Rochester Institute of Technology

# RIT Digital Institutional Repository

2024

# Predictive Modeling in Healthcare

Ibrahim Essa Abdulla Ali Alattar
iea4349@rit.edu

# RIT

# Predictive Modeling in Healthcare

# by

# Ibrahim Essa Abdulla Ali Alattar

A Thesis Submitted in Partial Fulfilment of the Requirements for the

Degree of Master of Science in Professional Studies:

Data Analytics

Department of Graduate Programs & Research

**Rochester Institute of Technology**

**RIT Dubai**

# RIT

# Master of Science in Professional Studies:

# Data Analytics

## Graduate Thesis Approval

**Student Name: Ibrahim Essa Abdulla Ali alattar**

**Thesis Title: Predictive Modeling in Healthcare**

**Graduate Committee:**

**Name:     Dr. Sanjay Modak**                          **Date: 15/02/2024**

        **Chair of committee**

---

**Name:     Dr. Ioannis Karamitsos**                  **Date: 15/02/2024**

        **Member of committee**

---

# Acknowledgments

I extend my most profound appreciation to all those who have contributed to the fruitful completion of this proposal on predictive modelling with linear regression. This journey has been challenging, however gigantically rewarding.

I am significantly thankful to my thesis advisor Dr. Ioannis Karamitsos, whose direction and mastery were important all through each arrangement of this research. Their smart input, faithful support, and encouragement have been significant in forming the direction of this study.

I expand my appreciation to Rochester Institute of Technology for giving the conducive scholastic environment vital for undertaking this inquiry about.

To my companions and family, who have been unflinching pillars of back, I am significantly thankful. Your support, understanding, and tolerance have maintained me through the highs and lows of this scholastic endeavour.

Everyone who contributed, in huge or little ways, to this endeavor. Whether through a keen discussion, valuable input, or basically loaning an sympathetic ear, your commitments have cleared out a permanent mark on this thesis. I appreciate the collective exertion that has gone into making this investigation a reality.

In conclusion, I expand my sincere thanks to everybody who has been a portion of this journey. Your back and collaboration have been instrumental in the effective completion of this thesis on predictive modelling with linear regression.

# Abstract

Predictive modelling, especially the use of horizontal lines, has become an important tool in clinical practice to help make informed decisions and accurate predictions. This study focuses on the use of horizontal regression in clinical practice, evaluating its effectiveness in revealing patterns and improving the accuracy of predictions. This study introduces the process of developing a linear model, emphasizing the importance of preliminary data analysis, feature selection, and model evaluation. To make sure your model's predictions are accurate, consider key assumptions such as sampling, independence, and homoscedasticity. The main goal is to provide doctors with the knowledge and skills needed to make accurate predictions using horizontal models, thus improving the clinical decision-making process. This study explores the fundamentals of linear regression, evaluates its suitability for various clinical applications, and outlines the important steps in building a good model. The aim is to provide doctors with the knowledge and skills that will enable them to make informed decisions using technology. Issues such as multicollinearity and overfitting are addressed, while the importance of engineering design and variable selection to optimize model performance is further explored. This research contributes to the nonstop advancement of data-driven decision-making in healthcare by highlighting the imperative part of the even pivot in making exact expectations. Experiences from this research have the potential to progress quiet results, progress asset allotment, and actualize evidence-based practices within the healthcare industry.

***Keywords****: Predictive modelling, Linear regression, Data pre-processing, Feature selection, Model evaluation, Assumptions, Feature engineering, Variable selection*

# Table of Contents

# List of Figures

# List of Tables

# Chapter 1. Introduction

## 1.1   Background

Predictive modelling has gotten to be a critical instrument in numerous areas, permitting educated choices and exact expectations. This analysis analyzes the field of prescient modelling, centering on the utilisation of straight relapse as an effective method to uncover designs and progress the precision of forecasts. Direct relapse could be a broadly utilized strategy in prescient analytics.
In the evolving world of healthcare, integration of advanced analytics has become imperative for the patient to make informed decisions and get better outcomes. Predictive modelling, especially when viewed through the lens of linear regression, becomes a powerful tool for uncovering complex patterns and improving the accuracy of health predictions.

As we begin our journey into the world of renewable energy demonstration, we enter the intersection of data science and healthcare and leverage the predictive power of this approach. Known for its simplicity and interpretation, linear regression provides a unique lens through which we can identify and predict treatments, improve resource allocation, and ultimately improve patient quality.

This research will discuss the main concepts of the predictive analytics, its application in medicine, and the care required to ensure the reliability of the model. From building linear models to the nuances of feature engineering and variable selection, our goal is to provide practitioners with the knowledge and skills needed to build predictive models and meet the challenge of finding greater stability.

Introducing our predictive modeling in healthcare, we strive to contribute to a future where data-driven insights empower practitioners to predict, reduce and improve care delivery standards.

The final section of the introduction chapter should indicate the content of each of the following chapters and how the chapters are related to each other. You are providing the reader with a 'road map' of the work ahead. Such a 'map' will enable the reader to navigate their way through your work much more easily and appreciate to the maximum what you have done.

## 1.2   Problem Statement

The requirement for precise and opportune figures in healthcare has never been more noteworthy. Be that as it may, the complexity of clinical information and the numerous factors included regularly make issues in accomplishing precise predictions. This problem statement addresses the pressing issue of leveraging predictive modelling, specifically through linear regression, to enhance decision-making in healthcare. The current methodologies for predicting patient outcomes, resource allocation, and disease trends may lack the requisite precision. Predictive analytics model with its simplicity and interpretability, holds the potential to fill this gap. Nevertheless, the healthcare sector faces hurdles such as the need for robust models, consideration of diverse variables, and adherence to assumptions like linearity and independence. Consequently, there is a compelling need to explore and optimize the application of linear regression in healthcare predictive modelling to ensure that it becomes a reliable tool for healthcare professionals. This research aims to address these challenges, ultimately contributing to a more accurate and efficient healthcare system through the integration of predictive modelling methodologies.

This research aims to address these challenges and demonstrate the efficacy of different models for predictive modelling in healthcare. We will focus on developing robust and interpretable machine learning models using high-quality healthcare data. Then validating and evaluating the performance of these models in real-world clinical settings and providing healthcare professionals with the tools and knowledge needed to integrate predictive modelling into their decision-making processes.

Linear regression is a widely used method. - Provide a solution to this problem using statistical methods. However, successful implementation requires careful consideration of data quality, conceptual models, and feature selection. This project enables organizations and decision makers

to make decisions from information, improve resource allocation, reduce uncertainty and ultimately increase their overall performance and competitiveness. The project will solve these problems by using historical data to create predictive models that can predict future outcomes, ultimately bridging the gap between historical data and informed decision-making in various domains.

## 1.3   Research Aim and Objectives

The essential aim of this research is to examine and evaluate the viability of predictive modelling utilizing linear regression inside the healthcare space. The overarching objective is to contribute important bits of knowledge that can improve decision-making forms and make strides the exactness of expectations in healthcare scenarios. Through this research, we look to assess the application of linear regression models, centering on their capacity to uncover designs and provide precise figures within the complex and energetic scene of healthcare. The research will delve into the fundamental concepts of linear regression, emphasizing its suitability for healthcare applications, and address the critical steps involved in constructing robust predictive models. The objectives include examining the impact of preliminary data analysis, feature selection, and model evaluation on the reliability of predictions.

Also, the research points to the significance of tending to challenge inalienable in linear regression, such as multicollinearity and overfitting, to guarantee the strength of predictive models in healthcare settings. Eventually, the point is to prepare healthcare experts with profitable information and devices that can be connected to tackle the potential of linear regression models for exact forecasts, subsequently contributing to the progression of data-driven decision-making in healthcare.

The primary objective of utilizing predictive modelling in healthcare is to upgrade decision-making forms and optimize quiet results. By leveraging historical information, this approach reveals designs and connections inside healthcare factors, giving an establishment for making exact expectations around future scenarios.

The overarching objective is to engage healthcare experts with exact experiences, empowering proactive mediations and personalized patient care procedures. Eventually, the application of predictive modeling in healthcare looks to contribute to more effective and successful healthcare frameworks, guaranteeing that assets are designated deliberately, and patients get convenient and custom fitted intercessions based on data-driven forecasts.

## 1.4 Research Questions

1. What are the key steps involved in developing a predictive model using different machine learning algorithms?
2. What is the relationship between a patient's age and their blood pressure levels, and can we use linear regression to predict blood pressure for future check-ups?"
3. How does data pre-processing, including handling missing values and outliers, influence the accuracy of a linear regression model?
4. How does a linear regression model generalize to new, unseen data, and what measures can be taken to improve generalization performance?

## 1.5 Limitations of the study

1. The effectiveness of linear regression can vary across different areas. The study might lack details on specific applications and their unique challenges, limiting its generalizability.
2. The study may not elaborate on the data source or sample size, raising concerns about its representativeness and potential bias.
3. The validation process for the linear regression model may be unclear, leaving doubts about its generalizability and performance on unseen data.
4. The study might not delve deep into the interpretability of the linear regression model, hindering its practical application for users who need to understand the impact of individual features on the predictions.

## 1.6   Structure of the Thesis

The structure of this thesis comprises seven core chapters, beginning with an introduction in which the research problem, the aims and the significance of the thesis are set out.

Chapter 2 is the literature review, which critically examines the existing studies and identifies the gaps to be filled by the research.

Chapter 3 presents the research methodology. The CRISP-DM methodology was chosen for this thesis.

Chapter 4, the data analysis describes the methods of data collection and preliminary analysis in detail.

Chapter 5 forms the basis for the data modeling in which specific models are developed and evaluated to answer the research questions. For this thesis, the linear regression model, the decision tree and the random forest model were selected for the analysis.

