



THE UNIVERSITY OF
WAIKATO
Te Whare Wānanga o Waikato

Research Commons

<https://researchcommons.waikato.ac.nz/>

Research Commons at the University of Waikato

Copyright Statement:

The digital copy of this thesis is protected by the Copyright Act 1994 (New Zealand).

The thesis may be consulted by you, provided you comply with the provisions of the Act and the following conditions of use:

- Any use you make of these documents or images must be for research or private study purposes only, and you may not make them available to any other person.
- Authors control the copyright of their thesis. You will recognise the author's right to be identified as the author of the thesis, and due acknowledgement will be made to the author where appropriate.
- You will obtain the author's permission before publishing any material from the thesis.

A Decision Support System for Predicting the Complications of Diabetes Mellitus: A Design Science Research Approach

A thesis
submitted in fulfilment
of the requirements for the degree

of
Doctor of Philosophy in Management Systems

at
The University of Waikato

by
Madurapperumage Anuradha Erandathi



THE UNIVERSITY OF
WAIKATO
Te Whare Wānanga o Waikato

2024

Abstract

Health information systems (HIS) serve as the cornerstone of modern healthcare, seamlessly weaving data into actionable insights and empowering professionals to make informed decisions and elevate patient care. Decision support systems became a prominent research area in the discipline of HIS, where clinical decision support systems (CDSSs) with the ability to diagnose and prognosis of disease are recognised as the most common and beneficial information systems. Moreover, the high prevalence and adverse effects of diabetes all over the globe evidently make it vital to predict diabetes and its complications.

The current study aims to resolve an issue at Te Whatu Ora, by generating a CDSS which can predict the complications of diabetes mellitus (CoDM) while answering the knowledge gap. Although a rich data set of diabetes patients' is maintained at Te Whatu Ora, their involvement in decision-making is unsatisfactory. This study created a CDSS to resolve the issue at Te Whatu Ora while considering two perspectives of the question: design and data analysis. The system design followed design science research methodologies (DSRM) while selecting suitable techniques in the steps of the empirical cycle to confirm their applicability in the domain. The data analysis perspective of the study focused on survival analysis methods due to their appropriateness in fulfilling the identified research gaps.

The created CDSS is the primary outcome of the research, which resolves the real-world issue while addressing the recognised research gaps. The solution's design perspective confirms the applicability of adopting design science research approaches in the context of a systematic solution-design process. The data analytics perspective confirms the appropriateness of survival techniques in the domain while validating the system's performance.

The outcome of this research has significant academical, and managerial implications. The implemented CDSS is capable of providing a chronological risk percentage for 10 CoDM in a cohort of New Zealand. The systematic procedure adopted in the research contributes to the existing knowledge gaps while answering the design, implementation, and evaluation stages. The managerial implications of the study expand through policy-makers, resource-allocators,

and healthcare administrators to doctors, nurses, and patients. The predicted risk of CoDM may be beneficial in managing the patients by issuing early warnings, starting treatment plans, conducting diagnosis tests, recommending dietary/exercise routines and more from the perspective of patient care. The visualisation of the statistical details and the survival curves in the system may assist in managing diabetes data repositories efficiently.

This study creates a cohort-specific risk prediction model based on the New Zealand cohort. The interoperability challenges in the CDSS could occur due to the variety of practices in real life. Future research studies in this domain can concentrate on increasing the accuracy of the models with more features of a rich dataset while protecting the patients' confidentiality. Additionally, the cohort specificity of the CDSS can be avoided with the engagement of a global dataset.

Academic Engagement

- Conference paper – “Predicting Diabetes Mellitus and Its Complications through a Graph-Based Risk Scoring System”. Erandathi, M. A., Wang, W. Y. C., and Mayo, M. 2020. "Predicting Diabetes Mellitus and Its Complications through a Graph-Based Risk Scoring System," in: Proceedings of the 4th International Conference on Medical and Health Informatics. Kamakura City, Japan: Association for Computing Machinery, pp. 1–7.
- Research presentation on Doctoral consortium at University of Waikato – 2020.
- Conference paper – “A Systematic Review on Extracting Predictors for Forecasting Complications of Diabetes Mellitus”. Madurapperumage, A., Wang, W. Y. C., and Michael, M. 2021. "A Systematic Review on Extracting Predictors for Forecasting Complications of Diabetes Mellitus," in: Proceedings of the 5th International Conference on Medical and Health Informatics. Kyoto, Japan: Association for Computing Machinery, pp. 327–330.
- Participation in New Zealand Information Systems Doctoral Consortium – NZISDC (2021).
- Journal article – “Clustering the countries for quantifying the status of Covid-19 through time series analysis”. Erandathi, M., Chung Wang, W.Y. and Hsieh, C.-C. , 2022, "Clustering the countries for quantifying the status of Covid-19 through time series analysis", Information Discovery and Delivery, Vol. 50 No. 3, pp. 297-311. <https://doi.org/10.1108/IDD-03-2021-0034>.
- Conference paper – “Prevalence of sociodemographic factors in a cohort of diabetes mellitus: a retrospective study”. Erandathi, M., Wang, W. Y. C., Mayo, M., and Shafuiu, I. 2022. "Prevalence of Sociodemographic Factors in a Cohort of Diabetes Mellitus: A Retrospective Study," in: Proceedings of the 6th International Conference on Medical and Health Informatics. Virtual Event, Japan: Association for Computing Machinery, pp. 193–197.
- Journal article – “Comprehensive Factors for Predicting the Complications of Diabetes Mellitus: A Systematic Review”. Erandathi, M. A., Wang, W. Y. C., Mayo, M., and Lee, C. C. 2024. "Comprehensive Factors for Predicting the Complications of Diabetes Mellitus: A Systematic Review," *Current diabetes reviews*, 10.2174/0115733998271863231116062601.

Acknowledgements

I wish to convey my profound gratitude to the individuals and organisations instrumental in facilitating this scholarly expedition. I extend sincere appreciation to my esteemed supervisors, namely Dr William Wang, Dr Michael Mayo, and Dr Emily Grout, whose guidance and expertise have significantly shaped the trajectory of this research endeavour. I am indebted to these distinguished academics' invaluable support and mentorship.

My gratitude also extends to the dedicated staff of Te Whatu Ora, particularly Dr. Ibrahim Shafiu and his team, for their substantial support, playing a pivotal role in the success of this research initiative.

I express heartfelt thanks to the AHEAD scholarship programme and Sabaragamuwa University of Sri Lanka for providing an exceptional opportunity. Furthermore, my appreciation encompasses the University of Peradeniya, which has played a foundational role in shaping my academic journey.

I extend my deepest gratitude to my beloved mother for her unwavering support, and I express appreciation to my family for their understanding and encouragement. Special thanks to my husband for his consistent and steadfast support during the most challenging moments.

Finally, I extend my heartfelt gratitude to the universe for orchestrating the circumstances that enabled me to navigate challenges, find the courage to overcome fears, and receive encouragement to pursue and fulfill my dreams.

Table of Contents

Chapter 1	Introduction	1
1.1	Research Motivation	6
1.2	Research Focus and Objectives	10
1.3	Research Methodology.....	13
1.4	Research Impact.....	16
1.5	Structure of the Thesis.....	19
Chapter 2	Reviewing the State of the Art.....	22
2.1	Data Analytics in Healthcare Management	23
2.2	Diabetes Mellitus	27
2.3	Applications of Health Information Systems in T2DM and its Complications.....	28
2.4	State-of-the-Art of the Clinical Decision Support System	34
2.5	Clinical Decision Support Systems in Predicting the Complications of Diabetes.....	37
2.6	Rationale of the Study	39
2.7	The Research Questions	43
Chapter 3	Research Methodology.....	46
3.1	Introduction	46
3.2	The Rationale for the Approach.....	49
3.3	Reviewing the Process of Design Science	53
3.3.1	Six Steps of DSRM	53
3.3.2	Design Science Research Methodology Introduced by Wieringa	55
3.3.2.1	Design Framework.....	56
3.3.2.2	The Engineering Cycle	59
3.3.2.3	Problem Investigation	61
3.3.2.4	Solution Design.....	61
3.3.2.5	Solution Validation	63
3.3.2.6	Solution Implementation.....	64
3.3.2.7	Implementation Evaluation	64
3.4	Identifying Situational Requirements and Artefacts.....	66
3.4.1	State of the Art of CDSS	67
3.4.2	Identifying the Feature Set	69
3.4.3	Artefact of the Study.....	70
3.5	The Design of the Artefact	72
3.5.1	Data Collection.....	72
3.5.1.1	Ethical Consideration.....	74
3.5.2	Data Pre-Processing	75
3.5.3	Exploratory Data Analysis	76
3.5.3.1	Exploring the Socio-demographic Details.....	77
3.5.3.2	Exploring the Dataset with Empirical Standards	79
3.5.4	Feature Selection	81
3.5.5	Model Selection	82
3.5.5.1	Non-Parametric Techniques in Analysis	85
3.5.5.2	Semi-Parametric Techniques in Survival Analysis.....	88

3.5.6	Web Portal Implementation as the Deployment.....	90
3.6	Evaluation of the Artefact.....	93
3.6.1	Evaluation Process of the Study.....	96
3.6.2	Design Evaluation.....	98
3.6.3	Algorithm Evaluation	100
3.6.4	Implementation Evaluation	102
Chapter 4	<i>Situational Awareness.....</i>	107
4.1	Introduction	107
4.2	Identifying the Contextual Requirements	107
4.2.1	The Social Context of the Research Study	108
4.2.2	Identified Stakeholders and their Goals.....	110
4.2.3	Features Selection through a Systematic Review	110
4.2.3.1	Methodology of the systematic review.....	112
4.2.3.2	Risk Factors for Predicting Complications of Diabetes Mellitus	115
4.3	Conceptual Problem Framework	121
4.4	Summary.....	122
Chapter 5	<i>The Design of the Artefact.....</i>	125
5.1	Proposed Artefact of the Study.....	125
5.2	Solution Design	126
5.3	Solution Implementation.....	134
5.3.1	Data Collection.....	136
5.3.2	Data Pre-Processing	138
5.3.3	Exploratory Data Analysis	141
5.3.3.1	Exploring the Prevalence of Factors in a Cohort of Diabetes Mellitus	141
5.3.3.2	Exploring the Data Set with Empirical Standards	142
5.3.4	Feature Selection	145
5.3.5	Model Selection	145
5.3.5.1	Non-Parametric Techniques in Survival Analysis.....	146
5.3.5.2	Semi-Parametric Techniques in Survival Analysis.....	148
5.3.6	Implementation of the Web Portal.....	152
Chapter 6	<i>Evaluation of the Artefact</i>	165
6.1	Design Evaluation	165
6.2	Algorithm Evaluation.....	169
6.3	Implementation Evaluation	170
6.3.1	Evaluation through ISO 25010	170
6.3.2	Evaluation through User Feedback.....	174
6.4	Summary.....	179
Chapter 7	<i>Research Findings and Interpretations.....</i>	181
7.1	RQ1: How can a CDSS be designed by utilising DSRM?	182
7.1.1	RQ1.1: What is the utilized design process applied for creating the CDSS?	182
7.1.1.1	Exploring the Prevalence of Sociodemographic Cohort of Diabetes Mellitus	184
7.1.1.2	Exploratory Data Analysis.....	188
7.1.2	RQ1.2: How Can a Designed CDSS be Evaluated with Existing Standards?.....	207
7.1.2.1	Results of the Evaluation through ISO 25010.....	208
7.1.2.2	The Results of the Evaluation of User Feedback	212

7.2	RQ2: How can the onset of CoDM be predicted using a longitudinal data set?.....	224
7.2.1	RQ2.1: How can the existing techniques for predicting CoDM be utilised in a cohort of New Zealand?	225
7.2.1.1	Non-Parametric Survival Curves.....	225
7.2.1.2	Semi-Parametric Survival Curves.....	233
7.2.2	RQ2.2: How Accurate is the Prediction of CoDM?	242
7.3	Summary.....	245
Chapter 8	<i>Discussion and Recommendation</i>	246
8.1	Research Implications.....	246
8.1.1	Academical implications	247
8.1.2	Managerial Implications.....	251
8.2	Limitations of the Study.....	253
8.3	Recommendations for Future Research	254
References	255	
Appendix A:	<i>Ethical Approval.....</i>	292
Appendix B:	<i>Requirement Specification Document</i>	293
Appendix C:	<i>Software Quality Measurement.....</i>	296
Appendix D:	<i>Questionnaire of Users' Feedback</i>	298

List of Tables

Table 3-1 : Comparison of the stages in design science research methodology.	52
Table 3-2 : Comparison of primary and secondary data collection methods.	73
Table 4-1 : Frequencies of the selected features (Madurapperumage et al., 2021).	117
Table 5-1 : Table of system requirements.	132
Table 5-3 : Data dictionary of collected datasets.	137
Table 5-4 : Table of ethnicities in the dataset and mapped ethnicity types in the study. ...	140
Table 6-1 : Table of user requirement and the expected effects of the artefact.	167
Table 6-2 : System requirements from the ISO 25010 standard.	171
Table 6-3 : Mapping the questions of the questionnaire with standard characteristics ...	177
Table 7-1 : Characteristics of diabetes patients in the Waikato region (N=2656)	185
Table 7-2 : Age distribution of diabetes cohort in identified sociodemographic groups. ...	186
Table 7-3 : Percentage of incident rates in CoDM	198
Table 7-4 : Table of user feedback for the system evaluation.....	222
Table 7-5 : Algorithm evaluation results.....	242
Table 7-6 : Accuracy of created Cox models.....	244

List of Figures

Figure 1-1 : Research focus diagram.	11
Figure 1-2 : Diagram of research questions.....	13
Figure 1-3 : Diagram illustrates the thesis structure	21
Figure 3-1 : Design science framework introduced by Hevner et al (2004).	57
Figure 3-2 : A design science framework introduced by Wieringa (2014).....	58
Figure 3-3 : Engineering cycle introduced by Wieringa (2014).....	60
Figure 3-4 : Diagram of presenting the extraction of research artefact through the research questions.	71
Figure 3-5 : FEDS (Framework for Evaluation in Design Science) (Venable et al., 2012).....	94
Figure 3-6 : Diagram of the evaluation process of the study.	98
Figure 4-1 : Flow chart of article selection of the systematic review.	114
Figure 4-2 : Percentages of selected frequently used features for predicting the complications of diabetes (adapted from	119
Figure 4-3 : Conceptual research framework adopted from the Wieringa (2014).	122
Figure 4-4 : Diagram to illustrate the ultimate research artefact of the study.	124
Figure 5-1 : Snippet of the requirement specification document where the functional requirement are extracted.	128
Figure 5-2 : Snippet of the use cases in the requirement specification document.	128
Figure 5-3 : Use case diagram for the functional requirements of the system.	134
Figure 5-4 : Empirical cycle of DSRM (Wieringa, 2014).	135
Figure 5-5 : Sample code for modelling the Kaplan-Meier in to the cohort of T2DM.	147
Figure 5-6 : Sample code for calling the functions in non-parametric survival analysis for T2DM.	147
Figure 5-7 : Code snippet used to generate the Kaplan Meier models for each complication.	148
Figure 5-8 : Code snippet for using Kaplan Meier model in complications.	148
Figure 5-9 : Code snippet used for creating Cox model for E1122 with demographic details.	149
Figure 5-10 : Saving created data model of E1122.	150
Figure 5-11 : Predicted results of survival of E1122 using Cox models with demographic details.	150
Figure 5-12 : Code snippet used to create the Cox model for E1122 with all the features..	151
Figure 5-13 : Saving the Cox model for E1122 which created for predicting the survival with all the features.	151
Figure 5-14 : Code snippet of predicting the survival rate for E1122 with all the features..	152
Figure 5-15 : Home screen of the NZTPCD - Part1.....	153
Figure 5-16 : Home screen of NZTPCD -Part 2- Visualisation of the distribution of gender and ethnicity.....	154
Figure 5-17 : Home screen of NZTPCD- Part 3- Distribution of Māori/Non-Māori and age on diagnosis.....	154
Figure 5-18 : Home screen of NZTPCD- Part 4 - Survival curve of diabetes.....	155
Figure 5-19 : NZTPCD user interface with the selection of E1122 Complication - Part1.	156
Figure 5-20 : NZTPCD user interface with the selection of E1122 Complication- Part 2.	156
Figure 5-21 : NZTPCD user interface with the selection of E1122 complication- Part 3.	157

Figure 5-22 : NZTPCD user interface with the selection of E1122 - Part 4.	157
Figure 5-23 : NZTPCD user interface for prediction with demographic details.	158
Figure 5-24 : NZTPCD user interface with the prediction of laboratory values.	159
Figure 5-25 : NZTPCD user interface for prediction of the survival rate with demographic details.	160
Figure 5-26 : NZTPCD user interface for the prediction of hazard of the complication using demographic details- Part 1.	160
Figure 5-27 : NZTPCD user interface for prediction of complication using demographic details – Part 2.	161
Figure 5-28 : NZTPCD User Interface for predicting the survival using laboratory values....	162
Figure 5-29 : NZTPCD user interface for the prediction results for laboratory-based survival prediction - Part 1.	163
Figure 5-30 : NZTPCD user interface for prediction results of laboratory-based prediction - Part 2.	163
Figure 6-1 : External characteristics used for evaluating the system with ISO 25010 standard.(ISO/IEC 25010, 2011)	172
Figure 6-2 : Characteristics used for build the questionnaire.	176
Figure 7-1 : Diagram of research questions.	181
Figure 7-2 : Distribution of complications of diabetes mellitus and the ethnicities of the cohort.	187
Figure 7-3 : E119 diagnosis age distribution of patients in the diagnosis table.	190
Figure 7-4 : Dispersion of age at diagnosis of E119 by Māori/non-Māori and gender.	191
Figure 7-5 : Boxplot of age on diagnosis of E119 by gender.	192
Figure 7-6 : Box plot of age on diagnosis of E119 by the characteristic of Māori.	192
Figure 7-7 : Violin plot of patient’s age on diagnosis categorised by Māori/non-Māori.	192
Figure 7-8 : Available records of test results in the test result table.	194
Figure 7-9 : The number of laboratory records collected through the years.	195
Figure 7-10 : Histogram of age at diagnosis of E119.	196
Figure 7-11 : Histogram of time between E119 to any-complication.	196
Figure 7-12 : Histogram of diagnosis year of E119.	197
Figure 7-13 : Histogram of number of complications recorded over time.	197
Figure 7-14 : Box plot of age at diagnosis of E119.	197
Figure 7-15 : Box plot of diagnosis year of complication.	197
Figure 7-16 : Incident rates of each complication of diabetes categorised by ethnicity.	199
Figure 7-17 : Incident rates of each complication categorised according to the characteristic of Māori/non-Māori.	199
Figure 7-18 : Incident rates of each complication of diabetes according to gender.	200
Figure 7-19 : Chart for representing the number of patients in CoDM by ethnicity	201
Figure 7-20 : Chart for representing the number of patients in CoDM by gender	202
Figure 7-21 : Bar charts represent the age at diagnosis of E119 categorised by gender and Māori/non-Māori characteristic.	203
Figure 7-22 : Bar chart of the age at diagnosis of E119 categorised by Māori/non-Māori and gender.	203
Figure 7-23 : Horizontal bar chart of age at diagnosis of E119 categorised by ethnicity.	203
Figure 7-24 : Violin plot of age at diagnosis of E119 categorised by Māori/non-Māori.	204
Figure 7-25 : Violin plot of age at diagnosis of E119 categorised by ethnicity.	204

Figure 7-26 : Chart representing the degree of influence of block of requirements on external characteristics.	209
Figure 7-27 : The degree of influence of the blocks of requirements on the external sub-characteristics.....	210
Figure 7-28 : The degree of influence of the requirements on the external characteristics.	211
Figure 7-29 : Chart represents the age distribution of participants in the questionnaire....	213
Figure 7-30 : Chart represents the gender distribution of the participants in the questionnaire.	214
Figure 7-31 : Chart represents the responses of users for the easiness of learning to use the CDSS.....	214
Figure 7-32 : Chart represents the user friendliness of the interfaces of CDSS.....	215
Figure 7-33 : Chart illustrates the responses for the ability to complete the tasks quickly using the CDSS.....	215
Figure 7-34 : Chart represents the received responses for the having clear icons and menu options in the interface of CDSS.....	216
Figure 7-35 : Chart illustrates the responses for the recall the steps to perform the tasks in CDSS.....	216
Figure 7-36 : Chart illustrates the satisfaction of the usability of the CDSS.	217
Figure 7-37 : Chart illustrates the responses whether the CDSS meets the users’ needs and expectations.	217
Figure 7-38 : Chart represents the responses of accessibility and usability of CDSS with assistive technologies.	218
Figure 7-39 : Chart for the responses for encountering accessibility challenges or limitations of the CDSS.	218
Figure 7-40 : Chart represents the responses for the visual appealing of the CDSS.	219
Figure 7-41 : Chart illustrates the responses for pleasant user experience of the CDSS.....	219
Figure 7-42 : Charts represents the responses of having well organised and intuitive layouts in the CDSS.	220
Figure 7-43 : Chart illustrates the responses of having consistent design elements and navigation in CDSS.....	220
Figure 7-44 : Chart illustrates the responses of having a good response time from the system.	221
Figure 7-45 : Chart illustrates the responses of the reliability of the CDSS.....	221
Figure 7-46 : Kaplan-Meier curve of E119 cohort by gender.....	226
Figure 7-47 : Kaplan-Meier curve of E119 cohort by Māori/non-Māori.....	226
Figure 7-48 : Kaplan-Meier curve of E119 cohort by age categories.....	226
Figure 7-49 : Kaplan-Meier curve of E119 cohort by ethnicities.	226
Figure 7-50 : Survival curve of E1122.	227
Figure 7-51 : Survival curve of E1129.	227
Figure 7-52 : Survival curve of E1131.	227
Figure 7-53 : Survival curve of E1139.	227
Figure 7-54 : Survival curve of E1142.	227
Figure 7-55 : Survival curve of E1151.	227
Figure 7-56 : Survival curve of E1164.	228
Figure 7-57 : Survival curve of E1165.	228
Figure 7-58 : Survival curve of E1171.	228

Figure 7-59 : Survival curve of E1172.	228
Figure 7-60 : Survival curve of E1122 by gender.	229
Figure 7-61 : Survival curve of E1122 by the age of diagnosis.....	230
Figure 7-62 : Survival curve of E1122 by ethnic groups.....	231
Figure 7-63 : Survival curve of E1122 by Māori/non-Māori.	232
Figure 7-64 : Demographic details forum filled for predict the patient's survival rates.....	234
Figure 7-65 : Survival curve of E1121 results from Cox model using demographic details..	235
Figure 7-66 : Survival curve of E1129 results from Cox model using demographic details..	235
Figure 7-67 : Survival curve of E1131 results from Cox model using demographic details..	235
Figure 7-68 : Survival curve of E1139 results from Cox model using demographic details..	235
Figure 7-69 : Survival curve of E1142 results from Cox model using demographic details..	236
Figure 7-70 : Survival curve of E1151 results from Cox model using demographic details..	236
Figure 7-71 : Survival curve of E1164 results from Cox model using demographic details..	236
Figure 7-72 : Survival curve of E1165 results from Cox model using demographic details..	236
Figure 7-73 : Survival curve of E1171 results from Cox model using demographic details..	237
Figure 7-74 : Survival curve of E1172 results from Cox model using demographic details..	237
Figure 7-75 : Forum for predict the survival of a patient with all features.	238
Figure 7-76 : Survival curve of E1122 results from Cox model using all features.....	239
Figure 7-77 : Survival curve of E1129 results from Cox model using all features.....	239
Figure 7-78 : Survival curve of E1131 results from Cox model using all features.....	239
Figure 7-79 : Survival curve of E1139 results from Cox model using all features.....	239
Figure 7-80 : Survival curve of E1142 results from Cox model using all features.....	240
Figure 7-81 : Survival curve of E1151 results from Cox model using all features.....	240
Figure 7-82 : Survival curve of E1164 results from Cox model using all features.....	240
Figure 7-83 : Survival curve of E1165 results from Cox model using all features.....	240
Figure 7-84 : Survival curve of E1171 results from Cox model using all features.....	241
Figure 7-85 : Survival curve of E1172 results from Cox model using all features.....	241
Figure 7-86 : Visualising algorithm evaluation results.....	243

List of Abbreviations

AI	Artificial Intelligence
ARIC	Atherosclerosis Risk in Communities study
AUSDRISK	The Australian Type 2 Diabetes Risk Assessment Tool
BMI	Body Mass Index
CDSS	Clinical Decision Support System
CoDM	Complications of Diabetes Mellitus
CVD	Cardio Vascular Disease
DA	Data Analytics
DCCT/EDIC	Diabetes Control and Complications Trial/ Epidemiology of Diabetes Interventions and Complications Study
DCSI	Diabetes Complications Severity Index
DM	Diabetes Mellitus
DNeph	Diabetic Nephropathy
DNeu	Diabetic Neuropathy
DR	Diabetic Retinopathy
DSR	Design Science Research
DSRM	Design Science Research Methodology
E119	Type 2 Diabetes with No Complications.
E1122	Diabetic Nephropathy
E1129	Kidney Complications AKI
E1131	Background Retinopathy
E1139	Other Ophthalmic Complications
E1142	Diabetic Polyneuropathy
E1151	PVD
E1164	Hypoglycaemia
E1165	Poor Control – Hyperglycaemia
E1171	Microvascular and Other Specified Nonvascular Complications
E1172	Fatty Liver
EDA	Exploratory Data Analysis.
eGFR	estimated Glomerular Filtration Rate
EHR	Electronic Health Records
ESRD	End Stage Renal Disease
FEDS	Framework for Evaluation in Design Science
FINRISK	Finland Cardiovascular Risk Study
HbA1c	Glycated Haemoglobin
HCI	Human-Computer Interaction
HDSS	Health Decision Support System
HDL	High-Density Lipoprotein
HREC	Human Research Ethics Committee
IDF	International Diabetes Federation
ISO 25010	Systems and Software Quality Requirements and Evaluation (SQuaRE) by International Organization for Standardization.
KB-CDSS	Knowledge Based Clinical Decision Support System
LEA	Lower Extremity Amputation

LDL	Low Density Lipoprotein
LoV	Loss of Vision
MDSS	Medical Decision Support System
ML	Machine Learning
NKB-CDSS	Non-Knowledge Based Clinical Decision Support System
PRISMA	Preferred Reporting Items for Systematic Reviews and Meta-Analyses.
QDiabetes	Risk Prediction Algorithm Created Based on UK.
QRISK	Prediction Algorithm for CVD
SAT	Self-Assessment Tools
SVM	Support Vector Method
T1DM	Type 1 Diabetes Mellitus
T2DM	Type 2 Diabetes Mellitus
UEMs	Usability Evaluation Methods
USCDI	United States Core Data for Interoperability

Chapter 1 Introduction

The emerging global issue of diabetes mellitus and its complications has become a significant focus in health informatics, encompassing areas such as disease diagnosis, prognosis, causal analysis, and identification of risk factors (Sugandh et al., 2023). The development of computational techniques has greatly impacted health informatics, offering promising improvements in several areas due to both health benefits and potential cost savings (Elhoseny et al., 2019). Among these, the prognosis of diabetes complications through non-invasive computational methods has seen substantial evolution over the past two decades (Chawla et al., 2016). Despite continuous research on the prognosis of diabetes and its complications, the scarcity of contextual data and limitations of existing models underscore the critical need for a computerized system tailored for predicting diabetes complications within a New Zealand cohort (Robinson et al., 2012). This real-world issue highlights the essential requirement to effectively manage diabetes data to extract insights and facilitate prognosis of its complications, thereby confirming the contemporary necessity of developing an information system to support decision-making in the complications of diabetes mellitus. The potential benefits of such a decision support system in healthcare settings, particularly in forecasting diabetes complications and predicting survival rates, span a range of services. These include alleviating the health burden on individuals and communities, enhancing health indices, and achieving economic efficiencies through reduced reliance on expert opinions, costly tests, and treatments (Zaman et al., 2023; Zikos & DeLellis, 2018). Addressing this research gap by establishing a platform for decision support in managing diabetes cohorts in healthcare settings is recognized as a pivotal endeavour with profound implications for healthcare enhancement.

The immense growth of technology has influenced all aspects of life, resulting in a gigantic digital universe growing exponentially. The intensity of this issue has increased with the evolution of health informatics, which is “the development and assessment of methods and systems for the acquisition, processing, and interpretation of patient data with the help of knowledge from scientific research” (Imhoff, 2002, p. 179). Further, health informatics is an analytical study of monitoring patients, selecting treatments, creating clinical knowledge, and

maintaining healthcare organisations (Coiera, 2015). However, the evolution of health informatics exacerbates the extensive growth of the digital universe. For instance, the development of the Genomes project has drastically increased the size of the digital world by depositing gene data in NCBI's (National Center for Biotechnology Information) GenBank during the first six months, which is nearly two times that of all the existing records in the last thirty 30 years (Stein, 2010). Health informatics has evolved to utilise these data repositories to improve the management of healthcare organisations, patient care operations, and the tasks-arrangement of caregivers. Since health informatics applications have implications for enhancing the quality of life and reducing economic burden, health information systems such as decision-making systems, healthcare management systems, self-assessment tools (SATs), risk prediction systems, and minimal or non-invasive disease diagnosis systems have widely in demand.

“A health information system (HIS) is an information system for processing data, information, and knowledge in health care environments. It can be defined as an integrated effort to collect, process, report, and use health information and knowledge to influence policy-making, program action, and research” (Jakovljević, 2008, p. 603).

Due to the increasing expenditures of the health sector and the high prevalence of diseases around the globe, diagnosing diseases without laboratory tests and estimating the current and future risk of individuals for chronic diseases have become fertile research areas (Thanga Selvi & Muthulakshmi, 2020). The ability to extract invaluable knowledge from health data repositories has led the healthcare sector to store the data in electronic formats. With the tremendous growth of the volume of electronic health records (EHRs), researchers tend to create information systems by utilising health data repositories to achieve multiple goals. In healthcare, where EHRs are prolific, researchers and healthcare institutions have harnessed the potential of health data repositories to develop sophisticated health information systems, with a significant emphasis on clinical decision support systems (CDSSs).

“Clinical decision support system (CDSS) is an interactive software which is developed on the basis of expert systems in order to assist and support the decision-making of physicians, health-care staff, and other personnel involved in broader domains of health-care systems” (Shahsavarani et al., 2015, p. 300).

CDSSs have evolved significantly over the years, reshaping the landscape of healthcare delivery and patient care. The inception of CDSSs can be traced back to early rule-based expert systems in the 1970s, designed to assist healthcare providers in making clinical decisions. These early systems relied on explicit knowledge and predefined rules to offer recommendations to clinicians. However, as healthcare data became more complex and voluminous, the limitations of rule-based CDSSs became evident. The evolution of CDSSs moved towards more sophisticated approaches, integrating data-driven techniques, statistical models, and machine learning (ML) algorithms. These modern CDSSs can analyse vast patient data, extract valuable insights, and provide personalised recommendations. They consider clinical data, patient history, genetics, and even social determinants of health to offer a holistic view of patient care.

CDSSs now seamlessly integrate with EHRs, making real-time patient data accessible to clinicians assisting in diagnosis, treatment planning, and medication management. Furthermore, the advent of mobile health technologies and wearable devices has extended the reach of CDSSs beyond clinical settings, enabling patients to take a more active role in managing their health. In recent years, CDSSs have also ventured into predictive analytics and population health management, aiding healthcare providers in identifying at-risk populations and preventive care strategies. As the healthcare industry continues to advance, the future of CDSSs may involve further integration with artificial intelligence (AI), natural language processing, and big data analytics, offering even more sophisticated decision support to improve patient outcomes and streamline clinical workflows. According to Wiens and Shenoy (2018), the applications of ML techniques can cover a range of clinical tasks, from identification/diagnosis to prediction. Primarily, ML techniques are used for risk stratification, identifying risk factors, understanding pathogen-host interaction, and predicting the

emergence and spread of infectious diseases. The field of health informatics has been evolving due to these massive EHR repositories and the evolution of ML techniques and AI techniques. Although innovative systems have resulted from this evolution, there is no consensus on methods for analysing, extracting patterns, and visualising the resulting enormous quantity of data from an overabundance of digital machines and sensors. Because of this, researchers usually compare the performances of various ML techniques to select the most appropriate one. Combining ML with AI techniques reveals a new dimension where expert knowledge can assist in making health decisions. Although a vast range of possibilities is revealed with ML and AI techniques, statistical modelling provides a stable and straightforward solution that is more reliable in clinical decision-making. This dynamic evolution underscores the growing importance of CDSSs in enhancing the quality and efficiency of healthcare delivery.

A prominent categorisation of CDSSs is knowledge-based and non-knowledge-based CDSSs (Alther & Reddy, 2015; Greenes, 2014) . The knowledge-based CDSS (KB-CDSS) consists of a built-in reference table comprising details such as patients' data, treatment data, diagnosis details, etc. The non-knowledge-based CDSS (NKB-CDSS) learns from past experiences to make decisions and implements these lessons into its knowledge base (Alther & Reddy, 2015). Knowledge-based systems have been introduced to fulfil the requirement of adding expert knowledge for medical purposes.

“Knowledge-based system is a computer-based system, which uses and generates knowledge from data, information and knowledge. These systems are capable of understanding the information under process. They can make decisions based on the residing information/knowledge in the system. In contrast, the traditional computer systems do not know or understand the data/information they process” (Sajja & Akerkar, 2010, p. 3).

Further, “since the expert systems gather its knowledge from several medical specialists, the system has a broader scope and can be more helpful to the patients in comparison to just one

physician” (Zeki et al., 2012, p. 1). Additionally, the inherent qualities of computerised decision-making systems, such as providing solutions heuristically, analysing problems based on previous cases, rules or models, and capability of assembling the techniques in inference engines, made them suitable for engaging with healthcare management. Applications of computerised decision-making systems in this field rigorously focus on diminishing health expenditures while authorising patients for an effective management process. Although computerised decision-making systems have become prominent in the field, due to inherent qualities of KB-CDSSs, such as the level of transparency and the ability to explain the results, and appropriateness of used inference engines and their capabilities, utilised knowledge resources lead to a controversial situation of using pure KB-CDSSs in the healthcare sector. Additionally, the evolution of novice statistical, ML and AI methods in building CDSSs leans more towards the NKB-CDSSs in the healthcare sector.

Due to the financial strains and the high prevalence rate of diabetes mellitus (DM) all over the globe, a significant number of research studies have been conducted on various aspects of managing diabetes patients. Diabetes Mellitus (DM) is a life-threatening chronic ailment that impacted 463 million individuals worldwide in 2019, with projections anticipating a rise to 578 million people by the year 2030 (Saeedi et al., 2019). Diabetes occurs either when the pancreas does not produce enough insulin or the body is resistant to the insulin produced by the pancreas; the former condition, named type 1 diabetes (T1DM), is prevalent among 10% of patients, and the latter, type 2 diabetes (T2DM), covers 90% of patients (WHO, 2016) . Additionally, impaired glucose tolerance and gestational diabetes are the other two common types of DM. The burden of DM has further soared due to the fatal outcomes of its complications. Loss of vision (LoV), end-stage renal diseases (ESRD), cardiovascular disease (CVD), and lower extremity amputations (LEA) are the most frequent complications of diabetes mellitus (CoDM), all of which directly affect the quality of life. “80% of ESRD are caused by diabetes, hypertension or a combination of the two” (WHO, 2016, p. 30). Further, patients with DM have two or three times higher rates of CVD than patients without diabetes (Emerging Risk Factors Collaboration, 2010). According to a study conducted in 2010, 1.9% of visual impairment and 2.6% of blindness are caused by diabetic retinopathy (Bourne et al., 2013). Further, LEA is 10 to 20 times higher among patients with DM (WHO, 2016). Besides

the fatality of DM, its economic impact becomes intolerable. DM's direct global annual cost is more than US\$825 billion (Smolen et al., 2016). Due to the widespread prevalence of DM and the increase in per capita spending on diabetes, the International Diabetes Federation (IDF) has estimated global healthcare expenditures on diabetes more than tripled from 2003 to 2013 (International Diabetes Federation, 2013) . The spending on diabetes will continue to increase in low-income countries compared to high-income countries, which would be significantly impacted by their national gross domestic product (GDP). “Further, this burden increases because of the indirect costs associated with productivity loss, premature mortality and the negative impact of diabetes on nations’ GDP” (WHO, 2016, p. 14). Along with the financial demand in national healthcare expenditures, there is a clear distinction between the catastrophic medical spending of people with and those without diabetes. This difference is apparent in low-income countries (Smith-Spangler et al., 2012).

Furthermore, people with diabetes are vulnerable to viruses because it is harder to treat viruses in the presence of diabetes complications and where blood glucose levels fluctuate (International Diabetes Federation, 2020) . The recent pandemic status in the world is due to a virus disease (COVID-19), an infectious disease caused by a newly discovered coronavirus (WHO, 2020) . According to a study conducted with 52 intensive care patients infected with COVID-19 in China, 22% of non-survivors were diabetes patients (Yang et al., 2020). Another study found that 12% of patients admitted to hospitals with COVID-19 were diabetes patients in Wuhan, China (Zhang et al., 2020). The severity of DM has thrived through this outbreak. Due to this immense financial burden, high prevalence rate of fatal CoDM, and adverse effect of DM and its complications healthcare providers, such as diabetes centres, hospitals, and specialist clinics, tend to engage in finding solutions for preventing DM and its complications.

1.1 Research Motivation

This research is motivated by the pressing necessity to confront the challenges within healthcare settings posed by the chronic disease of diabetes mellitus (DM). A system for assisting the management of diabetes data in decision-making can be beneficial to the stakeholders in healthcare, including policymakers, resource allocators, healthcare administrators, doctors, nurses, and general practitioners. Additionally, the outcomes of a

decision support system in managing diabetes data, expands through a vast spectrum of functionalities from resource allocation, and information delivery for policy making, to assist on managing/ preventing of diseases, monitoring individual health, issuing early warnings, suggesting diagnosis test/treatments/dietary requirements/exercise routines and many more. Diabetes has reached epidemic proportions globally, with millions of individuals affected. Complications arising from diabetes, such as CVD, diabetic retinopathy (DR), diabetic neuropathy (DNeu), and diabetic nephropathy (DNeph), pose substantial health risks and economic burdens. Early and accurate prediction of these complications is critical for timely intervention, assist on policy making, resource allocation, personalised patient care, and improved patient outcomes. Therefore, having a decision-making system for managing the diabetes datasets is became evident.

In addition, this research seeks to harness the power of advanced data analytics, statistical modelling, and predictive algorithms to enable healthcare providers to identify the risk of having CoDM of type 2 diabetes patients. By integrating vast datasets from EHRs, laboratory results, and patient demographics, this study aims to create a comprehensive risk assessment tool as a CDSS. Its successful implementation can potentially revolutionise diabetes care by enhancing preventive strategies and the management of complications, thus reducing healthcare costs and improving the quality of life for individuals with diabetes.

Although several risk-scoring systems are created to estimate and predict DM and CoDM, they all comprise their advantages and disadvantages, and some limitations of the existing risk-scoring models need to be addressed. Many methods in the diagnosis and prognosis of diseases result in a binary outcome, such as positive or negative (Elhoseny et al., 2019). This might be less useful in predicting CoDM due to its irreversibility and fatality. In contrast, predicting risk as a percentage is more informative and convenient. Most of the existing CDSSs forecast the risk in the short, medium, and long-term periods, which will deliver limited benefits compared to the prediction of chronological risk (Aminian et al., 2020). Although various applications are available in disease prognosis, creating a risk-scoring system accommodating wide acceptability has always been challenging. Even though current risk scoring systems can achieve high accuracy for similar cohorts, the performance is

controversial with other datasets. Since the feature values and feature importance vary according to the ethnicities (Schwarz et al., 2009), thorough consideration is required to create a widely accepted system that can dynamically calculate the risk by contemplating specific characteristics of each cohort. This is primarily due to the prevalence of risk-scoring systems that produce results based on ethnic backgrounds (Aekplakorn et al., 2006; Griffin et al., 2000; Lei Chen et al., 2019; Lindström & Tuomilehto, 2003; Schmidt.M.I. et al., 2005; Schulze.M.B. et al., 2007). According to Noble et al. (2011), the risk scoring systems are selection-biased due to the datasets considered to create the models. Most existing systems use datasets that are previously assembled or cohort-specific. Since the accuracy of the risk-scoring systems is highly dependent on the feature values of the cohorts (Noble et al., 2011), it is crucial to consider specific feature values for each cohort. Another shortcoming the scholars have less attended is calculating risk scores with a static dataset. Because of the static nature of the dataset, the validity of risk scoring systems gets worse since the assessments are conducted by an old dataset where the fundamental characteristics of the cohort are different from the details used in the risk score. "Data used to create the CVD risk score are old and might not accurately predict the cardiovascular risk of contemporary hypertensive patients that benefit from more modern treatments and management" (Prieto-Merino et al., 2013, p. 492). To cope with this, one strategy used by the scholars is introducing the extended versions of the risk-scoring systems by adding new features or externally validating the risk-scoring models. Lindström and Tuomilehto (2003) created a FINDRISC risk-scoring system using a Finnish cohort; Gomez-Arbelaez et al. (2015) used the FINDRISC score to predict the risk of a Colombian population; Štiglic et al. (2016) used the same scoring system for the Slovenian working population; and Kulkarni et al. (2017) used FINDRISC for identifying the high-risk individuals among white and black ARIC study participants. In annual updates, the QDiabetes scoring system invented by Hippisley-Cox et al. (2009), adds new potential factors for enhancing accuracy (Hippisley-Cox et al., 2017). The system's accuracy can be negatively affected since the dataset's characteristics have remained unchanged while the cohort is evolving. Most of the existing systems introduced extended versions of the system to overcome the issues.

Moreover, visualising health data conveniently plays a significant role in assisting the decision-making in health care management. Various tasks can get the assistance of representing public health data, such as tracking the distribution of data geographically, analysing the prevalence of diseases, predicting outbreaks, and identifying high-risk populations (Oluwakemi & Kamran, 2016). Further, visualisation methods can assist different stakeholders, such as clinicians, doctors, policymakers, patients, and other healthcare providers, in visualising the health data conveniently while serving a significant role in self-management and decision-making. The existing risk-scoring systems for estimating and predicting the risks of CoDM have several shortcomings that need to be addressed. It is also vital to consider the practical use of risk-scoring systems in enhancing healthcare management by creating systems for assisting clinicians, health authorities, and patients. Additionally, data visualisation techniques can be used to strengthen health management tools.

A range of tasks in healthcare management, including decision-making, policy-making, examining therapeutic effects, risk prediction, issuing early warnings, and self-monitoring systems, are achievable through research solutions in practical operations. However, the solutions derived from academic investigations and their pragmatic applications often reside at opposite poles within the spectrum, yielding diminished practical advantages. Scholars emphasised the importance of including practical considerations in research design paradigms (Druckman, 2005). The novice healthcare management systems can reduce the economic burden by diagnosing and predicting diseases through analysing clinical data with minimal involvement of a physician, making effective decisions using extracted knowledge from stored knowledge bases, estimating the risk of diseases and predicting the occurrence of fatal conditions while enhancing the awareness of self-health needs among patients, provide clinical support decisions, and reveal the severity in the future through predictions. Consequently, healthcare management authorities tend to improve self-management practices among patients to empower them to utilise health resources. “The growing worldwide burden of chronic illness and the imperative to reduce healthcare costs will continue to create demand for technology-based self-management interventions.” (Knight & Shea, 2014, p. 95). The use of electronic health tools for this purpose is trending due to the

frequent use of high-tech devices and digital biomarkers. Using ML techniques on health data repositories to enhance healthcare management through sophisticated technologies is vital to reducing the burden of chronic diseases worldwide.

The pragmatic advancements of the CDSS are highly beneficial in healthcare enhancements, while the theoretical advancements improve academic knowledge in designing, implementing, and validating the CDSS. The existing systems adopt various theoretical approaches for the process of system implementation. Design science research methodologies (DSRM) (Gregorio et al., 2021; Kempainen et al., 2017; Mombini et al., 2020; Ulapane et al., 2023), human-computer interactions (HCI) (Horsky et al., 2013; Vitabile et al., 2019), evidence-based medicine (Gholamzadeh et al., 2023; Sim et al., 2001; Zheng, 2007), and information processing theory (Kilsdonk et al., 2016; O'Neill et al., 2005) are some popular methodologies used in the implementation of CDSSs. Although the system implementation is enriched with various methodologies, selecting the most appropriate method for the problem in hand is challenging. The system implementation processes adopted in the same methodological frameworks are also highly diverse based on factors such as the nature of the dataset, the purpose of the system, stakeholders of the system, functional and non-functional requirements, available technological aspects, and many more. The scarcity of a straightforward system implementation process leads the academicians to a prominent research area of articulating a systematic approach which can be used in similar instants. Moreover, the system evaluation phase is recognised as another milestone of the system implementation. The methods of evaluating the implemented system can be adopted by considering the factors of client requirements, empirical evaluation methods, algorithmic evaluations and more. The knowledge gap in applying different theories and methodologies in a suitable manner to articulate a scientifically sound system implementation process urges the study's motivation.

1.2 Research Focus and Objectives

The conspicuous research gap made a contemporary requirement of implementing a CDSS to predict the risk of CoDM. In this spirit, the fundamentals of this study have been selected to focus on a combination of three components: system design, data analytics, and the

pragmatic discipline of medical science. The focus of this research is illustrated in Figure 1.1. The research study diversified into two primary disciplines, namely design science and data analytics, as it concurrently applied these methodologies within medicine. Adopting a retrospective case study approach, this study refined its focus within design science and statistical modelling while applying it to chronic DM. The study is planned to use a chronological dataset of diabetes patients collected in a healthcare setting. The retrospective study approach has been naturally selected for the study due to the limitations of available time and resources. Weiringa’s approach to DSRM (Wieringa, 2014) is adopted here to analyse the data using survival analysis techniques to predict the risk of a selected set of CoDM.

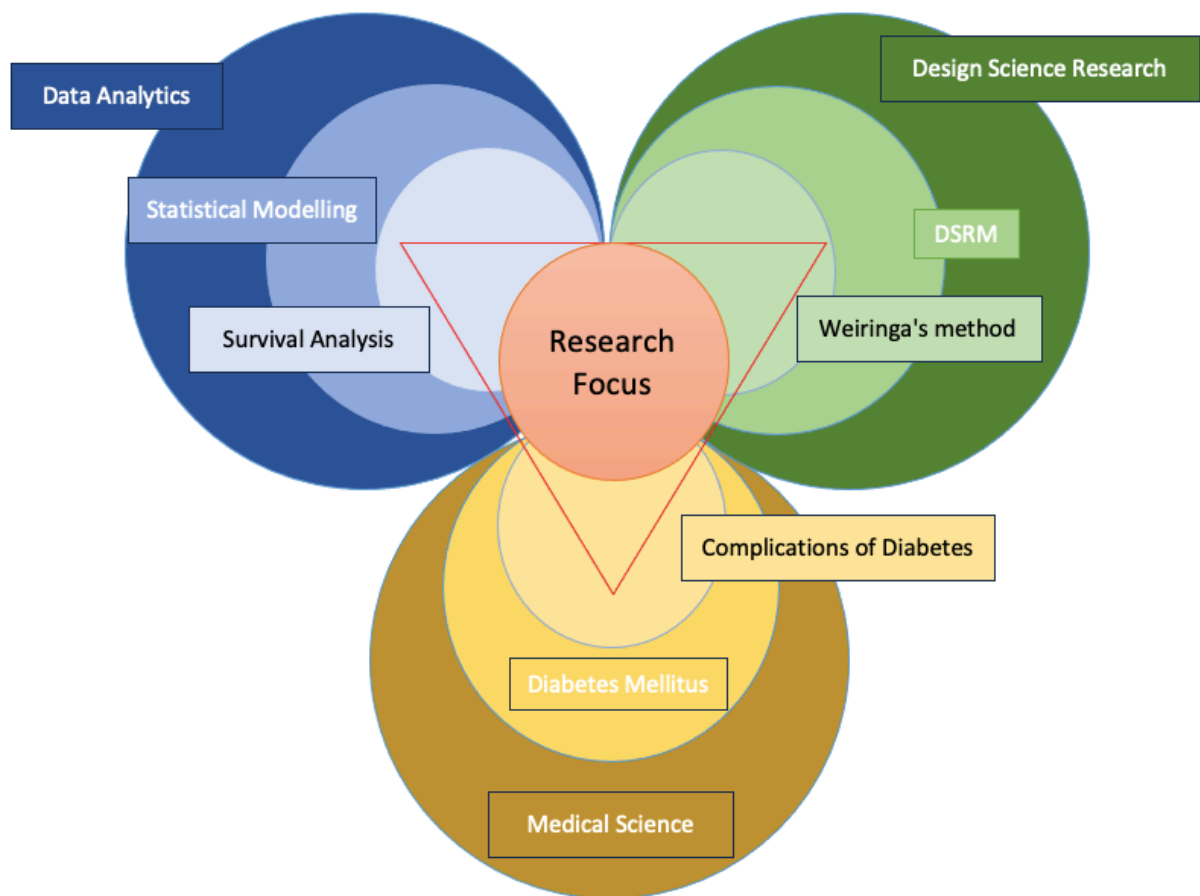


Figure 1-1 : Research focus diagram.

The risk-scoring method of this study is able to estimate and predict the survival rate of having CoDM chronologically for a decade from the diagnosis of diabetes. The ultimate research

focus is resolving a real-world issue by being aware of the situational status through a literature review and a thorough analysis of the real-world requirements. Since the literature review found a significant research gap in the field, the study focuses on filling the gaps in the literature through a practical solution. The contextual gap in the literature review motivates the research to find a local client to get the resources. The combination of extracted gaps in the literature and the real-world client's requirements made a solid research context. Ultimately the research has been focused on resolving a real-world issue in a healthcare sector through the aid of design science perspective and statistical modelling. This has been led to generate a CDSS, to fulfil the existing gaps while providing an innovative system for answering the real-world problem. The research focus has been selected to answer through two fundamentals: the design process and data analysis.

Following the aspects of the research focus, this study comprises with two major objectives:

1. Outcome a solid design process for developing a CDSS using DSRM.
2. Development of survival analysis model for CoDM.

The research questions arise from these two objectives and the relations to research objectives are further explained in section 2.7.

1. RQ1: How can a CDSS be designed by utilising DSRM?

RQ1.1: What is the utilised design process applied for creating the CDSS?

RQ1.2: How can a designed CDSS be evaluated with existing standards?

2. How can the onset of CoDM be predicted using a longitudinal data set?

RQ2.1: How can the existing techniques for predicting CoDM be utilised in a cohort of New Zealand?

RQ2.2: How accurate is the prediction of CoDM?

The research questions and their sub-research questions are illustrated in the figure 1.2.

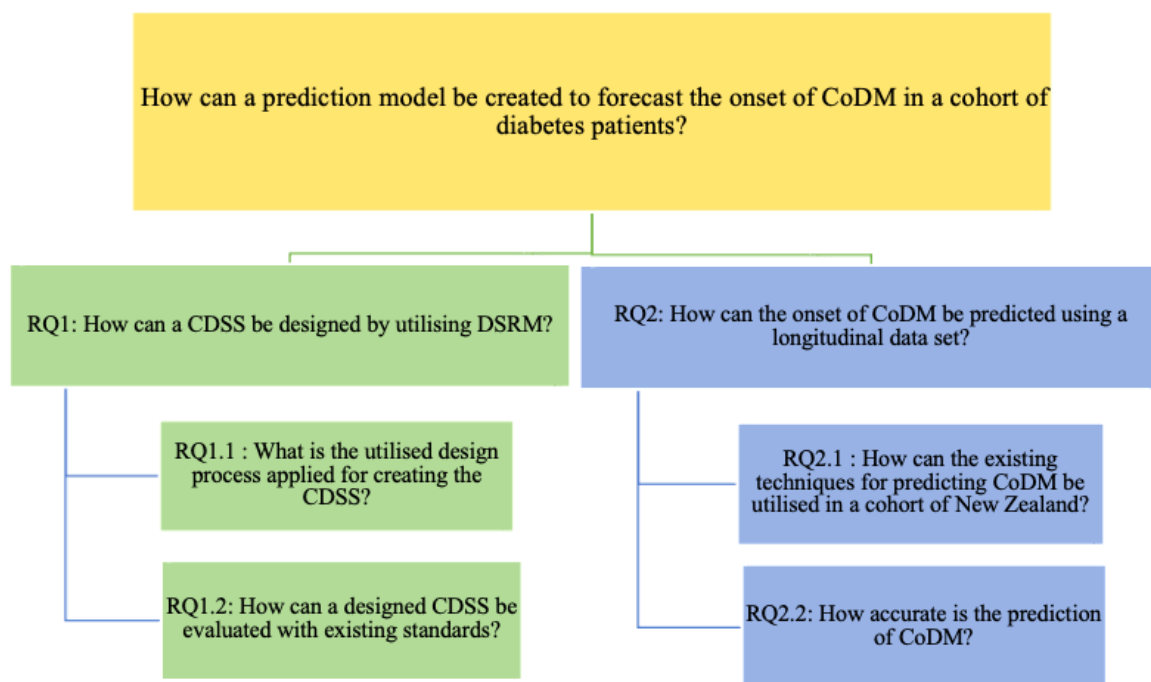


Figure 1-2 : Diagram of research questions

The research questions concentrate on developing a CDSS which can predict the risk of the onset of a selected set of CoDM while demonstrating a systematised process. Research question one focuses on the design perspective of the CDSS. The fundamental stages of the DSRM and the techniques that can be utilised in the healthcare domain are adopted in this section. The applicability and utilisation of the adopted concepts in the design process are confirmed by answering the first research question. The data analytics perspective of the CDSS is discussed in the second research question. The suitability of survival analysis techniques in the hand-in issue of healthcare sector and the accuracy of the designed models of the system are confirmed by the answers to the second research question. The ultimate focus of the research is to resolve a real-world issue with a CDSS which can predict CoDM while answering the recognised research gaps.

1.3 Research Methodology

This research study adopts a confirmatory research approach, synthesising design science research principles and data analytics principles to affirm their validity and applicability in addressing a real-world issue within the healthcare sector. The identified research gap is

addressed by developing a CDSS meticulously crafted following the tenets of DSRM. The CDSS incorporates survival analysis techniques derived from data analytics (DA) to manipulate and analyse the data, thereby facilitating the prediction of survival rates.

In pursuing DSRM, this study has adopted a specific method introduced by Wieringa (2014). The research framework is designed by incorporating key concepts from Weiringa's method, which is structured into three tiers: knowledge context, social context, and design science research. The knowledge context within this study encompasses a rich tapestry of theoretical foundations, including DSRM, software development theory, prediction algorithms, frameworks for design science, software development, data analytics techniques, software engineering techniques, and the formalism of data science. Although the knowledge context of the current study is enriched with the above-mentioned theories and concepts, the core of the study relies on the theories of DSRM and survival techniques. The solution's design process is followed by the DSRM to confirm its theoretical and pragmatic validity. Moreover, the DSRM approach in the design process of the solution explains the adopted methods and theories in design science research. In addition, the theories of survival analysis are confirmed through the study's findings. The social context comprises diverse stakeholders, including Te Whatu Ora as an organisation, decision-makers, policy-makers, general practitioners (GPs), nurses, diabetic patients, and the associated data infrastructures and prediction applications as adopted technologies. Within the design science component of the adopted framework, two primary focuses emerge: design improvement and knowledge enhancement. Design improvement is achieved by incorporating survival techniques and design research principles into implementing the CDSS. Knowledge enhancement is realised through the formalisation of the CDSS implementation process to address real-world issues using DSRM.

DSRM comprises five fundamental steps: Problem Investigation, Treatment Design, Treatment Validation, Treatment Implementation, and Implementation Evaluation. The stages of DSRM are broadly categorised into three main stages—**Designing, Implementation, and Evaluation**—in this study, considering their similarities and relevance in explaining the process. The design phase encompasses situational awareness and solution design, involving researching the existing issue, collecting resources and requirements, analysing

requirements, conducting a feasibility study, finalising solution functions, and selecting the best design. The chosen solution has been implemented using software development and data analysis techniques in the implementation stage. Further, the design implementation employed the Python programming language and the Streamlit application to create the CDSS as a web tool. Survival analysis techniques, specifically the Kaplan-Meier and Cox regression models, are applied to model the dataset for two distinct factor sets. The survival prediction models are implemented using Python. Moreover, when considering the formalisation of data science in modelling datasets, the data collection, pre-processing, exploratory data analysis, feature selection, and model selection are conducted to achieve the purposes of the study. The methodologies used, and their appropriateness in the context is explained in Chapter 3. The web portal is implemented to present the CDSS with a graphical user interface that the stakeholders can use.

The adopted methodology for the evaluation phase of the study comprises three main evaluation stages: design evaluation, algorithm evaluation, and implementation evaluation. Design evaluation ensures the validity of the proposed solution in the problem context. The design evaluation adopted a method , consisting of four components: expected effects, expected value, trade-offs, and sensitivity (Wieringa & Morali, 2012b). The algorithms developed for predicting the survival of each complication are evaluated through ten-fold cross-validation with a C-index and Brier score. Implementation evaluation is conducted as a summative assessment of the system's usability. Implementation evaluation is approached through two methods: ISO standards and user feedback. ISO-25010 standard measures eight characteristics: functional suitability, performance efficiency, compatibility, usability, reliability, security, maintainability, and portability. The evaluation process of the study conducted by Kadi et al. (2016) is adopted in the evaluation of the research outcome through the standards of ISO-25010. The evaluation of the outcome through user feedback is carried out by a distributed questionnaire among the potential stakeholders.

1.4 Research Impact

The research study holds substantial academical and managerial implications. The primary aim of this study is to address the identified research gap through a pragmatic solution of a CDSS. The developed CDSS effectively achieves this primary goal, serving as a tangible and practical tool that provides solution to the recognised real-world issue. The complete process of resolving the issue reveals the answers to the four research questions of the study, which broadly consider two perspectives of the solution: design and data analysis.

The design process of the CDSS significantly contributes to knowledge enhancement through its theoretical explication. Employing the DSRM in addressing a real-world health informatics issue demonstrates the practical applicability of design science theories in problem-solving. This design process serves as a theoretical augmentation of the DSRM. The used research framework and the derived process of generating a CDSS as a solution to a real-world issue exhibits the theoretical implication of this study. The design science research approach has been used in here mainly as three folded processes: Problem investigation, design and implementation and evaluation. The problem investigation phase comprises with a solid awareness of the real-world issue. The importance of a solid situational awareness process is identified in this study while defining the used processes in this stage. The vitality of regular client meetings, systematic review on the field, stakeholder identification, and brainstorming sessions through client meetings are emphasised through this study. Further, use of the result of data analysis stages, on communicating with the clients to collect a resourceful dataset has been recommended in this study. The design and implementation stage represents the solution designing, data collection, data pre-processing, exploratory data analysis, feature selection, and model selection stages with the most suitable techniques and their usage in the healthcare sector. Moreover, the process of implementation of web portal, its suitability in the domain and the techniques used in the implementation of the web portal contributes to the knowledge in the domain. The selected evaluation methods demonstrate their suitability for comparable designs, and enhancing knowledge in system evaluation. The adopted evaluation framework consists of three phases: design evaluation, algorithm evaluation, and implementation evaluation. The confirmation of the designed solution with a standard evaluation matrix to ensure the suitability of the solution are conclude the design

evaluation of the study with a thorough standardised process. Algorithm evaluation evaluate the prediction accuracy of created models while the implementation evaluation focuses on evaluating the final research outcome, the CDSS. The selected approach of evaluation, contributes to the system evaluation via a systematic process. The formative and summative evaluation processes cover the entire system evaluation, while emphasizing the accuracy and usability of the implemented system. The demonstrated process of solution design, implementation, and evaluation confirms the design science theories from a pragmatic perspective. Additionally, the formation of the study enhances the knowledge of adopting DSRM in healthcare issues. The adopted data analysis techniques confirmed the suitability of use of survival analysis in predicting the healthcare outcomes. The process of data analysis to achieve the prediction of survival confirms the theory of survival analysis. The employed data analytics techniques also contribute to knowledge enhancement. The chosen survival analysis techniques, particularly in predicting the survival rate of selected CoDM, affirm their applicability in the domain. The combined process of DSRM and survival techniques in the discipline of healthcare improved the design process. Besides the above-mentioned contributions of knowledge enhancement, a comprehensive feature set has been derived through a systematic review to predict complications associated with diabetes. This extracted feature set of the systematic review of the study, addresses the existing lack of consensus on predictive features for diabetes complications while expanding the knowledge of the field (Madurapperumage et al., 2021). Additionally, the outcome of the research, mitigates existing practical issues in the domain. The recognised practical issues such as fixed term prediction, binary/ordinary level risk presentation, limitations on the number of complications that can be predicted, and the existing contextual scarcity in New Zealand are addressed through the outcome of this study. The generated CDSS predict the survival of individuals using survival analysis techniques which provide the risk as a percentage in chronological manner. The CDSS is capable to predict the risk of 10 most common CoDM in their cohort, which is vastly beneficial than predicting one or rarely few complications as in existing systems. The contextual scarcity is addressed using a representative diabetes cohort from the Waikato region in New Zealand.

Moreover, the implemented CDSS offers an array of managerial implications. Policymakers and resource allocators in the sector of health care can utilise the CDSS to make informed decisions, by receiving an overarching view of the considered cohort. The generated CDSS offers a statistical overview of the cohort regarding complications, ethnicity, gender, and age groups, aiding in analysing policymaking and resource allocation. The health care administrators may use the system to assist in administrating the data of diabetes for the purpose of providing efficient and successful facilities to the patients. Doctors, GPs and nurses may use the CDSS to inform their decisions on specific laboratory tests, initiate treatments, advise patients on potential complications, suggest the necessary laboratory tests for diagnosing complications at onset, issuing early warnings, starting treatment plans, and recommending dietary/exercise routines. Disease diagnosis and prognosis are two primary areas that can actively contribute to preventive programmes. Predicting the chronological percentages of risk of having CoDM by considering their biomedical and behavioural data is highly beneficial for issuing early warnings and customisable preventive programmes throughout the patient's lifetime. Since the CoDM are fatal and permanent, the prediction of complications in advance is more beneficial than diagnosing them at the onset. The early prediction not only serves financial benefits but also reduces morbidity and mortality. The gravity of the comorbidities of DM, such as renal dysfunction, amputation, foot ulcers, stroke, blindness, and erectile dysfunction, can be communicated to the patient in advance, which might impact the changing attitudes towards the diseases. Moreover, forecasting the chronological percentage of risk will encourage patients to improve their dietary and physical activities, which is profoundly practised in controlling T2DM. Patients, in turn, can use the CDSS to assess their risk of diabetes complications by inputting their details into the system. The CDSS facilitates self-awareness of diseases and their prevention through continuous health monitoring. The aforementioned academic and managerial research implications contribute to enhancing the knowledge of the field while contributing to solving a real-world issue with a CDSS which also possess solid managerial implications.

1.5 Structure of the Thesis

The thesis is organised into distinct chapters to facilitate a systematic research study presentation. The introductory chapter provides an overview of the research problem, research focus and overall summary of the study. Chapter 2 critically examines the state of the art and identifies gaps and defines the current research aims to address. This critical examination of previous research helps situate the present study within the broader academic landscape, identifying gaps, controversies, and areas where further investigation is warranted. Chapter 2 emphasises the study's identified research gap through reviewing the literature. Chapter 3, research methodology chapter, outlines the research design, detailing the methods used for data collection and analysis and the methods used for implementation and evaluation. The justification of techniques and methods used in the study is explained in third chapter. The adopted DSRM, uses their approach of solving a problem in a cyclic process which comprises five stages; problem investigation, treatment design, design evaluation, treatment implementation, and implementation evaluation, which can be broadly considered as three main stages; problem investigation, design and implementation, evaluation. The chapters from 4 to 6 uses this structure to explain the process more systematically. Chapter 4 describes the situational awareness of the study. The investigation of the real-world issue is explained here while concluding the chapter with a framework for the issue. The designing of the artefact of the study is explained in the Chapter 5, from gathering requirements to the completion of the implementation of the system. The usage of solution design and implementation phase of the DSRM is described in the fifth chapter including the artefact designing, data collection, data pre-processing, data analysis, model selection, tool implementation and launching which cover the entire solution design, and solution implementation phase. The evaluation process of the study is included in Chapter 6, with an elaborated description of design, algorithm and implementation evaluations. The final phase of the adopted DSRM is described in the sixth chapter, with justifications, and used evaluation processes. Subsequently, Chapter 7 presents the key findings of the study, accompanied by relevant data visualisations and interpretations. How the results of conducted analysis being useful on answering the research questions of the study are explained in this seventh chapter while presenting the research findings. Chapter 8 offers insights and implications. Further, this summarises the study's main contributions, acknowledges limitations, and suggests

avenues for future research. The interconnections of the chapters play a vital role in structuring the thesis into a systematic format while navigating through the flow of the thesis. The situational awareness chapter (Chapter 4) connects with Chapter 2 to form and revise the research questions after understanding the real-world issue. The completed phases of setting the research questions and client's requirements lead to design the study's artefacts. The designing of artefacts describes the adopted process of system implementation which connects to chapters 2, 3, and 4 for clarifying the state-of-the art, justifying the methodological appropriateness, and exploring the real-world issue, respectively. The evaluation phase of the study refers back to Chapter 5, revising the methods used in the evaluation. The research findings and interpretation chapter describe the outcomes of the research by answering the extracted research questions of the study. The final chapter of the thesis explains the research implications with reference to the research background in the introduction chapter. The citations and references adhere to a consistent citation style, ensuring academic integrity. This structured format allows for a comprehensive explanation of the research topic, contributing to advanced knowledge in the respective field. The following diagram (Figure 1.2) illustrates the thesis structure.

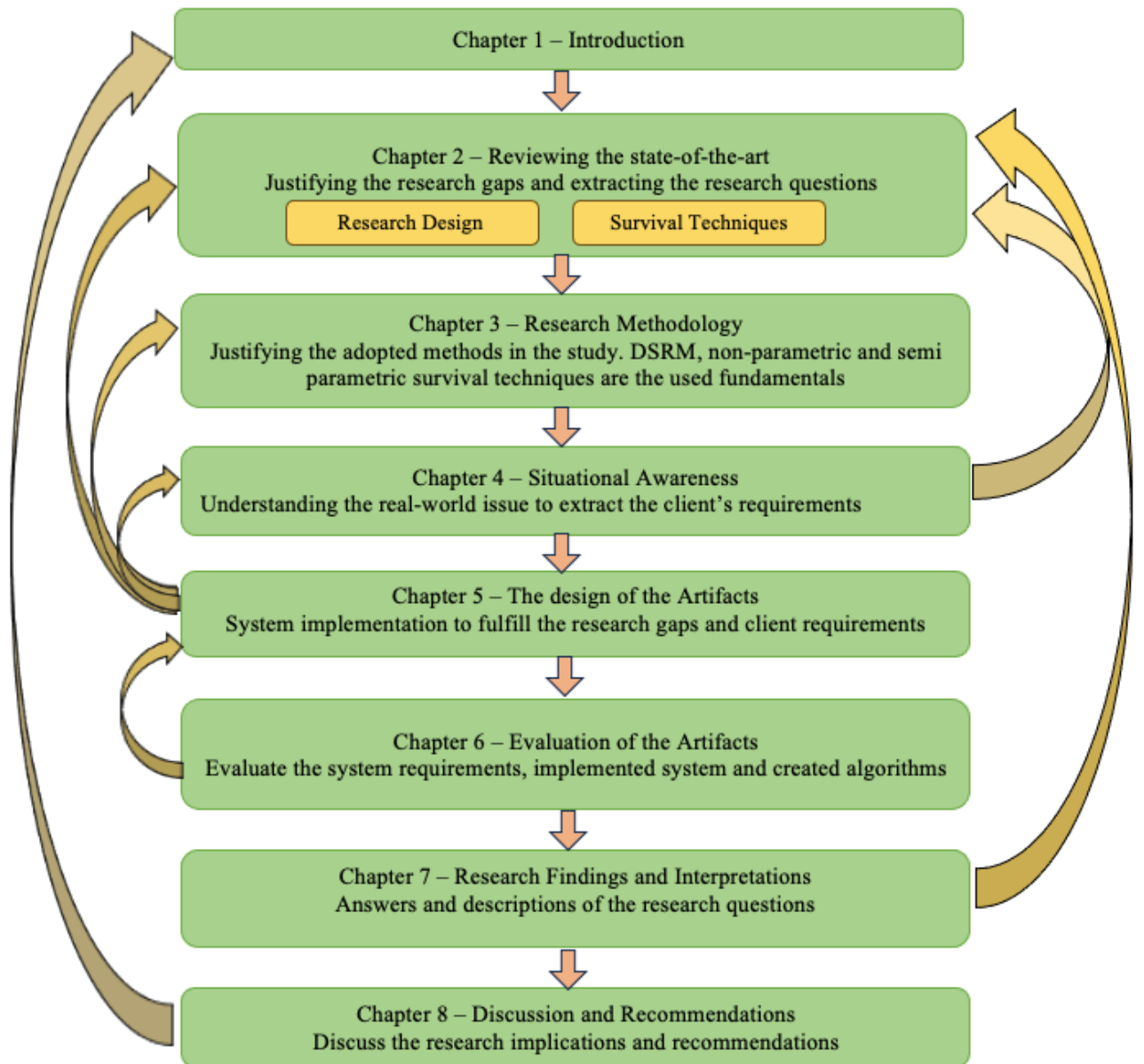


Figure 1-3 : Diagram illustrates the thesis structure

Chapter 2 Reviewing the State of the Art

The rapid growth of sensors, digital biomarkers, and smart devices makes health informatics a prominent research area. The eminence of health informatics has been triggered by three primary factors: massive datasets, soaring health expenditures, and a high prevalence rate of chronic diseases. According to the international data corporation, the digital universe may expand to 40,000 exabytes by 2020, where embedded and medical data will take 3400 exabytes of it (John Gantz & Reinsel, 2012), which is 6.8 million times larger than the amount of data added daily by Facebook users. Moreover, global spending on health was US\$7.8 trillion in 2017, which is nearly 0.4 of the GDP of the USA. Health informatics has evolved due to its ability to provide solutions for the severe burden of the healthcare sector. Due to the benefits achieved through health informatics, healthcare providers tend to develop more reliable tools and systems using diverse computerised techniques. The HIS profoundly resulted due to their capability of managing the health data repositories for efficient health care management (Haux, 2006). The vast number of types of data generated instantly from various resources leads to paying attention to reducing waste and inefficiency in three areas: clinical operations, research and development, and public health. Further, this can contribute to evidence-based medicine, genome analytics, pre-adjudication fraud analysis, device/remote monitoring, and patient profile analytics (Raghupathi & Raghupathi, 2014). McKinsey & Company (2011) state that using data analytics in the healthcare industry can reduce more than US\$300 billion annually in U.S. healthcare. Various commercial platforms have been created to cope with the requirements of the healthcare industry. Examples of such commercial platforms include IBM's Watson Health, Flatiron Health, MedeAnalytics, and Ayasdi. These platforms provide services such as sharing health-related data among hospitals and researchers, assisting with applications for managing clinical and health-related data, providing longitudinal patient data for risk assessment, deep learning methods for healthcare diagnosis, and providing healthcare analytics solutions with the aid of artificial intelligence (AI) techniques (Dash et al., 2019). The eminent use of HIS in managing health data repositories benefits in many aspects of healthcare management.

Patient empowerment systems have become a trend among healthcare providers, which results in the evolving self-management of electronic health tools. “The empowerment informatics (EI) framework suggests that patients living with chronic illnesses and collaborating nurses can use health-enabling technologies to support the relationships among patients’ behaviours (self-management), patients’ unique characteristics and context (health force), and patients’ individual goals” (Knight & Shea, 2014, p. 93). According to Alpay et al. (2011), little evidence exists for the effective use of electronic health in self-management. Therefore, systems that support health decision-making and self-assessment are profoundly beneficial for preventing the prevalence of chronic diseases and diminishing soaring health expenditures.

2.1 Data Analytics in Healthcare Management

Given the perspective of saving money and enhancing quality of life, various data analytics techniques have been utilised in the healthcare management sector. Many large companies worldwide, such as IBM, Microsoft, Google, and Intel, are moving to contribute to healthcare management through sophisticated platforms. Fundamentally, three forms of data analytics can be seen in healthcare: small data, predictive modelling expansion, and real-time analytics (Wills, 2014). Several studies (Hussain & Nguyen, 2014; Liu et al., 2012; Wang & Cao, 2014) focus on accomplishing the goals of these three forms of healthcare by applying statistical models, machine learning (ML) and AI techniques. Diseases diagnosing with non-invasive or minimally invasive methods, issuing early warnings on risky patients, experimenting with therapeutic outcomes, and predicting the prevalence of diseases are especially popular tasks of healthcare management which can be performed through statistical modelling, ML and AI techniques. Probabilistic analysis, regression modelling, survival analysis, and Naïve-Bayes’ analysis are popular statistical techniques used in the healthcare sector. Supervised, unsupervised and reinforced learning methods are used as data analysis techniques while applying data mining methods for database knowledge extraction in machine learning. Additionally, ML techniques in health data repositories emerge with the availability of clinically relevant datasets, where a broad range of clinical tasks can be accomplished, from identification/diagnostic to prediction tasks (Wiens & Shenoy, 2018). Disease prognosis has become a widespread research area that benefits patients and health authorities. Statistical

modelling and ML and AI techniques are repeatedly used with various enhancements in the healthcare sector. Logistic regression has been used for disease diagnosis (Li et al., 2022; Shariatnia et al., 2022) and prognosis fields (Sun et al., 2009; Urbanovych & Suslyk, 2018) over the years. The Cox proportional hazard model also became a popular statistical model in the field (Clair et al., 2022; Kim et al., 2020; Olfatbakhsh et al., 2022; Partridge et al., 2016; Rooney et al., 2013; Wilstrup & Cave, 2022), other parametric and non-parametric survival analysis techniques are also widely used in statistical data modelling (Chen et al., 2016; Derincek et al., 2020; Hajihosseini et al., 2016). In addition to the above-mentioned statistical models, various ML and AI techniques have recently been adopted in health informatics. Decision Trees, Support Vector Machines, K-nearest neighbours, Naïve Bayes, Neural Networks, clustering algorithms, Association Rules, Fuzzy logic systems, and ensemble techniques are the most common ML techniques applied in healthcare. Neural networks are used in various disease diagnosis and prognosis tasks. Costa et al. (2019) created an artificial neural network-based heart disease diagnosing system, and Oh et al. (2018) introduced a deep learning approach for diagnosing Parkinson's disease by analysing EEG signals using a convolutional neural network. A neural network-based early diagnosis method of Alzheimer's disease was invented by Peña-Bautista et al. (2019). A fuzzy-based cancer diagnosis system has been implemented by Govinda et al. (2017), and a neuro-fuzzy-based liver disease classification system has been created by (Vaidya et al., 2017) . Furthermore, classification techniques are the most prominent ML techniques among all, which can classify instances according to pre-defined knowledge. Classification techniques are frequently used in the prognosis and diagnosis of diseases (Elhoseny et al., 2019; Mariani et al., 2019), including dividing cohorts into groups by considering their characteristics (Monteiro-Soares et al., 2020) (Woodmansey et al., 2017), make rules to separate patients into various risk levels (Abdar et al., 2017; Gurumoorthy et al., 2018), and medical image classification (Elhassen, 2017; Rączkowski et al., 2019). Data analytics has been enhancing the healthcare sector by evolving its techniques. Due to the immense benefits of utilising information systems in the healthcare sector, their whole functionalities have been covered through the novice techniques, including disease diagnosis, prognosis, risk estimation, survival analysis, decision-making, policy creation, and many more.

Additionally, the realisation of the vitality of hidden knowledge in massive data sets, soaring cost, high demand for expert knowledge, expenses in the diagnosis of diseases through laboratory tests, the impact of undiagnosed chronic diseases and their complications on the health of individuals and community, leads to creating a whole spectrum of computerised systems for decision making in the healthcare sector. The HIS and its sub-sets, such as CDSS and SAT, are a few well-known decision-making systems in the spectrum (Winter et al., 2023). Ledley and Lusted (1959) analysed the reasoning process inherent in medical diagnosis to create a computerised decision-making system. The initiative of combining medical knowledge to develop decision support systems in medical/healthcare fields is rapidly growing while branching out to areas such as diagnosis, risk prediction, prognosis, finding associations between biomedical and behavioural features with diseases, and the therapeutic effect for disorders. CDSS, Medical decision support systems (MDSS), and health decision support systems (HDSS) have gradually evolved due to the benefits they deliver for managing healthcare. According to the mode of intervention in the decision process, MDSSs can be distinguished by a) indirect decision support systems or documentary assistance systems, b) systems for automatic reminders, and c) consulting systems (Mansoul et al., 2013). The terms “CDSS”, “MDSS”, and “HDSS” are used inconsistently by scholars to entitle the systems which can support the process of making decisions in healthcare sector. However, the term “CDSS” is widely used in the systems with the capability of assisting in making clinical, and health administrative decisions. Different CDSSs have been created to assist in the various steps of decision-making. FS-PSO-SVM is an expert system for diagnosing thyroid disease with the SVM approach. This system provides a guide for physicians to diagnose thyroid disease accurately (Chen et al., 2012). Since the differentiation between benign and malignant breast cancer by malignant mammographic findings is intricate, an expert system for the diagnosis of breast cancer (Ex-DBC) was created in 2011 (Keleş et al., 2011). This system assists in categorising patients into benign and malignant. Further, Zolnoori et al. (2011) have created an expert system to predict the possibility of having fatal asthma. An intelligent prediction and classification system has been designed for diagnosing chronic kidney disease with density-based feature selection and an ant colony-based optimisation (Elhoseny et al., 2019). The performances of decision-making systems are enhanced through various techniques,

such as feature selection, feature engineering, and handling class imbalances. Although the existing HDSSs perform satisfactorily, most are rarely in clinical practice.

Moreover, SATs have evolved parallel to CDSSs, where the SATs are beneficial for the patient's self-monitoring and self-evaluation. Richard Buckminster Fuller, an American architect, said, "If you want people to learn something, don't teach them. Instead, give them a tool, the use of which will make them think differently" (Duffy et al., 2008, p. 38). SAT is a helpful tool that can assist patients in assessing their health while being aware of medical conditions. SATs were created to fulfil a wide range of purposes, including the awareness of one's health, emergency alarm, diet recommendation, automatic medication serving, exercise recommendation, and making platforms for sharing with communities with the same disorders. The SATs focus on improving self-management tasks, such as medical, role, and emotional management, through six self-management skills: problem-solving, decision-making, resource utilisation, forming a patient-provider partnership, action planning, and self-tailoring (Lorig & Holman, 2003). Web-based SATs become prominent due to their high accessibility and profitability. Yoon et al. (2009) developed a web-based tool to assess the risk of six diseases: coronary heart disease, stroke, diabetes, colorectal, breast, and ovarian cancer. Further, risk assessment systems, early warning systems, and SATs are populating to empower patients while being aware of their health conditions. Although patient empowerment has become popular with SATs, using patient empowerment principles to design the e-health tool is inadequate. To overcome this, six essential components of empowerment have been extracted in literature by Alpay et al. (2011): a) communication, b) education and health literacy, c) Information, d) self-care, e) decision aids, and f) contact with fellow patients. The designers of SATs can create innovative tools by thoroughly attending to the principles of empowerment to make user-friendly and reliable systems. However, the poor performance and less engagement of the existing SATs unintentionally lead to less self-awareness and a high prevalence of chronic diseases, which creates an obvious research gap in making enhanced SATs.

2.2 Diabetes Mellitus

“Diabetes is an important public health problem, one of four priority noncommunicable diseases (NCDs) targeted for action by world leaders. Both the number of cases and the prevalence of diabetes have been steadily increasing over the past few decades” (WHO, 2016, p. 6). IDF shows that 537 million adults aged 20–79 are currently living with diabetes, which is predicted to rise to 643 million by 2030 and 783 million by 2045. Additionally, the economic impact of diabetes is non-neglectable. The direct expenditures for diabetes are estimated at US\$966 billion in 2021, while the total diabetes-related health expenditures will reach US\$1.03 trillion by 2030 and US\$1.05 trillion by 2045 (IDF, 2021). The high prevalence rate and soaring health expenditures of diabetes makes diabetes a prominent research field.

Diabetes mellitus (DM) is a severe chronic disease where the body's glucose levels rise due to either insufficient insulin production or inefficient use of insulin. Insulin is a hormone that plays a critical role in regulating blood sugar levels within the body. Hyperglycaemia, also known as high blood glucose levels, occurs due to the malfunctioning of insulin usage in the body. The blood glucose levels are used as clinical indicators of the diabetes. Four types of tests are available for diagnosis of diabetes. The equal or higher of 48 mmol/mol of glucose in the blood from the HbA1c test, or that of 11.1 mmol/L in the two-hour plasma glucose test, or of 7.0 mmol/L in the fasting plasma glucose, or 11.1 mmol/L in the random plasma glucose test are the standard criteria for diagnosing diabetes (IDF, 2021) .

