

# Trip Advisor Hotel Reviews

---

## **Problem Statement**

**Background Information:** Trip Advisor is a popular travel website where users can write reviews about hotels, restaurants, and other travel-related services. These reviews provide valuable insights into the quality, customer experience, and overall satisfaction of various establishments. Analyzing these reviews can help businesses understand customer preferences, improve their services, and make informed decisions.

**Dataset Information:** The dataset used in this project consists of Trip Advisor hotel reviews. It is stored in a CSV file format and contains several columns such as "Review," "Rating," and "Length." The "Review" column contains the text of the reviews written by customers, the "Rating" column represents the rating given by the customers (on a scale of 1 to 5), and the "Length" column indicates the length of each review.

**Problem Statement:** The objective of this project is to perform exploratory data analysis (EDA) on the Trip Advisor hotel reviews dataset. The analysis aims to gain insights into the sentiment, length, word count, and other characteristics of the reviews. By analyzing the reviews, we can understand customer sentiments, identify patterns or trends, and extract useful information to improve the quality of services provided by hotels.

**Approach:** The project starts by loading the dataset into a pandas DataFrame and examining the first few rows to get a glimpse of the data. The 'Rating' column is then transformed into binary labels ('Positive' and 'Negative') to facilitate sentiment analysis.

The EDA begins by analyzing the length of the reviews, including the character count, word count, mean word length, and mean sentence length. These metrics provide insights into the structure and complexity of the reviews. Visualizations, such as box

plots and density distributions, are created to explore the relationship between these features and sentiment.

Next, the reviews are preprocessed by removing stopwords (commonly used words with little sentiment value) and special characters to focus on meaningful words. This preprocessing step helps improve the accuracy of sentiment analysis.

The project then generates and visualizes the most frequently occurring words, bigrams, and trigrams in the reviews. This analysis provides insights into language patterns and commonly used phrases in customer reviews. By examining the sentiment associated with these frequent language patterns, we can better understand customer sentiments and preferences.

Overall, this project aims to provide a comprehensive analysis of Trip Advisor hotel reviews, enabling businesses to extract valuable insights and improve their services based on customer feedback.

## **Framework**

- 1) Importing Required Libraries:** Begin by importing the necessary libraries such as NumPy, Pandas, Matplotlib, Seaborn, NLTK, and other relevant modules. These libraries will be used for data manipulation, visualization, and natural language processing tasks.
- 2) Loading the Dataset:** Read the Trip Advisor hotel reviews dataset from a CSV file using the `pd.read_csv()` function. Print the first few rows of the DataFrame to inspect the structure and content of the dataset.
- 3) Transforming Rating into Binary Labels:** The 'Rating' column represents the customer ratings on a scale of 1 to 5. Transform this column into binary labels ('Positive' and 'Negative') based on a predefined threshold. In the provided code, the `ratings()` function is used to map ratings greater than 3 as 'Positive' and ratings between 1 and 3 as 'Negative'. Visualize the distribution of positive and negative ratings using a pie chart.
- 4) Exploratory Data Analysis (EDA):**
  - **Review Length Analysis:** Calculate and print the length of a sample review (character count). Create a new column 'Length' in the DataFrame to store the length of each review. Visualize the distribution of review lengths using box plots and density distributions.
  - **Word Count Analysis:** Count the number of words in a sample review and print the result. Implement the `word_count()` function to count the number of words in each review and create a new column 'Word\_count' to store the counts. Visualize the distribution of word counts using appropriate plots.
  - **Mean Word Length and Sentence Length Analysis:** Calculate the mean word length and mean sentence length for each review. Add the respective columns 'mean\_word\_length' and 'mean\_sent\_length' to the DataFrame. Visualize the distributions of these features using appropriate plots.
  - **Text Preprocessing:** Implement the `clean()` function to preprocess the reviews. This function converts the text to lowercase, removes special characters using regular expressions, and eliminates stopwords using NLTK's 'stopwords' corpus. Apply the `clean()` function to the 'Review' column and update the DataFrame with the cleaned reviews.

### **5) Most Frequently Occurring Words Analysis:**

- **Word Frequency Analysis:** Create a corpus by combining all the cleaned reviews. Count the frequency of each word using the Counter module. Visualize the top 10 most frequently occurring words using a bar plot.
- **Bigram Frequency Analysis:** Use the CountVectorizer with a specified n-gram range to generate bigrams from the cleaned reviews. Calculate the frequency of each bigram and create a DataFrame with the results. Visualize the top 10 most frequently occurring bigrams using a bar plot.

**6) Trigram Frequency Analysis:** Similarly, use the CountVectorizer to generate trigrams from the cleaned reviews. Calculate the frequency of each trigram and create a DataFrame. Visualize the top 10 most frequently occurring trigrams using a bar plot.

**7) Data Cleaning and Preparation:** Remove the unnecessary columns from the DataFrame, keeping only the relevant columns for further analysis. Print the information of the DataFrame to verify the changes.

**8) Implementation Considerations:** Provide additional information and considerations about the implementation, such as downloading NLTK resources (e.g., stopwords) and setting default styling and font size parameters for plots.

## **Code Explanation**

- 1) Importing Required Libraries:** This part of the code imports various libraries necessary for data analysis, visualization, and natural language processing tasks. The libraries include NumPy, Pandas, Matplotlib, Seaborn, NLTK, and other relevant modules. These libraries provide functions and tools to perform specific operations on data.
- 2) Loading the Dataset:** The code reads the Trip Advisor hotel reviews dataset from a CSV file using the `pd.read_csv()` function. This function loads the data into a Pandas DataFrame, which is a tabular data structure that allows easy manipulation and analysis of the data. The `data.head(10)` line prints the first 10 rows of the DataFrame to get an overview of the dataset's structure.
- 3) Transforming Rating into Binary Labels:** In this section, the code transforms the 'Rating' column of the dataset into binary labels. The `ratings()` function is defined to map ratings between 1 and 3 (inclusive) to 'Negative' and ratings between 4 and 5 (inclusive) to 'Positive'. This function is then applied to the 'Rating' column using the `data['Rating'].apply(ratings)` line. The resulting DataFrame now has 'Positive' and 'Negative' labels in place of the original ratings. The code further creates a pie chart to visualize the distribution of positive and negative ratings.
- 4) Exploratory Data Analysis (EDA):** This section focuses on analyzing various aspects of the dataset to gain insights into the reviews' characteristics.
  - **Review Length Analysis:** The code calculates the length of a sample review using the `len()` function and stores it in the 'Length' column of the DataFrame. This provides information about the number of characters in each review. The distribution of review lengths is visualized using box plots and density distributions.
  - **Word Count Analysis:** The code counts the number of words in a sample review by splitting the review text and counting the resulting list's length. This count is stored in the 'Word\_count' column of the DataFrame. The distribution of word counts is visualized using appropriate plots.
  - **Mean Word Length and Sentence Length Analysis:** The code calculates the mean word length and mean sentence length for each review. These values are added to the DataFrame as 'mean\_word\_length' and 'mean\_sent\_length' columns.

Visualizations, such as box plots or density distributions, can be used to understand the distributions of these features.

**5) Text Preprocessing:** In this section, the code preprocesses the review texts to clean and normalize them before further analysis. The `clean()` function is defined to convert the text to lowercase, remove special characters using regular expressions, and eliminate stopwords using NLTK's 'stopwords' corpus. The `clean()` function is applied to the 'Review' column using the `df['Review'].apply(clean)` line. This step ensures that the reviews are in a standardized format for subsequent analysis.

**6) Most Frequently Occurring Words Analysis:** This section aims to identify and visualize the most frequently occurring words, bigrams, and trigrams in the reviews.

- **Word Frequency Analysis:** The code creates a corpus by combining all the cleaned reviews. It then uses the Counter module to count the frequency of each word in the corpus. The top 10 most frequently occurring words are visualized using a bar plot.
- **Bigram and Trigram Frequency Analysis:** The code employs the CountVectorizer from scikit-learn to generate bigrams and trigrams from the cleaned reviews. It calculates the frequency of each bigram and trigram, creating separate DataFrames for each. The top 10 most frequently occurring bigrams and trigrams are visualized using bar plots.

**7) Data Cleaning and Preparation:** This part of the code removes unnecessary columns from the DataFrame, retaining only the relevant ones for further analysis. The `df.drop()` function is used to drop the specified columns. Finally, the `df.info()` line is included to print the information about the DataFrame, including the remaining columns and their data types.

**8) Implementation Considerations:** The code provides additional considerations for implementation, such as downloading NLTK resources (e.g., stopwords) using the `nltk.download()` function and setting default styling and font size parameters for plots using the `plt.rcParams[]` statements. These considerations ensure that the necessary resources are available and the plots generated are visually appealing.

## **Future Work**

The Trip Advisor Hotel Reviews project provides a foundation for analyzing and understanding hotel reviews. To further enhance the project and extract deeper insights, here are some future work suggestions:

**1. Sentiment Analysis:** Extend the analysis to perform sentiment analysis on the reviews. This can involve training a machine learning model to classify reviews as positive, negative, or neutral based on their content. Sentiment analysis can provide a more nuanced understanding of customer opinions and help identify areas for improvement in hotel services.

