

Librarian
Perpustakaan Sultanah Zanariah
UTM, Skudai
Johor

Sir,

CLASSIFICATION OF THESIS AS RESTRICTED
- *KNOWLEDGE DISCOVERY FOR LARGE DATABASES IN EDUCATION
INSTITUTES*
ROBAB SAADATDOOST

Please be informed that the above mentioned thesis entitled "*KNOWLEDGE DISCOVERY FOR LARGE DATABASES IN EDUCATION INSTITUTES*" be classified as RESTRICTED for a period of three (3) years from the date of this letter. The reasons for this classification are:

- (i) Data collected is confidential.
- (ii) Data given was on trust basis and cannot be revealed to public.

Thank you.
Sincerely yours,



DR.ALEX TZE HIANG SIM
N28-403-02,
UNIVERSITI TEKNOLOGI MALAYSIA,
81310 UTM SKUDAI JOHOR
+607-5532406

Note: This letter should be written by the supervisor, addressed to PSZ and a copy attached to the thesis

KNOWLEDGE DISCOVERY FOR LARGE DATABASES IN EDUCATION
INSTITUTES

ROBAB SAADATDOOST

A dissertation submitted in partial fulfillment of the
requirements for the award of the degree of
Master of Information Technology-Management

Faculty of Computer Science and Information Systems
Universiti Teknologi Malaysia

JUNE 2011

This dissertation is dedicated to my Parents who have never failed to give me every support, and my supervisor for giving all my need during the time we researched and for teaching me that even the largest task can be accomplished if it is done one step at a time.

ACKNOWLEDGEMENT

Praises to God for giving me the patience, strength and determination to go through and complete my study. I would like to express my deep and sincere gratitude to my supervisor; Dr. Alex Tze Hiang Sim, Her wide knowledge and her logical way of thinking have been of great value for me. Her understanding, encouraging and personal guidance have provided a good basis for the present thesis. I would like to express my appreciation to Assoc. Prof. Dr. Ali Selamat, Dr. Mohd Shahizan Othman, Dr. Roliana Ibrahim and Dr. Ab. Razak Che Hussin, for their support and guidance during the course of this study and Madam Lijah Rosdi for their friendly help.

During this work I have collaborated with many colleagues for whom I have great regard, and I wish to extend my warmest thanks to all those who have helped me with my work in Institute for Research and Planning in Higher Education (IRPHE).

ABSTRACT

This project presents the patterns and relations between attributes of Iran Higher Education (Iran Higher Education) data gained from the use of data mining techniques to discover knowledge and use them in decision making system of IHE. Large dataset of IHE is difficult to analysis and display, since they are significant for decision making in IHE. This study utilized the famous data mining software, Weka and SOM to mine and visualize IHE data. In order to discover worthwhile patterns we used clustering techniques and visualized the results. The selected dataset includes data of five medical university of Tehran as a small data set and Ministry of Science - Research and Technology's universities as a larger data set. Knowledge discovery and visualization are necessary for analyzing of these datasets. Our analysis reveals some knowledge in higher education aspect related to program of study, degree in each program, learning style, study mode and other IHE attributes. This study helps to IHE to discover knowledge in a visualize way; our results can be focused more by experts in higher education field to assess and evaluate more.

ABSTRAK

Disertasi ini merupakan penyelidikan terhadap pola dan hubungan antara atribut Institusi Pengajian Tinggi Iran dengan menggunakan teknik “data mining” dalam pengumpulan data untuk eksplorasi dan pencarian pengetahuan serta menggunakannya dalam sistem membuat keputusan IHE. Dataset IHE yang besar dan sukar untuk dianalisis dan paparan, kerana ia adalah penting untuk membuat keputusan dalam IHE. Penyelidikan ini menggunakan perisian “data mining” (pelombongan data) yang terkenal iaitu Weka dan SOM untuk melombong (mine) dan visualisasi data IHE. Dalam rangka mencari pola-pola yang berharga, penyelidik menggunakan teknik pengugusan, pengkelasan dan divisualisasikan hasilnya. Dataset yang dipilih merangkumi data dari lima universiti perubatan Teheran sebagai kumpulan data kecil dan Kementerian Sains - universiti Penyelidikan dan Teknologi sebagai data yang lebih besar ditetapkan. Penemuan pengetahuan dan visualisasi adalah perlu dalam proses untuk menganalisis dataset ini. Hasil analisis menunjukkan pengetahuan dalam aspek pendidikan tinggi adalah berkaitan dengan program pengajian, ijazah dalam setiap program, kaedah pembelajaran, mod pengajian dan atribut IHE yang lain. Penyelidikan ini membantu untuk IHE dalam mencari pengetahuan dengan cara pemvisualisasian dan hasil kajian boleh difokuskan lebih oleh pakar-pakar di bidang pengajian tinggi untuk menilai dan menilai lebih.