The type of diabetes occurring due to an autoimmune process in which the insulin-producing beta cells at the pancreas are attacked by the immune system, resulting in little or no insulin production is called T1DM. The most common type of diabetes is T2DM, accounting for over 90% of all diabetes worldwide. The insulin resistance of the body is the cause of T2DM. Initially, beta cells produce high amounts of insulin, eventually leading to inadequate insulin production due to the failures of beta cells. T2DM can be symptomless or may show less dramatic symptoms. Additionally, diagnosing the exact onset of T2DM is impossible to determine. As a result, one-third to one-half of people with T2DM in the population may be undiagnosed. Impaired glucose tolerance (IGT) and impaired fasting glycaemia (IFG) are the intermediate conditions from the normal range of blood glucose level to diabetes. Moreover,

the following common type of diabetes is gestational diabetes, including women with known T1DM, T2DM or rarer forms of diabetes before pregnancy.

T2DM is associated with many complications that can be classified into two principal categories: microvascular and macrovascular. Microvascular complications predominantly affect the small blood vessels, exemplified by diabetic retinopathy, a condition characterised by damage to retinal microvasculature that may lead to visual impairment or blindness. DNep, another microvascular complication, ensues from pathological changes in the renal microvasculature, progressively culminating in renal failure. DNeu, emanating from nerve damage in response to microvascular alterations, manifests as sensory disturbances, notably in the extremities. Conversely, macrovascular complications predominantly affect the larger blood vessels. Cardiovascular disease (CVD), including coronary artery disease, myocardial infarctions, and strokes, constitute the foremost macrovascular complications of diabetes. Furthermore, peripheral artery disease, typified by reduced blood flow in the limbs, and hypertension, a common comorbidity, are also prevalent among individuals with diabetes (Chawla et al., 2016; Papatheodorou et al., 2018). Although a spectrum of micro and macrovascular complications results from DM, the loss of vision, end-stage renal disease, cardiovascular events, and lower extremity amputation are emphasised in many situations due to their severity and significant impact on a patient's quality of life and overall health outcomes. Due to the vitality of early diagnosis of diabetes and its complications, risk prediction models, risk estimation tools, and decision support systems have resulted.

2.3 Applications of Health Information Systems in T2DM and its Complications

The rapid growth of the healthcare sector with the burgeoning of EHRs and enhancements in data analytics and AI techniques have resulted in many practical systems in diabetes and its complications. Diagnosis and prognosis of diseases through analysing datasets with various combinations of statistical models and ML and AI techniques have become popular due to their benefits. Notably, using information systems in the healthcare sector reduces the financial burden by diminishing the costs of laboratory tests, consultation fees for the specialist, and the need for specialised equipment and diagnosis tests. Furthermore, these systems can improve quality of life by issuing early warnings, empowering patients, and

customising treatment plans. Scholars have concentrated on implementing CDSSs to cope with these requirements. Computerised systems have been widely applied to fulfil a variety of functionalities in the healthcare management of diabetes and its complications. Moreover, the diagnosis of diabetes and its complications, predicting the onset of diabetes from pre-diabetes, predicting the onset of complications of diabetes, estimating the risk of diabetes and its complications with time, analysing the effect of a specific risk factor or group of risk factors on the onset of disease, and exploring the genetic contribution for diabetes are some of the topics that are being studied. As a result, non or minimal-invasive diagnosis tools, risk prediction models, clinical decision models, and self-assistant systems are widely created in diabetes and its complications.

Risk scoring systems are popular in the healthcare sector to estimate the vulnerability of individuals towards diseases. These automated tools provide the potential severity of a disease as a numerical outcome. The risk scoring systems are varied with considered features, datasets, outcomes, and focused purposes and techniques used to calculate the score. The risk scoring systems can deliver many benefits, including diagnosing diseases with minimally invasive methods, early detection of diseases, self-assessment, and quick responses for disease diagnosing. Due to the immense benefits of diagnosing T2DM through risk-scoring systems, a range of scoring tools has been implemented. Herman et al. (1995) created a simple questionnaire to identify undiagnosed diabetes by analysing their age, sex, history of delivery of a macrosomic infant, obesity, sedentary lifestyle, and family history of diabetes. The FINRISK model is created based on the details of a Finnish cohort (Lindström & Tuomilehto, 2003). Additionally, Griffin et al. (2000), Schmidt.M.I. et al. (2005), Aekplakorn et al. (2006), and Schulze et al. (2007) created risk scoring systems based on patients in Wessex, whites and African Americans in the U.S, Thailand, and Germany, respectively. Further, Chen Lei et al. (2010) created AUSDRISK to assess the risk of having T2DM based on a cohort who was born in Australia, New Zealand, and the United Kingdom. Since the extremity of feature values varied from one group to another, specific risk-scoring systems should be created to get high accuracy. Some scholars validate their models with external datasets to prove the model's capability in assessing patients in different cohorts. Kanaya et al. (2005) used two groups to create prediction rules and validate the model. Moreover, risk-scoring systems

enable access through web pages for self-assessment. AUSDRISK, QDiabetes, risk score by the American Diabetes Association, and Diabetes UK are some web-based risk assessment tools which can be used for self-assessment (ADA, 2020; DA, 2010; Diabetes UK DUK et al.; UKNHS, 2018). These risk-scoring systems deviate on the factors such as inconsistency of feature sets used to develop the model, nature of the used data sources, considered benchmark values for predicates, adopted techniques on predictions, and many more.

Due to the irreversibility of the outcomes of the CoDM, scholars have recognised the importance of forecasting the CoDM. Dogba et al. (2018) found the most crucial research topic for people with diabetes. According to that study, people care more about preventing and treating CoDM. These findings reveal the importance of engaging with the prediction of CoDM to patients and caregivers while emphasising the direction of the research interests. Further, the importance of awareness of diabetes and its complications was well established in a study conducted by Amoo et al. (2014). Due to these reasons, many research studies, which create models for diagnosing and predicting the risk of having CoDM, can be found in this field. Prognosticating the age of having complications, the risk of each complication at the current stage, forecasting the changes of risk factors, and chances of preventing complications by controlling the high-risk features are some of the pragmatic outcomes of the existing models. The onset of CoDM is predicted through a model created using logistic regression with a stepwise feature selection (Arianna Dagliati et al., 2018). The model is enhanced with appropriate methods for handling missing values and class imbalances. This model can predict the onset in fixed periods of 3,5,7 years from the first visit to hospital. However, due to the irreversibility and fatality of CoDM, fixed-time prediction is less valuable than continuous risk prediction. DCCT/EDIC Research Group (2005) investigated intensive diabetes treatment in cardiovascular disease. They predicted the time of onset of cardiovascular events in T1DM patients. The effect of intensive diabetes treatment on peripheral arterial calcification has been studied using the same dataset that resulted in the study mentioned above (Carter et al., 2007). Since this is an extension of the DCCT/EDIC study, it is considered DCCT2. The absolute 5-year risk of the onset of fatal/nonfatal CVD prediction model is created based on a Swedish cohort with T2DM (Cederholm et al., 2008). A risk equation is derived as the study's outcome, which can predict cardiovascular risk among

patients aged up to 70. Hippisley-Cox et al. (2007) introduced a risk score for predicting CVD (QRISK) for people who live in the United Kingdom. Nine clinical features were used in creating QRISK, where the Cox proportional hazard model utilises in prediction. QRISK2 is an extension of the QRISK study, where the validation was done with two ethnic groups in England and Wales (Hippisley-Cox et al., 2008). Due to the home advantage in both these algorithms, the researchers advised performing external validation. It becomes a requirement to validate the model externally to cope with the issue of handling different ethnicities. Kim et al. (2019) created a transferable machine learning (ML) model to predict the CoDM, which was evaluated through an external dataset. Further, with the need to quantify the severity of CoDM, the diabetes complications severity index (DCSI) was created. This index comprises seven complications: CVD, DNep, DR, peripheral vascular diseases, stroke, DNeu, and metabolic diseases (Young et al., 2008). Although DCSI is a powerful tool in predicting CoDM, it does not give the probability of risk as a percentage or numerical value. Instead, it provides the severity of the complication in three different levels. Further, the wide acceptability of the model has been limited because of the considered dataset at the creation of the model. Further, there is an update with the ICD 10 translation for DCSI in the 2017 (Glasheen et al., 2017).

Solving a real-world issue in a healthcare setting through an information system perspective has been popularly used for the last few decades. The existing applications of HIS in this field use various approaches to solve issues. The majority of the applications use empirical research methods in information systems. The popular risk scoring models, predictive models in DM, diagnosis and prognosis models are frequently use the empirical problem-solving methods to build the system as answers for the raised issues in healthcare settings. The case studies, field studies, field tests and laboratory studies are common empirical methods use in this field. However, solving real-world issues through an intervention of information systems falls into the design science approaches. The development of HIS through the aid of design science approach has become widely used due to their systematic process on solving the issues.

Although the HIS provide plenty of benefits, all of them are inherent with their limitations (Schwarz et al., 2009). The limitations of considered datasets, covariates, covariate values, function for risk calculation, and used techniques for implementation resulted in many restrictions for their performance. Many existing risk scoring systems estimate the risk based on the current feature values of the individual, as predictions are made by fixed periods such as the short, medium, and long terms (Amir Talaei-Khoei & Wilson, 2018) or frequently in 10 years (Exalto et al., 2013). Since CoDM is irreversible after the onset, this fixed-time prediction model serves fewer benefits in predicting CoDM. Furthermore, all the existing systems implemented their model based on a dataset of a specific cohort (Lei Chen et al., 2019; Lindström & Tuomilehto, 2003; Schulze.M.B. et al., 2007). As many researchers found, the feature values varied among ethnicities (Lagani et al., 2015). Some research studies externally validated their models (Kanaya et al., 2005) to test the ability to use their models widely. However, scholars have exaggerated the importance of creating unique risk-scoring systems for specific cohorts (Schwarz et al., 2009). Due to the lack of a standard risk-scoring system developed based on a cohort of New Zealand to the best of our knowledge, it will be worthwhile to create a CDSS by considering the characteristics of a New Zealand cohort. The ability of knowledge-based systems to find solutions with case-based, rule-based, and instance-based reasoning methods provides insight into the capability of using knowledge-based systems to make decisions dynamically (Alder et al., 2014). Although many models are available to support decision-making at clinics, for various reasons they do not populate as expected, for instance high time-consuming data entry, limited outcomes, less acceptability of the model, less accuracy of the decision, and limited expert knowledge embedded in the model.

Although scholars are inventing innovative solutions with high dedication, their practical usage is limited due to their inaccessibility, low popularity, or inconvenience in practice. As a result, healthcare providers promote various approaches such as CDSS, HDSS, SAT, self-monitoring systems, etc. Healthcare management is achieving several goals by promoting these systems among stakeholders such as patients, clinicians, doctors, and other health authorities (Dash et al., 2019). To make successful healthcare promotional strategies, people should be aware of the disease and its complications (Amoo et al., 2014). Self-awareness plays

a significant role in mitigating the gravity of CoDM from the perspectives of economics and the quality of life. Systems capable of evaluating patient profiles, giving treatment plans, suggesting diets and exercise routines, and assessing risks for each individual in CoDM become popular. Lagani et al. (2015) developed a system for determining the long-term risk of having CoDM, which can work in patient profiling and risk assessment. Although physicians have positively evaluated the long-term risk assessment (LTRA) system from this study as a better communication system with patients, it only represents the fluctuation of a patient's overall risk. Even though the variation of overall risk can deliver a message to patients, the visualisation of changes in high-risk features over time leads the patients to engage with maintaining their health firmly.

Due to the capability of changing patients' attitudes with powerful visualisation techniques, CDSS and SATs use precise visualisation methods. The visualisation of data in a convenient manner plays a significant role here. Policymakers and decision-makers can visualise the dataset to make firm decisions in allocating budgets for the health sector and designing preventive programmes. It can be concluded that risk scoring systems are inherent with several limitations that need to be improved to become widely accepted. Furthermore, it is vital to have multiple approaches to use scoring systems in day-to-day life. Hinz et al. (2014) created a temporal visualisation tool to represent the variation of HbA1c levels of a population over time, which uses parallel sets and Sankey diagrams for illustrating the trajectories of HbA1c levels for a general diabetic population. The representation of this study can be used to view the consolidated trends in the progression of diabetes from backwards and forwards across all diabetes control categories from uncontrolled to normal. The fewer participants diagnosed with diabetes and the considered period for the variation of HbA1c level makes this less useful in clinical practice as a model to view the HbA1c level trajectories of a general cohort. Moreover, due to the high use of mobile devices and the vast collection of functionalities, mobile devices tend to be used as personal health assistance. These devices can be used for collecting data, modelling and visualising data, interpreting data, and empowering and improving a patient's health. The massive number of wearable devices, such as smartphones, iPods, laptops, and notebook computers, enhances the opportunities for telemedicine and self-assessment. Burford et al. (2019) conducted research to get the

viewpoint of patients using iPads to monitor diabetes data. Users were highly appreciative of the data visualisation capability of apps; they find the apps are more useful when representing data with their history and patterns, and also the impressive charts form of chronological order helps them to be aware of their diet and control other risk factors. They confirm that the iPad tablet screen size is ideal for manipulating data and is of a quality that is easily viewed and comprehended. Since the requirements of diabetes patients indicate the need for good visualisation ways to self-monitor their health details, it is worthwhile to create visualising methods of diabetes data, including the clinical and biological factors.

2.4 State-of-the-Art of the Clinical Decision Support System

“Computerized clinical decision support systems, or CDSS, represent a paradigm shift in healthcare today” (Sutton et al., 2020, p. 1). Any computerised system to assist healthcare stakeholders in making decisions can be introduced as a sub sets of HIS, with the MDSS, HDSS, or CDSS. The CDSS are the most common type of HIS which involves in the decision making from policy makers, resorce allocators, health administrative to doctors, nurses and GPs (Winter et al., 2023). The rapid expansion of health records due to the EHRs, disease registries, patient surveys and information exchanges creates a contemporary requirement to utilise collected data in making decisions effectively. Additionally, adopting CDSS in the healthcare sector serves a range of functionalities, which reduces the cost of laboratory tests, expert knowledge, efficient diagnosis and prognosis, assisting in policy designing, resource allocation and many more. Although different categories of CDSS can be seen in the healthcare sector, the fundamental categorisation is based on five characteristics: system function, model for giving advice, style of communication, underlying decision-making process and human-computer interaction (Wasylewicz & Scheepers-Hoeks, 2019). “System function” based CDSS is focused on two fundamental questions: 1. What is true? 2. What to do? Disease diagnosis and recommendation of actions are their two respective purposes. The systems created for giving advice are alert systems, which can actively or passively advise the stakeholder. The way of communication is another characteristic of CDSS, which distinguishes two sub-models: the consulting and critiquing models. The human-computer interaction characteristic of CDSS concentrates on communication with humans. Although, initially, they were slow and mundane, the introduced techniques, such as embedded EHRs, voice

recognition, acoustic alarming, and pop-up messages, make it more interactive with humans. The final characteristic of CDSS is the underlying decision-making process. The basic flow chart method for decision-making and the other available mathematical models and statistical approaches are adopted in the decision-making process of the CDSS. Regression techniques (Coslovsky et al., 2015; Ng et al., 2016; Zlotnik et al., 2016), survival analysis models (Manzo et al., 2023; Todd et al., 2022), and Bayes network (Constantinou et al., 2016; Sesen et al., 2013; Zarikas et al., 2015) are popular statistical models used in decision-making of CDSS. ML and AI techniques are also widely adopted in the decision-making models (Araz et al., 2019; Hong et al., 2018; Olivia et al., 2018; Ong et al., 2012).

Moreover, a popular classification of CDSS is as a knowledge based clinical support system (KB-CDSS) and a non-knowledge based clinical support system (NKB-CDSS). A KB-CDSS is a clinical decision support system that relies on explicit domain knowledge, guidelines, rules, or expert-driven decision-making logic. It incorporates established medical knowledge to assist healthcare professionals in making informed decisions (Osheroff et al., 2007). An NKB-CDSS is a CDSS that leverages data-driven approaches, statistical modelling, ML, or other computational methods to derive insights from large datasets. It emphasises learning patterns and correlations from data without explicit incorporation of predefined rules (Kuhn & Johnson, 2013). The key characteristics of each type of CDSS leads to selecting the most appropriate type for fulfilling the research purpose. The KB-CDSSs are most suitable with the requirements of having an expert knowledge resource, making explicit rule sets for supporting decisions, and the necessity of interpreting the patient data. In contrast, NKB-CDSSs fit with the requirement of the system for data-driven decision-making, pattern recognition, adaptability, and learning. The novice CDSSs are leaning towards the non-knowledge-based systems for the above key reasons.

Further, the existing CDSSs are used different approaches to solve the healthcare issues. The implemented CDSSs commonly used the design science approaches or software development approaches. The used approach is mainly based on the nature of the solution. The risk scoring tools, and predictive models with user interactions, basically build with the design science approaches. The systems with high end technologies use software development approach to

implement the system. Although, CDSS can be categorised as a software product, the direct use of software development / design science approach in healthcare setting is not firmly established in the field. Therefore, a solid pragmatic process on adopting software development theories and design science theories in developing CDSS is a theoretical gap that needs to be addressed.

Design science research approach is a problem-solving paradigm (Hevner et al., 2004), whereas the software engineering approach “is an engineering discipline that is concerned with all aspects of software production from the early stages of system specification through to maintaining the system after it has gone into use” (Sommerville, 2011, p. 7). The design science principles focus more on solving a real-world issue, which may or may not necessarily ended up with a working product/ implemented solution. This characteristic express that the focal point of the design science research is to systematically solve an issue by paying a detailed attention on the artefacts in its context. The software engineering approach concentrates more on the engineering perspective of a solution. Further, design science principles created on more of a basis of information system and engineering aspects.

Design science research methodology invented on the purpose of systematizing the process of problem solving. The identification of the issue, designing the solution for the identified problem and validating the designed solution are the basic functions in design science researches. Although this methodology explains the process of problem solving in a general setting, the adaptation of this method in healthcare sector has not been well researched. The majority of researches on CDSSs reports their study from the data collection phase, to results of data analysis, where the research objectives are basically concentrate on the clinical flow, usability of the system, adoptability of the created system to a clinical setting, or the accuracy of results in decision making (Bozyel et al., 2024; Schoen et al., 2015; Shamsavarani et al., 2015). The design science principles are rarely used for developing the CDSS. Even the well-known CDSS are concentrate on grasping the concepts in clinical settings rather than concentrating on the process of solving the issue (Zaman et al., 2023; Zikos & DeLellis, 2018). The applicability of the design science research principles in a healthcare setting is lack in the literature. Therefore, a solid process of developing a CDSS with the aid of design science

research principles, made a clear theoretical gap. Further, it's vital to have a clear understanding of the techniques that can be used in the components of the DSRM. Since the healthcare setting is a different discipline from information systems, the approach of problem awareness, solution designing and evaluation consists of their unique challenges which may be mitigated through the inclusion of contextual knowledge. A more customised DSRM in solving a real-world issue in a healthcare sector, is recognised as a vital knowledge gap which needs to be filled. The processes that can be used for getting a clear understanding of the real-world situation, matching the designed solution with the client's requirements, designing the most feasible system to fulfil clients requirements, the process of implementing a CDSS to answer the real-world issue, and a systematic approach for evaluating the entire CDSS has been recognised as major concerns in the process of DS to be discussed. A systematic process of developing a CDSS by following the principles of design science research may be able to fill this gap while confirming the adoptability of DSRM in the field of healthcare.

2.5 Clinical Decision Support Systems in Predicting the Complications of Diabetes

CDSSs are widely used in predicting the risk of complications of diabetes. Due to the potential health and economic impact of diabetes and its complications, predicting them with a computer-aided system has mainly been used in the past two decades. The severity, irreversibility, and high prevalence of the CoDM directly impact all aspects of patients' lives. As mentioned above, although a series of complications result from chronic diabetes, the most severe and prevalent complications are LoV, ESRD, CVD and LEA (WHO, 2016). A clear understanding of the characteristics of existing CDSSs is important when exploring the state-of-the-art. The selected CDSS can be explored from the perspectives of the adopted decision-making mechanisms, considered factors in constructing the system, the potential tasks or functionalities of the system, the focus of the system, the complication/set of complications they predict, the prediction time, used dataset/cohort and many more. Since an overall comprehension of the existing CDSS is vital in realising the state-of-the-art, the following section provides a few examples of CDSS in diabetes and its complications.

A CDSS for screening diabetic retinopathy (DR) has been created based on fuzzy random forest techniques (Romero-Aroca et al., 2019). CDSSs, developed to predict DR using ML techniques, can be prominently found in the field (Bernardini et al., 2021). Additionally, ensemble techniques are used in creating the CDSS for predicting DR. This system uses logistic regression, decision trees, random forests, and neural networks as the decision-making process to construct the CDSS (Piri et al., 2017). Spain-based CDSSs were created to predict DR with the Cox proportional hazard model (Romero-Aroca et al., 2022). To create this model, nine risk factors are considered: current age, sex, body mass index (BMI), duration and treatment of DM, arterial hypertension, HbA1c, urine–albumin ratio and glomerular filtration. A diabetic foot risk stratification system has been developed as an evidence-based CDSS with surveys, focus groups, and an online web tool access (Schoen et al., 2015). Another CDSS for risk stratification for foot ulcers has been modelled with the Cox proportional hazard techniques (Schäfer et al., 2021). Another CDSS for diagnosing peripheral neuropathy with an integrated electronic medical record has been created using a fuzzy logic-based expert systems to make decisions (Kunhimangalam et al., 2014). ML techniques-based CDSSs were created for early prediction of diabetic nephropathy (Afrash et al., 2022). Another fuzzy-logic-based expert system was designed as a CDSS for predicting chronic kidney disease (Hamedan et al., 2020). An exciting approach to predicting diabetic nephropathy was constructed using casual knowledge graphs (Lyu et al., 2023).

The design process of the CDSSs in the area of diabetes and its complications is a vital aspect in reviewing the literature. Most of the existing systems in diabetes care also concentrate on solving the issue rather than focus on the process of problem-solving. As mentioned above the focus of the majority of CDSSs are on the clinical settings. The process of capturing the right question, design the most appropriate solution to the rightmost question and evaluation of the designed solution in the right context is rarely use in the research of CDSS. Most of the time, they consider that as a software product and use the software engineering principles in the implementation of the solution. Although the results are usable products, the process of developing the product is rarely presented as a research outcome. A well-known CDSS developed for diabetes care use the empirical software engineering method (Sim et al., 2017). The theoretical gap of adopting DSRM in solving a real-world issue in the healthcare

setting can fill through a systematised process of CDSS. Although a CDSS for diabetes care or its complications is a single scenario in a healthcare setting, the process will generate a more generalised design process.

The reviewing of the state-of-art of the existing CDSS in the prediction of CoDM reveals a clear research gap which needs to be addressed in order to improve the CDSS and enhance the domain knowledge. This study focuses on addressing the identified research gaps while mitigating the existing drawbacks. The research study's rationale and scope are described in the following section.

2.6 Rationale of the Study

Due to the increasing expenditures of the healthcare sector and the high prevalence of diseases around the globe, diagnosing diseases without laboratory tests and estimating the current and future risk of individuals for chronic diseases have become fertile research areas. The ability to extract invaluable knowledge from health data repositories leads to a swift move in the healthcare sector to store the data in electronic formats. With the tremendous growth of electronic health records (EHR), researchers tend to apply statistical modelling, data mining, and ML techniques on health data repositories to achieve multiple goals. According to Wiens and Shenoy (2018), the applications of ML techniques can cover a range of clinical tasks, from identification/diagnosis to prediction. Primarily, these techniques are used for risk stratification, identifying risk factors, understanding pathogen-host interaction, and predicting the emergence and spread of infectious diseases. The field of health informatics has been evolving due to these massive EHR repositories and the evolution of ML and AI techniques. Although innovative systems result from this evolution, there is no consensus on methods for analysing, extracting patterns, and visualising the resulting enormous data from the overabundance of digital machines and sensors. Because of this, researchers usually compare the performances of various ML techniques to select the most appropriate one. Further, combining the techniques of statistical models, ML with AI reveals a new dimension where statistical patterns and data driven decisions are assisted in making decisions. The NKB-CDSS became popular with the drastic use of statistical modelling and ML and AI techniques in the healthcare sector. The ability of the system to learn patterns and

relationships from the data without explicit incorporation of predefined rules or expert knowledge, and the use of statistical techniques in identifying correlations, trends, and predictive patterns makes the NKB-CDSS ideal for the predictive healthcare support systems. Moreover, CDSSs can address a range of tasks in healthcare management, including decision-making, policy-making, examining therapeutic effects, risk prediction, issuing early warnings, self-monitoring systems, etc. These systems can reduce the economic burden by diagnosing and predicting diseases through analysing clinical data without the involvement of a physician, making effective decisions using extracted knowledge from stored knowledge bases, estimating the risk of diseases and predicting the occurrence of fatal conditions while enhancing the awareness of self-health conditions among patients, provide clinical support decisions, and reveal the severity in future through predictions.

The high prevalence rate of DM and its complications and the serious financial crisis over the globe leads to conducting various research studies on managing diabetes. The predicted 643 million diabetic patients in 2030 (IDF, 2021) will worsen their economic burden. As IDF reported, the global expenditure of diabetes has been tripled from 2003 to 2013. (International Diabetes Federation, 2013). Further, the WHO stated that the indirect costs of diabetes associated with productivity loss, premature mortality, and negative impact on the GDP of the nations will lead to an unavoidable loss (WHO, 2016) .

The reviewing of the literature is directed to a contemporary requirement of a CDSS to predict the complications of diabetes. The existing drawbacks of the current CDSS made a clear research gap. Most of the CDSSs forecast the risk in discrete-time periods such as 3,5, 7, and 10 years, or as short, medium, and long terms, which delivers limited benefits compared to the prediction of chronological risk scores. (Bernardini et al., 2021; Chen et al., 2010; Hamedan et al., 2020; Kunhimangalam et al., 2014; Romero-Aroca et al., 2019; Saaristo et al., 2010; Schoen et al., 2015)Presentation of the risk of having complications of diabetes for patients is another important aspect of a CDSS. Presenting the result as a binary output or level-wise expression, such as low, medium, or high (Chen et al., 2010; Kunhimangalam et al., 2014; Romero-Aroca et al., 2022; Schoen et al., 2015; Singla et al., 2022), may not be as useful as presenting them in percentages. Therefore, it's a crucial characteristic of a CDSS to present

the risk of complications as percentages. Additionally, the ability of a CDSS is limited to one or a few complications of diabetes (Bernardini et al., 2021; Hamedan et al., 2020; Kunhimangalam et al., 2014; Romero-Aroca et al., 2019; Romero-Aroca et al., 2022; Schäfer et al., 2021; Schoen et al., 2015; Vartiainen et al., 2016). Due to the limited capabilities of the CDSS, their usefulness and popularity also become limited. Designing a CDSS by incorporating the data of multiple complications may address this gap in the current systems. Although current risk scoring systems can achieve high accuracy for similar cohorts, the performance is controversial with other datasets (Aekplakorn et al., 2006; Chen et al., 2010; King et al., 1999; Vartiainen et al., 2016; WICHAI AEKPLAKORN et al., 2006). Since the feature values and importance vary according to the ethnicities (Schwarz et al., 2009), thorough consideration is required to create a widely accepted system that can dynamically calculate the risk by contemplating the specific characteristics of each cohort. Although the CDSS can address this issue by enriching the dataset with the details of multiple cohorts, the feasibility is controversial. The data availability, accessibility, and consensus of laboratory results are common obstacles to developing an enriched dataset. Additionally, a contextual research gap has been recognised in New Zealand (Atlantis et al., 2017; Kenealy et al., 2005; Peiris et al., 2009; Pylypchuk et al., 2018). The diabetes population of New Zealand has significantly increased from 3.4% to 4.0% over the last decade, which results in 277,800 people around Aotearoa (NZ-MoH, 2022). The annual median of the estimated cost of diabetes per person in New Zealand was US\$3967.5 in 2021 (IDF, 2021). The impact of DM on the health sector in New Zealand, its irreversible effect on individual health, its enormous effect on the status of the country's health index, the high prevalence rate of diabetes, and high expenditure shows the need to design a CDSS with a local dataset. Although few studies have been conducted based on the diabetes population of Aotearoa, there are significant gaps when considering the existing CDSS based in New Zealand. Most studies look at the population of diabetes patients in Aotearoa without considering their ethnic diversities (Pylypchuk et al., 2021; Robinson et al., 2012). Additionally, studies are focused on one complication of diabetes, such as CVD (Elley et al., 2010; Pylypchuk et al., 2021; Pylypchuk et al., 2018; Robinson et al., 2012), nephropathy (Scott et al., 2006; Simmons, 1998; Simmons et al., 1994), and retinopathy (Hill et al., 2021; Rao et al., 2022). Moreover, the scarcity of a standard CDSS for predicting the complications of diabetes among the population of Aotearoa clearly shows the importance of

a standard, more functional and widely accepted CDSS for New Zealand. The identified issue of the existing CDSS leads to more concentration on a suitable data analytics method, which can mitigate the issues while improving the domain of CDSS with a stable data prediction method.

In addition to the recognised above-mentioned drawbacks of the existing CDSS, a theoretical research gap is identified by examining the utilised theories in CDSSs. Scholars used various designing theories and software development approaches in this field. Designing science principles is the most popular selection among the researchers of CDSSs. Further, a number of different approaches to design science research methodology (DSRM) have been invented by researchers (Wieringa, 2014) (Hevner et al., 2004; Nunamaker Jr et al., 1990; Takeda et al., 1990; Walls et al., 1992), which shows close similarities fundamentally. Although DSRM explains the process of solving an issue in a general context, its applicability and utilisation of techniques in solving a real-time problem in healthcare is inadequate. The principles of DSRM are focus more into the disciplines of engineering and information systems. Solving a real-world issue in a systematic manner is explained in the DSRM. The applications in healthcare setting, are information systems which consists of unique characteristics and challenges. Developing a CDSS by adopting the DSRM reveals a knowledge of applicability of principles in design science in the discipline of healthcare sector. Due to the scarcity of adopting DSRM in solving the issue of healthcare settings, a significant theoretical gap has been remained. The potential techniques use in each component of DSRM and the development process of CDSS using DSRM will fill the research gap while contributing for the extension of the knowledge. Solving a real-world issue by designing a CDSS in a structural and pragmatic manner can confirm the theories of design science while resulting in a useful application. Further, three basic steps of DSRM are emphasised in every approach: designing, implementation, and evaluation. The lack of consensus on techniques and standards of these fundamental steps leads to the development of a process which is widely acceptable and theoretically sound. The selected techniques and developed process result in a valuable confirmed method for similar approaches. Further, most scholars focused on specific purposes when developing a CDSS in relation to this field. Although prominent CDSSs can be found in the literature, a significant gap of the scarcity of a standard feature set to be used in the prediction of the

complications of diabetes has been recognised. Scholars used various feature sets as predictors. Some scholars develop CDSSs with specific features, such as biomedical and physiological features (ACCORD, 2008; Holman et al., 2008), clinical and demographic details (Knowler, 2002), lifestyle features (Hu et al., 2001), and genetic and molecular markers (Saxena et al., 2007; Scott et al., 2007). Therefore, extracting a standard set of features for predicting the complications is vital to filling the gap. Moreover, the evaluation phase of the DSRM is highly emphasised in all the approaches, although direct, standard evaluation criteria have not been standardised. The knowledge gap of the standard evaluation method creates a necessity for a suitable systemised evaluation process, at least for similar purposes. The above-mentioned practical and theoretical knowledge gaps lead this study to focus on addressing the existing issues and research gaps through a pragmatic CDSS for predicting the complications of diabetes, while resulting a systematic design process and prediction algorithms in survival analysis.

2.7 The Research Questions

The study's overarching objective is to design a decision support system that can assist in assessing the risk of diabetes patients to its complications. The research has aimed to create a CDSS by mitigating the issues as far as possible. The current CDSSs are more ethnicity specific. They provide efficient and accurate results with the same ethnic groups, leading to accurate CDSSs with ethnicity-specific models. Further, the CDSSs which can predict a single disease serve fewer purposes than the systems which can assist in predicting multiple complications. The fixed prediction periods of CDSS degrades its usability. The utilised factors for designing the CDSS are another important aspect of the usability and accuracy of the CDSS. Moreover, the contextual knowledge gap of CDSSs in the New Zealand domain made them eager to have a CDSS based on their local dataset. In addition, the scarcity of explanatory analysis among the ethnic groups of the New Zealand diabetes population creates an unneglectable space for a cohort-specific CDSS. Further, the scarcity of the use of DSRM in CDSS creates a theoretical gap in the application of DSRM in healthcare settings. A systematic process of adopting the principles of DSRM in an application of the healthcare sector, is vital to fill the research gap while confirming a solid theoretical framework for similar scenarios.

The newly suggested CDSS of this study tries to mitigate the existing pitfalls of CDSSs. The issues embedded with the data collection process, the identified problems in datasets of resource providers, data transformation and organising difficulties, and the ways of creating a functional dataset will be discussed. The inconsistencies of existing CDSSs confuse the researchers and designers about selecting the best feature sets, appropriate data analysis techniques, performance evaluation methods of the created systems, etc. This study aims to discuss the above-mentioned challenges to summarise and enhance the knowledge to overcome the obstacles. Further, the suitability of the data analytics and modelling techniques used in creating the proposed model and the accuracy of the model will be evaluated to justify the solution. The study reveals the existing drawbacks of designing CDSSs for a real-world client. The challenges of data acquisition, data transformation and organising to create a rich dataset which can provide the desired outcomes of a CDSS, the effect of inconsistencies of the existing CDSSs, and the impact of dissension of techniques in designing CDSSs are considered in this research to propose the solution. The current issues of the research domain direct the path of this research study into developing a better CDSS while confirming its design process. The proposed CDSS is answering to cover the existing problems. The functionalities and accuracy of the system prove the validity of the novel CDSS over the existing ones. Furthermore, knowledge of developing a CDSS for disease prognosis using survival analysis contributes to expanding the usage of DSRM in CDSSs.

The above-mentioned research objective will be achieved through two aspects of research objectives: designing the system and data analysis. The designing perspective of the system is engaged with the theory of design science, while the data analysis entirely leans towards the survival analysis methods. The design science approach of designing a multi-functional, widely acceptable CDSS will be discussed in this study to expand the knowledge of design science in CDSSs. In addition, adopting survival analysis techniques in predicting the complications of diabetes in a New Zealand cohort is explained through data analysis and modelling methods.

This study comprises with two major objectives (as aforementioned in Chapter 1.2):

1. Outcome a solid design process for developing a CDSS using DSRM.

2. Development of survival analysis model for CoDM.

To provide a clear image of research objectives and their classifications, we divided the main objective into two research questions and two additional sub-research questions.

The overarching research question of the study:

How can a prediction model be created to forecast the onset of CoDM in a cohort of diabetes patients?

The overarching research question has been divided into two sub-research questions to simplify the overarching question.

RQ1: How can a CDSS be designed by utilising DSRM?

RQ2: How can the onset of CoDM be predicted using a longitudinal data set?

The research study aims on answering these questions by creating a CDSS capable of predicting the selected CoDM from a longitudinal dataset. The direction of achieving the primary goal of the research through answering the research questions is provided in the thesis.

The above explained gaps in the predictive models in CDSSs, and specifically on the field of predicting complications of DM, leads to a more enhanced CDSS which can overcome the existing issues and is capable of delivering a user-friendly solution for a real-world issue. This study focuses on designing and implementing a CDSS to predict the complications of diabetes, using a longitudinal data set of a New Zealand cohort. The journey to achieve the primary outcome of the study answers four sub-research questions which eventually lead to the overarching research objective.

Chapter 3 Research Methodology

3.1 Introduction

The research methodology chapter explains the methods used in solving the identified real-world issue with justifications of their suitability. The chapter starts with explaining the base of fundamental research philosophy. A detailed explanation of placing the research in its rightmost niche in the perspectives of research philosophy and research paradigms are presented. The philosophical place of the study followed by describing the steps of the adopted research methodology and the required techniques in each of the steps in the process. First, the foundation of the research has been introduced with a detailed explanation of design science research principles and their adaptation to the current study. The selected DSRM introduced by Wieringa (2014) is explained and the use of the components of that in structuring the solution is described later on. The main three components of the adopted DSRM: problem investigation, artefact designing and implementation and the evaluation phases are described respectively to provide a sophisticated understanding and justification of the techniques and methods used in each step. The sub topics 3.4, 3.5 and 3.6 respectively described these three components, while creating a solid theoretical background for chapter 4, 5, and 6. The section 3.4 provide a detailed explanation of problem investigation phase of the DSRM which conclude the process of defining the ultimate artefacts of the study. The section 3.5 covers the entire design and implementation phase of the DSRM, through the explanations of the methods used in data collection, data pre-processing, exploratory data analysis, feature selection, model selection, and the web portal implementation. The selected techniques in each of these sub topics and their suitability for the current study are explained here. The evaluation phase of the DSRM is explained in section 3.6, which starts with a descriptive evaluation technique and their applicability in the current study, design evaluation, algorithm evaluation, and implementation evaluation respectively. The appropriateness of the selected evaluation methods, the customised evaluation process of the study and used techniques in each evaluation type is further described to provide a solid foundation for the evaluation phase of the study.

“Research methodology is a way to systematically solve the research problem” (Kothari, 2004, p. 8). Additionally, a clear understanding of the philosophical perspective of a research methodology not only drives the research in the correct direction but also emphasises its validity in the proper context. As stated above, this study aims to develop a clinical decision support system (CDSS) to assist in predicting the complications of diabetes. This study employs a Design science research approach to answer this research question. However, the main focus of the study is to produce a practical solution to solve the identified overarching research question while confirming the applicability of design science methodologies. The selected research approach of the study is described throughout this chapter to offer a clear understanding of the research design and methodology.

Research philosophy is an essential consideration in the design of a research study, as it shapes how the study is conducted and draws conclusions from the results. Research philosophy refers to the assumptions researchers make about reality, the relationship between the researcher and the research participants, and how research should be conducted (Holden & Lynch, 2004; Stern, 2004; Žukauskas et al., 2018). It is a set of beliefs that guide research practice and inform researchers’ choices, including their methods, the questions they ask, and the data they collect and analyse (Creswell & Miller, 1997; Kothari, 2004). The philosophical foundations precisely explain the overarching research study (Creswell, 2009). The research in IS, and IT can be presented by four principal philosophical foundations: ontology, epistemology, axiology, and methodology. Ontology deals with the nature of being, existence, or reality (Hirschheim et al., 1995). The most straightforward ontological questions are: What exists? How can existing things be known or understood? What is basic, and what is derived? What is reality? (Vaishnavi & Kuechler, 2015). Axiology is also known as the theory of value. Axiology is the analysis of values, which may consist of the values held by individuals or groups (Vaishnavi, 2007). In information systems research, the beliefs in the way of conducting research can be considered the axiological perspective of the research. Epistemology defines the relationship between the researcher and the research objectives. According to Hirschheim et al. (1995) epistemology is “the nature of human knowledge and understanding that can possibly be acquired through different types of inquiry and alternative methods of investigation” (p. 20). The fourth principle, methodology, explains

the ways of developing or constructing. It consists of the process, methods, artefacts, and guidelines for acquiring knowledge (Nunamaker Jr et al., 1990). Sound understanding of the study in perspectives, such as its existence or state of the art, what values can add to the field through the research, what are the objectives of the research are and how they relate to the researcher, and the logical manner of achieving the goals to fulfil the knowledge gap via conducting the research, creates a philosophically enriched research study.

Additionally, a research paradigm is a set of assumptions, values, and beliefs that shape how a researcher approaches a study (Tubey et al., 2015). There are several research paradigms, including positivism, interpretivism, and critical theory, each with assumptions and implications for conducting research. A clear philosophical understanding of research paradigms helps to conduct a research study by selecting the most suitable paradigm. Research paradigms are a collection of beliefs or a set of agreements on understanding the problems, the method of viewing the world, and the ways of conducting research (Creswell & Creswell, 2003; Rehman & Alharthi, 2016; Tubey et al., 2015). It is a framework that guides the research direction and the methods used to collect and analyse data (Oates et al., 2022). Although there are several different research paradigms, the main research paradigms involved in information systems and IT research are the interpretive or constructivist paradigm, the positivist or postpositivist paradigm, and the socio-technologist or developmentalist paradigm (Weber, 2010).

The positivistic paradigm is a leading research paradigm in recent information systems research (Orlikowski & Baroudi, 1991). The core concept of positivism is that the only way to gain knowledge about the world is through observation and experimentation and that this empirical evidence is the only reliable source of knowledge. In positivism, hypotheses define the causality between the variables while mainly using quantitative methods to verify them (Bailey, 2011). Positivism is called empirical science, scientific method, post positivism, and quantitative research (Levine et al., 1987). Interpretivism deals with a deep understanding and exploration of the phenomena relative to the researcher's world. The knowledge developed by interpretivism may be combined with the researcher's experience or may be created towards particular objects. Interpretive research is not focusing on discovering the universal or contextual truth but on keeping a space for understanding individuals'

interpretations of the phenomena. Interpretivism is also known as constructivism, social constructivism, and qualitative research (Levine et al., 1987). The interpretivist believes only a deep understanding and interpretation of a phenomenon can generate authentic knowledge.

Although positivism and interpretivism can adopt almost all the research concepts of the discipline of information systems research, scholars are arguing for a missing paradigm which can lead to building the theories in information systems research (Hevner, 2007; March & Storey, 2008; Weber, 2010). The socio-technologist or developmentalist paradigm is recommended by Gregg et al. (2001) to fill the gap in research paradigms in information systems research. Socio-technology focuses on the development of technology and its effect on individuals and organisations in a positive manner. In the paradigm of socio-technology, the construction and evaluation phases of technology and software are emphasised.

Design science research is derived from the engineering disciplines, where the primary goal is to solve a real-world problem (Hevner et al., 2004; Van der Merwe et al., 2020). The research in this discipline focuses on solving a real-world problem by designing, evaluating, and improving the artefacts (March & Storey, 2008). However, the entire design science research feels like a socio-technology or developmentalist paradigm. At the same time, positivism and interpretivism can be embedded in the research according to the relationship between the researcher and the research objectives (Gregg et al., 2001). This research study mainly focuses on the developmentalist paradigm while integrating positivism and interpretivism as required to strengthen the research study. The thesis explains the use of research paradigms in answering the study's research questions as needed. The following section describes placing the study in its right niche.

3.2 The Rationale for the Approach

Information systems is a field that studies the “effective design, delivery, use and impact of information technologies in organisations and society” (Avison & Fitzgerald, 2003, p. xi). The expanding of information systems over the boundaries of health institutions by managing the data of the field of healthcare is the health information systems (Haux et al., 2004). The

significance of information systems is altering the current situation into an effective, efficient, and preferable one. Developing a HIS to reduce the cost of health expenditures, lessen the consumption of time, improve the health indices of the community, and save the required expert knowledge in diagnosis and prognosis has direct enhancements in the artefacts of the healthcare sector. The ultimate product of this study is an effective information system to support clinical decisions by analysing a longitudinal dataset, where the outcome intends to improve the state of the art.

Research in information systems is fundamentally categorised into two distinct research paradigms: behavioural science and design science (March & Smith, 1995; Wieringa, 2014). Behavioural science has its roots in natural science research methods, while the origins of design science lie in engineering and the sciences of the artificial. Behavioural science aims to develop and justify theories that describe or predict organisational and human phenomena through analysing, designing, implementing, managing, and using information systems. In contrast, design science research concentrates on solving problems through innovations that define the ideas, practices, technical capabilities, and products by analysis, design, implementation, management, and use of information systems which can be effectively and efficiently accomplished (Hevner et al., 2004). The current study leans on design science instead of behavioural science because its sole purpose is to develop an innovative system for supporting clinical decisions rather than explaining a phenomenon through existing theories.

“Design science is the design and investigation of artefacts in context” (Wieringa, 2014, p. 3). Design science consists of two parts—design and investigation—which correspond to two types of research problems: design problems and knowledge questions. Although both research problems aim to enhance the existing organisational or human phenomena, design problem research calls for a change in the real world, whereas the knowledge question research asks for knowledge about the world as it is. Consequently, design problem research results from a design, whilst knowledge question research provides empirical knowledge (Wieringa, 2014). The overarching research question of this study intends to create a solution

for an existing real-world problem through an innovative solution, where the research question leans on the design problem rather than the knowledge question.

This study adopts the design science research methodology (DSRM) since the research aims to design an artefact in the healthcare context. The current study aims to develop a prototype of a clinical decision-making system to assist in predicting the complications of diabetes mellitus (CoDM). The prime intention of the study is to call for an alteration in the real world by analysing actual and hypothetical stakeholder goals, which divert the current investigation into design problem research part of design science. The research methodology is also placed in the DSRM. Although a consensus methodology is lacking in the information research discipline, highly accepted DSRMs have been introduced to information systems (Hevner et al., 2004; Nunamaker Jr et al., 1990; Simon, 1996; Takeda et al., 1990; Walls et al., 1992; Wieringa, 2014). The methodologies introduced by past scholars consisted of various steps and sub-components of the process, which led to a common approach to solving a real-world issue. Fundamentally, the methodological process of design science research can be presented in six main steps: problem identification and motivation, defining the objectives for a solution, design and development, demonstration, evaluation, and communication (Peppers et al., 2007). The following table represents some of the widely accepted methodologies in the field of information systems, with a mapping of the steps mentioned above in the fundamental process.

Fundamental steps of DSRM	(Wieringa, 2014)	(Hevner et al., 2004)	(Nunamaker Jr et al., 1990)	(Walls et al., 1992)	(Takeda et al., 1990)
Problem identification and motivation	Research Context 1. Knowledge goals 2. Improvement goals? 3. Current knowledge Research problem 4. Conceptual framework 5. Knowledge questions 6. Population	Important and relevant problems	Construct a conceptual framework	Meta requirements, Kernel theories	Problem enumeration
Definition of the objectives for a solution	Research Design and Validation 7. Object(s) of study 8. Treatment design 9. Measurement design	Implicit in relevance literature search process, artefact	Develop a system architecture, analyse and design the system	Design method, meta-design	Suggestions and development
Design and development	Inference Design and Validation 10. Inference design		Build the system experiment		
Demonstration	Research Execution 11. What has happened?		Observe		
Evaluation	Data Analysis 12. Descriptions 13. Statistical conclusions 14. Explanations 15. Generalisations 16. Answers	Evaluation	Evaluate the system	Testable design process/product hypothesis	Confirmatory evaluation
Communication	Research Context 17. Contribution to knowledge goal(s) 18. Contribution to improvement goal(s)?	Communication			

Table 3-1 : Comparison of the stages in design science research methodology.

3.3 Reviewing the Process of Design Science

Information systems is an applied research sector which combines the theory from disciplines such as economics, computer science, and social sciences to investigate solutions to a real-world problem at the intersection of information technology and organisations (Peppers et al., 2007). The methodology of examining the answers for the identified real-world problems is directed to the DSRM. As we explained in the above section, the basic six steps of DSRM are widely accepted as a systematic approach to solving a real-world issue. Many scholars invented novel methodologies in information systems research, which are advocated based on that fundamental manoeuvre. Due to the vitality of these six steps in the process of DSRM, we describe them here to review the process of design science.

3.3.1 Six Steps of DSRM

A widely accepted common design science framework is vital for conducting research in the discipline more scientifically while producing recognisable and reproducible results. Researchers from in and out of the domain of information systems, such as engineering (Archer, 1964; Fulcher & Hills, 1996), computer science (Preston & Mehandjiev, 2004; Takeda et al., 1990), and IS (March & Smith, 1995; Nunamaker Jr et al., 1990), define design science research framework guidelines. The most commonly used six activities in the design science research discipline provide direction for a more formative research approach.

Problem Identification

The research problem is defined in this step while justifying the validity of the solution. The researcher understands the research question in its context and recognises the current situation by atomising the problem. The outcome of the problem identification phase in DSRM can be identifying a problem, requirement, or concept that serves as the foundation for designing and implementing an artefact, model, construct, method, theory, or framework (Kuechler & Vaishnavi, 2008).

Defining Objectives of a Solution

Defining the research objectives by utilising the captured knowledge of the problem domain is the second step of DSRM. The main goal of this phase is to explore the identified problem from step one to set feasible objectives of the research by concerning the context and capturing knowledge. The deduced objectives can be quantitative or qualitative, where the former justifies the validity of the goals by comparing existing and novel solutions, and the latter describes the proposed artefact. Knowledge of state of the art and existing solutions should be used in this stage to infer the objectives from the identified problem statement rationally.

Design and Development

Artefacts of the study are developed in this step. The designed artefacts can be constructs, models, methods, or instantiations (Hevner et al., 2004). The design and development of the solution from the defined objectives of the study can be conducted by applying appropriate knowledge of theory. The designed solution should be capable of delivering the desired functionalities. This phase can be subdivided into more discrete activities (Eekels & Roozenburg, 1991) or considered an iterative search process (Hevner et al., 2004).

Demonstration

The developed artefacts should be utilised to demonstrate their usability in the problem context. The demonstration of artefacts varies concerning the proposed solution. Evidence of the accuracy of the artefacts, legitimacy of ideas, and various formal evaluation methods (Hevner & Chatterjee, 2010; Nunamaker Jr et al., 1990) can be used in the demonstration.

Evaluation

The defined artefacts of the study are appraised and evaluated at this stage. The study's objectives are compared with the observed results to assess the developed artefacts. Quantitative or qualitative techniques could be used to evaluate the artefacts based on the nature of the problem venue and the artefacts. Quantifiable system performance measures, such as accuracy, efficiency, response time, and the usage of resources, can evaluate the system quantitatively. In contrast, the feedback of stakeholders, satisfactory surveys, and

results of case studies can be utilised to prove the adequacy of the artefact qualitatively. Based on the impact of the evaluation phase, the researcher may reiterate the process to improve efficiency or communicate the research results.

Communication

It is vital to proclaim the research context, its artefacts, the importance of the study, design and implementation solutions, its effectiveness, and results to the public. Publishing a research paper, presenting the results at a conference, and deploying the system in the real world can be potential communication channels for the study. Communicating the research study is beneficial to the researcher because of its potentiality to receive constructive criticism and the audience being able to update with the state art with the expanded knowledge of the context.

Although these six steps are the backbone of design science research, a more sophisticated research methodology for conducting information systems and engineering research projects has been introduced by Wieringa (2014). The presented method advocated the core concepts of the fundamental approach while enhancing them to strategies in the design science process.

3.3.2 Design Science Research Methodology Introduced by Wieringa

DSRM has evolved over the past few decades as the combined effort of different disciplines, such as computer science, engineering, and information systems management. Although the research field is nourished and prominent, the utility of methodologies used in these disciplines is dispersed within a considerable range. An adapted method of the DSRM for information systems and software engineering has been introduced by Wieringa (2014). The presented concept combines a sensible framework, a systematic process of conducting design science, a detailed description of each step, and a checklist of the process that outstands the methodology itself. The following sections describe the research framework, design cycle, and the stages of the design cycle, respectively, to provide a detailed understanding of the adopted methodology.

3.3.2.1 Design Framework

The research framework guides the researcher in formulating research questions and provides an overall structure with a generalised framework for design science research (Wieringa, 2014). The framework introduced by Wieringa (2014) is compared to a widely accepted research framework which introduces Hevner et al. (2004) to emphasise its specialities. Hevner's framework comprises three components: environment, knowledge base, and information systems research. Additionally, it consists of five main stages: 1) Problem identification and motivation: The problem is identified and defined. The researchers also determine the cause for solving the problem and the research goals. 2) Solution design: The researchers design a solution to the problem. The solution can be a new artefact, a new process, or a new system. 3) Implementation and testing: The researchers implement the solution and test it in a controlled environment. The testing process helps to identify any errors or problems with the solution. 4) Evaluation: The researchers evaluate the solution based on its effectiveness in solving the problem. The evaluation process can involve user testing, expert feedback, and other types of assessments. 5) Communication: In this final stage, the researchers communicate their findings to the relevant stakeholders, including academic publications, presentations, and other forms of dissemination.

The following figure shows the design framework proposed by Hevner et al. (2004).

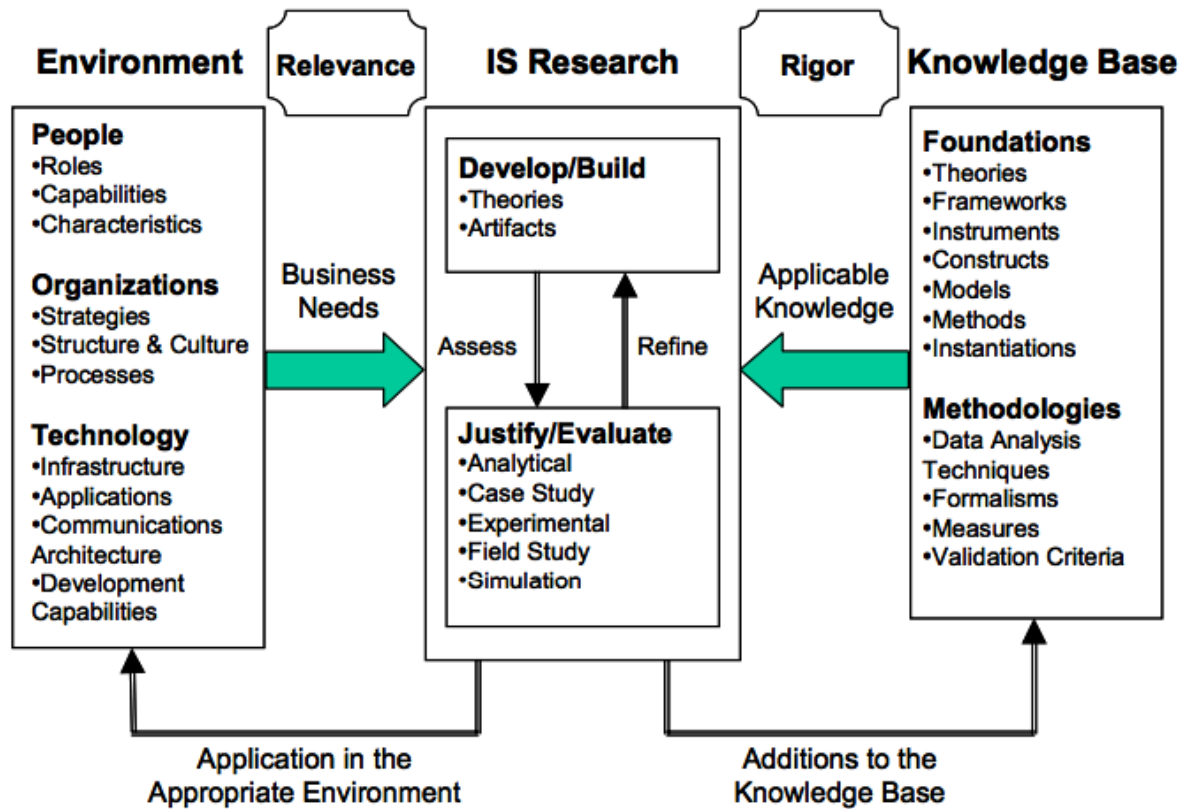


Figure 3-1 : Design science framework introduced by Hevner et al (2004).

The introduced framework divides the components of the environment into three sub-sections: people, organisations, and technology, while the knowledge base is categorised into foundations and methodologies. The information systems research section comprises two main activities: develop/build and justify/evaluate. The developed solutions are assessed through the selected evaluation method and are refined in the developing phase as required. The framework introduced (Wieringa, 2014) also has three main parts: social context, knowledge context, and design research. The following figure illustrates the oriented framework of the process.

Framework for design science

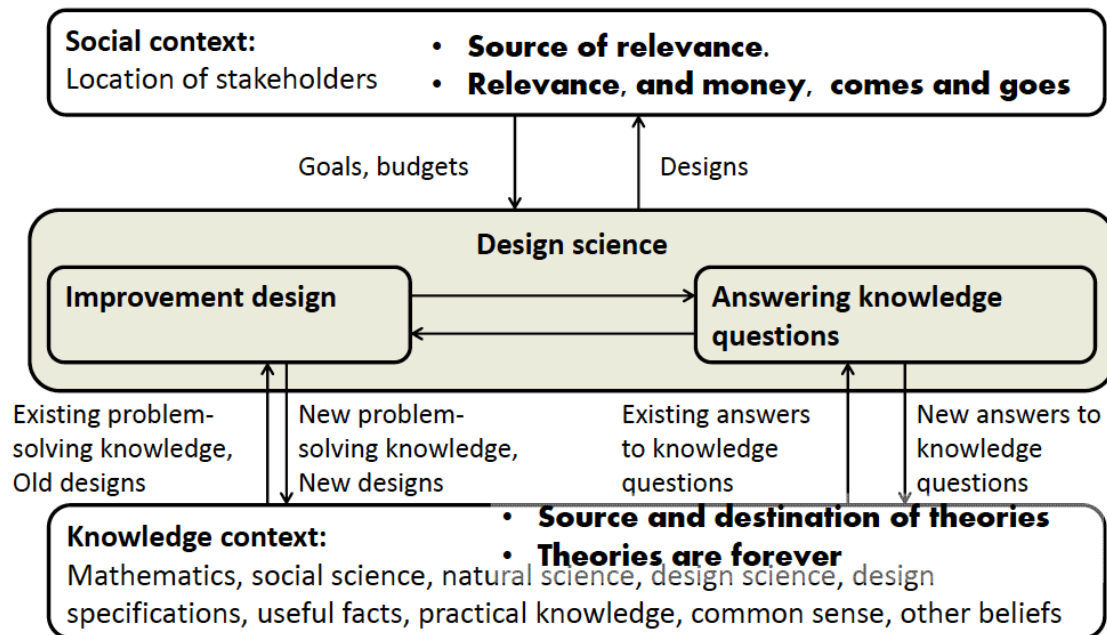


Figure 3-2 : A design science framework introduced by Wieringa (2014).

The problem context of an artefact has been extended into this framework's social and knowledge context. The stakeholders who may affect or are affected by the project create the social context. The users, operators, maintainers, instructors, and sponsors are included in the social context. The stakeholders request the requirements, the goals to be achieved, and the availability of the resources for the artefacts that need to be designed. The knowledge context contains the existing knowledge of the domain, theories, specifications of designs, valuable facts, personal experiences, and common sense.

Additionally, knowledge can be divided into prior and posterior knowledge, where the former is the knowledge before conducting the project, and the latter is the knowledge produced as the result of the project. Design science connects the knowledge context to suit the social context. The design science process can be an improved design or an investigation of a knowledge question. The enhanced design develops using the existing problem-solving knowledge and old strategies, which results in new problem-solving knowledge and new designs. Moreover, existing answers to knowledge questions investigate the answers to knowledge questions which produces new answers to knowledge questions. The developed

artefacts must be evaluated to ensure their validity in a social context. The efficiency of the artefacts, to what extent the artefacts meet the goals, and the accuracy and reliability of artefacts are some measures that can be used in evaluating the artefacts in a social context. Although these two frameworks share more common characteristics and behaviours, the framework introduced by Hevner et al. (2004) is more suitable for research on developing and evaluating innovative artefacts. In contrast, Wieringa (2014) framework is more suitable for research focused on requirements engineering for software systems. Since the outcome of this research study is a software system for predicting CoDM, and the framework emphasises its validity in engineering and computer science, the framework introduced by Wieringa (2014) is a good fit for the study.

3.3.2.2 *The Engineering Cycle*

The engineering cycle provides a structured and iterative approach to design and problem-solving, allowing for continuous improvement and refinement throughout the process. There is a vast range for the terminology used for the outcome; engineers prefer to use the term “solution”, social-scientists prefer “intervention”, and the stakeholders of the healthcare sector prefer to use “treatment”. In this study, we like to use the term “solution” since we are solving a real-world issue from an information system perspective. The engineering cycle is a well-structured problem-solving process consisting of five main tasks, including a design cycle as a sub-part. The five main stages of the engineering cycle are as follows:

1. Problem investigation: What phenomena must be improved? Why?
2. Treatment design: Design one or more artefacts that could treat the problem.
3. Treatment validation: Would these designs treat the problem?
4. Treatment implementation: Treat the problem with one of the designed artefacts.
5. Implementation evaluation: How successful has the treatment been? This may be the start of a new iteration through the engineering cycle. (Wieringa, 2014, p. 27)

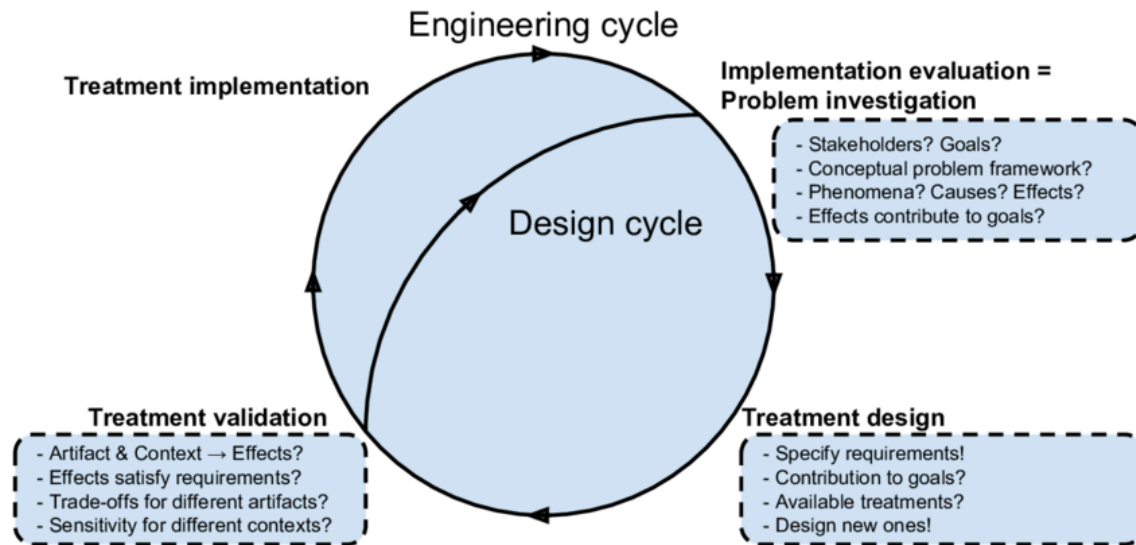


Figure 3-3 : Engineering cycle introduced by Wieringa (2014).

The engineering cycle describes solving a real-world issue with a viable solution deployed in its social context. In contrast, the design cycle focuses on systematising the design process of the conceptual solution, which may not be necessary to implement as a real-time implementation of the solution. The prime intention of the engineering cycle is to create a practical solution for an issue, whereas the design cycle focuses on the process of designing a solution. The study's outcome is a viable information system for assisting in decision-making for the survival of complications of diabetes. Therefore, the current study utilises the engineering cycle to conduct the research instead of limiting it with the design cycle.

Additionally, design research projects are broadly categorised into two types: problem-oriented and solution-oriented. Problem-oriented research projects investigate real-world implementations, whereas solution-oriented projects are technical research, which designs and validates artefacts. In the problem-oriented study, close attention is be given to the problem investigation. In contrast, solution-oriented projects focus more on treatment/solution design and validation. The nature of this study leans towards solution-oriented projects where the end goal is a practical solution. Therefore, the engineering cycle is a good fit for the current study. The stages of the engineering cycle describe the following with an emphasis on their usage in the study.

3.3.2.3 *Problem Investigation*

The first phase of the engineering cycle is problem investigation, where the authors describe the initial step of the development of information systems by focusing on understanding and defining the problem that needs to be solved by the information system. This phase involves gathering information about the current situation, identifying the stakeholders, and understanding their needs and requirements. The researcher understands the research context at this phase by gathering requirements through various channels, including stakeholder meetings, interviewing experts, systematic reviews, surveys, single-case mechanism experiments, statistical difference-making experiments, and observational case studies. The authors emphasise that involving stakeholders in the problem investigation phase is crucial to ensure their needs and expectations are considered in developing the information system. Sound domain knowledge is advantageous for grasping the client's requirements and identifying gaps in the existing knowledge. To systematise the process of problem investigation, Wieringa (2014) provides a set of checklists which is summarised as follows:

- Stakeholders? Goals?
- Conceptual problem framework?
- Phenomena? Causes, mechanisms, reasons?
- Effects? Positive/negative goal contribution?

This phase identifies, describes, explains, and evaluates the problem. A proper investigation of an existing real-world problem will result as the outcome of the problem investigation. This phase is crucial since it provides the foundation for the rest of the project. By thoroughly understanding the problem and the needs of the stakeholders, developers can ensure that the resulting information system meets the needs of the organisation and its stakeholders. A good research question can be derived by answering the checklist mentioned above. Further, it clears out the social context of the framework.

3.3.2.4 *Solution Design*

The main focus of the solution design/treatment design phase is to design the artefacts to treat the problem investigated from the previous stage. The solution should be designed

based on the established theoretical foundations. Researchers use various terminology interchangeably in multiple disciplines, such as solutions, treatments, interventions, or mitigations. In healthcare, the term “treatment” is frequently used to specify the designed solution. As mentioned above, we prefer to use the term “solution” in this study since we answer a real-world problem from a combination of computer science and information system perspectives. Artefacts are considered research outputs (March & Smith, 1995) or the outcomes of design science research projects (Hevner & Chatterjee, 2010). Artefacts are categorised into four broad categories: instantiations, methods, models, and constructs (Gregor & Hevner, 2013).

Further, new clarifying hypotheses, novel models for design and implementation, and processes and approaches for implementation are also considered artefacts (Ellis & Levy, 2010). Wieringa (2014) describes the artefact as something people implement for a practical purpose. The artefacts in the information systems and software engineering discipline include algorithms, methods, notations, techniques, and even conceptual frameworks (Wieringa, 2014). Artefacts are evaluated to prove their efficiency or effectiveness. The designed artefact should interact with the problem context in a way that is intended to treat the problem. Therefore, the “solution” is the interaction between the artefact and the problem context (Wieringa, 2014). There are four components attached to the solution design phase:

- Requirements and context assumptions
- (Requirements * context assumptions) contribute to stakeholder goal
- Available requirements
- Design new ones

The requirements are the characteristics of the solution expected from the stakeholders. The requirement collection and analysis are vital in designing a satisfactory artefact. A specification of an artefact should be released as a contribution argument that justifies the requirements choice. Moreover, the requirements can be classified based on various fundamentals, such as considering their contribution to stakeholder goals according to their priority and urgency (Gause & Weinberg, 1989; Lauesen, 2002). Additionally, classifying requirements as functional and non-functional is a popular way of specifying the requirements, where the former is the functions of the artefact, and the latter is the non-

functional properties of the artefact. Proving the existence of the property of an artefact is called the operationalisation of a property. The operationalisation of a functional attribute can be done by specifying the tests, whereas the existence of non-functional features can be proved with standard indicators.

3.3.2.5 Solution Validation

The validation phase justifies the implemented solution in the problem context to demonstrate its contribution to the stakeholder goals. Validation models are used in design science to simulate implementations. A wide range of validation models can be used to validate the artefacts, resulting in software engineering and information systems research, such as executable specifications, Matlab simulations, software prototypes, user-interface mock-ups, and even role-plays (Wieringa, 2014). There are various research methods to validate the implemented models, such as expert opinions, single-case mechanism experiments, technical action research, and statistical difference-making experiments. The solution validation phase consists of four components:

- (Artefacts * context) produce effects? Why? (Mechanisms)
- Effects satisfy requirements
- (Alternative artefacts * context) produce effects? Why? (Mechanisms)
- (artefacts * Alternative context) produce effects? Why? (Mechanisms)

The evaluation phase is an essential step in design science research, which can be used to validate the research outcome. March and Smith (1995) recognise that evaluation is a major step of design science research. The capability of designed artefacts to alternate the behaviour of systems, people, and organisations is recognised as artefact evaluation by Vaishnavi (2007). Moreover, evaluation is recognised as a critical phase demonstrating the designed artefacts' utility, quality, and effectiveness (Hevner et al., 2004). Evaluation methods of design science research are categorised based on various criteria: ex-ante and ex-post (Pries-Heje et al., 2008), quantitative and qualitative, formative and summative, naturalistic and artificial settings, goal dimension and structure dimension (Prat et al., 2014), and external and internal evaluation. The research artefacts and the process of developing artefacts are widely used to demonstrate the validity of the design science research projects.

3.3.2.6 Solution Implementation

The researcher implements a solution to the problem identified in the problem identification phase to treat the investigated issue. This phase usually involves the following steps:

1. Developing a prototype – The researcher builds a prototype which closely resembles the original solution to test the feasibility and efficacy of the intended solution.
2. Testing the prototype – The created prototype has to be tested to evaluate its effectiveness in resolving the identified issue. These tests ought to be made to assess the success of the implemented solution in achieving its goals.
3. Refining the solution – The researcher develops the solution to increase its efficacy based on the findings of the tests.
4. Implementing the solution – The researcher implements the refined solution in the social context of the research. The implementation should be carefully planned to ensure the solution is successfully incorporated into the current system.

The implemented solution in information systems research may be a software system, a protocol for solving an issue, a solution framework, or a hardware component which resolves an existing problem. The output of this stage is a validated and refined solution to the investigated issue. Specifically, a validated prototype, an advanced solution, an implementation plan, and an evaluation report result from the solution implementation phase of the process.

3.3.2.7 Implementation Evaluation

The implementation of a solution is evaluated after the answer has been applied in the original problem context. The same checklist items used to investigate the problem should be used to assess the improvement. This phase involves evaluating the implemented information system to determine whether it meets the objectives and requirements outlined in the earlier stages of the development process.

The researcher assesses the effectiveness of the implemented solution in this stage. This stage usually involves the following steps:

1. Setting the evaluation criteria – The suitable criteria to assess the effectiveness of the implemented solution is established here. The goals of the project and the theoretical foundations considered in the study have to be kept as the basis for the evaluation criteria.
2. Collecting data – The required data is collected to evaluate the solution. The collected data may be quantitative, such as performance metrics, or qualitative, such as user feedback.
3. Analysing data – The collected dataset is analysed to evaluate the validity of the implemented solution according to the established evaluation criteria.
4. Drawing conclusions – The conclusions can be made from the results of the analysed data about the validity of the implemented solution in its social context.
5. Reporting results – The results of the evaluation of the solution and the derived conclusions are reported as the final step in the solution evaluation phase. Additionally, the recommendation for further solution improvement can be reported here.

The goal of the implementation evaluation phase is to ensure that the implemented information system meets the objectives and requirements defined in the earlier stages, is accessible to the intended users, and is being used effectively and efficiently. By conducting thorough testing, training, and user acceptance testing, developers can ensure that the information system functions as intended and meets the organisation's and its stakeholders' needs.

This study adopts the methodology mentioned above by combining suitable techniques and methods to achieve the study's research goals successfully. The following sections describe how this methodology is adopted in the study. Although all the five main phases which are introduced in the methodology mentioned above were adopted, these phases were combined into three main sections: identifying the problem, solution design and implementation, and solution and implementation evaluation due to the clarity it brings to the study. The problem investigation is described in Section 3.4, the solution design and

solution implementation phases are described in Section 3.5, and the solution validation and implementation evaluation are described in Section 3.6.

3.4 Identifying Situational Requirements and Artefacts

This phase of the study outcomes a clear understanding of the state of the art, client requirements, and the scope of the study. This research study began by reviewing the academic literature associated with risk prediction models, machine learning (ML) models used in the diagnosis and prognosis of diseases, health decision support systems, and the models utilised artificial intelligence (AI) techniques in disease diagnosis. The reviewed literature provided a sound knowledge in the state-of-the-art decision-making support systems. The literature review was narrowed down to understand the current state of predicting DM. This motivated the study to focus on the prognosis of CoDM due to its necessity, as described in the literature review chapter (Chapter 2). This chapter summarises the derived reasons for developing a data analysis model for predicting CoDM from the conducted literature review, that includes the high prevalence rate of DM all over the globe, the increased risk of having one or more of CoDM during the lifetime of a diabetes patient, soaring health expenditures in diabetes, this resulted in a detrimental impact on health indices of individuals and the requirement of immoderate resources in the health care sector. Further, this chapter intends to be aware of the existing pitfalls of the current situation in the research context by deep diving into the findings of the literature review. The situational awareness phase will be enhanced with well structured client meetings and brain storming sessions. The client's requirement set has been planned to identify and clarify through the client meetings and brain storming sessions with the client. Additionally, the requirement specification documents will be used to collect the client requirements. Requirement specification documents are aim for "understand and define what functionalities are required from the software product"(Suárez-Figueroa et al., 2009, p. 966). Moreover, the feature set of the model will be collected after a thorough systematic review. The results of the systematic review may use to effectively communicate with the clients. The identified research gaps of the literature review are used in this chapter to make more scrutinised connections between the research gaps and research goals. The research objectives and artefacts are derived at the end of the chapter.

3.4.1 State of the Art of CDSS

A computerised system that can enhance medical decisions with targeted clinical knowledge, patient information, and other health information to improve healthcare delivery can be considered a clinical decision support system (CDSS) (Osheroff et al., 2012). Although CDSSs can be traced back to the early 1970s, the discipline can still be enhanced in three main pillars: high adaptation and practical use, best knowledge, and continuous improvement of knowledge in the CDSS (Wagholikar et al., 2012). Therefore, academicians and stakeholders in the healthcare sector have been trying to enhance CDSSs by introducing novel data analysis methods, ML and AI techniques, and feeding enhanced datasets.

The enhanced CDSSs are widely used in most healthcare sectors while assisting, guiding, predicting, and diagnosing medical conditions. However, some CDSSs' limitations keep it from reaching the peak of its wide acceptance. The following summarises the existing pitfalls of CDSSs.

- Data accessibility (data acquisition, data transformation, standardisation, temporal data gathering).
- Automating the decision process (ethical issues, safety issues, non-computerised mechanisms).
- Scarcity of consensus among the collaborators of the discipline (subject knowledge gap between the medical practitioners and system designers, lack of up-to-date knowledge of the practitioner).
- Hesitation in accepting novel systems.
- Lack of general agreement on the process of developing CDSSs (the rapidly changing techniques in the field of medicine should be utilised in new systems).
- Unavailability of standards in CDSSs (resulting in more true positives and fewer false negatives are better in some scenarios).

The CDSSs used in predicting diabetes or its complications extends another set of factors that need to be considered in improving the CDSSs of diabetes.

- Asymmetric data (unavailability of data, ethical issues in collecting medical data, scarcity of integrated data repositories)
- Diversity of the research focus (mainly focusing on a few well-known complications of diabetes, observing the effect of a particular factor for developing a complication, ethnicity specific studies)
- Limitation of functionality (only diabetes, capability of predicting few complications, time limitation).
- Considered factors in developing a CDSS (considered period, ethnicity-specific, reveal the effect of one particular aspect).
- The outcome of the CDSS (nature of predicting time, complication types).
- Scarcity of a standard scientific way of developing a CDSS.

The identified knowledge gap of existing CDSSs leads to the design of an advanced CDSS, which should minimise the current issues while providing satisfactory functionalities. The identified gaps divided the issue into two main areas: design issues and data analytics issues.

Developing a widely accepted, rigorous CDSS is challenging for the above reasons. Moreover, the procedures of designing, implementing, and deploying a CDSS make it even more complex. The scientific methodologies can be adapted to the process of developing a CDSS. Although many frameworks (Hevner et al., 2004; March & Smith, 1995; Nunamaker Jr et al., 1990; Wieringa, 2014) are established in the discipline of information systems in developing a computerised decision-assisting system, one standard method cannot be identified as fit for all situations due to the nature of the domain. However, the existing frameworks coincide with each other for leading the way in developing a system that is more straightforward and standardised. The following list shows the general activities for creating a decision-making system in DSRM.

- Requirement gathering
- Requirement analysis
- Design and prototyping
- System development
- Evaluating the system

- Deploying the system (optional)

All the listed steps have myriad ways of performing them, which must be selected appropriately. The process's most critical issue is the selection of the best techniques to make a better model. Since developing a CDSS is challenging due to the scarcity of knowledge in direct methods of choice for the development process, a clear knowledge gap has been derived from the literature on designing a CDSS for predicting the CoDM. The research context of the current study is a combination of the discipline of data analytics and design science research. As described in the chapter on literature review (Chapter 2), there is a significant requirement to assist the clinical practitioners of DM. The asymmetry of available information, lack of concrete methodology in design science research, and the contextual scarcity of a widely accepted model create a space for developing a CDSS. The study mainly focuses on developing a CDSS to assist in the prognosis of CoDM. The four highly prevalent complications—neuropathy, nephropathy, retinopathy, and cardiovascular diseases (CVDs)—are initially selected as the severe complications to prognosis to enhance the usability of the CDSS. Moreover, the suggested CDSS uses the patient's details for a decade to predict the survival and the hazard of the selected CoDM.

3.4.2 Identifying the Feature Set

A significant pitfall recognised in reviewing the state of the art of CDSSs is the considered factors for creating the CDSS. Since the accuracy of a model is highly dependent on the considered features, a suitable feature set is a vital requirement. The feature sets required to predict different complications are also varied. Therefore, selecting the most relevant feature set that predicts each complication is vital. Additionally, it was a requirement to request a functional dataset from client to prototype the model by considering the essential attributes. It was a crucial requirement to conduct a systematic review to identify the most suitable feature set for predicting the selected four significant complications of DM: neuropathy, nephropathy, retinopathy, and CVDs. Since the research topic is directed toward health informatics, perspective of the healthcare stakeholders may leverage the model. Therefore, the planned, systematic review confirms its vitality in enriching the study with a medical stance. The research papers used in the systematic review are published mostly by the stakeholders of medical sciences and health informatics. In this stage, we selected the

complications of diabetes, which has the highest prevalence rate according to the WHO (2016). The result of this systematic review is prepared into a journal article, which is published as a journal article. The systematic paper covers the research articles published in the last seven years, from 2015 to 2023 to extract the most recent researches. The feature set extracted from the literature leads us to request the features from the client. The systematic review revealed 59 features categorised into nine sub-categories: demographic, vital signs, lab orders/values, diagnosis, medication, problem list, family history, bio-sample data, and lifestyle features. The extracted feature set is used to effectively communicate with the client to extract the necessary features for the model. The results of the systematic review paper benefit from selecting feature sets, requesting data from client, being updated with the current situation of the field, and benchmarking the feature sets for future research.

3.4.3 Artefact of the Study

The state of the art of the considered niche clearly shows a requirement for a computerised information system for predicting CoDM by applying data analysis techniques on a diabetes dataset. The overarching research artefact is derived by solving this real-world issue by designing a CDSS which can predict the risk of CoDM through analysing a chronological dataset of diabetes patients. The derived artefact is broadly classified into two branches:

RQ1: How can a CDSS be designed by utilising DSRM?

RQ2: How can the onset of CoDM be predicted using a longitudinal data set?

This classification helps to clarify the overarching situational awareness. The sub-artefacts originated with this classification are

1. A conceptualised framework for the design process of a CDSS to predict the CoDM
2. Usage of data analysis techniques for predicting the CoDM

Moreover, the research questions above are further divided into sub-questions for understanding each derived artefact. Figure 3.4 illustrates the artefacts of the study, with the related research questions. The figure shows how the sub-research questions directly originate the study's artefacts. The proposed artefacts resolve the raised research questions while providing a real-time CDSS for predicting the CoDM. The research outcomes solve a real-world problem through a CDSS while contributing to expanding knowledge of developing

a CDSS in the perspectives of its designing process and its usage of data analysis techniques in disease prognosis.

