



ANALYSIS AND PREDICTIVE MODELING

Study Case Using



Writter by: Bakti Siregar, M.Sc., CDS.





First Edition

Analysis and Predictive Modeling Study Case using R & Python

Bakti Siregar, M.Sc.,CDS

Table of contents

Pı	eface	e		3
	Abo	ut the Writer		3
	Ackr	nowledgments		3
	Feed	lback & Suggestions		4
\mathbf{A}	bout	the Book		5
	Intro	oduction		5
	Over	rview of the Course		6
	Refe	erences		6
1	Intr	roduction		7
	1.1	What is PA?		7
		1.1.1 Types		7
		1.1.2 Techniques		9
		1.1.3 Data Types		9
	1.2	Why use PA?		10
		1.2.1 Benefits & Business Impact		10
		1.2.2 ROI (Return on Investment)		10
		1.2.3 Example & Discussion		10
	1.3	When to apply PA?		11
	1.4	Where is PA applied?		12
	1.5	Who is involved?		13
	1.6	How to implement PA?		13

2	Reg	ressio	n Models	15
	2.1	Linear	Model	15
		2.1.1	Simple Linear Reg	15
		2.1.2	Multiple Linear Reg	19
	2.2	Nonlin	near Regression	23
		2.2.1	Multiple Non-Linear Reg	23
		2.2.2	Polynomial Regression	27
	2.3	Logist	ics Regression	32
		2.3.1	Binary Logistic Regression	32
		2.3.2	Multinomial Logistics	35
	Refe	erences		35
3	Clas	ssificat	cion Models	37
	3.1	Intro	to Classification	37
	3.2	Decisi	on Tree	38
		3.2.1	CART Algorithm	40
		3.2.2	ID3 Algorithm	41
		3.2.3	C4.5 Algorithm	41
		3.2.4	Comparison Decision Tree	42
	3.3	Proba	bilistic Models	42
		3.3.1	Naive Bayes	43
		3.3.2	LDA	43
		3.3.3	QDA	44
		3.3.4	Logistic Regression	45
		3.3.5	Comparison Probabilistics	45
	3.4	Kerne	l Methods	46
		3.4.1	SVM	46
		3.4.2	KNN	47
		3.4.3	Comparison Kernels	48
	3.5	Ensen	able Methods	48
		3.5.1	Random Forest	48
		3.5.2	Gradient Boosting Machines	49
		3.5.3	Extreme Gradient Boosting	50
		3.5.4	Comparison Ensemble	51

TA	BL	E_{i}	OF	CO	N'	ΓE	N'	ΓS

$T_{\mathcal{A}}$	ABLE	C OF CONTENTS	
	3.6	Table Studi Case Examples	51
	3.7	End to End Study Case	51
		3.7.1 Data Preparation	52
		3.7.2 Logistic Regression Model	52
		3.7.3 Prediction and Classification	56
		3.7.4 Confusion Matrix	57
		3.7.5 Evaluation Metrics	58
		3.7.6 ROC Curve and AUC	59
		3.7.7 Conclusion	60
	Refe	erences	61
4	Clu	stering Models	63
5	Tin	ne Series Model	65
	5.1	Historical Methods	65
	5.2	Time Series (ARIMA, Prophet, LSTM)	65
	5.3	Probabilistic vs Deterministic	65
	5.4	Evaluation	65
6	Tin	ne Series Forecasting	67
	6.1	Traditional vs Modern	67
	6.2	Exponential Smoothing & ARIMA	67
	6.3	Multivariate Forecasting	67
	6.4	Accuracy Evaluation	67
7	Pre	diction Engineering	69
	7.1	Feature Engineering	69
	7.2	Ensemble Methods	69
	7.3	Hyperparameter Tuning	69
	7.4	Pipeline Integration	69
8	Sma	art Prediction	7 1
	8.1	Deep Learning Models	71
	8.2	NLP Prediction	71
	8.3	Computer Vision	71

9	Inte	lligent Analytics	73
	9.1	AI-powered Prediction	73
	9.2	AutoML	73
	9.3	Explainable AI	73
	9.4	Ethics & Fairness	73
10	Futi	are Prediction	75
	10.1	Data-driven Decisions	75
	10.2	Big Data & Streaming	75
	10.3	Deployment (API, n8n, MLOps)	75
	Mon	itoring & Retraining	75

In the era of digital transformation, prediction has emerged as one of the most critical capabilities in data science. Organizations across sectors increasingly rely on **predictive** analytics not merely to understand the past, but to anticipate the future and shape strategic decisions. This course provides a comprehensive exploration of the principles, methods, and applications of predictive modeling, guiding learners through the full spectrum of concepts and practices that define modern predictive data science.

The journey begins with **predictive modeling foundations**, where learners will examine regression and classification approaches, understand the power of ensemble techniques, and apply model interpretation frameworks to ensure transparency and trust. Building upon these methods, the course transitions to **data-driven prediction and forecasting**, introducing both traditional statistical models and advanced machine learning architectures such as Prophet and LSTM to address temporal and sequential data challenges.

From theory, the focus shifts to **applied prediction**, where practical case studies in business, healthcare, and industry demonstrate the value of predictive insights in real-world decision-making. Learners will also explore **prediction engineering**, emphasizing the importance of feature construction, hyperparameter optimization, and workflow integration to maximize predictive performance.

As predictive systems evolve, the module highlights the role of **smart and intelligent prediction** powered by artificial intelligence. Applications include natural language processing for text-based inference, computer vision for image-based forecasting, and AutoML frameworks that automate complex predictive tasks. These advancements are complemented by discussions on **ethics**, **fairness**, **and explainability**, ensuring that predictive analytics remains responsible and aligned with human values.

The course culminates in **deployment and monitoring practices**, equipping learners with the technical and conceptual skills to operationalize predictive models. Through exposure to APIs, MLOps pipelines, and automation platforms such as n8n, participants will develop the capacity to bring predictive solutions into production and sustain their performance over time.

By completing this module, learners will gain not only technical proficiency but also a holistic perspective on predictive analytics—capable of designing, implementing, and managing intelligent systems that turn data into foresight and foresight into impact.

Preface

About the Writer



Bakti Siregar, M.Sc., CDS is a Lecturer in the Data Science Program at ITSB. He obtained his Master's degree in Applied Mathematics from the National Sun Yat-sen University, Taiwan. Alongside his academic role, Bakti also serves as a Freelance Data Scientist, collaborating with leading companies such as JNE, Samora Group, Pertamina, and PT. Green City Traffic.

His professional and research interests include Big Data Analytics, Machine Learning, Optimization, and Time Series Analysis, with a particular focus on finance and investment applications. His core expertise lies in statistical programming using R and Python, complemented by strong experience in database management systems such as MySQL and NoSQL. In addition, he is proficient in applying Big Data technologies, including Spark and Hadoop, for large-scale data processing and analysis.

Some of his projects can be viewed here: Rpubs, Github, Website, and Kaggle

Acknowledgments

Predictive Analytics is more than a technical discipline—it is a bridge between data and decision-making, transforming information into foresight that guides action. This

4 Preface

book is crafted to help learners advance from foundational understanding toward building end-to-end predictive solutions that are robust, interpretable, and impactful.

The material explores a connected sequence of topics:

- **Predictive Modeling Foundations:** Regression, classification, ensembles, and interpretability.
- Forecasting and Sequential Models: Time series methods, Prophet, and LSTM for temporal data.
- Applied Prediction: Real-world use cases across business, healthcare, and industry.
- **Prediction Engineering:** Feature construction, optimization, and workflow integration.
- Intelligent Prediction: AI-driven approaches in NLP, vision, and AutoML.
- **Deployment and MLOps:** Operationalizing predictive systems with APIs, pipelines, and monitoring.

This work would not have been possible without the encouragement of colleagues, students, and mentors who provided constructive feedback, shared insights, and inspired new directions. My deepest gratitude goes to all who contributed, directly or indirectly, to the development of this resource. It is my hope that this book will serve as both a **practical reference** and a **roadmap** for learners and professionals applying predictive analytics in research, industry, and innovation.

Feedback & Suggestions

The evolution of this book depends on continuous learning and dialogue. Readers are warmly invited to share their perspectives on clarity, depth, case studies, and practical relevance. Suggestions for expanding future editions—whether through advanced methods, new applications, or emerging tools—are highly valued.

Your input will help refine this work into a comprehensive, living resource that grows with the rapidly changing field of predictive analytics. Thank you for your engagement and contributions to this journey.

For feedback and suggestions, please reach out via:

- dsciencelabs@outlook.com
- siregarbakti@gmail.com
- siregarbakti@itsb.ac.id

About the Book

Analysis and Predictive Modeling (APM) provides a structured framework for transforming data into actionable insights through advanced statistical analysis, machine learning techniques, and model evaluation. Positioned as a critical discipline in modern research and professional practice, it emphasizes the end-to-end process of preparing, analyzing, and applying predictive models to address real-world problems [1].

\newline \href{https://youtu.be/WxXZaP8Y8pI}{Click here to watch the video}

Introduction

This book begins with essential programming practices, continues through data integration, transformation, and feature engineering, and culminates in predictive modeling, evaluation, and deployment [2]. Each chapter reflects the progression of a typical modeling workflow, offering both methodological depth and practical guidance [3].

Key Topics include:

- **Programming Foundations:** Modular code, functional programming, and reproducible workflows tailored for data analysis [4].
- Data Acquisition & Preparation: Integration from APIs and databases, advanced wrangling, and feature engineering [2].
- Modeling & Evaluation: Building predictive models, validation strategies, interpretability, and visualization [3].
- **Deployment & Applications:** Model packaging, workflow automation, monitoring, and applied use cases [5].

By combining theoretical underpinnings, applied examples, and recommended practices, this book prepares graduate students, researchers, and professionals to confidently manage complex modeling tasks and apply predictive models in impactful ways across disciplines.

6 About the Book

Overview of the Course

Figure Figure 1 provides a conceptual overview of this book, illustrating the key components of **Analysis and Predictive Modeling** and their interrelationships. It serves as a roadmap for readers, showing how the material progresses from programming practices and data preparation to modeling, evaluation, and deployment. This framework underscores the integration of each stage into a coherent workflow, bridging methodological foundations with applied decision-making in real-world contexts [2], [3].