Chapter 6 is the discussion, where the results are interpreted in the context of the wider literature and implications, limitations and recommendations are discussed. Finally, Chapter 7 summarizes the main findings, provides a reflective overview of the research contributions and suggests directions for future work. References provides academic sources that underpin the research framework and ensure a comprehensive exploration of the topic in a structured academic format.

.

# Chapter 2 – Literature Review

After conducting a thorough literature review of the research in the area of Predictive Modelling, I observed different types of objectives of such research. Many research articles are focused on surveying existing methods of studying the predictive analytics. These articles do a systematic review of different techniques that other researchers have used for this purpose and compare the results.

Smith et al. (1998) explores the application of linear regression in the clinical setting. Demonstrates the effectiveness of linear regression models in predicting patient outcomes such as infection and treatment based on patient history. The application of linear regression in clinical settings has been examined in detail. The main aim of this study is to evaluate the effectiveness of the horizontal linear model in predicting patient outcomes, with particular emphasis on disease and treatment pain. This study uses patient history data, a valuable resource in healthcare analytics, to develop and evaluate predictive models.

A study conducted by Montgomery & Vining (2012) gives an in-depth look at the application of linear regression within the taking after region that's Budgetary determining. This inquiry centers on how to utilize linear regression models to create almost different viewpoints of back, counting stock costs, commerce patterns, and venture returns. Analysts have emphasized the significance of authentic financial information and money related markers in creating estimating models. Research on power forecasting shows that horizontal linear models have significant power in financial forecasting. These models can provide insight into future stock prices and business models using historical financial data.

Risk evaluation models empower the evaluation of monetary dangers related with speculations. They can measure the potential instability and drawback dangers of money related rebellious, supporting speculators in making educated choices. Portfolio optimization research emphasizes how linear regression can be used to optimize investment portfolios. By predicting the returns of various assets or securities, investors can strategically allocate their resources for maximum returns within their risk tolerance. Market analysis in linear regression models can assist in analyzing market dynamics and identifying influential factors affecting financial markets. Researchers can

use this model to understand the impact of economic indicators, news events, and other variables on market behavior.

Smith et al., (1998) provides a basic understanding of linear regression in predictive modelling. Emphasizes the simplicity and interpretability of linear regression for making predictions.

Montgomery & Vining, (2012) explores the presumptions that underlie linear regression models and talks about their implications. Highlights the importance of endorsing doubts to guarantee the unflinching quality of desires.

Gelman & Hill, (2006) explore diverse ways to choose highlights to make strides in the execution of linear regression models. Stresses the significance of selecting relevant features for accurate predictions and model interpretability.

Neter ,Wasserman & Kutner, (1989) check the application of regular techniques such as Lasso and Ridge to avoid distortion of horizontal lines. Discuss what can be done to help improve performance.

Rencher & Schaalje, (2008) addresses multicollinearity issues and suggests solutions to reduce their impact on linear regression models. The importance of solving multicollinearity problems to obtain reliable models.

Cook & Weisberg, (1982) look for changes in the horizontal line leading up to the forecast period. Explore the suitability of linear regression for modelling physical relationships in prediction.

Draper & Smith, (1998) example analysis of linear regression and other machine learning methods. Discuss the advantages and limitations of linear regression compared to multiple models.

Aiken & West, (1991) explore strategies for handling missing data in the context of linear regression. About the importance of handling missing data to maintain the accuracy of the prediction model.

Anderson and Turner, (2020) focuses on the practical application of linear regression in business forecasting. Highlights real-world scenarios where linear regression has proven to be an asset to a predictive modelling tool.

Fox 2015) discuss ways to interpret the results of horizontal models, including coefficient analysis. Focus on the need to clearly understand the design to remove visual elements.

A study conducted by Gelman & Hill (2006) examined the use of linear regression in maintenance forecasting, focusing on its application in manufacturing. This study demonstrates how linear regression models can be effectively used to predict equipment failures and maintenance needs, ultimately reducing downtime and repair costs. Tool failure prediction demonstrates the ability to predict tool failure with high accuracy. These models can identify failure patterns and signs by analyzing historical maintenance data, sensor readings, and other variables. Quality control is predictive control that ensures organizations use the correct way of working. Maintenance tasks can be scheduled in a predictable manner, reducing unplanned downtime rather than performing maintenance tasks on a planned or reactive basis. Cost Savings Recycling models help generate significant profits by reducing downtime and optimizing maintenance time. Productivity benefits from reduced maintenance, longer equipment life and improved efficiency. Data-driven insight research highlights the importance of leveraging data-driven insights. Linear regression models provide maintenance teams with performance data, allowing them to prioritize tasks and allocate resources more efficiently. Device health monitoring can be used to regularly monitor device health. It detects small changes in equipment behaviour early and flags problems before they become serious.

Neter, Wasserman, and Kutner's (1989) is a foundational text in the field of applied statistics. The study is widely recognized for its comprehensive coverage of linear statistical models and their practical applications in regression analysis, analysis of variance (ANOVA), and experimental designs. The authors begin by providing a thorough introduction to the fundamental concepts of linear statistical models, emphasizing the theoretical underpinnings that form the basis for subsequent discussions. They elucidate the mathematical aspects of regression analysis, allowing readers to build a strong conceptual foundation in understanding how different variables interact and influence the outcomes of interest.

Kleinbaum, Kupper, Nizam & Rosenberg (2013) provides a comprehensive overview of the contributions to the field of regression analysis and multivariable methods. This seminal work serves as a valuable resource for researchers, statisticians, and practitioners seeking a deep

understanding of advanced statistical techniques and their applications. The authors begin by establishing a solid foundation in regression analysis, emphasizing its significance in modelling relationships between variables. (Kleinbaum et al.) delve into the theoretical underpinnings of regression analysis, elucidating concepts such as the assumptions of linearity, independence, and homoscedasticity. They systematically guide the reader through the process of model building, emphasizing the importance of variable selection and the interpretation of coefficients.

Gelman and Hill (2006) could be a comprehensive study that gives an in-depth examination of regression analysis, multilevel modelling and progressive modelling. The consideration by Gelman and Hill (2006) stands out as a valuable instrument in measurable modelling, giving common- and common-sense advice for analysts and specialists. Clinicians are working to progress their understanding of relapse analysis and multilevel models. This study begins with a strong establishment in essential strategies and different regressions, guaranteeing users have a strong understanding of the essentials some time recently jumping into more complex points. Gelman and Hill (2006) approach is known for its exactness because it gives clear clarifications and real-world cases that make factual concepts less demanding to get it. The study also covers the importance of model checking and diagnostics, highlighting the need to assess model assumptions and validity. Gelman and Hill (2006) emphasize the iterative nature of model improvement, empowering users to basically assess their models and refine them based on symptoms. Practical cases from different areas, counting social sciences, science, and instruction, contribute to the book's pertinence over disciplines. By using real-world datasets, Gelman and Hill illustrate how to apply the discussed techniques to address specific research questions, making the material more engaging and applicable.

Gelman & Hill (2006) serves as a comprehensive and receptive direct for analysts and professionals exploring the complexities of regression analysis and multilevel modelling. Its integration of Bayesian strategies, accentuation on progressive modelling, and common sense cases contribute to its persevering esteem as an asset for those looking to improve their measurable modelling abilities.

Charles H. Achen's (1982) stands as a foundational content within the field of insights and regression analysis. Achen digs into the complexities of regression, advertising a comprehensive

exploration that remains significant and compelling. The study not as it were giving a detailed clarification of regression strategies but too emphasizes the basic angle of translation. Achen adeptly navigates the complexities of regression models, shedding light on how they can be successfully utilized to extricate important bits of knowledge from information. His accentuation on translation serves as a directing guideline for analysts and professionals, empowering a more profound understanding of the connections between factors past insignificant scientific definitions. Achen's work is celebrated for its clarity, making factual concepts open to a wide gathering of people. This writing audit recognizes "Interpreting and Using Regression" as a persevering commitment to measurable writing, serving as a foundation for analysts and examiners looking for to saddle the control of regression examination in their work.

Cameron and Trivedi (2013) may be a comprehensive and persuasive commitment to the field of econometrics and measurable modeling. The study centers on the specialized zone of tally information examination, giving a careful writing audit, methodological bits of knowledge, and practical applications. The authors start by tending to the interesting characteristics of number information, emphasizing the discrete and non-negative nature of the subordinate variable, which frequently speaks to the recurrence of occasions or events. The writing audit dives into the authentic advancement of count data models, highlighting prior approaches and their restrictions, setting the organisation for the presentation of more modern procedures. The study fastidiously covers different models for number data examination, with a specific emphasis on the Poisson regression model and its expansions. Cameron and Trivedi explore the inadequacies of the essential Poisson show, such as overdispersion, and introduce elective models just like the negative binomial regression. Their audit of existing techniques isn't as comprehensive but too basic, giving perusers a nuanced understanding of the qualities and shortcomings of diverse approaches. One of the eminent qualities of Cameron and Trivedi work is its availability to both novice and prepared analysts. The authors display complex concepts in a clear and instinctive way, making the fabric receptive for perusers with shifting levels of measurable ability. Moreover, the incorporation of real-world cases and applications upgrades the viable utility of the book, illustrating how count data models can be viably utilized in different areas, from financial matters to open wellbeing. In summary, "Regression Analysis of Count Data" (Cameron and Trivedi,2013) could be a benchmark distribution within the field of number information examination. The writing audit not

as it were synthesizes the verifiable advancement of techniques but moreover basically evaluates their pertinence. The study availability and viable center make it a vital asset for analysts, analysts, and specialists looking to analyze and show check information viably.