TABLE OF CONTENTS

CHAPTER	TITLE	PAGE
	DECLARATION	ii
	DEDICATION	iii
	ACKNOWLEDGEMENTS	iv
	ABSTRACT	v
	ABSTRAK	vi
	TABLE OF CONTENTS	vii
	LIST OF TABLES	xi
	LIST OF FIGURES	xiii
	LIST OF APPENDIXES	xvii
1	INTRODUCTION	1
	1.1 Introduction	1
	1.2 Problem Background	2
	1.3 Problem Statement	4
	1.4 Project Objectives	4
	1.5 Project Scope	5
	1.6 The Project Importance	5
	1.7 Summary	6
2	LITERATURE REVIEW	7
	2.1 Introduction	7
	2.2 Definition	8
	2.3 Knowledge Discovery Process	10

2.4	Advantages of Knowledge Discovery	11
2.5	Data Mining Techniques	12
2.6	Machine Learning Software	15
2.7	Data and Data visualization	20
2.8	Use of Data Mining Tools in Higher Education Institute	21
2.9	Higher Education Institute and Higher Education Attributes	26
2.10	Summary	30
3	RESEARCH METHODOLOGY	31
3.1	Introduction	31
3.2	Methodology	32
3.3	Make a user interface	36
3.4	Data Collection	37
3.5	Attributes	37
3.6	Pre processing	38
3.7	Filtering and Data visualizing	42
3.8	Apply Techniques on Data to Discover Knowledge	44
3.9	Analysis Findings	47
3.10	Using another data mining tool to compare differences and similarities between two tools	48
3.11	Definitions of attributes	48
3.12	Generalization	51
3.13	Summary	51
4	DATA ANALYSIS	53
4.1	Introduction	53
4.2	Visualize each attribute	53
4.3	SimpleKMeans Clustering Technique	63
4.4	Summary of findings by using Weka:	82
4.5	SOM Toolbox	88
4.6	The comparison between Weka and SOM	106
4.7	Summary	114
5	LARGE DATASET ANALYSIS	115
5.1	Introduction	115

5.2 Attribute visualizing	115
5.3 SimpleKMeans Clustering Technique	125
5.4 Summary of findings:	145
5.5 Summary	150
6 SUMMARY	151
6.1 Introduction	151
6.2 Findings and Achievements	152
6.3 Constraints and Challenges	152
6.4 Summary comparing results from Weka to SOM	153
6.5 Summary and Recommendation	153
REFERENCES	160
APPENDIX A	164
APPENDIX B	191
APPENDIX C	203

LIST OF TABLES

TABLE NO.	TITLE	PAGE
2.1	Weka Timeline	17
2.2	Attributes Description	26
2.3	Requirements	29
3.1	Attributes List	37
4.1	Test 1 With K=2	63
4.2	Test 2 With K=5	68
4.3	Table Of Most Common Degree	73
4.5	Table Of Most Common Study Mode	76
4.6	Test 3 With K=13	77
4.7	The Important Years In Each Cluster	79
4.8	Summary Of Findings	82
4.9	Table Of Most Common Degree	87
4.10	Table Of Most Common Study Mode	87
4.11	Table Of Important Clusters	88
4.12	University Code	89
4.13	Degree Code	89
4.14	Type Code	90
4.15	Study Mode Code	90
4.16	University Component Plane	94
4.17	Degree Component Plane	96
4.18	Type Component Plane	97