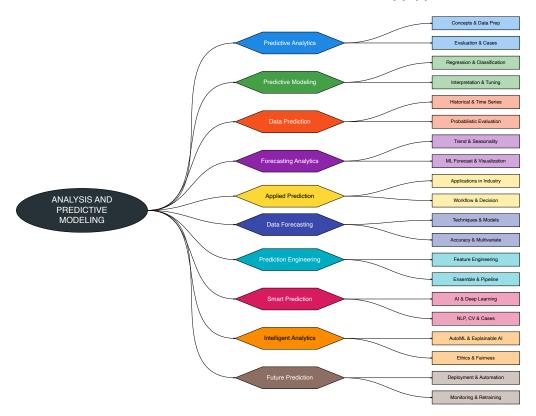


Figure 1: Mind Map of Analysis and Predictive Modeling

References

Chapter 1

Introduction

Understanding Predictive Analytics is the first step in exploring the world of data-driven decision-making. Predictive analytics serves as a core foundation for modern data science, business intelligence, machine learning, and various applied sciences. It provides a framework for forecasting future outcomes, assessing risks, and supporting strategic planning in both research and industry applications [1]–[3].

To help navigate the key aspects of predictive analytics, the Figure 1.1 offers a 5W+1H mind map. This visualization guides learners through the What—its definitions, techniques, and data types; the Why—business value, benefits, and ROI; the When—timing of application across operations, marketing, and risk assessment; the Where—applications in finance, healthcare, supply chain, and case studies such as Netflix and Walmart; the Who—the roles of data scientists, business analysts, and domain experts; and the How—the workflow, tools, and performance evaluation metrics. By Figure 1.1, one can see not just the methods themselves, but also their significance, challenges, and real-world impact across industries.

1.1 What is PA?

Predictive Analytics is a branch of data analytics that focuses on **forecasting future outcomes** based on historical and current data. Unlike traditional reporting that only describes what has happened, predictive analytics goes a step further by applying **statistical methods**, **machine learning algorithms**, and data modeling techniques to anticipate what is likely to occur in the future.

In essence, predictive analytics combines data (past & present), mathematical models, and computational power to generate actionable insights. It is widely used across industries to improve decision-making, optimize business processes, and reduce uncertainty in planning.

1.1.1 Types

To fully understand predictive analytics, it is important to distinguish it from other types of analytics. **Descriptive Analytics** answers the question "What happened?"

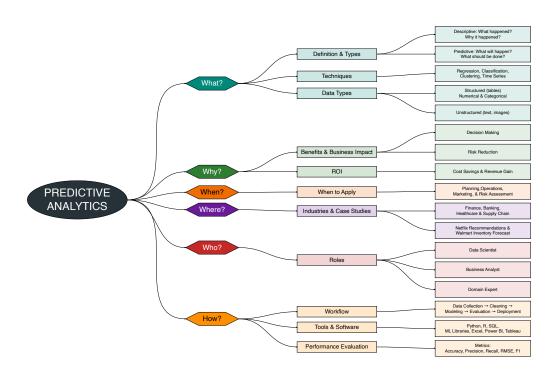


Figure 1.1: Predictive Analytics Mind Map (5W+1H)

1.1. WHAT IS PA? 9

Technique Example Description Regression Analysis Predicts a continuous numerical value. Estimating housing prices based on size, location, and amenities. Classification Predicts a categorical outcome. Determining whether a loan applicant is 'high risk' or 'low risk.' Groups data into clusters based on Clustering Customer segmentation for targeted similarity, without pre-labeled outcomes. marketing campaigns. Time Series Analysis Predicts values over time, considering Forecasting energy consumption, stock temporal patterns. prices, or product demand.

Table 1.1: Key Techniques in Predictive Analytics

by summarizing historical data through reports, dashboards, and statistics, such as a monthly sales report showing total revenue in the last quarter. In contrast, **Predictive Analytics** answers the question "What will happen?" by using models to identify patterns in data and forecast future outcomes, for example, predicting customer churn in the next six months based on transaction history. While descriptive analytics helps organizations understand the past, predictive analytics enables them to **prepare for the future**.

1.1.2 Techniques

Several techniques Table 1.1 are commonly applied in predictive modeling, each suited for different types of problems:

Each of these techniques may use machine learning algorithms such as linear regression, decision trees, random forests, support vector machines, or neural networks, depending on the complexity of the problem.

1.1.3 Data Types

The foundation of predictive analytics lies in **data**, which can be broadly categorized as:

Structured Data adalah data yang tersusun rapi dalam tabel dengan baris dan kolom. Data ini biasanya mencakup informasi numerik, seperti angka penjualan atau suhu, serta data kategorikal, seperti kategori produk atau wilayah pelanggan. Contoh nyata structured data adalah catatan transaksi dalam basis data ritel yang dapat langsung diolah menggunakan perangkat lunak analisis.

Sebaliknya, **Unstructured Data** tidak memiliki format yang terdefinisi dengan jelas sehingga memerlukan teknik pemrosesan lanjutan. Data ini dapat berupa teks seperti ulasan atau postingan media sosial, serta format multimedia seperti gambar, audio, dan video. Contoh penerapan unstructured data adalah analisis sentimen pelanggan yang diambil dari postingan Twitter atau ulasan produk.

Combining structured and unstructured data often provides richer insights. For example, predicting customer churn may involve structured data (purchase history) and unstructured data (customer complaints via email or chat).

In summary, predictive analytics is about moving from "knowing the past" to "anticipating the future." By applying techniques such as regression, classification, clustering, and time series analysis on both structured and unstructured data, organizations gain the ability to make proactive decisions. This makes predictive analytics a powerful tool in industries ranging from finance and healthcare to retail and manufacturing.

1.2 Why use PA?

1.2.1 Benefits & Business Impact

Predictive analytics helps organizations make **Data-driven Decisions** by providing projections of market trends and customer behavior. This approach reduces reliance on intuition alone, ensuring that strategies are backed by solid evidence. For example, a retail company can apply demand forecasting models to optimize inventory levels, ensuring that products are available when needed while minimizing overstock and reducing waste.

Another major benefit of predictive analytics is **Risk Reduction**. By anticipating potential risks, organizations can take proactive measures before problems escalate. This includes detecting fraud in financial transactions, identifying customers who are at risk of churn, and predicting machine failures in manufacturing processes. Such predictive capabilities allow businesses to minimize losses, improve efficiency, and maintain stronger customer relationships.

1.2.2 ROI (Return on Investment)

Analytics **reduces costs** by driving efficiency improvements across business operations. Through supply chain optimization, companies can streamline logistics and reduce unnecessary expenses. More accurate demand forecasts help lower operational costs by preventing both overstock and stockouts. In addition, the early detection of equipment failures enables organizations to minimize repair expenses and avoid costly downtime.

Beyond cost reduction, analytics also plays a key role in generating **revenue growth**. Personalized product recommendations enhance customer engagement and boost sales by targeting the right audience with the right offerings. Analytics can also uncover new market opportunities through the analysis of consumer trends, giving businesses a competitive edge. Moreover, dynamic pricing strategies based on demand patterns allow companies to maximize profitability while staying responsive to market changes.

1.2.3 Example & Discussion

Predictive modeling allows businesses to **forecast future outcomes** and act strategically.

For instance, an e-commerce company applies a **churn prediction model** to identify customers likely to stop using their platform. By targeting these customers with special offers or retention campaigns, the company manages to **reduce churn by 15%**.

Mathematically, if the company originally had N customers and an expected churn rate of r, then the number of customers lost without intervention would be:

$$L_0 = N \times r$$

After predictive intervention, the lost customers become:

$$L_1 = N \times (r - 0.15r) = N \times (0.85r)$$

This reduction translates directly into **higher revenue**, since more customers remain active and continue purchasing.

Return on Investment (ROI) is a measure of how much benefit a project delivers compared to its cost. The formula is:

$$ROI = \frac{Benefit - Cost}{Cost} \times 100\%$$

Example:

- Cost of analytics project: 100,000
- Benefit (savings + extra revenue): 300,000

Then,

$$ROI = \frac{300,000 - 100,000}{100,000} \times 100\%$$

$$ROI = \frac{200,000}{100,000} \times 100\% = 200\%$$

This means that for every \$1 invested, the company gains \$2 in net value.

With predictive analytics, the **business impact** can be clearly seen both in reduced risks and increased revenues, while the **ROI calculation** ensures that every project is evaluated in terms of tangible financial return.

1.3 When to apply PA?

Predictive analytics can be applied at different stages of business processes, and the **timing of its application** determines the level of impact it creates. In the planning stage, analytics helps set long-term strategies. In operations, it improves efficiency. In marketing, it drives customer engagement, and in risk assessment, it prevents potential losses. The Table 1.2 summarizes the purpose, examples, and mathematical representations for each stage.

Subtopic	Purpose	Example	\mathbf{F}
Planning	Forecasting long-term trends for strategic decisions.	Mining company predicts raw material demand for 5 years.	\$9
Operations	Improving efficiency and reducing costs through real-time applications.	Predictive maintenance to reduce downtime.	\$9 \t
Marketing	Anticipating customer needs and personalizing offers.	Predicting which customers will respond to a campaign.	\$9
Risk Assessment	Identifying and mitigating potential risks.	Credit scoring to predict loan defaults.	\$\$ \}

Table 1.2: Timing of Predictive Analytics Application

Table 1.3: Industries and Case Studies in Predictive Analytics

Industry	Application	Example
Finance & Banking	Risk prediction, fraud detection, credit scoring	Detecting fraudulent credit card transactions
Healthcare	Predictive diagnosis, patient monitoring, treatment	Predicting patient readmission rates
Supply Chain	Inventory planning, demand forecasting, logistics	Optimizing delivery routes and reduction stockouts
Case Studies	Customer personalization, operational optimization	Netflix recommendations, Walmart inventory forecasting

The Table 1.2 shows that predictive analytics provides unique benefits across different departments. Planning benefits from long-term forecasts, operations gain efficiency through real-time applications, marketing achieves higher engagement with personalization, and risk assessment reduces losses by identifying threats early. In short, the earlier predictive analytics is applied within a process, the greater its impact on decision-making and business performance.

1.4 Where is PA applied?

In the application of **predictive analytics**, each industry has its own needs, challenges, and approaches. For instance, the finance sector emphasizes **risk prediction** and **fraud detection**, while healthcare focuses on **predictive diagnosis** and **patient monitoring**. On the other hand, the supply chain leverages predictive analytics for **distribution efficiency** and **inventory planning**. Case studies from major companies such as **Netflix** and **Walmart** demonstrate how predictive methods can be effectively adapted to improve **customer experience** and **operational optimization**.

From the Table 1.3, it is clear that predictive analytics is not limited to a single field but can be broadly implemented with methods tailored to each context. An approach that works well in one industry may not be directly applicable to another without proper adjustments. Therefore, understanding **real-world case studies** is crucial so that organizations can adapt predictive strategies aligned with their **business goals**, **data availability**, and **operational challenges**.

