

Airbnb Project

Matthew Coleman, Austin Mac, Jeff Pittman, and Nick Reyes

2/28/2020

```

# levels(bnb$room_type)
#
# levels(bnb$neighbourhood_group)
#
# n_distinct(bnb$neighbourhood)
#
# n_distinct(bnb$neighbourhood_group)

# [1] "Entire home/apt" "Private room"      "Shared room"
# [1] "Bronx"           "Brooklyn"      "Manhattan"     "Queens"        "Staten Island"
# [1] 221
# [1] 5

```

1 Abstract

2 Introduction

3 Methods

3.1 Data

The dataset we will be using for our analysis is the dataset New York City Airbnb Open Data from Kaggle. This dataset contains the listing activity and metrics for Airbnb in New York City, New York during 2019. There are 48895 observations and 16 attributes to the dataset. The main features we are going to use for our analysis include the following:

- Price: Our main response variable. The price, in dollars, of the listing per night. Log-transformed to normalize distribution.

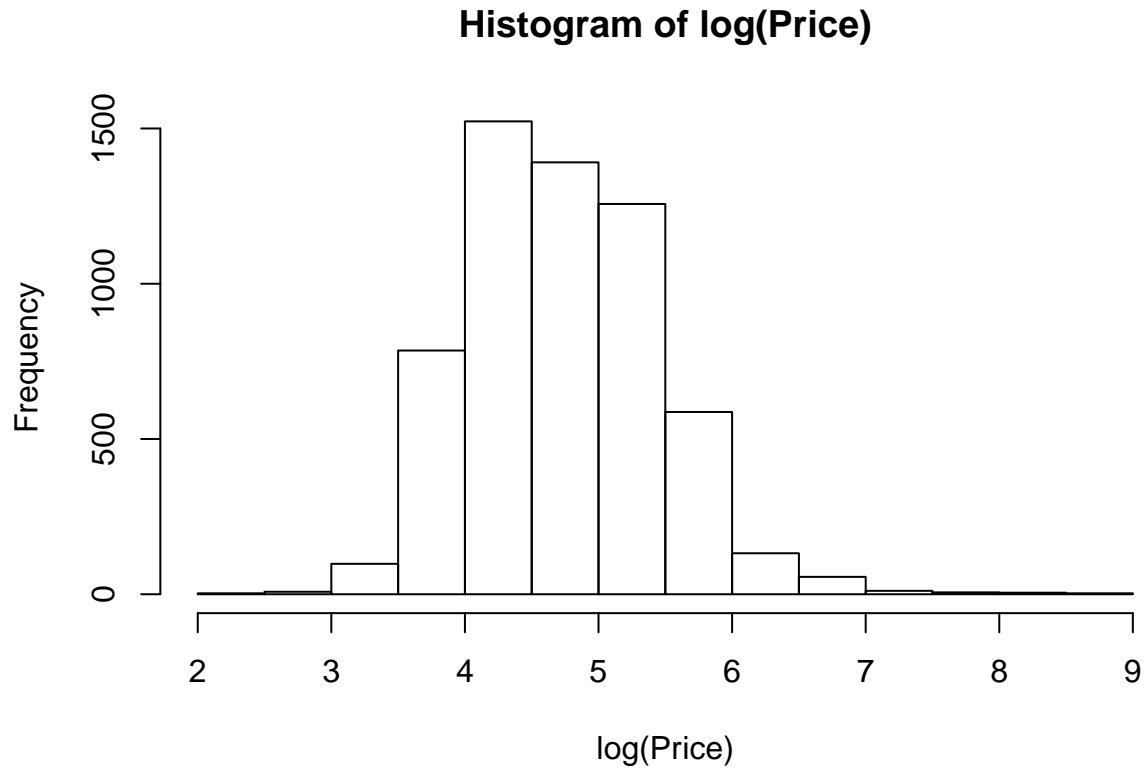


Figure 1: Log-transformed Price

- Price Above: Variable created from `price`, `price_above` is a binary variable of signaling whether a listings price is above the median listing price. 1 represents the price being above the median, and 0 represents the price being below the median.
- Neighbourhood: Categorical variable of the neighbourhood to which a listing belongs. This is a nested version of neighbourhood group, with 221 unique neighbourhood groups.
- Neighbourhood Group: Factor variable of the neighbourhood group to which the listing belongs. There are 5 neighbourhood groups in the dataset.
 - Plots of both neighbourhood and neighbourhood group are shown below:

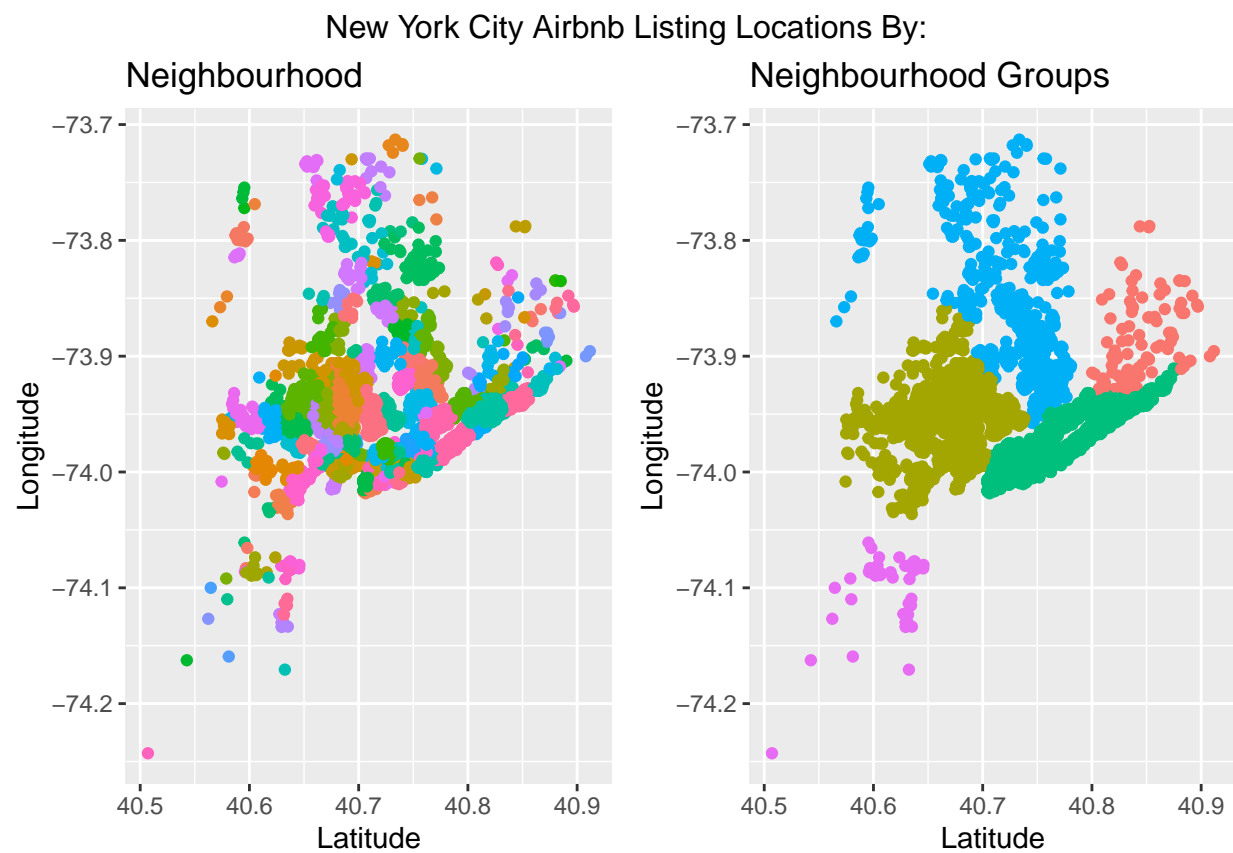


Figure 2: Neighbourhood and neighbourhood group



Figure 3: Justification for using neighbourhood group

- Latitude: Latitude coordinates of the listing.
- Longitude: Longitude coordinates of the listing.
- Room Type: The listing space type. Three types: *Entire home/apt*, *Private room*, *Shared room*.
- Minimum Nights: The minimum amount of nights someone can stay in the listing.
- Number of reviews: The number of reviews for the host.
- Reviews per Month: The number of reviews per month for the host. Formula: $\frac{\text{Number of Reviews}}{\text{Months Listed}}$.
- Calculated Host Listings Count: The number of listings per host.

All attributes were complete with the exception of `last_review`, which has the date of the last review, and `reviews_per_month`. Upon further exploration, the reviews per month feature was NA only when the host had no reviews. This resulted in us imputing 0's for NA values in the reviews per month column. Because the date of last review was unimportant to our analyses, we did not impute values for this column.

3.1.1 Assumptions.

Many of our machine learning methods are very computationally intensive, so we sampled 15% of the entire dataset, and then train-test split the 15% sample into 80% training 20% test dataset. To verify this was a viable practice, we plotted the distribution of our response variable, price, and verified the distribution is the similar to the distribution of the overall dataset. The histogram is very similar, and even contains some of the outliers we can see in the overall dataset, so we assumed our smaller dataset was representative of the population.

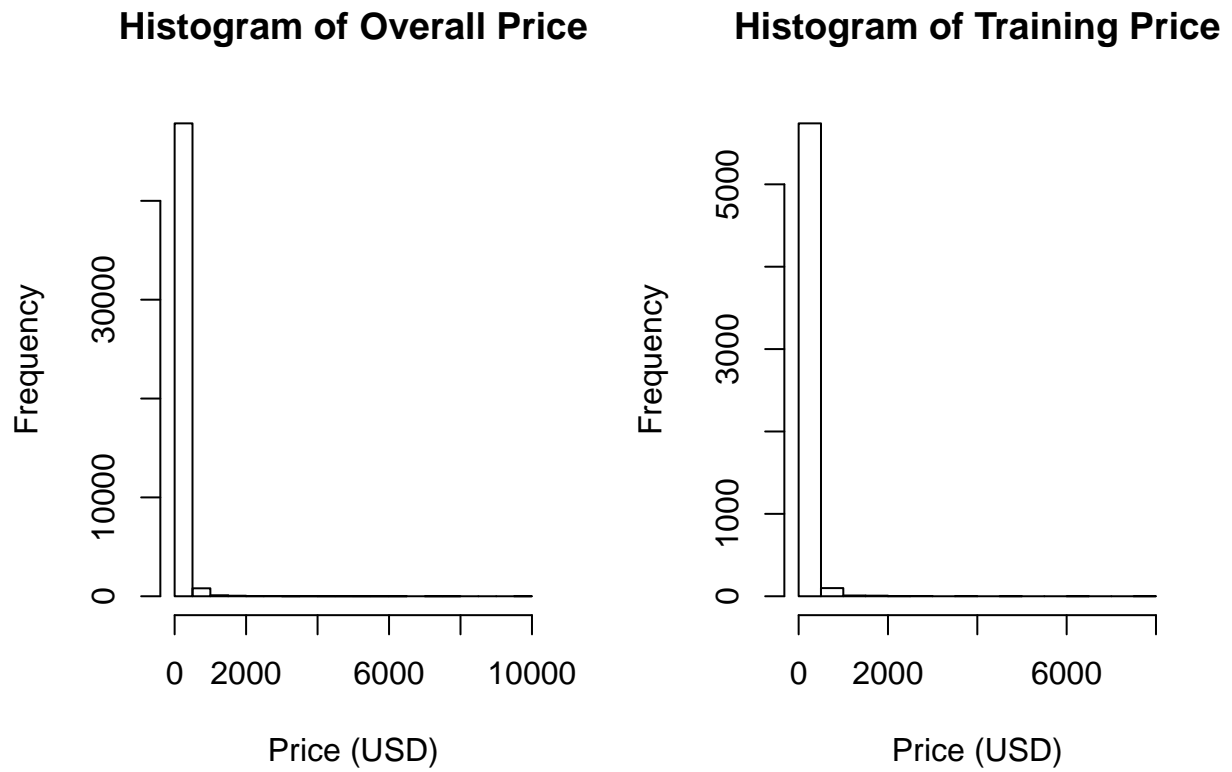


Figure 4: Training Data Histogram

We assessed the correlation between our variables with a correlation heatmap:

