# I2SL Statistical Learning

Nik Bear Brown

March 2024

# Contents

**16 Exploratory Data Analysis**                                            **31**

# Chapter 1

# Overview of Statistical Learning

## 1.1 Introduction to Statistical Learning

### 1.1.1 Definition and Scope

### 1.1.2 Historical Background

### 1.1.3 Importance and Applications in Various Fields

# Chapter 2

# Linear Regression

## 2.1  Introduction to Linear Regression

### 2.1.1  Definition and Importance

### 2.1.2  Historical Background

### 2.1.3  Applications in Various Fields

## 2.2  Linear Models for Regression

### 2.2.1  Theoretical Foundations

The Regression Equation

Assumptions Underlying Linear Regression Models

### 2.2.2  Simple Linear Regression

Estimating the Coefficients

Interpreting the Regression Coefficients

Assumptions of Simple Linear Regression

### 2.2.3  Multiple Linear Regression

Understanding Multiple Regression Outputs

The Use of Dummy Variables

Interactions Between Predictors

### 2.2.4  Assumptions of Linear Regression

Linearity

Homoscedasticity

Independence of Errors

Normal Distribution of Errors

Multicollinearity

## 2.3   Model Assessment and Validation

# Chapter 3

# Logistic Regression

## 3.1   Introduction to Logistic Regression

### 3.1.1   Definition and Overview

### 3.1.2   Comparison with Linear Regression

### 3.1.3   Applications in Various Fields

## 3.2   Theoretical Foundations of Logistic Regression

### 3.2.1   The Logistic Function

### 3.2.2   Odds and Log Odds

### 3.2.3   The Maximum Likelihood Estimation (MLE)

## 3.3   Binary Logistic Regression

### 3.3.1   Modeling Binary Outcomes

### 3.3.2   Interpreting the Coefficients

### 3.3.3   Assessing Model Fit and Accuracy

## 3.4   Assumptions of Logistic Regression

### 3.4.1   Requirement of Linearity in the Logit

### 3.4.2   Absence of Multicollinearity

### 3.4.3   Large Sample Size Requirement

## 3.5   Model Evaluation and Diagnostics

### 3.5.1   Confusion Matrix and Classification Accuracy

### 3.5.2   Receiver Operating Characteristic (ROC) Curve

### 3.5.3   Area Under the ROC Curve (AUC)

# Chapter 4

# Classification Techniques

## 4.1    Introduction to Classification Techniques

### 4.1.1    Definition and Importance

### 4.1.2    Overview of Classification in Machine Learning

### 4.1.3    Applications of Classification Techniques

## 4.2    Theoretical Foundations of Classification

### 4.2.1    Bayes Theorem and Decision Theory

### 4.2.2    The Concept of Decision Boundaries

### 4.2.3    Performance Metrics for Classification Models

## 4.3    Discriminant Analysis

### 4.3.1    Introduction to Discriminant Analysis

Historical Background

Basic Principles and Goals

### 4.3.2    Linear Discriminant Analysis (LDA)

Assumptions of LDA

Mathematical Formulation of LDA

Dimensionality Reduction with LDA

Multiclass Classification with LDA

### 4.3.3    Quadratic Discriminant Analysis (QDA)

Differences Between LDA and QDA

When to Use QDA Over LDA

Mathematical Formulation of QDA

### 4.3.4    Regularized Discriminant Analysis

# Chapter 5

# Resampling Methods

## 5.1   Introduction to Resampling Methods

### 5.1.1   Definition and Importance

### 5.1.2   Overview of Resampling in Statistical Analysis

### 5.1.3   Applications of Resampling Methods

## 5.2   Theoretical Foundations of Resampling Methods

### 5.2.1   Principles Behind Resampling

### 5.2.2   Advantages of Resampling Over Traditional Methods

### 5.2.3   Limitations and Considerations

## 5.3   Cross-Validation

### 5.3.1   Introduction to Cross-Validation

The Need for Cross-Validation

Types of Cross-Validation

### 5.3.2   K-Fold Cross-Validation

Implementation of K-Fold Cross-Validation

Advantages and Limitations

### 5.3.3   Leave-One-Out Cross-Validation (LOOCV)

Comparing LOOCV to K-Fold Cross-Validation

### 5.3.4   Stratified and Grouped Cross-Validation

When to Use Stratified vs. Grouped Cross-Validation

