# Predicting AirBnB Listing Price within New York City's Boroughs

**EAS 503 | Group 7 |** Gary Yu, Georgious Aland Feltama, Peter Pranata

# Agenda

# Presence of Airbnb in NYC

airbnb

NYC is the most visited city in the U.S

New York City has the 3rd largest active listings in the world
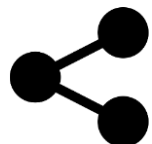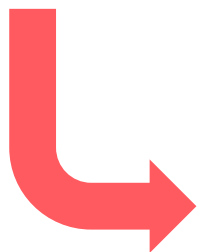
Approximately 48,000 active listings

1,143,036 reviews up to date

Bronx

Manhattan

Queens

Staten Island

Brooklyn

# Motivation and Objective(s)

Airbnb pricing varies a lot; sometimes its cheaper than hotel, other times no

What **factors** within an **Airbnb** listing **dictates** its **nightly price** within the respective **neighborhood groups**?

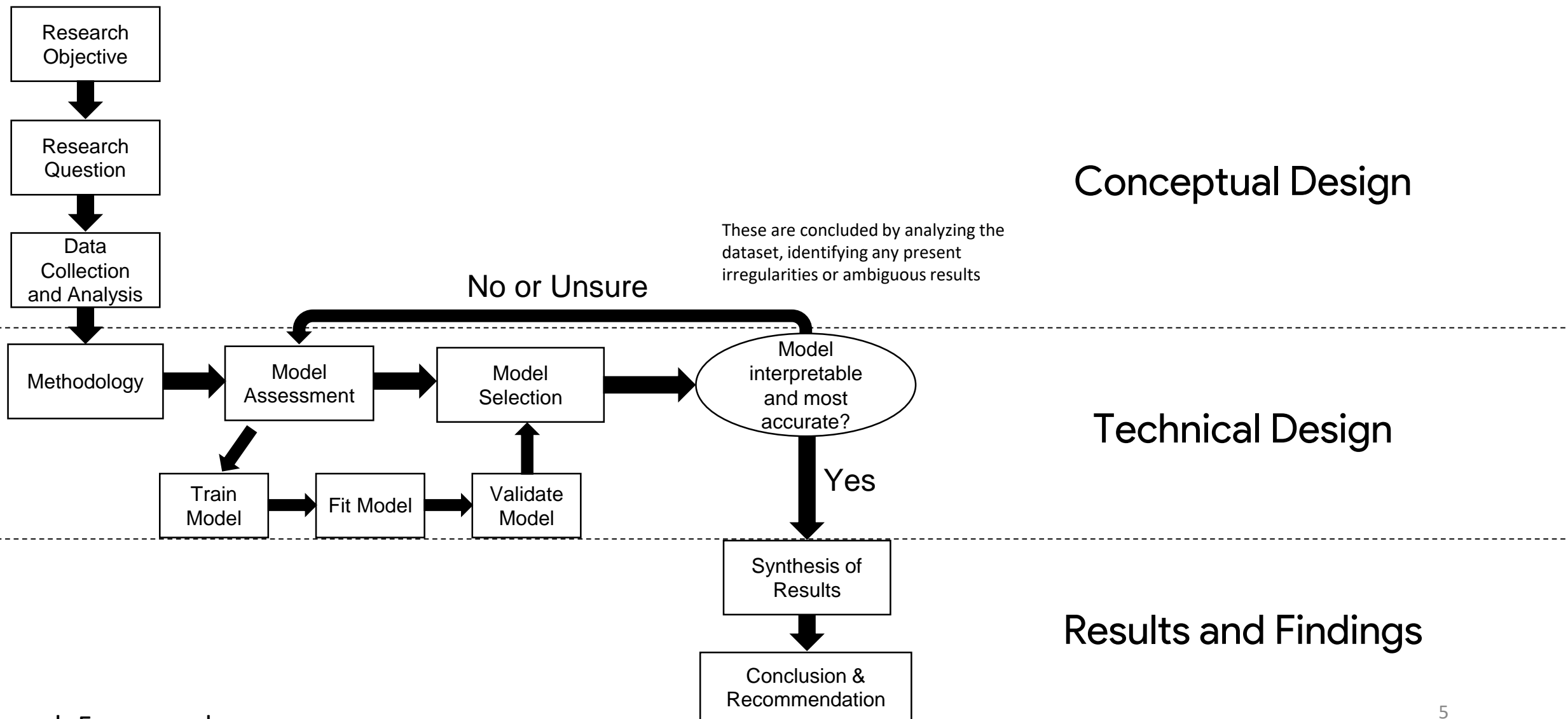Identify the relationship and differences between neighborhood groups