4.19	Study Mode Component Plane	98
4.20	Sum Component Plane	99
4.21	Program Component Plane	100
4.22	Year Component Plane	101
4.23	The Comparison Between Weka And Som Results	107
4.24	Clusters In The Som Analysis	111
4.25	Weka And Som Features	113
5.1	Test 1 With K=3	126
5.2	Test 2 With K=5	134
5.3	Test 3 With K=13	137
5.4	Test 4 With K=7	141
5.5	Summary Of Findings (Large Dataset)	145

LIST OF FIGURES

FIGURE NO.	TITLE	PAGE
2.1	Transition From Data To Wisdom	9
2.2	Kdd Process Steps	11
2.3	Download Histories For Weka	19
3.1	Methodology (The Overview)	34
3.2	Methodology (The Details)	35
3.3	The Flow Of Pre-Processing	39
3.4	Weka Tags	41
3.5	Out Of Memory Error	42
3.6	Visualize Attributes	43
3.7	The Sample Plot	44
3.8	Cluster Visualize	45
3.9	Cluster As Arff File	46
3.10	Accdb File	47
4.1	Visualize The Selected Attribute	54
4.2	Degree By University Class	55
4.3	Visualize All Attributes	56
4.4	Visualizing Section Of Weka	56
4.5	University – Program	57
4.6	Degree – Program	58
4.7	University - Degree	60
4.8	Program – Year	62
4.9	Degree – Sum Of Students	65
4.10	University – Sum Of Students	67

4.11	University – Program	69
4.12	Cluster – Program- Degree	70
4.13	University – Study Mode	71
4.14	Cluster – University – Degree	72
4.15	Cluster – Sum Of Students - Degree	73
4.16	Cluster – Sum Of Students – Type	74
4.18	University - Cluster	78
4.19	Cluster – Year	79
4.20	Cluster – Sum Of Students	80
4.21	Degree – Sum - Cluster	81
4.22	Matlab _ Som	91
4.23	Initialization And Training	92
4.24	Visualizing Of U-Matrix And Component Planes	93
4.25	U-Matrix	94
4.26	University Component Plane	95
4.27	Degree Component Plane	96
4.28	Type Component Plane	97
4.29	Study Mode Component Plane	98
4.30	Sum Component Plane	99
4.31	Program Component Plane	100
4.32	Year Component Plane	101
4.33	Relation Between University And Year	102
4.34	Relation Between Degree And Year	103
4.35	Relation Between Type And Year	104
4.36	Relation Between Study Mode And Year	105
4.37	Relation Between University And Degree	106
5.1	Study Mode – Sum Number Of Students	116
5.2	Study Mode - City	117
5.3	Year – Sum Number Of Students	118
5.4	Degree – Sum Number Of Students	119
5.5	City – Program	121
5.6	Learning Style – Study Mode	122
5.7	Year – Learning Style	123
5.8	City - Year	124

5.9	City – Program – Colour (Degree)	125
5.10	Cluster (K=3) – Colour (Year)	127
5.11	Instance Number – Year – Colour (Cluster)	128
5.12	City – University – Colour (Cluster)	129
5.13	Study Mode – University – Colour (Cluster)	130
5.14	Year – Province – Colour (Cluster)	131
5.15	Year – Type – Colour (Cluster)	132
5.16	University – Learning Style – Colour (Cluster)	133
5.17	Cluster – Cluster – Colour (Degree)	135
5.18	University – City – Colour (Cluster)	136
5.19	University – Cluster – Colour (Cluster)	138
5.20	Degree – Program – Colour (Study Field)	139
5.21	University – Cluster – Colour (Year)	140
5.22	Female Students – Cluster – Colour (Year)	141
5.23	Year – Cluster – Colour (Study Field)	142
5.24	Cluster – Program – Colour (Cluster)	143
5.25	Instance Number – Cluster – Colour (Cluster)	144

LIST OF APPENDIX

APPENDIX	TITLE	PAGE
A	The figures of discoveries	165
B	List of Clusters and Code Label of Program	193
C	Dataset and validation letter (IRPHE)	204

CHAPTER 1

INTRODUCTION

1.1 Introduction

Nowadays each organization deals with some data about their area and during time it increases, one of these organizations that includes big volume of data is higher education institute that has always held much data about universities, students and teachers. Thus, it is possible for us to discover some worthwhile relations or patterns that can be useful for making decision. For examples, planning the future development of a university and identifying the cluster of students who required more attentions. Management faces many challenges particularly in planning and for this purpose it needs some facts extracted from data. In our rapidly changing world, every year we accumulate data and add it to our data sets so after several years we will have a massive databank, in this environment every year our data volume increases so we need some tools to analysis this data for extracting some valuable outcome from it. Data mining has many techniques that can apply and facilitate analyzing of data.

