

JUNE, 2020 • CAPSTONE PROJECT

# Airbnb Price Prediction

Aisha Awan  
Daniela Collaguazo  
Darien Thayne  
Jessica Saenz  
Mulyono Kertajaya

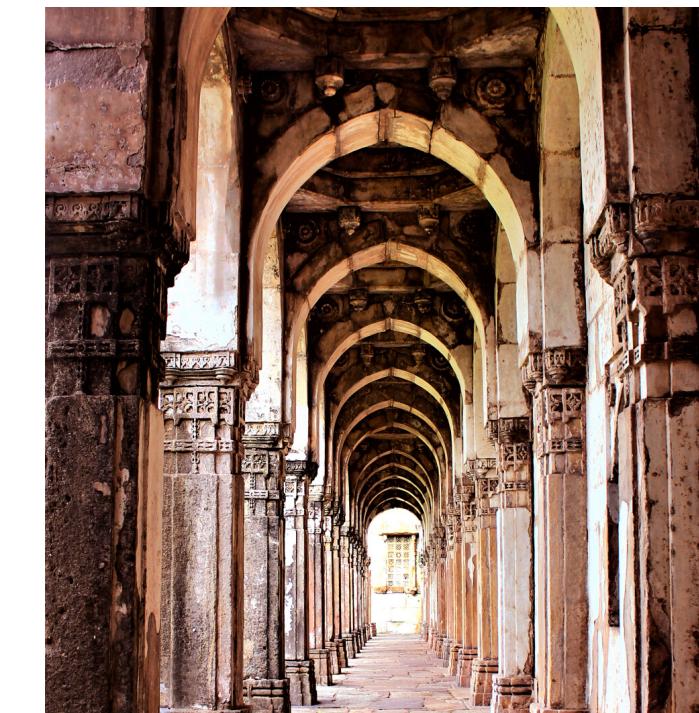
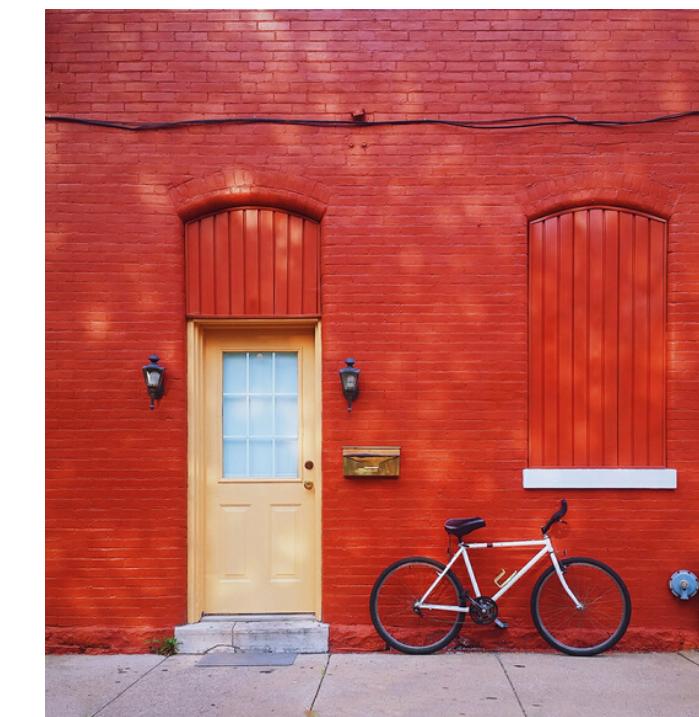


Georgetown Data Science Certificate

*GEORGETOWN UNIVERSITY*  
School of Continuing Studies



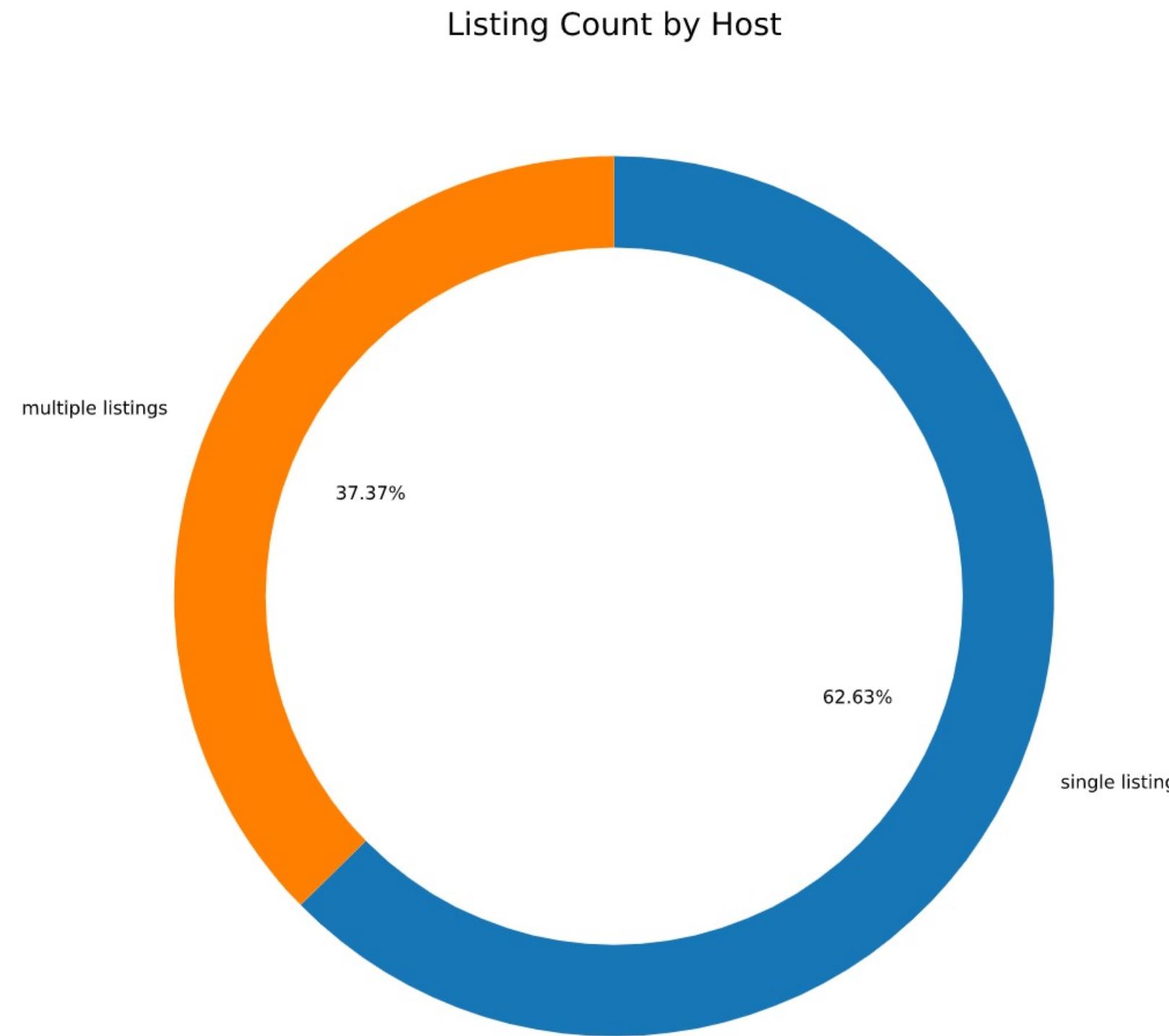
Jun, 2020 • Airbnb Prediciton Price



Airbnb began in 2008, it is a U.S based online marketplace that allows people to rent their properties or spare rooms to guests around the world.

The company's business model is to act as a broker that receives commission from each booking from both hosts and guests.

# Airbnb listing is now filled with professional hosts



We are trying to  
help new  
individual listers

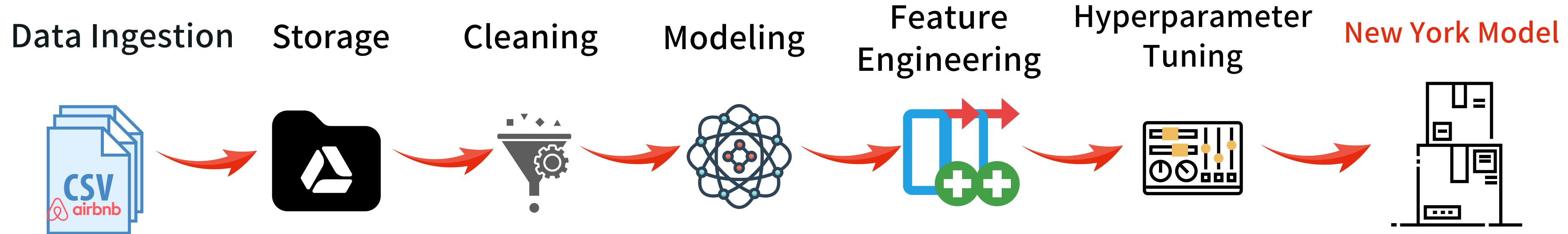


# Hypothesis

By understanding the type of property, location, booking availability per year, together with an analysis of the accommodations, we will be able to predict a price per night that will optimize revenue for a new Airbnb host.

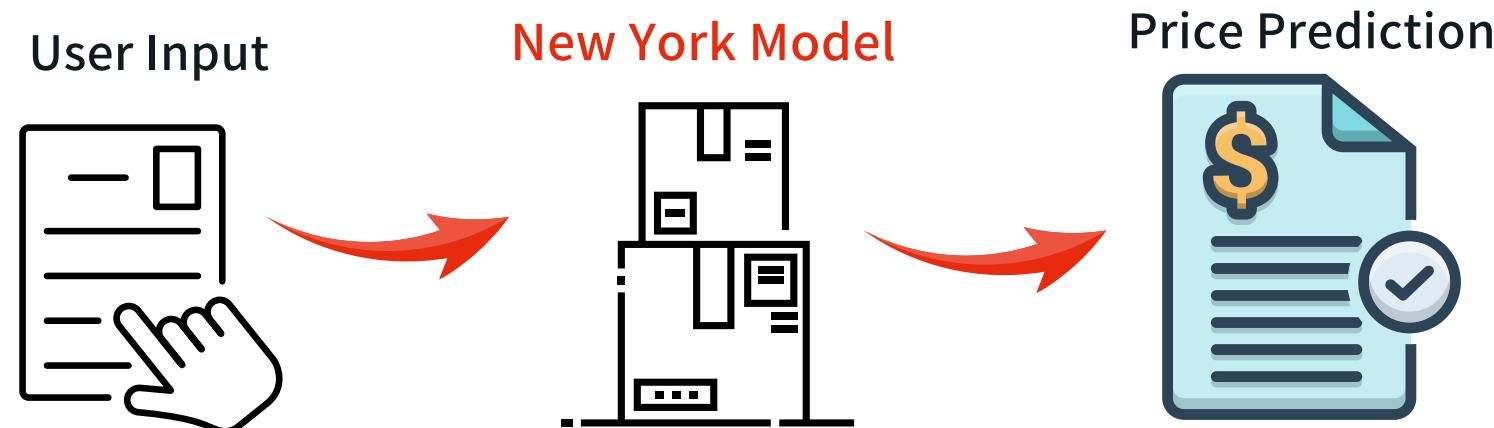
# The Architecture

## BUILD PHASE



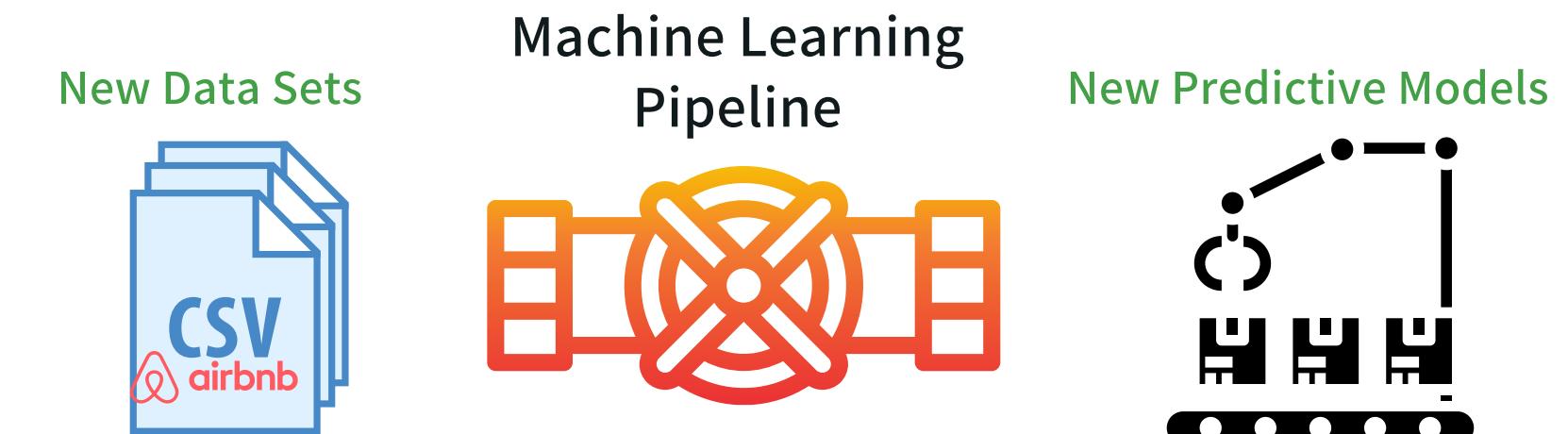
## OPERATIONAL PHASE (Option 1)

### Flask



## OPERATIONAL PHASE (Option 2)

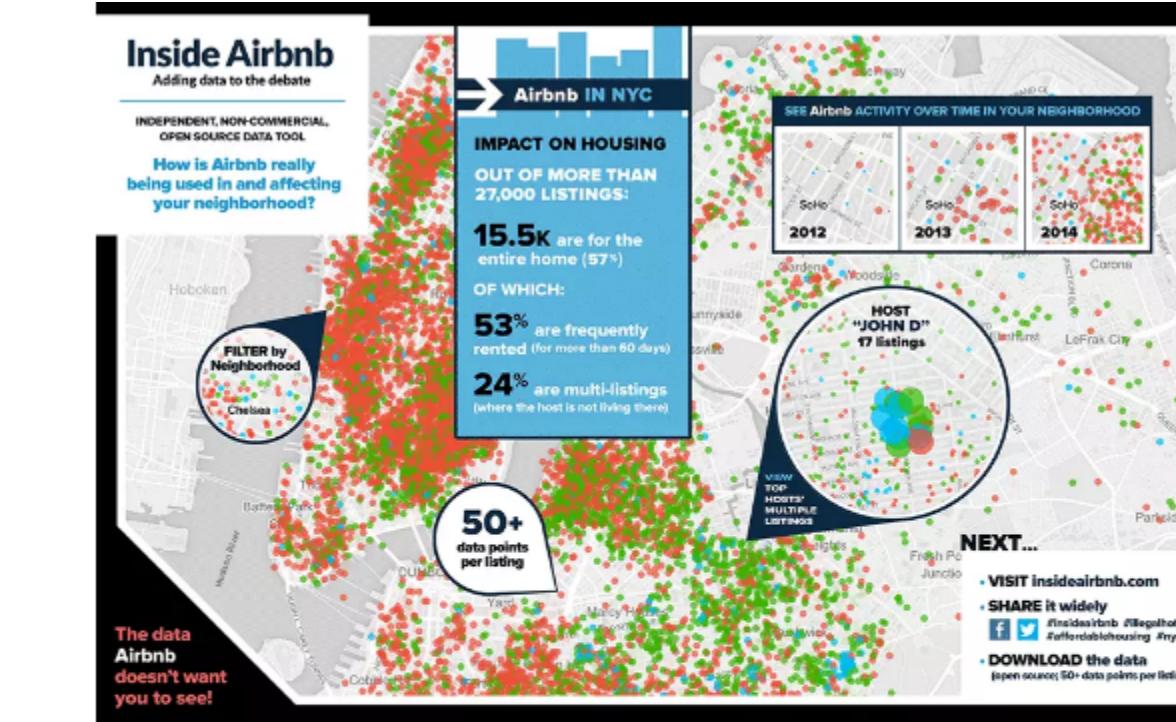
### Pipeline





# Data ingestion & Data Acquisition from AIRBNB OPEN DATA

- Listings -- this csv has 105 columns and 50k records for NY
- Reviews -- this csv has reviews by listing
- Calendars -- this has prices up to 1 year in the future



## New York City. Adding data to the debate.

Inside Airbnb is a set of independent tools and open data that allows you to explore how Airbnb is REALLY being used in cities around the world.

 Inside Airbnb



# Storage



You Retweeted

**Inside Airbnb** @InsideAirbnb · 04 Apr

Dear @InsideAirbnb community.

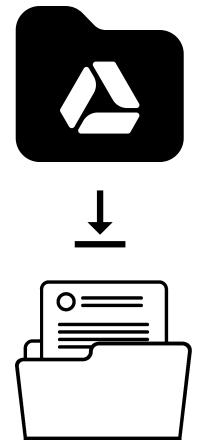
Between Jan and Mar, downloads (10TB) and associated costs increased by 500%, to an unsustainable level.

Downloads have been blocked for the time being. For those of you downloading bulk data, please stop. Or make a donation.



## DATA STORAGE

- Originals pulled from insideairbnb and stored on Google Drive
- Data pulled via code from Google Drive and stored locally
- Cleaned Data stored within separate cleaning stages during initial build stage.
- Simplified model is stored on flask hosting site.



# Data ingestion



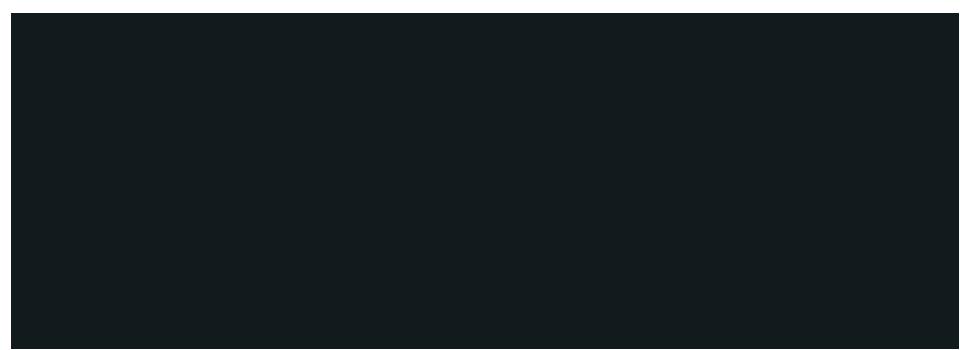
## REVIEWS

- Only used for data exploration and trying out text analysis but didn't work out.

## CALENDAR

We were planning to use this to see if the prices go up for different seasons or holidays.

