

Student Academic Performance Monitoring and Evaluation Using Data Mining Techniques

Emmanuel N. Ogor
Department of Natural Sciences
Turks & Caicos Islands Community College
Turks & Caicos Islands
ogor@tciiway.tc

Abstract

Assessment as a dynamic process produces data that reasonable conclusions are derived by stakeholders for decision making that expectedly impact on students' learning outcomes. The data mining methodology while extracting useful, valid patterns from higher education database environment contribute to proactively ensuring students maximize their academic output. This paper develops a methodology by the derivation of performance prediction indicators to deploying a simple student performance assessment and monitoring system within a teaching and learning environment by mainly focusing on performance monitoring of students' continuous assessment (tests) and examination scores in order to predict their final achievement status upon graduation. Based on various data mining techniques (DMT) and the application of machine learning processes, rules are derived that enable the classification of students in their predicted classes. The deployment of the prototyped solution, integrates measuring, 'recycling' and reporting procedures in the new system to optimize prediction accuracy.

Key Terms:-DM, DMT, KD, decision rules, assessment, performance monitoring, stakeholders

1. Introduction

Performance monitoring involves assessments which serve a vital role in providing information that is geared to help students, teachers, administrators, and policy makers take decisions.[1] The changing factors in contemporary education has led to the quest to effectively and efficiently monitor student performance in educational institutions, which is now moving away

from the traditional measurement and evaluation techniques to the use of DMT which employs various intrusive data penetration and investigation methods to isolate vital implicit or hidden information. Due to the fact that several new technologies have contributed and generated huge explicit knowledge, causing implicit knowledge to be unobserved and stacked away within huge amounts of data. The main attribute of data mining is that it subsumes Knowledge Discovery (KD) which according to [2] is a nontrivial process of identifying valid, novel, potentially useful and ultimately understandable patterns in data processes, thereby contributing to predicting trends of outcomes by profiling performance attributes that supports effective decisions making. This paper deploys theory and practice of data mining as it relates to student performance and monitoring Associate degree program in a Community College.

Technological developments and new programming techniques have improved understanding and use of Artificial Intelligence (AI). The isolation of hidden data and exposed relationships embedded within it, without a prior knowledge of the nature of any inherent relationship leading [3] to assert that data mining is a logical evolution of database technology with the development of enhanced query tools such as SQL, database managers are capable querying data more flexibly. Rules derived from various algorithms during the implementation of DMT in researches, support this opinion.

Recently educational institutions targets activities within its organizations with ERP tools to handle and store huge data available in educational processes for hidden patterns. The face value assessment of students at the point of entry can only be confirmed or dispelled by the dynamic follow-up monitoring of students'

performance during the course of study leading to serve as an indicator of the suitability and unsuitability of students before admission and during their course of study.

Fuzzy Set Theory is used in applications involving educational assessment and performance as it is regarded as efficient and effective in uncertain situations involving performance assessment. It is known that Expert Fuzzy scoring systems noted [4]; help teachers make assessment in less time and with a level of accuracy that compares favorably to the best teacher examiner.

This paper profiled students from factual and partly behavioral factors. The factual profile content such as gender, date of birth, race and others like different test results from the college each semester obtained from the student records. Performance profiling is dependent upon motivation, attitudes, peer influence, curriculum and by the continued real-time monitoring of student's performance using a simple rapid response system and as [5] noted predicts correctly which student may need some attention or reinforcements in the course of their education.. The model developed helps achieve a measurable student progress monitoring process that gives results quickly and meets a larger educational goal benefiting stakeholders in the educational system and the wider community.

2. Assessment and DM Rationale

The ideal goal of higher education is to continually maintain sustainable increasing graduation rates and growth with the most efficient procedures that allows for the accounting of input resources [6][7]. The degree of quality students' involves the pertinent issue of how to enhance and evaluate it through overt and covert processes. Hence, DM processes for knowledge is data which while dependent on quality, characteristics and preparation, supports and facilitates the thorough examination of the data's different aspects for knowledge discovery (KD) in tertiary processes. The result helps educational institutions to predict the degree of likelihood of a students' persistence, learning outcomes in terms of performance and by using clustering algorithms, meaningful learning outcome topologies are created [8]. Other studies have shown that some techniques are particularly beneficial for the various sub process. [9][10]