## 5.4   Bootstrap Methods

# Chapter 6

# Non-linear Models

## 6.1 Introduction to Non-linear Models

### 6.1.1 Definition and Importance

### 6.1.2 Contrast with Linear Models

### 6.1.3 Applications and Examples

## 6.2 Understanding Non-linearity in Data

### 6.2.1 Characteristics of Non-linear Relationships

### 6.2.2 Challenges in Modeling Non-linear Data

### 6.2.3 Tools for Identifying Non-linearity

## 6.3 Polynomial Regression

### 6.3.1 Introduction to Polynomial Regression

Why Polynomial Regression

Mathematical Foundation of Polynomial Regression

### 6.3.2 Implementing Polynomial Regression

Selecting the Degree of the Polynomial

Overfitting and Underfitting in Polynomial Regression

### 6.3.3 Advantages and Limitations of Polynomial Regression

## 6.4 Generalized Additive Models (GAM)

### 6.4.1 Introduction to Generalized Additive Models

From General Linear Models to GAM

Components and Formulation of GAM

### 6.4.2 Fitting GAM to Data

# Chapter 7

# Unsupervised Learning

## 7.1    Introduction to Unsupervised Learning

### 7.1.1    Definition and Overview

### 7.1.2    Contrast with Supervised Learning

### 7.1.3    Applications and Importance

## 7.2    Theoretical Foundations of Unsupervised Learning

### 7.2.1    Statistical Foundations

### 7.2.2    Dimensionality Reduction vs. Clustering

### 7.2.3    Metrics for Evaluating Unsupervised Learning

## 7.3    Clustering Methods

### 7.3.1    Overview of Clustering

Types of Clustering Methods

Choosing the Right Clustering Algorithm

### 7.3.2    K-Means Clustering

Algorithm and Implementation

Selecting the Number of Clusters

Strengths and Weaknesses

### 7.3.3    Hierarchical Clustering

Agglomerative vs. Divisive Hierarchical Clustering

Dendrogram Interpretation

Advantages and Limitations

## 7.4    Association Rules

# Chapter 8

# Handling Missing Data

## 8.1    Introduction

### 8.1.1    Importance of Handling Missing Data

### 8.1.2    Types of Missing Data

### 8.1.3    Impact of Missing Data on Analysis

## 8.2    Understanding Missing Data Mechanisms

### 8.2.1    Missing Completely at Random (MCAR)

### 8.2.2    Missing at Random (MAR)

### 8.2.3    Missing Not at Random (MNAR)

### 8.2.4    Imputation Techniques

## 8.3    Data Preprocessing Strategies

### 8.3.1    Identification of Missing Data

### 8.3.2    Deletion Methods

Listwise Deletion

Pairwise Deletion

### 8.3.3    Imputation Methods

Mean/Median Imputation

Mode Imputation

Regression Imputation

K-Nearest Neighbors (KNN) Imputation

Multiple Imputation

## 8.4    Advanced Techniques for Handling Missing Data

# Chapter 9

# Data Cleaning and Feature Selection

## 9.1   Introduction to Data Preprocessing

### 9.1.1   The Importance of Data Quality

### 9.1.2   Overview of Data Preprocessing Steps

### 9.1.3   Impact on Model Performance

## 9.2   Data Cleaning

### 9.2.1   Identifying and Handling Missing Values

Deletion vs. Imputation

Advanced Imputation Techniques

### 9.2.2   Detecting and Correcting Outliers

Statistical Methods

Proximity-Based Methods

### 9.2.3   Handling Duplicate Data

### 9.2.4   Normalization and Standardization

When and Why to Normalize

When and Why to Standardize

### 9.2.5   Dealing with Categorical Data

Encoding Techniques

Handling High Cardinality

## 9.3   Feature Selection

### 9.3.1   The Need for Feature Selection

### 9.3.2   Filter Methods

# Chapter 10

# Feature Engineering

## 10.1    Introduction to Feature Engineering

### 10.1.1    Definition and Importance

### 10.1.2    Role in Machine Learning and Data Science

### 10.1.3    Examples of Effective Feature Engineering

## 10.2    Principles of Feature Engineering

### 10.2.1    Understanding the Domain

### 10.2.2    Importance of Data Understanding in Feature Engineering

### 10.2.3    Balancing Complexity and Performance

## 10.3    Basic Techniques in Feature Engineering

### 10.3.1    Feature Creation

**Combining Features**

**Transformations and Normalizations**

### 10.3.2    Feature Extraction

**Principal Component Analysis (PCA)**

**Linear Discriminant Analysis (LDA)**

### 10.3.3    Feature Encoding

**One-Hot Encoding**

**Label Encoding**

**Encoding Categorical Variables with Many Categories**

## 10.4    Advanced Feature Engineering Techniques

### 10.4.1    Automated Feature Engineering

# Chapter 11

# Overfitting

## 11.1 