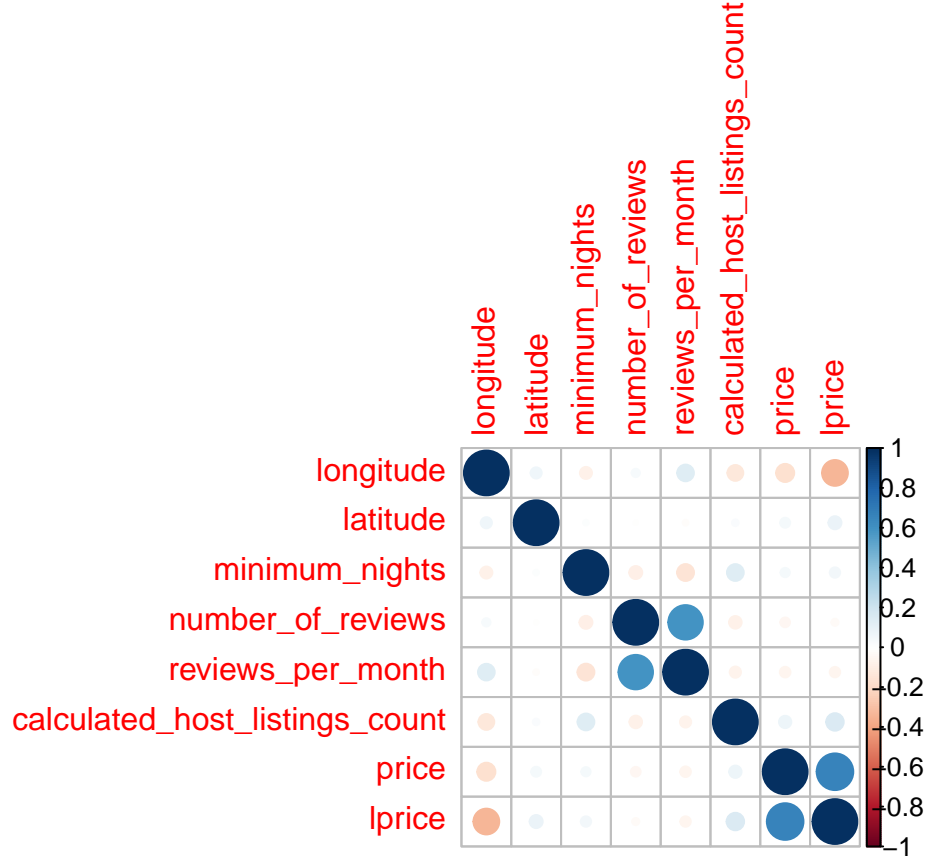


Figure 5: Feature Correlations

Reviews per month and number of reviews were highly correlated, so we decided to remove number of reviews to account for collinearity.

3.1.2 Sample Sizes

Our overall dataset is 48895. Taking the proposed 15% split on the data left us with an overall dataset of 7332 observations. The 80/20 train-test split left us with 5865 training samples and 1467 test samples.

3.2 Machine Learning Methods

3.2.1 Regression Methods

Methods used to predict the price of a listing:

- Ridge Regression:
 - Constraint optimization on the least squares criterion:

$$\hat{\beta}_{ridge} = \underset{\beta}{argmin} [||Y - XB||^2 + \lambda \sum_{j=1}^p \beta_j^2]$$

- Lasso Regression:

- Constraint optimization and model selection on the least squares criterion:

$$\hat{\beta}_{ridge} = \underset{\beta}{argmin} [||Y - XB||^2 + \lambda \sum_{j=1}^p |\beta_j|]$$

By using these two methods, we can try to reduce our estimates for the linear model by imposing some Bias on our estimates for β . Another benefit of using Lasso regression is that we can also perform model selection, making a simpler model.

- Tree Methods
 - Individual Trees: To compare the efficacy of ensemble tree methods, we will fit an individual regression tree on longitude and latitude, and then one tree on all variables of interest.
 - Bagging: We will fit an ensemble tree method which will grow large trees on bootstrapped data, resulting in high variance low bias. All of these trees predictions will be averaged to give the final prediction.
 - Random Forest: We will create multiple decision trees similar to bagging, but try to decorrelate each of the bootstrap trees through selecting $m = \frac{p}{3}$ variables.
 - Boosting: We will fit multiple (weak) trees sequentially, grown on information from the previously grown tree. Final prediction is a weighted prediction of the weak learners.

3.2.2 Classification Methods

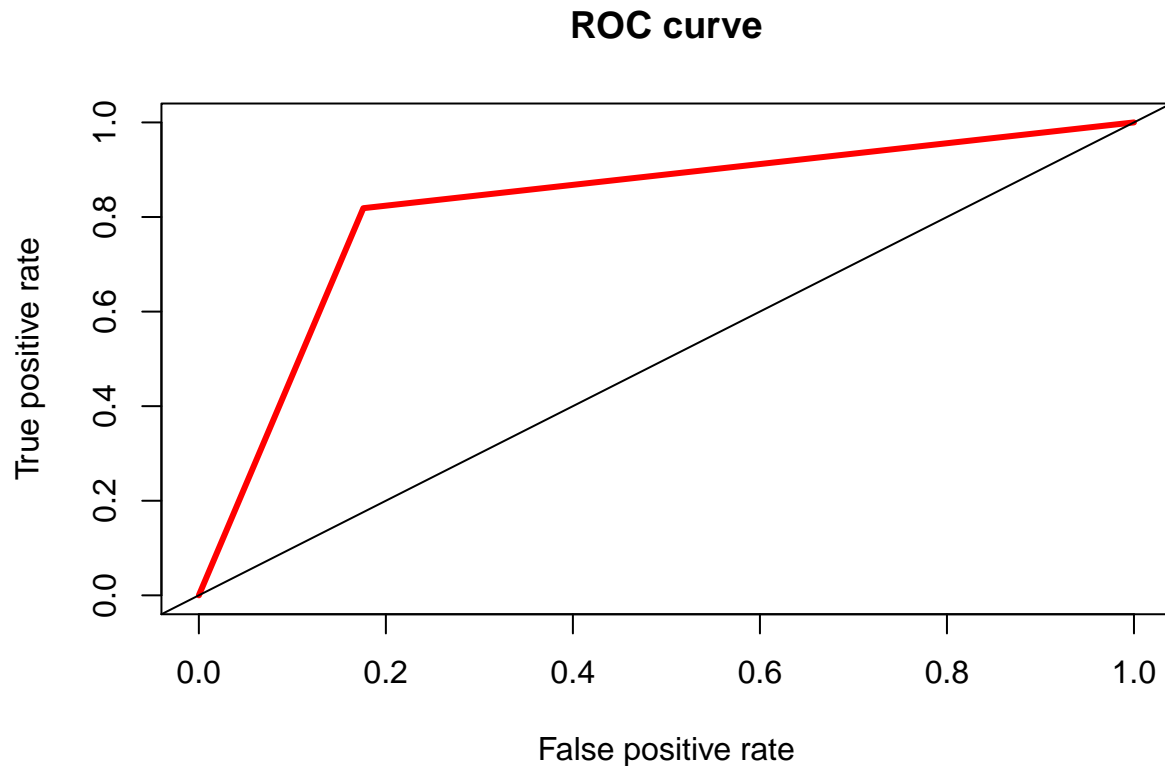
Methods used to predict whether a listings price is above the median:

- Logistic Regression: We will fit a logistic regression model on all variables of interest, using a binary classification output to predict whether a listing's price is above or below the median.
- LDA
- QDA
- Tree Methods
 - Individual Trees: To compare the efficacy of ensemble tree methods, we will fit an individual classification tree on longitude and latitude, and then one tree on all variables of interest.
 - Bagging: We will fit an ensemble tree method which will grow large trees on bootstrapped data, resulting in high variance low bias. All of these trees predictions will be chosen by majority voting for the final prediction.
 - Random Forest: We will create multiple decision trees similar to bagging, but try to decorrelate each of the bootstrap trees through selecting $m = \sqrt{p}$ variables. Final predictions will be through majority voting.
 - Boosting: We will fit multiple (weak) trees sequentially, grown on information from the previously grown tree. Final prediction is a weighted of the weak learners
- SVM: We will fit support vector machines with different kernels (Linear, Polynomial, Radial). In order to select the best possible support vector machines, we will use k-fold cross validation to tune the cost parameter to obtain the lowest misclassification rate.
- KNN: We will fit a K-Nearest Neighbours model with optimal K selected by cross validation.

4 Analysis and Discussion

4.1 Linear Models

```
## [1] 0.258828
##
## fit.pred  0  1
##          0 609 132
##          1 130 596
## [1] 0.1785958
```



4.2 Discriminant Analysis

```
## Call:
## lda(price_above ~ longitude + latitude + minimum_nights + calculated_host_listings_count +
##      availability_365 + reviews_per_month + room_type + neighbourhood_group,
##      data = train)
##
## Prior probabilities of groups:
##      0      1
## 0.5024723 0.4975277
##
## Group means:
##   longitude latitude minimum_nights calculated_host_listings_count
## 0 -73.93804 40.72435      5.993553      3.247031
## 1 -73.96637 40.73139      8.462303     12.267306
##   availability_365 reviews_per_month room_typePrivate room
```

```

## 0      102.3146      1.152922      0.7536478
## 1      122.8955      1.000398      0.1583276
## room_typeShared room neighbourhood_groupBrooklyn
## 0      0.042416016      0.4920258
## 1      0.004455106      0.3296779
## neighbourhood_groupManhattan neighbourhood_groupQueens
## 0      0.2857143      0.17848660
## 1      0.6028101      0.05962988
## neighbourhood_groupStaten Island
## 0      0.012894469
## 1      0.003427005
##
## Coefficients of linear discriminants:
##                                LD1
## longitude                    -7.2028485766
## latitude                     -2.3730643294
## minimum_nights               -0.0032140829
## calculated_host_listings_count -0.0009975322
## availability_365              0.0015398787
## reviews_per_month           -0.0226183788
## room_typePrivate room        -2.2892025706
## room_typeShared room        -2.7232244440
## neighbourhood_groupBrooklyn  -0.1307428980
## neighbourhood_groupManhattan  0.7403805159
## neighbourhood_groupQueens     0.1937494533
## neighbourhood_groupStaten Island -2.0321539095
##
## pred
## obs  0  1
##    0 598 141
##    1 128 600
## [1] 0.1833674
##
## Call:
## qda(price_above ~ longitude + latitude + minimum_nights + calculated_host_listings_count +
##      availability_365 + reviews_per_month + room_type + neighbourhood_group,
##      data = train)
##
## Prior probabilities of groups:
##      0      1
## 0.5024723 0.4975277
##
## Group means:
## longitude latitude minimum_nights calculated_host_listings_count
## 0 -73.93804 40.72435      5.993553      3.247031
## 1 -73.96637 40.73139      8.462303     12.267306
## availability_365 reviews_per_month room_typePrivate room
## 0      102.3146      1.152922      0.7536478
## 1      122.8955      1.000398      0.1583276
## room_typeShared room neighbourhood_groupBrooklyn
## 0      0.042416016      0.4920258
## 1      0.004455106      0.3296779
## neighbourhood_groupManhattan neighbourhood_groupQueens
## 0      0.2857143      0.17848660