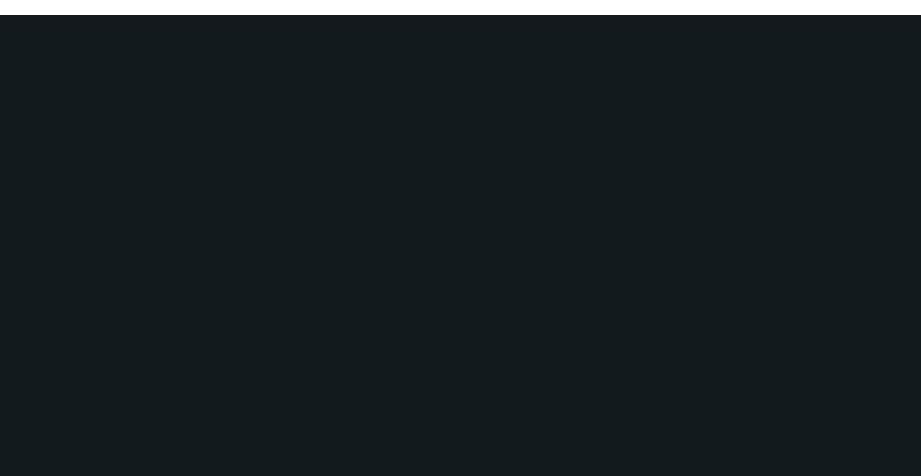
But then we realize the dates are for future dates. Most users probably don't setup their prices that advance.



## LISTINGS

- All listings that have reviews in last 1 year.
- Exclude records that are type of hotel, boutique hotel, resort, etc.
- Limit the price to between \$10 and \$800 to deal with outliers.

This listing.csv is the main data set





# Listings data

50796 rows × 105

columns

- Columns with same values
- Duplicates values
- Nan values
- Pictures and thumbnail columns

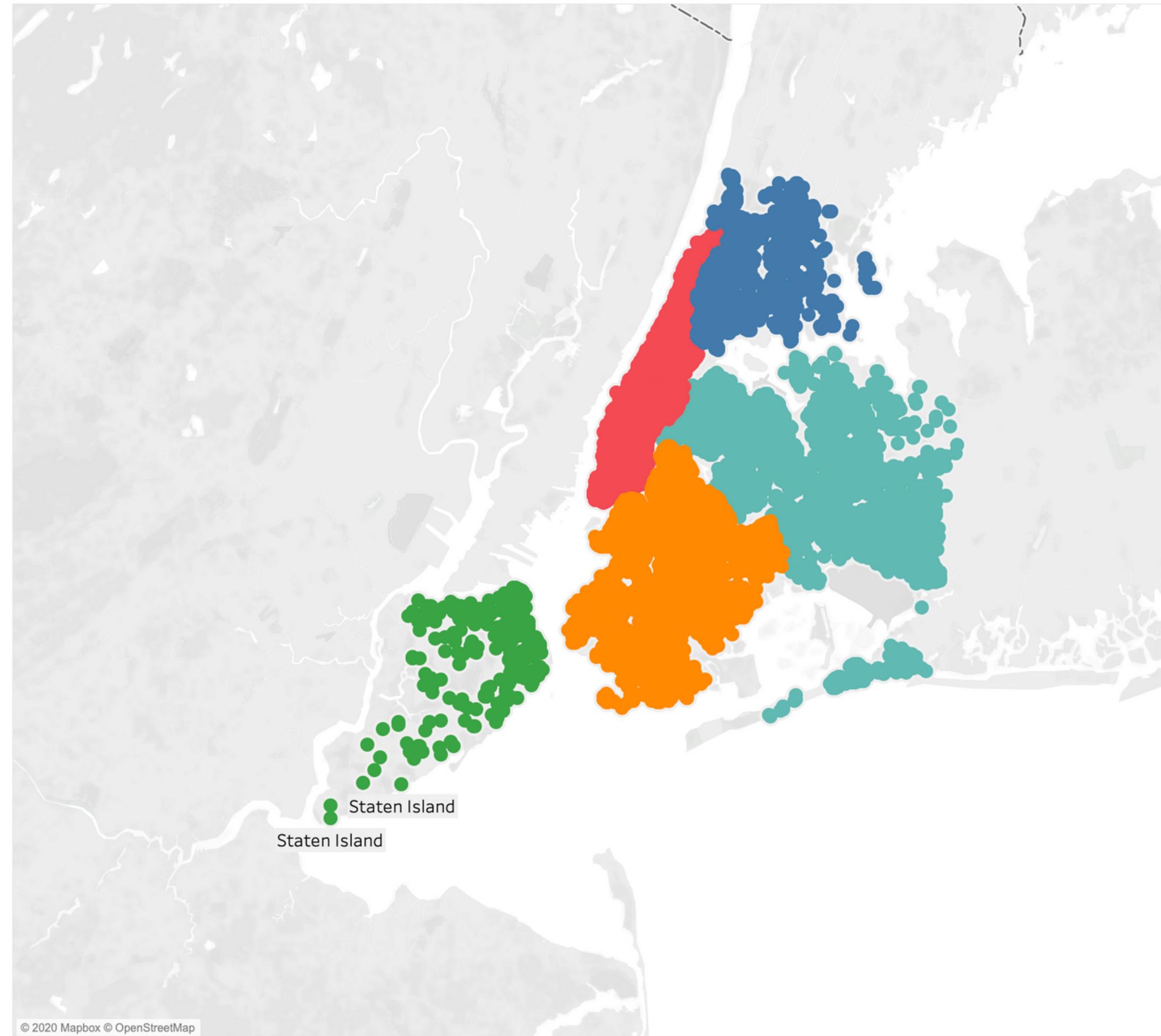
'id', 'listing\_url', 'name', 'host\_id', 'host\_url',  
'host\_name', 'host\_since', 'host\_location', 'host\_response\_time',  
'host\_response\_rate', 'host\_acceptance\_rate', 'host\_is\_superhost',  
    'host\_neighbourhood', 'host\_listings\_count',  
    'host\_total\_listings\_count', 'host\_verifications',  
    'host\_identity\_verified', 'street', 'neighbourhood',  
'neighbourhood\_cleansed', 'neighbourhood\_group\_cleansed', 'city',  
    'state', 'zipcode', 'market', 'smart\_location', 'latitude', 'longitude',  
    'is\_location\_exact', 'property\_type', 'room\_type', 'accommodates',  
    'bathrooms', 'bedrooms', 'beds', 'bed\_type', 'amenities', 'price',  
    'security\_deposit', 'cleaning\_fee', 'guests\_included', 'extra\_people',  
'minimum\_nights', 'maximum\_nights', 'minimum\_minimum\_nights',  
    'maximum\_minimum\_nights', 'minimum\_maximum\_nights',  
    'maximum\_maximum\_nights', 'minimum\_nights\_avg\_ntm',  
'maximum\_nights\_avg\_ntm', 'calendar\_updated', 'availability\_30',  
    'availability\_60', 'availability\_90', 'availability\_365',  
    'number\_of\_reviews', 'number\_of\_reviews\_ltm', 'first\_review',  
    'last\_review', 'review\_scores\_rating', 'review\_scores\_accuracy',  
    'review\_scores\_cleanliness', 'review\_scores\_checkin',  
    'review\_scores\_communication', 'review\_scores\_location',  
'review\_scores\_value', 'instant\_bookable', 'cancellation\_policy',  
    'calculated\_host\_listings\_count',  
    'calculated\_host\_listings\_count\_entire\_homes',  
    'calculated\_host\_listings\_count\_private\_rooms',  
'calculated\_host\_listings\_count\_shared\_rooms', 'reviews\_per\_month'

New York

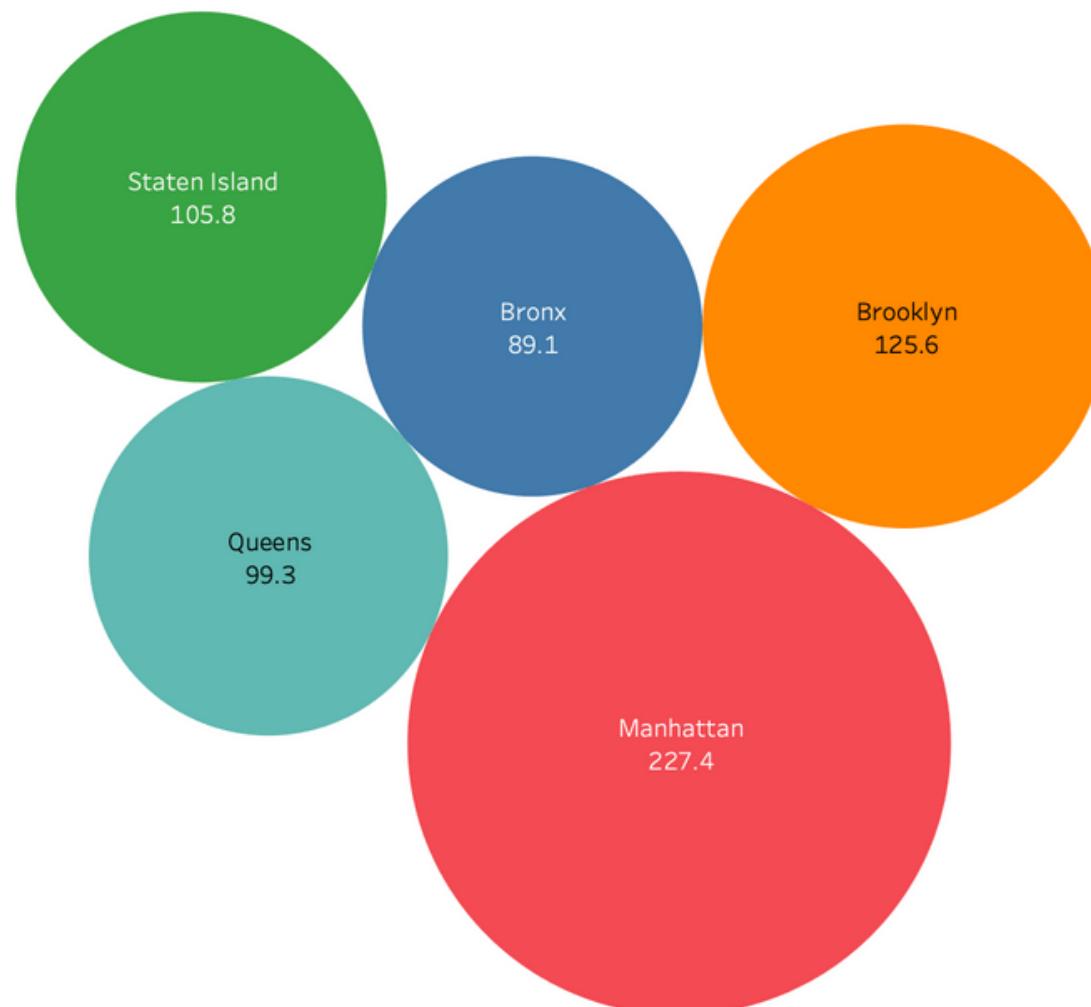
# Listings by Neighborhood

Neighbourhood Group

- Bronx
- Brooklyn
- Manhattan
- Queens
- Staten Island

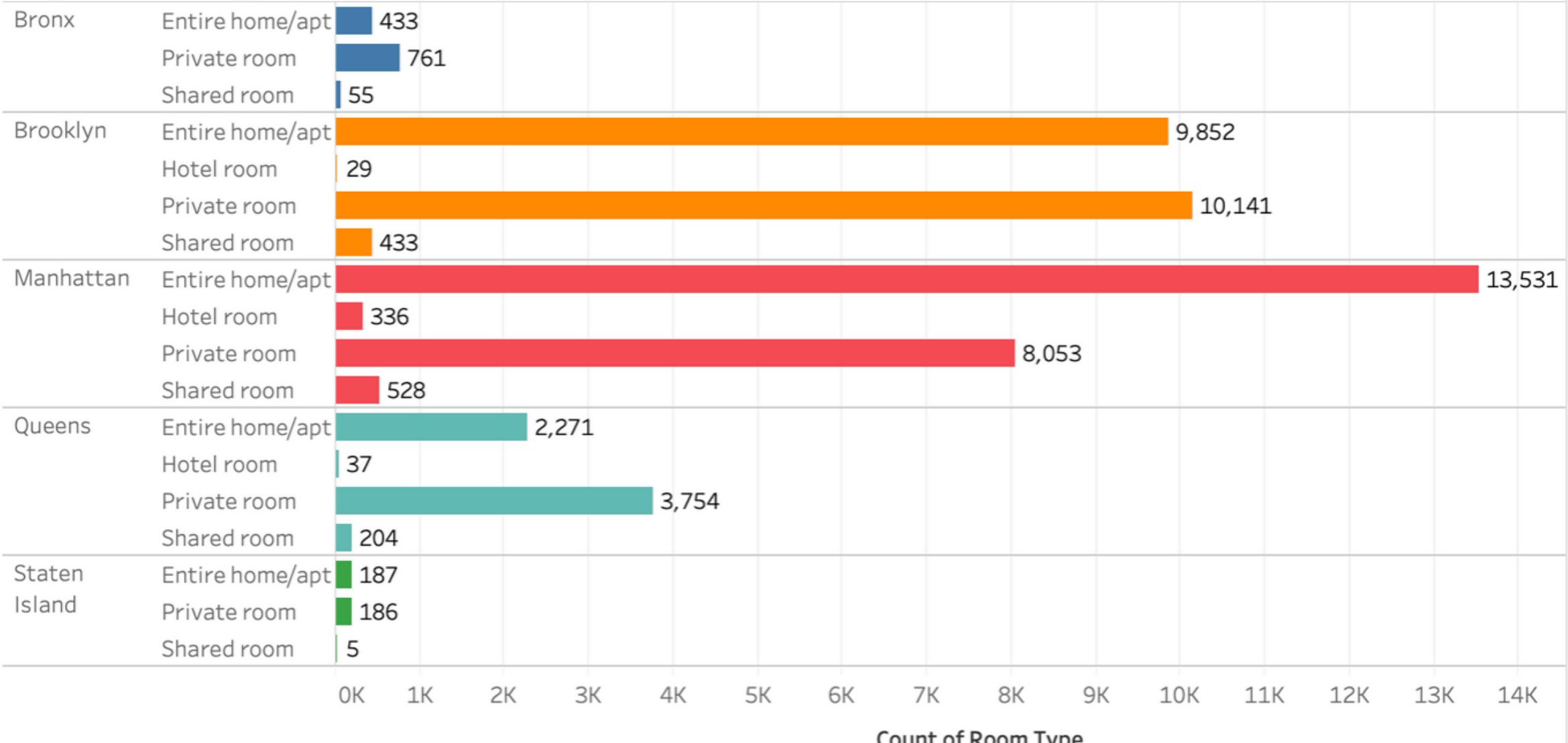


Average Price



Room Type

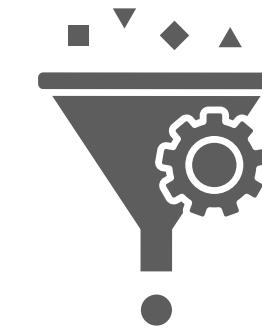
Neighbour.. Room Type



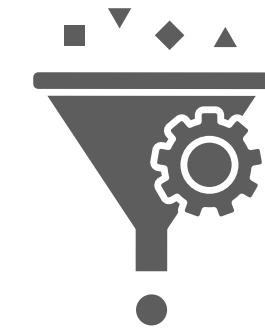
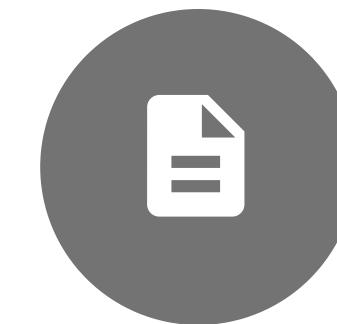
# Ingestion & Cleaning Data



- **Columns Removed**
  - < 50% Populated
  - 95% + rows with one value
  - Picture / Thumbnail URLs
  - Remove object fields
  - Remove other parts we didn't plan on using
- **Rows Removed**
  - No reviews within a year (non-active)
  - High price outliers (Hotels primarily)
- Features Columns added from Amenities text.
- Basic conversions (Boolean, Date, Category)



# Ingestion & Cleaning Data



## CLEANED DATA

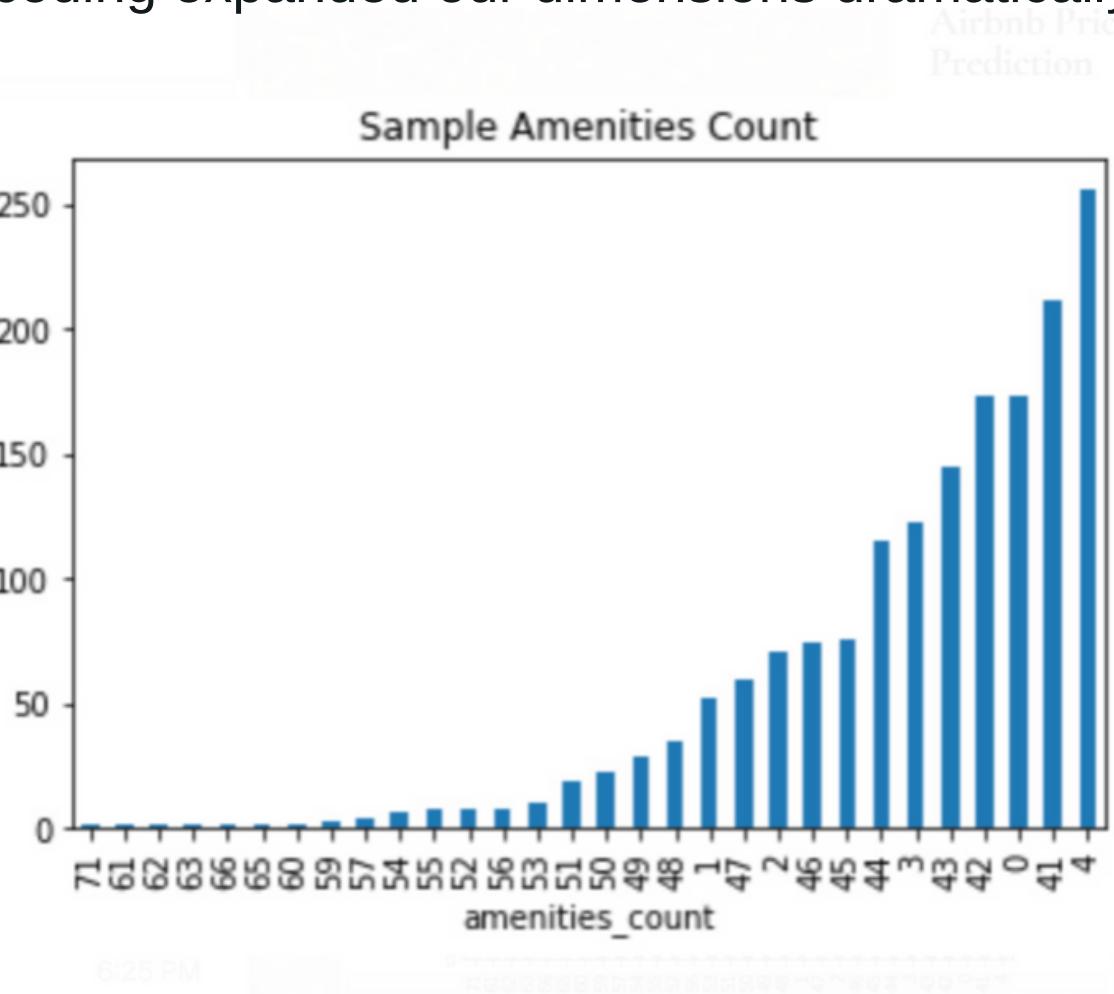
- 27926 rows × 45 columns
- 68 columns removed
- 22870 rows removed
- Feature Columns Engineered

# Feature engineering



## TROUBLESHOOTING

- First model used all the possible features (more than 120).
- Two features in particular, host verification and amenities have a lot of unique values where hot encoding expanded our dimensions dramatically.



- We also found out that the service fee and cleaning fee caused the model to behave differently than we would expect.
  - As we increased Cleaning Fees, the listing price also increased.
  - These were likely highly correlated with higher priced listing, but not predictive.

# Feature engineering

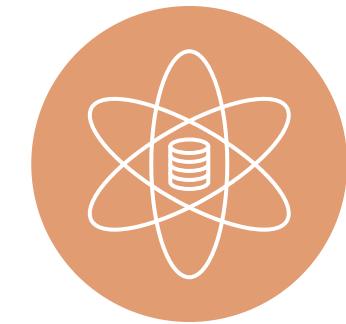


## OUR RESOLUTION

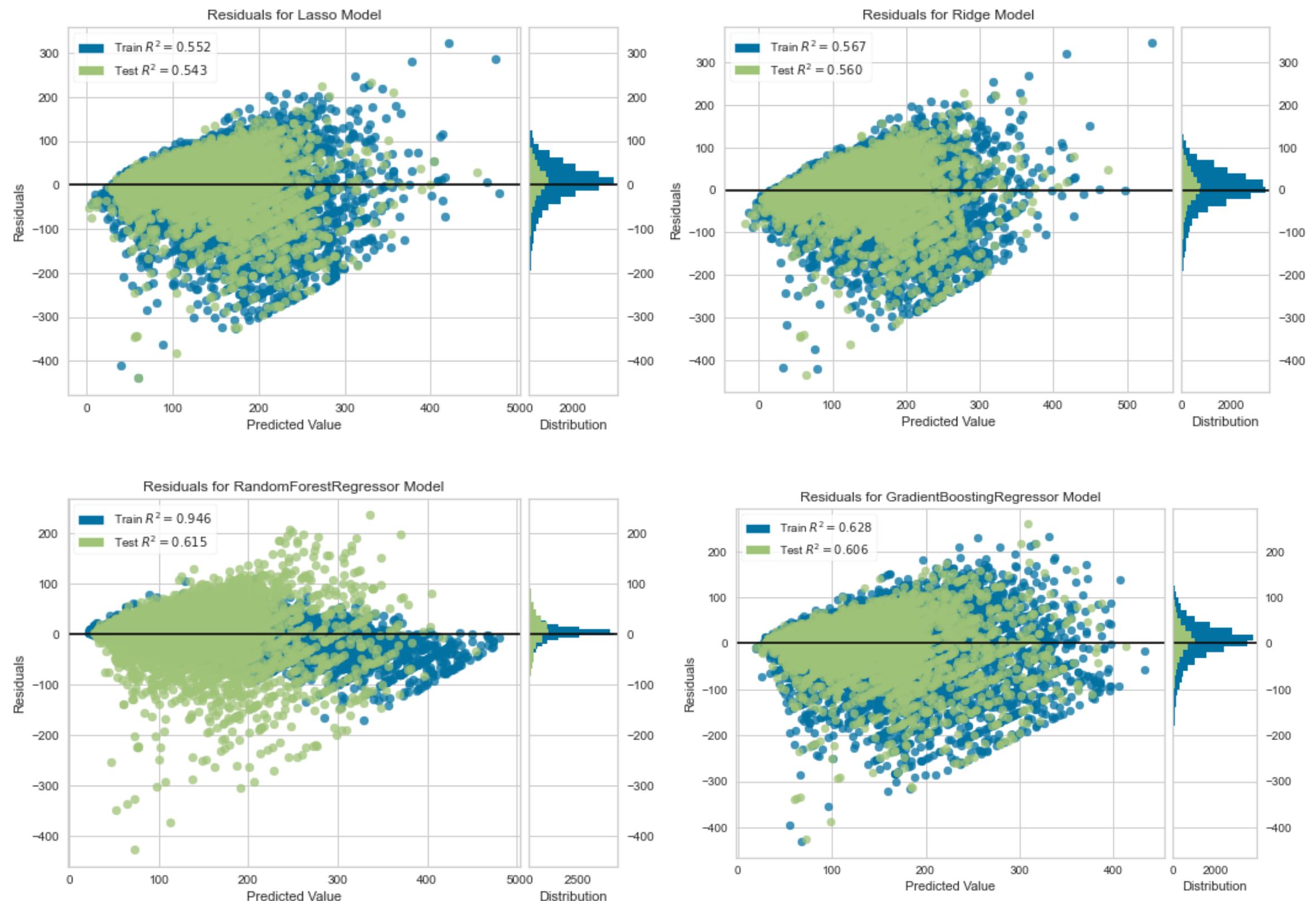
So we engineer some features to make the model simpler. Here are some of them:

- amenities , instead of using one hot encoding, we are counting them
- using yes or no for cleaning fee, instead of dollar amount
- host response type, property type, acceptance rate calculation, real bed, cancellation policy, listing counts

# Modeling: Base Models

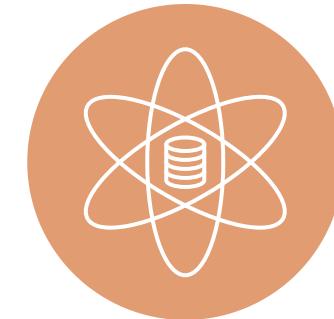


## Base Modeling



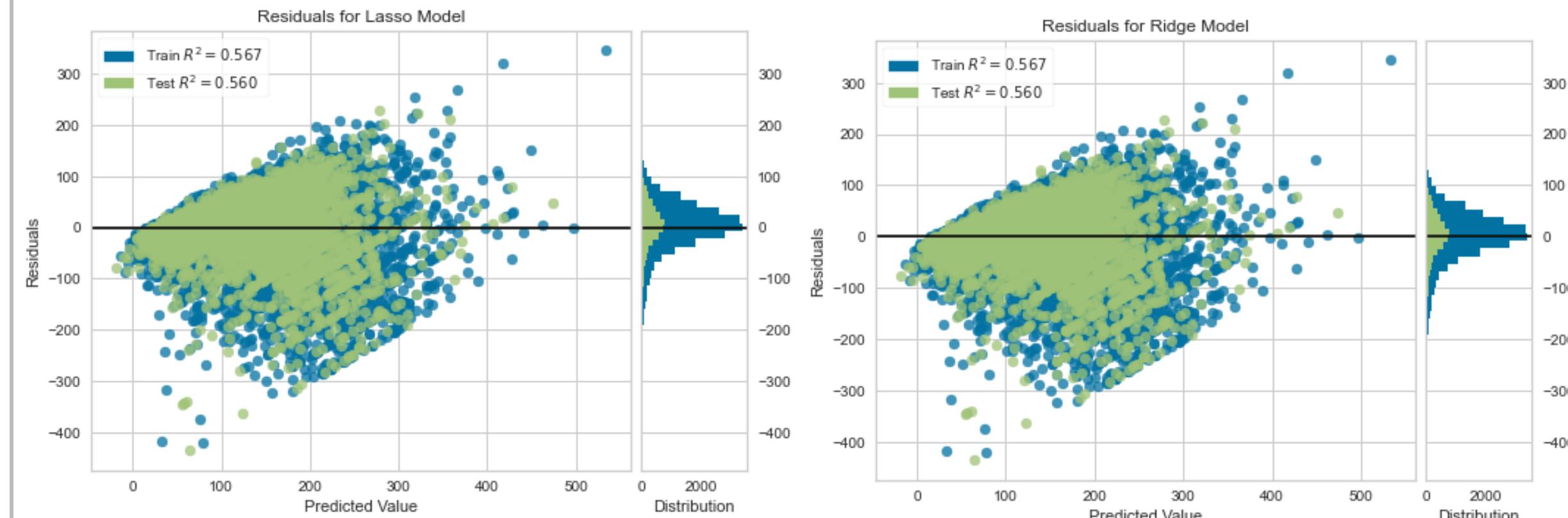
June, 2020 • Capstone Project

# Modeling: Hyper Parameter Tuning

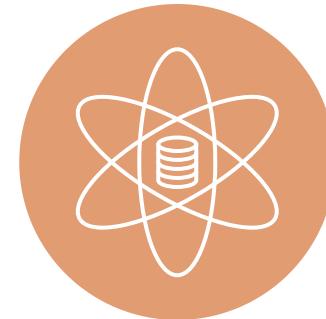


## Lasso & Ridge

- For Lasso & Ridge Regression models, we used Grid Search to find the optimal alpha to choose.
- For Lasso, Optimal alpha was .001
  - Applying this alpha improved  $R^2$ 
    - On Train Data, from .552 -> .567
    - On Test Data, from .543 -> .560
- For Ridge, Optimal alpha was 10
  - Applying this alpha had no impact on  $R^2$
- As you can see below, after optimization, the models performed identically.

