

Predicting AirBnB Rental Rates

Trevor Isaacson, Jonathan Olavarria, Jasmine DeMeyer

12/10/2021

Introduction

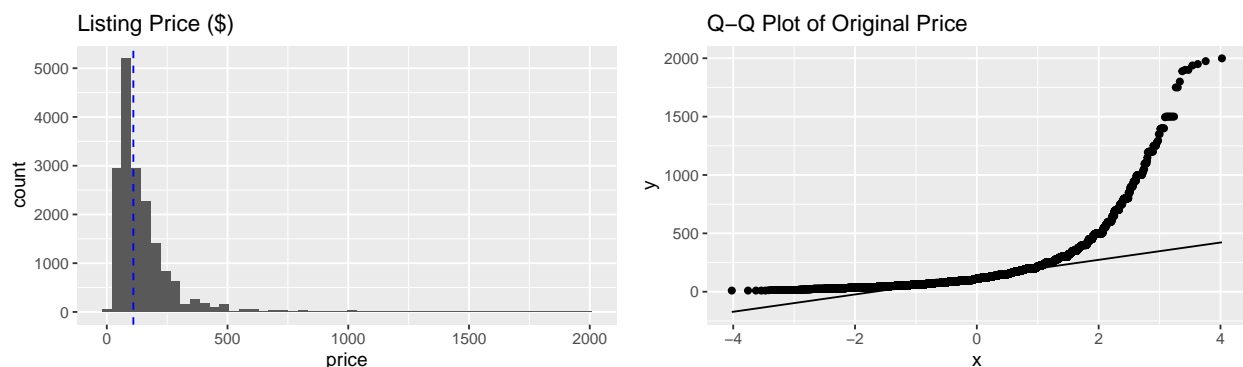
The Data

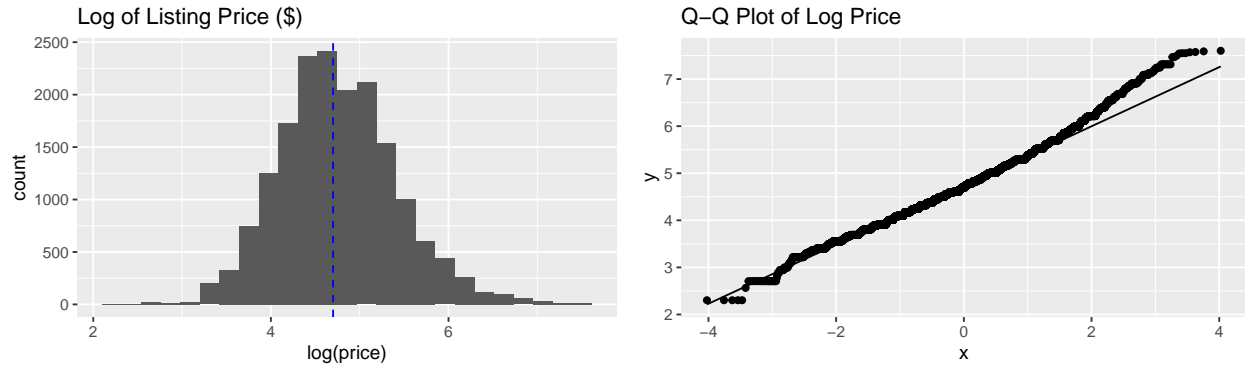
As a company, AirBnB is very open and transparent with the data they collect about their rental properties. They provide data about rental spaces in their system for cities and countries all over the world. Because of this, we were able to find a large dataset on Kaggle with AirBnB listings in major US cities including New York City, Los Angeles, San Francisco and others. The dataset available on Kaggle has over 74,000 entries and was used as a competition a few years ago. For the sake of time and processing, we trimmed our training data to about 17,500 entries and our test data to about 5,000 entries. We did this by taking a random sample of the provided training data. This allowed for easier access and faster processing while maintaining a large amount of data and individual AirBnB listings.