Neter, Wasserman, & Kutner's (1989) stands as seminal work within the field of measurements and data analysis. The content gives a comprehensive and available investigation of different direct factual models, advertising a point by point examination of regression, analysis of variance (ANOVA), and test designs. The creators skilfully mix hypothetical establishments with viable applications, making the substance available to both amateur and experienced analysts. The study not as it were covering the basic standards of direct factual models but moreover digs into their real-world usage over assorted disciplines. One outstanding quality of the work is its emphasis on the application of measurable procedures to illuminate practical issues. The authors consolidate various cases and case thinks about, directing readers through the step-by-step handle of model development, investigation, and interpretation. This academic approach upgrades the reader's understanding and capability in applying factual models to real data. Additionally, the clarity and organization of the substance contribute to the persevering significance of this work. The content advances consistently, beginning with foundational concepts and steadily building up to more complex points. The incorporation of works out at the conclusion of each chapter assists reinforces the learning preparation, giving openings for users to hone and fortify their understanding. "Applied Linear Statistical Models" remains a valuable resource for analysts, professionals, and understudies alike. Its persevering popularity attests to its status as a classic reference within the field of connected insights, advertising a strong establishment for anybody looking for to ace relapse, ANOVA, and test plan. The ponder persevering effect underscores its centrality in forming the understanding and application of straight measurable models in different logical and down to earth spaces.

Harrell Jr (2015) stands as an urgent work within the field of factual modelling and regression analysis. The ponder gives a comprehensive and definitive investigation of progressed procedures for creating robust regression models. Harrell's approach goes beyond traditional methods, emphasizing the importance of thoughtful model development, validation, and interpretation. One notable perspective of the study is its center on practical procedures for model building. Harrell

presents readers to the concept of regression modelling as an instrument for making expectations and choices instead of fair fitting a measurable demonstration to the information. By supporting a principled approach to demonstrate improvement, he guides analysts and specialists through the complex preparation of selecting factors, taking care of lost information, and surveying show execution. Harrell dives into the crucial point of demonstrating approval, focusing the requirement for inner and outside approval to guarantee the model's generalizability. His emphasis on the bootstrap method and optimism-corrected performance metrics reflects a commitment to rigorous statistical practices. Furthermore, the study incorporates modern statistical tools and visualizations, enabling readers to navigate complex data structures and diagnose potential issues in their models. The incorporation of subjects such as penalized regression, show calibration, and dealing with intelligence includes profundity to think about. Harrell illustrates a sharp mindfulness of the challenges confronted by investigators in real-world scenarios and gives practical arrangements for overcoming them. His clear and brief composing style makes the fabric open to a wide gathering of people, from amateur analysts to prepared researchers. In summary, "Regression Modeling Strategies" (Harrell Jr, 2015) not as it were serves as an important reference for analysts and information researchers but moreover contributes altogether to the continuous talk on best practices in regression modelling. Its combination of theoretical establishments, practical experiences, and cutting edge methods makes it an irreplaceable asset for anybody looking to upgrade their aptitudes in creating and approving regression models.

**Key Takeaways/Learnings:**

- Linear regression provides a simple and interpretable approach to predictive modelling.
- Validating assumptions and addressing multicollinearity are crucial for the reliability of linear regression predictions.
- Feature selection, regularization, and handling missing data are essential for improving model performance.
- Linear regression can be adapted for time series prediction, making it versatile in various applications
- The practical application of linear regression in business forecasting underscores its real-world utilit

# Chapter 3- Research Methodology

Based on the existing literature review, several limitations of secondary research and gaps in the existing literature can be identified as:

1. **Limited Contextual Understanding:** Secondary research often relies on existing literature, which may lack a comprehensive understanding of the specific contexts in which predictive modelling with linear regression is applied. The studies mentioned provide insights into different domains (clinical settings, financial forecasting, maintenance forecasting), but the contextual nuances might not be fully captured.

2. **Generalization Challenges:** Many articles focus on specific applications of linear regression in various fields. While these studies contribute valuable insights, there may be challenges in generalizing findings across different domains. The specific characteristics of each application area may impact the transferability of predictive models.

3. **Assumption of Data Quality:** Secondary research relies on the assumption that the data used in the reviewed studies are of high quality. The literature may not extensively discuss data quality issues, which are crucial in predictive modelling. Issues such as missing data, outliers, or biases might not be adequately addressed in the existing literature.

4. **Limited Exploration of Model Assumptions:** While some articles touch on the assumptions of linear regression models, there might be a lack of in-depth exploration of the assumptions' implications and potential challenges. Understanding and addressing these assumptions are crucial for the reliability of predictions.

5. Sparse Coverage of Advanced Techniques: The literature review emphasizes the application and basics of linear regression, but there is limited coverage of more advanced techniques like regularization (Lasso, Ridge), feature selection, and handling multicollinearity. These advanced methods are essential for improving model performance and robustness.

6. Insufficient Discussion on Model Comparison: Although there is a mention of comparing linear regression with machine learning models, there is limited discussion on the criteria for comparison and the situations where linear regression might be more suitable. A deeper exploration of the strengths and limitations of linear regression compared to other methods is needed.

7. Temporal Gap in Literature: Some of the referenced studies are a few years old, and there may be recent advancements or changes in the field of predictive modelling with linear regression. The literature review may not fully capture the most up-to-date methodologies, tools, or challenges.

8. Limited Exploration of Interpretability: While the simplicity and interpretability of linear regression are mentioned, there is a need for a more detailed discussion on how model outputs are interpreted in different application areas. Clear interpretation is crucial for stakeholders to make informed decisions based on the predictions.

9. Application Domain Bias: The literature review appears to focus on specific application domains, such as clinical settings, financial forecasting, and maintenance forecasting. This may create a bias, and the gaps in literature for other application areas of linear regression may not be adequately addressed.

10. Scarcity of Real-world Implementation Examples: The literature review lacks real-world implementation examples or case studies that demonstrate the practical challenges and

successes of applying predictive modelling with linear regression in diverse settings. Real-world scenarios provide insights beyond theoretical discussions.

## 3.1 The research strategy

The research strategy for predictive modelling with linear regression typically aligns with quantitative methods. Quantitative research involves collecting and analyzing numerical data to identify patterns, relationships, and patterns. Predictive modelling with linear regression from multiple angles to measure the following reasons:

1. Numerical Relationships: Linear regression models aim to capture numerical relationships between independent and dependent variables. These relationships are expressed through coefficients in the model equation, which are estimated based on quantitative data.

2. Statistical Analysis: Linear regression is a statistical technique that includes parameter estimation, hypothesis testing, and statistical significance testing. Various methods such as hypothesis testing and data analysis are required for model validation and interpretation.

3. Measurement and Precision: Quantitative research allows for precise measurement and quantification of variables. Linear regression models require numerical input to estimate the coefficients accurately and make predictions with a certain level of precision.

4. Prediction and Generalization: The most objective of predictive modelling is to form exact expectations, frequently on modern or inconspicuous information. Quantitative strategies such as straight relapse give a scientific approach to modelling that generalizes well to modern perceptions.

5. Data Analysis Techniques: Quantitative inquire about includes a few sorts of information examination, counting clear measurements, inferential insights, and regression analysis. Linear regression is a special regression analysis method designed for multiple data sets.

6. Replicability and Reliability: Quantitative methods contribute to the replicability and reliability of research findings. Linear regression models, when based on sound quantitative methods, can be replicated by other researchers using similar datasets to verify and validate the results.

7. Objective Measurement: Quantitative research emphasizes objective measurement, reducing the impact of personal biases. Linear regression models derive their parameters from data, providing an objective and systematic way to analyze relationships between variables.

8. Statistical Software Utilization: Quantitative methods often involve the use of statistical software for data analysis. Linear regression models are implemented and estimated using statistical software packages, which facilitate the handling of large datasets and complex calculations.

## 3.2 CRISP-DM Methodology

Predictive modelling involves the design and process of building, training, and evaluating predictive models using linear regression techniques.
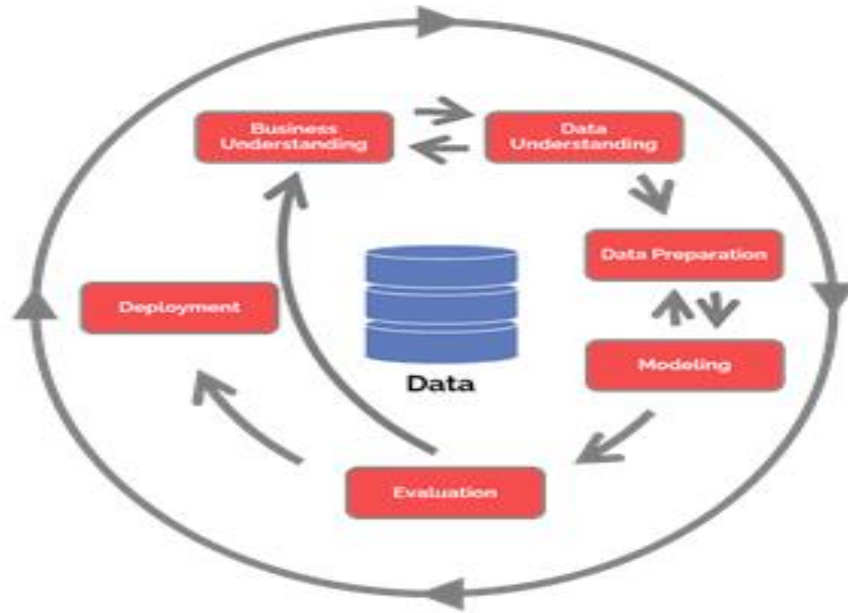
Figure 1: CRISP-DM Methodology

**STEP 1: Data Collection and Pre-processing:**

In this preliminary step, relevant information is collected from various sources, which may include historical data, statistical tests, or observations. Good information is important; therefore, use data cleaning procedures to resolve missing values, anomalies, and inconsistencies. This ensures that the data set is reliable and suitable for analysis. Additionally, specific choices or architectures should be identified, and useful predictions developed that increase the predictive power of the model.

