# Predicting Airbnb Rental Prices in the U.S.

T11: Minjie Bao, Yoojin Kathleen Jeong, Ria Stephanie Pribadi, Prerna Sharma, Alix Vermeulen, Qianyao Ye
TA: Karen Figueroa

## Results

- Predicted prices with three models: linear regression, multivariate adaptive regression splines, and XGBoost.
- XGBoost has the best performance with the highest R square and lowest RMSE and MAE.
- Room types and city are the key predictors.

## Background

Airbnb allows individual homeowners to rent their properties. With a wide variability of locations and price indicators, hosts often find themselves wondering if the price they offer is considered fair.

Offering right prices is very crucial for hosts, because Airbnb rental market is competitive. This project analyzed the factors that may have significant roles in affecting prices, which will provide comprehensive insights about Airbnb rental price valuation across the U.S.

## Data

Data is sourced from:
https://www.kaggle.com/datasets/kritikseth/us-airbnb-open-data.
- Contains 17 columns and 226030 rows.
- Key variables include host id, price, room type, city, and number of reviews.
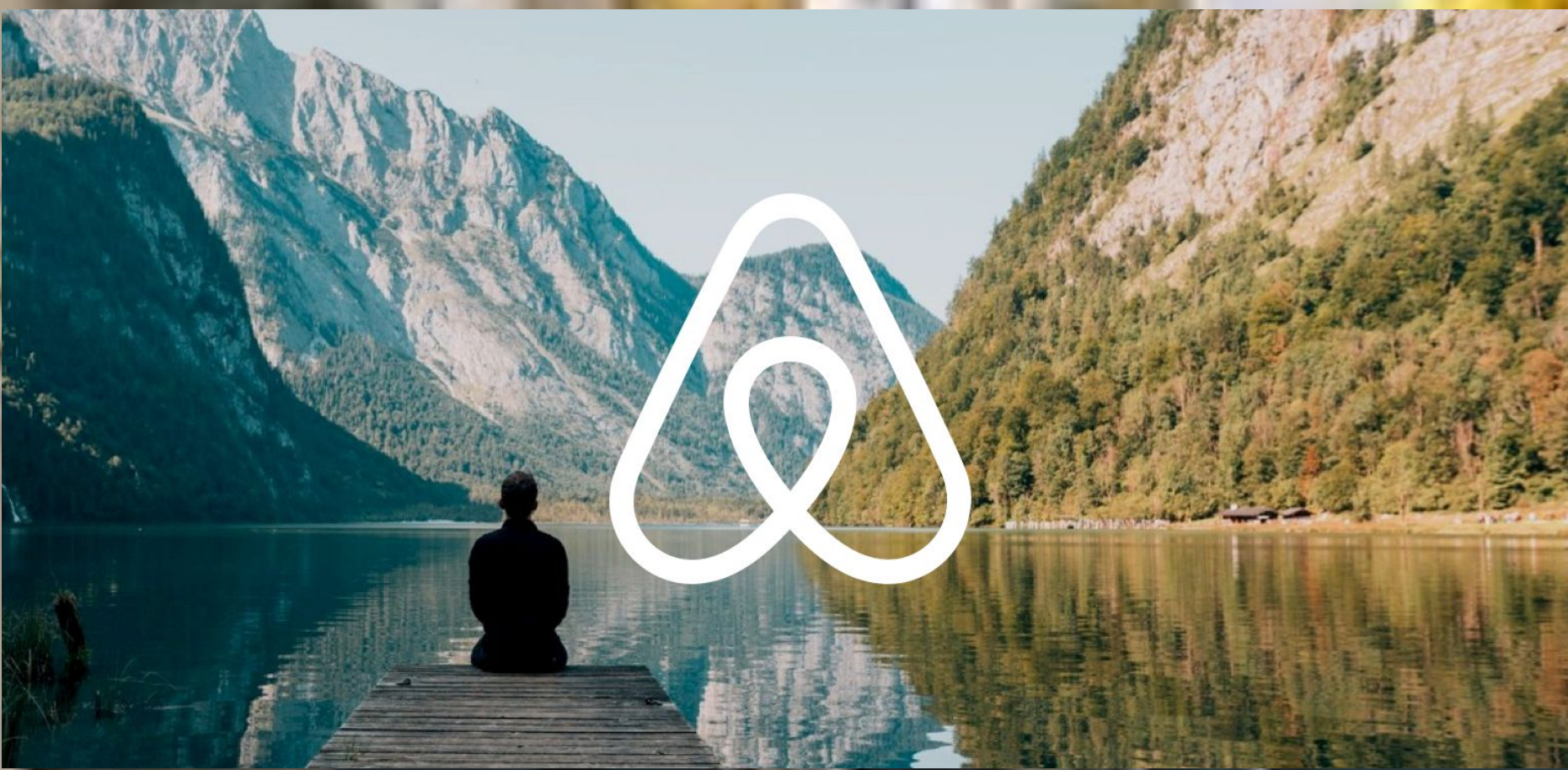- Gives insights regarding 2020 pricing as per location in the United States.