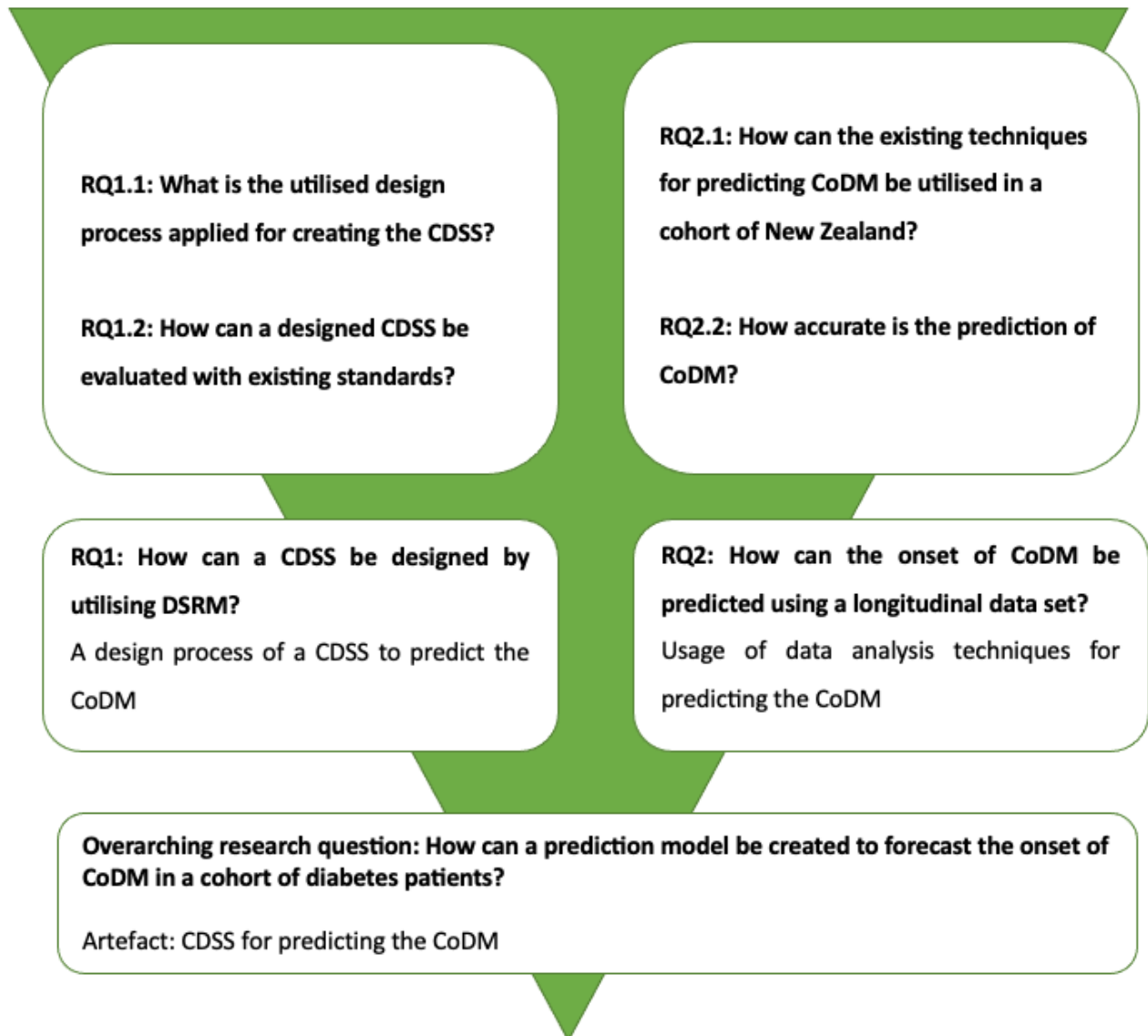


Figure 3-4 : Diagram of presenting the extraction of research artefact through the research questions.

3.5 The Design of the Artefact

The proposed solution is described in this chapter with the details of the methods and their suitability in each step of solution design. The data collection, data pre-processing, exploratory data analysis, feature selection, model selection, and the implementation of the CDSS as a web-based tool will be discussed in this section with their appropriateness for the considered purpose. The methods used in developing the artefacts of the study are described in the following sections consecutively.

In addition, Python is the programming language chosen to be used for the project. Python is a high-level programming language widely used in data analysis, web development, ML and many others (Idris, 2016; McKinney, 2022; Navlani et al., 2021). Python is a popular choice because of its versatility and interoperability. The rich ecosystem of Python (McKinney, 2012) and the available large community of support make Python a good choice for a study which consists of thorough data analysis and website development components. The entire CDSS development including data analysis, model creation, and web site implementation is done by Python.

3.5.1 Data Collection

Data collection is a systematic process of gathering the required information to find the answers to research questions, test the hypothesis, and evaluate the research outcomes. The data collection process should be performed with a clear understanding of the methods and tools of data collection (Goertzen, 2017). Further, the data collection method should be selected by considering factors such as research goals, the scope of the study, sample size, time, and type of data (Creswell, 2009). The accessibility and reliability of the data set, ethical safety, and confidentiality of the subjects must also be considered. Data can be qualitative or quantitative, and the data type should be chosen according to the requirement of the study. Data collection falls under two categories: primary and secondary data collection (Kothari, 2004). Primary data collection is gathering raw data by a researcher for a specific research purpose. Secondary data collection is the process of collecting data from an existing source. Therefore, the collector of the original dataset is not the researcher. The data collection types

mentioned above have pros and cons regarding ease, cost, and value. Table 3.2 illustrates the differences between primary and secondary data collection.

Criteria	Primary data	Secondary data
Definition	Researchers directly collect the data from the sources.	Researchers use the existing data.
Data	Real-time data.	Related to the past.
Process	Collection of data is very involved.	Quick and easy.
Cost	Expensive.	Relatively inexpensive.
Collection time	Long.	Short.
Source	Surveys, interviews, questionnaires, observations, experiments etc.	Books, journal articles, governmental publications, websites, internal records etc.
Reliability	High.	Relatively less.
Specific	Always pertinent to the researcher's need.	May or may not be specific to the researcher's need.

Table 3-2 : Comparison of primary and secondary data collection methods.

This study has to use a secondary dataset due to the inaccessibility, time, and high cost of collecting health-related data. Since the ultimate research objective is to predict the risk of CoDM, this research evidently falls into the quantitative data analysis approach. The preliminary dataset for designing a prediction model is a dataset with patients' diagnoses and test result values with the plan to collect them from the Te Whatu Ora Waikato. The potential communication channels, ease of accessibility, and availability of enriched datasets at the Te Whatu Ora Waikato make it a good fit as the data provider for the study. Although a range of quantitative data collection methods exists, the broad categorisation of it is as probability sampling, interviews, surveys, observations, and document review (Johnson & Turner, 2003). The selection of secondary-quantitative data set in the study leads to collecting the dataset by document review method. A database comprising the details of diabetes patients at the Te Whatu Ora is used to collect the desired dataset. Since the health-related data sets are sensitive and confidential, the data was collected through client meetings. The required data set was requested after collecting the client's requirements. The type of data, length of the dataset, and features needed to develop the solution were requested after a thorough discussion with the client. Since the study leans into health informatics, expertise-knowledge of DM can leverage the model development. This requirement has been fulfilled with an expert team in Te Whatu Ora and through a thorough literature review. The necessity of

expertise-knowledge was saturated from utilising the research articles of the systematic review. The chosen articles for this systematic review are well-known contributors to the field of health informatics. The extracted feature set from the systematic review was used to communicate with the client to request the dataset.

Additionally, another data collection phase has been planned to collect the data for the stage of implementation evaluation. In implementation evaluation, a primary data collection method has to be used. Since the user feedback on the system is collected to evaluate the implemented solution, the sole purpose of this data collection method is to use it for the current study. The evaluation method's use and suitability are explained in Section 3.6, which is the implementation evaluation. The selected evaluation technique requires feedback from the users to validate the design of the CDSS. Interviews, surveys, focus groups, case studies, and observations (Bachiochi & Weiner, 2004) are some primary data collection techniques for feedback gathering. However, users' feedback was gathered to evaluate the implemented CDSS through a questionnaire. Justification of the suitability of these evaluation methods and techniques can be found in Section 3.6, as mentioned above.

3.5.1.1 Ethical Consideration

Ethical consideration in data collection is paramount to the integrity and credibility of any research endeavour. Ethical approval for the study has been granted by the university's Human Research Ethics Committee (HREC), attesting to the adherence to established ethical guidelines. Throughout the data collection process, strict compliance with HREC protocols will be maintained to safeguard the privacy and confidentiality of participants. The WDHB responsibly provides datasets, employing procedures to anonymise and protect sensitive information. Furthermore, the collected data is securely stored in a dedicated Google Drive accessible only to the research team. To adhere to ethical norms and privacy regulations, the data will be retained for a period of five years from the thesis submission date, after which it will be automatically deleted in accordance with stipulated guidelines. This commitment to ethical considerations underscores the dedication to responsible and transparent research practices. The ethical approval received from the University of Waikato's HREC is attached in the appendix section (Appendix A).

3.5.2 Data Pre-Processing

It is crucial to have a quality dataset to generate accurate data models (Alasadi & Bhaya, 2017). The characteristics of a quality dataset include accuracy, completeness, consistency, timeliness, believability, and interpretability (Jiawei Han et al., 2011) . However, most real-world datasets are noisy, inconsistent, consist of missing values, and originate from multiple resources, which leads to a less quality dataset. Data pre-processing enhances the dataset's quality by using various techniques to generate the intended dataset without losing its information (García et al., 2015). The data repository at the Te Whatu Ora is a collection of data resources from multiple hospitals and general practitioners in the Waikato region. The data sources that combine to make the data repository of Te Whatu Ora are vastly different in their data storage format, structures, terminology, security, units of measurement, and equipment used in the same laboratory tests. Additionally, the dataset is highly susceptible to noisy, missing, and inconsistent data. Therefore, a comprehensive data pre-processing phase is required to generate a quality dataset for the study.

The received raw dataset has to transform into a useful, efficient, and quality dataset using pre-processing techniques. Four major tasks are involved in data pre-processing to create a quality dataset: data cleaning, integration, reduction, and transformation. Data cleaning is responsible for filling in the missing values, smoothing the noisy dataset, identifying and removing the outliers, and resolving inconsistencies (Alasadi & Bhaya, 2017). The missing values of the data set can handle by removing the tuples, filling them manually, filling the values with global constant, using central tendency measurement as filling value, using the statistical measures of the categorised class of the tuple, and using the most probable value to fill the missing value (Jiawei Han et al., 2011). The missing values of the received data set must be managed by applying contextual knowledge and statistical measures. Since most health-related datasets show the characteristic of Gaussian distributions, using statistical measures as the missing values minimally affects the statistical analysis results (D. M. C. MIT et al., 2016). The statistical measures were chosen as the filling value and were selected considering the existing outliers of the datasets. Additionally, the tuples consisting of missing values of more than the threshold value can be removed from the dataset.

Data smoothing can be done using various techniques such as binning, regression, and outlier analysis (García et al., 2015). Further, due to the capability of outliers to skew the results and impact the model's effectiveness with anomalies in the raw data, it's vital to detect and remove the outliers (Jiawei Han et al., 2011). The data set was smoothed by outlier analysis due to its effectiveness and minimal impact on the results. The inconsistencies of the dataset should be further resolved. More subjective and contextualised mechanisms can be used in determining inconsistencies (Famili et al., 1997).

Data integration has to be done when the data from multiple resources have to be collectively used in the analysis. Data integration was implemented in making the data frames by integrating the received datasets. Data reduction is the process of making a reduced representation of a dataset which is smaller in volume but can produce the same analytical results. There are three strategies of data reduction: dimensionality reduction, numerosity reduction, and data compression (Jiawei Han et al., 2011). Data transformation is another vital step in pre-processing, where the data are transformed or consolidated to make the data analysis process more efficient. The transformed data leads to efficient and precise data patterns, which create effective data models (Famili et al., 1997). Smoothing, attribute construction, aggregation, normalisation, discretisation, and concept hierarchy generation for nominal data are the most prominent strategies in data transformation (García et al., 2015). Due to the nature of the received dataset from the data provider, a combination of the data pre-processing techniques mentioned above was used accordingly to generate an accurate and efficient dataset. The methods used in each scenario will be explained.

3.5.3 Exploratory Data Analysis

Exploratory data analysis (EDA) is vital in any research analysis to obtain an overall understanding of the dataset. The fundamental aims of performing EDA in research are to detect mistakes, check assumptions, select preliminary data models, determine relationships among explanatory variables, and appraise the relationships between explanatory and outcome variables (Behrens, 1997). The EDA results give a detailed description of the dataset, including their hidden patterns, the nature of variables, existing anomalies, and the relationships among variables (Chatfield, 1986). Due to the importance of being familiar with

the dataset, a thorough phase of EDA was performed in this study. The initial exploration of socio-demographic details of the cohort may reveal a vital set of information of the cohort. The results of the exploration of socio-demographic detail can be used in problem awareness phase, to get an overall image of the cohort, which directly useful in client meetings. This initial exploration can be used to create a clearer communication channel where the stakeholders can facilitate with more satisfied solution. The overall knowledge of the social context is vital in defining the rightmost problem with a feasible designing solution.

3.5.3.1 Exploring the Socio-demographic Details

An overall perspective of a cohort is vital in understanding the existing patterns and behaviours of a dataset. An EDA phase of the prevalence of socio-demographic features was performed in this study as an entry point to the data analysis phase. The collected dataset from the initial stages of client meetings was used to conduct a short data exploration phase to understand the characteristics of the cohort. The gap between the literature and the pragmatic condition of the social context of the study can be filled with a sound initial data exploration phase. The results of this initial data exploration contribute to get a thorough knowledge of the context, which is highly beneficial in defining the correct real-world issue.

Scholars use a spectrum of sociodemographic features to reveal the information of cohorts. For example, level of education, civil status, age, and sex were used as sociodemographic factors in a study finding the association with HbA1c variability among T2DM patients (Mellergård et al., 2020). Willers et al. (2018) used sex, age, educational level, marital status, and region of birth as sociodemographic features when determining the sociodemographic determinants and health outcome variation in a cohort with T1DM. In a study on seeking health-related quality of life in diabetes associated with its social and clinical determinants, age, sex, occupation, education level, and marital status are considered as sociodemographic features to measure the quality of life of diabetes patients (Rodríguez-Almagro et al., 2018). According to Young-Hyman et al. (2016), sex and age are two crucial factors influencing the prevalence of diabetes either independently or as covariates. Further, ethnicity has been widely used as a feature to report and categorise the diabetes cohorts (IDF, 2021; WHO, 2016). Usage of the characteristics of different cohorts in different studies clearly expresses the importance of ethnicity-specific benchmarks. The studies conducted in multi-ethnic

countries usually consider ethnicity as a crucial sociodemographic factor for differentiating and specifying the characteristics of different ethnic groups (Adjei Boakye et al., 2018; Kyrou et al., 2020). Further, Mellergård et al. (2020), pointed out that sociodemographic factors associated with diabetes—age, ethnicity, gender—vary significantly between countries. The exploratory analysis of sociodemographic factors of a diabetes cohort provides insights for the healthcare management sector.

Furthermore, diabetes leads to a range of fatal macro and microvascular complications. The most common complications of diabetes are heart disease, stroke and hypertension, blindness, kidney disease, nervous system disorders, amputation, pregnancy problems, and other complications (Joseph et al., 2022). According to WHO (2016), diabetic retinopathy caused 2.6% of blindness in 2010, whereas 80% of end-stage renal diseases are consequences due to diabetes. Further, adults with diabetes have a two- to three-times higher rate of cardiovascular risk than those without diabetes.

Additionally, lower limb amputation incidents among those with diabetes have drastically increased over the past decade. The data shows that it reached 3.5 events per 1000 persons diagnosed with diabetes per year in 2016 (WHO, 2016). A study of the epidemiology of diabetes and its complications revealed that up to 65% of diabetes deaths are due to CVDs, 44% of new cases of end-stage renal diseases are due to diabetes, and 30%–50% of diabetes patients suffer from peripheral neuropathy (Deshpande et al., 2008). Therefore, the distribution of complications in a diabetes cohort is another way of understanding the underneath image of a cohort. This section of the study focuses on exploring the sociodemographic factors of a diabetes cohort in New Zealand while understanding the distribution of complications of diabetes. The results of this initial data exploration phase may use in effectively communicate with the clients. The patterns of the complications, the distribution of socio demographic details, and missing characteristics of the dataset retrieved through this EDA phase can be discuss with the clients to get a rich dataset which can lead to a more relevant system design. The revealed knowledge of this exploration phase directly contributes on defining the correct real-world issue and designing the solution to cater for the correct issue.

3.5.3.2 Exploring the Dataset with Empirical Standards

EDA is cross-classified in graphical, non-graphical, univariate, or multivariate (MIT et al., 2016). Graphical methods represent the data diagrammatically, while non-graphical methods summarise the data set using statistical measures. The univariate and multivariate analyses focus on exploring a single variable and exploring relationships between multiple variables, respectively. Univariate non-graphical methods statistically describe the dataset's variables (Behrens, 1997). The standard methods for categorical data are tabulating them with statistical measures, and that of nominal data is to describe the statistical features, such as central tendency measures, spread, skewness and kurtosis values of each variable (MIT et al., 2016). Although each measure reveals specific information, the most common statistical measures, such as mean, median, and standard deviation, are used in the study to understand the dataset.

The univariate graphical EDA methods visualise the characteristics of variables. Although univariate non-graphical EDA methods summarise the dataset well in a statistical manner, they are quantitative and objective, so univariate graphical EDA methods are essential in the overall understanding of the dataset. As the visual techniques are more qualitative and have a certain degree of subjective analysis, this is a vital step in EDA (MIT et al., 2016). Histograms, stem-and-leaf plots, boxplots, and quantile-normal (QN) plots are the widely used techniques here. However, due to the overlapping information revealed from histograms and stem-and-leaf plots, it was decided histograms would be used over stem-and-leaf plots. Additionally, “histograms are one of the best ways to quickly learn a lot about your data, including central tendency, spread, modality, shape and outliers” (Seltman, 2018, p. 78) . Moreover, since the QN plots are complicated and the box plots visualise the data dispersion in an understandable and straightforward manner, the study uses box plots over the QN plots.

Multivariate non-graphical EDA techniques were used to show the relationship between variables in the form of cross-tabulation or statistics. “Cross-tabulation is the basic bivariate non-graphical EDA technique” (Seltman, 2018, p. 90) . Further, univariate statistics by category is another well-known method that can be used to find the relationship between a

categorical and nominal variable. This method extracts the statistical measures of a quantitative variable concerning a related categorical variable. The attributes of health-related datasets contain interesting relationships among categorical and nominal variables, leading to performing the EDA method univariate statistics by category. The central tendency measures and the measures of spread are used here due to their simplicity, wide acceptancy, and capability of providing a sound understanding of the attributes. Correlation and covariance matrices are widely used to measure the relationship's strength and direction between two quantitative variables. Covariance measures the difference between two variables, which decided how the change of one variable reflected that of the other variable. Correlation is a measure which determines the degree of change of one variable and reciprocates the degree of change of the other variable (MIT et al., 2016). Although covariance and correlation are vital statistical measurements, correlation was selected for this study due to its tolerance to the change of scales. The multivariate analysis of categorical and numerical attributes can be done by performing a variety of univariate statistics for the numerical feature at each level of the categorical attribute. It shows the statistical differences among different categories of the cohort.

Multivariate graphical EDA methods illustrated the relationships among multiple variables graphically. Side-by-side box and scatter plots were widely used to graphically demonstrate variable changes (MIT et al., 2016). Because of the characteristics of these two techniques, they were both used in this study. Graphical representation of a relationship among multiple variables revealed new information about the dataset. The changes in statistical measurements of a quantitative variable with multiple categories of a categorical variable can be graphed using side-by-side box plots. It revealed a range of information on the data dispersion while comparing that with all categories. Moreover, the scatter plots can represent the relationship between two quantitative attributes, even accommodating one or two additional categorical variables by encoding them by colour and symbol type (Larose & Larose, 2019). Due to the potentiality of having hidden information in the relationship between categorical vs nominal variables and nominal vs nominal variables, both scatter plots and side-by-side box plots were used during the EDA of this study. The results of the EDA with empirical standards, may provide a range of information for system design. The most

prevalent complications among the cohort can be extracted as one of results, of this EDA phase. The system may design based on these extracted complications, which is beneficial in creating a more cohort specific system while providing vital information to discuss with the client. Further, identifying the sociodemographic details in the current cohort and their distributions may use in the CDSS, to give a statistical overview of each complication and general diabetes cohort for the stakeholders. The results of the EDA contribute to making decisions on designing the most feasible system while resulting a set of information to embed in the CDSS.

3.5.4 Feature Selection

Feature selection is an essential step in the process of model creation. Feature selection is selecting the most suitable feature set for the model to perform the desired tasks with the desired quality (Guyon & Elisseeff, 2003). Although real-world datasets consist of many attributes, the feature set of the data models can be a subset of the original feature set with more predictive power. The correct feature set benefits the model through low memory consumption, optimised dataset, and improved model performance (Dong & Liu, 2018). The feature selection technique (FST) chooses the appropriate variable subset, which can produce good results with less complexity (Kotsiantis, 2011). The primary feature selection techniques can be classified into three classes: filter methods, wrapper methods, and embedded methods. Still, the two most popular feature selection methods are filters and wrappers (Guyon & Elisseeff, 2003). Filter methods use statistical measures to select the feature set independently from the learning algorithm (Hall, 1999). Therefore, filter methods can apply as a pre-processing step. Wrapper methods combine the features iteratively to select the most suitable feature set for the model. The different combinations of feature sets create, evaluate, and compare with the model's performance to improve it. Filter methods are less computationally complex and faster than wrapper techniques (Zhang et al., 2019).

The preliminary feature set extraction uses the result of a systematic review, as mentioned above. Filter methods are used in this study as feature selection methods. The popularity, computationally less complexity, and quickness of filter methods in performing the feature selection in real-time datasets make them ideal for similar research projects like the current

study. Information gain and correlation coefficients can be filtering methods to select the features in the pre-processing phases. Those two filter methods are well suited to similar studies due to their characteristics, such as choosing the features due to their importance and their presentation of relationship with the target variable. The feature selection techniques are rarely used in this study due to the limited availability of features of the received dataset. The selected feature set should result in an effective model with high accuracy. However, inputting a feature set for the model without any weight biases is important. Feature scaling can be performed on the selected feature set before feeding them to the model to convert them into the same scale system. Feature scaling is altering the values of numeric columns into a standard scale without distorting differences in the range of values or losing information (Dong & Liu, 2018). The two most commonly used scaling methods are Z-score and Min-Max scaling. Since the attributes are normalised performing the above-mentioned pre-processing steps, feature scaling may not further improve the result. Therefore, feature scaling techniques are avoided in the study.

3.5.5 Model Selection

The sole purpose of this research is to predict the onset of complications in diabetes patients. Since the research goal is aligned with prognosing the onset of different complications of DM, the fundamental data analysis leans towards the survival analysis methods. "Survival analysis is important when the time between exposure and event is of clinical interest" (Ferreira & Patino, 2016, p. 77). Additionally, it is a set of statistical procedures for analysing a dataset where the target variable is time until an event occurs (David & Mitchel, 2012). The event can be death, the onset of a disease, recovery from surgery, failure of a part of a machine, divorce, or re-arresting.

Survival analysis is utilised in various fields, such as medicine, healthcare, engineering, biology, marketing, and social sciences (Tolley et al., 2016). Censoring is a crucial analytical problem in survival analysis that arises when the exact survival time is unknown. There are three significant reasons for censoring occurring in a dataset: the individual does not experience the event until the end of the study, the follow-up of the individual is lost, individual withdraw from the study (David & Mitchel, 2012). However, in this study, the

selected cohort for each complication is extracted by choosing the individuals who experienced that complication after the diagnosis of type 2 diabetes. Therefore, the censoring problem can be avoided in this study. The related terminology, notations and mathematical formula of fundamental survival analysis are explained here to provide a sound knowledge of survival analysis techniques.

T = survival time (T>=0)

t = any specific value of T

S(t) = survivor function

h(t) = hazard function

S(t) = P(T > t)

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

There is a clearly defined relationship between the survivor function and hazard function. If S(t) is known, the corresponding h(t) can be derived, and vice versa.

$$S(t) = \exp \left[- \int_0^t h(u) du \right]$$

$$h(t) = - \left[\frac{dS(t)/dt}{S(t)} \right]$$

Survival analysis has three primary goals: “1) estimate and interpret survivor and/or hazard functions from survival data 2) compare survivor and/or hazard functions 3) assess the relationship of explanatory variables to survival time” (David & Mitchel, 2012, p. 16). The first goal is achieved by calculating the survivor or hazard function using the existing survival dataset. The second goal is achieved through a comparison of two survivor functions. For example, to determine the effectiveness of a medical treatment using a placebo group. The third goal is achieved with mathematical modelling, such as Cox proportional hazard approach (David & Mitchel, 2012).

The survival analysis techniques are categorised into three classes based on the underlying assumptions about the survival distribution: parametric, non-parametric, and semi-parametric (Wang et al., 2019). If the survival time is assumed to follow a known distribution, a parametric survival model can be used in the analysis (David & Mitchel, 2012). Parametric models typically assume a specific functional form for the hazard function, which describes the instantaneous probability of an event occurring given that the subject has survived up to that time. This assumption allows for a more precise estimation of the survival probabilities and hazard rates but also requires that the assumed distribution accurately represents the actual distribution of survival times in the population (Jenkins, 2005). Further, parametric models assume some distributions over the covariates. However, if there is no prior knowledge or theoretical basis for assuming a specific distribution, or if the distribution is unknown or complex, non-parametric or semi-parametric models may be more appropriate for the analysis. Non-parametric models are a class of survival analysis techniques that do not make any assumptions about the underlying distribution of survival times in the population (Moore, 2016). Instead, they rely on the empirical estimation of the survival function, which describes an individual's probability of survival beyond a certain time. The semi-parametric methods assume some distributions on covariates, but they do not make any assumptions on the survival time (Emmert-Streib & Dehmer, 2019). Since the onset of complication cannot be assumed to illustrate a parametric distribution, parametric models have not been used in the survival analysis of the study.

“When modelling human or animal survival, it is hard to know what parametric family to choose, and often none of the available families has sufficient flexibility to model the actual shape of the distribution. Thus, in medical and health applications, nonparametric methods, which have the flexibility to account for the vagaries of the survival of living things, have considerable advantages”. (Moore, 2016, p. 25)

Further, using a semi-parametric model to fulfil a task which has to be done using a parametric model will not significantly change the statistical results since the semi-parametric model can still capture the underlying patterns and relationships in the data while allowing for greater flexibility and accommodating non-linear relationships, which a purely parametric model may

not capture (Kleinbaum & Klein, 1996). However, it is essential to carefully assess the assumptions and limitations of both models before deciding which one to use for a specific task and to correctly interpret and communicate the results to ensure their validity and reliability. When dealing with uncertainty in the data, it may be more appropriate to use a semi-parametric modelling approach, even if a parametric model is deemed the ideal choice based on the dataset's characteristics (Tolley et al., 2016). However, vice-versa is controversial. Semi-parametric models can provide greater flexibility and accommodate non-linear relationships, allowing for more robust and adaptable solutions in uncertain situations (David & Mitchel, 2012). Therefore, using a semi-parametric model is a safer choice.

The general survival curves of the cohort are worth visualising to compare them and have a good understanding of the characteristics of a cohort. Non-parametric modelling techniques have to be used to visualise these general curves since the purpose of the curves is to visualise the changes in survival curves, not to consider the impact of features. A series of graphs have visualised the general curves using non-parametric methods. The techniques used for non-parametric modelling and their suitability are described in the following section. Moreover, the proposed prediction algorithm of CDSSs consists of two components: 1) predicting through demographic details and 2) predicting based on a combination of demographic and laboratory values of the patients. The two algorithms of CDSSs are modelled using semi-parametric modelling techniques due to their flexibility and robustness. The impact of covariates on the survival rate is considered in these two components of the CDSS. Therefore, it is vital to use semi-parametric or parametric modelling techniques. However, as explained earlier, due to the potential uncertainty of the assumptions which has to be made in parametric modelling techniques, both components of the CDSS algorithm use semi-parametric modelling techniques. The used semi-parametric modelling techniques explain in the section after the following (section 3.5.5.2).

3.5.5.1 Non-Parametric Techniques in Analysis

Non-parametric survival analysis does not assume a particular distribution or functional form for the survival function. Non-parametric models help represent a cohort's survival time without considering the covariates' effect on the outcome (Stevenson, 2007). Further, non-

parametric survival analysis techniques are instrumental when the distribution of survival times is unknown or when the assumptions of parametric models cannot be met (David & Mitchel, 2012). They can also be used to explore the relationship between survival times and other covariates of interest, such as age, gender, or treatment status. However, they can be less potent than parametric techniques in detecting differences between groups, mainly when the sample size is small. This technique maps the survival time of a cohort using a simple step function (Jenkins, 2005). The non-parametric methods are very flexible, and the complexity of the model grows with the number of observations. However, it's hard to incorporate the effect of covariates with the non-parametric methods.

Moreover, since their survival function is not smooth, some resulting points of the analysis may not be realistic. There are three basic non-parametric techniques to describe the survival time: 1) the Kaplan-Meier method, 2) the life table method, and (3) the Nelson-Aalen method (Kleinbaum & Klein, 1996). The Kaplan-Meier method is a non-parametric technique used to estimate the survival function based on the observed survival times, and it takes censoring into account. The Kaplan-Meier method is the most straightforward survival analysis technique, which calculates the "survival probability as the number of subjects surviving divided by the number of patients at risk" (Goel et al., 2010, p. 275). The survival function of Kaplan-Meier is a probability distribution that expresses the proportion of individuals who survive beyond a certain time point. When the survival time distribution is unknown or parametric model assumptions cannot be met, the Kaplan-Meier technique is beneficial (Jenkins, 2005). The life-table method is another non-parametric model that also considers censoring. It is commonly used when survival times are grouped into intervals, such as age or time intervals. The life table method calculates the number of people at risk at each interval and then estimates the probability of surviving within each interval. It allows for incorporating censored observations and gives insights into the cumulative survival rate at different time intervals (Moore, 2016). The Nelson-Aalen method is a non-parametric statistical technique used to estimate the cumulative hazard function of a survival analysis dataset. It estimates the cumulative hazard function by calculating the cumulative sum of the hazard rates at each time point (Moore, 2016). The estimated hazard rate is the instantaneous rate at a particular

time, given that the individual has survived until that time. This method is advantageous when the hazard rate is not constant over time.

The Kaplan-Meier survival model was selected in this study to explore the received dataset due to the potential uncertainties embedded in health-related real-time datasets. Further, the flexibility, robustness, ease of use, and availability of many pre-defined libraries in performing the Kaplan-Meier survival technique made it a suitable choice for this study. Further, the ability of the Kaplan-Meier technique to generate clear visualisations of the survival curves of the cohort and the ability to use them to compare the characteristics of the cohort between different groups or periods stabilise the choice of the Kaplan-Meier technique. The most widely used non-parametric method is the Kaplan-Meier survival technique, introduced by Kaplan and Meier (Kaplan & Meier, 1958), also known as the product limit estimator, which estimates the survival time as the product of the number of failure times of the conditional probabilities of surviving to the next failure time. The survival function of the Kaplan-Meier method can be illustrated as follows.

$$\begin{aligned}\hat{S}(t_{(f)}) &= \prod_{i=1}^f \hat{Pr}[T > t_{(i)} | T \geq t_{(i)}] \\ &= \hat{S}(t_{(f-1)}) \\ &\quad \times \hat{Pr}(T > t_{(f)} | T \geq t_{(f)})\end{aligned}$$

This formula is also called the product limit formula since it multiplies the probability of surviving at previous failure time $t_{(f-1)}$ by the conditional probability of surviving at past time $t_{(f)}$, given survival to at least time $t_{(f)}$. The non-parametric survival techniques used in this study focus on exploring the characteristics of each cohort of the complications of diabetes. The exploration of cohorts can reveal valuable information regarding the existing cohorts. The cohorts were exploratorily analysed based on their demographic details, which may be beneficial in understanding the characteristics of cohorts. The exploration is done by graphing the non-parametric survival curves concerning the demographic information. The significance of the curves with strata will be analysed to check the existing differences in the cohort (Jager et al., 2008; Oomichi et al., 2006; Shang et al., 2021). Although the log-rank test, the Gehan-

Wilcoxon test, and the Tarone-Ware test are available in comparing the survival curves, the most commonly used method is the log-rank test (Li et al., 2015). The log-rank test is not only the most widely used technique for comparing the groups in survival analysis, “it has the considerable advantage that it does not require us to know anything about the shape of the survival curve or the distribution of survival times” (Bland & Altman, 2004, p. 1073). The log-rank test was used to compare the existing differences among the strata of the demographic details. The p-value also can be used to check the existence of differences among the curves. This measurement rejects or accepts the null hypothesis. There are few threshold values which can be used in here; 0.1,0.5,0.01, or 0.005, where 0.005 is the vast acceptance threshold in survival analysis (Bland & Altman, 2004). The valuable insights revealed through the comparisons are beneficial in model creation.

3.5.5.2 Semi-Parametric Techniques in Survival Analysis

The semi-parametric techniques in survival analysis are a category of methods which combine the advantages of both parametric and non-parametric techniques. The speciality of the semi-parametric methods is that they are flexible in estimating the underlying distribution of survival times while efficiently and precisely providing the estimates of the relevant covariates (Tolley et al., 2016). Although a range of semi-parametric techniques are available in survival analysis, the well-known methods are the Cox proportional hazard model and the accelerated failure time (AFT) (David & Mitchel, 2012). The Cox proportional hazard model is specifically used in medical and health-related research due to its ability to analyse the effect of multiple predictor variables on the hazard of an event while controlling for other factors that may impact survival. The Cox proportional hazard model does not assume the distribution of survival time, but it assumes that the hazard function is proportional across different levels of a set of covariates (Kleinbaum & Klein, 1996). The AFT is also a semi-parametric model which is widely used in the fields of economics and engineering (Moore, 2016). This method can be helpful when the researcher has a specific parametric form in mind for the distribution of survival times but wants to allow for flexible modelling of the effect of covariates on survival.

Further, AFT assumes a specific parametric form of the underlying distribution of survival time. The purpose of the algorithm of a CDSS is to make a prediction of the survival rate of an individual based on a series of covariates. The Cox proportional hazard model is a good fit in this scenario since it does not assume any parametric form of the survival times while accounting for the impact of covariates in the prediction. Moreover, the ease of interpreting hazard ratios resulting from Cox models and their less computational burden confirms using the Cox model for the study over AFT.

As mentioned earlier, there are two primary components of the algorithm in the CDSS: 1) one where prediction is made only with demographic details and 2) one which considers the demographic and laboratory test results of the patients for the prediction. Both of these components of the CDSS utilised the semi-parametric Cox model for the prediction of survival probabilities. The Cox model is selected due to its capability of considering the effect of covariates on the outcome while not assuming the parametric distribution of survival probability. Additionally, the Cox model is suitable for incidents which use continuous predictors or utilise multiple covariates at once (George et al., 2014). The Cox model is popular among health informatics researchers due to various reasons. A key reason is the “robustness” of the model, “so that the results from using the Cox model will closely approximate the results for the correct parametric model” (David & Mitchel, 2012, p. 110). The hazard function of the Cox model always provides non-negative hazard estimates, which protect the standards, where the hazard function must range between zero and infinity. Moreover, even though the baseline hazard ratio is unspecified, still the hazard ratio can be measured. When comparing the Cox model with logistic models, the Cox model uses more information, such as survival time and censoring, than logistic models, which considers binary outcomes and ignores survival times and censoring (David & Mitchel, 2012).

Semi-parametric methods help analyse the survival data of a cohort where the distribution shape of survival time is not assumed to follow a known shape. Since the Cox model has been selected here for the reasons mentioned above for analysing the dataset, the results of this method answer the remaining part of the first research question.

The formula of the Cox model is written in terms of the hazard model.

$$h(t, \mathbf{X}) = h_0(t)e^{\sum_{i=1}^p \beta_i X_i}$$

Here the $h(t, X)$ denotes an individual's hazard at time t with a set of explanatory variables represented as X . $X = (X_1, X_2, \dots, X_p)$. The survival formula of the Cox model consists of two quantities; $h_0(t)$, the baseline hazard function, which has a time sensitivity, and the exponential component, which involves X 's but does not have time sensitivity. The values of X 's are called time-independent covariates.

3.5.6 Web Portal Implementation as the Deployment

The outcome of this research is a CDSS which can predict a selected set of complications of diabetes. The CDSS is developed as a web portal due to its flexibility, cost-effectiveness, and its potential for the geographical dispersion of the stakeholders. The CDSS was decided to be developed as a web portal instead of a stand-alone system. Web portals have a range of advantages over stand-alone systems (Ginige & Murugesan, 2001):

1. Ease of access – Web portals can be easily accessed from any device through an Internet connection, though the user is not installed on the particular system in their device. This characteristic is vital for systems where the stakeholders are geographically dispersed.
2. Ease of update – The updates are easier in web portals than in stand-alone systems. Since the web portal updates are centralised, the alterations are readily available for all users.
3. Scalability – Web portals can quickly adapt to the growing number of users or data loads without a significant hardware upgrade or re-architecture.
4. Integration – Integrating web portals with other systems and data sources is more manageable than with stand-alone systems.
5. Lower costs – Development and maintenance is more effortless in web portals due to the availability of a large pool of developers familiar with web technologies.

This section describes the implementation of the front-end development of the web portal since all the other system components have already been discussed in the previous chapters. The process of front-end development of the web portal has been done according to the standards of usability, which ensure that the web portal is user-friendly and easy to navigate. Although the standards of the front-end development process—fast and responsive (even on slower Internet connection), responsive designs (optimised for various screen sizes and devices), and accessibility to all users (including individuals with disabilities)—are significant in web portal development (Pastore, 2012), the purpose of the planned stage of the web portal emphasises the standard of usability. The future versions of the web portal may consider improving its characteristics against more standards. The resulting web portal of this project serves as a risk prediction tool that anyone can access through the Internet. Since it is impossible to give training to the potential users of the web portal, this should be a user-friendly and easily navigable design. The process of the user interface designing of the web portal is described here. A software design life cycle (SDLC) was conducted in the web portal designing phase due to its ease of implementation, stability, and acceptance as a software design process. Among the available methodologies of SDLC, such as waterfall, Agile, rapid application development (RAD), and Spiral, the waterfall model was selected to implement the web portal. Waterfall methodology is a linear process with a strict sequence of phases to follow one after the other (Sommerville, 2011). The requirements of this web portal are almost understood, and due to the client's requirements not rapidly changing, the waterfall method can serve the primary purposes of the project. Moreover, the advantages of the waterfall method, such as its clear structure, well-defined phases of the process, and the ease of the process to manage and control, made it an ideal methodology for this project. The waterfall method consists of seven significant phases: 1) requirement gathering, 2) analysis, 3) design, 4) implementation, 5) testing, 6) deployment, and 7) maintenance (Sommerville, 2011). Developing the CDSS with these seven steps is described here.

The requirements of the study are gathered from resources such as reviewing the literature and meetings with clients. The background research conducted for this phase is discussed in Section 3.4. The gathered requirements are analysed to examine their feasibility. Moreover, the requirements are analysed to make the functionalities of the web portal. The

requirements and use cases of the web portal are mentioned in Section 3.4, and the design of the web portal is explained in a one-page standard to retain simplicity. The web portal has been designed as a service provider website, focusing on user-friendliness. The sketch of the website was created by considering the easy navigation, simple layout, and minimal options. The sole purpose of the web portal is to provide a user interface for the stakeholders, such as healthcare professionals, GPs, and diabetes patients, to assess the risk of having complications of diabetes.

The algorithm for predicting CoDM has been chosen from the survival analysis techniques. A detailed description of these algorithms is mentioned in Section 3.5.5. The whole project is developed using Python programming language as described above. The user interfaces are designed with the “Streamlit” library of Python. Streamlit is popular among data scientists and developers (Gopiseti et al., 2023; Zhou et al., 2023). It is an open-source library that is easy, fast, and flexible for creating interactive web applications. The conventional ways of creating websites, such as HTML, CSS, JavaScript, Angular, and Wix platform, are ideal for projects where the primary purpose is developing a website. The web portal of this study is created to display the results and provide a service for assessing the risk of complications of diabetes. The sole purpose of the web portal led to deploying the system as a simple, service-provider web portal which can be easily developed with Streamlit. Data visualisation, developing interactive graphs, and making simple, user-friendly interfaces are the primary concern of developing the web portal. The Streamlit library perfectly matches the project's scope while providing a free web hosting service. The whole website implementation has been done with Python with the Streamlit library. The testing phase of the project is described in the following section. The implemented website has been hosted in the Streamlit cloud environment. The maintenance of the website will be continued through a GitHub account, where all the resources and applications are stored. The Streamlit cloud is a free environment where developers can host their applications built by Streamlit quickly. Moreover, they provide a range of features such as automatic scaling, collaboration tools, and version control. However, the data and source code of the Streamlit cloud are open-sourced. Due to the agreed data confidentiality and possible threats in public host servers, alternative website host servers have to be used to launch this site. Most of the web site hosting services need

subscriptions: such as AWS, Google cloud hosting, and GoDaddy. Although some website host services are free: WordPress, Wix, and Squarespace some intolerant functional limitations are embedded with them. GitHub has been selected in this research due to its characteristics such as offering a set of hooks for continuous integration and automation, tags and semantic versioning, branching and forking (Perez-Riverol et al., 2016; Pipinellis, 2015; Utomo, 2020). Although Source Forge, GitLab, and Bitbucket are some of the alternatives for GitHub, the most widely used open-source platform has been selected in the study. Further, the ease of manipulating the account, the large community of support, the ability to version control, and the ease of integration with other tools and services (Pipinellis, 2015) made GitHub the ideal web-based platform for maintaining the web portal.

3.6 Evaluation of the Artefact

The evaluation phase is considered a crucial and essential step in design science research. The importance of evaluating artefacts is highly emphasised by many scholars due to the values it can bring to a study (Hevner, 2007; March & Storey, 2008; Vaishnavi, 2007; Wieringa, 2014). The evaluation phase of design science research is mainly concerned with design science outputs, which are theories and artefacts of the design science field (Pries-Heje et al., 2008). The research paradigms, such as positivism, interpretivism, and critical research paradigms, do not focus on designing and building new artefacts (Cleven et al., 2009). In contrast, the focal point of the research in design science is to solve a real-world issue with a novice artefact (March & Smith, 1995; Vaishnavi & Kuechler, 2004; Venable et al., 2012; Wieringa, 2014). Therefore, evaluation is vital in the discipline of design science research. The well-established framework of design science research introduced by Hevner et al. (2004) comprises two cycles: relevance and rigour, which can be considered as the evaluation phases to determine the relevance and rigorousness of the designed artefacts. A vast number of categorisations of evaluations can be seen in the literature on design science research: 1) formative vs summative (Remenyi & Sherwood-Smith, 2012; Wiliam & Black, 1996), 2) ex-ante vs ex-post (Bannister & Remenyi, 2000; Irani & Love, 2002; Keast, 2004), and 3) naturalistic vs artificial (Sun & Kantor, 2006; Venable, 2006). Moreover, a set of variables and values for the evaluation of design science research artefacts are categorised by Cleven et al. (2009), which has 12 variables with corresponding values: **Approach**: quantitative, qualitative; **Focus**:

technical, organisational, strategic; **Type**: construct, model, method, instantiation, theory; **Epistemology**: positivism, interpretivism; **Function**: knowledge function, control function, development function, legitimisation function; **Method**: action research, case study, field experiment, formal proofs, controlled experiment, prototype, survey; **Object**: artefact, artefact construct; **Ontology**: realism, nominalism; **Perspective**: economic, development, engineering, epistemological; **Position**: externally, internally; **Reference point**: artefact against research gap, artefact against the real world, research gap against the real world; **Time**: ex-ante, ex-post. Although the evaluation phase has been emphasised as a crucial step in design science research, a systematised method of choosing an evaluation strategy, the properties of evaluand and evaluation, and the specific evaluation period remain indeterminate. A well-structured framework (FEDS) for the assessment in design science research has been introduced by Venable et al. (2016), which provides clear guidance for evaluating the design science research to a fair extent. The FEDS framework introduces four basic strategies of evaluation based on two dimensions: functional purpose: formative and summative, and paradigm of the evaluation: artificial and naturalistic. Figure 3.5 illustrates the strategies of the FEDS framework.

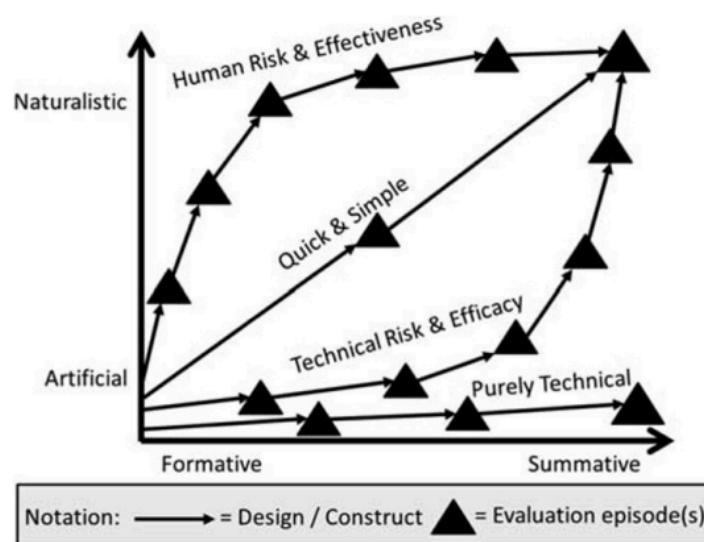


Figure 3-5 : FEDS (Framework for Evaluation in Design Science) (Venable et al., 2012).

The four strategies introduced in FEDS cover all the research in the design science research paradigm while keeping space for a potential hybrid approach. The formative evaluation is a process that contributes to improving the characteristics or performance of the evaluand

(William & Black, 1996). In comparison, summative evaluation intends to validate the generalisability of the evaluands for future applications or decide their appropriateness in the current application (Remenyi & Sherwood-Smith, 2012). Additionally, artificial evaluation is a method of proving the reliability of the artefacts through repeatability and falsifiability (Gummesson, 2000). Laboratory experiments, simulations, criteria-based analyses, theoretical arguments, and mathematical proofs are well-known techniques in artificial evaluation (Pries-Heje et al., 2008). Naturalistic evaluation is the method of testing the performance of the artefacts in their natural context. It explores the complexity of human practices in the relevant context (Nunamaker Jr et al., 1990). Case studies, focus groups, field experiments, surveys, ethnography, phenomenology, hermeneutic methods, and action research are popular naturalistic evaluation techniques (Venable et al., 2012). The evaluation strategy of the current study has been chosen by considering the characteristics of the dimensions of the FEDS framework. Moreover, the five-step process (Wieringa, 2014) that was adopted as the framework of the study also consisted of two evaluation phases: design evaluation and implementation evaluation. The current study used the strategy of FEDS, known as “technical risk and efficacy”, for implementing the evaluation phases. The quick and straightforward approach comprised low cost and quick project conclusion, which may result in high design risks (Venable et al., 2016). This research study intended to reduce the design risk of the project due to the possible burden it can bring to the study. The “human risk and effectiveness” evaluation strategy has not qualified for the current study due to its high user involvement requirement. The “purely technical artefact” method is suitable for purely technical research without the involvement of human users (Venable et al., 2016). The solution deployment should be far in the future for the studies suitable for this evaluation strategy. The current study leans towards technical research rather than interpreting complex human practices. Moreover, deploying the system in a natural setting and evaluating it with real users is not the focus of this research. The sole purpose and the available resources of the study lead to assessing the study through the strategy of “technical risk and efficacy”.

“The Technical Risk & Efficacy evaluation strategy emphasises artificial formative evaluations iteratively early in the process, but progressively moving towards summative artificial evaluations” (Venable et al., 2016, p. 82). The trajectory of the strategy has been constructed

with two evaluation episodes as described in the five-step DSRM: design evaluation and implementation evaluation. Moreover, these two phases mirror the two cycles, as mentioned earlier in the design research approach introduced by Hevner (2007). The rigour and relevance cycles of the process have similar purposes as the design and implementation evaluations. The relevance cycle evaluates the relevance of the research to its context, while the rigour cycle evaluates the added knowledge to the field (Hevner, 2007).

3.6.1 Evaluation Process of the Study

The approach to the evaluation phase of a particular design science research project has been identified by Venable et al. (2016) as a four-step process: “(1) explicate the goals of the evaluation, (2) choose the evaluation strategy or strategies, (3) determine the properties to evaluate, and (4) design the individual evaluation episode(s)” (Venable et al., 2016, p. 6). The study uses the FEDS evaluation strategy with the desired steps. First, the goals of the study’s evaluation are identified and categorised to generate the customised trajectory for the evaluation process of the study. The identified purposes in this study are categorised based on the evaluation period, such as ex-ante and ex-post. The objectives of the ex-ante phase are to validate the suitability of the artefact. The appropriateness of the suggested solution is evaluated at the beginning of the project to validate its relevance status and minimise human/social use risks. Moreover, assessing the suitability of the techniques used to build the solution is another goal set at this phase to reduce the project's technical risk and ethical constraints. The performance of the artefact and its user satisfaction are the evaluation goals in the ex-post stage. The overall efficacy and efficiency of the project are the intended goals at this phase. A series of formative evaluations are conducted during the project's implementation to test the suitability of techniques used in the prediction algorithm. This formative evaluation aims to check the efficiency of the methods used.

The technical risk and efficacy evaluation strategy has been selected as the evaluation strategy of the current study due to its appropriateness, as mentioned above. Two summative evaluation phases were conducted before and after the study to evaluate the validity of the designed artefact and to evaluate their performance of them at the beginning and after the study, respectively. Design evaluation focuses on assessing the appropriateness

of the suggested artefact as a real-world solution. Implementation evaluation concerns the performance of the artefact in the defined context at the end of the project. Several formative assessments have been conducted during the implementation of the artefact. The algorithms used to predict each complication are evaluated at the end of their execution to check their accuracy.

Determining the properties to evaluate is another vital step in assessing the project. The property set for the first goal of the ex-ante phase is its contribution to its social context. The dimensions of ISO 25010 standards are considered as the properties of the artefact to evaluate the CDSS at the ex-post phase of the evaluation. The dimensions used in ISO 25010 are considered here to measure the overall quality of the artefact through its functionality, reliability, usability, efficiency, maintainability, and portability (ISO/IEC 25010, 2011) . The accuracy of the techniques used for prediction is considered the potential property to evaluate in the phases of formative evaluation.

The last step of the evaluation process is to design the individual evaluation episodes. This study's evaluation was conducted in three types of episodes: design evaluation, algorithm evaluation, and implementation evaluation from the beginning to the end of the study. The design evaluation consists of a single episode. In contrast, the implementation evaluation phase consists of two evaluation episodes, and the algorithm evaluation is conducted after implementing each algorithm for predicting the complications of diabetes. The utilised evaluation process of this study is shown in the following figure, including its strategy and episodes.

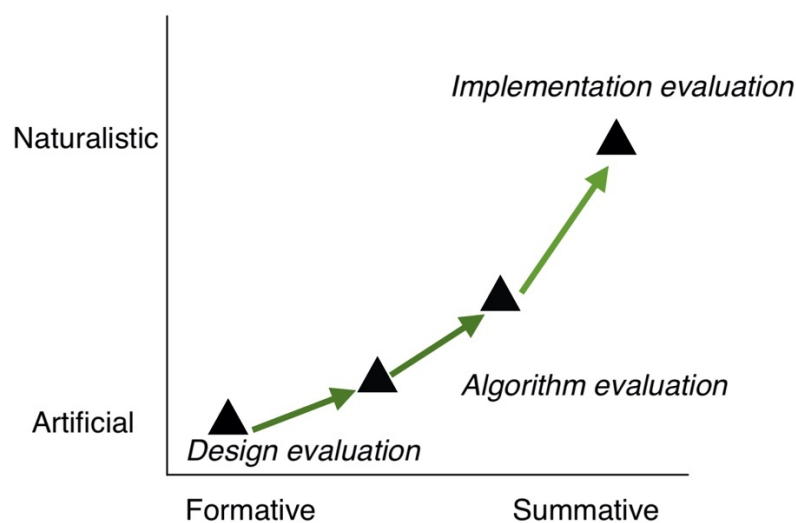


Figure 3-6 : Diagram of the evaluation process of the study.

3.6.2 Design Evaluation

A CDSS has been generated as the solution for predicting the complications of diabetes by analysing patients' datasets. Model validation is essential in empirical model designing (March & Smith, 1995). It is the process of determining how accurately a model can simulate the natural environment (Cleven et al., 2009). Clinical prediction model evaluation has evolved as a successful branch of model evaluation, systematised to improve the prediction models in the healthcare sector (Debray et al., 2013; Han et al., 2016; van Smeden et al., 2021). Furthermore, it introduces seven steps for developing a clinical prediction model, including model estimation, model performance evaluation, and internal validation.

The revealed requirements of the client and the existing research gaps in the field led to a CDSS being built to fulfil the client's requirements while addressing the current research gaps. The design evaluation phase was conducted at the beginning of the project to validate its suitability for the real-world issue. The goal of this was to reduce human social/user risk. The design evaluation phase leaned towards the summative evaluation since it concerned the overall suitability of the solution to the identified real-world issue. According to the method introduced by Wieringa and Morali (2012b), design validation is a process that is performed before implementing the actual system. This method concentrates on four components: 1) expected effects, 2) expected value, 3) trade-offs 4) sensitivity. The suggested solution is

validated against the gathered requirements for the above four components. The artefact's effects in a problem context are validated under the component of expected effects. The effect of the suggested solution in the problem context is validated through the requirements. The second component (expected values) is checking "how well will these effects satisfy the criteria?" (Wieringa & Morali, 2012b, p. 227). The utility, style, quality, and efficacy are considered essential criteria for evaluating the artefacts (Hevner et al., 2018; Hevner et al., 2004; March & Smith, 1995; Prat et al., 2015; Venable et al., 2016). Furthermore, Hevner et al. (2018) provide an extended version of the criteria as functionality, completeness, consistency, accuracy, performance, reliability, usability, fit with the organisation, and other relevant quality attributes. To what extent these criteria are satisfied through the effects of the artefacts is analysed here. The satisfaction of the effects of the artefacts is validated through the selected criteria. The base criteria for this study were chosen to match the client's requirements. Therefore, the specified criteria are user-friendliness, inexpensiveness, and accuracy. The above-mentioned standard criteria were embedded into the selected criteria. The next step was evaluating the design solution to consider the potential trade-offs and the comparison between them. The trade-offs were compared through brainstorming and discussion groups with the client. The last component of the process was to consider the system's sensitivity. It measured the system's ability to respond to the changes in the problem accurately and handle more stakeholders, data, etc. (Wieringa & Morali, 2012b). The factors such as data quality, algorithm, and clinical knowledge base are (Wieringa & Morali, 2012b). The factors, such as data quality, algorithm, and clinical knowledge base, were considered thoroughly to achieve the system's sensitivity. The data's accuracy, completeness, consistency, relevance, and validity were selected as the quality parameters for checking the data quality (Hasan & Padman, 2006). The data quality was checked with the Python programming code since the whole project was developed with Python. The algorithm of the design solution was considered to preserve the system's sensitivity.

Although the algorithm can be improved through various factors, such as feature selection, tuning threshold, regular updates and maintenance, incorporating expert knowledge, model calibration etc., the scope of this project selected the factors of feature selection, model tuning, and calibration as the relevant factors in enhancing algorithms for preserving the

sensitivity. The feature selection was covered with a thorough systematic review to fill the knowledge gap from the medical expertise and updated with the state-of-the-art in the field. Since the Cox model does not have traditional hyperparameters, the model tuning can be employed with regularisation methods. The contribution of the clinical knowledge base of the system plays a significant role in preserving the system's sensitivity (Olson et al., 2010). The root of developing specific risk prediction tools, such as UKPDS, FINDRISC, ASCVD Risk Estimator Plus etc., is characteristic of the model to be sensitive for the trained dataset. It has been recognised as a challenge to create a widely accepted model which is not sensitive to its knowledge base. The knowledge base of this study was trying to be sensitive through the used model types in the knowledge base.

Overall, the approach to evaluating the design solution as explained above is closer to the characteristics of artificial evaluation than natural evaluation. Due to time constraints, limitation of resource accessibility, and the potential complexity that the natural evaluation methods can bring to the study, the assessment of the design solution leaned towards the artificial evaluation method. Further, to be aligned with the above-mentioned evaluation strategy of the study, the beginning has to be more artificial than natural. The design evaluation phase used the criteria-based artificial evaluation strategy to reduce human social/user risk. Evaluating the suggested solution, related to its relevance to the user's requirements, minimises the project's human social/user risk.

3.6.3 Algorithm Evaluation

The algorithms developed for predicting the complications of diabetes were evaluated to check their ability to predict the survival rate of the individuals correctly. The selected survival analysis models have embedded evaluation criteria to validate their accuracy based on well-defined statistical techniques. Although the statistical methods for evaluating the models are well-accepted and established, there is some controversy in validating the Cox models. "Validating Cox models is not straightforward because event probabilities are estimated relative to an unspecified baseline function" (Royston & Altman, 2013, p. 1). However, the concordance index is the most popular criterion for evaluating the performance of the Cox model (David & Mitchel, 2012). "It measures the ability of the predictor to order the events

by estimating the fraction of correctly ordered pairs out of all comparable pairs in the dataset” (Alabdallah et al., 2022).

Moreover, the Brier score can be used to calibrate the models. It is not only used for model calibration but also for model discrimination. Since the C-index values were used in the study to evaluate the model's performance, the Brier score is used as a calibration method. The widely accepted calibration score for the survival analysis method is the Brier score (Heller, 2021). Additionally, the Brier score is a probability-based metric that assesses the calibration of predicted probabilities. In survival analysis, the primary goal is often to estimate the probability of an event occurring at a given time. The Brier score directly evaluates the accuracy of these predicted probabilities, making it relevant and interpretable in the context of Cox models (Gneiting & Raftery, 2007b). The Cox proportional hazard models in the current study basically use the concordance index and Brier score as their discrimination and calibration measurements. The C-index was used here to evaluate the performance of each prediction model. As mentioned earlier, the existing controversy of using the C-index for assessing the performance of the Cox models leads to a search for alternative methods of evaluating the performance of Cox models. Therefore, alternatives for determining the Cox model have recently been popular among academicians.

The internal and external ways of assessing the Cox models become prominent to customise the evaluation methods according to the user requirements. The external validation is checking the model accuracy with an external dataset. The internal validations are done with the collected dataset for the same study. The available resources, such as the availability of datasets, limited time frames, and associated costs, led this study to engage with internal validation methods. The evaluation method of the Cox models in this study used the concepts of cross-validation. Cross-validation is a well-established statistical method of evaluating and comparing algorithms, which divides the dataset into training and testing datasets, where the model is trained with the training dataset and tested against the testing dataset (Witten et al., 2005). This validation process is repeated multiple times with different subsets of the dataset used for training and testing to get a more reliable estimate of the dataset. The current study used C-index and Brier score to validate the performance of the models. The

model was trained and tested with two splits of the original dataset. The resulting survival curve of each dataset was checked with their log-rank values. The process was repeated 10 times to perform the 10-fold cross-validation techniques to evaluate the algorithm's performance. The algorithm evaluation was considered a formative evaluation method since it focuses on the improvements of the artefact.

3.6.4 Implementation Evaluation

Implementation evaluation was done at the end of the project to evaluate the developed system's quality. This summative evaluation phase evaluated the artefact at the end of the implementation process. Moreover, "summative evaluation is conducted to evaluate the efficacy of the final design or to compare competing design alternatives in terms of usability" (Hartson et al., 2001, p. 2). The evaluation of an artefact can be performed in various standard methods, broadly classified into user-based, expert-based, and model-based evaluation techniques, consisting of various pros and cons. "From the CDSS implementation point of view, authors combined quantitative and qualitative assessments to determine the system usability through questionnaires, ethnographic studies, group meetings, and individual interviews" (Souza-Pereira et al., 2021, p. 6). Moreover, the evaluation of IT artefacts is considered a combined mathematical and product quality evaluation process by Hevner et al. (2004). "IT artefacts can be evaluated in terms of functionality, completeness, consistency, accuracy, performance, reliability, usability, fit with the organisation, and other relevant quality attributes"(Hevner et al., 2004, p. 85). To comply with the literature, the implementation evaluation phase of this project was adopted as a combination of model-based and user-based techniques. The developed artefact was evaluated in two phases, where the product quality was measured as model-based standard evaluation and user feedback collection. The model-based evaluation method was selected for this study due to its benefits of evaluating the whole project, including its products and process. The evaluation of the artefact through user feedback was used to confirm the product quality further. The user-based and expert-based evaluation methods are more time-consuming and expensive (Dillon, 2001) than model-based techniques. A systematic review by Paz and Pow-Sang (2016) extracted the most frequent techniques for evaluating health informatics applications. They have concluded that survey/questionnaire, user testing-thinking aloud, user testing, and

heuristic evaluation as the most common usability evaluation methods of applications in health informatics. Although conducting a survey or questionnaire is appropriate to get user feedback, the influence of the implementor is unavoidable. Further, the high cost and time are other possible obstacles with surveys/questionnaires. The user testing-thinking aloud and heuristic evaluation methods comprise careful monitoring of users and high user involvement, respectively. The model-based evaluation technique was selected as the core technology for the implementation evaluation due to its suitability as less expensive and short time-constrained characteristics. The implementation evaluation is a summative evaluation process where the standards of ISO 25010 dimensions were adopted to evaluate the CDSS. ISO 25010 is a standard matrix for evaluating the quality of software products along with system quality in the used model (International Organization for Standardization, 2011) . The main criteria for evaluating the product quality of the model in this standard is categorised into eight characteristics: functional suitability, performance efficiency, compatibility, usability, reliability, security, maintainability, and portability. Additionally, these criteria are divided into 26 subcategories to evaluate the model further. The evaluation criteria express the existing characteristics of the developed model. ISO 25010 was utilised in this study to evaluate the product quality using the system requirements. System evaluation through ISO 25010 has been widely used in evaluating developed systems through various approaches. The evaluation of CDSSs through this standard matrix is commonly used in scenarios where the evaluation leans towards artificial methods rather than natural/ human-centred evaluation, which perfectly matches the purpose of the study. Quality evaluation of a cardiac decision support system was conducted using the standard ISO 25010 with a thorough evaluation of system qualities (Kadi et al., 2016). The characteristics and sub-characteristics of the system were evaluated against the standard measures using three simple mathematical expressions. The evaluation process of the study conducted by Kadi et al. (2016) was adopted in this research. Although the purpose of Kadi et al. (2016) is to recognise the most privileged characteristics and sub-characteristics of the ISO 25010 standard, their approach was used in a retrograded manner, where the purpose was to prove the existence of the standard characteristics in the developed system. Although the product quality can be validated through artificial techniques, the nature of evaluating the design science projects was a mixture of artificial and natural evaluating techniques. It was important to validate the

implemented artefact using both evaluating techniques to confirm its qualities. While the above-mentioned artificial evaluation technique evaluated the product quality with industrial standard measurements, the users' feedback was collected to evaluate it in a more naturalistic manner.

Since the resulting artefact of the study was a usable real-world solution, it was essential to evaluate the practical use of the system. Moreover, research focusing on solving practical problems should allow collaboration and cooperation in research to bring a naturalistic perspective to the solution (Siau & Rossi, 2011). The system analysis and design methods can be empirically evaluated through observations and propositions based on sensory experiences. Validating a system's user experience is a well-accepted method in empirical model evaluation. According to (Peppers et al., 2007), the practical use of the artefact can be differentiated into demonstration and assessment. The evaluation of the artefact is defined as demonstrating the utility of the artefact (Hevner et al., 2004). A CDSS developed for diagnosing a pulmonary disease was evaluated in two phases: algorithmic and clinical evaluation. The latter assessment was done with a questionnaire to assess user satisfaction, qualitative improvement, and proposal collection (Lee et al., 2010). Moreover, an architectural evaluation of a clinical guideline system application used three evaluation techniques: technical, functional and clinical evaluation, which used the user feedback of domain experts to evaluate the system's performance (Shalom et al., 2016). The naturalistic evaluation brings the values of usability, utility, and functionality of the system in a more human-centred perspective. Additionally, the criteria used to evaluate the human-centred perspective varied based on the focus of the study. The expert-based systems used expert knowledge to validate the accuracy of the system, and the improvements in the behaviour of users were focused on evaluating some of the CDSS, which used randomised control trial methods (Montgomery et al., 2000) etc. However, the user validation of the system has to be designed to fulfil the desired purposes of the study. A human-centred, more naturalistic approach was conducted to complete the implementation evaluation phase of the study with a more natural approach. "Human-computer interaction is a discipline concerned with the design, evaluation and implementation of interactive computing systems for human use and with the study of major phenomena surrounding them" (Hewett et al., 1992, p. 15). The

resulting artefact of the study was a usable software which interacted with humans to perform its functionalities. The technical evaluation of the artefact considered the standard criteria for evaluating the artefact. In contrast, the naturalistic evaluation phase focused on obtaining the user's perspective of the system design and modelling. The realistic evaluation phase was focused on the usability of the system. "Usability is considered a major success factor for current and future decision support systems" (Mucha et al., 2023, p. 92). Moreover, accuracy and usefulness are critical in evaluating computer-based diagnosis systems (Berner et al., 1994). The researchers alternatively use the term usability, utility, or usefulness to define the system's performance. The usability of health information technology can be performed in various ways, such as by conducting surveys, focus groups, interviews, thinking aloud, observations, heuristic evaluation, cognitive walkthroughs, review methods, etc. (Fitzpatrick, 1998; Maramba et al., 2019; Yen & Bakken, 2012). The human-centred evaluation of this study was conducted through the user feedback evaluation. Techniques such as thinking aloud, surveys, questionnaires, and interviews are more time- and resource-consuming. Moreover, these techniques are commonly used in analysing the behaviour of users. The purpose of the naturalistic evaluation of this study was to get the human perspective on the design and modelling of the system rather than the improvement or analysis of user behaviours due to the system. Therefore, the artefact of the study was evaluated by analysing the users' feedback. El-Halees (2014) uses opinion mining to test the software's usability. Additionally, the usability of the virtual reality applications was tested using a questionnaire technique with the participants to collect their responses on user experiences (Kamińska et al., 2022). The implementation evaluation phase of this study was completed with a technical and user involvement evaluation which confirmed the quality of the artefact through standard quality measurement and a human-computer interaction perspective, respectively.

In conclusion, this chapter justify the process of designing the proposed solution. Design science research methods and the principles of model creation in information systems are combined to develop the research solution. The research artefacts have been recognised, extracted, and designed by adopting design science methods. The data analysis, feature extraction, algorithm designing, and model creation are carefully selected based on their appropriateness of them in information systems. The DSRM introduced by Wieringa (2014)

was adopted to design the overall solution while selecting the survival analysis techniques as the appropriate method for analysing a dataset to predict the complications of diabetes. Moreover, a thorough justification of the used process of evaluation, is provided with the design, algorithm, and the final artefact evaluation.

Chapter 4 Situational Awareness

4.1 Introduction

The overarching research question of this study is “How can a prediction model be created to forecast the onset of complications of diabetes mellitus (CoDM) in a cohort of diabetes patients?” The study has deduced a solution for this overarching question: to build a CDSS capable of predicting the risk of CoDM. Although the literature review provides a sound knowledge of the existing research gaps in the field, conducting a thorough background exploration of the intended area is necessary. Moreover, the requirements, issues, challenges, and expectations vary from context to context. Therefore, explicating contextual information is vital in designing a solution for a real-world problem. This chapter describes the contextual details of the project while explaining their association with the designed artefact of the study. First, the current contextual issues and requirements are described here to explain the designed solution's contextual background and fundamentals. An extracted problem framework and a summary section at the end of the chapter provide a graphical overview of the current issue.

4.2 Identifying the Contextual Requirements

The existing issues and challenges of building a CDSS in general is discussed in Section 3.4. The gathered knowledge of a CDSS is utilised here to explore and understand the details of the field more pragmatically. The checklist recognised by Wieringa (2014) is adopted to guide the phase of exploring situational awareness of the study. Although the followed guideline is well established in the field of DSRM, a customised approach for solving a real-world issue in the healthcare sector is lack in the literature. Due to this reason, the situational awareness phase of this study has been thoroughly articulated to extract the most accurate client requirements. A series of well-planned client meetings were held to gather a sound understanding of the situation. The stakeholders, the challenges they faced, and their expectations were carefully discussed during these meetings to extract the contextual requirements. The social context of the real-world issue, the potential stakeholders of the system, and the feature set for creating the data models of prediction were mainly focused

in the situational awareness phase of this study. The brainstorming sessions and client meetings were scheduled for gathering the contextual requirements of the client. The system requirements were gathered through structured client meetings with a focus of requirement specifications. A requirement specification document has prepared and distributed among the participants to collect their requirements. The created requirement specification document is attached in the appendix B. Further, the results of EDA of the initial dataset have used for effectively communicating with the clients. The identified missing characteristics, the distributions of socio-demographic details, the categorisation of socio-demographic details are communicated with the client meetings to get their real-time feedback on the system requirements. Moreover, it has been identified that requesting the feature sets for modelling the data is vital aspect of the client meetings. Since the features for predicting the complications of diabetes are not well-established and its keep updating with new features, a systematic review is essential to get the most suitable feature set. The resulted feature set has been used to communicate with the client to request the dataset. The social context, potential stakeholders, their goals and expectations, and the possible set of features for building the models have been recognised as the output of this section which defines the foundation of the designed solution.

4.2.1 The Social Context of the Research Study

Data accessibility is one of the major obstacles identified in building a CDSS. It is crucial to have a quality dataset to construct a reliable CDSS with good performance. A quality dataset possesses the characteristics such as reliability, accuracy, completeness, relevance, and timeliness (Chen et al., 2014). Finding an excellent potential resource provider for collecting a quality dataset was challenging. However, Te Whatu Ora was selected as the client due to their availability of medical data repositories and the convenience of its reachability. Moreover, Te Whatu Ora is an institute of health administration which focuses on managing all health services of the people in the Waikato region of New Zealand. It serves 425,000 people while covering more than 21,000 km² (WDHB, 2022). Te Whatu Ora is a well-reputed organisation which maintains the medical records and details of the patients of Waikato in a confidential and formative manner. They have a separate team who work in the field of diabetes mellitus (DM). To understand the stakeholder requirements and the state of the art

of the area, client meetings were conducted with the data provider. The team members of the diabetes services in Te Whatu Ora explained the current state of diabetes in Aotearoa and the real-world issues regarding diabetes patients. The inconsistency of data structures, difficulties in data transformation, unavailability of medical records of diabetes patients who seek medical advice from general practitioners, and the problem of accessing the entire medical history of individuals were raised as common issues from the client's meetings. These asymmetric information issues were categorised into four pillars, which could be used to enhance the identified knowledge extension. The requirements, their feasibility, and the categorisation of asymmetric information issues into pillars were thoroughly discussed at the brain storming sessions with the client. These sessions were highly useful to recognise the true requirements of the client while understanding the context. The categorised pillars of issues in datasets of Te Whatu Ora are:

- The lack of accessibility of information
- The scarcity of consistent medical records of patients
- The lack of integration among medical data resources
- The lack of detailed longitudinal data records.

To understand the nature of the existing dataset, the data provider issued an exemplary dataset. After exploring the issued dataset, the state-of-the-art datasets in Te Whatu Ora were realised. An initial Exploratory data analysis phase has been conducted to explore the received dataset. The results of this data analysis phase revealed the missing characteristics of the dataset, the distributions of socio-demographic details, the features that are lack in the dataset, inconsistencies and irrelevant categorisations of the demographic details etc. We consistently communicated with the data provider to make modifications and improve the dataset. The ransomware attack at Te Whatu Ora in mid-May 2021 had slowed the data-receiving process. However, an adequate dataset for designing a model was collected that could satisfy the client's requirements. The social context—potential stakeholders, their requirements, available datasets, and the resulting answers of the outcome of this project—were determined by analysing the received dataset. The recognised stakeholders and their goals are described as follows to explain the status of the situation.

4.2.2 Identified Stakeholders and their Goals

Te Whatu Ora is the selected data provider for this study which makes it the primary stakeholder. Thorough discussions at the client's meetings revealed the potential stakeholders and their expectations from the study. The team members of the diabetes services in Te Whatu Ora, general practitioners, diabetes patients, and the policymakers at Te Whatu Ora are the identified potential stakeholders of the project. Additionally, the public can use this solution to get some insight into the onset of CoDM. Therefore, Te Whatu Ora, general practitioners, nurses, policy makers, resource allocators, and the patients of diabetes are selected as the main stakeholders of this system.

The sole purpose of the project's primary stakeholder is to have a mechanism to predict the onset of CoDM by utilising the existing data. However, the discussions lead to a broad spectrum of goals and expectations. Although a range of requirements and expectations were discussed, the feasibility and validity of them were checked carefully to extract the most relevant and realistic expectations. The following list indicates the selected main expectations of the stakeholders:

1. A clear understanding of the state-of-the-art of diabetes cohort in the Waikato district.
2. Identify the variations in demographic factors among diabetes cohorts in the Waikato district.
3. A concise image of dispersion, sociodemographic variations, and surviving of different complications among the cohort.
4. Accurately predicting the survival of CoDM in the cohort of diabetes.
5. Predicting the CoDM of individuals by feeding their details into the system.

4.2.3 Features Selection through a Systematic Review

The client meetings, and brain storming sessions reveal the requirement of a solid feature set to generate the prediction models. Since the used feature set in designing the system, plays a vital role of the system's accuracy, a thorough feature selection phase was conducted. A well-structured systematic review study was conducted to extract the most frequently used feature sets for predicting the above-mentioned highly prevalent four significant

complications. A detailed description of the planned methods for feature selection is included in Section 3.5.4. However, the real-world scenario led to a thorough systematic review of selecting the features for the proposed model.

Systematic reviews are commonly used for summarising research findings in health care (Gopalakrishnan & Ganeshkumar, 2013). Stakeholders in healthcare settings utilise periodic review articles to fulfil various purposes, such as academics seeking comprehensive literature overviews, clinicians relying on evidence-based guidelines, and doctors staying informed about the latest research developments. Further, this can be considered a justification for building clinical practice guidelines. The justification of the requirement of further research in the related topic may be useful for granting agencies (Moher et al., 2009). Primary care/family physicians use systematic reviews as decision-making tools. The identified gaps in existing prediction models recognised beneficial or harmful medical interventions, summarised through systematic reviews, are valuable for clinicians, researchers, the public, and policymakers (Gopalakrishnan & Ganeshkumar, 2013).

A well-known feature set is a contemporary requirement in the field of CoDM. Inaccurate and ineffective risk-scoring models result due to the lack of a standard risk factors in predicting the CoDM. (DCCT/EDIC & Braffett, 2016; Hippisley-Cox & Coupland, 2017). Therefore, extracting a standard set of features that can be used in multi-ethnic models is crucial. Moreover, electronic health records (EHRs) has been trending as an eminent data resource, which effectively use in health informatics (Häyrinen et al., 2008). Although EHRs are widely used for fulfilling various purposes in health informatics, due to the divergence of national health information stored in different countries, considerable variations can be seen in EHRs. United States Core Data for Interoperability (USCDI) makes a consensus to achieve interoperability. The data types used in USCDI are demographics, diagnoses, problem lists, family history, allergies, immunisation, medications, procedures, lab orders/values, vital signs, reports, and utilisation. There are a further nine data types under emerging data types: bio-sample data, genetic information, social data, patient-generated, community, geo-spatial, surveys, free text, and other data types (Gliklich et al., 2019). Due to the popularity and effectiveness of EHRs in health informatics, it is vital to focus on the most commonly available

features of EHRs when extracting features for creating prognosis or diagnosis models. This section of the study adopts the categories of data types in EHRs by considering the USCDI standardised method to categorise the features extracted from the literature. This portion of the research study focuses on extracting a frequently used feature set for predicting the selected set of complications of diabetes. A systematic review has been conducted to fulfil this aim with the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) standards.

4.2.3.1 Methodology of the systematic review

This systematic review has been performed according to PRISMA guidelines (Moher et al., 2009), a widely accepted scientific framework for reporting systematic reviews and meta-analyses. The four main PRISMA steps are identification, screening, eligibility, and inclusion. In the identification phase, researchers report the number of records identified through database searching and other sources. Duplicate records are removed in the screening phase while reporting the number of records that need to be screened and the number of records excluded. Studies should be filtered by excluding irrelevant articles through proper eligibility criteria. In the eligibility phase, the number of full-text articles assessed for eligibility and those excluded should be reported with reasons. Finally, the number of studies included in the systematic review should be specified at the inclusion (Moher et al., 2009). The repositories used for article selection, the methods used to extract the most relevant articles, excluding criteria, and reporting the results are milestones of conducting a more sound and scientific systematic review. The current systematic review study adopted the guidelines of the PRISMA method to conduct and report the research results. The research article identification has been done through a well-known article repository, "Web of Science"(Web Of Science, 2021) , which provides a consistent search interface to multiple databases of academic journals, conference proceedings, letters, and other related publications in various disciplines. The research articles were identified using a scientifically- structured search query:

Query (("risk*" or "risk model*" or "risk assess*" or "risk equation*" or "risk predict*") AND ("Diabetes*" or "Complication of diabetes" or "Complications of diabetes" or "comorbidities

of diabetes*" or "comorbidity of diabetes*" or "diabetic*") AND ("Statistical model*" or "Regression model" or "Cox*" or "Artificial Intelligence*" or "*model*" or "Time series analysis*" or "Machine Learning" or "Time series Forecasting"))

The search was filtered under publication year, document type, accessibility, and publication journal. Research papers published within the last eight years (2015–2023) were considered in this study. The latest electronic search was performed on 1st July 2023. The reviews, proceeding papers, meeting abstracts, editorial materials, book chapters, letters, and news items were filtered out from the search results due to the possible uncertainties and inconsistencies that they may have brought. Furthermore, only open-access articles were selected for this study. Fifteen journals with the highest impact factor were chosen as top-ranked journals among the resulting journals from the search query. The resulting articles were manually selected as eligible for the study by considering their relevance to the research topic. The relevancy criteria focused on the aim of the research study, the presence of considered risk factors, and the type of features they considered. Finally, articles were divided into four peers according to the complication type: DR, DNeu, DNep, and CVD. The flow diagram of the article selection of the study is illustrated in Figure 4.1.

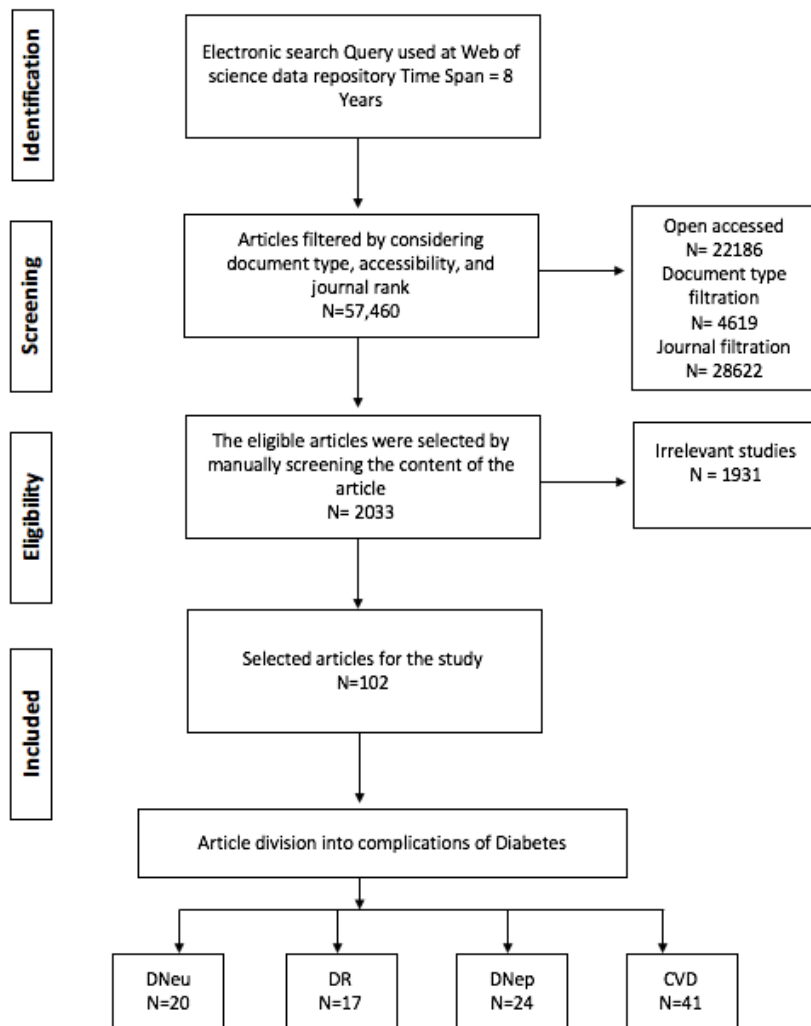


Figure 4-1 : Flow chart of article selection of the systematic review.

The feature list for each complication was extracted by reviewing the features used in selected studies. The selection criterion for the frequent features was their presence in at least 20% of selected articles. For example, to be selected as a frequent feature for DNeu, the feature should have been present in at least 4 of the 20 selected articles (20%) in that category. The frequency of each complication has been set as the threshold value to determine the features. The extracted frequent feature set has been categorised using the USCDI standards, resulting in nine categories to align with the standards of USCDI. Although the lifestyle features are unavailable in EHRs, they were included in a separate category as lifestyle features due to their importance highlighted in considerable number of research papers. The highly frequent feature list was created by selecting the common features for all

complications, with 20% or more in each complication. The resulted feature set is categorised into the nine categorisations of USCDI and presented with their frequencies in each complication below. Further, the frequencies of the features were calculated as percentages to visualise their utilisation in each complication. Moreover, the conducted systematic review has been published in *Current Diabetes Review*, with the title “Comprehensive factors for predicting the complications of diabetes mellitus: A systematic review” (Madurapperumage et al., 2021). A thorough description of the performed systematic review and its results with used article references can be seen in this published research paper.

4.2.3.2 Risk Factors for Predicting Complications of Diabetes Mellitus

The selected feature set for each complication are tabulated with their frequencies in Table 4.1 and their visual presentation can be seen in Figure 4.2.