3. Visualization and Articulation Methodology.

A total of 2215 records were initially identified and after a rigorous consideration of the impact of unknown and unavailable data, 1369 or 61.82% of the

data was selected for the data mining (DM) process. Missing and incomplete data were also included in the extracted data as learning from incomplete data is very possible and useful [11]. The operational data includes student demographic information and course assessment data of 1360 students from 2004/05 to 2005/06 academic year of the Business Studies Department of TCICC using five (5) different courses from three out of four consecutive semesters. The related and prerequisite courses selected are named for the purposes of data analysis as S1C1,.., S1C5; S2C1,.., S2C5; and S3C1,..,S3C5; with the nomenclature S_nC_nT and S_nC_nE represents the semester, course work test scores and course work examination scores respectively. Derived attributes includes S1_Avg_Performance Change = S1_Total - Class_Avg. Similar attributes applied to semester 2 (S2) and semester 3 (S3) respectively. Others include Overall_Gain_Performance_Avg= (S1_Avg + S2_Avg + S3_Avg)/3 and Avg_GainPerCourse = (Overall Gain Per Course)/Number of Courses Taken. Generally, quality derived attribute(s) output as shown in Table 1 plays a huge role in quality rule derivation in DM process as the research observed.

| Performance Ratio (PR) | Interpretation |
|------------------------|----------------|
| 0% =<PR <= 88% | Weak |
| 88% < PR < 94% | Fair |
| PR > = 94% | Strong |

Table 1 Performance Ratio

The working performance data is the overall grade point distribution of student performance is shown in Figure 1. The method used in this selection

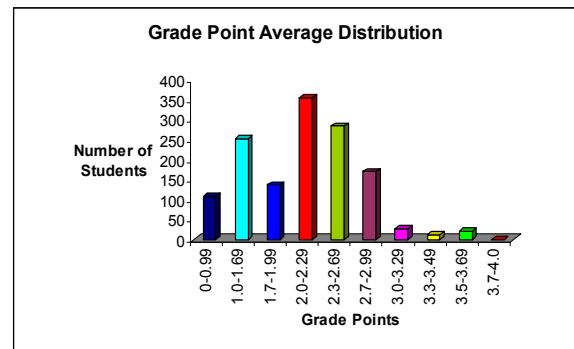


Figure 1 1st. 3-Semesters CGPA distribution

conforms to most methods of heuristic selection as data selected with reduced noise, approximates to the method of determining the most representative known as outstanding representatives [12]

4. Data Validation and Segmentation

The visualization and validation of the source data generates initial graphical objects and other clustering

models that may exist for this study. The identified data similarly remained consistent on an initial verification as the criteria laid down by the college authorities with respect to data handling, transactional quality, and integrity. Oversight mechanisms were strictly followed and adhered to ensuring the data gives the most accurate picture of students' performance especially test and examination scores. A final but simple validation with about 25% of the data is deployed on completion of the DM process as for comparative analysis, final model creation and error rate comparative analysis.

The segmentation by grouping performance level is a standard technique deployed by researchers is beneficial because in other studies DM segmented models usually out-perform aggregated models noted [13] as it preserves the overall interactive structures in a dataset while at the same time splits the analysis into segments of shorter duration. In this project attempts have been made to observe and compare individual, segmented and well aggregated student performance variables by analyzing the whole student base activities and then building one predictive model. This is achieved by building finer segments with latent positive values that effectively model student performance and allow for a homogeneous student performance prediction based solely on the given transaction or assessment and demographic data sets. [14]

5. Applicable Data Mining Techniques

The rule inductions and artificial neural network data mining techniques used in the project fall under the category of machine learning that uses high end modeling techniques for uncovering hidden patterns and/or predicting outcomes. Supervised knowledge discovery explains relationships found [15]. Hence for this project the data mining product used is the Clementine 10.0 software that parades in its arsenals C5.0, C&RT, ANN, CHAID, QUEST, Link Analysis, KMeans, and Kohonen etc. for both classification and estimation in supervised and unsupervised tasks respectively. The supervised modeling techniques for the classification and estimation of students' data and DM theories and techniques of analyses fully applied.

The C5.0 algorithm served as a good representative of the decision trees in this project as its statistical property of information gain helped to determine which of the several attributes best represent the division of the training sets. The processed neural network model reached an estimated accuracy of 99.57% and it contains 76 input layer neurons, 4 hidden layer neurons and 5 output layer neurons. The attributes with the highest relative importance is Credit Performance Per CGPA (CreditPerfCGPA) with value of 0.3116890

conforming to 0.35 thresholds.[16] Thus some of the resulting twenty five (25) rules from C5.0 used, contain the some of the following:-

Rule 1 for Good **(118, 1.0)**

if CGPA > 1 and CreditPerfCGPA <= 37.04 and OverallAvgTotal <= 81.78 and Avg_Gain_Per_Course > 0.40 and S1_Total > 65.1 and S2E_Avg > 33 and S2C4_T <= 37 then Good

Rule 6 for Satisfactory **(514, 1.0)**

if CGPA > 1 and CreditPerfCGPA > 37.04 and Avg_Gain_Per_Course > -0.44 and CreditPerfCGPA <= 53.57 then Satisfactory

Rule 3 for Fail **(65.341, 1.0)**

if CGPA <= 1 and CreditPerfCGPA > 81.25 and ExamAvg <= 44.53 then Fail

6. Data and Model Analyses

The prediction target processes summarized all the qualities of assessment and performance monitoring of students' which when expanded holds key information that answers questions using a matrix of 1369 student records x 78 derived variables .A preliminary statistical evaluation of the variables was carried out using some measures of central tendency and F-Test. During the first semester, students' performance was good but started to show signs of decreased performance in the second semester but during the second to third semester students' performance, though improving, also stabilized. Students' performance of note is that students in the class of "Excellent", "Good" and "Satisfactory" on the average maintained or remained consistent in their performance throughout the three semesters under study. The "Fail" class is consistently erratic and hence most likely group to keep close watch of along with the "Marginal" group, whose members require constant monitoring and advisement.

Clementine 10.0 Desktop Software provided an effective methodology to compare the various classifications for the training data and evaluating the test and validation datasets.

| | ANN | C&RT | C5.0 | CHAID |
|----------------|--------|--------|--------|--------|
| Training Set | 91.7% | 93.58% | 97.11% | 67.15% |
| Testing Set | 88.86% | 92.66% | 97.90% | 52.69% |
| Validation Set | 92.71% | 91.32% | 96.88% | 53.48% |
| Average | 91.42% | 92.52% | 97.30% | 57.77% |

Table 2 Matrix of Model Performance for Training, Testing and Validation Sets

To further understand groups, non-performance attributes, that latently affect performance and

therefore also worth monitoring in order to provide an overview of the entire dataset and identify main data pointer, a subjective selection of six (6) clusters was set for this study and two (2) distinct clusters-2 with 359 records and cluster-5 with 29 records out of the six (6) was identified. For instance cluster 2 members significantly foreign students performed better than their counterparts in cluster 5 predominantly local students; the mean credits passed are 44.441 and 8.207 respectively.

The attributes in the clusters were subjected to the Generalized Rule Induction (GRI) was also applied to demographics to produce association rules. This was applied on the attributes AgeClass, CreditAttemptClass, Campus, CreditPassClass and the CGPAClass. The result generated is a *two-way association* highlighting some patterns in the dataset shown as a web display of Figure 2.

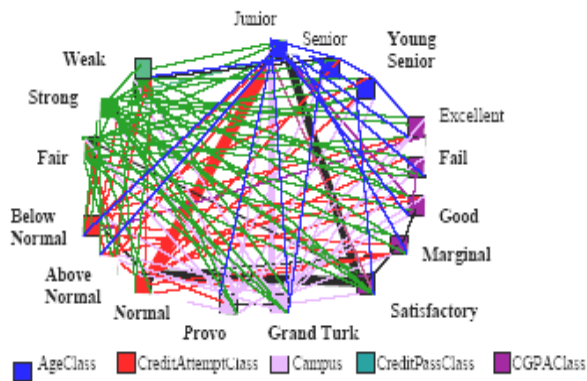


Figure 2 Web Display of Attributes

The achieved performance assessment monitoring modeling and classification led to a formalized prototype Performance Assessment Monitoring System (PAMS). This system trained with 60% of the data, tested with 20% data and to improve upon the accuracy, validated by 20% data. The eventual system prototype provide an “on-the-fly” and continuous “Just-In-Time” student performance assessment model for predicting performance with reasonable degree of accuracy, thereby enhancing monitoring of student academic pursuance and any other stakeholder’ interests, at any point, for any student during the student’s tenure at the educational institution.