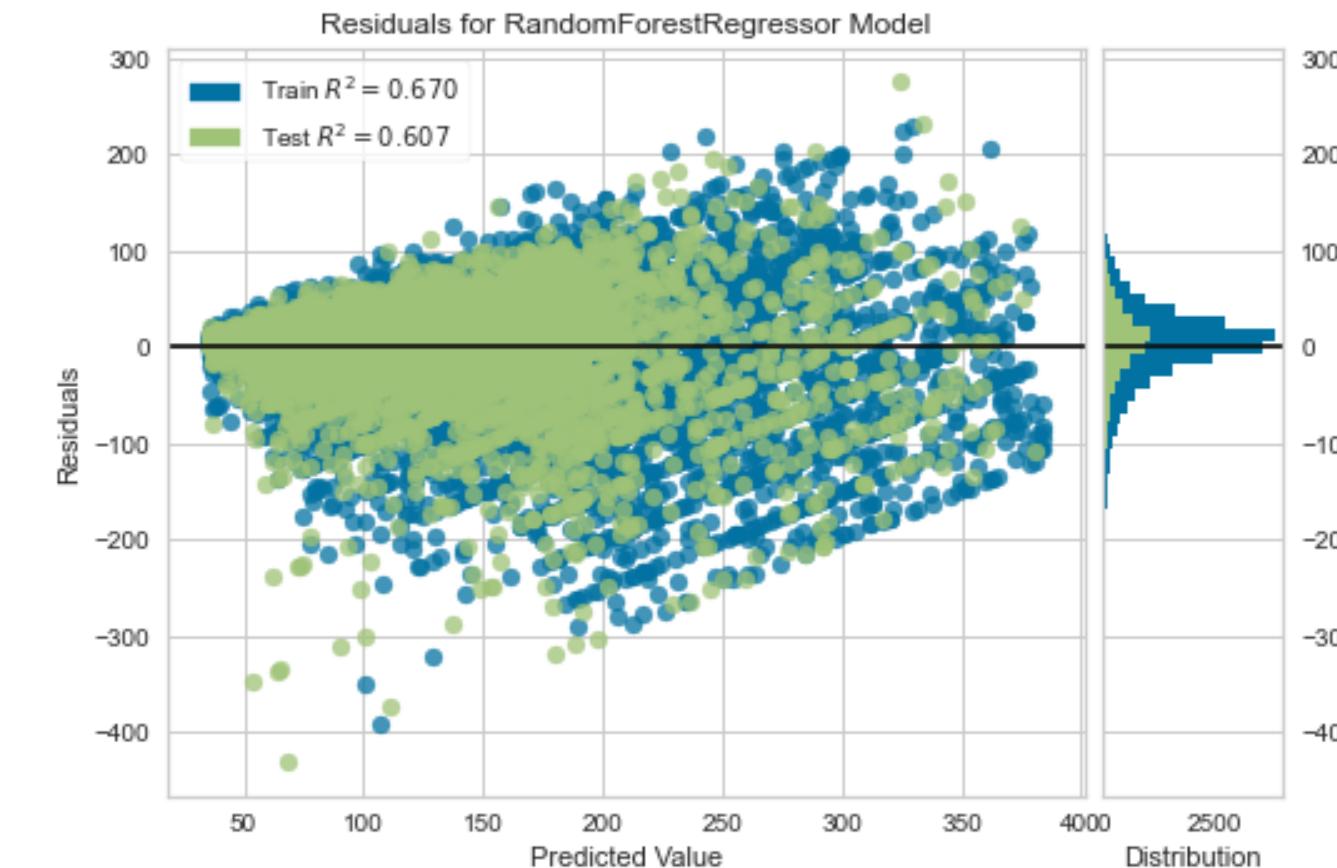


# Modeling: Hyper Parameter Tuning

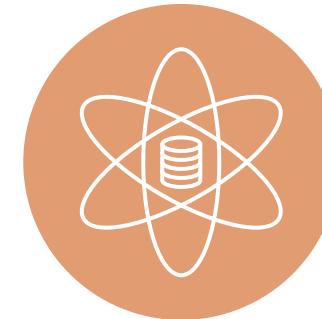


## Random Forest

- For Random Forest we were searching for the optimal number of estimators, which random state to use, and minimum sample splits.
- Base model showed high Train R<sup>2</sup> but much lower Test R<sup>2</sup>. This likely means the model is over-fit with the base model. This is the main issue we need to resolve if we want to use this model.
- Optimal parameters based on grid search we found were:
  - n\_estimators = 100
  - random\_state = none
  - min\_samples\_split\_100 = 100
- This still may be slightly over-fit still, but overall is much less than before while maintaining a relatively good test R<sup>2</sup>.

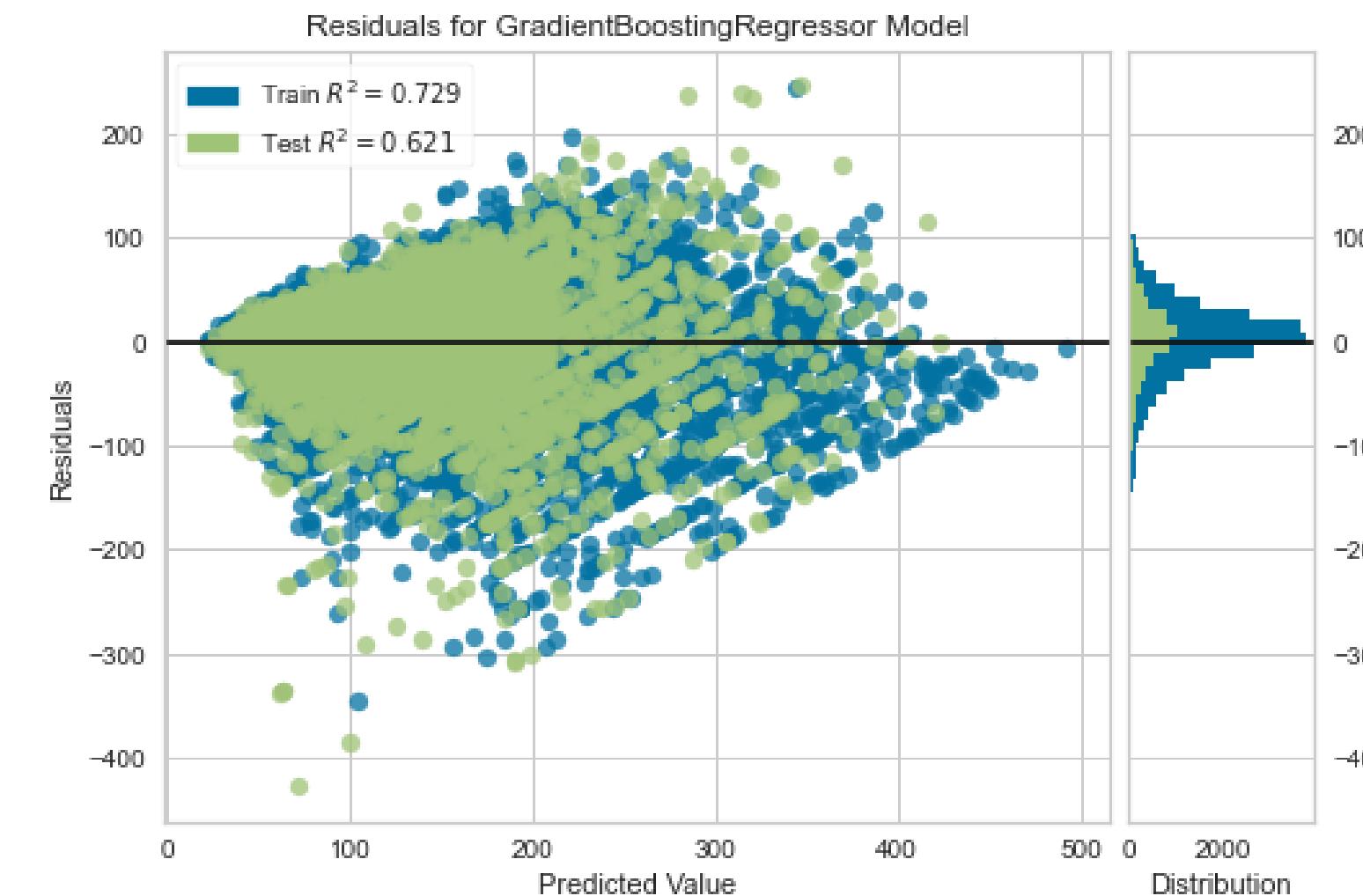


# Modeling: Hyper Parameter Tuning

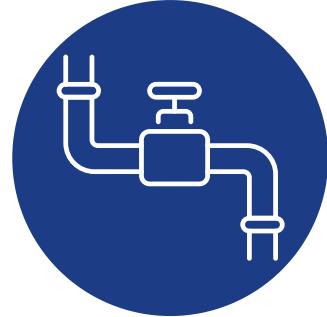


## Gradient Boost

- For the last model we tried, Gradient boost, the parameter we were looking to optimize was the max depth.
- Base model showed better  $R^2$  scores compared to the Lasso and Ridge, and was similar to the random forest regression scores on the test  $R^2$ .
- Optimal max depth from the grid search gave 6.
- Using this, we saw an increase in both  $R^2$  scores, but it may have been slightly over-fit.