```

```
## 1          0.6028101          0.05962988
##  neighbourhood_group Staten Island
## 0          0.012894469
## 1          0.003427005

##  pred
## obs    0    1
##  0 2199  748
##  1  380 2538
## [1] 0.1923274
```

4.3 Tree-based Methods

4.3.1 Classification and Regression Trees

Classification Tree for Price Above Median

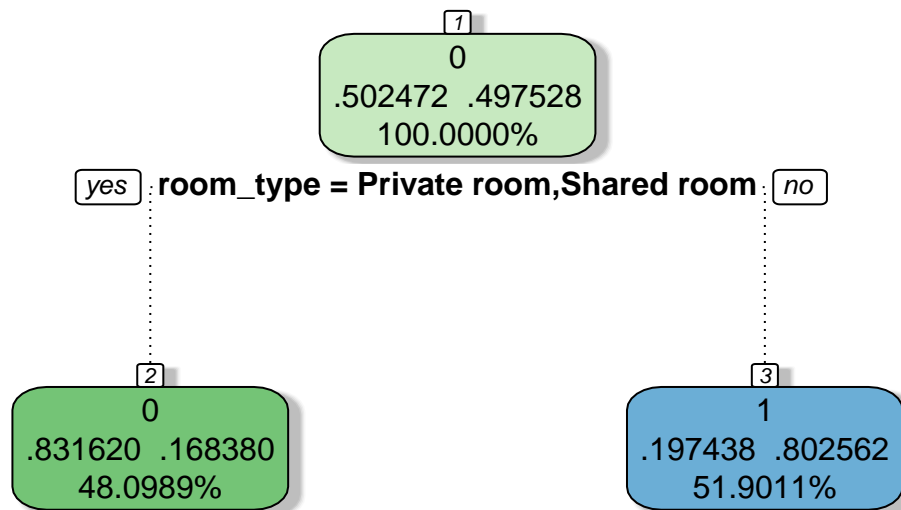
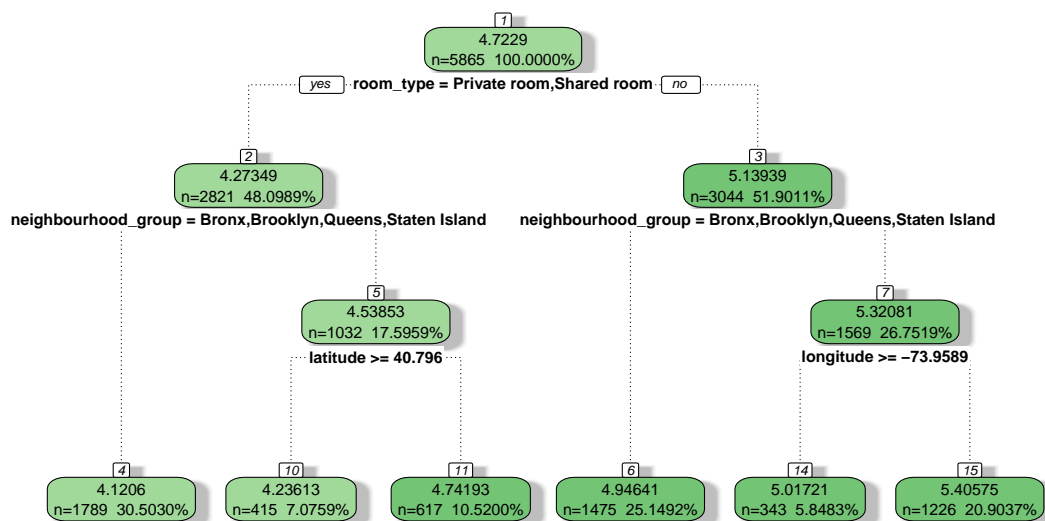


Figure 6: Classification Tree

The classification tree on all predictors split only on the type of room. We can see from the dendrogram that if a room is a private room or a shared room, the listing would be classified as “below the median price,” and if it is a whole apartment or home then it would be classified as “below the median price.”

Regression Tree for log(Price)



The regression tree on is more intricate than the classification tree. The main split is on the room type of the listing, and then the next two splits are made on the neighbourhood group. The last splits made are on the location of the the listing.

4.3.2 Random Forests

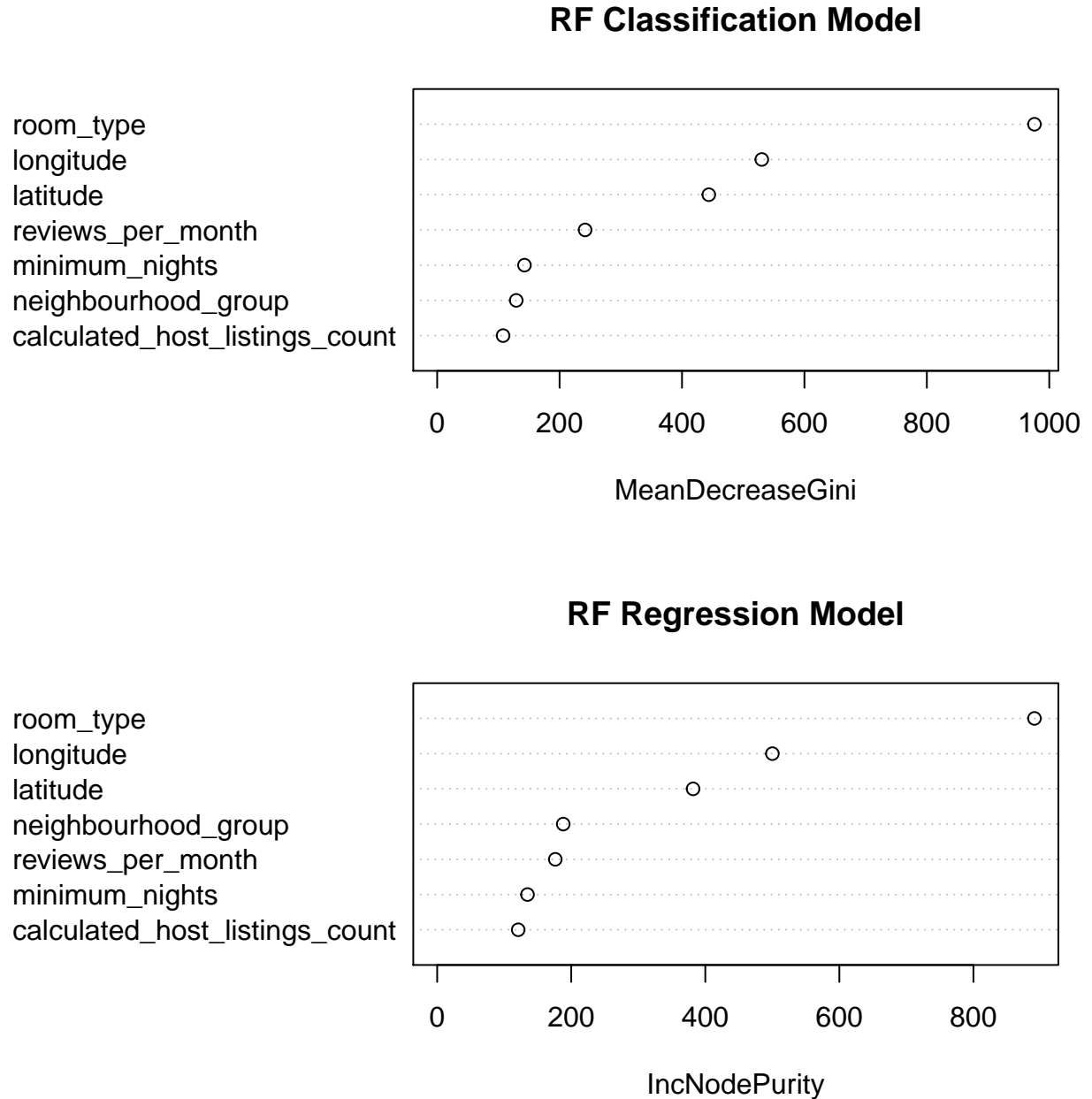


Figure 8: Variable Importances for RF Models

The random forest model importances showed that room type, longitude, latitude, and reviews per month were the most important variables for the decrease in gini impurity for classification forests. For the regression model, room type, longitude, latitude, and neighbourhood group were the most important variables for the increase in node purity.

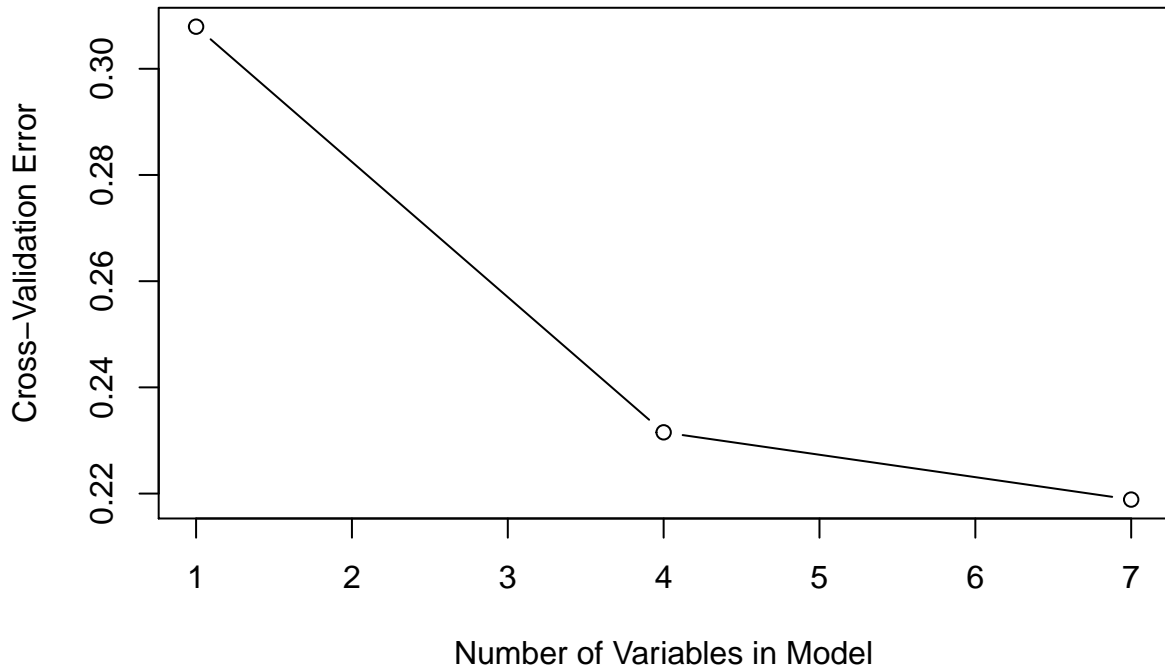


Figure 9: Cross Validation Error by # of Predictors

Through cross validation, we were able to determine 4 variables in the model leads to the greatest decrease in the cross-validation error while accounting for model complexity. To choose the 4 variables to use in the reduced random forest model, I used the criteria of greatest variable importance from above. This means we chose room type, longitude, latitude, and reviews per month for the classification model and room type, longitude, latitude, and neighbourhood group of a regression model.