## Business Impact

- Pricing is one of the most important decisions that hosts need to make to gain profits.
- The prediction exercise here provides a method to predict prices for properties around the same area with similar conditions.
- It will give the hosts a better idea of the property's market values.
- Airbnb can retain hosts, expand its business, and eventually increase revenues.
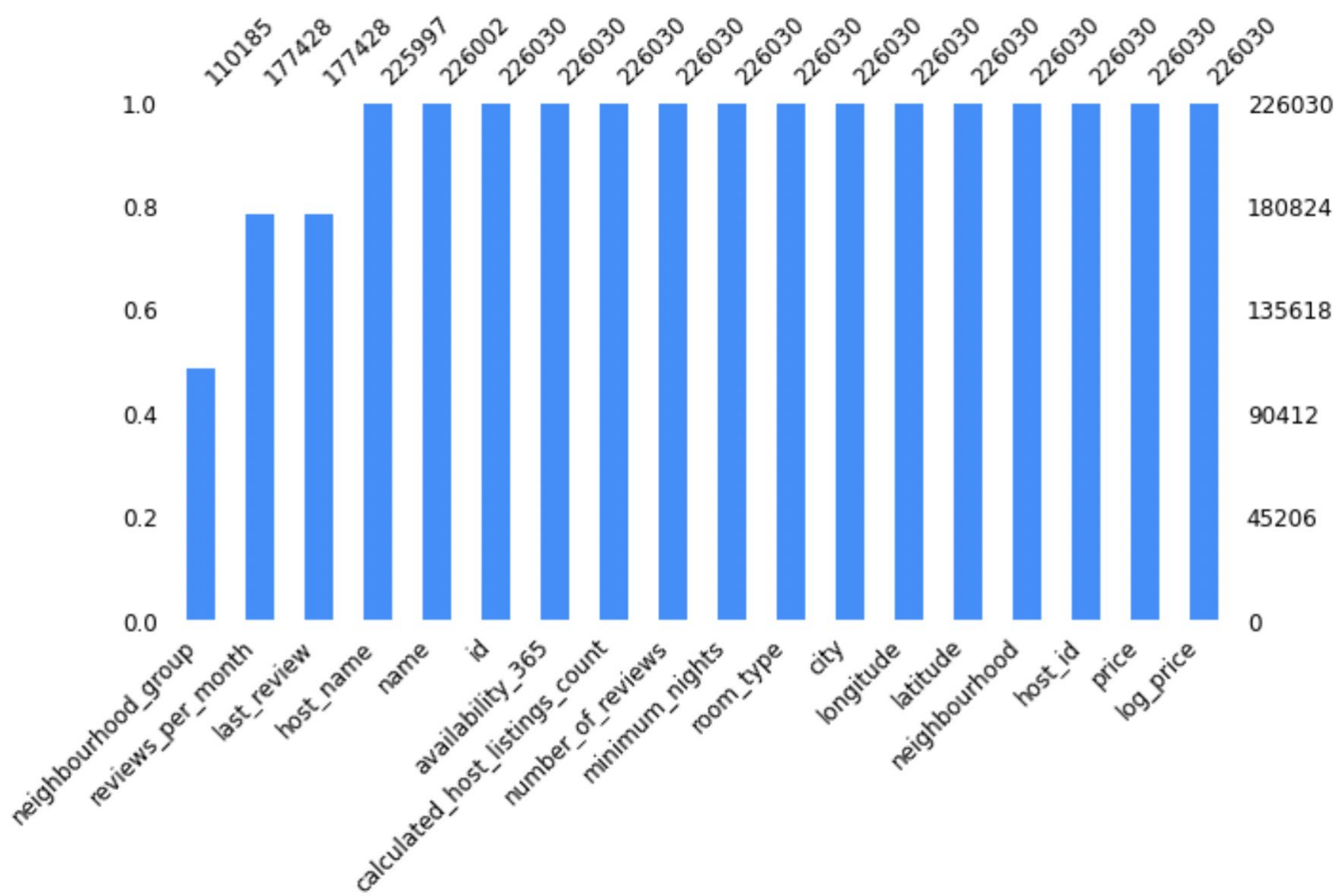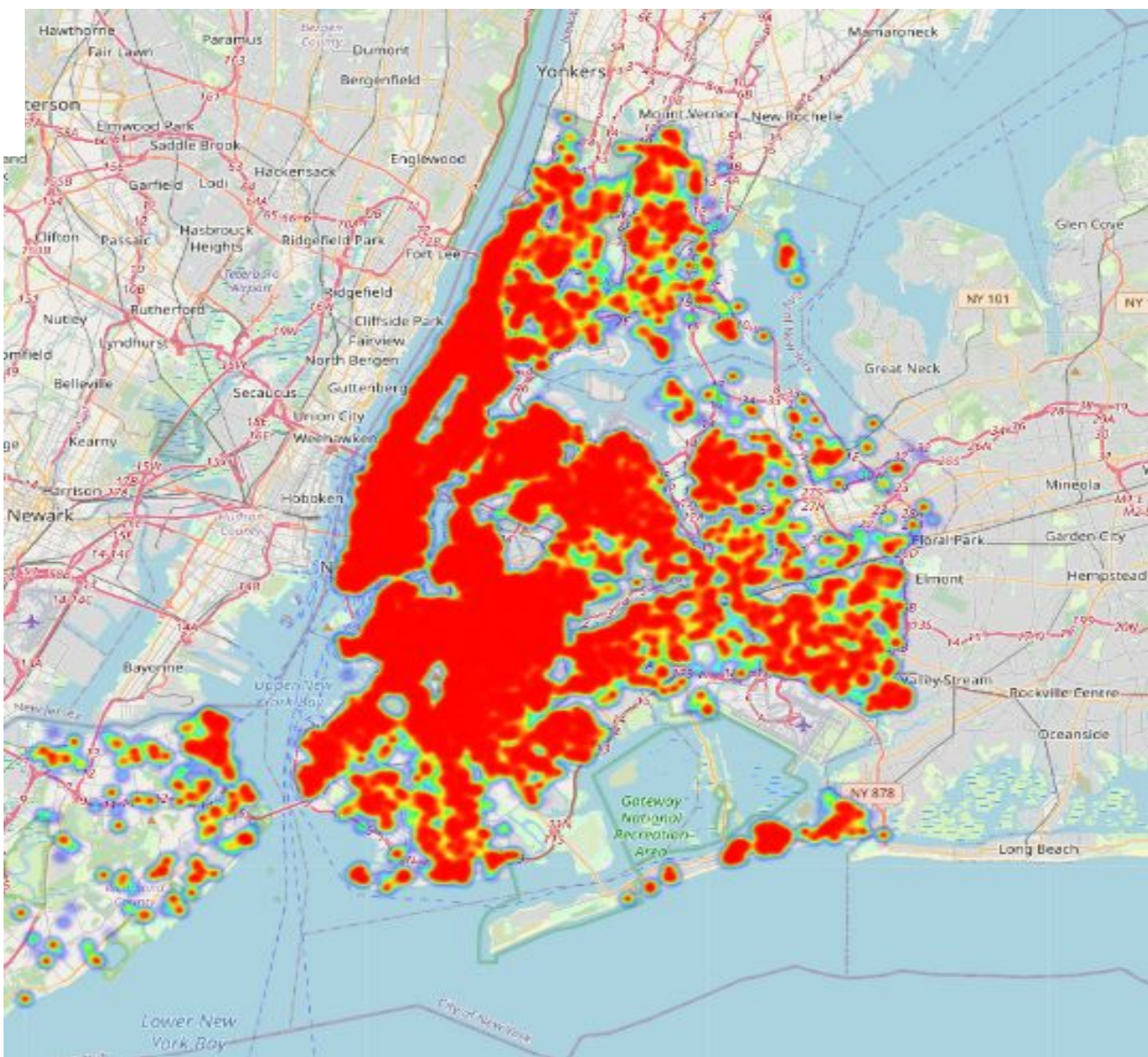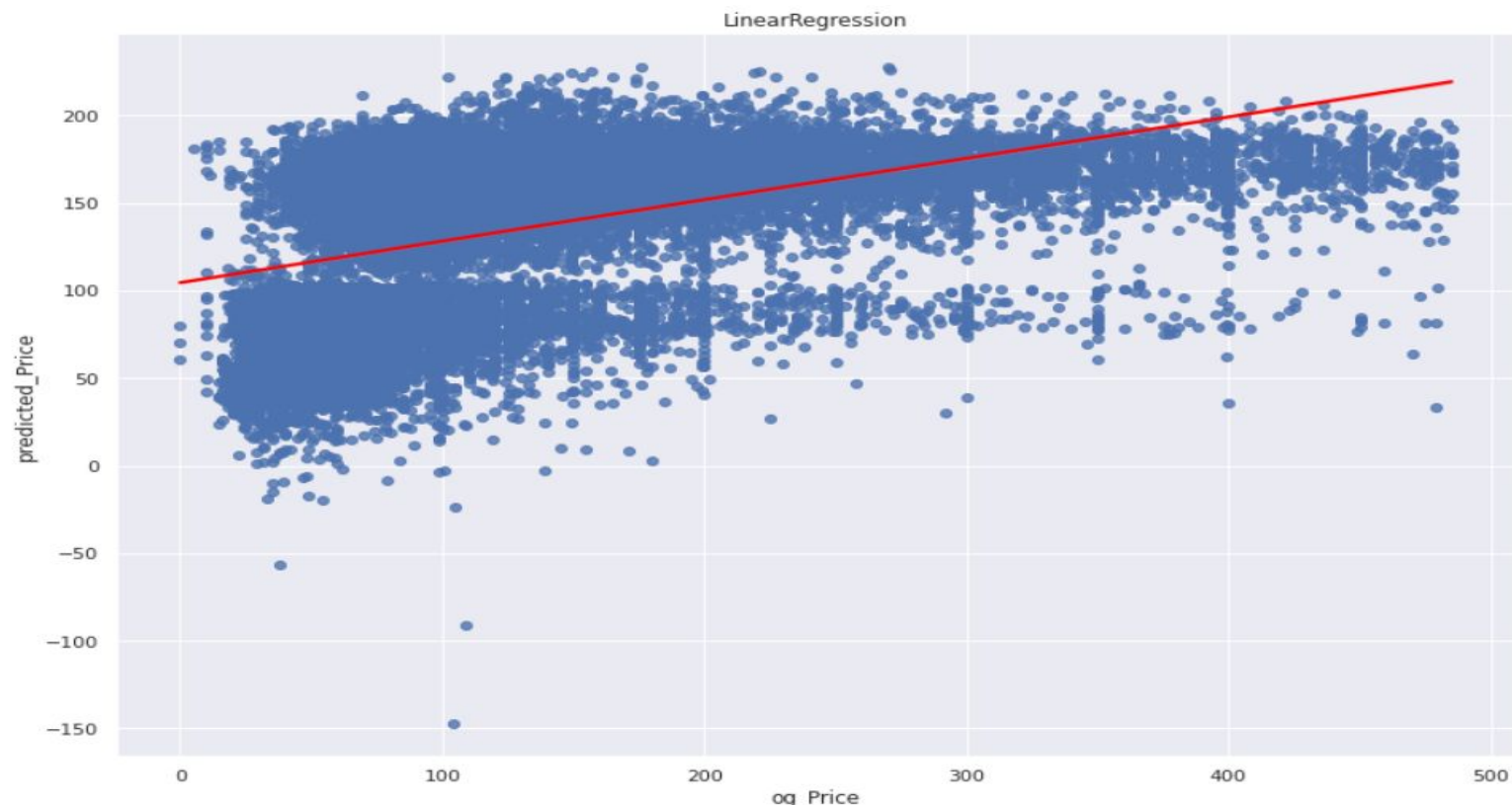
## Methodology

- Data preprocessing: Drop missing data and remove outliers
- EDA:
  - Create histograms and correlation matrix for numeric features
  - Visualize listings and prices distributions across states
- Data Modeling:
  - Use linear regression, XGBoost and multivariate adaptive regression splines to predict prices
  - Compare the model performances by RMSE, MAE and R2 Score

## Visualizations

Missing data histogram:



Model results table:

| Model | RMSE | MAE | R2 Score |
| --- | --- | --- | --- |
| Linear Regression | 79.3201 | 58.2739 | 0.2364 |
| MARS | 79.0596 | 58.0855 | 0.2414 |
| XGBoost | 71.2464 | 50.1792 | 0.3840 |

Linear regression scatter plot:



Highest Prices in New York map:



Correlation Matrix: