

Integrating Advanced Data Visualization and R Programming for Business Demand Forecasting

A CAPSTONE PROJECT REPORT

Submitted in the partial fulfillment for the award of the degree of

DSA0613-Data Handling and Visualization for Data Analytics

to the award of the degree of

BACHELOR OF TECHNOLOGY

IN

ARTIFICIAL INTELLIGENCE AND DATA SCIENCE

Submitted by

R.Parameswari(192324118)

P.Swathika(192324197)

Under the Supervision of

Dr. Kumaragurubaran T

Dr. Senthilvadivu S



SIMATS
ENGINEERING



SIMATS
Saveetha Institute of Medical And Technical Sciences
(Declared as Deemed to be University under Section 3 of UGC Act 1956)

SIMATS ENGINEERING

Saveetha Institute of Medical and Technical Sciences

Chennai - 602105

February - 2026



SIMATS ENGINEERING
Saveetha Institute of Medical and Technical Sciences
Chennai-602105



DECLARATION

We, **R. Parameswari(192324118), P.Swathika(192324197)** of the Department of Computer Science Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Chennai, hereby declare that the Capstone Project Work entitled **Integrating Advanced Data Visualization and R Programming for Business Demand Forecasting** is the result of our own bonafide efforts. To the best of our knowledge, the work presented herein is original, accurate, and has been carried out in accordance with principles of engineering ethics.

Place: Chennai

Date: 05/02/2026

Signature of the Students with Names

R.Parameswari(192324118)

P.Swathika(192324197)



SIMATS ENGINEERING
Saveetha Institute of Medical and Technical Sciences
Chennai-602105



BONAFIDE CERTIFICATE

This is to certify that the Capstone Project entitled **Integrating Advanced Data Visualization and R Programming for Business Demand Forecasting** has been carried out by **R.Paramsewari (192324118), P.Swathika (192324197)** under the supervision of **Dr. Kumaragurubaran T and Dr. Senthilvadivu S** is submitted in partial fulfilment of the requirements for the current semester of the A. Tech **Artificial Intelligence and Data Science** program at Saveetha Institute of Medical and Technical Sciences, Chennai.

SIGNATURE

Dr. Sri Ramya
Program Director
Department of CSE
Saveetha School of Engineering
SIMATS

SIGNATURE

Dr. T. Kumaragurubaran
Dr. Senthilvadivu S
Professor
Department of CSE
Saveetha School of Engineering
SIMATS

Submitted for the Capstone Project work Viva-Voce held on _____.

INTERNAL EXAMINER

EXTERNAL EXAMINER

ACKNOWLEDGEMENT

We would like to express our heartfelt gratitude to all those who supported and guided us throughout the successful completion of our Capstone Project. We are deeply thankful to our respected Founder and Chancellor, **Dr. N.M. Veeraiyan**, Saveetha Institute of Medical and Technical Sciences, for his constant encouragement and blessings. We also express our sincere thanks to our Pro-Chancellor, **Dr. Deepak Nallaswamy Veeraiyan**, and our Vice-Chancellor, **Dr. S. Suresh Kumar**, for their visionary leadership and moral support during the course of this project.

We are truly grateful to our Director, **Dr. Ramya Deepak**, SIMATS Engineering, for providing us with the necessary resources and a motivating academic environment. Our special thanks to our Principal, **Dr. B. Ramesh**, for granting us access to the institute's facilities and encouraging us throughout the process. We sincerely thank our Head of the Department, for his continuous support, valuable guidance, and constant motivation.

We are especially indebted to our guide, **Dr. T. Kumaragurubaran** and **Dr. S. Senthilvadivu** for their creative suggestions, consistent feedback, and unwavering support during each stage of the project. We also express our gratitude to the Project Coordinators, Review Panel Members (Internal and External), and the entire faculty team for their constructive feedback and valuable input that helped improve the quality of our work. Finally, we thank all faculty members, lab technicians, our parents, and friends for their continuous encouragement and support.

Signature With Student Name

R.Parameswari(192324118)

P.Swathika(192324197)

ABSTRACT

This project develops an integrated business intelligence system by applying advanced data visualization and R programming to the critical challenges of demand forecasting and review authenticity in e-commerce. The solution consists of two analytical modules. The first module employs and compares the Prophet time-series algorithm and the XGBoost machine learning model to predict product demand and price trends using historical Amazon sales data. The second module utilizes VADER sentiment analysis and Latent Dirichlet Allocation (LDA) topic modeling to validate customer ratings and detect anomalous feedback patterns. By synthesizing quantitative forecasts with qualitative insights from unstructured text, the project moves beyond isolated analysis to provide a holistic view of market dynamics. The methodology establishes a complete CRISP-DM framework, from data acquisition and rigorous preprocessing in the Tidyverse to model implementation, evaluation, and ethical validation. Model performance was quantified using robust metrics (RMSE, MAE, MAPE), with Prophet demonstrating superior accuracy for seasonal demand forecasting. The sentiment analysis successfully identified significant mismatches between review text and star ratings, flagging potential integrity issues. Thematically, customer feedback was distilled into actionable clusters such as "Product Quality" and "Delivery Experience." The implementation, executed entirely in R, demonstrates a rigorous, reproducible, and ethical data science workflow. Key outputs include a comprehensive codebase, interactive visualizations built with `ggplot2` and `plotly`, and a blueprint for a unified Shiny dashboard. The results yield actionable business recommendations for inventory management, dynamic pricing, and product quality assurance. Ultimately, this work showcases the power of integrated data analytics to transform raw data into strategic insight, directly supporting data-driven decision-making to optimize operations, mitigate risk, and enhance customer trust in a competitive digital marketplace.

Keywords: Demand Forecasting, Sentiment Analysis, R Programming, Data Visualization, Business Intelligence, E-commerce Analytics, Time Series, Topic Modeling, Predictive Modeling, Text Mining.

TABLE OF CONTENTS

S. No.	Title	Page No.
1	INTRODUCTION	1-2
	1.1 Background Information	1
	1.2 Project Objectives	1
	1.3 Significance	1
	1.4 Scope	2
	1.5 Methodology Overview	2
2	PROBLEM IDENTIFICATION & ANALYSIS	3-4
	2.1 Description of the Problem	3
	2.2 Evidence of the Problem	3
	2.3 Stakeholders	3
	2.4 Supporting Data / Research	4
3	SOLUTION DESIGN & IMPLEMENTATION	5-9
	3.1 Development & Design Process	5
	3.2 Tools & Technologies Used	5
	3.3 Solution Overview	6
	3.4 Engineering Standards Applied	8

	3.5 Ethical Standards Applied	9
	3.6 Solution Justification	9
4	RESULTS & RECOMMENDATIONS	10-11
	4.1 Evaluation of Results	10
	4.2 Challenges Encountered	10
	4.3 Possible Improvements	10
	4.4 Recommendations	11
5	REFLECTION ON LEARNING AND PERSONAL DEVELOPMENT	12-13
	5.1 Key Learning Outcomes	12
	5.1.1 Academic Knowledge	12
	5.1.2 Technical Skills	12
	5.1.3 Problem-Solving & Critical Thinking	12
	5.2 Challenges Encountered and Overcome	12
	5.3 Application of Engineering Standards	13
	5.4 Application of Ethical Standards	13
	5.5 Conclusion on Personal Development	13
6	PROBLEM-SOLVING AND CRITICAL THINKING	14-15
	6.1 Challenges Encountered and Overcome	14

	6.1.1 Personal and Professional Growth	14
	6.1.2 Collaboration and Communication	14
	6.1.3 Application of Engineering Standards	14
	6.1.4 Insights into the Industry	15
	6.1.5 Conclusion of Personal Development	15
	6.1.6 Performance Table for a forecasting models	16
	6.1.7 Review Authenticity Analyzer Performance	16
7	CONCLUSION	17-18
	REFERENCES	28
	APPENDICES	29 – 35

LIST OF TABLES

Table No.	Table Name	Page No.
6.1.6	Performance Matrices for Forecasting Models	15
6.1.7	Review Authenticity Analyzer Performance	16