Fundamental Concepts of Statistical Learning

### 11.1.1 Population vs. Sample

### 11.1.2 Bias-Variance Tradeoff

### 11.1.3 Supervised vs. Unsupervised Learning

### 11.1.4 Model Accuracy and Model Complexity

# Chapter 12

# Automated Machine Learning (AutoML)

## 12.1   Introduction to AutoML

### 12.1.1   Definition and Scope

### 12.1.2   The Evolution of AutoML

### 12.1.3   Importance and Impact on the Field of Machine Learning

## 12.2   The AutoML Pipeline

### 12.2.1   Overview of the AutoML Process

### 12.2.2   Data Preprocessing and Feature Engineering

### 12.2.3   Model Selection

### 12.2.4   Hyperparameter Optimization

### 12.2.5   Model Evaluation and Deployment

## 12.3   Key Components of AutoML

### 12.3.1   Data Cleaning Tools

### 12.3.2   Feature Engineering Automation

### 12.3.3   Automated Model Selection

### 12.3.4   Hyperparameter Tuning Techniques

Grid Search

Random Search

Bayesian Optimization

Evolutionary Algorithms

# Chapter 13

# Probability and Statistics

## 13.1    Introduction to Probability and Statistics

### 13.1.1    Definition and Importance

### 13.1.2    Role in Scientific Research and Data Analysis

### 13.1.3    Historical Evolution and Key Contributors

## 13.2    Probability Distributions

### 13.2.1    Overview of Probability Distributions

Definition and Significance

Discrete vs. Continuous Distributions

### 13.2.2    Key Probability Distributions

Uniform Distribution

Binomial Distribution

Normal Distribution

### 13.2.3    Properties of Probability Distributions

Mean, Variance, and Standard Deviation

Skewness and Kurtosis

## 13.3    Hypothesis Testing

### 13.3.1    Fundamentals of Hypothesis Testing

Null and Alternative Hypotheses

Type I and Type II Errors

### 13.3.2    Significance Levels and P-values

### 13.3.3    Commonly Used Hypothesis Tests

Z-test and T-test

# Chapter 14

# Tree-Based Methods

## 14.1   Introduction to Tree-Based Methods

### 14.1.1   Definition and Overview

### 14.1.2   Importance in Machine Learning

### 14.1.3   Types of Tree-Based Methods

## 14.2   Decision Trees

### 14.2.1   Fundamentals of Decision Trees

#### How Decision Trees Work

#### Criteria for Splitting

### 14.2.2   Building a Decision Tree

#### Algorithms for Tree Construction

#### Handling Overfitting in Decision Trees

### 14.2.3   Applications of Decision Trees

## 14.3   Ensemble Methods

### 14.3.1   Introduction to Ensemble Methods

### 14.3.2   Bagging

#### Bootstrap Aggregation

#### Random Forests

### 14.3.3   Boosting

#### Adaptive Boosting (AdaBoost)

#### Gradient Boosting

#### XGBoost, LightGBM, and CatBoost

## 14.4   Model Evaluation and Selection

# Chapter 15

# Support Vector Machines

## 15.1 Introduction to Support Vector Machines

### 15.1.1 Definition and Overview

### 15.1.2 Historical Background

### 15.1.3 Importance in Machine Learning

## 15.2 Theoretical Foundations of SVM

### 15.2.1 Linear SVM

Concept of Hyperplanes

Margin Maximization

### 15.2.2 Non-linear SVM

Kernel Trick

Types of Kernels

## 15.3 Mathematical Formulation of SVM

### 15.3.1 Optimization Problem

Objective Function

Constraints

### 15.3.2 Lagrange Multipliers

### 15.3.3 Dual Formulation

## 15.4 SVM for Classification

### 15.4.1 Binary Classification

Support Vectors and Decision Boundary

Interpretation of SVM Model Output

# Chapter 16

# Exploratory Data Analysis

## 16.1    Introduction to Exploratory Data Analysis

### 16.1.1    Definition and Scope

### 16.1.2    Importance in the Data Science Workflow

### 16.1.3    Goals and Principles of EDA

## 16.2    The Process of EDA

### 16.2.1    Understanding the Data Structure

### 16.2.2    Cleaning the Data

**Identifying and Handling Missing Values**

**Detecting and Removing Outliers**

### 16.2.3    Variable Identification

**Categorical vs. Continuous**

**Dependent vs. Independent Variables**

## 16.3    Univariate Analysis

### 16.3.1    Analyzing Continuous Variables

**Measures of Central Tendency**

**Measures of Dispersion**

### 16.3.2    Analyzing Categorical Variables

**Frequency Counts**

**Bar Charts and Pie Charts**

## 16.4    Bivariate and Multivariate