# Airbnb Analysis

Capstone Project Presentation

Link to Project Presentation

# Overview

Airbnb is an online marketplace for short term rentals. Airbnb allows people from all over the world to host their homes as someone's next stay.

Airbnb allows individuals to rent out their vacant properties to individuals looking for short term accommodation.
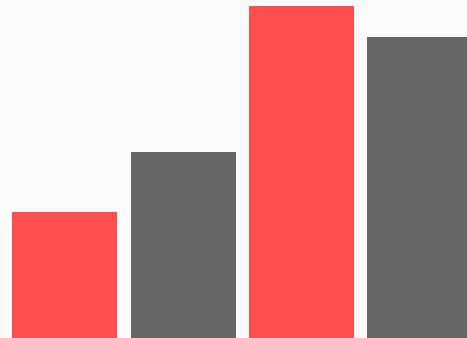
Properties can range from houses, apartments to single and shared rooms and are priced per night or per stay.

# The problem

With the ever changing market it can be challenging for Airbnb hosts to determine the optimal rent prices for their properties.

Since price depends on numerous factors ranging from property type to amenities offered, as well as location customer reviews and ratings.

Hosts can often misinterpret the prices for their neighbourhoods and miss the opportunity to good profits for their listings.

# The solution

Construct a data driven solution by using machine learning to predict rental prices for each property.

# Data Source

The file clean_data_bourgh.csv, contains data about Airbnb listings in Toronto, Canada. The dataset contains a total of 8387 records, where each row represents a unique listing and every column represents important data about the listing. The following is the description of some columns in this dataset.

1. **host_since**: The date that host listed their first Airbnb listing.
2. **host_response_rate**: How fast the host responded to customer inquiries.
3. **neighbourhood**: The Toronto neighbourhood of the listing.
4. **property_type**: Type of property (Entire home, private room, share room, hotel room)
5. **price**: Price of the listing property per day
6. **bourgh**: The Toronto bourgh where the of the listing property

# Exploration Data Analysis

The data exploration phase of the project is conducted in Python and Tableau. We primarily analyzed the data based on four segments, host details, location, reviews and amenities. Within the host details segment our goal was to determine whether factors like being a superhost, having a verified identity and the number of listings the host has in the city has any impact on average prices. In the location segment we derived insights like the most and least expensive neighbourhoods in the city, the average prices in each Toronto borough, etc. The reviews segment we dived a little deeper into understanding the impact customer ratings and reviews have on prices and popularity. Lastly, we looked into popular amenities and whether or not they are offered by most Airbnb listings and the impact they have on average prices.

# Project Analysis Phase

During the last segment the team was able to uncover some key insights on how different variables can affect average Airbnb prices in Toronto. The team analyzed the data by host, location, room type, property type, average customer reviews and rating, amenities offered, etc. On average we were able to see that there is a considerable difference in average prices based on the presence of each of these factors. Currently, the team is working towards analyzing the data through visualizations to derive any insights that can help us understand the impact different factors have on average price.

# Research Questions

1. Relation between price and Room Type
2. Relation between price and Property Type
3. Top five most popular amenities
4. Top five most expensive locations
5. Relation between price and amenities
6. Relation between price and location
7. Relation between price and customer reviews and ratings
8. Popular properties by number of reviews
9. Which month has the most bookings
10. Highest number of listings by boroughs
11. Top five expensive and least expensive neighbourhoods and boroughs
12. Relation between price and host response time
13. Relation between price and host response rate

# Data Cleaning

Cleaned the original listings.csv file using Pandas library and created clean_data_borough.csv for project analysis

1. Dropped 37 columns that are not relevant for project analysis
2. Scraped Toronto Postal Codes from Wikipedia to determine the relation between postal codes, neighbourhood and boroughs
3. Added postal codes and bourgh data to original dataset to match each neighbourhood with a bourgh

# ERD Diagram

# Machine Learning Model

Goal is to implement and compare results of these regression models to see which model helps us accurately predict prices

