

MODELING OF CARDIOVASCULAR
DISEASES (CVDS) AND DEVELOPMENT OF
PREDICTIVE HEART RISK SCORE



MIRZA RIZWAN SAJID

اونيورسيتي مليسيا قهغ

UNIVERSITI MALAYSIA PAHANG

DOCTOR OF PHILOSOPHY

UNIVERSITI MALAYSIA PAHANG

UNIVERSITI MALAYSIA PAHANG

DECLARATION OF THESIS AND COPYRIGHT

Author's Full Name : MIRZA RIZWAN SAJID

Date of Birth : 15 FEBRUARY 1983

Title : MODELING OF CARDIOVASCULAR DISEASES (CVDs)
AND DEVELOPMENT OF PREDICTIVE HEART RISK SCORE

Academic Session : SEMESTER 2 2020/2021

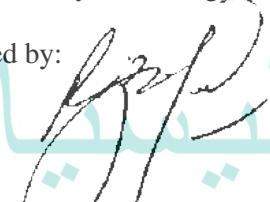
I declare that this thesis is classified as:


- CONFIDENTIAL (Contains confidential information under the Official Secret Act 1997)*
- RESTRICTED (Contains restricted information as specified by the organization where research was done)*
- OPEN ACCESS I agree that my thesis to be published as online open access (Full Text)

I acknowledge that Universiti Malaysia Pahang reserves the following rights:

1. The Thesis is the Property of Universiti Malaysia Pahang
2. The Library of Universiti Malaysia Pahang has the right to make copies of the thesis for the purpose of research only.
3. The Library has the right to make copies of the thesis for academic exchange.

Certified by:


(Student's Signature)


(Supervisor's Signature)

CP5573202
New IC/Passport Number
Date: 5 JULY 2021

DR. NOR YANTI MUHAMMAD
Name of Supervisor
Date: 5 JULY 2021

NOTE: * If the thesis is CONFIDENTIAL or RESTRICTED, please attach a thesis declaration letter.

SUPERVISOR'S DECLARATION

We hereby declare that we have checked this thesis and, in our opinion, this thesis is adequate in terms of scope and quality for the award of the degree of Doctor of Philosophy.

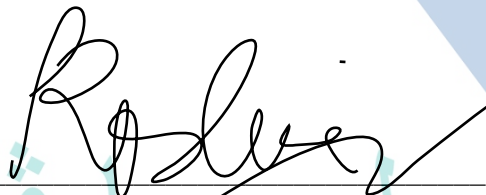


(Supervisor's Signature)

Full Name : DR. NORYANTI MUHAMMAD

Position : SENIOR LECTURER

Date : 5 JULY 2021



(Co-supervisor's Signature)

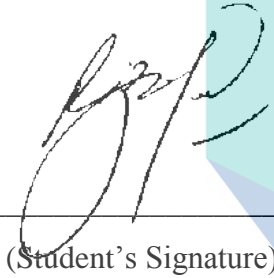
Full Name : DR. ROSLINAZAIRIMAH BINTI ZAKARIA

Position : ASSOCIATE PROFESSOR

Date : 5 JULY 2021

STUDENT'S DECLARATION

I hereby declare that the work in this thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at Universiti Malaysia Pahang or any other institutions.



(Student's Signature)

Full Name : MIRZA RIZWAN SAJID

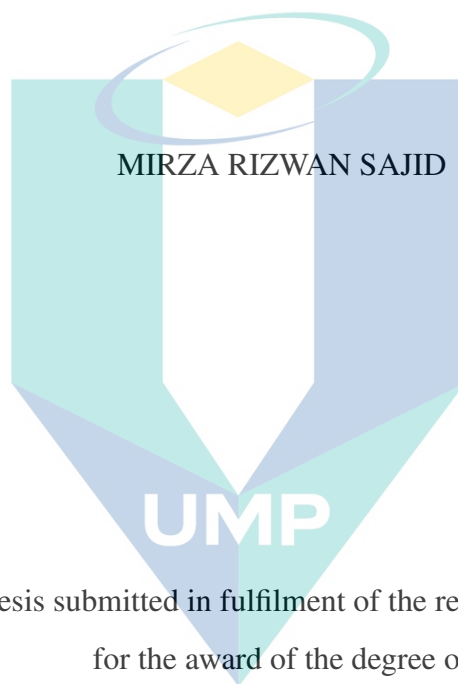
ID Number : PSS17002

Date : 5 JULY 2021

اونيورسيتي مليسيا قهغ

UNIVERSITI MALAYSIA PAHANG

MODELING OF CARDIOVASCULAR DISEASES (CVDs) AND DEVELOPMENT
OF PREDICTIVE HEART RISK SCORE



Thesis submitted in fulfilment of the requirements
for the award of the degree of
Doctor of Philosophy

اونيورسيتي ملايسيا قهغ

UNIVERSITI MALAYSIA PAHANG

Centre for Mathematical Sciences

UNIVERSITI MALAYSIA PAHANG

JULY 2021

ACKNOWLEDGEMENTS

Praise to Allah the Almighty for bestowing us with the best strengths and qualities to be used for the benefit of mankind. As a humble servant of Allah SWT, I am trying to express my gratitude in the form of little knowledge. I am a follower of His last Messenger Muhammad (PBUH), an eternal source of knowledge and guidance towards humanity's truth and is the gold standard for all intelligence.

