



# **10-Year Cardiovascular Disease (CVD) Risk Prediction of Sri Lankans: A Longitudinal Cohort Study**

**M.B Solangaarachchige**  
(Reg. No.: MS19805306)  
M.Sc. in IT

Supervisor: Mr Prasanna S. Haddela

December 2021

**Word Count 21,496**

**Faculty of Graduate Studies & Research  
Sri Lanka Institute of Information Technology**

# Acknowledgments

I'd like to express my gratitude to everyone whose encouragement and support helped me complete my thesis.

First and foremost, I want to thank my supervisor, Mr Prasanna S. Haddela, for his guidance, support, and encouragement in supporting the completion of this thesis. He has provided me with an abundance of invaluable advice and knowledge that will benefit me well in both my career and personal life. He was indeed the most patient and most understanding lecturer a student could ask for. Mr. Prasanna's guidance and teaching helped me understand and use machine learning techniques in this research. Thank you very much, Sir, for your unwavering support throughout this study.

Second, I'd like to express my heartfelt gratitude to Prof K.Chamila D.Mettananda for her invaluable guidance and unwavering support and assistance throughout this study. Prof Chamila's medical expertise was fundamental and extremely beneficial. Without her assistance, this research would not have been a success. Her guidance, support, and encouragement were priceless from the start of the project, when gathering the data, until the finish by validating the model.

I would also like to express my sincere gratitude to Vidyajyothi Professor H Janaka De Silva, Professor A R Wickremasinghe and Professor Anuradhani Kasthuriratne for believing in me and agreeing to provide a valuable dataset to make this project a success. Prof Janaka was extremely humble to discuss the feasibility of this study with me and my supervisor and ultimately decided to grant permission to use the dataset. Thank you, Sir, for making this possible.

I would like to thank our MSc in IT coordinator Mr. Samantha Rajapaksha for his support and valuable time and input throughout this research. His counsel to issues that developed over the two years of MSc degree and this thesis study is remarkable.

Lastly, I would like to express my sincere appreciation to my family for encouraging and supporting me throughout the study.

# Abstract

Cardiovascular diseases are one of the leading causes of mortality in the world. A cornerstone of preventive cardiology is identifying individuals at risk of cardiovascular diseases (CVD) at the earliest. Clinical guideline primarily recommends risk prediction models that are based on a limited number of predictors that perform poorly across all patient groups. Predicting cardiovascular risk is crucial for making treatment decisions, especially in the primary prevention of CVDs using a total risk approach. Despite the fact that several cardiovascular risk prediction models exist, only a handful are specifically designed for Asians, and none are generated from South Asians, including Sri Lankans. Machine learning (ML) and neural networks appear to be increasingly promising in supporting decision-making and forecasting from the huge amounts of data generated by the healthcare industry. This led us to develop a CVD model using Machine Learning to predict 10-year risk of developing a CVD in Sri Lankans. We investigated whether we could adopt ML to develop a model and whether there is an improvement in including non-traditional variables for the accuracy of CVD risk estimates and how to validate the ML model with existing WHO risk charts.

Using data on 2596 participants without CVD at baseline data collection of Ragama Medical Officer of Health (MOH) area in Sri Lanka, we developed a ML-based model for predicting CVD risk based on 75 available variables. However, the ratio of developing a CVD vs no CVD in 10 years was 7:93, which is extremely unbalanced. Therefore, at first, we derived a balanced dataset from the main dataset and build a ML model and it recorded an 80.56% accuracy. Secondly, to alleviate the dataset's imbalance, we adopted two techniques, which are 10-fold cross validation and stratified 10-fold cross validation (SKF) and trained six ML classification algorithms. They are Random Forest (RF), Decision Tree, AdaBoost, Gradient Boosting, K-Nearest Neighbor and 2D Neural Network. Out of these six algorithms RF model with SKF showed the highest accuracy in predicting a CVD event with an accuracy of 93.11%. Our ML model included predictors that are not usually considered in existing risk prediction models. Systolic blood pressure was the most important variable in this model. There were six non-traditional variables in the most ten important variable list and three of them were non-laboratory variables. To validate the model with existing WHO risk charts, we explored an experimental approach by developing a simple logistic regression function using the same techniques as the best selected model, with the seven traditional risk factors used in WHO risk charts and our

Random Forest model indicated the highest accuracy compared to the WHO model, with a difference of 26.20 %.

Our ML model improves the accuracy of CVD risk prediction in the Sri Lankan population. This approach justifies that the CVD prediction models also can be derived using ML for each subregion individually. Additionally, our research discovered novel CVD disease factors that may now be investigated in prospective studies.

**Keywords:** cardiovascular disease, risk assessment, models, machine learning, classification

# Table of Contents

Acknowledgments.....	ii
Abstract.....	iii
Table of Contents.....	v
List of Figures.....	viii
List of Tables.....	ix
List of Abbreviations.....	x
Chapter 1 Introduction.....	11
1.1 Cardiovascular Disease (CVD).....	12
1.2 What are cardiovascular diseases?.....	12
1.2.1 Types of CVD.....	12
1.2.2 The Risk Factors.....	13
1.3 Existing Cardiovascular Risk Prediction Models.....	13
1.3.1 WHO Risk Charts.....	13
1.4 Motivation.....	14
1.5 Outline.....	15
Chapter 2 Research Question & Objectives.....	17
2.1 Research Question.....	17
2.2 Objectives.....	17
2.2.1 General Objective.....	17
2.2.2 Sub Objectives.....	17
Chapter 3 Literature Review.....	18
Chapter 4 Materials & Methodology.....	24
4.1 Study Design.....	24
4.2 Study Setting.....	24
4.3 Study population.....	25
4.3.1 Exclusion Criteria.....	25
4.4 Machine Learning (ML).....	25
4.4.1 Supervised Learning.....	27
4.4.2 Machine Learning Life Cycle.....	28
4.4.3 The learning process adapted to the problem.....	29
4.5 Acquiring Data.....	30
4.6 Data Preparation.....	30
4.6.1 Data Exploration.....	30
4.6.2 Data Pre-Processing.....	34

