

Predictive Analysis of Dropouts in Information Technology Higher Education

U. G. N. Kumari
169317F

Dissertation submitted to the Faculty of Information Technology, University of
Moratuwa, Sri Lanka for the partial fulfillment of the requirements of Degree of Master
of Science in Information Technology

April 2019

Declaration

I declare that this thesis is my own work and has not been submitted in any form for another degree or diploma at any university or other institution of tertiary education. Information derived from the published or unpublished work of others has been acknowledged in the text and a list of references is given.

Name of Student

U.G.N.Kumari

Signature of Student

.....

Date:

Supervised by

Name of Supervisor:

Mr. S.Premaratne

Signature of Supervisor:

.....

Date:

Acknowledgment

I special gratitude pass to Mr. S. Premaratne, Senior Lecturer, Faculty of Information Technology, University of Moratuwa for his help, direction, and supervision was given to me all through my undertaking, making it a triumph.

Very special gratitude goes out to all down at Advanced Technological Institute, Labudawa for helping and providing the enormous support for data collection.

My genuine appreciation to Mrs. S.C.P. De Silva, Department of IT, Faculty of IT, ¹University Of Moratuwa for the animating talks, occupied hours we were cooperating before due dates.

I'm grateful to all the staff members of the Faculty of Information Technology, University of Moratuwa.

To wrap things up I might want to thank my family for supporting me profoundly throughout writing this thesis and my life in general.

Abstract

At this point attention on educational data mining methods have impact highly on predicting academic performance as the increased higher education dropout rates especially in information technology education has received huge attention in recent years due to the quality of higher education has been a topic of debate for many years. There is a huge necessity of mining educational data system and exact hidden knowledge to understand the factors affecting student dropouts and to understand the patterns that can lead to predict student performance at the entrance and improve and monitor performance of students enrolling in Information technology higher education by building early warning indicators based on factors affecting to dropouts and manage students drop out from the higher education. Data mining strategies have been utilized to effectively extract new, conceivably imperative Knowledge and diverse data mining techniques, for example, association, classification, clustering, prediction, sequential patterns, and decision trees are being utilized by numerous sorts of research. Identification of relevant attributes which affect to dropouts in ICT higher education is a leading concern in the field of education data mining as there are no significant studies that can be applied to understand the complex-inter correlated and distinct factors affecting to dropouts in ICT higher education Hence, the research has been conducted to identify the complex-inter correlated and distinct factors affecting to dropouts in ICT higher education. It is hypothesized that an experimental methodology can be adapted to generate a database that includes relevant information for extracting knowledge. The raw data will be preprocessed in terms of filling up missing values, transforming values in one form into another and relevant attribute/ variable selection. Thereby select student records, which can be used for classification prediction model construction. In constructing a classifier model different classifying algorithms can be applied and in this study evaluation of different classification algorithms is done to identify the most accurate algorithm. Finally, a predictive analyzing model will be building for student profile analyzing using the identified algorithm. Then this classification model will be used in developing an application to predict the students' dropouts. The overall research will be designed using the WEKA data mining tool and using java WEKA library for developing the application.

Table of Contents

Declaration.....	ii
Acknowledgments.....	iii
Abstract.....	iv
List of Figures.....	viii
List of Tables.....	ix
Chapter 1.....	1
1 Introduction.....	1
1.1 Introduction.....	1
1.2 Increased Higher Education Dropouts Rate is a Topic of Debate for Many Years.....	1
1.3 Problem Definition.....	3
1.4 Aim and Objectives.....	4
1.4.1 Aim.....	4
1.4.2 Objectives.....	4
1.5 Structure of the report.....	4
1.6 Chapter Summary.....	5
Chapter 2.....	6
2 Emerging use of Data Mining Techniques in Education Sector.....	6
2.1 Chapter Introduction.....	6
2.2. Different Data Mining Techniques and Increased Dropout Rates in Higher Education..._	6
2.2.1 Different Data Mining Techniques and Increased Dropout Rates in Higher Education	7
2.2.1 Different factors affecting to Dropouts.....	9
2.2.3 Current Status of Educational Data Mining.....	12
2.3 Chapter Summary.....	13
Chapter 3.....	14
3 Predictive Analysis Through Classification Using WEKA Data Mining Tool.....	14
3.1 Chapter Introduction.....	14
3.2 Correlation-Based Feature Selection.....	14
3.3 Brief Introduction to Classification Techniques Used in WEKA.....	14
3.3.1 Decision Tree Algorithms.....	15
3.3.2 Naïve Bayes Algorithms.....	16

