performance in distance learning using machine learning techniques, *Journal of Applied Artificial Intelligence*, 18(5), p.p. 411-426, 2004.

[9] C. Moucary, M. Khair and W. Zakhem, 2011. Improving student's performance using data clustering and neural networks in foreign-language based higher education. The Research Bulletin of Jordan ACM, 2(3), pp 27-34

[10] Y. Hongsuk, H. Jung. and S. Bae, 2017. Deep Neural Networks for traffic flow prediction. In Big Data and Smart aComputing (BigComp), IEEE International Conference on (pp. 328-331). IEEE.

[11] S. Yadav, and S. Pal. Data mining: A prediction for performance improvement of engineering students using classification. World of Computer Science and Information Technology Journal (WCSIT). (ISSN: 2221-0741), Vol. 2, No. 2, 51-56, 2012