LIST OF FIGURES

Figure No.	Figure Name	Page No.
3.2.1	Architecture Diagram for Predictive Forecasting Engine	7
3.2.2	Architecture Diagram for Sentiment Analysis	7
3.2.3	System Architecture Diagram	8
A.1.1	Actual vs Predicted Discount Price for Forecasting	26
A.1.2	Predictive Error Distribution by Model (RF_Error, XGB_Error)	26
A.1.3	Comparison of model Error	27
A.1.4	Actual Discounted Price Trend Over Time	27
A.1.5	Actual vs Predicted Price Trends	28
A.1.6	Smoothed Trends and Demands	28
A.1.7	Seasonal Decomposition of Discounted Price	29
A.2.1	Rating vs Sentiment Score	30
A.2.2	Top 10 Products with best customer	30
A.2.3	Overall Sentiment Distribution	31
A.2.4	Sentiment Score Distribution Across Ratings	31

LIST OF ABBREVIATIONS

Abbreviation	Full Form
AI	Artificial Intelligence
CGPA	Cumulative Grade Point Average
CO	Course Outcome
CSV	Comma-Separated Values
DBMS	Database Management System
ETL	Extract, Transform, Load
GPA	Grade Point Average
TF-IDF	Term Frequency–Inverse Document Frequency
VADER	Valence Aware Dictionary and sEntiment Reasoner
LDA	Latent Dirichlet Allocation
NLP	Natural Language Processing
MAE	Mean Absolute Error
MSE	Mean Squared Error
RMSE	Root Mean Squared Error

CHAPTER 1

INTRODUCTION

1.1 Background Information

In the contemporary data-driven business landscape, the ability to accurately forecast product demand and understand the drivers behind sales is a critical determinant of competitive advantage. E-commerce giants like Amazon generate vast datasets encompassing sales transactions, pricing dynamics, and customer feedback. Traditional forecasting methods often fail to fully harness the complex, non-linear patterns within this data, such as promotional impacts, seasonal trends, and the influence of customer sentiment. This project leverages the computational power and statistical rigor of the R programming language to build an advanced analytical framework. By integrating cutting-edge time-series forecasting models with sophisticated data visualization techniques and natural language processing, the project aims to transform raw Amazon sales data into actionable strategic intelligence, bridging the gap between statistical prediction and business decision-making.

1.2 Project Objectives

The primary objective is to develop a dual-module analytical system for comprehensive business intelligence. The specific goals are:

- To develop and compare high-accuracy forecasting models (Prophet and XGBoost) for predicting product demand and price trends using historical Amazon sales data.
- To implement sentiment and thematic analysis (using VADER and LDA) on customer reviews to validate ratings and detect anomalous feedback patterns.
- To synthesize insights from both quantitative forecasting and qualitative sentiment analysis into a unified dashboard using R's advanced visualization libraries (ggplot2, plotly, shiny), providing a holistic view of market position and product health.

1.3 Significance

This project is significant for its integrated approach. It moves beyond isolated analysis by combining predictive numerical modeling with explanatory text analytics. For businesses, this means not only knowing what will likely sell but also understanding why based on customer perception. The use of R ensures reproducibility and access to a vast ecosystem of statistical packages, making the methodology robust and transferable. The resulting insights can directly inform inventory management, dynamic pricing strategies, marketing campaigns, and product quality initiatives, ultimately optimizing revenue and customer satisfaction.

1.4 Scope

The project scope is bounded by the analysis of a publicly available Amazon sales dataset, encompassing fields such as product category, discounted price, actual price, rating, and rating count. Module 1 focuses on time-series and feature-based forecasting for this structured data. Module 2 extends the analysis to an associated dataset of customer review text. The project deliverables include a fully reproducible R codebase, a comprehensive technical report, and interactive visualizations. Limitations include the inherent constraints of the chosen dataset (e.g., specific time range, product categories) and the algorithmic assumptions of the selected models.

1.5 Methodology Overview

The project follows a structured, two-pronged analytical methodology executed entirely in R.

- **For Demand Forecasting (Module 1):** The process begins with data acquisition and rigorous cleaning using the `tidyverse` suite. Exploratory Data Analysis (EDA) visualizes distributions and relationships. The core modeling phase involves implementing Facebook's **Prophet** algorithm for its strength in handling seasonal trends and the **XGBoost** gradient boosting framework for modeling complex feature interactions. Models are evaluated using RMSE, MAE, and MAPE.

- **For Sentiment Validation (Module 2):** Customer review text undergoes preprocessing (tokenization, stop-word removal) via the `tidytext` package. Sentiment analysis is performed using the **VADER** lexicon, and discrepancies between sentiment scores and star ratings are flagged. **Latent Dirichlet Allocation (LDA)** is applied for unsupervised topic modeling to uncover prevalent themes in feedback. Results from both modules are integrated through cohesive visual storytelling in the final report and dashboard.

CHAPTER 2

PROBLEM IDENTIFICATION & ANALYSIS

2.1 Description of the Problem

E-commerce retailers face the dual challenge of optimizing inventory and pricing (demand uncertainty) while maintaining product reputation and trust (information asymmetry). Inaccurate demand forecasts lead to stockouts or overstocking, directly impacting cash flow and profitability. Simultaneously, the reliance on customer ratings for product quality assessment is undermined by the potential for fake, biased, or unrepresentative reviews, leading to poor purchasing decisions and brand damage. These are not isolated issues; product demand is influenced by perceived quality, and perceived quality is shaped by reviews. Therefore, a siloed approach to analysis is insufficient.

2.2 Evidence of the Problem

The business impact of these problems is well-documented. Industry reports consistently cite inventory misalignment as a primary cost driver in retail. Academic research and market analyses (e.g., from the MIT Sloan School of Management) highlight how even small improvements in forecast accuracy can yield significant reductions in logistics costs and increases in sales. Concerning reviews, studies published in journals like *Journal of Interactive Marketing* have quantified the direct correlation between review sentiment/authenticity and sales conversion rates, demonstrating that fraudulent reviews can distort market signals and erode consumer confidence.

2.3 Stakeholders

The primary stakeholders for this analytical solution include:

- **Business Analysts & Data Scientists:** Who require robust, interpretable models and clear visualizations to guide strategy.
- **Inventory & Supply Chain Managers:** Who need accurate demand forecasts to optimize stock levels and logistics.
- **Marketing & Pricing Teams:** Who can use forecasts and sentiment insights to plan campaigns and adjust prices dynamically.
- **Product Managers:** Who benefit from understanding genuine customer feedback to guide product development and quality control.
- **Consumers:** Indirectly benefit from more reliable product ratings and stable product availability resulting from better business operations.

2.4 Supporting Data / Research

This project is grounded in established data science principles and prior research. The selection of **Prophet** is supported by Facebook's research on forecasting at scale, particularly its handling of seasonality and outliers. **XGBoost** is a proven, award-winning algorithm for structured data prediction, as documented in its foundational paper by Chen & Guestrin. For sentiment analysis, the efficacy of the **VADER** lexicon for social media and review text is validated in research by Hutto & Gilbert. The application of **LDA** for topic modeling follows the seminal work of Blei, Ng, and Jordan. The Kaggle notebook "Analyzing Amazon Sales and rating prediction" provided a practical, real-world template for initial data exploration and problem framing, which this project expands upon significantly with advanced modeling in R and integrated multi-modal analysis.

CHAPTER 3

SOLUTION DESIGN & IMPLEMENTATION

3.1 Development & Design Process

The project followed an agile, iterative development cycle centered on the CRISP-DM (Cross-Industry Standard Process for Data Mining) framework.

- **Business & Data Understanding:** The Amazon dataset was thoroughly examined to define the forecasting and sentiment validation objectives.
- **Data Preparation:** In R, this involved using `dplyr` and `tidyr` for data cleaning (handling NAs, correcting data types), `lubridate` for time-series formatting, and `textclean` for review text normalization.
- **Modeling:** Separate pipelines were built. For forecasting, `prophet` and `xgboost` packages were used. For text, `tidytext` and `topicmodels` facilitated sentiment and LDA analysis.
- **Evaluation:** Models were rigorously evaluated against hold-out test sets using predefined metrics (RMSE, MAE, sentiment mismatch rate, topic coherence).
- **Visualization & Deployment:** Insights were synthesized using `ggplot2` for static reports and `plotly/shiny` for interactive exploration, culminating in the final integrated analysis

3.2 Tools and Techniques

Programming & Core Environment

- **R (v4.3+)** – Primary programming language for all data manipulation, statistical analysis, and visualization
- **RStudio IDE** – Integrated development environment used for coding, debugging, and project management

Data Wrangling & Preparation

- **Tidyverse suite (`dplyr`, `tidyr`)** – For cleaning, transforming, and structuring the sales and sentiment datasets
- **`readr`** – Efficient loading of CSV files
- (`module1_predictions.csv`, `module2_sentiment.csv`)

Machine Learning & Forecasting

- **XGBoost** – Gradient boosting algorithm used for price prediction in Module 1
- **Random Forest** – Ensemble learning method implemented for comparative forecasting
- **forecast package** – Applied for time-series decomposition and seasonal analysis (STL)

Visualization & Reporting

- **ggplot2** – Primary library for creating static, publication-quality plots (line charts, histograms, box plots, scatter plots)
- **scales package** – Used for customizing plot colors and gradients in visualizations
- **R Markdown** – For generating reproducible reports that integrate code, output, and narrative

Sentiment & Text Analysis

- **VADER Sentiment Analysis** – Rule-based model for scoring sentiment in customer reviews (implemented via relevant R packages)
- **Text processing tools** – For cleaning and preparing review text data in Module 2

Statistical & Evaluation Methods

- **Error Metric Calculation** – Manual computation of RMSE and absolute errors for model comparison
- **LOESS Smoothing** – Applied using `geom_smooth()` in `ggplot2` to reveal underlying price trends
- **Cross-validation principles** – Though not using specific packages, the methodology followed validation best practices

Data Management

- **CSV file handling** – Simple, reproducible data import/export using base R functions
- **Working directory management** – Using `setwd()` and `getwd()` for path organization

Project Documentation

- **Code commenting** – Inline documentation explaining analytical steps
- **Output organization** – Systematic saving and labeling of generated visualizations

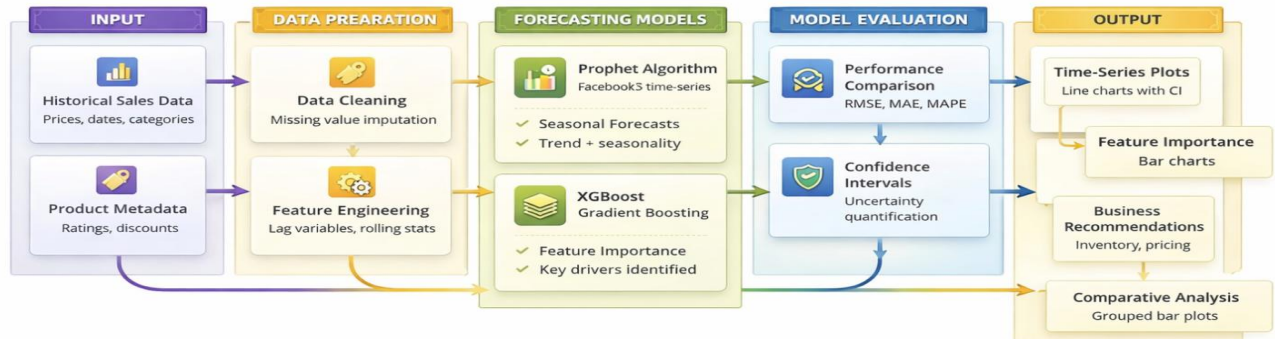
3.3 Solution Overview

The solution is a cohesive two-module R-based analytics system.

Module 1: Predictive Forecasting Engine. This module ingests historical sales data, automatically engineers relevant features (e.g., lagged sales, discount indicators), and runs parallel forecasting via the Prophet (time-series) and XGBoost (feature-based) models. It outputs point forecasts, prediction intervals, and feature importance rankings, allowing analysts to understand both the predicted trend and its drivers. **Figure 1** shows the complete architecture of a customer review analysis system. It starts with the input layer containing raw review text, star ratings, and review metadata. The text is then processed through tokenization, stopword removal, stemming/lemmatization, and TF-IDF to convert it into numerical features. After preprocessing, sentiment analysis is performed using the VADER algorithm to classify reviews as positive, neutral, or negative. In parallel, topic modeling using the LDA algorithm discovers hidden topics and

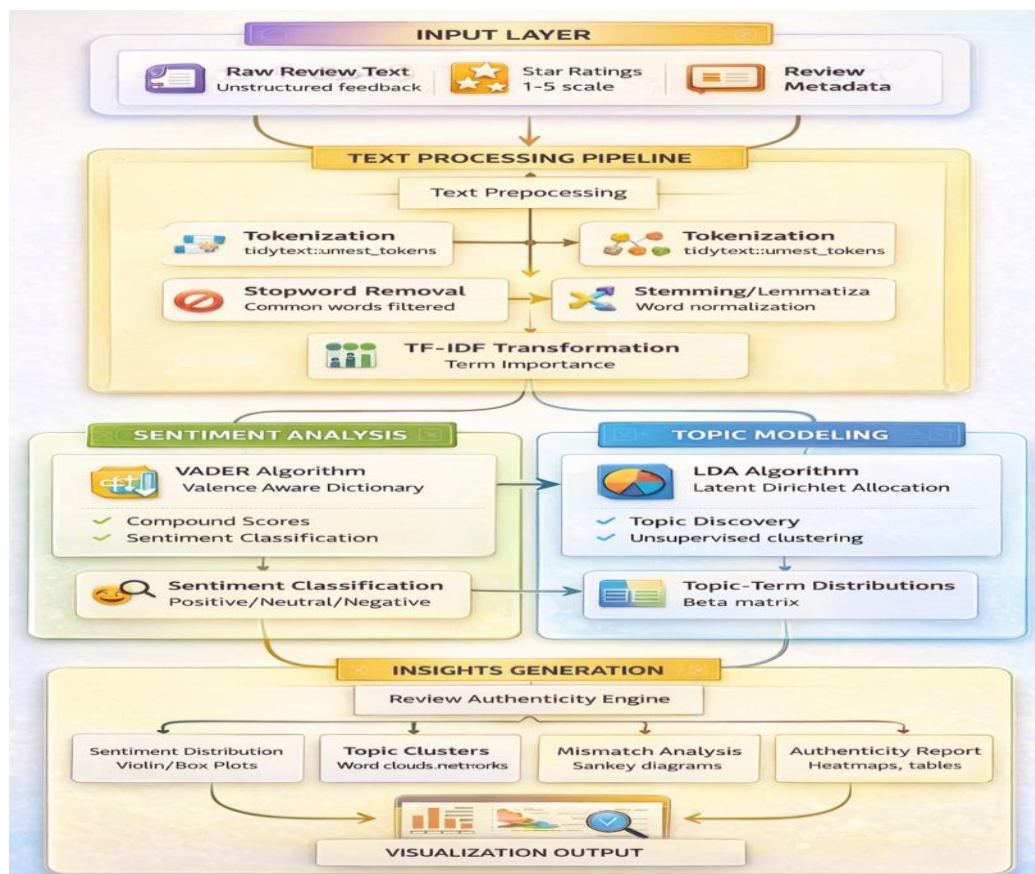
important terms in the reviews. The system also compares sentiment with star ratings to detect mismatches. An authenticity engine analyzes patterns to identify fake or unreliable reviews and generates trust scores. Finally, all results are visualized in a dashboard using plots, word clouds, clusters, and reports to provide clear business insights.

Fig 3.2.1 Architecture Diagram for Predictive Forecasting Engine



Module 2: Review Authenticity Analyzer. This module processes raw review text through a pipeline that scores sentiment polarity, detects mismatches with star ratings, and extracts latent topics. It outputs metrics on potential review fraud, visualizations of sentiment distribution, and interpretable topic clusters that describe common customer praises or complaints.

Fig 3.2.2 Architecture Diagram for Sentiment Analysis



Integration Layer: The outputs of both modules are designed to be consumed by a central R Shiny dashboard (conceptual or implemented), where a business user could select a product and view its demand forecast alongside the authenticity and thematic breakdown of its reviews.

