# Analyzing Airbnb rental prices and occupancy rate

## Problem Description:

The objective of this project is to analyze Airbnb rental prices and occupancy rates in New York City. The project uses the Airbnb dataset from Kaggle, which contains information on Airbnb listings in New York City, including the rental prices, occupancy rates, and other variables such as location, type of property, and number of bedrooms.

The goal of the project is to gain insights into the factors that affect Airbnb rental prices and occupancy rates in New York City, and to develop a predictive model that can forecast rental prices and occupancy rates based on various variables.

**Dataset**: The dataset used in this project is the New York City Airbnb Open Data from Kaggle. The dataset contains information on 48,895 Airbnb listings in New York City, including the rental prices, occupancy rates, and various other variables such as location, type of property, and number of bedrooms. The dataset covers a period of 11 years, from 2008 to 2019.

*The dataset is provided in a CSV file format, and it contains 16 columns, including the following:*

- id: Unique identifier for each listing

- name: Name of the listing

- host_id: Unique identifier for each host

- hostname: Name of the host

- neighbourhood_group: The borough in which the listing is located (Bronx, Brooklyn, Manhattan, Queens, or Staten Island)

- neighborhood: The specific neighborhood in which the listing is located

- latitude: The latitude of the listing's location

- longitude: The longitude of the listing's location

- room type: The type of room being listed (Private room, Shared room, or Entire home/apt)

- price: The price per night of the listing

- minimum nights: The minimum number of nights a guest is required to stay

- number_of_reviews: The number of reviews that the listing has received

- last_review: The date of the last review

- reviews_per_month: The number of reviews per month

- calculated_host_listings_count: The number of listings that the host has

- availability_365: The number of days the listing is available for booking in the next 365 Days.

Deliverables:

The deliverables for this project include:

- Exploratory data analysis of the Airbnb dataset to gain insights into the factors that affect rental prices and occupancy rates in New York City

- Data preprocessing and cleaning to ensure data quality and integrity

- Development of a predictive model that can forecast rental prices and occupancy rates based on various variables.

- Creation of visualizations and dashboards to communicate the findings and insights gained from the analysis

- Report generation and conclusion that summarizes the findings and insights, provides for future research, and discusses the limitations of the analysis.

The project aims to provide insights into the factors that affect Airbnb rental prices and occupancy rates in New York City, and to develop a predictive model that can help hosts and guests make informed decisions about renting and booking Airbnb listings.

# Framework/Approach:

To analyze Airbnb rental prices and occupancy rates in New York City, we will follow the following framework:

**Section 1: Data Loading and Preparation:**

•        Load the Airbnb dataset from Kaggle

•        Check for missing values and anomalies in the data

•        Filter out any irrelevant columns or rows that are not needed for the analysis

•        Convert data types to appropriate formats

**Section 2: Exploratory Data Analysis:**

•        Create visualizations to explore the data and gain insights into the factors that affect Airbnb rental prices and occupancy rates in New York City

•        Use descriptive statistics and summary metrics to gain a general understanding of the data

•        Identify any patterns or trends in the data

•        Conduct hypothesis testing to identify significant differences between different groups or categories

**Section 3: Data Preprocessing and Cleaning:**

•        Handle missing values by imputing or dropping them as appropriate

•        Correct any errors or inconsistencies in the data

•        Normalize or standardize the data as appropriate

•        Address any outliers or extreme values in the data

**Section 4: Feature Engineering and Selection:**

•        Create new features or variables that may be useful for the analysis

- Select the most relevant features or variables for the predictive model

- Use feature scaling or selection techniques as appropriate

**Section 5: Predictive Modeling:**

- Develop a predictive model that can forecast Airbnb rental prices and occupancy rates based on various variables

- Use machine learning algorithms such as linear regression, decision trees, or random forests to build the model

- Train the model on a subset of the data and evaluate its performance on a validation set

- Tune the hyperparameters of the model to optimize its performance

**Section 6: Data Visualization and Dashboard Creation:**

- Create visualizations that illustrate the findings and insights gained from the analysis

- Use interactive charts and graphs if possible to enable users to explore the data and gain a better understanding of Airbnb rental prices and occupancy rates in New York City

- Create a dashboard to summarize the findings and insights gained from the analysis

**Section 7: Report Generation and Conclusion:**

- Summarize the findings and insights gained from the analysis in a report

- Provide recommendations for future research and areas for further exploration

- Discuss the limitations of the analysis and possible sources of error in the data

This framework aims to provide a comprehensive and systematic approach to analyzing Airbnb rental prices and occupancy rates in New York City. By following this framework, we can ensure that our analysis is rigorous and thorough, and that we are able to gain meaningful insights into the factors that affect Airbnb rentals in the city.

# Code Explanation :

Here is the simple explanation for the code you can find at code.py file.