7. The Process and Function Analysis

A random control group and finally groups that may or may not conform to obviously known rules are considered. These analyses were based on the subjective selection then *with Associativity analysis*, to

expose latent patterns and determining the degrees of co-relationship within a dataset. This system is developed and based on Microsoft VB features along with the algorithms format as presented by the C5.0 rule set with the SQL capabilities.

8. The Modeling Process and Access

This involves processes that rely on application systems and a number of manual transactions. The rules are exported to a database management system manually or using OLAP, allowing for rule migration. After each monitoring process concluded by a stakeholder, and an output performance prediction obtained, report output or logs of predictions over a particular period(s) in time are obtained to enable the tracking, monitoring and comparative analysis of the particular students’ performance over a given period the student is under observation. The modeling process was effective as it integrated all the data objects and rules needed for performance prediction allowing for quality control.

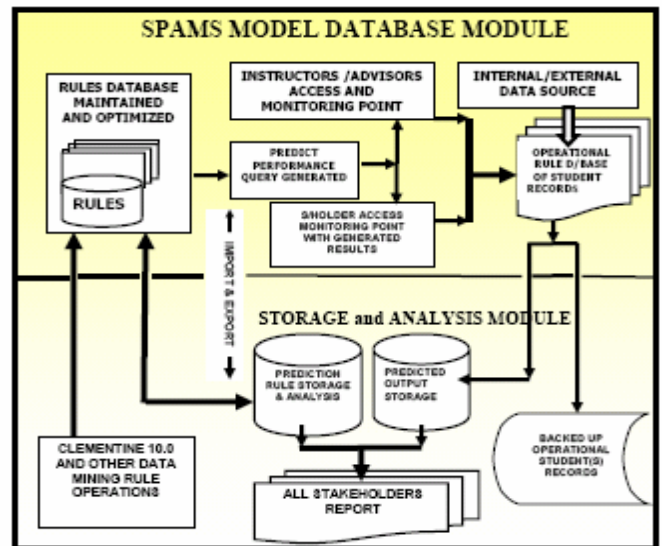


Figure 3 The Data Flow Within the SPAMS Processes

9. Implementation

The implementation of the prototype was carried out having among other integrated analytics features in Microsoft Access application for ease of use for end-users and data size. The system was deployed using ‘live’ student data never used previously deployed and the comparative output between the manual and system results. The startup application, displays a simple GUI data window in the form shown in Figure 4.