# Data Pipelines



## Testing Models on other cities

Now that we had done each stage of the data pipeline, we were then able to stitch together a pipeline procedure so we could expand our testing to other cities.

Because the Lasso and Ridge performed identically after optimization, we decided to just use Lasso with the optimized alpha alongside both the Random forest and Gradient Boost models.

# Additional Exploration: Topic Modeling



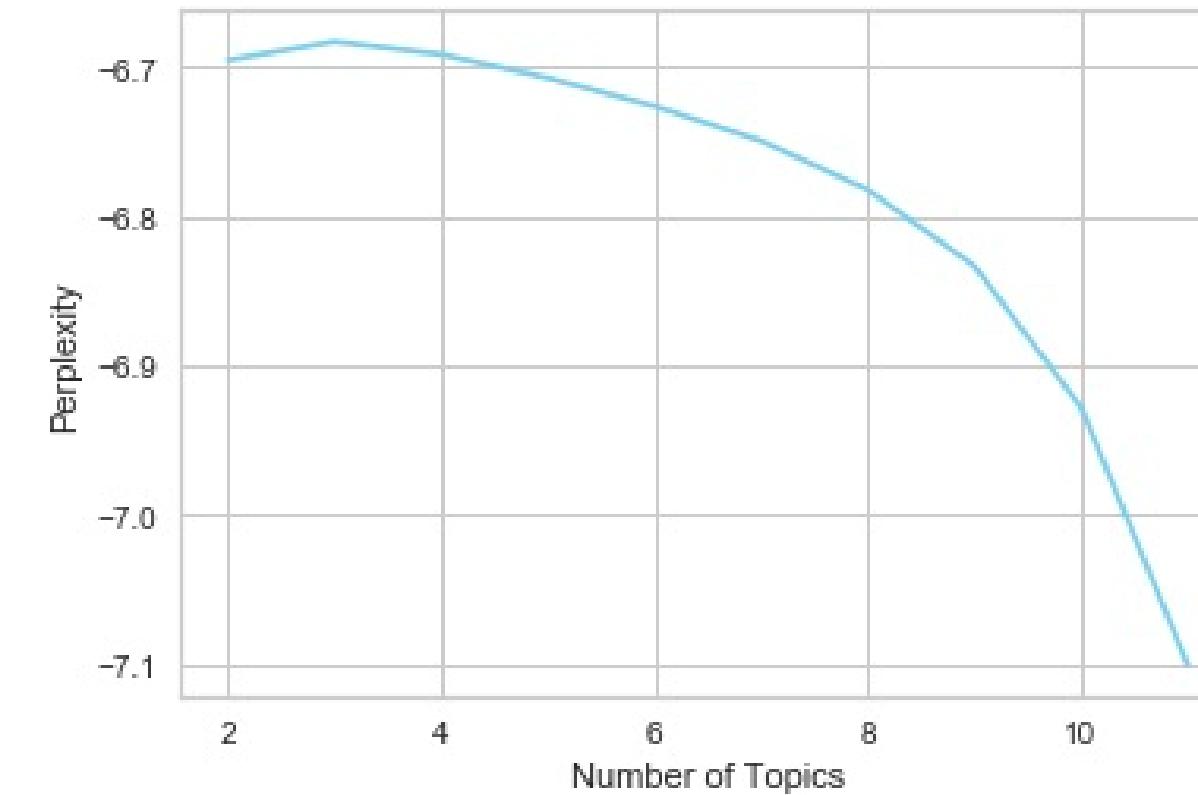
## Missing values in Airbnb dataset (TEXT)

- 4.8% missing in summaries
- 28.0% missing of space description
- 2.4% description
- 34.4% missing overviews on their neighborhood
- 59.6% missing notes from the host
- 34.4% missing information on transit
- 47.2% missing information on accessibility
- 40.2% missing information on interaction
- 39.1% missing on house\_rules
- 0% mising in name

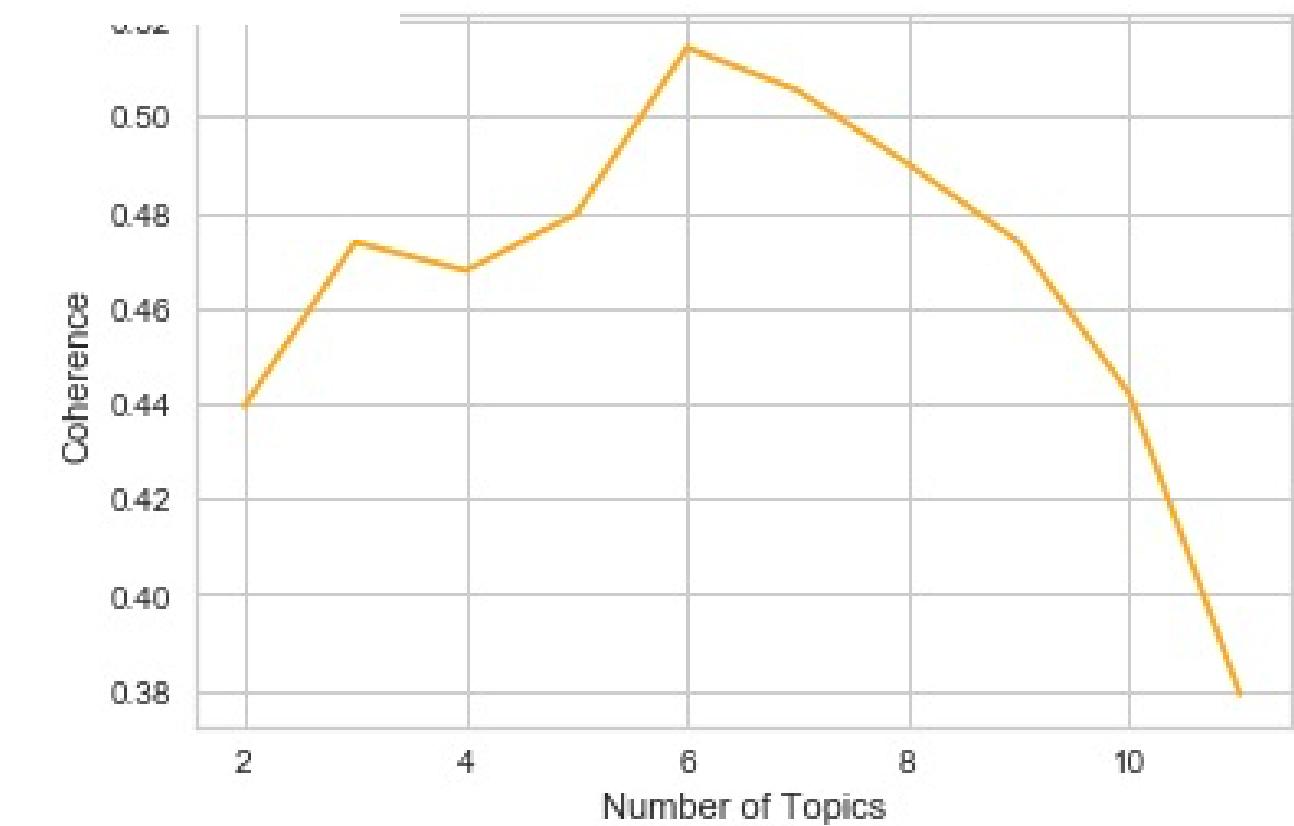
# Additional Exploration: Topic Modeling



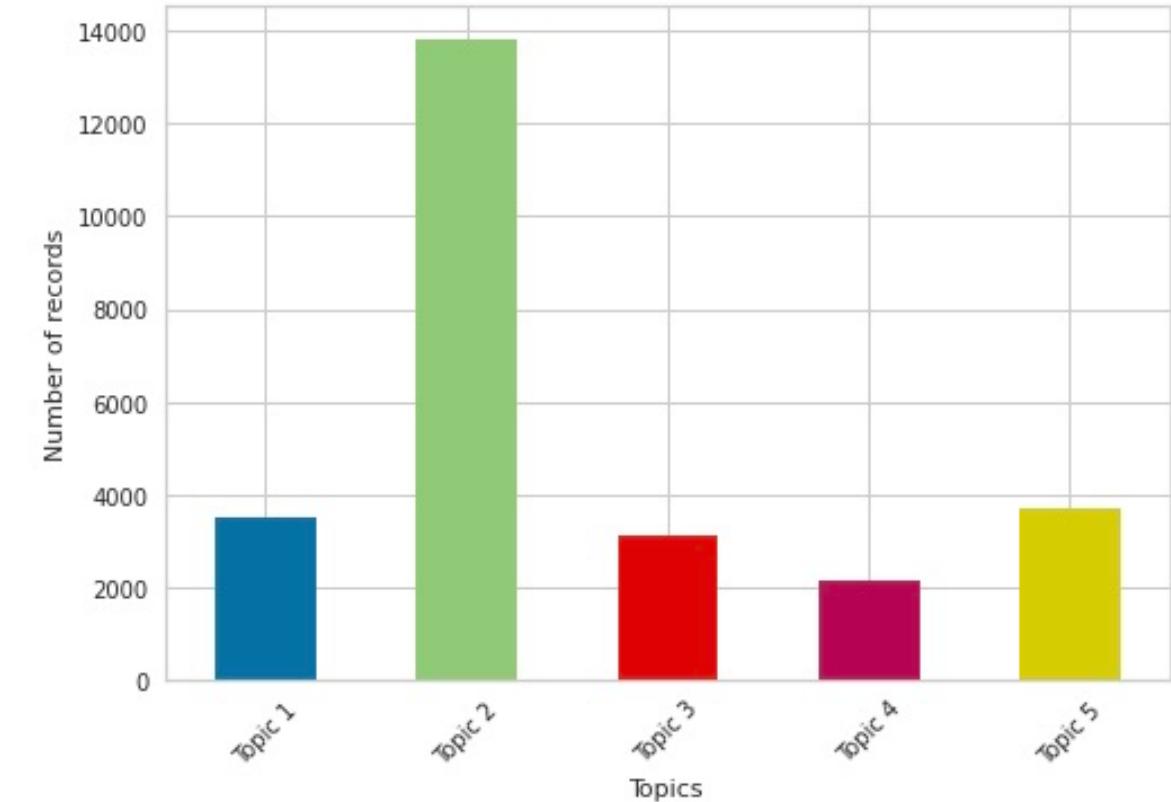
## LDA Text Analysis



[Visualization](#)