Chapter 5 Model Building/ Model Fitting .....	40
5.1 Initial Model Building.....	40
5.2 Machine Learning Techniques.....	40
5.2.1 Random Forest .....	41
5.2.2 AdaBoost.....	42
5.2.3 Decision tree .....	43
5.2.4 Gradient Boosting .....	44
5.2.5 K-Nearest Neighbor .....	46
5.2.6 Neural Network.....	47
5.3 Train & Test Model .....	48
5.3.1 Failure of k-Fold Cross-Validation .....	49
5.3.2 Fix Cross-Validation for Imbalanced Classification.....	50
5.3.3 Hyper Parameter Tuning.....	51
5.4 Evaluate the performance of the ML model .....	52
5.4.1 Confusion Matrix .....	52
5.4.2 Precision.....	53
5.4.3 Recall .....	54
5.4.4 F-Measure .....	54
5.4.5 Accuracy .....	54
5.4.6 ROC Curve and ROC AUC .....	55
5.4.7 Precision-Recall Curve and AUC .....	55
Chapter 6 Results .....	57
6.1 Results of the Initial Model.....	57
6.2 Results of the Main Model.....	59
6.2.1 10-Fold Cross Validation.....	59
6.2.2 Stratified 10-Fold Cross Validation .....	59
6.3 Validating the ML model predictions .....	62
Chapter 7 Discussion .....	64
7.1 Machine Learning Model.....	64
7.2 Improvement in including non-traditional variables.....	65
7.3 Validating the ML model predictions against WHO risk chart predictions.....	66
Chapter 8 Limitations and Future Work .....	67
Chapter 9 Conclusion.....	68
Bibliography .....	69
Appendix.....	73
Appendix 1: WHO Cardiovascular Disease Risk Charts for Southeast Asia .....	73

Appendix 2: RHS questioner .....	75
Appendix 3: Exploratory Data Analysis (EDA) of Baseline Dataset (2007).....	81
Appendix 4: Feature importance of the initial model .....	90
Appendix 5: Source code for splitting the dataset using k-fold cross validation.....	92
Appendix 6: Source code for splitting the dataset using stratified k-fold cross validation.....	93
Appendix 7: Source code of the Initial Model .....	94
Appendix 8: Feature importance of the main model.....	96

# List of Figures

Figure 4.1 Machine Learning vs Traditional Programming.....	26
Figure 4.2 The Machine Learning Life Cycle.....	28
Figure 4.3 The shape of preliminary database .....	31
Figure 4.4 The shape of baseline database.....	32
Figure 4.5 Database information.....	32
Figure 4.6 Visualization of Missing Values with a bar chart.....	33
Figure 4.7 Target Variable .....	33
Figure 4.8 Boxplot and distribution plot of average height .....	34
Figure 5.1 Random Forest with two trees .....	41
Figure 5.2 AdaBoost algorithm intuition .....	43
Figure 5.3 General structure of a decision tree .....	44
Figure 5.4 KNN Example .....	47
Figure 5.5 2-layer Neural Network.....	48
Figure 5.6 Splits of k-fold cross validation for CVD database .....	50
Figure 5.7 Splits of Stratified k-fold cross validation for CVD database .....	50
Figure 5.8 Example of a grid search .....	51
Figure 5.9 Confusion Matrix.....	53
Figure 6.1 The plot of feature importance of the initial model .....	58
Figure 6.2 The ten most important features of the initial model.....	58
Figure 6.3 Main ML Model comparison.....	60
Figure 6.4 ROC-AUC plot of Random Forest model in stratified 10-fold cross validation .....	60
Figure 6.5 PR AUC plot of Random Forest model in stratified 10-fold cross validation.....	61
Figure 6.6 The plot of feature importance of the main Random Forest model.....	61
Figure 6.7 The ten most important features of the main model .....	62
Figure 6.8 ROC-AUC plots of Random Forest and Logistic Regression models.....	63
Figure 6.9 PR AUC plots of Random Forest and Logistic Regression models .....	63



# List of Tables

Table 4.1 Data Transformation .....	38
Table 4.2 Feature Description .....	38
Table 6.1 Classification Report of the Initial Model.....	57
Table 6.2 10-Fold Cross Validation Results .....	59
Table 6.3 Stratified 10-Fold Cross Validation Results .....	59
Table 6.4 Validating the ML model predictions with WHO risk prediction charts.....	62

# List of Abbreviations

IT	Information Technology
UN	United Nation
WHO	World Health Organization
CVD	Cardiovascular Disease
CVDs	Cardiovascular Diseases
ML	Machine Learning
SKF	Stratified K-Fold
RF	Random Forest
2D	2 Dimensional
EDA	Exploratory Data Analysis
RHS	Ragama Health Study
CV	Cardiovascular
CVEs	Cardiovascular Events
kNN	k-Nearest Neighbor
KNN	K-Nearest Neighbor
ROC	Receiver Operating Characteristic
AUC	Area Under the Curve
PR	Precision-Recall
IHD	Ischemic Heart Disease
NCDs	Noncommunicable Diseases
NCD	Noncommunicable Disease
MOH	Medical Officer of Health
GN	Grama Niladhari