

Introduction

"Al-Ofoq Al-Lami" (The Shining Horizon) is a prestigious hotel chain with multiple properties spread across various locations. In an effort to enhance its operational efficiency and market position, the chain has initiated a comprehensive data analysis project. This project aims to analyze performance metrics across all hotels in the chain, refine marketing strategies, and increase overall revenue.

The primary objectives of this analysis are to:

- 1. Evaluate hotel performance within the chain
- 2. Identify areas for improvement in operational efficiency
- 3. Develop targeted marketing strategies based on data insights
- 4. Increase overall revenue and profitability for the entire chain

This report presents the methodology, findings, and recommendations derived from our extensive data analysis. The insights gathered will serve as a foundation for strategic decision-making, aimed at optimizing the performance of the chain.

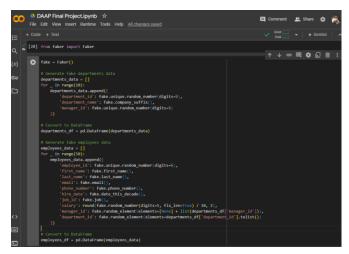
Table of Contents

- 1. Executive Summary
- 2. Project Overview
 - 2.1 Background
 - 2.2 Objectives
 - 2.3 Scope
- 3. Data Collection and Preparation
 - 3.1 Data Sources
 - 3.2 Data Cleaning and Preprocessing
 - 3.3 Data Integration
- 4. Methodology
 - 4.1 Tools and Technologies Used
 - 4.2 Data Analysis Techniques
 - 4.3 Sentiment Analysis Model
- 5. Data Visualization and Analysis
 - 5.1 Key Performance Indicators
 - 5.2 Customer Segmentation
 - 5.3 Booking Patterns and Trends
 - 5.4 Room Type Analysis
 - 5.5 Cancellation Analysis
- 6. Findings and Insights
- 7. Extra- Sentiment analysis model

1- Data cleaning, collecting and preparation: Tools: Google Collab, MySQL Server

In this phase of the Hotel Chain Project, we focused on preparing and cleaning the data to ensure its quality and consistency for analysis. We performed the following steps:

- 1. Sample Data Creation (Python Faker library):
 - We created sample data for the Employees table to complement the existing booking and review data.

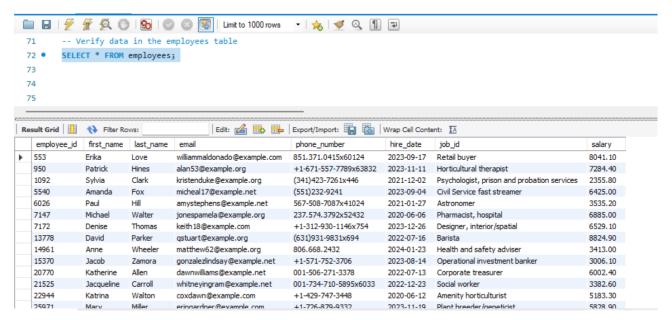


The data is then saved into a CSV file which was loaded into MySQL:

Employees Table Sample:

```
-- Create the employees table

CREATE TABLE IF NOT EXISTS employees (
employee_id INT AUTO_INCREMENT PRIMARY KEY
first_name VARCHAR(255) NOT NULL,
last_name VARCHAR(255) NOT NULL,
email VARCHAR(255),
phone_number VARCHAR(50) NOT NULL,
hire_date DATE NOT NULL,
job_id VARCHAR(255),
salary DECIMAL(10, 2) NOT NULL
);
```



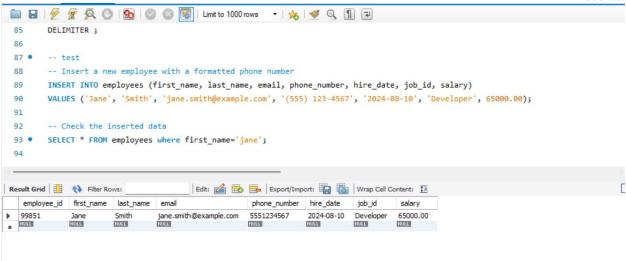
Creating a trigger on the employee table:

The trigger is for formatting employees' phone numbers correctly.