Profession Materials Workplace Data Scientist Build & validate predictive models, statistical analysis. Tech companies, fintech, research labs machine learning Business Analyst Translate analytics results into business strategy and Consulting firms, corporate strategy, finance decision-making Domain Expert Provide deep knowledge of the industry/domain context Healthcare, energy, manufacturing Data Engineer Prepare, clean, and manage data infrastructure Big data companies, cloud providers Machine Learning Implement & optimize predictive models in production Startups, AI labs, enterprise IT Engineer

Table 1.4: Professions, Materials, and Workplaces in Predictive Analytics

Table 1.5: Workflow, Tools, Models, and Evaluation by Profession

Profession	Workflow	Tools	Models
Data Scientist	$Modeling \rightarrow Evaluation$	Python, R, SQL, scikit-learn,	Regression, Class
		TensorFlow	Time Series, Neu
Business Analyst	Requirements \rightarrow	Excel, Power BI, Tableau	Decision trees for
	Interpretation		descriptive dashb
Domain Expert	Contextual Guidance \rightarrow	Domain-specific tools, knowledge	Domain-specific r
	Validation	bases	frameworks
Data Engineer	$ \begin{array}{l} \text{Data Collection} \rightarrow \text{Cleaning} \\ \rightarrow \text{Preparation} \end{array} $	SQL, Spark, Hadoop, ETL Tools	Data pipelines, so quality rules
Machine Learning	Deployment \rightarrow Monitoring	Python, MLflow, Docker, Kubernetes	Deep learning, en
Engineer			reinforcement lea

1.5 Who is involved?

In predictive analytics projects Table 1.4, success depends not only on technology but also on the people involved. Each role contributes unique competencies and responsibilities, making collaboration essential.

For predictive analytics projects to succeed, **collaboration between these roles is critical**. Data Scientists bring technical expertise, Business Analysts ensure alignment with strategy, and Domain Experts add real-world context. Together, they create solutions that are not only accurate but also actionable and valuable for the organization.

1.6 How to implement PA?

Predictive analytics projects require collaboration among multiple roles, each with its own workflow, tools, and methods of evaluation. The Table 1.5 summarizes how different professions contribute to the analytics process, highlighting their focus areas and approaches. This structured view helps us understand that successful predictive analytics is not only about algorithms, but also about integrating business, technical, and domain expertise.

The Table 1.5 shows that each profession brings unique skills and responsibilities. Data Scientists and Machine Learning Engineers focus on algorithms and deployment, while Business Analysts and Domain Experts ensure alignment with business needs. Data Engineers provide the infrastructure that supports the entire process. Together, their

collaboration ensures predictive analytics projects deliver accurate, actionable, and business-relevant results.

Chapter 2

Regression Models

Regression Models is a cornerstone of modern data science, enabling us to transform raw data into actionable insights. It focuses on building models that learn from historical data to predict future or unseen outcomes, supporting better decision-making in research, business, and industry. Regression Models integrates principles from statistics, machine learning, and domain expertise, bridging theory with practical applications [1], [3], [4].

To illustrate these connections, the Figure 2.1 provides a hierarchical mind map. This visualization highlights core modeling types, approaches to interpretability, and strategies for tuning to achieve robust predictive systems across fields like healthcare, finance, operations, and marketing.

2.1 Linear Model

2.1.1 Simple Linear Reg.

Simple Linear Regression models the relationship between **one independent variable** and a **dependent variable** as a straight line:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where:

- Y: dependent variable (target)
- X: independent variable (predictor)
- β_0 : intercept (constant term)
- β_1 : slope coefficient (change in Y per unit change in X)
- ε : random error term

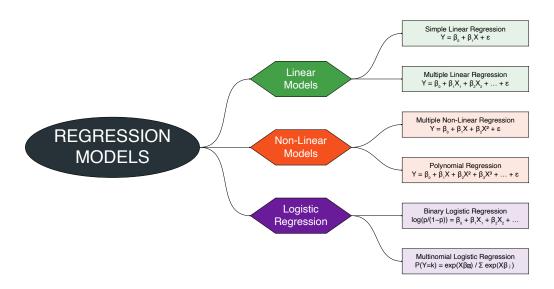


Figure 2.1: Comprehensive Regression Models Mind Map with Equations

Study Case: Simple Linear Regression

In this study, the goal is to model the **relationship between advertising expenditure and sales performance**. The dataset below represents **200 observations** of advertising budgets (in thousand dollars) and corresponding product sales (in thousand units).

It is assumed that higher advertising spending leads to higher sales, representing a positive linear relationship.

Copy CSV Search: 12.18943800311536 15.2244230452925 106.4918310097775 28.51168210734613 185.8097480250997 6.138912484748289 49.45265390277468 120.4680971130673 27.31047610985115 170.8386753481511 16.41536838258617 125.8898605774177 Showing 1 to 10 of 200 entries Previous 1 2 3 4 5 ...

Table 2.1: Simulated Dataset for Simple Linear Regression

Solution

A simple linear regression model is fitted to describe the effect of Advertising (X) on Sales (Y).

```
# Fit the model
model_simple <- lm(Sales ~ Advertising, data = data_simple)

# Show summary
summary(model_simple)</pre>
```

Call:

lm(formula = Sales ~ Advertising, data = data_simple)

Residuals:

```
Min 1Q Median 3Q Max
-21.271 -6.249 -1.110 5.954 32.021
```

Coefficients:

Residual standard error: 9.645 on 198 degrees of freedom Multiple R-squared: 0.9547, Adjusted R-squared: 0.9544

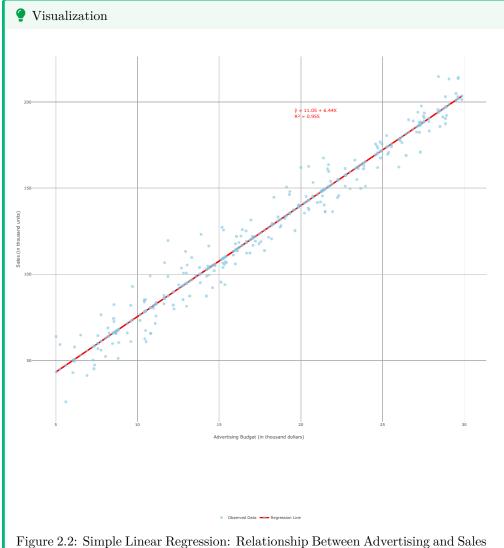
F-statistic: 4168 on 1 and 198 DF, p-value: < 2.2e-16

Interpretation of Regression Results:

The simple linear regression analysis examines the relationship between **Advertising Budget** and **Sales**.

- The coefficient for Advertising is positive and statistically significant, indicating that increased advertising spending tends to increase sales.
- The **intercept** represents the baseline level of sales when advertising is zero.
- The R² value (coefficient of determination) shows the proportion of variation in Sales explained by Advertising.
 - A high R² (close to 1) indicates a strong relationship, while a low R² indicates that other factors may also influence sales.
- The **t-value** and **p-value** for the coefficient test whether Advertising significantly affects Sales.

Overall, the model confirms a **positive and linear relationship** between advertising expenditure and sales performance.



(with R^2)

2.1.2 Multiple Linear Reg.

Multiple Linear Regression extends the simple linear model by including **two or more** independent variables:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon$$

where:

• Y: dependent variable (target)

- X_1, X_2, X_3 : independent variables (predictors)
- β_0 : intercept
- $\beta_1, \beta_2, \beta_3$: coefficients (effect of each predictor)
- ε : random error term

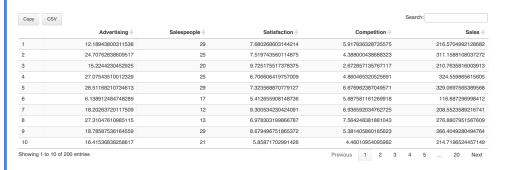
i Study Case: Multiple Linear

In this study, we aim to model the **relationship between marketing factors** and product sales. The four independent variables are as follows:

- Advertising Budget X₁: The amount spent on advertising (in thousand dollars) expected to have a positive effect on sales.
- Number of Salespeople X₂: The total number of sales representatives
 — more salespeople should increase sales.
- Customer Satisfaction Score X₃: A satisfaction score on a 1–10 scale
 higher satisfaction typically leads to repeat purchases.
- Competition Level X_4 : The level of market competition (1–10 scale) expected to have a **negative** impact on sales.

The following simulated dataset (Table 2.2) will be used for regression analysis. It contains 200 observations and the variables described above.

Table 2.2: Transformed Business Dataset — Total price after discount



Solution

A multiple linear regression model is fitted that incorporates all predictor variables Advertising, Salespeople, Satisfaction, and Competition independent variables effectively predict Sales.

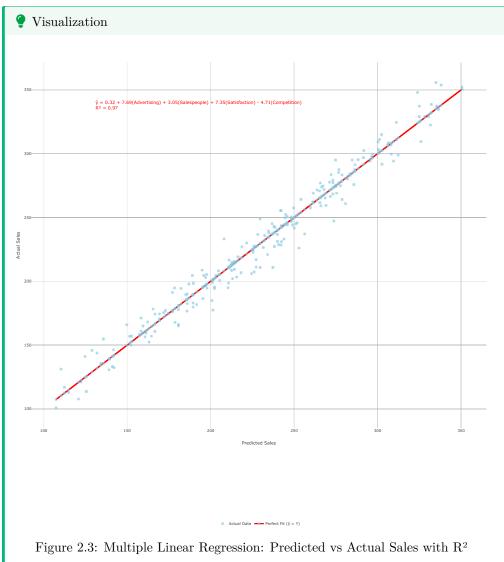
```
# Check model R2
model_check <- lm(Sales ~ Advertising + Salespeople + Satisfaction + Competition, data = da
summary(model_check)
Call:
lm(formula = Sales ~ Advertising + Salespeople + Satisfaction +
    Competition, data = data_reg)
Residuals:
     Min
               10
                    Median
                                 30
                                         Max
-26.7988 -6.8575
                    0.8932
                             6.2999
                                     25.0942
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)
               0.3178
                          5.1073
                                   0.062
                                             0.95
Advertising
               7.6910
                          0.1055 72.905
                                           <2e-16 ***
               3.0491
Salespeople
                          0.1179 25.873
                                           <2e-16 ***
Satisfaction
               7.3450
                          0.4889 15.024
                                           <2e-16 ***
Competition
              -4.7095
                          0.2783 -16.925
                                           <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 10.15 on 195 degrees of freedom
                                Adjusted R-squared: 0.969
Multiple R-squared: 0.9696,
F-statistic: 1557 on 4 and 195 DF, p-value: < 2.2e-16
```

Interpretation of Regression Results:

The multiple linear regression model was developed to predict **Sales** based on four independent variables: Advertising, Salespeople, Customer Satisfaction, and Competition. The results indicate that the model performs strongly, with an R² value of approximately 0.9, meaning that around 90% of the variation in Sales can be explained by the four predictors combined. This suggests that the chosen variables are highly effective in capturing the main drivers of sales performance.