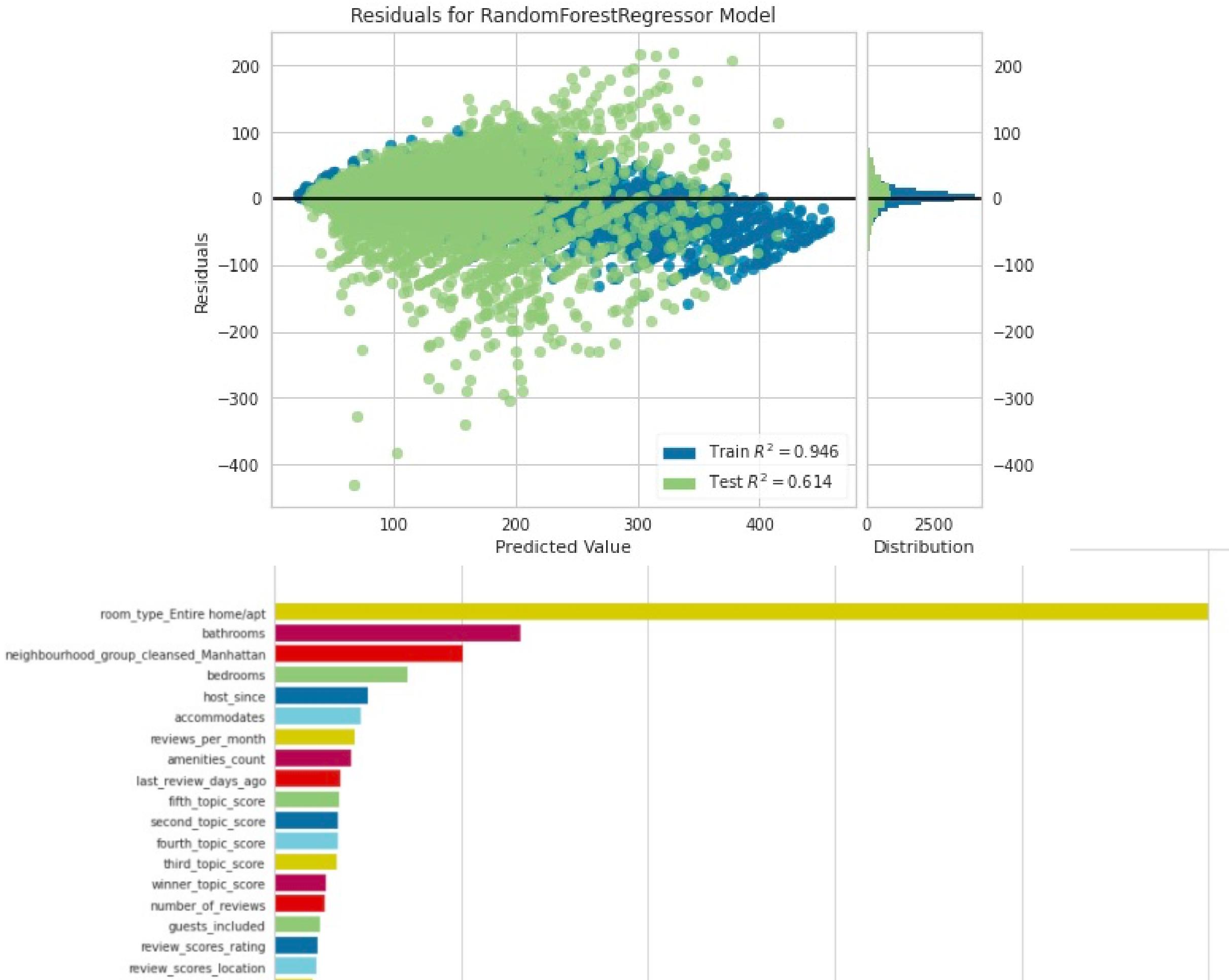
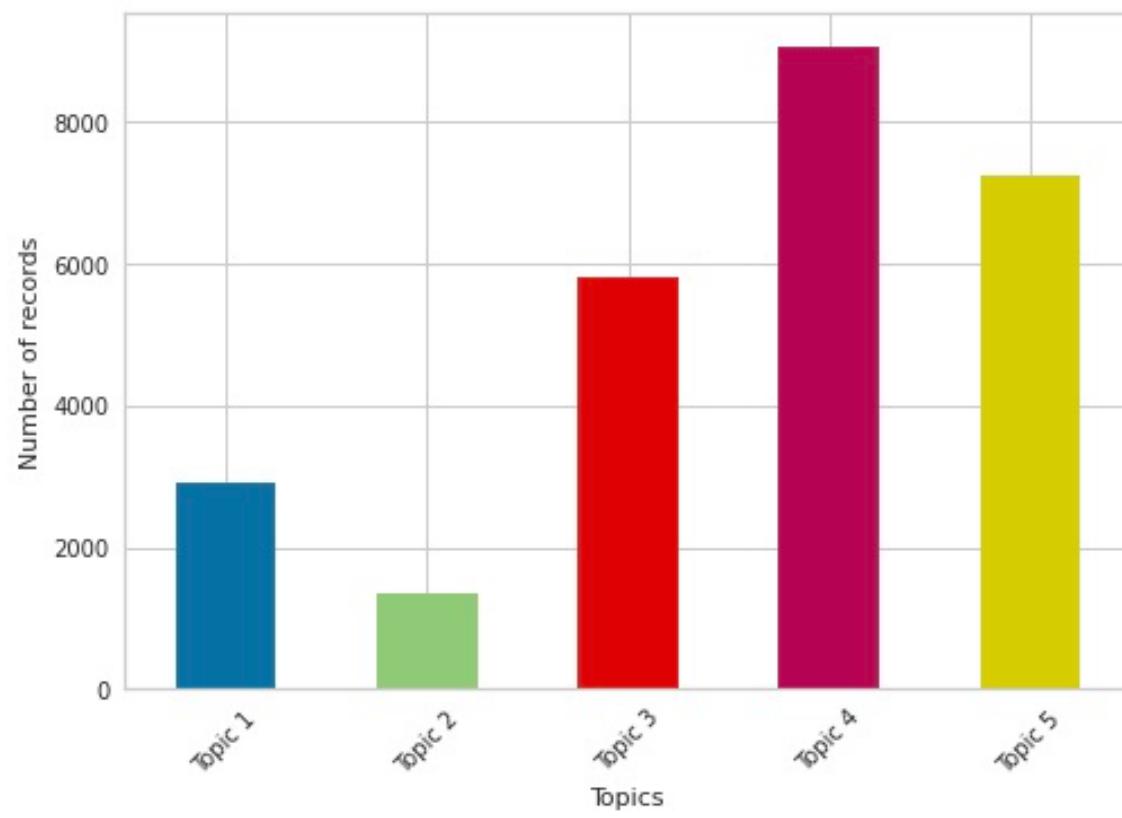
# Additional Exploration: Topic Modeling



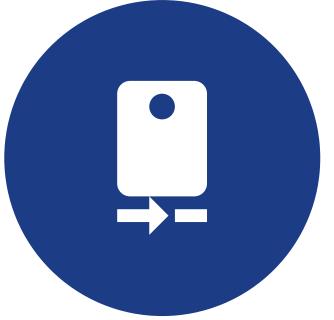
Topic # 2



# Modeling with TM scores



# Application End Product



June, 2020 • Capstone Project

## FLASK MODEL

- To simplify the usage of the model in a live system, we decided to use a simplified model (6 features).
- Though this app is interesting, and works as expected, without further testing in the field, we won't be able to verify its accuracy in the real world.

## Link to Flask App

### Predict Airbnb Price

Are you renting an entire home?  Yes  No

Is this in Manhattan?  Yes  No

How many people will it accommodate?

Number of bedrooms

Number of bathrooms

How many guests are included?

Show Me the Money

Suggested Listing Price is: \$

Entire Home:

In Manhattan:

Accomodations:

Number of Bedrooms:

Number of Bathrooms:

Total Guests:

# Further Research



Address heteroscedasticity of our model and R2.

Better understand of why the topic scores are important for prediction

Do a bigram or trigram analysis instead of single tokens because text is all about context

Use other algorithms to calculate frequency like TF-IDF to have a more detailed understanding of the text

Picture analysis