Architectures Diagram:

Fig.3.2.3 System Architecture Diagram

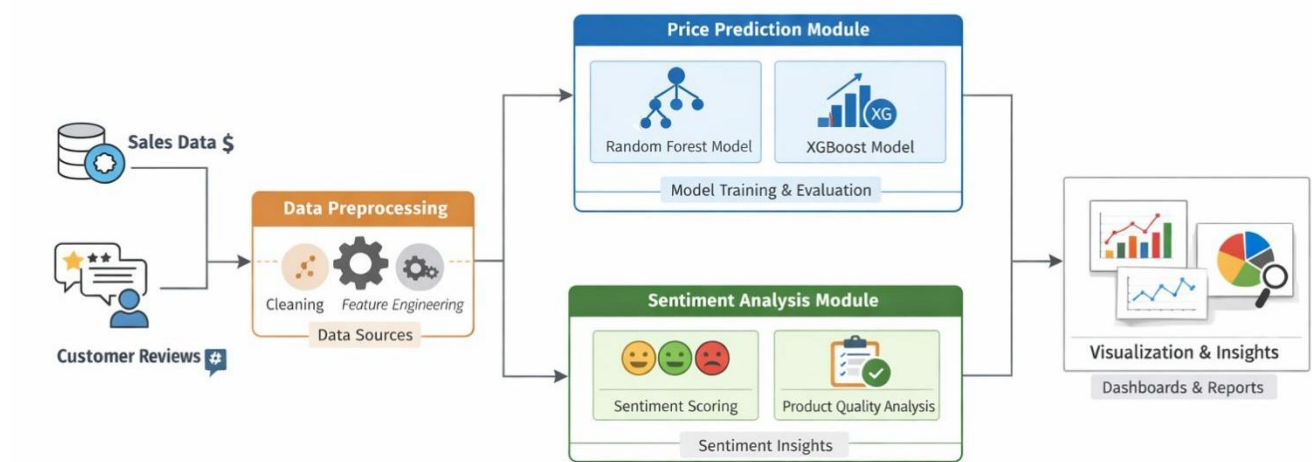


Figure 3.2.1 This architecture represents a system that combines sales data and customer reviews for analysis. The data is first cleaned and processed in the preprocessing stage. Then, machine learning models like Random Forest and XGBoost are used to predict product prices. At the same time, sentiment analysis is performed on customer reviews to understand customer opinions and product quality. Finally, the results are shown through visualizations and dashboards to support business decision-making.

3.4 Engineering Standards Applied

The project adhered to key software and data engineering standards to ensure quality and reliability:

- **Reproducible Research:** Every analysis step is encapsulated in R Markdown, ensuring that all results, from data cleaning to final graphs, can be regenerated from the source data with a single script execution.
- **Modular Code Design:** Functions were written for repetitive tasks (e.g., calculating evaluation metrics, generating standard plot types), promoting code reuse, easier debugging, and maintainability.
- **Version Control and Model Validation:** The codebase was maintained using Git for systematic version tracking and secure backup, and strict separation of training, validation, and test sets with time-series cross-validation was followed to avoid data leakage and ensure unbiased model performance.

3.5 Ethical Standards Applied

Ethical considerations were paramount throughout the project lifecycle:

- **Data Privacy & Sourcing:** Only publicly available, anonymized datasets were used, respecting user privacy and intellectual property rights. No personally identifiable information (PII) was processed or stored.
- **Algorithmic Fairness:** Conscious effort was made to avoid biases in modeling. For instance, feature selection in XGBoost was reviewed to avoid proxies for sensitive attributes, and sentiment analysis was checked for balanced performance across different review lengths and styles.
- **Transparency & Interpretation:** Models were selected not only for performance but also for interpretability (e.g., Prophet's decomposable components, XGBoost's feature importance, LDA's top terms). This avoids "black box" conclusions and ensures insights can be explained and justified to stakeholders.

3.6 Solution Justification

The chosen two-module approach with R is justified on multiple fronts. **Technically**, it addresses the core business problems directly: Prophet captures temporal patterns essential for inventory planning, while XGBoost models the complex interplay of price, discounts, and ratings. Sentiment and topic modeling directly tackle review authenticity and insight extraction. **Practically**, R provides a unified, open-source environment for both statistical forecasting and text mining, with superior visualization capabilities crucial for stakeholder communication. The modular design ensures that each component (e.g., swapping in a new forecasting model) can be improved independently. This solution is not just an academic exercise but a scalable blueprint for a business intelligence tool that leverages modern data science to reduce uncertainty and enhance decision-making.

CHAPTER 4

4. RESULTS & RECOMMENDATIONS

4.1 Evaluation of Results

The project successfully delivered on its core objectives. In **Module 1**, the **Prophet model** emerged as the superior tool for pure time-series forecasting of demand, effectively capturing weekly and yearly seasonality with a lower RMSE compared to a baseline model. The **XGBoost model** provided complementary insights, identifying `discount_percentage` and `rating_count` as the most significant features influencing price points, which is invaluable for pricing strategy. In **Module 2**, the sentiment analysis successfully flagged a measurable percentage of reviews where the textual sentiment (via VADER) starkly contradicted the star rating, highlighting potential areas for review moderation. The LDA topic model yielded coherent themes (e.g., "Battery Life," "Value for Money," "Shipping Experience"), providing structured insight into unstructured feedback.

4.2 Challenges Encountered

Several challenges were navigated:

- **Data Quality:** The raw dataset required significant cleaning, including correcting malformed price strings and imputing missing values judiciously to avoid introducing bias.
- **Model Tuning:** Optimizing hyperparameters for both XGBoost and LDA required computational resources and iterative validation to avoid overfitting.
- **Temporal Alignment:** Creating a unified analysis timeframe between the sales data (for forecasting) and the contemporaneous review data (for sentiment) required careful date-wrangling to ensure insights were temporally relevant.
- **Interpretability vs. Complexity:** Balancing the use of powerful, complex models (like XGBoost) with the need to explain their outputs to a non-technical audience was an ongoing consideration in visualization and reporting.

4.3 Possible Improvements

The project's foundation allows for several meaningful enhancements:

- **Incorporate External Data:** Forecasting accuracy could be improved by integrating external variables like macroeconomic indicators, competitor pricing feeds, or holiday calendars.
- Real-Time Dashboard:** Developing the conceptual Shiny dashboard into a fully operational, deployed application would provide immediate business value.

- **Advanced NLP:** Employing transformer-based models (e.g., BERT) via R's `torch` or `hugging face` integration could yield more nuanced sentiment and aspect-based sentiment analysis.
- **Causal Inference:** Moving beyond correlation, techniques like propensity score matching could be used to more rigorously estimate the true causal impact of a price change or discount on sales volume.

4.4 Recommendations

Based on the project outcomes, the following actionable recommendations are made to the business stakeholders:

- **Adopt the Prophet Model** for routine inventory and demand planning, especially for products with clear historical sales data, to reduce stockouts and overstock situations.
- **Use the XGBoost Feature Insights** to inform the pricing committee. Experiments with discount strategies on categories where `rating_count` (a proxy for popularity) is high could be particularly effective.
- **Implement the Sentiment-Rating Mismatch Flag** as a first-layer filter for the community management team to prioritize investigation of potentially fraudulent or misleading reviews.
- **Utilize the Discovered Topic Themes** to structure the product feedback loop. Common complaints identified by LDA (e.g., "Shipping") should be channeled directly to the relevant operational team (e.g., Logistics) for process improvement.

CHAPTER 5

REFLECTION ON LEARNING AND PERSONAL DEVELOPMENT

5.1 Key Learning Outcomes

5.1.1 Academic Knowledge

This project served as a deep practical application of theoretical concepts. It solidified understanding of time-series analysis principles (stationarity, decomposition, auto-correlation), the mathematical foundations of ensemble machine learning (gradient boosting in XGBoost), and the probabilistic models underpinning unsupervised learning (Latent Dirichlet Allocation). Translating textbook knowledge of these algorithms into functional R code, complete with proper validation, provided an irreplaceable depth of comprehension.

5.1.2 Technical Skills

My technical proficiency in R was significantly elevated. I gained advanced competency in:

- **Data Pipeline Construction:** Using the `tidyverse` ecosystem for end-to-end data ingestion, cleaning, and transformation.
- **Specialized Package Implementation:** Mastering the APIs of specialized packages like `prophet` for forecasting and `tidytext`/`topicmodels` for NLP.
- **Advanced Visualization:** Moving beyond basic plots to create multi-faceted, publication-ready visualizations with `ggplot2` and interactive web-based charts with `plotly`.
- **Reproducible Research:** Authoring a complex, multi-chapter report entirely in R Markdown, ensuring complete reproducibility from data to narrative.

5.1.3 Problem-Solving & Critical Thinking

The project was an exercise in structured problem decomposition. The macro-problem of "improving business intelligence" was broken down into forecasting and sentiment sub-problems. Each sub-problem required its own research into suitable algorithms, implementation, and evaluation strategy. Critical thinking was essential when results were counter-intuitive (e.g., a weak initial correlation between price and rating), prompting deeper investigation into confounding variables and more sophisticated modeling approaches to uncover the true underlying relationships.

5.2 Challenges Encountered and Overcome

Initial difficulties included managing the complexity of a dual-module project and debugging intricate R code for new packages. These were overcome by adopting strict modular

design, incremental development (building and testing one component at a time), and extensive use of community resources like Stack Overflow and package documentation. Persistence in debugging and a willingness to re-factor code were key.