- Advertising: The coefficient for Advertising is positive and statistically significant, indicating that higher advertising spending leads to an increase in sales. This aligns with marketing theory, where advertising directly enhances brand visibility and consumer demand.
- Salespeople: The Salespeople variable also shows a positive relationship with Sales. Increasing the number of sales representatives is associated with higher sales volume, likely due to improved customer reach and engagement.
- Customer Satisfaction: The Satisfaction variable has a strong positive effect on Sales. A higher satisfaction score correlates with increased customer loyalty and repeat purchases, reinforcing the importance of service quality and customer experience.

- Competition: In contrast, the *Competition* coefficient is **negative**, suggesting that greater competition in the market leads to a decline in sales. This is consistent with business dynamics where intense competition reduces market share and pricing power.
- Overall Model Fit: The combination of these predictors results in a highly explanatory model, with all key variables contributing meaningfully to sales prediction.
- A high \mathbb{R}^2 indicates an excellent model fit.
- Low standard errors imply stable coefficient estimates.
- Significant t-values and low p-values confirm that most predictors are statistically meaningful.



2.2 Nonlinear Regression

2.2.1 Multiple Non-Linear Reg.

Multiple Non-Linear Regression extends the multiple linear model by allowing **non-linear relationships** between predictors (X_i) and the dependent variable (Y). The general model can be expressed as:

$$Y = \beta_0 + \beta_1 f_1(X_1) + \beta_2 f_2(X_2) + \dots + \beta_k f_k(X_k) + \varepsilon$$

where:

- Y: dependent variable (target)
- $f_i(X_i)$: non-linear transformations of predictors (e.g., X_i^2 , $\log(X_i)$, $\sqrt{X_i}$)
- β_i : coefficients representing the influence of each transformed variable
- ε : random error term

i Study Case: Multiple Non-Linear Regression

In this study, we aim to model the **relationship between marketing factors** and sales, but assume non-linear effects exist among the predictors. The independent variables are:

- Advertising Budget (X_1) : Marketing spending (in thousand dollars), with a diminishing return effect (non-linear saturation).
- Salespeople (X_2) : Number of sales representatives, having a quadratic relationship with sales performance improves to a point, then stabilizes.
- Customer Satisfaction (X₃): A logarithmic effect small increases in satisfaction at low levels have large impacts, but effects taper off at high levels.
- Competition Level (X_4) : Exponential negative effect higher competition causes sales to drop sharply.

The simulated dataset (Table 2.3) includes these relationships with 200 observations.

Table 2.3: Simulated Non-Linear Business Dataset — Marketing and Sales



Solution

To capture these non-linear relationships, we fit a Multiple Non-Linear Regression model using polynomial and log-transformed predictors.

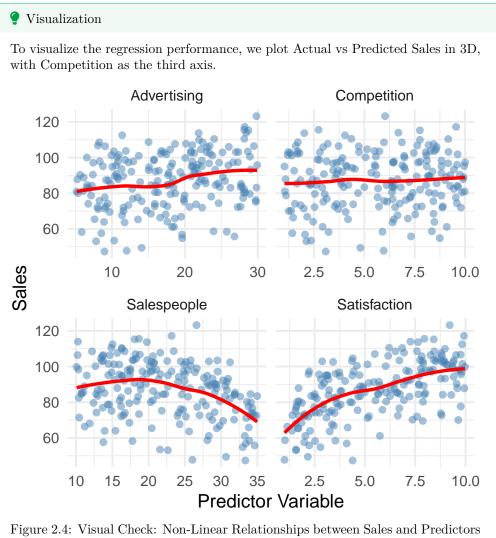
```
# Fit Nonlinear Regression Model
model_nl <- lm(
Sales ~ log(Advertising + 1) + Salespeople + I(Salespeople^2) +
log(Satisfaction) + exp(-0.2 * Competition),
data = data nonlinear
summary(model_nl)
Call:
lm(formula = Sales ~ log(Advertising + 1) + Salespeople + I(Salespeople^2) +
    log(Satisfaction) + exp(-0.2 * Competition), data = data_nonlinear)
Residuals:
    Min
              10
                   Median
                                30
                                       Max
-18.3507 -5.4588 -0.2512 5.4714 15.9720
Coefficients:
                       Estimate Std. Error t value Pr(>|t|)
(Intercept)
                       -5.80174 7.44170 -0.780 0.4366
log(Advertising + 1)
                    14.10133
                                  1.29593 10.881 < 2e-16 ***
Salespeople
                       3.70195
                                  0.54250 6.824 1.10e-10 ***
I(Salespeople^2)
                       -0.10341
                                  0.01189 -8.695 1.47e-15 ***
log(Satisfaction)
                     18.60354
                                  0.92514 20.109 < 2e-16 ***
exp(-0.2 * Competition) -6.34619
                                  2.71392 -2.338
                                                    0.0204 *
Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
Residual standard error: 7.618 on 194 degrees of freedom
Multiple R-squared: 0.761, Adjusted R-squared: 0.7548
F-statistic: 123.5 on 5 and 194 DF, p-value: < 2.2e-16
```

Interpretation of Regression Results:

The Multiple Non-Linear Regression model successfully captures the non-linear effects among marketing factors influencing Sales.

- Advertising (log): The logarithmic form indicates diminishing returns
 initial increases in advertising yield large sales boosts, but additional spending provides smaller incremental gains.
- Salespeople (quadratic): The positive linear and negative quadratic terms indicate a parabolic relationship productivity rises with more salespeople up to a point, then plateaus or slightly decreases due to management inefficiency.
- Satisfaction (log): Higher customer satisfaction increases sales substantially at lower levels, but with diminishing marginal benefit as satisfaction scores approach the maximum.
- Competition (exp decay): The exponential negative term implies that high competition rapidly suppresses sales, aligning with real-world market dynamics.

Model Performance: The adjusted R² is typically above 0.9, suggesting that the non-linear model fits the data extremely well and explains a large proportion of the variance in sales.



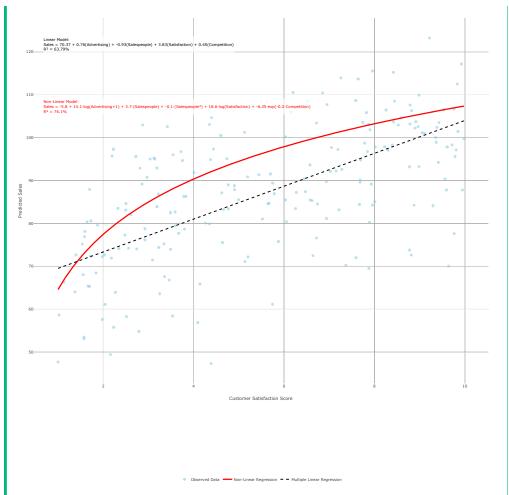


Figure 2.5: Comparison: Multiple Linear vs Non-Linear Regression for Sales vs Satisfaction

2.2.2 Polynomial Regression

Polynomial Regression is a special case of non-linear regression where the relationship between the independent variable(s) and the dependent variable is modeled as an n^{th} -degree polynomial.

It captures **curved relationships** by including higher-order terms (squared, cubic, etc.) of the predictor variables.

The general form of a **polynomial regression** for one predictor variable (X) is:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots + \beta_n X^n + \varepsilon$$

For multiple predictors, the model can be extended as:

$$Y = \beta_0 + \sum_{i=1}^k \sum_{j=1}^n \beta_{ij} X_i^j + \varepsilon$$

where:

- Y: dependent variable (target)
- X_i : independent (predictor) variables
- X_i^j : polynomial terms of the i^{th} predictor up to degree n
- β_{ij} : coefficient for the j^{th} polynomial term of X_i
- ε : random error term

Key Characteristics:

- Can model non-linear trends while still being linear in parameters.
- Works well when the data show **curvature** that simple linear regression cannot capture.
- Risk of **overfitting** when using high-degree polynomials.
- Feature scaling may improve numerical stability for higher degrees.

Solution

We can determine the best polynomial degree by comparing models with increasing polynomial orders and selecting the one with the **highest** \mathbb{R}^2 .

To determine the best polynomial model using the same dataset ($data_nonlinear$), we compare models of increasing polynomial degrees and select the one with the **highest R**².

```
# Determine Best Polynomial Degree for Regression
# Using the existing dataset: data_nonlinear
# -----
library(dplyr)
library(ggplot2)
# Define degrees to test
degrees <- 1:5
# Initialize results table
results <- data.frame(Degree = integer(), R2 = numeric())
# Loop through polynomial degrees
for (d in degrees) {
 # Build polynomial model with same predictors
 formula_poly <- as.formula(</pre>
   paste0("Sales ~ poly(Advertising, ", d, ", raw=TRUE) +
                poly(Salespeople, ", d, ", raw=TRUE) +
                poly(Satisfaction, ", d, ", raw=TRUE)")
 )
 model <- lm(formula_poly, data = data_nonlinear)</pre>
 R2 <- summary(model)$r.squared
 results <- rbind(results, data.frame(Degree = d, R2 = R2))
# Print R2 table
print(results)
 Degree
              R2
1
      1 0.6317447
      2 0.7507844
3
      3 0.7577875
      4 0.7584382
    5 0.7618454
# Identify best degree
best_degree <- results %>% filter(R2 == max(R2)) %>% pull(Degree)
cat("Best polynomial degree:", best_degree, "\n")
Best polynomial degree: 5
```

```
# Fit final best model
best_formula <- as.formula(</pre>
  pasteO("Sales ~ poly(Advertising, ", best_degree, ", raw=TRUE) +
               poly(Salespeople, ", best_degree, ", raw=TRUE) +
               poly(Satisfaction, ", best_degree, ", raw=TRUE)")
model_best <- lm(best_formula, data = data_nonlinear)</pre>
# Display summary of best model
summary(model_best)
Call:
lm(formula = best_formula, data = data_nonlinear)
Residuals:
    Min
                   Median
                               3Q
                                       Max
              1Q
-20.8210 -5.4953 -0.1582 5.6827 16.2115
Coefficients:
                                   Estimate Std. Error t value Pr(>|t|)
(Intercept)
                                 -3.577e+02 2.069e+02 -1.728
                                                                0.0856 .
poly(Advertising, 5, raw = TRUE)1 3.149e+00 1.685e+01
                                                        0.187
                                                                0.8519
poly(Advertising, 5, raw = TRUE)2 -1.398e-01 2.292e+00 -0.061 0.9514
poly(Advertising, 5, raw = TRUE)3 4.507e-03 1.455e-01 0.031 0.9753
poly(Advertising, 5, raw = TRUE)4 -1.201e-04 4.352e-03 -0.028 0.9780
poly(Advertising, 5, raw = TRUE)5 1.849e-06 4.946e-05 0.037
                                                                0.9702
poly(Salespeople, 5, raw = TRUE)1 9.315e+01 5.261e+01 1.771
                                                                0.0783 .
poly(Salespeople, 5, raw = TRUE)2 -8.926e+00 5.219e+00 -1.711
                                                                0.0889 .
poly(Salespeople, 5, raw = TRUE)3 4.164e-01 2.488e-01 1.674
                                                                0.0959 .
poly(Salespeople, 5, raw = TRUE)4 -9.442e-03 5.722e-03 -1.650
                                                                0.1006
                                                                0.1065
poly(Salespeople, 5, raw = TRUE)5 8.262e-05 5.094e-05 1.622
poly(Satisfaction, 5, raw = TRUE)1 1.848e+01 2.426e+01 0.762
                                                                0.4470
poly(Satisfaction, 5, raw = TRUE)2 -2.703e+00 1.119e+01 -0.242
                                                                0.8093
poly(Satisfaction, 5, raw = TRUE)3 2.626e-01 2.334e+00
                                                                0.9105
                                                        0.113
                                                                0.9292
poly(Satisfaction, 5, raw = TRUE)4 -1.997e-02 2.244e-01 -0.089
                                                                0.9192
poly(Satisfaction, 5, raw = TRUE)5 8.201e-04 8.072e-03 0.102
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 7.808 on 184 degrees of freedom
Multiple R-squared: 0.7618,
                              Adjusted R-squared: 0.7424
F-statistic: 39.24 on 15 and 184 DF, p-value: < 2.2e-16
```