Suggests which factors to improve if you want to increase listing price

Provide insight on potential gaps in listing locations and business opportunities

# Research Framework

# Data Source

**New York City Airbnb Open Data**

https://www.kaggle.com/dgomonov/new-york-
city-airbnb-open-data

Approximately 48,000 rows and 16 columns

**Data Preparation**

**Data Manipulation and Model Assessment**

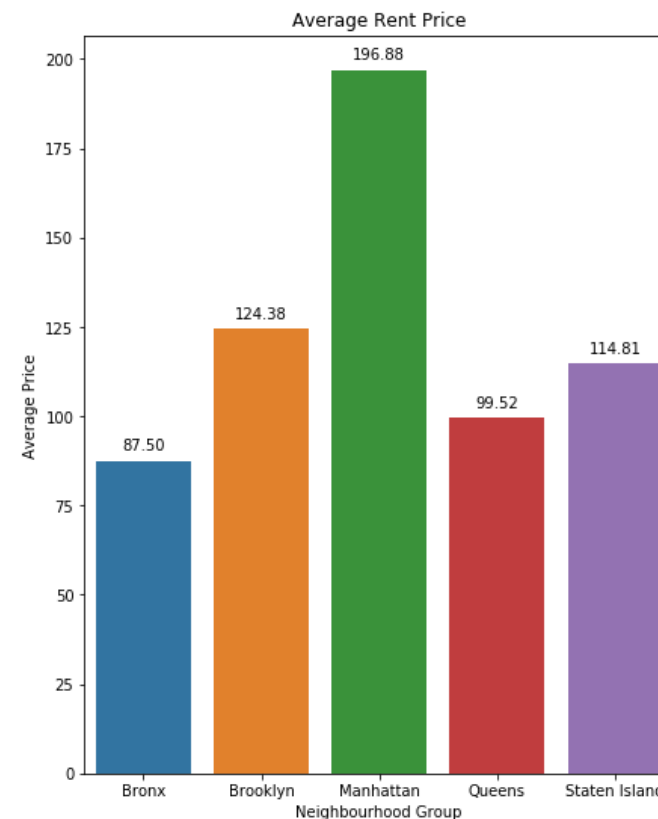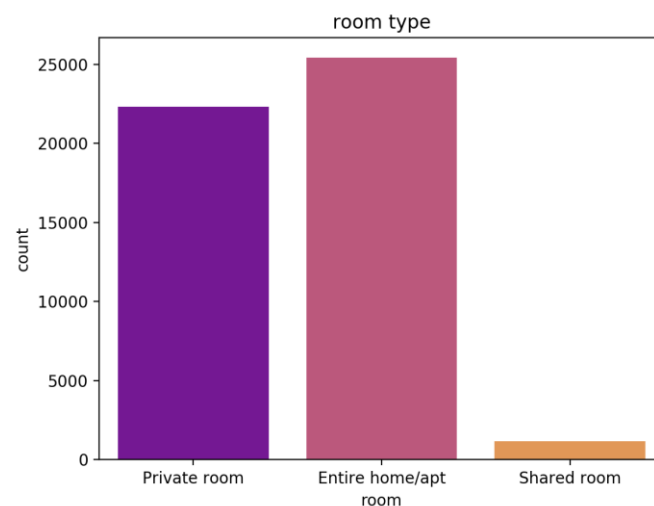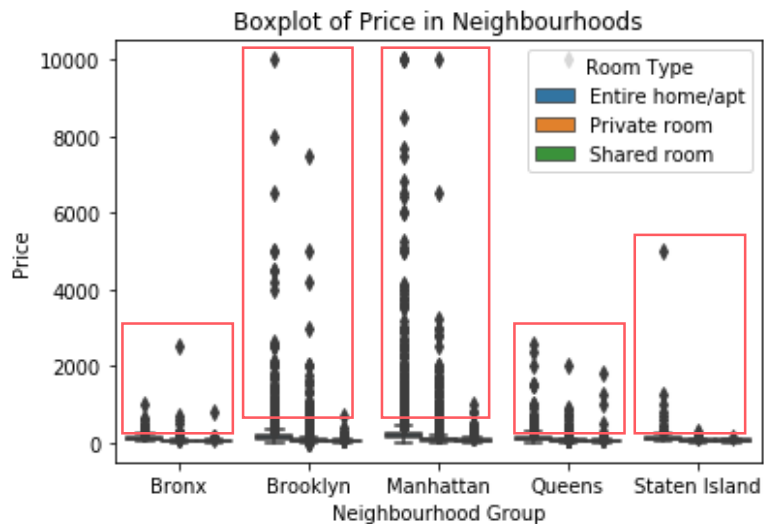| sqlite3 | seaborn |
| pandas | matplotlib |
| numpy | sklearn |

# Initial Data Exploration

Converted data types to its appropriate format

Removed rows with NA values

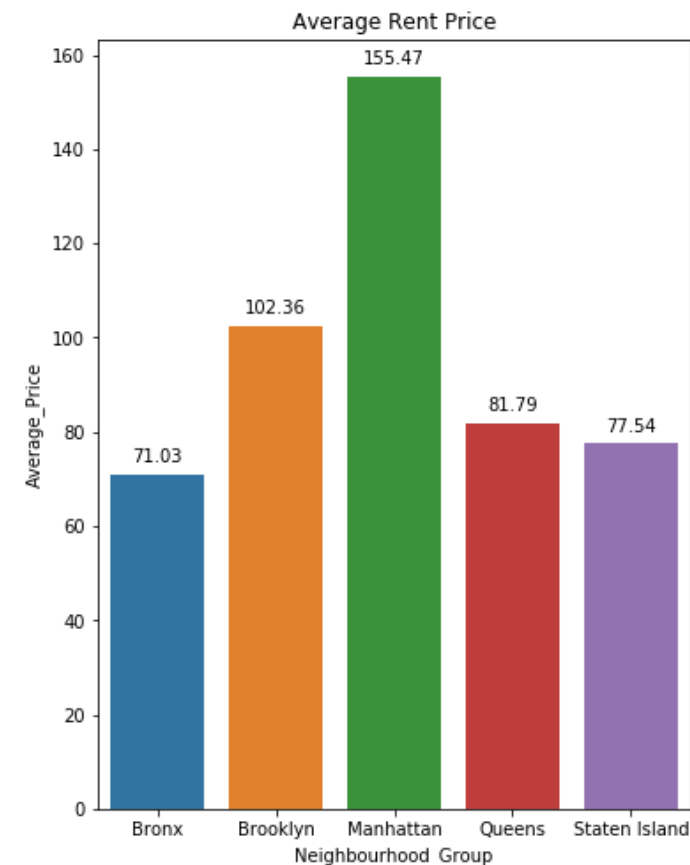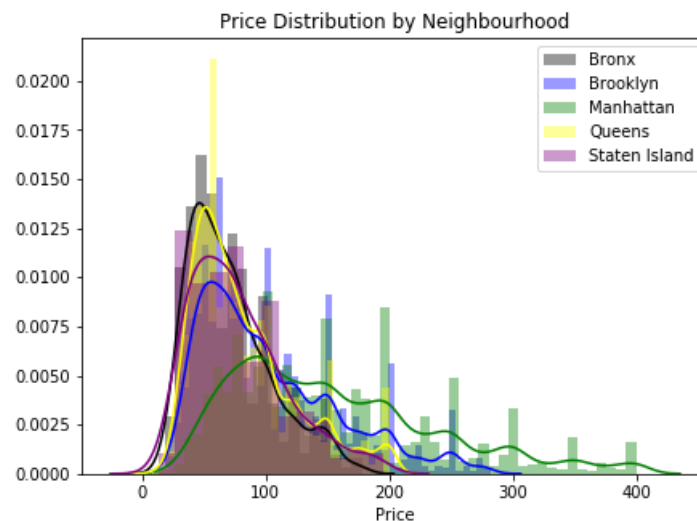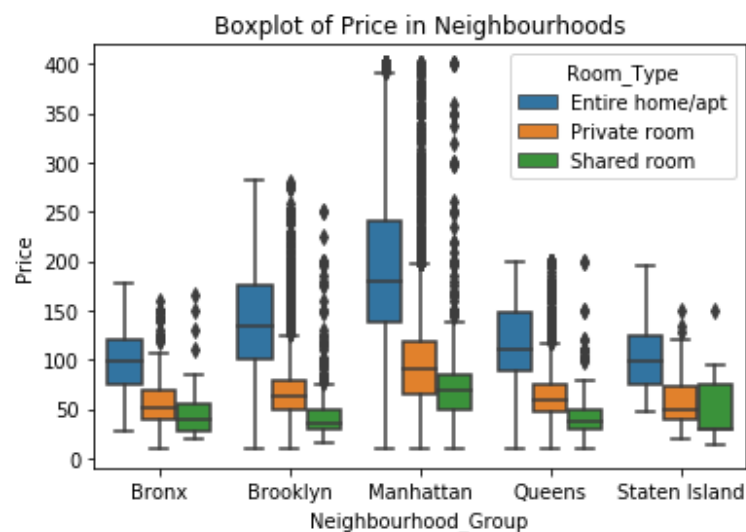Generated basic plots to understand data

