# Predicting AirBnB Rental Prices

Group 11: Trevor Isaacson, Jonathan Olavarria, Jasmine
DeMeyer

12/10/2021

# Motivation

- You are looking for some additional income and decide renting on AirBnB is the best option
- How much should you rent your extra space for?

# Data

- In general, AirBnB data is very open and be easily accessed
- The original dataset is from a past Kaggle competition
  - Contained over 74,000 individual listings
- For sake of time and processing power, we took a random sample of 17,500 from those 74,000 listings
- They also provided a testing file
- Since the competition is over, we will compile our final predictions on that file using our best model

# Data

- Original data consists of 30 variables
- Variables are about the property, property location, the host and host reviews
- After cleaning and eliminating variables, our data consisted of 22 variables
- Property:
  - property_type, room_type, accommodates, bedrooms, beds, bed_type, bathrooms
- Location:
  - latitude, longitude, city
- Host:
  - cancellation_policy, cleaning_fee, host_has_profile_pic, host_identify_verified, etc

# Baseline Regression

```
linear = lm(price ~ ., data = training)
```
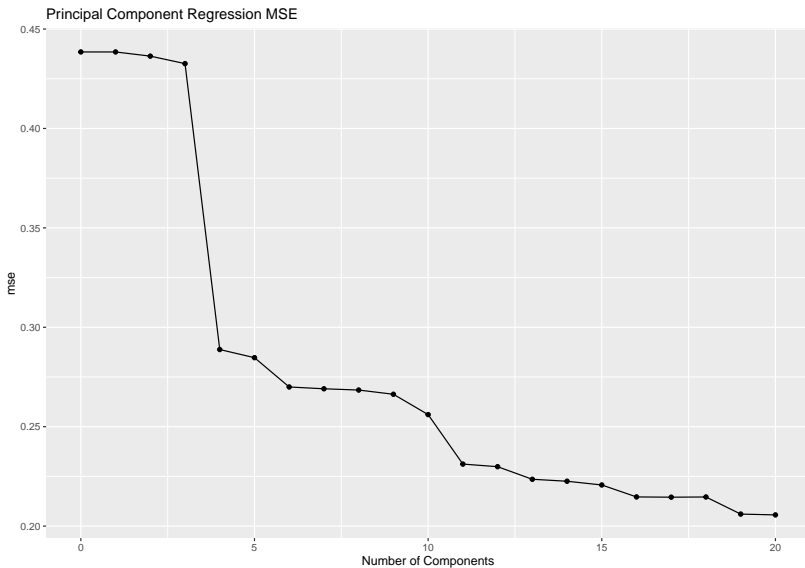
```
## [1] "MSE of Testing Set:  0.165"
```

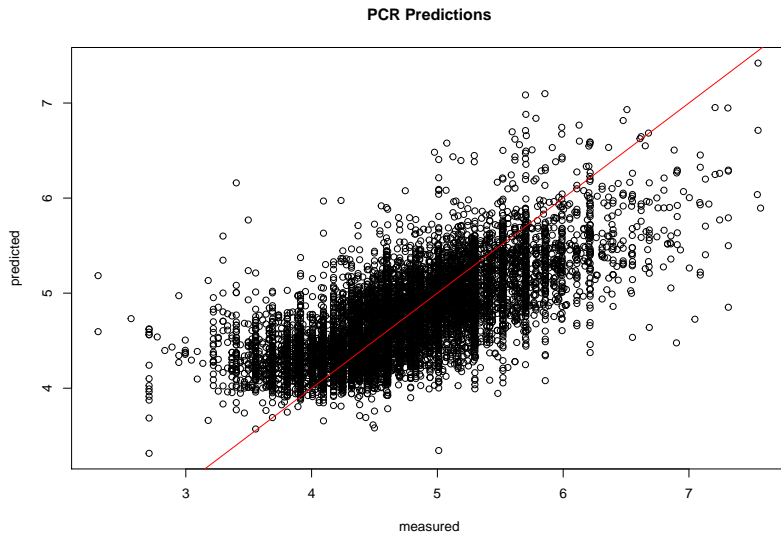# Regression Splines/Generalized Additive Models

# PCR and PLS

- ▶ 10 Fold Cross-Validation was performed for number of components ranging from 1 to 20.
- ▶ The Cross-Validation MSE was used to pick optimal number of components for both models.
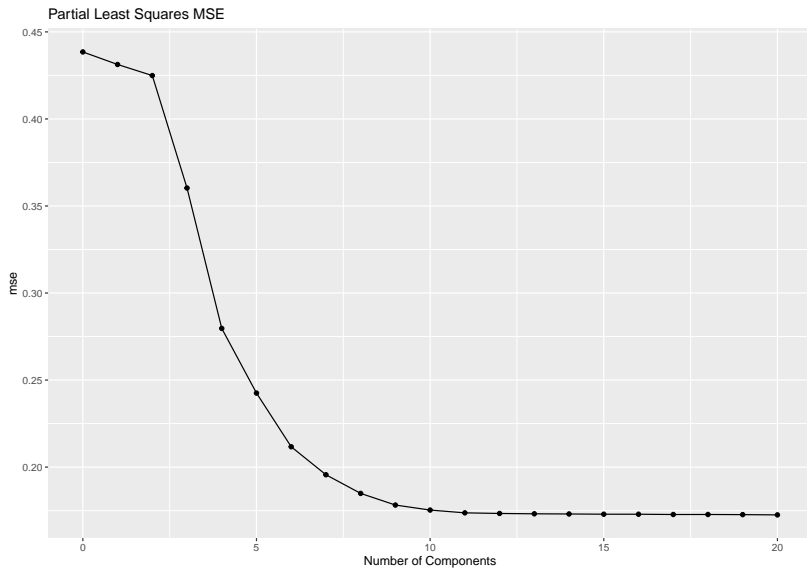
# PCR
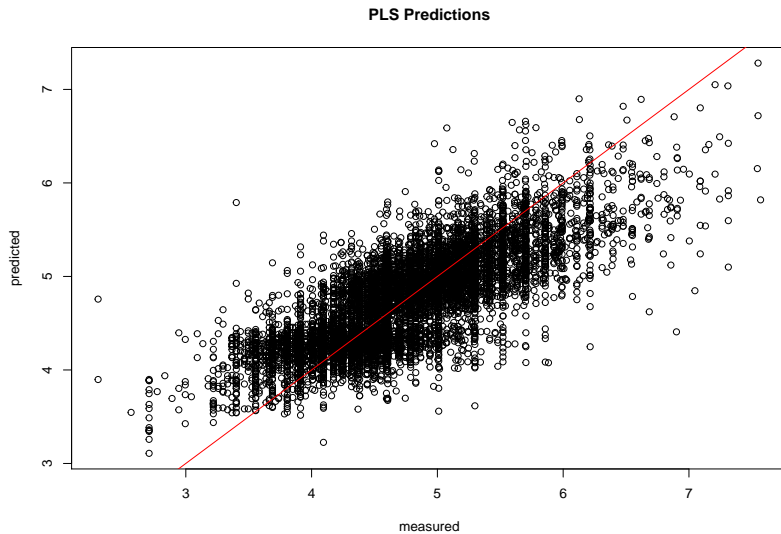


Principal Component Regression MSE

# PCR Predictions



PCR Predictions

# PLS



Partial Least Squares MSE

# PLS Predictions



PLS Predictions

# PCR and PLS Summary

```
##                         PCR     PLS
## Components          15.0000 10.0000
## Test MSE             0.1765  0.2192
## % Variance Explained 99.7000 99.9000
```