The misclassification rate for the random forest model with all the predictors included was 0.1731425, as opposed to the 0.1799591 misclassification rate of the reduced model. While the RF model with all the predictors is more accurate than the smaller model, the smaller model is simpler and more likely to be scalable in different scenarios. The mean squared prediction error of the full model, 0.222965, is also lower than the 0.2435506 MSPE of the smaller model. As with the classification forest, the simpler model is more scalable at the cost of prediction error. Another downside of the larger model is the possibility of overfitting to the training dataset.

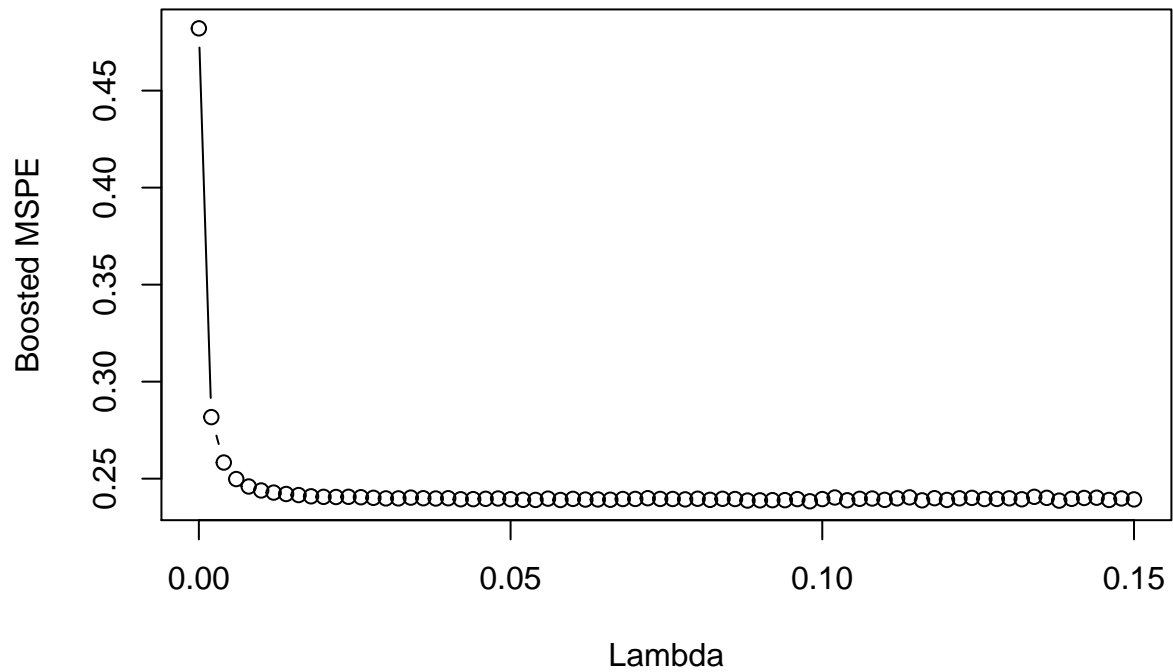
4.3.3 Bootstrap-Aggregating (Bagging)

The bagging model was similar to the random forest model with a misclassification rate of 0.1820041 and a MSPE of 0.2426593.

4.3.4 Boosting

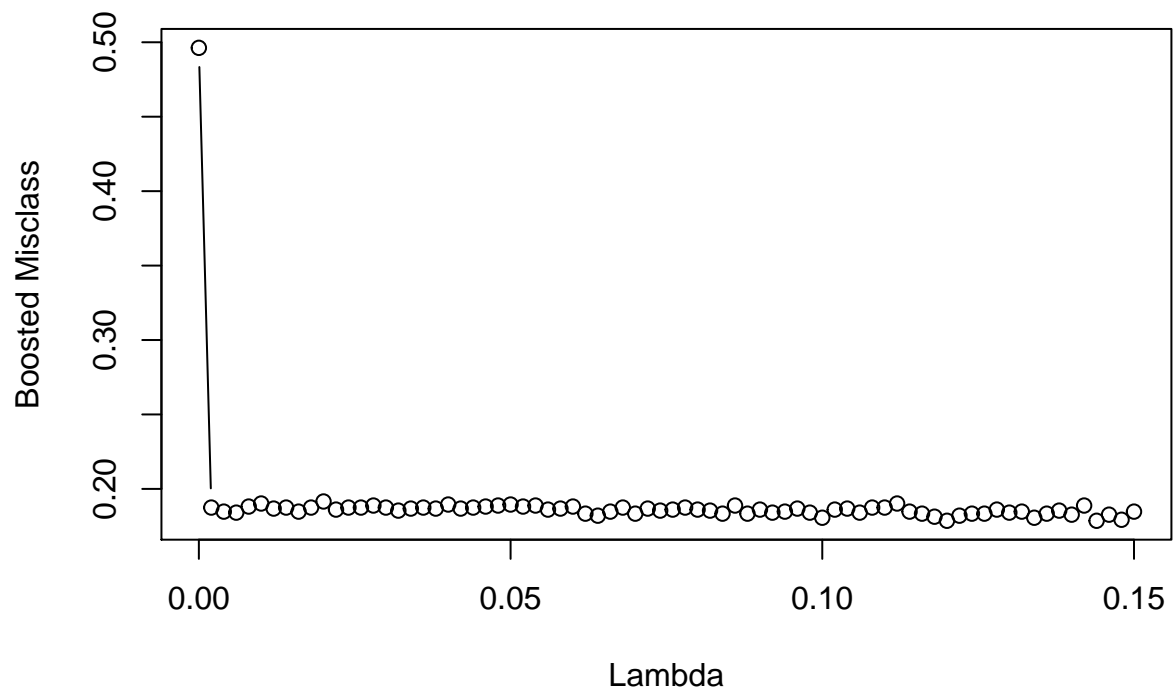
The first parameter we tuned in the boosting model was the shrinkage parameter, λ . The lambda which results in the lowest MSPE will be the lambda used in the final boosting model.

MSPE vs. Lambdas



There was not a discernable optimal lambda for test MSE, so we decided to use the default value of $\lambda = 0.1$. This made model parameter selection the easiest.

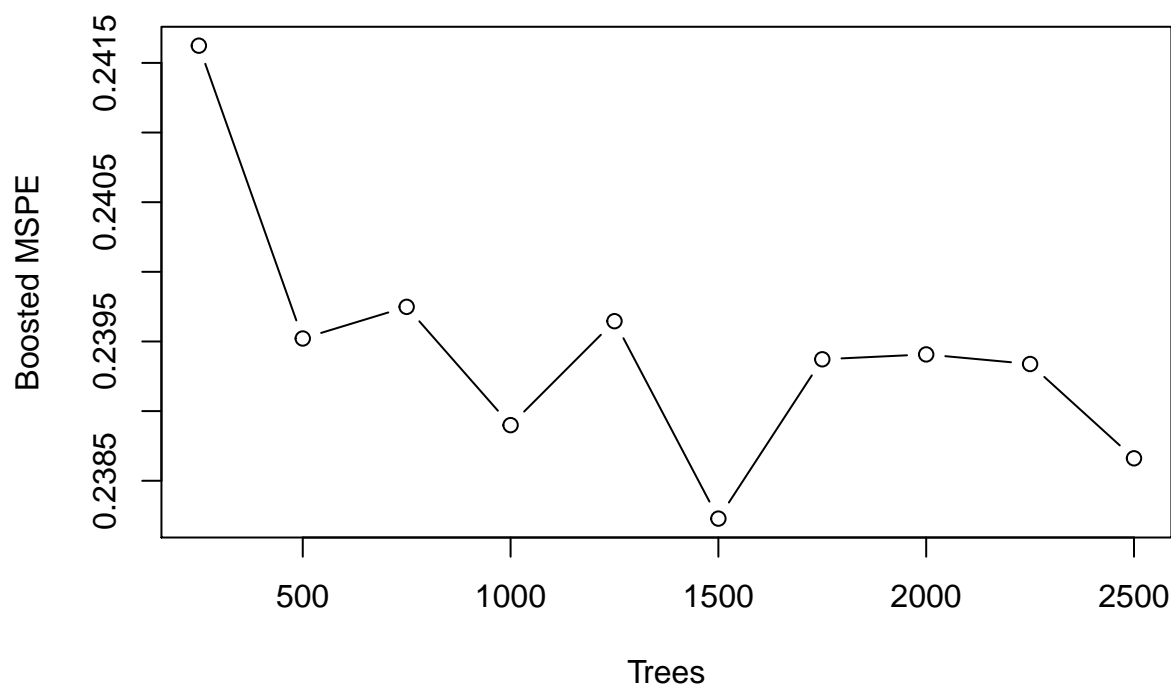
Misclassification. vs. Lambdas



As with the regression boosting model, the misclassification error rate does not seem to be at a minimum for any value of lambda, so we will also choose the default value for λ .

The other parameter tuned for the boosting model was the number of trees used.

MSPE vs. Number of Trees



The MSPE was the lowest for the model with 1500 trees, so we implemented this into the final boosting model.

4.3.5 Tree Method Summary

All The tables for tree methods put together. This will allow us to evaluate the efficacy of our tree models and determine which is the best.