Analysis

### 16.4.1    Correlation Analysis

**Pearson Correlation**

# Chapter 17

# Model Interpretability

## 17.1 Introduction to Model Interpretability

### 17.1.1 Definition and Importance

### 17.1.2 Overview of Methods in Model Interpretability

### 17.1.3 The Role of Interpretability in Machine Learning

## 17.2 The Need for Model Interpretability

### 17.2.1 Ethical and Legal Considerations

### 17.2.2 Building Trust in AI Systems

### 17.2.3 Debugging and Improving Models

## 17.3 Basics of Model Interpretability

### 17.3.1 Transparent vs. Post-hoc Interpretability

### 17.3.2 Local vs. Global Interpretability

### 17.3.3 Interpretability Techniques Overview

## 17.4 Introduction to SHAP

### 17.4.1 Background and Theoretical Foundations

Game Theory and Shapley Values

From Shapley Values to SHAP

### 17.4.2 Advantages of SHAP over Other Methods

## 17.5 SHAP in Practice

### 17.5.1 SHAP for Tree-based Models

TreeSHAP Algorithm

# Chapter 18

# Multiple Testing

## 18.1    Introduction to Multiple Testing

### 18.1.1    Definition and Importance

### 18.1.2    The Problem with Multiple Comparisons

### 18.1.3    Real-world Scenarios and Examples

## 18.2    Theoretical Foundations

### 18.2.1    Probability Theory and Error Rates

### 18.2.2    Type I and Type II Errors

### 18.2.3    Family-Wise Error Rate (FWER)

### 18.2.4    False Discovery Rate (FDR)

## 18.3    Controlling the Family-Wise Error Rate

### 18.3.1    Bonferroni Correction

### 18.3.2    Holm-Bonferroni Method

### 18.3.3    Šidák Correction

## 18.4    Controlling the False Discovery Rate

### 18.4.1    Benjamini-Hochberg Procedure

### 18.4.2    Benjamini-Yekutieli Procedure

### 18.4.3    Control of FDR in Practice

## 18.5    Advanced Topics in Multiple Testing

### 18.5.1    Post-hoc Analysis

### 18.5.2    Power Analysis in the Context of Multiple Testing

# Chapter 19

# Deep Learning

## 19.1   Introduction to Deep Learning

### 19.1.1   Definition and Importance

### 19.1.2   Historical Overview

### 19.1.3   Applications in Various Fields

## 19.2   Multilayer Perceptrons (MLPs)

### 19.2.1   Basic Structure and Architecture

**Input Layer**

**Hidden Layers**

**Output Layer**

### 19.2.2   Activation Functions

**Sigmoid**

**ReLU (Rectified Linear Unit)**

**Hyperbolic Tangent (tanh)**

### 19.2.3   Training MLPs

**Backpropagation Algorithm**

**Gradient Descent Optimization**

## 19.3   Convolutional Neural Networks (CNNs)

### 19.3.1   Fundamental Concepts

**Convolutional Layers**

**Pooling Layers**

**Fully Connected Layers**

### 19.3.2   Architectural Variants

# Chapter 20

# Generative Adversarial Networks (GANs)

## 20.1    Introduction to GANs

### 20.1.1    Definition and Importance

### 20.1.2    Brief History

### 20.1.3    Applications in Various Fields

## 20.2    Discriminative versus Generative Models

### 20.2.1    Discriminative Models

Definition and Characteristics

Examples: Logistic Regression, Support Vector Machines

### 20.2.2    Generative Models

Definition and Characteristics

Examples: Naive Bayes, Hidden Markov Models

## 20.3    Generative Adversarial Networks (GANs)

### 20.3.1    Basic Concept and Architecture

Generator

Discriminator

Training Process

### 20.3.2    Loss Functions

Generator Loss

Discriminator Loss

### 20.3.3    Variants of GANs

Conditional GANs

# Chapter 21

# Transformer Neural Networks

## 21.1 Introduction to Transformer Neural Networks

### 21.1.1 Motivation for Transformers

### 21.1.2 Overview of Transformer Architecture

### 21.1.3 Advantages over Recurrent and Convolutional Models

## 21.2 Attention is All You Need

### 21.2.1 Transformer Architecture

Self-Attention Mechanism

Positional Encoding

Feed-Forward Networks

Layer Normalization and Residual Connections

### 21.2.2 Training Procedure

Masked Self-Attention

Position-wise Feed-Forward Networks

Optimizer and Learning Rate Scheduling

## 21.3 BERT Neural Network

### 21.3.1 Introduction to BERT

### 21.3.2 BERT Architecture

Transformer Encoder Structure

Pre-training and Fine-tuning

### 21.3.3 BERT Variants

BERT Base vs. BERT Large

RoBERTa, DistilBERT, ALBERT, etc.