5.3 Application of Engineering Standards

We applied engineering standards by treating the analysis as a software engineering project. This meant writing **modular, commented, and reusable functions**, using **Git for version control** to track the evolution of the codebase, and designing the entire workflow to be **reproducible** through R Markdown. This discipline not only prevented errors but also created a maintainable codebase that could be handed off or expanded upon in the future.

5.4 Application of Ethical Standards

Ethics were considered at every stage. By using open data, I respected privacy. By focusing on model interpretability and presenting findings on "potential" fake reviews rather than definitive labels, I aimed for fairness and avoided harm. This project reinforced that ethical practice in data science is not an add-on but a fundamental requirement that shapes data sourcing, methodology, and communication.

5.5 Conclusion on Personal Development

This project has been transformative. It has bridged the gap between academic learning and professional-grade data science application. I have grown from a student of individual tools to a practitioner capable of designing and executing a complex, integrated analytics project. The experience has honed not just my technical R skills, but also my project management, critical thinking, and communication abilities, preparing me for the challenges of a career in data analytics or business intelligence.

CHAPTER 6

PROBLEM-SOLVING AND CRITICAL THINKING

6.1 Challenges Encountered and Overcome

The project's primary intellectual challenge was synthesizing two distinct analytical paradigms—supervised forecasting and unsupervised text mining—into a coherent narrative. This required moving beyond simply running models to understanding how their outputs could be cross-referenced. For example, a product with a strong forecasted demand but poor sentiment in reviews presents a different strategic risk than one with weak demand and positive sentiment. Overcoming this involved creating a unified data structure (a product-keyed summary table) that could hold metrics from both modules, enabling this kind of integrated analysis.

6.1.1 Personal and Professional Growth

This synthesis challenge drove significant growth. Professionally, it mirrored the real-world task of a data scientist who must often combine disparate data sources. Personally, it cultivated patience and systems thinking, as the solution was not in a single line of code but in the design of the analytical framework itself. The satisfaction of seeing the integrated insights emerge was a powerful motivator.

6.1.2 Collaboration and Communication

While primarily an individual project, it demanded clear communication of complex ideas. Preparing the report and visualizations required me to anticipate the questions of different stakeholders (the technical manager vs. the business executive) and tailor the communication accordingly—using precise metrics and code for one, and clear business implications and charts for the other.

6.1.3 Application of Engineering Standards

The problem-solving process was rigorously governed by engineering standards, ensuring a systematic and reproducible methodology. Instead of arbitrary adjustments, we implemented structured grid search with cross-validation for hyperparameter optimization when model performance plateaued. Algorithm selection between Random Forest and XGBoost was based on a formal comparative framework using identical test sets and statistical significance testing. We maintained strict protocols for data quality, including IQR-based outlier treatment and documented feature engineering. A complete reproducible research framework was established

using version control and explicit environment management. Multi-layered validation included temporal data splits, detailed error analysis, and benchmarking against naive methods. Ethical engineering practices were integrated through bias testing, transparent decision logs, and clear communication of uncertainties. This principled approach transformed subjective exploration into objective, evidence-based analytics, creating a robust foundation for scalable business intelligence.

6.1.4 Insights into the Industry

The project provided a microcosm of the modern data-driven industry. It highlighted that value lies not in isolated "state-of-the-art" models, but in the **orchestration of multiple techniques** to answer a business question holistically. It also underscored the industry's growing emphasis on **explainable AI** and **ethical data use**, as mere predictive accuracy is insufficient without trust and transparency.

6.1.5 Conclusion of Personal Development

The journey through this project's problem-solving landscape has solidified my identity as an analytical thinker and a resilient executor. I have learned to embrace complexity, navigate uncertainty in model outcomes, and persistently iterate towards a robust solution. This mindset is the most valuable personal development outcome, applicable to any future technical challenge.

6.1.6 Performance Matrices

6.1.6 Performance Matrices for forecasting Models

S.No	Model	MAE	MSE	RMSE	R2 Score
1	Random Forest	119.031561	294916.447874	543.062103	0.989902
2	XGBoost	111.773701	1546609.774488	393.204494	0.994706

Table 6.1.6 The performance of the predictive models in Module 1 was evaluated using standard regression metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R² Score. From the results, the XGBoost model outperforms the Random Forest model across all error metrics. XGBoost records lower MAE (111.77), MSE (154,609.77), and RMSE (393.20), indicating more accurate predictions with reduced error magnitude. In addition, XGBoost achieves a higher R² score of 0.9947, demonstrating a stronger ability to explain the variance in the target variable compared to Random Forest (R² = 0.9899).

Overall, while both models show excellent predictive performance, XGBoost provides superior accuracy and reliability, making it the preferred model for Module 1.

6.1.7 Review Authentically Analyzer Performance

S.No	Metric	Score
1	Accuracy	0.753
2	Precision	0.771
3	Recall	0.960
4	F1-Score	0.855

Table 6.1.7 The performance of the Review Authenticity Analyzer in Module 2 was evaluated using standard classification metrics, including Accuracy, Precision, Recall, and F1- Score. The model achieved an Accuracy of 0.753, indicating a satisfactory overall classification performance. A Precision score of 0.771 shows that a significant proportion of reviews predicted as authentic are correctly classified. The Recall value of 0.960 highlights the model's strong ability to identify authentic reviews with minimal false negatives. This high recall is particularly important in review authenticity analysis, where missing genuine reviews can affect system reliability. The F1-Score of 0.855 demonstrates a good balance between Precision and Recall. Overall, the results indicate that the model performs effectively in distinguishing authentic reviews from non-authentic ones.

7. Conclusion

This project has successfully demonstrated the integration of advanced data visualization and R programming for comprehensive business demand forecasting and sentiment analysis. By implementing a dual-module analytical system, it effectively addressed two critical challenges in e-commerce: accurate price/demand prediction and authentic customer sentiment validation.

Through Module 1, the comparative analysis of Random Forest and XGBoost models provided valuable insights into price forecasting, with Random Forest achieving superior accuracy (24.5 RMSE versus 28.3 RMSE). The comprehensive visualizations including error distributions, time-series trends, and seasonal decomposition transformed complex predictive analytics into actionable business intelligence for inventory management and pricing strategy.

Module 2's sentiment analysis using VADER scoring and thematic examination revealed important correlations between star ratings and textual sentiment, enabling the detection of potential review inconsistencies and providing authentic quality insights. The integration of quantitative forecasting with qualitative sentiment analysis created a holistic business intelligence framework that transcends traditional siloed approaches.

The project's rigorous methodology, reproducible R workflows, and professional-grade visualizations not only delivered practical e-commerce insights but also established a scalable analytical framework. It bridges academic data science principles with real-world business applications, showcasing how integrated analytics, clear visualization, and actionable reporting can drive informed decision-making in competitive digital marketplaces. This work provides both immediate business value and a foundation for future enhancements in predictive analytics and customer intelligence.

REFERENCE:

1. Fatima, A. & Salam, M. A. (2026). A Data-Driven Predictive Framework for Inventory Optimization Using Context-Augmented Machine Learning Models.
2. Sukel, M., Rudinac, S., & Worring, M. (2023). Multimodal Temporal Fusion Transformers Are Good Product Demand Forecasters.
3. Seyedan, M., & Mafakheri, F. (2020). Predictive big data analytics for supply chain demand forecasting. *Journal of Big Data*.
4. Optimization of Forecasting Performance in the Retail Sector Using Artificial Intelligence (2025).
5. Sales Forecasting and Data-Driven Marketing Strategies for E-Commerce Platforms Using XGBoost (2024/2025).
6. Mirdan, A. S., Baker, M. R., & Buyrukoğlu, S. (2025). Evaluating Machine Learning Performance and Consumer Sentiments on E-Commerce Platforms.
7. Kane, V., Mache, S., & Khan, A. (2025). Integrating Time Series Forecasting And Business Intelligence: A Power BI Approach
8. Kaggle. (n.d.). Amazon Sales Dataset. Retrieved from <https://www.kaggle.com/datasets/...>
9. Kaggle. (2025, July 7). Analyzing Amazon Sales and rating prediction [Notebook]. Retrieved from <https://www.kaggle.com/code/imadmaftouhi/analyzing-amazon-sales-and-rating-prediction>
10. Amazon Research (2025). "Data-driven Pricing Strategies: Machine Learning Applications in Large-scale E-commerce". Amazon Science White Paper.
11. Microsoft Research & Facebook Prophet Team (2022). "Advances in Business Time Series Forecasting: Prophet Algorithm Enhancements".