Table 1: Error Rates for Tree Models

Methods	MSPE	Methods	Misclassification
Tree	0.25993	Tree	0.187457
Bagging	0.242659	Bagging	0.182004
Boosting	0.240589	Boosting	0.184731
Random Forest	0.222965	Random Forest	0.173142
Reduced Random Forest	0.243551	Reduced Random Forest	0.179959

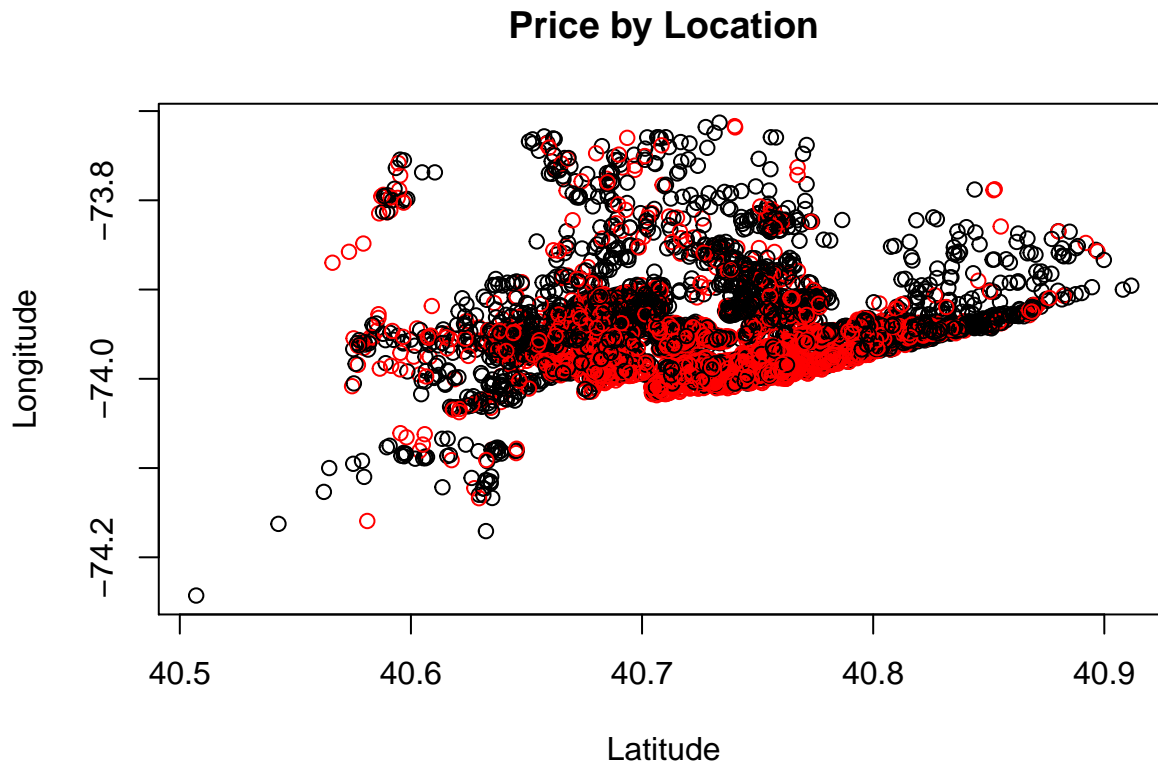
The table with all the tree method error rate shows there is not a clear “best model” for the data. The random forest with all predictors has the lowest test MSE and misclassification error rate, but suffers from the possibility of overfitting the data and being too complex of a model. Because all of the methods are approximately equal in terms of predictive performance, a simpler model such as a simple tree or the reduced random forest may be best. Both of these models are simpler, and would therefore be more scalable in larger-data environments.

4.4 SVM

In order to gain a rudimentary understanding of our data, we plotted a scatterplot Airbnb properties in New York coded by price above or below the median. It appears that the price of properties is noticeably higher along waterfront properties.

```
train$price_above <- as.factor(train$price_above)
```

```
plot(rbind(train, test)$latitude,  
     rbind(train, test)$longitude,  
     col = rbind(train, test)$price_above,  
     main = "Price by Location",  
     xlab = "Latitude",  
     ylab = "Longitude")
```



4.4.1 Best Linear Kernel SVM

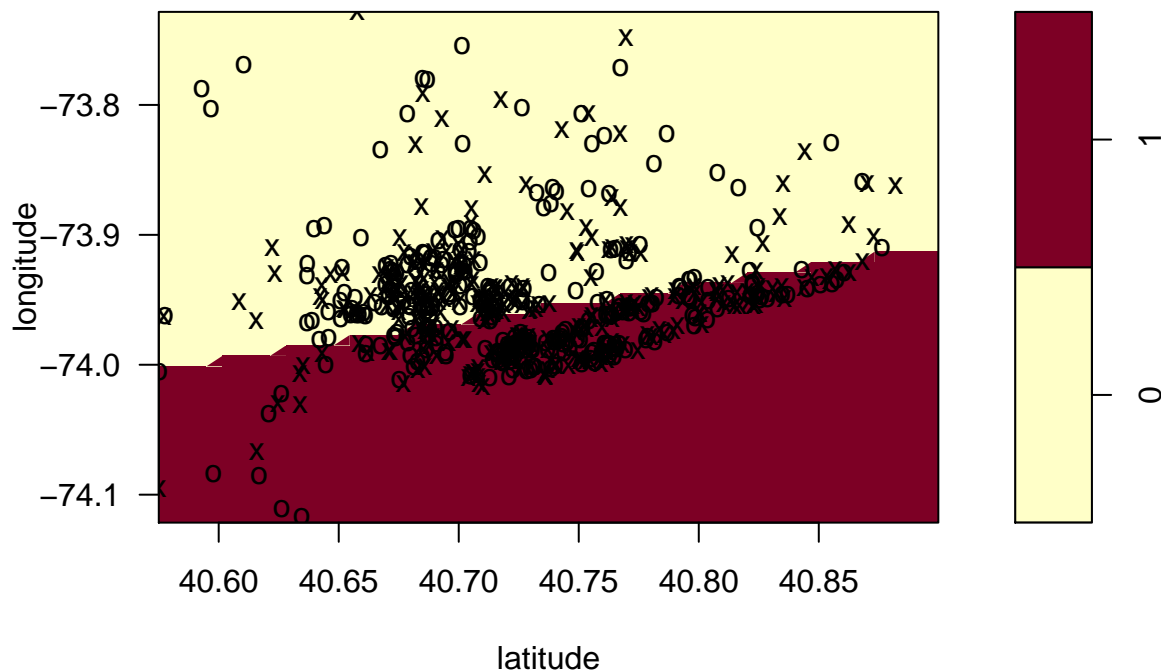
- Misclass. Rate: 0.1867757
- Cost Parameter: 0.06
- Support Vectors: 4598

Below is a preliminary plot of a linear kernel SVM model fit off of latitude and longitude. After including other predictors including `neighbourhood_group`, `minimum_nights`, `room_type`, `number_of_reviews`, `calculated_host_listings_count`, `availability_365`, and using 10-fold cross validation to select the best cost parameter, we obtain a cost parameter with value 0.06. This gives us a test misclassification rate of 0.1867757.

```
# determine approximate best cost parameter  
# tune.linear <- tune(svm, price_above ~ latitude+longitude, data = train, kernel = "linear", ranges =  
# increase precision of cost parameter
```

```
# tune.linear <- tune(sum, price_above ~ longitude
#                               +latitude
#                               +neighbourhood
#                               +minimum_nights
#                               +room_type
#                               +minimum_nights
#                               +number_of_reviews
#                               +calculated_host_listings_count
#                               +availability_365, data = train, kernel = "linear", ranges = list(cost = seq(.05,
plot(svm(price_above ~ longitude+latitude, data = train, kernel = "linear", cost = 0.06), test[,c("price_above", "longitude", "latitude")])
```

SVM classification plot



```
best.linear <- svm(price_above ~ longitude
+latitude
+neighbourhood_group
+minimum_nights
+room_type
+number_of_reviews
+calculated_host_listings_count
+availability_365, data = train, kernel = "linear", cost = 0.06)
pred.linear <- predict(best.linear, test)

# confusion matrix
conf.linear <- table(obs = test$price_above, pred = pred.linear)
(acc.linear <- 1 - sum(diag(conf.linear)/sum(conf.linear)))

## [1] 0.1874574
```

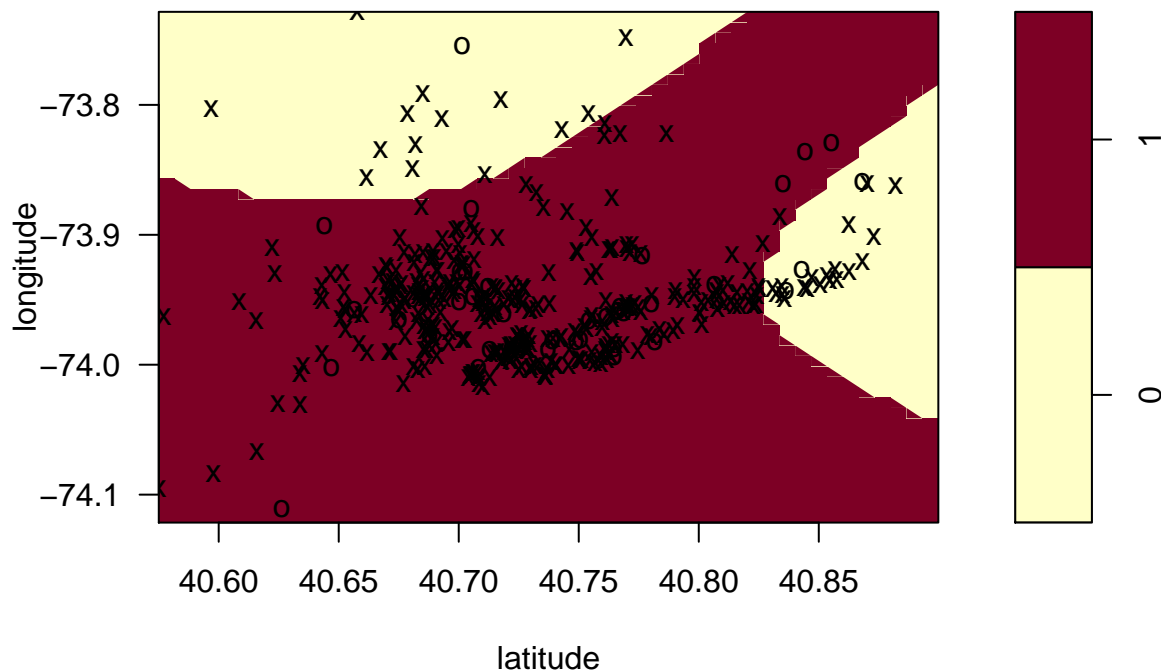
4.4.2 Best Polynomial Kernel SVM

- Misclass. Rate: 0.5132924
- Cost Parameter: 100
- Degree: 3
- Support Vectors: 5737

Below is a plot of a polynomial kernel SVM decision boundary fit off of `longitude` and `latitude`. Fitting the model with identical predictors as our linear kernel SVM and using 10-fold cross validation, we obtain an optimal cost parameter value of 100. Our test misclassification rate is 0.5132924.

```
# tune.poly <- tune(svm, price_above ~ longitude
#                   +latitude
#                   +neighbourhood_group
#                   +minimum_nights
#                   +room_type
#                   +minimum_nights
#                   +number_of_reviews
#                   +calculated_host_listings_count
#                   +availability_365, data = train, kernel = "polynomial", ranges = list(cost = c(10
plot(svm(price_above ~ longitude+latitude, data = train, kernel = "polynomial", cost = 100), test[,c("p
```

SVM classification plot



```
best.poly <- svm(price_above ~ longitude
+latitude
+neighbourhood_group
+minimum_nights
+room_type
+number_of_reviews
+calculated_host_listings_count
+availability_365, data = train, kernel = "polynomial", cost = 100)
```

```

pred.poly <- predict(best.poly, test)
# (MSE.poly <- mean((as.numeric(test$price_above) - as.numeric(pred.poly))^2))

# confusion matrix
conf.poly <- table(obs = test$price_above, pred = as.factor(pred.poly))
(acc.poly <- 1 - sum(diag(conf.poly)/sum(conf.poly)))

## [1] 0.1738241