Data mining has many definitions and almost all of them point to the discovery of patterns, and the analysis of some relations between variables in data. It does not limit to collecting and managing data; it also includes analysis. In this study, we intend to use historical data as the basis of discovering hidden relations. We intend to perform data mining techniques to discover knowledge. There are some examples, such as:

- Mining of statistical data of one university to discover successful students(Venus Shokorniaz & akbari, 2008)(Venus Shokorniaz & akbari, 2008) .
- Mining on students and discovering groups of students those are available from the data and their relations (Yghini, Akbari, & Sharifi, 2008)(Yghini, Akbari, & Sharifi, 2008) .

In this project, we applied data mining techniques on data related to Iran Higher Education Institute to discover some relations and patterns that are useful in decision making system of higher education.

1.2 Problem Background

We have chosen this topic because of government and management of universities need to plan before an event occurrence. We face huge data and need to analysis them to reach some knowledge. For this purpose we need some techniques that data mining helps us on this way, data mining has two common techniques that are classification and clustering. In this project we study about these techniques and choose one of them in our project.

Clustering is a data mining technique that is a division of data elements into groups of similar objects without advance knowledge of the group definitions. In addition it is a tool for data analysis, which solves classification problems. In clustering, there are strong associations between members of each group and according to the type of clustering; Clustering algorithms have 4 types: exclusive, overlapping, Hierarchical, and Probabilistic. We may find some associations between different groups. Some of these associations are strong and some of them are weak. For example exclusive algorithm has weak association and overlapping has strong (Berkhin, 2006). Clustering is a discovery tool that may discover associations and patterns in data which is not previously obvious. In short: clustering attempts to find some groups of elements, based on some similarities (Ong, 2000). One of the cluster analyses is SOM (self organizing method) that is one of the most important algorithms in data visualization and exploration. Visualization transforms from the invisible to the visible (Alhoniemi, et al., 2002,2003). SOM is a particular type of neural network used in clustering. It maps high dimensional input onto two dimensional.

Classification is a data mining technique that predicts data elements' group, for example we can predict the weather of a day will be sunny, rainy or cloudy. In classification we have predefined classes that classification is a task to assign instances to these classes opposite of clustering that we don't have knowledge about group definitions. In clustering we cluster elements based on their attribute on the contrary in classifying we classify elements into groups by recognizing pattern.

Our concern in this study is finding a way to discover exciting knowledge for universities management to achieve an appropriate plan to improve society. Each society can be affected by higher education; its economic, political and scientific improvement can be resulted by advanced higher education, when we develop our programs and management part in universities so its output will be successful graduated students that can be helpful in every part of one society.

1.3 **Problem Statement**

The main problem is the volume of data that we have in higher education and government needs to analysis and discover fast and correct knowledge from them. “How we can discover it?” is our significant question. Our endeavour was proposing a methodology to discover knowledge that reveals some patterns and relations between data. We have huge databases about universities but most of time we cannot use them efficiently because we don’t have appropriate pattern for our executive system whereas we have data to this purpose.

We encounter many challenges in higher education such as allocating budget to universities in start of new semester without delay, accepting accurate number of students for every semester, finding ways for effective teaching. In this study our intension is discover some findings that help us in this way to improve higher education decisions and policies.

1.4 **Project Objectives**

Objectives of this study are:

- To study for an appropriate data mining approach suitable for analysing Higher Education Institute of Iran.
- To analyse data on medical universities to discover patterns usable for managerial decision making.
- To generalise steps for discovery on larger number of universities.
- To suggest an appropriate software for the analyses of this research.