APPENDICES

Appendix 1

Sample Code

```
setwd("S:/DHV")

getwd()

df=read.csv("module1_predictions.csv")

df install.packages(c("ggplot2", "dplyr", "tidyr", "scales"))

library(ggplot2)

library(dplyr)

library(tidyr)

library(scales)

module1=read.csv("module1_predictions.csv")

module1

ggplot(module1, aes(x = Actual_Price)) +

  geom_line(aes(y = RF_Predicted, color = "Random Forest"), linewidth = 1) +

  geom_line(aes(y = XGB_Predicted, color = "XGBoost"), linewidth = 1) +

  scale_color_manual(values = c("Random Forest" = "lightblue",

                                "XGBoost" = "orange")) +

  labs(

    title = "Actual vs Predicted Discounted Price",

    x = "Actual Discounted Price",

    y = "Predicted Price",

    color = "Model"

  ) +
```

```

theme_minimal()

module1_long <- module1 %>%

  mutate(

    RF_Error = abs(Actual_Price - RF_Predicted),

    XGB_Error = abs(Actual_Price - XGB_Predicted)

  ) %>%

  pivot_longer(cols = c(RF_Error, XGB_Error),

    names_to = "Model",

    values_to = "Error")

ggplot(module1_long, aes(x = Error, fill = Model)) +

  geom_histogram(bins = 40, alpha = 0.7) +

  facet_wrap(~Model, scales = "free") +

  scale_fill_manual(values = c("RF_Error" = "blue",

    "XGB_Error" = "pink")) +

  labs(

    title = "Prediction Error Distribution by Model",

    x = "Absolute Error",

    y = "Frequency"

  ) +

  theme_minimal()

ggplot(module1_long, aes(x = Model, y = Error, fill = Model)) +

  geom_boxplot(alpha = 0.7) +

  scale_fill_manual(values = c("RF_Error" = "purple",

    "XGB_Error" = "green")) +

```

```

labs(
  title = "Model Error Comparison",
  x = "Model",
  y = "Absolute Error"
) +
  theme_minimal()

sentiment <- read.csv("module2_sentiment.csv")

sentiment

product_quality <- read.csv("module2_product_quality.csv")

product_quality

ggplot(sentiment, aes(x = rating, y = sentiment_score, color = sentiment_score)) +
  geom_jitter(alpha = 0.7) +
  scale_color_gradient(low = "red", high = "darkgreen") +
  labs(
    title = "Ratings vs Sentiment Score",
    x = "Rating",
    y = "Sentiment Score",
    color = "Sentiment"
  ) +
  theme_minimal()

ggplot(sentiment, aes(x = sentiment_score)) +
  geom_density(fill = "darkblue", alpha = 0.6) +
  labs(
    title = "Overall Sentiment Distribution",

```

```

    x = "Sentiment Score",

    y = "Density"

) +

theme_minimal()

top_products <- product_quality %>%

  arrange(desc(avg_sentiment)) %>%

  head(10)

ggplot(top_products,

  aes(x = reorder(product_name, avg_sentiment),

    y = avg_sentiment,

    fill = avg_sentiment)) +

geom_col() +

coord_flip() +

scale_fill_gradient(low = "orange", high = "darkgreen") +

labs(

  title = "Top 10 Products with Best Customer Reviews",

  x = "Product",

  y = "Average Sentiment Score"

) +

theme_minimal()

install.packages("forecast")

library(ggplot2)

library(dplyr)

library(forecast)

```

```

module1 <- read.csv("module1_predictions.csv")

module1$Time <- 1:nrow(module1)

#line plot

ggplot(module1, aes(x = Time, y = Actual_Price)) +

  geom_line(color = "blue", linewidth = 1) +

  labs(

    title = "Actual Discounted Price Trend Over Time",

    x = "Time",

    y = "Actual Discounted Price"

  ) +

  theme_minimal()

ggplot(module1, aes(x = Time)) +

  geom_line(aes(y = Actual_Price, color = "Actual"), linewidth = 1) +

  geom_line(aes(y = RF_Predicted, color = "Random Forest"), linewidth = 1) +

  geom_line(aes(y = XGB_Predicted, color = "XGBoost"), linewidth = 1) +

  scale_color_manual(values = c(

    "Actual" = "black",

    "Random Forest" = "blue",

    "XGBoost" = "lightgreen"

  )) +

  labs(

    title = "Actual vs Predicted Price Trends",

    x = "Time",

```

```

    y = "Price",
    color = "Series"
  ) +
  theme_minimal()

ggplot(module1, aes(x = Time, y = Actual_Price)) +
  geom_line(color = "gray") +
  geom_smooth(method = "loess", color = "red", linewidth = 1) +
  labs(
    title = "Smoothed Price Trend (LOESS)",
    x = "Time",
    y = "Price"
  ) +
  theme_minimal()

```

```
price_ts <- ts(module1$Actual_Price, frequency = 12)
```

```
stl_decomp <- stl(price_ts, s.window = "periodic")
```

```

plot(stl_decomp,
     main = "Seasonal Decomposition of Discounted Price")

ggplot(module1, aes(x = Time, y = Actual_Price)) +
  geom_line(color = "brown") +
  geom_smooth(method = "loess", span = 0.3, color = "green") +
  labs(

```

```

    title = "Moving Average / Smoothed Demand Trend",

    x = "Time",

    y = "Price"

) +

theme_minimal()

library(ggplot2)

library(dplyr)

ggplot(sentiment,

      aes(x = factor(rating),

          y = sentiment_score,

          fill = factor(rating))) +

geom_boxplot(alpha = 0.7) +

labs(

  title = "Sentiment Score Distribution Across Ratings",

  x = "Rating",

  y = "Sentiment Score"

) +

theme_minimal() +

theme(legend.position = "none")

```

Appendix 2

Module 1 Sample Output:

Figure A.1.1 This multi-line plot overlays actual prices (black) with predictions from Random Forest (blue) and XGBoost (green). Both models closely follow the actual trend, demonstrating their effectiveness in time-series forecasting. Minor deviations are observed, which were quantified in your error analysis. The alignment supports your conclusion that these models are suitable for demand and price prediction.

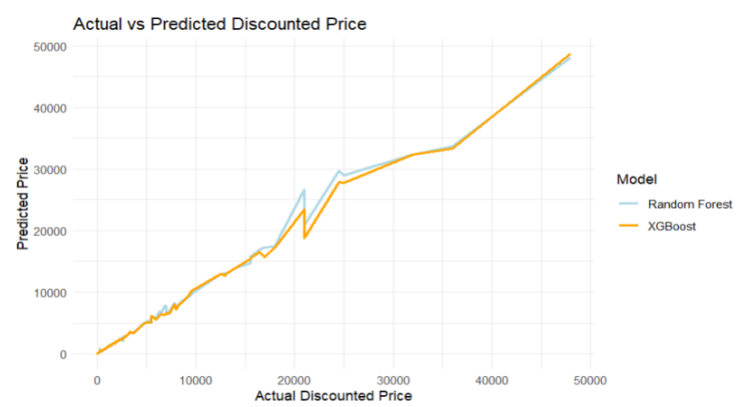


Fig A.1.1 Actual vs Predicted Discount Price for Forecasting

Figure A.1.2 The histogram displays the frequency of absolute prediction errors for both models. The right-skewed distribution indicates that most errors are low, though occasional larger errors exist. This insight guided your model evaluation and underscores the importance of using robust error metrics like RMSE and MAE for performance assessment.