Attribute Category	Attributes	Neuropathy (n=20)	Retinopathy (n=17)	Nephropathy (n=25)	CVD (n=42)
Demographic	Age	95.0	70.6	100.0	90.5
	Gender	95.0	70.6	80.0	73.8
	Ethnicity	30.0	5.9	40.0	33.3
Vital Signs	BMI	55.0	47.1	80.0	76.2
	SBP	60.0	52.9	80.0	83.3
	DBP	50.0	41.2	68.0	54.8
	Heart Rate	15.0	5.9	4.0	26.2
	Pulse Pressure	0.0	5.9	8.0	26.2
	Weight	30.0	5.9	4.0	21.4
	Height	30.0	0.0	8.0	4.8
	Waist Circumference	25.0	5.9	16.0	11.9
	Retinal arterial calibre (CRAE)	0.0	11.8	0.0	0.0

Table (Continue)					
Attribute Category	Attributes	Neuropathy (n=20)	Retinopathy (n=17)	Nephropathy (n=25)	CVD (n=42)
Lab Orders/ Values	HBA1c	70.0	47.1	76.0	71.4
	HDL	45.0	35.3	56.0	59.5
	Total Cholesterol	45.0	35.3	44.0	54.8
	LDL	40.0	23.5	36.0	54.8
	Triglycerides	30.0	23.5	44.0	52.4
	EGFR	55.0	23.5	84.0	31.0
Lab Orders/Values	Fasting Blood glucose	0.0	17.6	20.0	11.9
	Serum Creatinine	5.0	5.9	28.0	11.9
	Urine albumin-to-creatinine ratio	10.0	11.8	36.0	9.5
	Serum uric acid	0.0	5.9	20.0	0.0
	Haemoglobin	0.0	0.0	16.0	0.0
	Serum Albumin	10.0	0.0	16.0	7.1
Diagnoses	Duration of diabetes	50.0	47.1	64.0	52.4
	Age at diagnoses	0.0	11.8	16.0	0.0
	Prior Amputations	30.0	5.9	0.0	2.4
Medication	Statin	20.0	23.5	16.0	31.0
	Metformin	20.0	5.9	8.0	4.8
	Blood pressure lowering treatment	10.0	11.8	28.0	38.1
	Insulin	45.0	23.5	8.0	26.2
	ACE	5.0	17.6	36.0	26.2
	Beta-blocker	5.0	17.6	8.0	19.0
	Aspirin	5.0	0.0	8.0	14.3
	Calcium channel blocker	5.0	17.6	0.0	16.7
	Lipid Lowering drugs	5.0	5.9	16.0	19.0
	Oral medication for diabetes	5.0	0.0	0.0	16.7
	ARB	5.0	11.8	4.0	14.3

Table (Continue)					
Attribute Category	Attributes	Neuropathy (n=20)	Retinopathy (n=17)	Nephropathy (n=25)	CVD (n=42)
Medication	Diuretic	0.0	5.9	8.0	11.9
	Antidiabetic medication	0.0	11.8	0.0	14.3
	Fibrates	0.0	11.8	0.0	0.0
	Sulfonylureas	15.0	0.0	4.0	4.8
Problem List	Hypertension	30.0	17.6	28.0	31.0
	Myocardial Infarction	0.0	0.0	0.0	28.6
	CKD	15.0	5.9	8.0	19.0
	No-albuminuria/ Microalbuminuria /Microalbuminuria /Proteinuria	15.0	5.9	4.0	11.9
	Atrial fibrillation	5.0	0.0	0.0	19.0
	CVD	0.0	29.4	28.0	9.5
	Diabetic Neuropathy	15.0	23.5	8.0	2.4
	Diabetes	45.0	17.6	12.0	9.5
	Renal disease requiring dialysis	15.0	0.0	0.0	0.0
	Retinopathy	15.0	11.8	20.0	4.8
	Family History	Family history of CVD	60.0	29.4	32.0
Bio-sample Data	TNFR-1	5.0	0.0	16.0	0.0
Lifestyle Features	Smoking History	80.0	52.9	60.0	81.0
	Alcohol usage	15.0	23.5	12.0	28.6
	Exercise	10.0	5.9	4.0	14.3
	Income	0.0	11.8	4.0	4.8
	Education level	10.0	11.8	4.0	9.5

Table 4-1: Frequencies of the selected features (Madurapperumage et al., 2021).

Among the chosen attributes, age, gender, ethnicity, weight, height, BMI, smoking history, HbA1c, SBP, eGFR, DBP, HDL, LDL, total cholesterol, triglyceride, use of insulin, duration of

diabetes, and family history of cardiovascular (CVD), and diabetes was recognised as the feature subset which could be used in risk prediction of all four complications. According to the selected feature list, predicting CVD got the most significant number of features (n=29), while DR had the most minor features (n=19). DNeu and DNeph result with 24 features. Age and gender remained the two most frequent features for DNeu and DR, while gender was the third-highest priority for CVD and DNeph. Although the absence of a percentage of a factor in one complication represented its infrequency, it does not mean the invalidity of it in predicting that complication. For example, the percentage of the feature of "urine albumin to creatinine ratio" only showed in DNeph, but it had been used for all other complications less frequently. Some features were extracted due to their frequency in one complication, which was entirely unrelated to other complications. "Retinal arterial calibre" was a feature that had been selected due to its frequency in DR, which cannot be used anywhere else. "Renal disease requiring dialysis", "Myocardial infraction" and "Fibrates" were a few other features that were selected based on one complication. Moreover, the terminology used in different articles varied hugely. Due to the requirement to extract the terms in the papers as they were, some feature terms may overlap. For example, "Metformin" and "Oral medication for diabetes" are two feature values. Further, the term "Antidiabetic medication" was also included under the medications of diabetes to maintain the authenticity of the words.

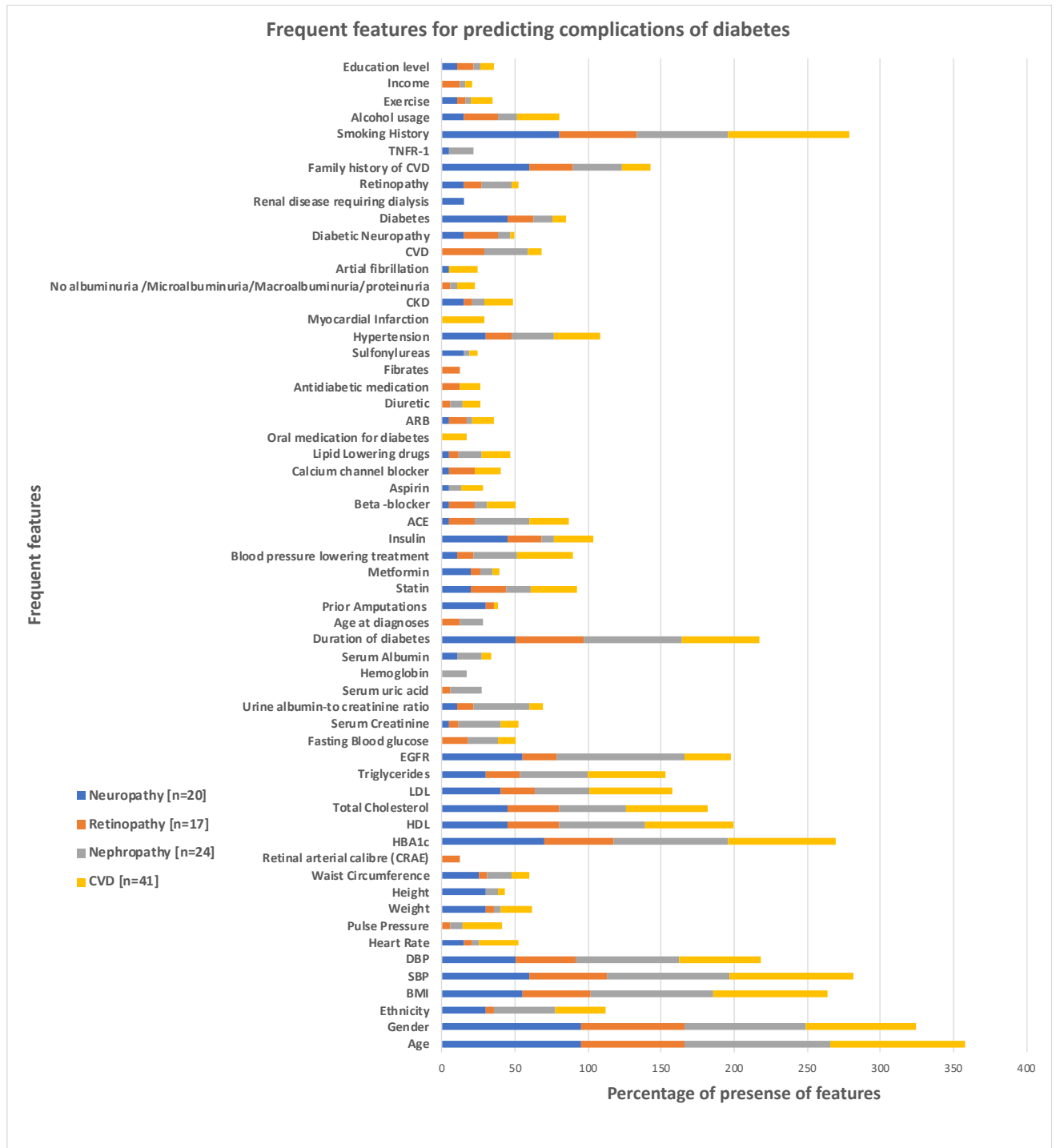


Figure 4-2 : Percentages of selected frequently used features for predicting the complications of diabetes (adapted from Madurapperumage et al., 2021).

As a result of this systematic review, 59 features were identified as the most common features. Nineteen features were recognised as common in all four considered complications, which are age, gender, ethnicity, weight, height, BMI, smoking history, HbA1c, SBP, eGFR, DBP, HDL, LDL, total cholesterol, triglyceride, use of insulin, duration of diabetes, family history of CVD, and diabetes. Although the frequent features used in literature were considered, some features may have had significant importance in predicting CoDM but were not used as frequently as the other traditional features used in the literature. The genetic risk factors were an excellent example of features not commonly used to design risk prediction models but possessed significant prediction power. Recently, there has been a trend in using genetic factors (Bebu et al., 2021) and biomarkers (Chan et al., 2021) for predicting CoDM. Single nucleotide polymorphisms (SNP) are used vastly as a genetic factor for predicting CVD (Bebu et al., 2021). Further, the biomarkers TNFR-1, TNFR-2, and KIM1 have been used in predicting diabetes nephropathy. Moreover, diagnosis and prognosis of diabetes retinopathy have been carried out with image processing techniques. Therefore, the features used in this complication are varied from the rest. Retinal arterial calibre (CRAE) (Coca et al., 2017; Deal et al., 2018), arteriolar tortuosity, and fractal dimensions are commonly used in predicting DR with image processing techniques (Sandoval-Garcia et al., 2021). Further, the effect of retinal venular tortuosity and fractal dimensions in predicting incident retinopathy was explored to understand their prediction capability (Forster et al., 2021). The feature categories used terminology and how features are used vary among scholars. The usage of features in prediction models is hard to generalise since data availability, the nature of data sources, the focus of the risk model, and selected machine learning (ML) techniques are hugely different.

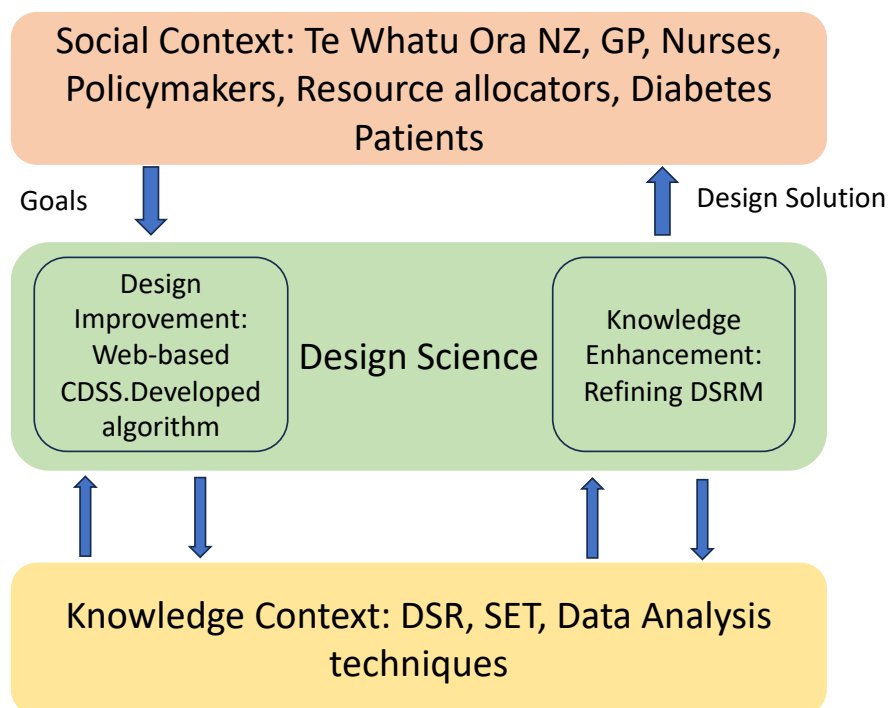
Although the extracted, detailed feature set of the systematic review had been requested from Te Whatu Ora, a set of obstacles limited the availability of the dataset. The limitations of existing datasets of Te Whatu Ora, the inadequate interconnections of data repositories in Te Whatu Ora, the limitations of gathered features of the patients, shortages of staff during Covid-19, the accessibility, and the changes of policies that occurred after the ransomware attack of Te Whatu Ora all impacted on collecting the desired dataset with the requested feature set. However, the received dataset from Te Whatu Ora consisted of only 10 features,

some demographic details, and lab values. The requested feature set consisted of 59 attributes categorised into nine groups. The significant limitation of the availability of the desired features may have had an adverse effect on the model. Therefore, the model accuracy may have been compromised with the limited number of available features. The solution for the practical issue is design by considering the collected requirements and available features. The systematic approach for the solution design was made by adopting the framework mentioned above by Wieringa (2014) and explained in the following section.

4.3 Conceptual Problem Framework

The explored contextual details are embedded with the research question to designing the ideal solution. A conceptual problem framework is a structured approach to identify, analyse, solve, and systematically visualise complex problems. It assists in visualising the issue at hand in a well-structured and more decomposed manner. The adopted research framework of the study is illustrated in Figure 4.3. It consisted of three tiers: the social context, the knowledge base, and the component of design science. The social context of the study provided the potential stakeholders and existing societal issues which needed to be addressed. The knowledge context comprised established theories, formalities, methodologies and frameworks regarding design and data science. The design theory in design science research was adopted and the theory of survival analysis to design the solution to address the issues that arose from the social context. The developed solution consisted of an improved design of a CDSS and an enhancement of existing knowledge. The implemented CDSS can predict a selected range of CoDM by utilising statistical methods on a longitudinal dataset. Additionally, the implemented solution and the knowledge extracted from the process of implementing the CDSS expanded the existing knowledge of designing a CDSS. The output of the research answers the issues of social context through these improved designs and knowledge expansion.

Research framework



SET – Software Engineering Techniques

Figure 4-3 : Conceptual research framework adopted from the Wieringa (2014).

4.4 Summary

Situational awareness is essential for understanding the state of the art of the real-world issue. A series of client meetings and brain storming sessions were conducted to grasp the feasible content of the project. Additionally, the extracted knowledge from the literature review was utilised in situational awareness. The existing pitfalls of the domain, the necessary improvements of the context, the issues with the techniques used to create the models, and the differences of methods used in model creation and their pros and cons were explored using the literature review. The real-world problem investigation was done with our client, Te Whatu Ora. As stated above, the requirement gathering meetings were conducted to understand the user requirements and the nature of existing data sources. The revealed stakeholders from the meetings are Te Whatu Ora, general practitioners, nurses, policymakers, resource allocators, and patients. The client meetings held with a team at Te

Whatu Ora, who are specifically work in the division of diabetes mellitus. The team consist with a general practitioner, data base management specialist, divisional lead, and a system analyst. The stakeholders had several goals to achieve through a CDSS. Te Whatu Ora was eager to reduce the expenses of diabetes and its complications while understanding the nature of diabetes and CoDM among the cohort in Aotearoa. They have been seeking a statistical analysis of the cohort and a system which could support the prognosis of CoDM using computerised systems to reduce the expenses of expert knowledge, expensive tests, and medications. Further, they were enthusiastic about increasing the health indices of the cohort by prognosing irreversible and chronic diseases. The advancement of health indices directly benefits enhancing the individual and societal health indices, which are the main intentions of Te Whatu Ora. Moreover, the CDSS may assist general practitioners with the prognosis of CoDM. The high prevalence rate and the irreversible adverts in individual health made prognosing CoDM by a computerised system more convenient. Additionally, patients with DM may benefit from a CDSS to be conscious of their health. The awareness of upcoming health issues may motivate the patients to lead healthy lifestyles while closely monitoring their well-being. The awareness of patients highly affects in increasing the health indices and protecting them from future health adverts.

The client meetings, literature review and the brainstorming sessions leads to the utmost artefact of this study which is a CDSS to predict the complications of diabetes mellitus. The artefact created based on three identified categories of issues: asymmetry of data, knowledge gap, and contextual scarcity. The above-mentioned issues in the datasets of Te Whatu Ora made the information asymmetry. The lack of accessibility of information, the lack of consistent medical records of patients, the lack of integration among medical data resources, and the lack of explicit longitudinal data records are the severe issue extracted from the client meetings and brainstorming sessions. The existing knowledge gap of the field including the pragmatic drawback of existing decision-making systems, the scarcity of a standard feature set for predicting CoDM, the lack of consensus of adopting DSRM to solve a healthcare issue, and the gap in using data analysis techniques for solving a healthcare issue were identified through the literature review and the conducted systematic review. Additionally, the results of exploratory data analysis phases reveal the nature of the datasets at Te Whatu Ora, while

providing a clear understanding of the context. The lack of risk prediction models created for differentiate the ethnic groups in New Zealand, and the scarcity of CDSS used by stakeholders of healthcare in New Zealand made the component of contextual scarcity. The suggested artefact of the study and the issues directed to build the artefact is visualised in Figure 4.4.

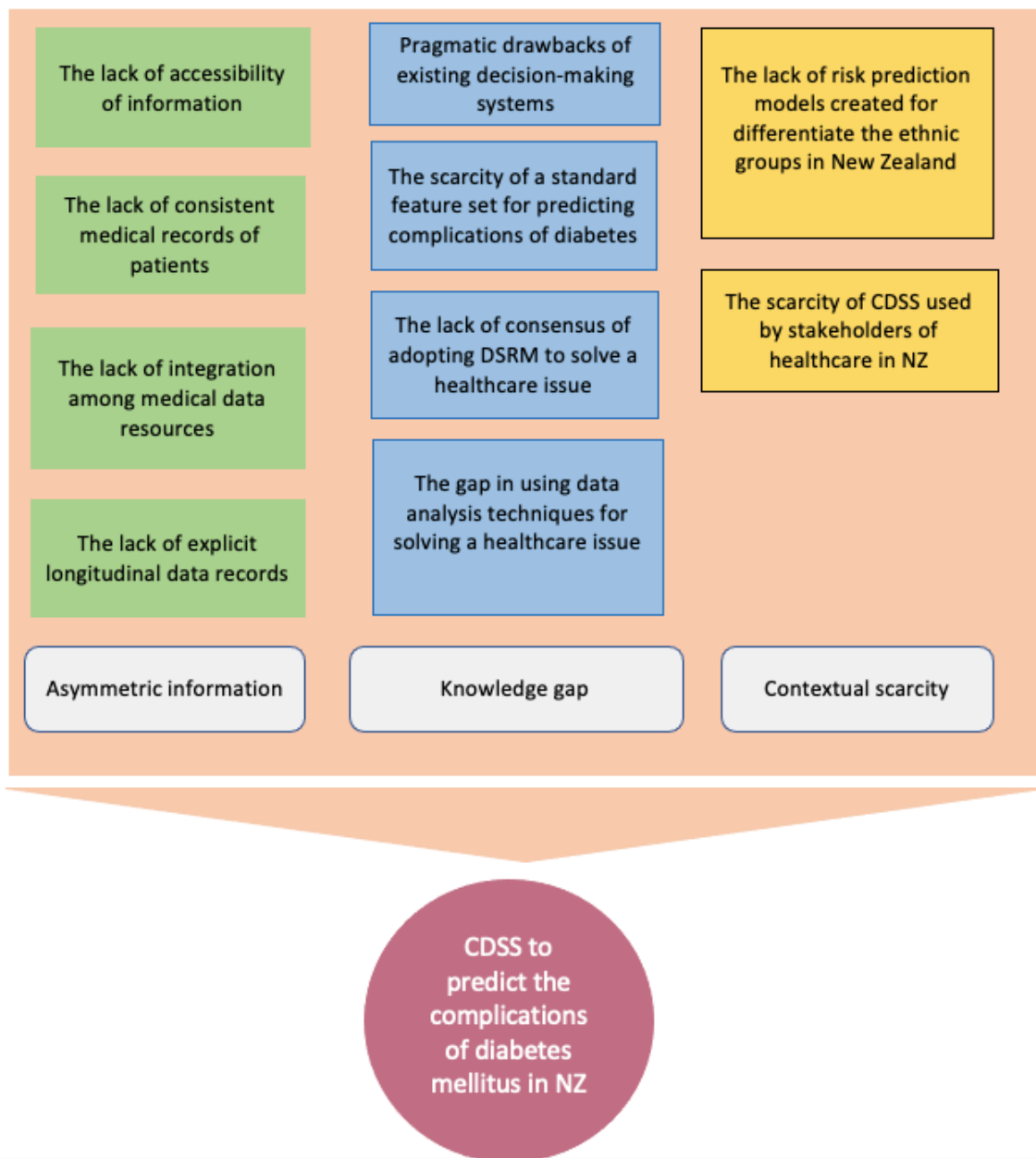


Figure 4-4 : Diagram to illustrate the ultimate research artefact of the study.

Chapter 5 The Design of the Artefact

This research study adopts the methodology introduced by Wieringa (2014) due to its pertinence to the nature of the study. The research framework and the design cycle explained in the previous chapter are directly adopted for systematically performing this research study. This section describes how the well-established concepts of design science research methodology (DSRM) were used for conducting this research. Moreover, this chapter presents a detailed description of the solution implementation process while introducing the resulting conceptual framework for solving similar approaches.

5.1 Proposed Artefact of the Study

This study focuses on solving an issue in the healthcare industry through a design science approach. Since diabetes mellitus (DM) is a global health issue with a high prevalence rate and vast expenditures, it became an emerging field of using computerised systems in designing models for analysing, exploring, and predicting diabetes. Further, the complications caused by diabetes have worsened the status of diabetes in societal and individual health indices. Predicting the risk of complications of diabetes becomes vital due to its benefits in the healthcare industry. The capability of issuing early warnings on complications enhances individuals' health and reduces expenditures. Additionally, the exploratory analysis of complications of diabetes in different cohorts extends the knowledge context, which benefits policymakers and healthcare resource allocators.

The overarching artefact of this study is a clinical decision support system (CDSS) which can predict the onset of a selected set of complications of DM using a diabetes cohort in New Zealand. This study generates the following outcomes during the journey of creating the CDSS.

- i. Design science approach for creating a CDSS
- ii. A systematic review of feature selection
- iii. Knowledge of the methodology of designing a CDSS

- iv. Exploratory analysis of socio-demographic factors among diabetes cohort in New Zealand
- v. Statistical regression models to predict the complications of diabetes
- vi. A web portal-based CDSS for predicting the complication of diabetes mellitus (CoDM).

The artefact of this study was designed, implemented, and evaluated by adopting the methodologies of design science due to the reasons mentioned earlier. The approach introduced by (Wieringa, 2014) is especially adopted here to conduct a well-structured process to innovate a widely accepted solution and a knowledge of the design process of the artefact. The outcomes of the study mentioned earlier serve to solve a real-world problem while extending the knowledge of design science.

5.2 Solution Design

The proposed artefact has been designed to provide the desired services to answer the identified real-world issues. The designed artefact should interact with the problem context in a way that is intended to treat the problem. Therefore, the "solution" is the interaction between the artefact and the problem context (Wieringa, 2014). The identified problem framework from the previous chapter is utilised here to explicate the fundamental requirements, their feasibility, applicability, usability, and relevance with the suggested artefact.

The solution designing phase dedicated for carefully examining the requirements, refining them by analysing their feasibility, and validity, while defining the scope of the study. The issues with designing perspective of the solution were gathered from the literature review, where the focus was to identify the existing pitfalls in the designing perspective and mitigating them to result in an innovative solution. The major design concern of the study is to confirm the adoptability of DSRM in healthcare settings. The selected DSRM was followed through the process of formulating the study starting from problem investigations, design artefacts to implementation evaluation of the project. The problem investigation phase is conducted in the literature review and the situational awareness chapter to define the correct real-world

issue. The methods used in this phase may guide the future researchers to adopt similar approach of extracting the real-world issue. The situational awareness phase is mainly composed with client meetings and brainstorming sessions, which provided a sound understanding of the context. Additionally, the conducted exploratory data analysis phase and the systematic review clarifies the state-of-art of Te Whatu Ora, and their requirements, while agreeing on the feasibility of the solution. This was a cyclic process for gathering and clarifying the contextual awareness. The processes used in the situational awareness phase of this study can be adopted in similar situations in healthcare sector. The phases of design of artefact and the evaluation describe the techniques and processes used in each step of the phases and provide a systematic process on adopting DSRM concepts throughout the project. The techniques used in the stages of design artefact: requirement gathering, and analysis, data collection, data pre-processing, exploratory data analysis, feature selection, model selection, implementation of the web-portal, and the phase of evaluation provide a solid knowledge of using DSRM in the healthcare setting while fulfilling the identified knowledge gap of the study.

The requirements of the solutions were mainly gathered through the conducted client's meetings and brain storming sessions. The diabetes team at Te Whatu Ora, explained their stakeholders' requirements and expectations through this project. The brainstorming sessions conducted with the stakeholders further clarified the requirements and the scope of the project. The functional and non-functional requirements are gathered and clarified in the brainstorming sessions with the client. The use case diagrams, conceptual frameworks, and prototypes were used to communicate with the client while defining a solid scope for the project. As mentioned in the section 4.2, a requirement specification document was finalised by combining the gathered requirements from the client meetings. (Appendix B). The following figure illustrate a part of requirement specification document which express the functional requirements.

3. Requirements Collection

3.1. Functional Requirements

- **What should the system do?** (Use MoSCoW framework to prioritize)
 - **Must Have:**
 - Predict the risk of the complications of diabetes using patient data (age, lab results, medical history, etc.).
 - Provide a user-friendly tool for the public.
 - **Should Have:**
 - Statistical overview of the cohort of diabetes.
 - Graphical visualization of survival rates of the cohort.
 - **Could Have:**
 - Comparison of individual health against the cohort.
 - **Won't Have** (for now):
 - Integration with mobile apps for patients (optional in the future).

Figure 5-1 : Snippet of the requirement specification document where the functional requirements are extracted.

Further, the use cases of this requirement specification clarify the clients' expectations and goals. The following figure illustrate the gathered use cases.

3.4. User Stories/Use Cases

- Wants to input a patient's data and receive a risk prediction for complications of diabetes so that I can plan the patient's treatment.
- Wants to get the patients risk for the complications of diabetes as a percentage over the time.
- Wants to visualize the distribution of socio-demographic details of the cohort.
- Wants to see the changes in survival rate of patients over the time.
- Wants to get the survival of patients for the most common complications of diabetes in the cohort.
- Wants to download the patients reports for each complication/all of the selected complications.

Figure 5-2 : Snippet of the use cases in the requirement specification document.

The gathered requirements in the situational awareness stage are analysed to design a solution for the identified research questions. Further, the requirements are more generalised to match the necessities of the general stakeholders of the healthcare sector through the series of client meetings. The necessities are clarified through the involvement

of potential stakeholders and validated against the existing literature. The contribution of the novice system to fulfilling the requirements of stakeholders can be listed as follows:

1. A clear understanding of the territory of diabetes patients
2. Variations in socio-demographic factors among diabetes cohorts in New Zealand
3. Understanding the distribution and evolution of CoDM
4. Predict the CoDM using a survival dataset
5. Developing a web portal to present the CDSS to predict the CoDM.

Additionally, the literature review of the study reveals another set of requirements for making a CDSS:

1. Abstract knowledge of the process of implementing and evaluation of a CDSS
2. A solid feature set for predicting the CoDM
3. A widely accepted CDSS to predict the risk of CoDM regardless of the considered cohort
4. Data analysis model for accurate prediction.

All the gathered requirements from the literature, client meetings, and brainstorming sessions were combined to finalise the project's scope. The limitations of data availability of cohorts of different countries made it impossible to achieve the goal of a widely accepted CDSS with a good training dataset. Although there is a significant research gap for a CDSS which can be sensitive to the ethnic diversities of the diabetes cohorts, the practical issue of obtaining health-related data for diverse population groups still made it an unachievable target. Moreover, the following assumptions were made about the social context of the artefact:

1. The gathered stakeholders' requirements represent the necessity of a CDSS in diabetes.
2. The collected data is a representative sample of the diabetes population of New Zealand.

The selected feasible requirements were converted to the functionalities of the CDSS. The system's selected functional and non-functional requirements are presented in Table 5.1. The requirement set of the CDSS is presented with their descriptions in this table, to provide a clear understanding of the functions of suggested system whereas the stakeholder interaction with the functions is visually represented in the figure 5.1.

No	Functional/ Non-Functional	Requirement Description	The functionality of the system
1	Functional	Understand the dispersion of diabetes patients based on demographic details.	Visualise the gender dispersion among the diabetes cohort.
2	Functional		Visualise the age distribution among the diabetes cohort.
3	Functional		Visualise the dispersion of patients by considering their Māori/non-Māori characteristics.
4	Functional		Visualise the ethnic diversity among diabetes patients.
5	Functional	Understand the differences in survival rates among diabetes patients based on demographic details.	Visualise the survival curves of diabetes based on gender dispersion.
6	Functional		Visualise the survival curves of diabetes based on the age groups.
7	Functional		Visualise the survival curves of diabetes Māori/non-Māori ethnic groups.
8	Functional		Illustrate the survival curves of diabetes based on the ethnic groups.
9	Functional	Understand the socio-demographic dispersion of patients with each complication.	Visualise the gender dispersion among the cohort of each complication.
10	Functional		Visualise the age distribution among the cohort of each complication.
11	Functional		Visualise the dispersion of patients by considering their Māori/non-Māori characteristics in each complication.
12	Functional		Visualise the ethnic diversity among the cohort of each complication.

Table (Continue)			
No	Functional/ Non-Functional	Requirement Description	The functionality of the system
13	Functional	Understand the differences in survival rates among diabetes patients for each complication based on demographic details.	Visualise the survival curves of each complication based on gender dispersion.
14	Functional		Visualise the survival curves of each complication based on the age groups.
15	Functional		Visualise the survival curves of each complication of Māori/non-Māori ethnic groups.
16	Functional		Illustrate the survival curves of each complication based on the ethnic groups.
17	Functional	Understand the survival rates of each complication	Visualise the survival curves for each complication.
18	Functional	Analyse the survival rates of the cohort for complications based on demographic details.	Visualise the survival curve of the cohort in each complication based on demographic details.
19	Functional	Analyse the survival rates of the cohort for complications based on demographic and laboratory details.	Visualise the survival curve of the cohort of each complication based on demographic and laboratory details.
20	Functional	Predict the survival of individuals in each complication.	Predict the survival of an individual for each complication based only on the demographic details.

Table (Continue)			
No	Functional/ Non-Functional	Requirement Description	The functionality of the system
21	Functional	Predict the survival of individuals in each complication.	Visualise the resulting demographic details-based survival curves of the individual with the cohort.
22	Functional	Predict the survival of individuals in each complication.	Predict the survival of an individual for each complication base on a combination of demographic and laboratory details.
23	Functional		Visualise the resulting all-details-based survival curves of the individual with the cohort.
24	Functional	Generate a report for individual.	Generate a report for each complication based on the demographic details.
25	Functional		Generate a report for each complication based on the demographic and laboratory details.
26	Functional		Generate a report for all complications based on the demographic details.
27	Functional		Generate a report for all complications based on the demographic and laboratory details.
28	Non-Functional	A user-friendly approach to predicting the CoDM.	Implement clear graphical user interfaces for the web portal.
29	Non-Functional	Accurate predictions results.	Use statistically significant methods for predicting the CoDM.
30	Non-Functional	Efficient system for predicting the CoDM.	Using a web portal to launch the CDSS makes the system more efficient than a stand-alone system.

Table 5-1 : Table of system requirements.

The selected requirement set was mapped into a use case diagram, as shown below. It represented the functional requirements that were chosen for building the CDSS. The tabulated requirement set in the table 5.1 was divided into four major use cases: exploratory data analysis, data visualisation, predicting the CoDM, and generating reports. The categorised four use cases were assigned with the related sub sets of use cases which covers the entire functional requirements in the Table 5.1. Exploratory data analysis was conducted to provide the state-of-the-art of the cohorts of diabetes. The socio-demographic data visualisation and survival curves of diabetes and its complications were visualised as graphical representations. Moreover, the predicted survival curves were visualised based on the demographic and all details. The survival curves of individuals were displayed with their cohort to compare the survival rates. The prediction of CoDM was made with two sets of features: demographic details and a combined set of demographic and laboratory values. The report generation was also implemented as two phases, where the report of one chosen complication or all complications could be taken as the outcome.

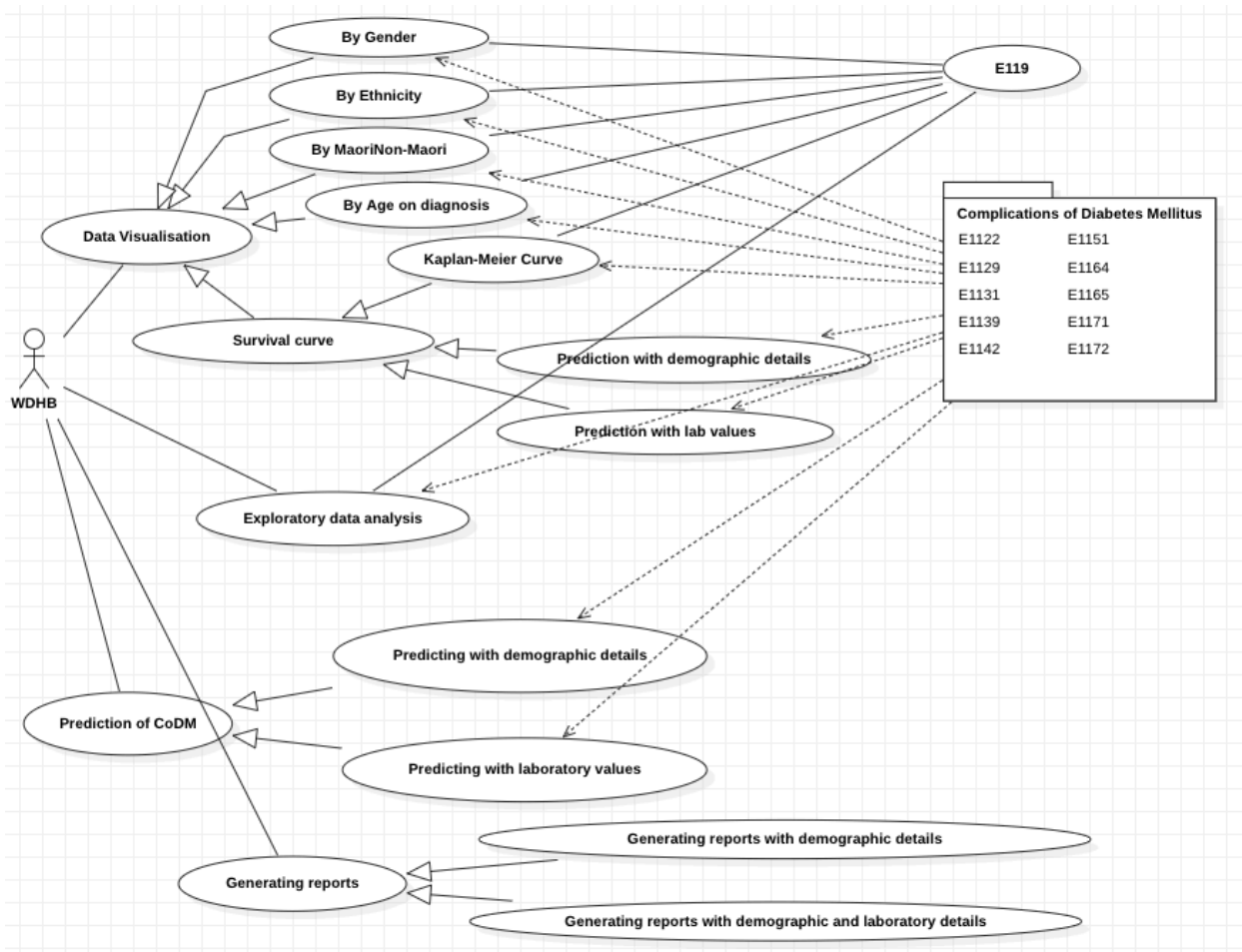


Figure 5-3 : Use case diagram for the functional requirements of the system.

The suggested solution with the functionalities mentioned above was implemented to achieve the goal of this project. The procedure of solution implementation has been described in the following section. The suggested solution has validated through a feasible study to confirm the validity of the design with the user requirements. The selected feasibility study is included in the section 6.1, since it belongs to the evaluation phase of the study.

5.3 Solution Implementation

The methodologies of the design science research approach have been used to identify, design, and evaluate the suggested CDSS. The implementation of the solution adopted the widely accepted empirical cycle due to its obvious appropriateness of system implementation. The empirical cycle is

a "rational way to answer scientific knowledge questions" (Wieringa, 2014, p. 109). Figure 5.2 illustrates the adapted empirical cycle for the study. The CDSS was implemented by following the steps in the empirical cycle to build a scientifically sound solution to the research question.

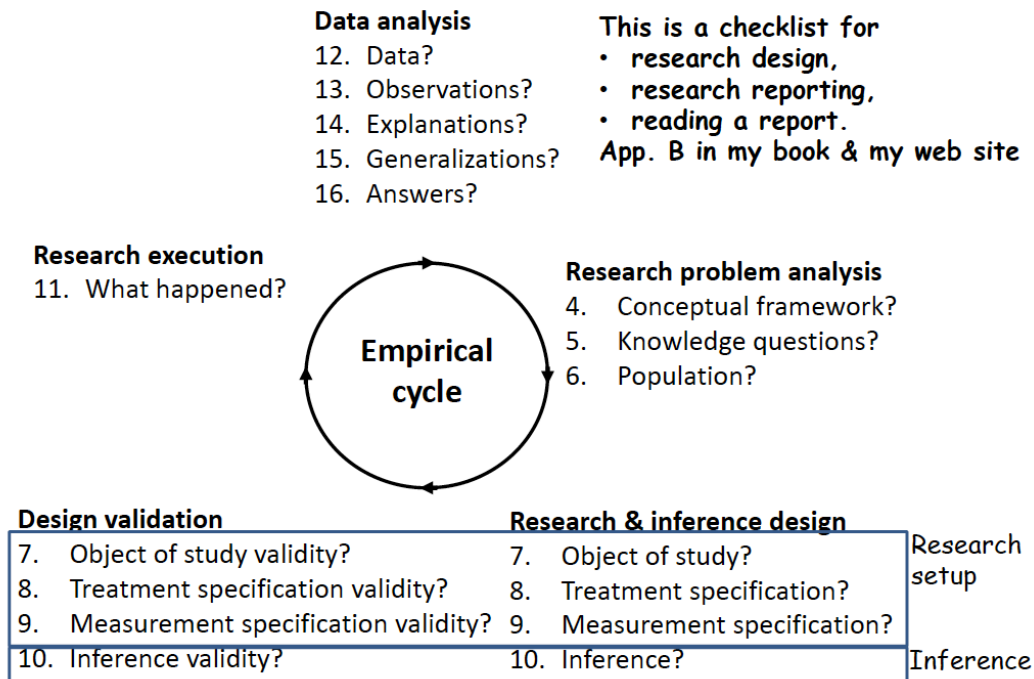


Figure 5-4 : Empirical cycle of DSRM (Wieringa, 2014).

A concise system summary is explained here to provide an overview of the system before going over the details. The principles of design science research were adopted to craft this study. The framework and design process introduced by (Wieringa, 2014) was used for systematically conducting the whole project. The entire system was implemented using Python programming language due to its suitability, as explained in chapter 3. Survival analysis techniques were selected as the data analysing method. Moreover, the web portal implementation was done with the aid of the Streamlit library of Python. The built solution was evaluated through a technical risk and efficacy strategy which consisted of summative and formative evaluation episodes. The steps of the solution implementation process are described as follows.

5.3.1 Data Collection

The implementation process of the designed artefact begins with collecting the required dataset. The dataset for this study was collected from Te Whatu Ora. The members of the diabetes team at the Te Whatu Ora actively contributed on the data collection process. The discussed requirements and the results of brain storming sessions are closely monitored during this phase. For example the requirement specification document (Appendix B) was thoroughly analysed at the data collection phase to match the clients' requirements at the dataset. Further, As mentioned in the solution design phase the extracted requirements from the client meetings and brain storming sessions are inclined with the functional and non-functional requirements of the solution. The document review method was used by the team members of the diabetes service group to collect the necessary dataset. The data bases of Te Whatu Ora were browsed and reviewed to collect the dataset. The data extraction and integration were done through a set of storages of data repositories to meet the requirements of the requested dataset. The data set was collected as a secondary dataset, via reviewing the governmental databases, which leads to use document review method of data collection being the most suitable fit. As explained earlier, the data-collection phase was conducted simultaneously with the requirement-gathering stage through a few client meetings. A data repository of diabetic patients was maintained by Te Whatu Ora, which contained data from various hospitals and practitioners in the Waikato region. Although the laboratory values and general demographic data were available in the existing repositories, most lifestyle and biomedical features were not included. The gathered data set from Te Whatu Ora is a rich dataset representing diabetes patients in Aotearoa. It was assumed that the compiled dataset could be used to generalise the solution. The patients' collected lab data was considered to be standardised on the same scale with similar potential errors. Most of the laboratories' equipment and used scales were identical throughout the Aotearoa. Therefore, the laboratory values collected from different sources were considered to have a common consensus. Further, due to the general legal and ethical issues arising with medical data collection, the identification details of patients were kept encoded. The ethical considerations were seriously considered and did not breach any violations of confidentiality or security of the data set. Since the dataset of this study was not a primary dataset, a few unavoidable effects

adversely impacted the research. The selection of the group of complications of diabetes and the features used to predict the survival rates of the complications had to be done according to the availability of the dataset. The received dataset from the WDHB consisted of two flat files—diagnosis and test results—which spanned a decade, from 2010 to 2020. The diagnosis table consisted of the details of the diagnosis history of the patients. The test result table comprised the laboratory test results of the patients. A concise description of the received dataset is included as follows.

Dataset name	Diagnosis	Test result
Description	The diagnosis history of the patients is included in the diagnosis table.	The test result table consists of the patient's laboratory test values.
No of records	273,126	1,048,575
No of columns	12	8
Column names and data types in the dataset	Diagnosis Patient Refno <int> Diagnosis ICD10 Code <String> Diagnosis Description <String> Diagnosis Create DateTime <Date> Diagnosis Seq <int> Patient Refno <int> Patient_DOB <Date> Patient_Age_OnDiagnosis <int> Patient_Ethnicity <String> Patient_Maori_NonMaori <String> Patient_Gender <Enum> Patient_DeceasedFlag <Enum>	Patient Refno <int> Collection_DateTime <Date> Result_DateTime <Date> Order_Number <int> Specimen_ID <int> Test_Description <String> Test_Id <int> ResultValue <String> ResultFlag <Enum> ResultUnit <String> Result_RefRange <String> Result_Extended <String> Result_Comments <String> Result_Status_Description <String>

Table 5-2 : Data dictionary of collected datasets.

The data collection was conducted via a series of discussions with the assigned team members. The amount of existing missing values of some desired features, human mistakes in data extraction, and clarifications of the test values were some of the major concerns of the data collection phase. Strong and collaborative discussions had to be conducted to grasp the desired dataset and its meanings to collect it with the intended qualities. Characteristics of the collected dataset are included here to get an overall understanding of the received dataset. The diagnosis table consists of 273126 records, including 142,067 males, 111,085 females and 22 unknown values. The dataset is comprised with 160193 non-Māori patients and 92981 Māori patients. Additionally, the dataset consists of 24 ethnicities, where the majority is NZ European/Pakeha with 120,630 patients and Māori as 90311. The dispersion of age in this dataset is up to 105 years with the mean age is 66.63. The exploratory data analysis section includes a detailed description of the collected dataset.

5.3.2 Data Pre-Processing

The dataset that was received from Te Whatu Ora was two flat files, which consisted of laboratory details of the patients and the diagnosis history of the patients. The received dataset was no different from a real-world dataset, which usually shows inaccurate data types, outliers, high dimensionality, missing values, unstructured data, unit values transformation, and other errors (Alasadi & Bhaya, 2017). The received row dataset was thoroughly processed to eliminate the common misconceptions and mistakes using the above-mentioned pre-processing techniques to obtain a quality data set. The data pre-processing phase started with a data cleaning process by filling the missing values. Contextual knowledge was used when deciding the methods of data filling. Each patient's missing laboratory data values were filled in by their nearest recorded laboratory value. Further, the missing values of a particular laboratory test were filled with the median value of the cohort. The median value was selected as a suitable statistical filling value here due to its ability to reduce the effect of the recognised outliers. The patients with no laboratory value for a specific test for the whole decade were removed from the dataset by removing the entire tuple of details of the patients. The missing dates were filled with recognised patterns in the dataset. The tuples with ambiguous/unmatched data/noisy values were identified

and removed from the dataset. The time periods were calculated as the number of dates. The existing inconsistencies of the dataset, such as different units in HbA1C, LDL, HDL, and eGFR, resolved by selecting only the tuples with the most frequently used units. The flat files received from Te Whatu Ora were integrated to create the datasheets that fulfilled the requirements of the data analysis phase of the study. The data integration was completed using the common unique feature (patient id) that both files used to identify the tuples. The dataset of this study was further transformed by normalising the values such as diagnosis period and age at diagnosis. Additionally, the dataset was discretised according to ethnicity, age group, gender, disease, and Māori/non-Māori. The received dataset had values of 24 ethnic groups, where some groups had fewer patients. Therefore, the ethnic groups of the study were recreated to match the existing popular ethnic group categorisation used in Stats NZ, which is New Zealand's official data agency (Stats New Zealand, 2023) . Table 5.4 presents the mapping of 24 ethnic groups in the received dataset into six groups, which rearrange the use of ethnicity feature for the rest of the study.

Existing ethnicities of the raw dataset	The new set of ethnic groups	
NZ European/Pākehā	European	
Other European		
European No Further Definition		
Māori	Māori	
Cook Island Māori		
Tongan	Pacific people	
Samoan		
Tokelauan		
Fijian		
Niuean		
Other Pacific Peoples		
Pac People No Further Definition		
Indian		Asian
Chinese		
Other Asian		
South-East Asian		
Asian No Further Definition		
African (Or Cultural Group of African Origin)	Middle Eastern/Latin American/ African	
Latin American/Hispanic		
Middle Eastern		
Not Stated	Other ethnicities	
Not Specified		
Response Unidentifiable		
Other Ethnicity		

Table 5-3 : Table of ethnicities in the dataset and mapped ethnicity types in the study.

Additionally, creating age groups is a requirement in comparing the survival curves. Therefore, the ages of the patients have been categorised into five groups <25, 25–49, 50–64, 65–74, and ≥75. The age categorisation was taken from the systematic review results, where the age categories were extracted from the literature. The pre-processed data set was critically analysed

using exploratory data analysis methods to understand the dispersion, existing patterns, outliers, etc.

5.3.3 Exploratory Data Analysis

Exploratory analysis is critical in understanding the dataset with an initial investigation. The process assisted in discovering the patterns, identifying anomalies, and providing summary statistics, which helped in understanding the overall characteristics of the dataset. The exploratory data analysis phase of this study was conducted in two phases, where the first phase aimed to analyse the prevalence of socio-demographic factors of a portion of the dataset, and the second phase was targeted for analysing the whole dataset with detailed analysis with fundamental four tiers: univariate graphical, univariate non-graphical, multivariate graphical, and multivariate non-graphical.

5.3.3.1 Exploring the Prevalence of Factors in a Cohort of Diabetes Mellitus

The socio-demographic factors were analysed through a short dataset collected from the initial client meetings. The considered data set of the study consisted of 2,656 patients from 2018 to 2020. Population details of the Waikato region were collected from the official website of Stats NZ (Stats New Zealand, 2023). The exploratory analysis of the dataset was conducted using the Python programming language. The sociodemographic factors used in this study were age, sex, and ethnicity. Due to the retrospective design study and the existing limitations of data issuing policies of Te Whatu Ora, the sociodemographic features had to be restricted. Data pre-processing techniques were used in data cleaning and data transformation. The patients with missing attributes, such as age and the noisy values, were removed from the dataset to maintain consistency. In contrast, values with a negligible amount of data were removed to avoid the unnecessary complexity of the results. The data frames were created to fulfil the requirement of exploratory analysis by extracting and aggregating the necessary attributes. Although the sociodemographic factors had to be limited due to the study's retrospective nature, the selected features were matched with existing scholarly works.

Diabetes complications were extracted from the dataset to understand the distribution of micro and macrovascular complications of diabetes among the cohort. Among the unique disease sets of the dataset, the diseases with the highest number of patients were chosen to extract the complications of diabetes among the cohort. The results of exploratory data analysis were presented using tables and a Sankey diagram, to visualise the results. The complications with more than 80 patients were selected to be visualised in the diagram to simplify the representation.

The analysed results of the socio-demographic details of the cohort are included in Section 7.2.3 – Exploring the prevalence of sociodemographic factors in a cohort of diabetes mellitus.

5.3.3.2 Exploring the Data Set with Empirical Standards

The empirical exploratory analysis begins with the univariate non-graphical methods to illustrate the sample distribution and "to make some tentative conclusions about what population distribution(s) is/are compatible with the sample distribution" (Seltman, 2018, p. 63) . The exploratory data analysis phase was divided into three components to clearly explain the procedure and results, where each component consisted of four general types of exploratory data analysis (EDA): univariate non-graphical, univariate graphical, multivariate non-graphical, and multivariate graphical. The three components mentioned above analysed the original diagnosis table, the original test-result table, and the main data frame created to fulfil the system's functionalities. The combined results of all three components of the EDA phase revealed the important statistical perspective while showing a clear understanding of the dataset. Furthermore, an exploration of socio-demographic factors in the diabetes cohort is included at the end of the section to demonstrate an overall view of the cohort.

The diagnosis table consisted of 273,126 rows and 13 columns and was analysed using the above-mentioned EDA methods. The categorical variables of the diagnosis table, such as gender,

ethnicity, Māori/non-Māori, and the patient's age on the diagnosis of E119, were first analysed through statistical measurement, followed by histograms and box plots to visualise their characteristics. The cross-tabulation was done with the attribute combinations such as gender vs Māori/non-Māori and gender vs ethnicity. The age on the diagnosis of E119 was analysed based on gender, ethnicity, and Māori/non-Māori to explain how the statistical measurements of age on the diagnosis of E119 change in the categories. Moreover, correlation was calculated to understand the relationships among selected variables. Correlation matrices were represented to deeply understand the relationships among variables such as the age of the patient vs age at diagnosis of diabetes, and age of the patient vs age at diagnosis of the selected complications of diabetes.

Additionally, side-by-side box plots and violin plots were used to visualise the results of the above analysis. The meaningful data visualisations of the results are included in the results section. The test result table consisted of the laboratory test values of the patients. The results id, their description, units of measurement, and the value are the considered columns here. The analysis was done here to answer questions such as 1) What are the most frequent test results in the dataset? 2) What units are used to measure the laboratory tests, and 3) What tests can be combined to create the same test values? The gathered test results were scattered into different categories, representing the same tests with the same measuring units recognised from the univariate analysis of the test result table. They were rearranged to resolve the issue. The final test result table included HbA1c, cholesterol, triglycerides, HDL, LDL, and eGFR attributes. The used measurement units were mmol/mol, mmol/l, mmol/l, mmol/l, mmol/l, and mL/min/1.73m², respectively. The number of test details for each year was cross-tabulated here to get information regarding the test result details available in each year. The information revealed from the analysis helped request data sets from the client meetings. Moreover, the number of tests conducted each year was visualised using a bar plot to illustrate the distribution of test results along with the considered period. The visualisation of the count of all test results each year produced a messy representation due to the vast number of existing test results. The issue was overcome by visualising the top test ids that had the greatest number of test results.

After analysing the two received data files, the most vital data frame of the study was created by filtering the diabetes patients diagnosed with E119, which is type 2 diabetes, without any complications. The created cohort was the foundational diabetes cohort for creating the prediction models. The constructed data frame was named "t2dm" and pickled for later use. The produced data frame consisted of new columns such as "Diagnosis ICD10 Code _Comp", "Diagnosis Descripton_comp", "days", "event", "Year_E119", and "Year_Comp". The diagnosed ICD codes of the complications of the filtered patients, the description of the selected complications, the number of days between the diagnosis of E119 and its complication, whether the patient was diagnosed with the complication or not, the year of diagnosis of E119, and the year of diagnosis of complications were the respective descriptions of the columns. The "t2dm" data frame was thoroughly analysed for more insights about the data set. The attributes of the created t2dm data frame were analysed separately in graphical and non-graphical methods to analyse the characteristics of the considered cohort. Additionally, "Patient_Gender" vs "Patient_Maori_NonMaori", "Patient_Ethnicity" vs "Patient_Gender", "Diagnosis ICD10 Code_comp" vs "Patient_Gender", "Diagnosis ICD10 Code_comp" vs "Patient_Maori_NonMaori", "Patient_Ethnicity" vs "Diagnosis ICD10 Code_comp" were cross-tabulated to analyse them against each other. Correlation and covariance matrices were used to get the relationship between "Patient_Age_OnDiagnosis" vs "days". The pair plots, horizontal and vertical bar charts, side-by-side box plots, and violin plots were used to visualise the multivariate EDA results.

The analysis of the cohorts of each complication were done as the last step of the EDA process. The pickled "t2dm" data frame was used to create a separate data frame for each complication. Although this was not necessary to fulfil this purpose, individual data frames were required in making the models. The characteristics of each cohort of complications were analysed through the created data frames. The gender, age on the diagnosis of E119, ethnicity, and Māori/non-Māori were analysed for each data frame. Moreover, the survival curves of complications were developed using the Kaplan-Meier function to present the survival rates of the complications.

The characteristics of each cohort of complications were also compared through developed survival curves using the Kaplan-Meier survival technique. The results revealed the differences among the ethnicities, gender, age groups, and Māori/non-Māori characteristics of cohorts.

5.3.4 Feature Selection

The features used in this study were extracted in a method initiated with a critical systematic review followed by a feasibility test. Although the results of the systematic review paper (Madurapperumage et al., 2021) extracted the most frequently used feature sets from the literature, there was a practical issue finding the dataset to match this requirement. Since the above-described systematic review had been conducted to extract the most frequently used features for predicting the most common complications of diabetes, the extracted features set was sent over to the client to check their availability in the existing data repositories of Te Whatu Ora. However, since a secondary dataset was collected in this study, some unavoidable limitations were embedded in the dataset. Therefore, the feature set for this research was selected, as mentioned earlier, by considering the existing features in the collected dataset. Although the extracted feature set from the literature review was composed of 59 elements divided into nine categories, the received dataset from Te Whatu Ora consisted of four demographic values: gender, age, ethnicity, and Māori/non-Māori, and six laboratory values: HbA1c, cholesterol, HDL, LDL, triglyceride, and eGFR. The laboratory values were in different units of measures. However, the values in the standard units of measures such as mmol/mol, mmol/l, mmol/l, mmol/l, ml/min/1.73 m² were considered, respectively.

5.3.5 Model Selection

The survival analysis techniques were used in this study to predict the survival of complications of individuals. The research methodology chapter explained the suitability of the data analysis and model selection methods. The Kaplan-Meier survival technique was used in this study to visualise the survival curves of each cohort without considering the covariates. Additionally, the

Cox model was selected to predict the survival rate of complications with the aid of demographic details and laboratory values. The created model is described in the following.

5.3.5.1 Non-Parametric Techniques in Survival Analysis

The Kaplan-Meier survival techniques were used in exploring the general survival rates of the cohorts without modelling them with covariates. The whole diabetes cohort and the cohorts of each complication were modelled using the Kaplan-Meier survival technique. Further, this non-parametric survival technique was used to differentiate the cohorts according to their demographic details. Gender, age, ethnicity, and Māori/non-Māori characteristics were used to explore the diversity of the cohorts. The resulting survival curves are included in the results and findings section. The resulting Kaplan-Meier curves were beneficial in identifying the existing differences in cohorts. The Kaplan-Meier techniques were applied to the created datasets of each cohort by importing the library of "KaplanMeierFitter" from the lifeline package. The datasets were prepared to have an "event" column which stated whether the particular tuple had experienced the event or not, and a "days" column, which had values of the number of days from the diagnosis of diabetes to the event's occurrence in days. The survival curves of each cohort and the differences between each cohort based on their demographic values were graphed to visualise the existing characteristics of the cohorts. The resulting Kaplan-Meier curves are included in the result section. The Kaplan-Meier curves were created for each complication with the stratification of demographic details. The curves were compared with a log-rank test to check the statistical differences among the strata of each demographic value sector. The results of the log-rank tests of demographic strata are included in the results section. Moreover, the code used for generating the Kaplan-Meier curve of the cohort of T2DM is included as follows.

```

In [9]: from lifelines import KaplanMeierFitter
from lifelines.statistics import pairwise_logrank_test
from lifelines.statistics import multivariate_logrank_test
from lifelines.utils import restricted_mean_survival_time

kmf = KaplanMeierFitter()
x= t2dm_final_col['days']
y= t2dm_final_col['event']

def plot_km(col):
    ax = plt.subplot(111)
    for r in t2dm_final_col[col].unique():
        ix = t2dm_final_col[col] == r
        kmf.fit(x[ix], y[ix], label =r)
        kmf.plot(ax=ax)

def print_logrank(col):
    log_rank = pairwise_logrank_test(t2dm_final_col['days'], t2dm_final_col[col], t2dm_final_col['event'])
    return log_rank.summary

```

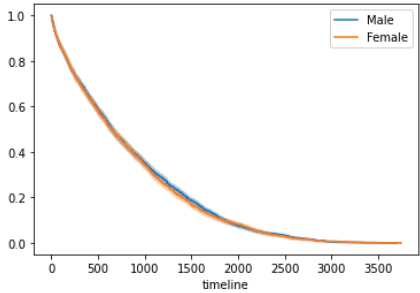
Figure 5-5 : Sample code for modelling the Kaplan-Meier in to the cohort of T2DM.

The used data frame for this code was “t2dm_final_col” which consisted of the details of all T2DM patients in the cohort. Two functions—plot_km(col), print_log_rank(col)—were used to plot the Kaplan Meier curves with different demographic details and measured their log-rank values respectively. Figure 5.4 shows the usage of the created function.

```

In [14]: plot_km('Patient_Gender')

```



```

In [15]: print_logrank('Patient_Gender')
Out[15]:

```

	test_statistic	p	-log2(p)
Female Male	0.792966	0.373205	1.421961

```

In [16]: print_wilcoxon_logrank('Patient_Gender')
Out[16]: <bound method StatisticalResult.print_summary of <lifelines.StatisticalResult: multivariate_Wilcoxon_test>
t_0 = -1
null_distribution = chi squared
degrees_of_freedom = 1
test_name = multivariate_Wilcoxon_test

---
test_statistic  p  -log2(p)
0.63 0.43 1.22>

```

Figure 5-6 : Sample code for calling the functions in non-parametric survival analysis for T2DM.

Additionally, the Kaplan-Meier models for each complication were created and saved for future system use. The following code generates the Kaplan Meier models for each complication.

```
In [130]: for k in ICD10_Code_comp_mapper.keys():
comp_id_en = ICD10_Code_comp_mapper[k]
temp = t2dm_final_col[t2dm_final_col["Diagnosis ICD10 Code_comp"] == comp_id_en]

# Initialize the Kaplan-Meier model
kmf = KaplanMeierFitter()

# Fit the model to your data
kmf.fit(temp['days'], temp['event'])

# Save the data frames
df_files = f'Trained_Models/Kaplan/Data_Frames/comp_{k}_DF.pkl'
open_file = open(df_files, "wb")
pickle.dump(temp, open_file)

# pickle.dump(kmf, open(df_files, 'wb'))

# save the model to disk
filename = f'Trained_Models/Kaplan/comp_{k}_model.sav'
pickle.dump(kmf, open(filename, 'wb'))
```

Figure 5-7 : Code snippet used to generate the Kaplan Meier models for each complication.

The following code was used to plot the curves and measure the log-rank values of the survival curves of each complication.

```
In [253]: def plot_km_comp(data_f,col):
x= data_f['days']
y= data_f['event']
name=data_f['Diagnosis ICD10 Code_comp'].unique()
ax = plt.subplot(111)
for r in data_f[col].unique():
ix = data_f[col] == r
kmf.fit(x[ix], y[ix], label =r)
kmf.plot(ax=ax)

df_name = str(name)
col1= str(col)
plt.title("Title : " + df_name + col1)
plt.savefig(f'Trained_Models/Kaplan/{name}_{col}.png')

def print_logrank_comp(data_f,col):
log_rank = pairwise_logrank_test(data_f['days'], data_f[col], data_f['event'])
return log_rank.summary
```

Figure 5-8 : Code snippet for using Kaplan Meier model in complications.

5.3.5.2 Semi-Parametric Techniques in Survival Analysis

The Cox model was used to predict the survival rate of individuals using the considered covariates. The survival of individuals in each cohort of complications was predicted through demographic details and a combination of demographic and laboratory values. The Cox models

were created for each complication with the two sets of covariates mentioned above. The models, based only on demographic details, used four demographic values: gender, age, ethnicity, and Māori/non-Māori values. The next model type used a combination of demographic and laboratory values, such as the four demographic values of the previous model and HbA1c, cholesterol, HDL, LDL, triglycerides, and eGFR. The models' accuracy was validated using the 10-fold cross-validation method, where the training and testing data sets were made out of 80/20 proportions of the original dataset. The models' performances were measured using the C-index and Brier scores. The Cox models were implemented by importing the library of "CoxPHFitter" from the lifelines package. The function "predict_survival_function" to measure the mean survival of the cohort was used. The function "predict_partial_hazard" to predict the partial hazard rates for new or existing data points based on a fitted Cox proportional hazards model was also used. The code used to create the Cox model with demographic values for predicting the survival of E1122 is shown in Figure 5.7 to demonstrate the algorithms used for the system.

```
In [8]: from lifelines import CoxPHFitter
cph_E1122 = CoxPHFitter()
cph_E1122.fit(Cox_E1122, duration_col = 'days', event_col = 'event')
cph_E1122.print_summary()
```

model	lifelines.CoxPHFitter										
duration col	'days'										
event col	'event'										
baseline estimation	breslow										
number of observations	893										
number of events observed	893										
partial log-likelihood	-5167.51										
time fit was run	2023-03-21 09:21:18 UTC										

	coef	exp(coef)	se(coef)	coef lower 95%	coef upper 95%	exp(coef) lower 95%	exp(coef) upper 95%	cmp to	z	p	-log2(p)
Patient_Age_OnDiagnosis	0.01	1.01	0.00	0.01	0.02	1.01	1.02	0.00	3.76	<0.005	12.54
Patient_Ethnicity	0.06	1.06	0.04	-0.02	0.14	0.98	1.15	0.00	1.42	0.15	2.70
Patient_Maori_NonMaori	-0.12	0.89	0.08	-0.28	0.04	0.76	1.04	0.00	-1.48	0.14	2.83
Patient_Gender	-0.00	1.00	0.07	-0.14	0.13	0.87	1.14	0.00	-0.03	0.97	0.04

Concordance	0.56
Partial AIC	10343.03
log-likelihood ratio test	22.74 on 4 df
-log2(p) of ll-ratio test	12.77

Figure 5-9 : Code snippet used for creating Cox model for E1122 with demographic details.

The created model "cph_E1122" is permanently saved for future usage using the following code.

```
In [12]: # Save the model to disk
pk.dump(cph_E1122, open(f'/Users/anuradhamadurapperuma/Documents/thesis_code/Thesis_codes/Trained_Models/CPH_Basic/
```

Figure 5-10 : Saving created data model of E1122.

Additionally, Figure 5.9 shows the predicted result of survival using the cph_E1122.

```
In [14]: t = range(0, 3500)
survival_func = cph_E1122.predict_survival_function(Cox_E1122.iloc[0:893], times=t)

In [15]: survival_func

Out[15]:
```

	2504	2510	2519	2521	2523	2526	2536	2537	2547	2549	2552	2556	2557	2562	2563
0.0	0.997334	0.997926	0.997735	0.997588	0.997818	0.998031	0.997981	0.998306	0.997432	0.997981	0.997981	0.998468	0.997846	0.998002	0.998002
1.0	0.997334	0.997926	0.997735	0.997588	0.997818	0.998031	0.997981	0.998306	0.997432	0.997981	0.997981	0.998468	0.997846	0.998002	0.998002
2.0	0.997334	0.997926	0.997735	0.997588	0.997818	0.998031	0.997981	0.998306	0.997432	0.997981	0.997981	0.998468	0.997846	0.998002	0.998002
3.0	0.996000	0.996888	0.996601	0.996382	0.996727	0.997046	0.996971	0.997459	0.996148	0.996970	0.996971	0.997701	0.996769	0.997002	0.997003
4.0	0.994669	0.995851	0.995469	0.995177	0.995636	0.996062	0.995963	0.996612	0.994865	0.995961	0.995963	0.996935	0.995692	0.996003	0.996004
...
3495.0	0.000193	0.001290	0.000699	0.000436	0.000914	0.001812	0.001543	0.004380	0.000264	0.001540	0.001543	0.007357	0.001000	0.001646	0.001650
3496.0	0.000192	0.001287	0.000697	0.000435	0.000911	0.001807	0.001539	0.004370	0.000264	0.001536	0.001539	0.007342	0.000997	0.001642	0.001645
3497.0	0.000191	0.001283	0.000695	0.000434	0.000909	0.001803	0.001535	0.004360	0.000263	0.001532	0.001535	0.007327	0.000994	0.001637	0.001641
3498.0	0.000191	0.001280	0.000692	0.000432	0.000906	0.001798	0.001531	0.004350	0.000262	0.001528	0.001531	0.007312	0.000991	0.001633	0.001637
3499.0	0.000190	0.001276	0.000690	0.000431	0.000903	0.001793	0.001527	0.004340	0.000261	0.001524	0.001527	0.007297	0.000988	0.001629	0.001632

3500 rows x 893 columns

Figure 5-11 : Predicted results of survival of E1122 using Cox models with demographic details.

Moreover, the code for the creation of Cox model for predicting the survival of E1122 with laboratory values is present in Figure 5.10.

```
In [37]: from lifelines import CoxPHFitter
cph_E1122_All = CoxPHFitter()
cph_E1122_All.fit(E1122_Cox_Final_InDep, duration_col = 'days', event_col = 'event')
cph_E1122_All.print_summary()
```

model	lifelines.CoxPHFitter											
duration col	'days'											
event col	'event'											
baseline estimation	breslow											
number of observations	893											
number of events observed	893											
partial log-likelihood	-5156.82											
time fit was run	2023-03-21 09:49:43 UTC											
	coef	exp(coef)	se(coef)	coef lower 95%	coef upper 95%	exp(coef) lower 95%	exp(coef) upper 95%	cmp to	z	p	-log2(p)	
Patient_Age_OnDiagnosis	0.01	1.01	0.00	0.01	0.02	1.01	1.02	0.00	3.49	<0.005	10.99	
Patient_Ethnicity	0.07	1.07	0.04	-0.01	0.15	0.99	1.16	0.00	1.68	0.09	3.42	
Patient_Maori_NonMaori	-0.13	0.88	0.08	-0.29	0.03	0.75	1.03	0.00	-1.62	0.11	3.25	
Patient_Gender	0.02	1.02	0.07	-0.12	0.16	0.88	1.17	0.00	0.26	0.80	0.33	
HbA1c_ResultValue	0.00	1.00	0.00	-0.00	0.01	1.00	1.01	0.00	1.04	0.30	1.74	
Cholesterol_ResultValue	0.13	1.14	0.09	-0.05	0.31	0.95	1.37	0.00	1.41	0.16	2.66	
Triglyceride_ResultValue	-0.05	0.96	0.05	-0.13	0.04	0.87	1.04	0.00	-1.00	0.32	1.66	
HDL_ResultValue	-0.47	0.63	0.15	-0.77	-0.17	0.46	0.85	0.00	-3.05	<0.005	8.76	
LDL_ResultValue	-0.05	0.96	0.10	-0.25	0.15	0.78	1.17	0.00	-0.45	0.65	0.61	
EGFR_ResultValue	-0.00	1.00	0.00	-0.01	-0.00	0.99	1.00	0.00	-2.51	0.01	6.38	
Concordance	0.57											
Partial AIC	10333.64											
log-likelihood ratio test	44.13 on 10 df											
-log2(p) of ll-ratio test	18.29											

Figure 5-12 : Code snippet used to create the Cox model for E1122 with all the features.

The created model was stored permanently for use in the future predictions of the system.

```
In [41]: # Save the model to disk
pk1.dump(cph_E1122_All, open(f'/Users/anuradhamadurapperuma/Documents/thesis code/Thesis_codes/Trained_Models/CPH_A
```

Figure 5-13 : Saving the Cox model for E1122 which created for predicting the survival with all the features.