1.5 **Project Scope**

The study is analyse huge volume of data in institute of research and planning for higher education (IRPHE) in Iran. We choose two dataset; the smaller dataset is medical universities of Tehran and the larger one is Ministry of Science - Research and Technology's universities. We proposed a methodology that can be used in similar institute in Iran with large volume of data.

1.6 **The Project Importance**

Knowledge discovery is necessary for most of our plantings. Nowadays planning for higher education has significant impact on developing of one society, successful planning needs to analysis some huge and historical data that is available in higher education institutes; most of the time there is not any correct and precise analysis whereas this analysis can be helpful for managers, researchers to plan, report and discover some knowledge.

The other importance of this project is collection of data with large volume that relate to many universities during many years so it has high probability to discover some exciting knowledge. This study can be used in every university, because it faces much data about students, teachers, staff and financial resources and most of the time, this data includes information and worthwhile patterns.

The methodology that will be suggested is helpful for higher education institute in collecting data and to discover some knowledge for improving

management decisions. Furthermore it can be useful for researchers that study and research about higher education. This knowledge leads us to development in technical, scientific and economic aspects according to type of data that we will analyse.

1.7 Summary

This chapter discussed the overview of this study which is knowledge discovery in higher education institute with statistical data and other associated issue in data mining techniques that helps us in discovering relations and patterns that are useful in decision making in higher education. Higher education usually does not have appropriate methodology to analysis huge volume of data and thus most of time it ignores much knowledge that can be very helpful in many decisions to improve quality of education. There are four objectives that were achieved successfully in order to improve higher education decision making systems and their policies that certainly affect on society. At first we studied knowledge discovery and then analysed Iran higher education data to find some relations and patterns. Finally we proposed a methodology to knowledge discovery from higher education data.

REFERENCES

- Al-Noukari, M., & Al-Hussan, W. (2008). *Using Data Mining Techniques for Predicting Future Car market Demand; DCX Case Study*. Paper presented at the Information and Communication Technologies: From Theory to Applications, 2008. ICTTA 2008. 3rd International Conference on.
- Aldana, W. A. (2000). A Brief History of Data Mining.
- Alhoniemi, E., Himberg, J., Hollm'en, J., Laine, S., PasiLehtim'aki, Raivio, K., et al. (2002,2003). *SOM in data mining*.
- Becerra-Fernandez, I., Zanakis, S. H., & Walczak, S. (2002). Knowledge discovery techniques for predicting country investment risk. *Computers & Industrial Engineering*, 43(4), 787-800.
- Bellinger, G., Castro, D., & Mills, A. (2010). Data, Information, Knowledge, and Wisdom.
- Berkhin, P. (2006). A Survey of Clustering Data Mining Techniques.
- Cios, k. j., pedrycz, w., swiniarski, r., & kurgan , l. A. (2007). Data Mining- A knowledge discovery approach.
- Delavari, N., & Beikzadeh, M. R. (2004). A new model for using data mining in higher educational system.
- Delavari, N., Beikzadeh, M. R., & Amnuaisuk, S. (2005). Application of enhanced analysis model for data mining processes in higher educational system.
- Famili, F. (2009, 27-28 Oct. 2009). *Knowledge discovery and management in life sciences: Impacts and challenges*. Paper presented at the Data Mining and Optimization, 2009. DMO '09. 2nd Conference on.
- Fangjun, W. (2010, 6-7 March 2010). *Apply Data Mining to Students' Choosing Teachers Under Complete Credit Hour*. Paper presented at the Education