# Outliers & Extreme Values

Extreme values present for all neighborhood groups

**Sample 99th percentile of each neighborhood**

Removed rows containing a price of 0 (system error)



Boxplot of Price in Neighbourhoods



Price Distribution by Neighbourhood



Average Rent Price

Preliminary Analysis

# k-Level Categorical Features

'Neighborhood_Groups' = 5 Levels

'Room_Type' = 3 Levels

"One-Hot Encoding"

- Converts multiple level categories into dummy columns

- Binary response (1 or 0)



Features       Dummy Features

# Model Exploration
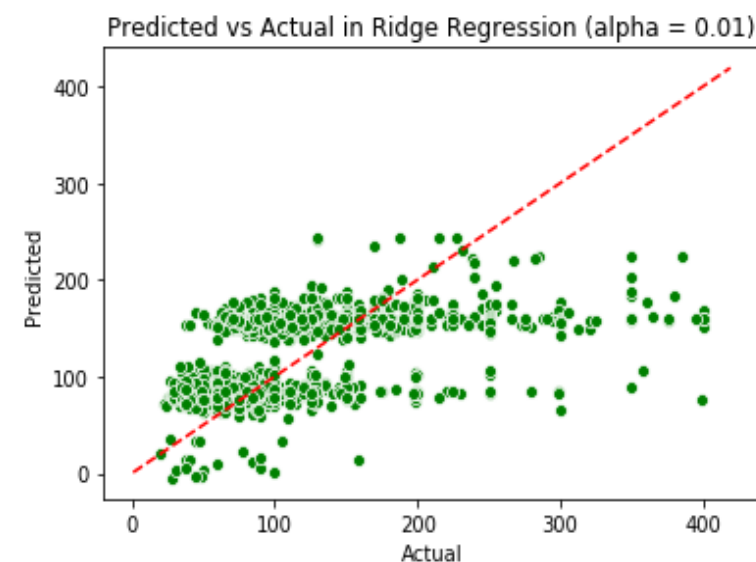
Source: https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

Train Model using **Random Forest**

↓

Cross-Validate & Evaluate

↓

Synthesize Results

Train data using **Ridge Regression**

↓

Cross-Validate & Evaluate

↓

Synthesize Results

Model Fitting

# Regression Models



**Sampled 10% of entire data (~4,000 rows)**

75% Training set
25% Testing set

Model Fitting

# Accuracy and Validation

|              | Accuracy  | R-sq     |
|--------------|-----------|----------|
| **Linear**       | 55.740437 | 0.326937 |
| **Ridge (a=0.01)** | 55.740380 | 0.326938 |
| **Ridge (a=100)**  | 55.138659 | 0.325496 |
| **Lasso**        | 55.700691 | 0.326868 |

No significant differences!

