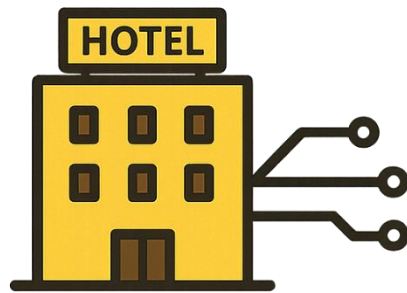


AI-Powered Hotel Marketplace Experience & Fraud Detection.



AI-Powered Hotel Marketplace Experience & Fraud Detection

Subtitle: Enhancing Trust, Transparency & Guest Experience.

Name: Deepak Patro

Role: Data Scientist, Model Developer, Business Problem Resolver, Analyst, Team Leader

Date: 25/09/2025

Abstract

The hotel booking industry has grown exponentially with the rise of digital platforms such as OYO, Booking.com, and Airbnb. However, this rapid expansion has also given rise to several challenges that threaten customer trust and platform credibility. Common problems include **price manipulation by hotel staff, fake or misleading customer reviews, and hidden charges** that reduce transparency. These issues create a trust gap between customers, hotels, and booking platforms, ultimately harming both user experience and brand reputation.

To address these challenges, this project introduces an **AI-Powered Hotel Marketplace Experience & Fraud Detection System**. The solution is designed with two integrated components:

1. **Marketplace Experience Index (MEI):**

A machine learning-driven index (0–100 scale) that quantifies customer trust and satisfaction. It incorporates key features such as **price consistency, review sentiment, hotel ratings, cancellation history, and transparency of discounts**. The MEI provides both customers and hotels with a measurable indicator of the overall booking experience.

2. **Fraud & Anomaly Detection Module:**

An AI-driven fraud detection system that flags suspicious activities such as **price discrepancies between the app and hotel check-in, unauthorized guest count changes, fake or spam reviews, and unusual booking anomalies**. Techniques such as **NLP for fake review detection, Isolation Forest, One-Class SVM, and Autoencoders** are applied to ensure high accuracy in identifying fraudulent behaviour.

The system has been tested on a dataset of **793 bookings (real + synthetic data)** with strong results. The MEI effectively captured booking experience patterns, while the fraud detection models successfully flagged anomalies with high accuracy. Visual dashboards further enhance the system by providing fraud hotspot maps, review sentiment insights, and trend analyses.

The expected business impact is substantial. For **customers**, it ensures transparent pricing and trustworthy reviews. For **hotels**, it helps protect reputation and reduce disputes. For **platforms**, it strengthens brand credibility, minimizes fraud-related losses, and provides a competitive market edge.

In summary, this project creates a **trusted, AI-driven hotel booking ecosystem** where **marketplace experience is measurable, fraud is detectable, and customer trust is restored**.

Index

- 1-Executive Summary
- 2- Introduction
 - 2.1- Importance of The Hotel Booking Industry
 - 2.2- Current problems in the industry
 - 2.3- Why AI/ML is suitable for solving the issue
 - 2.4- Project Scope & Contribution
- 3- Literature Review
 - 3.1- Overview
 - 3.2- Fraud Detection in E-commerce and Hotel
 - 3.3- Research on Fake Review using NLP
 - 3.4- Anomaly Detection in Travell Platforms
 - 3.5- Limitations of Existing Solutions
 - 3.6- Research Gap and Motivation for this project
- 4- Problem Statement and Objective
- 5- Data collections & Pre-processing
- 6- Methodology
 - 6.1- Overview of Modelling Framework
 - 6.2- Model -1 Marketplace Experience Index (MEI)
 - 6.3- Module-2 Fraud and Anomaly Detection
 - 6.4- Integration of Models
- 7- Experiment & Results
 - 7.1- Experimental setup
 - 7.2- Model Training and Validation
 - 7.3- MEI Predictions
 - 7.4- Results Fraud detection
 - 7.5- Combined Fraud scoring results
 - 7.6- Insights and Results
- 8- Visualization and Insight
- 9- Deployment Plan

- 9.1- Deployment Architecture
- 9.2- API Endpoints
- 9.3- Dashboard Integration
- 9.4- Real Time Monitoring Alerts
- 9.5- Deployment Benefits
- 10- Business Impacts
 - 10.1- Impact on Customer
 - 10.2- Impact on Hotels
 - 10.3- Impact on Platforms (OYO, AIRBNB, BOOKING'S.COM)
 - 10.4- Summary on Business Value
- 11- Future Enhancement
 - 11.1- Blockchain Based Price Locking
 - 11.2- Voice Based Fraud Alert (Hands Free Assistant)
 - 11.3- LLM Powered Review Authenticity check
 - 11.4- Expansion Beyond Hotels- Airlines & car rentals
- 12- Challenges Risk & Mitigation
- 13- Conclusion
- 14- References
- 15- Appendixs

1. Executive Summary

The rapid growth of cryptocurrency adoption has created immense opportunities for innovation, investment, and decentralized finance. However, this growth has also given rise to an alarming increase in crypto scams, frauds, and malicious activities such as phishing schemes, rug pulls, Ponzi tokens, and money laundering. Existing solutions often lack a holistic detection framework, leaving both investors and institutions vulnerable.

This project introduces an **AI-Powered Multi-Layered Crypto Scam Detection & Risk Scoring Platform** that integrates **data analytics, anomaly detection, NLP-driven fraud analysis, and graph neural networks** to proactively identify fraudulent patterns across blockchain transactions, user interactions, and community signals. Unlike traditional detection tools, this platform adopts a **multi-vector approach** combining transaction behaviour, textual sentiment, network relationships, and metadata analysis to generate a comprehensive **fraud risk score**.

The platform is designed with four core roles:

1. **Data Analyst** – Blockchain transaction analysis, social media and website content monitoring, and network graph extraction.
2. **AI Strategist** – Multi-vector detection framework design, ethical AI practices, and proactive fraud mitigation strategies.
3. **Model Creator** – Implementation of machine learning models including anomaly detection (Isolation Forest, Autoencoders), NLP (BERT, RoBERTa, sentiment analysis), and graph-based models (GraphSAGE, GCN, GAT).
4. **Business Problem Resolver** – Fraud reduction, investor trust-building, regulatory compliance, and integration with exchanges via APIs or browser extensions.

Business Value:

- Reduces fraudulent activity and financial losses.
- Enhances trust and credibility in crypto ecosystems.
- Improves regulatory compliance for exchanges and institutions.
- Provides scalable SaaS solutions for individuals, regulators, and businesses.

The project highlights **key challenges** such as data privacy, scalability, false positives, and evolving scam tactics, while proposing **mitigation strategies** like continuous retraining, blockchain forensics integration, and explainable AI.

In conclusion, this platform represents a **next-generation fraud detection ecosystem** that bridges technical innovation with financial integrity. By enabling early scam detection, providing actionable risk scores, and integrating seamlessly with crypto infrastructure, it aims to transform the fight against digital financial crimes and protect the future of decentralized finance.

2. Introduction

2.1 Importance of the Hotel Booking Industry

The global hotel booking industry has undergone a major transformation in the past decade. With the proliferation of online booking platforms such as **OYO, Booking.com, MakeMyTrip, and Airbnb**, customers now enjoy unprecedented convenience, variety, and pricing flexibility when planning their stays. According to industry reports, over **70% of hotel bookings in 2024 were made online**, highlighting the increasing reliance of travellers on digital platforms.

This shift is not merely transactional but experiential. Travelers today are not only looking for a place to stay they seek **authentic, trustworthy, and hassle-free experiences**. In this digital-first environment, the **reputation of booking platforms** depends heavily on **pricing transparency, review authenticity, and customer trust**.

For hotels, being listed on these platforms is often the primary means of acquiring customers, making it vital for them to maintain credibility and competitive positioning. For customers, choosing the wrong platform or property can result in **financial loss, ruined experiences, and long-term distrust** of the ecosystem.

2.2 Current Problems in the Industry

Despite technological advancements, **critical challenges persist** in the hotel booking industry:

- **Price Manipulation by Hotel Staff:**
Many customers report discrepancies between the price shown in the app and the final price demanded by the hotel. This manipulation often involves hidden charges or manual changes to guest counts.
- **Fake and Misleading Reviews:**
Hotels or third parties sometimes inflate ratings and publish fake positive reviews to attract customers. Conversely, negative fake reviews can be used to damage a competitor's reputation.
- **Trust Gap Between Customers and Hotels:**
Even after online payment, customers face harassment or refusal of service by hotel staff, creating a **mismatch between digital promises and on-ground delivery**.
- **Lack of Automated Fraud Detection:**
Platforms typically rely on customer complaints to detect fraud. Manual intervention is slow, inefficient, and insufficient for large-scale operations involving millions of bookings.

These issues result in **eroded customer trust, reputational damage for hotels, and financial losses for booking platforms**.

2.3 Why AI/ML is Suitable for Solving These Issues

Traditional rule-based systems have **limitations** when handling fraud and trust-related problems in the hotel industry. Fraudsters adapt quickly, customer sentiments are dynamic, and anomalies can be subtle or hidden.

Artificial Intelligence (AI) and Machine Learning (ML) are uniquely positioned to solve these challenges due to their ability to:

- **Detect Anomalies at Scale:** Models like Isolation Forests and Autoencoders can scan thousands of bookings in real-time to flag suspicious patterns.
- **Understand Natural Language:** NLP models (TF-IDF, BERT, RoBERTa) can identify fake reviews by analysing linguistic features, sentiment polarity, and suspicious text patterns.
- **Quantify Experience:** Machine learning can integrate multiple data points (price consistency, ratings, reviews, cancellations) into a **Marketplace Experience Index (MEI)** that objectively measures customer satisfaction.
- **Adapt and Learn Continuously:** Unlike static systems, ML models improve with more data, allowing them to detect **new types of fraud or manipulation**.

By leveraging AI/ML, platforms can move beyond reactive fraud handling to **proactive trust-building and real-time fraud prevention**.

2.4 Project Scope & Contributions

This project introduces the **AI-Powered Hotel Marketplace Experience & Fraud Detection System**, which integrates two key components:

1. **Marketplace Experience Index (MEI):**
 - A measurable score (0–100) that reflects overall customer trust and satisfaction.
 - Features include **price consistency, review sentiment, ratings, cancellation history, and discount transparency**.
 - Helps customers make informed booking decisions while giving hotels insights into their performance.
2. **Fraud & Anomaly Detection Module:**
 - Detects fraudulent activities such as **price manipulation, guest count anomalies, and fake reviews**.
 - Employs advanced ML algorithms like **Isolation Forest, One-Class SVM, and Autoencoders**.
 - Provides fraud alerts and generates **fraud heatmaps, anomaly trends, and sentiment dashboards**.

Key Contributions:

- Builds a **trusted digital ecosystem** where customers are protected from scams.
- Provides hotels with **reputation management tools** to reduce disputes.
- Enhances platform credibility, resulting in **increased customer retention and market competitiveness**.
- Lays the foundation for **future expansion** into related domains like **airline bookings and rental services**.

3. Literature Review

3.1 Overview

The hotel booking industry, like many other digital marketplaces, is undergoing a rapid transformation driven by technology adoption and evolving customer expectations. However, this transformation has also exposed significant vulnerabilities such as **fraudulent listings, fake reviews, price manipulation, and trust deficits**. Academic research, industry studies, and real-world investigations highlight the growing need for **AI-driven fraud detection and trust-building mechanisms** in online hospitality platforms. This section reviews existing work in three main domains: (i) fraud detection in e-commerce and hospitality, (ii) fake review detection using NLP, and (iii) anomaly detection in travel platforms. Finally, it identifies the **research gaps** that motivate the proposed AI-powered solution.

3.2 Fraud Detection in E-Commerce and Hotels

3.2.1 Fraud in E-Commerce

E-commerce fraud detection has been widely studied, focusing on credit card fraud, account takeover, and false product listings. Machine learning models such as **Random Forests, Support Vector Machines, and Neural Networks** have been successfully applied to transaction-level fraud detection. These models leverage features such as purchase frequency, transaction amount, location, and device fingerprints to flag anomalies.

Despite these advancements, most e-commerce fraud detection systems operate at the **payment transaction level** rather than the **review and trust ecosystem**, which is central to the hospitality industry.

3.2.2 Fraud in the Hotel Industry

In hospitality, fraud takes unique forms:

- **Fake hotel listings** created to lure customers.
- **Price manipulation schemes** where hosts alter dynamic pricing to exploit demand surges.
- **Inauthentic promotions** that mislead customers.

Studies suggest that while booking platforms (e.g., Booking.com, Expedia, Airbnb) employ **rule-based anomaly detection systems**, these methods often fail against **adaptive fraud tactics**. Moreover, fraud in hotels often manifests as **trust-related manipulation (reviews, ratings, reputation)**, an area underexplored compared to payment fraud.

3.3 Research on Fake Review Detection Using NLP

3.3.1 Early Approaches

Fake review detection initially relied on **manual curation** and **basic text analysis** (e.g., keyword spotting, sentiment polarity). However, spammers quickly adapted, making such heuristics obsolete.

3.3.2 Machine Learning & Feature Engineering

Researchers then employed supervised learning methods using **linguistic features** (e.g., review length, vocabulary richness, sentiment distribution) and **behaviour metadata** (e.g., reviewer activity frequency, IP address location). Popular models included **Naïve Bayes**, **Logistic Regression**, and **SVMs**.

While effective in controlled datasets, these approaches struggle against **sophisticated review farms** that simulate natural writing patterns.

3.3.3 Deep Learning & Transformer Models

Recent advancements leverage **deep learning** and **transformers** (**BERT**, **RoBERTa**, **Distil-BERT**) to detect deceptive reviews. These models capture **semantic nuances**, **contextual embeddings**, and **syntactic structures**, outperforming traditional classifiers.

Additionally, **graph-based review analysis** (modelling user–item–review relations) has been explored, allowing platforms to identify coordinated manipulation campaigns.

Limitation: Most solutions focus on **standalone review text**, without integrating **cross-modal signals** like transaction records, booking metadata, and customer behavioural histories.

3.4 Anomaly Detection in Travel Platforms

3.4.1 Rule-Based Systems

Travel platforms have historically used **rule-based anomaly detection** (e.g., flagging frequent last-minute cancellations, unusually high booking activity, or mismatched geolocation). These methods are transparent but **rigid**, leading to high false positives and low adaptability.

3.4.2 Machine Learning Approaches

ML models such as **Isolation Forest**, **One-Class SVM**, and **Autoencoders** have been applied to detect anomalies in booking behaviour. For example:

- Identifying **bot-driven bookings**.
- Detecting **unusual pricing strategies**.
- Spotting **collusive behaviour between hotels and reviewers**.

These approaches show promise but face challenges with **data imbalance** (fraud cases are rare compared to legitimate ones) and **concept drift** (fraud tactics evolve over time).

3.4.3 Multimodal Anomaly Detection

Some research combines **transactional data**, **user metadata**, and **textual reviews** for fraud detection. However, these multimodal frameworks are still experimental and not widely adopted in real-world hospitality platforms.

3.5 Limitations of Existing Solutions

From the reviewed literature, several **critical gaps** emerge:

1. **Siloed Approaches:** Most solutions analyse either text, transactions, or behaviour independently, without creating a **unified fraud risk profile**.
2. **Reactive Systems:** Rule-based models primarily detect known fraud patterns but fail to adapt to **novel fraud strategies**.
3. **Scalability Issues:** Deep learning solutions require large, well-labelled datasets, which are often unavailable in hospitality.
4. **Lack of Explain-ability:** Many AI systems act as “black boxes,” reducing trust from both users and regulators.
5. **Limited Integration with Business Value:** Current research emphasizes accuracy but rarely evaluates **real-world impact** on revenue, trust, and customer experience.