### **Implementation Guide:**

- Collect additional labeled data for sentiment analysis, consisting of hotel reviews with corresponding sentiment labels (positive, negative, neutral).
- Preprocess the data by cleaning and normalizing the reviews.
- Split the dataset into training and testing sets.
- Train a machine learning model (e.g., Naive Bayes, Support Vector Machines, or deep learning models like LSTM) on the training set to predict sentiment labels.
- Evaluate the model's performance using the testing set.
- Apply the trained model to predict the sentiment of new hotel reviews.

**2. Topic Modeling:** Explore topic modeling techniques to identify the underlying topics or themes present in the reviews. Topic modeling algorithms like Latent Dirichlet Allocation (LDA) can uncover latent topics by grouping related words together. This analysis can reveal the key aspects of the hotel experience that customers frequently mention.

### **Implementation Guide:**

- Preprocess the reviews by removing stopwords, special characters, and performing lemmatization.
- Convert the preprocessed text into a document-term matrix.
- Apply the LDA algorithm to the document-term matrix to extract the underlying topics.
- Determine the optimal number of topics using techniques such as coherence scores or visual inspection of topic distributions.

- Interpret the generated topics by examining the most representative words for each topic.
- Analyze the prevalence of different topics across the reviews and identify key areas of interest or concern.

**3. Aspect-Based Sentiment Analysis:** Dive deeper into the reviews by performing aspect-based sentiment analysis. Instead of analyzing the overall sentiment of the reviews, this technique focuses on extracting sentiments related to specific aspects of the hotel experience (e.g., room cleanliness, staff behavior, amenities). This can provide more detailed insights into the strengths and weaknesses of the hotel.

#### **Implementation Guide:**

- Identify a list of key aspects or categories to analyze (e.g., room, service, location, food).
- Preprocess the reviews and tokenize them into sentences.
- Use dependency parsing or rule-based techniques to extract sentences that mention the desired aspects.
- Apply sentiment analysis techniques to determine the sentiment polarity of the aspect-specific sentences.
- Aggregate the sentiment scores for each aspect to get an overall sentiment assessment.
- Analyze the sentiment distribution across different aspects and identify areas for improvement or areas of excellence.

**4. Review Summarization:** Implement an automatic review summarization feature to generate concise summaries of the reviews. This can help hotel managers quickly grasp the main points and sentiments expressed by customers without manually reading each review in detail.

#### **Implementation Guide:**

- Preprocess the reviews by cleaning and tokenizing the text.
- Apply sentence segmentation to divide the reviews into individual sentences.
- Use extractive summarization techniques (e.g., TextRank algorithm) to identify the most important sentences that capture the essence of each review.
- Generate a summary by combining the selected sentences from multiple reviews.



- Evaluate the quality of the generated summaries using metrics such as ROUGE scores or human evaluation.
- Fine-tune the summarization model based on feedback and iterate for improved results.

**5. User Review Clustering:** Perform clustering analysis on the reviews to group similar reviews together based on their content. This can help identify distinct customer segments or preferences and tailor hotel services accordingly.

#### **Implementation Guide:**

- Preprocess the reviews by cleaning, tokenizing, and vectorizing the text.
- Apply a clustering algorithm (e.g., K-means, DBSCAN, or hierarchical clustering) to cluster the reviews into groups.
- Explore different feature representations such as TF-IDF or word embeddings (e.g., Word2Vec, GloVe) to capture the semantic meaning of the text.
- Visualize the clusters using techniques like t-SNE or PCA to understand the distribution of reviews in the feature space.
- Analyze the characteristics of each cluster to identify common themes or preferences among the customer segments.

By incorporating these future work suggestions into the Trip Advisor Hotel Reviews project, you can gain deeper insights, improve customer satisfaction, and make data-driven decisions to enhance the hotel experience. Remember to adapt the implementation steps based on the specific requirements of each task and experiment with different techniques to achieve optimal results.

## **Concept Explanation**

Alright, buckle up and get ready for a fun and friendly explanation of the algorithm used in the Trip Advisor Hotel Reviews project! Imagine you're a detective trying to uncover the secrets hidden within the reviews of hotels. You need a special tool to help you analyze and understand all the information. That's where the algorithm comes in!

The algorithm used in this project is called Natural Language Processing (NLP). It's like having a super smart language detective that can make sense of the words people use in their hotel reviews. Let's break it down step by step with a funny example:

### **Step 1: Preprocessing**

Just like you need to tidy up your detective's office before getting to work, the algorithm preprocesses the reviews. It cleans up the text by removing unnecessary stuff like punctuation and converting everything to lowercase. It's like giving your detective a clean slate to work with.

**Example:** Imagine a review that says, "OMG!!! This hotel is AMAZING!!!!" The algorithm would clean it up to "omg this hotel is amazing."

### **Step 2: Tokenization**

Tokenization is like giving your detective a magnifying glass to zoom in on each individual word in the reviews. It breaks down the text into smaller pieces called tokens, which are usually words. It's like handing your detective a collection of clues to investigate.

**Example:** Using our previous cleaned-up review, the tokens would be ["omg", "this", "hotel", "is", "amazing"].

### **Step 3: Stopword Removal**

Stopword removal is like getting rid of the boring and irrelevant words that don't contribute much to the investigation. These words, like "is," "the," and "and," are like distractions that your detective can do without.

**Example:** After removing the stopwords, our tokens would be ["omg", "hotel", "amazing"].

#### **Step 4: Lemmatization**

Lemmatization is like your detective finding the root form of each word. It helps to group similar words together, making it easier to understand their meaning. It's like your detective discovering that "amazing" and "amaze" are related.

**Example:** After lemmatization, our tokens would be ["omg", "hotel", "amaze"].

#### **Step 5: Analysis and Visualization**

Now that your detective has all the clean and meaningful words, it's time to analyze and visualize the data. The algorithm can do all sorts of cool things like counting the number of words, finding the most frequently occurring words, and even detecting patterns and sentiments in the reviews. It's like your detective using special tools to uncover interesting insights and create visual charts and graphs.

**Example:** The algorithm might reveal that "amaze" is mentioned a lot, indicating that customers find the hotels to be amazing.

So, there you have it! The algorithm used in the Trip Advisor Hotel Reviews project is like your trusty language detective, helping you understand and make sense of all the reviews. It cleans, organizes, and analyzes the text, giving you valuable insights and making your investigation a whole lot easier and fun!

Remember, algorithms are like helpful sidekicks to humans, assisting us in understanding the vast amount of data we encounter. So, embrace the power of algorithms and embark on your detective journey to uncover the secrets hidden within hotel reviews! Happy investigating!

## **Exercise Questions**

**1. Question: How can you determine the average length of the hotel reviews in the dataset? Provide the code snippet to calculate it.**

**Answer:** To determine the average length of the hotel reviews, we can calculate the mean length of the 'Review' column. Here's the code snippet to calculate it:

```
average_length = data['Review'].str.len().mean()

print(f"The average length of the hotel reviews is: {average_length}")
```

**2. Question: Can you explain the purpose and significance of tokenization in natural language processing?**

**Answer:** Tokenization is the process of breaking down a text into smaller units, usually words or tokens. Its purpose is to facilitate analysis and understanding of the text data. By tokenizing the text, we can treat each word as an individual entity, enabling us to perform tasks such as counting words, identifying patterns, and extracting meaningful information. Tokenization is a crucial step in natural language processing as it forms the foundation for further analysis and processing of textual data.

**3. Question: How can you determine the most frequently occurring words in the hotel reviews? Provide the code snippet to find the top 10 most common words.**

**Answer:** To find the most frequently occurring words in the hotel reviews, we can use the Counter class from the collections module. Here's the code snippet to find the top 10 most common words:

```
from collections import Counter

corpus = ' '.join(data['Review']).split()

word_freq = Counter(corpus)

most_common = word_freq.most_common(10)
```

```
print("Top 10 Most Common Words:")
```

```
for word, count in most_common:
```

```
    print(f'{word}: {count} occurrences')
```

**4. Question: How can you visualize the distribution of review lengths for positive and negative ratings? Explain the steps involved.**

**Answer:** To visualize the distribution of review lengths for positive and negative ratings, we can use a box plot. Here are the steps involved:

- Split the dataset into two groups based on ratings: positive and negative.
- Calculate the lengths of the reviews in each group.
- Create a box plot with the review lengths as the y-axis and the ratings as the x-axis.
- Set the labels and titles for the plot.
- Display the plot.

This visualization allows us to compare the distribution of review lengths between positive and negative ratings, providing insights into how review lengths may differ based on sentiment.

**5. Question: How can you identify the most frequently occurring bigrams (pairs of consecutive words) in the hotel reviews? Provide the code snippet to find the top 10 most common bigrams.**

**Answer:** To identify the most frequently occurring bigrams in the hotel reviews, we can use the `CountVectorizer` class from the `sklearn.feature_extraction.text` module. Here's the code snippet to find the top 10 most common bigrams:

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
cv = CountVectorizer(ngram_range=(2, 2))
```

```
bigrams = cv.fit_transform(data['Review'])
```

```
bigram_freq = bigrams.sum(axis=0)
```

```
bigram_freq_sorted = sorted([(count, bigram) for bigram, count in
zip(cv.get_feature_names(), bigram_freq.tolist())], reverse=True)[:10]
```

```
print("Top 10 Most Common Bigrams:")
```

```
for count, bigram in bigram_freq_sorted:
```

```
    print(f"{bigram}: {count} occurrences")
```

By using the `ngram_range=(2, 2)` parameter, we specify that we want to extract bigrams. The code then calculates the frequency of each bigram and sorts them in descending order to obtain the top 10 most common bigrams.