# Forecasting Financial Schedules for Cancer Patients

P.M.U.A.Jayatilaka

149214T

Faculty of Information Technology, University of Moratuwa, Sri Lanka

# Forecasting Financial Schedules for Cancer Patients

P.M.U.A.Jayatilaka

149214T

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

**May 2017**

# Declaration

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 Information derived from the published or unpublished work of others has been acknowledged in the text and the list of references is given

Name of the Student                                    Signature of the student

P.M.U.A.Jayatilaka

                                                                Date:

Supervised by                                          Signature of the supervisor

Mr. S.C.Premarathne

                                                                Date:

# Dedication

I dedicate this thesis to my parents who have always been my nearest and reverse nearest neighbors and have been so close to me that I found them with in me whenever I needed. It is their unconditional love that motivates me to set higher targets.

# Acknowledgement

# Abstract

During the time a patient is admitted in hospital rich sources of clinical, bio medical, contextual, and environmental data about patients have been available in medical and health sciences. These clinical sources of information are marked increasing in both volume and variety. Due to continuous increasing of the size of health care data, certain complexity raised in it. The hidden patterns among patient data can be extracted by applying different techniques. The techniques and tools are very helpful as they provide health care professionals with significant knowledge towards a decision.

This research has been conducted to analyze cancer patient's medical records in Sri Lanka in an effective manner to predict the future cost estimation for medicine. It is hypothesized that analyzing cancer patient's medical record can be done through machine learning data mining techniques. By doing so we can help the patients to schedule their financial matters in coming years. In Sri Lanka almost 23,000 individuals diagnosed with a new invasive cancer each year. The costs of the diagnosis and treatment for the cancer are considerable important. These costs have increased over the past years and are expected to increase in the future. It has been claimed that most of the families face difficulties with a diagnosis of cancer have financial issues. To respond to this challenge, we developed an interactive application which utilizes a multiple linear regression model to forecast financial schedule for cancer patients

This solution take data set collected from Apeksha Hospital -Maharagama as the input and predict the factors associate with the question. Having received the input this approach prepossess the data set to remove the anomalies. Then build the data model to predict the future cost for the patient. Linear regression technique is used to develop for prediction of future cost estimation. This total solution is build using WEKA data mining software. (WEKA GUI and WEKA API). According to the evaluation of the predictive model the correlation coefficient value is 0.5725 and the root mean squared error is Rs. 553183.25.

# Table of Content

# List of Figures

# List of Tables