The following code shows the predicted survival rate with the created Cox model with all the features.

```
In [30]: E1122_predict = cph_E1122_All.predict_survival_function(E1122_Cox_Final_InDep)
E1122_predict
```

```
Out[30]:
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	
2.0	0.997346	0.997993	0.997601	0.997828	0.998224	0.998171	0.997985	0.998650	0.997412	0.998193	0.998558	0.998242	0.998097	0.998043	0.997697	0
4.0	0.994693	0.995986	0.995202	0.995655	0.996447	0.996342	0.995969	0.997299	0.994824	0.996384	0.997114	0.996482	0.996193	0.996085	0.995393	0
6.0	0.990719	0.992976	0.991608	0.992399	0.993783	0.993598	0.992946	0.995272	0.990948	0.993673	0.994948	0.993844	0.993339	0.993150	0.991942	0
7.0	0.988068	0.990967	0.989209	0.990225	0.992003	0.991765	0.990928	0.993917	0.988361	0.991861	0.993500	0.992082	0.991433	0.991190	0.989638	0
8.0	0.986741	0.989961	0.988009	0.989137	0.991112	0.990848	0.989918	0.993239	0.987067	0.990955	0.992776	0.991200	0.990478	0.990209	0.988485	0
...
3189.0	0.000967	0.005264	0.001886	0.003414	0.009632	0.008388	0.005148	0.029364	0.001149	0.008870	0.023041	0.010087	0.006908	0.005998	0.002423	0
3260.0	0.000661	0.003949	0.001337	0.002501	0.007469	0.006455	0.003857	0.024203	0.000793	0.006847	0.018741	0.007841	0.005260	0.004532	0.001742	0
3319.0	0.000361	0.002499	0.000774	0.001524	0.004983	0.004255	0.002436	0.017794	0.000439	0.004535	0.013491	0.005252	0.003409	0.002901	0.001030	0
3352.0	0.000123	0.001106	0.000292	0.000631	0.002422	0.002024	0.001074	0.010285	0.000153	0.002176	0.007509	0.002571	0.001573	0.001310	0.000404	0
3745.0	0.000009	0.000151	0.000027	0.000073	0.000416	0.000330	0.000145	0.002696	0.000012	0.000362	0.001795	0.000449	0.000238	0.000188	0.000041	0

702 rows x 893 columns

Figure 5-14 : Code snippet of predicting the survival rate for E1122 with all the features

Further, the results of Cox models in predicting the complications are included in the result section.

5.3.6 Implementation of the Web Portal

The web portal of this project was created by adopting the waterfall method of SDLC for the reasons mentioned above. The requirements were gathered as the first step of the SDLC through the clients' meetings, brainstorming sessions, and literature review. The gathered requirements were analysed to ensure feasibility by continuously discussing with the client and brainstorming ideas. The selected feasible requirements are described in Section 5.2. The use case diagram in Section 5.2 clearly illustrates the functionalities of the suggested solution. The sole purpose of the artefact, which is the prediction of complications of diabetes, was done with the aforementioned suitable prediction methods. The web portal was the tool which combined the outputs of all the steps mentioned above of the solution implementation. The web portal was designed to present the tool as user-friendly, which should have been self-explanatory to the users. We developed the tool in a simple manner with no separate tabs, accordions, multi-step forums, or hidden or non-discoverable features, making it a complex and non-user-friendly tool.

The web portal we created used a simple sidebar to provide the full tool functionalities, which could hardly be missed. Two selection boxes were used to select the complication and prediction method, with carefully selected wording to keep it simple. One button was scheduled to pop up after selecting the user's purpose to print the resulting survival rates as a table and a graph. The main body of the website was used to show the statistics of the cohorts, and when the options for prediction were selected, the forums were loaded there. The prediction results were also displayed in the main body of the website. Additionally, after submitting a report request, the created report was available to download through a simple link at the end of the page. The whole website was implemented using the Streamlit library of the Python language. Since the entire project was conducted through Python, Streamlit was a good option for creating a simple website. The website was hosted through a GitHub account. The link for the website package at the GitHub account is attached here (<https://github.com/AnuradhaMadurapperuma/NZTPCD---web-site/tree/main>). The designed interfaces of the website are presented as follows.

The home screen was comprised of a general description of DM and their statistics in the region. The home screen was presented in two screen shots of the website as follows.

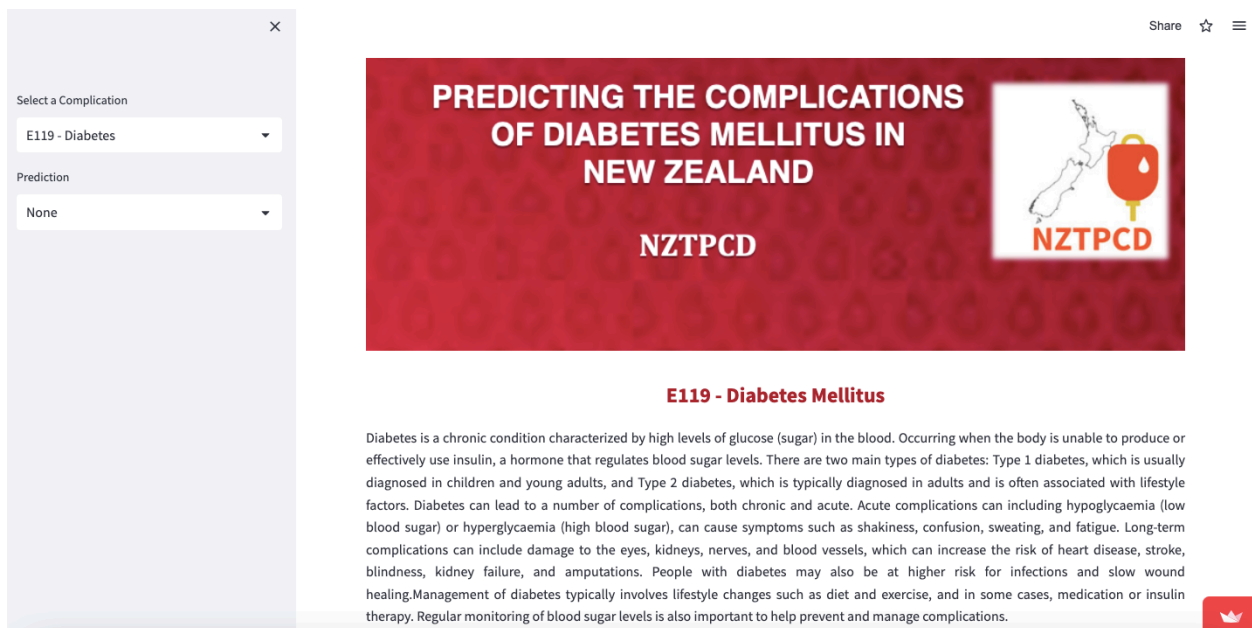


Figure 5-15 : Home screen of the NZTPCD - Part1.

The following two screenshots present the visualisation of statistics of the considered diabetes cohort. The visualisation of gender and ethnicity distribution among the cohort is presented in the following diagram.

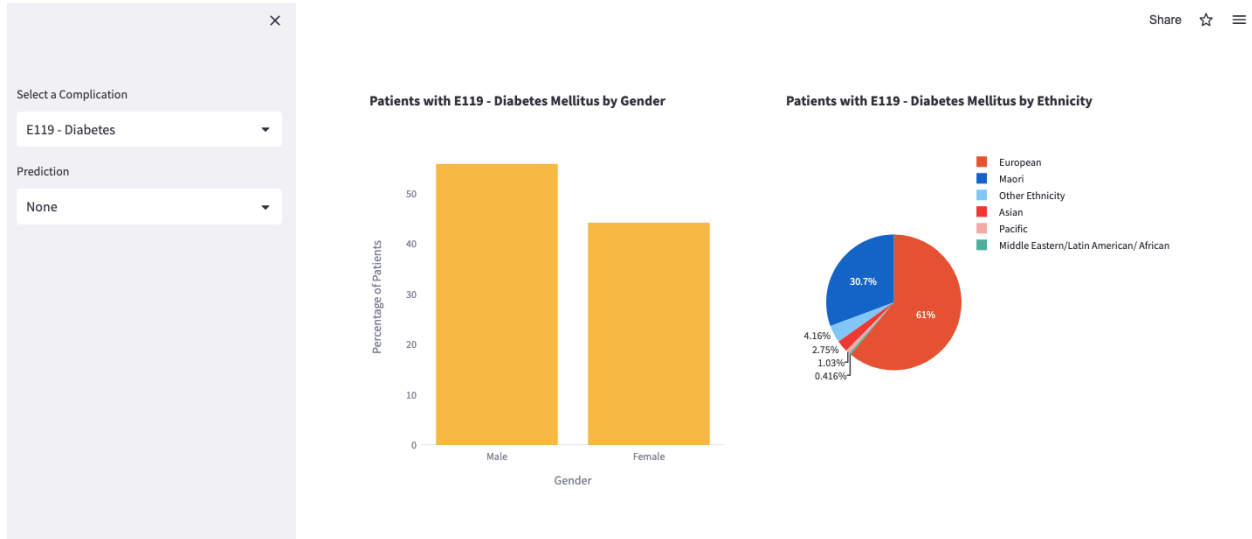


Figure 5-16 : Home screen of NZTPCD -Part 2- Visualisation of the distribution of gender and ethnicity.

Distribution of the factors of Māori/non-Māori and the age on the diagnosis are presented in the following screenshot.

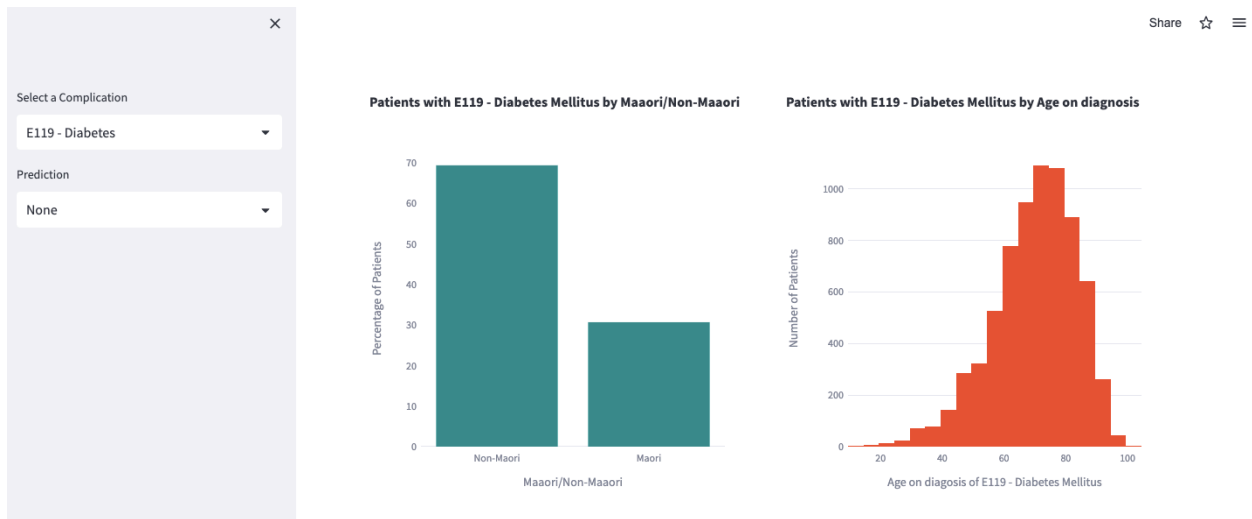


Figure 5-17 : Home screen of NZTPCD- Part 3- Distribution of Māori/Non-Māori and age on diagnosis.

The general survival curve of the cohort for diabetes mellitus was included in the end of the home screen of NZTPCD.

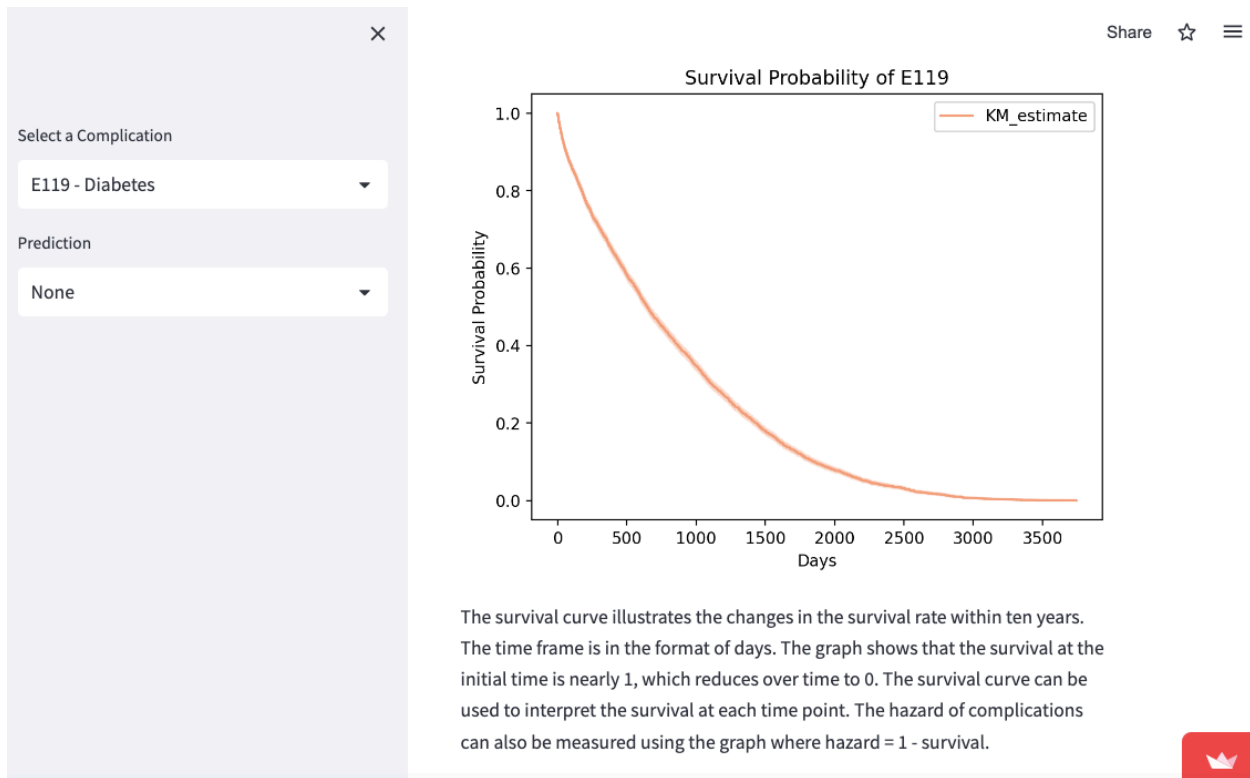


Figure 5-18 : Home screen of NZTPCD- Part 4 - Survival curve of diabetes.

The default setting of the website presented the details of diabetes. The drop-down list on the left-hand side can be used to select the complication and the type of prediction. There are 10 options for the complications: E1122 – Diabetic Nephropathy, E1122 – Diabetic Nephropathy, E1131 – Background Retinopathy, E1139 – Other Ophthalmic Complications, E1142 – Diabetic Polyneuropathy, E1151 -PVD, E1164 – Hypoglycaemia, E1165 – Poor Control – Hyperglycaemia, E1171 – Microvascular and other specified nonvascular complications, E1172 – Fatty Liver, and All Complications. The content of the website will dynamically change with the selected complication type. A general description of the complication, the same set of statistical representations, and the survival curve for the particular complication takes place at the right-hand side of the website. A sample of the website with the complication of E1122, is presented as follows to express the fundamental concept of the website.



E1122 - Diabetic Nephropathy

Diabetic nephropathy is a complication of diabetes that affects the kidneys, and it is one of the leading causes of end-stage renal disease (ESRD) worldwide. In diabetic nephropathy, high levels of blood sugar damage the small blood vessels in the kidneys, leading to progressive kidney damage and impaired kidney function. Over time, this can lead to proteinuria (excretion of protein in the urine), hypertension, and ultimately, kidney failure. Management of diabetic nephropathy involves controlling blood sugar levels, blood pressure, and other risk factors that can worsen kidney function. Treatment options may include lifestyle modifications, medications, and in severe cases, kidney transplant or dialysis. Early diagnosis and treatment can help slow the progression of the disease and reduce the risk of complications.

Figure 5-19 : NZTPCD user interface with the selection of E1122 Complication - Part1.

The representation of the distribution of gender and ethnicity of the cohort for the E1122 complication on the website is shown in the Figure 5.18.

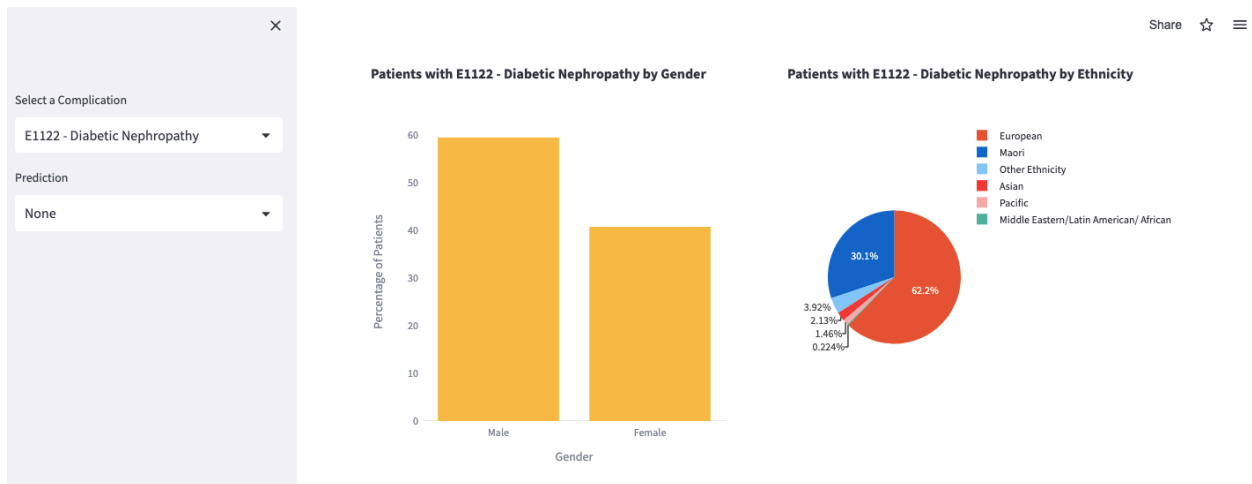


Figure 5-20 : NZTPCD user interface with the selection of E1122 Complication- Part 2.

Figure 5.19 shows the distribution of Māori/non-Māori and the age of diagnosis of E1122, which are represented on the website.

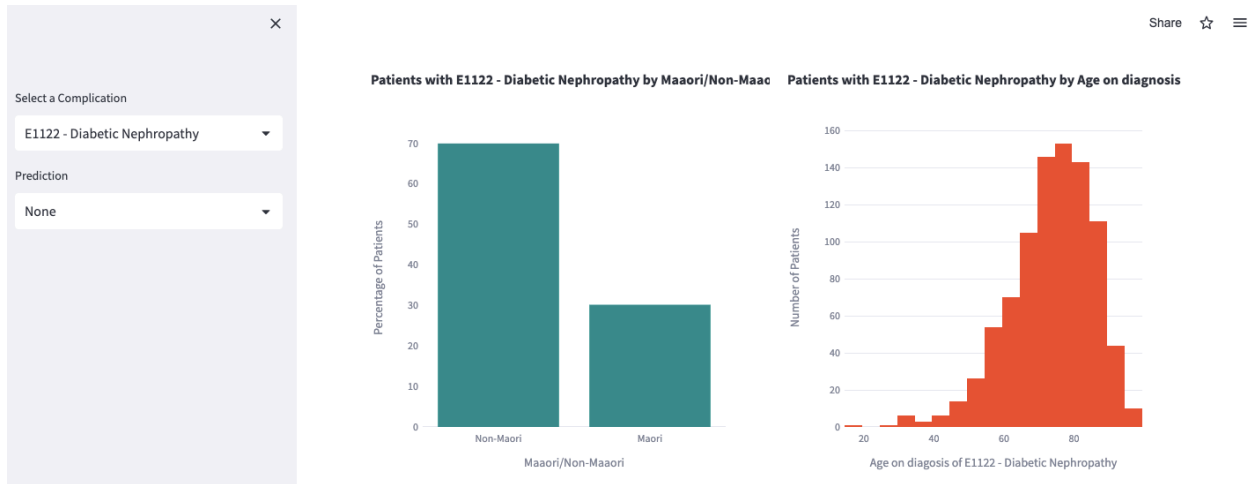
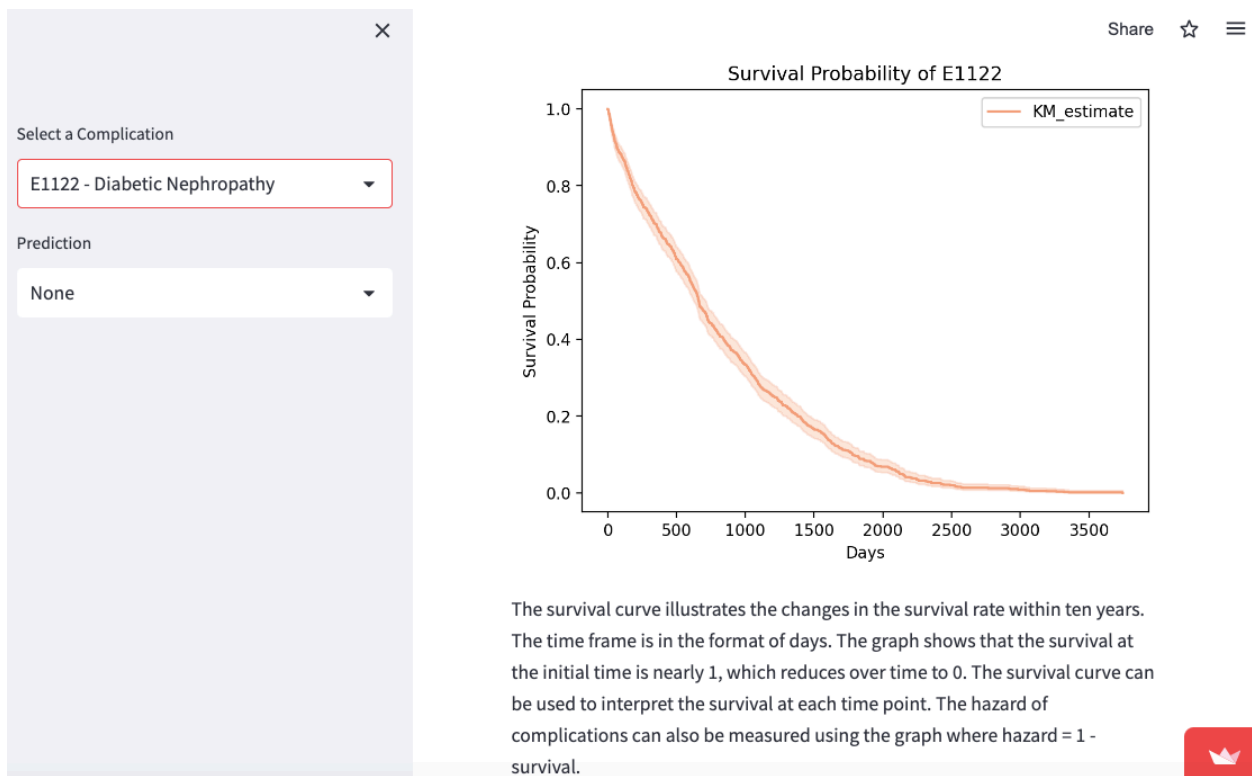


Figure 5-21 : NZTPCD user interface with the selection of E1122 complication- Part 3.

The survival curve of E1122 included at the end of the web page is illustrated in Figure 5.20.



The survival curve illustrates the changes in the survival rate within ten years. The time frame is in the format of days. The graph shows that the survival at the initial time is nearly 1, which reduces over time to 0. The survival curve can be used to interpret the survival at each time point. The hazard of complications can also be measured using the graph where hazard = 1 - survival.

Figure 5-22 : NZTPCD user interface with the selection of E1122 - Part 4.

Additionally, a drop-down menu on the left-side of the website provides the option to predict the survival of individual for each complication with two options: demographic details and laboratory values. The forms for each option were loaded into the right-sided context space of the website. The following figures illustrate the content of the website with the prediction options. The illustration of the interface for predictions based on the demographic details of the individuals is included as follows.

The screenshot shows a web interface for predicting survival. On the left, a sidebar contains a close button (X), a 'Select a Complication' dropdown menu with 'E1122 - Diabetic Nephropathy' selected, a 'Prediction' dropdown menu with 'Demographic Details' selected, and a 'Want to create a report?' section with a 'Create Report with demographic data' button. The main content area is titled 'E1122 - Diabetic Nephropathy' and contains a form with the following fields: 'Enter your demographic details' header, 'Gender' dropdown (Male), 'Maaori' dropdown (Maaori), 'Ethnicity' dropdown (European), 'Enter your age at diagnosis of Diabetes' text input (0), 'Enter your current age' text input (0), and a 'Submit' button. In the top right corner, there are 'Share', a star icon, and a hamburger menu icon. A red button with a white crown icon is visible in the bottom right corner.

Figure 5-23 : NZTPCD user interface for prediction with demographic details.

The demographic form consisted of gender, ethnicity, Māori/non-Māori, age at diagnosis of diabetes, and current age.

Figure 5.22 shows the interface for predicting the survival of E1122 with laboratory values.

Figure 5-24 : NZTPCD user interface with the prediction of laboratory values.

The prediction with laboratory values consisted of 11 features: gender, ethnicity, Māori/Non-Māori, age at diagnosis of diabetes, current age, HbA1c (mmol/mol), Cholesterol (mmol/l), Triglyceride (mmol/l), HDL (mmol/l), LDL (mmol/l), and eGFR (mL/min/1.73m²) value. The unit of measurement for each laboratory value and an example value was given with a help button for each value button. The following figures show examples of prediction for a hypothetical patient.

The completed form for a hypothetical patient is included as follows.

Share ☆ ☰

E1122 - Diabetic Nephropathy

Enter your demographic details

Gender

Maaori

Ethnicity

Enter your age at diagnosis of Diabetes

Enter your current age

Close X

Select a Complication

Prediction

Want to create a report?

Figure 5-25 : NZTPCD user interface for prediction of the survival rate with demographic details.

The results predicted by the system are included in the following diagrams. The hazard of the complication for 10 years from now was tabulated and visualised in the first diagram (figure 5.24) while the second diagram (figure 5.25) presents the survival curves of the individual against the cohort.

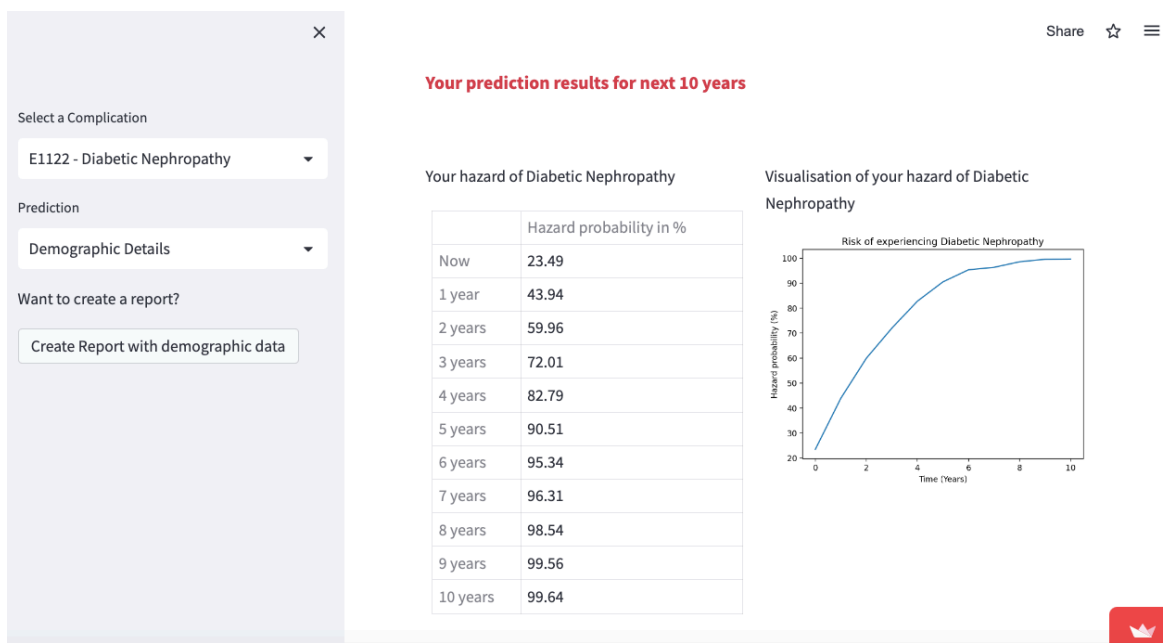


Figure 5-26 : NZTPCD user interface for the prediction of hazard of the complication using demographic details- Part 1.

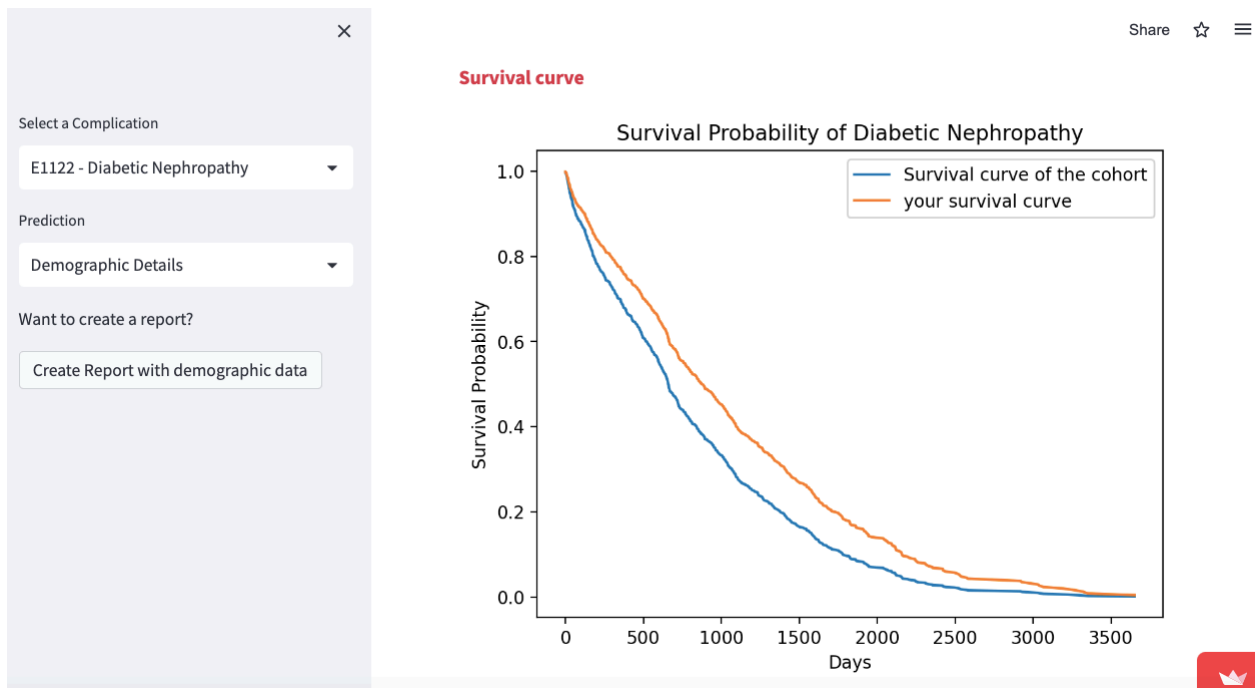


Figure 5-27 : NZTPCD user interface for prediction of complication using demographic details – Part 2.

The following diagrams present the prediction with laboratory values of a hypothetical patient. The first diagram is an example of a completed form with laboratory values.

The screenshot shows the NZTPCD User Interface for predicting survival using laboratory values. The interface is titled "E1122 - Diabetic Nephropathy" and is divided into two main sections: a left sidebar for report configuration and a main form for patient details and lab values.

Left Sidebar (Report Configuration):

- Select a Complication:** E1122 - Diabetic Nephropathy
- Prediction:** Lab Vales
- Want to create a report?:** Create Report with lab data
- Select a Complication:** E1122 - Diabetic Nephropathy
- Prediction:** Lab Vales
- Want to create a report?:** Create Report with lab data

Main Form (Enter your details here):

- Gender:** Female
- Maaori:** Non-Maaori
- Ethnicity:** European
- Enter your age at diagnosis of Diabetes:** 45
- Enter your current age:** 46
- HbA1c:** 50.00
- Cholesterol:** 5.50
- Triglyceride:** 1.60
- HDL:** 1.56
- LDL:** 2.60
- eGFR:** 49.00
- Submit**

Figure 5-28 : NZTPCD User Interface for predicting the survival using laboratory values.

The prediction results generated for the above hypothetical patient are included in the following two diagrams.

×

Select a Complication

E1122 - Diabetic Nephropathy

Prediction

Lab Vales

Want to create a report?

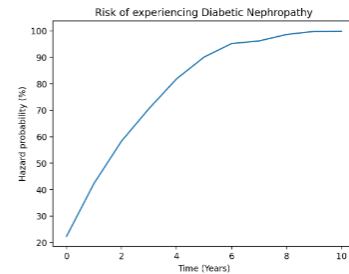
Create Report with lab data

Your prediction results for next 10 years

Your hazard of Diabetic Nephropathy

	Hazard probability in %
Now	22.41
1 year	42.40
2 years	58.37
3 years	70.65
4 years	81.86
5 years	90.08
6 years	95.20
7 years	96.19
8 years	98.63
9 years	99.80
10 years	99.86

Visualisation of your hazard of Diabetic Nephropathy



Share ☆ ≡

Figure 5-29 : NZTPCD user interface for the prediction results for laboratory-based survival prediction - Part 1.

The survival curves resulted from the laboratory-based prediction.

×

Select a Complication

E1122 - Diabetic Nephropathy

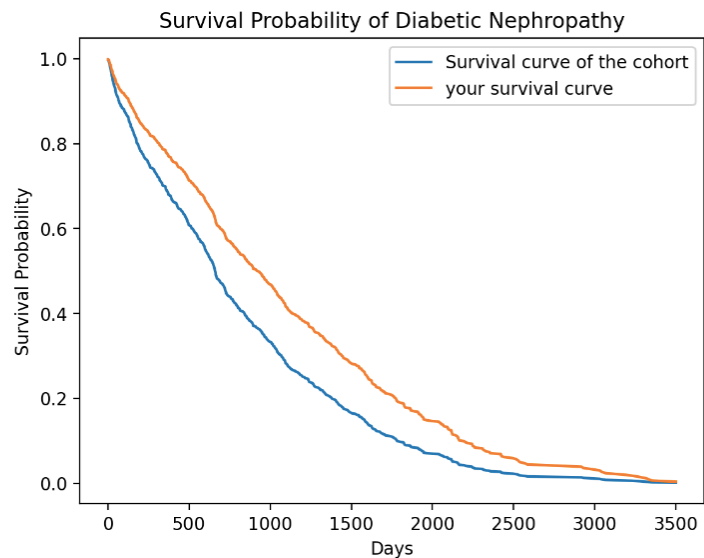
Prediction

Lab Vales

Want to create a report?

Create Report with lab data

Survival curve



Share ☆ ≡

Figure 5-30 : NZTPCD user interface for prediction results of laboratory-based prediction - Part 2.

A similar content on the website is resulted with the selected complications and prediction types while

website predicting the survival with given values.

The selection of the option “All Complications” from the drop-down menu with prediction options behave differently than selecting a regular complication. The content of the website for the option of “All Complications” shows the results of the predicted hazard rates and survival curves for all the complications. Moreover, the prediction results can be downloaded as a pdf using the bottom most option of the website, which appeared with the prediction results. Four types of reports can be downloaded through the system: prediction report of a complication created with demographic details, prediction report of a complication created with all details and the same two types of reports with the prediction results of all complications.

Chapter 6 Evaluation of the Artefact

The study's suggested artefact is evaluated in three phases: design evaluation, algorithm evaluation and implementation evaluation. The well-known FEDS framework (Venable et al., 2016) is adopted here to select the appropriate evaluation approach for the study and conduct the evaluation process more systematically due to its vitality in the process of design science research methodology (DSRM). The chosen strategy and the adopted framework of evaluation of this study are included in the research methodology section. In summary, the selected evaluation approach is the technical risk and efficacy approach due to its suitability for the current study. The initial design evaluation phase checks the validity of the suggested artefact in its context. The artefact's relevance will be monitored in this phase to reduce the human/user risk of the system. The implemented algorithms for predicting the complications are evaluated to prove their accuracy using statistical measurements. The implemented system is assessed for product quality using the standard ISO 25010 matrix while checking its usability through user feedback. The details of each evaluation phase are included as follows.

6.1 Design Evaluation

Design evaluation is the first component of the evaluation process of this study. The design evaluation phase has been conducted with the approach introduced by Wieringa (2016). The method comprised four components: expected effects, expected values, trade-offs, and sensitivity. These four components were utilised in evaluating the suggested design solution in its context. The proposed solution's relevance was checked to reduce the system's potential social/user risk. The proposed solution was validated against the requirement set to check the effect of the solution in its context. This step was conducted to study the "expected effect" component in the adopted validation method. Table 6.1 represents the expected effect of the artefact through the user requirements.

Requirement	Suggest solution	Effect
Understand the dispersion of diabetes patients based on demographic details.	Visualise the dispersion of diabetes patients based on their gender, age, ethnicity, and Māori/non-Māori.	Clear representation of the characteristics of diabetes patients helps to make decisions on policymaking, resource allocation, designing workshops and awareness programmes, etc.
Understand the differences in survival rates among diabetes patients based on demographic details.	Visualise the survival curves of diabetes patients based on demographic characteristics.	The resulting survival curves illustrate the differences in survival rates based on the demographic details of the cohort. This affects the decisions made by stakeholders in the healthcare sector for analysing the severity of diabetes among the participants.
Understand the socio-demographic dispersion of patients with each complication.	Analyse and visualise the dispersion of demographic details of the cohorts with the complications.	The stakeholders can grasp information on the prevalence of demographic details among the cohorts of different complications to assist in making decisions.
Understand the differences in survival rates among diabetes patients for each complication based on demographic details.	Survival curves for each complication, which are based on their demographic details.	The resulting survival curves illustrate the contrast of survival curves based on the demographic details to understand the patterns and variations of survival rates among different cohorts of complications.
Understand the survival rates of each complication.	Visualisation of survival curves for each complication.	The survival curves provide information on the survival rates of each complication, which is beneficial in rating the complications on their severity and adversity on the health indices.
Analyse the survival rates of the cohort for complications based on demographic details.	Survival curves are created by considering the demographic details of the whole cohort.	The survival rate of the cohort for each complication can be visualised. This would be beneficial in understanding the changes in the survival of complications among the cohort.
Analyse the survival rates of the cohort for complications based on demographic and laboratory details.	A series of survival curves are built with the combined features set.	The resulting survival curves can be utilised in deciding the complications' differences.

Table (Continue)		
Requirement	Suggest solution	Effect
Predict the survival of individuals in each complication based on demographic details.	Cox model for predicting the survival and hazard of individuals using their demographic details.	The ability to predict an individual's survival chance assists in conducting laboratory tests, advising pre-cautionary actions, and deciding medications.
Predict the survival of individuals in each complication based on demographic and laboratory details.	Cox model for predicting the survival of an individual through a combined feature set.	
Generate a report for individual.	A report on the survival chances of individuals is created.	Communication between the patient and general practitioners is getting easy and productive.
A user-friendly approach to predicting the CoDM.	A web portal with clear graphical interfaces has been implemented.	The stakeholders can easily engage with the web portal without specific training.
Accurate predictions results.	Use statistically effective methods for predicting CoDM.	The prediction results can be used to assist decision-making.
Efficient system for predicting the CoDM.	Using a web portal to launch the CDSS makes the system more efficient than a stand-alone system.	The web portal efficiently provides the services as a CDSS.

Table 6-1 : Table of user requirement and the expected effects of the artefact.

The expected values were validated as the next step of evaluating the design solution. The expected values from the system were gathered during the client meetings. The client's expectation of the perspective of values was having a computerised system with less cost, and which was user-friendly and accurate. The characteristics of the design solution matched these clients' expectations. The CDSS was designed as a web portal rather than a stand-alone system to satisfy the client's expected values, which saved the system's cost by avoiding the installation cost and minimising the system specification and expert knowledge. The graphical user interfaces were designed without complex widgets and tools to keep them simple and user-friendly. The system's accuracy was achieved through the statistical accuracy of the algorithms. Since evaluating the solution design was to meet the client's requirements with the suggested overall solution, this was a summative evaluation phase of the solution design. The recognition and comparison of trade-offs were conducted as the next step of evaluating the design solution. The recognised potential alternatives were a stand-alone system, mobile application, electronic

health record (EHR) integrated system, and smart devices and wearables. The stand-alone system required installation and was vulnerable to updating the versions (Sommerville, 2011). The mobile applications were embedded with characteristics such as platform compatibility, limitations of screen size, varying compatibility with devices and performance, and the cost of development and maintenance (Mohammadzadeh & Safdari, 2014). The accessibility, confidentiality, and limitations of the existence of EHRs were unavoidable barriers to building an EHR-integrated system (Jardim, 2013). Although smart devices and wearable systems are the trend in the healthcare industry, their high cost of developing and maintaining, privacy and security concerns, and their suitability for the sole purpose of the study made these systems controversial. The last component of the evaluation phase was conducted to check the system's sensitivity.

With regard to sensitivity, the following question was asked: "Would the treatment still be effective and useful if the problem changes?" (Wieringa & Morali, 2012a, p. 228). The potential changes in a similar situation were considered here as the characteristics to evaluate the solution, such as changes in the data set and predicting a different complication or disease (Wieringa & Morali, 2012a, p. 228). The considered dataset was thoroughly pre-processed to enhance the quality of the dataset. The data pre-processing techniques preserved the data's accuracy, completeness, consistency, relevance, and validity. Since the dataset was a sophisticated resource created to fulfil a particular set of tasks, the alterations of the dataset may have affected the solution. Additionally, the system's algorithm was designed and maintained by considering the factors such as a relevant and informative feature set, tuning the Cox model for threshold values, and model calibration with the Brier score. The feature set was selected through a thorough systematic review which revealed the most commonly used attribute set for predicting the complications of diabetes. The created models of this research were sensitive to the demographic and laboratory details of the individuals. Therefore, it could be argued that the models were sensitive to their demographic information, such as gender, ethnicity, and Māori/non-Māori. The sensitivity could not be reached for a widely acceptable level, as the models were not trained with high resourceful datasets. Anyway, the sensitivity of the suggested

solution was adequate since it covered the cohort in Aotearoa the same as other models (UKPDS, FINDRISK) do. The solution proposed for the overarching research question was evaluated through the above-mentioned four criteria: expected effects, expected values, trade-offs, and sensitivity. The evaluation results revealed that the suggested solution covered most of the evaluation criteria while showing that the proposed solution of the study was the most appropriate solution for resolving the recognised real-world issue.

6.2 Algorithm Evaluation

Algorithm evaluation is the formative evaluation phase of the study's evaluation trajectory which consists of episodes for each algorithm used for predicting the complications. The accuracy of the survival analysis techniques used in this study to analyse and predict the survival of individuals were validated through this phase. The Kaplan-Meier survival techniques are used in analysing the survival of a cohort without using any parameter interference. The resulting survival curves were used for visualising the survival rate of the cohort. The prediction of survival was conducted through time-independent Cox models. The prognosis was conducted through two branches of feature, one with only the demographic characteristics (gender, age, ethnicity, Māori/non-Māori), and the other model with a combination of demographic and laboratory features (gender, age, ethnicity, Māori/non-Māori, HbA1c, Cholesterol, HDL, LDL, Triglycerides, eGFR). These algorithms were evaluated using the Akaike information and C-index values. The C-index values of each model were recorded to see the accuracy of each model in the prediction. The results are tabulated in the result section. The interpretation of the C-index ranges from 0 to 1, where 1 indicates perfect discrimination power and 0.5 indicates random prediction, which has no discrimination ability. As mentioned above, due to the existence of the dilemma of using the C-index in evaluating the Cox models, cross-validation and log-rank tests were used as the initial steps of assessing the accuracy of the models. The internal evaluation of the models with training and testing data sets provided each model's C-index and Brier score with the 10-fold cross-validation results. The results of the cross-validation of each model are represented in the result section.

6.3 Implementation Evaluation

The implementation evaluation is the summative evaluation process of the system to check the quality of the developed system. This phase consists of two episodes: technical and user feedback evaluation. The system has been technically evaluated through a standard quality measurement system: ISO 25010 (International Organization for Standardization, 2011). Additionally, the usability of the system has been confirmed with user feedback.

6.3.1 Evaluation through ISO 25010

The technical assessment of the implemented system was conducted against the standard product quality measurements, ISO 25010, to check the system quality of the product. The adopted method validated the requirements of the developed system against the standard characteristics and sub-characteristics of ISO 25010 to prove the existence of standard qualities in the developed system. Since the evaluation process was done with the existing system requirements, first, a list of current requirements of the system was created. Second, the system requirements had to be categorised into blocks. The degree of influence of the requirements to each category and sub-category in the quality of the product was measured using mathematical formulas introduced by Kadi et al. (2016). The requirements of this system were listed with their categorised blocks of requirements to evaluate against the standard matrix. The following table represents the system's requirements, acronyms, and categorisations.

Requirement block	Requirement No	Requirement
Data accessibility (DA)	DA1	installation server.
	DA2	Operating system (OS) type.
	DA3	OS version
	DA4	target user
	DA5	Cost
	DA6	Internet connection.
Personal information (PI)	PI1	age
	PI2	gender
	PI3	ethnicity
	PI4	Māori/non-Māori
Quantitative data (QD)	QD1	HbA1c
	QD2	Cholesterol
	QD3	HDL
	QD4	LDL
	QD5	Triglyceride
	QD6	eGFR
User's actions (UA)	UA1	Add information
	UA2	Assess the survival through demographic data
	UA3	Assess the survival through demographic + laboratory data
	UA4	Report downloads
	UA5	Image/graph downloads
CDSS's functionalities (DF)	DF1	Visualise the statistical details of the diabetes cohort
	DF2	Visualise the statistical details of the cohort of each complication
	DF3	Notes on complications and diabetes
	DF4	Predicting the selected CoDM using demographic details.
	DF5	Predicting the selected CoDM using demographic+ laboratory details.

Table 6-2 : System requirements from the ISO 25010 standard.

The considered characteristics and sub-characteristics are represented in Figure 6.1.

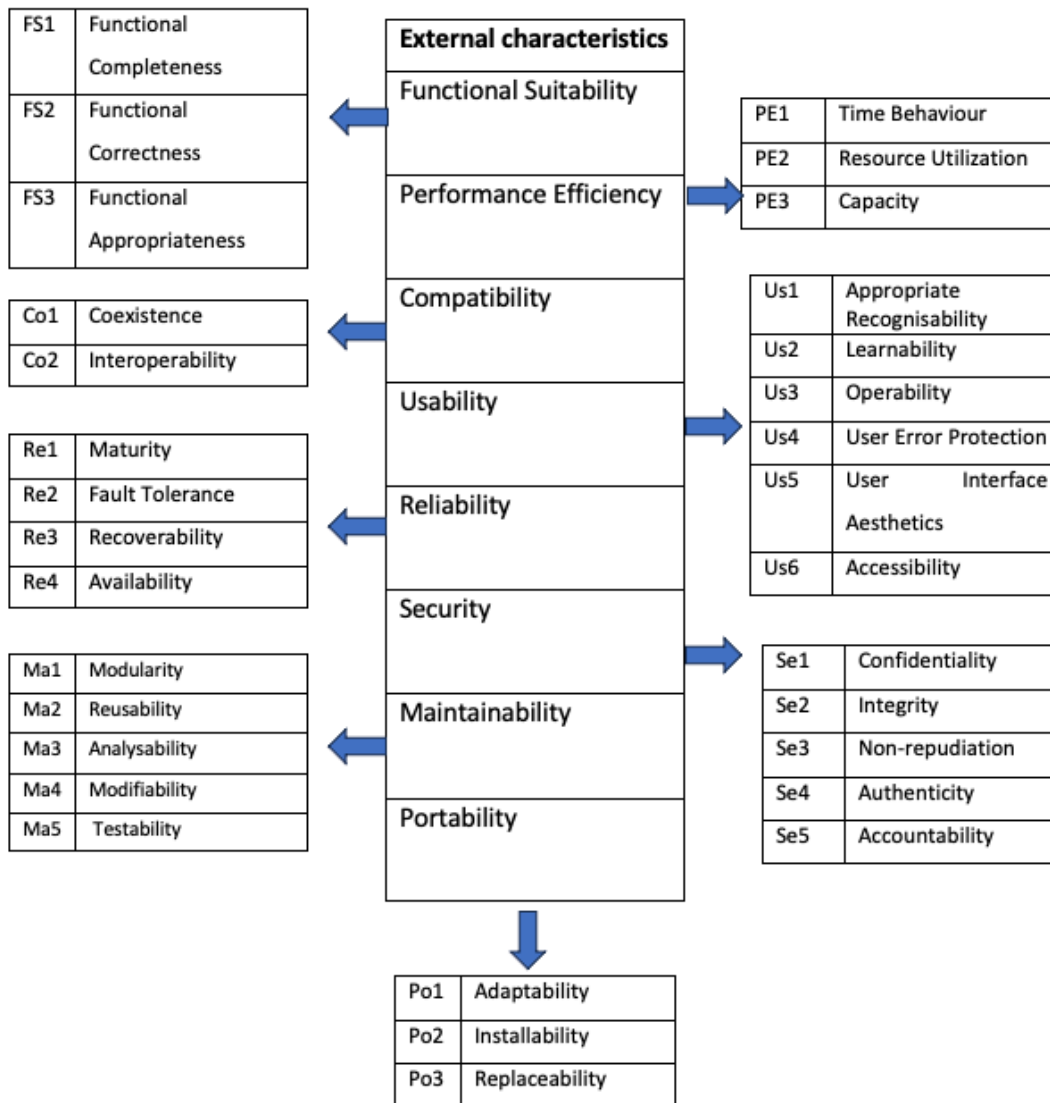


Figure 6-1 : External characteristics used for evaluating the system with ISO 25010 standard.(ISO/IEC 25010, 2011)

The influence of the requirements on the above-mentioned external characteristics and sub-characteristics is presented as a separate chart in the appendix section (Appendix C). The presence or absence of a requirement on each sub-category is represented as + and –, respectively, on the chart in appendix C.

The mathematical formulas used are as follows for calculating the degree of influence. The first equation calculates the degree of leverage of a block of requirements B on an external characteristic EC: (DI (EC, B))

$$DI (EC, B) = \Sigma DI(EC, R) / N(R) \quad (1)$$

N(R) = total requirements in block B.

The degree of influence of a block of requirements B on an external sub-characteristic EsC: (DI (EsC, B)) is determined through the following equation.

$$DI (EsC, B) = \Sigma DI (EsC, R) / N(R) \quad (2)$$

N(R) = total requirements in block B.

Additionally, the degree of influence of each requirement R on an external characteristic EC: (DI (EC, R)) can be calculated using the following equation.

$$DI (EC, R) = N (EsC, R) / N(EsC) \quad (3)$$

N (EsC, R) = number of sub-characteristics EsC of the external characteristic EC which are influenced by the requirement R

N(EsC) = total number of sub-characteristics of EC.

The calculated degree of influence of the block of requirement on external characteristics, sub-characteristics and the degree of leverage of the requirements on the external elements was

used to measure the product quality of the developed system. The results of the implementation evaluation are included in the result section, with detailed interpretations of the obtained results.

6.3.2 Evaluation through User Feedback

The artefact has been evaluated through users' feedback to embed the user perspective on the design and modelling of the system. The user-centred evaluation techniques use a confirmatory method to support the product quality of the artefact. The product quality of the artefact has been evaluated by adopting standard criteria. The created system of this study was a web-based clinical decision support system (CDSS), which could predict the survival of a selected set of complications while providing a statistical overview and informative visualisations of the cohort. Therefore, the usability evaluation methods (UEMs) commonly used in evaluating web applications were considered here. Usability has several definitions in various research fields. Those above-mentioned adopted standards for technical evaluation: ISO 25010 defined usability as the "degree to which specified users can use a product or system to achieve specified goals with effectiveness, efficiency and satisfaction in a specified context of use" (ISO/IEC 25010, 2011, p. 1). The UEMs of web applications were categorised into empirical and inspection methods. Empirical methods focus on collecting and analysing usage data from real users. In contrast, inspection methods focused on "reviewing the usability aspects of web artefacts, which are commonly user interfaces, concerning their conformance with a set of guidelines." (Fernandez et al., 2011, p. 4). The empirical methods were more resource-consuming and could be used to analyse the real participants' usage and behaviour. In contrast, inspection evaluation was a simple, cost-effective method that could be easily adopted in collecting user feedback on the system's design. We selected the inspection method for evaluating the system's usability by considering the above characteristics. Although the inspection method of evaluation uses a wide variety of techniques to perceive the system's usability, including questionnaires, interviews, focus groups etc., the questionnaire method was adopted due to the advantage it brought to the study. The questionnaire method was used by past scholars to collect the user perspective of the web systems (Cao et al., 2004; Chiew & Salim, 2003; Zaharias, 2006). Preparing the questionnaire, selecting the participants, and analysing the responses were important steps in conducting a user

feedback evaluation phase. The methods of conducting user feedback evaluation introduced by several studies was adopted (Chiew & Salim, 2003; EL-firjani et al., 2017; Tullis & Stetson, 2006). The questionnaire for evaluating the website usability was developed by considering various evaluation metrics. Successful task completion, mistakes, and time on task were considered by EL-firjani et al. (2017). Twenty usability criteria were extracted from the literature and sub-categorised in a study which intended to create a website usability evaluation tool (WEBUSE) (Chiew & Salim, 2003). Further, a system usability scale was introduced by Brooke (1996), which consists of 20 questions. The questionnaire developed by IBM (Computer System Usability Questionnaire) consists of 19 questions. The questionnaire developed in this study considered the characteristics emphasised in the literature, such as content, organisation and readability, navigation and links, user interface design, and performance and effectiveness (Chiew & Salim, 2003); navigation, learnability, accessibility, consistency, visual design, interactivity, content and resources, media use, learning strategies design, instructional feedback, instructional assessment, and learner guidance and support (Zaharias, 2006). The questionnaire of the study developed based on the commonly selected standard set of qualities: consistent, navigation, learnability, readability, comprehensive and content, organisation, accessibility, user interface design, performance, and user satisfaction. The qualities and the developed questions to check the quality are highly depend on the existing literature. (Chiew & Salim, 2003; EL-firjani et al., 2017; Tullis & Stetson, 2006; Zaharias, 2006). The evaluation of standard qualities of the web-portal of CDSS is specifically adopted a method introduced by Chiew and Salim (2003), which illustrates in the figure 6.2. Four major qualities were examined here: Content, organization and readability, Navigation and Link, User interface design and Performance and effectiveness. The potential questions for testing these qualities also extracted through a thorough literature review. The questions of the developed questionnaire and the mapping of questions into the qualities are presented in the Table 6.3.

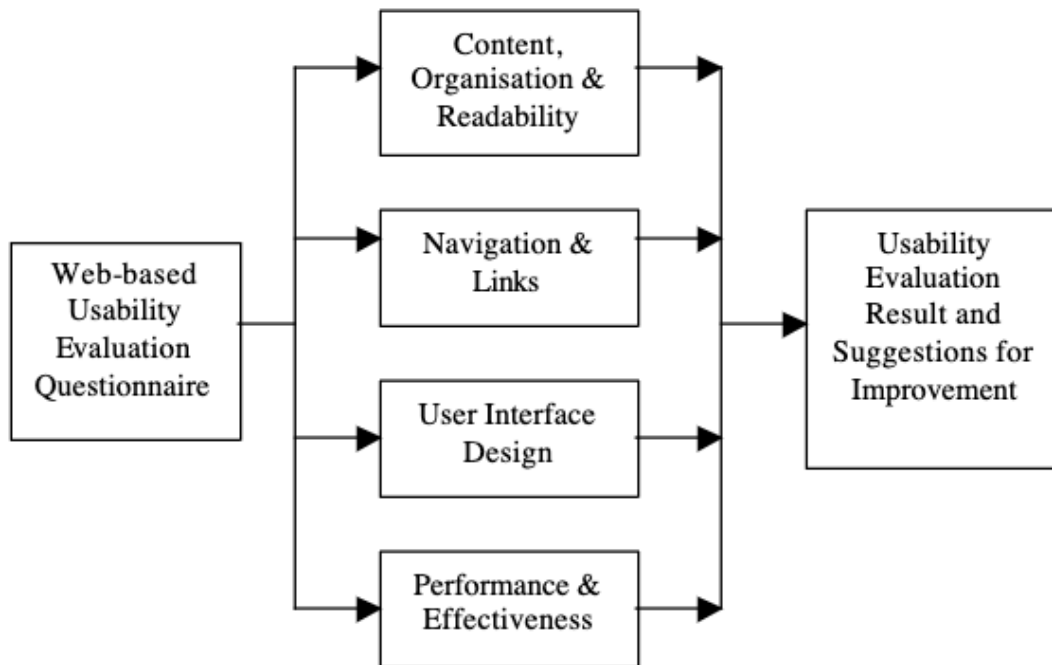


Figure 6-2 : Characteristics used for build the questionnaire.

Standard Characteristic	Question in the survey	Response type
Personal info	What is your age?	18-24,25-34,35-44,45-54,55-64,65-74, 75 or older
Personal Info	What is your gender?	Female, Male, prefer not to say
Personal Info	What is your occupation?	Text data field
Content, organization and readability	How easy was learning to use this CDSS?(EL-firjani et al., 2017)	Likert scale (1-5)
Content, organization and readability	I found the CDSS user interfaces easy to understand(EL-firjani et al., 2017)	Likert scale (1-5)
Content, organization and readability	I was able to complete the tasks quickly using the CDSS. (Chiew & Salim, 2003)	Likert scale (1-5)
Content, organization and readability	The CDSS had clear icons, labels, and menu options that were easily remembered.(EL-firjani et al., 2017)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)

Table (Continue)		
Standard Characteristic	Question in the survey	Response type
Navigation and Link	I easily recall the steps to perform the tasks in the CDSS (EL-firjani et al., 2017)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Performance and effectiveness	Overall, I am satisfied with the usability of the CDSS. (Chiew & Salim, 2003)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Performance and effectiveness	The CDSS meets my needs and expectations.(Alasadi & Bhaya, 2017)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Performance and effectiveness	The CDSS was accessible and usable with assistive technologies. (Chiew & Salim, 2003)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Performance and effectiveness	I encountered accessibility challenges or limitations while using the CDSS.	Boolean response (Yes, No)
Performance and effectiveness	If the answer for Q12 is yes, please specify the encountered challenges and limitations	Open response textual response
User interface design	The CDSS was visually appealing(Chiew & Salim, 2003)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Performance and effectiveness	The CDSS provided a pleasant user experience. (Chiew & Salim, 2003)	Likert scale (1-5)
User interface design	The CDSS had a well-organised and intuitive layout.(EL-firjani et al., 2017)	Likert scale (1-5)
Navigation and Link	The CDSS had consistent design elements and navigations.(EL-firjani et al., 2017)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Performance and effectiveness	The CDSS has a good system response time.(Chiew & Salim, 2003)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Content, organization and readability	The CDSS were reliable.(EL-firjani et al., 2017)	Likert scale (Strongly disagree, disagree, neutral, Agree, Strong agree)
Open ended question to gather additional comments.	Please provide any additional comments suggestions or concerns about the usability of the CDSS, and any features you think are missing or would be advantageous compared to existing prediction models. (Self-defined)	Open response textual response

Table 6-3 : Mapping the questions of the questionnaire with standard characteristics

Further, the developed questionnaire for this study is attached in the appendix D.

The participants for the study were selected to cover the potential stakeholders. The selection of participants was personalised to match the research purposes. The purposive sampling method was selected due to its characteristics, such as ability to deliberately focus on individuals with specific traits or experiences relevant to the research question, thereby ensuring that the sample is composed of participants who are most likely to provide valuable insights for the study. Moreover, Purposeful sampling is a commonly employed technique in qualitative research to identify and select cases that provide rich information, maximizing the effective use of limited resources (Palinkas et al., 2015). The limitation of resources, the purpose of conducting the user feedback evaluation, and the flexibility of selecting participants of this method were well suited to the problem at hand. Since the goal of this evaluation was to gather specific insights rather than to generalize findings, purposive sampling was employed. The selection process focused on identifying individuals within the community who were expected to provide the most relevant and insightful information for the study (Campbell et al., 2020). The foundations of selecting a set of users for a usability test of websites was familiarised with similar systems, possessing good computer literacy, and covering the potential type of stakeholders. Forty students were randomly selected in WEBUSE to test the website's usability at the Faculty of Computer Science and Information Technology, University of Malaya (Chiew & Salim, 2003). Five testers with good computer literacy were selected by EL-firjani et al. (2017), and 123 employers of the company were selected to compare the questionnaire of website usability evaluation (Tullis & Stetson, 2006). Participants for evaluating the developed CDSS in this study were selected by considering the potential stakeholders in the healthcare sector. Since the convenience sampling method was a non-probabilistic sampling method and the purpose of this evaluation was not to analyse the overall opinion of the population, the sample size could be selected based on availability and convenience. Responses were gathered from seven participants to evaluate the usability of the CDSS. Two general practitioners, three nursing lecturers, and two research officers volunteered to evaluate the system through the questionnaire. The two general practitioners and the three

nursing lecturers brought the subject knowledge and contextual awareness to the study while the researchers providing an external perspective of the system. Their computer literacy, knowledge of working with similar systems and expectations of the improvements of existing systems provide a vital perspective of the implemented system. The career experiences and the advantageous age range support to strengthen their responses and opinion of the CDSS.

The questionnaire was created as a Google form due to its popularity, simplicity, and availability. The collected responses were analysed to check the usability of the CDSS. The questionnaire used for this purpose was attached in the appendix section (Appendix D). The questionnaire consisted of 20 questions created to get the user information, their evaluation of the CDSS, and recommendations. The questionnaire was designed with a Likert scale of five responses such as “Strongly disagree”, “Disagree”, “Neutral”, “Agree”, and “Strongly agree”. Moreover, some responses were collected through a scaler chart which ranged from 1 to 5, where 1 represented the least and 5 represented the most. The only personal details recorded here were the participants' age, gender, and occupation. The system's usability was checked by requesting responses about its interfaces, used icons, navigation pathways, ease of performing tasks, etc. Finally, an open-ended question was asked to collect the users' suggestions and recommendations for the system.

6.4 Summary

The evaluation is a critical step in design science research emphasised through decades. The design science research approaches describe the evaluation as evaluating the suitability of the proposed design to solve the real-world issue and evaluating the implemented solution as the system that performs in the context. The concept is also known as evaluating as being rigour and relevant. The current study's evaluation process adopted the technical risk evaluation strategy introduced in the FEDS framework. The evaluation is conducted in three sections, where the suitability of the suggested solution is evaluated in solution design, the accuracy of the algorithms is checked in the algorithm evaluation phase, and the final implemented solution is evaluated using two episodes: technical and user feedback evaluation. The design evaluation phase

confirmed the validity of the suggested solution for the problem at hand. The algorithm evaluation provided statistically significant utilised algorithms. The implementation evaluation phase evaluated the product quality via a standard technical process while confirming its usability by collecting user feedback through a questionnaire. The evaluation of the above-mentioned three pillars of the study systematically validated the overarching research question.

Chapter 7 Research Findings and Interpretations

This chapter provides the results of the research study with their interpretations. The results gained through the exploratory data analysis, survival analysis, the results of generated algorithms, and the findings of implementation evaluation are included here to answer the research questions. The results are presented according to their relevance in answering the research questions. The research questions and the sub-research questions are included here to provide a clarified image of the study.

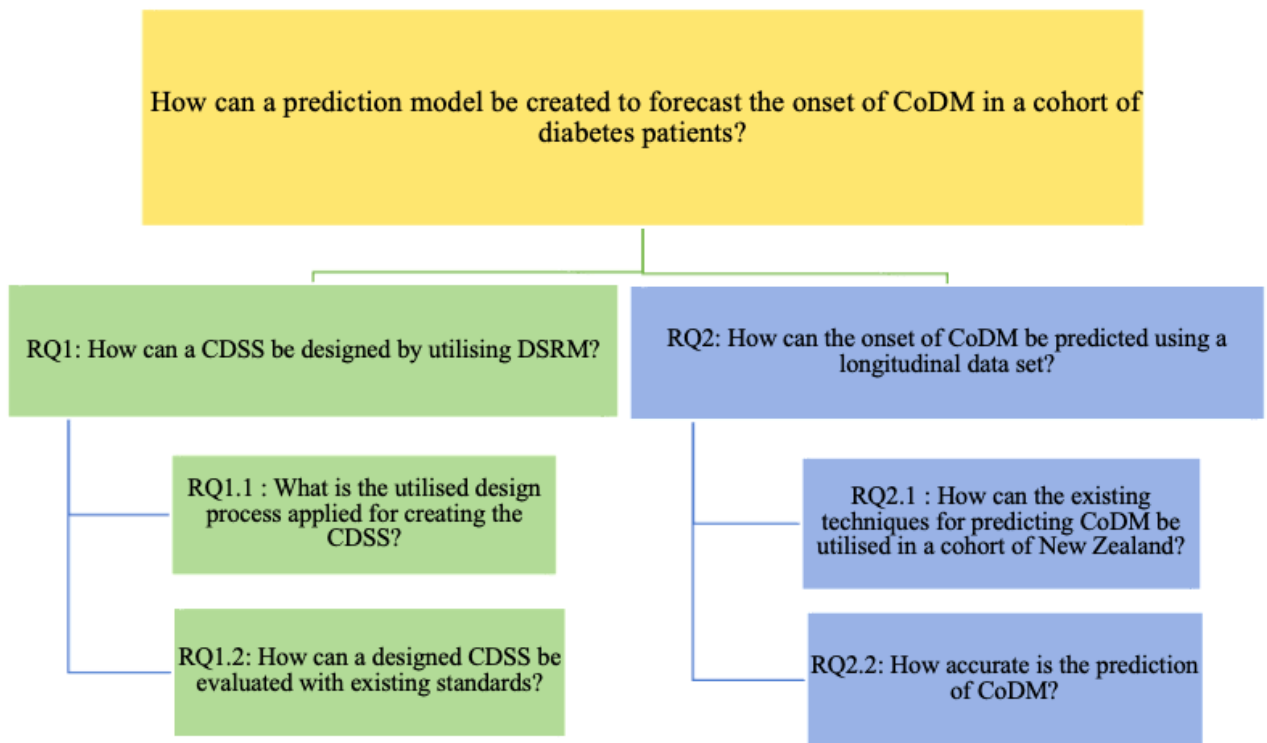


Figure 7-1 : Diagram of research questions.

A set of research objectives were achieved on the way to implement the clinical decision support system (CDSS). The objectives and their contribution are described in the following section. The

overarching research question of this study is divided into two main categories based on its design perspective and data analysis perspective. Further, the epistemological and evaluation perspectives of both categories are considered in the research study while making two sub-research questions from each category. The following section presents the research findings respective to the above research questions.

7.1 RQ1: How can a CDSS be designed by utilising DSRM?

The design perspective of the information system is focused on in this section of the study. A thorough literature review directs the research into its most appropriate philosophical ground with a clear understanding of the best suitable research approach. The design perspective of CDSS is explained through two sub-research questions, where the former focuses on the design process, and the latter concentrates on the evaluation of the designed solution.

7.1.1 RQ1.1: What is the utilized design process applied for creating the CDSS?

The recognised real-world problem and the existing knowledge gap in the domain led the current study to create an improved CDSS as the solution. Further, the literature and philosophical understanding of the research concentrated on a confirmatory, design science-embedded solution. Design science research methodology (DSRM) specialises in scientifically resolving real-world issues with its fundamental guide: designing the information system, implementation of the designed system, and evaluation (Wieringa, 2014). Although DSRM provides a scientific guideline for resolving a real-world issue with a design, it has been challenging to adopt those concepts in an information system in healthcare management. Awareness of the existing issue, understanding the requirements of the client, and recognising the pitfalls of current information systems in the domain has been fulfilled in chapters 2 and 4. The extracted research gap from the literature review in Chapter 2 made the fundamental requirement of designing a CDSS, while the problem in-hand specifies the extended expectations, functional and non-functional requirements of the client. A clear image of situational awareness has been provided in Chapter 4 to understand the contextual requirement of the research study and situational understanding

of healthcare management systems in Te Whatu Ora Waikato. The recognised requirements and drawbacks of existing CDSSs are addressed by developing a CDSS through the process of developing this study. The DSRM approach introduced by Wieringa (2014) has been used in this study with the justifications of each steps. The steps of the empirical design cycle of the adopted method are justified in Chapter 3 and explained in Chapter 5. Additionally, the situational awareness phase of the study is described in chapter 4, with the identified milestones of the conducted process of the study. The use of literature review and the contextual knowledge at the client meetings, conducted systematic review in the extraction of updated feature set for prediction models, and the use of the results of exploratory data analysis for effectively communicating with the brain storming sessions with the clients were explained in chapter 4, while providing a guideline for similar approaches. The initial step of requirement gathering and data collection methods in a healthcare setting, data pre-processing techniques, exploratory data analysis, feature selection, model selection and solution implementation are described throughout Chapter 5 to provide a solid design process for a CDSS. The decision-making for the system design, mapping the client requirements with feasible deliverables, and conceptualising the design framework with identified social and knowledge contexts in a healthcare setting are results of the extended knowledge outcomes through the study. A significantly identified research gap of the study is applying the knowledge of DSRM and utilising those techniques on the process of solving the real-world issue. The process of developing the solution and adopting techniques in each step of the process confirms the applicability of DSRM in developing a CDSS for predicting the complications of diabetes. Further, this approach to resolving the real-world issue demonstrates the applicability and validity of DSRM in the healthcare sector, which fills another recognised knowledge gap.

The above-mentioned outcomes of chapters 2,3,4 and 5 fill the aforementioned knowledge gaps while providing vital statistical inferences from the data analysis phase. The exploratory data analysis (EDA) phase has been recognised as a milestone of the adopted empirical cycle to analyse real-world problems, to build strong situational awareness, and define the scope of the system. The EDA phase of the study was conducted in two stages, due to the availability of datasets. The

first phase was conducted with the initially collected dataset, to gain a good grasp on the existing cohort. The sociodemographic factors in the diabetes cohort were analysed in the initial phase. The second phase of EDA was conducted over the entire dataset which was collected at the later stages of the client meetings. The statistical analysis of the entire dataset is presented in the later section. The outcomes of the two stages of EDA are vital in refining the functional and non-functional requirements and the scope of the CDSS. Due to the importance of justifying the scope of the CDSS, the results of the two phases of explanatory data analysis are presented as follows. For example, the initial exploration of socio-demographic factors was hugely benefit on effectively communicating with the clients to clarify the availability of datasets, their characteristics, and handling the missing or inconsistent values etc. The results of the EDA of the entire dataset with empirical standards, express the nature of the cohort, while revealing the statistical overview for illustrate on the final web portal.

7.1.1.1 Exploring the Prevalence of Sociodemographic Cohort of Diabetes Mellitus

Examining the sociodemographic factors within a cohort stands as a crucial phase, unveiling significant dimensions of societal health. The healthcare sector strategically employs the outcomes of exploratory analyses of sociodemographic attributes for diverse purposes, including formulating healthcare policies, allocating resources effectively, prescribing essential medications, and more. Substantial evidence substantiates that comprehending the prevalence of sociodemographic factors such as age, ethnicity, and gender yields pivotal insights. Consequently, this section of the study endeavours to elucidate this knowledge through an analysis of the sociodemographic particulars within a diabetes cohort in New Zealand. The examination of this cohort, specifically those with diabetes complications, sheds light on the frequency of complications among the diabetes patient population. The objective of this section is to present an initial exploration of the dataset profile, visualizing patterns in sociodemographic details from the samples and their correlation with diabetes complications (Erاندathi et al., 2022). Scrutinizing the sociodemographic features of a nation's cohort constitutes a crucial stride

in comprehending behaviours and prevalent patterns, directly contributing to the enhancement of healthcare sector management.

In the Waikato region, the prevalence of diabetes stands at 5.79%. Table 7.1 provides a tabulated overview of the characteristics of all patients in the cohort. The selected cohort exhibits a slightly higher representation of males (52%) compared to females (48%). Furthermore, the cohort encompasses 22 ethnic groups, classified into five categories: European, Māori, Pacific, Asian, and others for result simplification. Predominantly, the patients in this cohort belong to the European ethnic group, followed by the Māori population.

Factors	Categories	N=2656	Percentage	Patient/ Population percentage
Gender	Male	1380	51.96	0.61
	Female	1275	48.00	0.55
Ethnicity	European (NZ European/Pākehā, Other European, European No Further Definition)	1649	62.09	0.48
	Māori	744	28.01	0.68
	Pacific (Cook Island Māori, Samoan, Fijian, Tongan, Other Pacific Peoples, Niuean, Pac People No Further Definition)	92	3.46	0.44
	Asian (Indian, Other Asian, Chinese, South-East Asian, Asian No Further Definition)	131	4.93	0.30
	Other (Not Stated, African, Other Ethnicity, Middle Eastern, Latin American/Hispanic, Response Unidentifiable)	40	1.51	0.38

Table 7-1 : Characteristics of diabetes patients in the Waikato region (N=2656)

The age distribution for each patient group is delineated in Table 7.2. The average age at which diabetes is diagnosed consistently tends to be lower for females. However, the age variability among females was more pronounced across all ethnic groups and the non-Māori subset. In comparison, both the mean age at diagnosis and the age range were lower among Māori individuals than their non-Māori counterparts. The Pacific ethnic group exhibited the earliest mean age at diagnosis, while Europeans demonstrated the highest mean age at diagnosis. Notably, Māori individuals had the lowest age at diagnosis recorded at 13, whereas Europeans recorded the highest at 103.

	Categories	Mean (Standard deviation)	Minimum value(0th)	First quartile (25 th)	Median (50 th)	Third quartile (75 th)	Maximum value (100 th)
All patient		67.4(14.2)	13.0	58.0	69.0	78.0	103.0
Gender	Male	69.3(13.1)	23.0	61.0	71.0	79.0	97.0
	Female	66.6(15.1)	13.0	57.0	68.0	78.0	103.0
Māori	Male	61.5(13.1)	23.0	53.0	62.0	70.0	95.0
	Female	60.3(14.5)	13.0	52.0	62.0	70.0	91.0
Non-Māori	Male	71.8(12.1)	33.0	64.0	73.0	80.0	97.0
	Female	69.7(14.5)	20.0	61.0	72.0	80.0	103.0
Ethnicity	European	71.7(12.5)	20.0	64.0	73.0	80.0	103.0
	Māori	60.2(13.9)	13.0	52.0	61.0	69.0	93.0
	Pacific	60.0(15.0)	26.0	50.0	59.0	71.0	92.0
	Asian	61.0(15.0)	29.0	49.0	63.0	72.0	94.0
	Other	66.7(14.2)	24.0	60.3	68.0	75.0	92.0

Table 7-2 : Age distribution of diabetes cohort in identified sociodemographic groups.

The distribution of microvascular and macrovascular complications within the T2DM cohort is visually depicted in Figure 7.2, segregated by ethnicities. The left section of the diagram delineates the ethnic composition, while the right side illustrates the prevalence of common complications within the cohort. The Sankey diagram highlights the most prevalent complications, accounting for at least 99 patients (3.73%) in the cohort. Hypertension emerged as the most frequent complication, succeeded by Type 2 Diabetes Mellitus (T2DM) with insulin resistance features, kidney complications, ophthalmic issues, and atherosclerotic heart disease. Approximately half of the cohort presented with hypertension, with 23.5%, 18.9%, and 13.1%

experiencing kidney complications, ophthalmic issues, and neuropathy, respectively. These findings align with the global distribution of complications as reported by the WHO (2016).

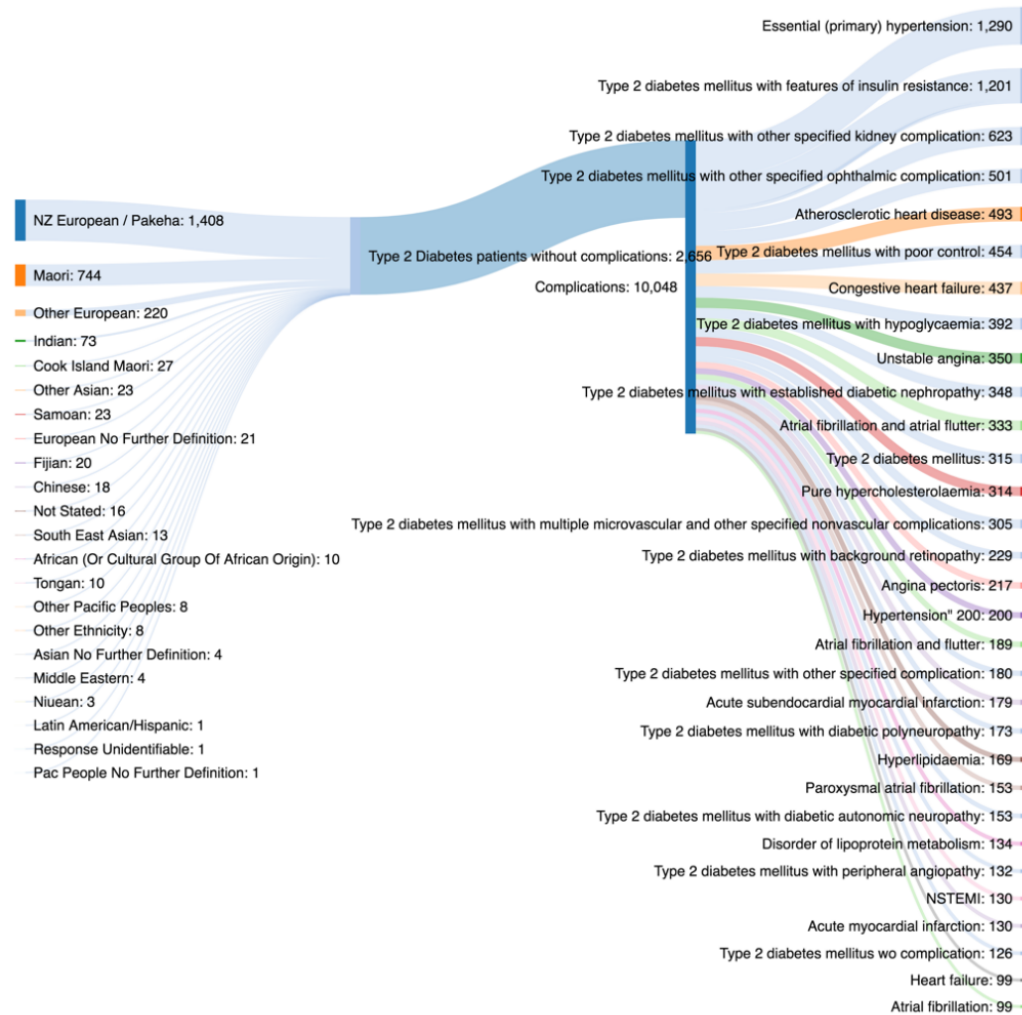


Figure 7-2 : Distribution of complications of diabetes mellitus and the ethnicities of the cohort.

Examining the demographic factors within a cohort is essential for gaining profound insights into a nation's characteristics. Given the substantial impact of Diabetes Mellitus (DM) on New Zealand's healthcare sector, an analysis of the demographic landscape unveils diverse perspectives of the nation. The findings indicate that the Māori ethnic group constitutes a higher percentage of the patient population (0.68%) in the Waikato region compared to other ethnic groups. While there are more males (0.61%) diagnosed with diabetes than females (0.55%),

females exhibit a lower average age of diagnosis than males. Furthermore, the diagnosis age for diabetes among Māori individuals (60.2) is lower than non-Māori individuals (70.3). The earliest age of diagnosis is observed among Māori (60.2) and Pacific (60.0) populations, with the highest average age recorded among Europeans (71.7). These research findings align with the overall statistics of New Zealand (NZ-MoH, 2022). Hypertension emerges as the most prevalent complication among diabetes patients, affecting 48.6% of the cohort (n=1290).

Moreover, this analysis provides insights into the epidemiology of a patient group in the Waikato region, leading to exploratory analysis and an in-depth examination of DM complications. Notably, limited research exists on the distribution of diabetes complications within a New Zealand cohort. Previous studies on diabetes complications have explored various thematic areas, including epidemiology and pathogenesis, micro, macro, and miscellaneous vascular complications, and treatment options (Papatheodorou et al., 2018). While research commonly focuses on epidemiology, the expansion of sociodemographic factors, and their association with specific clinical attributes like HbA1c, there is a dearth of scholarly works specifically addressing the distribution of DM complications. The vulnerability of the Māori population in the Waikato region to diabetes is highlighted, and females are diagnosed at an earlier age than males. These findings carry implications for the healthcare sector, offering valuable guidance for policymaking, resource allocation, and the identification of diabetes patient clusters. Stakeholders can focus on highly vulnerable patient classes, directing awareness programs and research studies towards the identified risk groups. This study presents a cross-sectional diabetes cohort in the Waikato region, characterized by a high density of diabetes cases. The outcomes of this stage have been employed to engage with the client, clarifying requirements and deliverables.