Firstly, I am immensely grateful to my main supervisor, Dr. Noryanti Muhammad, and co-supervisor Associate Professor Dr. Roslinazairimah Zakaria for their academic guidance and unlimited support throughout my humble struggle to complete this thesis. Dr. Noryanti's patience, advice, thought-provoking ideas, objective-oriented approaches, and unconditional support have been invaluable to me. I am also thankful to Dr. Ahmad Shahbaz from Punjab Institute of Cardiology for his valuable help in data collection, especially for the standardized ethical considerations.

I am highly indebted to my dearest parents, wife, and my children (Muhammad Shahaan Alam and Emaan Fatima) for their moral support and patience throughout this work. It was really tough for them to manage their things alone, which were supposed to be my duties. Special thanks to my sisters, their husbands, and brother, whose prayers have always been sources of strength and inspiration for me.

I also owe my thankfulness to all my teachers, friends and colleagues, especially Dr. Waris Ali Khan, Dr. Nadeem Shafique Butt, Dr. Asif Hanif, Asim Butt, Arshad Ali Khan, Dr. Syed Ahmad Chan Bukhari, Jawad Asif, Dr. Waqas Sami, Dr. Fayyaz Ahmad, Mehfooz Ullah Dar, Muhammad Bilal, Dr. Asad Ullah Khan and Muhammad Ilyas for their encouragement, appreciation and support during my research.

My final thanks to the University of Gujrat (UoG) for providing me the opportunity to pursue my studies at Universiti Malaysia Pahang (UMP). I am also thankful to the Centre for Mathematical Sciences, College of Computing & Applied Sciences for offering me an excellent academic and conducive environment, which had enabled me to study smoothly.

UNIVERSITI MALAYSIA PAHANG

ABSTRAK

Penyakit kardiovaskular (CVDs) adalah penyebab utama kematian dengan 31% kematian global. Tujuan kajian ini adalah untuk membangunkan model lintasan yang sah secara statistik yang mempertimbangkan kemungkinan lintasan bukan linear dan ciri binari endogenos, dan pengantara keduanya bagi status CVDs. Kajian ini menumpukan pembangunan pelbagai bentuk model ramalan risiko tempatan dan penggunaan skor risiko jantung yang ringkas dengan menggunakan ciri bukan makmal dan algoritma pembelajaran mesin (*machine learning*) (*ML*). Walau bagaimanapun, penukaran bentuk algoritma *ML* yang kompleks kepada model statistik yang ringkas menjadi perhatian utama didalam kajian ini. Kajian kawalan kes yang sesuai dengan jantung dilakukan di Institut Kardiologi Punjab, Pakistan, di mana 460 individu sebagai sampel dipilih melalui persampelan bersistematik. Kaedah *warp-partial least square* digunakan untuk mengganggar pelbagai lapisan model lintasan yang dihipotesiskan. Model ini menganggarkan pekali *warp* menggunakan keseluruhan aliran linear yang terdapat dalam segmen linear daripada hubungan bukan linear. Model yang dibangunkan ini merupakan laluan pintasan yang novel di mana ciri demografi dan sosioekonomi menjadi pemacu utama bagi ciri tingkah laku, yang membawa kepada status CVD secara langsung dan tidak langsung melalui metabolic sindrom. Dalam membangunkan model ramalan risiko, dua algoritma *ML*, iaitu sokongan linear mesin vektor (*linear support vector machine*) dan *artificial neural network* mengatasi model konvensional iaitu analisis regresi logistik (*Logistic Regression Analysis*) (*LRA*). Prestasi model yang dibangunkan ini dinilai melalui pelbagai matriks yang ditetapkan dengan menggunakan pengesahan bersilang 10-kali ganda (*10-fold cross-validation*). Kemudian, satu novel metodologi dibangunkan dan digunakan untuk mengira skor risiko jantung yang ringkas berdasarkan ciri bukan makmal yang dipanggil *Non-Laboratory based Heart Risk Score (NLHRS)*. Metodologi ini menyusun dan mengumpul algoritma *ML* yang terbaik dan digunakan sebagai asas untuk mengira indeks pemberat ciri relatif (*relative feature weights*). Indeks pemberat ini disebut sebagai *NLHRS*, yang selanjutnya digunakan sebagai kovariat di model *LRA* ringkas untuk menganggarkan kemungkinan CVD berlaku. Perubahan yang berlaku iaitu dari model algoritma yang bersifat kotak-hitam kompleks kepada model statistik ringkas yang menghasilkan model yang tidak memerlukan sistem berautomatik untuk pelaksanaannya. *NLHRS* yang berasaskan algoritma *ML* dan model yang bersangkutan yang dibangunkan ini menunjukkan prestasi yang lebih baik daripada model sedia ada yang berdasarkan skor risiko semi-kuantitatif dari segi penilaian diskriminasi dan penentuan. Akhirnya, keupayaan model ramalan *NLHRS* juga diuji dan disesuaikan mengikut strata penduduk. Kajian ini menyimpulkan beberapa perkara. Pertama, penggunaan pendekatan kaedah yang fleksibel dalam anggaran dapat memodelkan ciri binari bagi status CVD dan lintasan tidak linear didalam lintasan model yang kompleks. Model lintasan CVD yang dianggarkan dapat digunakan sebagai strategi penangguhan penyakit dalam pengaturan klinikal. Kedua, model algoritma *ML* menawarkan model ramalan risiko yang lebih baik dan konsisten berbanding model yang berasaskan *LRA*. *NLHRS* dan model yang berkaitan dengannya yang dibangunkan merupakan hasil metodologi novel yang memberikan bentuk skor risiko yang sah dan ringkas dan dapat digunakan tanpa sistem berautomatik.