Very similar capability and accuracy

**Next: Ensemble Regressor Methods**

Model Fitting

# Training on Random Forest

**Due to computing power constraint**

**Sampled 10% of entire data (~4,000 rows)**

75% Training set
25% Testing set

Ran model using standard &
automatic predictors/parameters
- 1000 trees

**Accuracy: 62.73%**
- Accuracy is calculated by 100 - MAPE

We'll try to improve it

Model Fitting

# Variable Importance in RF

# Hyperparameter Tuning

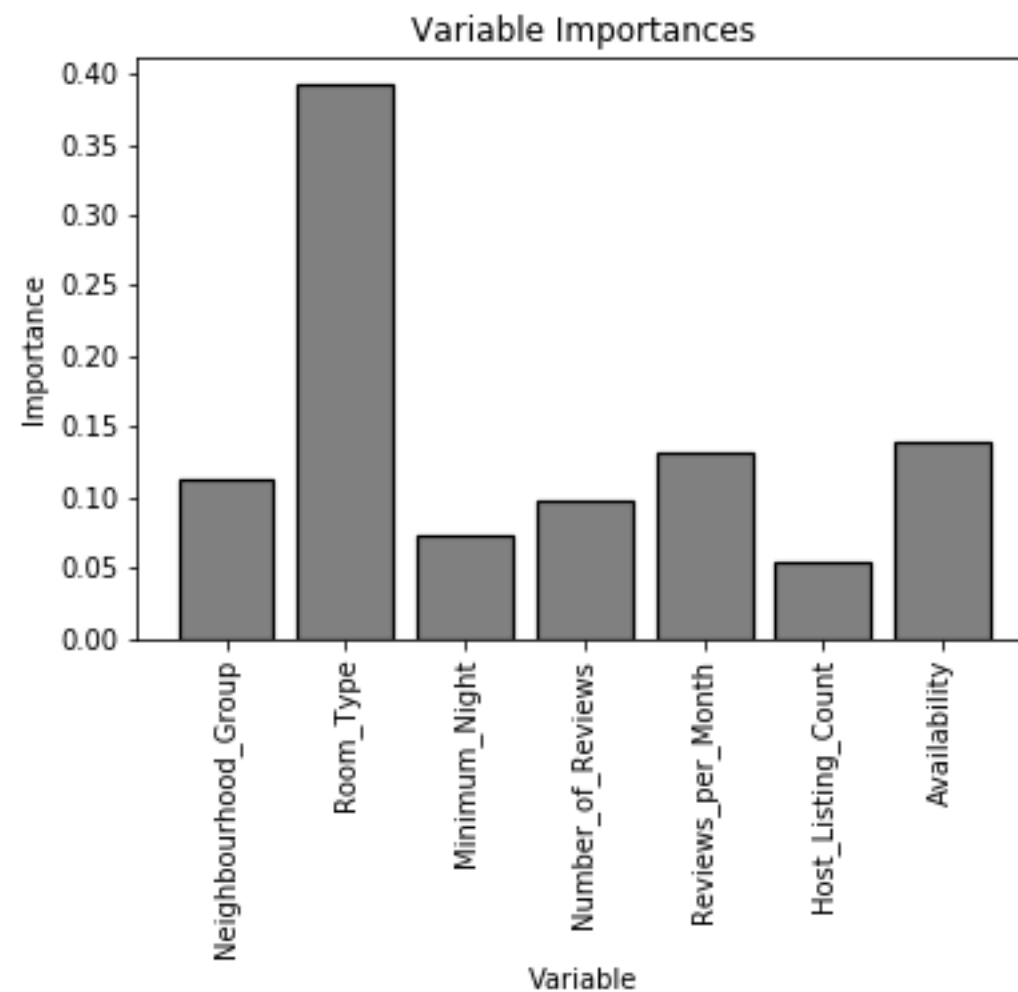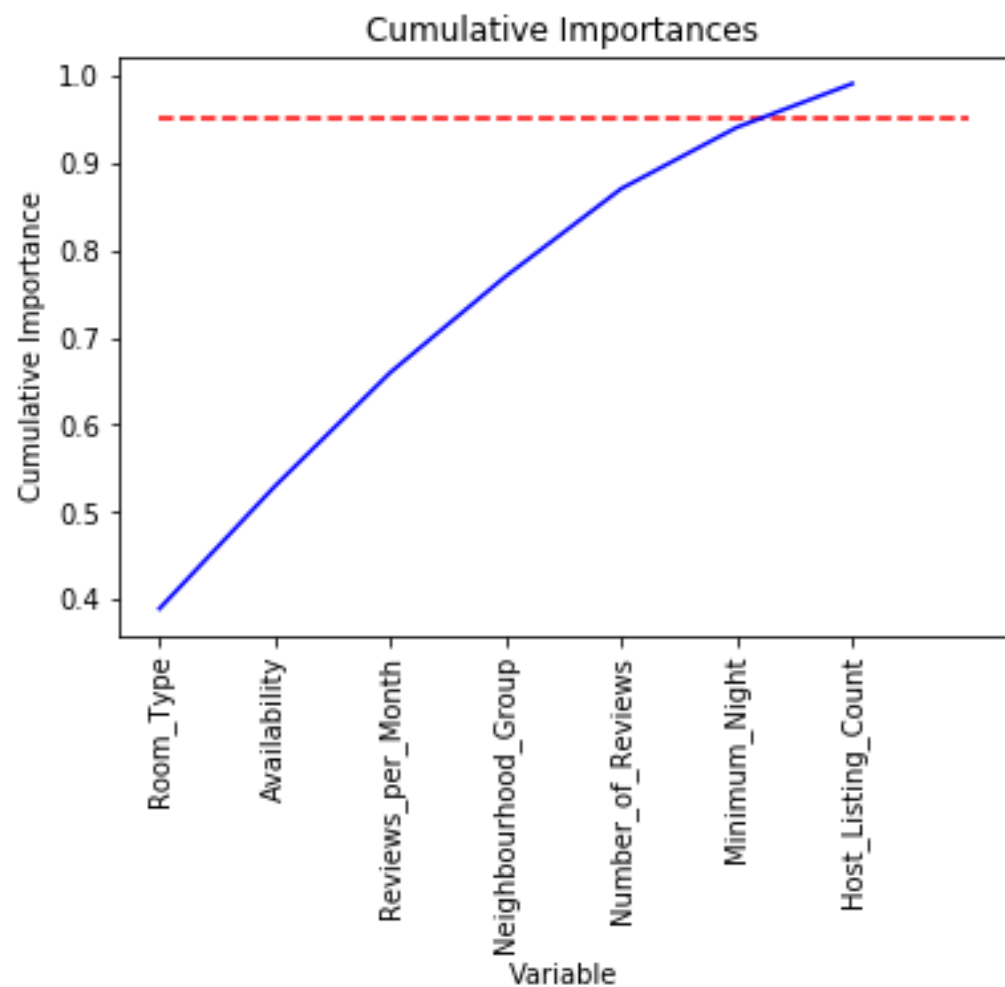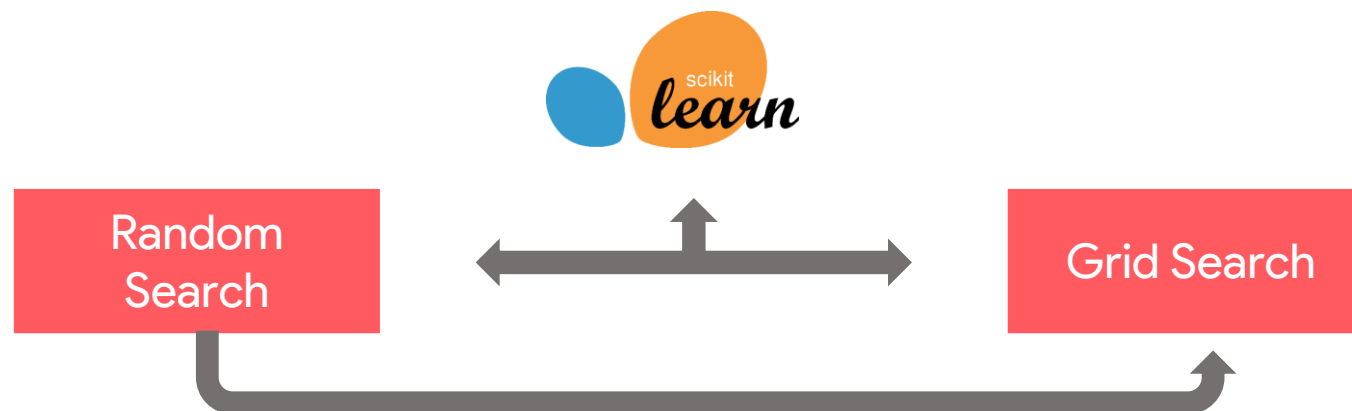"It is a parameter whose value is set before the learning process begins"

scikit learn

| Random Search | | Grid Search |

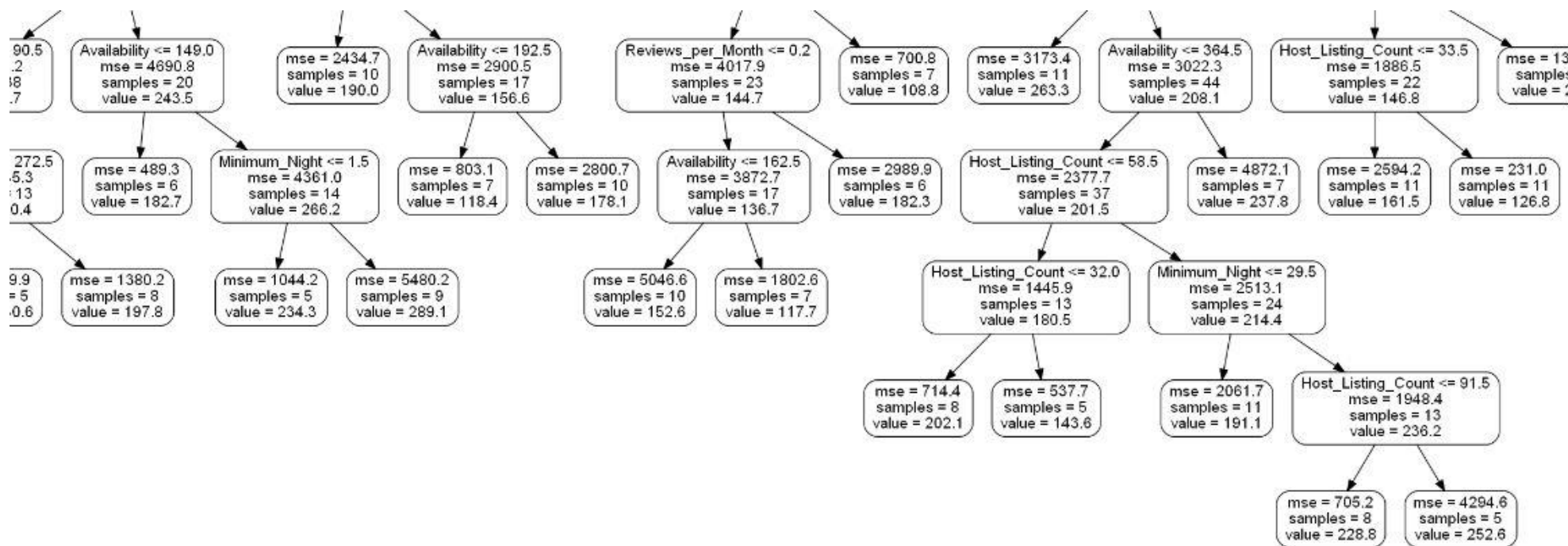**Train model with random combinations of the parameters and see which iterations/set is the best**

**Utilize Random Search to narrow down parameters for Grid Search**

**Using the selected optimized parameters obtained from Random Search, re-train model using every single combination of hyperparameter values**

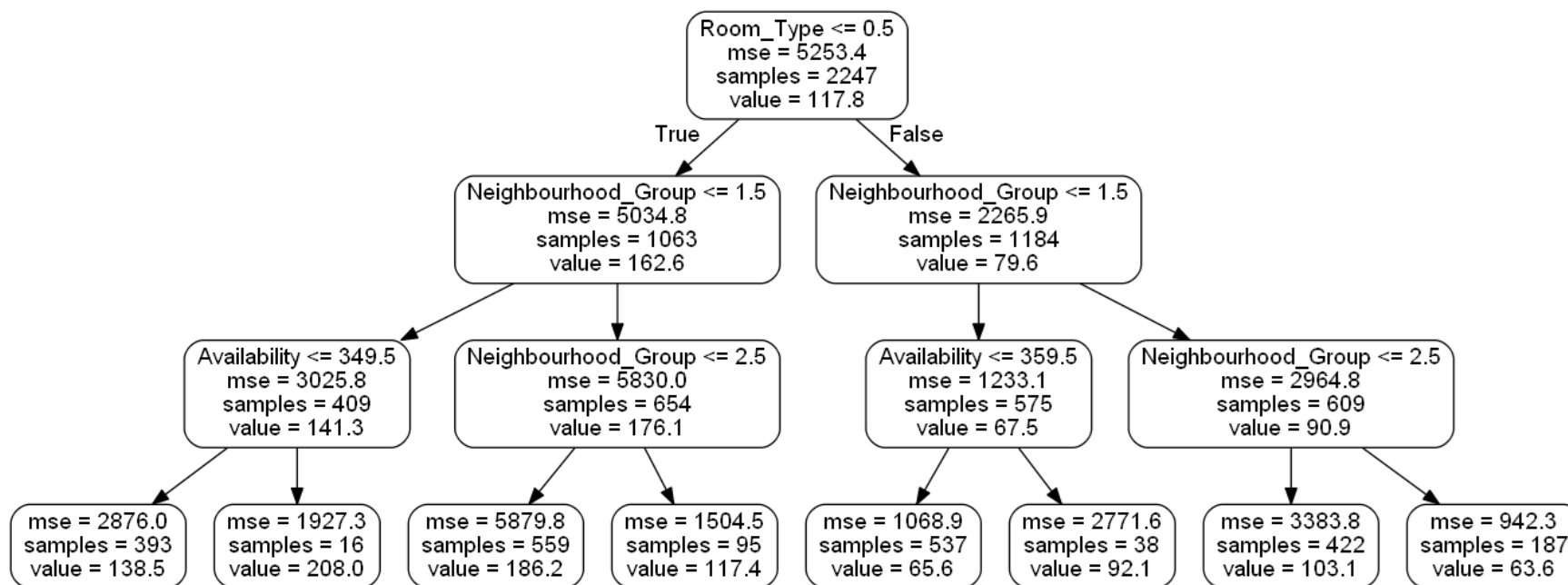**Best Grid RF Accuracy: 64.89%
Improvement of 2.16%**