7.1.1.2 Exploratory Data Analysis

The results of the second stage of EDA are reported under the three most vital data resources: diagnosis table, test result table, and t2dm data frame. The three main data resources considered here are the details of the patient's diagnosis history, details of laboratory data of patients, and a data frame created as the foundational data frame for the study, which has the

dataset of the patients with E119. The results are further divided into four major EDA categories: univariate non-graphical, univariate graphical, multivariate non-graphical, and multivariate graphical. A summary of the results is included in the following section.

EDA results of Diagnosis Table

- **Univariate Non-Graphical EDA**

The diagnosis table consisted of the details of 142,067 males, 111,085 females, and 22 unknown gender values which followed the percentage values as 56.5%, 43.8% and 0.008%, respectively. According to the national diabetes registry (Health New Zealand, 2023), male percentage was higher than the female in 2021 which aligned with the research findings. The table comprised the details of 26,182 unique patients. The ethnicities of the diagnosis table comprised 24 ethnicities, where the majority of the patients belonged to NZ European/ Pākehā, Māori, other European, Indians and so on. The ethnicity division of the VDR has four classes: Māori, Pacific, Indian, European or other. Due to the existing taxonomical differences, the comparison of them against the results of the research does not take place. The mean age of diagnosis of E119 of the dataset was 66.63 with a standard deviation of 15.18. Moreover, the categorisation of patients as Māori or non-Māori revealed that 37% and 63% of patients were there for each class. This categorisation could not be seen in the VDR; therefore, the comparison could not be placed. However, according to Romana et al. (2022) Māori people are 1.8 times more likely to have diabetes than non-Māori people. The statistics of this study are about the dataset that was received from the WDHB. Sometimes the people who report to the WDHB, the availability of the data, and the regional specifications may affect the statistics revealed from the study.

- **Univariate Graphical EDA**

Patients' age on diagnosis of E119 was visualised as a histogram to see the age distribution. The following figure represented the age of the diagnosis of diabetes without any complications. The ages were scattered between 50 to 80 years, which showed a significant increase after 40 years. The average age for diagnosing T2DM was around 46 years (Koopman et al., 2005) among the cohort in the USA; according to the NIDDK (NIDDK, 2023), the most common age for diagnosis of diabetes was 45 years or older. The study's data set showed similar results as the age

distribution started to show significantly high values after the age of 40. The following histogram shows the age distribution among the patients in the diagnosis table.

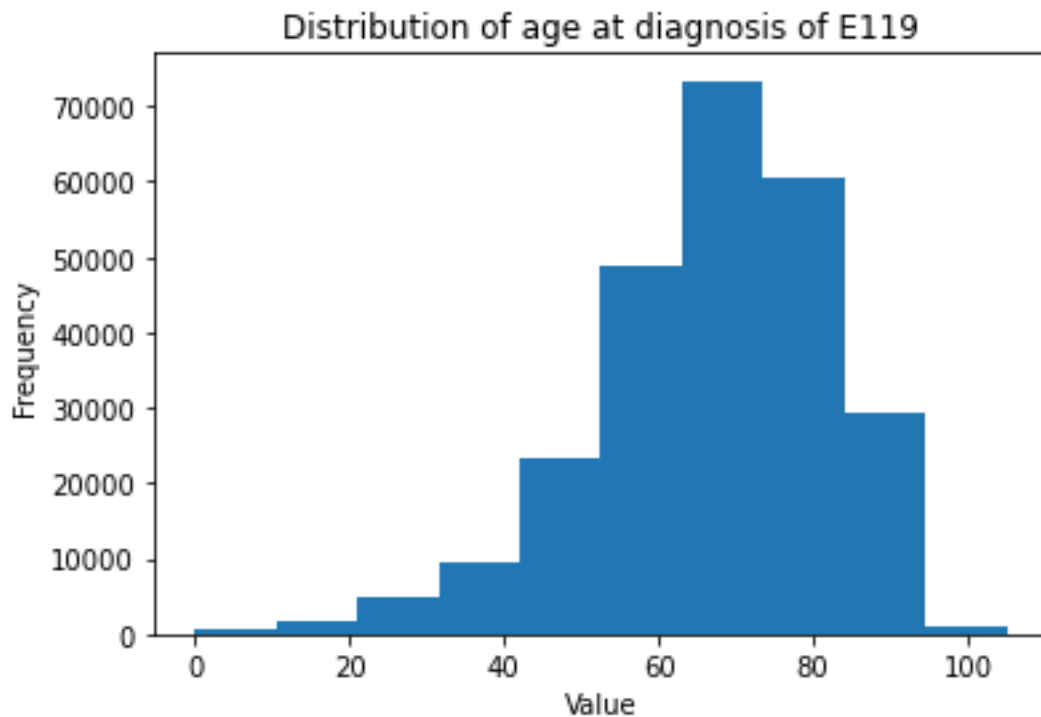


Figure 7-3 : E119 diagnosis age distribution of patients in the diagnosis table.

- **Multivariate Non-Graphical EDA**

The multivariate analysis revealed that in this dataset, there were more diabetes in the male population than females, regardless of their Māori or non-Māori ethnicity specification. The mean age of diagnosing E119 in females was 65.87, whereas that of males was 67.22. The average age value at diagnosing E119 of Māori was 60.44, whereas that of non-Māori was 70.22. Although a conflict in the percentage of Māori/non-Māori was detected due to the limitations of data accessibility, this result showed that the health indices among Māori and non-Māori aligned with the prominent general statistics.

- **Multivariate Graphical EDA**

A graphical representation of “Patient_Gender”, “Patient_Age_OnDiagnosis”, and “Patient_Maori_NonMaori” attributes are included in Figure 7.4 to visualise the relationships more clearly. Other than that, a few resulting box plots and violin plots are included to visualise

the relationships among the attributes of the diagnosis table. The following figure illustrates the age dispersion at diagnosis of E119, categorised into Māori/non-Māori by gender.

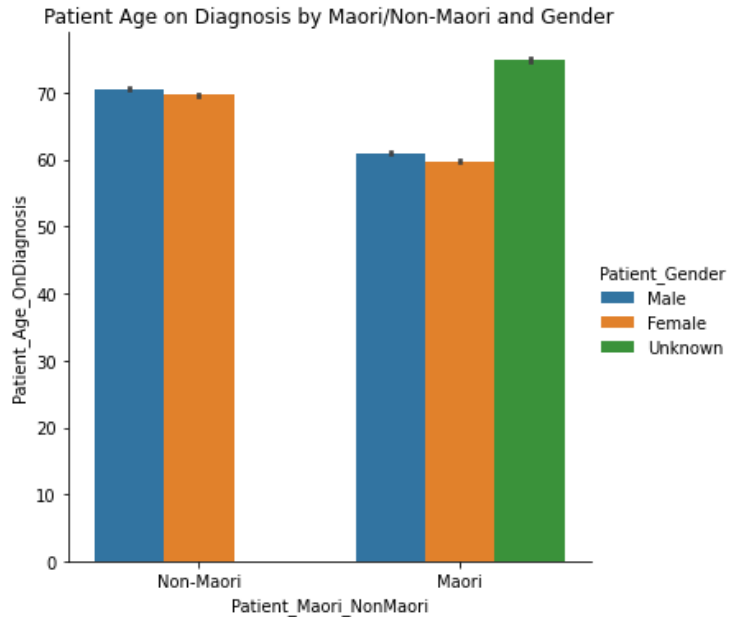


Figure 7-4 : Dispersion of age at diagnosis of E119 by Māori/non-Māori and gender.

The age of diagnosis of E119 for non-Māori males and females was around 71 and 69, respectively, whereas that of the male and female Māori population was around 60. The graph illustrates that the Māori population showed younger age than the non-Māori population at diagnosis of diabetes regardless of gender.

The following two box plots represent the age of diagnosis of E119 categorised by their gender and Māori/non-Māori characteristics. The minimum age of diagnosis of E119 was higher for males than females. Additionally, the mean age and minimum age of diagnosing E119 was higher for non-Māori than for Māori. This expresses that the status of health of Māori was comparatively lower than non-Māori.

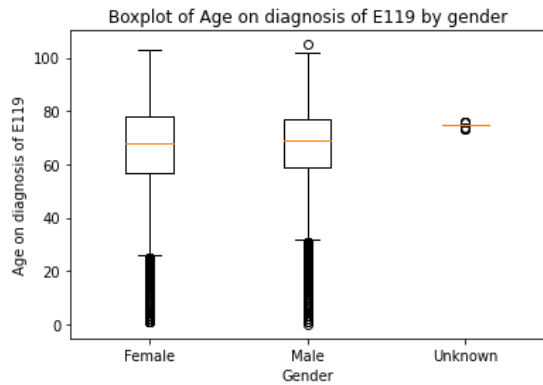


Figure 7-5 : Boxplot of age on diagnosis of E119 by gender.

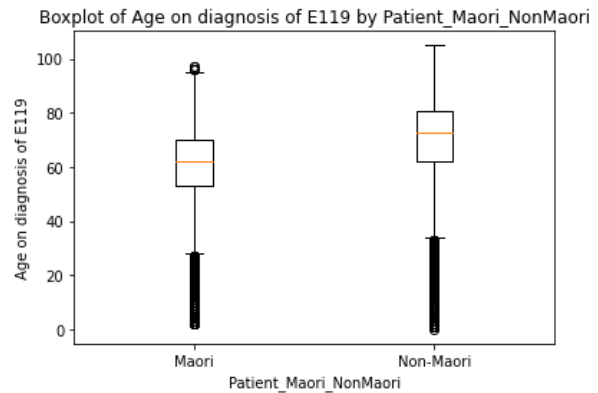


Figure 7-6 : Box plot of age on diagnosis of E119 by the characteristic of Māori.

The following violin plot shows the age distribution of the population among Māori and non-Māori ethnicities. A wider dispersion of population can be seen in Māori at a lower age of diagnosing, whereas comparatively a narrower and higher age of diagnosis of E119 can be seen among the population of non-Māori.

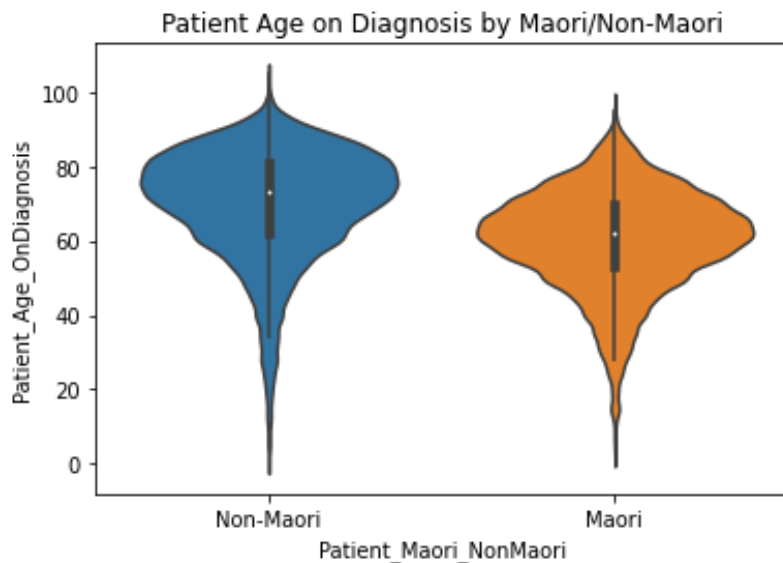


Figure 7-7 : Violin plot of patient's age on diagnosis categorised by Māori/non-Māori.

EDA Results of the Test Result Table

The interesting features of the test result table used in this phase of EDA are “Test_Id”, “Test_Description”, “ResultUnit”, and “Collection_DateTime”. The “Test_Id” column represented a unique identification number for the conducted laboratory test, “Test_Description” column described the laboratory test, “ResultUnit” column showed the measuring unit used for the particular test, and the column “Collection_DateTime” indicated the date and time of the laboratory test. Although the column with the result of laboratory test is vital, the value column was not considered an interesting feature, since the focal point of this phase was not about the individual values.

- **Univariate Non-Graphical EDA**

The dataset consisted of 1,048,575 tuples of data with a set of 23,883 unique patients. There were details of 121 unique laboratory tests, which had 40 different measurement units. The test details were collected from 2009 to 2019, where the least number of details could be seen in 2009 and 2019.

- **Univariate Graphical EDA**

The distribution of collected data over the considered decade was visualised using a histogram to check the availability of potential biases in the dataset. The following graph represents the resulting visualisation of the distribution of collected data over the years.

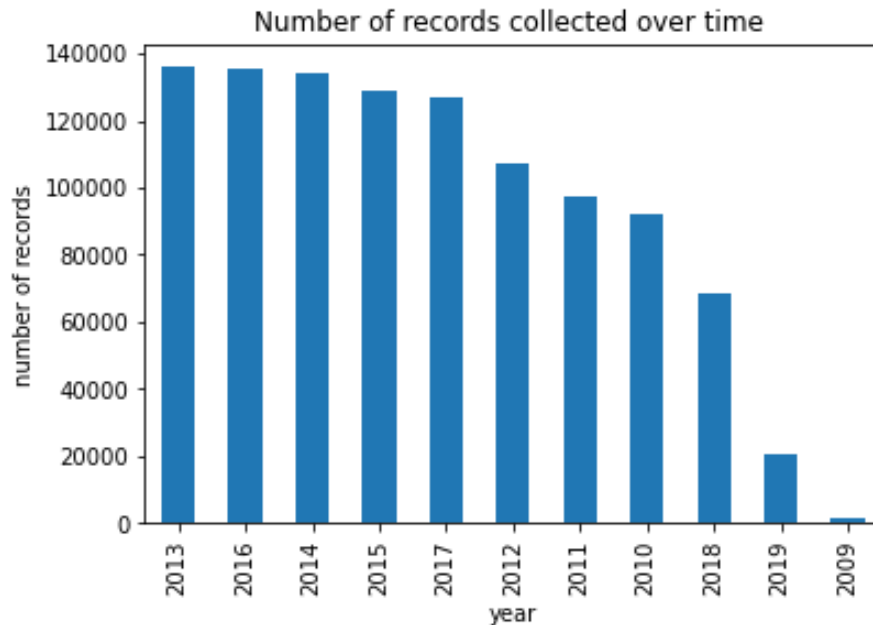


Figure 7-8 : Available records of test results in the test result table.

- **Multivariate Non-Graphical EDA**

The cross-tabulated results of “Test_Id” vs collected year of data provided information on the number of test results available on each year. The highest test result values of each year can be seen in Haemoglobin values whereas the eGFR shows the lowest values. The highest number of records for the Haemoglobin records are in the year 2013. The issues recognised through the EDA phase, such as few records on some selected features were communicated with clients to request more data values.

- **Multivariate Graphical EDA**

The number of tests conducted in each year has been visualised using a bar plot. The resulting bar plot is complicated due to the presence of a large number of test types. Therefore, the visualisation of the 20 most frequent test results each year has been illustrated in the following graph.

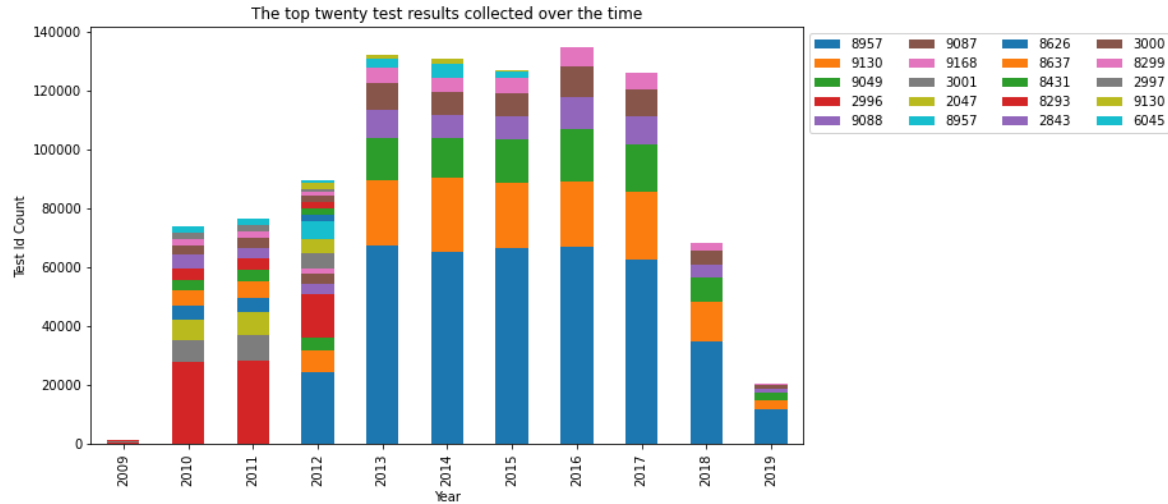


Figure 7-9 : The number of laboratory records collected through the years.

The figure shows the count of records in each laboratory test collected each year. The Haemoglobin values are the most significant laboratory test values in the dataset. Additionally, the least number of records of test results can be seen in 2009.

EDA Results of the Created Fundamental Data Frame of Patients with E119

This study's fundamental data frame was created by filtering only the patients with E119 (Diabetes mellitus without any complication). The data frames created to fulfil the various tasks of the study were created by slicing and filtering the above-mentioned fundamental data frame (t2dm). The created data frame was pickled for future purposes (t2dm.pkl). Due to the importance of this fundamental data frame, the following EDA phases describe the results of the EDA of the data frame t2dm. The diabetes cohort used in this study is the cohort of this data frame.

- **Univariate Non-Graphical EDA**

There are 7,206 unique patients in the t2dm data frame diagnosed as type 2 diabetes patients without any complications at the diagnosis time of t2dm. The complication with the highest number of patients was E1129 (Kidney Complications AKI), whereas the least complication was E1151 (Peripheral vascular disease - PVD), where the respective percentages were 20.2% and 1.9%. The selected most frequent 10 complications to predict the survival of the cohort were

E1129, E1172(Fatty Liver), E1122 (Diabetic Nephropathy), E1139 (Other Ophthalmic Complications), E1164 (Hypoglycaemia), E1165 (Poor Control – Hyperglycaemia), E1171(Microvascular and other specified nonvascular complications), E1131 (Background Retinopathy), E1142 (Diabetic Polyneuropathy), E1151 in the order of their number of patients. Moreover, this data frame comprised 60.9% of Europeans, 30.6% of Māori, 4.2% of Asians and so on. The results may mislead due to the characteristics of the resulting t2dm data frame. The data set had 55.8% of males and 44.15% of females. The data frame consisted of more non-Māori population (4996) than Māori population (2210). Moreover, the mean age at diagnosis of E119 was 67.5. The average days of diagnosis with any complication was 830 days, whereas the mean value was 648. The mean of the diagnosis year of E119 is 2014, while that of the complication diagnosis was 2017. The existing imbalances of the dataset may have an unavoidable effect on the study's results.

- **Univariate Graphical EDA**

The following two histograms represent the distribution of patients' age at diagnosis of E119 and the number of days between the non-complication to any complication at the left and right-hand sides, respectively. The histogram of age at diagnosis of E119 shows a normal distribution. In contrast, the one which shows the dispersion of time between the non-complication to any-complication shows a left-skewed distribution. The later distribution explains that the majority of the population stack on the left side and decreases over time.

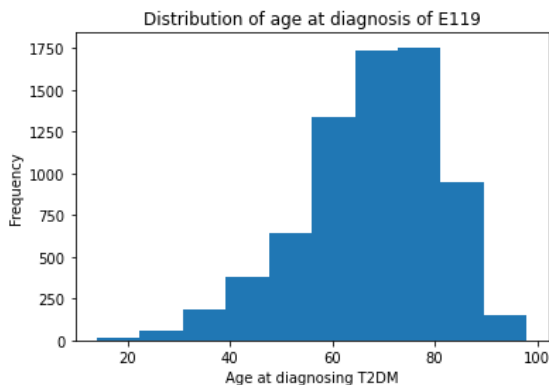


Figure 7-10 : Histogram of age at diagnosis of E119.

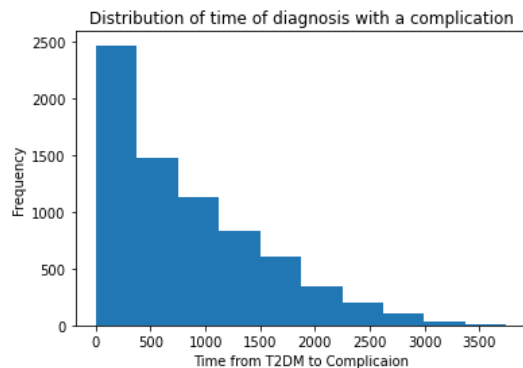


Figure 7-11 : Histogram of time between E119 to any-complication.

The histogram illustrating the diagnosis year of E119, and the diagnosis year of complications has been included as follows. The histogram of E119 values shows a reasonable symmetry with a mean in the year 2015. The histogram of the diagnosis year of complications is right-skewed, where the complication details are increased with the period.

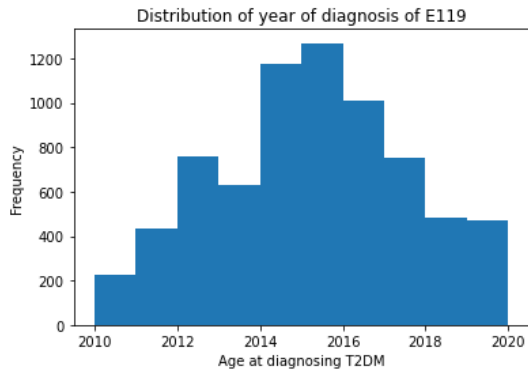


Figure 7-12 : Histogram of diagnosis year of E119.

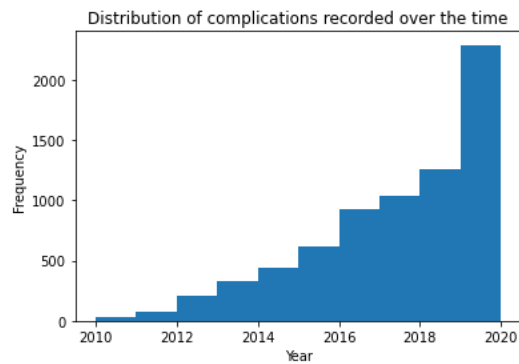


Figure 7-13 : Histogram of number of complications recorded over time.

The resulting box plots of the age at diagnosis of E119 and diagnosis year of complications are presented as follows. The box plot of age at diagnosis of E119, indicates the points up to age 35 are as outliers, which aligned with the knowledge of onset of type 2 diabetes. Additionally, the mean age of diagnosis of T2DM is indicated as 67.5 years. The complications occurred before 2012 are considered as outliers in the box plot, although they might be valid observations since the patients can be immediately diagnosed with complications after diagnosis of E119.

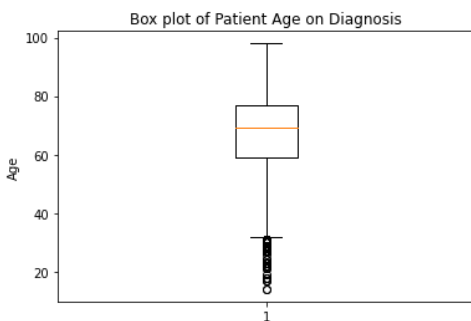


Figure 7-14 : Box plot of age at diagnosis of E119.

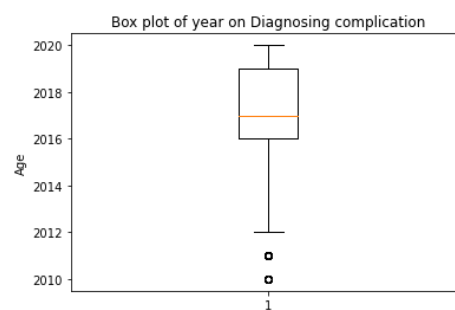


Figure 7-15 : Box plot of diagnosis year of complication.

- **Multivariate Non-Graphical EDA**

The cross-tabulated results of “Patient_Gender” vs “Patient_Maori_NonMaori” summarises that a higher number of the non-Māori male population could be seen in the dataset. In contrast, the Māori male population showed the least. Further, except for the “Māori” and “Middle Eastern/Latin American/African” ethnicities, all other showed a higher number of males than females. The population of each selected complication also followed the same patterns as above by having more males than females. The percentages of the incident rate of complications, along with their ethnicity, Māori/non-Māori, and gender categorisation, are represented in the following table.

Categorical characteristics		E1122	E1129	E1131	E1139	E1142	E1151	E1164	E1165	E1171	E1172
Ethnicities	European	12.63	21.07	5.05	12.11	4.57	2.21	10.97	8.58	9.04	13.77
	Māori	12.17	19.41	5.97	9.19	4.43	1.49	10.41	12.53	8.01	16.38
	Asian	11.67	16.33	8.33	11.67	6.33	2.00	12.00	11.67	8.67	11.33
	Pacific	9.60	14.65	11.62	9.60	4.04	1.52	9.09	13.64	8.59	17.68
	Other ethnicities	17.57	25.68	2.70	8.11	4.05	1.35	6.76	8.11	12.16	13.51
	Middle Eastern/Latin American/African	6.67	13.33	10.00	23.33	16.67	0.00	3.33	0.00	10.00	16.67
Māori/non-Māori	Māori	12.17	19.41	5.97	9.19	4.43	1.49	10.41	12.53	8.01	16.38
	Non-Māori	12.39	20.39	5.46	11.89	4.69	2.12	10.76	8.84	8.98	13.68
Gender	Female	11.41	19.08	5.94	12.01	4.37	1.67	11.25	10.06	8.23	16.00
	Male	13.17	21.10	5.42	10.44	4.85	2.16	10.29	9.99	9.12	13.47

Table 7-3 : Percentage of incident rates in CoDM

Although the results of incident rates were extracted through multivariate non-graphical methods of EDA, the resulting figures are presented as three graphs due to their clarity.

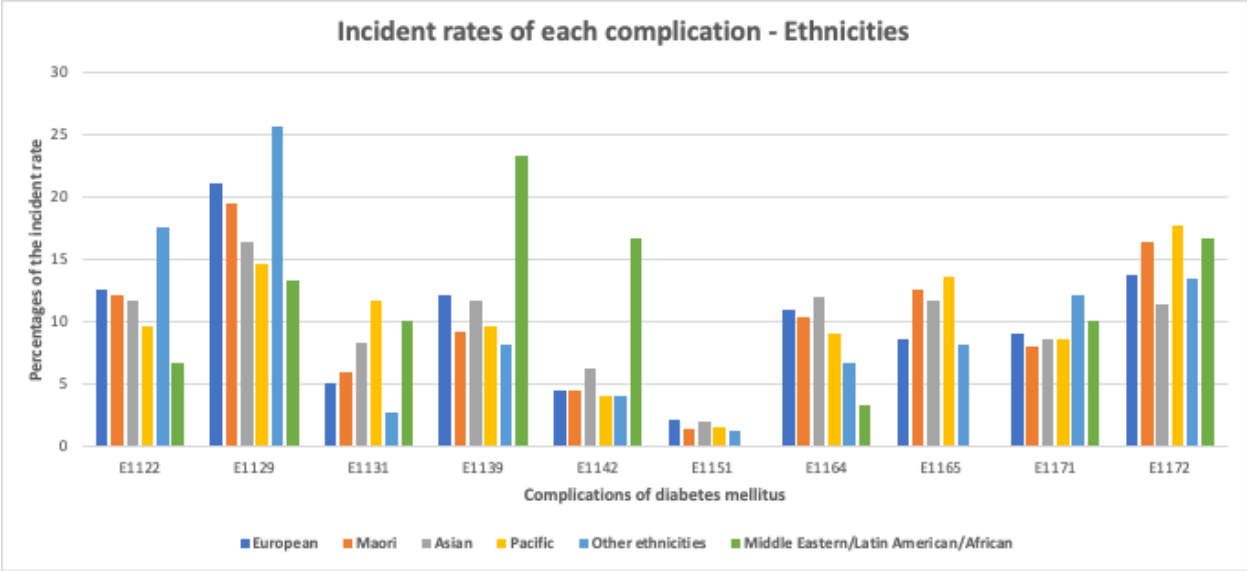


Figure 7-16 : Incident rates of each complication of diabetes categorised by ethnicity.

The highest percentage values of the graph belong to complication E1129, and the lowest percentage can be seen in E1151. The incident percentages of the “other ethnicities” group in E1122 and E1129 show the highest percentages compared with other ethnicities. The same characteristic can also be seen in the “Middle Eastern/Latin American/African” group in E1139 and E1142.

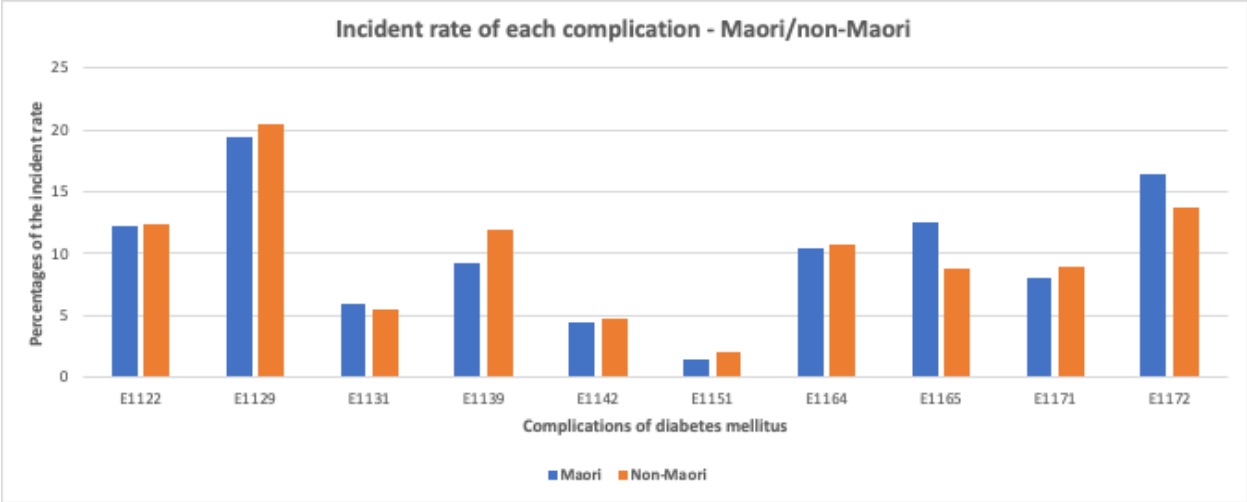


Figure 7-17 : Incident rates of each complication categorised according to the characteristic of Māori/non-Māori.

The incident percentages of the Māori population are higher in E1131, E1165, and E1172 compared to the non-Māori population, while the percentages are much similar in E1122, E1142,

and E1164. Distinctive values of the percentages can be seen in E1172, E1165 and E1139, where the non-Māori population shows higher values in E1172 and E1165 while Māori shows a high value in E1139.



Figure 7-18 : Incident rates of each complication of diabetes according to gender.

The incident rates of complications are different from the original values. As mentioned above, the male population shows higher values for all the complications than the females. However, the incident rate shows that the male population has higher incident percentages only in E1122, E1129, E1142, E1151 and E1171.

- **Multivariate Graphical EDA**

The following figure represents the number of patients in each complication according to ethnicity. The graph shows that the highest number of patients are from the ethnicities of European and Māori while the least is from the ethnicity categorisation known as "other ethnicity". Further, the change in the number of records for each complication can be seen in the figure, which shows the records in ascending order.

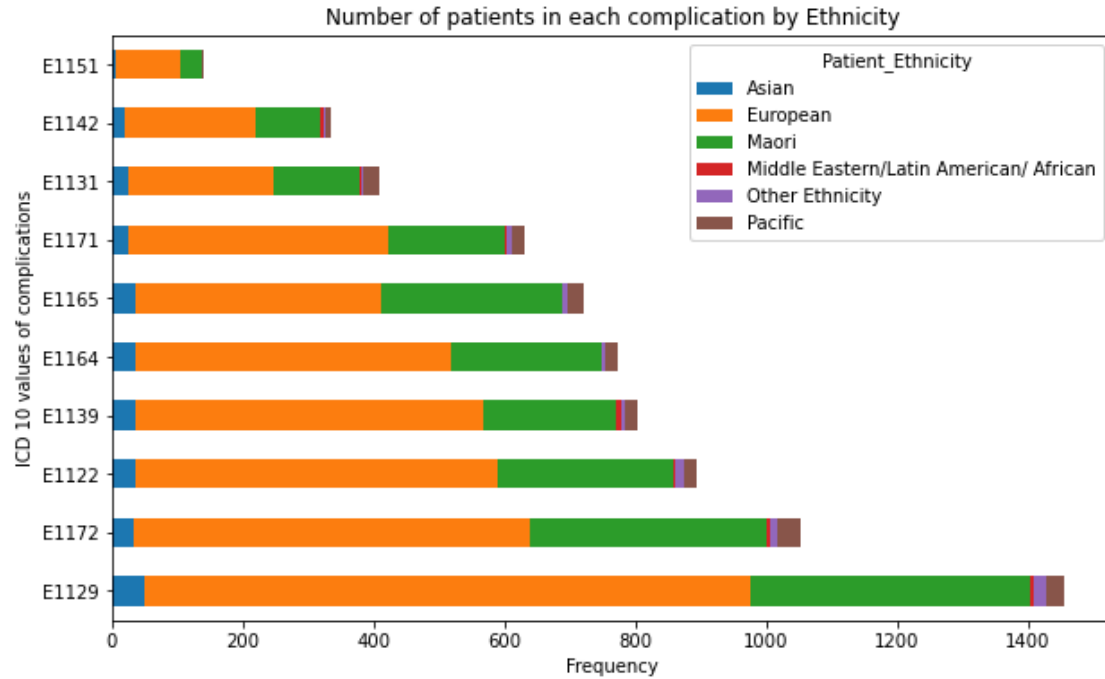


Figure 7-19 : Chart for representing the number of patients in CoDM by ethnicity

The following horizontal bar chart represents the number of patients in each complication with respect to their gender. It shows the gender distribution among the selected complications graphically.

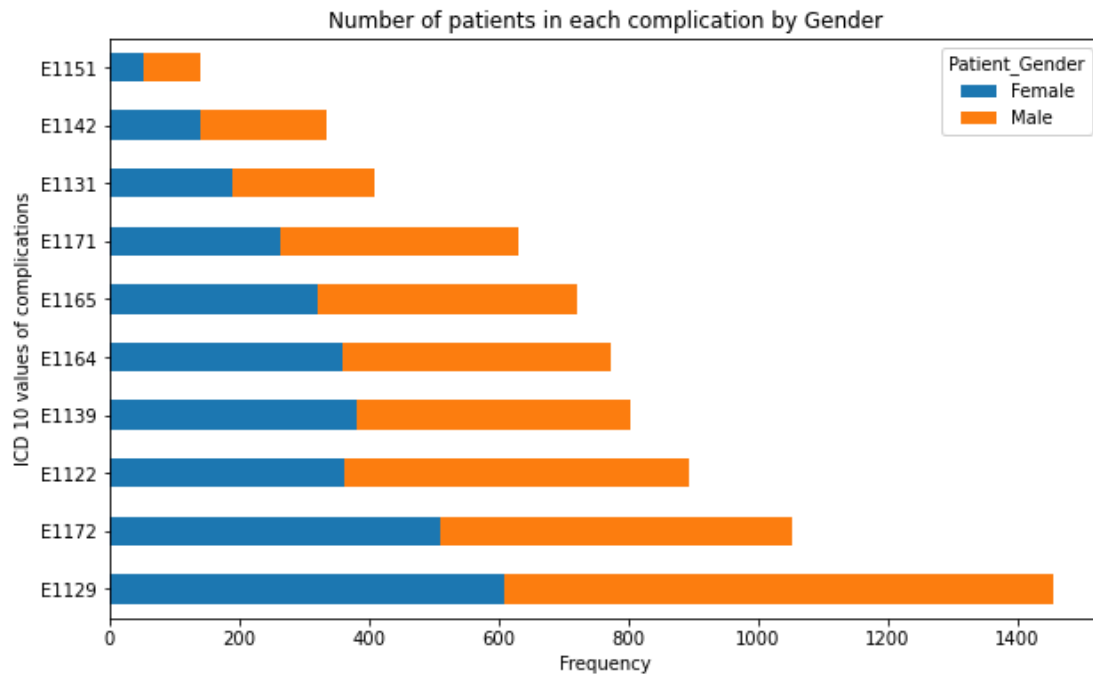


Figure 7-20 : Chart for representing the number of patients in CoDM by gender

The following two figures are representations of the distribution of age at diagnosis of E119 according to gender and Māori/Non-Māori characteristics. The graphs show that the age value is always higher in the non-Māori population regardless of gender value. It expresses important information on the health index of the Māori and non-Māori populations. Additionally, the age at diagnosis is less in the Māori females than males and vice versa for the non-Māori population. Although both of these graphs represent the same information about the cohort, their illustrations are important in visualising the differences between the cohort according to different categorisations.

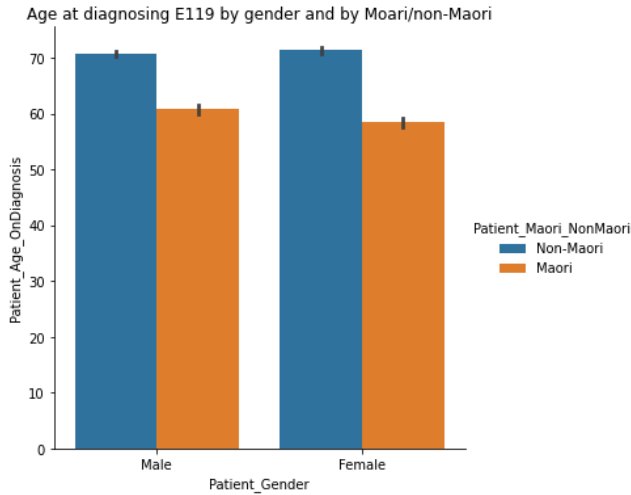


Figure 7-21 : Bar charts represent the age at diagnosis of E119 categorised by gender and Māori/non-Māori characteristic.

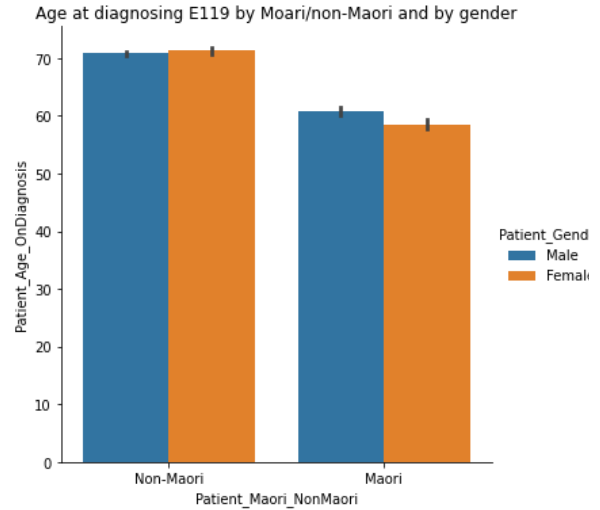


Figure 7-22 : Bar chart of the age at diagnosis of E119 categorised by Māori/non-Māori and gender.

The distribution of age at diagnosis of E119 among the ethnic groups has been visualised using a horizontal bar chart, where the Europeans showed the highest age of diagnosis of E119 and the lowest belonged to the “Middle Eastern/Latin American/African”. The Māori population showed a low age at diagnosis of E119 compared to the Europeans.

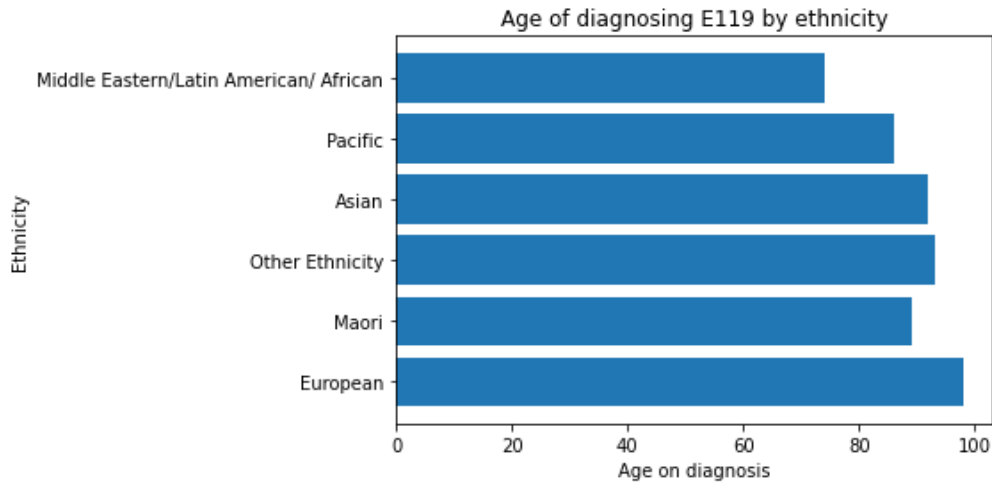


Figure 7-23 : Horizontal bar chart of age at diagnosis of E119 categorised by ethnicity.

The following two violin plots represent the age at diagnosis of E119 with the Māori/non-Māori and ethnicity values. The plot on the left side shows that the age distribution of Māori is more symmetric than non-Māori. The graph on the right-side expresses that the highest mean age

values belong to the ethnicity type “other ethnicity”. In contrast, the European, Asian, Māori, Pacific, and “middle eastern/Latin American/African” show a decreasing means of age.

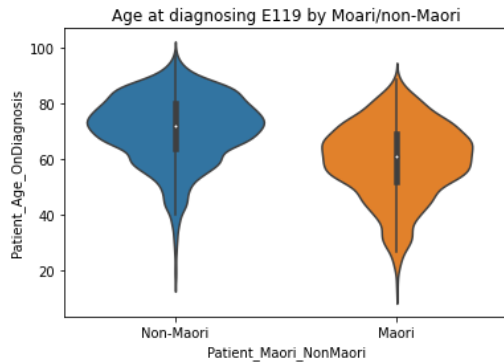


Figure 7-24 : Violin plot of age at diagnosis of E119 categorised by Māori/non-Māori.

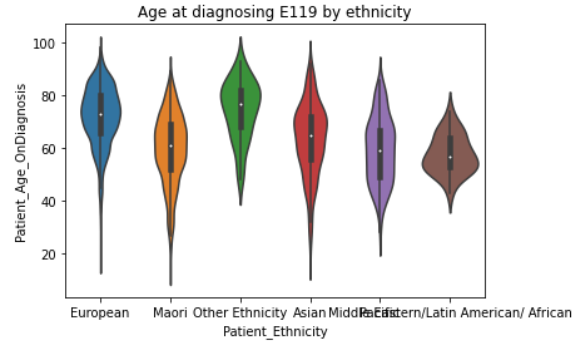


Figure 7-25 : Violin plot of age at diagnosis of E119 categorised by ethnicity.

Moreover, the data frames created to build the prediction models for the complications were analysed to reveal specific information regarding the complications. The categorical attributes of the dataset values—gender, age at diagnosis of E119, and Māori/non-Māori—were considered in the EDA of each data frame.

The outcomes of the above EDA results are used to refine the scope of the suggested CDSS. The initial step of design science research methodology (DSRM), which is the “designing of the solution” phase, consists of the literature review, situational awareness, requirement gathering, data collection, data pre-processing, EDA, and refining the scope of the CDSS. The “design implementation” phase starts with model selection and implementation. A justification of selected models included in Chapter 3 while the implementation of models is included in Chapter 5. Moreover, the presentation of the CDSS as an application is also discussed in chapters 3 and 5. The suitability of the selection of a CDSS as a web portal is described in Chapter 3, while the implementation process is discussed in Chapter 5. The resulting code of the website is another vital outcome of the study.

The first research question, lean towards the design perspective of the CDSS, which also divided into two sub questions where the former one focus on the design process and the later one

concentrate on the evaluation. The adopted design process on this study is started from initiating the phase of problem investigation. Similar to any research study, the state-of-art of the field has been captured via a thorough literature review. The design process of the study answers the first research question. The process of design from situation awareness, to design and implementation of the artefacts, which use in this study provide clear guidelines on adopting DSRM on solving a real-world issue in healthcare sector. This process refines the theory of DSRM while confirming their applicability in a healthcare setting. The situational awareness chapter describes the methods used in capturing the right real-world issue. The existing challenges, and obstacles in healthcare settings, and the expectations of stakeholders in a reputed healthcare organization in New Zealand are summaries as the outcome of chapter 4 to define the real-world issue with a systematically developed research framework. The used approach in the situational awareness phase provides a guideline for similar scenarios. The vitality of using structured client meetings, brainstorming sessions with the knowledge of the current status of the field, and potential deliverables, the use of results of systematic review on requesting datasets, and using EDA results for effectively communicating with the client to gather a strong dataset are emphasised as recommended milestones at situational awareness. The phase of design and implementation describe the stages and the techniques use in each step while providing a systematic procedure. The designing of artifacts uses a combined knowledge of literature review and situational awareness to design the rightmost artifact. The identified pragmatic issues in CDSS from the literature and situational awareness, leads to create an improved CDSS which avoid the practical issues as much as possible. The issue of predicting limited number of complications of the existing systems, are mitigated with analysing ten most common complications of the dataset. The drawback of short and fixed prediction periods in current CDSSs is overcome in the improved CDSS by using survival techniques on data analysis phase. The limitations on output of prediction result such as binary or ordinary output overtake with presenting the results as a risk percentage. Predicting the results by considering limited number of features, is handled with providing a standard feature set and using Cox models with two different feature sets. Although the issue of cohort specificity is emphasized through the literature review, the existing obstacles prevent this study from creating a globally accepted

system. The pragmatic issues of existing CDSS defines the expectations of the improved version of CDSS. The solution design and implementation phase continue with the empirical cycle of the system design. The data collection phase explains the use of document review method and its suitability due to the limitation of resources, obstacles in data collection in healthcare settings, challenges and necessary sensitivity of ethical considerations. Data pre-processing stage provide the nature of datasets in the discipline of healthcare, vital attentive points such as having different instruments, test results with non-standard measurement units, importance of considering outliers, filling the missing values with statistical measurements and cohort specific values. The exploratory data analysis phase reveals vital information on the hand-in dataset. The results included in chapter 7.1.1 represent how the initial EDA results use in effective communicating with stakeholders and defining the scope of the project. The results of second phase of EDA use to implement the web-portal. Moreover, the process of EDA followed through the study refines the process of EDA in a healthcare setting while providing the common factors and discretization of factors in a real-world scenario. The feature selection phase leads to a significant knowledge gap where a standard set of common risk factors are lack in prediction of diabetes mellitus and CoDM. The conducted systematic review extracted a standard feature set to predict the most common CoDM. The selected feature set is beneficial in designing solutions to CoDM while using that as a benchmark for future researches. The model selection phase mapped the data analysis techniques with the intended goal of the study. The use of predictive techniques, and the suitability of survival techniques to fulfil the purpose of the research, guide the data analysis process. Although the data analysis perspective is considered in the second research question, the process selection provides insights in the design process of the CDSS. Finally, the implementation of the web-portal, describe the research presentation. The existing solution presentation methods, the advantageous of using a web-portal over stand-alone system, selected libraries in implementing the web-portal their suitability, launching a web-portal in a set up to match with the criteria of the study are guiding the future researchers into a systematized path while refining the DSRM in a healthcare sector.

The findings described in this section explain the process of designing a CDSS, from their initial stage of finding the research gaps, awareness of the situation, conceptualising the framework with correct entities to create the solution, fundamental data analysis steps in designing a practical tool, such as data collection, pre-processing, feature selection, and model selection, and the final stage of designing a GUI for present the CDSS as a web-based system. The fundamentals of DSRM are adopted throughout the study design process, with descriptions of their appropriateness and confirming their usage. The theory of DSRM refines through this provided design process of the suggested CDSS. The above-mentioned systematic design process guide the researchers in healthcare sector to solve a real-world issue while adopting theories in DSRM. Further, due to the vitality of the evaluation phase in the DSRM, it has been separated as the second sub-question of the study.

7.1.2 RQ1.2: How Can a Designed CDSS be Evaluated with Existing Standards?

The system evaluation plays a vital role in the information research domain (Wieringa, 2016). The validity of the proposed solution is confirmed with a systematised evaluation process. Although system evaluation is considered a crucial step in information system research, adopting the existing evaluation methods are controversial. The evaluation method adopted in validating the CDSS is explained in Chapter 6. In summary, the FEDS framework's technical risk and efficacy method has been selected to evaluate the CDSS (Venable et al., 2016). The evaluation phase has been divided into three major components: design, algorithm, and implementation evaluation. The design and implementation evaluation are the initial and final evaluation of the CDSS, respectively, while algorithm evaluation comprises episodes of evaluating each model. The design evaluation phase is focused on evaluating the requirements, investigating the feasibility of requirements, mitigating the risk of mismatches with the client's requirements, and enriching the design with the potential functionalities to enhance the system's usability. The utilised design evaluation phase comprised four components: 1) Expected effects, 2) Expected value, 3) Trade-offs 4) Sensitivity (Wieringa, 2016). The given components shape the design to select the most suitable functionalities by considering their overall characteristics. Using existing standards in the design evaluation phase reveals information on configuring the requirements, selecting the

functionalities, mapping the stakeholder needs to the system requirements, and validating the design. The revealed knowledge of the design evaluation process is important as a practical guideline to evaluate the design at the commencement of the study. The above-mentioned process and techniques used in the design evaluation phase validated the current system while providing a systematised approach for the “design evaluation” phase which can be considered a guideline for similar approaches.

Further, the implementation evaluation is the system's summative evaluation phase. The appropriateness of the selection of techniques and their implementation are described in chapters 3 and 6. The adopted techniques of the “implementation evaluation” phase reveals the necessity of two standard ways to validate the system. The standard software evaluation method of ISO-25010 (International Organization for Standardization, 2011; ISO/IEC 25010, 2011) is primarily used to evaluate the implemented system with a standard set of criteria. Further, a user-involved evaluation method (Chiew & Salim, 2003; EL-firjani et al., 2017; Tullis & Stetson, 2006) has been performed to validate the CDSS to check the users' perspective. The following section describes the results of the implementation evaluation phase. The results will be presented as two sub sections, where the first section describes the results of the evaluation of created CDSS against the standard ISO 25010, and the second sub section present the results of the user feedback collected through the questionnaire.

7.1.2.1 Results of the Evaluation through ISO 25010

The implementation evaluation results of the CDSS conducted through the standard matrix ISO 25010 are included in this section. The evaluation process is justified and described in the chapter 3 and 6 respectively. The overall product quality of the system has been evaluated using the ISO 25010 standards (International Organization for Standardization, 2011; ISO/IEC 25010, 2011).. The influence of the requirements on external characteristic and sub-characteristics are evaluated in three phases, such as evaluating the degree of influence of a block of requirement on external characteristics, sub-characteristics, and the degree of influence of requirements in external characteristics. The results of these three phases are included as follows. The used blocks

of requirements are data accessibility (DA), personal information (PI), quantitative data (QD), user's actions (UA), CDSS's functionalities (DF). The standard external characteristics are functional suitability, performance efficiency, compatibility, usability, reliability, security, maintainability, and portability. The following graph visualises the influence of a block of requirements on external characteristics.

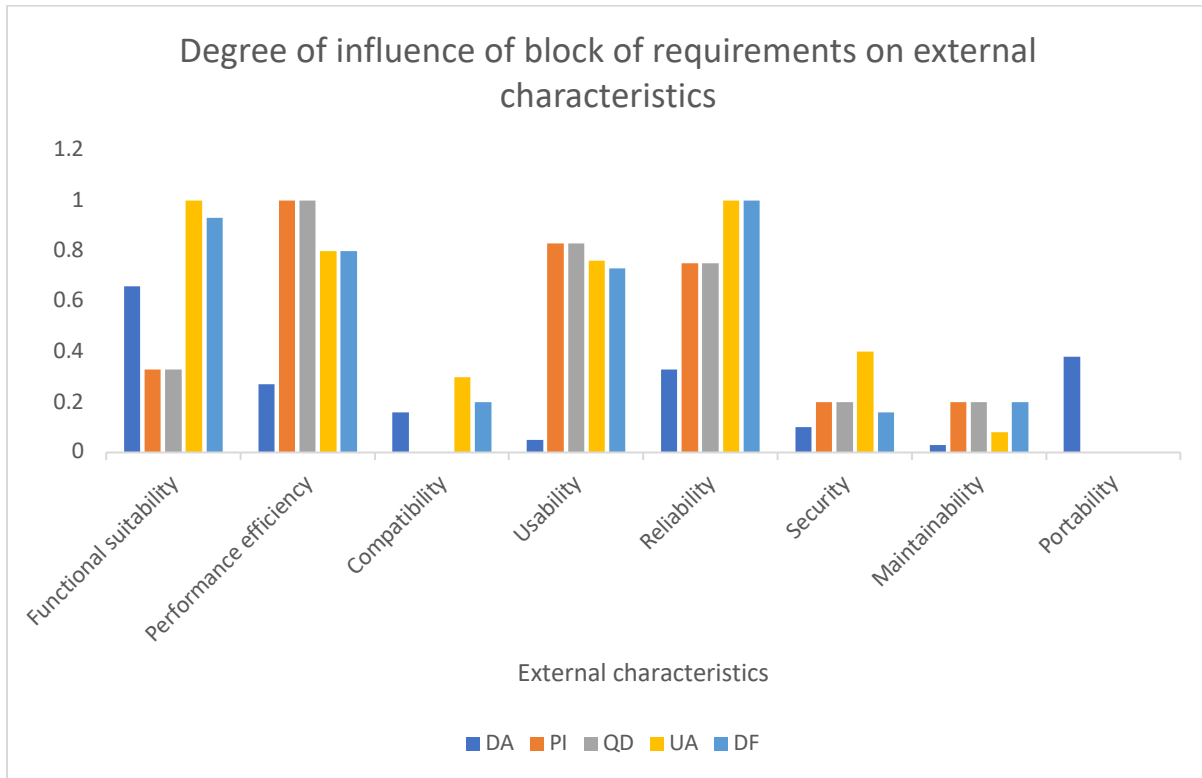


Figure 7-26 : Chart representing the degree of influence of block of requirements on external characteristics.

According to the results, the external characteristics, such as performance efficiency, usability, and reliability of the system, presented a significant strength of the system. Additionally, the functional suitability also showed a good representation. Compatibility and portability were the characteristics that showed the minimum contribution to the system. Portability and compatibility are two characteristics which are focused on having different versions of the system to be compatible with all platforms and devices. Since the created system was a web-based CDSS, these two characteristics were minimally affected by the quality of the system. Moreover, the graph shows the presence of at least one block of requirement in all external characteristics

which expresses the existence of all the standard external characteristics in the developed system.

The following two figures show the degree of influence of a block of requirements for the sub-characteristics and that of requirements to the sub-characteristics.

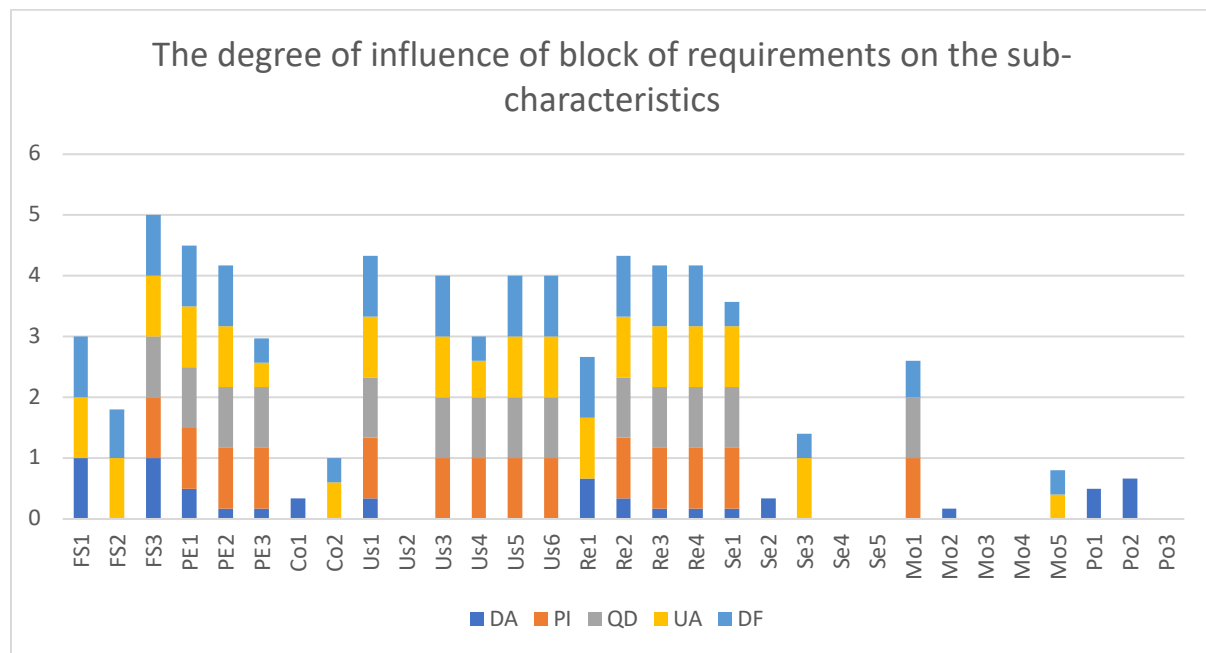


Figure 7-27 : The degree of influence of the blocks of requirements on the external sub-characteristics.

The above graph shows the presence of external sub-characteristics of the system. Functional appropriateness, time behaviour, resource utilisation, appropriate recognizability, fault tolerance, recoverability, and availability characteristics were the sub-characteristics that had the most decisive influence gained from the selected requirements of the system. Authenticity, accountability, analysability, modifiability, and replaceability were the qualities revealed to have minimum influence gained through the requirements. Most of these qualities were tolerable in the setting in which this research was developed. However, if the results of the above graph were statistically analysed, the results expressed that the existence of qualities in each requirement was higher than the absence of the qualities in the requirements. Since five blocks of requirements were considered, the threshold value of influence on each sub-characteristic was

set to three. All 31 sub-characteristics were considered here. Sixteen sub-characteristics were influenced by at least three blocks of requirements. The results illustrate that more than 50% of the sub-characteristics were influenced by 3/5 blocks of requirements.

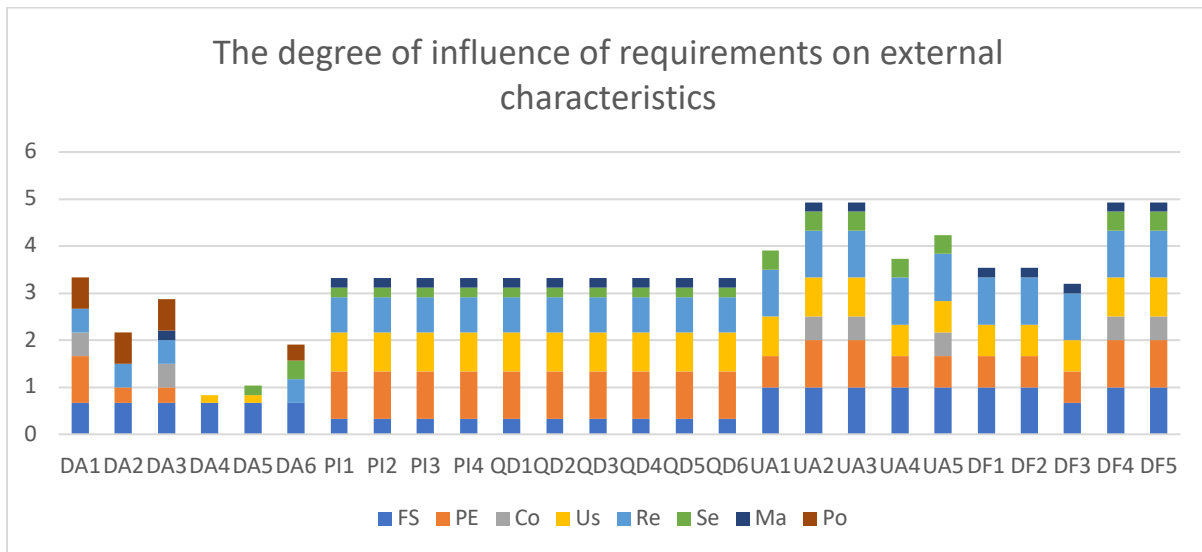


Figure 7-28 : The degree of influence of the requirements on the external characteristics.

The above graph shows the influence of each requirement on external characteristics. The results illustrate that reliability, performance efficiency, and usability can be seen in most requirements, while portability and maintainability were limited to specific requirements. However, this graph shows that most of the system's requirements covered at least six sub-characteristics out of the eight.

The implementation evaluation phase revealed that the developed CDSS possessed external qualities, such as functional suitability, performance efficiency, compatibility, usability, reliability, security, maintainability, and portability. The user's actions and CDSS functionality, respectively, showed 100% and 93% influence on the quality of functional suitability. Personal information and quantitative data represented 100% influence on performance efficiency, while user action and CDSS functionality showed 100% influence on reliability. The average influence of requirements

on each external characteristic expressed the system's strength. The highest strength of the system was reliability, whereas the least of that is portability. Performance efficiency and functional suitability were also significantly higher, which followed usability, security, compatibility, and manageability, respectively, showing the deduction of average influence in the system. The results revealed that the developed artefact of the research showed its product quality against the standard quality matrix ISO 25010.

7.1.2.2 The Results of the Evaluation of User Feedback

The second method used in implementation evaluation is justified and executed in chapter 3 and 6 respectively. The results of user feedback are included in this section to prove the usability of the system. The process of evaluation of the system is included in the chapter 6 while presenting their results in here. The user feedback on the generated CDSS was collected through a questionnaire distributed among the potential stakeholders. Seven volunteer stakeholders responded to the questionnaire: two general practitioners, three nursing lecturers, and two research officers. Most participants were "25 to 34", whereas two belonged to the age group "35 to 44". Four participants were females, two males, and one preferred not to reveal their gender. The fourth question asked about the easiness of learning to use the CDSS, which had the majority of responses to number four, where the four participants among seven evaluated the learning easiness as 4 and 5 (the highest level of easiness). Furthermore, 57.1% of responses for the understandability of the user interfaces were received for number 4 in the linear scaler. Since numbers 1 and 2 did not have any responses received, the systems' understandability was evaluated as having an acceptable level. The sixth questionnaire question checked the quickness of doing tasks using the CDSS, which evaluated to the highest possible number by 57.1%. Moreover, 85.7% of the participants agreed that the created system had clear icons, labels, and menu options. The ease of recalling the steps to perform tasks through the system received 57.1% of responses as "Agree" and 14.3% as "Strongly agree". However, one participant recorded it as "Neutral" and another responded as "Strongly disagree". The satisfaction with the usability of the system received 57.2% positive responses, whereas 14.3% received negative responses.

The 10th question checked the level of satisfying user needs and expectations, where 71.4% selected the " Agree " response. Additionally, the participants responded to the system's accessibility and usability of assistive technologies with 57.1% on the value of "Agree". There were no records of any accessibility challenges or limitations of the CDSS. The visually appealing quality of the system was evaluated as "Agree" by 28.6%, "Neutral" by 42.9%, and "Disagree" by 28.6%. The responses to the quality of visually appealing revealed the requirement to enhance the appearance of the system. However, most participants voted for an excellent user experience of the CDSS, which only encountered number 3 and above. The 16th question checked the layout design as a well-organised and intuitive structure which received 42.9% on number 4 and 28.6% on number 5. Furthermore, 85.7% of the users voted for using consistent design elements and system navigation as "Agree"; one participant voted as "Strongly Agree". The system's response time also received only positive responses, such as 57.1% on "Agree" and 42.9% on "Strongly agree". The system's reliability was evaluated as "Agree" by 57.1% and "Neutral" by 42.9%.

The visualisation of the results of the questionnaire are included as follows to provide a clear understanding of the questions and their responses.

Age Distribution of the Participants.

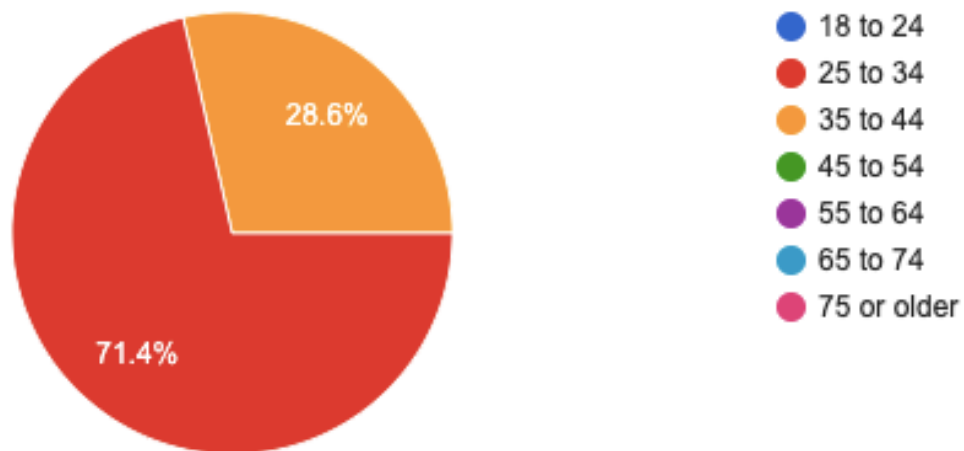


Figure 7-29 : Chart represents the age distribution of participants in the questionnaire.

Gender Distribution of the Participants.

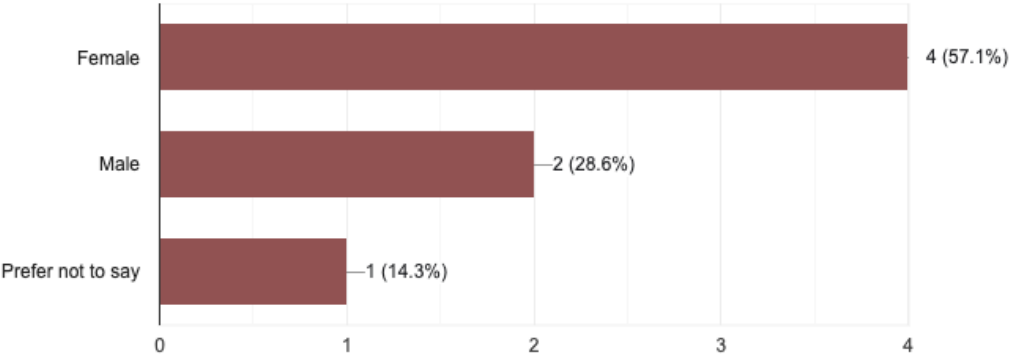


Figure 7-30 : Chart represents the gender distribution of the participants in the questionnaire.

Visualisation of the Responses to the Easiness of Learning to use the CDSS.

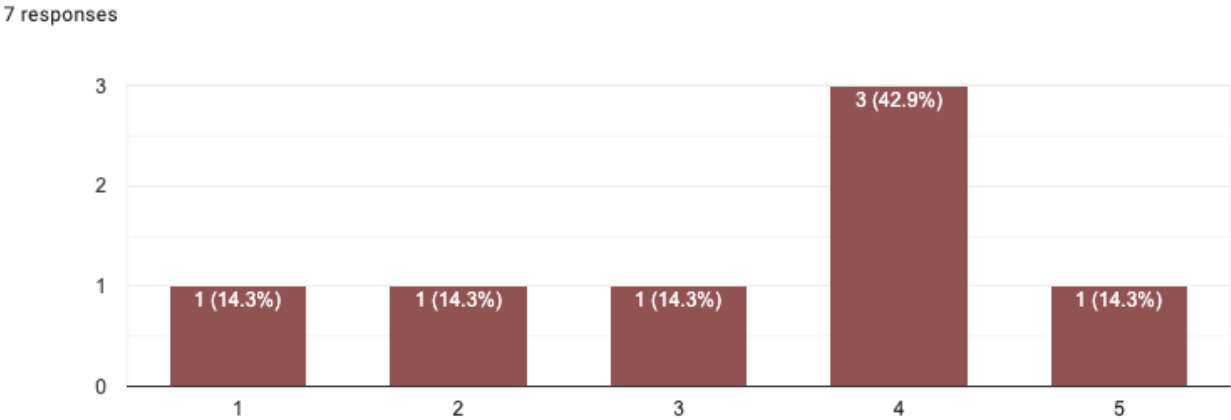


Figure 7-31 : Chart represents the responses of users for the easiness of learning to use the CDSS.

Responses for the user friendliness of the interfaces of CDSS.

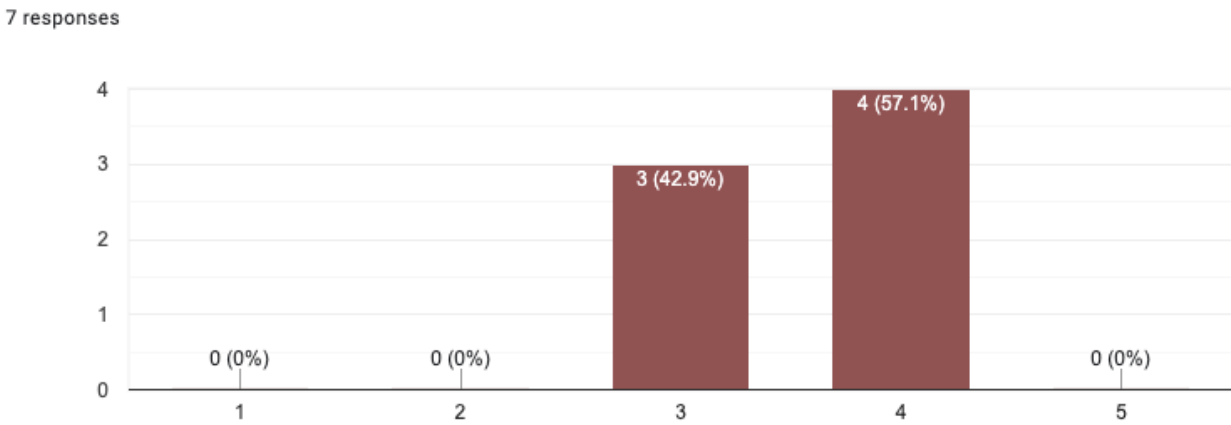


Figure 7-32 : Chart represents the user friendliness of the interfaces of CDSS.

Visualisation of the Responses of the Ability to Complete the Tasks Quickly Using the CDSS.

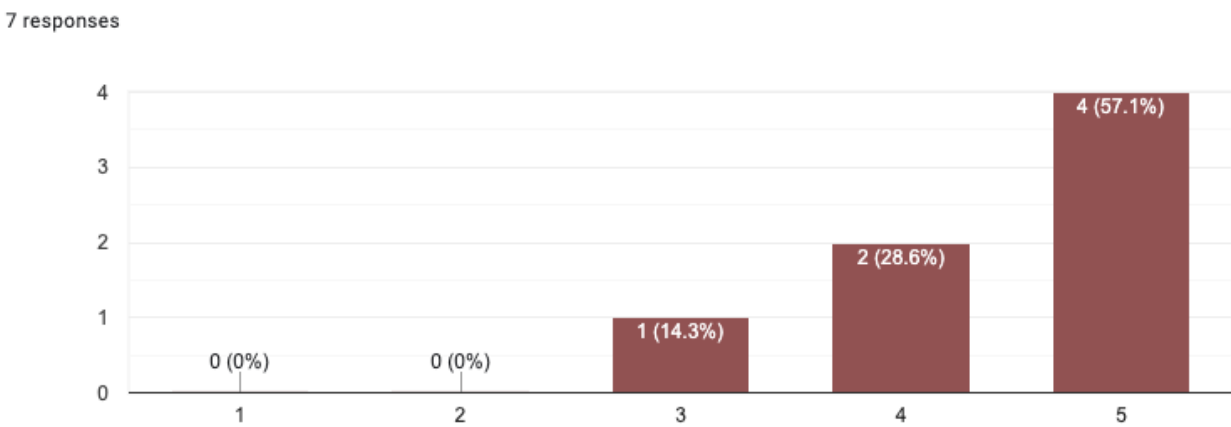


Figure 7-33 : Chart illustrates the responses for the ability to complete the tasks quickly using the CDSS.

Representation of the Responses of Having Clear Icons and Menu Options in the Interface of CDSS.

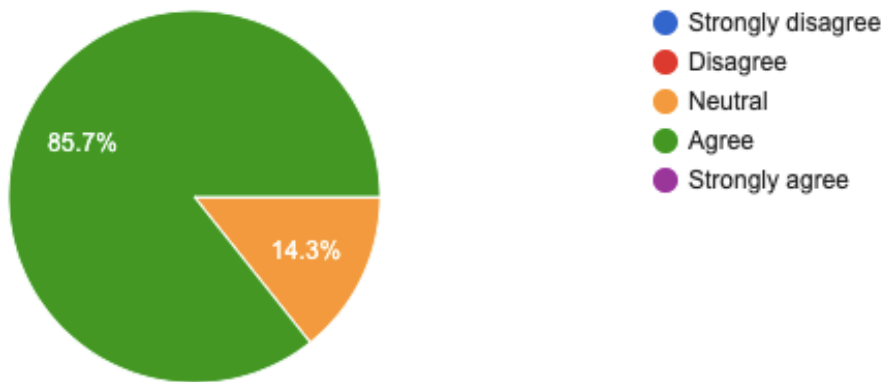


Figure 7-34 : Chart represents the received responses for the having clear icons and menu options in the interface of CDSS.

Visualisation of the Responses of Recalling the Steps to Perform the Tasks in CDSS.

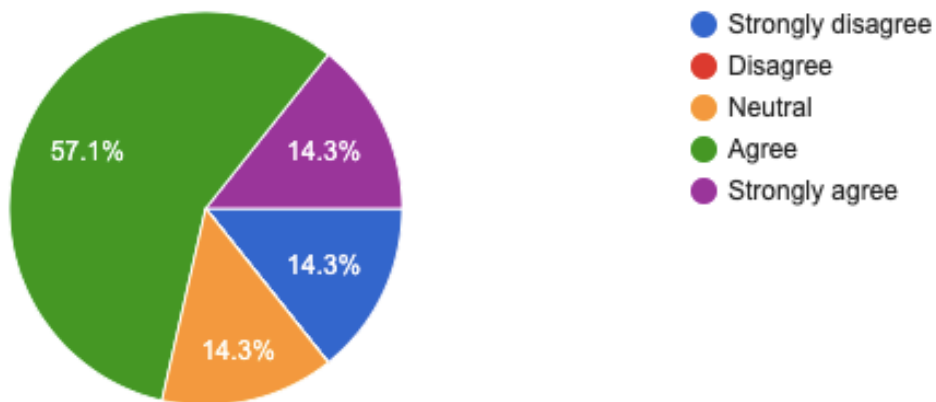


Figure 7-35 : Chart illustrates the responses for the recall the steps to perform the tasks in CDSS.

Illustration of the Responses of Satisfaction of the Usability of the CDSS.

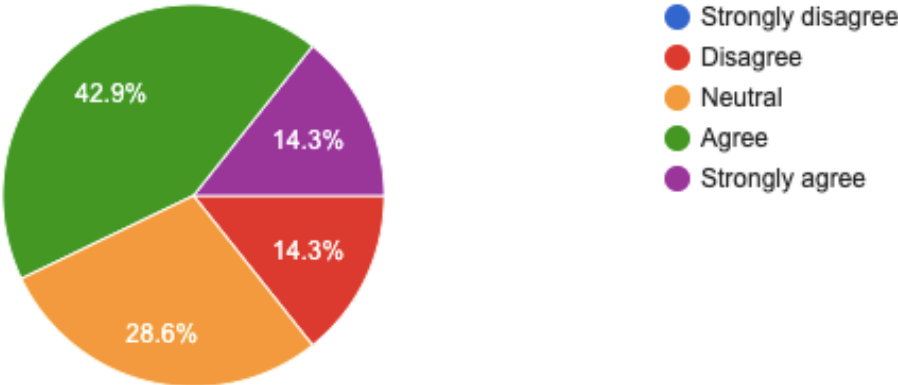


Figure 7-36 : Chart illustrates the satisfaction of the usability of the CDSS.

Visualisation of the Responses of the Ability to Meets the Users' Needs and Expectations.

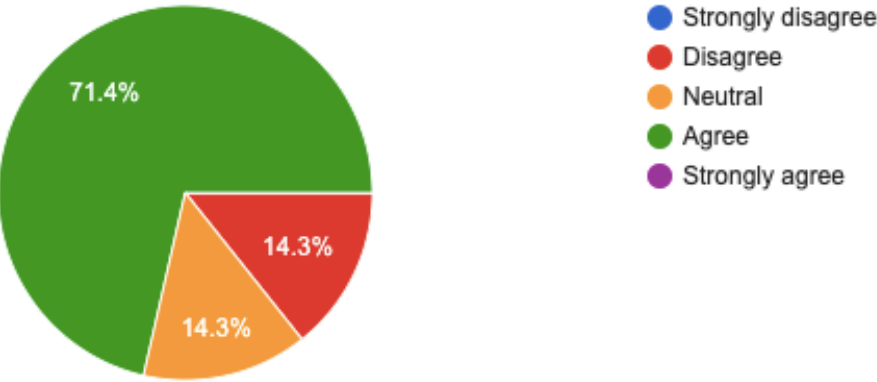


Figure 7-37 : Chart illustrates the responses whether the CDSS meets the users' needs and expectations.

Illustration of the Responses of the Accessibility and Usability of CDSS with Assistive Technologies.

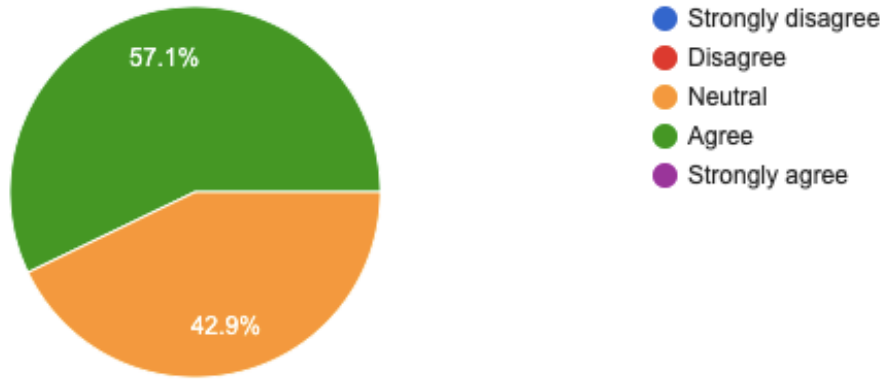


Figure 7-38 : Chart represents the responses of accessibility and usability of CDSS with assistive technologies.

Illustration of the Responses for Encountering Accessibility Challenges or Limitations of the CDSS.

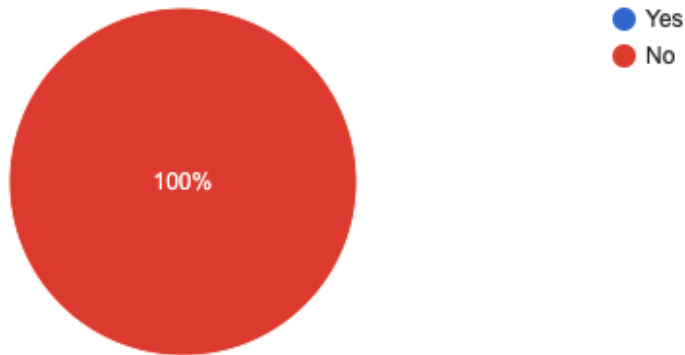


Figure 7-39 : Chart for the responses for encountering accessibility challenges or limitations of the CDSS.

Representation of the Responses for Visually Appealing Quality of the CDSS.

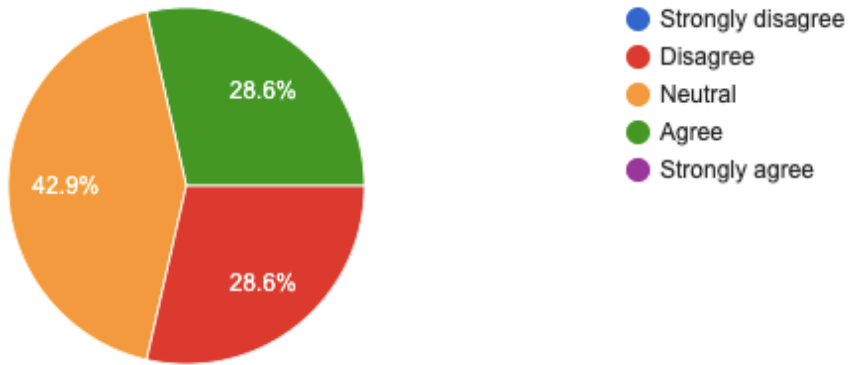


Figure 7-40 : Chart represents the responses for the visual appealing of the CDSS.

Illustration of the Responses of Having Pleasant User Experience in Using the CDSS.

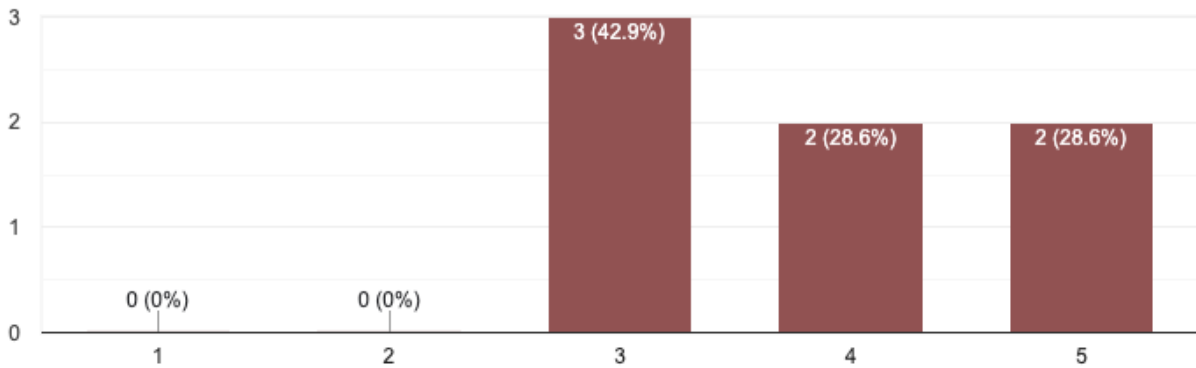


Figure 7-41 : Chart illustrates the responses for pleasant user experience of the CDSS.

Visualisation of the Responses of Having Well Organised and Intuitive Layouts in the CDSS.

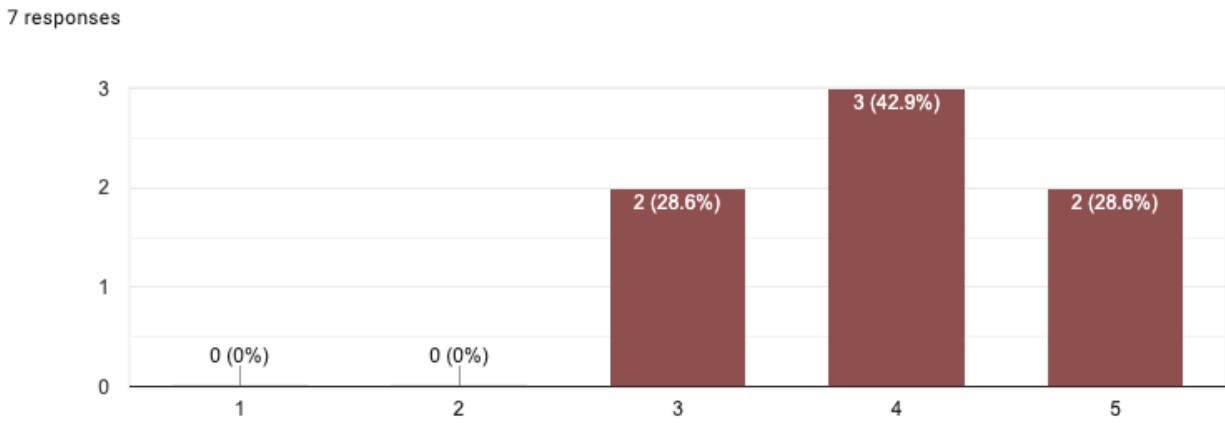


Figure 7-42 : Charts represents the responses of having well organised and intuitive layouts in the CDSS.

Represent of Having Consistent Design Elements and Navigation in the CDSS.

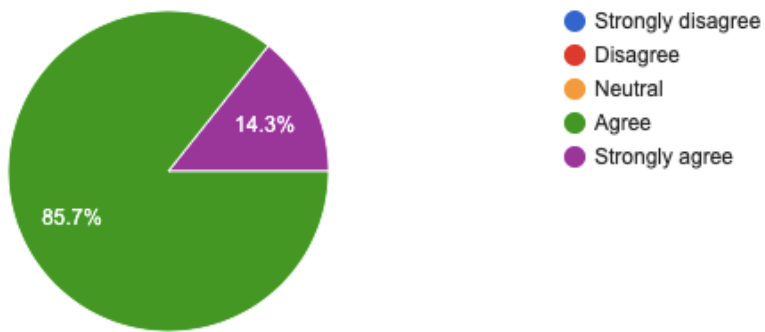


Figure 7-43 : Chart illustrates the responses of having consistent design elements and navigation in CDSS.

Illustration of the Responses of Having a Good Response time from the system.

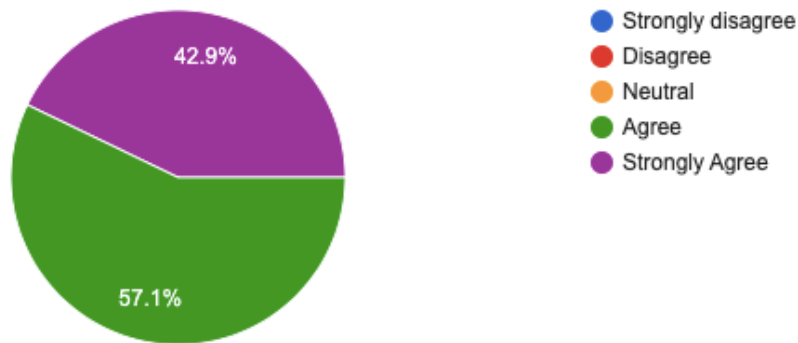


Figure 7-44 : Chart illustrates the responses of having a good response time from the system.

Visualisation of the Responses of the Reliability of the CDSS.

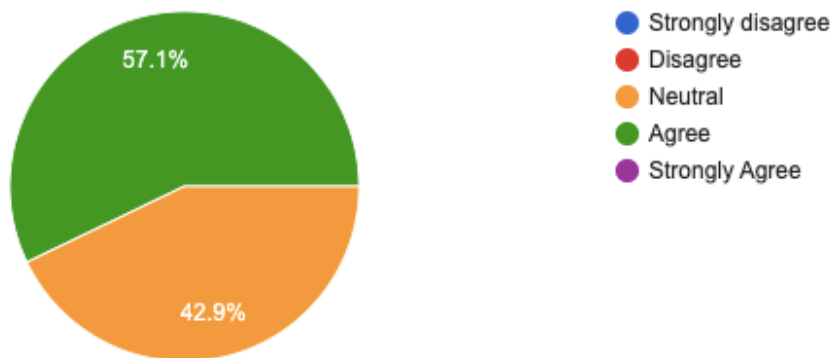


Figure 7-45 : Chart illustrates the responses of the reliability of the CDSS.

Moreover, the participants mentioned valuable suggestions, criticisms, compliments, and recommendations for the system. The following table summarises the user comments.

Suggestions	Criticism	Compliment	Recommendation
Include a description about the specific domain that the system can use.	Peculiar behaviour of the system for the predictions which are more than 10 years difference in between diagnosis of diabetes and the current age.	System has great potentials of educating the patients and making shared-decisions.	Suitable for chronic disease treatment centres.
Provide measuring units for the laboratory values.	Mobile platform use has distorted the appearance.	Easy to navigate.	Since the system has the potential to predict the whole set of complications this is recommend for shared decision making.
Explain the term “survival probability”.	Predicting hypoglycaemia only based on demographic details may not be effective.	App is very easy to use.	
Specify the clear guidance on defining ethnicity.	Returning to home page is not clear.	Helpful for medical professionals.	
Draw more clearer survival curves with explained legends.	The graphics used in head banner of the tool are not in high resolution.	A good system for predicting the complications holistically.	
Use the colour schema of Health NZ			

Table 7-4 : Table of user feedback for the system evaluation.

The collected feedback provided great insights into the usability of the generated CDSS. Overall, the design, navigation, performing tasks, and user experience were positively characterised. At

the same time, visually appealing quality is suggested to enhance through high-resolution banners, different colour schemas, and platform-independent graphs. The capability of predicting accurate survival rates was questioned when the prediction was not made within a decade of diagnosing diabetes.

The above section provides a clear explanation on the evaluation of a CDSS with existing standards. The adopted evaluation methods of the study in design evaluation and implementation evaluation are described with the gained results of the study. Furthermore, the suitability of the adopted techniques and their usage in the study are explained to answer the research question. The chapter 6 is dedicated for the evaluation phase of the DSRM, where the 6.1 and 6.3 chapters describe evaluation process in the design perspective of the study. The answer for the RQ1.2 is developed from above mentioned chapters while providing the real-time results obtained in the study to confirm their applicability in a similar scenario. The summary of the used evaluation techniques present here. The design evaluation was done by focusing on four elements: evaluating the requirements, investigating the feasibility of requirements, mitigating the risk of mismatches with the client's requirements, and enriching the design with the potential functionalities to enhance the system's usability. Moreover, a four-peered design evaluation phase was used in the study to reveal information on configuring the requirements, selecting the functionalities, mapping the stakeholder needs to the system requirements, and validating the design. The implementation evaluation was carried out using two methods to confirm the validity of implemented system in the real world: ISO-25010 and through collecting the user feedback. The evaluation through ISO-25010 focused on a quantitative approach to validating the product quality. A set of external characteristics of the implemented system was mapped with the standards to evaluate the existence of standard quality measures in the system. The procedure followed to evaluate the CDSS through ISO-25010 explained the practical usage of standards in the evaluation while confirming its suitability in the healthcare domain. The use of ISO-25010 in the implementation evaluation of a CDSS confirms the adaptability of standards in information system research while guiding future research on evaluation through standards. Additionally, a summative evaluation of the system through user feedback was performed in this research to

confirm the product quality in a qualitative evaluation method. User-centred evaluation methods are always encouraged to be performed by the system implementation to bring the users' perspective of the system. Selecting the user feedback collection method, potential participants for the evaluation, developing the questionnaire to cover the whole spectrum of characteristics that need to be evaluated, and expressing the result of gathered user feedback are explained in the implementation evaluation phase of the research. The knowledge extracted from the literature and the procedure developed to use in the evaluation of CDSS through user feedback enhances the knowledge of implementation evaluation in theoretical and practical aspects. The used process of these two approaches of evaluation is building a systematic approach which is able to validate the current system while resulting in a systematic process of implementation evaluation. The identified knowledge gap of utilisation of evaluation methods and systematised process of evaluation of the CDSS can be filled with the aforementioned outcomes of the study.

7.2 RQ2: How can the onset of CoDM be predicted using a longitudinal data set?

The data analysis perspective of the research is explained through this research question. The backbone of this study is accurately predicting the onset of complications through suitable data analysis methods. The appropriateness of the survival analysis techniques in predicting the complications of diabetes using a longitudinal dataset is trying to confirm through this research question. Survival analysis techniques were selected to analyse the longitudinal data set of Te Whatu Ora to model them into data prediction models. Survival techniques are prominently used in various fields, including medicine, healthcare, engineering, biology, marketing, and social sciences. Moreover, the selection of survival techniques has to be done with a conscious understanding of the hand-in issue. The following two research questions are derived to answer a vital aspect of the CDSS which is data analysis.

7.2.1 RQ2.1: How can the existing techniques for predicting CoDM be utilised in a cohort of New Zealand?

The selection of suitable survival analysis techniques for achieving the desired purposes depends on the nature of the dataset. The survival analysis of this study was conducted using non-parametric and semi-parametric methods. Kaplan-Meier survival techniques were selected as the non-parametric method to analyse the survival of the cohort in selected complications of diabetes. The resulting survival curves of the cohorts are highly beneficial in comparing them under demographic value changes, such as gender, age groups, ethnicity, and Māori/non-Māori. The outcomes of the nonparametric survival curves delivered crucial information about the cohort, which is directly applicable in the management of healthcare settings, including policy-making, resource allocation, decision-making, and forecasting. The confirmation of the suitability of selected survival techniques in this study provides a guideline for achieving similar purposes while filling the knowledge gap of model selection and implementation for similar approaches.

7.2.1.1 *Non-Parametric Survival Curves*

The survival curves generated from the Kaplan-Meier estimator are used in observing the statistical differences between the strata of cohorts with demographic changes. The resulting graphs reveal important information about cohorts in Aotearoa. The cohort of E119 were graphed by differentiating according to their gender, Māori/non-Māori, age, and ethnicity. The following graphs illustrate the characteristics of E119.

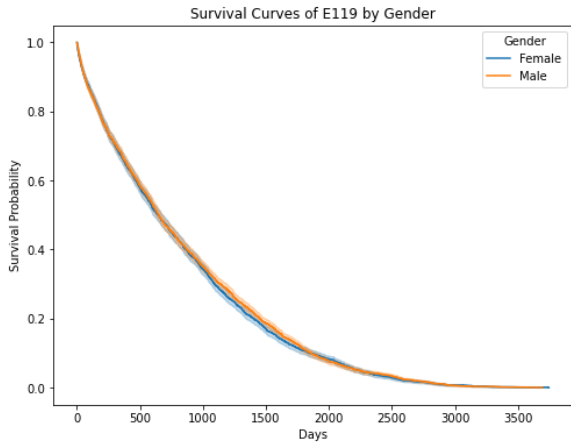


Figure 7-46 : Kaplan-Meier curve of E119 cohort by gender.

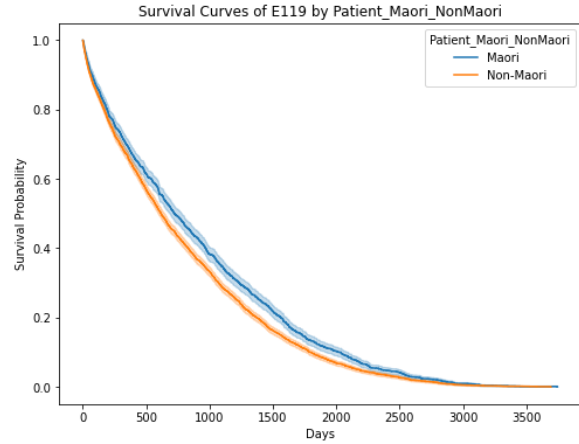


Figure 7-47 : Kaplan-Meier curve of E119 cohort by Māori/non-Māori.

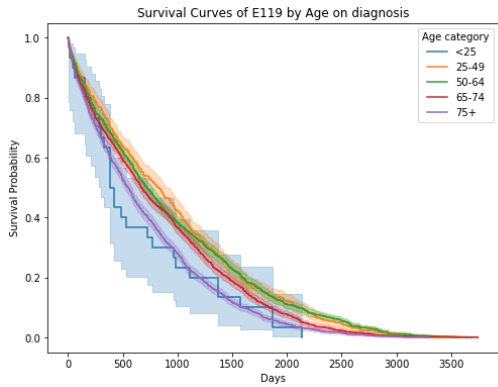


Figure 7-48 : Kaplan-Meier curve of E119 cohort by age categories.

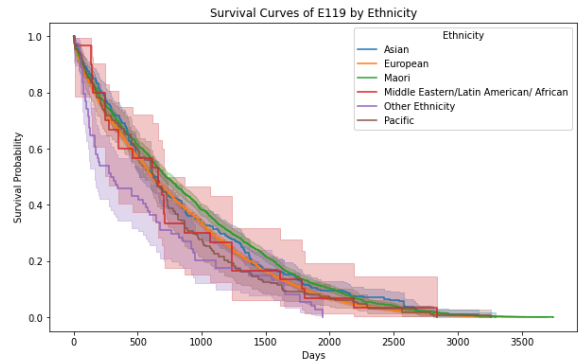


Figure 7-49 : Kaplan-Meier curve of E119 cohort by ethnicities.

Additionally, the resulting Kaplan-Meier curves of the cohorts of selected complications are included below to present the general survival probabilities of each cohort.

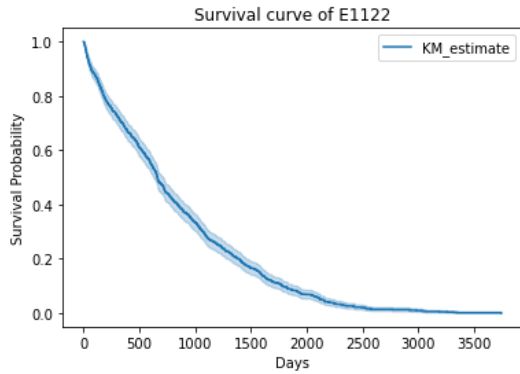


Figure 7-50 : Survival curve of E1122.

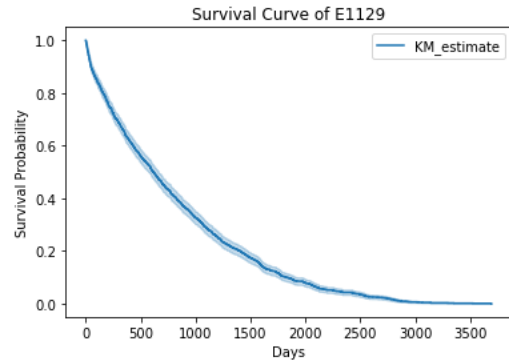


Figure 7-51 : Survival curve of E1129.

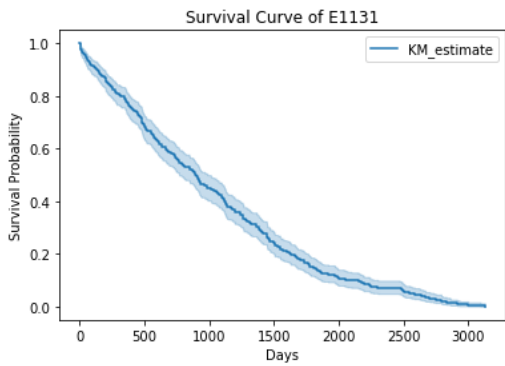


Figure 7-52 : Survival curve of E1131.

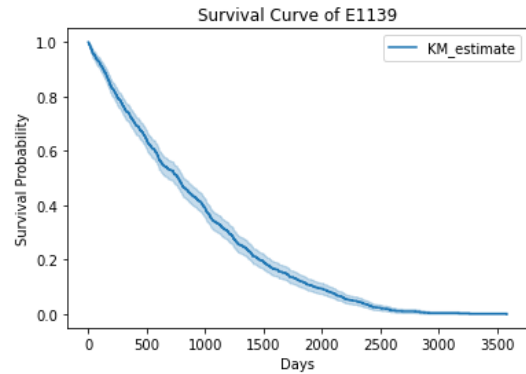


Figure 7-53 : Survival curve of E1139.

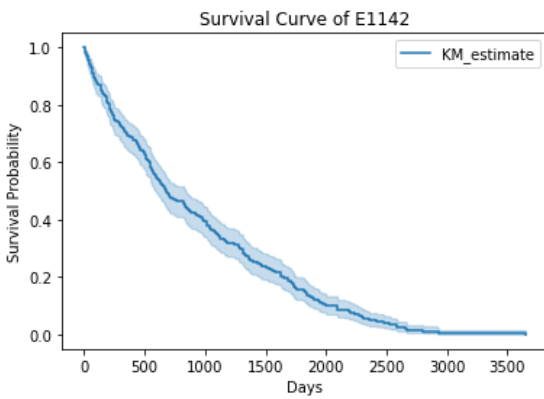


Figure 7-54 : Survival curve of E1142.

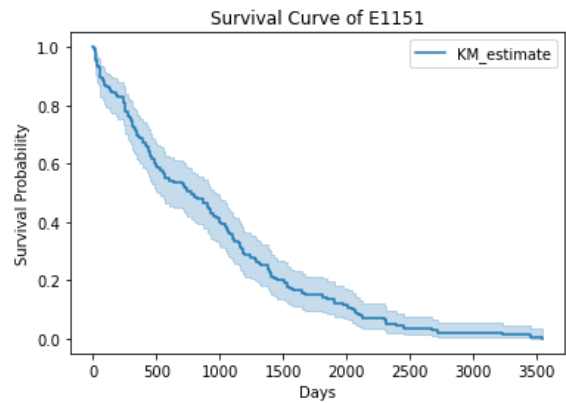


Figure 7-55 : Survival curve of E1151.

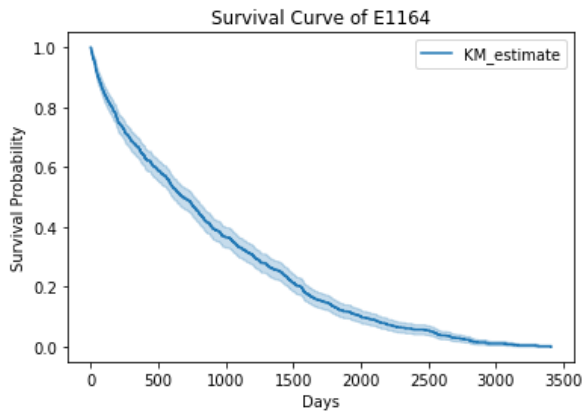


Figure 7-56 : Survival curve of E1164.

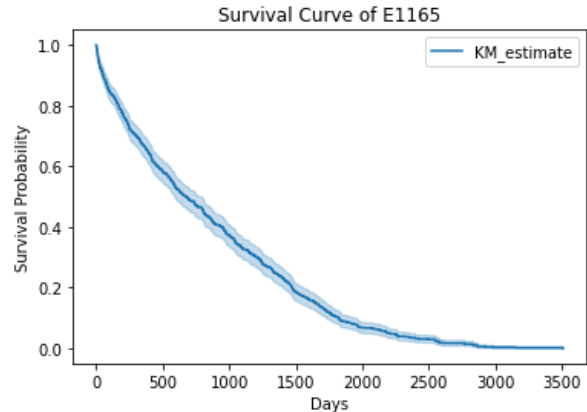


Figure 7-57 : Survival curve of E1165.

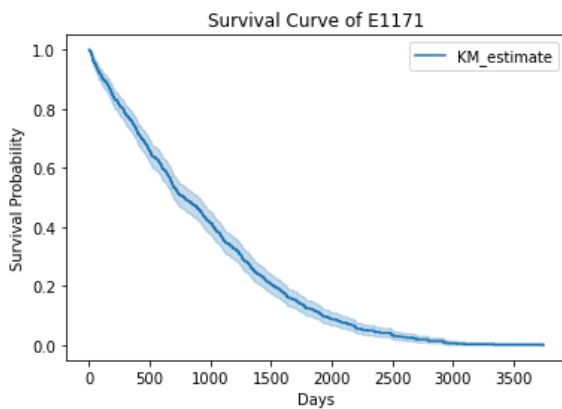


Figure 7-58 : Survival curve of E1171.

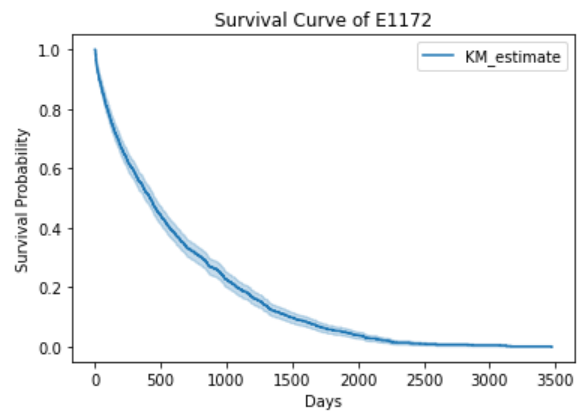


Figure 7-59 : Survival curve of E1172.

The survival curves can be used to analyse the characteristics of the cohorts of complications. The survival curves of complications such as E1122, E1139, E1151, E1165, and E1172 reached 0 survival probability around 2500 days. The E1131 and E1164 cohorts presented the survival curves as steeped functions throughout the period. Moreover, reaching 50% of survival is around 750 days in E1122, E1129, E1139, E1142, E1164, E1165, and E1171 complications. The shortest time of reaching 50% of survival belonged to E1172, whereas E1131 and E1151 took around 1000 days to reach a 50% survival rate.

Moreover, the survival curves of each complication categorised by gender, ethnicity, Māori/non-Māori, and age categories were created to visualise the details of each cohort. The graphs of E1122 are included here as examples.

The survival curve of the cohort of E1122, which was stratified according to gender, is as follows. The two survival curves did not proportionally change over time. The survival rate was similar among females and males around the survival probabilities such as 1–0.65, 0.37–0.2, and 0.1–0.0. Among the day's 400–900, the curve of females resided above and between 1250–1800 beneath the curve of males.

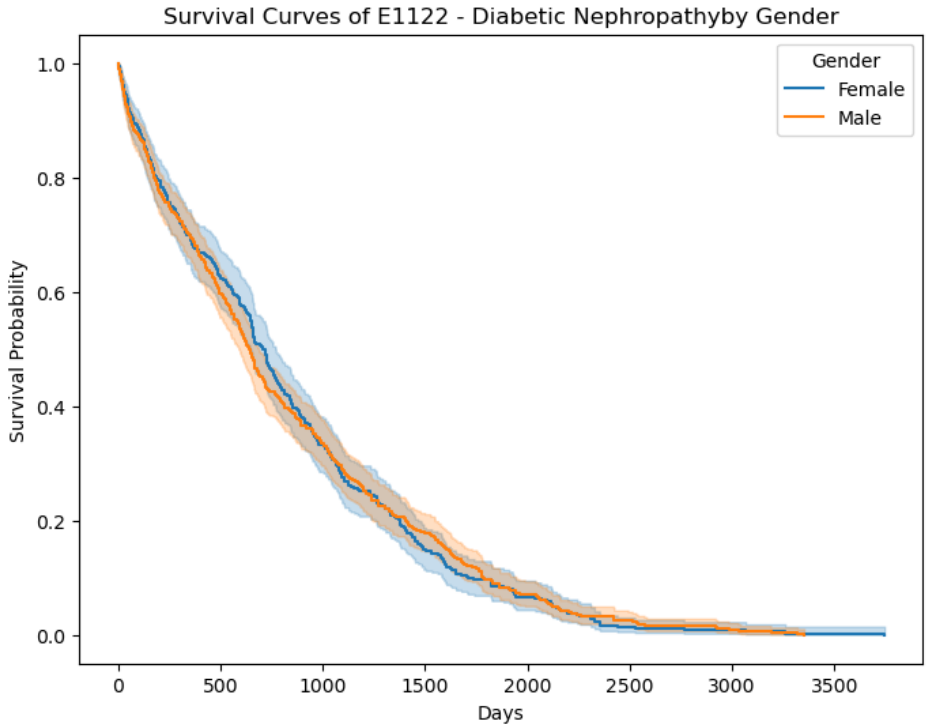


Figure 7-60 : Survival curve of E1122 by gender.

The following graph shows the survival changes among the different age categories. The curve of the age category of '75+' is a steeper curve than the curve of the age category of '25-49'. The patients who were diagnosed with diabetes at young ages had a better survival chance for E1122 than those who were diagnosed later in life.

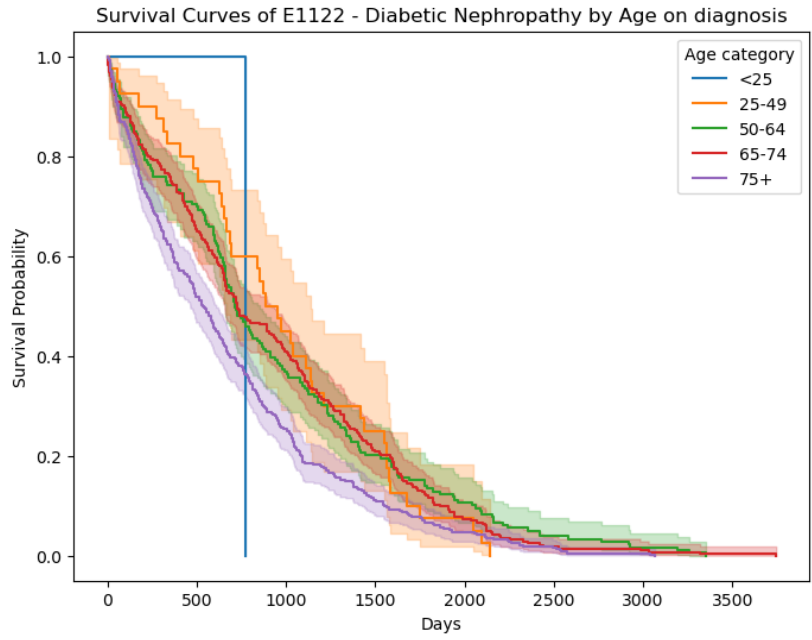


Figure 7-61 : Survival curve of E1122 by the age of diagnosis.

The following graph shows the ethnic differences among the considered cohort of patients. The graph shows that every ethnicity reaches 0 survival probability around 2250 days, except “Pacific” and “Middle Eastern/Latin American/ African”. The limited number of patients in those ethnicities may result in confusion. The time taken to reach 50% of survival probability was increased from the ethnicity of “other ethnicities”, “European”, “Asian”, “Pacific”, “Middle Eastern/Latin American/ African”, and “Māori, respectively.

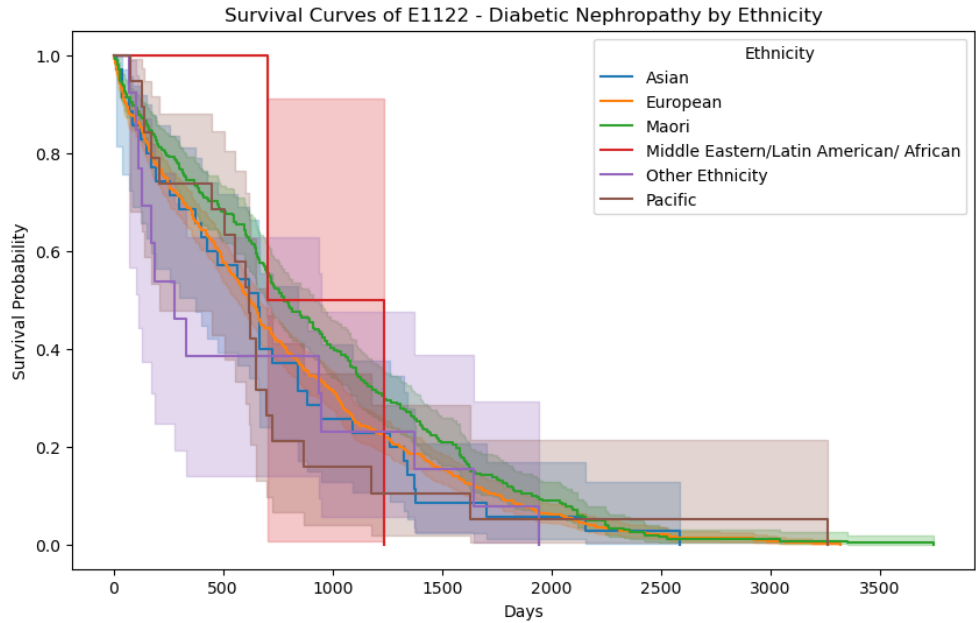


Figure 7-62 : Survival curve of E1122 by ethnic groups.

The following curve shows the difference in survival among the Māori and non-Māori populations. A better survival chance can be seen among Māori than the non-Māori population. When the days reached 2250 the survival curves overlapped. Additionally, the confidence interval of the survival curve of “Māori” had a higher width than “non-Māori” while leading to a higher prediction uncertainty in the population of “Māori”.

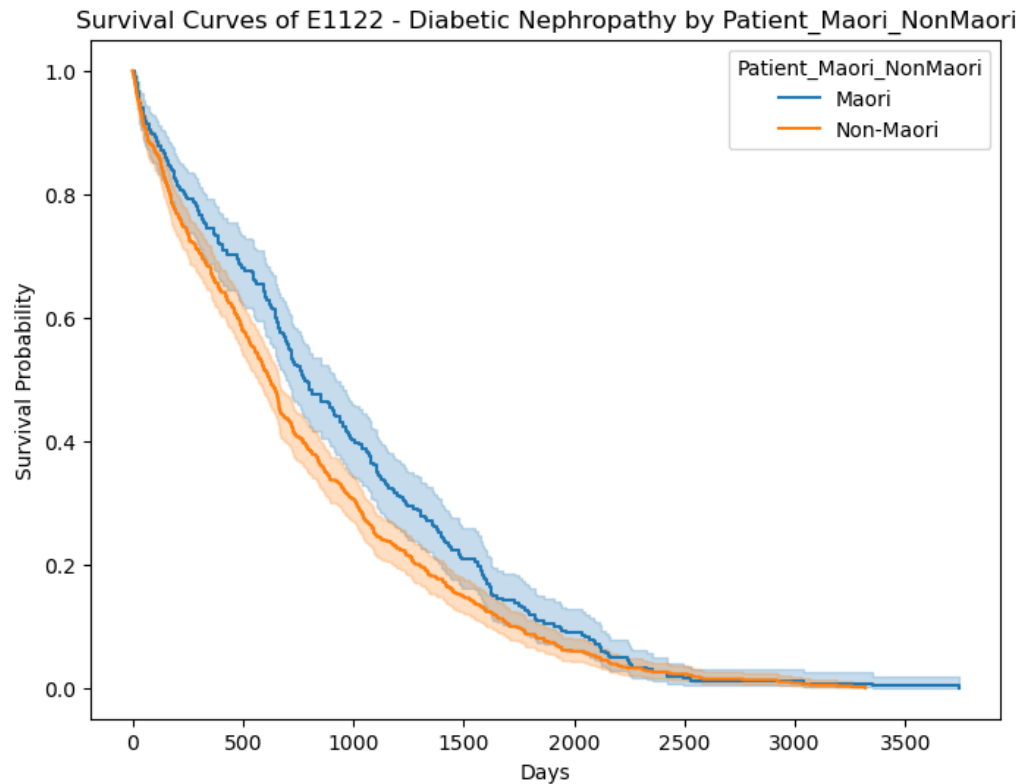


Figure 7-63 : Survival curve of E1122 by Māori/non-Māori.

The resulting survival curves of each complication concerning the demographic strata were compared to analyse the differences among the characteristics of the cohorts. The resulting survival curves of each cohort on gender, age, ethnicity, and Māori/non-Māori values were created and tested with a log-rank test for checking the presence of significant differences among the cohort. The gender curves overlapped and were similar in shape, although their log-rank values were statistically significantly different. The survival curves of the gender value of E1131 and E1142 showed that the curve of female patients resided beneath the curve of male patients, which presented that females have better chances of surviving on E1131 and E1142 than males. The age curves created from cohorts also showed a significant difference in each age category of all cohorts. Additionally, most showed the same characteristic of having a steeper curve for adults than young patients. The patients for the age category '<25' were rarely present in the data frame, which led to some abnormalities in the category of '<25'. The age curves for most cohorts started and ended within the considered range, except for E1139, E1142, E1151, and E1172. The

age curves of '75+' of E1142 and E1172 ended around 2000 days, and E1139 and E1172 ended around 2250 days, which interpreted that the survival duration of those complications was shorter for the age groups of 75+ patients in the cohort.

The ethnicity survival curves of the cohorts generally followed the same ethnicity orders as the E119 cohort showed. However, Pacifica showed fewer survival chances at the initial diagnosis of diabetes in E1131 and E1139 than other ethnicities. Additionally, E1151 and E1171 showed abnormalities of the survival curves instead of laying as the highest chances of surviving. The survival curves of European ethnicity were always beneath the survival curve of Māori. However, the curves of Pacifica resided above the survival curves of Māori at the initial diagnosis time of diabetes in E1129 and E1151.

7.2.1.2 Semi-Parametric Survival Curves

Moreover, the selected semi-parametric model is the Cox proportional hazard model. Two types of models are developed using the Cox models by considering different feature sets: 1) demographic details and 2) demographic details with laboratory values. The collected dataset of the New Zealand cohort was pre-processed with feature selection methods to predict survival based on demographic and laboratory features. The CDSS can predict the patient's survival of complications of diabetes based on these two feature sets.

The following section presents the resulting survival curves of the selected complications predicted using the demographic details of a hypothetical patient. The survival curve for each complication was graphed with the survival curve of the cohort to provide easy comparison. The entered details of the patients are presented as a forum in the following diagram.

Enter your demographic details

Gender

Female

Maaori

Non-Maaori

Ethnicity

European

Enter your age at diagnosis of Diabetes

47

Enter your current age

48

Submit

Figure 7-64 : Demographic details forum filled for predict the patient's survival rates.

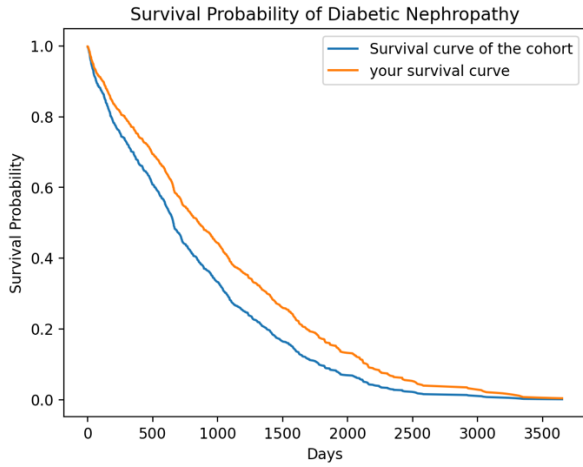


Figure 7-65 : Survival curve of E1121 results from Cox model using demographic details.

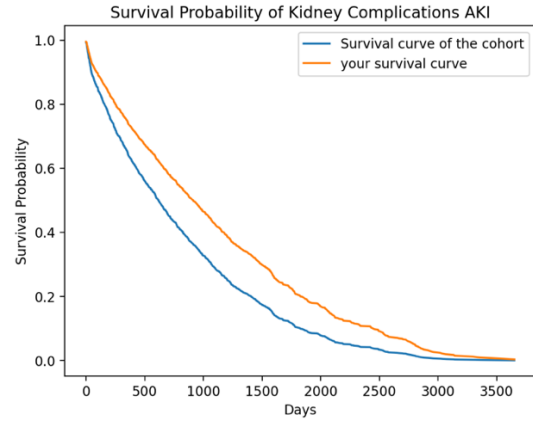


Figure 7-66 : Survival curve of E1129 results from Cox model using demographic details.

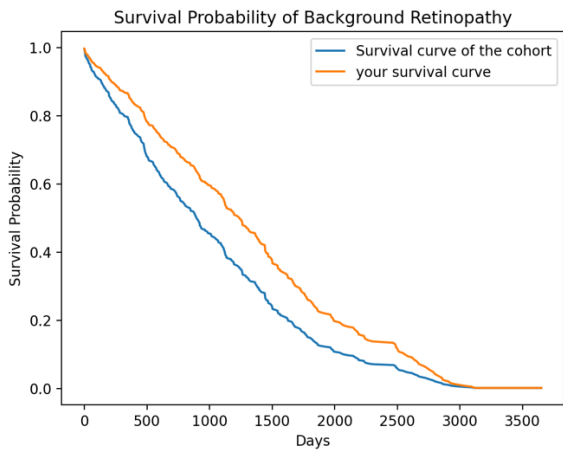


Figure 7-67 : Survival curve of E1131 results from Cox model using demographic details.

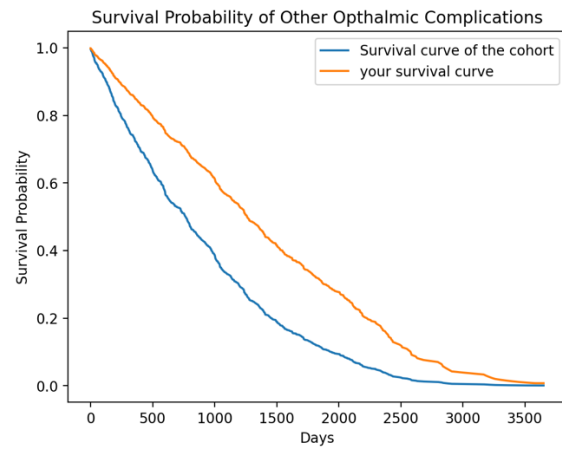


Figure 7-68 : Survival curve of E1139 results from Cox model using demographic details.

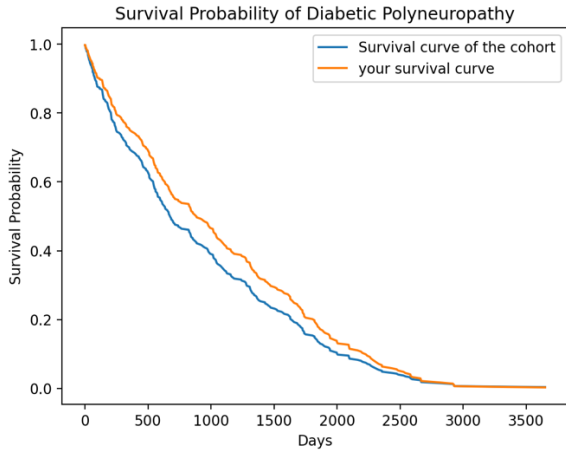


Figure 7-69 : Survival curve of E1142 results from Cox model using demographic details.

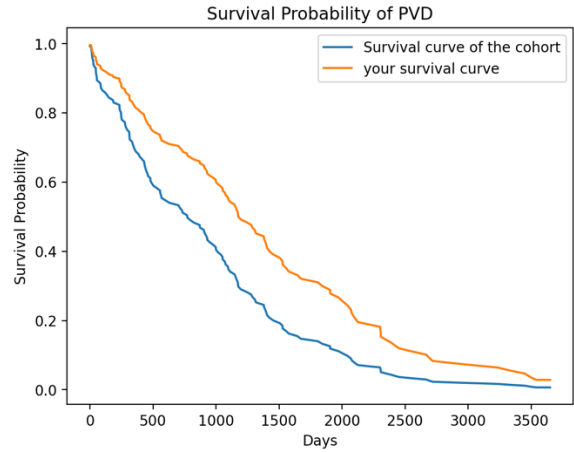


Figure 7-70 : Survival curve of E1151 results from Cox model using demographic details.

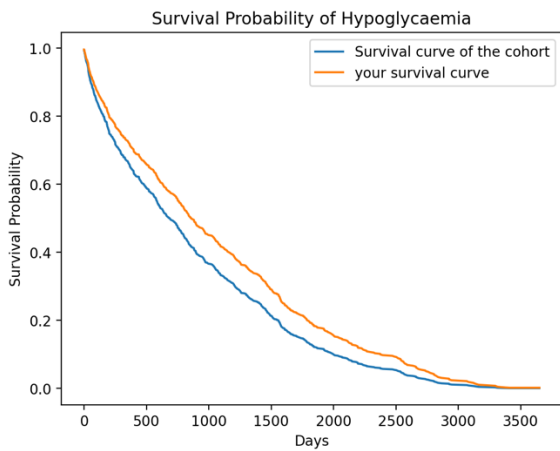


Figure 7-71 : Survival curve of E1164 results from Cox model using demographic details.

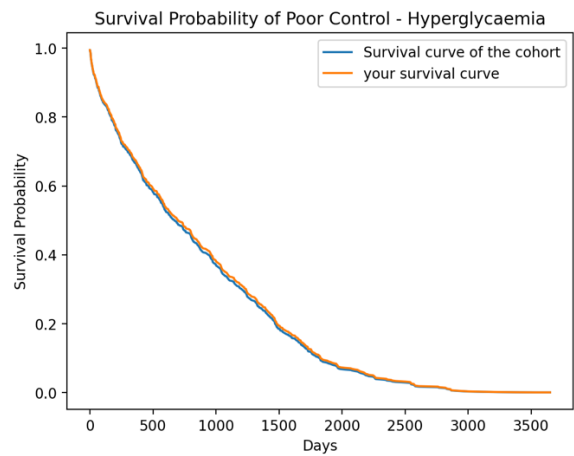


Figure 7-72 : Survival curve of E1165 results from Cox model using demographic details.

Survival Probability of Microvascular and other specified nonvascular complications

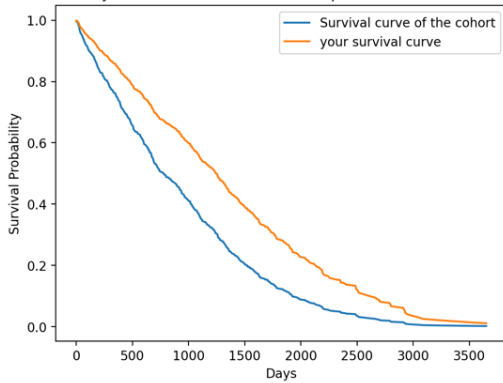


Figure 7-73 : Survival curve of E1171 results from Cox model using demographic details.

Survival Probability of Fatty Liver

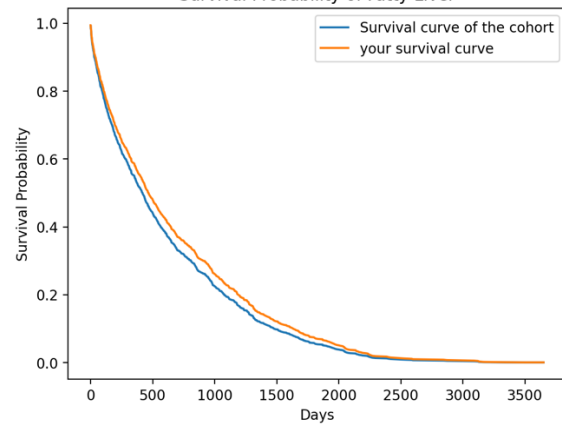


Figure 7-74 : Survival curve of E1172 results from Cox model using demographic details.