Syntax:

```
CREATE TRIGGER before_employee_insert
BEFORE INSERT ON employees
FOR EACH ROW
BEGIN
    -- Remove non-numeric characters from phone_number
    SET NEW.phone_number = REGEXP_REPLACE(NEW.phone_number, '[^0-9]', '');
END$$
```

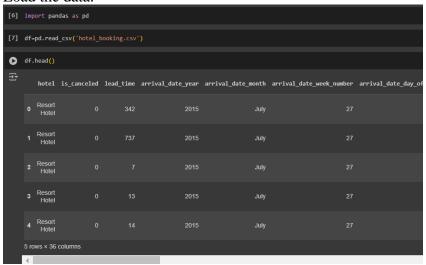
Results:



2. Booking Data Cleaning:

- o Converted date columns to proper datetime format for easier analysis.(Using excel)
- Load the data into python for cleaning.

Load the data:



Handle missing data:

```
### Check for missing values print(df.isnull().sum())
# fill missing values for numrics df.fillna(df.median(numeric_only=True), inplace=True)
### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

### Drop rows with missing values df.dropna(axis=0, inplace=True)

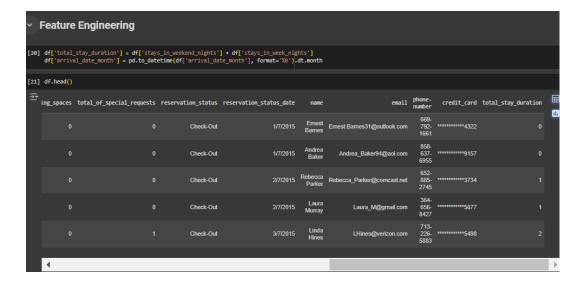
### Drop rows with missing v
```

Results:

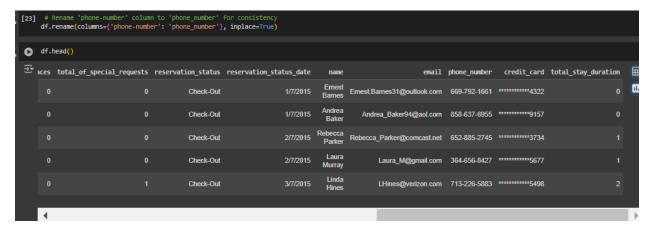
```
print(df.isnull().sum())
⊕ hotel
    lead_time
arrival_date_year
     arrival_date_month
     arrival_date_week_number
     arrival_date_day_of_month
stays_in_weekend_nights
     stays_in_week_nights
     adults
     children
    market_segment
distribution_channel
     is_repeated_guest
     previous_cancellations
     previous_bookings_not_canceled
     assigned_room_type
    booking_changes
deposit_type
     agent
     days_in_waiting_list
     adr
    required_car_parking_spaces
total_of_special_requests
     reservation_status
     reservation_status_date
     email
     phone-number
     credit_card
```

Handling duplicates:
Handling duplicates:
[18] # Drop duplicate rows
df.drop_duplicates(inplace=True)

Feature engineering (create columns that might be useful later)



Renaming columns (Formatting):



Save the cleaned data to CSV:



- 3. Review Data Cleaning:
- Will be used for sentiment analysis

This data cleaning process has prepared our datasets for integration into the SQL database and subsequent analysis in Power BI. The cleaned data will allow for more accurate and meaningful insights into hotel performance, customer satisfaction, and operational efficiency.