3.3.3 KNearestNeighbors Algorithms.....	17
3.3. Rule-Based Algorithms.....	18
3.4 Model Evaluation.....	19
3.5 Java WEKA API.....	20
3.6 Chapter Summary	20
Chapter 4.....	21
4 J48 Decision Tree Algorithm to Predict Student Dropouts.....	22
4.1 Chapter Introduction.....	22
4.2 Research Questions.....	22
4.3 Input	22
4.3.1 Selection of specific dataset	22
4.3.2 Dropout data.....	22
4.4 Output.....	23
4.5 Process.....	23
4.5 Users.....	23
4.6 Chapter Summary	24
Chapter 5.....	25
5 Design.....	25
5.1 Chapter Introduction.....	25
5.2 Analysis and Design.....	25
5.2.1 Collection of Student Data	25
5.2.2 Preprocessing	25
5.2.3 Data Selection and Transformation	27
5.2.4 Feature Extraction: Identifying important factors affecting to dropouts	27
5.3 Model Building Using Various Data Mining Techniques.....	28
5.4 Deployment of GUI Application to predict Student Dropouts	28
5.5 Chapter Summary.....	29
Chapter 6.....	30
6 Implementation of the Dropout Prediction System	30
6.1 Chapter Introduction	30
6.2 Implementation	30
6.2.1 Building up classifier models and evaluation for accuracy	30

6.2.2 Development of dropout prediction system.....	31
6.2.2.1 Development of GUI of dropout prediction system	31
6.2.2.2 Buildup classifier model.....	31
6.2.2.3 Classify new data	31
6.3 Chapter Summary.....	31
Chapter 7.....	32
7 Introduction.....	32
7.1 Chapter Introduction.....	32
7.2 Results.....	32
7.2.1 Performance evaluation of classifiers Objectives	32
7.2.1.1 Decision Tree algorithms (J48) Classifier.....	32
7.2.1.2 NaiveBayers Classifier	33
7.2.1.3 KNearestNeighbor(Lazy-IBK) Classifier.....	33
7.2.1.4 RuleBased(ZeroR).....	33
7.2.2 Accuracy comparison by increasing the number of instances in test dataset...	34
7.2.2.1 Decision Tree algorithms (J48) Classifier.....	34
7.2.2.2 NaiveBayers Classifier	34
7.2.2.3 KNearestNeighbor(Lazy-IBK) Classifier.....	35
7.2.2.4 rule-based(ZeroR).....	35
7.3 Implemented the Predictive analysis model.....	38
7.4 Chapter Summary.....	42
Chapter 8.....	43
8 Conclusion and Further work	43
8.1 Chapter Introduction	43
8.2 Conclusion	43
8.3 Future Work	44
8.4 Limitation of the proposed solution.....	44
8.5 Chapter Summary.....	45
References.....	46
Appendix	x
Appendix A – Selected Source Code	xi

List of Figures

Figure 5.1 - High-Level design of the proposed system.....	29
Figure 7.1 –Accuracy comparison by increasing the number of instances in test dataset of Decision Tree (J48) Classifier.....	34
Figure 7.2 –Accuracy comparison by increasing the number of instances in test dataset of Classifier.....	34
Figure 7.3 –Accuracy comparison by increasing the number of instances in test dataset of KNearestNeighbor (Lazy-IBK).....	35
Figure 7.4 –Accuracy comparison by increasing the number of instances in test dataset of rule-based (ZeroR) Classifier.....	35
Figure 7.5 - Decision Tree Diagram generated by j48 algorithm	37
Figure 7.6 - Variable importance graph of the Decision Tree J48 classifier	38
Figure 7.7a – Browsing a file from the local file system	39
Figure 7.7b - User selects a file from the local file system	39
Figure 7.8 - Building the model using the Decision Tree J48 classifier	40
Figure 7.9: Classify new data item.....	41

List of Tables

Table 2.1 – Factors that influence the dropouts in higher education system	11
Table: 3.1 Confusion matrix for performance evaluation.....	19
Table 5.1 - Student information data formats	26
Table 7.1 – Re-evaluate model summary of Decision Tree algorithms (J48) Classifier.....	32
Table 7.2 – Re-evaluate model summary of naive Bayes	33
Table 7.3 – Re-evaluate model summary of KNearestNeighbor - Lazy.IBK Classifier.....	33
Table 7.4 – Re-evaluate model summary of Rule-based - ZeroR Classifier.....	33
Table 7.5 - Performance Evaluation Matrix of all Classifiers.....	36