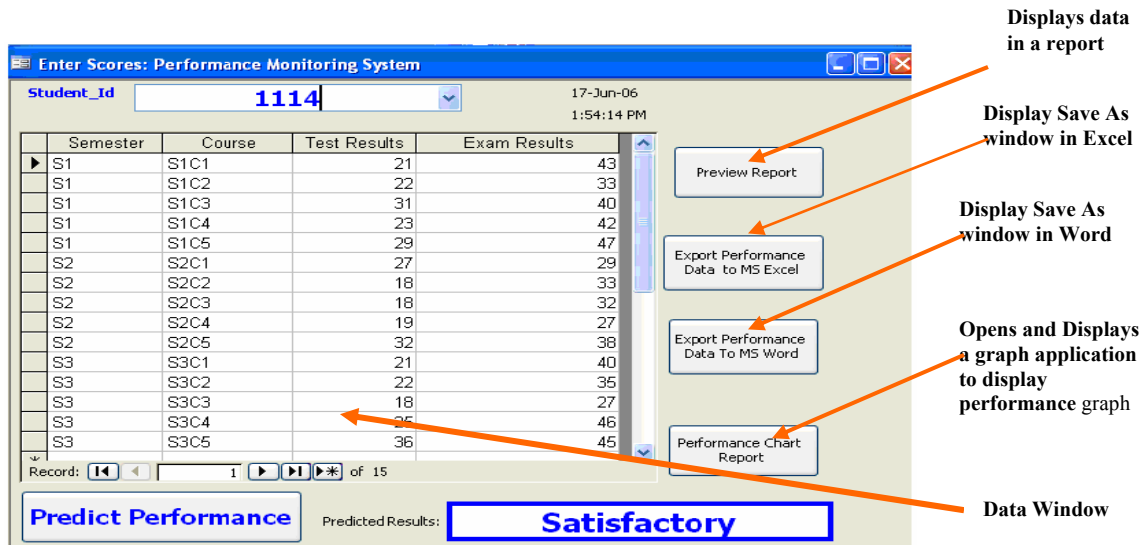


Figure 4 : SPAM System Windows Model

| CGPAClass | Frequency | Class % |
|-------------------|-----------|----------------|
| Fail | 2 | 4.4 |
| Marginal | 7 | 15.2 |
| Satisfactory | 27 | 58.7 |
| Good | 8 | 17.4 |
| Excellent | 2 | 4.4 |
| Misclassification | 0 | 0 |
| No Predictions | 0 | 0 |
| TOTAL | 46 | 100.00% |

Table 4 Summary | SPAMS Manual Prediction

| CGPAClass | Frequency | Class % |
|-------------------|-----------|----------------|
| Fail | 2 | 4.4 |
| Marginal | 4 | 8.7 |
| Satisfactory | 28 | 60.9 |
| Good | 8 | 17.40 |
| Excellent | 2 | 4.4 |
| Misclassification | 2 | 4.4 |
| No Prediction | 0 | 0 |
| TOTAL | 46 | 100.00% |

Table 5 Summary of non-SPAMS Prediction

Further analysis: Marginal class showed a significant difference with 7 students or 15.2% in the manual method compares to 4 students or 8.7% correct predictions using the SPAMS. The misclassification was due to some identified overlapping rules in the Marginal and Satisfactory classes. Hence the SPAMS attained an accuracy rating of 95.6%. The graphical output of the system model clearly shows performance trend. For as shown in figures below, there is significant difference in the trend between “Excellent” and “Fail” performing students and Tables 4 and 5

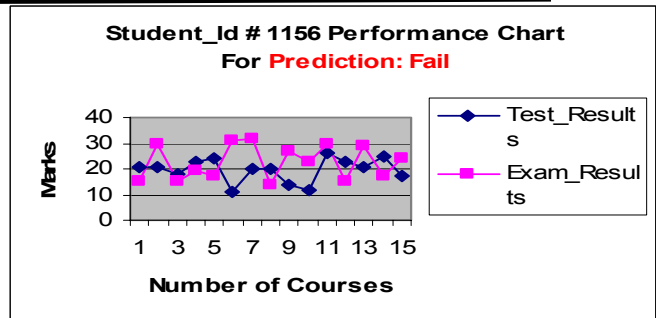


Figure 5: Sample Chart for a “Fail”

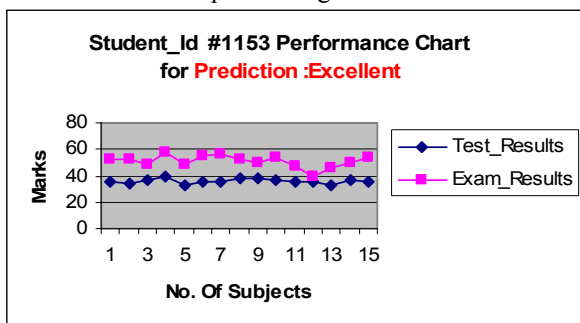


Figure 4: Sample Chart for Excellent

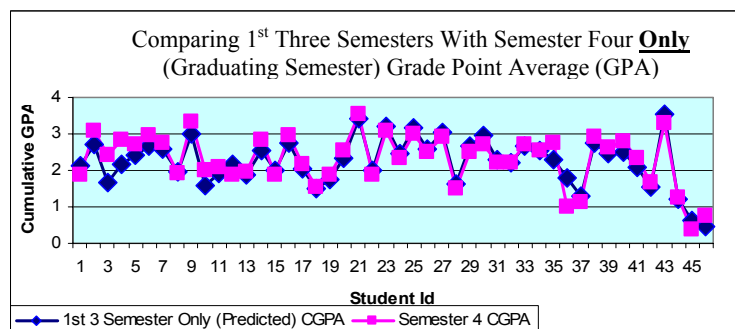


Figure 6: Predicted CGPA versus Actual Graduating GPA

10 Conclusion

The result of this project indicates that DMT capabilities provided effective monitoring tools for student academic performance with overall 94% success rating and fine tuning derived variables improves rules quality producing improved performance. The various reporting tools that this system offers serve mainly to compare changes over time in performances as may be affected by the different rules that are available plus other well chosen variables exposes systematic structures required to improve performance monitoring. OLAP implementation with dynamic reporting capabilities and efficiency is perceived as better solution and recommended for very large student databases in Oracle or MS SQL Server database environment