```

4.4.3 Best Radial Kernel SVM

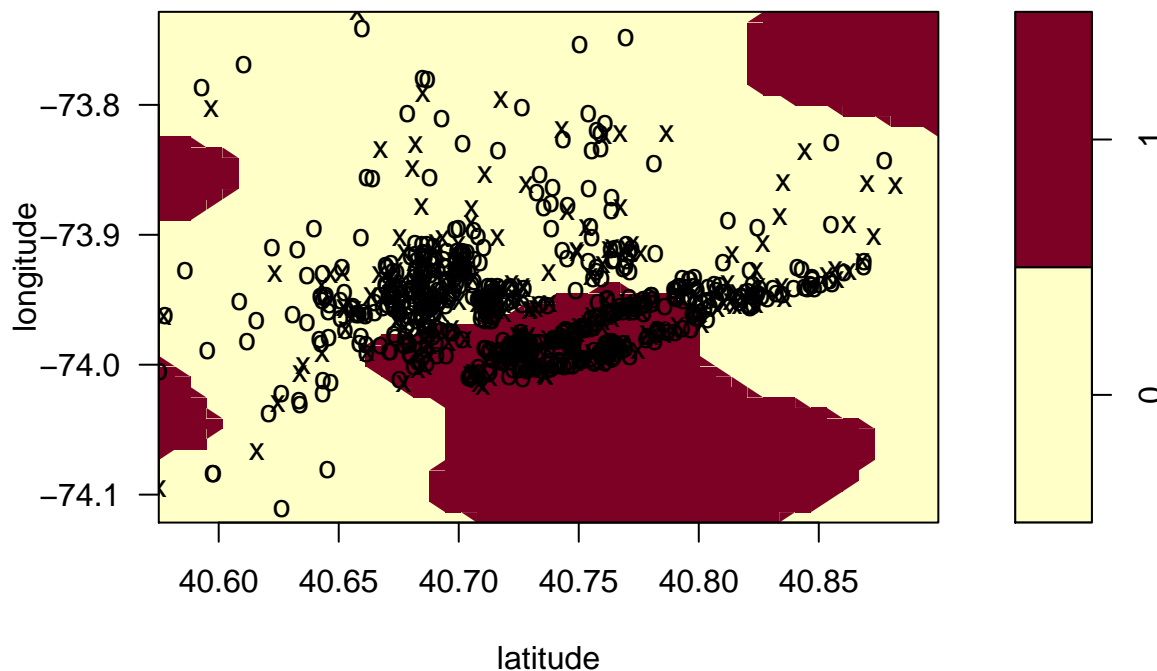
- Misclass. Rate: 0.1785958
- Cost Parameter: 8
- Support Vectors: 3615 Below is a plot of the decision boundary of a radial kernel SVM fit on **latitude** and **longitude**. Fitting a model with identical parameters as the previous two models, we obtain a cost parameter value of 8 through 10-fold cross validation. Our test misclassification rate is 0.1785958.

```

# tune.rad <- tune(svm, price_above ~ longitude+latitude, data = train, kernel = "radial", ranges = lis
tune.rad <- tune(svm, price_above ~ longitude
                +latitude
                +neighbourhood_group
                +minimum_nights
                +room_type
                +number_of_reviews
                +calculated_host_listings_count
                +availability_365, data = train, kernel = "radial", ranges = list(cost = seq(8, 11,
plot(svm(price_above ~ longitude+latitude, data = train, kernel = "radial", cost = 8), test[,c("price_al

```

SVM classification plot



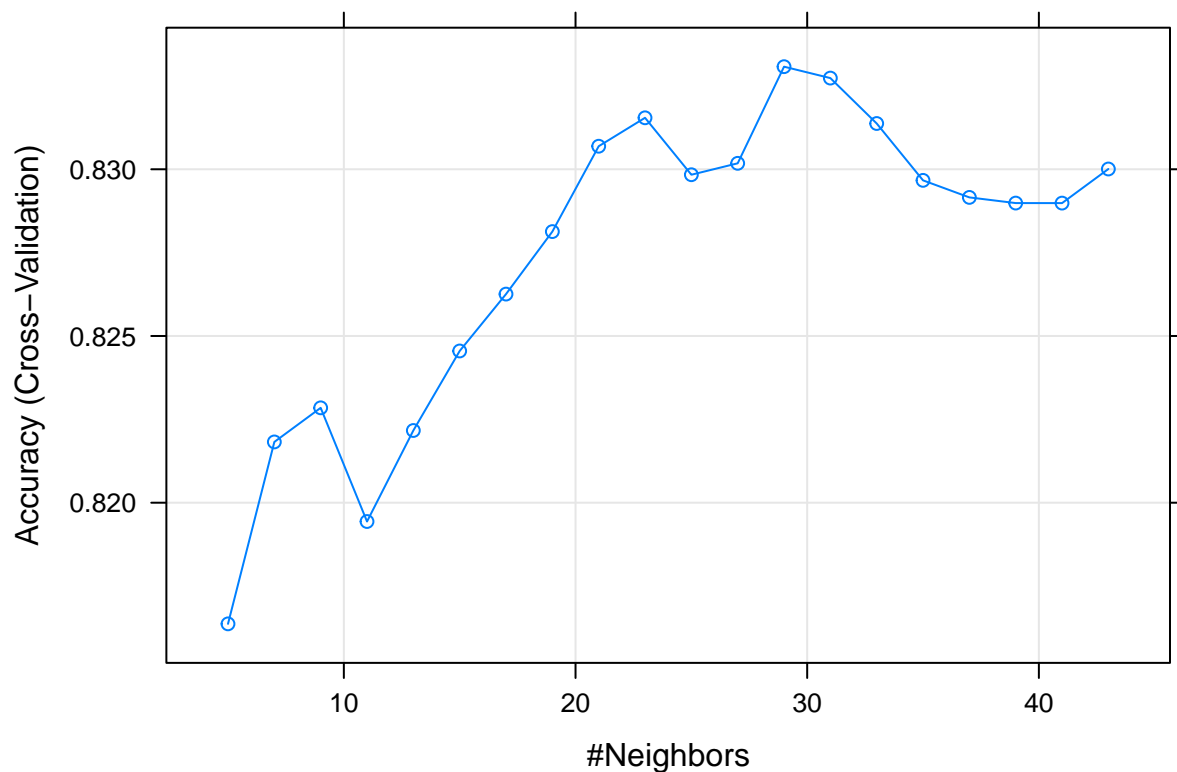
```
best.rad <- svm(price_above ~ longitude
               +latitude
               +neighbourhood_group
               +minimum_nights
               +room_type
               +number_of_reviews
               +calculated_host_listings_count
               +availability_365, data = train, kernel = "radial", cost = 8)
pred.rad <- predict(best.rad, test)
conf.rad <- table(obs=test$price_above, pred=pred.rad)
(acc.rad <- 1 - sum(diag(conf.rad)/sum(conf.rad)))
```

```
## [1] 0.1792774
```

4.5 KNN

Using cross validation to find k such that the cross validated error is minimized, we find that $k = 27$ gives us the minimum cross validated error of 0.1642808.

```
set.seed(3)
knn.caret <- train(price_above ~ neighbourhood_group+latitude+longitude+room_type+minimum_nights+number_of_reviews, data = train, method = "knn", plot = TRUE)
plot(knn.caret)
```



```
pred.knn <- predict(knn.caret, test)
(conf.matrix <- table(pred = pred.knn, obs = test$price_above))
```

```
##      obs
## pred  0   1
```

```
##      0 605 109
##      1 134 619

1 - sum(diag(conf.matrix))/sum(conf.matrix)

## [1] 0.1656442
```

4.6 Regularization

4.6.1 Ridge Regression

```
X <- data.matrix(num.feats)
#X <- X[, -c(10,2,4,5,6,9,13,14)]
y <- log(train$price)
Ridge <- glmnet(x = X, y = y, alpha = 0)
RidgeCV <- cv.glmnet(x = X, y = y, alpha = 0, lambda = Ridge$lambda, nfolds = 10)
lambda.ind <- which.min(RidgeCV$cvm)
lambda.best <- Ridge$lambda[lambda.ind]
lambda.best

## [1] 0.07057983
Ridge$beta[, lambda.ind]

##              longitude              latitude
##      -5.022231e-01          1.262035e-01
##      minimum_nights      number_of_reviews
##      2.251623e-05          2.646160e-05
##      reviews_per_month calculated_host_listings_count
##      4.951531e-04          2.390493e-04
##      price              lprice
##      2.811029e-04          8.425852e-01

xTrain <- num.feats %>% dplyr::select(-c(price,lprice)) %>% as.matrix()
yTrain <- num.feats$price
# Fit model using the best lambda from above
proRidgeT <- glmnet(x = xTrain, y = yTrain, alpha = 0, lambda = lambda.best)
# Use testing data and fitted model to predict
XTest <- test %>% dplyr::select(longitude, latitude, minimum_nights, number_of_reviews,
                              reviews_per_month, calculated_host_listings_count) %>% as.matrix()
yTest <- log(test$price)
yPredRidge <- proRidgeT$a0 + XTest%*%proRidgeT$beta
# Compute mean-squared prediction error
proMSPE.Ridge <- mean((yTest - yPredRidge)^2)
proMSPE.Ridge

## [1] 0.4183765
```

4.6.2 Lasso

```
X <- data.matrix(bnb)
X <- X[, -c(10,2,4,5,6,9,13,14)]
y <- log(bnb$price)
```

```

#X <- num.feats %>%
Lasso <- glmnet(x = X, y = y, alpha = 1)

LassoCV <- cv.glmnet(x = X, y = y, alpha = 1, lambda = Lasso$lambda, nfolds = 10)
lambda.indL <- which.min(LassoCV$cvm)
lambda.bestL <- Lasso$lambda[lambda.indL]
lambda.bestL

## [1] 0.001305949

Lasso$beta[, lambda.indL]

##              id              host_id
##      0.000000e+00      4.351725e-11
##      latitude      longitude
##      5.490008e-01      -1.671619e+00
##      minimum_nights      number_of_reviews
##      -4.791877e-04      -3.442020e-04
##      calculated_host_listings_count      availability_365
##      2.413531e-04      3.356195e-04
##      price_above
##      1.045644e+00

xTrain <- num.feats %>% dplyr::select(-c(price,lprice)) %>% as.matrix()
yTrain <- num.feats$price
# Fit model using the best lambda from above
proLassoT <- glmnet(x = xTrain, y = yTrain, alpha = 1, lambda = lambda.best)
# Use testing data and fitted model to predict
XTest <- test %>% dplyr::select(longitude, latitude, minimum_nights, number_of_reviews,
                                reviews_per_month, calculated_host_listings_count) %>% as.matrix()

yTest <- log(test$price)
yPredLasso <- proLassoT$a0 + XTest*%proLassoT$beta
# Compute mean-squared prediction error
proMSPE.Lasso <- mean((yTest - yPredLasso[x])^2)
proMSPE.Lasso

## [1] 0.4289998

```

5 Conclusion

There were many considerations to be taken into account when selecting our final model. First, we had to consider the accuracy of the model; if our test error rate was too high, there would be no reason for to implement a poorly-constructed model. Next was model complexity: we could use a model with multiple predictors and parameters, decreasing the error rate of our model, but increasing the variance, or we could construct a simpler model with more bias, but lower variability. Choosing a simpler model would allow us to use our model in a larger setting and achieve similar test error rates. We can compare the test error rates of our different models and determine a best model.