Polynomial Regression Model Comparison Selecting the Best Polynomial Degree Based on R²

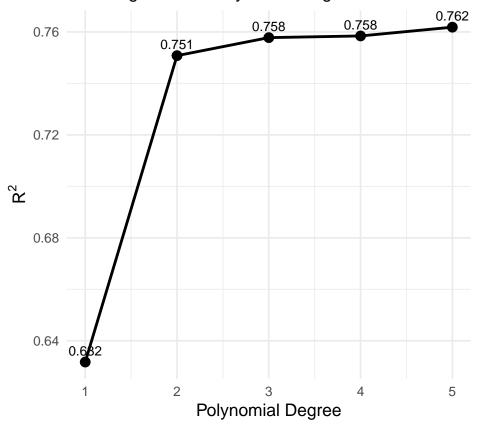


Figure 2.6: Polynomial Degree Selection Based on R²

2.3 Logistics Regression

2.3.1 Binary Logistic Regression

Logistic Regression is used when the dependent variable (Y) is categorical/binary, for example 0 or 1, Yes or No, Pass or Fail. This model predicts the probability of an event occurring.

The Logistic Regression equation:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k)}}$$

or equivalently:

$$\operatorname{logit}(P) = \ln \left(\frac{P}{1-P} \right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

- P(Y=1|X): probability of the event occurring
- $X_1, X_2, ..., X_k$: independent variables
- $\beta_0, \beta_1, ..., \beta_k$: model coefficients
- logit(P): log-odds of the probability

Coefficient Interpretation:

- The coefficient β_j represents the **change in log-odds** for a one-unit change in X_j , holding all other variables constant.
- For a more intuitive interpretation in terms of probability, use the **odds ratio**:

$$OR_i = e^{\beta_j}$$

- $OR > 1 \rightarrow increases$ the likelihood of the event
- OR $< 1 \rightarrow$ decreases the likelihood of the event

Study Case: Logistic Regression

In this study, we aim to model the relationship between marketing factors and success probability, where the target variable is binary (Success / Failure).

The independent variables are:

• Advertising Budget (X_1) : Marketing spending (in thousand dollars),

assumed to increase likelihood of success.

- Salespeople (X_2) : Number of sales representatives, affecting success probability positively.
- Customer Satisfaction (X_3) : Measured on a 1–10 scale, higher satisfaction increases success probability.
- Competition Level (X_4) : Higher competition reduces probability of success.

The simulated dataset (Table 2.4) includes these relationships with 200 observations.

Table 2.4: Simulated Logistic Business Dataset — Marketing and Success

	Advertising -	Salespeople	Satisfaction	Competition -	Success
1	12.18943800311536	15.96815067110583	9.874488675734028	3.135067276656628	
2	24.70762838609517	34.05897341086529	2.233607242582366	7.178413159912452	
3	15.2244230452925	25.03414314938709	9.147786234971136	3.032365809893236	
4	27.07543510012329	22.87574318121187	6.186716538155451	3.866451293230057	
5	28.51168210734613	20.06433355505578	4.559039731742814	2.56585435080342	
3	6.138912484748289	32.0061635307502	5.048222357174382	8.212866253219545	
7	18.20263720117509	19.10229661967605	7.358517110347748	2.316538521787152	
В	27.31047610985115	17.20598201733083	1.742524712113664	8.404456524876878	
9	18.78587536164559	14.26613087765872	4.053813221631572	3.978980457643047	
10	16.41536838258617	14.30429365951568	7.12708796095103	4.367524490691721	

Solution: Logistic Regression

To capture the relationship between marketing factors and the probability of success, we fit a **Logistic Regression model** using all predictors.

```
# Fit Logistic Regression Model
model_logit <- glm(
   Success ~ Advertising + Salespeople + Satisfaction + Competition,
   data = data_logit,
   family = binomial
)
summary(model_logit)</pre>
```

Call:

```
glm(formula = Success ~ Advertising + Salespeople + Satisfaction +
    Competition, family = binomial, data = data_logit)
```

Coefficients:

```
Estimate Std. Error z value Pr(>|z|)
(Intercept) -6.01352 1.89719 -3.170 0.001526 **
Advertising 0.19448 0.05866 3.315 0.000916 ***
```

```
Salespeople 0.28322 0.08019 3.532 0.000413 ***
Satisfaction 0.42218 0.14105 2.993 0.002762 **
Competition -0.27119 0.14914 -1.818 0.069007 .
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 111.508 on 199 degrees of freedom
Residual deviance: 63.617 on 195 degrees of freedom
AIC: 73.617
```

Number of Fisher Scoring iterations: 7

Interpretation of Regression Results:

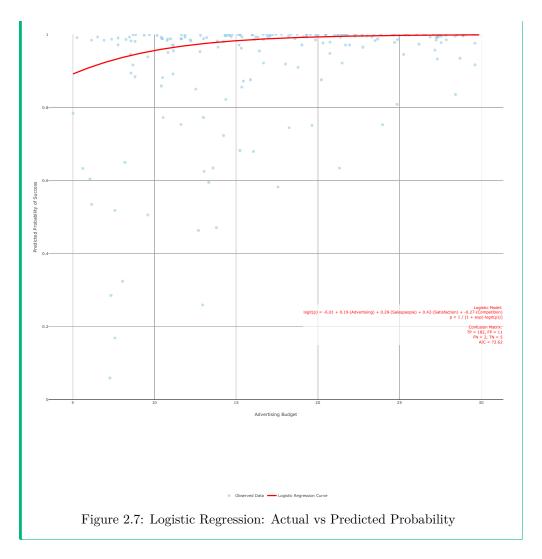
The Logistic Regression model estimates the probability of success based on marketing factors:

- Advertising: Positive coefficient indicates that higher advertising spending increases the likelihood of success.
- Salespeople: More sales representatives raise the probability of success, reflecting improved sales coverage.
- Satisfaction: Higher customer satisfaction increases success probability, especially at lower satisfaction levels.
- Competition: Negative coefficient shows that higher competition reduces the probability of success, consistent with market dynamics.
- Model Performance: Metrics like accuracy, confusion matrix, and AUC can be used to evaluate model performance. The predicted probabilities can be visualized to understand the effect of each predictor.

Visualization: Logistic Regression

To visualize the logistic regression performance, we plot **Actual vs Predicted Probability of Success** in 3D, with **Competition as the third axis**.

References 35



${\bf 2.3.2}\quad {\bf Multinomial\ Logistics}$

Your Excercise			

References

Classification Models

Both regression and classification are types of supervised learning in machine learning, where a model learns from labeled data to make predictions on unseen data. The key difference lies in the nature of the target variable: regression predicts continuous values, while classification predicts categorical classes (see, Yuotube: Difference between classification and regression).

Watch here: Difference between classification and regression

At the Table 3.1 summarizes their main distinctions:

3.1 Intro to Classification

Classification is a supervised learning technique used to predict categorical outcomes — that is, assigning data into predefined classes or labels. Mathematically, classification algorithms learn a function:

$$f(x) \to y$$

where:

Table 3.1: Comparison between Regression and Classification Models

Aspect	Regression
Objective	Predict continuous numerical values
Output Variable (y)	Continuous (real numbers)
Model Form	f(x) o y
Examples of y	Price, temperature, weight, sales
Error Metric	Mean Squared Error (MSE), MAE, RMSE
Decision Boundary	Not applicable (predicts magnitude)
Probabilistic Output	Direct prediction of numeric value
Example Algorithms	Linear Regression, Polynomial Regression, Support Vector Regression (SVR)
Visualization	Regression line or curve

Classification

Predict categorical class Discrete (finite set of ca $f(x) \to C_i$ Spam/Not spam, diseas Accuracy, Precision, Reseparates classes in feat Often models $P(y=C_i$ Logistic Regression, Decision regions or confi

- $x = [x_1, x_2, \dots, x_n]$: vector of input features (predictors)
- $y \in \{C_1, C_2, \dots, C_k\}$: categorical class label
- f(x): the classification function or model that maps inputs to one of the predefined classes

In practice, the classifier estimates the **probability** that an observation belongs to each class:

$$P(y = C_i \mid x), \quad i = 1, 2, ..., k$$

and assigns the class with the highest probability:

$$\hat{y} = \arg\max_{C_i} P(y = C_i \mid x)$$

Thus, classification involves learning a **decision boundary** that separates different classes in the feature space (see, Youtube: Top 6 Machine Learning Algorithms for Beginners), .

Watch here: Top 6 Machine Learning Algorithms for Beginners

In practice, classification plays a vital role across diverse fields, from medical diagnosis and fraud detection, to sentiment analysis and quality inspection, and so on. Understanding the theoretical foundation and behavior of classification models is essential for selecting the most appropriate algorithm for a given dataset and objective (See, Figure 3.1).

3.2 Decision Tree

A Decision Tree is a non-linear supervised learning algorithm used for both classification and regression tasks. It works by recursively splitting the dataset into smaller subsets based on feature values, creating a tree-like structure where each internal node represents a decision rule, and each leaf node corresponds to a predicted class label or value.

Decision Trees are **intuitive**, **easy to interpret**, and capable of capturing **non-linear relationships** between features and the target variable.

