# GENERAL APPROACH FOR CHURN PREDICTION WITH GENETIC ALGORITHM OPTIMIZED K-NEAREST NEIGHBOR FRAMEWORK

Tennakoon Mudiyanselage Nipunika Priyadarshani Tennakoon

(118231G)

Dissertation submitted in partial fulfillment of the requirements for the degree Master of Science

Department of Computer Science & Engineering

University of Moratuwa

Sri Lanka

June 2015

## Declaration

I declare that this is my own work and this dissertation does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my thesis, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature:                                                          Date:

Name:       University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

The above candidate has carried out research for the Masters Dissertation under my supervision.

Signature of the supervisor:                                        Date:

Name of the supervisor:

# Abstract

Customer churn has become one of the most significant topics in today's business. It has become a major challenge for a business with the evolving market and low barriers to switch between the service providers. It has identified that, retaining the old customers is more profitable for a company than acquiring new customers. That motivates business personnel for churn prediction. Service providers can get necessary measures to retain their customers if they could gain prior knowledge on the probable churns in the customer base. But, churn prediction is considered a difficult task.

Various attempts have been made in predicting churn and churn related information. Different data mining techniques had been used in developing churn models. Regression analysis, decision tree based methods and neural network based methods were among the most commonly used techniques. The most successful models suffered from low interpretability which is a main consideration in a churn model while some of the models were domain specific.

K nearest neighbor classifier is one of the best algorithms to be used in classifications. But it has been rarely used in churn prediction. Genetic algorithms are considered an optimization technique which could be used in optimizing performance of other algorithms. Genetic Algorithm Optimized K Nearest Neighbor (gaKnn) is a framework that has tested for its high accuracy. Hence, we developed a Tool based on the gaKnn framework which could be used for churn prediction. We also incorporated two voting mechanisms; Bayesian weights and class confidence weights (*ccw*) to weight the prediction in order to address misclassification issues occur due to class skew.

# Acknowledgements

I sincerely acknowledge the work carried out by the WEKA community and publishing it as an open source code.

I would like to dedicate my sincere thank to my supervisor Dr. Amal Shehan Perera, Senior Lecturer, Department of Computer Science and Engineering for his dedicated support for the success of this research. This would not have become a success without your support from the initial stage of the research.

I would like to thank the entire academic and the non academic staff of the Department of Computer Science and Engineering for their kindness extended to me in every aspect.

Last but not least, I thank my family and all my friends who supported me for the success of this piece of work. Your support was so precious.

University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

## Table of Contents

## List of abbreviations

| Abbreviation | Description |
|---|---|
| AUC | Area Under Curve |
| CCW | Class Confidence Weight |
| CDR | Call Detail Record |
| CSV | Comma Separated Values |
| DMEL | Data Mining by Evolutionary Learning |
| FN | False Negative |
| FP | False Positive |
| GA | Genetic Algorithm |
| gaKnn | Genetic Algorithm Optimized K nearest neighbor |
| GUI | Graphical User Interface |
| JGAP | Java Genetic Algorithms Package |
| KNN | K nearest Neighbor |
| NN | Neural Network |
| PR | Precision- recall |
| ROC | Receiver Operating Characteristic |
| TN | True Negative |
| TP | True Positive |

# Table of figures

# List of tables