- Technology and Computer Science (ETCS), 2010 Second International Workshop on.
- Fayyad, U. M., & Uthurusamy, R. (1996). Data mining and knowledge discovery in databases.
- Ghanbari, M. (2007, 4-6 March 2007). *Visualization Overview*. Paper presented at the System Theory, 2007. SSST '07. Thirty-Ninth Southeastern Symposium on.
- Guruler, H., Istanbulu, A., & Karahasan, M. (2010). A new student performance analysing system using knowledge discovery in higher educational databases. *Computers & Education*, 55(1), 247-254.
- Hilderman, R. J., Liangchun, L., & Hamilton, H. J. (1997, 3-8 Nov 1997). *Data visualization in the DB-Discover system*. Paper presented at the Tools with Artificial Intelligence, 1997. Proceedings., Ninth IEEE International Conference on.
- Himberg, J. (1999). SOM based cluster visualization and its application for false coloring.
- Holmes, G., Donkin, A., & Witten, I. H. (1994). Weka: A machine learning workbench.
- Jakulin, A. (2010).
- Koua, E. L. (2005). USING SELF-ORGANIZING MAPS FOR INFORMATION VISUALIZATION AND KNOWLEDGE DISCOVERY IN COMPLEX GEOSPATIAL DATASETS. *Seminar on Data and Information Management*
SS 2005 2D, 3D and High-dimensional Data and Information Visualization.
- Lam, S. B. (2007). DATA MINING WITH CLUSTERING AND CLASSIFICATION.
- LAROSE, D. T. (2005). DISCOVERING KNOWLEDGE IN DATA.
- Mannila, H. (1996, 18-20 Jun 1996). *Data mining: machine learning, statistics, and databases*. Paper presented at the Scientific and Statistical Database Systems, 1996. Proceedings., Eighth International Conference on.
- Mierle, K., K., L., Roweis, S., & Wilson, G. (2005). Mining student CVS repositories for performance indicators.
- Mitchell, T. M. (1997). Machine learning.

- Ong, C. S. (2000). KNOWLEDGE DISCOVERY IN DATABASES: AN INFORMATION RETRIEVAL PERSPECTIVE.
- Pazzani, M. J. (2000). Knowledge discovery from data? *Intelligent Systems and their Applications, IEEE, 15(2)*, 10-12.
- Salazar, A., Gosalbez, J., Bosch, I., Miralles, R., & Vergara, L. (2004, 28 June-1 July 2004). *A case study of knowledge discovery on academic achievement, student desertion and student retention*. Paper presented at the Information Technology: Research and Education, 2004. ITRE 2004. 2nd International Conference on.
- Saraee, M. H., & Theodoulidis, B. (1995, 1 Feb 1995). *Knowledge discovery in temporal databases*. Paper presented at the Knowledge Discovery in Databases, IEE Colloquium on (Digest No. 1995/021 (A)).
- Siraj, F., & Abdoulha, M. A. (2009). *Uncovering Hidden Information Within University's Student Enrollment Data Using Data Mining*. Paper presented at the Modelling & Simulation, 2009. AMS '09. Third Asia International Conference on.
- Tan, Steinbach, & Kumar. (2006). Introduction to Data Mining.
- Tsai, P. S. M., & Chen, C.-M. (2001). Discovering knowledge from large databases using prestored information. *Information Systems, 26(1)*, 1-14.
- Vembu, S., & Baumann, S. (2004). A Self-Organizing Map Based Knowledge Discovery for Music Recommendation Systems. *Computer Music Modeling and Retrieval Second International Symposium, CMMR 2004, Esbjerg, Denmark, May 26-29, 2004, 3310*, 119-129.
- Venus Shokorniaz, & akbari, A. H. a. (2008). Mining of statistical data of one university to discover successful students.
- Vesanto, J., & Alhoniemi, E. (2000). Clustering of the Self-Organizing Map. *IEEE TRANSACTIONS ON NEURAL NETWORKS, 11*.
- Vesanto, J., Himberg, J., Alhoniemi, E., & Parhankangas, J. (1999). Self-organizing map in Matlab: the SOM Toolbox. *Proceedings of the Matlab DSP Conference 1999, Espoo, Finland, November 16-17, pp. 35-40, 1999*.
- Witten, I. H., & Frank, E. (2005a). Data Mining :Practical Machine Learning Tools and Techniques,.
- Witten, I. H., & Frank, E. (2005b). Data Mining: Practical machine learning tools and techniques, 2nd Edition.

- Yghini, M., Akbari, A., & Sharifi, M. (2008). Mining on students and discovering groups of students those are available from the data and their relations.
- Zanakis, S. H., & Becerra-Fernandez, I. (2005). Competitiveness of nations: A knowledge discovery examination. *European Journal of Operational Research*, 166(1), 185-211.