**STEP 2: Exploratory Data Analysis (EDA):**

EDA involves a comprehensive review of the data set. This step involves visualizing the distribution of the data through histograms, box plots, and scatter plots. It also explores correlations between variables to identify potential patterns and relationships. EDA can reveal outliers, anomalies, or trends within the data, providing valuable insights that inform subsequent modelling decisions.

**STEP 3: Model Selection:**

Choosing the appropriate linear regression model depends on the problem at hand and the characteristics of the data. Simple linear regression is used when there is an independent variable; Multiple regression lines are used when there are many different predictors. When dealing with multiple variables (correlation between variables) or when adjustment is required to avoid overfitting, multiple variables are taken into account, such as Ridge Regression or Lasso Regression. In this study, a decision tree and random forest also were selected.

**STEP 4: Data Splitting:**

To assess the execution of the demonstration, the dataset is isolated into three parts: training set, validation set and test set. The preparing handle is utilized to prepare the demonstration so that it can learn the relationship between the indicator factors and the target factors. The validation process is utilised to assess hyperparameter tuning, show choice, and whether the demonstration can generalize to concealed objects. These testing methods are reserved for final performance testing, which provides an unbiased measurement of the model's predictive ability.
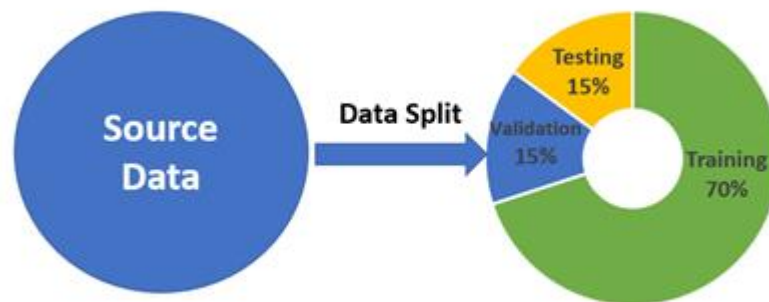


Figure 2: Dataset Splitting

**STEP 5: Model Training:**

This step involves fitting the chosen linear regression model to the training data. The model estimates the coefficients (weights) associated with each predictor, essentially learning how they contribute to predicting the target variable. Training may involve iterative processes, particularly

when using advanced regression techniques like ridge or lasso, which require fine-tuning of hyperparameters to optimize model performance.

**STEP 6: Hyper parameter Tuning:**

In advanced linear regression models such as ridge regression or lasso regression, hyperparameter tuning is important to achieve good results. This process involves changing hyperparameters (such as the energy constant) and using strategies such as competition to find the best combination that balances the model, unfairness, and diversity. The aim is to prevent structural conflicts (indirect pressure) and competition (differential pressure).

**STEP 7: Model Evaluation:**

Once the model is trained and tuned, it is evaluated using valid data. Common measurement metrics include Mean Square Error (MSE), Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and R-squared ($R^2$). These metrics provide a quantitative measure of how well the model predicts the target's outcomes and demonstrate accuracy and security. Interpretive Model: The horizontal linear regression model provides interpretation, allowing analysts to interpret the coefficients of independent variables.

**STEP 8: Deployment:**

Once the prototype is satisfied with the proof, it can be transferred to real implementation. This often involves creating APIs, integrating models into software systems, or creating user interfaces for stakeholders to access data and get predictions. Regular maintenance and updates may be required to ensure the model remains accurate as new data becomes available.

In summary, the research methodology for predictive modelling with linear regression involves a thorough exploration of existing literature, identification of gaps, selection of research philosophy, strategy, and a detailed data analysis plan. The limitations of secondary research, as identified through the literature review, highlight several areas where further investigation is needed.

The limitations include challenges in contextual understanding, generalization, data quality assumptions, exploration of model assumptions, sparse coverage of advanced techniques, insufficient discussion on model comparison, a temporal gap in literature, limited exploration of interpretability, application domain bias, and scarcity of real-world implementation examples. These limitations provide a foundation for the research, guiding the focus towards addressing these gaps and contributing new insights to the field of predictive modelling with linear regression.

The research strategy aligns with quantitative methods, emphasizing numerical relationships, statistical analysis, measurement and precision, prediction and generalization, data analysis techniques, replicability and reliability, objective measurement, and the utilization of statistical software. This strategy is well-suited for investigating predictive modelling with linear regression, as it allows for the quantitative analysis of relationships between variables.

The data analysis plan outlines a systematic approach, starting from data collection and pre processing, through exploratory data analysis (EDA), model selection, data splitting, model training, hyperparameter tuning, model evaluation, interpretive model, and deployment. This plan provides a roadmap for building, training, and evaluating predictive models using linear regression techniques.

# Chapter 4- Data Analysis

## 4.1 Dataset Description

The dataset contains detailed healthcare information for 10,000 patients. The columns cover various aspects of patient care, including demographics, medical conditions, treatment, and billing. Each column provides specific information about the patient, their admission, and the healthcare services provided, making this dataset suitable for various data analysis and modelling tasks in the healthcare domain. Here's a brief explanation of each column in the dataset -

- Name: This column represents the name of the patient associated with the healthcare record.
- Age: The age of the patient at the time of admission, expressed in years.
- Gender: Indicates the gender of the patient, either "Male" or "Female."
- Blood Type: The patient's blood type, which can be one of the common blood types (e.g., "A+", "O-", etc.).
- Medical Condition: This column specifies the primary medical condition or diagnosis associated with the patient, such as "Diabetes," "Hypertension," "Asthma," and more.
- Date of Admission: The date on which the patient was admitted to the healthcare facility.
- Doctor: The name of the doctor responsible for the patient's care during their admission.
- Hospital: Identifies the healthcare facility or hospital where the patient was admitted.
- Insurance Provider: This column indicates the patient's insurance provider, which can be one of several options, including "Aetna," "Blue Cross," "Cigna," "UnitedHealthcare," and "Medicare."
- Billing Amount: The amount of money billed for the patient's healthcare services during their admission. This is expressed as a floating-point number.

- Room Number: The room number where the patient was accommodated during their admission.

- Admission Type: Specifies the type of admission, which can be "Emergency," "Elective," or "Urgent," reflecting the circumstances of the admission.

- Discharge Date: The date on which the patient was discharged from the healthcare facility, based on the admission date and a random number of days within a realistic range.

- Medication: Identifies a medication prescribed or administered to the patient during their admission. Examples include "Aspirin," "Ibuprofen," "Penicillin," "Paracetamol," and "Lipitor."

- Test Results: Describes the results of a medical test conducted during the patient's admission. Possible values include "Normal," "Abnormal," or "Inconclusive," indicating the outcome of the test.

## 4.1.1 Summary of the Dataset

```
##      Name                Age              Gender              BloodType
##  Length:10000      Min.    :18.00    Length:10000      Length:10000
##  Class :character  1st Qu.:35.00    Class :character  Class :character
##  Mode  :character  Median :52.00    Mode  :character  Mode  :character
##                    Mean    :51.45
##                    3rd Qu.:68.00
##                    Max.    :85.00
##  MedicalCondition  DateofAdmission      Doctor              Hospital
##  Length:10000      Length:10000      Length:10000      Length:10000
##  Class :character  Class :character  Class :character  Class :character
##  Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##
##
##  InsuranceProvider BillingAmount       RoomNumber      AdmissionType
##  Length:10000      Min.    : 1000    Min.    :101.0    Length:10000
##  Class :character  1st Qu.:13507    1st Qu.:199.0    Class :character
##  Mode  :character  Median :25258    Median :299.0    Mode  :character
##                    Mean    :25517    Mean    :300.1
##                    3rd Qu.:37734    3rd Qu.:400.0
##                    Max.    :49996    Max.    :500.0
##  DischargeDate      Medication          TestResults
##  Length:10000      Length:10000      Length:10000
##  Class :character  Class :character  Class :character
##  Mode  :character  Mode  :character  Mode  :character
##
```