ABSTRACT

Cardiovascular diseases (CVDs) are the leading cause of death, with 31% of global mortality. The purpose of this study is two folds such as the development of a statistically valid path model which considered the possible non-linear paths, mediators, and binary endogenous feature of CVDs status. Further, it focuses on the development of various forms of local risk prediction models and simple heart risk scores using non-laboratory features and machine learning (*ML*) algorithms. However, the conversion of a complex form of *ML* algorithms into a simple statistical model is the prime concern. A gender-matched case-control study was conducted in Punjab Institute of Cardiology, Pakistan, in which a sample of 460 individuals was selected through systematic sampling. The warp-partial least square method was utilized to estimate the multi-layer hypothesized path model. This model estimated warped coefficients using the overall linear trend found in linear segments of non-linear relationships. This model found novel pathways in which demographic and socioeconomic features are the main drivers of behavioral features, leading to CVDs status directly and indirectly through metabolic syndrome. In developing risk prediction models, two *ML* algorithms, linear support vector machine and artificial neural network outperformed the existing conventional logistic regression analysis (LRA) model. The performance of the models was assessed through various established matrices using 10-fold cross-validation. A novel methodology was used to compute simple heart risk scores called non-laboratory based heart risk score (*NLHRS*). The methodology is proposed as stacking ensemble *ML* and the best *ML* algorithms are used as a base learner to compute relative feature weights. The index of these weights is referred to as *NLHRS*, which was further used as a covariate in the simple LRA model to estimate the likelihood of CVDs. This conversion from a complex black-box nature of *ML* algorithms into simple statistical models yielded such models, which do not require automated systems for their implementation. *ML*-based *NLHRS* and their associated models outperformed the existing semi-quantitative risk score-based model in terms of discrimination and calibration assessments. Finally, the predictive capability of valid *NLHRS* models has also been tested and adjusted for different strata of the population. Firstly, the study concludes that the adoptions of the flexible approach in estimation can model the binary feature of CVDs and non-linear paths in the complex path models. The estimated CVDs path model can be implemented as a disease delay strategy in clinical settings. Secondly, the *ML* models offer better and consistent risk prediction models as compared to LRA-based model. The *NLHRS* and their associated models which are the outputs of novel methodology provide valid and simple forms of risk scores and can be used without automated systems.