The original dataset contained 30 variables about each listing. Due to high correlations and lack of relevancy, our final dataset consisted of twenty-two variables. Those twenty-two variables can be split into four categories: property, location, host and host reviews.

Property includes:

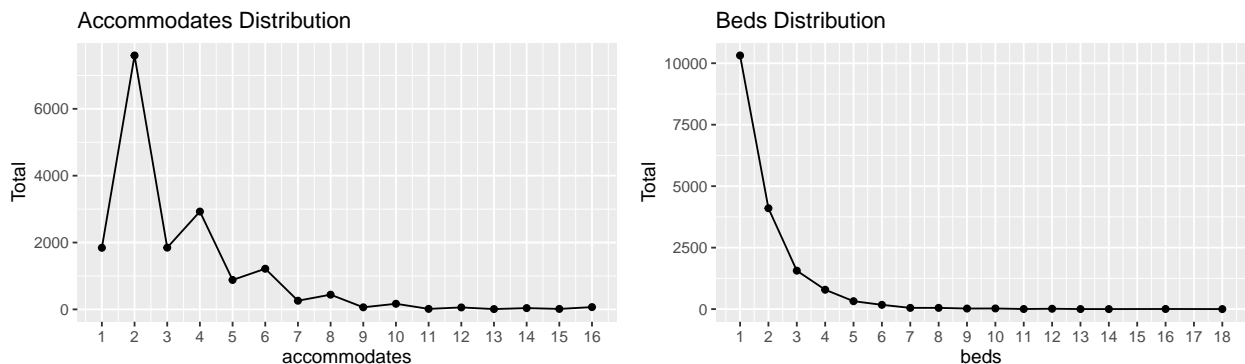
- price: listing price
 - Because the original price data is very heavily skewed, we needed to log transform the prices. As shown in the histogram, we have a very heavy right tail because there are a some listings with very high prices compared to the median price of \$110 (blue line). This non-normal shape and distribution is clearly evident in the Q-Q plot. The observations clearly curve away from the line depicting how skewed the distribution is.

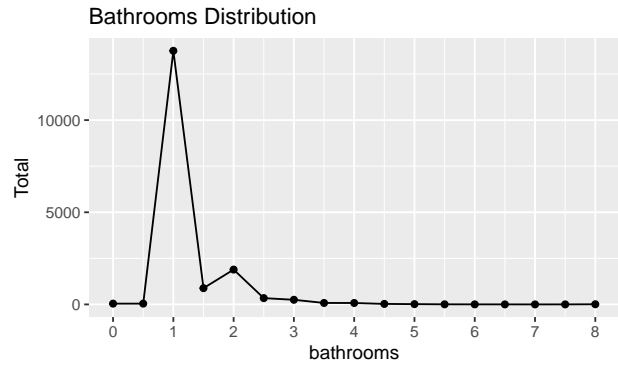




By applying a log transformation, we now have a more normal shaped distribution. Our Q-Q plot shows some evidence of a right tail but this can be expected since the original distribution is very right tailed.