11 Recommendations for Future Work

The encouraging results obtained on application of knowledge discovery, begs for a comprehensive strategic implementation, an integration of the results of other research efforts in areas such as Instructor assessment and performance, curriculum, course relevance, student attitude, demographics, etc and its impact on the student learning process must be determined and integrated into any prototype. Learning process must be determined and integrated into any prototype performance, course relevance, student attitude, demographics, etc and its impact on the student learning process must be evaluated and integrated into any future performance monitoring prototype. DMT has a potential in performance monitoring of High school and other levels education offering historical perspectives of students' performances. The results may both complement and supplement tertiary education performance monitoring and assessment implementations.

12 References

- [1] Council N. "Knowing What Student Knows. The Science and Design of Educational Assessment". National Academic Press. Washington, D.C. 2001
- [2] Frawley, W.J., Piatetsky-Shapiro, G and Matheus, C. J.), "Knowledge Discovery databases: An overview In": Piatetsky-Shapiro and Frawley, W. J. (eds) Knowledge Discovery in Databases, AAAI/MIT .1991. pp 1-27
- [3] Rubenking N. (2001), "Hidden Messages" PC Magazine, May 22, 2001.
- [4] Nolan J "A Prototype Application of Fuzzy Logic and Expert Systems in Education Assessment". AAAI/IAAI Proceedings pp 1134-1139. 1998
- [5] Luan J. "Data Mining And Knowledge Management" Presentation at AIR Conference, Long Beach, CA. June 2001 pp1-17
- [6] Yorke M., "Leaving Early: Undergraduate Non-completion in Higher Education". Philadelphia: Palmer Press. 1999
- [7] Kash Barker, Theodore Trafalis, and Teri Reed Rhoads "Learning From Student Data". Proceedings of the 2004 Systems and Information Engineering Design Symposium. Mathew H. Jones, Stephen D. Patek, and Barbara E. Towney eds. 2004. pp79-86
- [8] Luan, J. "Chapter 2: Data Mining and Its Application in Higher Education. Knowledge Management – Building a Competitive Advantage in Higher Education." Serban, A. & Luan, J. (eds.) Jossey-Bass. 2002
- [9] Waiyamai K, "Improving the Quality of Graduate Students by Data Mining". Department of Computer Engineering, Faculty of Engineering, Kasetsart University, Bangkok, Thailand. 2003.
- [10] Naeimeh Delavari and Mohammad Reza Beikzadeh and Somnuk Phon-Amnuaisuk, "Application of Enhanced Analysis Model for Data Mining Processes in Higher Educational System". ITHET 6th Annual International Conference. Juan Dolio, Dominican Republic. July 7 – 9, 2005.pp F4B-1 -6
- [11] Lakshminarayan, K., Harp, S.A., Goldman, R. and Samad, T. "Imputation of Missing Data Using Machine Learning Techniques". Proceedings of the Second International Conference on Knowledge Discovery and Data Mining . Portland, OR, pp. 1996. 140-145.
- [12] Michalski, R.S., "A Planar Geometrical Model for Representing Multi-Dimensional Discrete Spaces and Multiple-Valued Logic Functions," ISG Report No. 897, Department of Computer Science, University of Illinois, Urbana,1978
- [13] Allenby, G.M. and Rossi P.E. "Marketing Models of Customer Heterogeneity". Journal of Econometrics. 1999.
- [14] Shavelson R and Towne, L "Scientific Research In Education". National Research Council, Washington, D.C.; National Academy Press. 2002
- [15] Westphal, C and Blaxton, T., Data Mining Solutions: Methods and Tools for Solving Real-World Problems. New York: Wiley Computer Publishing, 1998.
- [16] Luan J. "Data Mining and Knowledge Management in Higher Education – Potential Applications". Presentation at AIR Forum Toronto. Canada. 2002.