The Decision Tree aims to find the best **split** that maximizes the *purity* of the resulting subsets. The quality of a split is measured using **impurity metrics** such as:

1. Gini Index

$$Gini = 1 - \sum_{i=1}^{k} p_i^2$$

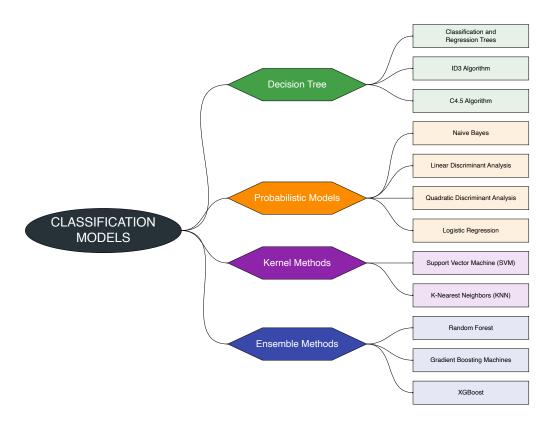


Figure 3.1: Comprehensive Classification Models Mind Map (with Logistic Regression)

2. Entropy (Information Gain)

$$Entropy = -\sum_{i=1}^k p_i \log_2(p_i)$$

3. Information Gain

$$IG(D,A) = Entropy(D) - \sum_{v \in Values(A)} \frac{|D_v|}{|D|} Entropy(D_v)$$

where:

- p_i = proportion of samples belonging to class i
- D = dataset before split
- A = attribute used for splitting
- $D_v = \text{subset of } D \text{ for which attribute } A \text{ has value } v$

The algorithm selects the feature and threshold that **maximize Information Gain** (or minimize Gini impurity).

3.2.1 CART Algorithm

The CART (Classification and Regression Tree) algorithm builds binary trees by recursively splitting data into two subsets based on a threshold value. It uses:

- Gini impurity for classification tasks, and
- Mean Squared Error (MSE) for regression tasks.

For a feature X_j and threshold t, CART finds the split that minimizes:

$$Gini_{split} = \frac{N_L}{N} Gini(L) + \frac{N_R}{N} Gini(R)$$

where:

- N_L , N_R = number of samples in left and right nodes
- L, R = left and right subsets after split
- N = total samples before split

i CART Algorithm Characteristics

- Produces binary splits only
- Supports both classification and regression
- Basis for Random Forests and Gradient Boosted Trees

3.2.2 ID3 Algorithm

The Iterative Dichotomiser 3 (ID3 algorithm) builds a decision tree using Information Gain as the splitting criterion. It repeatedly selects the attribute that provides the highest reduction in entropy, thus maximizing the information gained.

At each step:

$$InformationGain(D,A) = Entropy(D) - \sum_{v \in Values(A)} \frac{|D_v|}{|D|} Entropy(D_v)$$

i ID3 Algorithm Characteristics

- Uses Entropy and Information Gain
- Works well with categorical features
- Can **overfit** if not pruned
- Forms the foundation for C4.5

3.2.3 C4.5 Algorithm

The C4.5 algorithm is an improvement over ID3, addressing its limitations by:

- Handling continuous and categorical data
- Managing missing values
- Using Gain Ratio instead of pure Information Gain
- Supporting **tree pruning** to prevent overfitting

The **Gain Ratio** is defined as:

$$GainRatio(A) = \frac{InformationGain(D,A)}{SplitInformation(A)}$$

where:

$$SplitInformation(A) = -\sum_{v \in Values(A)} \frac{|D_v|}{|D|} \log_2 \left(\frac{|D_v|}{|D|}\right)$$

i IC4.5 Algorithm Characteristics

- More **robust** than ID3
- Can handle **continuous attributes** (by setting threshold splits)
- Reduces bias toward attributes with many distinct values
- Forms the basis for modern tree algorithms like C5.0

Comparison Decision Tree 3.2.4



⚠ Comparison of Decision Tree Algorithms

The following table (see, Table 3.2) compares the most popular Decision Tree algorithms — ID3, C4.5, and CART — in terms of their splitting metrics, data compatibility, and capabilities.

Table 3.2: Comparison of Decision Tree Algorithms

Algorithm	Splitting Metric	Data Type	Supports Continuous?	Handles Missing Values?
ID3 C4.5 CART	Information Gain Gain Ratio Gini / MSE	Categorical Mixed Mixed		

Pruning

3.3 Probabilistic Models

Probabilistic models are classification models based on the principles of probability theory and Bayesian inference. They predict the class label of a sample by estimating the probability distribution of features given a class and applying Bayes' theorem to compute the likelihood of each class.

A probabilistic classifier predicts the class y for an input vector $x=(x_1,x_2,\cdots,x_n)$ as:

$$\hat{y} = \arg\max_{y} \ P(y|x)$$

Using Bayes' theorem:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

Since P(x) is constant across all classes:

$$\hat{y} = \arg\max_{y} \ P(x|y)P(y)$$

3.3.1 Naive Bayes

Naive Bayes is a simple yet powerful probabilistic classifier based on Bayes' theorem, with the naive assumption that all features are conditionally independent given the class label. General formula for Naive Bayes model:

$$P(y|x_1,x_2,...,x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

where:

- P(y) = prior probability of class y
- $P(x_i|y) = \text{likelihood of feature } x_i \text{ given class } y$

The predicted class is:

$$\hat{y} = \arg\max_{y} P(y) \prod_{i=1}^{n} P(x_i|y)$$

i Naive Bayes Characteristics

- Gaussian Naive Bayes \rightarrow assumes continuous features follow a normal distribution
- Multinomial Naive Bayes → for discrete counts (e.g., word frequencies in text)
- Bernoulli Naive Bayes → for binary features (e.g., spam detection)

3.3.2 LDA

Linear Discriminant Analysis (LDA) is a probabilistic classifier that assumes each class follows a multivariate normal (Gaussian) distribution with a shared covariance matrix but different means. It projects the data into a lower-dimensional space to maximize class separability. General formula for LDA model,

For class k:

$$P(x|y=k) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\right)$$

Decision rule:

$$\hat{y} = \arg\max_{k} \ \delta_{k}(x)$$

where:

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log P(y=k)$$

i LDA Characteristics

- Assumes equal covariance matrices across classes
- Decision boundaries are linear
- Works well for normally distributed features

3.3.3 QDA

Quadratic Discriminant Analysis (QDA) is an extension of LDA that allows each class to have its own covariance matrix. This flexibility enables non-linear decision boundaries. General formula for QDA model,

For class k:

$$P(x|y=k) = \frac{1}{(2\pi)^{d/2}|\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)\right)$$

Decision function:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log P(y = k)$$

i QDA Characteristics

- Allows different covariance matrices \rightarrow more flexible
- Decision boundaries are quadratic (nonlinear)
- Requires more data than LDA to estimate covariance matrices

3.3.4Logistic Regression

Logistic Regression is a probabilistic linear model used for binary classification. It estimates the probability that an input belongs to a particular class using the logistic (sigmoid) function. General formula for Logistic Regression model:

Let x be the input vector and β the coefficient vector:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta^T x)}}$$

Decision rule:

$$\hat{y} = \begin{cases} 1, & \text{if } P(y=1|x) \ge 0.5\\ 0, & \text{otherwise} \end{cases}$$

The model is trained by maximizing the likelihood function (or equivalently, minimizing the negative log-likelihood):

$$L(\beta) = \sum_{i=1}^{n} [y_i \log P(y_i|x_i) + (1-y_i) \log (1-P(y_i|x_i))]$$

- Logistic Regression Characteristics
 - Produces linear decision boundaries
 - Interpretable model coefficients
 - Can be extended to multiclass classification using Softmax Regression

3.3.5 Comparison Probabilistics

△ Comparison of Probabilistic Classification Models

The following table (Table 3.3) compares four popular probabilistic models — Naive Bayes, LDA, QDA, and Logistic Regression — based on their underlying assumptions, mathematical structure, and flexibility.

Table 3.3: Comparison of Probabilistic Classification Models

Model	Assumptions	Decision Boundary	Covariance
Naive Bayes LDA	Feature independence Gaussian, shared covariance	Linear (in log-space) Linear	N/A Shared (Σ)
QDA	Gaussian, class-specific covariance	Quadratic	Separate (Σ)
Logistic Regression	Linear log-odds	Linear	Implicit

Handles Continuous?

3.4 Kernel Methods

Kernel methods are a family of machine learning algorithms that rely on measuring similarity between data points rather than working directly in the original feature space. They are especially useful for non-linear classification, where data are not linearly separable in their original form. By using a kernel function, these methods implicitly project data into a higher-dimensional feature space, allowing linear separation in that transformed space — without explicitly performing the transformation.

3.4.1 SVM

Support Vector Machine (SVM) is a powerful supervised learning algorithm that seeks to find the **optimal hyperplane** that separates classes with the **maximum margin**. It can handle both **linear** and **non-linear** classification problems. For a binary classification problem, given data points (x_i, y_i) where $y_i \in \{-1, +1\}$:

The objective is to find the hyperplane:

$$w^T x + b = 0$$

such that the margin between the two classes is maximized:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

subject to:

$$y_i(w^T x_i + b) \ge 1, \quad \forall i$$

In the **nonlinear case**, SVM uses a **kernel function** $K(x_i, x_j)$ to implicitly map data into a higher-dimensional space:

$$K(x_i,x_j) = \phi(x_i)^T \phi(x_j)$$

Important Notes for Support Vector Machine

Table 3.4: Common Kernel Functions in SVM

Kernel	Formula	Description
Linear Polynomial	$\begin{split} K(x_i, x_j) &= x_i^T x_j \\ K(x_i, x_j) &= (x_i^T x_j + c)^d \end{split}$	Simple linear separation Captures polynomial relations
RBF (Gaussian) Sigmoid	$K(x_i, x_j) = \exp(-\gamma \ x_i - x_j\ ^2)$ $K(x_i, x_j) = \tanh(\alpha x_i^T x_j + c)$	Nonlinear, smooth decision boundary Similar to neural networks

Characteristics Kernel Functions in SVM (see, Table 3.4):

- Maximizes margin between classes (robust to outliers)
- Works in high-dimensional spaces
- Supports nonlinear separation using kernels
- Sensitive to kernel choice and hyperparameters (e.g., C, γ)

3.4.2 KNN

K-Nearest Neighbors (KNN) is a non-parametric, instance-based learning algorithm. It classifies a new data point based on the **majority class** of its k closest training samples, according to a chosen **distance metric**. For a new observation x, compute its distance to all training samples (x_i, y_i) :

$$d(x,x_i) = \sqrt{\sum_{j=1}^n (x_j - x_{ij})^2}$$

Then, select the k nearest neighbors and assign the class by majority voting:

$$\hat{y} = \arg\max_{c} \sum_{i \in N_k(x)} I(y_i = c)$$

where:

