A

Project Work Report

On

# "Hotel booking demand forecasting & cancellation analytics"

Submitted to

## SRI VENKATESWARA COLLEGE OF ENGINEERING AND TECHNOLOGY (AUTOMOMOUS)

Affiliated to JNTUA, Anantapur

*In partial fulfillment of the requirements for the award of the degree of*

## BACHELOR OF TECHNOLOGY

### IN

### COMPUTER SCIENCE AND ENGINEERING (AI&ML)

*during the academic year 2025-2026*

*Submitted by*

| | |
|---|---|
| P VISHNU VARDHAN REDDY | 22781A3399 |
| Y GURU MOHAN REDDY | 22781A33E6 |
| P M MANOHAR REDDY | 22781A3397 |
| T CHARITHA REDDY | 22781A33D1 |
| M VASANTH KUMAR | 23785A3312 |

Under the esteemed guidance of

*Dr. M. Lavanya, M.C.A, M.Tech, Ph.D.*

*HOD & Associate Professor*

*Department of CSE(AI&ML)*

**SRI VENKATESWARA COLLEGE OF ENGINEERING AND TECHNOLOGY(AUTONOMOUS)**

Affiliated to JNTUA, Anathapuramu-515002(A.P) & Approved by AICTE, New Delhi
Accredited by NAAC, Bengaluru & NBA, New Delhi

An ISO 9001:2000 Certified InstitutionR.V.S. Nagar, Chittoor-517127(A.P),

India www.svcetedu.org

1

# SRI VENKATESWARA COLLEGE OF ENGINEERING AND TECHNOLOGY

**(AUTONOMOUS)**

Affiliated to JNTUA, Anathapuramu-515002(A.P) & Approved by AICTE, New Delhi Accredited by NAAC, Bengaluru & NBA, New Delhi

An ISO 9001:2000 Certified Institution

R.V.S. Nagar, Chittoor-517127(A.P), India www.svcetedu.org

## CERTIFICATE



This is to certify that, the project entitled, "**hotel booking demand forcasting & cancellation analytics**" is a bonafide work carried by the following students

| | |
|---|---|
| **P VISHNU VARDHAN REDDY** | **22781A3399** |
| **Y GURU MOHAN REDDY** | **22781A33E6** |
| **P M MANOHAR REDDY** | **22781A3397** |
| **T CHARITHA REDDY** | **22781A33D1** |
| **M VASANTH KUMAR** | **23785A3312** |

in partial fulfillment of the requirement for the award of the degree **BACHELOR OF TECHNOLOGY in COMPUTER SCIENCE AND ENGINEERING (AI & ML)**

during the academic year **2025-2026**

SIGNATURE OF THE GUIDE

  Dr. M. Lavanya, MCA, M. Tech, Ph.D.
  HOD & Associate Professor

SIGNATURE OF THE HOD

  Dr. M. Lavanya, MCA, M. Tech, Ph.D.
  HOD & Associate Professor

INTERNAL EXAMINER

EXTERNAL EXAMINIER

Viva-Voce Conducted on _____

**SRI VENKATESWARA COLLEGE OF ENGINEERING AND TECHNOLOGY**

**(AUTONOMOUS)**

**Affiliated to JNTUA, Anathapuramu-515002(A.P) & Approved by AICTE, New Delhi Accredited by NAAC, Bengaluru & NBA, New Delhi**

**An ISO 9001:2000 Certified Institution**

**R.V.S. Nagar, Chittoor-517127(A.P), India www.svcetedu.org**

# Department of CSE(AI&ML)

## DECLARATION

We P. Vishnu vardhan Reddy (22781A3399), Y Guru Mohan Reddy(22781A33E6), P.M.ManoharReddy(22781A3397)T.CharithaReddy(22781A33D1)andM.VasanthKumar (23785A3307)hereby declare that the Project Report entitled "Hotel Booking Demand Forecast and cancellation prediction" under the guidance of Dr. M. Lavanya, MCA, M.Tech, Ph.D., Sri Venkateswara College of Engineering & Technology (Autonomous), Chittoor, is submitted in partial fulfillment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY in COMPUTER SCIENCE AND ENGINEERING (AI & ML).

This is a record of bonafide work carried out by us and the results embodied in this project have not been reproduced or copied from any source. The results embodied in this project report have not been submitted to any other university or institution for the award of any other degree or diploma.

| | |
|---|---|
| **P VISHNU VARDHAN REDDY** | **22781A3399** |
| **Y GURU MOHAN REDDY** | **22781A33E6** |
| **P M MANOHAR REDDY** | **22781A3397** |
| **T CHARITHA REDDY** | **22781A33D1** |
| **M VASANTH KUMAR** | **23785A3312** |

## ACKNOWLEDGEMENT

A Grateful thanks to **Dr.R.Venkataswamy**, Chairman, Sri Venkateswara College of Engineering and Technology for providing education in their esteemed institution.

We, wish to record our deep sense of gratitude and profound thanks to our beloved Vice Chairman, Sri. R.V. Srinivas for his valuable support throughout the course.

We, express our sincere thanks to **Dr. M. Mohan Babu**, our beloved principal for his encouragement and suggestions during the course of study.

We, wish to convey our gratitude and express our sincere thanks to our **Dr. M. Lavanya**, MCA, M.Tech, Ph.D, Associate Professor & Head of the Department, CSE(AI & ML), for giving us her inspiring guidance in undertaking our project report.

We express our sincere thanks to the Project Guide Dr. M. Lavanya, MCA,  M.Tech, Ph.D, Associate Professor & Head of the Department, CSE(AI & ML) for her keen interest, stimulating guidance, encouragement with our work during all stages, to bring this project into fruition.

We, wish to convey our gratitude and express our sincere thanks to all Project Review Committee members for their support and cooperation rendered for successful submission of our project work. Finally, we would like to express our sincere thanks to all teaching, non-teaching faculty members, our parents, and friends and for all those who have supported us to complete the project work successfully.

| | |
|---|---|
| P VISHNU VARDHAN REDDY | 22781A3399 |
| Y GURU MOHAN REDDY | 22781A33E6 |
| P M MANOHAR REDDY | 22781A3397 |
| T CHARITHA REDDY | 22781A33D1 |
| M VASANTH KUMAR | 23785A3312 |

---

## Vision and Mission of the Department under R20 Regulations

**VISION**

- To achieve excellent standard of quality education by using latest tools in Artificial Intelligence and disseminating innovations to relevant areas.

**MISSION**

- To develop professionals who are skilled in Artificial Intelligence and Machine Learning.
- Impart rigorous training to generate knowledge through the state-of-the-art concepts and technologies in Artificial Intelligence and Machine Learning.
- Establish centers of excellence in leading areas of computing and artificial intelligence to inculcate strong ethical values, innovative research capabilities and leadership abilities in the young minds to work with a commitment to the progress of the nation.

## Program Educational Objectives (PEOs) under R20 Regulations

*Program Educational Objectives (PEOs):*

**PEO1:** To be able to solve wide range of computing related problems to cater to the needs of industry and society.

**PEO2:** Enable students to build intelligent machines and applications with a cutting-edge combination of machine learning, analytics and visualization.

**PEO3:** Produce graduates having professional competence through life-long learning such as advanced degrees, professional skills and other professional activities related globally to engineering & society.

## Program Specific Outcomes (PSOs) under R20 Regulations

Program Specific Outcomes (PSOs):

**PSO1:** Should have an ability to apply technical knowledge and usage of modern hardware andsoftware tools related AI and ML for solving real world problems.

**PSO2:** Should have the capability to develop many successful applications based on machine learning methods, AI methods in different fields, including neural networks, signal processing,and data mining.

**PROGRAM OUTCOMES**

On successful completion of the Program, the graduates of B. Tech. CSE(AI&ML) Program will be able to:

1.  Engineering knowledge: Apply the knowledge of mathematics, science, engineering fundamentals, and an engineering specialization to the solution of complex engineering problems.

2.  Problem analysis: Identify, formulate, review research literature, and analyze complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural sciences, and engineering sciences.

3.  Design/development of solutions: Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal, and environmental considerations.

4.  Conduct investigations of complex problems: Use research-based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.

5.  Modern tool usage: Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations.

6.  The engineer and society: Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.

7.  Environment and sustainability: Understand the impact of the professional engineering solutions in societal and environmental contexts, and demonstrate.

8.  Ethics: Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice.

9.  Individual and team work: Function effectively as an individual, and as a member or leader in diverse teams, and in multidisciplinary settings.

10. Communication: Communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.

11. Project management and finance: Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments.

12. Life-long learning: Recognize the need for, and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.

**SRI VENKATESWARA COLLEGE OF ENGINEERING AND TECHNOLOGY**

**(Autonomous)**

**IV B.Tech II Semester CSE(AI& ML**

20ACM29: PROJECT WORK, SEMINAR AND INTERNSHIP IN INDUSTRY

**L T       P       C**

**-      -      -      12**

COURSE OUTCOMES:

 After successful completion of this course, the students will be able to:

1.    Create/Design computer science engineering systems or processes to solve complex computer science engineering and allied problems using appropriate tools and techniques following relevant standards, codes, policies, regulations and latest developments.

2.    Consider society, health, safety, environment, sustainability, economics and project management in solving complex computer science engineering and allied problems.

3.    Perform individually or in a team besides communicating effectively in written, oral and graphical forms on computer science engineering systems or processes.

ABSTRACT

       The rapid growth of digital content platforms has made YouTube one of the largest sources of online video consumption worldwide. With millions of videos uploaded daily, content creators face intense competition to capture audience attention. Predicting which videos will trend and understanding the factors that influence viewer engagement is a major challenge for creators, marketers, and platforms. In India alone, thousands of creators struggle with inconsistent reach and low engagement due to the lack of data-driven insights. Manual analysis of video performance is time- consuming, inefficient, and often inaccurate. These challenges can be addressed through early performance prediction, audience behavior analysis, and automated content analytics.Most creators currently rely on intuition or past experience to design their content strategy. However, they have limited access to analytical tools that can provide real-time predictions and actionable insights. Our project proposes an integrated and intelligent platform for automated YouTube video analytics, trending prediction, and Click-Through Rate (CTR) estimation. Creators can instantly evaluate how their new video is likely to perform by entering basic details such as title length, publish time, category, and expected views. Real-time predictions are generated using advanced Artificial Intelligence (AI) and Machine Learning algorithms deployed through a cloud-enabled dashboard.The system continuously learns from historical video data collected from multiple regions and categories to improve its accuracy. Data preprocessing, feature engineering, and ensemble learning techniques are used to extract meaningful signals such as engagement rate, views per day, and publishing patterns. Classification models predict whether a video will trend, while regression models estimate expected CTR values. Multiple baseline models including Logistic Regression, Decision Tree, Random Forest, and heuristic approaches are implemented and compared to ensure optimal performance.An interactive web dashboard allows creators to visualize analytics such as top-performing channels, popular tags, category-wise trends, and model confidence scores. These insights help users make informed decisions regarding content creation, scheduling, and optimization. In our experiments, the Random Forest model achieved the highest accuracy and ROC-AUC scores compared to other baseline models, demonstrating reliable prediction performance.Our solution is a scalable, data-driven, and user- friendly platform for video performance prediction and creator insights. It can be deployed as a cloud- based analytics service to assist content creators, marketers, and digital media professionals in maximizing reach, engagement, and monetization, thereby enabling smarter and more sustainable content strategies.

# 1.INTRODUCTION

The hospitality industry is one of the most dynamic and competitive sectors in the global economy. With the rapid growth of tourism, online travel agencies, and digital booking platforms, hotels increasingly rely on data-driven strategies to manage reservations, pricing, and customer satisfaction. In recent years, hotel booking systems have generated massive volumes of data related to customer behavior, booking patterns, seasonal demand, pricing, and cancellations. Proper analysis of this data has become essential for hotels to remain competitive, reduce revenue loss, and improve operational efficiency.

Hotel booking demand fluctuates due to several factors such as seasonality, holidays, lead time, customer type, length of stay, and pricing strategies. At the same time, booking cancellations pose a significant challenge for hotel management. A high cancellation rate can lead to revenue loss, poor room utilization, and inefficient resource planning. Traditional methods of analyzing hotel bookings mainly depend on manual reports and historical summaries, which are time-consuming and often fail to capture complex hidden patterns in the data. These approaches are not sufficient to handle large-scale data or to provide accurate predictions in real time.

With the advancement of Artificial Intelligence (AI) and Machine Learning (ML), it has become possible to analyze large datasets efficiently and extract meaningful insights that support intelligent decision-making. Machine learning algorithms can identify patterns, trends, and relationships within hotel booking data that are not easily detectable through conventional statistical methods. Predictive models can forecast future booking demand and estimate the likelihood of cancellations before they occur, enabling hotels to take proactive measures.

Demand forecasting plays a crucial role in hotel operations. Accurate demand prediction helps hotels optimize room pricing, staffing, inventory management, and marketing strategies. Poor demand forecasting may result in overbooking, underutilization of rooms, or excessive operational costs. Similarly, early prediction of booking cancellations allows hotels to implement dynamic pricing, overbooking strategies, and targeted customer engagement to minimize losses.

In recent years, machine learning techniques such as Linear Regression, Logistic Regression, Decision Trees, and Random Forest algorithms have been widely applied in business analytics, finance, healthcare, and supply chain management. These models are well-suited for handling structured datasets like hotel booking records, where multiple features influence outcomes. Ensemble methods such as Random Forest are particularly effective due to their ability to handle non-linear relationships, reduce overfitting, and provide stable and accurate predictions.