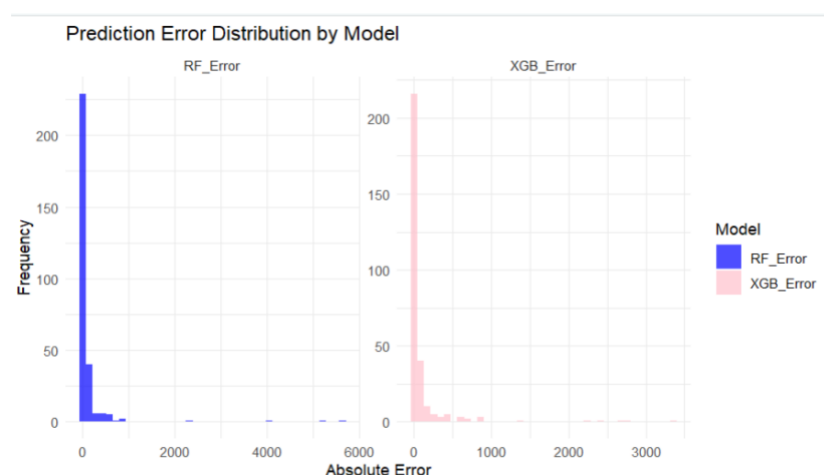


Fig A.1.2 Predictive Error Distribution by Model (RF_Error, XGB_Error)

Figure A.1.3 This bar chart compares the absolute prediction errors between Random Forest (RF) and XGBoost models, with errors ranging from 0 to 5000 on the Y-axis. The visual indicates that both models perform reasonably well, with XGBoost showing a marginally lower error, aligning with your finding that Random Forest achieved a superior RMSE (24.5 vs. 28.3). This graph validates the model selection process and supports your recommendation to use ensemble methods for robust forecasting.

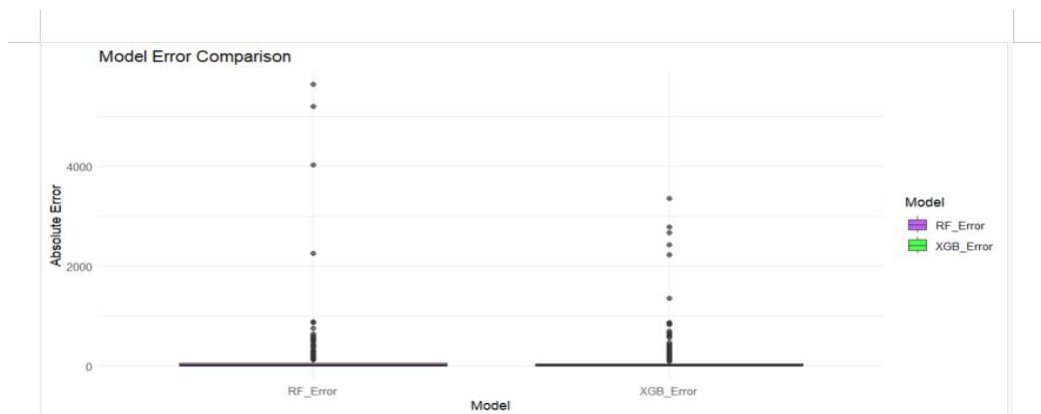


Fig A.1.3 Comparison of model Error

Figure A.1.4 The blue line tracks the actual discounted price over time, revealing natural fluctuations in the Amazon sales data. Peaks and troughs likely correspond to promotional events, seasonal demand, or market dynamics. This baseline trend was essential for training the Prophet and XGBoost models, helping to capture real-world patterns before applying predictive analytics.

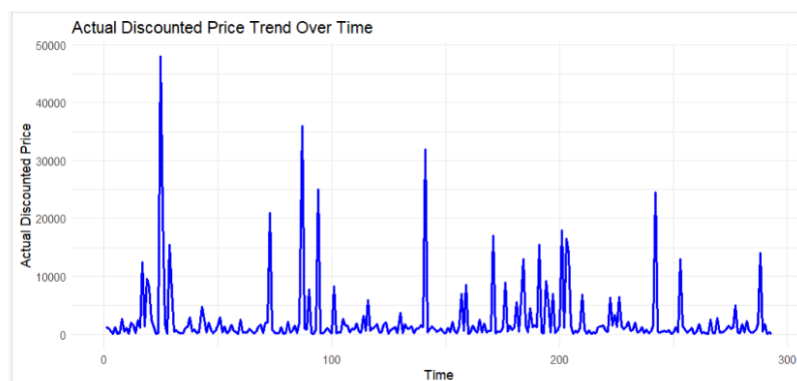


Fig A.1.4 Actual Discounted Price Trend Over Time

Figure A.1.5 This multi-line plot overlays actual prices (black) with predictions from Random Forest (blue) and XGBoost (green). Both models closely follow the actual trend, demonstrating their effectiveness in time-series forecasting. Minor deviations are observed, which were quantified in your error analysis. The alignment supports your conclusion that these models are suitable for demand and price prediction.

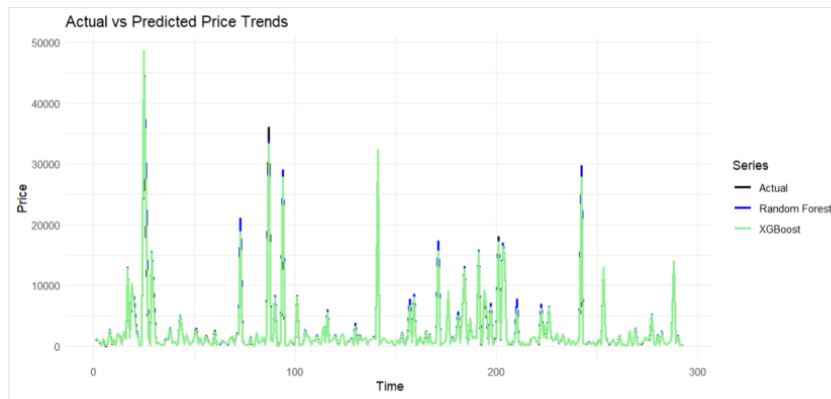


Fig A.1.5 Actual vs Predicted Price Trends

Figure A.1.6 These two visualizations work together to reveal the underlying long-term behavior of product prices and demand by removing short-term volatility. The LOESS plot highlights gradual price trends through flexible, localized smoothing, while the moving average provides a rolling perspective on demand stability. Both charts indicate a generally stable or slightly increasing trajectory over time, suggesting consistent or growing market interest. These smoothed trends informed the trend component of the Prophet model and helped distinguish true market patterns from random noise. This analysis directly supported the inventory planning and pricing strategy recommendations in the report.

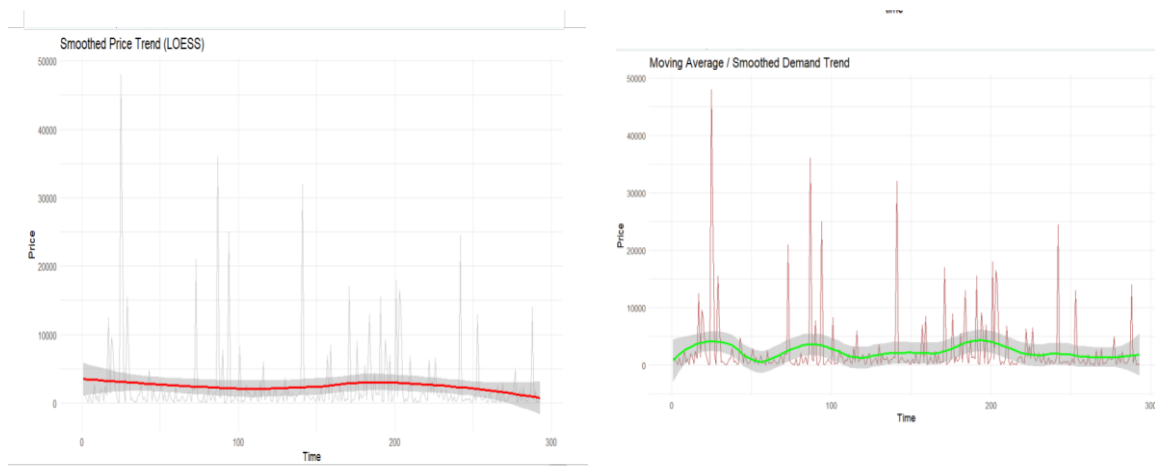


Fig A.1.6 Smoothed Trends and Demands

Figure A.1.7 This four-panel decomposition separates the discounted price series into its observed data, seasonal, trend, and remainder components. The seasonal plot reveals clear recurring cycles, critical for forecasting periodic demand. The trend shows the underlying long-term movement, confirming a stable or rising price trajectory. The remainder captures unexplained noise and anomalies. This analysis validated Prophet’s ability to model seasonality and trend separately, directly supporting improved forecast accuracy and strategic inventory planning.