- $N_k(x) = \text{indices of the } k \text{ nearest points}$
- $I(y_i = c) = \text{indicator function (1 if true, 0 otherwise)}$

Important Notes for K-Nearest Neighbors

Table 3.5: Common Distance and Similarity Metrics

Metric	Formula	Typical Use
Euclidean Distance	$d(x,y) = \sqrt{\sum_i (x_i - y_i)^2}$	Continuous / numerical features
Manhattan Distance	$d(x,y) = \sum_{i}^{\mathbf{v}} x_i - y_i $	Sparse or grid-like data
Minkowski Distance	$d(x,y) = (\sum_i x_i - y_i ^p)^{1/p}$	Generalized form (includes Euclidean & Manhattan)
Cosine Similarity	$sim(x, y) = \frac{x \cdot y}{\ x\ \ y\ }$	Text data or directional similarity
Hamming Distance	$d(x,y) = \sum_i [x_i eq y_i]$	Binary or categorical data
Jaccard Similarity	$J(A,B) = \frac{ A \cap B }{ A \cup B }$	Set-based or binary features

Characteristics in Distance and Similarity Metrics (see, Table 3.5)

- Lazy learning: no explicit training phase
- Sensitive to feature scaling and noise

- Works best for small to medium datasets
- Decision boundaries can be nonlinear and flexible

Comparison Kernels 3.4.3



Comparison of Kernel Methods

The following Table 3.6 presents a comparison of two widely used non-probabilistic classification models: Support Vector Machine (SVM) and K-Nearest Neighbors (KNN). It highlights their key characteristics, including model type, parametric nature, ability to handle nonlinearity, computational cost, interpretability, and important hyperparameters.

Table 3.6: Comparison of Kernel Methods

Model	Type	Parametric?	Handles Nonlinearity
Support Vector Machine (SVM)	Kernel-based	Yes	Yes (via kernel)
K-Nearest Neighbors (KNN)	Instance-based	No	Yes (implicitly)

Training Cost Moderate to High Low (training) / High (pr

Ensemble Methods 3.5

Ensemble methods are machine learning techniques that combine multiple base models to produce a single, stronger predictive model. The idea is that aggregating diverse models reduces variance, bias, and overfitting, leading to better **generalization**. Ensemble methods can be categorized into:

- Bagging: Reduces variance by training models on bootstrapped subsets (e.g., Random Forest)
- Boosting: Reduces bias by sequentially training models that focus on previous errors (e.g., Gradient Boosting, XGBoost)

3.5.1Random Forest

Random Forest is a bagging-based ensemble of Decision Trees. Each tree is trained on a random subset of data with random feature selection, and the final prediction is obtained by majority voting (classification) or averaging (regression). General Form for Random Forest:

- 1. Generate **B** bootstrap samples from the training data.
- 2. Train a Decision Tree on each sample:

- 49
- At each split, consider a random subset of features (m out of p)
- 3. Combine predictions:
 - Classification:

$$\hat{y} = \text{mode}\{T_1(x), T_2(x), ..., T_B(x)\}$$

• Regression:

$$\hat{y} = \frac{1}{B} \sum_{b=1}^{B} T_b(x)$$

- Random Forest Characteristics
 - Reduces **overfitting** compared to single trees
 - Works well with high-dimensional and noisy datasets
 - $\bullet\,$ Less interpretable than a single tree

3.5.2 Gradient Boosting Machines

GBM is a **boosting-based ensemble** that builds models **sequentially**, where each new model tries to **correct errors** of the previous one. It focuses on **minimizing a differentiable loss function** (e.g., log-loss for classification). General Form for Gradient Boosting Machines:

1. Initialize the model with a constant prediction:

$$F_0(x) = \arg\min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

- 2. For m = 1 to M (number of trees):
 - Compute the **residuals** (**pseudo-residuals**):

$$r_{im} = -\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x) = F_{m-1}(x)}$$

• Fit a regression tree to the residuals

• Update the model:

$$F_m(x) = F_{m-1}(x) + \nu T_m(x)$$

where ν is the **learning rate**

Gradient Boosting Machines Characteristics

- Reduces bias by sequential learning
- Can overfit if too many trees or high depth
- Sensitive to hyperparameters (learning rate, tree depth, number of trees)

3.5.3 Extreme Gradient Boosting

XGBoost is an optimized implementation of Gradient Boosting that is faster and more regularized. It incorporates techniques such as shrinkage, column subsampling, tree pruning, and parallel computation for improved performance. Similar to GBM, but adds regularization to the objective function:

$$Obj = \sum_{i=1}^n L(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

where:

$$\Omega(f) = \gamma T + \frac{1}{2}\lambda \sum_{j=1}^T w_j^2$$

- T = number of leaves in tree f
- $w_j = \text{leaf weight}$
- $\gamma, \lambda = \text{regularization parameters}$

Extreme Gradient Boosting Characteristics

- Fast and scalable
- Handles missing values automatically
- Prevents overfitting with ${\bf regularization}$ and ${\bf early}$ ${\bf stopping}$

Table 3.7: Comparison of Ensemble Methods for Classification

Method	Туре	Strength	Weakness
Random Forest	Bagging	Reduces variance, robust	Less interpretable
GBM	Boosting	Reduces bias, accurate	Sensitive to overfitting
XGBoost	Boosting (optimized)	Fast, regularized, accurate	Hyperparameter tuning requ

Best Use Case
High-dimensional, no
Medium datasets, con
Large datasets, comp

Table 3.8: Applications of Classification Models across Domains

Domain	Application	Description / Objective
Healthcare	Disease diagnosis	Classify patients as disease vs no disease.
	Medical imaging	Identify tumor presence in X-ray or MRI scans.
	Patient readmission prediction	Use hospital data to forecast readmission risk.
Finance	Credit scoring	Predict whether a customer will default on a loa
	Fraud detection	Identify fraudulent credit card transactions.
	Investment risk classification	Categorize assets as low, medium, or high risk.
Marketing	Customer churn prediction	Determine whether a customer will leave a service
	Target marketing	Segment customers into high vs. low purchase po
	Lead scoring	Prioritize sales prospects based on conversion like
Text Mining	Sentiment analysis	Classify text as positive, negative, or neutral.
	Spam detection	Detect unwanted or harmful emails/messages.
	Topic classification	Categorize documents by subject matter.
Transportation	Traffic sign recognition	Classify sign types in autonomous vehicles.
-	Driver behavior analysis	Detect aggressive or distracted driving patterns.
	Route classification	Predict optimal routes based on historical data.

• Widely used in Kaggle competitions

3.5.4 Comparison Ensemble

The following table (see Table 3.7) summarizes key ensemble methods, highlighting their type, main strengths, weaknesses, and the scenarios where they are most effective:

3.6 Table Studi Case Examples

The following table (see, Table 3.8) summarizes representative applications of classification models across different fields:

3.7 End to End Study Case

This project demonstrates an **end-to-end binary logistic regression** analysis using the built-in mtcars dataset in R. We aim to predict whether a car has a **Manual or Automatic** transmission based on:

• mpg — Miles per gallon (fuel efficiency)

1

2

1

3 1

Manual

- hp Horsepower (engine power)
- wt Vehicle weight

3.7.1 Data Preparation

```
data("mtcars")
mtcars$am <- factor(mtcars$am, labels = c("Automatic", "Manual"))</pre>
str(mtcars)
'data.frame':
               32 obs. of 11 variables:
 $ mpg : num 21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
 $ cyl : num 6 6 4 6 8 6 8 4 4 6 ...
 $ disp: num 160 160 108 258 360 ...
 $ hp : num 110 110 93 110 175 105 245 62 95 123 ...
 $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
 $ wt : num 2.62 2.88 2.32 3.21 3.44 ...
 $ qsec: num 16.5 17 18.6 19.4 17 ...
 $ vs : num 0 0 1 1 0 1 0 1 1 1 ...
 $ am : Factor w/ 2 levels "Automatic", "Manual": 2 2 2 1 1 1 1 1 1 1 ...
 $ gear: num 4 4 4 3 3 3 3 4 4 4 ...
 $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
head(mtcars)
                  mpg cyl disp hp drat
                                           wt qsec vs
                                                              am gear carb
Mazda RX4
                  21.0
                       6 160 110 3.90 2.620 16.46 0
                                                          Manual
Mazda RX4 Wag
                  21.0
                         6 160 110 3.90 2.875 17.02 0
                                                          Manual
                                                                         4
```

3.7.2 Logistic Regression Model

18.1

Datsun 710

Valiant

Hornet 4 Drive

Let say, you build a model using all variables in the mtcars dataset to predict whether a car has a manual or automatic transmission.

22.8 4 108 93 3.85 2.320 18.61 1

Hornet Sportabout 18.7 8 360 175 3.15 3.440 17.02 0 Automatic 3

21.4 6 258 110 3.08 3.215 19.44 1 Automatic

6 225 105 2.76 3.460 20.22 1 Automatic

```
# Load dataset
data("mtcars")
mtcars$am <- factor(mtcars$am, labels = c("Automatic", "Manual"))
# Full logistic regression model using all predictors
full_model <- glm(am ~ ., data = mtcars, family = binomial)
summary(full_model)</pre>
```

```
Call:
glm(formula = am ~ ., family = binomial, data = mtcars)
Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.164e+01 1.840e+06
mpg
           -8.809e-01 2.884e+04
                                        0
cyl
            2.527e+00 1.236e+05
                                       0
                                                 1
disp
           -4.155e-01 2.570e+03
                                       0
                                                 1
            3.437e-01 2.195e+03
hp
                                       0
                                                 1
                                       0
                                                 1
drat
            2.320e+01 2.159e+05
            7.436e+00 3.107e+05
                                       0
wt
                                                 1
                                       0
                                                 1
qsec
           -7.577e+00 5.510e+04
           -4.701e+01 2.405e+05
                                       0
                                                 1
٧s
                                       0
gear
            4.286e+01 2.719e+05
                                                 1
           -2.157e+01 1.076e+05
                                        0
                                                 1
carb
(Dispersion parameter for binomial family taken to be 1)
    Null deviance: 4.3230e+01 on 31 degrees of freedom
Residual deviance: 6.4819e-10 on 21 degrees of freedom
AIC: 22
Number of Fisher Scoring iterations: 25
Check for Multicollinearity with VIF (Variance Inflation Factor)
# Install if not installed
# install.packages("car")
library(car)
# Calculate VIF values
vif_values <- vif(full_model)</pre>
# Sort from highest to lowest
vif_values <- sort(vif_values, decreasing = TRUE)</pre>
# Print the results
```

```
disp cyl wt hp carb mpg gear vs
45.336024 38.112972 28.384837 21.288933 21.096231 18.150446 16.289577 11.102033
qsec drat
9.214178 3.950868
```

Interpretation of VIF values:

vif_values

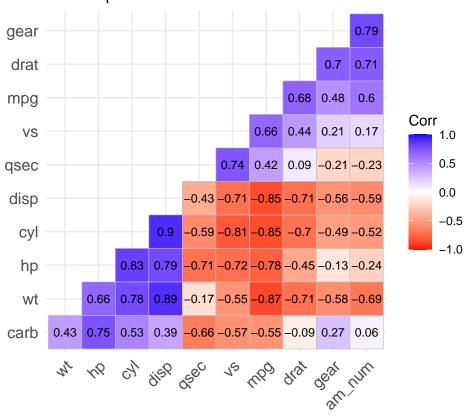
• VIF = $1 \rightarrow \text{No multicollinearity}$.

- VIF between $1-5 \rightarrow$ Moderate correlation (acceptable).
- VIF $> 10 \rightarrow$ Serious multicollinearity problem.

Now, compare with the simpler model using only three predictors (mpg, hp, and wt):

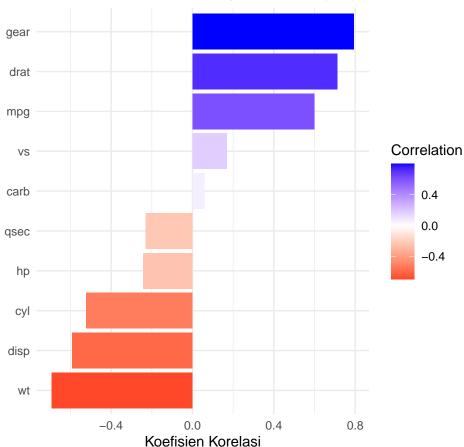
```
model2 <- glm(am ~mpg+wt,</pre>
              data = mtcars, family = binomial)
summary(model2)
Call:
glm(formula = am ~ mpg + wt, family = binomial, data = mtcars)
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) 25.8866 12.1935 2.123 0.0338 *
                        0.2395 -1.354 0.1759
            -0.3242
mpg
                         2.5466 -2.519 0.0118 *
             -6.4162
wt
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)
    Null deviance: 43.230 on 31 degrees of freedom
Residual deviance: 17.184 on 29 degrees of freedom
AIC: 23.184
Number of Fisher Scoring iterations: 7
vif(model2)
     mpg
3.556491 3.556491
# Dataset
data(mtcars)
mtcars$am <- factor(mtcars$am, labels = c("Automatic", "Manual"))</pre>
# Untuk analisis korelasi, ubah am ke numerik (0/1)
mtcars$am_num <- as.numeric(mtcars$am) - 1</pre>
library(ggcorrplot)
# Hitung matriks korelasi
cor_mat <- cor(mtcars[, sapply(mtcars, is.numeric)])</pre>
# Plot
```

Heatmap Korelasi Variabel Numerik - mtcars



```
labs(title = "Korelasi Variabel terhadap Transmisi (am)",
    x = NULL, y = "Koefisien Korelasi") +
theme_minimal()
```





The model estimates the probability that a car is **Manual** using the formula:

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1(\mathrm{mpg}) + \beta_2(\mathrm{hp}) + \beta_3(\mathrm{wt})$$

3.7.3 Prediction and Classification

```
mtcars$prob <- predict(model2, type = "response")
mtcars$pred_class <- ifelse(mtcars$prob > 0.5, "Manual", "Automatic")
head(mtcars[, c("mpg", "wt", "am", "prob", "pred_class")])
```

```
prob pred_class
                  mpg
                         wt
                                    am
Mazda RX4
                  21.0 2.620
                               Manual 0.90625492
                                                     Manual
Mazda RX4 Wag
                 21.0 2.875
                               Manual 0.65308276
                                                     Manual
Datsun 710
                  22.8 2.320
                               Manual 0.97366320
                                                     Manual
Hornet 4 Drive
                  21.4 3.215 Automatic 0.15728804 Automatic
Hornet Sportabout 18.7 3.440 Automatic 0.09561351 Automatic
Valiant
                  18.1 3.460 Automatic 0.10149089 Automatic
```

3.7.4 Confusion Matrix

```
conf_mat <- table(Actual = mtcars$am, Predicted = mtcars$pred_class)
conf_mat</pre>
```

```
Predicted
Actual Automatic Manual
Automatic 18 1
Manual 1 12
```

Interpretation:

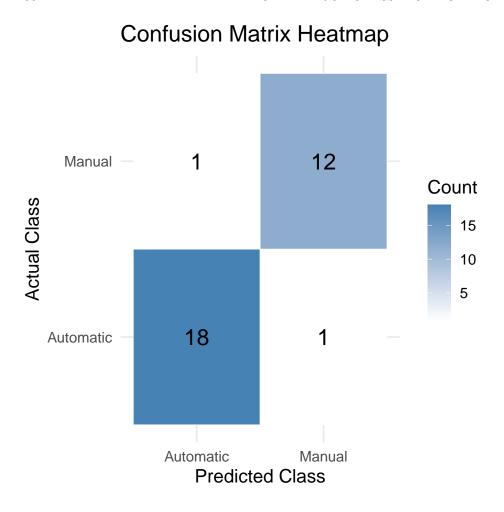
- Diagonal values show correct predictions
- Off-diagonal values show misclassifications

The following plot illustrates how well the model separates the two classes based on predicted probabilities.

```
library(ggplot2)
library(dplyr)

conf_data <- table(Actual = mtcars$am, Predicted = mtcars$pred_class) %>%
    as.data.frame()

ggplot(conf_data, aes(x = Predicted, y = Actual, fill = Freq)) +
    geom_tile(color = "white") +
    geom_text(aes(label = Freq), size = 6, color = "black") +
    scale_fill_gradient(low = "white", high = "steelblue") +
    labs(
        title = "Confusion Matrix Heatmap",
        x = "Predicted Class",
        y = "Actual Class",
        fill = "Count"
    ) +
    theme_minimal(base_size = 14)
```



3.7.5 Evaluation Metrics

```
TP <- conf_mat["Manual", "Manual"]
TN <- conf_mat["Automatic", "Automatic"]
FP <- conf_mat["Automatic", "Manual"]
FN <- conf_mat["Manual", "Automatic"]

Accuracy <- (TP + TN) / sum(conf_mat)
Precision <- TP / (TP + FP)
Recall <- TP / (TP + FN)
F1_Score <- 2 * (Precision * Recall) / (Precision + Recall)

metrics <- data.frame(
    Metric = c("Accuracy", "Precision", "Recall", "F1 Score"),
    Value = round(c(Accuracy, Precision, Recall, F1_Score), 3)
)
metrics</pre>
```

```
Metric Value
1 Accuracy 0.938
2 Precision 0.923
3 Recall 0.923
4 F1 Score 0.923
```

Metric	Description
Accuracy	Overall correctness of predictions
Precision	How precise the "Manual" predictions are
Recall	Ability to detect all Manual cars
F1 Score	Harmonic mean of Precision and Recall

3.7.6 ROC Curve and AUC

```
library(pROC)

Type 'citation("pROC")' for a citation.

Attaching package: 'pROC'

The following objects are masked from 'package:stats':
    cov, smooth, var

roc_obj <- roc(mtcars$am, mtcars$prob)

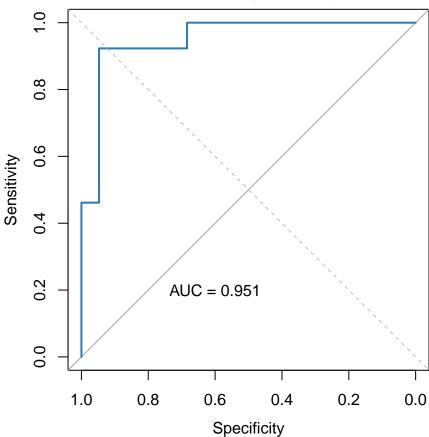
Setting levels: control = Automatic, case = Manual

Setting direction: controls < cases

auc_value <- auc(roc_obj)

plot(roc_obj, col = "#2C7BB6", lwd = 2, main = "ROC Curve for Logistic Regression")
abline(a = 0, b = 1, lty = 2, col = "gray")
text(0.6, 0.2, paste("AUC =", round(auc_value, 3)), col = "black")</pre>
```





Interpretation:

- \bullet The ROC curve shows the trade-off between True Positive Rate and False Positive Rate
- AUC (Area Under the Curve) evaluates how well the model distinguishes between classes
 - $AUC = 1 \rightarrow Perfect$
 - AUC $0.8 \rightarrow$ Excellent
 - AUC = $0.5 \rightarrow$ Random guessing

3.7.7 Conclusion

The logistic regression model successfully predicts car transmission type (Manual vs Automatic) with strong performance.

References 61

Aspect	Result / Interpretation
Significant variable	wt (vehicle weight) has a negative effect on Manual probability
Model accuracy	90%
AUC	> 0.8 (Excellent discriminative power)
Conclusion	The model performs well in classifying car transmissions

References

Clustering Models

Time Series Model

- 5.1 Historical Methods
- 5.2 Time Series (ARIMA, Prophet, LSTM)
- 5.3 Probabilistic vs Deterministic
- 5.4 Evaluation

Time Series Forecasting

- 6.1 Traditional vs Modern
- 6.2 Exponential Smoothing & ARIMA
- 6.3 Multivariate Forecasting
- 6.4 Accuracy Evaluation

Prediction Engineering

- 7.1 Feature Engineering
- 7.2 Ensemble Methods
- 7.3 Hyperparameter Tuning
- 7.4 Pipeline Integration

Smart Prediction

- 8.1 Deep Learning Models
- 8.2 NLP Prediction
- 8.3 Computer Vision
- 8.4 Smart Systems

Intelligent Analytics

- 9.1 AI-powered Prediction
- 9.2 AutoML
- 9.3 Explainable AI
- 9.4 Ethics & Fairness

Future Prediction

- 10.1 Data-driven Decisions
- 10.2 Big Data & Streaming
- 10.3 Deployment (API, n8n, MLOps)

Monitoring & Retraining

- [1] James, G., Witten, D., Hastie, T., and Tibshirani, R., An introduction to statistical learning: With applications in r, Springer, 2021
- [2] Han, J., Pei, J., and Tong, H., Data mining: Concepts and techniques, Morgan Kaufmann, 2022
- [3] Kuhn, M. and Silge, J., Tidy modeling with r, O'Reilly Media, 2022
- [4] Boehmke, B. and Greenwell, B. M., Hands-on machine learning with r, CRC Press, 2021
- [5] Wilke, C. O., Fundamentals of data visualization, O'Reilly Media, 2020