# Chapter 22

# Natural Language Processing (NLP)

## 22.1 Introduction to Natural Language Processing

### 22.1.1 Definition and Scope of NLP

### 22.1.2 Importance and Applications

## 22.2 WordNet

### 22.2.1 Definition and Purpose

### 22.2.2 WordNet Structure

### 22.2.3 Applications in NLP

## 22.3 Collocations

### 22.3.1 Definition and Examples

### 22.3.2 Identification Methods

### 22.3.3 Role in NLP

## 22.4 Text Mining and Natural Language Processing

### 22.4.1 Text Mining vs. NLP

### 22.4.2 Text Processing Techniques

### 22.4.3 NLP Applications in Text Mining

## 22.5 Python Natural Language Tools

### 22.5.1 Overview of Python NLP Libraries

### 22.5.2 NLTK (Natural Language Toolkit)

### 22.5.3 spaCy

# Chapter 23

# Data Visualization

## 23.1 Introduction to Data Visualization

### 23.1.1 Definition and Importance

### 23.1.2 Role in Data Analysis and Communication

## 23.2 Add Content to Data Visualization

### 23.2.1 Enhancing Visualization with Additional Content

### 23.2.2 Interactive Visualizations

## 23.3 Data Types, Graphical Marks, and Visual Encoding Channels

### 23.3.1 Understanding Data Types

### 23.3.2 Graphical Marks

### 23.3.3 Visual Encoding Channels

Position

Color

Size

Shape

Texture

## 23.4 Edward Tufte

### 23.4.1 Background and Contributions

### 23.4.2 Tufte's Principles of Data Visualization

## 23.5 Hans Rosling

### 23.5.1 Rosling's Work in Data Visualization

# Chapter 24

# Grammar of Graphics

## 24.1   Introduction to Grammar of Graphics

### 24.1.1   Definition and Concept

### 24.1.2   Importance in Data Visualization

## 24.2   Grammar of Graphics in R

### 24.2.1   Overview of ggplot2 Package

### 24.2.2   Components of ggplot2 Grammar

Data

Aesthetic Mapping

Geometric Objects

Facets

Statistics

Coordinates

Themes

## 24.3   Grammar of Graphics in Python

### 24.3.1   Introduction to Plotnine

### 24.3.2   Comparison with ggplot2

## 24.4   Applications of Grammar of Graphics

### 24.4.1   Data Exploration and Analysis

### 24.4.2   Statistical Graphics

### 24.4.3   Publication-Quality Plots

## 24.5   Case Studies

# Chapter 25

# Python Review

## 25.1    Introduction to Python

### 25.1.1   What is Python?

### 25.1.2   Why Python?

### 25.1.3   Python in Various Domains

## 25.2    Intro to Python Data Structures

### 25.2.1   Lists

### 25.2.2   Tuples

### 25.2.3   Dictionaries

### 25.2.4   Sets

## 25.3    Data Visualization with matplotlib

### 25.3.1   Introduction to matplotlib

### 25.3.2   Basic Plotting with matplotlib

### 25.3.3   Advanced Plot Customization

### 25.3.4   Plotting Data Structures

## 25.4    Jupyter Markdown

### 25.4.1   Markdown Basics

### 25.4.2   Markdown for Jupyter Notebooks

### 25.4.3   Markdown Syntax and Formatting

## 25.5    Hands-On Python Exercises

### 25.5.1   Practice Problems

# Chapter 26

# R Review

## 26.1   Introduction to R

### 26.1.1   What is R?

### 26.1.2   Why R?