TABLE OF CONTENTS

DECLARATION

TITLE PAGE

ACKNOWLEDGEMENTS

ii

ABSTRAK

iii

ABSTRACT

iv

TABLE OF CONTENTS

v

LIST OF TABLES

xi

LIST OF FIGURES

xiii

LIST OF SYMBOLS

xv

LIST OF ABBREVIATIONS

xvii

LIST OF APPENDICES

xx

CHAPTER 1 INTRODUCTION

1

1.1 Overview

1

1.2 Problem Statement

4

1.3 Research Questions

6

1.4 Objectives of the Study

6

1.5 Research Scope

7

1.6 Research Activities

7

1.7 Overall Framework of Thesis

11

CHAPTER 2 LITERATURE REVIEW

12

2.1 Introduction

12

2.2 Section A: Basic Terminology and Epidemiology of CVDs

12

2.2.1 Types of Cardiovascular Disease

13

2.2.2 Risk Factors and its Types

13

2.2.3	Metabolic Syndrome and its Components	15
2.2.4	Epidemiology of Disease	15
2.3	Section B: Underpinning Theories and Theoretical Framework of Study	18
2.3.1	Relationship of Demographic Factors, Family History with Behavioral Risk factors	21
2.3.2	Relationship of Demographic Factors, Family History with CVDs	24
2.3.3	Relationship of Socio-economic Factors with Behavioral Risk Factors	26
2.3.4	Relationship of Behavioral Risk Factors and CVDs	29
2.3.5	Relationship of Behavioral Risk Factors and MS	31
2.3.6	Relationship of MS and CVDs	32
2.3.7	MSA as Mediator in CVDs Model	33
2.3.8	Hypothesized Path Model	35
2.3.9	Main Statistical Techniques Used in Causal Relationships	37
2.3.10	Gap Analysis for Theory-driven Path Models	39
2.4	Section C: Risk Prediction Models and Heart Risk Scores	40
2.4.1	Issues in Existing Risk Prediction Models	44
2.4.2	Traditional Algorithms to Machine Learning (<i>ML</i>)	47
2.4.3	Brief Introduction of Machine Learning <i>ML</i> Algorithms	49
2.4.4	Black-box Nature of <i>ML</i> Algorithms	52
2.4.5	Risk Stratification of Risk Prediction Models	54
2.4.6	Gap Analysis of Risk Prediction Models	55
2.5	Summary of the Chapter	56
2.6	Concluding Remarks	57
CHAPTER 3 METHODOLOGY FOR DEVELOPMENT OF PATH MODEL		58
3.1	Introduction	58
3.2	Pre-test	58
3.3	Pilot Study	59
3.4	Main Study	59
3.4.1	Research Design	60

3.4.2	Inclusion and Exclusion Criteria of Cases and Controls	60
3.4.3	Confounding	61
3.4.4	Target Population	62
3.4.5	Sampled Population and Selection of Hospital	63
3.4.6	Sample Size Estimation	64
3.4.7	Sampling Technique	67
3.4.8	Data Collection Methods and Tools	68
3.4.9	Research Ethics	69
3.4.10	Data Management	69
3.4.11	Flow of Data Analysis	70
3.5	Description of Statistical Techniques	70
3.5.1	Reliability Analysis	70
3.5.2	Exploratory Factor Analysis	71
3.6	Structural Equation Modeling	71
3.6.1	Path Analysis and WarpPLS	72
3.7	Summary of the Chapter	74
3.8	Concluding Remarks	75

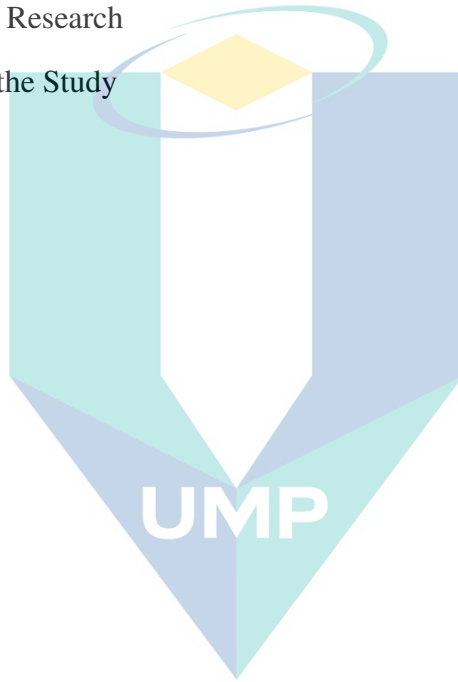
CHAPTER 4 METHODOLOGY FOR DEVELOPMENT OF RISK PREDICTION MODELS

4.1	Introduction	76
4.2	Baseline Study	76
4.3	Methodology for Development of Quantitative and <i>NLHRS</i> -based Risk Prediction Models	78
4.3.1	Development of Quantitative Risk Prediction Models	80
4.3.2	Development of <i>NLHRS</i> and their Risk Prediction Models	84
4.4	Description of Statistical Techniques	88
4.4.1	Confirmatory Factor Analysis (CFA)	88
4.4.2	Multiple Linear Regression Analysis	90
4.4.3	Moderation Analysis	91

