# Project Overview

This project involves a comprehensive analysis of Airbnb reviews, leveraging machine learning and natural language processing (NLP) to extract insights from demographic property, location and textual data. The goal is to understand customer experiences, preferences, and broader themes from market, pricing,reviews and to present these insights effectively using data visualizations and PowerPoint presentations.

It involves:

- Text analysis to highlight key factors influencing guest satisfaction.
- Predictive modeling to forecast trends that could help hosts optimize their offerings.
- Ultimately, the project intends to enhance the guest experience and improve service quality for Airbnb hosts.

The translation module enhances the project's ability to handle multilingual review data, particularly beneficial for platforms like Airbnb where host offerings and user-generated content can vary widely in language/region. By ensuring reviews are translated into English, the subsequent analysis (sentiment analysis, topic modeling, etc.) can be standardized and applied uniformly, improving the reliability and comparability of insights across different linguistic demographics.

## Libraries and Setup:

- Libraries: The code imports several libraries necessary for handling data and performing machine learning and NLP tasks. These include:
  - Pandas and NumPy for data manipulation.
  - Scikit-learn, WordCloud  for machine learning tasks like analysis, clustering and dimensionality reduction.
  - NLTK, Langdetect,Spacy for natural language processing, such as tokenizing and lemmatizing text.
  - Matplotlib and Seaborn for data visualization.
  - Python-pptx for creating PowerPoint presentations programmatically.
- Colors and Styles: Set a consistent color palette and font styles for visualizations using Matplotlib and Seaborn to ensure that all plots have a uniform appearance.

# Data Extraction and Initial Processing:

1. Connection to MongoDB: The code connects to a MongoDB database and retrieves data from the `listingAndReviews` collection.
2. DataFrame Creation: The retrieved data is converted into a Pandas DataFrame for easier manipulation.
3. Data Extraction: Text data is loaded from external json files into pandas DataFrames for analysis.
4. Data Cleaning: Implements a series of text preprocessing steps including converting text to lowercase, removing punctuation, numbers, and stopwords. Special attention is given to French text, employing SpaCy for lemmatization to standardize the text further.
5. Reviews/Comments  are carefully extracted and translated to English for NLP.

# Data Cleaning and Transformation:

1. Transformation: Each record in the DataFrame is passed through the `transform_record` function, which cleans and formats various fields, including converting numerical values and dates, and flattening nested structures.
2. Handling NaN Values: Columns with NaN values are identified, and specific strategies are applied to impute missing data for certain columns like `bedrooms`, `beds`, `bathrooms`, `security_deposit`, and `cleaning_fee`.
3. Handling Specific Records: Filters out specific records (e.g., a listing with a particular `_id`) and adjusts bedroom counts in specific scenarios (e.g., studios).
4. An outlier in terms of pricing is removed.
5. Text Vectorization: Converts cleaned textual data into a TF-IDF matrix, which numerically represents the importance of words in the dataset.
6. Dimensionality Reduction with SVD: Applies Singular Value Decomposition (SVD) to reduce the high-dimensional TF-IDF matrix into a two-dimensional space for visualization and further analysis.

Here's a concise summary of the main tasks performed and insights derived:

# Statistical Analysis:

1. Correlation Analysis: The code calculates correlation matrices for various numerical features including prices, bedrooms, and review scores. This helps in understanding the relationships between different features of the listings.
2. Amenities Analysis: A frequency analysis of amenities listed across properties is performed, identifying the most common amenities and visualizing their frequencies.
3. Total Cost Calculation: Adds a derived column for 'total cost', which sums up the price, cleaning fee, and extra people charges. This feature is then correlated and visualized against other features like 'accommodates' to explore financial dynamics.
4. Review Scores Analysis: Investigates the distribution of various review score metrics and their correlations with pricing and total cost. This includes generating histograms for individual review categories and a correlation heatmap to visualize these relationships.
5. The visualizations provide a clear depiction of various aspects of the Airbnb listings, such as which cities have the most listings, which amenities are most common in popular listings, and how pricing varies with the number of beds and property types across different locations.
6. The analyses also delve into less common scenarios, such as listings without bedrooms or beds, highlighting anomalies or special cases in the dataset.
7. Finally, the integration of results into a presentation format ensures that these insights are communicated effectively, making them accessible and actionable for decision-makers or stakeholders interested in the Airbnb market dynamics.

This comprehensive approach not only aids in understanding the current state of Airbnb listings but also supports strategic decisions related to property management, pricing strategies, and marketing based on the characteristics that most influence listing popularity and pricing.

- Variance Computation: The variances of numeric columns are calculated to understand the dispersion of data points. This information is formatted and added to the presentation to provide insights into the variability of each numerical feature.

## Visualization:

- Heatmaps: Used to show the correlation matrices, providing a color-coded visualization of how closely different variables are related.

- Bar Plots: Employed to showcase the frequency of common amenities and visualize the mean values of certain features.
- Scatter Plots: Created to display relationships between continuous variables, such as total cost versus the number of people accommodated.
- Histograms: Utilized to examine the distribution of review scores, illustrating the frequency of different scores received by the listings.
- Word Clouds: Generates word clouds for visual summarization of the most frequent and significant words in the data set.
- Cluster Visualization: Visualizes the data clusters in two-dimensional space post-SVD, offering a graphical representation of the text data segmentation.
- Elbow and Silhouette Analysis: Visual plots are used to determine the optimal number of clusters, enhancing the meaningfulness of the cluster groups.
- City Distribution: Visualizes the distribution of listings across the top cities, highlighting the concentration of listings in certain areas.
- Amenity Popularity: Analyzes and displays the top amenities that correlate with higher booking frequencies, inferred from the average number of reviews.
- Cost Analysis:
  - Highest Total Cost: Identifies and displays the record with the highest total cost.
  - Price vs. Accommodates Including Extra Fees: Visualizes the relationship between the total cost (including extra people fees) and the capacity of listings.
- Property Price Analysis:
  - Aggregates data by city, property type, and number of beds to analyze the average prices in various configurations.
  - Filters and visualizes data specifically for apartments in Barcelona to explore how the number of beds affects pricing in this city and property type.

## Amenity Analysis:

1. Amenity Extraction: Splits the 'amenities' string into individual amenities for each listing, allowing for more detailed analysis.
2. Popularity of Amenities: Calculates the average number of reviews per amenity, using the number of reviews as a proxy for booking frequency. This analysis helps identify which amenities are most associated with frequently booked listings.

3. Visualization of Top Amenities: Creates a bar plot showing the top ten amenities based on booking frequency, providing insights into what features make a listing more attractive to guests.

## Geographical Data Extraction:

1. City and Country Extraction: Pulls city and country information from the 'address' field of each listing, refining the dataset for potential geographical analyses or visualizations.
2. Relevant Data Display: Prepares a final DataFrame that contains key geographical and temporal data, allowing for straightforward access to city, country, and review dates.

## Data Preparation and Preprocessing for NLP:

1. Text Vectorization: Converts cleaned textual data into a TF-IDF matrix, which numerically represents the importance of words in the dataset.
2. Dimensionality Reduction with SVD: Applies Singular Value Decomposition (SVD) to reduce the high-dimensional TF-IDF matrix into a two-dimensional space for visualization and further analysis.

## Clustering and Visualization:

1. Elbow Method: Uses the elbow method to determine the optimal number of clusters for KMeans by plotting the sum of squared errors (SSE) against the number of clusters.
2. KMeans Clustering: Implements KMeans clustering to categorize the documents into clusters based on their textual content, reduced to two principal components by SVD.
3. Visualization of Clusters: Plots the clustered data points in the two-dimensional SVD space, showing how documents group together based on their content similarities.
4. Association Matrix: Constructs an association matrix that maps the relationship between the dominant topics (from LDA) and the clusters (from KMeans), offering insights into how well the topics align with the clusters identified by KMeans.

5. Silhouette Scores for Clustering: Employs silhouette analysis to validate the cohesiveness and separation of the identified clusters, ensuring the clustering results are robust and informative.

## Topic Modeling:

1. LDA Implementation: Utilizes Latent Dirichlet Allocation (LDA) to model topics in the dataset, extracting themes or topics from the textual data.
2. Displaying Topics: Shows the top terms associated with each topic, providing insights into the thematic structure of the dataset.

## Insights:

- The analysis effectively uncovers underlying themes and clusters in the textual data, which can be used to understand common topics or concerns expressed in reviews or descriptions.
- The combination of clustering and topic modeling provides a robust framework for extracting actionable insights from unstructured text, which is crucial for content-driven strategies and improving user engagement or satisfaction.

## Sentiment Analysis:

1. Sentiment Extraction: Applies sentiment analysis to each text snippet to determine the overall sentiment (positive, negative, or neutral) expressed in the text.
2. Aggregation and Visualization: Aggregates sentiment scores by property ID and visualizes the average sentiment scores per host, providing insights into the overall guest satisfaction.
3. Polarity and Subjectivity Analysis: Implements sentiment analysis to compute polarity and subjectivity scores for the reviews, providing insights into the overall sentiment trends within the comments.

## Advanced Text Analysis:

1. N-Gram Analysis: Performs bi-gram and trigram analysis to identify common two-word and three-word phrases, offering deeper insights into common themes or issues mentioned in the text.

2. Frequency Analysis: Counts the occurrences of each word and phrase, helping to pinpoint key topics or concerns in the reviews.
3. Silhouette Analysis: You evaluate the optimal number of clusters for KMeans clustering based on silhouette scores, which measure how similar an object is to its own cluster compared to other clusters. The optimal number of clusters is then used to categorize the comments, providing a quantitative measure to guide the clustering decision.
4. Cluster Visualization: While not fully implemented here, this would typically involve visualizing the clusters in reduced dimensional space to assess how distinct the groupings are, further supported by your use of SVD for dimensionality reduction.
5. Language Detection: Uses the `langdetect` library to identify the language of the input text. This functionality is crucial for deciding whether translation is needed based on the target language (English in this case).

## Overall Insights:

- The combination of these techniques provides a comprehensive understanding of the text data, revealing underlying themes, sentiments, and the structure of the text.
- The visualizations and summarizations make the data more accessible and interpretable, which is crucial for making informed decisions based on textual analysis.
- This approach is particularly valuable in contexts where text data is abundant, such as customer feedback, reviews, or any other user-generated content, enabling stakeholders to glean actionable insights from unstructured data.

## Data Integration and Reporting:

- PowerPoint Presentation: Integrates analytical results into a dynamic PowerPoint presentation, including visual and textual summaries of key findings like topic distributions, clustering results, and sentiment analysis.
- Comprehensive Reporting: Each step from data loading to analysis is structured to culminate in a report that provides actionable insights, effectively communicated through well-designed slides.

## Additional Functionalities:

- Handling Multiple Languages: Special provisions for handling French reviews, including custom stopword removal and lemmatization, highlight the project's capability to adapt to multilingual datasets.
- The topic modeling and subsequent analysis was repeated for french reviews to find regional and language bias in reviews.
- For french reviews the clusters were not clearly detected.

## Conclusion:

This project showcases a robust application of various data science techniques to analyze textual data comprehensively. From preprocessing to detailed analysis with LDA and clustering, followed by effective visualization and presentation of results, the project aims to provide a deep understanding of customer feedback in a structured and visually appealing manner.

This module significantly expands the project's scope by incorporating multilingual support, allowing for a more inclusive and comprehensive analysis of global user feedback. It is a vital component for projects involving international data, ensuring that language barriers do not hinder analytical accuracy or depth.