Table 2: Error Rates for Tree Models

Methods	MSPE	Methods	Misclassification
Tree	0.25993	Tree	0.187457
Bagging	0.242659	Bagging	0.182004
Boosting	0.240589	Boosting	0.184731
Random Forest	0.222965	Random Forest	0.173142
Reduced Random Forest	0.243551	Reduced Random Forest	0.179959

5.1 Tree Error Rates

5.2 Final Models and Error

5.2.1 Regression

```
##
## Call:
## lm(formula = log(price) ~ latitude + longitude + minimum_nights +
##     reviews_per_month + neighbourhood_group + room_type + calculated_host_listings_count,
##     data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0405 -0.3215 -0.0597  0.2442  4.2232
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -1.398e+02  2.024e+01  -6.911 5.33e-12
## latitude      -1.108e+00  1.990e-01  -5.570 2.67e-08
## longitude     -2.569e+00  2.274e-01 -11.295 < 2e-16
## minimum_nights -1.050e-03  3.460e-04  -3.034 0.00242
## reviews_per_month 1.065e-03  4.426e-03   0.241 0.80983
## neighbourhood_groupBrooklyn -7.379e-02  6.004e-02  -1.229 0.21912
## neighbourhood_groupManhattan 3.154e-01  5.540e-02   5.694 1.30e-08
## neighbourhood_groupQueens 5.742e-02  5.854e-02   0.981 0.32670
## neighbourhood_groupStaten Island -6.991e-01  1.064e-01  -6.571 5.41e-11
## room_typePrivate room -7.491e-01  1.390e-02 -53.876 < 2e-16
## room_typeShared room -1.166e+00  4.439e-02 -26.266 < 2e-16
## calculated_host_listings_count 5.089e-04  1.986e-04   2.562 0.01042
##
## (Intercept)          ***
## latitude             ***
## longitude            ***
## minimum_nights      **
## reviews_per_month
## neighbourhood_groupBrooklyn
## neighbourhood_groupManhattan ***
## neighbourhood_groupQueens
## neighbourhood_groupStaten Island ***
## room_typePrivate room ***
## room_typeShared room ***
## calculated_host_listings_count *
## ---
```



```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5081 on 5853 degrees of freedom
## Multiple R-squared:  0.4828, Adjusted R-squared:  0.4818
## F-statistic: 496.7 on 11 and 5853 DF,  p-value: < 2.2e-16
```

The final mean squared prediction error, 0.2572045, of the model is slightly higher than some of the more complex models. This choice was made for our final model because the model is simpler and more efficient with larger datasets. While the accuracy is slightly decreased, it is not a large enough difference where we would consider another, more accurate model. Another reason for choosing the linear model is because of the ability to interpret the coefficients. Some of the more significant coefficients can be interpreted below:

Note: All interpretations are given all other predictors are held constant.

- **Latitude:** A one unit increase in latitude results in a -66.9781261% decrease in the price of the listing.
- **Longitude:** A one unit increase in longitude results in a -92.3387881% decrease in the price of the listing.
- **Neighbourhood Group:**
 - Being in the *Brooklyn* neighbourhood group results in the price of the listing being -7.1133265% lower than a listing in the Bronx.
 - Being in the *Manhattan* neighbourhood group results in the price of the listing being 37.0807524% higher than a listing in the Bronx.
 - Being in the *Queens* neighbourhood group results in the price of the listing being 5.9100539% higher than a listing in the Bronx.
 - Being in the *Staten Island* neighbourhood group results in the price of the listing being -50.2967568% lower than a listing in the Bronx.
- **Room Type**
 - Being a *Private Room* results in the price of the listing being -52.7208126% lower than whole apartment or house listing.
 - Being a *Shared Room* results in the price of the listing being -68.8389105% lower than whole apartment or house listing.
- **Calculated Host Listing Count:** A one unit increase in longitude results in a 0.050903% decrease in the price of the listing.

5.2.2 Classification

Table 3: Confusion Matrix

	0	1
0	19251	3776
1	5210	20647

The logistic regression model ended up being the best choice for our classification model. As with the regression model, the misclassification error rate, 0.1838229, was slightly higher than some of the more complex models, but we decided to choose the simpler model at a small penalty to misclassification rate. We can interpret the most important coefficients as follows:

Note: All interpretations are given all other predictors are held constant.

- **Latitude:** A one unit increase in latitude results in a -4.699 decrease in the log-odds of the listing price being above the median.
- **Longitude:** A one unit increase in longitude results in a -13.08 decrease in the log-odds of the listing price being above the median.
- **Minimum Nights:** A one unit increase in the minimum amount of nights results in a decrease of -.002933 of the listing price being above the median.
- **Neighbourhood Group:**

- Being in the *Brooklyn* neighbourhood group results in the log-odds of the price being below the median to be reduced by -.1586, compared to being a listing in the Bronx.
- Being in the *Manhattan* neighbourhood group results in the log-odds of the price being below the median to be increased by 1.411, compared to being a listing in the Bronx.
- Being in the *Queens* neighbourhood group results in the log-odds of the price being below the median to be increased by .5048, compared to being a listing in the Bronx.
- Being in the *Staten Island* neighbourhood group results in the log-odds of the price being below the median to be reduced by -3.371, compared to being a listing in the Bronx.
- **Room Type**
 - Being a *Private Room* results in the log-odds of the price being below the median to be reduced by -3.078, compared to a whole apartment or house listing.
 - Being a *Shared Room* results in the log-odds of the price being below the median to be reduced by -4.019, compared to a whole apartment or house listing.
- **Calculated Host Listing Count:** A one unit increase in longitude results in a .00375 increase in the price of the listing.

6 Appendix

```
# Anything below here before the abstract has to be here so we can run code in the analysis and have ou
library(tidyverse)
library(dplyr)
library(randomForest)
library(class)
library(tree)
library(gbm)
library(caret)
library(rpart.plot)
library(rattle)
library(knitr)
library(fastAdaboost)
library(ggpubr)
library(MASS)
library(kableExtra)
library(glmnet)
library(faraway)
library(ROCR)
library(e1071)
#library(train)
```

Read in the CSV and check the dimensions of the data.

```
bnb <- read.csv('AB_NYC_2019.csv')

bnb <- as_tibble(bnb)

dims <- dim(bnb)

#sprintf('Our dataset has %d observations and %d attributes', dims[1],dims[2])
# [1] "Our dataset has 48895 observations and 16 attributes"
```

Since observations which have a price of 0 will not be useful to our analysis, and are likely to be representative of a bad data point, we will remove these observations.

```
bnb <- bnb[(bnb$price!=0),]
bnb[is.na(bnb$reviews_per_month), 'reviews_per_month'] <- 0
bnb$price_above <- ifelse(bnb$price > median(bnb$price), 1, 0)
```

Train-test split the data

```
set.seed(123)
train.ind <- sample(1:nrow(bnb), size = .8*nrow(bnb))
small.ind <- sample(1:nrow(bnb), size = .15*nrow(bnb))

train.big <- bnb[train.ind,]
test.big <- bnb[-train.ind,]

small <- bnb[small.ind,]
train.small <- sample(1:nrow(small), size = .8*nrow(small))
train <- small[train.small,]
test <- small[-train.small,]
```

We can get the column names of the columns which contain NA's with the following code:

```
#colnames(train)[apply(train, 2, anyNA)]
```

We can see that there are NA reviews in the `reviews_per_month`, the number of reviews per month. We can also see upon visual inspection `last_review`, the date of the last review the host received, also contains empty values. We will not be using `last_review` in our analysis, so we will not worry about imputing values here.

We believe the reason there are NA's in the `reviews_per_month` column is because the hosts have 0 reviews overall. We further explore this claim the below:

```
#with(train, sum((is.na(reviews_per_month)) & (number_of_reviews!=0)) )
# [1] 0
```

We can see there are no cases where the `number_of_reviews` and `reviews_per_month`. As a result, we will impute 0 where `reviews_per_month` is NA.

```
train[is.na(train$reviews_per_month), 'reviews_per_month'] <- 0
test[is.na(test$reviews_per_month), 'reviews_per_month'] <- 0

#sum(is.na(train$reviews_per_month))
# [1] 0
#sum(is.na(test$reviews_per_month))
# [1] 0
```

We can assess the correlation between numeric features with a correlation heatmap:

```
num.feats <- train %>% dplyr::select(longitude, latitude, minimum_nights, number_of_reviews,
                                   reviews_per_month, calculated_host_listings_count, price)
num.feats$price <- log(num.feats$price)
feat.corr <- cor(num.feats)
corrplot::corrplot(feat.corr)

pairs(num.feats)

# levels(bnb$room_type)
#
# levels(bnb$neighbourhood_group)
```

```
#
# n_distinct(bnb$neighbourhood)
#
# n_distinct(bnb$neighbourhood_group)

# [1] "Entire home/apt" "Private room"      "Shared room"
# [1] "Bronx"           "Brooklyn"        "Manhattan"      "Queens"          "Staten Island"
# [1] 221
# [1] 5
```

There are 221 neighborhoods covered in the overall data, but only 5 neighbourhood groups. We will further investigate whether we need to use the neighbourhood, or whether we would like to use the neighbourhood groups for simplicity of our model.

We will determine whether we should use the neighbourhood by seeing if there is a large disparity in mean price by calculating the mean price for the neighbourhood. If there seems to be large disparities within the neighbourhood group for mean pricing, we will attempt to use neighbourhood itself.

```
hist(log(train$price), main='Histogram of log(Price)', xlab = 'log(Price)')

n <- ggplot(data = train, aes(x = latitude, y = longitude, color = neighbourhood)) +
  geom_point() + theme(legend.position="none") + xlab('Latitude') + ylab('Longitude') +
  ggtitle('Neighbourhood')

ng <- ggplot(data = train, aes(x = latitude, y = longitude, color = neighbourhood_group)) +
  geom_point() + theme(legend.position="none") + xlab('Latitude') + ylab('Longitude') +
  ggtitle('Neighbourhood Groups')

fig <- ggarrange(n, ng, nrow = 1, ncol = 2)

annotate_figure(fig, top = 'New York City Airbnb Listing Locations By:')

mean_price <- train %>% group_by(neighbourhood) %>% summarise(mean_price = mean(price),
                                                             latitude = median(latitude), longitude =
                                                             median(longitude))

#plot(mean_price$latitude, mean_price$longitude, col = mean_price$mean_price)

ggplot(data = mean_price, aes(x = latitude, y = longitude, color = mean_price)) +
  geom_point() + scale_color_gradient(low="blue", high="red", name = 'Mean Price (USD)') + xlab('Latitude') +
  ylab('Longitude') + ggtitle('Mean Price by Neighborhood')
```

It does not seem there are any large disparities in pricing, and all the neighbourhood groups seems to be similar to their nearby neighbours. To reduce the complexity of our model, we will use the neighbourhood group.

```
par(mfrow = c(1,2))
hist(bnb$price, main='Histogram of Overall Price', xlab = 'Price (USD)')
hist(train$price, main='Histogram of Training Price', xlab = 'Price (USD)')

num.feats <- train %>% dplyr::select(longitude, latitude, minimum_nights, number_of_reviews,
                                   reviews_per_month, calculated_host_listings_count, price)
num.feats$price <- log(num.feats$price)
feat.corr <- cor(num.feats)
corrplot::corrplot(feat.corr)
```

```

set.seed(123)
ols.price <- lm(log(price) ~ latitude + longitude + minimum_nights + reviews_per_month +
                neighbourhood_group + room_type + calculated_host_listings_count, data = train)
ols.pred <- predict(ols.price, test)

ols.mspe <- mean((log(test$price)-ols.pred)^2)
ols.mspe

#summary(ols.price)