4.5	Classification and Prediction Techniques	92
4.5.1	Bivariate Odds Ratio Analysis	93
4.5.2	Logistic Regression Analysis (LRA)	93
4.5.3	Decision Tree Learning	95
4.5.4	Artificial Neural Networks	98
4.5.5	Support Vector Machine	100
4.5.6	Cross-validation of Risk Prediction Models and Heart Risk Scores	102
4.5.7	Comparative Performance Analysis of Classifiers	103
4.6	Features Importance Approaches	104
4.7	Validation of <i>NLHRS</i> -based Risk Prediction Models	108
4.7.1	Brier Score (BS)	108
4.7.2	Spiegelhalter's Z-statistic	109
4.7.3	Discrimination	109
4.7.4	Calibration	110
4.8	Identification of Thresholds of <i>NLHRS</i> through QUEST Algorithm	110
4.9	Findings from the Baseline Study	111
4.9.1	Assessment of Relationship Between CE and CVM	112
4.9.2	Moderation Analysis through Urbanization for Relationship Between CE and CVM	113
4.9.3	Interaction Graphs	116
4.9.4	Discussion	118
4.9.5	Role of Baseline Study in the Main Study	121
4.10	Summary of the Chapter	121
4.11	Concluding Remarks	122
CHAPTER 5 MODELING OF CARDIOVASCULAR DISEASES		124
5.1	Introduction	124
5.2	Section-I: Findings from Pre-test and Pilot Study	124
5.2.1	Results of Pre-test	125
5.2.2	Results of Pilot Study	125

5.3	Section-II: Path Modeling of Cardiovascular Diseases	129
5.3.1	Response Rate	130
5.3.2	Basic Characteristics of Study Individuals	130
5.3.3	Path Analyses	131
5.3.4	Assessment of Linear and Nonlinear Relationships	133
5.3.5	Assessment of Mediation through Indirect Effects	135
5.4	Discussions	135
5.5	Concluding Remarks	142
CHAPTER 6 DEVELOPMENT OF RISK PREDICTION MODELS AND PREDICTIVE HEART RISK SCORE		145
6.1	Introduction	145
6.2	Section-I: Development of Quantitative Risk Prediction Models	146
6.2.1	Development of Quantitative Risk Prediction Models through LRA and <i>ML</i> Algorithms	147
6.3	Section-II: Development of <i>NLHRS</i> and Their Risk Prediction Models	155
6.3.1	Implementation of Approach-I for Computation of <i>NLHRS</i>	156
6.3.2	Implementation of Approach-II for Computation of <i>NLHRS</i>	157
6.3.3	<i>NLHRS</i> -based Risk Prediction Models	163
6.3.4	Validation Process of <i>NLHRS</i> -based Risk Prediction Models	166
6.3.5	Identification of Risk Thresholds for Valid Heart Risk Scores	170
6.4	Section-III: Risk Stratification of Valid <i>NLHRS</i> -based Risk Prediction Models	178
6.4.1	Moderated Logistic Regression for ANN-RS and CVDs Status	180
6.4.2	Moderated Logistic Regression for SVM-RS and CVDs Status	182
6.4.3	Finalized Pooled and Stratified Risk Prediction Models	184
6.4.4	Discussions	186
6.5	Concluding Remarks	196
CHAPTER 7 CONCLUSION		198
7.1	Introduction	198

7.2	Research Questions and Objectives Revisited	198
7.2.1	Research Question 1	198
7.2.2	Research Question 2	199
7.2.3	Research Question 3	200
7.2.4	Research Question 4	202
7.2.5	Research Question 5	203
7.3	Contributions of the Study	204
7.4	Significance of Research	205
7.5	Limitations of the Study	206
7.6	Future Works	207
REFERENCES		209



اونيورسيتي ملايسيا قهغ

UNIVERSITI MALAYSIA PAHANG

LIST OF TABLES

Table 2.1	Different Criteria for Measuring Metabolic Syndrome	16
Table 2.2	Trade-off Between Accuracy and Interpretability of Models	53
Table 3.1	Sample Size Requirements in Different Statistical Techniques	65
Table 3.2	Main Contents of Questionnaire	69
Table 3.3	Model Fitness Indicators	74
Table 4.1	Performance Matrices for Classifiers	104
Table 4.2	Input-hidden-output Connection Weights Methodology	107
Table 4.3	Interpretation of C-Statistic for Different Thresholds	110
Table 4.4	Levels of Urbanization	114
Table 4.5	Descriptive Statistics of CE,CVM and Confounder in the Overall Sample and Different Forms of Urbanization	115
Table 4.6	Regression and Moderation Analysis	116
Table 5.1	EFA for Dietary Habits in Cases and Controls	128
Table 5.2	Basic Characteristics of Study Individuals	131
Table 5.3	Direct Path Coefficients of the Path Model	133
Table 5.4	Assessment of the Mediation Impact in the Path Model	135
Table 5.5	Overall Hypothesis Testing	144
Table 6.1	Input Features	147
Table 6.2	Point and Interval Estimates of Bivariate Odds Ratio	149
Table 6.3	Feature Importance using Forward Stepwise Logistic Regression Analysis	150
Table 6.4	Performance Assessment of Baseline LRA Model through 10-fold Cross-validation	151
Table 6.5	Performance Assessment of <i>ML</i> -based Risk Prediction Models	154