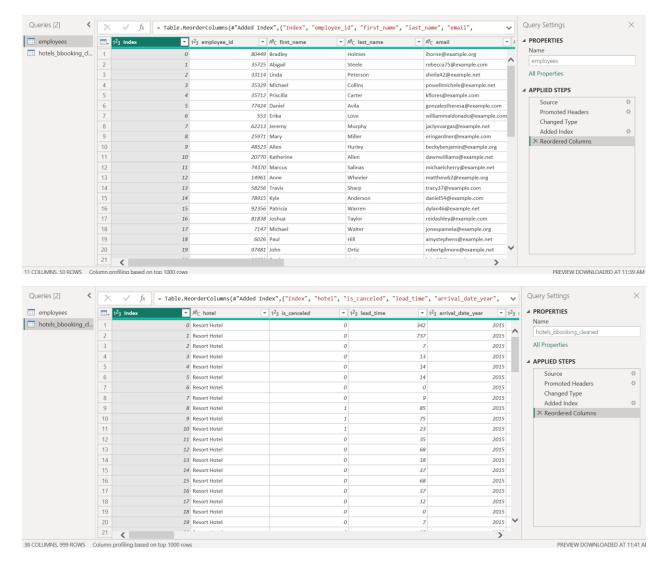
By creating sample data for departments and employees, we've also set the stage for potential analysis of staff performance and departmental efficiency, which could provide valuable insights for hotel management.

2-Data Visualization

Tools: Power BI

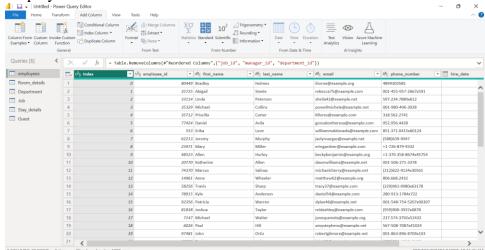
After loading the two tables (Cleaned Booking data + Employees) we are going to drill down the tables using power Bi's DAX and Built-in functions:

First, lets add index column to each table, we will use it as a PK to manage the relationship between the tables.

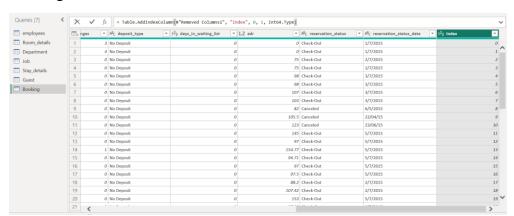


The new tables:

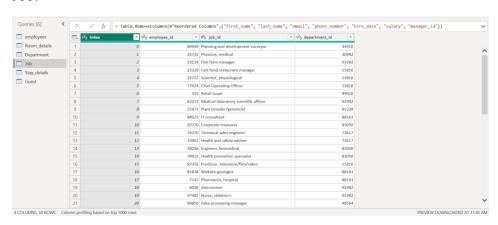
Employee:



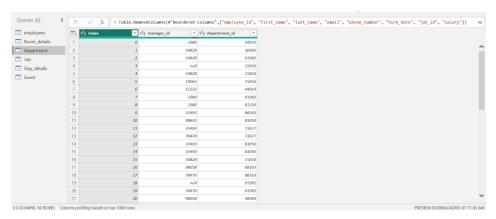
Booking:



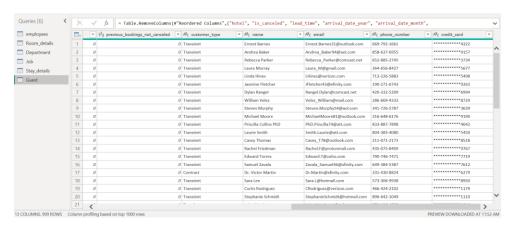
Job:



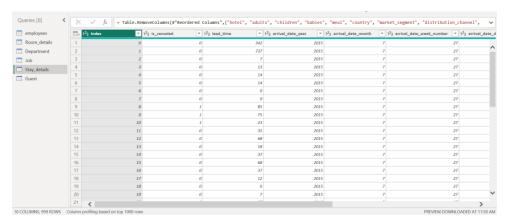
Department:



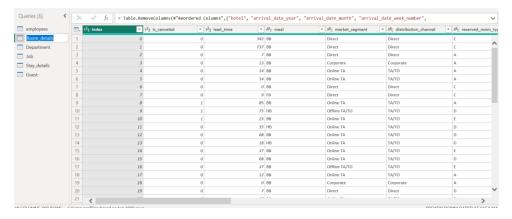
Guest:



Stay details:



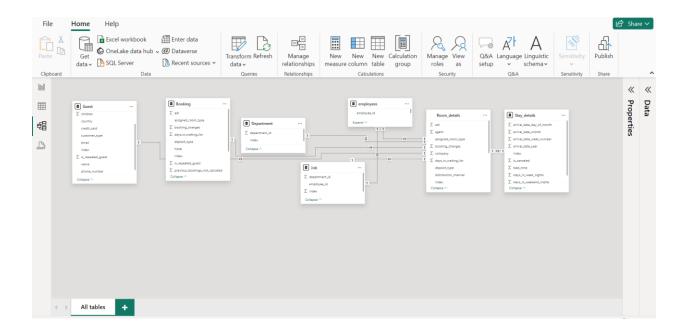
Room_Details:



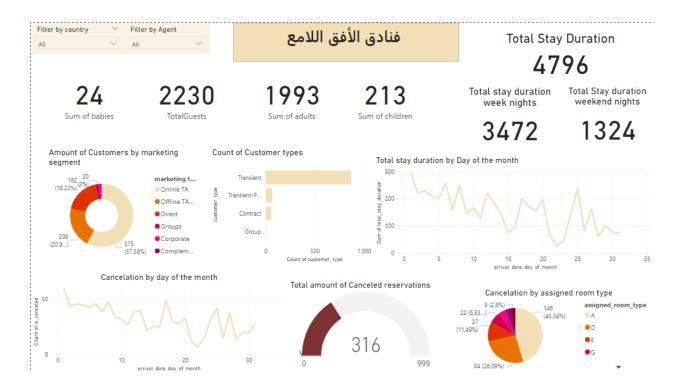
Summary

- Employee Table: Contains core employee information.
- **Job Table:** Contains job-related details.
- **Department Table:** Contains department-related details.
- **Booking Table:** Contains booking-related information.
- Guest Table: Contains guest-specific details.
- Stay Details Table: Contains stay-specific details.
- Room Details Table: Contains room-related details.

Model View:



3- Visualization analysis:



Observations:

- 1. Correlation between Guest Types and Stay Duration:
 - The total number of guests (2230) is composed of 1993 adults, 213 children, and 24 babies.
 - The total stay duration of 4796 nights is a combination of 3472 weeknight stays and 1324 weekend night stays.
 - This suggests that the adult guests are likely responsible for the majority of the weeknight stays, while the guests with children (213 children + 24 babies) may be contributing more to the weekend night stays.
- 2. Segmentation and Booking Behavior:
 - The largest marketing segment is "Offline TA" (57.54% of customers), followed by "Direct" (20%) and "Online TA" (13.22%).
 - This implies that the hotel has a strong relationship with offline travel agencies, and a significant portion of its business comes from this channel.

• However, the presence of "Direct" and "Online TA" segments suggests that the hotel is also attracting direct bookings and leveraging online travel agencies.

3. Cancellation Patterns and Room Type:

- The dashboard shows that the total amount of canceled reservations is 316.
- Looking at the "Cancellation by assigned room type" chart, we can see that the majority of canceled reservations (84 or 26.58%) are for room type "A", followed by room types "E" (146 or 46.20%) and "D" (45 or 14.24%).
- This could indicate that certain room types are more prone to cancellations, potentially due to factors like pricing, availability, or guest preferences.
 Understanding these patterns can help the hotel optimize its room inventory and pricing strategies.

4. Seasonal Trends in Stay Duration:

- The "Total Stay Duration by Day of the Month" chart shows some fluctuations in the total stay duration throughout the month.
- These fluctuations could be related to seasonal factors, such as weekends, holidays, or specific events that influence the hotel's occupancy and stay patterns.
- Analyzing these trends over time can help the hotel anticipate demand and adjust its operations and pricing accordingly.

5. Cancellation Trends and Day of the Month:

- The "Cancellation by day of the month" chart displays the number of cancellations per day of the month.
- Identifying any patterns or peaks in cancellations could help the hotel understand the underlying factors driving these cancellations, such as changes in travel plans, last-minute bookings, or specific events.
- This information can be used to implement proactive measures to reduce cancellations, such as adjusting cancellation policies, improving communication with guests, or offering incentives for guests to maintain their reservations.

4- Extra Credit:

Tools: Google colab

Sentiment analysis model:

In this task, we developed a machine learning application to perform sentiment analysis on hotel reviews using Python in Google Colab. The goal was to classify hotel reviews as either positive (1) or negative (0) based on the text content.

Steps Undertaken:

1- Data Loading and Exploration:

We began by loading the (hotel_Reviews.csv) dataset and exploring its structure using basic pandas functions. This included checking for missing values and reviewing the first few rows of the dataset to understand the content.

2-Data Preprocessing:

The dataset had a 'Review' column containing the review text and a 'Liked' column with binary values (1 for positive, 0 for negative). The text data in the 'Review' column was then preprocessed:

Converted all text to lowercase to ensure uniformity.

Removed punctuation and numbers to focus only on the words that carry meaning.

3- Feature Extraction:

We converted the cleaned text data into numerical form using TF-IDF (Term Frequency-Inverse Document Frequency) vectorization. This technique transforms the text into a matrix of features that the machine learning model can understand.

4-Model Building and Training:

We split the dataset into training and testing sets using an 80-20 split.

A Logistic Regression model was then trained on the TF-IDF features extracted from the training set.

```
Model:

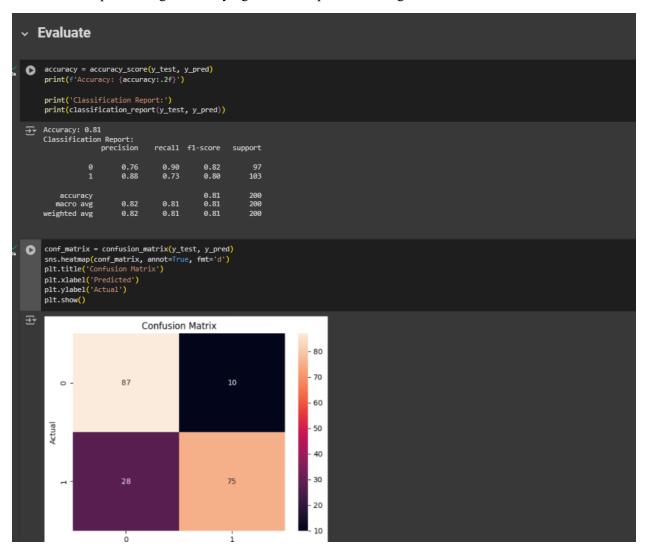
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = LogisticRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)

[13] def predict_sentiment(review):
    # Preprocess the review text
    review = review.lower() # Convert to lowercase
    review = "'.join([c for c in review if c.isalnum() or c.isspace()]) # Remove punctuation and numbers
    review_tfidf = tfidf.transform([review]).toarray() # Convert to TF-IDF

# Predict sentiment using our trained model
prediction = model.predict(review_tfidf)
    return prediction[0] # Return 1 for positive, 0 for negative
```

5-Model Evaluation:

The model's performance was evaluated on the testing set using metrics such as accuracy, classification report (precision, recall, F1-score), and a confusion matrix. These metrics provided insights into how well the model was performing in classifying reviews as positive or negative.



6-Sentiment Prediction for New Reviews:

A custom function, predict sentiment, was created to take a new review as input, preprocess it, and use the trained model to predict its sentiment. This allows for real-time sentiment analysis on new reviews, outputting either a positive (1) or negative (0) classification.

```
Predicating:

[14] new_review = "The hotel was great, very clean and the staff was friendly!"
    sentiment = predict_sentiment(new_review)
    print(f'The sentiment of the review is: {"Positive (1)" if sentiment == 1 else "Negative (0)"}')

The sentiment of the review is: Positive (1)

Predicating:

[14] new_review = "The hotel was great, very clean and the staff was friendly!"
    sentiment of the review is: {"Positive (1)" if sentiment == 1 else "Negative (0)"}')

Predicating:

[14] new_review = "The hotel was great, very clean and the staff was friendly!"
    sentiment of the review is: {"Positive (1)" if sentiment == 1 else "Negative (0)"}')

The sentiment of the review is: Negative (0)

The
```

7-Conclusion:

The final product is a machine learning model capable of performing sentiment analysis on hotel reviews, which can be used to predict whether a review is positive or negative. This application could be valuable for understanding customer feedback, improving services, and enhancing overall customer satisfaction in the hospitality industry.