**Section 1: Data Loading and Preparation:**

This section loads the Airbnb dataset from Kaggle, checks for missing values and anomalies in the data, filters out any irrelevant columns or rows that are not needed for the analysis, and converts data types to appropriate formats.

**Section 2: Exploratory Data Analysis:**

This section creates visualizations to explore the data and gain insights into the factors that affect Airbnb rental prices and occupancy rates in New York City, uses descriptive statistics and summary metrics to gain a general understanding of the data, identifies any patterns or trends in the data, and conducts hypothesis testing to identify significant differences between different groups or categories.

**Section 3: Data Preprocessing and Cleaning:**

This section handles missing values by imputing or dropping them as appropriate, corrects any errors or inconsistencies in the data, normalizes or standardizes the data as appropriate, and addresses any outliers or extreme values in the data.

**Section 4: Feature Engineering and Selection:**

This section creates new features or variables that may be useful for the analysis, selects the most relevant features or variables for the predictive model, uses feature scaling or selection techniques as appropriate, and prints the selected feature names.

**Section 5: Predictive Modeling:**

This section splits the data into training and test sets, trains three different regression models (linear regression, decision tree regression, and random forest regression), and prints the evaluation metrics ($R^2$ and RMSE) for each model.

**Section 6: Data Visualization and Dashboard Creation:**

This section creates visualizations that illustrate the findings and insights gained from the analysis using Plotly Express, and creates a dashboard using Dash to summarize the findings and insights gained from the analysis.

**Section 7: Report Generation and Conclusion:**

This section summarizes the findings and insights gained from the analysis in a report, provides recommendations for future research and areas for further exploration, and discusses the limitations of the analysis and possible sources of error in the data.

The code uses several Python libraries, including pandas for data loading and manipulation, seaborn and matplotlib for data visualization, scipy for hypothesis testing, scikit-learn for data preprocessing, feature selection, and machine learning, Plotly Express and Dash for interactive data visualization and dashboard creation.

Overall, this code provides a comprehensive framework for analyzing Airbnb rental prices and occupancy rates in New York City, from data loading and preparation to predictive modeling and data visualization, and can be easily adapted and extended for similar projects in other cities or regions.

# Future Work :

**Section 1: Data Collection and Preparation:**

This section involves collecting more recent and comprehensive data on Airbnb rentals in New York City, cleaning and preprocessing the data as before, and conducting a comparative analysis of changes in rental prices and occupancy rates over time.

**Section 2: Exploratory Data Analysis:**

This section involves conducting a more detailed and granular analysis of the factors that affect Airbnb rental prices and occupancy rates in different neighbourhoods, using advanced visualization and machine learning techniques to identify patterns, trends, and outliers in the data, and conducting deeper hypothesis testing to validate the findings.

**Section 3: Feature Engineering and Selection:**

This section involves creating and selecting more sophisticated and complex features or variables that capture the nuances and complexities of Airbnb rentals in New York City, using advanced feature selection and dimensionality reduction techniques to identify the most important and relevant features for the predictive model.

**Section 4: Predictive Modeling:**

This section involves building more advanced and sophisticated machine learning models that can predict Airbnb rental prices and occupancy rates with greater accuracy and precision, using ensemble methods, deep learning, or other advanced techniques, and evaluating the performance of the models using more rigorous and robust metrics.

**Section 5: Data Visualization and Dashboard Creation:**

This section involves creating more interactive and engaging visualizations and dashboards that enable users to explore the data and gain insights in real-time, using cutting-edge tools and technologies such as virtual reality, augmented reality, or gamification.

**Section 6: Report Generation and Conclusion:**

This section involves summarizing the key findings and insights gained from the analysis in a more comprehensive and detailed report, providing more actionable recommendations and insights for policymakers, researchers, and other stakeholders, and discussing the implications and impact of the research on the wider community.

**To implement this future work, you can follow these steps:**

1.      Collect more recent and comprehensive data on Airbnb rentals in New York City, either by scraping the web or by obtaining it from a public or private source.

2.      Clean and preprocess the data as before, but use more sophisticated and advanced techniques such as natural language processing, computer vision, or machine learning to handle the complexities and nuances of the data.

3.      Conduct a more detailed and granular analysis of the data using advanced visualization and machine learning techniques, and identify patterns, trends, and outliers in the data that can help to explain or predict rental prices and occupancy rates.

4.      Create more advanced and sophisticated machine learning models using ensemble methods, deep learning, or other advanced techniques, and evaluate their performance using more rigorous and robust metrics such as precision, recall, or F1 score.

5.      Create more interactive and engaging visualizations and dashboards using cutting-edge tools and technologies such as virtual reality, augmented reality, or gamification, and enable users to explore the data and gain insights in real-time.

6.      Summarize the key findings and insights gained from the analysis in a more comprehensive and detailed report, provide more actionable recommendations and insights for policymakers, researchers, and other stakeholders, and discuss the implications and impact of the research on the wider community.

By following these steps, you can build on the existing project and take it to the next level, creating more value and impact for the community and the world.