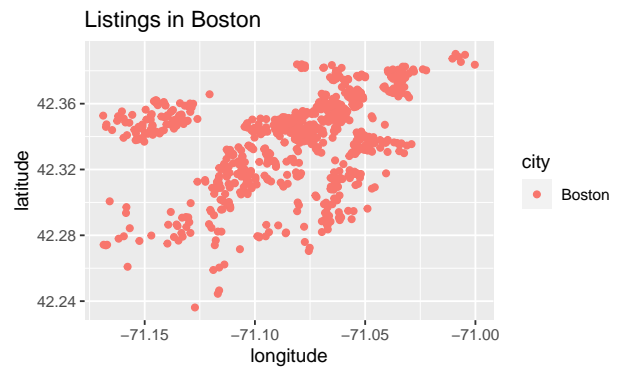
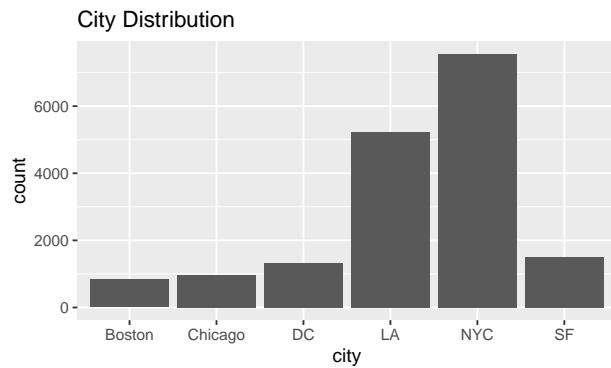
- `property_type`: defines the type of property listed
 - There are 21 different types ranging from apartments, houses, and condos to boats, cabins, hostels and even castles
- `room_type`: defines type of rental within the property
 - Includes entire home/apt, private room and shared room
- `accommodates`: number of people the property can comfortably accommodate
- `bedrooms`: number of bedrooms within the property
- `beds`: number of beds within the property
- `bed_type`: type of bed available
 - This includes a Real Bed, Futon, Pull-out Sofa, Airbed or Couch
 - Only 463 listings have something other than a Real Bed
- `bathrooms`: number of bathrooms within the property





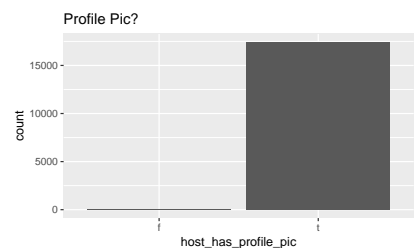
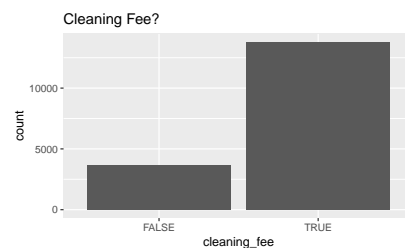
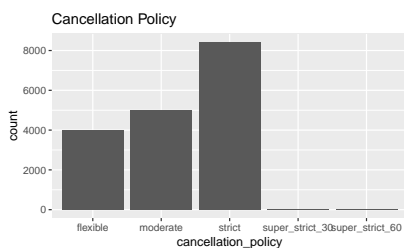
Location includes:

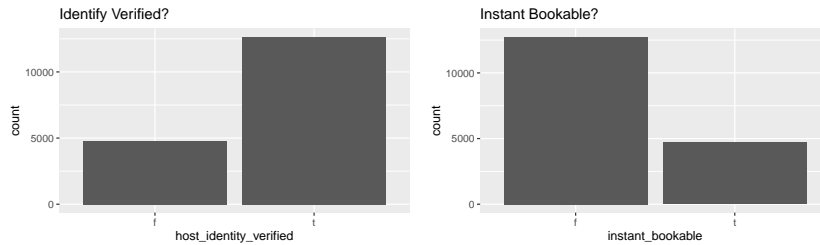
- city: Location of listing
- latitude and longitude: latitude and longitude coordinates of the listing



Host includes:

- cancellation_policy: strictness of cancellation policy set by the host
 - Levels include strict, moderate, flexible, super_strict_30 and super_strict_60
- cleaning_fee: TRUE/FALSE determines if host charges a cleaning fee
- host_has_profile_pic: TRUE/FALSE determines if the host has uploaded a picture to their profile
- host_identify_verified: TRUE/FALSE determines if the host's identity has been verified by AirBnB
- instant_bookable: TRUE/FALSE determines if the property can be booked in short notice
- host_response_rate: how often does the host reply to potential clients?





Host Reviews:

- number_of_reviews: Number of reviews the host has received
- review_scores_rating: average review rating for the host and property
- first_review_year: year of the first review
- last_review_year: year of most recent review
- host_since_year: year the property was first listed on AirBnB

As a group, we felt these twenty-two predictors were all relevant and important in helping predict price.

Models

There are have a lot of machine learning methods discussed this past semester and we wanted to incorporate some of our favorites into our research. Thus, we have included linear regression, splines, general additive models, PCR, PLS, trees, bagging, random forests and bagging.