Model Validation

# Optimized Tree Diagram

# Pruned Tree



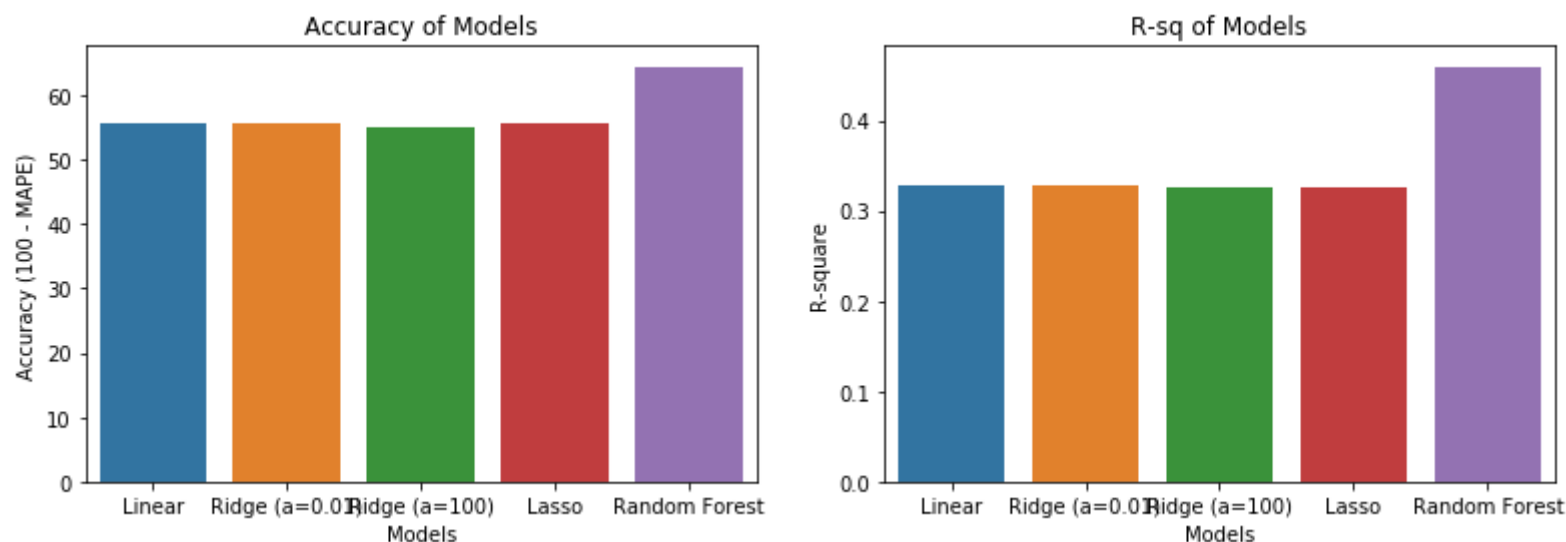Number of Trees = 10
Max Depth = 3

**Accuracy: 62.53%**
**R-sq: 43.99%**

**Easier to interpret despite the slight loss of accuracy**

# Synthesis of Results



**Random Forest** yields the best **Accuracy** and **R-sq** compared to the regression models

There is no significant differences between the regression models

Best Model Accuracy & R-sq:
- **Accuracy: 64.89%**
- **R-sq: 45.96%**

Results

# Conclusion

More features is needed to increase model accuracy

Room type has the highest impact in determining the listing price of Airbnb in NYC

Hosts does not necessarily increase their price based on higher reviews
- Having more reviews doesn't mean the listing will be more expensive
- Number of reviews is a small factor which means price is inelastic to reviews

airbnb ➡ **Utilize model to provide pricing suggestions for their customers and business partners**
**Given the basic features (room type, location, etc), they can predict future outcome**

Results