# Exercise Questions :

**1.       What are some potential factors that influence Airbnb rental prices in New York City, and how could you investigate them using exploratory data analysis?**

Answer: Some potential factors that could influence Airbnb rental prices in New York City include the neighborhood, the property type, the number of bedrooms or bathrooms, the minimum nights stay, the availability ratio, the number of reviews, and the host experience or reputation. To investigate these factors using exploratory data analysis, we could create visualizations such as scatterplots, histograms, or boxplots, to compare the distribution and variation of rental prices across different categories or groups. We could also use summary statistics such as mean, median, or standard deviation, to quantify the central tendency and variability of rental prices across different categories or groups. We could conduct hypothesis testing using t-tests, ANOVA, or chi-square tests, to identify significant differences or associations between rental prices and different categorical or continuous variables.

**2.       What are some potential challenges or limitations of building a predictive model for Airbnb rental prices in New York City, and how could you address them using machine learning techniques?**

Answer: Some potential challenges or limitations of building a predictive model for Airbnb rental prices in New York City include the high dimensionality and sparsity of the data, the heterogeneity and non-linearity of the relationships between the features and the target variable, the presence of outliers or extreme values, the potential overfitting or underfitting of the models, and the bias or randomness of the training and test sets. To address these challenges using machine learning techniques, we could use feature selection or dimensionality reduction methods, such as PCA or LASSO, to reduce the number of features and capture the most important and relevant information. We could use ensemble methods, such as bagging, boosting, or stacking, to combine multiple models and reduce the variance and bias of the predictions. We could use regularization methods, such as Ridge or Lasso regression, to penalize the model complexity and prevent overfitting or underfitting. We could use cross-validation or bootstrapping techniques, to assess the generalization performance of the models and reduce the impact of bias or randomness in the data.

**3.       What are some potential insights or recommendations that could be derived from the analysis of Airbnb rental prices and occupancy rates in New York City, and how could you communicate them effectively to different stakeholders?**

Answer: Some potential insights or recommendations that could be derived from the analysis of Airbnb rental prices and occupancy rates in New York City include the identification of the most popular and profitable neighbourhoods for Airbnb rentals, the factors that most influence rental prices and occupancy rates, the impact of different policies or regulations on the Airbnb market, the potential effects of demographic or economic changes on the Airbnb market, and the comparisons with other cities or regions. To communicate these insights effectively to different stakeholders, we could use different visualizations and dashboards tailored to their needs and preferences, such as heatmaps, line charts, pie charts, or bar charts, that highlight the most relevant and interesting findings. We could also use storytelling techniques and narratives, that frame the analysis in a compelling and persuasive way, and that provide concrete and actionable recommendations or implications for different stakeholders.

**4.      What are some potential ethical or legal issues that could arise from the analysis of Airbnb rental prices and occupancy rates in New York City, and how could you address them appropriately?**

Answer: Some potential ethical or legal issues that could arise from the analysis of Airbnb rental prices and occupancy rates in New York City include the privacy and security of the data, the potential discrimination or bias in the analysis or the models, the compliance with the regulations or policies related

# Concept Explanation :

The algorithm used in this project for predictive modeling is called the Random Forest Regression algorithm. Now, let's dive into the details of this algorithm in a fun and friendly way!

Imagine you're in a forest, and you want to predict the type of tree that's in front of you. But the forest is so big and diverse that you don't know where to start. So, you decide to ask some of your forest friends for help.

Your friends are different animals in the forest, and each one has a different perspective and expertise. Some are good at identifying the shape and size of the leaves, some are good at identifying the texture and color of the bark, and some are good at identifying the type of soil and the climate of the area.

So, you decide to ask each of your friends to give you their opinion on the type of tree based on their expertise. Then, you combine all their opinions into a final decision.

This is basically what the Random Forest Regression algorithm does. It combines the opinions of multiple "decision trees" to make a final prediction.

A decision tree is like a flowchart that asks a series of yes-or-no questions about the features of the data (like the size of the tree, the number of branches, etc.), and assigns a value to the target variable (like the type of tree) based on the answers. Each decision tree is trained on a random subset of the data, and produces a different prediction.

The Random Forest algorithm then combines the predictions of all the decision trees, using a technique called "bagging", which takes the average of all the predictions. This results in a more accurate and robust prediction, since it reduces the variance and bias of the individual decision trees.

For example, if we want to predict the rental price of an Airbnb listing based on features like the neighborhood, the number of bedrooms, the number of reviews, and the host experience, we can use a Random Forest Regression algorithm to combine the opinions of multiple decision trees and make a final prediction.

In summary, the Random Forest Regression algorithm is like asking multiple friends in a forest for their opinions on the type of tree, and combining their opinions to make a final decision. It's a powerful and flexible algorithm that can handle complex and high-dimensional data, and can produce accurate and robust predictions for a wide range of applications.