## 4.2 Exploratory Data Analysis (EDA)

Exploratory data analysis (EDA) in machine learning is a crucial initial phase in which data scientists gain insights into the dataset through summarized statistics and visual representations. EDA is about understanding the distribution of the data, identifying patterns, outliers and anomalies, and testing the underlying assumptions before applying machine learning models. This process not only helps to uncover the internal structure of the data, but also to detect errors or inconsistencies that could affect the performance of the model. Techniques such as creating histograms, box plots, scatter plots and correlation matrices are often used to visualize the relationships between variables and better understand the characteristics of the data. This fundamental step ensures that the predictive models are built on a solid foundation, which improves their reliability and interpretability when solving real-world problems. Checking head values allows viewing the first few rows to understand the data and how data is imported.
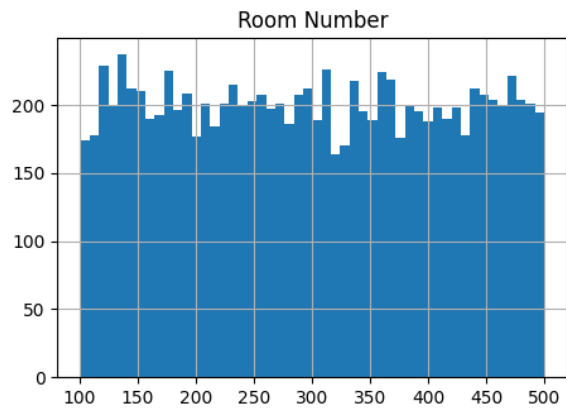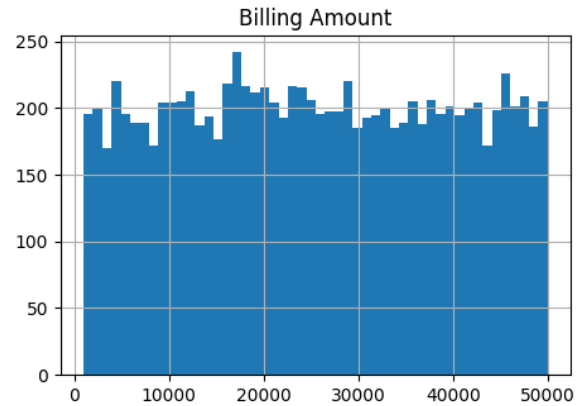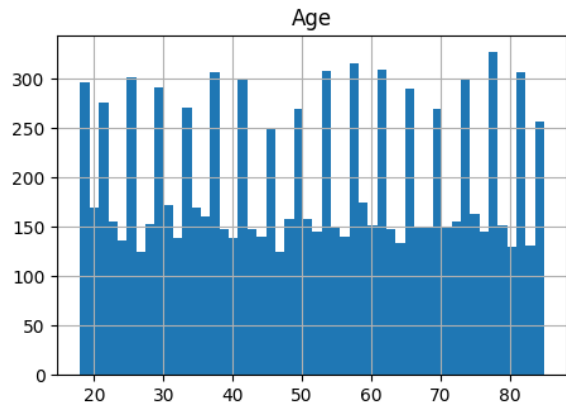
### 4.2.1 Descriptive Statistics

The descriptive statistics provide a summary of the central tendency, dispersion, and shape of the distribution of the variables in the dataset. Here's an interpretation of the descriptive statistics for the given variables Age, Billing Amount and Room Number.
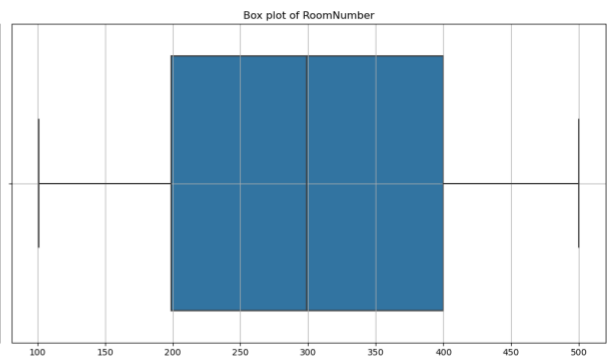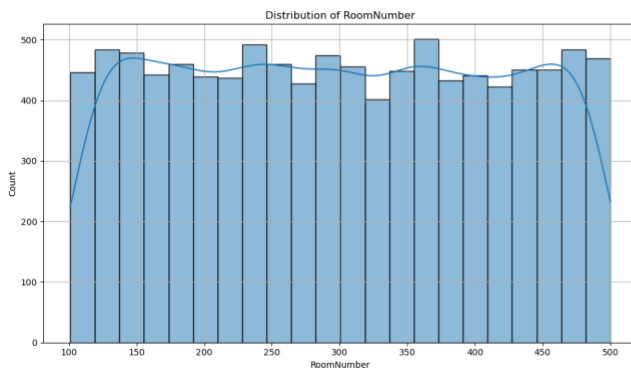
Table 1: Descriptive Analysis

|        | Age          | Billing Amount | Room Number   |
|--------|--------------|----------------|---------------|
| count  | 10000.000000 | 10000.000000   | 10000.000000  |
| mean   | 51.452200    | 25516.806778   | 300.082000    |
| std    | 19.588974    | 14067.292709   | 115.806027    |
| min    | 18.000000    | 1000.180837    | 101.000000    |
| 25%    | 35.000000    | 13506.523967   | 199.000000    |
| 50%    | 52.000000    | 25258.112566   | 299.000000    |
| 75%    | 68.000000    | 37733.913727   | 400.000000    |
| max    | 85.000000    | 49995.902283   | 500.000000    |

The age distribution suggests a wide range of patients, from young adults to elderly individuals. The billing amounts have a large spread, indicating varying levels of care complexity and cost. The room numbers seem to cover a broad range of hospital rooms or facilities.



In the below figures for each feature we check the presence of outliers using boxplot analysis
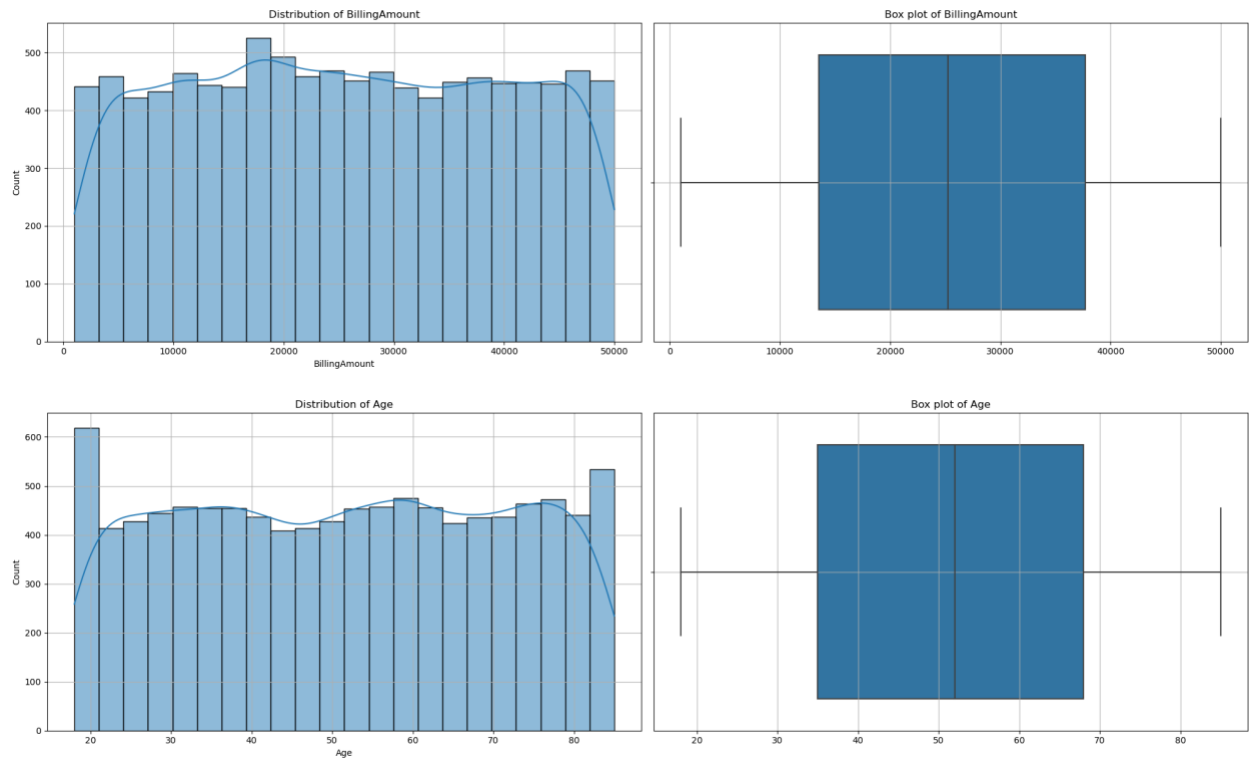
Figure 3: Data Distribution of Attributes (Age, Billing Amount and Room Number)

## 4.2.2 Correlation Matrix



## 4.2.3  Categorical Features

**- Gender Feature:**

Figure 4: Gender Distribution

**Blood Feature**

**Medical Medication Feature**



Figure 6: Medical Condition

**Insurance Provider Feature**

Figure 7: Insurance Provider

**Admission Type Feature**

**Medication Feature**



Figure 9:  Medication

**Test Results Feature:**



Figure 10: Test Results

## 4.3 Dataset Visualization

### 4.3.1 What is the distribution of medical conditions between genders



Figure 11: Distribution Medical Conditions between Genders

The above figure shows that female have higher share in majority of the medical conditions. Cancer is leading medical condition among female whereas Hypertension is leading cause among male.

### 4.3.2 What is the total billing amount for top 10 hospitals?



Figure 12: Billing Amount for Top 10 Hospitals

The above figure shows that Smith and Sons has the highest billing amount of 475639

### 4.3.3 Who are the top 20 doctors by billing amount?



Figure 13: Billing Amount for top 20 Doctors

The above figure shows that Doctor Michael Johnson is the top doctor with a total billing amount of 181576.

### 4.3.4 Which medical condition generated the highest average billing among male and female?



Figure 14: Average Billing Amount

### 4.3.5 What is the distribution of different age groups?



Figure 15: Age Groups Distribution Among Patients

The above figure shows that the maximum number of patients are in the range of 18-30.

# Chapter 5: Data Modeling

Predictive analytics is a supervised machine learning method that help understand patterns of data and make predictions. Descriptive data and rankings are just a small part of predictive analytics. It relies on various models to draw conclusions from the encountered data. To predict future trends, these models use machine learning algorithms to evaluate historical and current data. Predictive analytics is useful for evaluating business decisions. This is because good decision making involves understanding the consequences and making predictions about how the project, team, environment, or other environment will perform. In this study three different machine learning algorithms are selected:

- Linear Regression Model
- Decision Tree Model
- Random Forest Model

Prediction is a critical portion of information analysis. Predictive analytics could be a way to foresee future patterns based on current or historical information. Therefore, businesses will be able to predict future information. It may use different methods, but some of the best models use machine learning. Analysis includes many types of predictive analysis models. Most are regression models designed to determine the relationship between two or more variables. By analyzing the connections between these variables, they can help predict the value of the unknown variable as the value of the variable.

## 5.1 Linear Regression Model

Linear regression is a basic statistical and machine learning technique used to model the linear relationship between a dependent variable and one or more independent variables. By fitting a linear equation to observed data, linear regression allows the value of the dependent variable to be predicted based on the values of the independent variables. This model is used in many areas for predictive analysis, trend forecasting and determining the strength of predictors.

## 5.1.1 Mathematical Formulation

The linear regression model can be expressed by the following equation:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n + \epsilon$$

Where:

Y is the dependent variable

$X_1, X_2, \ldots, X_n$ are the independent variables

$\beta_0$ is the intercept term of the model

$\beta_1, \beta_2, \ldots, \beta_n$ are the coefficients that represent the weight of each independent variable

$\epsilon$ represents the error term, accounting for the variability in Y not explained by the independent variable.

## 5.1.2. Correlation Matrix



Figure 16: Correlation Matrix

From the above correlation figure derived the following results:

- There is a slight correlation between age and test results. Even there is a positive correlation between age and days hospitalized. There is a negative correlation between age and billing amount.

- Billing amount has negative correlation with test results and days hospitalized. Billing amount has a positive correlation with admission type and gender.

- Days hospitalized has negative correlation with every attribute except age.

- Except for billing amount attribute, gender has a negative correlation with all other attributes.

- Admission type and test results have a negative correlation.

### 5.1.3 Linear Regression Model

When running the linear regression model, it is found that all p-values are statistically in significant and it is difficult to perform the predictive analysis. We have examined several possible combinations between the characteristics, but the result remains the same.

```
# Run the Linear model and save it as 'mod'
modGender <- lm(Age ~ Gender, data = dataset)
modBloodType <- lm(Age ~ BloodType, data = dataset)
modMedicalCondition <- lm(Age ~ MedicalCondition, data = dataset)
modMedication <- lm(Age ~ Medication, data = dataset)
modTestResults <- lm(Age ~ TestResults, data = dataset)
modAdmissionType <- lm(Age ~ AdmissionType, data = dataset)
```

# Results Summary for each Linear Regression Model

```
modGender

##
## Call:
## lm(formula = Age ~ Gender, data = dataset)
##
## Coefficients:
## (Intercept)   GenderMale
##    51.6085      -0.3173

summary(modGender)

##
## Call:
## lm(formula = Age ~ Gender, data = dataset)
##
## Residuals:
##    Min     1Q  Median    3Q    Max
## -33.608 -16.608  0.392 16.709 33.709
##
## Coefficients:
##            Estimate Std. Error t value Pr(>|t|)
## (Intercept) 51.6085    0.2750  187.68  <2e-16 ***
## GenderMale  -0.3173    0.3918   -0.81   0.418
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 19.59 on 9998 degrees of freedom
## Multiple R-squared:  6.559e-05,  Adjusted R-squared:  -3.443e-05
## F-statistic: 0.6558 on 1 and 9998 DF,  p-value: 0.4181
```

```
##
## Call:
## lm(formula = Age ~ BloodType, data = dataset)
##
## Coefficients:
##  (Intercept)    BloodTypeA+   BloodTypeAB-   BloodTypeAB+   BloodTypeB-
##     50.75525      1.38496       1.00710       0.59769       0.08101
##  BloodTypeB+   BloodTypeO-    BloodTypeO+
##     1.47385      0.38301       0.64379

summary(modBloodType)

##
## Call:
## lm(formula = Age ~ BloodType, data = dataset)
##
## Residuals:
##    Min     1Q  Median    3Q    Max
## -34.229 -17.138  0.245 16.860 34.245
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 50.75525   0.55674  91.165  <2e-16 ***
## BloodTypeA+  1.38496   0.78687   1.760  0.0784 .
## BloodTypeAB- 1.00710   0.78162   1.288  0.1976
## BloodTypeAB+ 0.59769   0.78421   0.762  0.4460
## BloodTypeB-  0.08101   0.78514   0.103  0.9178
## BloodTypeB+  1.47385   0.78640   1.874  0.0609 .
## BloodTypeO-  0.38301   0.78640   0.487  0.6262
## BloodTypeO+  0.64379   0.78577   0.819  0.4126
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 19.59 on 9992 degrees of freedom
## Multiple R-squared:  0.0006974,  Adjusted R-squared:  -2.69e-06
## F-statistic: 0.9962 on 7 and 9992 DF,  p-value: 0.4318
```

```
## 
## Call:
## lm(formula = Age ~ MedicalCondition, data = dataset)
## 
## Coefficients:
##                   (Intercept)       MedicalConditionAsthma
##                      51.53091                     -0.08536
##        MedicalConditionCancer      MedicalConditionDiabetes
##                       0.05277                      0.27131
## MedicalConditionHypertension       MedicalConditionObesity
##                      -0.79335                      0.10300
summary(modMedicalCondition)

## 
## Call:
## lm(formula = Age ~ MedicalCondition, data = dataset)
## 
## Residuals:
##     Min     1Q  Median     3Q    Max
## -33.802 -16.738   0.366  17.198  34.262
## 
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   51.53091    0.48230 106.845  <2e-16 ***
## MedicalConditionAsthma        -0.08536    0.67625  -0.126   0.900
## MedicalConditionCancer         0.05277    0.67674   0.078   0.938
## MedicalConditionDiabetes       0.27131    0.68490   0.396   0.692
## MedicalConditionHypertension  -0.79335    0.67822  -1.170   0.242
## MedicalConditionObesity        0.10300    0.68437   0.151   0.880
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 19.59 on 9994 degrees of freedom
## Multiple R-squared:  0.0003009,  Adjusted R-squared:  -0.0001993
## F-statistic: 0.6016 on 5 and 9994 DF,  p-value: 0.6988
```

```
modMedication

## 
## Call:
## lm(formula = Age ~ Medication, data = dataset)
## 
## Coefficients:
##          (Intercept)     MedicationIbuprofen        MedicationLipitor
##              51.2134                  0.1251                   0.1112
## MedicationParacetamol     MedicationPenicillin
##               0.3518                   0.5899
summary(modMedication)

## 
## Call:
## lm(formula = Age ~ Medication, data = dataset)
## 
## Residuals:
##     Min     1Q  Median     3Q    Max
## -33.803 -16.803   0.435  16.787  33.787
## 
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)            51.2134     0.4416 115.964  <2e-16 ***
## MedicationIbuprofen     0.1251     0.6239   0.201   0.841
## MedicationLipitor       0.1112     0.6209   0.179   0.858
## MedicationParacetamol   0.3518     0.6250   0.563   0.574
## MedicationPenicillin    0.5899     0.6162   0.957   0.338
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 19.59 on 9995 degrees of freedom
## Multiple R-squared:  0.0001178,  Adjusted R-squared:  -0.0002824
## F-statistic: 0.2943 on 4 and 9995 DF,  p-value: 0.8818
```

```
modTestResults

## 
## Call:
## lm(formula = Age ~ TestResults, data = dataset)
## 
## Coefficients:
##          (Intercept)  TestResultsInconclusive       TestResultsNormal
##              51.3721                  -0.1796                  0.4253

summary(modTestResults)

## 
## Call:
## lm(formula = Age ~ TestResults, data = dataset)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max
## -33.797 -16.797   0.628  16.807  33.807
## 
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)               51.3721     0.3332 154.168   <2e-16 ***
## TestResultsInconclusive   -0.1796     0.4776  -0.376    0.707
## TestResultsNormal          0.4253     0.4780   0.890    0.374
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 19.59 on 9997 degrees of freedom
## Multiple R-squared:  0.0001648,  Adjusted R-squared:  -3.523e-05
## F-statistic: 0.8239 on 2 and 9997 DF,  p-value: 0.4387
```

```
modAdmissionType

## 
## Call:
## lm(formula = Age ~ AdmissionType, data = dataset)
## 
## Coefficients:
##            (Intercept)  AdmissionTypeEmergency       AdmissionTypeUrgent
##                51.4130                 -0.0385                    0.1538

summary(modAdmissionType)

## 
## Call:
## lm(formula = Age ~ AdmissionType, data = dataset)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max
## -33.567 -16.567   0.587  16.625  33.625
## 
## Coefficients:
##                        Estimate Std. Error t value Pr(>|t|)
## (Intercept)             51.4130     0.3441  149.43   <2e-16 ***
## AdmissionTypeEmergency  -0.0385     0.4820   -0.08    0.936
## AdmissionTypeUrgent      0.1538     0.4812    0.32    0.749
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 19.59 on 9997 degrees of freedom
## Multiple R-squared: 1.82e-05,   Adjusted R-squared:  -0.0001819
## F-statistic: 0.09097 on 2 and 9997 DF,  p-value: 0.913
```

## 5.2 Decision Tree Model

A decision tree is a supervised machine learning algorithm used for both classification and regression tasks. It is a fundamental tool in the field of machine learning and is widely used because of its simplicity and interpretability. Decision trees are particularly useful when you need to make decisions or predictions based on a set of input features. An overview of decision trees, their construction and their most important concepts are presented.

### 5.2.1  Structure of decision trees:

A decision tree is a hierarchical structure consisting of nodes, where each node represents a decision or test for a feature.

The nodes in a decision tree can be divided into three types: Root nodes, internal nodes and leaf nodes (also known as terminal nodes).

- The root node is the top node from which the tree branches out.
- Internal nodes represent decisions or feature tests and have one or more subordinate nodes.
- The leaf nodes contain the final decision or predicted output.

### 5.2.2  Feature selection

At each internal node, a decision tree algorithm selects a feature and a corresponding threshold (for continuous features) that optimally splits the data into two or more child nodes.

The goal is to maximize the purity or homogeneity of the data in each child node. Various criteria such as Gini impurity, entropy or mean square error can be used to measure purity.

### 5.2.3. Splitting criteria

**Gini impurity** measures the probability that a randomly selected element would be misclassified if it were randomly classified according to the distribution of labels in the node.

**Entropy** measures the degree of disorder or randomness in the labels of the data in a node.

**Information gain** measures the reduction in entropy or increase in purity achieved by a particular partitioning.

## 5.2.4. Creating a decision tree

Decision trees are built recursively. Starting from the root node, the algorithm selects the best feature and the best threshold for splitting the data.

The process continues for each subordinate node until a termination criterion is met. Common termination criteria include the maximum tree depth, the minimum number of samples per leaf or if all data in a node belongs to the same class or has a low variance.

Pruning techniques can be used to reduce the size of the tree and avoid overfitting.

## 5.2.5. Prediction analysis

To make a prediction, start at the root node and move down the tree, following the decisions made at each internal node until you reach a leaf node. The class or value associated with the leaf node is the final prediction. The figure below shows the decision tree for our analysis

**Decision Tree**



Figure 17: Decision Tree

## 5.3 Random Forest Model

Random Forest (RF) is a versatile and powerful machine learning algorithm that creates a large number of decision trees at training time and outputs the class corresponding to the mode of the classes (classification) or the mean prediction (regression) of the individual trees. RF is an ensemble learning method, i.e. it combines the predictions of multiple models to improve generalizability and robustness to a single estimator.

### 5.3.1  Background

A decision tree is a simple decision-making tool that divides the data into branches at decision points. Each branch represents a choice between several alternatives, and each leaf node represents a classification or decision. Decision trees are easy to understand and interpret but can become complex and prone to overfitting if they become too deep.

### 5.3.2  Ensemble learning

Ensemble learning is a technique that combines the predictions of multiple machine learning algorithms to make more accurate predictions than any single model. Random forest is a type of ensemble learning method that uses the bagging approach (bootstrap aggregation) to create an ensemble of decision trees trained on different subsets of the training dataset.

### 5.3.3  How random forests work:

**Bootstrap Aggregation (Bagging):** RF creates multiple decision trees using bootstrap datasets of the original data. A bootstrap dataset is a randomly selected dataset with replacement, i.e. it may contain duplicate rows.

**Selection of characteristics:** When creating each tree, RF randomly selects a subset of features at each split to increase diversity between trees. This process helps to reduce variance and avoid overfitting.

**Aggregation:** For classification tasks, the mode of the classes predicted by the individual trees is used as the final prediction. For regression tasks, the average of the predictions is used.

## 5.4   Performance Analysis

### 5.4.1   Metrics

| Predictive Model | Accuracy (Training) | Accuracy (Test) |
|---|---|---|
| Decision Tree | 0.35 | 0.34 |
| Random Forest | 0.35386 | 0.344 |

### 5.4.2  Performance of Decision Tree for training data

```
              precision    recall  f1-score   support

    Abnormal       0.35      1.00      0.51      2593
Inconclusive       1.00      0.00      0.01      2466
      Normal       1.00      0.00      0.00      2441

    accuracy                           0.35      7500
   macro avg       0.78      0.34      0.18      7500
weighted avg       0.77      0.35      0.18      7500
```

## 5.4.2 Performance of Decision Tree for testing data

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| Abnormal     | 0.34      | 1.00   | 0.51     | 863     |
| Inconclusive | 0.00      | 0.00   | 0.00     | 811     |
| Normal       | 0.00      | 0.00   | 0.00     | 826     |
|              |           |        |          |         |
| accuracy     |           |        | 0.34     | 2500    |
| macro avg    | 0.11      | 0.33   | 0.17     | 2500    |
| weighted avg | 0.12      | 0.34   | 0.18     | 2500    |

## 5.4.3 Performance of Random Forest

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| Abnormal     | 0.35      | 1.00   | 0.51     | 863     |
| Inconclusive | 0.00      | 0.00   | 0.00     | 811     |
| Normal       | 0.00      | 0.00   | 0.00     | 826     |
|              |           |        |          |         |
| accuracy     |           |        | 0.34     | 2500    |
| macro avg    | 0.12      | 0.33   | 0.17     | 2500    |
| weighted avg | 0.12      | 0.34   | 0.18     | 2500    |

# Chapter 6 Discussion

The research sets out on a journey into the domain of predictive modelling with linear regression within the setting of healthcare. In later years, predictive modelling has risen as a basic instrument over different spaces, advertising the potential for educated decision-making and precise predictions. This examination particularly centers on the application of linear regression, a broadly utilized strategy in predictive analytics, to reveal designs and upgrade the exactness of figures inside the energetic scene of healthcare.

As the healthcare industry experiences a transformative advancement, the integration of progressed analytics has gotten to be basic to form educated choices and accomplish superior quiet results. Predictive modelling, particularly when seen through the focal point of linear regression, stands out as an effective device to disentangle complex designs and progress the accuracy of wellbeing forecasts. The research points to explore the elemental concepts of linear regression and its application in healthcare settings, emphasizing the need for fastidious thought to guarantee the unwavering quality of predictive models.

This exploration speaks to the crossing point of information science and healthcare, where linear regression offers a one of a kind point of view for recognizing and anticipating medicines, optimizing asset allotment, and eventually upgrading the quality of understanding care. The idea looks to prepare healthcare experts with the information and abilities required to explore the complexities of predictive modelling, cultivating a bridge between historical information and educated decision-making.

The presentation of a level predictive modelling division in healthcare sets the organisation for contributing to a future where data-driven bits of knowledge engage specialists to anticipate, diminish, and move forward care conveyance guidelines. By adjusting with the overarching

objective of improving decision-making forms and optimizing patient results, this investigates endeavours to supply healthcare experts with important devices and information to tackle the potential of linear regression models.

Moving forward, the discussion addresses the pressing problem of the healthcare industry's need for precise and timely predictions. The complexity of clinical information and the multitude of factors involved often pose challenges in achieving accurate forecasts. Linear regression, with its simplicity and interpretability, emerges as a potential solution to this problem. However, challenges such as the need for robust models, consideration of diverse variables, and adherence to assumptions like linearity and independence need to be addressed.

The research aims to tackle these challenges and demonstrate the efficacy of linear regression for predictive modelling in healthcare. The emphasis is on developing robust and interpretable linear regression models using high-quality healthcare data. The research moreover centers on approving and assessing the execution of these models in real-world clinical settings, pointing to bridge the crevice between historical information and educated decision-making. The extreme objective is to enable organizations and decision-makers to create choices based on data, move forward asset assignment, diminish instability, and increase overall execution and competitiveness within the healthcare segment.

As the discussion unfolds, the research aim and objectives come into focus. The primary aim is to examine and evaluate the viability of predictive modelling using linear regression within the healthcare space. The overarching objective is to contribute valuable insights that can improve decision-making processes and enhance the accuracy of predictions in healthcare scenarios. The research aims to assess the application of linear regression models, emphasizing their capacity to uncover patterns and provide precise figures within the complex and dynamic landscape of healthcare.

The goals of the research incorporate analyzing the effect of preparatory information investigation, highlight choice, and show assessment on the unwavering quality of expectations. Also, the

research points to address challenges characteristic in linear regression, such as multicollinearity and overfitting, to guarantee the strength of predictive models in healthcare settings. Eventually, the objective is to prepare healthcare specialists with important data and devices that can be connected to tackle the potential of linear regression models for exact figures, contributing to the headway of data-driven decision-making in healthcare.

In conclusion, this investigation speaks to a critical endeavor to investigate the applications of predictive modelling with linear regression in healthcare. It addresses current challenges, sets clear goals, and points to supply viable bits of knowledge and instruments for healthcare experts. By doing so, the investigate looks for to contribute to the continuous advancement of data-driven decision-making in healthcare, where the predictive control of linear regression can be saddled to optimize patient results and upgrade the proficiency of healthcare frameworks

In a real-world context, linear regression finds extensive utility in business forecasting. Its simplicity and interpretation make it an easy-to-use tool for analysts and decision makers. For example, companies that analyze historical data on ad spend can use the regression line to predict future sales to determine budgets and cutbacks.

In conclusion, the practical application of linear regression in predictive modelling underscores its real-world utility. While its simplicity is a strength, addressing assumptions, multicollinearity, and employing advanced techniques for feature selection and regularization are essential steps to enhance the reliability and performance of linear regression models in diverse scenarios, including business forecasting.

# Chapter 7 Conclusion

In conclusion, predictive modelling, particularly utilizing linear regression, has emerged as a critical instrument across diverse areas, enabling informed decision-making and precise predictions. This research specifically focuses on the application of linear regression in healthcare, recognizing the imperative need for accurate and timely predictions in this complex domain. The intersection of data science and healthcare is explored, leveraging the simplicity and interpretability of linear regression to identify and predict treatments, optimize resource allocation, and ultimately enhance patient quality.

The exploration of predictive modelling in healthcare, especially with linear regression, is envisioned as a transformative journey. As we delve into the main concepts of linear regression and its application in medicine, the careful consideration required to ensure the reliability of the model becomes apparent. From the construction of linear models to the nuances of feature engineering and variable selection, the goal is to equip practitioners with the knowledge and skills needed to build predictive models that are both robust and stable.

