# An Investigation on Factors Affecting the Further Education Desires of Youth in Sri Lanka
# &
# A Decision Support Tool for Data Mining

By
D.R.T. Jayasundara
06/8105

This dissertation is submitted in partial fulfillment of the requirements for the M.Sc. Degree in Operational Research.

Department of Mathematics
University of Moratuwa
Sri Lanka

October 2009

# Declaration

I officially state that this dissertation titled "An Investigation on Factors Affecting the Further Education Desires of Youth in Sri Lanka and a Decision Support Tool for Data Mining" is entirely my own work except where explicitly mentioned. It has not been submitted nor being currently submitted for any other degree.

Candidate : Ms. D. R. T. Jayasundara     Signature:.. *UOM Verified Signature*

Date : 01/03/2010

Supervisor : Dr. M.D.T. Attygalle     Signature:.. *UOM Verified Signature*
Senior Lecturer
Department of Statistics,     Date: 01/03/2010
Faculty of Science,
University of Colombo.

Supervisor : Mr. M. Firdhous     Signature:.. *UOM Verified Signature*
Senior Lecturer
Department of Information
Technology,     Date : 01/03/2010.
Faculty of Information
Technology,
University of Moratuwa.

Coordinator Mr. T.M.J.A. Cooray     *UOM Verified Signature*
Senior Lecturer     Signature:.........
Department of Mathematics,     Date : 01/05/2010
Faculty of Engineering,
University of Moratuwa.

To My Beloved Husband &Son Geethaka

With

Heartfelt Love and Admiration

# Abstract

Education plays an important role in the development of a country. Sri Lanka has given priority to widen access to education and has taken measures but whether everyone gets the equal opportunity to education is still open to question. This has been considered in an island wide National Youth Survey, conducted in the year 2000 which has been a joint undertaking involving United Nations Development Program (UNDP) and six Sri Lankan and German Institutions. The main aim of the survey has been to collect up-to-date and reliable information about opinions, values, perceptions, concerns, grievances and aspirations of the young generation in Sri Lanka. Further it has aimed at the identification and better understanding of the main commitment and ideas to solve them. It has been intended to provide a scientific database to assist policy makers and development organizations. The data set collected in this survey has data on specific segments related to youth in the age spectrum 15-29 years. This set of data is subjected to a social research. In this study, statistical evidence is sought to follow the proceedings of data mining.

A statistical modeling approach has been used in the analysis. The type of Further Educational Desire of a person has been found to be mainly related with the Type of Current Activity in terms of current employment status and so on, Educational Level, Province, Gender, Social Class and Age Group. Moreover, sufficient statistical evidence has been available to say that even, the Financial Situation in Past and the Major Problems Occurred in Education have an effect on developing their further educational desires. The importance of these findings is that they can be useful to understand the facilities and opportunities to be provided for students to assist in their education and to achieve their ultimate educational and career goals. Besides that, the information could be useful to assist in educational reforms and policy making.

A decision support tool has been developed using the rule based method in data mining using the information furnished by the inferential statistical analysis. The social characteristics "Type of Current Activity", "Educational Level", "Province", "Gender", "Social Class" and "Age Group" were found to be directly influential in predicting the type of further educational desires of youth and they were used in developing the decision support tool. It predicts what type of educational desires can be there with the individuals in a selected sample of youth. These predictions can be observed by any individual education provider or an organization before launching their educational centers targeting the selected population.

*M.Sc. in Operational Research, Project Dissertation, 2009*

i

# Acknowledgement

I wish to express my appreciation to the course coordinator, Mr. T.M.J.A. Cooray, the Head of the Department of Mathematics, Faculty of Engineering, University of Moratuwa, who has given his best assistance and advice towards the success of this project.

Words cannot express my gratitude to my supervisors Dr. Dilhari Attygalle and Mr. M. Firdhous, without whose expert advice, guidance and help this study would not have been easily get succeeded.

I specially thank Ms. Geethanjalee Henadheera for making me to do an important analysis on Youth Education in Sri Lanka by providing the data set.

I thankfully mention the support given by the staff members of Social Policy Analysis & Research Centre, Faulty of Arts, University of Colombo by giving valuable information on the data set and the questionnaire.

My Sincere gratitude goes to Prof. Roshini Sooriyarachchi, Ms. A. Karunaratne and Mr. R. A. B. Abeygunawardane, senior lecturers in the Department of Statistics, University of Colombo for guiding me to complete this project work punctually, overcoming all obstacles.

Last but not least, my heartfelt thanks go to all my friends who were always with me in-need and helped me in many ways.

# Table of Contents

## List of Tables

## List of Figures

*M.Sc. in Operational Research, Project Dissertation, 2009*