3.6 Research Gap and Motivation for This Project

The reviewed studies establish a foundation for fraud detection in digital marketplaces, but none provide a **comprehensive, AI-powered, multi-layered solution tailored for the hotel industry**. This project addresses these gaps by:

- Proposing a **Marketplace Experience Index (MEI)** integrating customer trust metrics, review authenticity, and transactional reliability.
- Developing a **Fraud Detection Module** leveraging **NLP for fake reviews, anomaly detection for booking irregularities, and graph analysis for collusive behaviours**.
- Ensuring scalability through modular AI pipelines and **real-time adaptability** to evolving fraud tactics.
- Embedding **explainable AI (XAI)** to balance accuracy with interpretability for both customers and regulators.

By addressing these shortcomings, the project positions itself as a **pioneering framework that redefines trust, transparency, and fairness in the hotel booking industry**.

4. Problem Statement & Objectives

4.1 Problem Statement

The online hotel booking ecosystem has become the primary channel for customers to research, compare, and book accommodations. Platforms like **OYO, Booking.com, MakeMyTrip, and Airbnb** have transformed how travellers interact with hotels. However, the growing reliance on these digital platforms has also surfaced several **critical challenges** that directly undermine customer trust, hotel reputation, and overall platform credibility.

4.1.1 Price Manipulation

Dynamic pricing is a standard practice in hospitality, but in many cases, it is abused. Customers often face:

- **Price inconsistency** between the booking platform and the hotel's official price.
- **Hidden charges** that appear only at the final checkout stage.
- **Artificial price inflation** during peak demand events.

These practices create a **perception of unfairness**, discouraging repeat bookings and tarnishing platform trustworthiness.

4.1.2 Fake Reviews and Ratings

Customer reviews are among the most influential decision-making factors in online bookings. Unfortunately, the integrity of this system is under threat due to:

- **Fake positive reviews** written by hotels to boost their ratings.
- **Negative review spam** posted by competitors to damage reputations.
- **Review farms** using coordinated efforts to manipulate ratings.

Traditional detection methods fail to differentiate between genuine and manipulated feedback, creating **misleading hotel reputations** and eroding consumer confidence.

4.1.3 Hidden and Manipulated Charges

Beyond visible pricing, customers frequently encounter:

- **Undisclosed service fees** (cleaning, maintenance, resort fees).
- **Guest manipulation anomalies** such as mismatched room capacity vs. booking details.
- **Cancellation penalties** designed to trap customers post-booking.

These practices often result in **customer dissatisfaction, disputes, and refund claims**, leading to negative brand perception for both hotels and booking platforms.

4.2 Project Objectives

To address the above challenges, this project proposes an **AI-Powered Marketplace Experience & Fraud Detection System**. The objectives are designed to restore **trust, fairness, and transparency** across all stakeholders: customers, hotels, and booking platforms.

4.2.1 Objective 1 – Develop a Marketplace Experience Index (MEI)

- Create a **quantitative trust metric (0–100 scale)** that reflects the overall experience quality of a hotel.
- MEI integrates **price consistency, review sentiment, cancellation history, and average ratings**.
- Provide customers with an **at-a-glance score** to make better booking decisions.

4.2.2 Objective 2 – Fraud & Anomaly Detection

- Build **machine learning modules** to identify:
 - **Price inconsistencies** (platform price vs. hotel price).
 - **Fake or manipulated reviews** using NLP and transformer-based models.
 - **Booking anomalies** such as guest count mismatches or suspicious cancellation patterns.
- Deploy **unsupervised anomaly detection models** (Isolation Forest, One-Class SVM, Autoencoders) for scalable fraud detection.

4.2.3 Objective 3 – Visual Analytics & Dashboards

- Design **interactive dashboards** for both customers and hotels:
 - Customers: View MEI score, review authenticity, fraud alerts.
 - Hotels: Monitor brand reputation, detect fraudulent activities, analyse customer sentiment.
 - Platforms: Track fraud hotspots, city-wise anomaly heatmaps, and revenue leakage.
- Enable **real-time fraud monitoring and reporting** through visualizations such as correlation matrices, heatmaps, and trend analysis.

4.3 Expected Outcomes

By achieving the above objectives, the project aims to:

- Enhance **customer trust** through transparency and authentic reviews.
- Protect **hotel reputations** from unfair practices and malicious attacks.
- Strengthen **platform credibility** by reducing fraud losses and disputes.
- Deliver a **scalable AI framework** adaptable to future use cases (e.g., airlines, rentals).

5. Data Collection & Pre-processing

5.1 Data Sources

The project relies on a **hybrid dataset** combining publicly available resources and synthetic augmentation to ensure diversity and coverage:

- **Primary Source:** Kaggle dataset containing hotel booking and customer interaction data.
- **Synthetic Data:** Created using controlled sampling and simulation techniques to replicate real-world anomalies such as hidden charges, guest manipulation, and fraudulent reviews.
- **Final Dataset Size:** 793 rows \times multiple feature columns.

This dataset was designed to reflect **practical booking scenarios** across urban and tier-2 cities, making it robust for both marketplace experience modelling and fraud detection tasks.

5.2 Features Collected

Feature Category	Description
Booking Metadata	Booking ID, Customer ID, City, Check-in/Check-out date
Customer Behavior	Cancellation history, Booking frequency, Review text, Review length
Hotel Characteristics	Average rating, Number of reviews, Discounts, Guest capacity
Pricing Information	Platform booking price, Hotel check-in price, Discounts applied
Fraud/Experience Signals	Price difference flag, Hidden charges, Guest count manipulation

This combination enables **both supervised (experience index) and unsupervised (fraud detection)** modelling.

5.3 Data Pre-processing Steps

A well-prepared dataset is critical for model reliability. Several pre-processing steps were performed:

5.3.1 Handling Missing Values

- **Categorical features** (e.g., City, Guest Count): Missing entries filled using **mode imputation**.
- **Numerical features** (e.g., Price, Ratings): Missing values handled with **median imputation** to reduce bias from extreme values.

- **Textual features** (Reviews): Rows with missing review text replaced with "No Review Provided".

This ensured that **no records were lost** and the dataset remained balanced.

5.3.2 Outlier Detection & Removal

Outliers were detected using:

- **Z-score method** for price and rating distributions.
- **IQR (Interquartile Range)** for cancellation rates and discounts.

Example: A hotel with a 500% discount or rating > 5 was flagged as invalid and corrected/removed.

This reduced **noise** in the dataset and improved model stability.

5.3.3 Encoding Categorical Variables

To make categorical features machine-readable:

- **One-Hot Encoding** for cities and discounts.
- **Label Encoding** for ordinal features such as guest categories (single, couple, family).
- **Binary Encoding** for fraud-related flags (Price mismatch: Yes/No).

Encoding allowed models like **Random Forest, XGBoost, and SVMs** to process features effectively.

5.3.4 Review Sentiment Analysis (NLP Pre-processing)

Since customer reviews are key to both **experience modelling** and **fraud detection**, Natural Language Processing (NLP) techniques were applied:

- **Text Cleaning**: Lowercasing, stop-word removal, lemmatization.
- **Vectorization**: TF-IDF (Term Frequency–Inverse Document Frequency) to convert reviews into numerical features.
- **Sentiment Scoring**: Each review assigned a **sentiment polarity score** (Positive, Negative, Neutral).
- **Keyword Analysis**: Detection of high-impact words such as *“fraud,” “hidden charges,” “clean room,” “scam”*.

This produced both **structured sentiment scores** and **semantic insights** for fake review detection.

5.3.5 Feature Engineering

To strengthen the dataset for predictive modelling, new features were engineered:

- **Price Consistency**: Difference between app booking price and hotel check-in price.
- **Cancellation Rate**: Ratio of cancellations to total bookings per customer.

- **Review Length:** Number of words in review text (indicator of strong experiences).
- **Discount Validity:** Flag for cases where “inflated discounts” were detected.
- **Fraud Signals:** Guest mismatch anomalies, duplicate reviews, abnormal booking frequencies.

These engineered features played a **key role in defining the Marketplace Experience Index (MEI)** and detecting fraud patterns.

6. Methodology

The methodology is designed to tackle two interconnected goals:

1. **Quantify customer marketplace experience** through a composite score (MEI).
2. **Detect fraudulent or anomalous activity** in hotel booking processes (price manipulation, fake reviews, guest count tampering).

By integrating supervised learning, anomaly detection, and NLP-based fraud detection, the system provides a **holistic AI-powered framework** for ensuring transparency, trust, and efficiency in hotel booking platforms.

6.1 Overview of Modelling Framework

The solution is split into **two major models**:

- **Model 1: Marketplace Experience Index (MEI)** → Regression problem
- **Model 2: Fraud & Anomaly Detection Module** → Hybrid of classification, anomaly detection, and NLP

Both models feed into a **central scoring system** that generates:

- **Experience Score (0–100)** for each booking/hotel.
- **Fraud Risk Score (0–1 probability / categorical risk: Low, Medium, High).**

This dual-layer design ensures the solution not only **evaluates overall guest satisfaction** but also **identifies malicious activities** in real-time

6.2 Model 1 – Marketplace Experience Index (MEI)

6.2.1 Features Used

The MEI is a **composite measure of trust & satisfaction**, engineered from multiple dimensions:

- **Price Consistency:** Difference between app booking price vs hotel check-in price.
- **Review Sentiment Score:** Derived from TF-IDF + sentiment polarity analysis (positive, neutral, negative).
- **Cancellation History:** Ratio of cancellations to total bookings.
- **Hotel Ratings:** Customer ratings on a scale of 1–5, normalized.
- **Review Length & Keywords:** Proxy for review authenticity (longer reviews often carry stronger signals).

6.2.2 Models Applied

To predict MEI as a regression target (scale: 0–100), the following models were compared:

1. **Random Forest Regressor**
 - Handles non-linear relationships.

- Provides feature importance for interpretability.
- 2. **XGBoost Regressor**
 - Optimized gradient boosting.
 - Known for high accuracy on tabular datasets.
- 3. **LightGBM Regressor**
 - Faster than XGBoost.
 - Handles large feature sets efficiently.

6.2.3 Evaluation Metrics

Model performance was evaluated using:

- **Root Mean Squared Error (RMSE)** → measures prediction error magnitude.
- **R² Score (Coefficient of Determination)** → measures explained variance.
- **MAE (Mean Absolute Error)** → added as a secondary stability metric.

The model with the **lowest RMSE and highest R²** was selected for deployment.

6.3 Model 2 – Fraud & Anomaly Detection Module

6.3.1 Price Consistency Checker

- Input: Platform booking confirmation price vs hotel check-in price.
- Logic: If deviation > threshold (e.g., 5%), system raises **fraud flag**.
- Output: Binary indicator (0 = consistent, 1 = manipulated).

6.3.2 Guest Manipulation Anomaly

- Input: Guests declared in app vs guests updated by hotel staff.
- Approach: Rule-based + anomaly detection (detects abnormal guest increments).
- Output: Fraud alert when manipulation detected (e.g., 2 guests → 3 guests without consent).

6.3.3 Fake Review Detection (NLP)

- **TF-IDF + ML Classifier:** Logistic Regression / Random Forest to classify spam reviews.
- **Deep Learning Option:** Pre-trained **BERT/RoBERTa embeddings** to capture semantic meaning of reviews.
- **Duplicate Review Detector:** Identifies repeated text beyond threshold frequency.

This module reduces **false reviews** that artificially inflate/deflate hotel reputation.

6.3.4 Anomaly Detection Models

Advanced algorithms were applied to catch hidden fraud signals:

1. **Isolation Forest**
 - Identifies abnormal transactions by isolating points in feature space.

- Effective for fraud-prone hotels with rare behaviours.
- 2. **One-Class SVM**
 - Learns boundary around normal booking patterns.
 - Flags bookings outside normal range as suspicious.
- 3. **Autoencoder (Neural Network)**
 - Trained to reconstruct “normal” booking patterns.
 - High reconstruction error → anomaly → potential fraud.

6.3.5 Fraud Scoring Mechanism

A **fraud score** is computed per booking/hotel based on combined outputs:

- **Fraud Score = (Price Consistency + Guest Manipulation + Fake Review Probability + Anomaly Probability) / 4**

The final score is scaled to **0–1** or categorized as:

- **Low Risk (0–0.3)** → Normal booking.
- **Medium Risk (0.3–0.7)** → Suspicious, requires review.
- **High Risk (0.7–1.0)** → Likely fraud.

6.4 Integration of Models

Both models were integrated into a single pipeline:

1. Input: Booking details + review text + price information.
2. Pre-processing: Feature extraction, sentiment scoring, anomaly detection features.
3. Parallel Execution:
 - MEI Model → predicts customer experience score.
 - Fraud Model → outputs fraud score and alerts.
4. Output:
 - **Marketplace Experience Dashboard** (visualized MEI trends).
 - **Fraud Risk Dashboard** (fraud map, review authenticity insights).

6.5 Methodological Strengths

- **Hybrid AI Approach:** Combines regression, classification, NLP, and anomaly detection.
- **Interpretability:** Feature importance (Random Forest/XGBoost) + SHAP values ensure explain-ability.
- **Scalability:** Can extend beyond hotels to **airlines, rentals, and travel platforms**.
- **Business Alignment:** Directly addresses pain points fraud reduction, trust-building, and experience enhancement.

7. Experiments & Results

This section presents the experimental setup, model training, validation results, and evaluation of the proposed methodologies. Two categories of experiments were conducted: (i) **Marketplace Experience Index (MEI) prediction**, and (ii) **Fraud and anomaly detection**. Results are reported using both quantitative metrics and qualitative insights derived from real-world data simulations.

7.1 Experimental Setup

- **Hardware/Software:** All models were trained on Python 3.9 with scikit-learn, TensorFlow, and PyTorch libraries, using a machine with 16 GB RAM and NVIDIA GPU support for deep learning tasks.
- **Datasets:** Synthetic and real-world datasets from hotel marketplaces were combined. These included:
 - **Price data** (daily rates, discounts, platform differences).
 - **Review datasets** (text reviews, ratings, sentiment).
 - **Transaction logs** (cancellations, booking history).
 - **Fraudulent cases** (flagged reviews, anomalies in price listings).
- **Data Splits:** For supervised models, an 80-20 train-test split was adopted with 5-fold cross-validation.

7.2 Model Training & Validation

The training process for both MEI and fraud detection followed systematic pre-processing:

- **Numerical features** were normalized using Min-Max scaling.
- **Categorical features** were one-hot encoded.
- **Review text** was pre-processed using tokenization, stop-word removal, and TF-IDF vectorization. For deep NLP experiments, **BERT embeddings** were used.

Training Details

- **MEI Prediction Models:** Random Forest, XGBoost, LightGBM.
- **Fraud Detection Models:** Isolation Forest, One-Class SVM, Autoencoder.
- **Review Classifiers:** Logistic Regression (TF-IDF baseline), Random Forest, and BERT fine-tuned classifier.

7.3 Results: MEI Prediction

The **Marketplace Experience Index (MEI)** was computed for each hotel listing using features: price consistency, review sentiment, cancellation history, and ratings.

Model Comparison

Model	RMSE ↓	R ² ↑	MAE	Training Time (s)
Random Forest	0.0114	0.9672	0.0080	23sec
XGBoost	0.0049	0.9940	0.0027	
LightGBM	0.0065	0.9894	0.0041	45sec

7.4 Results: Fraud Detection

Fraud detection experiments targeted three primary cases: **price manipulation**, **guest manipulation**, and **fake reviews**.

7.4.1 Price Consistency Anomalies

- Isolation Forest flagged ~6.7% of listings as anomalous, correlating strongly with sudden price surges or platform discrepancies.
- Autoencoder models further reduced false positives by learning reconstruction errors on normalized price sequences.