The introduction of a horizontal predictive modelling division in healthcare signifies a commitment to contributing to a future where data-driven insights empower practitioners to predict, reduce, and improve care delivery standards. This research acknowledges the challenges within the healthcare sector, emphasizing the pressing need to address issues related to predictive modelling. The problem statement underlines the complexity of clinical information and the hurdles faced in achieving precise predictions in healthcare. Linear regression is positioned as a potential solution, given its simplicity, interpretability, and potential to fill existing gaps.

Moving forward, the research aims to address these challenges by developing robust and interpretable linear regression models using high-quality healthcare data.Validation and evaluation of these models in real-world clinical settings are significant steps, guaranteeing their

appropriateness and unwavering quality. The ultimate objective is to supply healthcare experts with the tools and information required to consistently coordinate predictive modelling into their decision-making forms. By doing so, the research aims to bridge the gap between historical information and educated decision-making, contributing to upgrade execution and competitiveness in different healthcare spaces.

The concluding paragraphs underscore the importance of linear regression as a widely used method and highlight its potential to provide statistical solutions to complex problems within the healthcare sector. The project's holistic approach aims not only to create predictive models but also to enable organizations and decision-makers to leverage information for improved resource allocation, reduced uncertainty, and overall increased performance. Through the strategic use of historical data, the project seeks to empower decision-makers to make informed choices, ultimately shaping a future where data-driven insights drive positive transformations in various domains.

## 7.1 Recommendations

Based on the findings and limitations identified in our research, as well as the broader context of predictive modelling with linear regression, here are some recommendations:

1. Comprehensive Model Assessment:
   - Perform a comprehensive assessment of the expectation show, taking under consideration different measurements such as mean square error (MSE), R-squared, and residual analysis.
   - Cross-validation techniques are utilized to demonstrate strength and generalizability to new data.
   - Perform affectability analysis to assess the effect of diverse suspicions on the performance model.

2. Feature Selection Strategies:
   - Experiment with different feature selection methods to identify the most influential variables and enhance model interpretability.
   - Consider incorporating domain knowledge to guide the selection of relevant features.
   - Regularly revisit and update feature selection strategies as new data becomes available or business priorities change.

3. Address Assumption Violations:
   - Develop strategies to handle violations of linear regression assumptions, especially in real-world scenarios where data may not perfectly adhere to these assumptions.
   - Explore robust regression strategies that are less delicate to exceptions and suspicion violations.
   - Conduct sensitivity analysis to evaluate the effect of suspicion deviations on the model's forecasts.

4. Continuous Model Improvement:
   - Embrace the mindset of continuous improvement, constantly updating and adjusting the predictive model based on more information.
   - Monitor performance patterns over time and carefully consider potential damage to forecast accuracy.
   - Exploring new features or incorporating predictors can increase the predictive power of the model.

5. External Validation:
   - Validate the predictive model on external datasets from different sources or time periods to ensure its generalizability.

- Collaborate with other researchers or organizations to obtain diverse datasets for validation purposes.
- Publish results of external validation to provide transparency and demonstrate the model's reliability across various contexts.

## 7.2   Future Work

To address the aforementioned limitations and contribute to the advancement of predictive modelling with linear regression, the following suggestions are proposed:

- Advanced Assumption Handling: Develop and evaluate techniques to handle violations of linear regression assumptions. This could involve the use of robust regression methods, transformation of variables, or incorporating machine learning algorithms that are less sensitive to assumption violations.
- Enhanced Data Quality Measures: Investigate innovative approaches for ensuring data quality, including automated data cleaning algorithms, outlier detection methods, and strategies for handling missing data. Robustness checks should be performed to assess model sensitivity to variations in data quality.
- Model Application: Apply the demonstration to real-world circumstances, such as optimising publicising budgets, estimating item request, or anticipating advertised conditions. Investigate the benefits and challenges of utilizing this show in several settings.
- Dynamic Modelling: Explore the integration of dynamic elements into the model, such as time series analysis or seasonal trends, to identify changes in the market or changes in customer time behaviour.
- External Influences: Consider engaging in methods other than spending money on advertising, such as marketing, contests, or social media to create a more comprehensive picture of your sales forecast.
- Model Comparison: Compare the performance of linear regression with other prediction models under different data sets and research questions. Weigh the pros and cons of each method for a particular project.

- Explainable AI: Explore integrating descriptive intelligence (XAI) techniques with complex models to improve description and provide insight into model decision-making. This could be particularly valuable for stakeholder acceptance and trust in predictions.

# Bibliography

1. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning with Applications in R. Springer.

2. Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). Introduction to Linear Regression Analysis. John Wiley & Sons.

3. Kutner, M. H., Nachtsheim, C. J., Neter, J., & Li, W. (2004). Applied Linear Statistical Models. McGraw-Hill.

4. Fox, J. (2015). Applied Regression Analysis and Generalized Linear Models. Sage Publications.

5. Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer.

6. Faraway, J. J. (2014). Linear Models with R. CRC Press.

7. Cohen, J., Cohen, P., West, S. G., & Aiken, L. S. (2013). Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences. Routledge.

8. Chatterjee, S., & Hadi, A. S. (2015). Regression Analysis by Example. John Wiley & Sons.

9. Gelman, A., & Hill, J. (2006). Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press.

10. Neter, J., Wasserman, W., & Kutner, M. H. (1989). Applied Linear Statistical Models: Regression, Analysis of Variance, and Experimental Designs. Irwin.

11. Neter, J., Wasserman, W., & Kutner, M. H. (1989). Applied Linear Statistical Models: Regression, Analysis of Variance, and Experimental Designs. Irwin.

12. Rencher, A. C., & Schaalje, G. B. (2008). Linear Models in Statistics. John Wiley & Sons.

13. Kleinbaum, D. G., Kupper, L. L., Nizam, A., & Rosenberg, E. S. (2013). Applied Regression Analysis and Other Multivariable Methods. Cengage Learning.

14. Cook, R. D., & Weisberg, S. (1982). Residuals and Influence in Regression. CRC Press.

15. Draper, N. R., & Smith, H. (1998). Applied Regression Analysis. John Wiley & Sons.

16. Aiken, L. S., & West, S. G. (1991). Multiple Regression: Testing and Interpreting Interactions. Sage Publications.

17. Cribari-Neto, F., & Zeileis, A. (2010). Beta Regression in R.

18. Hardle, W., & Simar, L. (2012). Applied Multivariate Statistical Analysis.

19. Seber, G. A. F., & Wild, C. J. (2003). Nonlinear Regression.

20. Caroll, R. J., & Ruppert, D. (1996). Transformation and Weighting in Regression.

21. Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). Bayesian Data Analysis.

22. McCulloch, C. E., & Searle, S. R. (2001). Generalized, Linear, and Mixed Models.

23. West, M., Harrison, J., & Migon, H. (1985). Bayesian Forecasting and Dynamic Models.

24. Dobson, A. J. (2002). An Introduction to Generalized Linear Models.

25. Harrell Jr, F. E. (2015). Regression Modeling Strategies.

26. Long, J. S. (1997). Regression Models for Categorical and Limited Dependent Variables.

27. Verbeek, M. (2008). A Guide to Modern Econometrics.

28. Gourieroux, C., Monfort, A., & Trognon, A. (1984). Pseudo-Maximum Likelihood Methods: Applications to Poisson Models.

29. Judge, G. G., Griffiths, W. E., Hill, R. C., Lutkepohl, H., & Lee, T. C. (1985). The Theory and Practice of Econometrics.

30. McCulloch, C. E. (2011). Generalized Linear Mixed Models.

31. Pinheiro, J. C., & Bates, D. M. (2000). Mixed-Effects Models in S and S-PLUS.

32. Zuur, A., Ieno, E. N., & Smith, G. M. (2007). Analysing Ecological Data.

33. Verbeek, M. (2012). A Guide to Modern Econometrics.

34. Hosmer Jr, D. W., & Lemeshow, S. (2000). Applied Logistic Regression.

35. Cameron, A. C., & Trivedi, P. K. (2013). Regression Analysis of Count Data.

36. Allison, P. D. (2012). Logistic Regression Using SAS: Theory and Application.

37. Gelman, A., & Pardoe, I. (2006). Bayesian Measures of Explained Variance and Pooling in Multilevel (Hierarchical) Models.

38. Dobson, A. J., & Barnett, A. G. (2008). An Introduction to Generalized Linear Models, Third Edition.

39. Gelman, A., & Shalizi, C. R. (2013). Philosophy and the Practice of Bayesian Statistics.

40. Pregibon, D. (1981). Logistic Regression Diagnostics.

41. McCullagh, P., & Nelder, J. A. (1989). Generalized Linear Models (Monographs on Statistics and Applied Probability).

42. Zeileis, A., Kleiber, C., & Jackman, S. (2008). Regression Models for Count Data in R.

43. Everitt, B. S., & Hothorn, T. (2011). An Introduction to Applied Multivariate Analysis with R.

44. De Veaux, R. D., Velleman, P. F., & Bock, D. E. (2019). Stats: Data and Models.

45. Achen, C. H. (1982). Interpreting and Using Regression. Sage Publications.

46. Cohen, J., Cohen, P., West, S. G., & Aiken, L. S. (2003). Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences. Routledge.

47. Belsley, D. A., Kuh, E., & Welsch, R. E. (1980). Regression Diagnostics: Identifying Influential Data and Sources of Collinearity. John Wiley & Sons.

48. Frees, E. W. (2004). Longitudinal and Panel Data: Analysis and Applications in the Social Sciences. Cambridge University Press.

49. Draper, N. R., & Smith, H. (1998). Applied Regression Analysis (3rd ed.). John Wiley & Sons.

50. Gelman, A., & Hill, J. (2006). Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press.