```

```

#initial LR model
set.seed(123)
logit.fit <- glm(price_above ~longitude + latitude+ minimum_nights+ calculated_host_listings_count+
                reviews_per_month + room_type + neighbourhood_group, data=train,
                family=binomial('logit'))

fit.pred <- predict(logit.fit, test, type="response")
fit.pred <- ifelse(fit.pred>.5, 1, 0)
table(fit.pred,test$price_above)
mean(fit.pred!=test$price_above)

pred <- prediction(fit.pred, test$price_above)
perf <- performance(pred, "tpr","fpr")
plot(perf, col=2, lwd=3, main="ROC curve")
abline(0,1)

auc = performance(pred, "auc")@y.values

```

```

set.seed(123)
class.tree <- rpart(as.factor(price_above)~latitude + longitude + minimum_nights + reviews_per_month +
                neighbourhood_group + room_type + calculated_host_listings_count, data = train)

tree.class.prediction <- predict(class.tree, test, type = 'class')
tree.class.misclass <- mean(test$price_above !=tree.class.prediction)

fancyRpartPlot(class.tree, digits = 6, main = 'Classification Tree for Price Above Median', sub = '')

```

```

set.seed(123)

#Pruning tree did not improve tree

# prune <- prune.tree(loc.tree, best = 4, newdata = test)
# plot(prune)
# text(prune)
# tree.class.misclass

```

```

set.seed(123)
tree.reg <- rpart(log(price) ~ latitude + longitude + minimum_nights + reviews_per_month +
                neighbourhood_group + room_type + calculated_host_listings_count, data = train)

reg.prediction <- predict(tree.reg, test)
tree.mspe <- mean((log(test$price)-reg.prediction)^2)

fancyRpartPlot(tree.reg, digits = 6, sub = '', main = 'Regression Tree for log(Price)')

```

```

set.seed(123)
rf.class <- randomForest(as.factor(price_above) ~ latitude + longitude + minimum_nights + reviews_per_month +
                        neighbourhood_group + room_type + calculated_host_listings_count,
                        data = train, importance = TRUE)

#randomForest::importance(rf.class)

rf.class.pred <- predict(rf.class, test)
rf.misclass <- mean(test$price_above != rf.class.pred)

#sprintf('The misclassification rate for the classification random forest is %f', rf.misclass)
#varImpPlot(rf.class)

```

```

set.seed(123)
rf.reg <- randomForest(log(price) ~ latitude + longitude + minimum_nights + reviews_per_month +
                      neighbourhood_group + room_type + calculated_host_listings_count,
                      data = train, importance = TRUE)

rf.reg.pred <- predict(rf.reg, test)
rf.mspe <- mean((log(test$price) - rf.reg.pred)^2)

#sprintf('The mean squared prediction error for the regression random forest is %f', rf.mspe)

```

There still does not seem to be much of an improvement over the regression tree. We can try to re-evaluate the random forest model through cross validation and seeing if we can select important features.

The variable importance plot shows us that `room_type`, `longitude`, `latitude`, and `reviews_per_month` are the most important variables.

```

par(mfrow = c(2,1))
varImpPlot(rf.class, main = 'RF Classification Model')
varImpPlot(rf.reg, main = 'RF Regression Model')

set.seed(123)
rf.cv.trainx <- train %>% dplyr::select(latitude, longitude, minimum_nights, reviews_per_month,
                                       neighbourhood_group, room_type, calculated_host_listings_count)
rf.cv.trainy <- log(train$price)
cv.rf <- rfcv(rf.cv.trainx, rf.cv.trainy, cv.fold = 5)
plot(cv.rf$n.var, cv.rf$error.cv, type = 'b', xlab = 'Number of Variables in Model', ylab = 'Cross-Validated Error')

```

As we can see, the cross validation error is the lowest when we use the most predictors. Despite this, There does not seem to be much of a decrease after there are 4 variables in the model, so we will try to fit a model with 4 variables.

We will try fitting the 4 most important variables from the regression random forest, and seeing whether this model is better, or the same, as our more complex model.

```

set.seed(123)
rf.reg.reduced <- randomForest(log(price) ~ latitude + longitude + neighbourhood_group + room_type,
                              data = train, importance = TRUE)

rfr.reg.pred <- predict(rf.reg.reduced, test)
rfr.reg.mspe <- mean((log(test$price) - rfr.reg.pred)^2)

#sprintf('The mean squared prediction error for the reduced regression random forest is %f', rf.red.mspe)

```

By reducing the number of predictors, we were able to slightly increase the MSE, while creating a much

simpler model.

Lets try bagging with the smaller subset of variables

The prediction error is about the same as it is for a random forest.

```
set.seed(123)
bag.class <- randomForest(as.factor(price_above) ~ latitude + longitude + reviews_per_month + room_type,
                          data = train, mtry = 4 , importance = TRUE)

bag.class.pred <- predict(bag.class,test)
bag.misclass <- mean(test$price_above!=bag.class.pred)

# sprintf('The misclassification rate for the bagged classification model is %f', bag.misclass)

set.seed(123)
boost.mod <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                 n.trees = 1000, cv.folds = 5, distribution = 'gaussian')
boost.pred <- predict(boost.mod, test, n.trees = 1000)

boost.mspe <- mean((log(test$price)-boost.pred)^2)

#sprintf('The mean squared prediction error for boosting is %f', boost.mspe)
```

As with all our other models, this one is about the same. We can try different values of the shrinkage and see if we can find a best model for cross validation error

```
set.seed(123)
lambdas <- seq(0,.15, .002)
b.mspe.list <- NULL

for(lambda in lambdas){
  boost.l.mod <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                    n.trees = 1000, shrinkage = lambda, distribution = 'gaussian')
  boost.l.pred <- predict(boost.l.mod, test, n.trees = 1000)

  b.mspe.list <- append(b.mspe.list,mean((log(test$price)-boost.l.pred)^2))
}
plot(lambdas,b.mspe.list, type = 'b', ylab = 'Boosted MSPE', xlab = 'Lambda',
     main = 'MSPE vs. Lambdas')
```

There does not seem to be a discernable lambda from the plot.

```
set.seed(123)
best.lambda <- lambdas[which.min(b.mspe.list)]

best.boost <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                 n.trees = 1000, distribution = 'gaussian')
best.boost <- predict(boost.mod, test, n.trees = 1000)

b.boost.mspe <- mean((log(test$price)-best.boost)^2)

#sprintf('The mean squared prediction error for bagging with the optimal lambda is is %f', b.boost.mspe)

set.seed(123)
tree.err <- NULL
```

```

ntrees <- seq(250,2500,250)

for(ntree in ntrees){
  boost.t.mod <- gbm(log(price) ~ latitude + longitude +
                    reviews_per_month + room_type, data = train, n.trees = ntree,
                    shrinkage = best.lambda, distribution = 'gaussian')
  boost.t.pred <- predict(boost.t.mod, test, n.trees = ntree)

  tree.err <- append(tree.err, mean((log(test$price) - boost.t.pred)^2))
}
#tree.err

plot(ntrees, tree.err, type = 'b', ylab = 'Boosted MSPE', xlab = 'Trees',
     main = 'MSPE vs. Number of Trees')
best.tree <- ntrees[which.min(tree.err)]

```

```

set.seed(123)
lambdas <- seq(0, .15, .002)
b.mis.list <- NULL

for(lambda in lambdas){
  boost.m.mod <- gbm(as.character(price_above) ~ latitude + longitude +
                    reviews_per_month + room_type, data = train,
                    n.trees = 1000, shrinkage = lambda, distribution = 'bernoulli')
  boost.m.pred <- predict(boost.m.mod, test, n.trees = 1000, type = 'response')
  boost.m.pred <- ifelse(boost.m.pred >= .51, 1, 0)

  b.mis.list <- append(b.mis.list, mean(test$price_above != boost.m.pred))
}
plot(lambdas, b.mis.list, type = 'b', ylab = 'Boosted Misclass', xlab = 'Lambda',
     main = 'Misclassification. vs. Lambdas')

best.class.lambda = 0.1

```

```

set.seed(123)

# class.boost <- train(as.factor(price_above) ~ latitude + longitude + reviews_per_month + room_type,
#                      method = 'gbm', data = train, verbose = FALSE)

class.boost <- gbm(as.character(price_above) ~ latitude + longitude + reviews_per_month + room_type,
                  n.trees = best.tree, data = train, distribution = 'bernoulli')

boost.class.pred <- predict(class.boost, test, n.trees = 1000, type = 'response')
boost.class.pred <- ifelse(boost.class.pred >= .51, 1, 0)

boost.misclass <- mean(test$price_above != boost.class.pred)
#boost.misclass

```

```

set.seed(123)
boost.reg.final <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                      n.trees = 1000, distribution = 'gaussian')
boost.pred <- predict(boost.reg.final, test, n.trees = 1000)

```



```

boost.mspe <- mean((log(test$price)-boost.pred)^2)

#sprintf('The mean squared prediction error for boosting is %f', boost.mspe)

set.seed(123)
lambdas <- seq(0,.15, .002)
b.mspe.list <- NULL

for(lambda in lambdas){
  boost.l.mod <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                    n.trees = 1000, shrinkage = lambda, distribution = 'gaussian')
  boost.l.pred <- predict(boost.l.mod, test, n.trees = 1000)

  b.mspe.list <- append(b.mspe.list,mean((log(test$price)-boost.l.pred)^2))
}
plot(lambdas,b.mspe.list, type = 'b', ylab = 'Boosted MSPE', xlab = 'Lambda',
     main = 'MSPE vs. Lambdas')

set.seed(123)
lambdas <- seq(0,.15, .002)
b.mspe.list <- NULL

for(lambda in lambdas){
  boost.l.mod <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                    n.trees = 1000, shrinkage = lambda, distribution = 'gaussian')
  boost.l.pred <- predict(boost.l.mod, test, n.trees = 1000)

  b.mspe.list <- append(b.mspe.list,mean((log(test$price)-boost.l.pred)^2))
}
plot(lambdas,b.mspe.list, type = 'b', ylab = 'Boosted MSPE', xlab = 'Lambda',
     main = 'MSPE vs. Lambdas')

set.seed(123)
lambdas <- seq(0,.15, .002)
b.mspe.list <- NULL

for(lambda in lambdas){
  boost.l.mod <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                    n.trees = 1000, shrinkage = lambda, distribution = 'gaussian')
  boost.l.pred <- predict(boost.l.mod, test, n.trees = 1000)

  b.mspe.list <- append(b.mspe.list,mean((log(test$price)-boost.l.pred)^2))
}
plot(lambdas,b.mspe.list, type = 'b', ylab = 'Boosted MSPE', xlab = 'Lambda',
     main = 'MSPE vs. Lambdas')

set.seed(123)
lambdas <- seq(0,.15, .002)
b.mspe.list <- NULL

for(lambda in lambdas){
  boost.l.mod <- gbm(log(price) ~ latitude + longitude + reviews_per_month + room_type, data = train,
                    n.trees = 1000, shrinkage = lambda, distribution = 'gaussian')
  boost.l.pred <- predict(boost.l.mod, test, n.trees = 1000)

```

```
b.mspe.list <- append(b.mspe.list, mean((log(test$price)-boost.l.pred)^2))  
}  
plot(lambdas, b.mspe.list, type = 'b', ylab = 'Boosted MSPE', xlab = 'Lambda',  
     main = 'MSPE vs. Lambdas')
```