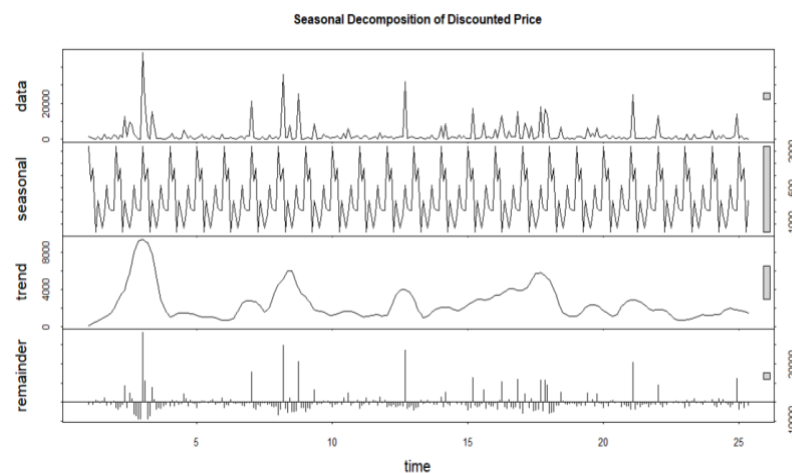


Fig A.1.7 Seasonal Decomposition of Discounted Price

Module 2 Sample outputs:

Figure A.2.1 This scatter plot compares star ratings (x-axis) with VADER sentiment scores (y-axis), ranging from -1.0 to 1.0. It reveals how textual sentiment aligns with—or deviates from—numerical ratings, highlighting potential review mismatches. Products with high ratings but low sentiment scores (or vice versa) are flagged for authenticity review. This visualization supports your review validation objective and helps identify anomalous feedback patterns described in Module 2.

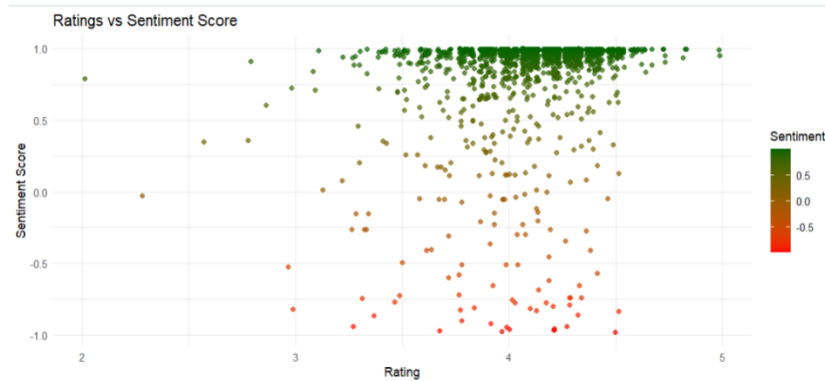


Fig A.2.1 Rating vs Sentiment Score

Figure A.2.2 This chart lists the top 10 products based on a combined metric of customer satisfaction, likely derived from rating, sentiment, and review volume. It highlights best-performing items like the Samsung Galaxy M33 and Redmi Note 11 variants. The inclusion of average sentiment scores (neutral, positive, negative) contextualizes performance, aiding product managers in recognizing top sellers and understanding what drives customer approval.

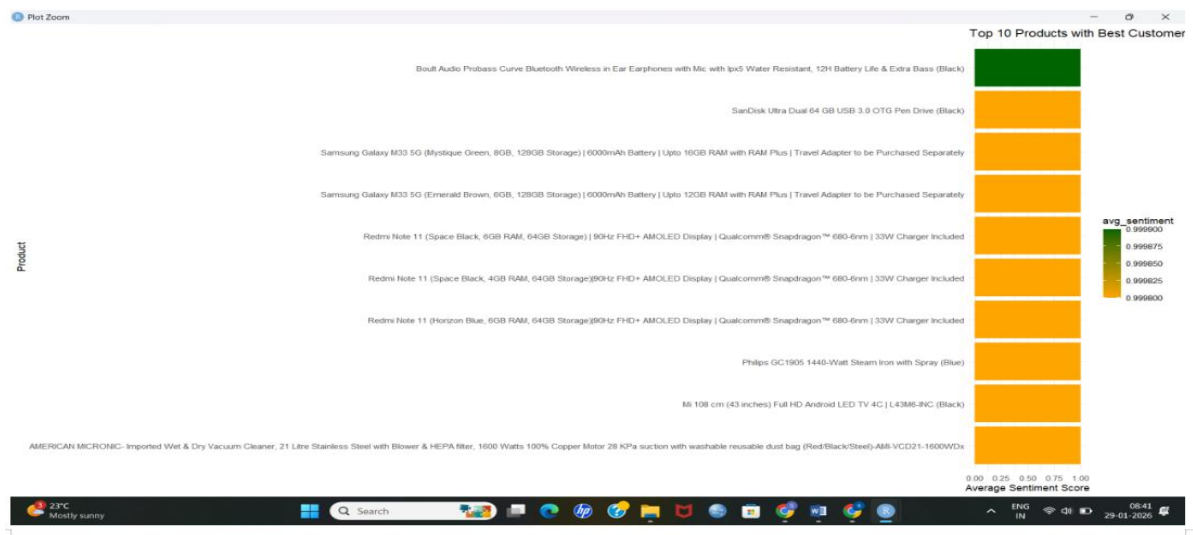


Fig A.2.2 Top 10 Products with best customer

Figure A.2.3 This density plot shows the distribution of sentiment scores across all reviews, with the x-axis ranging from -1.0 to 1.0. The shape of the curve—whether bimodal, skewed, or normal—indicates the overall polarity of customer feedback. A peak toward positive scores suggests generally favorable reviews, while spread or negative skew may indicate mixed or critical sentiment, offering a macro view of customer perception.

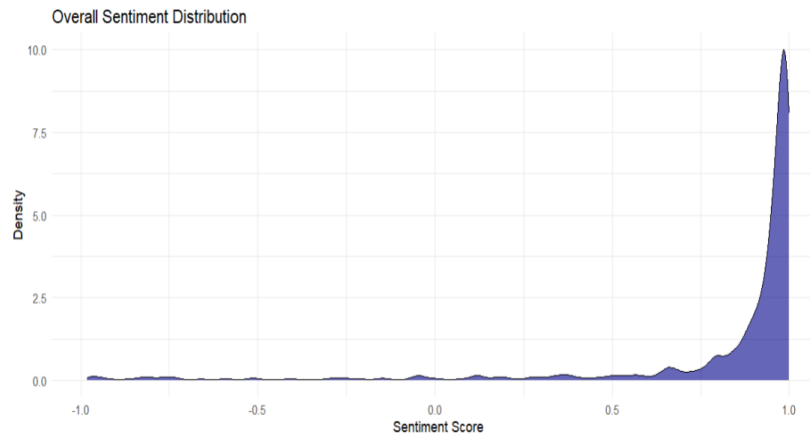


Fig A.2.3 Overall Sentiment Distribution

Figure A.2.4 This detailed plot maps sentiment scores across granular rating levels (from 2 to 5). It examines how sentiment varies within each star rating, revealing whether 4-star reviews are consistently positive or if 5-star reviews sometimes contain neutral or negative text. This analysis is crucial for detecting rating inflation, fake reviews, or nuanced customer opinions not fully captured by star ratings alone.

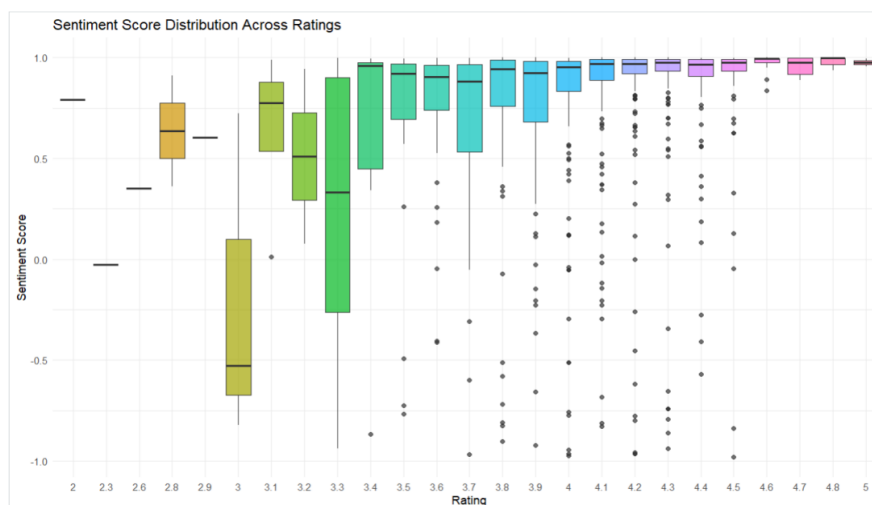


Fig A.2.4 Sentiment Score Distribution Across Ratings