### 26.1.3   R in Various Domains

## 26.2   Intro to R Data Structures

### 26.2.1   Vectors

### 26.2.2   Matrices

### 26.2.3   Data Frames

### 26.2.4   Lists

## 26.3   Data Visualization with ggplot

### 26.3.1   Introduction to ggplot

### 26.3.2   Basic Plotting with ggplot

### 26.3.3   Advanced Plot Customization

### 26.3.4   Plotting Data Structures

## 26.4   Jupyter Markdown

### 26.4.1   Markdown Basics

### 26.4.2   Markdown for Jupyter Notebooks

### 26.4.3   Markdown Syntax and Formatting

## 26.5   Hands-On R Exercises

### 26.5.1   Practice Problems

# Chapter 27

# Data Munging

## 27.1    Introduction to Data Munging

### 27.1.1    Definition and Importance

### 27.1.2    Role of Data Munging in Data Analysis

### 27.1.3    Challenges in Data Munging

## 27.2    Data Cleaning Techniques

### 27.2.1    Handling Missing Values

### 27.2.2    Removing Duplicate Data

### 27.2.3    Standardizing and Normalizing Data

### 27.2.4    Dealing with Outliers

## 27.3    Data Transformation

### 27.3.1    Data Reshaping

### 27.3.2    Variable Transformation

### 27.3.3    Feature Engineering

## 27.4    Data Integration

### 27.4.1    Combining Data Sources

### 27.4.2    Joining and Merging Datasets

### 27.4.3    Reshaping Data for Integration

## 27.5    Data Reduction

### 27.5.1    Dimensionality Reduction Techniques

**Principal Component Analysis (PCA)**

# Chapter 28

# Case Studies and Applications of Statistical Learning

## 28.1 Introduction

### 28.1.1 Overview of Statistical Learning

### 28.1.2 Importance of Case Studies in Understanding Applications

## 28.2 Application in Computational Biology

### 28.2.1 Genomic Data Analysis

### 28.2.2 Protein Structure Prediction

### 28.2.3 Drug Discovery

## 28.3 Application in Finance

### 28.3.1 Stock Price Prediction

### 28.3.2 Portfolio Optimization

### 28.3.3 Credit Scoring

## 28.4 Application in Healthcare

### 28.4.1 Disease Diagnosis

### 28.4.2 Medical Image Analysis

### 28.4.3 Patient Outcome Prediction

## 28.5 Application in Marketing

### 28.5.1 Customer Segmentation

### 28.5.2 Market Basket Analysis

# Chapter 29

# Bayesian Statistical Methods

## 29.1    Introduction to Bayesian Statistics

### 29.1.1   Overview of Bayesian Inference

### 29.1.2   Comparison with Frequentist Statistics

### 29.1.3   Importance and Applications

## 29.2    Bayesian Probability

### 29.2.1   Bayes' Theorem

### 29.2.2   Prior, Likelihood, and Posterior Distributions

### 29.2.3   Conjugate Priors

## 29.3    Bayesian Modeling

### 29.3.1   Parameter Estimation

### 29.3.2   Model Comparison and Selection

### 29.3.3   Hierarchical Modeling

## 29.4    Markov Chain Monte Carlo (MCMC)

### 29.4.1   Gibbs Sampling

### 29.4.2   Metropolis-Hastings Algorithm

### 29.4.3   Hamiltonian Monte Carlo (HMC)

## 29.5    Bayesian Computation

### 29.5.1   Computational Techniques

### 29.5.2   Simulation Methods

### 29.5.3   Approximate Bayesian Computation (ABC)

# Chapter 30

# Survival Analysis and Censored Data

## 30.1    Introduction to Survival Analysis

### 30.1.1    Definition and Scope

### 30.1.2    Key Concepts: Survival Time, Hazard, Censoring

## 30.2    Types of Censoring

### 30.2.1    Right Censoring

### 30.2.2    Left Censoring

### 30.2.3    Interval Censoring

### 30.2.4    Informative vs. Non-Informative Censoring

## 30.3    Survival Probability and Hazard Function

### 30.3.1    Kaplan-Meier Estimator

### 30.3.2    Nelson-Aalen Estimator

### 30.3.3    Hazard Ratio

## 30.4    Parametric Survival Models

### 30.4.1    Exponential Distribution

### 30.4.2    Weibull Distribution

### 30.4.3    Log-Normal Distribution

### 30.4.4    Parametric Regression Models

## 30.5    Non-Parametric Survival Models

### 30.5.1    Cox Proportional Hazards Model

# Chapter 31

# Time Series Analysis and Forecasting

# Chapter 32

# Real-World Implementations

## 32.1 GNS Healthcare

# References

# Acknowledgements