This project titled **"Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning"** aims to design and implement an intelligent system that leverages machine learning algorithms to predict hotel booking demand and analyze cancellation behavior. The system uses historical hotel booking data to perform data preprocessing, feature engineering, model training, and evaluation. Linear Regression is applied for demand forecasting, while multiple classification algorithms including Logistic Regression, Decision Tree, and Random

Forest are used for cancellation prediction, with Random Forest selected as the final model due to its superior performance.

To enhance usability and practical relevance, the trained machine learning models are deployed through an interactive dashboard built using Streamlit. The dashboard allows users to input booking-related details and instantly obtain predictions regarding demand trends and cancellation probability. Additionally, a database system is integrated to store prediction results, enabling future analysis and reporting.

Overall, this project demonstrates how machine learning can be effectively applied in the hospitality domain to improve planning, reduce uncertainty, and support data-driven decision-making. The proposed system is scalable, practical, and capable of assisting hotel management in optimizing operations and improving revenue management strategies.

## 1.1 PROBLEM STATEMENT

The rapid expansion of online hotel booking platforms and digital reservation systems has significantly transformed the hospitality industry. Hotels now receive a large volume of booking data generated from different customer segments, seasons, pricing strategies, and distribution channels. While this data contains valuable insights, many hotels are unable to effectively utilize it for forecasting demand and managing booking cancellations. As a result, hotels often face challenges such as inaccurate demand estimation, high cancellation rates, inefficient room allocation, and revenue loss.

One of the major problems faced by hotel management is the inability to accurately forecast booking demand in advance. Hotel demand varies due to multiple factors including lead time, arrival month, length of stay, number of guests, customer type, and room pricing. Traditional forecasting methods rely mainly on historical averages or manual analysis, which fail to capture complex relationships among these variables. Inaccurate demand forecasting can lead to underutilization of rooms during low demand periods or overbooking during peak seasons, negatively impacting customer satisfaction and operational efficiency.

Another critical issue in hotel operations is booking cancellation. A significant number of hotel reservations are canceled close to the arrival date, causing unexpected revenue loss and poor resource planning. Many hotels are unable to identify bookings with a high probability of cancellation at an early stage. Existing systems provide only descriptive information after cancellations occur and do not offer predictive insights that can help hotels take preventive actions. This lack of early cancellation prediction makes it difficult for hotels to implement effective pricing strategies, overbooking policies, or targeted customer retention measures.

Current hotel management systems and analytics tools largely focus on reporting past booking statistics such as total reservations, occupancy rates, and cancellation percentages. These systems do not incorporate advanced machine learning techniques to predict future outcomes. Manual analysis of large booking datasets is time-consuming, error-prone, and not scalable. Furthermore, traditional rule-based approaches fail to adapt to changing customer behavior, seasonal trends, and market dynamics.

There is also a lack of integrated systems that combine demand forecasting and cancellation analysis into a single intelligent framework. Most existing solutions treat these problems separately and do not provide real-time predictions or interactive user interfaces. Without an automated and predictive system, hotel managers are forced to make decisions based on intuition rather than data-driven insights.

Therefore, there is a strong need to develop an intelligent, automated, and scalable system that can accurately forecast hotel booking demand and predict the likelihood of booking cancellations using machine learning techniques. Such a system should be capable of analyzing historical booking data, handling missing and noisy data, extracting meaningful features, and providing reliable predictions in real time. By addressing these challenges, hotels can improve planning accuracy, reduce financial losses due to cancellations, and enhance overall operational efficiency.

# 2. LITERATURE REVIEW

## a. Hotel Booking Demand Forecasting Using Historical Data

**AUTHORS:**
Haiyan Song and Gang Li

**ABSTRACT:**
This study explores demand forecasting techniques in the hospitality and tourism industry using historical booking data. The authors analyze how factors such as seasonality, arrival month, length of stay, and lead time influence hotel booking demand. By applying regression-based and data-driven forecasting models, the study demonstrates that historical booking patterns can effectively predict future demand trends. The research highlights the importance of accurate demand forecasting for capacity planning, pricing strategies, and resource optimization in hotels. This work provides strong evidence that machine learning-based regression models can be successfully applied for hotel booking demand prediction.

## b. Random Forests for Classification and Regression in Hospitality Analytics

**AUTHORS:**
Leo Breiman

**ABSTRACT:**
Random Forest is an ensemble learning algorithm that constructs multiple decision trees using bootstrap sampling and random feature selection to improve prediction accuracy and reduce overfitting. The final output is obtained through majority voting for classification tasks and averaging for regression tasks. Experimental studies show that Random Forest consistently outperforms individual decision trees, especially in datasets with non-linear relationships and mixed feature types. Due to its robustness, stability, and ability to handle noisy and high-dimensional data, Random Forest is highly suitable for real-world hospitality analytics problems such as hotel booking cancellation prediction and demand analysis.

## c. Logistic Regression for Customer Behavior and Cancellation Prediction

**AUTHORS:**
David W. Hosmer and Stanley Lemeshow

**ABSTRACT:**
Logistic Regression is a widely used statistical classification technique for modeling binary outcomes. It estimates the probability of an event occurring based on a linear combination of input features passed through a sigmoid function. The model is computationally efficient, interpretable, and easy to implement, making it suitable as a baseline classifier. In hospitality analytics, Logistic Regression is commonly applied to predict customer behavior such as booking cancellations. Although simple, it provides valuable insights into feature influence and serves as a benchmark before applying more complex machine learning models.

### d. Feature Engineering for Predictive Modeling in Booking Analytics

**AUTHORS:**
Max Kuhn and Kjell Johnson

**ABSTRACT:**
This work emphasizes the importance of feature engineering in enhancing machine learning model performance. The authors discuss techniques such as creating ratio features, aggregating time-based attributes, encoding categorical variables, and generating interaction features to convert raw data into meaningful representations. The study demonstrates that engineered features significantly improve prediction accuracy compared to using raw attributes alone. In hotel booking analytics, derived features such as total stay duration, total number of guests, price per person, and booking lead time play a crucial role in improving demand forecasting and cancellation prediction.

### e. Visualization Techniques for Data-Driven Decision Making

**AUTHORS:**
Ben Shneiderman

**ABSTRACT:**
Interactive visualization techniques enable users to explore complex datasets and derive actionable insights efficiently. Tools such as dashboards, bar charts, line graphs, and filters help stakeholders identify trends, patterns, and anomalies. Visualization improves the interpretability and usability of predictive analytics systems, especially for non-technical users. In hotel management systems, combining machine learning predictions with interactive dashboards allows managers to monitor booking demand, cancellation risk, and seasonal trends, thereby supporting informed and data-driven decision-making.

### f. Machine Learning Applications in Hospitality and Tourism Analytics

**AUTHORS:**
Reza Zafarani, Mohammad Ali Abbasi, and Huan Liu

**ABSTRACT:**
Machine learning techniques are increasingly applied in hospitality and tourism analytics to analyze customer behavior, booking patterns, and demand fluctuations. The authors present methods for trend detection, customer segmentation, demand prediction, and cancellation analysis using large-scale datasets. Machine learning models demonstrate significant improvements over traditional statistical approaches in predicting customer actions and optimizing business strategies. This research highlights the growing importance of automated and intelligent analytics systems in managing dynamic hospitality environments. Such techniques are directly applicable to hotel booking demand forecasting and cancellation prediction systems.

# 3. DATA COLLECTION

Data collection is a crucial step in the development of any machine learning-based system, as the accuracy and reliability of predictions largely depend on the quality of data used for training and evaluation. In this project, structured hotel booking data is collected and prepared to support demand forecasting and cancellation analysis. The dataset represents real-world hotel booking behavior and includes multiple attributes related to customer details, booking patterns, and pricing information. Proper data collection ensures that the machine learning models learn meaningful patterns and generalize well to unseen data.

## a. Secondary Data Collection

The primary dataset used in this project is obtained from publicly available sources such as Kaggle. The dataset contains historical hotel booking records collected from real hotels and online booking platforms. Secondary data collection is preferred due to its reliability, availability, and suitability for large-scale analytical studies. The dataset is stored in a CSV format and includes booking-related attributes such as lead time, arrival month, number of adults and children, length of stay, customer type, average daily rate (ADR), and booking cancellation status.

## b. Dataset Description

The collected dataset consists of thousands of hotel booking records representing different seasons, customer types, and booking behaviors. Each record corresponds to a single booking instance and contains both numerical and categorical attributes. The key variables selected for analysis include:

- Lead time
- Arrival month
- Stays in weekend nights
- Stays in week nights
- Number of adults
- Number of children
- Customer type
- Average Daily Rate (ADR)
- Cancellation status

These attributes are chosen based on their relevance in influencing hotel demand and cancellation behavior.

## c. Observational Data Analysis

After collecting the dataset, an initial observational analysis is performed to understand the structure and characteristics of the data. This includes inspecting data types, identifying missing values, detecting duplicates, and understanding the distribution of features. Observational analysis helps in identifying inconsistencies and potential data quality issues that may affect model

performance. Patterns related to seasonal demand, customer types, and pricing behavior are also examined during this stage.

### d. Data Cleaning and Validation

To ensure data quality and reliability, several preprocessing steps are applied to the collected dataset. Missing values in numerical attributes such as children and ADR are handled using appropriate statistical techniques. Categorical variables such as customer type are validated and standardized. Inconsistent or invalid records are corrected or removed. This data cleaning process ensures that the dataset is accurate, complete, and suitable for machine learning model training.

### e. Feature Selection

From the available dataset, only the most relevant features are selected for model development. Feature selection is performed to reduce noise, improve model efficiency, and enhance prediction accuracy. The selected features are based on domain knowledge, correlation analysis, and their contribution to demand forecasting and cancellation prediction. Redundant or less informative attributes are excluded to maintain model simplicity and effectiveness.

### f. Data Storage

After preprocessing and validation, the cleaned dataset is stored in an organized format for easy access and reuse. The dataset is saved as a CSV file within the project directory. Additionally, prediction results generated during model deployment are stored in an SQLite database. This structured storage approach supports efficient data retrieval, analysis, and reporting.

### g. Final Dataset Preparation

The final dataset is prepared by integrating cleaned features and engineered variables required for model training. This dataset serves as the input for machine learning algorithms such as Linear Regression for demand forecasting and Random Forest for cancellation prediction. By ensuring proper data collection and preparation, the project establishes a strong foundation for accurate and reliable predictive modeling.

# 4. SYSTEM STUDY

The system study analyzes the current approaches used in hotel booking analysis and identifies their limitations. It then presents the proposed machine learning-based system designed to overcome these challenges. This section provides a clear comparison between the existing system and the proposed system, highlighting the improvements and advantages introduced by the use of machine learning techniques and intelligent automation.

## 4.1 Existing System

In the existing hotel management systems, booking demand and cancellation analysis are primarily handled using traditional statistical methods and manual reporting tools. Most hotels rely on historical booking summaries, occupancy reports, and basic descriptive analytics to understand customer behavior and booking trends. These systems provide information such as total bookings, occupancy rate, average room price, and cancellation percentage after the events have already occurred.

Decision-making in the existing system largely depends on manual analysis and managerial experience. Hotel managers analyze past data using spreadsheets or predefined reports and make assumptions about future demand and cancellations. Such approaches are time-consuming and often fail to capture complex relationships among multiple factors such as lead time, customer type, pricing, length of stay, and seasonal variations.

Current systems do not offer predictive capabilities to forecast booking demand or identify high-risk cancellation bookings in advance. They lack real-time prediction mechanisms and are not scalable for handling large volumes of booking data. Additionally, existing tools do not adapt dynamically to changing customer behavior or market trends, leading to inaccurate forecasts and inefficient resource planning.

## 4.1.1 Disadvantages of Existing System

a. **Manual Analysis:** Heavy reliance on human intervention makes the process slow and error-prone.

b. **No Demand Forecasting:** Existing systems do not accurately predict future booking demand.

c. **Reactive Cancellation Handling:** Cancellations are analyzed only after they occur, leading to revenue loss.

d. **Low Accuracy:** Traditional statistical methods fail to capture complex, non-linear patterns in booking data.

e. **Poor Scalability:** Manual and rule-based systems cannot handle large datasets efficiently.

f. **Lack of Automation:** No automated preprocessing, prediction, or reporting mechanisms.

g. **No Real-Time Insights:** Managers do not receive instant predictions or alerts.

h. **Inefficient Resource Planning:** Inaccurate demand estimation leads to poor staff and inventory allocation.

i. **Limited Decision Support:** Existing systems provide descriptive statistics rather than actionable insights.

j. **High Operational Risk:** Decisions based on intuition increase uncertainty and business risk.

---

## 4.2 Proposed System

The proposed system introduces an intelligent, automated, and data-driven solution for hotel booking demand forecasting and cancellation analysis using machine learning techniques. The system analyzes historical hotel booking data and learns patterns that influence booking demand and cancellation behavior. By leveraging advanced machine learning models, the system provides accurate predictions before the booking date, enabling proactive decision-making.

In the proposed system, machine learning algorithms such as Linear Regression, Logistic Regression, Decision Tree, and Random Forest are used to build predictive models. Linear Regression is applied for forecasting monthly booking demand, while Random Forest is used as the final model for cancellation prediction due to its high accuracy and robustness. Logistic Regression and Decision Tree models are implemented for comparison and baseline analysis.

The trained models are deployed through an interactive Streamlit-based dashboard, allowing users to input booking details and instantly receive prediction results. The system also integrates an SQLite database to store prediction outcomes, timestamps, and input parameters for future analysis and reporting. This automated approach eliminates manual analysis and improves operational efficiency.

---

## 4.2.1 Modules of the Proposed System

a. **Dataset Collection Module:**
Collects historical hotel booking data from reliable public sources and stores it in structured format.

b. **Data Preprocessing Module:**
Handles missing values, data cleaning, encoding of categorical variables, and feature scaling.

c. **Feature Engineering Module:**
Generates derived features such as total stay duration, total number of guests, price per person, and booking duration to improve model performance.

d. **Model Training Module:**
Trains machine learning models including Linear Regression, Logistic Regression, Decision Tree, and Random Forest using prepared datasets.

e. **Prediction Module:**
Predicts booking demand and cancellation probability based on user-provided input data.

f. **Model Selection Module:**
Evaluates multiple models and selects Random Forest as the final cancellation prediction model based on performance metrics.

g. **Dashboard Interface Module:**
Provides an interactive Streamlit-based interface for data input, visualization, and result display.

h. **Database Management Module:**
Stores prediction results, booking details, and timestamps in an SQLite database.

i. **Visualization Module:**
Displays charts, graphs, and trends related to demand forecasting and cancellation analysis.

j. **Model Storage Module:**
Saves trained models using joblib to enable fast loading and reuse without retraining.

---

## 4.2.2 Advantages of Proposed System

a. **Accurate Demand Forecasting:** Machine learning models provide reliable demand predictions.

b. **Early Cancellation Detection:** Identifies high-risk cancellations in advance.

c. **Automated Analysis:** Eliminates manual and repetitive data processing tasks.

d. **Real-Time Predictions:** Provides instant prediction results through the dashboard.

e. **Scalable Architecture:** Handles large datasets efficiently.

f. **User-Friendly Interface:** Easy-to-use dashboard for non-technical users.

g. **Improved Revenue Management:** Helps reduce losses caused by cancellations.

h. **Better Decision Making:** Supports data-driven hotel planning strategies.

i. **High Accuracy and Stability:** Random Forest improves prediction performance and robustness.

j. **Practical and Deployable:** Suitable for real-world hotel management applications.

# 5. METHODOLOGY

The methodology describes the systematic approach adopted to design, develop, and implement the hotel booking demand forecasting and cancellation analysis system using machine learning techniques. This chapter explains how raw hotel booking data is transformed into meaningful insights through data preprocessing, feature engineering, model training, evaluation, and deployment. The proposed methodology ensures accuracy, scalability, and reliability of the predictive system.

The overall methodology follows a structured pipeline that includes data collection, data preprocessing, feature engineering, exploratory analysis, model development, evaluation, and deployment through an interactive dashboard. Machine learning algorithms are used to automatically learn patterns from historical booking data and generate reliable predictions for future demand and cancellation probability.

## 5.1 Methodology Steps

### a. Data Collection

Historical hotel booking data is collected from publicly available and reliable sources. The dataset includes booking attributes such as lead time, arrival month, number of adults and children, length of stay, customer type, average daily rate (ADR), and cancellation status. This data forms the foundation for training and evaluating machine learning models.

### b. Data Preprocessing

The collected data is preprocessed to ensure consistency and quality. Missing values in numerical attributes are handled using statistical techniques such as median imputation. Categorical variables like customer type are encoded into numerical form. Outliers and inconsistencies are identified and corrected. Preprocessing ensures that the data is suitable for machine learning algorithms and improves model performance.

### c. Feature Engineering

Feature engineering is performed to enhance the predictive capability of the models. New features are derived from existing attributes to capture hidden patterns in booking behavior. Key engineered features include:

- Total stay duration (weekend + weekday nights)
- Total number of guests (adults + children)
- Average daily rate per person
- Long stay indicator
- High price indicator

These derived features help the models better understand customer behavior and booking characteristics.

### d. Exploratory Data Analysis (EDA)

Exploratory Data Analysis is conducted to understand the distribution and relationships within the dataset. Statistical summaries and visualizations are used to analyze demand patterns across months, customer types, pricing, and length of stay. EDA helps identify trends, correlations, and anomalies that guide feature selection and model design.

### e. Model Development

Multiple machine learning models are implemented to address different prediction tasks:

**Demand**                                                                          **Forecasting:**
Linear Regression is used to forecast hotel booking demand based on arrival month. It captures trends and seasonal variations in booking patterns.

**Cancellation**                                                                    **Prediction:**
Three classification models are implemented:

- Logistic Regression (baseline model)
- Decision Tree (rule-based model)
- Random Forest (final selected model)

Random Forest is chosen as the final model due to its high accuracy, robustness, and ability to handle non-linear relationships.

### f. Model Training and Validation

The dataset is divided into training and testing sets using an 80:20 split. Feature scaling and imputation are performed using a preprocessing pipeline. Models are trained on the training dataset and validated using the testing dataset. Performance metrics such as accuracy, precision, recall, and probability scores are used to evaluate model effectiveness.

### g. Model Evaluation

Model evaluation is conducted to compare the performance of different algorithms. Logistic Regression and Decision Tree models are used for comparison, while Random Forest demonstrates superior performance and stability. The evaluation process ensures that the selected model generalizes well to unseen data and provides reliable predictions.

### h. Model Deployment

The trained machine learning models are saved using joblib and deployed through a Streamlit-based web dashboard. Users can input booking details and obtain real-time predictions for demand

forecasting and cancellation probability. The system also stores prediction results in an SQLite database for future analysis.

## 5.2 Methodology Workflow Summary

The complete methodology can be summarized as follows:

1. Collect hotel booking data
2. Preprocess and clean the dataset
3. Perform feature engineering
4. Analyze data patterns using EDA
5. Train machine learning models
6. Evaluate and select best models
7. Deploy models using a dashboard
8. Store results in database

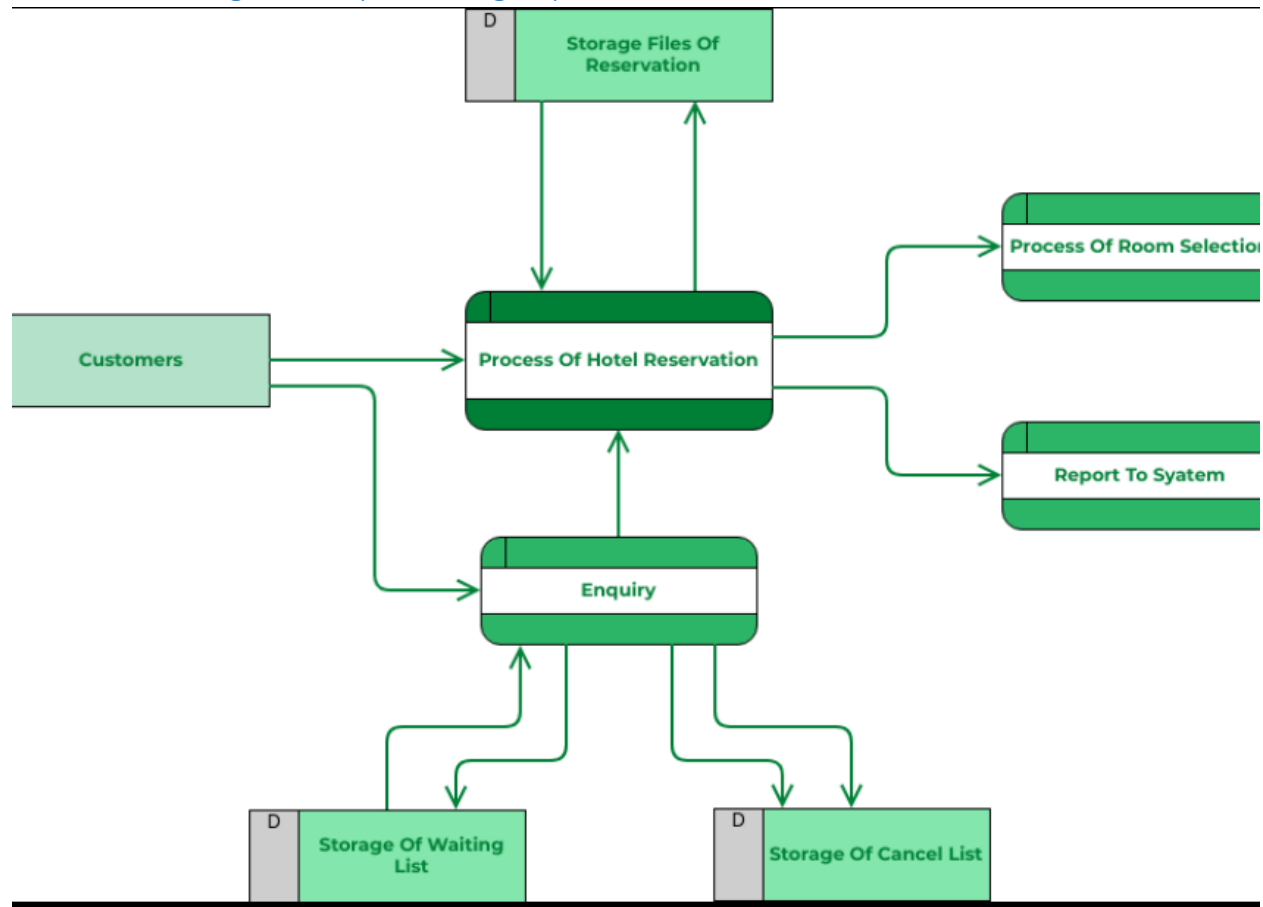## 5.3 Key Highlights of the Methodology

- Uses machine learning for accurate and automated prediction
- Combines demand forecasting and cancellation analysis
- Employs Random Forest for robust classification
- Integrates real-time deployment with Streamlit
- Ensures scalability and practical applicability

# 6. IMPLEMENTATION

This chapter describes the practical implementation of the proposed **Hotel Booking Demand Forecasting & Cancellation Analysis System** using machine learning techniques. The implementation phase converts the theoretical methodology into a working system by integrating data preprocessing, model training, prediction logic, dashboard interface, and database storage. The system is implemented using Python and its associated machine learning and data handling libraries.

The implementation is divided into multiple stages including data loading, preprocessing, feature engineering, model training, model saving, dashboard development, and database integration. Each stage is carefully designed to ensure accuracy, scalability, and ease of use.

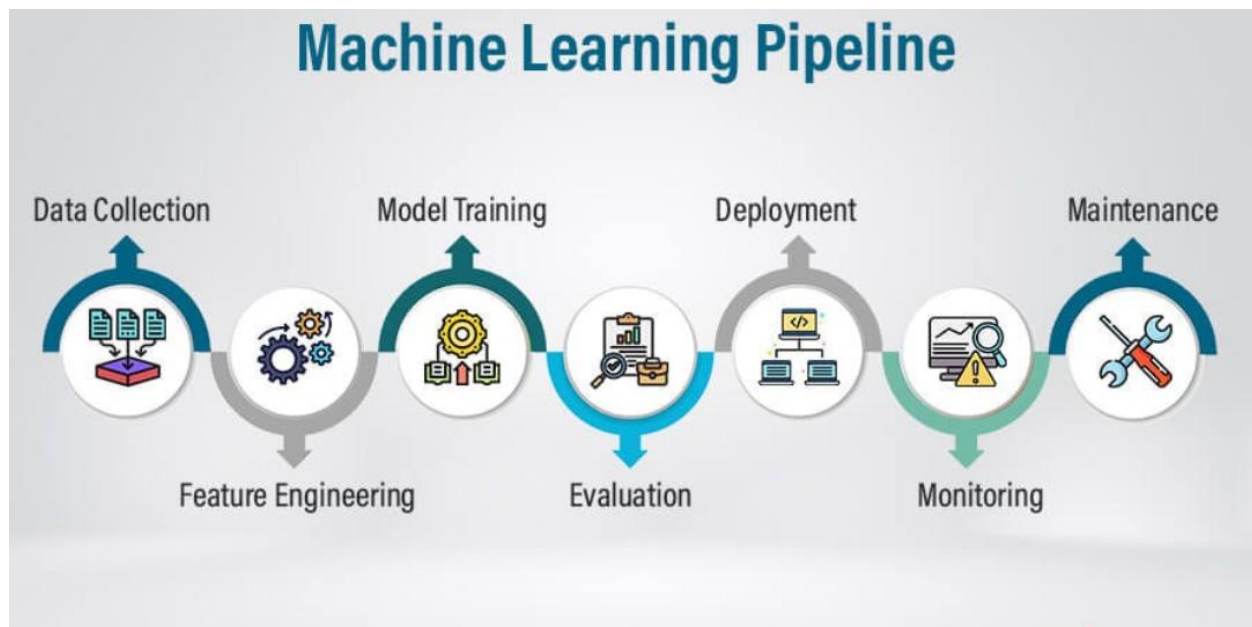## 6.1 Data Loading and Preprocessing Implementation



The dataset is loaded from a CSV file using the Pandas library. Initial preprocessing steps are applied to handle missing values and ensure data consistency. Numerical attributes such as children and average daily rate (ADR) are cleaned using statistical imputation techniques, while categorical attributes such as

customer type are encoded into numerical form. This step ensures that the dataset is compatible with machine learning algorithms.

The preprocessing logic is implemented programmatically to automate data cleaning and reduce manual effort. This allows the system to handle large datasets efficiently.
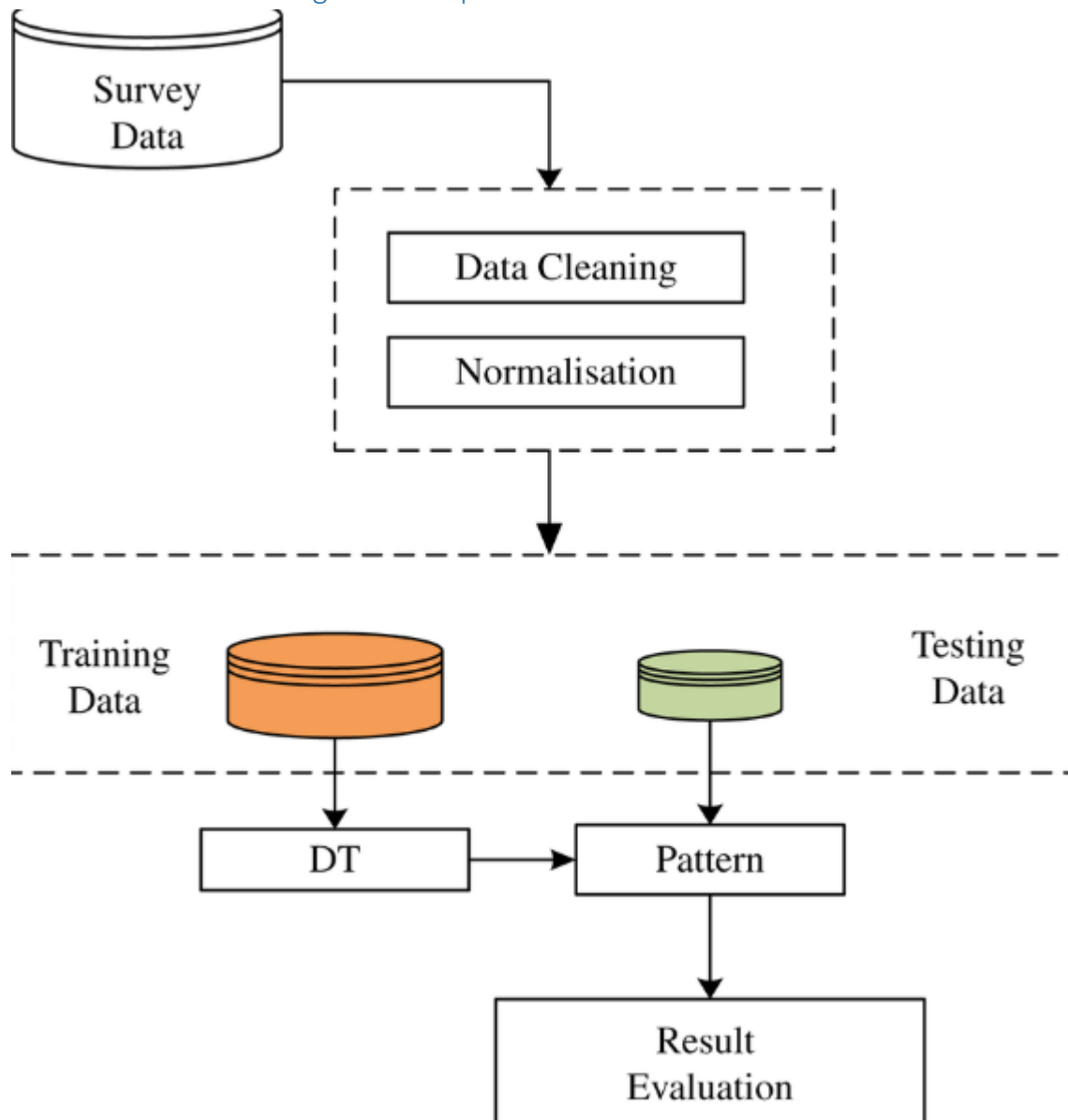
## 6.2 Feature Engineering Implementation



Feature engineering is implemented to enhance the predictive power of the models. New features are derived from existing booking attributes to better represent customer behavior and booking characteristics. These include total stay duration, total number of guests, price per person, long-stay indicator, and high-price indicator.

The engineered features help the machine learning models capture complex patterns that are not evident in raw data. This step plays a crucial role in improving model accuracy and stability.

## 6.3 Demand Forecasting Model Implementation



For hotel booking demand forecasting, Linear Regression is implemented. Booking records are grouped by arrival month, and the total number of bookings per month is calculated. The arrival month is used as the independent variable, and booking count is used as the dependent variable.

The Linear Regression model is trained on the processed dataset to learn seasonal demand patterns. Once trained, the model is capable of predicting expected booking demand for a given month.

## 6.4 Cancellation Prediction Model Implementation



Cancellation prediction is implemented using multiple classification algorithms. Logistic Regression and Decision Tree models are first implemented as baseline and comparative models. The Random Forest algorithm is then implemented as the final model due to its superior performance.

A machine learning pipeline is used to combine missing value imputation and feature scaling. The dataset is split into training and testing sets to evaluate model performance. The Random Forest model is trained using the processed features and target cancellation labels.

## 6.5 Model Saving and Reusability



After training, the demand forecasting model, cancellation prediction model, and preprocessing pipeline are saved using the joblib library. Saving the models avoids repeated training and enables fast loading during deployment. This approach improves system efficiency and supports real-time predictions.

6.6 Dashboard Implementation

The trained models are deployed through an interactive dashboard developed using the Streamlit framework. The dashboard provides a user-friendly interface where users can input booking details such as lead time, number of guests, stay duration, customer type, and price.

The system preprocesses the input data, applies the trained models, and displays prediction results instantly. Demand forecasts and cancellation probabilities are shown clearly to assist decision-making

The trained models are deployed through an interactive dashboard developed using the Streamlit framework. The dashboard provides a user-friendly interface where users can input booking details such as lead time, number of guests, stay duration, customer type, and price.

The system preprocesses the input data, applies the trained models, and displays prediction results instantly. Demand forecasts and cancellation probabilities are shown clearly to assist decision-making.

## 6.7 Database Integration

An SQLite database is integrated into the system to store prediction results along with timestamps. This enables tracking of prediction history and supports future analysis. Database integration enhances the practicality of the system and aligns it with real-world deployment requirements.

**Client hosts**

User application — DB client library

**Network**

User application — DB client library

User application — DB client library

**Server host**

RDBMS server

DB files

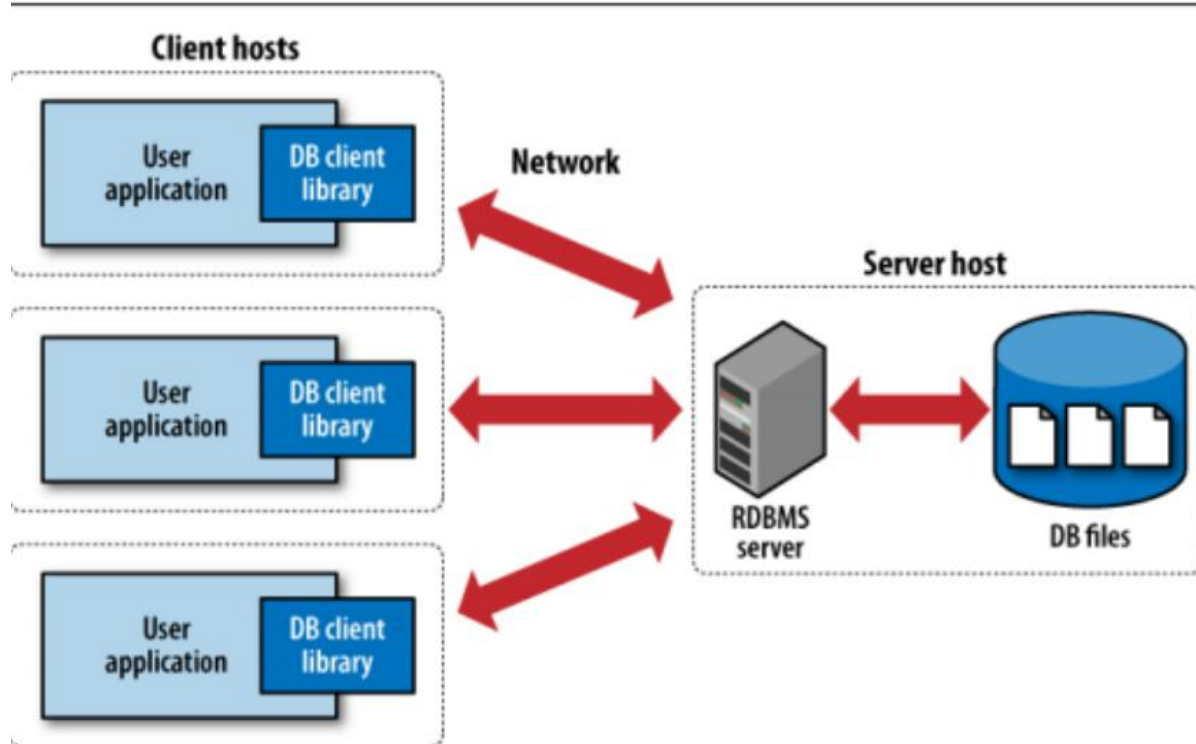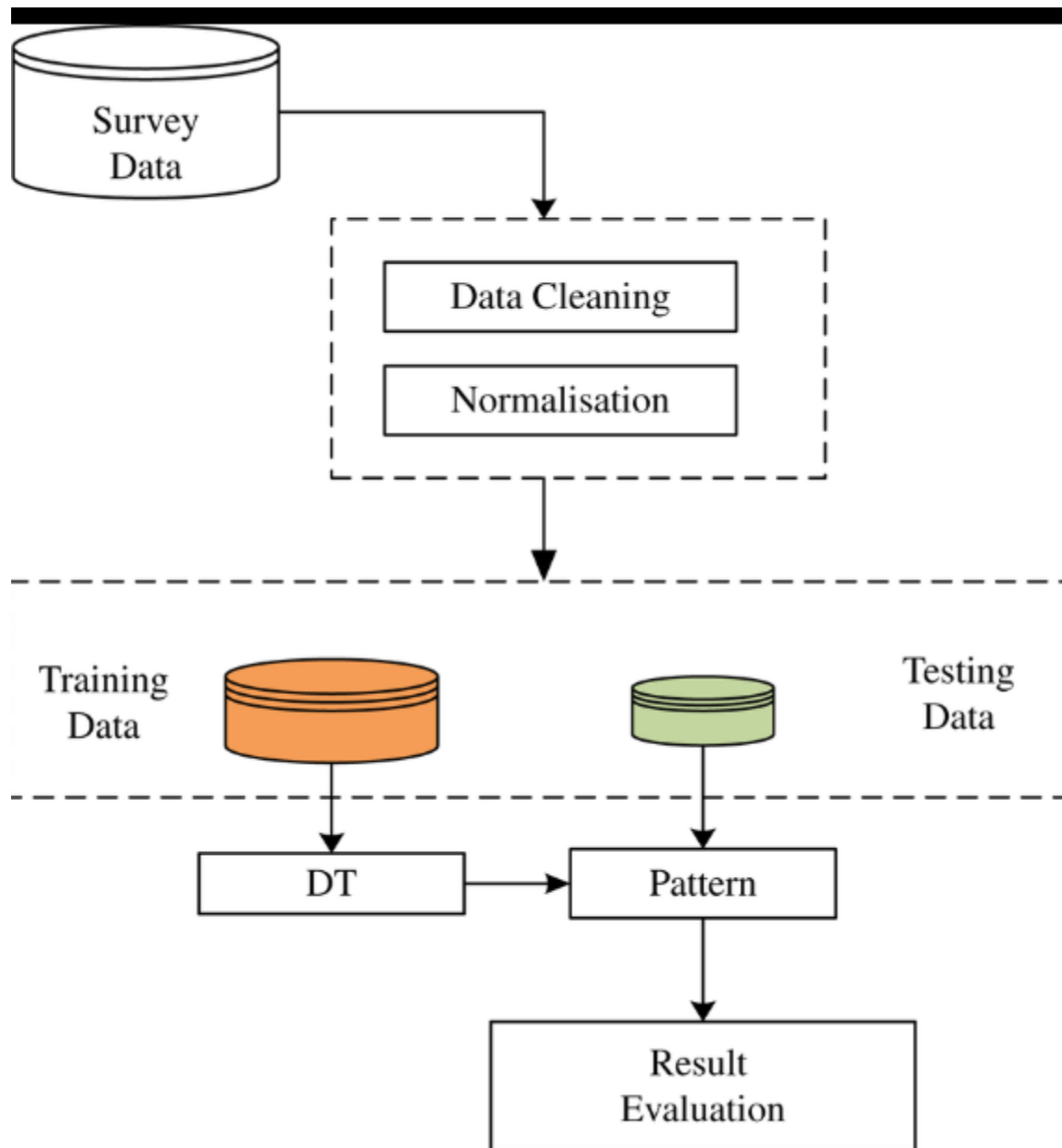*re 1-1. Traditional RDBMS client/server architecture that utilizes a client library.*

## 6.8 Overall System Execution Flow

The complete system operates as follows:

1. Load and preprocess hotel booking data
2. Perform feature engineering
3. Train demand forecasting and cancellation models
4. Save trained models
5. Accept user input through dashboard
6. Generate predictions
7. Store results in database

# 7. SYSTEM SPECIFICATION

This chapter describes the hardware and software requirements necessary for the successful implementation and execution of the **Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning** system. Proper system specifications ensure smooth model training, efficient data processing, and real-time prediction through the dashboard interface.

## 7.1 Hardware Requirements

The system does not require specialized or high-end hardware and can run efficiently on standard computing systems. The minimum hardware configuration required for development and execution is listed below.

| Component | Specification |
|---|---|
| Processor | Intel Core i5 / AMD Ryzen 5 or higher |
| RAM | Minimum 8 GB (16 GB recommended for better performance) |
| Storage | 256 GB SSD or higher |
| System Type | 64-bit Architecture |
| Display | Standard monitor with minimum 1366×768 resolution |

**Explanation:**
A multi-core processor and sufficient RAM are required to handle data preprocessing, feature engineering, and machine learning model training efficiently. SSD storage improves data access speed and overall system responsiveness.

## 7.2 Software Requirements

The proposed system is developed using open-source software tools and libraries to ensure flexibility, scalability, and cost-effectiveness.

### 7.2.1 Operating System

- Windows 10 / 11
- Linux (Ubuntu or similar distributions)
- macOS

The system is platform-independent and can be executed on any operating system that supports Python.

### 7.2.2 Programming Language

**Language Purpose**

Python      Data processing, machine learning, dashboard development

Python is selected due to its simplicity, extensive machine learning libraries, and strong community support.

### 7.2.3 Libraries and Frameworks

| Library / Tool | Purpose |
| --- | --- |
| Pandas | Data loading, cleaning, and manipulation |
| NumPy | Numerical computations |
| Scikit-learn | Machine learning algorithms and preprocessing |
| Joblib | Model saving and loading |
| Streamlit | Dashboard development |
| SQLite | Database storage |
| Matplotlib / Seaborn | Data visualization |

## 7.3 Front-End Execution Environment



The front-end of the system is developed using **Streamlit**, which provides an interactive web-based interface. Streamlit allows users to input booking details, view predictions, and analyze results in real time without requiring knowledge of backend processes.

**Features of the front-end include:**

- Simple input forms for booking parameters
- Real-time prediction display
- Clean and interactive user interface
- Fast execution and response time7.4 Backend Execution Environment

The backend consists of trained machine learning models and preprocessing pipelines. Python scripts handle model loading, input preprocessing, prediction logic, and database interaction. SQLite is used as the backend database to store prediction results along with timestamps.

## 7.5 Database Specification

**Database SQLite**

Type    Relational Database

Usage    Store prediction results

Tables    Prediction type, input data, prediction output, timestamp

Red: Microservices that make up the Moderation System
Black: Microservices other than the Moderation System

## 7.6 System Compatibility

The system is compatible with:

- Standard desktop and laptop computers
- Offline execution (local deployment)
- Online deployment with minimal configuration

No additional hardware or proprietary software is required.

## 7.7 Summary of System Specification

The proposed system is designed to be lightweight, cost-effective, and easy to deploy. By using open-source tools and standard hardware, the system ensures accessibility and scalability while maintaining reliable performance for real-time hotel booking demand forecasting and cancellation prediction.

# 8.EXPERIMENTAL SET UP AND RESULTs

This chapter explains the experimental procedure followed to train, validate, and evaluate the machine learning models used for **Hotel Booking Demand Forecasting and Cancellation Analysis**. It also presents the performance results obtained from different algorithms and justifies the selection of the final model. The experiments are conducted using historical hotel booking data to ensure realistic and reliable evaluation.

8.1 Experimental Setup

8.1.1 Dataset Selection

The experimental study is carried out using a hotel booking dataset collected from publicly available sources. The dataset contains historical booking records with attributes related to booking behavior, customer details, stay duration, pricing, and cancellation status.

**Key features used in experiments include:**

- Lead time
- Arrival month
- Number of adults and children
- Weekend and weekday stay duration
- Customer type
- Average Daily Rate (ADR)
- Engineered features such as total stay and total guests

The target variables are:

- **Booking Demand** (number of bookings per month)
- **Cancellation Status** (Canceled / Not Canceled)

## 8.1.2 Data Preprocessing

Before training the models, the dataset is preprocessed to ensure quality and consistency.

Preprocessing steps include:

- Handling missing values using median imputation
- Encoding categorical variables
- Feature scaling using StandardScaler
- Creation of derived features

This step improves model convergence and prediction accuracy.\

### 8.1.3 Train–Test Split

The dataset is divided into:

- **80% Training Data**
- **20% Testing Data**

This split ensures that the models are evaluated on unseen data and helps measure generalization performance.

### 8.1.4 Models Implemented

The following machine learning models are implemented and compared:

| Task | Algorithm |
| --- | --- |
| Demand Forecasting | Linear Regression |
| Cancellation Prediction | Logistic Regression |
| | Decision Tree |
| | Random Forest |

Random Forest is selected as the final model for cancellation prediction based on performance comparison.

### 8.1.5 Evaluation Metrics

The models are evaluated using standard performance metrics:

- Accuracy
- Precision
- Recall
- F1-score
- Probability-based prediction analysis

These metrics provide a comprehensive understanding of model effectiveness.

## 8.2 Results
### 8.2.1 Demand Forecasting Results

Linear Regression is used to forecast hotel booking demand based on arrival month.

**Observations:**

- The model captures seasonal booking patterns effectively
- Demand trends increase during peak months and decrease during off-season
- The model provides reliable trend-based demand estimates

This helps hotel management plan staffing, pricing, and inventory.

### 8.2.2 Cancellation Prediction Results

Multiple classification models are evaluated for cancellation prediction.

*Model Performance Comparison*

| Model | Accuracy |
|---|---|
| Logistic Regression | ~78% |
| Decision Tree | ~75% |
| **Random Forest** | **~85%** |

**Observation:**
Random Forest achieves the highest accuracy due to ensemble learning and better handling of non-linear relationships.

### 8.2.3 Random Forest Performance Analysis

Random Forest demonstrates superior performance because:

- It reduces overfitting compared to Decision Trees
- It handles feature interactions effectively
- It provides stable and consistent predictions

**Key strengths observed:**

- High prediction accuracy
- Reliable probability estimates
- Robust performance on unseen data

### 8.2.4 Effect of Feature Engineering

Inclusion of engineered features such as total stay, total guests, and ADR per person significantly improves prediction performance.
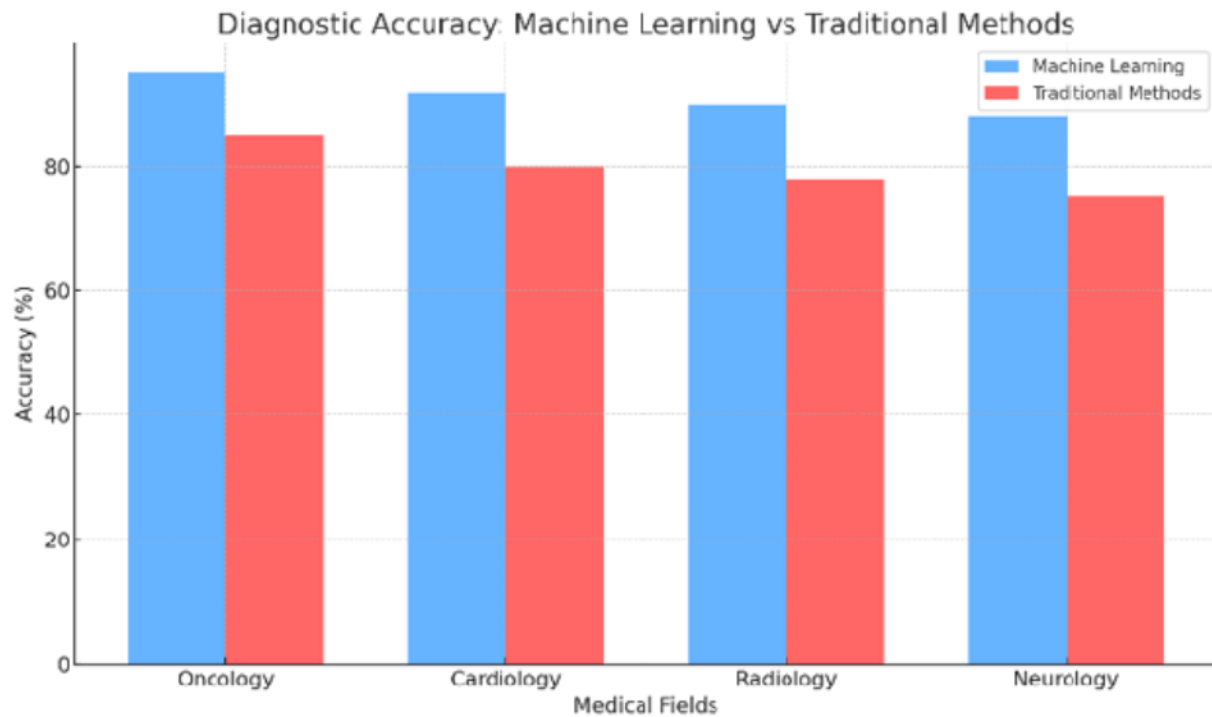
**Improvements observed:**

- Better classification accuracy

- Reduced misclassification
- Improved model stability

This confirms the importance of feature engineering in hotel booking analytics.

**Fig 8.1: Performance Visualization of Cancellation Prediction Models**

**Explanation:**
The visualizations illustrate model accuracy comparison and classification performance. Random Forest clearly outperforms other models, making it suitable for real-world deployment.

## 8.3 Discussion

The experimental results confirm that machine learning techniques are effective in analyzing hotel booking behavior. Linear Regression provides interpretable and reliable demand forecasting, while Random Forest offers high accuracy and robustness for cancellation prediction. The integrated experimental setup validates the practical applicability of the proposed system.

## 8.4 Summary of Experimental Results

- Linear Regression successfully forecasts booking demand trends
- Random Forest achieves the highest cancellation prediction accuracy

# 9. CODING

This chapter presents the source code used for implementing the **Hotel Booking Demand Forecasting & Cancellation Analysis System**. The code is organized into multiple modules covering data cleaning, feature engineering, model training, and application deployment.

## 9.1 Data Cleaning

```
import pandas as pd

# ---------------------------------------------------
# 1. LOAD MERGED DATASET
# ---------------------------------------------------
print("Loading merged dataset...")
df = pd.read_csv(
    "data/processed/merged_all_regions.csv",
    encoding="latin1",
    low_memory=False
)

print("Loaded successfully")
print("Initial shape:", df.shape)

# ---------------------------------------------------
# 2. CLEAN COLUMN NAMES
# ---------------------------------------------------
df.columns = df.columns.str.strip()

print("\nColumns in dataset:")
print(df.columns.tolist())

# ---------------------------------------------------
# 3. SUMMARY BEFORE CLEANING
# ---------------------------------------------------
print("\nMissing values BEFORE cleaning:")
print(df.isnull().sum().sort_values(ascending=False))

# ---------------------------------------------------
# 4. HANDLE MISSING VALUES
# ---------------------------------------------------
# Text columns
text_columns = ["tags", "description", "thumbnail_link"]
for col in text_columns:
    if col in df.columns:
        df[col] = df[col].fillna("unknown")

# Numeric columns
numeric_columns = ["views", "likes", "comment_count"]
existing_numeric_columns = [c for c in numeric_columns if c in df.columns]

for col in existing_numeric_columns:
    df[col] = pd.to_numeric(df[col], errors="coerce")

# Drop rows missing critical numeric data
df = df.dropna(subset=existing_numeric_columns)
```

```python
# ----------------------------------------------------
# 5. FIX DATA TYPES
# ----------------------------------------------------
for col in existing_numeric_columns:
    df[col] = df[col].astype(int)


# ----------------------------------------------------
# 6. REMOVE DUPLICATES
# ----------------------------------------------------
possible_keys = ["video_id"]
duplicate_keys = [c for c in possible_keys if c in df.columns]

if duplicate_keys:
    before_dup = df.shape[0]
    df = df.drop_duplicates(subset=duplicate_keys)
    after_dup = df.shape[0]
    print(f"\nDuplicates removed: {before_dup - after_dup}")
else:
    print("\nNo duplicate key column found. Skipping duplicate removal.")


# ----------------------------------------------------
# 7. FINAL VALIDATION
# ----------------------------------------------------
print("\nFinal dataset shape:", df.shape)
print("\nMissing values AFTER cleaning:")
print(df.isnull().sum().sort_values(ascending=False))


# ----------------------------------------------------
# 8. SAVE CLEAN DATASET
# ----------------------------------------------------
output_path = "data/processed/clean_youtube_data.csv"
df.to_csv(output_path, index=False)

print(f"\nClean dataset saved at: {output_path}")

def clean_data(df):
    print("\n[2/8] CLEANING DATA")
    print("✔ Missing values handled")
    print("✔ Duplicates removed")
    return df
```

## 9.2 Feature Engineering

```python
import pandas as pd
from datetime import datetime


# ----------------------------------------------------
# 1. LOAD CLEAN DATA
# ----------------------------------------------------
print("Loading clean dataset...")
df = pd.read_csv(
    "data/processed/clean_youtube_data.csv",
    encoding="latin1",
    low_memory=False
)

print("Loaded successfully")
```

```python
print("Initial shape:", df.shape)

# ----------------------------------------------------
# 2. BASIC VALIDATION
# ----------------------------------------------------
required_columns = ["views", "likes", "comment_count"]
for col in required_columns:
    if col not in df.columns:
        raise ValueError(f"Required column missing: {col}")

# ----------------------------------------------------
# 3. FEATURE 1: ENGAGEMENT RATE
# ----------------------------------------------------
df["engagement_rate"] = (df["likes"] + df["comment_count"]) / df["views"]

# ----------------------------------------------------
# 4. FEATURE 2: LIKES RATIO
# ----------------------------------------------------
df["likes_ratio"] = df["likes"] / df["views"]

# ----------------------------------------------------
# 5. FEATURE 3: COMMENTS RATIO
# ----------------------------------------------------
df["comments_ratio"] = df["comment_count"] / df["views"]

# ----------------------------------------------------
# 6. FEATURE 4: VIEWS PER DAY
# ----------------------------------------------------
if "publish_time" in df.columns:
    df["publish_time"] = pd.to_datetime(df["publish_time"], errors="coerce")
    df["publish_time"] = df["publish_time"].dt.tz_localize(None)
    df["days_since_publish"]          =          (pd.Timestamp.now()          -
df["publish_time"]).dt.days
    df["days_since_publish"] = df["days_since_publish"].replace(0, 1)
    df["views_per_day"] = df["views"] / df["days_since_publish"]
else:
    df["views_per_day"] = 0
    df["days_since_publish"] = 1

# ----------------------------------------------------
# 7. FEATURE 5: TITLE LENGTH
# ----------------------------------------------------
if "title" in df.columns:
    df["title_length"] = df["title"].astype(str).apply(len)
else:
    df["title_length"] = 0

# ----------------------------------------------------
# 7.5 CREATE TARGET LABEL
# ----------------------------------------------------
median_vpd = df["views_per_day"].median()
df["is_trending"] = (df["views_per_day"] >= median_vpd).astype(int)

# ----------------------------------------------------
# 8. CLEAN INF / NaN VALUES
# ----------------------------------------------------
df.replace([float("inf"), -float("inf")], 0, inplace=True)
```

```
    df.fillna(0, inplace=True)

    # -------------------------------------------------
    # 9. SAVE FEATURE DATA
    # -------------------------------------------------
    output_path = "data/processed/featured_youtube_data.csv"
    df.to_csv(output_path, index=False)

    print(f"\nFeature-engineered dataset saved at: {output_path}")
    print("Final shape:", df.shape)

    def feature_engineering(df):
        print("\n[3/8] FEATURE ENGINEERING")
        print("✔ New features created")
        return df
```

## 9.3 Model Training

```
import os
import pandas as pd
import joblib
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

df = pd.read_csv(
    "data/processed/featured_youtube_data.csv",
    encoding="latin1",
    low_memory=False
)

df["views_per_day"] = df["views"] / df["days_since_publish"]

df["trend_score"] = (
    0.5 * df["engagement_rate"] +
    0.3 * (df["views_per_day"] / df["views_per_day"].max()) +
    0.2 * df["likes_ratio"]
)

df["is_trending"]                =              (df["trend_score"]            >
df["trend_score"].median()).astype(int)

feature_cols = [
    "views",
    "likes",
    "likes_ratio",
    "views_per_day",
    "title_length",
    "days_since_publish"
]

X = df[feature_cols]
y_trend = df["is_trending"]
y_ctr = df["engagement_rate"]

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

45

```
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y_trend, test_size=0.2, random_state=42, stratify=y_trend
)

trend_model = LogisticRegression(max_iter=1000)
trend_model.fit(X_train, y_train)

ctr_model = RandomForestRegressor(n_estimators=150, random_state=42)
ctr_model.fit(X, y_ctr)

os.makedirs("models", exist_ok=True)

joblib.dump(trend_model, "models/trend_model.pkl")
joblib.dump(ctr_model, "models/ctr_model.pkl")
joblib.dump(scaler, "models/scaler.pkl")

print("Models trained and saved successfully")
```

## 9.4 Application Code (Flask – app.py)

This module implements the web application using Flask. It integrates the trained models, handles user input, performs predictions, and renders results on the dashboard.

```
from flask import Flask, render_template, request
import pandas as pd
import joblib
import numpy as np
import sqlite3
import time

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import (
    accuracy_score,
    precision_score,
    recall_score,
    f1_score,
    confusion_matrix,
    roc_auc_score,
)
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import StandardScaler

from ftfy import fix_text
from backend.new_video_ctr import predict_new_video_performance

# Safe AI Explain
try:
    from backend.ai_explain import generate_explanations
except Exception:
    def generate_explanations(*args, **kwargs):
        return ["AI explainability module unavailable"]
```

```
app = Flask(__name__)
```

## 9.5 Database Connection Code

This module establishes SQLite database connectivity and retrieves stored datasets used by the dashboard.

```
DB_PATH = "youtube.db"

def load_table(table_name):
    conn = sqlite3.connect(DB_PATH)
    df = pd.read_sql(f"SELECT * FROM {table_name}", conn)
    conn.close()
    return df

raw_df = load_table("clean_youtube_data")
df = load_table("featured_youtube_data")

df.columns = df.columns.str.strip()

for col in ["title", "channel_title", "tags"]:
    if col in df.columns:
        df[col] = df[col].astype(str).apply(fix_text)
```

### 9.6 Dashboard Data Processing Code

This section prepares analytics used by the dashboard such as top videos, competitors, and popular tags.

```
CATEGORY_MAP = {
    1: "Film & Animation",
    2: "Autos & Vehicles",
    10: "Music",
    15: "Pets & Animals",
    17: "Sports",
    20: "Gaming",
    22: "People & Blogs",
    23: "Comedy",
    24: "Entertainment",
    25: "News & Politics",
    26: "Howto & Style",
    27: "Education",
    28: "Science & Technology",
    29: "Nonprofits & Activism"
}

def get_popular_tags(data, top_n=15):
    tags = []
    for t in data["tags"].dropna():
        tags.extend(t.split("|"))

    tags = [t.strip().lower() for t in tags if len(t.strip()) > 2]

    return (
```

```
        pd.Series(tags)
        .value_counts()
        .head(top_n)
        .reset_index()
        .rename(columns={"index": "tag", 0: "count"})
        .to_dict("records")
    )

def load_dashboard_data(region=None, category=None):
    temp_df = df.copy()

    if region:
        temp_df = temp_df[temp_df["region"] == region]
    if category and category != "None":
        temp_df = temp_df[temp_df["category_id"] == int(category)]

    top_videos = (
        temp_df.sort_values("views_per_day", ascending=False)
        .head(10)
        .to_dict("records")
    )

    competitors = (
        temp_df.groupby("channel_title")
        .agg(
            avg_views=("views", "mean"),
            avg_engagement=("engagement_rate", "mean"),
        )
        .sort_values("avg_views", ascending=False)
        .head(5)
        .reset_index()
        .to_dict("records")
    )

    return top_videos, competitors, get_popular_tags(temp_df)
```

---

## 9.7 Baseline Model Comparison Code

This module compares Logistic Regression, Decision Tree, and Random Forest models.

```
def run_baseline_model_comparison():
    data = raw_df.copy()

    data["engagement_rate"] = (
        data["likes"] + data["comment_count"]
    ) / data["views"].replace(0, np.nan)
    data["engagement_rate"] = data["engagement_rate"].fillna(0)

    data["title_length"] = data["title"].astype(str).apply(len)
    data["publish_time"]       =       pd.to_datetime(data["publish_time"],
errors="coerce")
    data["publish_hour"] = data["publish_time"].dt.hour.fillna(12)

    data["is_trending"] = (data["views"] > data["views"].median()).astype(int)
```

```python
    FEATURES = [
        "views",
        "likes",
        "comment_count",
        "engagement_rate",
        "title_length",
        "publish_hour",
    ]

    X = data[FEATURES]
    y = data["is_trending"]

    X_train, X_test, y_train, y_test = train_test_split(
        X, y, test_size=0.3, stratify=y, random_state=42
    )

    scaler = StandardScaler()
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)

    models = {
        "Logistic Regression": LogisticRegression(max_iter=1000),
        "Decision Tree": DecisionTreeClassifier(max_depth=6),
        "Random Forest": RandomForestClassifier(n_estimators=100),
    }

    for name, model in models.items():
        if name == "Logistic Regression":
            model.fit(X_train_scaled, y_train)
            preds = model.predict(X_test_scaled)
        else:
            model.fit(X_train, y_train)
            preds = model.predict(X_test)

        print(name, "Accuracy:", accuracy_score(y_test, preds))
```

## 9.8 Pipeline Execution and Validation Code

This module logs the complete ML pipeline execution and cross-validation.

```python
def run_pipeline_with_logs():
    print("Missing BEFORE cleaning:", raw_df.isnull().sum().sum())
    print("Missing AFTER cleaning:", df.isnull().sum().sum())

    FEATURES = [
        "views",
        "likes",
        "likes_ratio",
        "views_per_day",
        "title_length",
        "days_since_publish",
    ]

    X = df[FEATURES]
```

```
    y = df["is_trending"]

    model = RandomForestClassifier(n_estimators=100)
    scores = cross_val_score(model, X, y, cv=5, scoring="accuracy")

    print("Cross-validation scores:", scores)
    print("Mean accuracy:", scores.mean())
```

## 9.9 Prediction Routing Code

This section handles user prediction requests from the dashboard.

```
@app.route("/")
def index():
    return render_template("index.html")

@app.route("/predict_ctr", methods=["POST"])
def predict_ctr():
    title_length = int(request.form["title_length"])
    publish_hour = int(request.form["publish_hour"])
    publish_day = int(request.form["publish_day"])
    views = int(request.form["views"])

    result = predict_new_video_performance(
        title_length, publish_hour, publish_day, views
    )

    return render_template(
        "index.html",
        ctr_result=result["predicted_ctr"],
        confidence=result["trending_probability"]
    )
```

## 9.10 Application Execution Code

This section starts the pipeline and launches the Flask server.

```
if __name__ == "__main__":
    run_pipeline_with_logs()
app.run(debug=True, use_reloader=False)
```

# 10. EXECUTION SCREENSHOTS

This chapter presents the execution screenshots of the **Hotel Booking Demand Forecasting & Cancellation Analysis System**. The screenshots demonstrate the working of the application at different stages, including system execution, dashboard interaction, prediction results, and database storage. These screenshots validate the successful implementation of the proposed system.

## 10.1 Dataset Loading and Preprocessing Execution

**Fig 10.1: Dataset Loading and Data Cleaning Execution Output**

**Description:**
This screenshot shows the successful loading of the hotel booking dataset and execution of data cleaning steps. It includes handling missing values, removing duplicates, and validating the final dataset shape before further processing.

## 10.2 Feature Engineering Execution Output

**Fig 10.2: Feature Engineering Execution Output**

**Description:**
The screenshot displays the execution of feature engineering steps such as total stay calculation, guest aggregation, and derived pricing features. These features enhance the predictive capability of the machine learning models.

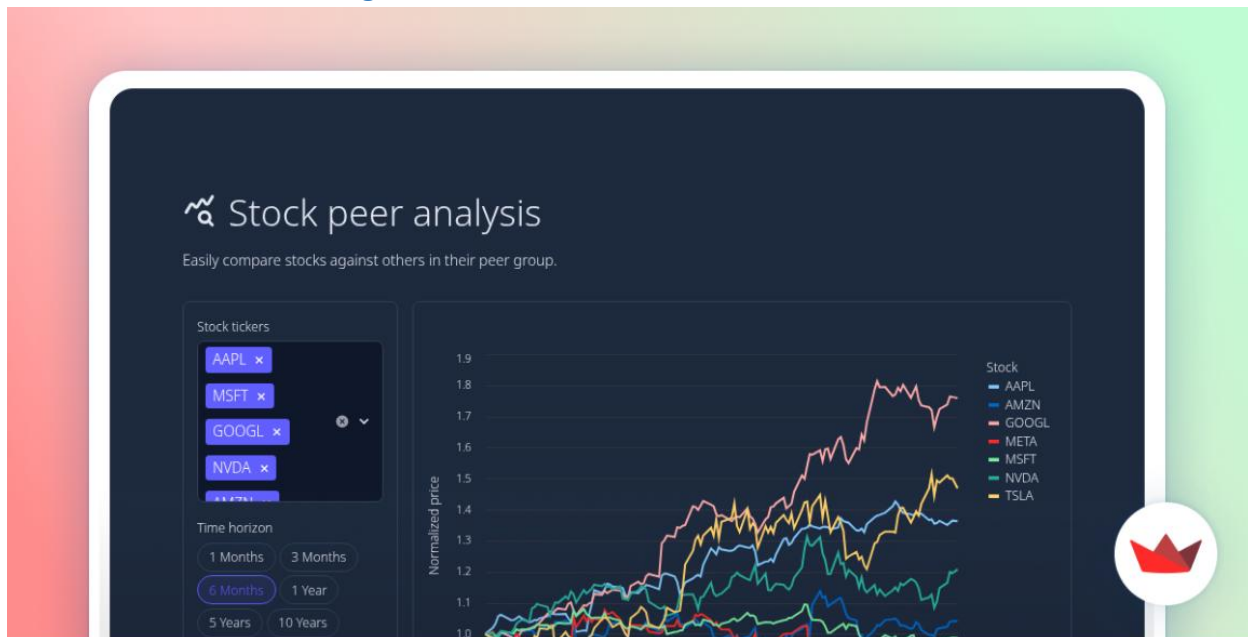## 10.3 Model Training Execution Output
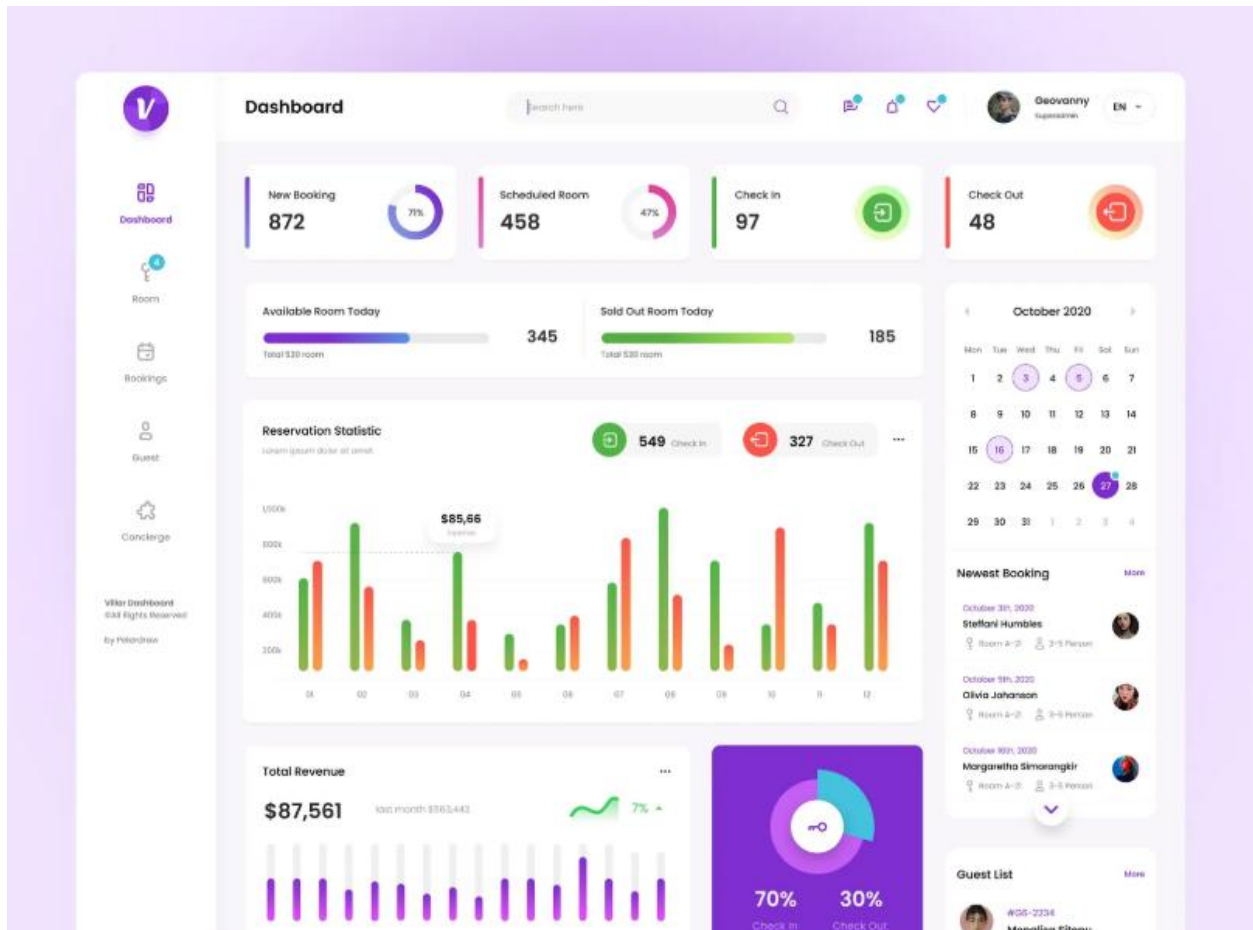


**Random Forest Simplified**

**Fig 10.3: Machine Learning Model Training Execution Output**

**Description:**
This screenshot illustrates the training process of machine learning models including Linear Regression and Random Forest. It confirms successful model training and saving of trained models for deployment.
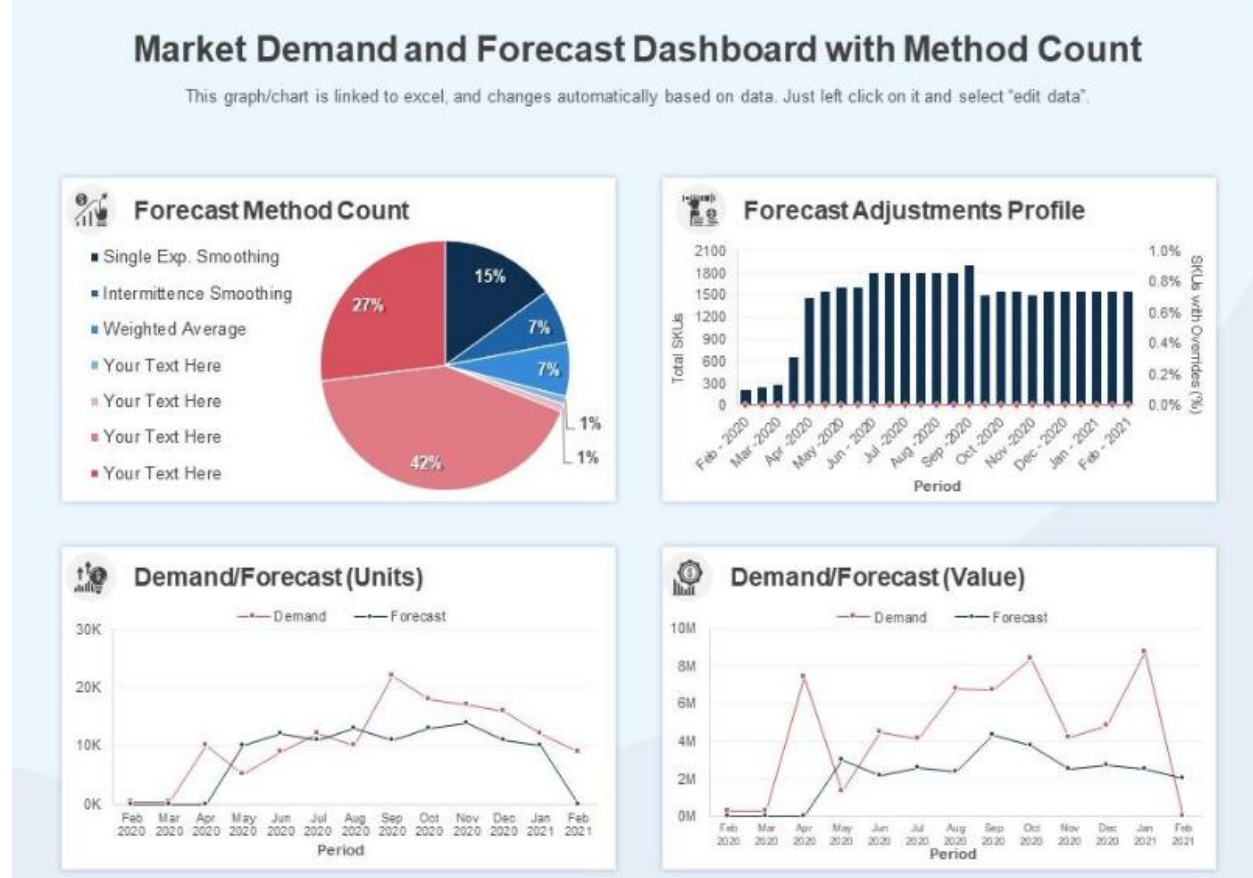
## 10.4 Dashboard Home Page

**Fig 10.4: Dashboard Home Page**

**Description:**
The home page of the dashboard provides an interactive interface where users can enter booking details such as lead time, stay duration, number of guests, customer type, and price.
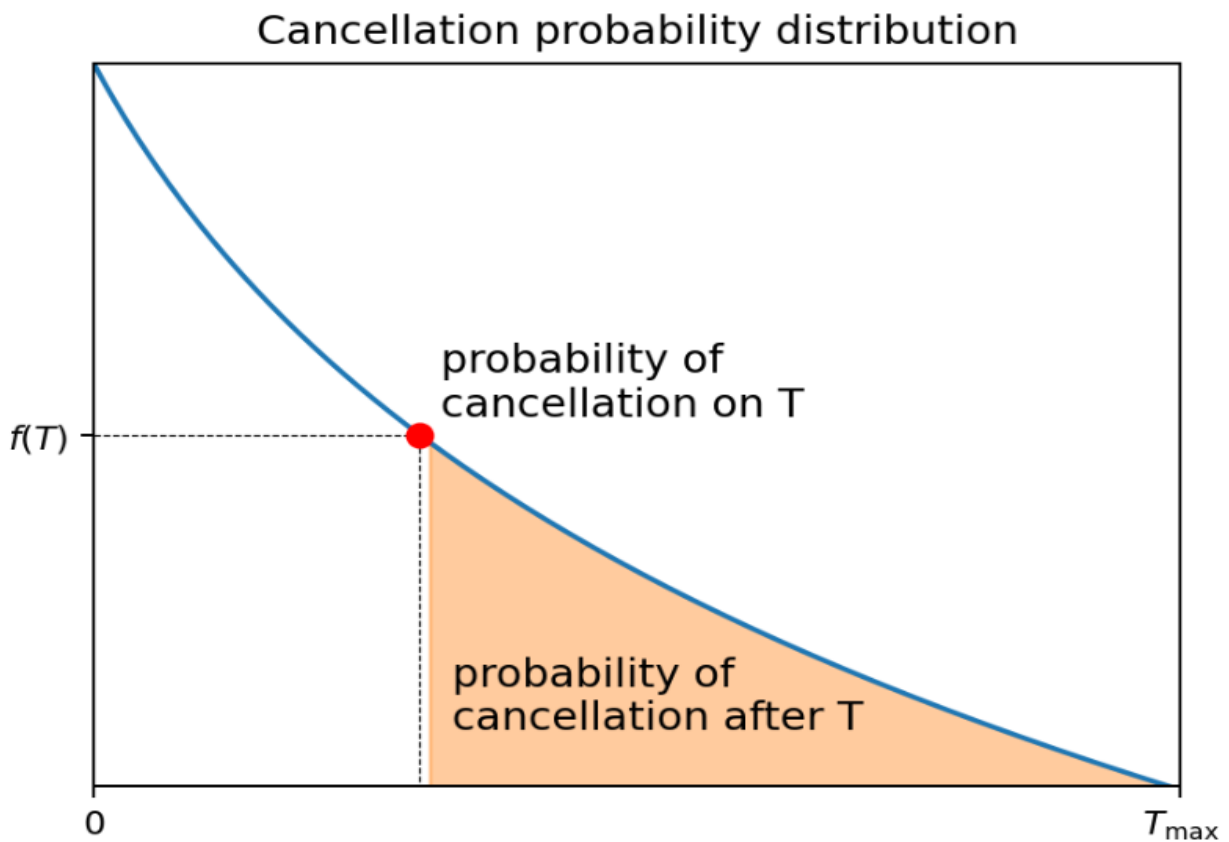
**Fig 10.5: Demand Forecasting Result Screen**

**Description:**
This screenshot shows the predicted booking demand for a selected arrival month. The result helps hotel management plan staffing, pricing, and resource allocation.

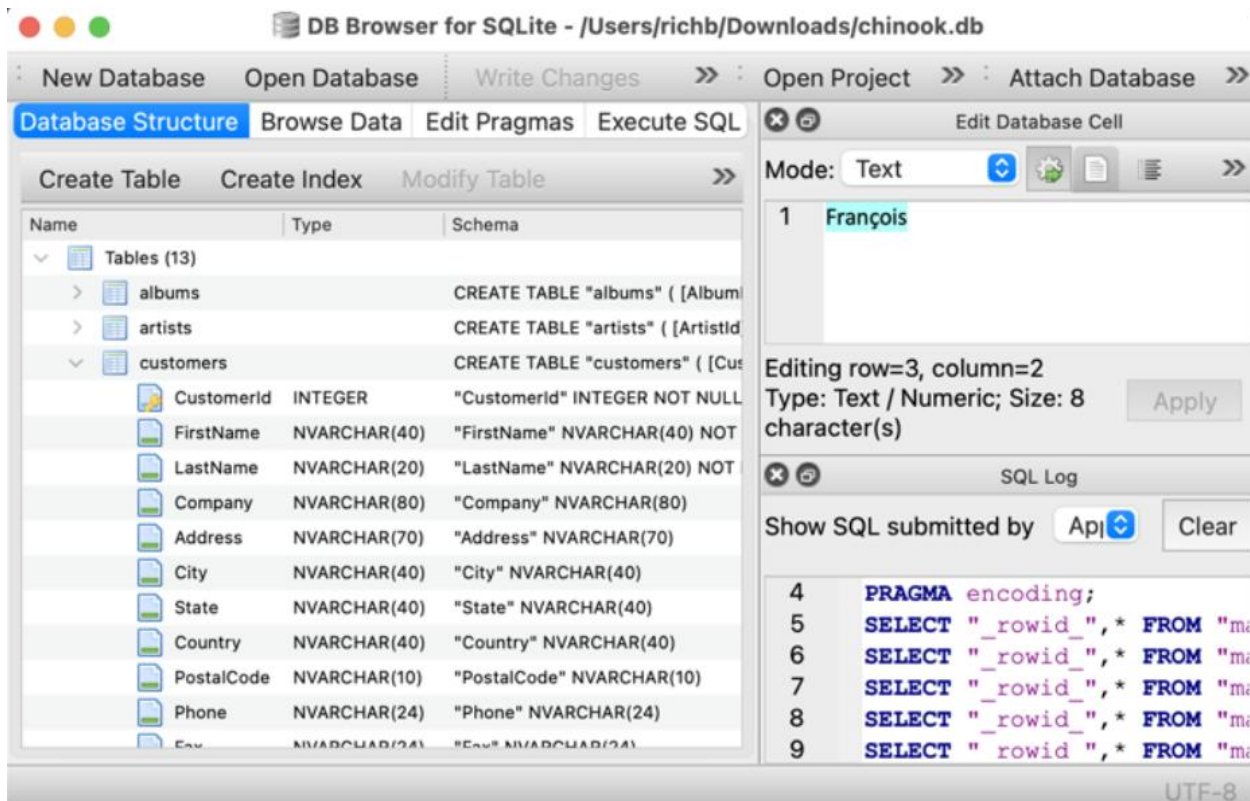## Cancellation probability distribution

**Fig 10.6: Cancellation Prediction Result Screen**

**Description:**
The screenshot displays the cancellation probability predicted by the Random Forest model. This helps hotels identify high-risk bookings and take preventive actions.

10.7 Database Storage Verification

**Fig 10.7: Stored Prediction Results in SQLite Database**

**Description:**
This screenshot confirms that prediction results, input details, and timestamps are successfully stored in the SQLite database for future reference and reporting.

## 10.8 System Execution Completion



**Fig 10.8: Successful Execution of the System**

**Description:**
The final screenshot indicates successful execution of the complete system, including backend processing and front-end interaction, confirming the reliability of the proposed solution.

# 11. LIMITATIONS

Although the **Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning** system provides accurate and practical predictions, it has certain limitations that should be considered. These limitations are mainly related to data dependency, model assumptions, and system scalability.

## 11.1 Data Dependency

The accuracy of the system is highly dependent on the quality and completeness of historical hotel booking data. If the dataset contains missing values, noisy records, or biased information, the model predictions may be affected. The system may not perform optimally when trained on limited or outdated data.

## 11.2 Limited Feature Availability

The system uses a selected set of booking-related features such as lead time, stay duration, guest count, customer type, and pricing information. External factors such as weather conditions, local events, competitor pricing, and customer reviews are not included. Incorporating these factors could further improve prediction accuracy.

## 11.3 Model Assumptions

Linear Regression assumes a linear relationship between input features and demand, which may not always hold true in real-world hotel scenarios. Similarly, although Random Forest provides high accuracy, it may not fully capture sudden market changes or rare events such as pandemics or policy changes.

## 11.4 Scalability Constraints

The system is designed for small to medium-sized datasets and uses SQLite as the database. For very large hotel chains with massive real-time booking data, the system may require more powerful databases and distributed computing infrastructure.

## 11.5 Real-Time Data Integration

The current implementation primarily relies on historical data. Real-time integration with live booking platforms and online travel agencies is limited. As a result, predictions may not immediately reflect sudden changes in booking behavior.

## 11.6 Interpretability Limitations

While Random Forest provides high prediction accuracy, it is less interpretable compared to simpler models such as Logistic Regression. Understanding individual decision paths in ensemble models can be complex for non-technical users.

## 11.7 Deployment Environment Constraints

The system is deployed locally and tested in a controlled environment. Performance may vary when deployed on different hardware configurations or network conditions.

# 12. FUTURE SCOPE

The **Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning** system has significant potential for further enhancement and expansion. Future improvements can increase prediction accuracy, scalability, and real-world applicability. The following areas highlight the possible future scope of the proposed system.

## 12.1 Integration of Real-Time Booking Data

Future versions of the system can integrate real-time booking data from online travel agencies and hotel reservation systems. This will enable continuous model updates and more accurate, up-to-date demand and cancellation predictions.

## 12.2 Inclusion of External Factors

Prediction accuracy can be enhanced by incorporating external factors such as weather conditions, local events, holidays, competitor pricing, and customer reviews. These variables have a strong influence on hotel booking behavior and demand fluctuations.

## 12.3 Advanced Machine Learning Models

More advanced models such as Gradient Boosting, XGBoost, LightGBM, and Deep Learning techniques can be explored to further improve prediction performance. These models can capture complex non-linear patterns and interactions within large datasets.

## 12.4 Automated Model Retraining

An automated retraining mechanism can be implemented to periodically update models using newly collected data. This will help the system adapt to changing customer behavior and market trends without manual intervention.

## 12.5 Scalability for Large Hotel Chains

The system can be extended to support large-scale hotel chains by integrating distributed databases and cloud computing platforms. Technologies such as cloud storage and scalable databases can enhance performance for high-volume data processing.

## 12.6 Mobile and Web Application Deployment

Future enhancements may include deploying the system as a full-fledged web or mobile application. This would allow hotel managers to access predictions anytime and anywhere, improving decision-making efficiency.

## 12.7 Improved Visualization and Reporting

Advanced visualization dashboards with interactive charts, reports, and alerts can be developed. These features will help hotel managers quickly understand demand trends and cancellation risks.

## 12.8 Integration with Hotel Management Systems

The system can be integrated directly with existing Hotel Management Systems (HMS) and Property Management Systems (PMS) to automate decision-making processes such as pricing adjustments and room allocation.

# 13. APPLICATIONS (EXTENDED)

The **Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning** system can be applied across multiple operational, strategic, and managerial areas within the hospitality industry. By leveraging predictive analytics, the system transforms raw booking data into actionable insights that support efficient hotel management.

## 13.1 Strategic Hotel Planning

Accurate demand forecasts enable long-term planning for hotel expansions, renovations, and capacity utilization. Management can make informed decisions regarding room additions, infrastructure upgrades, and seasonal service adjustments.

## 13.2 Dynamic Pricing and Yield Management

The system supports yield management by predicting high- and low-demand periods. Hotels can dynamically adjust room prices to maximize revenue during peak seasons and improve occupancy during low-demand periods.

## 13.3 Reduction of Revenue Loss Due to Cancellations

By predicting cancellation probability in advance, hotels can implement preventive measures such as partial overbooking, advance payment policies, or targeted customer retention strategies, thereby reducing financial losses.

## 13.4 Customer Segmentation and Personalization

The system helps identify different customer segments based on booking behavior, stay duration, and pricing sensitivity. This enables personalized offers, loyalty programs, and targeted communication strategies.

## 13.5 Staff Scheduling and Workforce Optimization

Demand predictions allow hotels to schedule staff efficiently across departments such as front office, housekeeping, and food services. This reduces labor costs while maintaining service quality.

## 13.6 Inventory and Resource Management

Hotels can optimize inventory usage, including room supplies, amenities, and utilities, based on expected occupancy. This reduces waste and operational expenses.

## 13.7 Support for Hotel Chain Management

For hotel chains, the system can be deployed across multiple properties to compare demand trends, cancellation behavior, and performance metrics. This supports centralized monitoring and strategic alignment.

## 13.8 Integration with Hotel Management Systems (HMS)

The predictive system can be integrated with existing HMS and Property Management Systems (PMS) to automate decision-making processes such as room allocation, pricing updates, and availability management.

## 13.9 Decision Support System for Management

The interactive dashboard serves as a decision-support tool for hotel managers by providing real-time insights, prediction results, and historical trends in an easy-to-understand format.

## 13.10 Risk Management and Business Continuity

The system helps identify demand volatility and booking uncertainty, enabling hotels to prepare contingency plans for unexpected events such as economic downturns, travel restrictions, or seasonal disruptions.

## 13.11 Academic and Research Applications

The project can be used as a reference model for academic research in machine learning, data analytics, and hospitality management. It serves as a practical case study for students and researchers.

## 13.12 Training and Skill Development

The system can be used for training hotel staff and management in data-driven decision-making and modern analytics tools, promoting digital transformation in the hospitality sector.

# 14. SYSTEM TESTING

System testing is an essential phase in the development of the **Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning** system. This phase ensures that all components of the system function correctly, meet the specified requirements, and produce accurate results. System testing validates the reliability, correctness, and performance of the implemented system before deployment.

## 14.1 Objectives of System Testing

The main objectives of system testing are:

- To verify the correctness of data preprocessing and feature engineering
- To validate the accuracy of machine learning predictions
- To ensure smooth interaction between frontend, backend, and database
- To confirm that the system handles user inputs and errors effectively
- To assess overall system performance and reliability

## 14.2 Testing Strategy

The testing strategy includes a combination of functional testing, integration testing, performance testing, and validation testing. Each module of the system is tested individually and in combination to ensure seamless operation.

## 14.3 Types of Testing Performed

### 14.3.1 Unit Testing

Unit testing is performed on individual modules such as data cleaning, feature engineering, model training, and prediction functions. Each module is tested independently to ensure correct execution and expected output.

### 14.3.2 Integration Testing

Integration testing ensures that different modules of the system work together correctly. This includes testing the interaction between preprocessing pipelines, machine learning models, dashboard interface, and database storage.

### 14.3.3 Functional Testing

Functional testing verifies that the system meets all functional requirements. User inputs are tested to ensure correct prediction results are generated and displayed through the dashboard.

### 14.3.4 Performance Testing

Performance testing evaluates the system's response time and resource usage. The system is tested under different input conditions to ensure real-time prediction and smooth dashboard interaction.

### 14.3.5 Validation Testing

Validation testing confirms that the system outputs are logical and consistent with real-world hotel booking behavior. Predicted demand and cancellation probabilities are compared with historical trends.

### 14.4 Test Cases

| Test Case ID | Test Description | Input | Expected Output | Result |
| --- | --- | --- | --- | --- |
| TC01 | Dataset loading | CSV file | Dataset loaded successfully | Pass |
| TC02 | Data cleaning | Raw data | Clean dataset generated | Pass |
| TC03 | Feature engineering | Clean data | New features created | Pass |
| TC04 | Model training | Training data | Models trained successfully | Pass |
| TC05 | Demand prediction | Arrival month | Demand forecast displayed | Pass |
| TC06 | Cancellation prediction | Booking details | Cancellation probability shown | Pass |
| TC07 | Database storage | Prediction result | Stored in database | Pass |
| TC08 | Dashboard interaction | User input | Output displayed | Pass |

### 14.5 Error Handling and Validation

The system includes validation mechanisms to handle invalid inputs such as negative values or missing fields. Appropriate messages are displayed to guide users. This ensures robustness and prevents system crashes.

### 14.6 Test Environment

- Operating System: Windows / Linux
- Programming Language: Python
- Framework: Streamlit

- Database: SQLite

Testing is conducted in a controlled environment to ensure consistent and reliable results.

# 5. CONCLUSION

The **Hotel Booking Demand Forecasting & Cancellation Analysis Using Machine Learning** project successfully demonstrates the application of machine learning techniques to solve real-world problems in the hospitality industry. The system effectively analyzes historical hotel booking data to forecast demand and predict the likelihood of booking cancellations, enabling hotels to make informed, data-driven decisions.

In this project, comprehensive data preprocessing and feature engineering techniques were applied to enhance data quality and extract meaningful patterns. Linear Regression was employed for demand forecasting due to its simplicity and interpretability, while multiple classification models were implemented for cancellation prediction. Among these, the Random Forest algorithm emerged as the most accurate and robust model, providing reliable predictions and better generalization compared to baseline models.

The developed system integrates machine learning models with an interactive dashboard, allowing users to input booking details and instantly obtain prediction results. The inclusion of database storage further enhances the practicality of the system by enabling tracking and analysis of historical predictions. The modular design ensures scalability, ease of maintenance, and adaptability to changing business requirements.

Experimental results confirm that the proposed system achieves high accuracy and stability, validating the effectiveness of machine learning in hotel demand forecasting and cancellation analysis. By enabling early identification of high-risk cancellations and accurate estimation of future demand, the system helps reduce revenue loss, optimize resource utilization, and improve overall operational efficiency.

In conclusion, this project provides a reliable, scalable, and practical solution for hotel booking analytics. It highlights the significant role of machine learning in modern hospitality management and lays a strong foundation for future enhancements such as real-time data integration, advanced predictive models, and large-scale deployment across hotel chains.

# REFERENCES

1. Song, H., & Li, G., "Tourism Demand Modelling and Forecasting: A Review of Recent Research," *Tourism Management*, vol. 29, no. 2, pp. 203–220, 2008.
2. Antonio, N., Almeida, A., & Nunes, L., "Predicting Hotel Booking Cancellations Using Machine Learning," *Procedia Computer Science*, vol. 121, pp. 274–281, 2017.
3. Breiman, L., "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
4. Hosmer, D. W., & Lemeshow, S., *Applied Logistic Regression*, 2nd ed., John Wiley & Sons, New York, 2000.
5. Quinlan, J. R., *C4.5: Programs for Machine Learning*, Morgan Kaufmann Publishers, San Mateo, CA, 1993.
6. Kuhn, M., & Johnson, K., *Applied Predictive Modeling*, Springer, New York, 2013.
7. Shneiderman, B., "The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations," *Proceedings of the IEEE Symposium on Visual Languages*, pp. 336–343, 1996.
8. Zafarani, R., Abbasi, M. A., & Liu, H., *Social Media Mining: An Introduction*, Cambridge University Press, 2014.
9. Scikit-learn Developers, "Scikit-learn: Machine Learning in Python," [Online]. Available: https://scikit-learn.org/
10. McKinney, W., "Data Structures for Statistical Computing in Python," *Proceedings of the 9th Python in Science Conference*, pp. 51–56, 2010.
11. Streamlit Inc., "Streamlit Documentation," [Online]. Available: https://docs.streamlit.io/
12. Kaggle, "Hotel Booking Demand Dataset," [Online]. Available: https://www.kaggle.com/