1. Linear Regression
2. Support Vector Regression
3. HistGradientBoostingRegressor
4. RandomForestRegressor
5. ExtraTreesRegressor
6. XGBoostRegressor

```python
In [ ]: from sklearn.model_selection import train_test_split
        from sklearn.preprocessing import StandardScaler
        from sklearn.linear_model import LinearRegression

In [ ]: # Numerical Features to be included in ML
        new_data_df=data_df[["host_is_superhost", "host_identity_verified", "instant_bookable","accommodates", "bedrooms","minimum_nights
                    "toiletries", "high_end_electronics","ac_heater","internet","bbq","home_appliances","coffee_machine","long_t
                    "parking", "kitchen", "elevator","gym","price"]]

In [ ]: # Processng Categorical Features
        for cat_feature in ["host_response_time","neighbourhood_cleansed","bedrooms","property_type","room_type"]:
            new_data_df=pd.concat([new_data_df, pd.get_dummies(data_df[cat_feature])], axis=1)

In [ ]: #Create Target and Feature Variables
        X=new_data_df.drop(["price"], axis=1)
        y=new_data_df["price"]

        #Split the preprocessed data into a training and testing dataset
        X_train, X_test, y_train, y_test= train_test_split(X,y, random_state=78)

In [ ]: # Create a StandardScaler Instance
        scaler=StandardScaler()

        #Fit the StandardScaler
        X_scaler=scaler.fit(X_train)

        #Scale the data
        X_train_scaled=X_scaler.transform(X_train)
        X_test_scaled=X_scaler.transform(X_test)

In [ ]: # Run the Linear Regression model
        model=LinearRegression()
        model.fit(X_train_scaled, y_train)

In [ ]: #Making predictions using the testing data
        predictions=model.predict(X_test_scaled)
        predictions

In [ ]: results=pd.DataFrame({"Prediction" : predictions, "Actual" : y_test}).reset_index(drop=True)
        results.head(20)
```

# Model Results

### Linear Regression Model

```
RMSE train: 0.392
RMSE test: 0.398
R^2 train: 0.685
R^2 test: 0.666
```

### Support Vector Regression Model

```
RMSE train: 0.287
RMSE test: 0.396
R^2 train: 0.831
R^2 test: 0.669
```

### XGBoost Regressor

```
RMSE train: 0.289
RMSE test: 0.353
R^2 train: 0.828
R^2 test: 0.738
```

### HistGradientBoosting Regressor

```
RMSE train: 0.295
RMSE test: 0.352
R^2 train: 0.822
R^2 test: 0.739
```
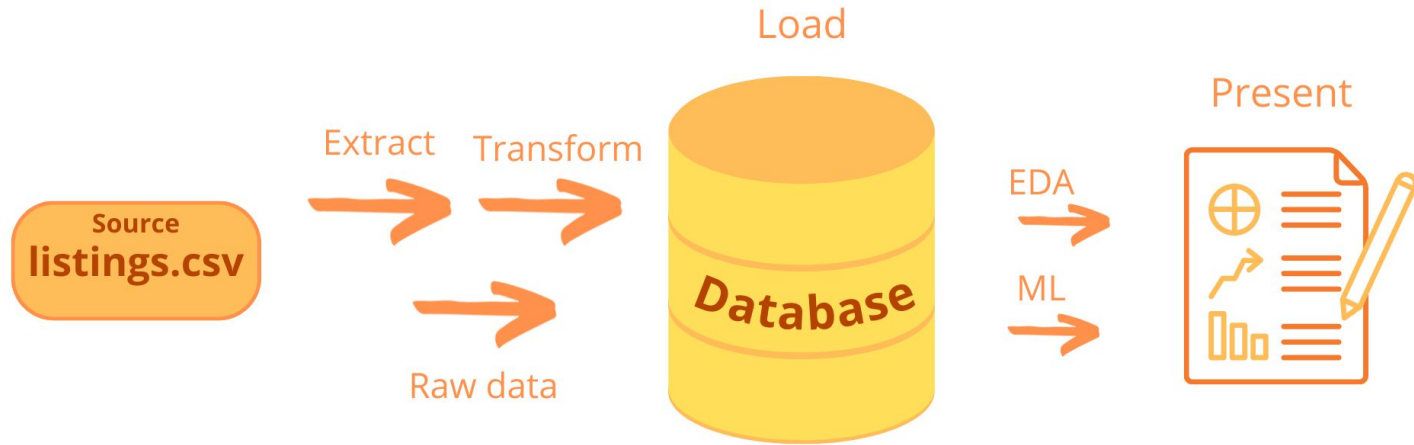
### RandomForest Regressor

```
RMSE train: 0.151
RMSE test: 0.360
R^2 train: 0.953
R^2 test: 0.727
```

### ExtraTrees Regressor

```
RMSE train: 0.000
RMSE test: 0.370
R^2 train: 1.000
R^2 test: 0.711
```

# ETL Pipeline

# Story Outline



Page 1

Page 2

Page 3

Page 4

Toronto Airbnb Analysis

# Dashboard Drafts

## Price vs. Reviews Analysis



## Price vs. Neighbourhood Analysis

# Other Visualizations