Additionally, a set of Cox models were created to predict the survival of the patients using a combination of demographic and laboratory values. Although the laboratory values were not comprehended enough as was extracted from the systematic review, the study process wanted to demonstrate the possibilities of using different features for predictions and exploring their validity in the domain. The results gained through the Cox models, which used all the considered features, are attached as follows. Further, the following diagram represent the used values for creating the survival curves, through the CDSS.

Gender

Female

Maaori

Non-Maaori

Ethnicity

European

Enter your age at diagnosis of Diabetes

47

Enter your current age

48

HbA1c ? Cholesterol ?

48.00 5.20

Triglyceride ? HDL ?

1.50 1.60

LDL ? eGFR ?

2.60 49.00

Submit

Figure 7-75 : Forum for predict the survival of a patient with all features.

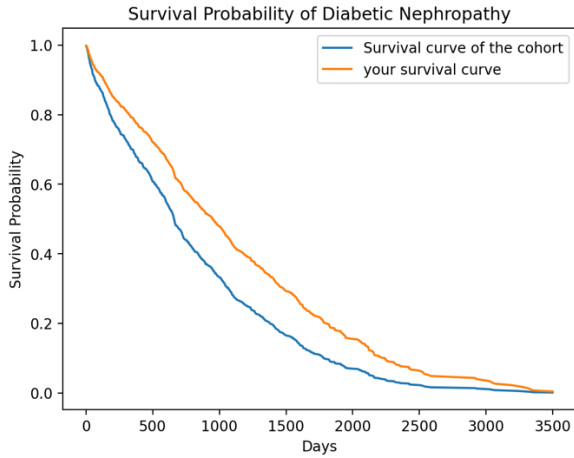


Figure 7-76 : Survival curve of E1122 results from Cox model using all features.

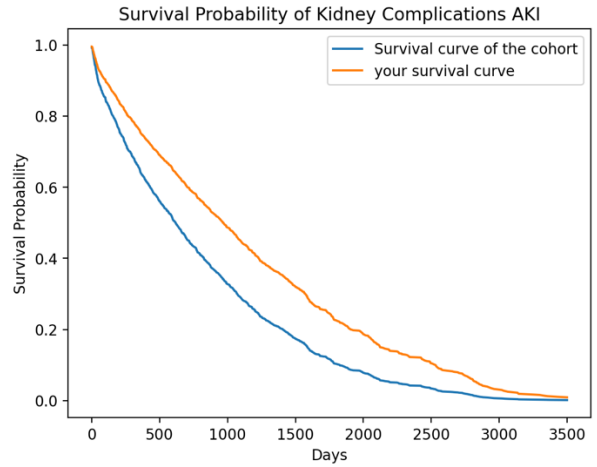


Figure 7-77 : Survival curve of E1129 results from Cox model using all features.

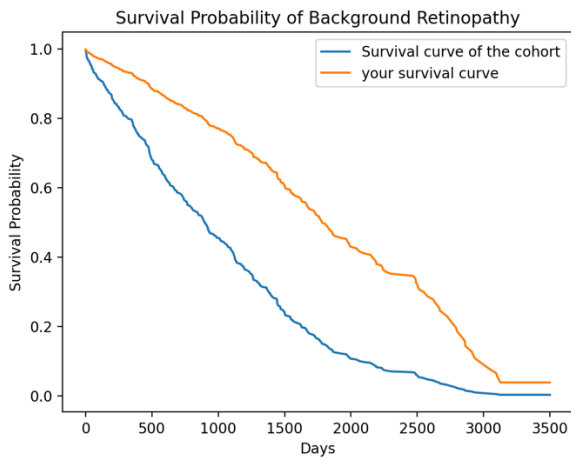


Figure 7-78 : Survival curve of E1131 results from Cox model using all features.

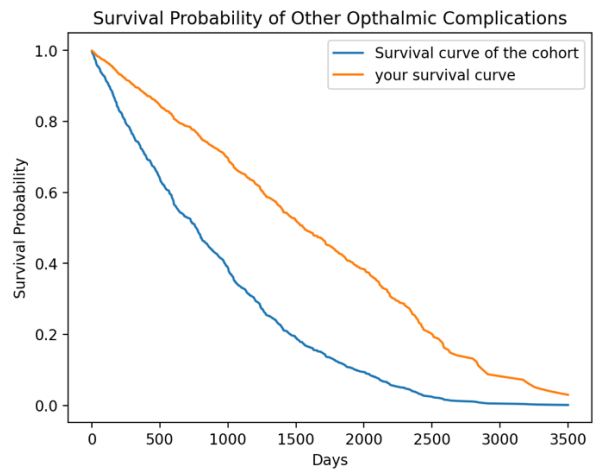


Figure 7-79 : Survival curve of E1139 results from Cox model using all features.

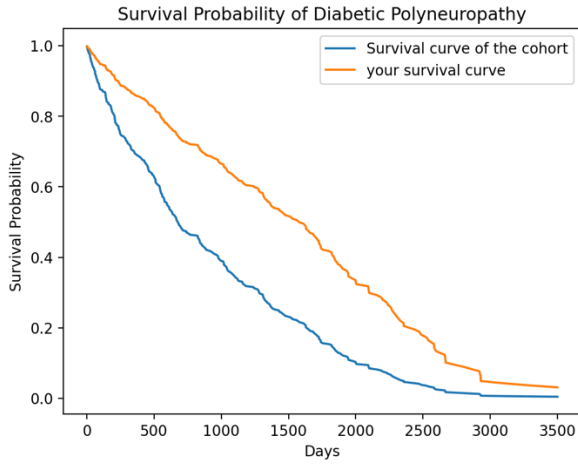


Figure 7-80 : Survival curve of E1142 results from Cox model using all features.

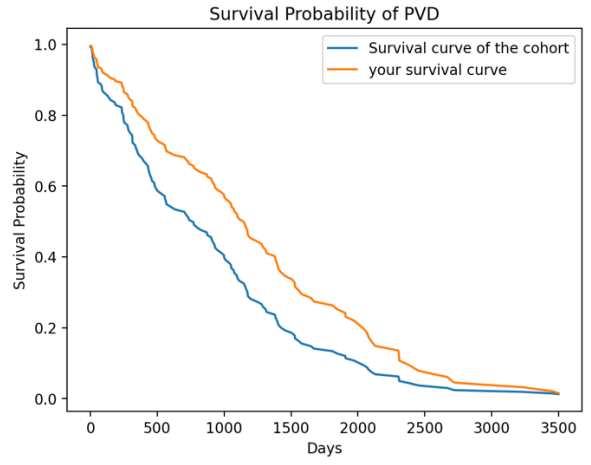


Figure 7-81 : Survival curve of E1151 results from Cox model using all features.

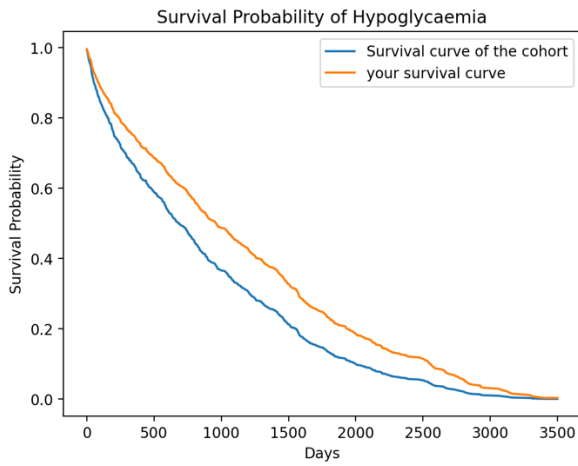


Figure 7-82 : Survival curve of E1164 results from Cox model using all features.

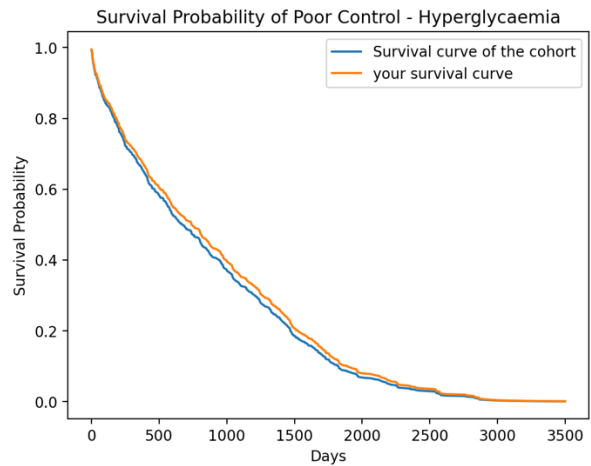


Figure 7-83 : Survival curve of E1165 results from Cox model using all features.

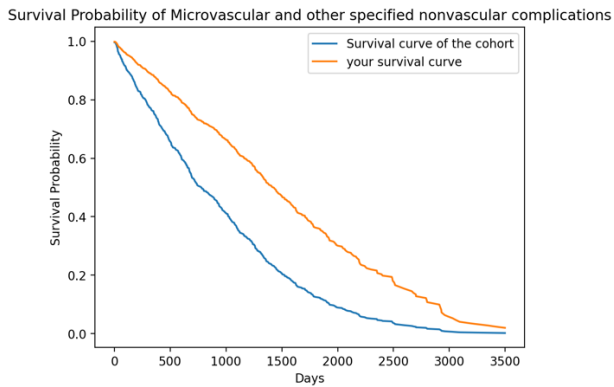


Figure 7-84 : Survival curve of E1171 results from Cox model using all features.

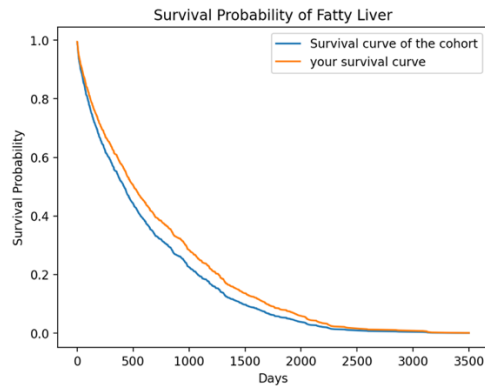


Figure 7-85 : Survival curve of E1172 results from Cox model using all features.

The resultant survival curves can be used to support the decision making on patients' health management while comparing their health with the survival curves of cohort. Further, the prediction capability of the CDSS delivers a package of practical outcomes. The resulting personalised survival rate for each individual based on two types of features can assist general practitioners in making decisions, such as conducting laboratory tests to diagnose complications, commencing treatments and medications, and advising patients on healthier lifestyles. Moreover, the patients can use their survival curves to be aware of their health conditions and be motivated to lead a life to mitigate the risk of complications. The utilisation of existing survival techniques for predicting the complications of diabetes primarily serves practical benefits to the individual and society to enhance health indices by considering the given results.

The developed models enhance knowledge of developing algorithms for forecasting through survival analysis. It confirms the usage of knowledge of using survival analysis to predict complications of diabetes from a longitudinal dataset. The selection of the models confirms the applications of Kaplan-Meier and Cox proportional hazard techniques in predicting the survival rate of individuals. The resulting survival curves fill the existing contextual knowledge gap of cohorts in New Zealand while the selected techniques confirm the applicability of survival analysis techniques.

7.2.2 RQ2.2: How Accurate is the Prediction of CoDM?

The evaluation of all components in a system makes the entire system more validated and credible. Further, the algorithm's accuracy is a significant factor in validating the practical applicability of the system. The trained models of each complication have been validated through C-index values and Brier scores. Moreover, the models' C-index values and Brier scores were assessed with a cross-validation method to validate the created models. The selection of evaluation methods and scores in survival analysis has its own controversialities. The assessed accuracies of the models are beneficial in confirming the credibility of the predictions achieved for each complication. The practical applicability of implemented CDSS can be verified through the accuracy of the trained models. Further, evaluating models explains the suitability of adapted methods in the considered domain.

Results of Algorithm Evaluation

The following section describes the results of the algorithm evaluation phase. The Cox models' C-index values, created based on demographic and laboratory details, are calculated separately through 10-fold cross-validation techniques. The following table represents the resulting C-index values from the 10-fold cross-validation of the Cox models.

Complication of diabetes	Prediction based on demographic details	Prediction based on demographic and laboratory details
E1122	0.583	0.609
E1129	0.575	0.583
E1131	0.551	0.575
E1139	0.604	0.600
E1142	0.523	0.593
E1151	0.580	0.583
E1164	0.505	0.580
E1165	0.588	0.646
E1171	0.602	0.588
E1172	0.599	0.651

Table 7-5 : Algorithm evaluation results.

Moreover, a figure that represents the C-index values of the above-mentioned two types of Cox models is included in the following. The results show that the models created using the details of demographic and laboratory values have higher C-index values most of the time, except E1172. This makes it clear that the models created with more covariates of laboratory values fit better than the curve results using only the demographic details.

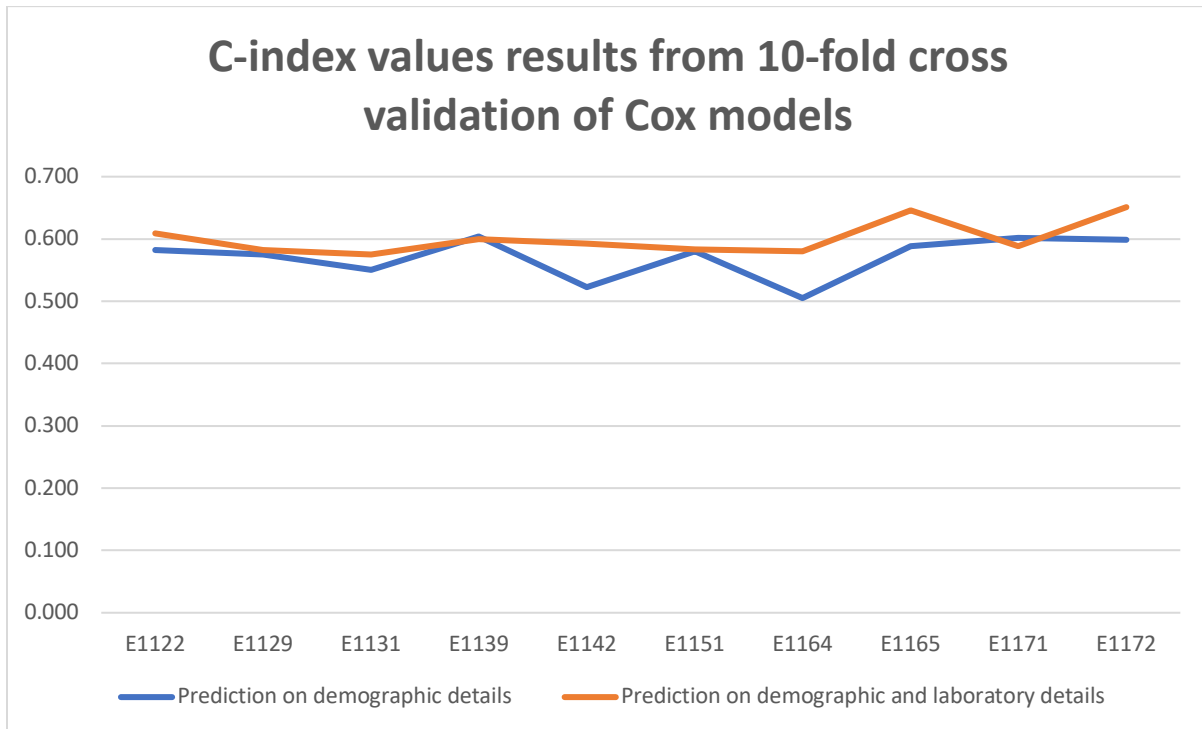


Figure 7-86 : Visualising algorithm evaluation results.

Additionally, the C-index value of E1164 with the demographic detail prediction model was close to random prediction. E1164 was a situation that occurred due to the effect of medication on diabetes. Since this is a common complication among diabetes patients, we need their history of usage of medication to predict the E1164. Although some ambiguities can be seen in a few occurrences of the graph, most of the time, the C-index values were around 0.55–0.6. The C-index values might be enhanced if the models are trained with the extracted feature sets.

The Brier scores were also used to evaluate the Cox models with cross-validation techniques. The resulting Brier scores are presented in the following table.

Complication of diabetes	Prediction based on demographic details	Prediction based on demographic and laboratory details
E1122	0.3193	0.2955
E1129	0.2456	0.3193
E1131	0.2631	0.3119
E1139	0.3233	0.2519
E1142	0.2625	0.3420
E1151	0.2792	0.2714
E1164	0.2600	0.2792
E1165	0.2915	0.3740
E1171	0.3286	0.3219
E1172	0.2545	0.3531

Table 7-6 : Accuracy of created Cox models.

The resulting Brier scores show that the Cox models accurately predicted the complications since the values were closer to 0 than 1. This shows that the accuracy of Cox models is in the acceptable range.

The algorithm evaluation provided information on the accuracy of the created models for predicting the selected complications of diabetes. The models' accuracy presented the prediction's success, where the models can be used to assist in making personalised decisions. Further, comparison of Cox models with and without including the laboratory values in the feature set, show a general differentiation. The high accuracy shows from the model trained with the laboratory values, which express the importance of laboratory values in predicting the complications of diabetes. The selected evaluation matrices and their applicability of the domain confirmed the validity of the models. The findings of this section express the methods of evaluation and the resultant accuracies, which can be considered as an extension of knowledge and practical implication of the study, while providing a solid answer for RQ2.2.

7.3 Summary

In conclusion, the identified theoretical and practical research gaps are answered through the above-mentioned four research questions, by focusing on designing and data analytics perspectives. The overarching research question is “How can a prediction model be created to forecast the onset of CoDM in a cohort of diabetes patients?”. The overarching question achieved its answer through the solutions of those four sub research questions. The adopted design process of the study is a vital outcome of the research, which describes the systematic use of concepts of DSRM in building applications for solving real-world issues in the healthcare sector. The applicability and utilisation of the concepts of DSRM in the phases of problem investigation, designing, implementing, and evaluating a system fills the identified knowledge gap in designing perspectives of the system. The techniques used for the steps of problem identification, situational awareness, data collection, data pre-processing, feature selection, refining the scope of the CDSS, implementation of the CDSS as a web portal, and evaluation of the CDSS are validated in this research. The process is systematised for similar approaches, filling the knowledge gap of designing perspectives. Further, the resulting feature set of the systematic review conducted in the feature selection stage of the study fills the knowledge gap of having a standard feature set as predicates for the complications of diabetes. The design evaluation phase of the study results in the suitable evaluation methods for a CDSS while confirming the validity of the implemented system and applicability of the extracted process in the domain. The model selection, implementation, and evaluation phases are focused on the data analytics perspective of the CDSS. The applicability of survival analysis techniques in predicting the complications of diabetes among a New Zealand cohort is confirmed through this research. The data analytics perspective of the research answers the knowledge gap of selection and usage of survival analysis techniques in a CDSS. In addition, the resulted CDSS mitigates the drawbacks of the exiting CDSS while filling the identified practical research gaps. The fixed period of prediction, presentation of risk, limitations of the capability of prediction, and the contextual gap in New Zealand are answered with the implemented CDSS. The resulted CDSS serves the purpose of predicting the complications of diabetes while improving its characteristics by answering the identified issues.

Chapter 8 Discussion and Recommendation

The following section discusses the research implications from academical and managerial standpoints. The contribution of the research outcome on the enhancement of the domain, the limitations of the study, and finally, the directions of future research are discussed in this chapter. This research was undertaken to accomplish the primary goal of implementing a CDSS to resolve a real-world issue. Te Whatu Ora has been experiencing issues managing diabetes patients' details on fruitful domains. The existing knowledge gaps in the literature and the hands-on practical problems, and the benefits it can bring into the decision-making led the study to investigate the solution of a CDSS. The overarching research study is two-fold: the design of the information system and the data analysis approach adopted to predict the complications. Design science research methodology (DSRM) has been used in the study to gain situational awareness of the issue, design, and evaluate the system, while survival analysis techniques are adopted to predict the survival of patients with different complications. The final outcome of the study is a clinical decision support system (CDSS) to predict the survival of diabetes patients for selected complications of diabetes.

8.1 Research Implications

Developing a CDSS by employing DSRM principles for predicting complications of diabetes via survival analysis holds significant research implications. Foremost among these is the profound practical impact of the CDSS effectively meeting the specific requirements of its end-users. Beyond its practical utility, the study engenders academic implications. The following sections explain the academic and managerial implications resulting from the research study.

8.1.1 Academical implications

The identified knowledge gaps of the domain regarding the design perspective of the CDSS are mainly focused on confirming the applicability of DSRM concepts in the healthcare sector. The utilisation of techniques in the stages of the empirical cycle, and evaluating the CDSS while the followed procedure of the study refined the theory of DSRM. This study adopted the method of DSRM introduced by Wieringa (2014). The basic three stages of this study to build the CDSS demonstrate the systematic process of solving the real-world issue, through problem investigation, design and implementation, and evaluation. Although, different approaches of this process have been introduced by past scholars (Hevner et al., 2004; March & Smith, 1995; Nunamaker Jr et al., 1990; Vaishnavi, 2007), the fundamentals are lying in the same ground. The adopted approach for this study initiates the process with design stage by analysing the real-world issue (Wieringa, 2014), reviewing literature (Nunamaker Jr et al., 1990), and gathering requirements (Weber, 2010) via brainstorming and client's meetings. The used techniques on exploring the real-world issue, the expectations of stakeholders, functional and non-functional requirements gathered through the above-mentioned techniques are providing a sound process of problem investigation. Moreover, the used approaches in situational awareness phase of the study, have directed to a pragmatic guideline on solving similar real-time issues in the healthcare sector. The use of the knowledge of literature and contextual understanding is vital when solving a real-world issue. The client meetings and brainstorming sessions contributes vastly in understanding the real requirements. The selected techniques on effectively communicating with the client are proved their validity on the process. The use of the results of EDA in clarifying requirements, conducting feasibility study, and defining the scope highly contributed in awaring of the real situation. Moreover, a requirement of a thorough systematic review on features for the prediction models was identified as a vital stage on the situational awareness stage. The recognised feature set of predicting the common complications of diabetes through the systematic review were actively contributed on understanding the situation of the health board and defiing the scope of the system. The situational awareness phase of this study provides a guideline for resolving similar issues with proven success. Additionally, the problem has been

frame worked by concerning the existing framework development strategies (March & Storey, 2008; Wieringa, 2014). The developed framework of the study consists of three components such as social context, knowledge context and design science approach. The developed framework is a vital outcome of the study which directs to generate frameworks for solving real-world issues in the healthcare sector. The design and implementation phase explains the process by most suitable techniques and their implementation on solving the hand-in issue. Each step in the phase of solution design and implementation contribute to refine the theory of DSRM. The data collection, data pre-processing, feature selection, model selection, and web-portal implementation provide a guidance on adopting different techniques in each step while systematising the process of designing and implementation in a healthcare sector. Additionally, the other major outcomes of this phase are the feature set extracted through the systematic review, and the implemented web-portal of CDSS. The existing issues of lack of a standard feature set for predicting the CoDM is overcome with the identified feature set while fulfilling the knowledge gap in the disciplinary. Clinicians, general practitioners, policymakers, and other stakeholders in the industry of healthcare can use the findings of this research to update their knowledge of the domain. The most frequently used feature set provide state-of-the-art and updated knowledge, beneficial in disease diagnosis and decision-making. Clinicians and healthcare professionals can utilise the identified risk factors as valuable tools to assess the likelihood of complications in diabetic patients, enabling personalised and targeted intervention strategies. The importance of each feature in predicting the complications provides a direction for clinicians to investigate them to make informed decisions. The identified set of features of the CoDM is highly useful for researchers who predict CoDM with statistical and computerised prediction models. By systematically extracting and analysing frequently used risk factors from the existing literature, the review sheds light on the critical determinants contributing to the development and progression of complications associated with diabetes. The feature set can be used as a quick reference in feature selection and feature engineering phases in model creation. Identifying these risk factors not only advances our understanding of the multifaceted nature of diabetes complications but also provides a comprehensive framework for future research endeavours. Further, academics can use the feature set to validate and prove the credibility of

their feature selection against not only individuals but also against a thorough literature review, which enhances the model's acceptability. The findings of the systematic review can adapt to future risk prediction, disease diagnosis and prognosis models. Additionally, the overall outcome of the systematic review provides a thorough state of the art, which keeps the clinicians and other stakeholders updated with the domain knowledge.

The implementation stage also covers the entire process of resulting the CDSS as a web-portal. Implementing the CDSS as a web-based information system provides a spectrum of knowledge from the perspective of explaining, including developing a computer algorithm, integrating the algorithm with user interfaces, assisting in making healthcare management decisions, visualising the statistics of the general cohort, and illustrating the appealing results on survival of patients. The techniques used, their suitability and applicability in the domain of healthcare fill a significant research gap while refining the theory of DSRM. The evaluation phase of the study provides a solid background in evaluation techniques and evaluation methods which fits with the similar problem domains. The evaluation strategy used in this study is adopted after a careful selection of evaluation methods. The FEDS method of evaluation Field (Venable et al., 2016) created the fundamentals of the evaluation phase, while the selected methods proved their validity in the domain. The evaluation strategy of this study is comprised with three components such as design evaluation (Wieringa & Morali, 2012b), algorithm evaluation (David & Mitchel, 2012) (Heller, 2021), and implementation evaluation (Kadi et al. (2016) (International Organization for Standardization, 2011) (Cao et al., 2004; Chiew & Salim, 2003; Zaharias, 2006). The selected methods in each component of the evaluation and their usage in the domain answer the research gaps identified in the study. The recognised knowledge gaps such as scarcity of the confirmed applicability of DSRM in the healthcare sector, the utilisation of techniques in the process of developing a CDSS, the evaluation strategies of the CDSS are answers through the first research question.

Further, the academical implication on data analysis component of the study is also a vital aspect.

Although survival analysis techniques are used in the healthcare sector for fulfilling various purposes, their suitability in CDSS is not confirmed. The model selection phase of this study explains the suitability of survival analysis techniques in the given domain, while the phase of model evaluation illustrates the statistically significant result. Survival analysis techniques are applied in the disciplines of social sciences, medicine, healthcare, engineering, biology, and marketing (Tolley et al., 2016), but it is rarely used in CDSS. Although, non-parametric, semi-parametric and parametric models of survival analysis are applied in the discipline of healthcare, it is challenging to select the most suitable model type to fulfil the purposes. The selection of the Kaplan-Meier model for visualising the cohort details illustrates the usage of non-parametric models in healthcare. Kaplan-Meier techniques basically used for comparing the cohorts (Jager et al., 2008; Oomichi et al., 2006; Shang et al., 2021). Kaplan-Meier method is used in this study to visualise the survival curves based on demographic values, which is ideal for comparing the cohort against age, ethnicity, Māori/non-Māori or gender. The results generated through the use of Kaplan-Meier techniques enhances the visually appealing quality of CDSS. Further, the created two types of data modelling for survival prediction, by Cox proportional hazard models, confirms their applicability in predicting the CODM with given feature sets. Selection of models by considering the availability of data set is expressed in this research by generated models. The use of Cox models in healthcare (Chien et al., 2009; Jia et al., 2019; Sim et al., 2022) is adopted in this research to fulfill the client's requirements while confirming their applicability in CDSS. The used evaluation methods (David & Mitchel, 2012; Gneiting & Raftery, 2007a) of the generated models reveals the suitable evaluation techniques and the accuracy of the models in the study. Additionally, the results of Cox models confirms the high accuracy of models with laboratory details while confirming the effect of laboratory details in risk prediction.

In conclusion, the academic implication of this study can be summarised as confirming the applicability of DSRM concepts by utilising survival analysis techniques for solving a real-world issue. The designed approach to achieve the CDSS in this research contributes to the knowledge of designing solutions for similar scenarios in healthcare management. This study provides a solid foundation on adopting DSRM in solving a real-world issue in the discipline of healthcare. The

procedure of resolving the issue provide a systematic process while refining the theory of DSRM in the field of healthcare. Further, the process of designing, implementing, and evaluating the CDSS results in a systematic approach to solving similar issues. The resulted system review fulfill the knowledge gap in predicting CoDM. Another significant implication is the confirmation of the applicability of survival techniques to analyse the dataset for achieving the desired goals in the healthcare sector. Further, the data analysis techniques adopted in study results a solid prediction method, with a proven accuracy. Thus the considered two aspects of CDSS contributes the academical advancements in the discipline of design sciences in a significant extent.

8.1.2 Managerial Implications

From a practical standpoint, the CDSS holds immense promise for improving diabetes care and management. The potential for revolutionary impact lies within the application of this system, which holds the capability to transform the diabetic data management. The managerial implications achieved through the mitigation of existing issues in CDSS result a more advanced and powerful CDSS. The CDSS enhanced the diabetes data management in various ways while delivering the stakeholders' expectations. Furthermore, the innovatively crafted CDSS proves to be a valuable asset, offering contextual enhancements that cater to the diverse needs of stakeholders within the expansive healthcare sector. The nurses, and general practitioners can use the CDSS to make the decisions regarding the risks of the complications of patients, conducting the diagnosis tests, initiating treatment plans, and issuing early warning for their patients. The resource allocators and policy makers can get statistical overview of the cohort, such as the distributions of socio demographic details, and survival curves, for making decisions on focusing at allocating resources to match with the cohort distributions, and assisting on policy making. Additionally, the patients of diabetes and the public can use the CDSS to assess the risk of complications of diabetes, where they can make informative decisions on their risk levels and plan their life style accordingly. The stakeholders can proactively act on their roles with the assistance of this innovative CDSS.

The data analysis perspective of the CDSS vastly contribute on improving the CDSS. The state-of-art of the domain leads to concentrate on selecting a suitable prediction method for predicting the risk of CoDM, and on evaluating the selected technique. The identified research gap of the existing CDSS leads the study to select a suitable prediction technique to mitigate the existing drawbacks. The nature of the prominently used prediction methods results practical issues in CDSS such as fixed time prediction, less useful risk presentation, limitation on the predicted number of complications, use of cohort details, and the contextual gap. The existing risk prediction methods provide the results in fixed periods such as 3,5, 7 or 10 years while some of the studies predict the one-time prediction (Bernardini et al., 2021; Chen et al., 2010; Hamedan et al., 2020; Kunhimangalam et al., 2014; Romero-Aroca et al., 2019; Saaristo et al., 2010; Schoen et al., 2015). Due to the irreversibility and severity of the CoDM, predicting the risk in a chronological manner has been recognised as a significant issue to handle. The data analysis method used in this study has been selected to mitigate this drawback. The used Cox proportional hazard model provides the risk of individuals along with the time, which would be beneficial in recognising the onset of CoDM sooner to the onset. The presentation of the risk is another vital aspect identified from the literature. In literature, the risk score is presented as high, moderate , or low risk (Bernardini et al., 2021; Kunhimangalam et al., 2014; Romero-Aroca et al., 2019; Schoen et al., 2015) , or as a numerical value (Chen et al., 2010), or as a binary outcome (Romero-Aroca et al., 2022). The selected survival technique of this study calculates the risk score as a percentage, which avoid the risk presenting drawback of the literature. The previous studies always concentrate on predicting a single CoDM (Bernardini et al., 2021; Hamedan et al., 2020; Kunhimangalam et al., 2014; Romero-Aroca et al., 2019; Romero-Aroca et al., 2022; Schäfer et al., 2021; Schoen et al., 2015; Vartiainen et al., 2016). or rarely predict multiple CoDM. Since a CDSS with the capability of predicting multiple CoDM is more beneficial in decision making, the resulted CDSS has been designed with the capability of predicting the most prevalent ten CoDM in the cohort. The selected data analysis technique for solving the above-mentioned issue is Cox proportional hazard regression model. The results of the Cox proportional models mitigate the above-mentioned data analysis related drawbacks of the existing CDSS while confirming their suitability in the domain. The methods used for the stages of data analytical steps, such as data

pre-processing, feature selection, model selection, and model evaluation, are carefully selected in this study to demonstrate the applicability of the techniques. The justification of selected methods presents the suitability of the techniques while selected model and their results of evaluation confirms the applicability of survival techniques on the domain. The use of Cox proportional hazard model for the prediction of onset of CoDM, present their appropriateness on the similar domain. Further, its suitability for predicting the risk as a survival rate along with time, mitigate the identified gaps and confirms the applicability of survival techniques on CDSS.

The theoretical foundations of this Clinical Decision Support System (CDSS) provide a strong base for continually improving predictive models, making risk assessments more accurate and reliable. While the focus is on diabetes, the insights gained can apply to other chronic diseases, helping us better predict, prevent, and manage health outcomes. As this research solidifies, it leads to the creation of healthcare decision systems that support more personalized and informed patient care. This CDSS not only represents a theoretical advancement but also plays a key role in the evolving healthcare field. It's expected to make a lasting impact on both the theory and practice of healthcare technology, especially in managing diabetes in Aotearoa.

8.2 Limitations of the Study

Although the research study holds great promise in advancing CDSSs as a solution to the identified real-world issue, it is essential to acknowledge and navigate through the raised multifaceted limitations. The CDSS's potential for generalisability may face considerable challenges due to the inherently variable nature of clinical practices and the diverse demographics of patients encountered across disparate healthcare settings.

From a survival analysis standpoint, the CDSS is confronted with challenges related to the inherent assumptions and the prevalent issue of collected datasets. These factors can exert constraints on the system's predictive accuracy, influencing the reliability of the outcomes it generates. Finally, there are ethical and practical concerns to consider. Protecting patient privacy when handling sensitive data is a major challenge.

8.3 Recommendations for Future Research

As we continue to develop and implement a CDSS (Clinical Decision Support System) for predicting diabetes complications using survival analysis, many opportunities for future research arise. One important focus is to strengthen survival models by including more factors. These should go beyond standard clinical markers to also cover genetic, lifestyle, and environmental influences, recognizing their role in diabetes outcomes and aiming for a more complete model.

Future research should also look at improving risk predictions by using advanced machine learning (ML) techniques. These techniques can reveal complex patterns in data and help make predictions more accurate. Another key area is integrating real-time patient data into the CDSS, which could lead to more personalized and timely care. Developing adaptive CDSSs, which adjust as patient conditions change, could greatly improve the system's ability to provide effective interventions for diabetes. In addition, the design of CDSS interfaces should be user-friendly, ensuring that healthcare professionals can easily use these tools in their daily work. A focus on making the system intuitive will help it fit smoothly into clinical routines and increase its effectiveness.

Beyond technical improvements, future research must address ethical concerns, such as patient privacy, data security, and the risk of bias in the predictive algorithms. Understanding these issues is essential for building trust in CDSSs and ensuring they are used responsibly and fairly in healthcare. Finally, because healthcare challenges are interdisciplinary, future research should encourage collaboration between data scientists, clinicians, and policymakers. This teamwork is crucial to successfully integrating CDSSs into different healthcare settings and addressing the unique challenges of each environment.

In summary, these research directions will help continue to advance CDSS technology for predicting diabetes complications. Through ongoing innovation, this work has the potential to improve both the field of healthcare and the quality of care for people living with diabetes.

References

- Abdar, M., Zomorodi-Moghadam, M., Das, R., & Ting, I.-H. (2017). Performance analysis of classification algorithms on early detection of liver disease. *Expert systems with applications*, 67, 239-251.
- ACCORD, A. t. C. C. R. i. D. S. G. (2008). Effects of intensive glucose lowering in type 2 diabetes. *New England Journal of Medicine*, 358(24), 2545-2559.
- ADA, A. D. A. (2020). *Our 60-Second Type 2 Diabetes Risk Test*. <https://www.diabetes.org/risk-test>
- Adjei Boakye, E., Varble, A., Rojek, R., Peavler, O., Trainer, A. K., Osazuwa-Peters, N., & Hinyard, L. (2018). Sociodemographic factors associated with engagement in diabetes self-management education among people with diabetes in the United States. *Public Health Reports*, 133(6), 685-691.
- Aekplakorn, W., Bunnag, P., Woodward, M., Sritara, P., Cheepudomwit, S., Yamwong, S., Yipintsoi, T., & Rajatanavin, R. (2006). A risk score for predicting incident diabetes in the Thai population. *Diabetes care*, 29(8), 1872-1877.
- Afrash, M. R., Rahimi, F., Kazemi-Arpanahi, H., Shanbezadeh, M., Amraei, M., & Asadi, F. (2022, 2022/01/01/). Development of an intelligent clinical decision support system for the early prediction of diabetic nephropathy. *Informatics in Medicine Unlocked*, 35, 101135. <https://doi.org/https://doi.org/10.1016/j.imu.2022.101135>
- Alabdallah, A., Ohlsson, M., Pashami, S., & Rögngvaldsson, T. (2022). The Concordance Index decomposition--A measure for a deeper understanding of survival prediction models. *arXiv preprint arXiv:2203.00144*.
- Alasadi, S. A., & Bhaya, W. S. (2017). Review of data preprocessing techniques in data mining. *Journal of Engineering and Applied Sciences*, 12(16), 4102-4107.
- Alder, H., Michel, B. A., Marx, C., Tamborrini, G., Langenegger, T., Bruehlmann, P., Steurer, J., & Wildi, L. M. (2014). Computer-based diagnostic expert systems in rheumatology: where do we stand in 2014? *International journal of rheumatology*, 2014.

Alpay, L., van der Boog, P., & Dumaij, A. (2011). An empowerment-based approach to developing innovative e-health tools for self-management. *Health informatics journal*, 17(4), 247-255.

Alther, M., & Reddy, C., K. (2015). *Clinical Decision Support Systems*.

Aminian, A., Zajichek, A., Arterburn, D. E., Wolski, K. E., Brethauer, S. A., Schauer, P. R., Nissen, S. E., & Kattan, M. W. (2020, Apr). Predicting 10-Year Risk of End-Organ Complications of Type 2 Diabetes With and Without Metabolic Surgery: A Machine Learning Approach. *Diabetes care*, 43(4), 852-859. <https://doi.org/10.2337/dc19-2057>

Amir Talaei-Khoei, & Wilson, J. M. (2018). Identifying people at risk of developing type 2 diabetes: A comparison of predictive analytics techniques and predictor variables. *International Journal of Medical Informatics*, 119, 22-38.

Amoo, T., Green, B. O., & Raghupathi, V. (2014). The awareness of diabetes, its complications, and preventive measures in a developing country. *International Journal of Healthcare Management*, 7(4), 247-250.

Araz, O. M., Olson, D., & Ramirez-Nafarrate, A. (2019, 2019/02/01/). Predictive analytics for hospital admissions from the emergency department using triage information. *International Journal of Production Economics*, 208, 199-207. <https://doi.org/https://doi.org/10.1016/j.ijpe.2018.11.024>

Archer, L. B. (1964). Systematic method for designers. *Design*, 56-59.

Arianna Dagliati, Lucia Sacchi, Valentina Tibollo, & Chiovato, L. (2018). Machine Learning Methods to Predict Diabetes Complications. *Journal of Diabetes Science and Technology*, 12(2), 295-302.

Atlantis, E., Joshy, G., Williams, M., & Simmons, D. (2017). Diabetes Among Māori and Other Ethnic Groups in New Zealand. In S. Dagogo-Jack (Ed.), *Diabetes Mellitus in Developing Countries and Underserved Communities* (pp. 165-190). Springer International Publishing. https://doi.org/10.1007/978-3-319-41559-8_10

Avison, D., & Fitzgerald, G. (2003). *Information systems development: methodologies, techniques and tools*. McGraw-Hill.

- Bachiochi, P. D., & Weiner, S. P. (2004). Qualitative data collection and analysis. *Handbook of research methods in industrial and organizational psychology*, 161-183.
- Bailey, B. (2011). Case studies: A security science research methodology.
- Bannister, F., & Remenyi, D. (2000). Acts of faith: instinct, value and IT investment decisions. *Journal of information Technology*, 15(3), 231-241.
- Bebu, I., Keshavarzi, S., Gao, X., Braffett, B. H., Canty, A. J., Herman, W. H., Orchard, T. J., Dagogo-Jack, S., Nathan, D. M., & Lachin, J. M. (2021). Genetic risk factors for CVD in type 1 diabetes: the DCCT/EDIC study. *Diabetes care*, 44(6), 1309-1316.
- Behrens, J. T. (1997). Principles and procedures of exploratory data analysis. *Psychological methods*, 2(2), 131.
- Bernardini, M., Romeo, L., Mancini, A., & Frontoni, E. (2021). A Clinical Decision Support System to Stratify the Temporal Risk of Diabetic Retinopathy. *Ieee Access*, 9, 151864-151872.
- Berner, E. S., Webster, G. D., Shugerman, A. A., Jackson, J. R., Algina, J., Baker, A. L., Ball, E. V., Cobbs, C. G., Dennis, V. W., & Frenkel, E. P. (1994). Performance of four computer-based diagnostic systems. *New England Journal of Medicine*, 330(25), 1792-1796.
- Bland, J. M., & Altman, D. G. (2004). The logrank test. *bmj*, 328(7447), 1073.
- Bourne, R. R., Stevens, G. A., White, R. A., Smith, J. L., Flaxman, S. R., Price, H., Jonas, J. B., Keeffe, J., Leasher, J., & Naidoo, K. (2013). Causes of vision loss worldwide, 1990–2010: a systematic analysis. *The lancet global health*, 1(6), e339-e349.
- Bozyel, S., Şimşek, E., Koçyiğit, D., Güler, A., Korkmaz, Y., Şeker, M., Ertürk, M., & Keser, N. (2024). Artificial intelligence-based clinical decision support systems in cardiovascular diseases. *Anatolian Journal of Cardiology*, 28(2), 74.
- Brooke, J. (1996). SUS -- a quick and dirty usability scale. In (pp. 189-194).
- Burford, S. J., Park, S., & Dawda, P. (2019). Small Data and Its Visualization for Diabetes Self-Management: Qualitative Study. *JMIR DIABETES*, 4(3), e10324.

- Campbell, S., Greenwood, M., Prior, S., Shearer, T., Walkem, K., Young, S., Bywaters, D., & Walker, K. (2020, Dec). Purposive sampling: complex or simple? Research case examples. *J Res Nurs*, 25(8), 652-661. <https://doi.org/10.1177/1744987120927206>
- Cao, J., Crews, J. M., Nunamaker, J. F., Burgoon, J. K., & Lin, M. (2004). User experience with Agent99 Trainer: A usability study. 37th Annual Hawaii International Conference on System Sciences, 2004. Proceedings of the,
- Carter, R. E., Lackland, D. T., Cleary, P. A., Yim, E., Lopes-Virella, M. F., Gilbert, G. E., & Orchard, T. J. (2007). Intensive treatment of diabetes is associated with a reduced rate of peripheral arterial calcification in the diabetes control and complications trial. *Diabetes care*, 30(10), 2646-2648.
- Cederholm, J., Eeg-Olofsson, K., Eliasson, B., Zethelius, B., Nilsson, P. M., & Gudbjörnsdottir, S. (2008). Risk prediction of cardiovascular disease in type 2 diabetes: a risk equation from the Swedish National Diabetes Register. *Diabetes care*, 31(10), 2038-2043.
- Chan, L., Nadkarni, G. N., Fleming, F., McCullough, J. R., Connolly, P., Mosoyan, G., El Salem, F., Kattan, M. W., Vassalotti, J. A., & Murphy, B. (2021). Derivation and validation of a machine learning risk score using biomarker and electronic patient data to predict progression of diabetic kidney disease. *Diabetologia*, 64, 1504-1515.
- Chatfield, C. (1986, 1986/01/01/). Exploratory data analysis. *European Journal of Operational Research*, 23(1), 5-13. [https://doi.org/https://doi.org/10.1016/0377-2217\(86\)90209-2](https://doi.org/https://doi.org/10.1016/0377-2217(86)90209-2)
- Chawla, A., Chawla, R., & Jaggi, S. (2016, Jul-Aug). Microvascular and macrovascular complications in diabetes mellitus: Distinct or continuum? *Indian J Endocrinol Metab*, 20(4), 546-551. <https://doi.org/10.4103/2230-8210.183480>
- Chen, H., Hailey, D., Wang, N., & Yu, P. (2014). A review of data quality assessment methods for public health information systems. *International journal of environmental research and public health*, 11(5), 5170-5207.
- Chen, H.-L., Yang, B., Wang, G., Liu, J., Chen, Y.-D., & Liu, D.-Y. (2012). A three-stage expert system based on support vector machines for thyroid disease diagnosis. *Journal of medical systems*, 36(3), 1953-1963.

- Chen, L., Magliano, D. J., Balkau, B., Colagiuri, S., Zimmet, P. Z., Tonkin, A. M., Mitchell, P., Phillips, P. J., & Shaw, J. E. (2010, Feb 15). AUSDRISK: an Australian Type 2 Diabetes Risk Assessment Tool based on demographic, lifestyle and simple anthropometric measures. *Med J Aust*, 192(4), 197-202. <https://doi.org/10.5694/j.1326-5377.2010.tb03507.x>
- Chen Lei, Magliano Dianna J, Balkau Beverley, Colagiuri Stephen, Zimmet Paul Z, Tonkin Andrew M, Mitchell Paul, Phillips Patrick J, & Shaw Jonathan E. (2010). AUSDRISK: an Australian Type 2 Diabetes Risk Assessment Tool based on demographic, lifestyle and simple anthropometric measures. *Medical Journal of Australia*, 192(4), 197-202.
- Chen, Y., Jiang, L., Gao, B., Cheng, Z. Y., Jin, J., & Yang, K. H. (2016, Jun). Survival and disease-free benefits with mastectomy versus breast conservation therapy for early breast cancer: a meta-analysis. *Breast Cancer Res Treat*, 157(3), 517-525. <https://doi.org/10.1007/s10549-016-3830-z>
- Chien, K., Cai, T., Hsu, H., Su, T., Chang, W., Chen, M., Lee, Y., & Hu, F. B. (2009, 2009/03/01). A prediction model for type 2 diabetes risk among Chinese people. *Diabetologia*, 52(3), 443-450. <https://doi.org/10.1007/s00125-008-1232-4>
- Chiew, T. K., & Salim, S. S. (2003). Webuse: Website usability evaluation tool. *Malaysian Journal of Computer Science*, 16(1), 47-57.
- Clair, L., Anderson, H., Anderson, C., Ekuma, O., & Prior, H. J. (2022, Jun). Cardiovascular disease and the risk of dementia: a survival analysis using administrative data from Manitoba. *Can J Public Health*, 113(3), 455-464. <https://doi.org/10.17269/s41997-021-00589-2>
- Cleven, A., Gubler, P., & Hüner, K. M. (2009). Design alternatives for the evaluation of design science research artifacts. Proceedings of the 4th International Conference on Design Science Research in Information Systems and Technology,
- Coca, S. G., Nadkarni, G. N., Huang, Y., Moledina, D. G., Rao, V., Zhang, J., Ferket, B., Crowley, S. T., Fried, L. F., & Parikh, C. R. (2017). Plasma biomarkers and kidney function decline in early and established diabetic kidney disease. *Journal of the American Society of Nephrology*, 28(9), 2786-2793.
- Coiera, E. (2015). *Guide to health informatics*. CRC press.

Constantinou, A. C., Fenton, N., Marsh, W., & Radlinski, L. (2016). From complex questionnaire and interviewing data to intelligent Bayesian network models for medical decision support. *Artificial Intelligence in Medicine*, 67, 75-93.

Coslovsky, M., Takala, J., Exadaktylos, A. K., Martinolli, L., & Merz, T. M. (2015, 2015/06/01). A clinical prediction model to identify patients at high risk of death in the emergency department. *Intensive Care Medicine*, 41(6), 1029-1036.
<https://doi.org/10.1007/s00134-015-3737-x>

Costa, W., Figueiredo, L., & Alves, E. (2019). Application of an Artificial Neural Network for Heart Disease Diagnosis. XXVI Brazilian Congress on Biomedical Engineering,

Creswell, J. W. (2009). Research designs: Qualitative, quantitative, and mixed methods approaches. *Callifornia: Sage*.

Creswell, J. W., & Creswell, J. (2003). *Research design*. Sage publications Thousand Oaks, CA.

Creswell, J. W., & Miller, G. A. (1997). Research Methodologies and the Doctoral Process. *New directions for higher education*, 99, 33-46.

DA, D. A. (2010). *Welcome to the diabetes risk calculator*.
<https://www.diabetesaustralia.com.au/risk-calculator>

Dash, S., Shakyawar, S. K., Sharma, M., & Kaushik, S. (2019). Big data in healthcare: management, analysis and future prospects. *Journal of Big Data*, 6(1), 54.

David, G. K., & Mitchel, K. (2012). *Survival analysis: a Self-Learning text*. Springer.

DCCT/EDIC, & Braffett, B. (2016). Coprogression of cardiovascular risk factors in type 1 diabetes during 30 years of follow-up in the DCCT/EDIC study. *Diabetes care*, 39(9), 1621-1630.

DCCT/EDIC Research Group. (2005). Intensive diabetes treatment and cardiovascular disease in patients with type 1 diabetes. *New England Journal of Medicine*, 353(25), 2643-2653.

Deal, J. A., Sharrett, A. R., Rawlings, A. M., Gottesman, R. F., Bandeen-Roche, K., Albert, M., Knopman, D., Selvin, E., Wasserman, B. A., Klein, B., & Klein, R. (2018). Retinal signs and

- 20-year cognitive decline in the Atherosclerosis Risk in Communities Study. *Neurology*, 90(13), e1158-e1166. <https://doi.org/10.1212/WNL.0000000000005205>
- Debray, T. P., Moons, K. G., Ahmed, I., Koffijberg, H., & Riley, R. D. (2013). A framework for developing, implementing, and evaluating clinical prediction models in an individual participant data meta-analysis. *Statistics in medicine*, 32(18), 3158-3180.
- Derincek, A., Guler, U. O., Uysal, M., & Ozalay, M. (2020, Mar). Spinal Metastatic Disease: Survival Analysis of 146 Patients and Evaluation of 4 Different Preoperative Scoring Systems. *Clin Spine Surg*, 33(2), E81-e86. <https://doi.org/10.1097/bsd.0000000000000858>
- Deshpande, A. D., Harris-Hayes, M., & Schootman, M. (2008). Epidemiology of diabetes and diabetes-related complications. *Physical therapy*, 88(11), 1254-1264.
- Diabetes UK DUK, University of Leicester UoL, & NHSTrust, U. H. o. L. *Diabetes UK* <https://riskscore.diabetes.org.uk/start>
- Dillon, A. (2001). The evaluation of software usability. *Encycl Hum Factors Ergon*.
- Dogba, M. J., Dipankui, M. T., Chipenda Dansokho, S., Légaré, F., & Witteman, H. O. (2018). Diabetes-related complications: Which research topics matter to diverse patients and caregivers? *Health Expectations*, 21(2), 549-559.
- Dong, G., & Liu, H. (2018). *Feature engineering for machine learning and data analytics*. CRC Press.
- Druckman, D. (2005). *Doing Research: Methods of Inquiry for Conflict Analysis*. SAGE.
- Duffy, F. D., Lynn, L. A., Didura, H., Hess, B., Caverzagie, K., Grosso, L., Lipner, R. A., & Holmboe, E. S. (2008). Self-assessment of practice performance: Development of the ABIM Practice Improvement Module (PIMSM). *Journal of Continuing Education in the Health Professions*, 28(1), 38-46.
- Eekels, J., & Roozenburg, N. F. (1991). A methodological comparison of the structures of scientific research and engineering design: their similarities and differences. *Design studies*, 12(4), 197-203.

- EL-firjani, N. F., Elberkawi, E. K., & Maatuk, A. M. (2017). METHOD FOR WEBSITE USABILITY EVALUATION: A COMPARATIVE ANALYSIS.
- El-Halees, A. M. (2014). Software Usability Evaluation Using Opinion Mining. *J. Softw.*, 9(2), 343-349.
- Elhassen, S. S. M. (2017). *Characterization of Multiple Sclerosis on the Brain Magnetic Resonance Images Using Texture Analysis* Sudan University of Science and Technology].
- Elhoseny, M., Shankar, K., & Uthayakumar, J. (2019). Intelligent Diagnostic Prediction and Classification System for Chronic Kidney Disease. *Scientific reports*, 9(1), 1-14.
- Elley, C. R., Robinson, E., Kenealy, T., Bramley, D., & Drury, P. L. (2010). Derivation and validation of a new cardiovascular risk score for people with type 2 diabetes: the New Zealand diabetes cohort study. *Diabetes care*, 33(6), 1347-1352.
- Ellis, T. J., & Levy, Y. (2010). A guide for novice researchers: Design and development research methods. Proceedings of Informing Science & IT Education Conference (InSITE),
- Emerging Risk Factors Collaboration. (2010). Diabetes mellitus, fasting blood glucose concentration, and risk of vascular disease: a collaborative meta-analysis of 102 prospective studies. *The Lancet*, 375(9733), 2215-2222.
- Emmert-Streib, F., & Dehmer, M. (2019). Introduction to survival analysis in practice. *Machine Learning and Knowledge Extraction*, 1(3), 1013-1038.
- Erandathi, M., Wang, W. Y. C., Mayo, M., & Shafuii, I. (2022). Prevalence of sociodemographic factors in a cohort of diabetes mellitus: a retrospective study. Proceedings of the 6th International Conference on Medical and Health Informatics,
- Exalto, L. G., Biessels, G. J., Karter, A. J., Huang, E. S., Katon, W. J., Minkoff, J. R., & Whitmer, R. A. (2013). Risk score for prediction of 10 year dementia risk in individuals with type 2 diabetes: a cohort study. *The Lancet Diabetes & Endocrinology*, 1(3), 183-190.
- Famili, A., Shen, W.-M., Weber, R., & Simoudis, E. (1997). Data preprocessing and intelligent data analysis. *Intelligent data analysis*, 1(1), 3-23.

Fernandez, A., Abrahão, S., & Insfran, E. (2011, 2011//). A Web Usability Evaluation Process for Model-Driven Web Development. *Advanced Information Systems Engineering*, Berlin, Heidelberg.

Ferreira, J. C., & Patino, C. M. (2016). What is survival analysis, and when should I use it? *Jornal Brasileiro de Pneumologia*, 42, 77-77.

Fitzpatrick, R. (1998). Strategies for evaluating software usability. *Methods*, 353(1).

Forster, R. B., Garcia, E. S., Sluiman, A. J., Grecian, S. M., McLachlan, S., MacGillivray, T. J., Strachan, M. W., Price, J. F., & investigators, E. T. D. S. (2021). Retinal venular tortuosity and fractal dimension predict incident retinopathy in adults with type 2 diabetes: the Edinburgh Type 2 Diabetes Study. *Diabetologia*, 64, 1103-1112.

Fulcher, A., & Hills, P. (1996). Towards a strategic framework for design research. *Journal of Engineering Design*, 7(2), 183-193.

García, S., Luengo, J., & Herrera, F. (2015). *Data preprocessing in data mining*. Springer.

Gause, D. C., & Weinberg, G. M. (1989). *Exploring requirements: quality before design* (Vol. 7). Dorset House New York.

George, B., Seals, S., & Aban, I. (2014, 2014/08/01). Survival analysis and regression models. *Journal of Nuclear Cardiology*, 21(4), 686-694. <https://doi.org/10.1007/s12350-014-9908-2>

Gholamzadeh, M., Abtahi, H., & Safdari, R. (2023). The Application of Knowledge-Based Clinical Decision Support Systems to Enhance Adherence to Evidence-Based Medicine in Chronic Disease. *Journal of Healthcare Engineering*, 2023.

Ginige, A., & Murugesan, S. (2001). Web engineering: an introduction. *IEEE MultiMedia*, 8(1), 14-18. <https://doi.org/10.1109/93.923949>

Glasheen, W. P., Renda, A., & Dong, Y. (2017). Diabetes complications severity index (DCSI)—update and ICD-10 translation. *Journal of Diabetes and its Complications*, 31(6), 1007-1013.

- Gliklich, R. E., Leavy, M. B., & Dreyer, N. A. (2019). Tools and Technologies for Registry Interoperability, Registries for Evaluating Patient Outcomes: A User's Guide, Addendum 2 [Internet].
- Gneiting, T., & Raftery, A. E. (2007a). Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical association*, 359-378.
- Gneiting, T., & Raftery, A. E. (2007b). Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical association*, 102(477), 359-378.
- Goel, M. K., Khanna, P., & Kishore, J. (2010). Understanding survival analysis: Kaplan-Meier estimate. *International journal of Ayurveda research*, 1(4), 274.
- Goertzen, M. J. (2017). Introduction to quantitative research and data. *Library Technology Reports*, 53(4), 12-18.
- Gomez-Arbelaez, D., Alvarado-Jurado, L., Ayala-Castillo, M., Forero-Naranjo, L., Camacho, P. A., & Lopez-Jaramillo, P. (2015). Evaluation of the Finnish Diabetes Risk Score to predict type 2 diabetes mellitus in a Colombian population: A longitudinal observational study. *World journal of diabetes*, 6(17), 1337.
- Gopalakrishnan, S., & Ganeshkumar, P. (2013). Systematic Reviews and Meta-analysis: Understanding the Best Evidence in Primary Healthcare. *Journal of family medicine and primary care*, 2(1), 9-14. <https://doi.org/10.4103/2249-4863.109934>
- Gopiseti, L. D., Kummera, S. K. L., Pattamsetti, S. R., Kuna, S., Parsi, N., & Kodali, H. P. (2023). Multiple Disease Prediction System using Machine Learning and Streamlit. 2023 5th International Conference on Smart Systems and Inventive Technology (ICSSIT),
- Govinda, K., Singla, K., & Jain, K. (2017). Fuzzy based uncertainty modeling of Cancer Diagnosis System. 2017 International Conference on Intelligent Sustainable Systems (ICISS),
- Greenes, R. (2014). *Clinical decision support: the road to broad adoption*. Academic Press.
- Gregg, D. G., Kulkarni, U. R., & Vinzé, A. S. (2001). Understanding the philosophical underpinnings of software engineering research in information systems. *Information systems frontiers*, 3(2), 169-183.

- Gregor, S., & Hevner, A. R. (2013). Positioning and presenting design science research for maximum impact. *MIS quarterly*, 337-355.
- Gregorio, J., Reis, L., Peyroteo, M., Maia, M., da Silva, M. M., & Lapão, L. V. (2021). The role of Design Science Research Methodology in developing pharmacy eHealth services. *Research in Social and Administrative Pharmacy*, 17(12), 2089-2096.
- Griffin, S. J., Little, P. S., Hales, C. N., Kinmonth, A. L., & Wareham, N. J. (2000). Diabetes risk score: towards earlier detection of Type 2 diabetes in general practice. *Diabetes/Metabolism Research and Reviews*, 16(3), 164-171.
[https://doi.org/10.1002/1520-7560\(200005/06\)16:3<164::aid-dmrr103>3.0.co;2-r](https://doi.org/10.1002/1520-7560(200005/06)16:3<164::aid-dmrr103>3.0.co;2-r)
- Gummesson, E. (2000). *Qualitative methods in management research*. Sage.
- Gurumoorthy, S., Muppalaneni, N. B., & Gao, X.-Z. (2018). Classification and Analysis of EEG Using SVM and MRE. In *Computational Intelligence Techniques in Diagnosis of Brain Diseases* (pp. 33-46). Springer.
- Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar), 1157-1182.
- Hajihosseini, M., Kazemi, T., & Faradmal, J. (2016, Jul). Multistate Models for Survival Analysis of Cardiovascular Disease Process. *Rev Esp Cardiol (Engl Ed)*, 69(7), 714-715.
<https://doi.org/10.1016/j.rec.2016.04.009>
- Hall, M. A. (1999). *Correlation-based feature selection for machine learning* [The University of Waikato].
- Hamedan, F., Orooji, A., Sanadgol, H., & Sheikhtaheri, A. (2020, 2020/06/01/). Clinical decision support system to predict chronic kidney disease: A fuzzy expert system approach. *International Journal of Medical Informatics*, 138, 104134.
<https://doi.org/https://doi.org/10.1016/j.ijmedinf.2020.104134>
- Han, K., Song, K., & Choi, B. W. (2016). How to develop, validate, and compare clinical prediction models involving radiological parameters: study design and statistical methods. *Korean journal of radiology*, 17(3), 339-350.

- Hartson, H. R., Andre, T. S., & Williges, R. C. (2001). Criteria for evaluating usability evaluation methods. *International journal of human-computer interaction*, 13(4), 373-410.
- Hasan, S., & Padman, R. (2006). Analyzing the effect of data quality on the accuracy of clinical decision support systems: a computer simulation approach. *AMIA Annu Symp Proc*, 2006, 324-328.
- Haux, R. (2006). Health information systems—past, present, future. *International Journal of Medical Informatics*, 75(3-4), 268-281.
- Haux, R., Ammenwerth, E., Winter, A., & Brigl, B. (2004). *Strategic information management in hospitals: an introduction to hospital information systems*. Springer Science & Business Media.
- Häyrinen, K., Saranto, K., & Nykänen, P. (2008, May). Definition, structure, content, use and impacts of electronic health records: a review of the research literature. *Int J Med Inform*, 77(5), 291-304. <https://doi.org/10.1016/j.ijmedinf.2007.09.001>
- Health New Zealand. (2023). *Virtual Diabetes Register web tool*. <https://tewhatuora.shinyapps.io/virtual-diabetes-register-web-tool/>
- Heller, G. (2021). The added value of new covariates to the brier score in cox survival models. *Lifetime data analysis*, 27(1), 1-14.
- Herman, W. H., Smith, P. J., Thompson, T. J., Engelgau, M. M., & Aubert, R. E. (1995). A new and simple questionnaire to identify people at increased risk for undiagnosed diabetes. *Diabetes Care*, 18(3), 382-387.
- Hevner, A., & Chatterjee, S. (2010). Design science research in information systems. In *Design research in information systems* (pp. 9-22). Springer.
- Hevner, A., Prat, N., Comyn-Wattiau, I., & Akoka, J. (2018). A pragmatic approach for identifying and managing design science research goals and evaluation criteria. AIS SIGPrag Pre-ICIS workshop on "Practice-based Design and Innovation of Digital Artifacts",
- Hevner, A. R. (2007). A three cycle view of design science research. *Scandinavian journal of information systems*, 19(2), 4.

- Hevner, A. R., March, S. T., Park, J., & Ram, S. (2004). Design science in information systems research. *MIS quarterly*, 75-105.
- Hewett, T. T., Baecker, R., Card, S., Carey, T., Gasen, J., Mantei, M., Perlman, G., Strong, G., & Verplank, W. (1992). *ACM SIGCHI curricula for human-computer interaction*. ACM.
- Hill, S., Mullins, P., Murphy, R., Schmiedel, O., Vaghefi, E., Ramke, J., & Squirrell, D. (2021). Risk Factors for Progression to Referable Diabetic Eye Disease in People With Diabetes Mellitus in Auckland, New Zealand: A 12-Year Retrospective Cohort Analysis. *The Asia-Pacific Journal of Ophthalmology*, 10(6), 579-589.
<https://doi.org/10.1097/apo.0000000000000464>
- Hinz, E., Borland, D., Shah, H., West, V. L., & Hammond, W. E. (2014). Temporal visualization of diabetes mellitus via hemoglobin a1c levels. Proceedings of the 2014 Workshop on Visual Analytics in Healthcare,
- Hippisley-Cox, J., & Coupland, C. (2017). Development and validation of QDiabetes-2018 risk prediction algorithm to estimate future risk of type 2 diabetes: cohort study. *bmj*, 359, j5019.
- Hippisley-Cox, J., Coupland, C., & Brindle, P. (2017). Development and validation of QRISK3 risk prediction algorithms to estimate future risk of cardiovascular disease: prospective cohort study. *bmj*, 357, j2099.
- Hippisley-Cox, J., Coupland, C., Robson, J., Sheikh, A., & Brindle, P. (2009). Predicting risk of type 2 diabetes in England and Wales: prospective derivation and validation of QDScore. *bmj*, 338, b880.
- Hippisley-Cox, J., Coupland, C., Vinogradova, Y., Robson, J., May, M., & Brindle, P. (2007). Derivation and validation of QRISK, a new cardiovascular disease risk score for the United Kingdom: prospective open cohort study. *bmj*, 335(7611), 136.
- Hippisley-Cox, J., Coupland, C., Vinogradova, Y., Robson, J., Minhas, R., Sheikh, A., & Brindle, P. (2008). Predicting cardiovascular risk in England and Wales: prospective derivation and validation of QRISK2. *bmj*, 336(7659), 1475-1482.

- Hirschheim, R., Klein, H. K., & Lyytinen, K. (1995). *Information systems development and data modeling: conceptual and philosophical foundations* (Vol. 9). Cambridge University Press.
- Holden, M. T., & Lynch, P. (2004). Choosing the appropriate methodology: Understanding research philosophy. *The marketing review*, 4(4), 397-409.
- Holman, R. R., Paul, S. K., Bethel, M. A., Matthews, D. R., & Neil, H. A. W. (2008). 10-year follow-up of intensive glucose control in type 2 diabetes. *New England Journal of Medicine*, 359(15), 1577-1589.
- Hong, W. S., Haimovich, A. D., & Taylor, R. A. (2018). Predicting hospital admission at emergency department triage using machine learning. *PloS one*, 13(7), e0201016.
- Horsky, J., Phansalkar, S., Desai, A., Bell, D., & Middleton, B. (2013). Design of decision support interventions for medication prescribing. *International Journal of Medical Informatics*, 82(6), 492-503.
- Hu, F. B., Manson, J. E., Stampfer, M. J., Colditz, G., Liu, S., Solomon, C. G., & Willett, W. C. (2001). Diet, lifestyle, and the risk of type 2 diabetes mellitus in women. *New England Journal of Medicine*, 345(11), 790-797.
- Hussain, T., & Nguyen, Q. T. (2014). Molecular imaging for cancer diagnosis and surgery. *Advanced drug delivery reviews*, 66, 90-100.
- IDF, I. D. F. (2021). *IDF Diabetes Atlas*.
- Idris, I. (2016). *Python data analysis cookbook*. Packt Publishing Ltd.
- Imhoff, M. (2002). Health informatics. In *Evaluating Critical Care* (pp. 255-269). Springer.
- International Diabetes Federation, I. (2013). *IDF Diabetes Atlas*.
- International Diabetes Federation, I. (2020). *COVID-19 and diabetes*.
<https://www.idf.org/aboutdiabetes/what-is-diabetes/covid-19-and-diabetes.html>

- International Organization for Standardization. (2011). *Systems and software Quality Requirements and Evaluation (SQuaRE) – System and software quality models*.
- Irani, Z., & Love, P. E. (2002). Developing a frame of reference for ex-ante IT/IS investment evaluation. *European Journal of Information Systems, 11*, 74-82.
- ISO/IEC 25010, I. S. M. S. (2011). *ISO/IEC 25010:2011, Systems and software engineering — Systems and software Quality Requirements and Evaluation (SQuaRE) — System and software quality models*.
- Jager, K. J., Van Dijk, P. C., Zoccali, C., & Dekker, F. W. (2008). The analysis of survival data: the Kaplan–Meier method. *Kidney international, 74*(5), 560-565.
- Jakovljević, B. (2008). Health Information SystemHealth information system. In W. Kirch (Ed.), *Encyclopedia of Public Health* (pp. 603-607). Springer Netherlands.
https://doi.org/10.1007/978-1-4020-5614-7_1425
- Jardim, S. V. (2013). The electronic health record and its contribution to healthcare information systems interoperability. *Procedia technology, 9*, 940-948.
- Jenkins, S. P. (2005). Survival analysis. *Unpublished manuscript, Institute for Social and Economic Research, University of Essex, Colchester, UK, 42*, 54-56.
- Jia, X., Baig, M. M., Mirza, F., & GholamHosseini, H. (2019, 2019/04/09). A Cox-Based Risk Prediction Model for Early Detection of Cardiovascular Disease: Identification of Key Risk Factors for the Development of a 10-Year CVD Risk Prediction. *Advances in preventive medicine, 2019*, 8392348. <https://doi.org/10.1155/2019/8392348>
- Jiawei Han, Micheline Kamber, & Pei, J. (2011). *Data Mining: Concepts and Techniques* (3rd edition ed.).
- John Gantz, & Reinsel, D. (2012). THE DIGITAL UNIVERSE IN 2020: Big Data, Bigger Digital Shadow s, and Biggest Growth in the Far East.
- Johnson, B., & Turner, L. A. (2003). Data collection strategies in mixed methods research. *Handbook of mixed methods in social and behavioral research, 10*(2), 297-319.

- Joseph, J. J., Deedwania, P., Acharya, T., Aguilar, D., Bhatt, D. L., Chyun, D. A., Di Palo, K. E., Golden, S. H., Sperling, L. S., Lifestyle, A. H. A. D. C. o. t. C. o., Health, C., Council on Arteriosclerosis, T., Biology, V., Cardiology, C. o. C., & Hypertension, C. o. (2022). Comprehensive management of cardiovascular risk factors for adults with type 2 diabetes: a scientific statement from the American Heart Association. *Circulation*, *145*(9), e722-e759.
- Kadi, I., Idri, A., & Ouhbi, S. (2016). Quality evaluation of cardiac decision support systems using ISO 25010 standard. 2016 IEEE/ACS 13th International Conference of Computer Systems and Applications (AICCSA),
- Kamińska, D., Zwoliński, G., & Laska-Leśniewicz, A. (2022). Usability Testing of Virtual Reality Applications—The Pilot Study. *Sensors*, *22*(4), 1342.
- Kanaya, A. M., Fyr, C. L. W., De Rekeneire, N., Shorr, R. I., Schwartz, A. V., Goodpaster, B. H., Newman, A. B., Harris, T., & Barrett-Connor, E. (2005). Predicting the development of diabetes in older adults: the derivation and validation of a prediction rule. *Diabetes care*, *28*(2), 404-408.
- Kaplan, E. L., & Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American statistical association*, *53*(282), 457-481.
- Keast, R. L. (2004). Participatory Evaluation: A Missing Component in the Social Change Equation for Public Services. Social Change in the 21st Century Conference; Centre for Social Change Research; Queensland University of Technology,
- Keleş, A., Keleş, A., & Yavuz, U. (2011). Expert system based on neuro-fuzzy rules for diagnosis breast cancer. *Expert systems with applications*, *38*(5), 5719-5726.
- Kemppinen, J., Korpela, J., Elfvingren, K., & Polkko, J. (2017). Improving the productivity and efficiency of an integrated mental and addiction care—an application of the theory of constraints and five-focusing step to evaluation of adult ADHD patients: An application of the theory of constraints and the five-focusing step to evaluation of the adult ADHD patients. *Finnish Journal of eHealth and eWelfare*, *9*(1), 18-30.
- Kenealy, T., Arroll, B., & Petrie, K. J. (2005, Oct). Patients and computers as reminders to screen for diabetes in family practice. Randomized-controlled trial. *J Gen Intern Med*, *20*(10), 916-921. <https://doi.org/10.1111/j.1525-1497.2005.0197.x>

- Kilsdonk, E., Peute, L., Riezebos, R. J., Kremer, L. C., & Jaspers, M. W. (2016). Uncovering healthcare practitioners' information processing using the think-aloud method: From paper-based guideline to clinical decision support system. *International Journal of Medical Informatics*, 86, 10-19.
- Kim, E., Caraballo, P. J., Castro, M. R., Pieczkiewicz, D. S., & Simon, G. J. (2019). Towards more accessible precision medicine: building a more transferable machine learning model to support prognostic decisions for micro-and macrovascular complications of type 2 diabetes mellitus. *Journal of medical systems*, 43(7), 185.
- Kim, H. C., Lee, J. S., Lee, E. Y., Ha, Y. J., Chae, E. J., Han, M., Cross, G., Barnett, J., Joseph, J., & Song, J. W. (2020, Dec). Risk prediction model in rheumatoid arthritis-associated interstitial lung disease. *Respirology*, 25(12), 1257-1264.
<https://doi.org/10.1111/resp.13848>
- King, P., Peacock, I., & Donnelly, R. (1999, Nov). The UK prospective diabetes study (UKPDS): clinical and therapeutic implications for type 2 diabetes. *Br J Clin Pharmacol*, 48(5), 643-648. <https://doi.org/10.1046/j.1365-2125.1999.00092.x>
- Kleinbaum, D. G., & Klein, M. (1996). *Survival analysis a self-learning text*. Springer.
- Knight, E. P., & Shea, K. (2014). A patient-focused framework integrating self-management and informatics. *Journal of Nursing Scholarship*, 46(2), 91-97.
- Knowler, W. C. (2002). Diabetes Prevention Program Research Group: Reduction in the incidence of type 2 diabetes with life-style intervention or metformin. *N. Engl. J. Med.*, 346, 393-403.
- Koopman, R. J., Mainous, A. G., Diaz, V. A., & Geesey, M. E. (2005). Changes in age at diagnosis of type 2 diabetes mellitus in the United States, 1988 to 2000. *The Annals of Family Medicine*, 3(1), 60-63.
- Kothari, C. R. (2004). *Research methodology: Methods and techniques*. New Age International.
- Kotsiantis, S. (2011). Feature selection for machine learning classification problems: a recent overview. *Artificial Intelligence Review*, 42(1), 157-176.

- Kuechler, B., & Vaishnavi, V. (2008). On theory development in design science research: anatomy of a research project. *European Journal of Information Systems*, 17(5), 489-504.
- Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling* (Vol. 26). Springer.
- Kulkarni, M., Foraker, R. E., McNeill, A. M., Girman, C., Golden, S. H., Rosamond, W. D., Duncan, B., Schmidt, M. I., & Tuomilehto, J. (2017). Evaluation of the modified FINDRISC to identify individuals at high risk for diabetes among middle-aged white and black ARIC study participants. *Diabetes, Obesity and Metabolism*, 19(9), 1260-1266.
- Kunhimangalam, R., Ovallath, S., & Joseph, P. K. (2014). A clinical decision support system with an integrated EMR for diagnosis of peripheral neuropathy. *Journal of medical systems*, 38, 1-14.
- Kyrou, I., Tsigos, C., Mavrogianni, C., Cardon, G., Van Stappen, V., Latomme, J., Kivelä, J., Wikström, K., Tsochev, K., & Nanasi, A. (2020). Sociodemographic and lifestyle-related risk factors for identifying vulnerable groups for type 2 diabetes: a narrative review with emphasis on data from Europe. *BMC Endocrine Disorders*, 20, 1-13.
- Lagani, V., Chiarugi, F., Manousos, D., Verma, V., Fursse, J., Marias, K., & Tsamardinos, I. (2015). Realization of a service for the long-term risk assessment of diabetes-related complications. *Journal of Diabetes and its Complications*, 29(5), 691-698.
- Larose, C. D., & Larose, D. T. (2019). *Data science using Python and R*. John Wiley & Sons.
- Lauesen, S. (2002). *Software requirements: styles and techniques*. Pearson Education.
- Ledley, R. S., & Lusted, L. B. (1959). Reasoning foundations of medical diagnosis. *Science*, 130(3366), 9-21.
- Lee, Y., Chae, Y., & Jeon, S. (2010). Integration and evaluation of clinical decision support Systems for Diagnosis Idopathics Pulmonary Fibrosis (IPF). *Healthcare informatics research*, 16(4), 260-272.
- Lei Chen, Dianna J Magliano, Beverley Balkau, Stephen Colagiuri, Paul Z Zimmet, Andrew M Tonkin, Paul Mitchell, Patrick J Phillips, & Shaw, J. E. (2019). AUSDRISK: an Australian

- Type 2 Diabetes Risk Assessment Tool based on demographic, lifestyle and simple anthropometric measures. *MJA*, 192(4), 197-202.
- Levine, A., Sober, E., & Wright, E. O. (1987). Marxism and methodological individualism. *New Left Review*, 162, 67-84.
- Li, H., Han, D., Hou, Y., Chen, H., & Chen, Z. (2015). Statistical inference methods for two crossing survival curves: a comparison of methods. *PLoS one*, 10(1), e0116774.
- Li, Y., Lu, F., & Yin, Y. (2022, Jul 5). Applying logistic LASSO regression for the diagnosis of atypical Crohn's disease. *Sci Rep*, 12(1), 11340. <https://doi.org/10.1038/s41598-022-15609-5>
- Lindström, J., & Tuomilehto, J. (2003). The diabetes risk score: a practical tool to predict type 2 diabetes risk. *Diabetes care*, 26(3), 725-731.
- Liu, M., Ray, M., Zhang, D., Rundensteiner, E. A., Dougherty, D. J., Gupta, C., Wang, S., & Ari, I. (2012). Realtime healthcare services via nested complex event processing technology. Proceedings of the 15th International Conference on Extending Database Technology,
- Lorig, K. R., & Holman, H. R. (2003). Self-management education: history, definition, outcomes, and mechanisms. *Annals of behavioral medicine*, 26(1), 1-7.
- Lyu, K., Tian, Y., Shang, Y., Zhou, T., Yang, Z., Liu, Q., Yao, X., Zhang, P., Chen, J., & Li, J. (2023, 2023/03/01/). Causal knowledge graph construction and evaluation for clinical decision support of diabetic nephropathy. *Journal of Biomedical Informatics*, 139, 104298. <https://doi.org/https://doi.org/10.1016/j.jbi.2023.104298>
- Madurapperumage, A., Wang, W. Y. C., & Michael, M. (2021). A Systematic Review on Extracting Predictors for Forecasting Complications of Diabetes Mellitus Proceedings of the 5th International Conference on Medical and Health Informatics, Kyoto, Japan. <https://doi.org/10.1145/3472813.3473211>
- Mansoul, A., Atmani, B., & Benbelkacem, S. (2013). A hybrid decision support system: application on healthcare. *arXiv preprint arXiv:1311.4086*.
- Manzo, G., Pannatier, Y., Duflot, P., Kolh, P., Chavez, M., Bleret, V., Calvaresi, D., Jimenez-del-Toro, O., Schumacher, M., & Calbimonte, J.-P. (2023, 2023/04/01/). Breast cancer

- survival analysis agents for clinical decision support. *Computer methods and programs in biomedicine*, 231, 107373. <https://doi.org/https://doi.org/10.1016/j.cmpb.2023.107373>
- Maramba, I., Chatterjee, A., & Newman, C. (2019). Methods of usability testing in the development of eHealth applications: a scoping review. *International Journal of Medical Informatics*, 126, 95-104.
- March, S. T., & Smith, G. F. (1995). Design and natural science research on information technology. *Decision support systems*, 15(4), 251-266.
- March, S. T., & Storey, V. C. (2008). Design science in the information systems discipline: an introduction to the special issue on design science research. *MIS quarterly*, 725-730.
- Mariani, M. C., Tweneboah, O. K., & Bhuiyan, M. A. M. (2019). Supervised machine learning models applied to disease diagnosis and prognosis. *AIMS Public Health*, 6(4), 405.
- McKinney, W. (2012). *Python for data analysis: Data wrangling with Pandas, NumPy, and IPython*. " O'Reilly Media, Inc."
- McKinney, W. (2022). *Python for data analysis*. " O'Reilly Media, Inc."
- McKinsey & Company, I. (2011). *Big Data: The Next Frontier for Innovation, Competition & Productivity*.
- Mellergård, E., Johnsson, P., & Eek, F. (2020). Sociodemographic factors associated with HbA1c variability in type 2 diabetes: a prospective exploratory cohort study. *BMC Endocrine Disorders*, 20, 1-8.
- MIT, MIT Critical, Komorowski, M., Marshall, D. C., Saliccioli, J. D., & Crutain, Y. (2016). Exploratory data analysis. *Secondary analysis of electronic health records*, 185-203.
- MIT, D. M. C., Salgado, C. M., Azevedo, C., Proença, H., & Vieira, S. M. (2016). Missing data. *Secondary analysis of electronic health records*, 143-162.
- Mohammadzadeh, N., & Safdari, R. (2014). Patient monitoring in mobile health: opportunities and challenges. *Med Arch*, 68(1), 57-60. <https://doi.org/10.5455/medarch.2014.68.57-60>

- Moher, D., Liberati, A., Tetzlaff, J., & Altman, D. G. (2009, Aug 18). Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Ann Intern Med*, 151(4), 264-269, w264. <https://doi.org/10.7326/0003-4819-151-4-200908180-00135>
- Mombini, H., Tulu, B., Strong, D., Agu, E., Nguyen, H., Lindsay, C., Loretz, L., Pedersen, P., & Dunn, R. (2020). Design of a machine learning system for prediction of chronic wound management decisions. Designing for Digital Transformation. Co-Creating Services with Citizens and Industry: 15th International Conference on Design Science Research in Information Systems and Technology, DESRIST 2020, Kristiansand, Norway, December 2–4, 2020, Proceedings 15,
- Monteiro-Soares, M., Russell, D., Boyko, E. J., Jeffcoate, W., Mills, J. L., Morbach, S., Game, F., & Foot, I. W. G. o. t. D. (2020). Guidelines on the classification of diabetic foot ulcers (IWGDF 2019). *Diabetes/Metabolism Research and Reviews*, 36, e3273.
- Montgomery, A. A., Fahey, T., Peters, T. J., MacIntosh, C., & Sharp, D. J. (2000). Evaluation of computer based clinical decision support system and risk chart for management of hypertension in primary care: randomised controlled trial. *bmj*, 320(7236), 686-690.
- Moore, D. F. (2016). *Applied survival analysis using R* (Vol. 473). Springer.
- Mucha, H., Robert, S., Breitschwerdt, R., & Fellmann, M. (2023). Usability of clinical decision support systems. *Zeitschrift für Arbeitswissenschaft*, 77(1), 92-101.
- Navlani, A., Fandango, A., & Idris, I. (2021). *Python Data Analysis: Perform data collection, data processing, wrangling, visualization, and model building using Python*. Packt Publishing Ltd.
- Ng, C. J., Liao, P. J., Chang, Y. C., Kuan, J. T., Chen, J. C., & Hsu, K. H. (2016, Jun). Predictive factors for hospitalization of nonurgent patients in the emergency department. *Medicine*, 95(26), e4053. <https://doi.org/10.1097/md.0000000000004053>
- NIDDK, N. I. o. D. a. D. a. K. D. (2023). *National Institute of Diabetes and Digestive and Kidney Diseases*. Retrieved April 27 from <https://www.niddk.nih.gov/>
- Noble, D., Mathur, R., Dent, T., Meads, C., & Greenhalgh, T. (2011). Risk models and scores for type 2 diabetes: systematic review. *bmj*, 343, d7163.

- Nunamaker Jr, J. F., Chen, M., & Purdin, T. D. (1990). Systems development in information systems research. *Journal of management information systems*, 7(3), 89-106.
- NZ-MoH. (2022). *Diabetes*. Retrieved 24/08 from <https://www.health.govt.nz/search/results/diabetes>
- O'Neill, E. S., Dluhy, N. M., & Chin, E. (2005). Modelling novice clinical reasoning for a computerized decision support system. *Journal of Advanced Nursing*, 49(1), 68-77.
- Oates, B. J., Griffiths, M., & McLean, R. (2022). *Researching information systems and computing*. Sage.
- Oh, S. L., Hagiwara, Y., Raghavendra, U., Yuvaraj, R., Arunkumar, N., Murugappan, M., & Acharya, U. R. (2018). A deep learning approach for Parkinson's disease diagnosis from EEG signals. *Neural Computing and Applications*, 1-7.
- Olfatbakhsh, A., Heidari, L., Omid, Z., Hashemi, E. O., Ansari, M., Mozaffarian, S., & Haghghat, S. (2022, Sep 1). Long-term Survival and Prognostic Factors of Breast Cancer. *Arch Iran Med*, 25(9), 609-616. <https://doi.org/10.34172/aim.2022.96>
- Olivia, D., Nayak, A., & Balachandra, M. (2018, 2018//). Machine Learning Based Electronic Triage for Emergency Department. *Applications and Techniques in Information Security*, Singapore.
- Olson, C. A., Tooman, T. R., & Alvarado, C. J. (2010, Oct). Knowledge systems, health care teams, and clinical practice: a study of successful change. *Adv Health Sci Educ Theory Pract*, 15(4), 491-516. <https://doi.org/10.1007/s10459-009-9214-y>
- Oluwakemi, O., & Kamran, S. (2016, 12/28). Beyond simple charts: Design of visualizations for big health data. *Online journal of public health informatics*, 8(3). <https://doi.org/10.5210/ojphi.v8i3.7100>
- Ong, M. E. H., Lee Ng, C. H., Goh, K., Liu, N., Koh, Z. X., Shahidah, N., Zhang, T. T., Fook-Chong, S., & Lin, Z. (2012, 2012/06/21). Prediction of cardiac arrest in critically ill patients presenting to the emergency department using a machine learning score incorporating heart rate variability compared with the modified early warning score. *Critical Care*, 16(3), R108. <https://doi.org/10.1186/cc11396>

- Oomichi, T., Emoto, M., Tabata, T., Morioka, T., Tsujimoto, Y., Tahara, H., Shoji, T., & Nishizawa, Y. (2006). Impact of glycemic control on survival of diabetic patients on chronic regular hemodialysis: a 7-year observational study. *Diabetes care*, 29(7), 1496-1500.
- Orlikowski, W. J., & Baroudi, J. J. (1991). Studying information technology in organizations: Research approaches and assumptions. *Information systems research*, 2(1), 1-28.
- Osheroff, J. A., Teich, J. M., Levick, D., Saldana, L., Velasco, F. T., Sittig, D. F., Rogers, K. M., & Jenders, R. A. (2012). *Improving outcomes with clinical decision support: an implementer's guide*. Himss Publishing.
- Osheroff, J. A., Teich, J. M., Middleton, B., Steen, E. B., Wright, A., & Detmer, D. E. (2007, Mar-Apr). A roadmap for national action on clinical decision support. *J Am Med Inform Assoc*, 14(2), 141-145. <https://doi.org/10.1197/jamia.M2334>
- Palinkas, L. A., Horwitz, S. M., Green, C. A., Wisdom, J. P., Duan, N., & Hoagwood, K. (2015, Sep). Purposeful Sampling for Qualitative Data Collection and Analysis in Mixed Method Implementation Research. *Adm Policy Ment Health*, 42(5), 533-544. <https://doi.org/10.1007/s10488-013-0528-y>
- Papatheodorou, K., Banach, M., Bekiari, E., Rizzo, M., & Edmonds, M. (2018, 2018/03/11). Complications of Diabetes 2017. *Journal of diabetes research*, 2018, 3086167. <https://doi.org/10.1155/2018/3086167>
- Partridge, A. H., Hughes, M. E., Warner, E. T., Ottesen, R. A., Wong, Y. N., Edge, S. B., Theriault, R. L., Blayney, D. W., Niland, J. C., Winer, E. P., Weeks, J. C., & Tamimi, R. M. (2016, Sep 20). Subtype-Dependent Relationship Between Young Age at Diagnosis and Breast Cancer Survival. *J Clin Oncol*, 34(27), 3308-3314. <https://doi.org/10.1200/jco.2015.65.8013>
- Pastore, S. (2012). Website development and web standards in the ubiquitous world: Where are we going. *WSEAS Transactions on Computers*, 11(4), 309-318.
- Paz, F., & Pow-Sang, J. A. (2016). A systematic mapping review of usability evaluation methods for software development process. *International Journal of Software Engineering and Its Applications*, 10(1), 165-178.

- Peffers, K., Tuunanen, T., Rothenberger, M. A., & Chatterjee, S. (2007). A design science research methodology for information systems research. *Journal of management information systems*, 24(3), 45-77.
- Peiris, D., Joshi, R., Webster, R., Groenestein, P., Usherwood, T., Heeley, E., Turnbull, F., Lipman, A., & Patel, A. (2009). An electronic clinical decision support tool to assist primary care providers in cardiovascular disease risk management: development and mixed methods evaluation. *Journal of medical Internet research*, 11(4), e1258.
- Peña-Bautista, C., Durand, T., Oger, C., Baquero, M., Vento, M., & Cháfer-Pericás, C. (2019). Assessment of lipid peroxidation and artificial neural network models in early Alzheimer Disease diagnosis. *Clinical biochemistry*, 72, 64-70.
- Perez-Riverol, Y., Gatto, L., Wang, R., Sachsenberg, T., Uszkoreit, J., Leprevost, F. d. V., Fufezan, C., Ternent, T., Eglen, S. J., Katz, D. S., Pollard, T. J., Kononov, A., Flight, R. M., Blin, K., & Vizcaíno, J. A. (2016). Ten Simple Rules for Taking Advantage of Git and GitHub. *PLOS Computational Biology*, 12(7), e1004947. <https://doi.org/10.1371/journal.pcbi.1004947>
- Pipinellis, A. (2015). *GitHub essentials* (Vol. 2). Packt Publishing.
- Piri, S., Delen, D., Liu, T., & Zolbanin, H. M. (2017). A data analytics approach to building a clinical decision support system for diabetic retinopathy: Developing and deploying a model ensemble. *Decision support systems*, 101, 12-27.
- Prat, N., Comyn-Wattiau, I., & Akoka, J. (2014). Artifact evaluation in information systems design-science research—a holistic view.
- Prat, N., Comyn-Wattiau, I., & Akoka, J. (2015). A taxonomy of evaluation methods for information systems artifacts. *Journal of management information systems*, 32(3), 229-267.
- Preston, M., & Mehandjiev, N. (2004). A framework for classifying intelligent design theories. Proceedings of the 2004 ACM workshop on Interdisciplinary software engineering research,
- Pries-Heje, J., Baskerville, R., & Venable, J. R. (2008). Strategies for design science research evaluation.

- Prieto-Merino, D., Dobson, J., Gupta, A. K., Chang, C. L., Sever, P. S., Dahlöf, B., Wedel, H., Pocock, S., Poulter, N., & On Behalf of the, A.-B. I. (2013, 2013/08/01). ASCORE: an up-to-date cardiovascular risk score for hypertensive patients reflecting contemporary clinical practice developed using the (ASCOT-BPLA) trial data. *Journal of Human Hypertension*, 27(8), 492-496. <https://doi.org/10.1038/jhh.2013.3>
- Pylypchuk, R., Wells, S., Kerr, A., Poppe, K., Harwood, M., Mehta, S., Grey, C., Wu, B. P., Selak, V., Drury, P. L., Chan, W. C., Orr-Walker, B., Murphy, R., Mann, J., Krebs, J. D., Zhao, J., & Jackson, R. (2021, 2021/06/12/). Cardiovascular risk prediction in type 2 diabetes before and after widespread screening: a derivation and validation study. *The Lancet*, 397(10291), 2264-2274. [https://doi.org/https://doi.org/10.1016/S0140-6736\(21\)00572-9](https://doi.org/https://doi.org/10.1016/S0140-6736(21)00572-9)
- Pylypchuk, R., Wells, S., Kerr, A., Poppe, K., Riddell, T., Harwood, M., Exeter, D., Mehta, S., Grey, C., & Wu, B. P. (2018). Cardiovascular disease risk prediction equations in 400 000 primary care patients in New Zealand: a derivation and validation study. *The Lancet*, 391(10133), 1897-1907.
- Rączkowski, Ł., Możejko, M., Zambonelli, J., & Szczurek, E. (2019). ARA: accurate, reliable and active histopathological image classification framework with Bayesian deep learning. *Scientific reports*, 9(1), 1-12.
- Raghupathi, W., & Raghupathi, V. (2014). Big data analytics in healthcare: promise and potential. *Health information science and systems*, 2(1), 3.
- Rao, B. N., Quinn, N., Januszewski, A. S., Peto, T., Brazionis, L., Aryal, N., O'Connell, R. L., Li, L., Summanen, P., Scott, R., O'Day, J., Keech, A. C., & Jenkins, A. J. (2022, 2022/04/01/). Retinopathy risk calculators in the prediction of sight-threatening diabetic retinopathy in type 2 diabetes: A FIELD substudy. *Diabetes research and clinical practice*, 186, 109835. <https://doi.org/https://doi.org/10.1016/j.diabres.2022.109835>
- Rehman, A. A., & Alharthi, K. (2016). An introduction to research paradigms. *International Journal of Educational Investigations*, 3(8), 51-59.
- Remenyi, D., & Sherwood-Smith, M. (2012). *IT investment: making a business case*. Routledge.
- Robinson, T., Elley, C. R., Wells, S., Robinson, E., Kenealy, T., Pylypchuk, R., Bramley, D., Arroll, B., Crengle, S., & Riddell, T. (2012). New Zealand Diabetes Cohort Study cardiovascular

- risk score for people with Type 2 diabetes: validation in the PREDICT cohort. *Journal of primary health care*, 4(3), 181-188.
- Rodríguez-Almagro, J., García-Manzanares, Á., Lucendo, A. J., & Hernández-Martínez, A. (2018). Health-related quality of life in diabetes mellitus and its social, demographic and clinical determinants: A nationwide cross-sectional survey. *Journal of clinical nursing*, 27(21-22), 4212-4223.
- Romana, J., Law, M., Murphy, R., Morunga, E., & Broadbent, E. (2022). Illness perceptions and diabetes self-care behaviours in Māori and New Zealand Europeans with type 2 diabetes mellitus: a cross-sectional study. *The New Zealand Medical Journal (Online)*, 135(1561), 31-35.
- Romero-Aroca, P., Valls, A., Moreno, A., Sagarra-Alamo, R., Basora-Gallisa, J., Saleh, E., Baget-Bernaldiz, M., & Puig, D. (2019). A clinical decision support system for diabetic retinopathy screening: creating a clinical support application. *Telemedicine and e-Health*, 25(1), 31-40.
- Romero-Aroca, P., Verges, R., Maarof, N., Vallas-Mateu, A., Latorre, A., Moreno-Ribas, A., Sagarra-Alamo, R., Basora-Gallisa, J., Cristiano, J., & Baget-Bernaldiz, M. (2022). Real-world outcomes of a clinical decision support system for diabetic retinopathy in Spain. *BMJ Open Ophthalmology*, 7(1), e000974.
- Rooney, J., Byrne, S., Heverin, M., Corr, B., Elamin, M., Staines, A., Goldacre, B., & Hardiman, O. (2013). Survival analysis of irish amyotrophic lateral sclerosis patients diagnosed from 1995-2010. *PLoS one*, 8(9), e74733. <https://doi.org/10.1371/journal.pone.0074733>
- Royston, P., & Altman, D. G. (2013). External validation of a Cox prognostic model: principles and methods. *BMC medical research methodology*, 13, 1-15.
- Saaristo, T., Moilanen, L., Jokelainen, J., Korpi-Hyövälti, E., Vanhala, M., Saltevo, J., Niskanen, L., Peltonen, M., Oksa, H., Cederberg, H., Tuomilehto, J., Uusitupa, M., & Keinänen-Kiukaanniemi, S. (2010, 2010/12/01/). Cardiometabolic profile of people screened for high risk of type 2 diabetes in a national diabetes prevention programme (FIN-D2D). *Primary Care Diabetes*, 4(4), 231-239. <https://doi.org/https://doi.org/10.1016/j.pcd.2010.05.005>
- Saeedi, P., Petersohn, I., Salpea, P., Malanda, B., Karuranga, S., Unwin, N., Colagiuri, S., Guariguata, L., Motala, A. A., & Ogurtsova, K. (2019). Global and regional diabetes

prevalence estimates for 2019 and projections for 2030 and 2045: Results from the International Diabetes Federation Diabetes Atlas. *Diabetes research and clinical practice*, 157, 107843.

Sajja, P. S., & Akerkar, R. (2010). *Knowledge-Based Systems for Development*. TMRF E-Book.

Sandoval-Garcia, E., McLachlan, S., Price, A. H., MacGillivray, T. J., Strachan, M. W., Wilson, J. F., & Price, J. F. (2021). Retinal arteriolar tortuosity and fractal dimension are associated with long-term cardiovascular outcomes in people with type 2 diabetes. *Diabetologia*, 64(10), 2215-2227.

Saxena, R., Voight, B. F., Lyssenko, V., Burt, N. P., de Bakker, P. I., Chen, H., Roix, J. J., Kathiresan, S., Hirschhorn, J. N., & Daly, M. J. (2007). Genome-wide association analysis identifies loci for type 2 diabetes and triglyceride levels. *Science*, 316(5829), 1331-1336.

Schäfer, Z., Mathisen, A., Svendsen, K., Engberg, S., Rolighed Thomsen, T., & Kirketerp-Møller, K. (2021). Toward machine-learning-based decision support in diabetes care: A risk stratification study on diabetic foot ulcer and amputation. *Frontiers in medicine*, 7, 601602.

Schmidt, M. I., Duncan, B. B., Bang, H., Pankow, J. S., Ballantyne, C. M., Golden, S. H., Folsom, R., & Chambless, L. E. (2005). Identifying Individuals at High Risk for Diabetes The Atherosclerosis Risk in Communities study. *Diabetes care*, 28(8), 2013-2018.

Schoen, D. E., Glance, D. G., & Thompson, S. C. (2015, 2015/12/12). Clinical decision support software for diabetic foot risk stratification: development and formative evaluation. *Journal of Foot and Ankle Research*, 8(1), 73. <https://doi.org/10.1186/s13047-015-0128-z>

Schulze, M. B., Hoffmann, K., Boeing, H., Linseisen, J., Rohrmann, S., Möhlig, M., Pfeiffer, A. F., Spranger, J., Thamer, C., & Häring, H.-U. (2007). An accurate risk score based on anthropometric, dietary, and lifestyle factors to predict the development of type 2 diabetes. *Diabetes care*, 30(3), 510-515.

Schulze, M. B., Hoffmann, K., Boeing, H., Linseisen, J., Rohrmann, S., Möhlig, M., Pfeiffer, A. F. H., Spranger, J., Thamer, C., Häring, H., Fritsche, A., & Joost, H. (2007). An Accurate Risk Score Based on Anthropometric, Dietary, and Lifestyle Factors to Predict the Development of Type 2 Diabetes. *Diabetes care*, 30(3), 510-515.

- Schwarz, P. E., Li, J., Lindstrom, J., & Tuomilehto, J. (2009). Tools for predicting the risk of type 2 diabetes in daily practice. *Hormone and metabolic research*, 41(02), 86-97.
- Scott, A., Toomath, R., Bouchier, D., Bruce, R., Crook, N., Carroll, D., Cutfield, R., Dixon, P., Doran, J., & Dunn, P. (2006). First national audit of the outcomes of care in young people with diabetes in New Zealand: high prevalence of nephropathy in Maori and Pacific Islanders. *The New Zealand Medical Journal*, 119(1235).
- Scott, L. J., Mohlke, K. L., Bonnycastle, L. L., Willer, C. J., Li, Y., Duren, W. L., Erdos, M. R., Stringham, H. M., Chines, P. S., & Jackson, A. U. (2007). A genome-wide association study of type 2 diabetes in Finns detects multiple susceptibility variants. *Science*, 316(5829), 1341-1345.
- Seltman, H. J. (2018). *Experimental Design and Analysis*.
- Sesen, M. B., Nicholson, A. E., Banares-Alcantara, R., Kadir, T., & Brady, M. (2013). Bayesian networks for clinical decision support in lung cancer care. *PloS one*, 8(12), e82349.
- Shahsavarani, A. M., Azad Marz Abadi, E., Hakimi Kalkhoran, M., Jafari, S., & Qaranli, S. (2015). Clinical decision support systems (CDSSs): state of the art review of literature. *International Journal of Medical Reviews*, 2(4), 299-308.
- Shalom, E., Shahar, Y., & Lunenfeld, E. (2016). An architecture for a continuous, user-driven, and data-driven application of clinical guidelines and its evaluation. *Journal of Biomedical Informatics*, 59, 130-148.
- Shang, J., Wang, Q., Zhang, H., Wang, X., Wan, J., Yan, Y., Gao, Y., Cheng, J., Li, Z., & Lin, J. (2021). The relationship between diabetes mellitus and COVID-19 prognosis: a retrospective cohort study in Wuhan, China. *The American Journal of Medicine*, 134(1), e6-e14.
- Shariatnia, S., Ziaratban, M., Rajabi, A., Salehi, A., Abdi Zarrini, K., & Vakili, M. (2022, Mar 29). Modeling the diagnosis of coronary artery disease by discriminant analysis and logistic regression: a cross-sectional study. *BMC Med Inform Decis Mak*, 22(1), 85.
<https://doi.org/10.1186/s12911-022-01823-8>

- Siau, K., & Rossi, M. (2011). Evaluation techniques for systems analysis and design modelling methods—a review and comparative analysis. *Information Systems Journal*, 21(3), 249-268.
- Sim, I., Gorman, P., Greenes, R. A., Haynes, R. B., Kaplan, B., Lehmann, H., & Tang, P. C. (2001). Clinical decision support systems for the practice of evidence-based medicine. *Journal of the American Medical Informatics Association*, 8(6), 527-534.
- Sim, L. L. W., Ban, K. H. K., Tan, T. W., Sethi, S. K., & Loh, T. P. (2017). Development of a clinical decision support system for diabetes care: A pilot study. *PloS one*, 12(2), e0173021.
- Sim, R., Chong, C. W., Loganadan, N. K., Adam, N. L., Hussein, Z., & Lee, S. W. H. (2022). Comparison of a chronic kidney disease predictive model for type 2 diabetes mellitus in Malaysia using Cox regression versus machine learning approach. *Clinical Kidney Journal*, 16(3), 549-559. <https://doi.org/10.1093/ckj/sfac252>
- Simmons, D. (1998). Diabetic nephropathy in New Zealand Maori and Pacific Islands people. *Nephrology*, 4, S72-S75.
- Simmons, D., Shaw, L. M., Scott, D. J., Kenealy, T., & Scragg, R. K. (1994). Diabetic nephropathy and microalbuminuria in the community: The South Auckland Diabetes Survey. *Diabetes care*, 17(12), 1404-1410.
- Simon, H. A. (1996). *The Sciences of the Artificial, reissue of the third edition with a new introduction by John Laird*. MIT press.
- Singla, R., Aggarwal, S., Bindra, J., Garg, A., & Singla, A. (2022, Jan-Feb). Developing Clinical Decision Support System using Machine Learning Methods for Type 2 Diabetes Drug Management. *Indian J Endocrinol Metab*, 26(1), 44-49. https://doi.org/10.4103/ijem.ijem_435_21
- Smith-Spangler, C. M., Bhattacharya, J., & Goldhaber-Fiebert, J. D. (2012). Diabetes, its treatment, and catastrophic medical spending in 35 developing countries. *Diabetes care*, 35(2), 319-326.
- Smolen, J., Burmester, G., & Combeet, B. (2016). NCD Risk Factor Collaboration (NCD-RisC). Worldwide trends in diabetes since 1980: a pooled analysis of 751 population-based

studies with 4· 4 million participants. *Lancet* 2016; 387: 1513–30—In this Article, Catherine Pelletier.

Sommerville, I. (2011). *Software engineering* (ed.). *America: Pearson Education Inc.*

Souza-Pereira, L., Ouhbi, S., & Pombo, N. (2021). Quality-in-use characteristics for clinical decision support system assessment. *Computer methods and programs in biomedicine*, 207, 106169.

Stats New Zealand. (2023). *Stats New Zealand*. <https://www.stats.govt.nz/>

Stein, L. D. (2010). The case for cloud computing in genome informatics. *Genome biology*, 11(5), 207.

Stern, E. (2004). Philosophies and types of evaluation research. *Evaluation and impact of education and training: The value of learning. Third report on vocational training research in Europe: Synthesis report. Luxembourg: Office for Official Publications of the European Communities.*

Stevenson, M. (2007). *An Introduction to Survival Analysis*.

Štiglic, G., Fijačko, N., Stožer, A., Sheikh, A., & Pajnikihar, M. (2016). Validation of the Finnish Diabetes Risk Score (FINDRISC) questionnaire for undiagnosed type 2 diabetes screening in the Slovenian working population. *Diabetes research and clinical practice*, 120, 194-197.

Suárez-Figueroa, M. C., Gómez-Pérez, A., & Villazón-Terrazas, B. (2009). How to write and use the ontology requirements specification document. On the Move to Meaningful Internet Systems: OTM 2009: Confederated International Conferences, CoopIS, DOA, IS, and ODBASE 2009, Vilamoura, Portugal, November 1-6, 2009, Proceedings, Part II,

Sugandh, F., Chandio, M., Raveena, F., Kumar, L., Karishma, F., Khuwaja, S., Memon, U. A., Bai, K., Kashif, M., Varrassi, G., Khatri, M., & Kumar, S. (2023, Aug). Advances in the Management of Diabetes Mellitus: A Focus on Personalized Medicine. *Cureus*, 15(8), e43697. <https://doi.org/10.7759/cureus.43697>

Sun, Q. F., Ding, J. G., Xu, D. Z., Chen, Y. P., Hong, L., Ye, Z. Y., Zheng, M. H., Fu, R. Q., Wu, J. G., Du, Q. W., Chen, W., Wang, X. F., & Sheng, J. F. (2009, Jul). Prediction of the prognosis of

- patients with acute-on-chronic hepatitis B liver failure using the model for end-stage liver disease scoring system and a novel logistic regression model. *J Viral Hepat*, 16(7), 464-470. <https://doi.org/10.1111/j.1365-2893.2008.01046.x>
- Sun, Y., & Kantor, P. B. (2006). Cross-Evaluation: A new model for information system evaluation. *Journal of the American Society for Information Science and Technology*, 57(5), 614-628.
- Sutton, R. T., Pincock, D., Baumgart, D. C., Sadowski, D. C., Fedorak, R. N., & Kroeker, K. I. (2020). An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ digital medicine*, 3(1), 17.
- Takeda, H., Veerkamp, P., & Yoshikawa, H. (1990). Modeling design process. *AI magazine*, 11(4), 37-37.
- Thanga Selvi, R., & Muthulakshmi, I. (2020). An Extensive Survey on Recent Machine Learning Algorithms for Diabetes Mellitus Prediction. In S. Balaji, Á. Rocha, & Y.-N. Chung, *Intelligent Communication Technologies and Virtual Mobile Networks Cham*.
- Todd, J., Gepp, A., Stern, S., & Vanstone, B. J. (2022, 2022/05/01/). Improving decision making in the management of hospital readmissions using modern survival analysis techniques. *Decision support systems*, 156, 113747. <https://doi.org/https://doi.org/10.1016/j.dss.2022.113747>
- Tolley, H. D., Barnes, J. M., & Freeman, M. D. (2016). Chapter 10 - Survival Analysis. In M. D. Freeman & M. P. Zeegers (Eds.), *Forensic Epidemiology* (pp. 261-284). Academic Press. <https://doi.org/https://doi.org/10.1016/B978-0-12-404584-2.00010-0>
- Tubey, R. J., Rotich, J. K., & Bengat, J. K. (2015). Research Paradigms.
- Tullis, T., & Stetson, J. (2006, 06/27). A Comparison of Questionnaires for Assessing Website Usability.
- UKNHS, U. N. H. S. (2018). *Welcome to the QDiabetes®-2018 risk calculator*. <https://qdiabetes.org/>

- Ulapane, N., Forkan, A. R. M., Jayaraman, P. P., Schofield, P., Burbury, K., & Wickramasinghe, N. (2023). Using Task Technology Fit Theory to Guide the Codesign of Mobile Clinical Decision Support Systems.
- Urbanovych, A., & Suslyk, H. (2018). Prognosis of ischaemic heart disease in patients with newly diagnosed dm type 2 by logistic regression. *Wiad Lek*, 71(3 pt 2), 691-694.
- Utomo, P. (2020). Building serverless website on GitHub pages. IOP Conference Series: Materials Science and Engineering,
- Vaidya, H., Chaudhari, M. S., & Ingale, H. (2017). Neuro Fuzzy Based Liver Disease, Classification. *International Journal of Advance Research and Innovative Ideas in Education*.
- Vaishnavi, V., & Kuechler, B. (2004, 01/01). Design Science Research in Information Systems. *Association for Information Systems*.
- Vaishnavi, V. K. (2007). *Design science research methods and patterns: innovating information and communication technology*. Auerbach Publications.
- Vaishnavi, V. K., & Kuechler, W. (2015). *Design science research methods and patterns: innovating information and communication technology*. Crc Press.
- Van der Merwe, A., Gerber, A., & Smuts, H. (2020). Guidelines for conducting design science research in information systems. ICT Education: 48th Annual Conference of the Southern African Computer Lecturers' Association, SACLA 2019, Northern Drakensberg, South Africa, July 15–17, 2019, Revised Selected Papers 48,
- van Smeden, M., Reitsma, J. B., Riley, R. D., Collins, G. S., & Moons, K. G. (2021). Clinical prediction models: diagnosis versus prognosis. *Journal of clinical epidemiology*, 132, 142-145.
- Vartiainen, E., Laatikainen, T., Peltonen, M., & Puska, P. (2016, 2016/06/01/). Predicting Coronary Heart Disease and Stroke: The FINRISK Calculator. *Global Heart*, 11(2), 213-216. <https://doi.org/https://doi.org/10.1016/j.gheart.2016.04.007>
- Venable, J. (2006). The role of theory and theorising in design science research. Proceedings of the 1st international conference on design science in information systems and technology (desrist 2006),

- Venable, J., Pries-Heje, J., & Baskerville, R. (2012). A comprehensive framework for evaluation in design science research. *Design Science Research in Information Systems. Advances in Theory and Practice: 7th International Conference, DESRIST 2012, Las Vegas, NV, USA, May 14-15, 2012. Proceedings 7,*
- Venable, J., Pries-Heje, J., & Baskerville, R. (2016, 2016/01/01). FEDS: a Framework for Evaluation in Design Science Research. *European Journal of Information Systems, 25(1), 77-89.* <https://doi.org/10.1057/ejis.2014.36>
- Vitabile, S., Marks, M., Stojanovic, D., Pillana, S., Molina, J. M., Krzyszton, M., Sikora, A., Jarynowski, A., Hosseinpour, F., & Jakobik, A. (2019). Medical data processing and analysis for remote health and activities monitoring. In *High-Performance Modelling and Simulation for Big Data Applications: Selected Results of the COST Action IC1406 cHiPSet* (pp. 186-220). Springer International Publishing Cham.
- Wagholikar, K. B., Sundararajan, V., & Deshpande, A. W. (2012). Modeling paradigms for medical diagnostic decision support: a survey and future directions. *Journal of medical systems, 36(5), 3029-3049.*
- Walls, J. G., Widmeyer, G. R., & El Sawy, O. A. (1992). Building an information system design theory for vigilant EIS. *Information systems research, 3(1), 36-59.*
- Wang, P., Li, Y., & Reddy, C. K. (2019). Machine learning for survival analysis: A survey. *ACM Computing Surveys (CSUR), 51(6), 1-36.*
- Wang, Y., & Cao, K. (2014). A proactive complex event processing method for large-scale transportation internet of things. *International Journal of Distributed Sensor Networks, 10(3), 159052.*
- Wasylewicz, A., & Scheepers-Hoeks, A. (2019). Clinical decision support systems. *Fundamentals of clinical data science, 153-169.*
- WDHB. (2022). *Waikato District Health Board*. Retrieved 20/10 from <https://www.waikatodhb.health.nz/about-us/snapshot-of-waikato-dhb/>

- Web Of Science. (2021). Retrieved 2020/12/10 from https://apps.webofknowledge.com/WOS_GeneralSearch_input.do?product=WOS&search_mode=GeneralSearch&SID=C3fvHdcit4mReEsqQnW&preferencesSaved=
- Weber, S. (2010). Design science research: Paradigm or approach?
- WHO, W. H. O. (2016). *GLOBAL REPORT ON DIABETES*.
- WHO, W. H. O. (2020). WHO Information Network for epidemics.
- WICHAI AEKPLAKORN, PONGAMORN BUNNAG, MARK WOODWARD, PIYAMITR SRITARA, SAYAN CHEEPUDOMWIT, SUKIT YAMWONG, TADA YIPINTSOI, & RAJATA RAJATANAVIN. (2006). A Risk Score for Predicting Incident Diabetes in the Thai Population. *Diabetes care*, 29(8), 1872-1876.
- Wiens, J., & Shenoy, E. S. (2018). Machine learning for healthcare: on the verge of a major shift in healthcare epidemiology. *Clinical Infectious Diseases*, 66(1), 149-153.
- Wieringa, R., & Morali, A. (2012a). Technical action research as a validation method in information systems design science. International Conference on Design Science Research in Information Systems,
- Wieringa, R., & Morali, A. (2012b). Technical action research as a validation method in information systems design science. Design Science Research in Information Systems. Advances in Theory and Practice: 7th International Conference, DESRIST 2012, Las Vegas, NV, USA, May 14-15, 2012. Proceedings 7,
- Wieringa, R. J. (2014). *Design science methodology for information systems and software engineering*. Springer.
- Wieringa, R. J. (2016). Design science research in information systems and software systems engineering. ClbSE,
- William, D., & Black, P. (1996). Meanings and consequences: a basis for distinguishing formative and summative functions of assessment? *British educational research journal*, 22(5), 537-548.

- Willers, C., Iderberg, H., Axelsen, M., Dahlström, T., Julin, B., Leksell, J., Lindberg, A., Lindgren, P., Looström Muth, K., & Svensson, A.-M. (2018). Sociodemographic determinants and health outcome variation in individuals with type 1 diabetes mellitus: A register-based study. *PloS one*, *13*(6), e0199170.
- Wills, M. J. (2014). Decisions Through Data: Analytics in Healthcare. *Journal of Healthcare Management*, *59*(4), 254-262.
https://journals.lww.com/jhmonline/Fulltext/2014/07000/Decisions_Through_Data_Analytics_in_Healthcare.5.aspx
- Wilstrup, C., & Cave, C. (2022, Jul 25). Combining symbolic regression with the Cox proportional hazards model improves prediction of heart failure deaths. *BMC Med Inform Decis Mak*, *22*(1), 196. <https://doi.org/10.1186/s12911-022-01943-1>
- Winter, A., Ammenwerth, E., Haux, R., Marschollek, M., Steiner, B., & Jahn, F. (2023). *Health Information Systems: Technological and Management Perspectives*. Springer Nature.
- Witten, I. H., Frank, E., Hall, M. A., Pal, C. J., & DATA, M. (2005). Practical machine learning tools and techniques. *Data Mining*,
- Woodmansey, C., McGovern, A. P., McCullough, K. A., Whyte, M. B., Munro, N. M., Correa, A. C., Gatenby, P. A., Jones, S. A., & de Lusignan, S. (2017). Incidence, demographics, and clinical characteristics of diabetes of the exocrine pancreas (type 3c): a retrospective cohort study. *Diabetes care*, *40*(11), 1486-1493.
- Yang, X., Yu, Y., Xu, J., Shu, H., Liu, H., Wu, Y., Zhang, L., Yu, Z., Fang, M., & Yu, T. (2020). Clinical course and outcomes of critically ill patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered, retrospective, observational study. *The Lancet Respiratory Medicine*.
- Yen, P.-Y., & Bakken, S. (2012). Review of health information technology usability study methodologies. *Journal of the American Medical Informatics Association*, *19*(3), 413-422.
- Yoon, P. W., Scheuner, M. T., Jorgensen, C., & Khoury, M. J. (2009). Developing Family Healthware, a family history screening tool to prevent common chronic diseases. *Preventing Chronic Disease*, *6*(1).

- Young, B. A., Lin, E., Von Korff, M., Simon, G., Ciechanowski, P., Ludman, E. J., Everson-Stewart, S., Kinder, L., Oliver, M., & Boyko, E. J. (2008). Diabetes complications severity index and risk of mortality, hospitalization, and healthcare utilization. *The American journal of managed care*, *14*(1), 15.
- Young-Hyman, D. L., Peterson, C. M., Fischer, S., Markowitz, J. T., Muir, A. B., & Laffel, L. M. (2016). Depressive symptoms, emotion dysregulation, and bulimic symptoms in youth with type 1 diabetes: Varying interactions at diagnosis and during transition to insulin pump therapy. *Journal of Diabetes Science and Technology*, *10*(4), 845-851.
- Zaharias, P. (2006). A usability evaluation method for e-learning: focus on motivation to learn. CHI'06 Extended Abstracts on Human Factors in Computing Systems,
- Zaman, S. B., De Silva, N., Goh, T. Y., Evans, R. G., Singh, R., Singh, R., Singh, A., Singh, P., & Thrift, A. G. (2023). Design and development of a clinical decision support system for community health workers to support early detection and management of non-communicable disease. *BMJ Innovations*, *9*(1).
- Zarikas, V., Papageorgiou, E., & Regner, P. (2015). Bayesian network construction using a fuzzy rule based approach for medical decision support. *Expert Systems*, *32*(3), 344-369.
- Zeki, T. S., Malakooti, M. V., Ataeipoor, Y., & Tabibi, S. T. (2012). An expert system for diabetes diagnosis. *American Academic & Scholarly Research Journal*, *4*(5), 1.
- Zhang, J. j., Dong, X., Cao, Y. Y., Yuan, Y. d., Yang, Y. b., Yan, Y. q., Akdis, C. A., & Gao, Y. d. (2020). Clinical characteristics of 140 patients infected by SARS-CoV-2 in Wuhan, China. *Allergy*.
- Zhang, T., Zhu, T., Xiong, P., Huo, H., Tari, Z., & Zhou, W. (2019). Correlated differential privacy: Feature selection in machine learning. *IEEE Transactions on Industrial Informatics*, *16*(3), 2115-2124.
- Zheng, K. (2007). *Design, implementation, user acceptance, and evaluation of a clinical decision support system for evidence-based medicine practice*. Carnegie Mellon University.
- Zhou, C.-M., Wang, Y., Xue, Q., Yang, J.-J., & Zhu, Y. (2023, 2023/05/31). Predicting early postoperative PONV using multiple machine-learning- and deep-learning-algorithms.

BMC medical research methodology, 23(1), 133. <https://doi.org/10.1186/s12874-023-01955-z>

Zikos, D., & DeLellis, N. (2018). CDSS-RM: a clinical decision support system reference model. *BMC medical research methodology*, 18(1), 137.

Zlotnik, A., Alfaro, M. C., Pérez, M. C. P., Gallardo-Antolín, A., & Martínez, J. M. M. (2016, 2016/05//). Building a Decision Support System for Inpatient Admission Prediction With the Manchester Triage System and Administrative Check-in Variables. *Computers, informatics, nursing : CIN*, 34(5), 224-230.
<https://doi.org/10.1097/cin.0000000000000230>

Zolnoori, M., Zarandi, M. H. F., & Mostafa, M. (2011). Fuzzy rule-base expert system for evaluation possibility of fatal asthma. *Journal of Health Informatics in Developing Countries*, 5(1).

Žukauskas, P., Vveinhardt, J., & Andriukaitienė, R. (2018). Philosophy and paradigm of scientific research. *Management culture and corporate social responsibility*, 121, 139.

Appendix A: Ethical Approval

The University of Waikato
Private Bag 3105
Gate 1, Knighton Road
Hamilton, New Zealand

Human Research Ethics Committee
Roger Moltzen
Telephone: +64021658119
Email: humanethics@waikato.ac.nz



THE UNIVERSITY OF
WAIKATO
Te Whare Wānanga o Waikato

28 April 2020

Madurapperumage Anuradha Erandathi
School of Management
By email: am448@students.waikato.ac.nz

Dear Anuradha

HREC(Health)2020#26 : Chronological Risk Estimation and Prediction in Health Informatics through a Knowledge-Based System: A Case Study of Complication of Diabetes Mellitus

Thank you for submitting your application HREC(Health)2020# 26 for ethical approval. The committee noted the project is low risk and will access anonymised medical records from two Sri Lankan medical centres to assess how to improve medical decisions and creating an online information hub for diabetes patients to assess their own risk factors for diabetes mellitus.

Please contact the committee by email (humanethics@waikato.ac.nz) if you wish to make changes to your project as it unfolds, quoting your application number with your future correspondence. Any minor changes or additions to the approved research activities can be handled outside the monthly application cycle.

We wish you all the best with your research.

Regards,

A handwritten signature in black ink, appearing to be 'RM'.

Emeritus Professor Roger Moltzen MNZM
Chairperson
University of Waikato Human Research Ethics Committee

Appendix B: Requirement Specification Document

Clinical Decision Support System (CDSS) Requirements Specification Document

Project Name: CDSS for Predicting Diabetes

Client: Te Whatu Ora

Date: 2020.10.27

Prepared By: Anuradha Madurapperuma

1. Meeting Overview

- **Objective:** Collect requirements for developing a CDSS to predict diabetes based on patient datasets.
 - **Participants:**
 - Team members and team lead of diabetes-related data management.
 - System administrator
 - **Meeting Date/Time:** 2020.10.16
 - **Location:** Te Whatu Ora, Hamilton.
-

2. Project Scope

- **Overview:**

A brief description of the CDSS project, its goals, and how it will be used by clinicians to predict diabetes based on historical patient data.
 - **Goals:**
 - Develop a system that uses past patient data to predict the likelihood of complications of diabetes.
-

3. Requirements Collection

3.1. Functional Requirements

- **What should the system do?** (Use MoSCoW framework to prioritize)
 - **Must Have:**
 - Predict the risk of the complications of diabetes using patient data (age, lab results, medical history, etc.).
 - Provide a user-friendly tool for the public.
 - **Should Have:**
 - Statistical overview of the cohort of diabetes.
 - Graphical visualization of survival rates of the cohort.
 - **Could Have:**
 - Comparison of individual health against the cohort.
 - **Won't Have** (for now):
 - Integration with mobile apps for patients (optional in the future).

3.2. Non-Functional Requirements

- **Performance:**
 - Accurate predictions on the complications.
- **Security:**
 - Provide security for patient data.
 - Do not publish the dataset even with the anonymous patient details.
- **Usability:**
 - Easy-to-navigate interface for all stakeholders.
 - Minimal training required for usage.
 - Inexpensive or cost free tool.

3.3. Data Requirements

- **Available Dataset Features:**
 - Demographics (age, gender, ethnicity, Moari/non-Maori)
 - Lab Results (HbA1c, cholesterol, HDL, LDL, triglyceride, eGFR)
- **Missing/Additional Data:**
 - eGFR values are missing after 2012.
 - Is there any way to collect BMI values of the diabetes patients?
- **Data Privacy and Compliance:**
 - What legal or ethical considerations are there for handling patient data?

3.4. User Stories/Use Cases

- Wants to input a patient's data and receive a risk prediction for complications of diabetes so that I can plan the patient's treatment.
- Wants to get the patients risk for the complications of diabetes as a percentage over the time.
- Wants to visualize the distribution of socio-demographic details of the cohort.
- Wants to see the changes in survival rate of patients over the time.
- Wants to get the survival of patients for the most common complications of diabetes in the cohort.
- Wants to download the patients reports for each complication/all of the selected complications.

4. Pain Points and Challenges

- **Current Challenges:**
 - Difficulty predicting diabetes early with current methods.
 - Lack of integration between patient data systems and prediction tools.
- **Client Expectations:**
 - Accurate and early prediction of the complications of diabetes.
 - User-friendly CDSS to use as a prediction model and information system.

5. Stakeholder Input

- **What do stakeholders want from the system?**
 - Risk prediction from basic demographic details and laboratory values.
 - Data visualization of the socio-demographic details to understand the overview of the cohort.
 - **Ethical/Regulatory Concerns:**
 - The privacy and security of the dataset.
-

6. Success Criteria

- **How will success be measured?**
 - Level of matching functionality with the requirements.
 - Usability of the CDSS.
 - The accuracy of the predictions of complications of diabetes.
 - Ease of use for stakeholders.
 - Feedback from the stakeholders.

7. Next Steps/Action Items

- **Follow-Up Actions:**
 - Clarifications on dataset.
 - Completing the dataset by fixing the inconsistencies, missing data, and noise handling.
 - System design, implementation, and evaluation.
- **Additional Meetings/Discussions:**
 - The feedback gathering at the end of the system implementation.

Appendix C: Software Quality Measurement

Requirements	Functional Suitability			Performance Efficiency			Compatibility		Usability						Reliability				Security					Maintainability					Portability		
	Functional Completeness	Functional Correctness	Functional Appropriateness	Time Behaviour	Resource Utilization	Capacity	Coexistence	Interoperability	Appropriate Recognisability	Learnability	Operability	User Error Protection	User Interface Aesthetics	Accessibility	Maturity	Fault Tolerance	Recoverability	Availability	Confidentiality	Integrity	Non-repudiation	Authenticity	Accountability	Modularity	Reusability	Analysability	Modifiability	Testability	Adaptability	Installability	Replaceability
DA1	+	-	+	+	+	+	+	-	-	-	-	-	-	+	-	+	-	-	-	-	-	-	-	-	-	-	-	+	+	-	
DA2	+	-	+	+	-	-	-	-	-	-	-	-	-	+	+	-	-	-	-	-	-	-	-	-	-	-	-	+	+	-	
DA3	+	-	+	+	-	-	+	-	-	-	-	-	-	+	+	-	-	-	-	-	-	-	-	+	-	-	-	+	+	-	
DA4	+	-	+	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
DA5	+	-	+	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	+	+	-	-	-	-	-	-	-	-	-	-	
DA6	+	-	+	-	-	-	-	-	-	-	-	-	-	+	-	-	+	+	+	-	-	-	-	-	-	-	-	-	+	-	
PI1	-	-	+	+	+	+	-	-	+	-	+	+	+	-	+	+	+	+	+	-	-	-	-	+	-	-	-	-	-	-	
PI2	-	-	+	+	+	+	-	-	+	-	+	+	+	-	+	+	+	+	+	-	-	-	-	+	-	-	-	-	-	-	
PI3	-	-	+	+	+	+	-	-	+	-	+	+	+	-	+	+	+	+	+	-	-	-	-	+	-	-	-	-	-	-	
PI4	-	-	+	+	+	+	-	-	+	-	+	+	+	-	+	+	+	+	+	-	-	-	-	+	-	-	-	-	-	-	
QD1	-	-	+	+	+	+	-	-	+	-	+	+	+	-	+	+	+	+	+	-	-	-	-	+	-	-	-	-	-	-	

Appendix D: Questionnaire of Users' Feedback

**PREDICTING THE COMPLICATIONS
OF DIABETES MELLITUS IN
NEW ZEALAND**

NZTPCD



**Evaluation of Web-based NZTPCD (New Zealand Tool for Predicting the Complications of Diabetes)
- Gathering User Feedback**

Thank you for participating in evaluating the Web-based NZTPCD (New Zealand Tool for Predicting the Complications of Diabetes). Your feedback is essential in assessing the usability and effectiveness of this web tool. Please take a few minutes to complete the following form, providing your honest opinions and insights.

This form is automatically collecting emails from all respondents. [Change settings](#)

Q1. What is your age? *

- 18 to 24
- 25 to 34
- 35 to 44
- 45 to 54
- 55 to 64
- 65 to 74
- 75 or older

Q2. What is your gender? *

Female

Male

Prefer not to say

Q3. What is your occupation? *

Short answer text
.....

Q4. How easy was learning to use this clinical decision support system (CDSS)? (1 represents the least and 5 represents the most) *

1

2

3

4

5

Q5. I found the CDSS user interfaces easy to understand. (1 represents the least and 5 represents the most) *

1

2

3

4

5

Q6. I was able to complete tasks quickly using the CDSS. (1 represents the least and 5 represents the most) *

1	2	3	4	5
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q7. The CDSS had clear icons, labels, and menu options that were easily remembered. *

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly agree

Q8. I easily recalled the steps to perform tasks in the CDSS. *

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly agree

Q9. Overall, I am satisfied with the usability of the CDSS. *

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly agree

Q10. The CDSS met my needs and expectations. *

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly agree

Q11. The CDSS was accessible and usable with assistive technologies. *

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly agree

Q12. I encountered accessibility challenges or limitations while using the CDSS. *

Yes

No

Q13. If the answer for Q12 is yes, please specify the encountered challenges or limitations. (If *
not fill the field with "NULL")

Long answer text
.....

Q14. The CDSS was visually appealing. *

Strongly disagree

Disagree

Neutral

Agree

Strongly agree

...

Q15. The CDSS provided a pleasant user experience. (1 represents the least and 5 represents *
the most)

1

2

3

4

5

Q16. The CDSS had a well-organised and intuitive layout. (1 represents the least and 5 represents the most)

1	2	3	4	5
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q17. The CDSS had consistent design elements and navigation.

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly agree

Q18. The CDSS had a good system response time.

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly Agree

Q19. The CDSS was reliable.

- Strongly disagree
- Disagree
- Neutral
- Agree
- Strongly Agree

Q20. Please provide any additional comments suggestions or concerns about the usability of CDSS, and any features you think are missing or would be advantageous compared to existing predictive models.

Long answer text