Table 6.6	Percentage Change in Performance Matrices of <i>ML</i> -based Models from Baseline LRA Model	154
Table 6.7	Overall Performance Comparison of Quantitative Risk Prediction Models	155
Table 6.8	Performance Assessment of <i>ML</i> Models for the Development of <i>NLHRS</i>	157
Table 6.9	Percentage Change in Performance Matrices of <i>ML</i> -based Models from LRA Model	158
Table 6.10	Overall Comparison of <i>ML</i> Models with LRA Model for Developing <i>NLHRS</i>	158
Table 6.11	Input-hidden-output Nodes Connection Weights	160
Table 6.12	Absolute Product of Input-hidden and Hidden-output Connections	160
Table 6.13	Relative Contribution of Each Feature to the Output Node through Each Hidden Node	161
Table 6.14	Relative Importance of Each Feature in the Overall Artificial Neural Network	161
Table 6.15	Linear Support Vector Machines-based Feature Importance	163
Table 6.16	Performance Assessment of <i>NLHRS</i> -based Risk Prediction Models	166
Table 6.17	Percentage Change in Performance Matrices of <i>ML</i> -based <i>NLHRS</i> Models from RFS Model	166
Table 6.18	Summary of Validity Assessment of <i>NLHRS</i> -based Risk Prediction Models	168
Table 6.19	Categories of Risk based on Valid <i>NLHRS</i>	175
Table 6.20	Gender-wise Identification of Risk Thresholds for ANN-RS	176
Table 6.21	Gender-wise Identification of Risk Thresholds for SVM-RS	176
Table 6.22	Risk Stratification: Discrimination and Calibration of <i>NLHRS</i> -based Risk Prediction Models by Area of Living	179
Table 6.23	Moderation Analysis for ANN-RS by Area	181
Table 6.24	Moderation Analysis for SVM-RS	183

LIST OF FIGURES

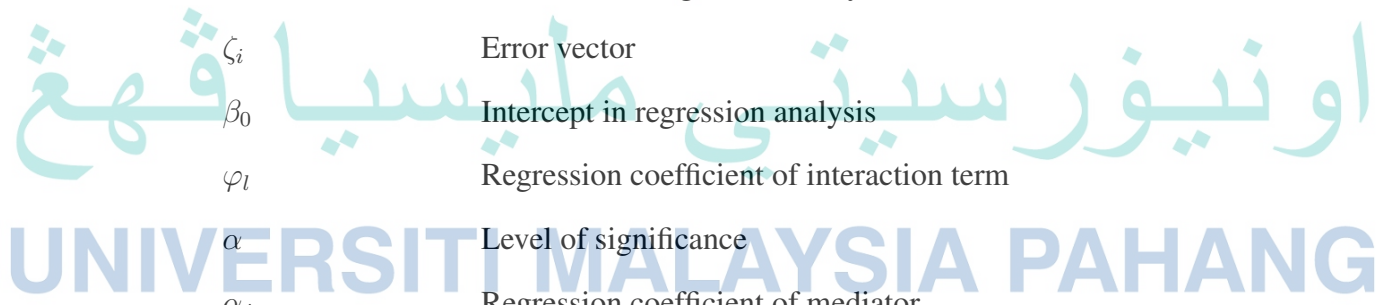
Figure 1.1	Research Activities	10
Figure 1.2	Overall Framework of Thesis	11
Figure 2.1	Causes of Deaths in the World	18
Figure 2.2	Wider Determinants of Health Model	20
Figure 2.3	Kastorini Model	36
Figure 2.4	Hypothesized Model of Study	37
Figure 2.5	Types of Risk Prediction Models	42
Figure 2.6	Example of Artificial Neural Network (ANN)	51
Figure 2.7	Example of Support Vector Machine (SVM)	52
Figure 2.8	Black box Nature of <i>ML</i> Algorithms	53
Figure 3.1	Political Map of Pakistan	63
Figure 4.1	Development of Baseline Quantitative Risk Prediction Model	82
Figure 4.2	Development of Machine Learning based Quantitative Risk Prediction Models	84
Figure 4.3	Development of Risk Factor Score (RFS) or Semi-quantitative Risk Score	86
Figure 4.4	Transformation of <i>ML</i> Model into Simple Statistical Model	88
Figure 4.5	Development and Validation of <i>NLHRS</i> -based Risk Prediction Model	89
Figure 4.6	Features Importance Approaches	106
Figure 4.7	Validation Process	108
Figure 4.8	PPLUA as Moderator for CE and CVM	117
Figure 4.9	Urbanization Status as Moderator for CE and CVM	117
Figure 4.10	Level of Urbanization as Moderator for CE and CVM	118
Figure 5.1	Finalized Path Model	134

Figure 6.1	Computation of Product of Input-hidden Connections and Hidden-output Connections for Age Groups	162
Figure 6.2	RFS and Estimated Risk of CVDs	165
Figure 6.3	ANN-RS and Estimated Risk of CVDs	165
Figure 6.4	SVM-RS and Estimated Risk of CVDs	165
Figure 6.5	AUC for RFS based <i>NLHRS</i> Model	168
Figure 6.6	AUC for ANN-RS based <i>NLHRS</i> Model	168
Figure 6.7	AUC for SVM-RS based <i>NLHRS</i> Model	168
Figure 6.8	Calibration Plot for Different Types of <i>NLHRS</i> Based Risk Prediction Models	170
Figure 6.9	Distributional Behavior of ANN-RS in Cases and Controls	171
Figure 6.10	Distributional Behavior of SVM-RS in Cases and Controls	172
Figure 6.11	Decision Tree for Identification of Thresholds of ANN-RS in Overall Sample	173
Figure 6.12	Decision Tree for Identification of Thresholds of SVM-RS in Overall Sample	174
Figure 6.13	Thresholds of ANN-RS and SVM-RS and their Associated Risk	175
Figure 6.14	Gender-wise Decision Tree for Thresholds of ANN-RS	177
Figure 6.15	Gender-wise Decision Tree for Thresholds of SVM-RS	178
Figure 6.16	Area as Moderator for ANN-RS and Predicted Logit of CVDs	182
Figure 6.17	Area as Moderator for SVM-RS and Predicted Logit of CVDs	184

LIST OF SYMBOLS

A	Attributes for classification
B	Matrix of regression coefficients of y_i variables on y_i variables where $i \neq i$
C	Set of classes in the data set
C	Correlation Matrix
c	Class
c_t	Cost function in kernels
f	Map from n-dimensional to m-dimensional
F	Weight vector
H	Hosmer and Lemeshow Test
H_0	Null Hypothesis
H_1	Alternative Hypothesis
$H(S)$	Entropy of data set S
$H(t)$	Entropy of subset t
I	Identity matrix
k_r	Kernel
$k - 1$	A method of cross-validation
\log_2	The logarithm of base 2
M	Margin
N	Size of population
N_i	Number of items
n	Sample size
P	Probability for favorable events
p	Proportion
Q	Bias
R^2	Coefficient of determination

S	Data set
S^2	Sample Variance
SE	Standard error
s_p	Split point in decision tree
T	Subsets creating from splitting set S
t	Target class
v_i	Training example
W_O	Moderator
x	First input feature
x'	Second input feature
x_j	Independent variable
y_i	Dependent variable
Z	Z-statistic
z_i	Logit function or logit index
σ	Free parameter in kernel function
Γ	Matrix of regression coefficients of y_i on x_i
ϕ_{ij}	Covariance between exogenous variables
ϵ	Error term in regression analysis
ζ_i	Error vector
β_0	Intercept in regression analysis
φ_l	Regression coefficient of interaction term
α	Level of significance
α_j	Regression coefficient of mediator
ρ	Reliability coefficient
ρ_O	Regression coefficient of moderator
γ_{11}	Regression coefficients of x_j on y_i
β_j	Regression coefficient
β_{ii}	Regression coefficient of y_i with another y_i where $i \neq i$
λ	Eigen value



LIST OF ABBREVIATIONS

AFVIF	Average Full Collinearity Variance Inflation Factor
AMI	Acute Myocardial Infarction
AMOS	Analysis of Moments Structures
ANN	Artificial Neural Network
ANN-RS	Artificial Neural Network based Risk Score
ANOVA	Analysis of Variance
AO	Abdominal Obesity
APC	Average Path Coefficient
ARS	Average R-Square
AUC	Area under the Curve
BMI	Body Mass Index
BS	Brier Score
CB-SEM	Covariance based Structure Equation Modeling
CE	Combined Exposure
CART	Classification Analysis and Regression Trees
CI-TC	Corrected Item Total Correlation
CFA	Confirmatory Factor Analysis
CFI	Comparative Fit Index
CHD	Coronary Heart Disease
CI	Confidence Interval
CVDs	Cardiovascular Diseases
CVM	Cardiovascular Mortality
DM	Diabetes Mellitus
DT	Decision Trees
DV	Dependent Variable
FFQ	Food Frequency Questionnaire

FN	False Negative
FH	Family History
FP	False Positive
GFI's	Goodness of Fit Indices
GI	Gini Index
GoF	Goodness of Fit
HDL	High-density Lipoprotein
HICs	High-income Countries
HTN	Hypertension
ICT	Islamabad Capital Territory (Capital city of Pakistan)
ICR	Internal Consistency Reliability
IDF	International Diabetes Federation
ID3	Iterative Dichotomizer
IV	Independent Variable
KMO	Kaiser-meyer-olkin Test
KS	Kolmogorov Smirnov
LMICs	Low-middle-income countries
LL	Log likelihood
LRA	Logistic Regression Analysis
ML	Machine Learning
MLE	Maximum Likelihood Estimation
MLP	Multilayer Perceptron
MS	Metabolic Syndrome
MSA	Metabolic Syndrome Abnormalities
NDH	Negative Dietary Habits
NCDs	Non-communicable Diseases
NLHRS	Non-laboratory based Heart Risk Score
OLS	Ordinary Least Square
OR	Odds Ratio

PA	Physical Activity
PI	Prognostic Index
PI_A	Prognostic Index based on ANN-RS
PI_S	Prognostic Index based on SVM-RS
PI_R	Prognostic Index based on RFS
PLS	Partial Least Square
PLS-SEM	Partial Least Square based Structure Equation Modeling
PPLUA	Percentage Population Living in Urban Areas
RBF	Radial Basis Function
RFS	Risk Factors based Risk Score
RMSE	Root Mean Square Error
RMSEA	Root Mean Square Error Approximation
ROC	Receiver Operating Characteristics
RR	Risk Ratio
RT	Regression Trees
SDH	Social Determinants of Health
SEM	Structure Equation Modeling
SE	Socio-economic
SFW	Subjective financial Well-being
ST	Sleep Satisfaction
SPSS	Statistical Packages for Social Sciences
SS	Self-reported Subjective Stress
SVM	Support Vector Machine
SVM-RS	Support Vector Machine based Risk Score
TN	True Negative
TP	True Positive
WC	Waist Circumference
WHO	World Health Organization
WHR	Waist to Hip Ratio

LIST OF APPENDICES

APPENDIX A	Disease Terminology	240
APPENDIX B	List of Experts	242
APPENDIX C	Definition and Coding Scheme of Variables Used in Study	243
APPENDIX D	Informed Consent	247
APPENDIX E	Questionnaire	248
APPENDIX F	Ethical Review Committee Approval	254
APPENDIX G	Algorithms Scheme or Configuration for <i>ML</i> Models	256
APPENDIX H	WarpPLS Model Output and Examples of Linear and Nonlinear Relationships	258
APPENDIX I	Weka Output	260
APPENDIX J	Normality Assessment for ANN-RS and SVM-RS	263
APPENDIX K	List of Publications and Certification	265

اونیورسیتی ملیسیا قہق

UNIVERSITI MALAYSIA PAHANG

CHAPTER 1

INTRODUCTION

1.1 Overview

Cardiovascular diseases (CVDs) are a group of diseases, usually referred to as conditions that involve narrowed or blocked blood vessels leading to heart attack and other related problems. Recent estimates of the World Health Organization (WHO) showed that 31% of cause-specific mortality occurs due to this disease group (Organization et al., 2018a). Global cardiovascular mortality (CVM) has established an exponential trend where 40.8% increase was observed from 1990 to 2013 (Roth et al., 2015). This rise has made it the most important and the largest cause of noncommunicable diseases (NCDs) at over 50% (McAloon et al., 2016). However, the nature of disease burden is diverse in different regions of the world. It usually occurs in low-middle-income countries (LMICs), contributing to 80% of the annual deaths (Organization et al., 2019). These statistics of CVDs related deaths reflect the enormity of disease, which is continuously growing. Therefore, it has become a public health challenge, especially for LMICs and needs to be duly addressed.

Global burden of disease (GBD) reported that age-standardized CVDs mortality in high-income countries (HICs) is decreased by 21% in the last two decades (Lozano et al., 2012). The adoption of population and individual-based preventive strategies recommended by WHO is the main reason for this substantial decline in HICs (Bonita et al., 2013). In population-based strategies, the focus is on developing such public health policies, which transform the distribution of primary risk factors of CVDs in the population. However, individual-based strategies tend to focus on the early assessment of the likelihood of CVDs events, identifying high-risk individuals and their management through modifiable risk factors. Besides, it is more suitable in clinical settings as they