In order to both train and test using the training data file, we needed to split the 17,500 total observations into roughly a 50/50 split. To do this, we took a random sample of 9000 observations and made that the training set. From now on, this random sample will be referred to as the training set. The remaining 7500 observations became our testing set for determining the the performance of each method.

Regression

To begin, we started with a simple multiple linear regression model. We wanted to give ourselves a baseline mean squared error value and because linear regression is the easiest to apply and interpret, we determined this was the best place to start. The model was fit using all twenty-two variables and the training set.

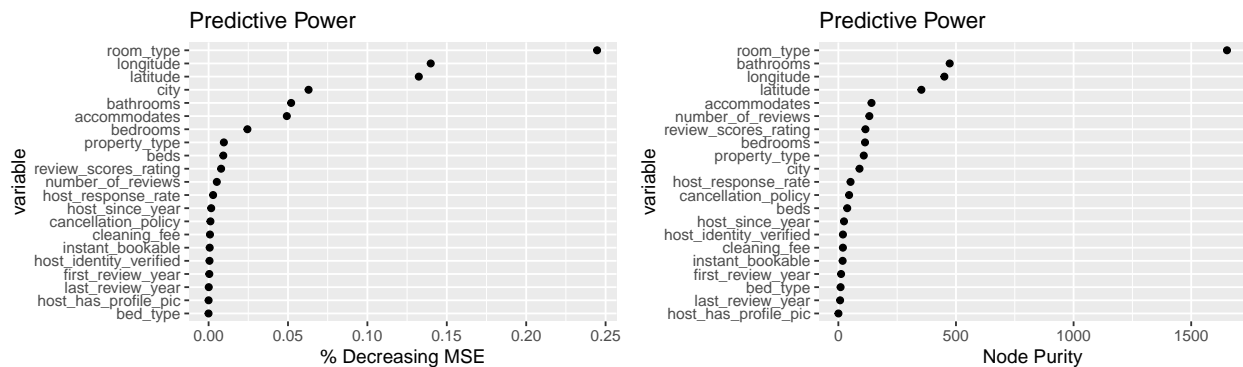
```
## [1] "R^2: 0.632 , Adjusted R^2: 0.63"
## [1] "MSE of Testing Set: 0.1652"
## [1] "Leave One Out Cross Validation: 0.1829"
## [1] "k-fold Cross Validation: 0.3063"
```

As shown in the results above, our linear regression model was able to set a good baseline for future methods with a mean squared error of 0.1652. With our linear regression model, we also applied some cross validation. For Leave One Out Cross Validation, we achieved a mean squared error 0.1829 and applying k-fold cross validation with $k = 10$, we achieved a mean square error of 0.3063. Clearly, Leave One Out Cross Validation performed better than the k-fold cross validation.