# Regression Trees

```
## 
## Regression tree:
## tree(formula = price ~ ., data = training)
## Variables actually used in tree construction:
## [1] "room_type" "longitude" "bathrooms" "city"       "bec
## Number of terminal nodes:  8
## Residual mean deviance:  0.1885 = 1695 / 8992
## Distribution of residuals:
##     Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -2.5050 -0.2999 -0.0196  0.0000  0.2558  2.8310

## [1] "Test MSE of Initial Tree:  0.1926"
```
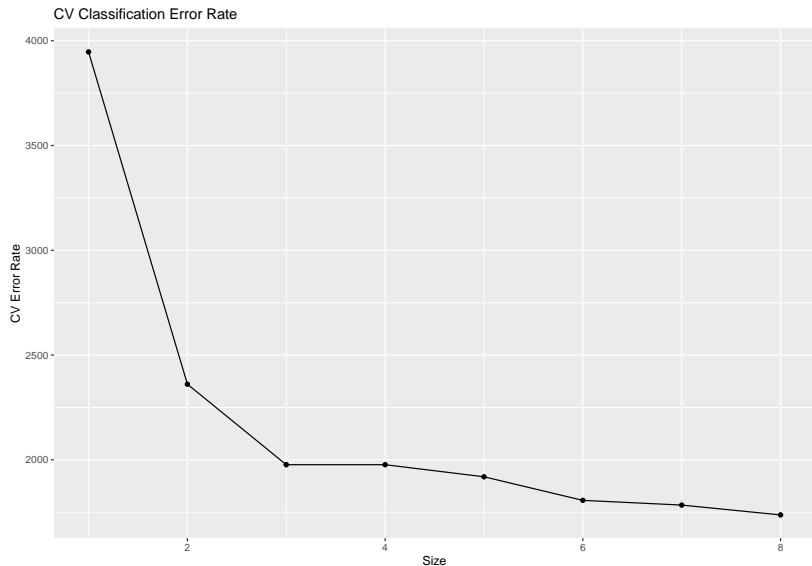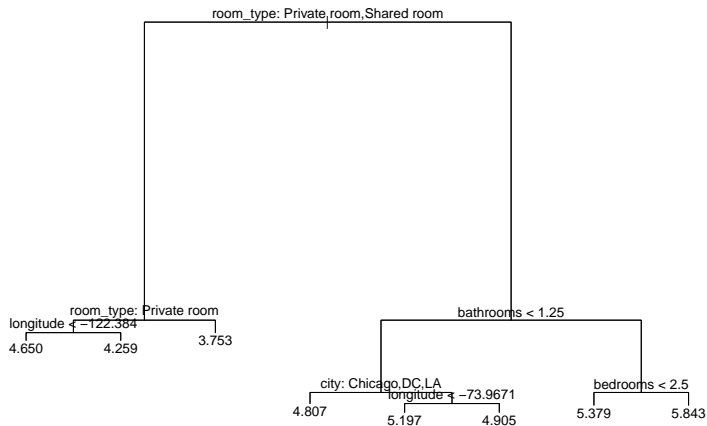
# Regression Trees

# Regression Trees



room_type: Private room,Shared room

room_type: Private room

longitude < −122.384

4.650     4.259     3.753

bathrooms < 1.25

city: Chicago,DC,LA

4.807

longitude < −73.9671

5.197     4.905

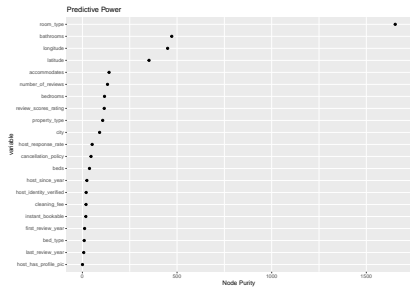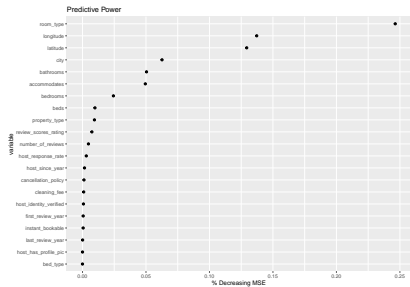bedrooms < 2.5

5.379     5.843

# Bagging

```r
bag_fit <- randomForest(price ~ ., data = training, mtry =
bag_predict = predict(bag_fit, testing, type = "response")
bag_MSE = round(mean((testing$price - bag_predict)^2), 4)
print(paste("Test MSE of Bagging: ", bag_MSE))
```

```
## [1] "Test MSE of Bagging:  0.1294"
```
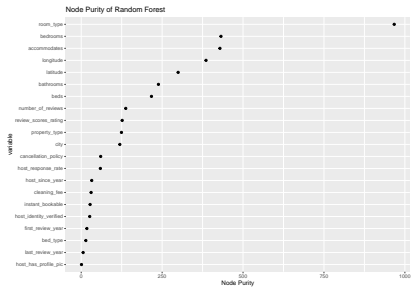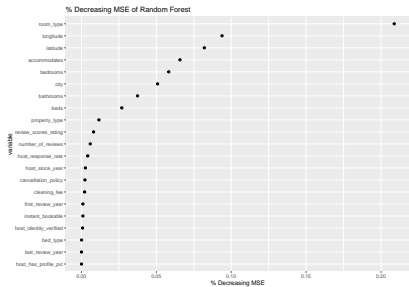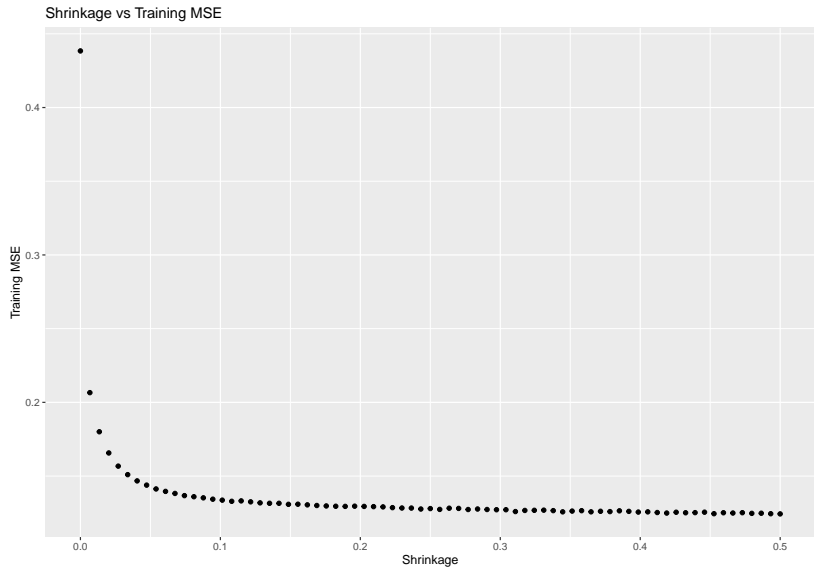
# Bagging

# Random Forests

```
rf_fit <- randomForest(price ~ ., data = training, mtry = s

## [1] "Test MSE of Random Forest:  0.1299"
```

# Random Forests
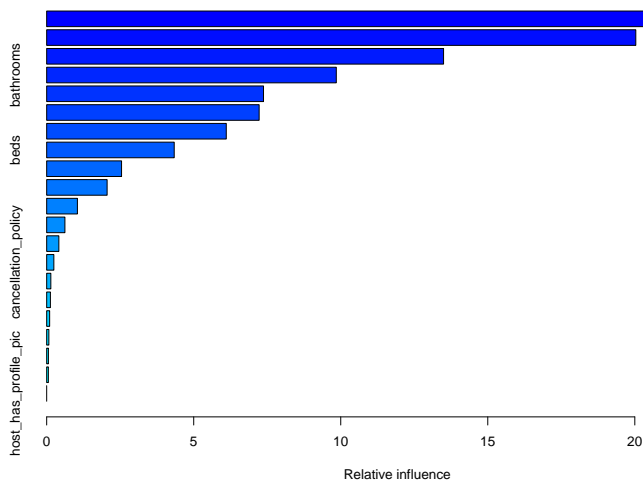
# Boosting



Shrinkage vs Training MSE

# Boosting



```
##                                                    var     rel.in
## property_type                             property_type 24.0900283
```

# MSE Table

# Final Model

# Going Forward

- Our data has data from multiple cities across the country
- Can we apply this to a certain city and see similar results?
- Is this accurate enough to help AirBnB hosts in selected cities?
  - Using current data, can this model help hosts correctly adjust their rates?

# Questions?

References