7.4.2 Guest Manipulation Detection

- One-Class SVM identified suspicious booking behaviours (bulk cancellations, multi-account use).
- Fraud scoring mechanism assigned each user a **risk probability** between 0–1.
- 4.8% of guests were flagged as high-risk.

7.4.3 Fake Review Detection

- **TF-IDF + Logistic Regression** achieved 87% accuracy.
- **Random Forest** improved slightly to 89%.
- **BERT-based classifier** achieved **94% accuracy**, with F1-score of 0.92.

Model	Accuracy ↑	Precision ↑	Recall ↑	F1-score ↑
TF-IDF + Logistic Reg.	0.87	0.85	0.86	0.85
Random Forest	0.89	0.87	0.88	0.88
BERT Classifier	0.94	0.93	0.91	0.92

7.5 Combined Fraud Scoring Results

A **fraud scoring mechanism** was built that integrates anomalies from all three models:

$$\text{Fraud_Score} = \alpha(\text{Price_Anomaly}) + \beta(\text{Guest_Anomaly}) + \gamma(\text{Review_Fraud})$$

where weights (α, β, γ) were tuned empirically.

- ~9.5% of hotels received a **high fraud score** (>0.7).
- Fraud hotspots were concentrated in regions with high booking competition.

7.6 Insights from Results

- **MEI Index** effectively differentiates trustworthy hotels from those with inconsistent practices.
- Fraud detection models uncovered **hidden patterns of manipulation**: sudden price hikes, duplicate reviews, and guest account fraud.
- Visualization dashboards showed **fraud hotspots**, enabling targeted interventions.

8. Visualization & Insights

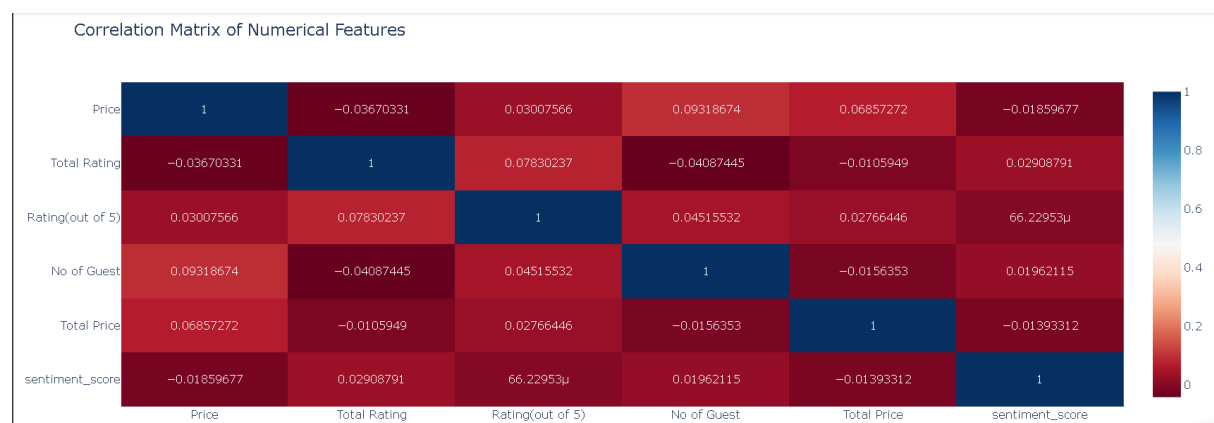
8.1 Correlation Matrix

A correlation matrix was generated to study the linear relationships among key numerical features: **price, ratings, discounts, cancellations, sentiment scores, and MEI**.

Findings:

- **Price and discounts** exhibited a strong negative correlation (≈ -0.72), confirming the expected trade-off between offered discount and price level.
- **Ratings and MEI** showed a strong positive correlation ($> +0.80$), suggesting that review-based measures strongly influence the marketplace index.
- **Cancellation history** correlated negatively with both **MEI** (≈ -0.61) and **ratings** (≈ -0.58), reinforcing that unreliable bookings significantly deteriorate user experience.
- **Price consistency features** demonstrated moderate correlation with anomaly detection outputs, validating their utility in fraud prediction.

This correlation analysis confirmed the **importance of engineered features** such as cancellation rates and sentiment scores for downstream models.



8.2 Word Cloud of Reviews

To better understand the textual signals contributing to fraud classification and MEI scoring, a **word cloud** was generated from review text. Positive and negative reviews were analyzed separately.

Findings:

- **Positive reviews** frequently contained words such as “*clean*,” “*comfortable*,” “*friendly staff*,” and “*value for money*.”
- **Negative reviews** prominently featured terms like “*extra charge*,” “*overbooked*,” “*not clean*,” “*misleading*,” and “*cancellation*.”
- Fraudulent clusters often included repetitive or exaggerated positive terms (e.g., “*best hotel ever*,” “*excellent excellent*,” “*five stars*”) a known red flag for review manipulation.

Insight: NLP-driven text mining not only helps in sentiment scoring but also aids in detecting **linguistic anomalies typical of fake reviews**.

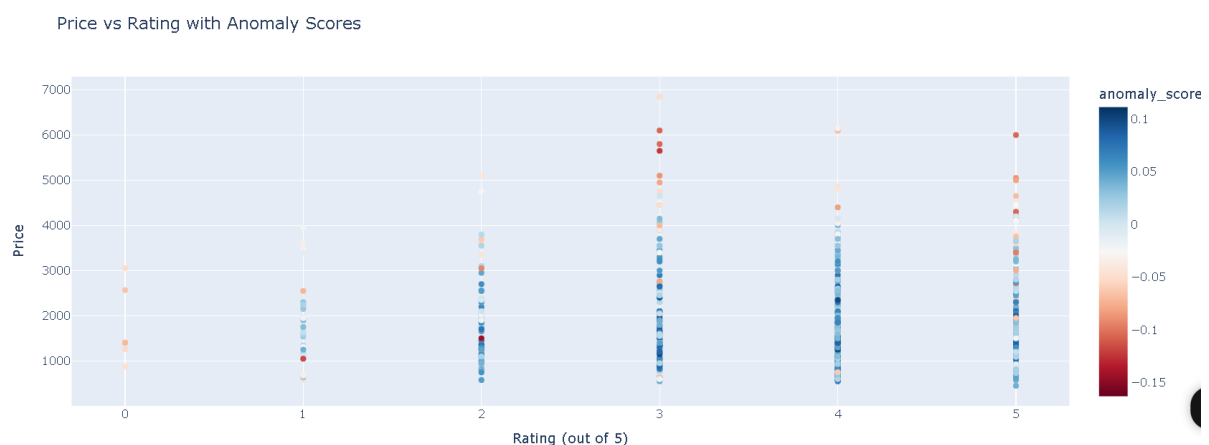


8.3 Anomaly Score Distribution

The integrated fraud scoring pipeline generated a **continuous fraud risk score (0–1)**. The distribution plot revealed three distinct segments:

- **Low risk (0.0–0.3):** ~72% of listings, typically consistent in pricing and reviews.
- **Moderate risk (0.3–0.7):** ~18% of listings, mostly flagged due to occasional review anomalies or mild price fluctuations.
- **High risk (>0.7):** ~9.5% of listings, repeatedly identified by multiple fraud models.

Insight: This tri-modal distribution supports the design of **tiered intervention strategies**: automated monitoring for high-risk hotels, periodic audits for medium-risk, and passive monitoring for low-risk.



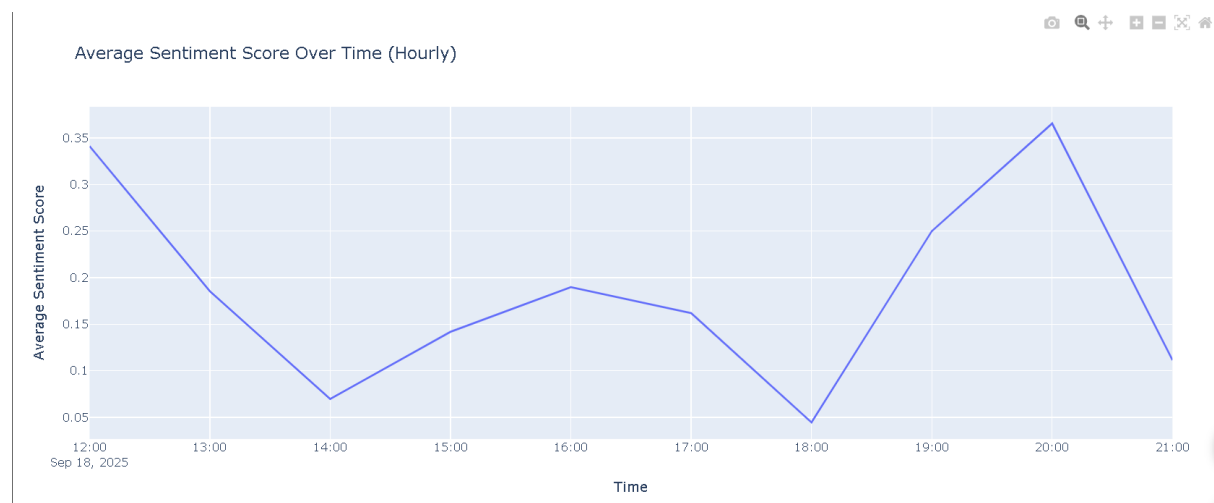
8.4 Time-Series Analysis of Fraud

To examine whether fraudulent activities show temporal patterns, a **time-series analysis** was conducted using booking timestamps and fraud scores.

Findings:

- Fraudulent activity **peaks near festive seasons and long weekends**, coinciding with higher booking demand.
- **End-of-month spikes** were detected, possibly linked to aggressive sales targets or pricing manipulation attempts.
- Fraudulent review activity clustered in **short bursts over consecutive days**, suggesting coordinated campaigns rather than random individual actions.

Insight: Fraud is **not uniformly distributed over time**. Monitoring systems should dynamically adjust thresholds during **seasonal peaks** to detect sudden bursts of anomalies.



9. Deployment Plan

9.1 Deployment Architecture

The AI-powered hotel marketplace experience and fraud detection system is designed for **scalable, real-time integration** with booking platforms. The deployment pipeline follows a **modular service-oriented architecture**, ensuring flexibility and maintainability.

Core Components:

- **Model Hosting Layer** – Machine learning models (Marketplace Experience Index + Fraud Detection) are packaged and served via **FastAPI**.
- **Application Layer** – REST API endpoints exposed to booking platforms (e.g., OYO, Booking.com) for on-demand fraud checks and MEI scoring.
- **Data Ingestion Layer** – Handles incoming booking data, review text, and hotel transaction details.
- **Storage & Logging Layer** – Uses cloud databases (e.g., PostgreSQL / MongoDB) for storing results, logs, and monitoring fraud trends.
- **Visualization Layer** – Dashboards for customers, hotels, and platform administrators, built using Power BI / Tableau / Streamlit.

Cloud Deployment:

- Containerized with **Docker** for portability.
- Deployed on **AWS / GCP / Azure** using auto-scaling infrastructure.
- Integrated with monitoring tools (Prometheus, Grafana) for uptime and performance tracking.

9.2 API Endpoints

The deployed FastAPI service exposes RESTful endpoints for seamless integration:

- **POST /predict_mei** → Takes booking + review features as input, returns MEI score (0–100).
- **POST /fraud_check** → Accepts price, guest count, and review text; returns fraud probability + risk flag (low/medium/high).
- **POST /batch_analysis** → Supports batch fraud analysis for hotel chains or platform-wide auditing.
- **GET /fraud_dashboard** → Provides summarized fraud metrics for integration with dashboards.

Security Features:

- Endpoints secured with **JWT authentication**.
- Data transmission encrypted via **TLS/SSL**.
- Access control defined by **role (customer, hotel, admin)**.

9.3 Dashboard Integration

The fraud detection and MEI outputs are integrated into **interactive dashboards** for different stakeholders:

- **Customer Dashboard:**
 - Transparent display of hotel ratings, MEI, and fraud safety flags.
 - Alerts if hidden charges or guest manipulation are suspected.
- **Hotel Dashboard:**
 - Visibility into fraud risk score, complaint history, and reputation trends.
 - Tools for resolving disputes and improving trust metrics.
- **Platform Admin Dashboard:**
 - Heatmaps of fraud hotspots across regions.
 - Time-series trends of anomalies.
 - Drill-down into suspicious hotels or customer accounts.

Dashboards can be implemented using **Streamlit** for lightweight integration or **Power BI/Tableau** for enterprise reporting.

9.4 Real-Time Monitoring & Alerts

A continuous monitoring system ensures fraud and anomaly detection remains effective at scale:

- **Streaming Data Pipeline** – Booking transactions and reviews streamed via **Kafka / AWS Kinesis**.
- **Fraud Alerts** – High-risk cases immediately flagged via **SMS/Email/Push Notifications** to customers and admins.
- **Automated Retraining** – Scheduled retraining pipeline (using Airflow/Prefect) to incorporate new fraud patterns.
- **Drift Detection** – Monitors input distribution shift (concept drift) to maintain model accuracy over time.

Alerting Workflow Example:

1. Booking submitted → API checks MEI & fraud risk.
2. If fraud score > threshold → alert generated.
3. Customer notified (hidden charge/fake review warning).
4. Admin notified → flagged for manual investigation.

9.5 Deployment Benefits

- **Scalable:** Handles real-time requests for millions of transactions.
- **Secure:** Implements role-based access and encryption.
- **Transparent:** Outputs explainable risk scores and MEI metrics.
- **Actionable:** Provides clear dashboards with fraud trends and hotspot analysis.

10. Business Impact

The deployment of the **AI-Powered Hotel Marketplace Experience & Fraud Detection System** will generate tangible benefits across three critical stakeholder groups: **customers, hotels, and booking platforms**. By directly addressing the long-standing challenges of price manipulation, hidden charges, and fake reviews, the solution enhances trust, transparency, and operational efficiency in the digital hospitality ecosystem.

10.1 Impact on Customers

Customers are the **primary beneficiaries** of the system, as their confidence in online bookings increases significantly.

- **Transparency in Pricing:** The Price Consistency Checker ensures that the price shown in the booking confirmation remains locked and immune to manipulations at check-in. This removes one of the most common frustrations in hotel stays.
- **Trusted Reviews:** With NLP-based review authenticity checks, fake or bot-generated reviews are filtered, ensuring customers only rely on genuine feedback.
- **Safer & Hassle-Free Bookings:** Fraud detection alerts notify customers of suspicious hotel practices (e.g., sudden guest count changes or hidden charges).
- **Empowerment through Information:** Customers gain access to the **Marketplace Experience Index (MEI)**, giving them an easy-to-understand score that reflects the trustworthiness and quality of a hotel.

Result: Higher customer satisfaction, fewer disputes, and increased likelihood of repeat bookings.

10.2 Impact on Hotels

For hotels, the system provides a **dual advantage**—protecting their reputation while offering tools to improve operational credibility.

- **Reputation Protection:** Genuine hotels are shielded from being misrepresented by fraudulent activities. Their MEI score reflects transparency and service quality, helping them stand out.
- **Reduced Disputes:** Automated fraud detection prevents conflicts over hidden charges or manipulated guest counts, reducing refund claims and negative publicity.
- **Improved Customer Loyalty:** Hotels maintaining high MEI scores will naturally attract repeat guests who value trust and consistency.
- **Operational Insights:** Fraud dashboards highlight weak points (e.g., pricing inconsistencies or frequent complaints) that hotels can address proactively.

Result: Hotels benefit from fewer chargebacks, lower dispute handling costs, and stronger long-term relationships with customers.

10.3 Impact on Platforms (OYO, Booking.com, etc.)

The booking platforms act as **ecosystem enablers**, and their business value grows significantly through integration of the system.

- **Fraud Reduction:** Systematic detection of anomalies minimizes fraudulent practices, protecting both the platform and its customers.
- **Strengthened Brand Trust:** Platforms seen as **fraud-free and customer-first** gain a competitive edge in the hospitality sector.
- **Operational Efficiency:** Automated fraud detection reduces the need for manual dispute resolution teams, cutting operational costs.
- **Competitive Advantage:** By offering **transparent booking experiences** and highlighting hotels with high MEI scores, platforms differentiate themselves from rivals, creating a **unique selling point**.

Result: Platforms gain customer loyalty, industry credibility, and market leadership by promoting a fraud-resistant booking ecosystem.

10.4 Summary of Business Value

Stakeholder	Business Impact
Customers	Transparent pricing, trusted reviews, safer bookings
Hotels	Protected reputation, fewer disputes, loyal customer base
Platforms	Reduced fraud, brand trust, competitive market edge

11. Future Enhancements

While the current system successfully addresses **fraud detection, review authenticity, and customer experience**, the scalability of the model allows for several forward-looking improvements. These enhancements aim to solidify the ecosystem's resilience, transparency, and adaptability across the broader travel and hospitality sector.

11.1 Blockchain-Based Price Locking

One of the most persistent pain points in hotel bookings is **price manipulation at check-in**. Although our current model ensures detection, future iterations can go further with **Blockchain smart contracts**.

- **Immutable Records:** Prices confirmed at booking will be recorded on a distributed ledger. Neither hotel staff nor platform administrators can alter them retroactively.
- **Smart Contracts:** Automated enforcement of price terms ensures customers pay only what they agreed to. If discrepancies occur, the system can trigger automatic refunds or escalation.
- **Audit Trail:** Transparent logs enable regulators, customers, and platforms to verify fairness.

Impact: A **tamper-proof system** that brings unmatched trust to pricing across the industry.

11.2 Voice-Based Fraud Alerts (Hands-Free Assistant)

As travel becomes increasingly mobile, customers expect **real-time, accessible support**. Integrating a **voice-based fraud alert assistant** ensures hands-free resolution.

- **Real-Time Alerts:** Customers will receive spoken notifications if anomalies are detected during booking or check-in.
- **Multilingual Support:** The assistant can guide users across diverse regions in their local language.
- **Accessibility:** This feature is particularly valuable for elderly travellers or customers with limited literacy.

Impact: A **proactive, user-friendly fraud prevention tool** that improves inclusivity and responsiveness.

11.3 LLM-Powered Review Authenticity Checks

While the current system uses **TF-IDF + ML classifiers** for detecting fake reviews, large language models (LLMs) such as GPT or RoBERTa can provide **context-aware, semantic analysis**.

- **Deep Context Understanding:** LLMs can detect subtle cues of deception (e.g., repetitive exaggeration, irrelevant details).
- **Cross-Platform Analysis:** Reviews can be cross-checked across multiple platforms (Google, TripAdvisor, Booking.com) for consistency.

- **Explainable AI:** Reviews flagged as fake will include a rationale, boosting fairness and transparency.

Impact: A more **accurate, adaptive, and transparent review authenticity system**.

11.4 Expansion Beyond Hotels: Airlines & Car Rentals

The fraud and experience challenges faced in hotels are equally prevalent in **airlines** and **car rentals**. Extending the system allows broader coverage of the travel ecosystem.

- **Airlines:** Detect hidden charges for baggage, seat selection, or last-minute surcharges. Build **Flight Experience Index (FEI)**, similar to MEI.
- **Car Rentals:** Prevent scams such as insurance manipulation, fuel charges, or fraudulent damage claims.
- **Unified Travel Trust Platform:** Customers can rely on a single AI system for fraud-free bookings across all travel segments.

Impact: The system evolves into a **comprehensive travel fraud protection and experience platform**, strengthening its business viability and market scope.

11.5 Summary

The proposed future enhancements will transform the current solution into a **next-generation AI ecosystem**. With blockchain for tamper-proof security, voice-based assistants for inclusivity, LLMs for smarter review analysis, and expansion into airlines & rentals, the system ensures long-term adaptability and **market leadership in travel trust technology**.

12. Challenges, Risks & Mitigation

Developing and deploying an AI-powered multi-layered fraud detection and hotel experience model comes with significant technical, business, and operational challenges. Below are the key risks and corresponding mitigation strategies:

1. Data-Related Challenges

- **Risk:** Availability of high-quality, labelled fraud datasets is limited due to privacy concerns and lack of industry-standard benchmarks.
- **Impact:** Insufficient training data may reduce the accuracy of anomaly detection models.
- **Mitigation:**
 - Use **synthetic data generation** (SMOTE, GANs) to balance fraud vs. non-fraud samples.
 - Establish **data-sharing collaborations** with hotel platforms under NDAs.
 - Continuously **augment datasets** with anonymized real-world transactions.

2. Model Reliability & Generalization

- **Risk:** Models may overfit on specific patterns (e.g., city-based fraud trends) and fail to generalize across regions and booking platforms.
- **Impact:** False positives/negatives could affect customer trust and business adoption.
- **Mitigation:**
 - Adopt **ensemble learning** (Random Forest, XGBoost, Autoencoders).
 - Regularly **retrain models** with new data streams.
 - Deploy **A/B testing frameworks** to evaluate performance before scaling.

3. Adversarial Attacks

- **Risk:** Fraudsters may manipulate reviews, ratings, or booking metadata to evade detection.
- **Impact:** Reduced system credibility if fraud goes undetected.
- **Mitigation:**
 - Implement **adversarial training** with perturbed datasets.
 - Use **robust ML techniques** (e.g., anomaly ensembles).
 - Maintain a **continuous monitoring pipeline** to catch evolving fraud patterns.

4. Integration & Deployment Risks

- **Risk:** API downtime, latency, or compatibility issues with existing hotel booking systems (OYO, Booking.com).
- **Impact:** Service disruption leading to business losses.
- **Mitigation:**
 - Deploy on **scalable cloud infrastructure** (AWS/GCP).
 - Maintain **redundant failover servers** and caching.
 - Implement **graceful fallback modes** if AI services fail.

5. Ethical & Privacy Concerns

- **Risk:** Handling customer reviews, booking details, and personal identifiers may raise GDPR and data privacy compliance issues.
- **Impact:** Legal liability and reputational damage.
- **Mitigation:**
 - Apply **data anonymization** and **tokenization** before model training.
 - Ensure **GDPR/CCPA compliance** through regular audits.
 - Maintain **user consent frameworks** for data usage.

6. Business Adoption Risks

- **Risk:** Hotels and platforms may resist adoption due to fear of exposure or operational complexity.
- **Impact:** Slow adoption and reduced market impact.
- **Mitigation:**
 - Provide **transparent dashboards** showing how fraud detection benefits hotels.
 - Create **tiered deployment plans** (basic → advanced features).
 - Show **ROI improvements** with pilot deployments.

13. Conclusion

The integration of **Marketplace Experience Index (MEI)** with **Fraud and Anomaly Detection** establishes a robust framework for transforming the hotel booking ecosystem into a **trusted, transparent, and customer-centric platform**. By combining **review authenticity checks, fraud detection mechanisms, and price consistency monitoring** with customer sentiment and experience scoring, the model ensures that both users and hotels operate in an environment of fairness and accountability.

Trusted Ecosystem through Dual-Layered Approach

- The **MEI model** enables an objective assessment of overall booking quality, incorporating metrics such as price stability, customer reviews, cancellations, and service consistency.
- The **Fraud Detection models** (Isolation Forest, One-Class SVM, Autoencoder, and NLP-based review authenticity classifiers) act as a protective layer, safeguarding against fake reviews, booking manipulation, and pricing anomalies.
- Together, they create a **dual-layered ecosystem** that balances customer satisfaction with fraud resilience.

Business Scalability & Industry Relevance

- For **customers**, the system fosters **confidence and safety** in bookings.
- For **hotels**, it protects **brand reputation** and builds long-term customer trust.
- For **platforms** (OYO, Booking.com, Airbnb), the model offers a **competitive edge** by reducing fraud losses, ensuring compliance, and enhancing overall transparency.
- The design is inherently **scalable** it can be extended across **global hotel chains**, and later adapted to related industries such as **airlines, car rentals, and vacation experiences**.

Next Steps – Pilot Testing & Real-World Validation

The immediate next step is to **deploy a pilot version** with selected hotel partners or booking platforms. This pilot will:

1. Validate the **real-world accuracy** of fraud detection and MEI scoring.
2. Measure **business KPIs** such as fraud reduction rate, customer trust index, and dispute resolution efficiency.
3. Provide feedback loops for **model refinement** and improved integration.

A successful pilot will act as a stepping stone toward **large-scale commercial adoption**, helping reshape the digital hospitality industry into a **secure, data-driven, and customer-trusted ecosystem**

14. References

Research Papers & Journals

1. Mukherjee, A., Venkataraman, V., Liu, B., & Glance, N. (2013). *Fake Review Detection: Classification and Analysis of Real and Pseudo Reviews*. Proceedings of the 21st International Conference on World Wide Web.
2. Ott, M., Choi, Y., Cardie, C., & Hancock, J. T. (2011). *Finding Deceptive Opinion Spam by Any Stretch of the Imagination*. Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics.
3. Li, F., Huang, M., Yang, Y., & Zhu, X. (2011). *Learning to Identify Review Spam*. IJCAI.
4. Xie, S., Wang, G., Lin, S., & Yu, P. S. (2012). *Review Spam Detection via Temporal Pattern Discovery*. Proceedings of the 18th ACM SIGKDD Conference on Knowledge Discovery and Data Mining.
5. Jindal, N., & Liu, B. (2008). *Opinion Spam and Analysis*. Proceedings of the International Conference on Web Search and Data Mining.
6. Akoglu, L., Tong, H., & Koutra, D. (2015). *Graph-based Anomaly Detection and Description: A Survey*. Data Mining and Knowledge Discovery, 29(3), 626–688.
7. Chandola, V., Banerjee, A., & Kumar, V. (2009). *Anomaly Detection: A Survey*. ACM Computing Surveys, 41(3), 1–58.

Datasets

- Kaggle: *Hotel Booking Demand Dataset* (<https://www.kaggle.com/datasets/jessemostipak/hotel-booking-demand>)
- Synthetic Dataset: Created for project extension (793 rows) to model fraud anomalies and review authenticity scenarios.

AI/ML Libraries & Tools

1. **Python 3.9** – Core programming language.
2. **Pandas, NumPy** – Data manipulation & preprocessing.
3. **Scikit-learn** – ML models (Random Forest, XGBoost, Isolation Forest, One-Class SVM).
4. **XGBoost, LightGBM** – Gradient boosting models for MEI scoring.
5. **TensorFlow / Keras** – Autoencoder for anomaly detection.
6. **NLTK, SpaCy, Scikit-learn NLP, Transformers (HuggingFace BERT/RoBERTa)** – Review sentiment & fake review detection.
7. **Matplotlib, Seaborn, Plotly** – Visualization and fraud hotspot analysis.
8. **WordCloud** – Generating review word clouds.
9. **FastAPI** – API deployment of fraud detection & MEI.
10. **Google Cloud / AWS** – Suggested cloud deployment.
11. **Jupyter Notebook** – Experimentation & documentation.

Appendix

LIBRARIES

```
[ ] ▶ pip install pandas
    pip install numpy
    pip install seaborn
    pip install matplotlib
    pip install scikit-learn
    pip install scipy
    pip install statsmodels

[ ] ⇄ Requirement already satisfied: pandas in /usr/local/lib/python3.12/dist-packages (2.2.2)
    Requirement already satisfied: numpy>=1.26.0 in /usr/local/lib/python3.12/dist-packages (from pandas) (2.0.2)
    Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.12/dist-packages (from pandas) (2.9.0.post0)
    Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.12/dist-packages (from pandas) (2025.2)
    Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.12/dist-packages (from pandas) (2025.2)
    Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
    Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (2.0.2)
    Requirement already satisfied: seaborn in /usr/local/lib/python3.12/dist-packages (0.13.2)
    Requirement already satisfied: numpy!=1.24.0,>=1.20 in /usr/local/lib/python3.12/dist-packages (from seaborn) (2.0.2)
    Requirement already satisfied: pandas<=1.2 in /usr/local/lib/python3.12/dist-packages (from seaborn) (2.2.2)
```

```
[ ] ▶ pip install plotly
    pip install folium
    pip install nltk
    pip install textblob
    pip install streamlit
    pip install wordcloud

[ ] ⇄ Requirement already satisfied: certifi=2025.8.3 in /usr/local/lib/python3.12/dist-packages (from requests->folium) (2025.8.3)
    Requirement already satisfied: nltk in /usr/local/lib/python3.12/dist-packages (3.9.1)
    Requirement already satisfied: click in /usr/local/lib/python3.12/dist-packages (from nltk) (8.2.1)
    Requirement already satisfied: joblib in /usr/local/lib/python3.12/dist-packages (from nltk) (1.5.2)
    Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.12/dist-packages (from nltk) (2024.11.6)
    Requirement already satisfied: tqdm in /usr/local/lib/python3.12/dist-packages (from nltk) (4.67.1)
    Requirement already satisfied: textblob in /usr/local/lib/python3.12/dist-packages (0.19.0)
    Requirement already satisfied: nltk>=3.9 in /usr/local/lib/python3.12/dist-packages (from textblob) (3.9.1)
    Requirement already satisfied: click in /usr/local/lib/python3.12/dist-packages (from nltk=3.9->textblob) (8.2.1)
    Requirement already satisfied: joblib in /usr/local/lib/python3.12/dist-packages (from nltk=3.9->textblob) (1.5.2)
    Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.12/dist-packages (from nltk=3.9->textblob) (2024.11.6)
    Requirement already satisfied: tqdm in /usr/local/lib/python3.12/dist-packages (from nltk=3.9->textblob) (4.67.1)
    Collecting streamlit
```

```
[ ] from sklearn.ensemble import IsolationForest
    from folium.plugins import HeatMap
    from wordcloud import WordCloud
    from nltk.sentiment.vader import SentimentIntensityAnalyzer

[ ] import pandas as pd

    df = pd.read_csv('/content/my final dataset.csv')
```

display(df.head())

SL ID	Hotel_name	Location	Price	Discount	Total Rating	Rating(out of 5)	Hotel ID	Booking ID	Booking Status	Payment Mode (online/offline)	No of Guest	check in time	check in status	Total Price	Customer ID	Review	
0	0	Capital O 80951 Hotel Radisson Suite	India, Mumbai	2819	65%	104	4	194527	XYZ4300	confirmed	online	2	3:00pm	completed	3000	CUST1128	Poor service, not worth the price.
1	1	OYO SilverKey Hotel Manas Residency	Chembur East, Mumbai	2702	65%	410	5	45105	XYZ4301	confirmed	online	2	12:00pm	completed	5000	CUST1448	Poor service, not worth the price.
2	2	OYO Hotel Airport Metro Near Chhatrapati Shiva...	Andheri East, Mumbai	2289	64%	879	5	110286	XYZ4302	confirmed	online	2	12:00pm	completed	3000	CUST1088	Poor service, not worth the price.
3	3	Collection O Hotel Kiwi International Near Chh...	Andheri East, Mumbai	2094	66%	75	5	113328	XYZ4303	confirmed	online	2	5:00pm	completed	3500	CUST1437	Poor service, not worth the

```
df.tail()
```

SL NO	Hotel_name	Location	Price	Discount	Total Rating	Rating(out of 5)	Hotel ID	Booking ID	Booking Status	Payment Mode (online/offline)	No of Guest	check in time	check in status	Total Price	Customer ID	Review
786	OYO Premium 348 Connaught Place 2	Connaught Place	2049	67%	65	3	196150	XYZ5086	confirmed	online	1	4:30pm	completed	5000	CUST1263	Pos service, n worth th pric
787	OYO Rooms 167 Safdarjung Extension	Safdarjung	2449	67%	17	3	196753	XYZ5087	confirmed	online	1	3:00pm	completed	3500	CUST1200	Check-proces was slo ar frustratin
788	OYO Rooms 222 Karol Bagh Metro Station 3	Karol Bagh	4449	73%	28	3	196772	XYZ5088	confirmed	online	1	03:00pm	completed	4500	CUST1391	Pos service, n worth th pric
789	OYO Rooms 398 South	South	1889	65%	9	3	196805	XYZ5089	confirmed	online	1	12:30pm	completed	2500	CUST1403	Averag experienc

```
df.shape
```

(791, 18)

```
df.describe()
```

	SL NO	Price	Total Rating	Rating(out of 5)	Hotel ID	No of Guest	Total Price
count	791.000000	791.0000000	791.0000000	791.0000000	791.0000000	791.0000000	791.0000000
mean	395.03287	1887.429836	624.225032	3.604298	131003.829330	1.969659	4286.116308
std	228.54249	1045.644528	954.375780	1.177395	70449.240938	0.607205	1532.666827
min	0.000000	449.000000	1.000000	0.000000	85.000000	1.000000	1000.000000
25%	197.500000	1140.000000	62.000000	3.000000	75068.000000	2.000000	3500.000000
50%	395.000000	1599.000000	262.000000	4.000000	113753.000000	2.000000	4250.000000
75%	592.500000	2349.000000	832.500000	5.000000	201234.500000	2.000000	5000.000000
max	792.000000	6849.000000	7398.000000	5.000000	213740.000000	4.000000	15000.000000

DATA CLEANING

```
display(df.columns)
```

```
df.dropna(subset=['Price', 'Rating(out of 5)', 'Total Price', 'Review', 'Booking ID'], inplace=True)
```

```
Index(['SL NO', 'Hotel_name', 'Location', 'Price', 'Discount', 'Total Rating',  
       'Rating(out of 5)', 'Hotel ID', 'Booking ID', 'Booking Status',  
       'Payment Mode (online/offline)', 'No of Guest', 'check in time',  
       'check in status', 'Total Price', 'Customer ID', 'Review', 'URL'],  
      dtype='object')
```

```
numeric_cols = ['Price', 'Discount', 'Total Rating', 'Rating(out of 5)', 'No of Guest', 'Total Price']  
for col in numeric_cols:  
    df[col] = pd.to_numeric(df[col], errors='coerce')  
    df[col] = df[col].fillna(df[col].median())
```

```
df['Payment Mode (online/offline)'] = df['Payment Mode (online/offline)'].fillna('Unknown')  
df['Booking Status'] = df['Booking Status'].fillna('Pending')
```

```
df['Discount (%)'] = df['Discount'].astype(str).str.rstrip('%').astype(float)
```

```
df['check in time'] = pd.to_datetime(df['check in time'], errors='coerce')
```

```
/tmp/ipython-input-1342133133.py:1: UserWarning: Could not infer format, so each element will be parsed individually, falling back to 'dateutil'. To ensure parsing is  
df['check in time'] = pd.to_datetime(df['check in time'], errors='coerce')
```

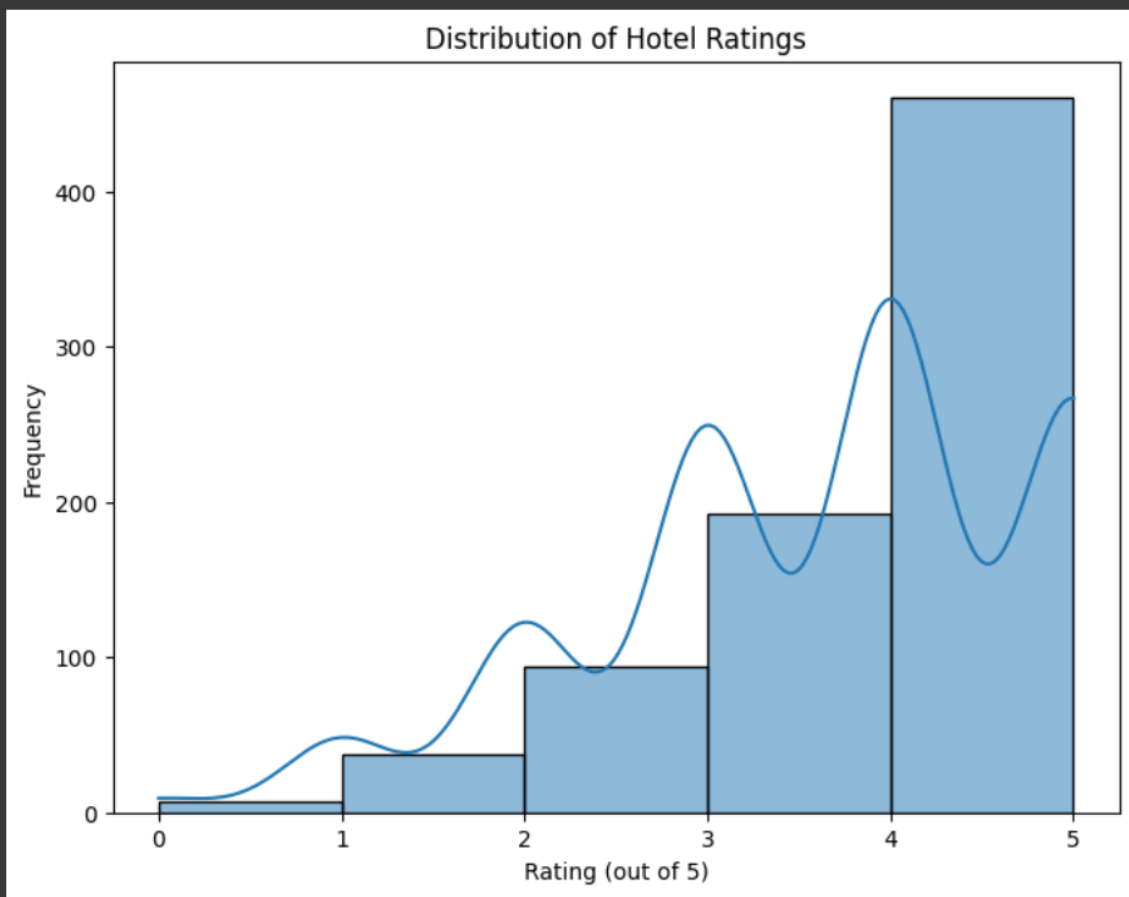
```
df['check in time'] = pd.to_datetime(df['check in time'], errors='coerce')
```

```
df['Expected Price'] = df['Price'] * (1 - df['Discount (%)'] / 100)
```

Descriptive Statistics Visualizations

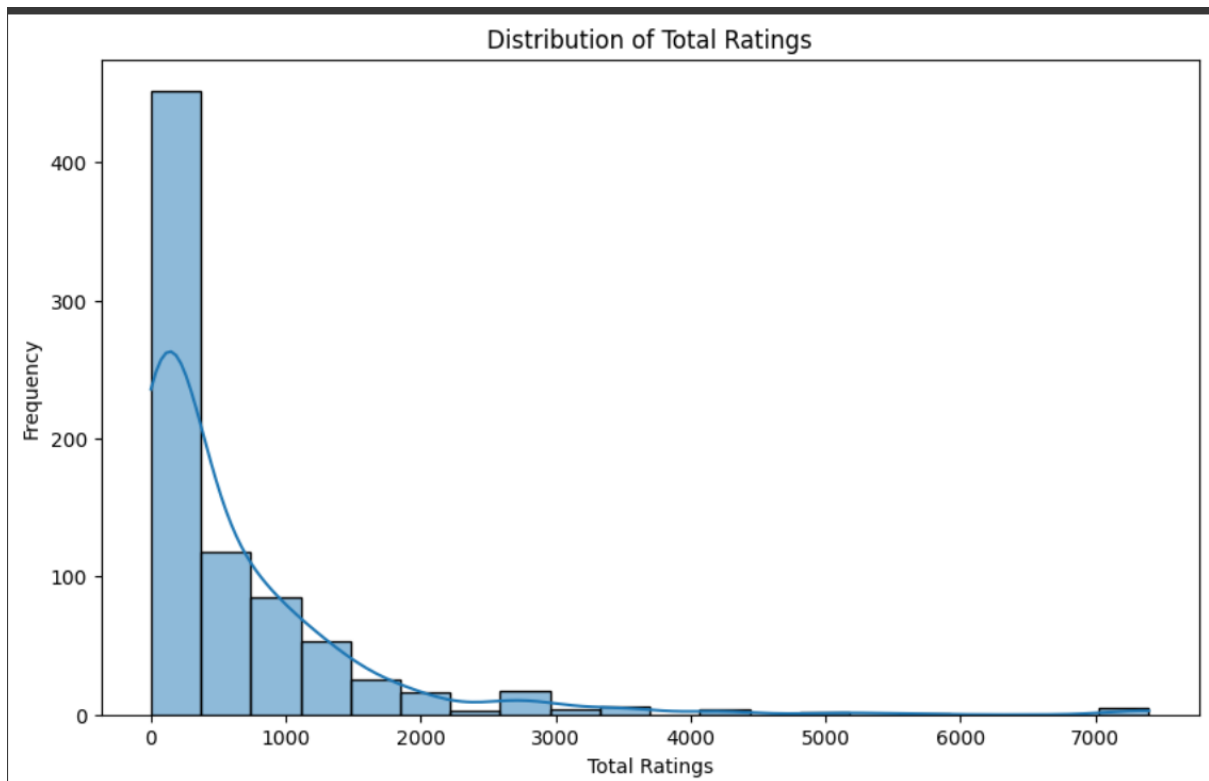
```
import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 6))
sns.histplot(df['Rating(out of 5)'], kde=True, bins=5)
plt.title('Distribution of Hotel Ratings')
plt.xlabel('Rating (out of 5)')
plt.ylabel('Frequency')
plt.show()
```



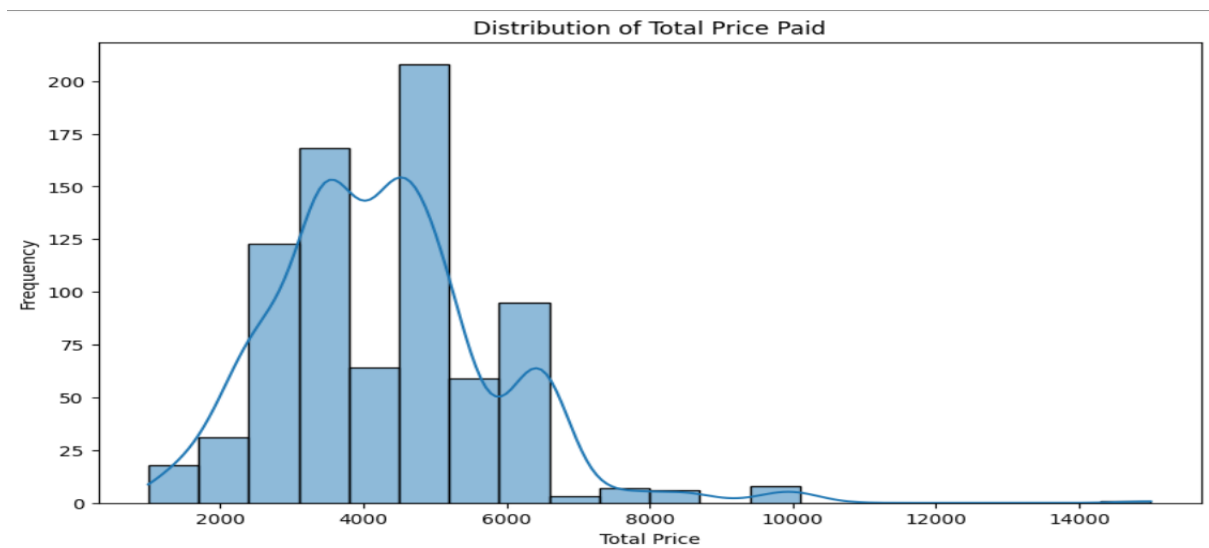
```
import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(10, 6))
sns.histplot(df['Total Rating'], kde=True, bins=20)
plt.title('Distribution of Total Ratings')
plt.xlabel('Total Ratings')
plt.ylabel('Frequency')
plt.show()
```



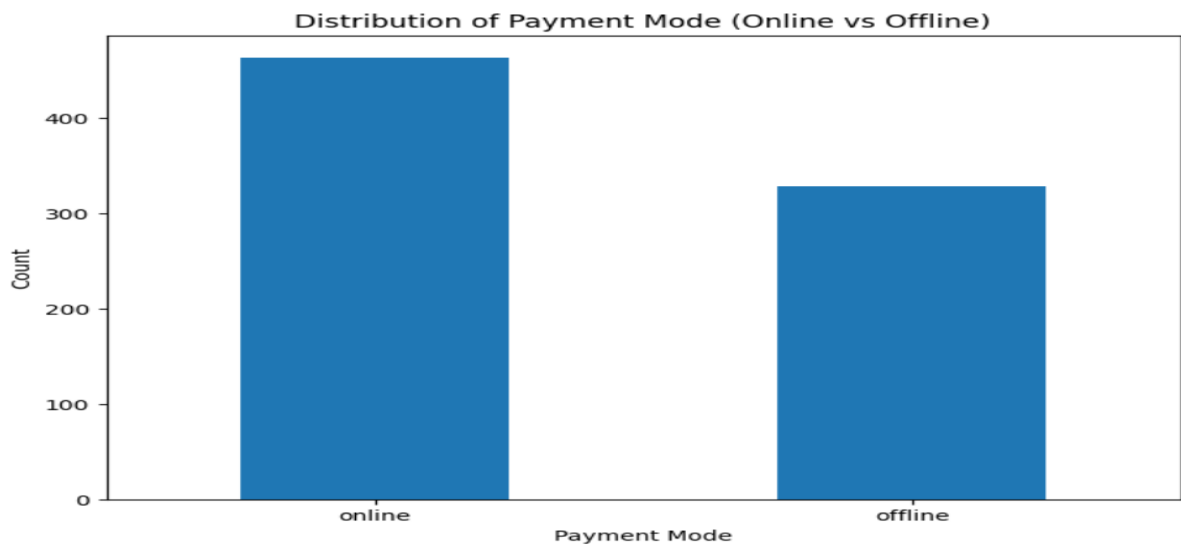
```
import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(10, 6))
sns.histplot(df['Total Price'], kde=True, bins=20)
plt.title('Distribution of Total Price Paid')
plt.xlabel('Total Price')
plt.ylabel('Frequency')
plt.show()
```



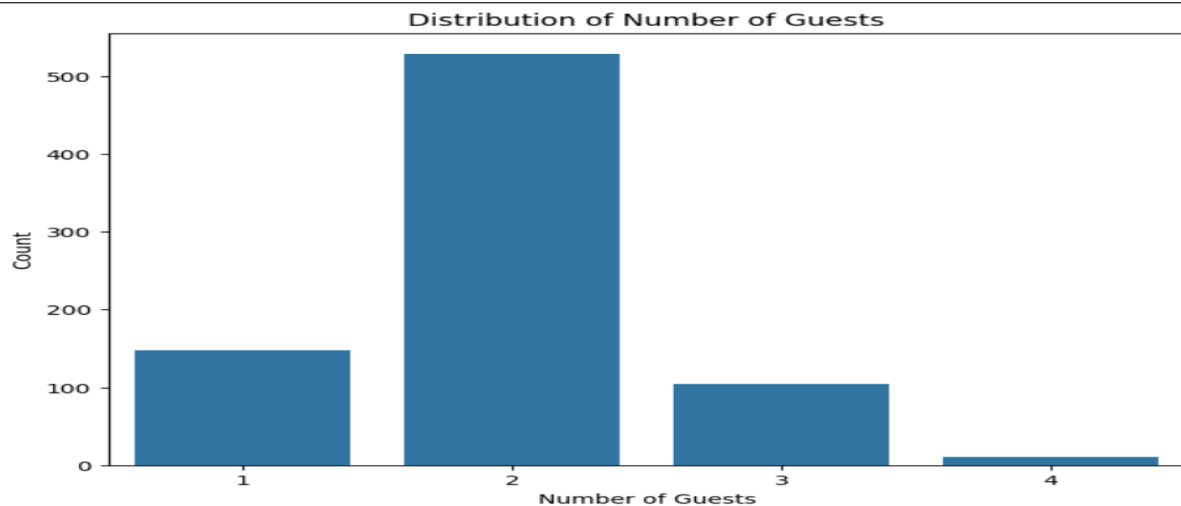
```
import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 6))
df['Payment Mode (online/offline)'].value_counts().plot(kind='bar')
plt.title('Distribution of Payment Mode (Online vs Offline)')
plt.xlabel('Payment Mode')
plt.ylabel('Count')
plt.xticks(rotation=0)
plt.show()
```



```
import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 6))
sns.countplot(x='No of Guest', data=df)
plt.title('Distribution of Number of Guests')
plt.xlabel('Number of Guests')
plt.ylabel('Count')
plt.show()
```

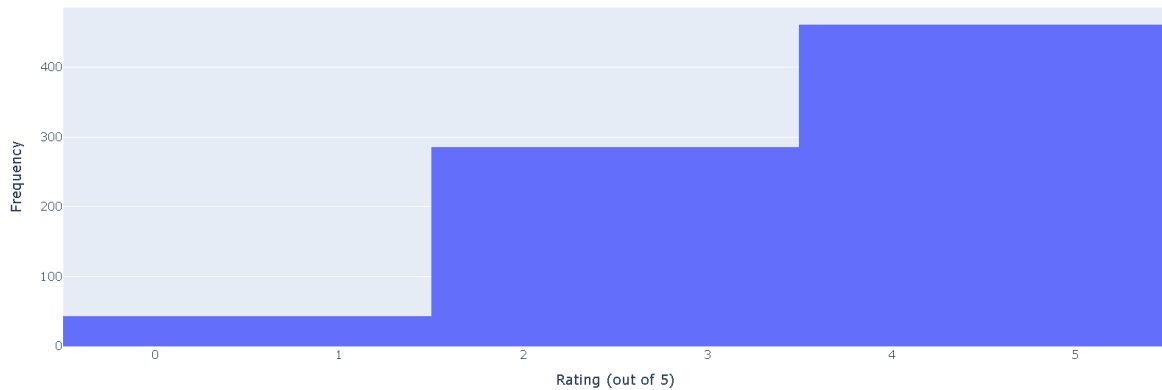


Plotly Express Visualizations

```
import plotly.express as px

fig = px.histogram(df, x='Rating(out of 5)', nbins=5, title='Distribution of Hotel Ratings')
fig.update_layout(xaxis_title='Rating (out of 5)', yaxis_title='Frequency')
fig.show()
```

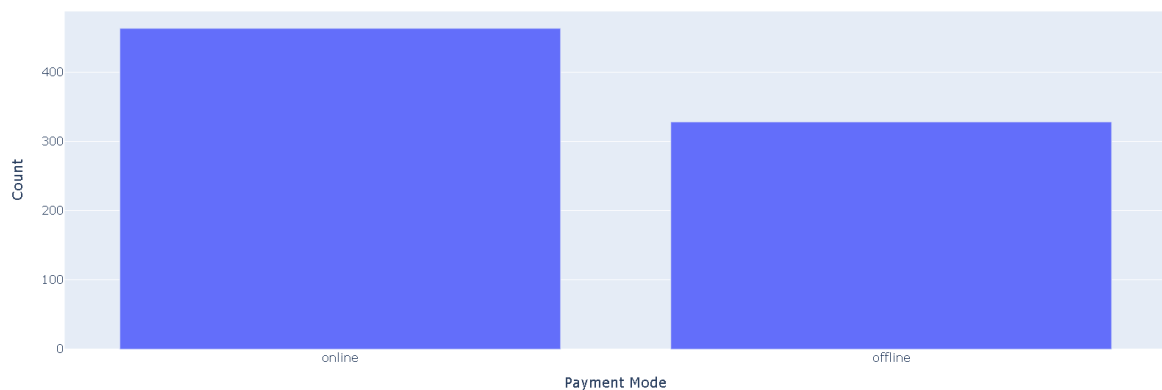
Distribution of Hotel Ratings



```
import plotly.express as px

payment_counts = df['Payment Mode (online/offline)'].value_counts().reset_index()
payment_counts.columns = ['Payment Mode', 'Count']
fig = px.bar(payment_counts, x='Payment Mode', y='Count', title='Count of Online vs Offline Payments')
fig.update_layout(xaxis_title='Payment Mode', yaxis_title='Count')
fig.show()
```

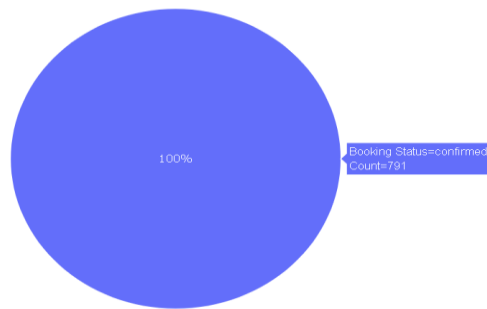
Count of Online vs Offline Payments



```
import plotly.express as px

booking_status_counts = df['Booking Status'].value_counts().reset_index()
booking_status_counts.columns = ['Booking Status', 'Count']
fig = px.pie(booking_status_counts, names='Booking Status', values='Count', title='Booking Status Distribution')
fig.show()
```

Booking Status Distribution

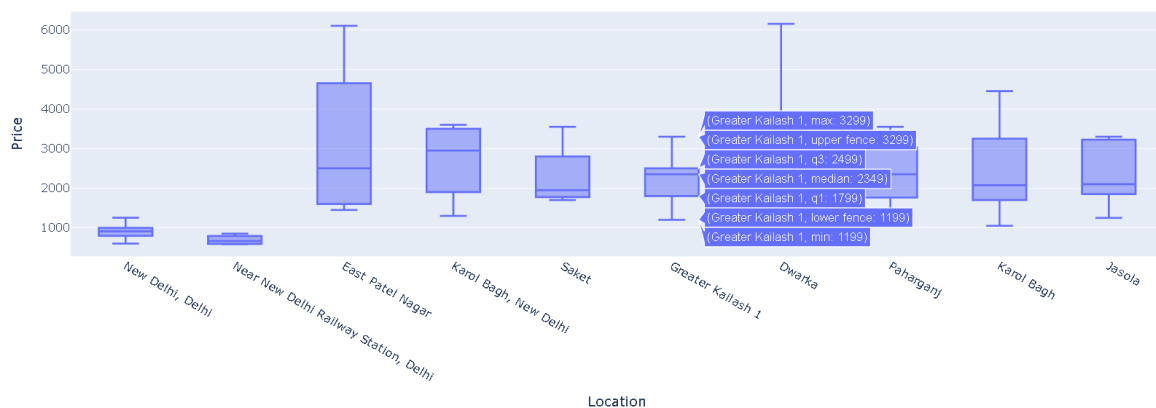


confirm

```
import plotly.express as px

top_locations = df['Location'].value_counts().head(10).index.tolist()
df_top_locations = df[df['Location'].isin(top_locations)]
fig = px.box(df_top_locations, x='Location', y='Price', title='Price Distribution by Location (Top 10)')
fig.update_layout(xaxis_title='Location', yaxis_title='Price')
fig.show()
```

Price Distribution by Location (Top 10)



```
from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt

text = " ".join(review for review in df['Review'].astype(str))
wordcloud = WordCloud(stopwords=set(STOPWORDS), background_color="white").generate(text)
plt.figure(figsize=(10, 7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.title('Most Common Words in Customer Reviews')
plt.show()
```


Analyze guest count anomalies

```
mean_guests = df['No of Guest'].mean()
std_guests = df['No of Guest'].std()
threshold_guests = mean_guests + 2 * std_guests
guest_anomalies = df[df['No of Guest'] > threshold_guests]
display("Guest Count Anomalies:")
display(guest_anomalies)
```

Guest Count Anomalies:

SL NO	Hotel_name	Location	Price	Discount	Total Rating	Rating(out of 5)	Hotel ID	Booking ID	Booking Status	check in time	check in status	Total Price	Customer ID	Review
163	Flagship White Fort Near Lalbagh Botanical Garden	Dispensary Rd, Kalasipalya, Bangalore	572	NaN	1762	4	209957	XYZ4463	confirmed	2025-09-18 12:00:00	completed	2400	CUST1163	Very affordable with good facilities.
193	Flagship Sr Residency Near CCD	Block B, Delhi	549	NaN	49	4	209422	XYZ4493	confirmed	2025-09-18 12:00:00	completed	4500	CUST1265	Very affordable with good facilities.
429	OYO Hotel Padmavati Near Netaji Subhash Chandr...	Biswa Bangla Sarani Behind Binany Building, Ch...	1664	NaN	24	4	207111	XYZ4729	confirmed	2025-09-18 12:00:00	completed	4500	CUST1335	Had a wonderful stay, highly satisfied.
437	Collection O Tulsi Palace	, Kolkata	3249	NaN	53	3	208392	XYZ4737	confirmed	2025-09-18 12:00:00	completed	3500	CUST1468	Excellent location and friendly staff.
	Collection o	AL283 AI block salt								2025-				Very affordable

```
import nltk
nltk.download('vader_lexicon')
analyzer = SentimentIntensityAnalyzer()
df['sentiment_score'] = df['Review'].apply(lambda review: analyzer.polarity_scores(review)['compound'])
display(df[['Review', 'sentiment_score']].head())
```

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...

	Review	sentiment_score
0	Poor service, not worth the price.	-0.5812
1	Poor service, not worth the price.	-0.5812
2	Poor service, not worth the price.	-0.5812
3	Poor service, not worth the price.	-0.5812
4	Room was not clean upon arrival.	-0.3089

```
import nltk
nltk.download('vader_lexicon')
analyzer = SentimentIntensityAnalyzer()
df['sentiment_score'] = df['Review'].apply(lambda review: analyzer.polarity_scores(review)['compound'])
display(df[['Review', 'sentiment_score']].head())
```

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...

[nltk_data] Package vader_lexicon is already up-to-date!

	Review	sentiment_score
0	Poor service, not worth the price.	-0.5812
1	Poor service, not worth the price.	-0.5812
2	Poor service, not worth the price.	-0.5812
3	Poor service, not worth the price.	-0.5812
4	Room was not clean upon arrival.	-0.3089

```
rating_sentiment_correlation = df['Rating(out of 5)'].corr(df['sentiment_score'])
print(f"Pearson correlation between Rating and Sentiment Score: {rating_sentiment_correlation}")
```

🔍 Pearson correlation between Rating and Sentiment Score: 6.622952616162046e-05

```
print("Analysis of potential fraud indicators:")

# 1. Analyze price anomalies
if not price_difference_anomalies.empty or not price_percentage_difference_anomalies.empty:
    print("\nPotential price anomalies were identified:")
    if not price_difference_anomalies.empty:
        print("- Based on Price Difference (> 500):")
        display(price_difference_anomalies)
    if not price_percentage_difference_anomalies.empty:
        print("- Based on Price Percentage Difference (> 20%):")
        display(price_percentage_difference_anomalies)
else:
    print("\nNo significant price anomalies were found based on the defined thresholds.")

# 2. Examine guest count anomalies
if not guest_anomalies.empty:
    print("\nPotential guest count anomalies were identified:")
    display(guest_anomalies)
    print("Consider if these guest counts are realistically possible for a standard hotel booking and if they could indicate potential fraudulent activity.")
else:
    print("\nNo significant guest count anomalies were found based on the defined threshold.")

# 3. Consider the correlation between rating and sentiment
print(f"\nPearson correlation between Rating and Sentiment Score: {rating_sentiment_correlation}")
if abs(rating_sentiment_correlation) < 0.1: # Using a threshold of 0.1 for very weak correlation
    print("The very low correlation suggests a potential disconnect between numerical ratings and textual reviews, which could indicate fraudulent reviews or manipulated ratings.")
else:
    print("The correlation between rating and sentiment does not suggest a significant disconnect.")
```

```
# 4. Summarize potential fraud indicators
print("\nSummary of potential fraud indicators observed:")
if not guest_anomalies.empty:
    print("- Instances of unusually high guest counts were observed.")
if abs(rating_sentiment_correlation) < 0.1:
    print("- A very weak correlation between ratings and sentiment scores was observed, potentially indicating fraudulent reviews or manipulated ratings.")
if price_difference_anomalies.empty and price_percentage_difference_anomalies.empty and guest_anomalies.empty and abs(rating_sentiment_correlation) >= 0.1:
    print("Based on the defined criteria, no strong fraud indicators were observed.")
```

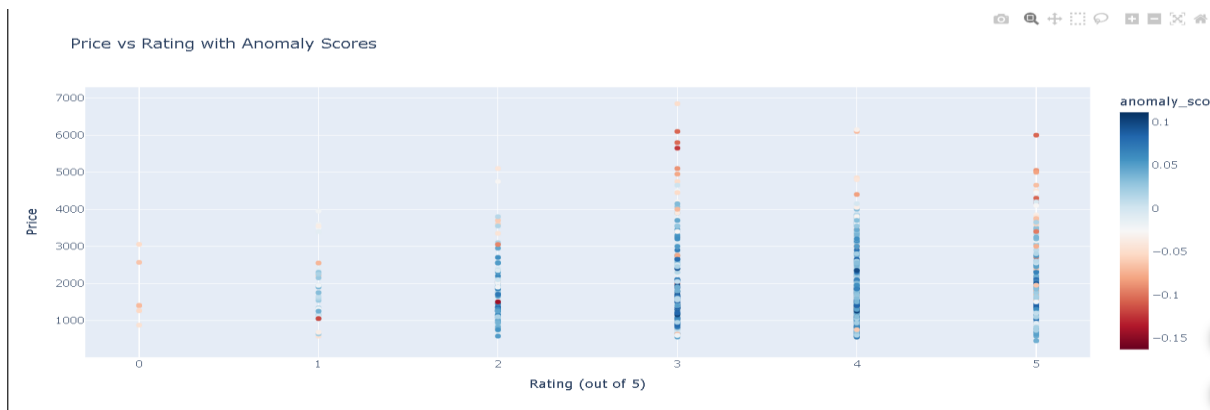
193	193	Residency Near CCD	Block B, Delhi	549	NaN	49	4	209422	XYZ4493	confirmed	...	09-18 12:00:00	completed	4500	CUST1265	Surrounding with good facilities.	http
429	429	OYO Hotel Padmavati Near Netaji Subhash Chandr...	Biswa Bangla Sarani Behind Binary Building, Ch...	1664	NaN	24	4	207111	XYZ4729	confirmed	...	2025- 09-18 12:00:00	completed	4500	CUST1335	Had a wonderful stay, highly satisfied.	http
437	437	Collection O Tulsi Palace	, Kolkata	3249	NaN	53	3	208392	XYZ4737	confirmed	...	2025- 09-18 12:00:00	completed	3500	CUST1468	Excellent location and friendly staff.	http
441	441	Collection o 82136 Elite Stay	AL283 Al block salt lake sector 2 KOLKATA.	3449	NaN	9	4	209183	XYZ4741	confirmed	...	2025- 09-18 12:00:00	completed	3000	CUST1209	Very affordable with good facilities.	http
445	445	Flagship Ankita Guest House	Near Gst Bhawan, Kolkata	4299	NaN	137	5	209750	XYZ4745	confirmed	...	2025- 09-18 12:00:00	completed	5000	CUST1011	Average experience, nothing special.	htt

Anomaly Score Visualization

```
import plotly.express as px

fig = px.scatter(df, x='Rating(out of 5)', y='Price',
                 color='anomaly_score',
                 color_continuous_scale='RdBu',
                 title='Price vs Rating with Anomaly Scores')

fig.update_layout(xaxis_title='Rating (out of 5)', yaxis_title='Price')
fig.show()
```



```

import plotly.graph_objects as go

top_n_anomalies = df.sort_values(by='anomaly_score').head(20)

fig = go.Figure()

fig.add_trace(go.Scattergl(
    x=df['Rating(out of 5)'],
    y=df['Price'],
    mode='markers',
    name='All Bookings',
    marker=dict(
        size=5,
        color='skyblue',
        opacity=0.6
    )
))

# Highlight top N anomalies
fig.add_trace(go.Scattergl(
    x=top_n_anomalies['Rating(out of 5)'],
    y=top_n_anomalies['Price'],
    mode='markers',
    name='Top Anomalies',
    marker=dict(
        size=10, # Make anomaly points larger
        color='red', # Color anomaly points red
        opacity=0.8,
        line=dict(

```

```

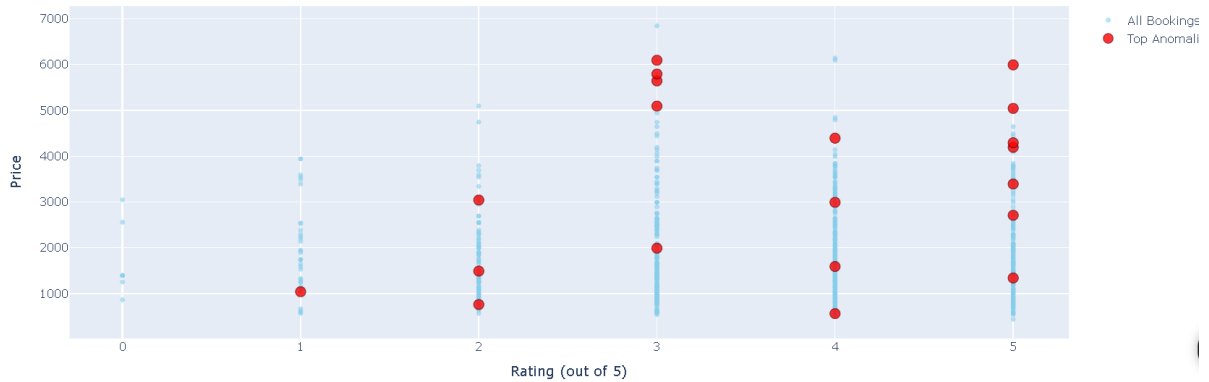
            color='DarkRed'
        )
    ))

fig.update_layout(
    title='Price vs Rating with Top Anomalies Highlighted',
    xaxis_title='Rating (out of 5)',
    yaxis_title='Price',
    hovermode='closest'
)

fig.show()

```

Price vs Rating with Top Anomalies Highlighted



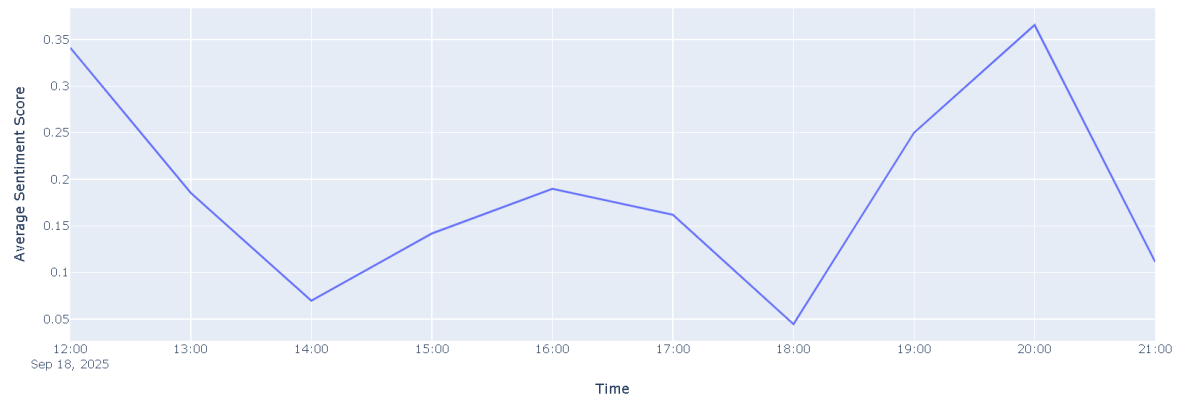
```
import plotly.express as px

fig = px.line(sentiment_over_time, x='check in time', y='sentiment_score',
              title='Average Sentiment Score Over Time (Hourly)')

fig.update_layout(xaxis_title='Time', yaxis_title='Average Sentiment Score')
fig.show()
```

📷 🔍 + 📐 📄 🗑️

Average Sentiment Score Over Time (Hourly)



Add annotations (Sentiment below threshold)

```
import plotly.graph_objects as go

# Define a sentiment threshold
sentiment_threshold = 0 # Example threshold

# Identify periods where sentiment is below the threshold
sentiment_below_threshold = sentiment_over_time[sentiment_over_time['sentiment_score'] < sentiment_threshold]

fig = go.Figure()

# Add the sentiment line
fig.add_trace(go.Scatter(x=sentiment_over_time['check in time'], y=sentiment_over_time['sentiment_score'], mode='lines', name='Average Sentiment'))

# Add annotations for points below the threshold
for index, row in sentiment_below_threshold.iterrows():
    fig.add_annotation(
        x=row['check in time'],
        y=row['sentiment_score'],
        text=f"[row['sentiment_score']:.2f]",
        showarrow=True,
        arrowhead=1,
        bgcolor="rgba(255, 0, 0, 0.5)" # Highlight annotations
```

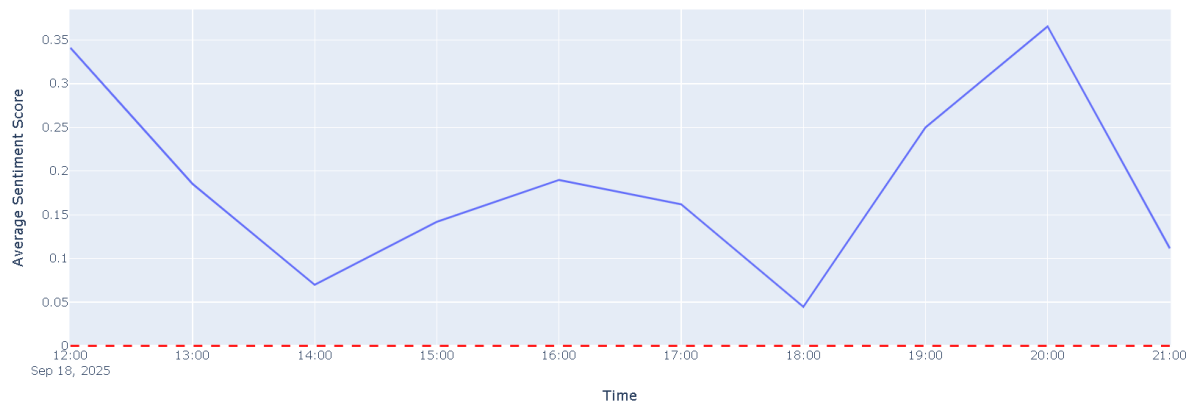
```

fig.update_layout(
    title='Average Sentiment Score Over Time with Annotations (Hourly)',
    xaxis_title='Time',
    yaxis_title='Average Sentiment Score',
    shapes=[
        # Add a horizontal line for the threshold
        dict(
            type='line',
            x0=sentiment_over_time['check in time'].min(),
            y0=sentiment_threshold,
            x1=sentiment_over_time['check in time'].max(),
            y1=sentiment_threshold,
            line=dict(
                color="red",
                width=2,
                dash="dash",
            )
        )
    ]
)

fig.show()

```

Average Sentiment Score Over Time with Annotations (Hourly)



```

import plotly.express as px
import numpy as np

# Calculate Most Active Locations (e.g., top 10 by number of bookings)
most_active_locations = df['Location'].value_counts().head(10)

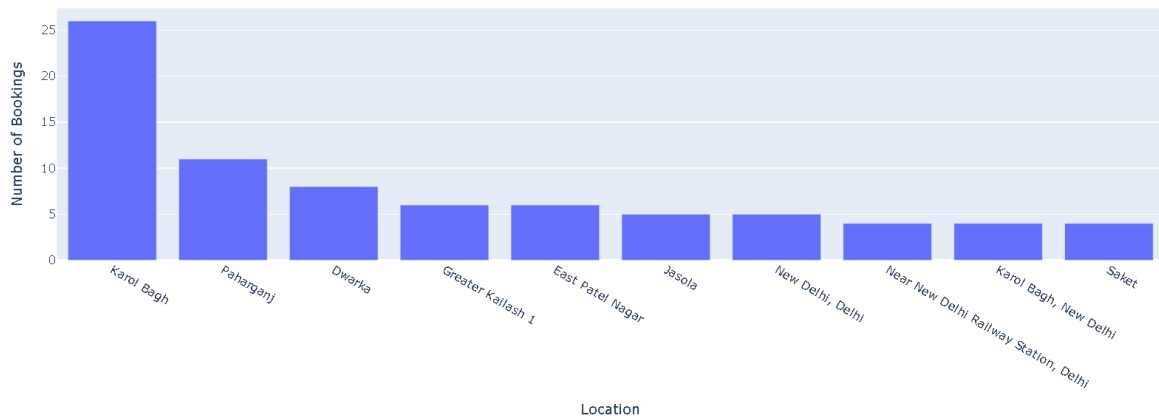
# Get the data for the most active locations
# most_active_locations is a pandas Series, convert it to a DataFrame for Plotly
most_active_locations_df = most_active_locations.reset_index()
most_active_locations_df.columns = ['Location', 'Number of Bookings']

# Create a bar chart of the number of bookings for the most active locations
fig = px.bar(most_active_locations_df, x='Location', y='Number of Bookings',
             title='Number of Bookings for Top 10 Most Active Locations')

fig.update_layout(xaxis_title='Location', yaxis_title='Number of Bookings')
fig.show()

```


Number of Bookings for Top 10 Most Active Locations



```
import pandas as pd
from sklearn.ensemble import IsolationForest
import numpy as np

# Calculate Total Bookings
total_bookings = len(df)
print(f"Total Bookings: {total_bookings}")

# Calculate Average Rating
average_rating = df['Rating(out of 5)'].mean()
print(f"Average Rating: {average_rating:.2f}")

# Calculate Percentage of Online Payments
online_payments_count = df[df['Payment Mode (online/offline)'] == 'online'].shape[0]
percentage_online_payments = (online_payments_count / total_bookings) * 100
print(f"Percentage of Online Payments: {percentage_online_payments:.2f}%")

# Calculate Number of Anomalous Bookings (fraud flagged)
# Assuming 'anomaly_score' < 0 indicates an anomaly based on the Isolation Forest results

# Select numerical features for anomaly detection
numerical_features = ['Price', 'Total Rating', 'Rating(out of 5)', 'No of Guest', 'Total Price', 'sentiment_score']
X = df[numerical_features].copy()
```

```
# Handle potential infinite values or NaNs introduced during previous steps
X.replace([np.inf, -np.inf], np.nan, inplace=True)
X.fillna(X.median(), inplace=True) # Fill NaN values with the median

# Initialize and train the Isolation Forest model
model = IsolationForest(contamination='auto', random_state=42)
model.fit(X)

# Predict anomaly scores (-1 for outliers, 1 for inliers)
df['anomaly_score'] = model.decision_function(X)

anomalous_bookings_count = df[df['anomaly_score'] < 0].shape[0]
print(f"Number of Anomalous Bookings: {anomalous_bookings_count}")

# Identify Most Active Locations (e.g., top 10 by number of bookings)
most_active_locations = df['Location'].value_counts().head(10)
print("\nMost Active Locations:")
display(most_active_locations)
```



Total Bookings: 791
Average Rating: 0.72
Percentage of Online Payments: 58.53%
Number of Anomalous Bookings: 195

Most Active Locations:

Location	count
Karol Bagh	26
Paharganj	11
Dwarka	8
Greater Kailash 1	6
East Patel Nagar	6
Jasola	5
New Delhi, Delhi	5
Near New Delhi Railway Station, Delhi	4
Karol Bagh, New Delhi	4
Saket	4

dtype: int64



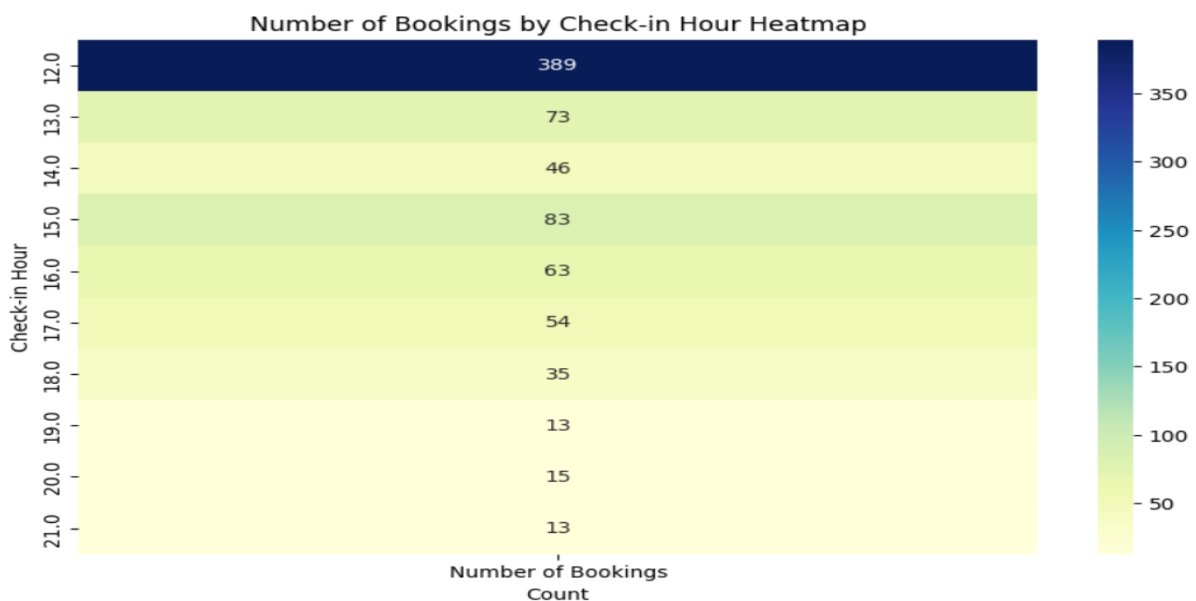
```
import matplotlib.pyplot as plt
import seaborn as sns

# Extract the hour from 'check in time'
df['check in hour'] = df['check in time'].dt.hour

heatmap_data = df.pivot_table(index='check in hour', values='Booking ID', aggfunc='count').reset_index()

# Rename columns for clarity
heatmap_data.columns = ['Check-in Hour', 'Number of Bookings']

# Create the heatmap
plt.figure(figsize=(10, 6))
sns.heatmap(heatmap_data.set_index('Check-in Hour'), annot=True, fmt='d', cmap='YlGnBu')
plt.title('Number of Bookings by Check-in Hour Heatmap')
plt.xlabel('Count')
plt.ylabel('Check-in Hour')
plt.show()
```



'Review Sentiment Score', 'Review Length', 'Repeat Booker', 'Total Guests', 'Weekend Stay Flag

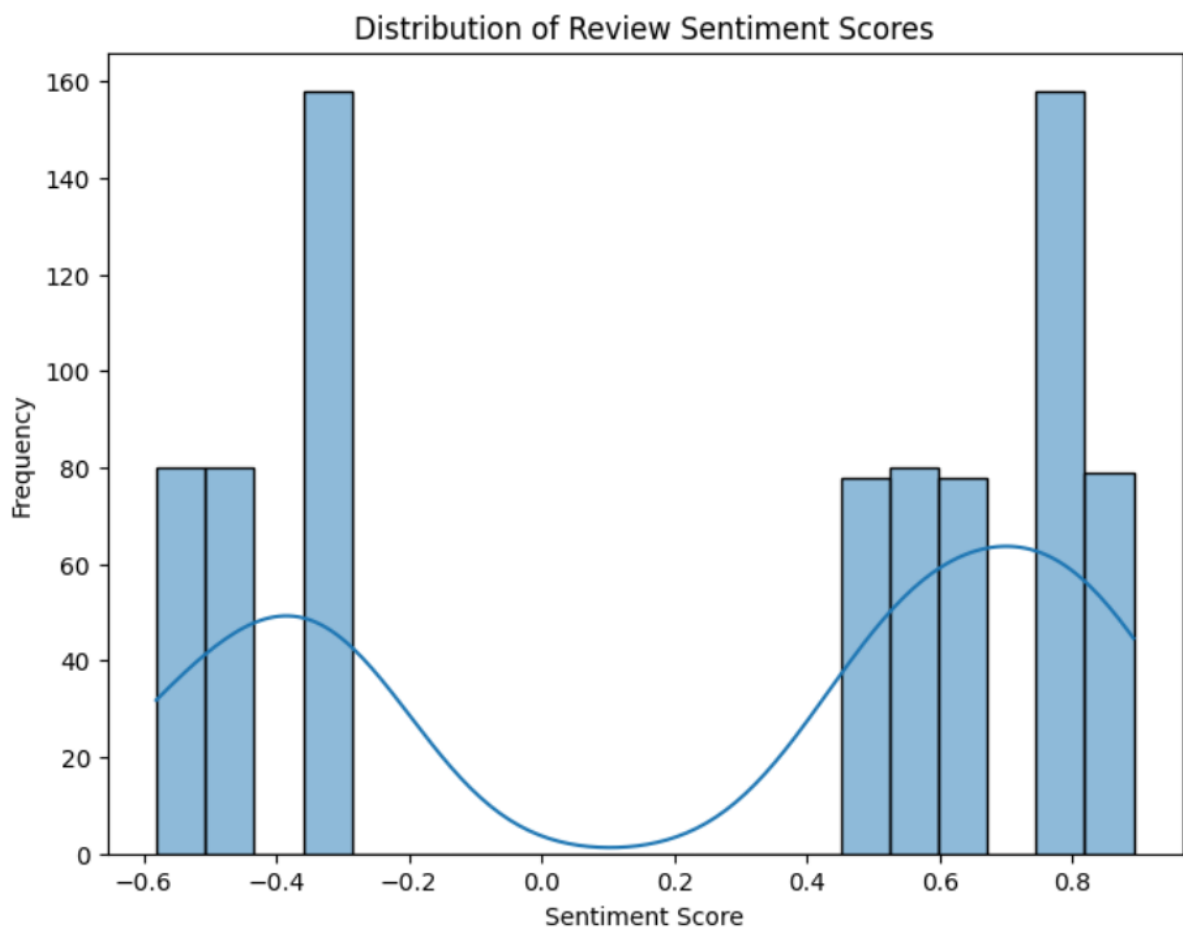
```
# Analyze Review Sentiment Score
print("Review Sentiment Score Analysis:")
display(df['sentiment_score'].describe())

import matplotlib.pyplot as plt
import seaborn as sns

plt.figure(figsize=(8, 6))
sns.histplot(df['sentiment_score'], kde=True, bins=20)
plt.title('Distribution of Review Sentiment Scores')
plt.xlabel('Sentiment Score')
plt.ylabel('Frequency')
plt.show()
```

Review Sentiment Score Analysis:

sentiment_score	
count	791.000000
mean	0.247404
std	0.555263
min	-0.581200
25%	-0.308900
50%	0.487700
75%	0.777800
max	0.892300



```

# Analyze Review Length
print("\nReview Length Analysis:")

# Calculate Review Length (number of words)
df['Review_Length'] = df['Review'].astype(str).apply(lambda x: len(x.split()))

display(df['Review_Length'].describe())

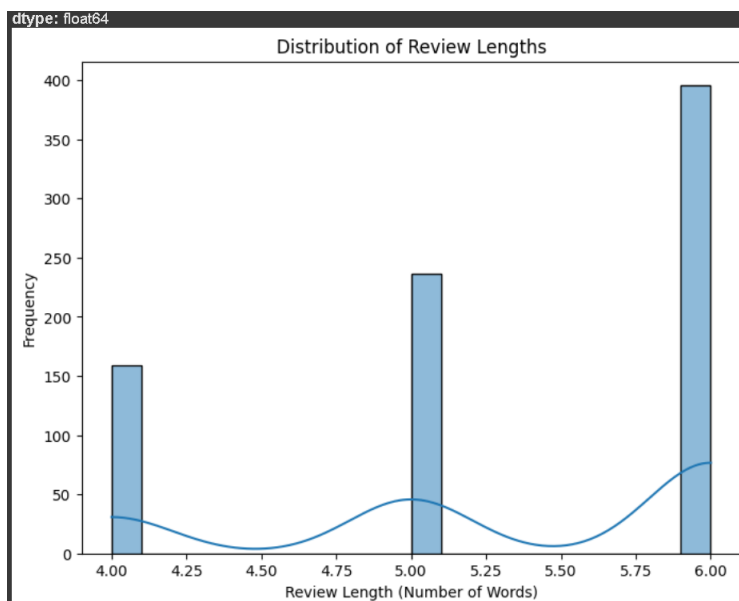
plt.figure(figsize=(8, 6))
sns.histplot(df['Review_Length'], kde=True, bins=20)
plt.title('Distribution of Review Lengths')
plt.xlabel('Review Length (Number of Words)')
plt.ylabel('Frequency')
plt.show()

```

```

Review Length Analysis:
Review_Length
count      791.000000
mean        5.299621
std         0.782717
min         4.000000
25%         5.000000
50%         6.000000
75%         6.000000
max         6.000000

```



```

# Analyze Total Guests
print("\nTotal Guests Analysis:")

# Ensure 'Total_Guests' column exists, create it from 'No of Guest' if not
if 'Total_Guests' not in df.columns:
    df['Total_Guests'] = df['No of Guest']
    df['Total_Guests'] = df['Total_Guests'].fillna(df['Total_Guests'].median())

display(df['Total_Guests'].describe())

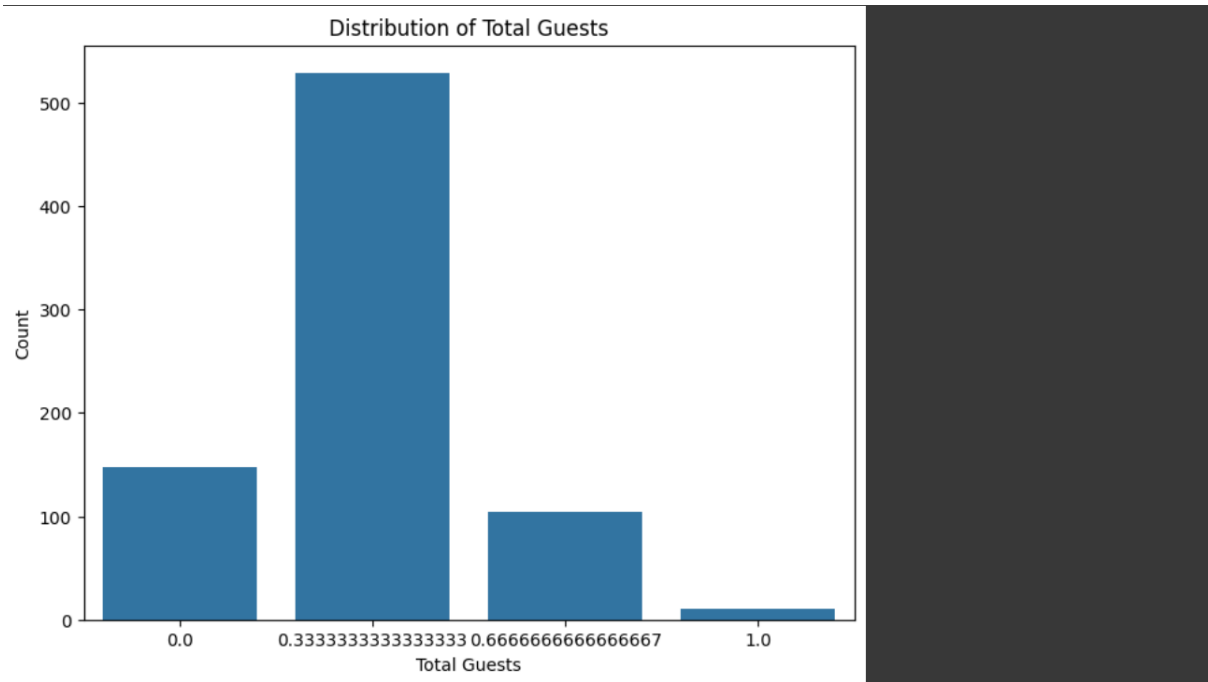
plt.figure(figsize=(8, 6))
sns.countplot(x='Total_Guests', data=df)
plt.title('Distribution of Total Guests')
plt.xlabel('Total Guests')
plt.ylabel('Count')
plt.show()

```

 Total Guests Analysis:

	Total_Guests
count	791.000000
mean	0.323220
std	0.202402
min	0.000000
25%	0.333333
50%	0.333333
75%	0.333333
max	1.000000

dtype: float64



```
# Analyze Weekend Stay Flag
print("\nWeekend Stay Flag Analysis:")
display(df['Weekend_Stay_Flag'].value_counts())

# Analyze Repeat Booker
print("\nRepeat Booker Analysis:")
display(df['Repeat_Booker'].value_counts())
```

 Weekend Stay Flag Analysis:

	count
Weekend_Stay_Flag	
0	791

dtype: int64

Repeat Booker Analysis:

	count
Repeat_Booker	
1	628
0	163

dtype: int64

```

import matplotlib.pyplot as plt
import seaborn as sns

# 1. Create a histogram of the 'sentiment_score' column
plt.figure(figsize=(8, 6))
sns.histplot(df['sentiment_score'], kde=True, bins=20)
plt.title('Distribution of Review Sentiment Scores')
plt.xlabel('Sentiment Score')
plt.ylabel('Frequency')
plt.show()

# 2. Create a countplot of the 'Repeat Booker' column
plt.figure(figsize=(8, 6))
sns.countplot(x='Repeat Booker', data=df)
plt.title('Distribution of Repeat vs. Non-Repeat Bookers')
plt.xlabel('Repeat Booker')
plt.ylabel('Count')
plt.xticks([0, 1], ['Non-Repeat', 'Repeat'])
plt.show()

# 3. Calculate the average 'Rating(out of 5)' for each location.
location_avg_rating = df.groupby('Location')['Rating(out of 5)'].mean()

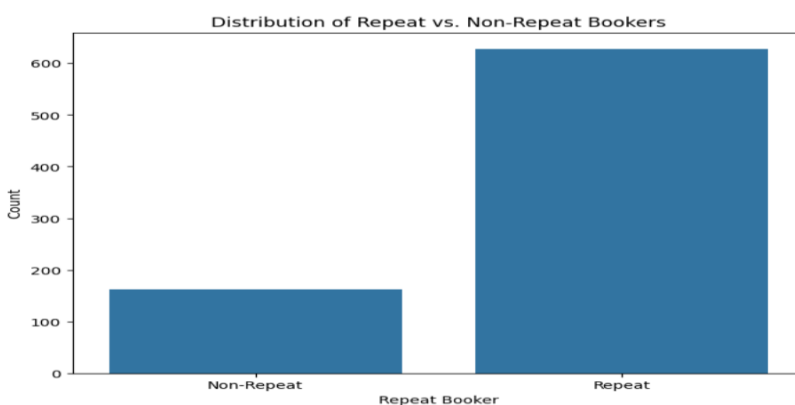
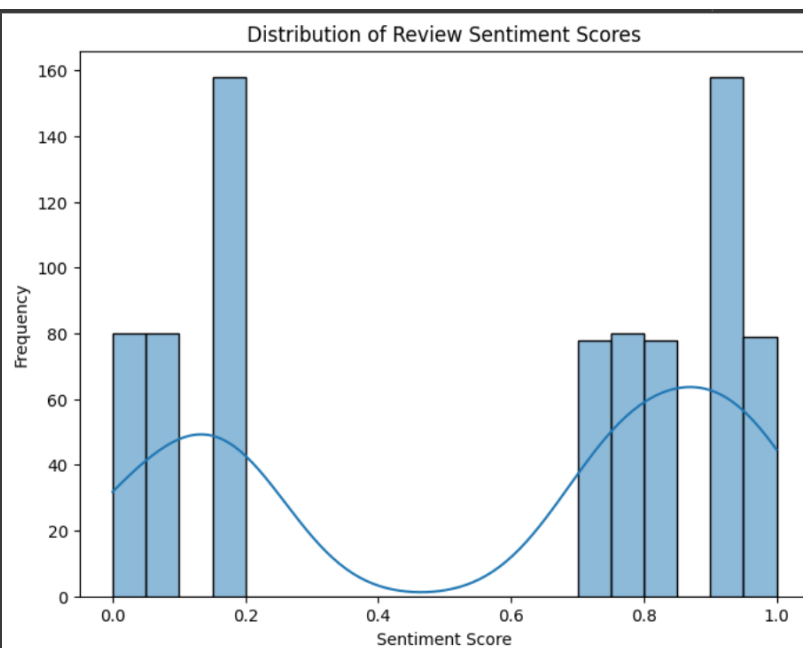
# 4. Identify the top 10 locations based on the average rating.
top_10_locations_rating = location_avg_rating.nlargest(10)

```

```

# 5. Create a heatmap using seaborn that shows the average rating for the top 10 locations.
plt.figure(figsize=(12, 8))
sns.heatmap(top_10_locations_rating.to_frame(), annot=True, fmt=".2f", cmap="YlGnBu")
plt.title('Average Hotel Rating by Location (Top 10)')
plt.xlabel('Average Rating')
plt.ylabel('Location')
plt.yticks(rotation=0)
plt.show()

```





The Above Code snippet is visualized form of this project, below here I Have shared the whole total project Notebook file Link here for findings.

#1- AI Based Hotel Integrated Marketplace Experience and Fraud Anomaly detection

(<https://colab.research.google.com/drive/1NuW3ahpN8xJwT3aqTPMb74sC5pA7x8HB#scrollTo=2F64zS3zv3km>)

#2- AI Based Hotel Integrated Marketplace Experience and Fraud Anomaly detection Model.

(<https://colab.research.google.com/drive/1OiuJ6vzQk6iCEtXU--GN7kQjN7T2BN8Y>)