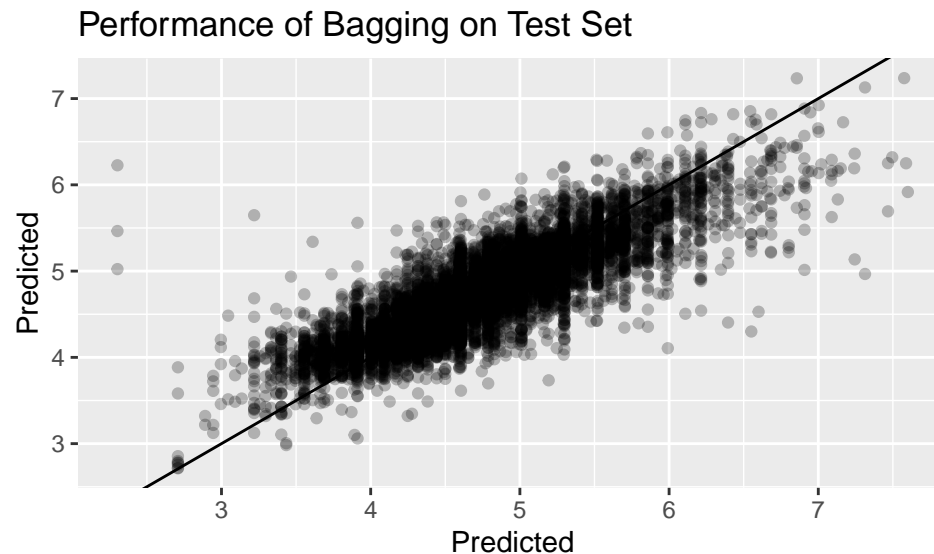
Bagging

Because decision trees suffer from high variance, we then moved onto bagging or bootstrap aggregation to see if we could lower the MSE of our trees. Although bagging decision trees can be slow as it has to average hundreds or thousands of trees together, it won't lead to any overfitting. We performed bagging on the training set to predict price and found that room type, bathrooms and latitude/longitude were the most important variables. Our bagging model was able to obtain a MSE of 0.1293 and was found to be one of our best methods.

```
## [1] "Test MSE of Bagging: 0.1293"
```



The plot below shows an Actual vs Predicted plot for the predicted values using bagging. In general, the data points move positively along the diagonal line with a slight horizontal tilt in the data points compared to the horizontal line. There are a few outliers especially on the right side of the plot. It appears bagging is predicting slightly higher than the actual values for prices less than 5 and slightly lower than the actual values for prices more than 5. However, this plot shows the overall effectiveness for bagging.

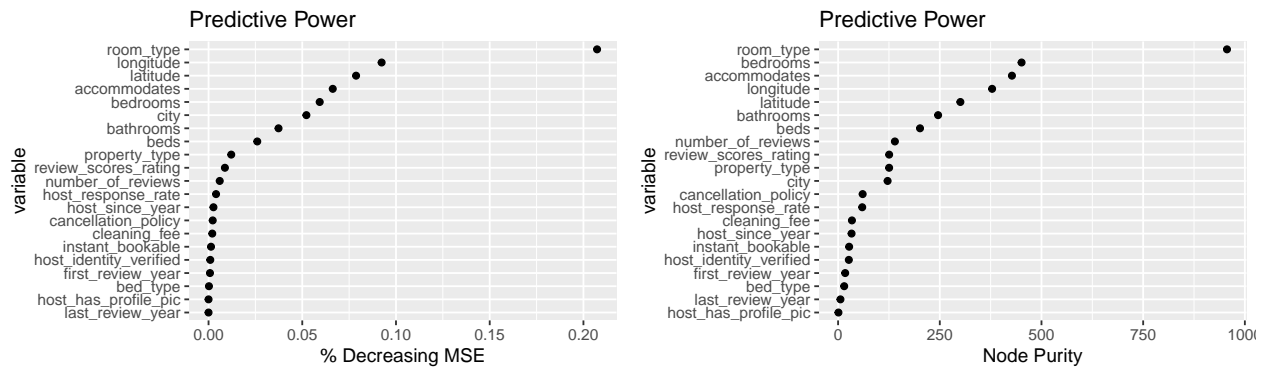


Random Forests

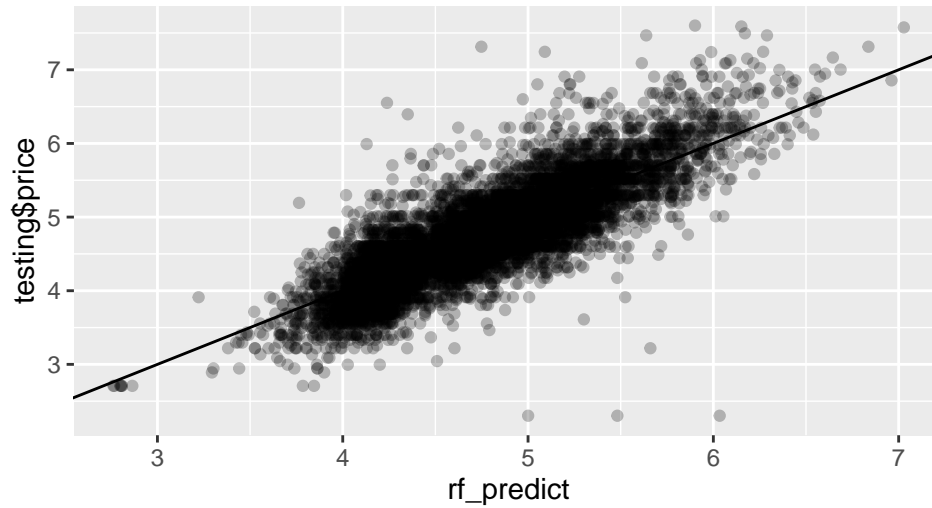
Next, we wanted to determine if random forests could provide an improvement over the bagged trees and see if decorrelating the trees will lead to better MSE result. The choice of m or the size of the predictor

subset was set to the square root of the total predictor set.

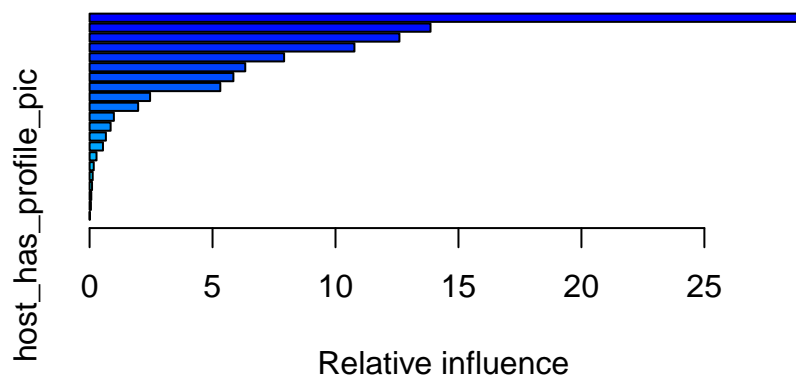
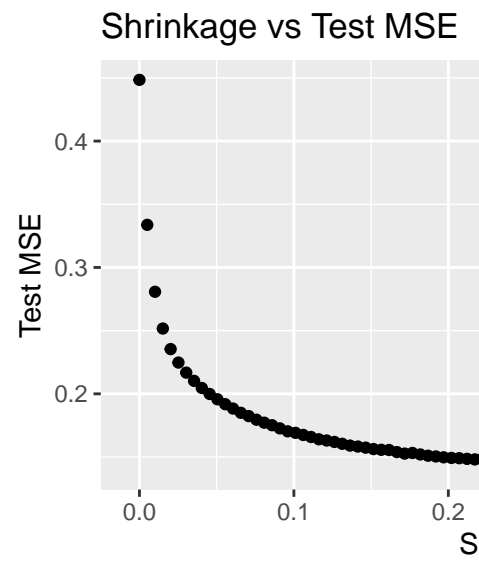
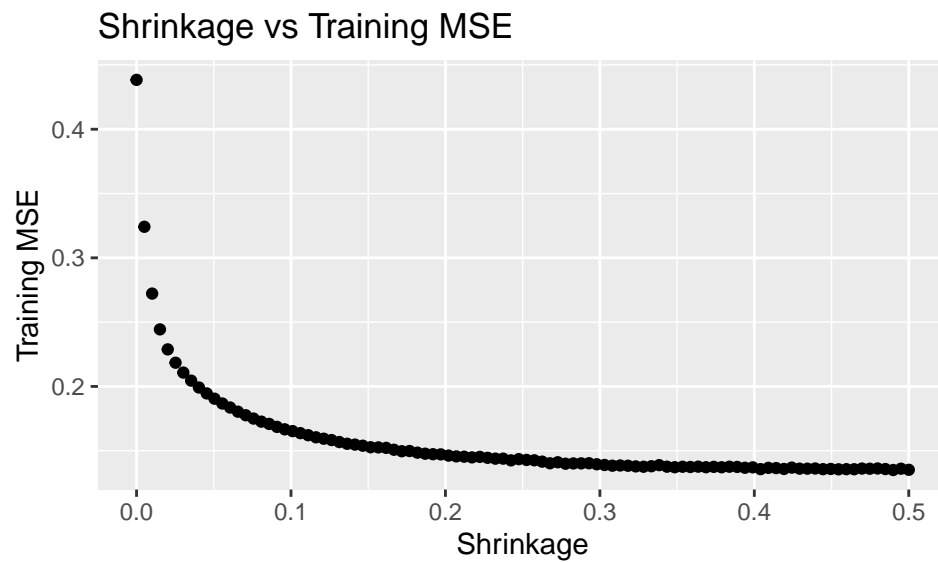
```
## [1] "Test MSE of Random Forest: 0.1305"
```



Performance of Bagging on Test Set



Boosting



##	var	rel.inf
## property_type	property_type	29.065131471
## room_type	room_type	13.863836420
## bathrooms	bathrooms	12.599142177
## bedrooms	bedrooms	10.769605688
## accommodates	accommodates	7.910744021
## longitude	longitude	6.330597804
## latitude	latitude	5.847626717
## beds	beds	5.315361788
## review_scores_rating	review_scores_rating	2.460158299
## city	city	1.974567122
## number_of_reviews	number_of_reviews	0.987549731
## bed_type	bed_type	0.855576511
## last_review_year	last_review_year	0.663648178


```
## host_response_rate      host_response_rate 0.544699245
## cancellation_policy     cancellation_policy 0.282229574
## instant_bookable        instant_bookable 0.169424186
## host_since_year         host_since_year 0.127083521
## cleaning_fee            cleaning_fee 0.099457369
## first_review_year       first_review_year 0.069712030
## host_identity_verified  host_identity_verified 0.053887259
## host_has_profile_pic    host_has_profile_pic 0.009960889
```

```
## [1] "Testing MSE for Boosted Model: 0.1311"
```

MSE Results

City Prediction

Our original data set is from a past and completed Kaggle competition. Thus, the data was split into training and testing files. The testing file has no response column included because that is how they would determine a winner. Even though our testing results can't be compared to the actual prices, we can still use it for testing our final model.

```
## [1] 4984
```

```
##   property_type      room_type accommodates bathrooms bed_type
## 1   Apartment      Private room           2          1.5 Real Bed
## 2   Apartment      Private room           2          1.0 Real Bed
## 3   Apartment Entire home/apt            4          1.0 Real Bed
## 4   Apartment Entire home/apt            3          1.0 Real Bed
## 5     House Entire home/apt            4          1.0 Real Bed
## 6   Apartment Entire home/apt            4          1.0 Real Bed
##   cancellation_policy cleaning_fee city host_has_profile_pic
## 1           moderate          TRUE   LA                   t
## 2           moderate          TRUE   DC                   t
## 3           flexible          FALSE  NYC                   t
## 4           strict           FALSE  NYC                   t
## 5           moderate          TRUE  NYC                   t
## 6           moderate          TRUE  NYC                   t
##   host_identity_verified host_response_rate instant_bookable latitude
## 1                      t                 0.00                t 34.10388
## 2                      f                 0.00                t 38.90633
## 3                      f                 1.00                t 40.67878
## 4                      t                 0.00                f 40.80747
## 5                      f                 1.00                t 40.69085
## 6                      t                 0.86                t 40.68004
##   longitude number_of_reviews review_scores_rating bedrooms beds
## 1 -118.34929             1             100             1      1
## 2  -77.00622             5              96             1      1
## 3  -73.95585             2             100             1      2
## 4  -73.96030            12              93             1      1
## 5  -73.79916             6             100             1      3
## 6  -73.92150             2             100             1      2
##   first_review_year last_review_year host_since_year price
```

## 1	Greater2014	Greater2014	Less2014	0
## 2	Greater2014	Greater2014	Greater2014	0
## 3	Greater2014	Greater2014	Greater2014	0
## 4	Less2014	Greater2014	Less2014	0
## 5	Greater2014	Greater2014	Greater2014	0
## 6	Greater2014	Greater2014	Less2014	0

References: