



UNIVERSITÀ  
DEGLI STUDI  
DI MILANO

---

LA STATALE

DEPARTMENT OF ECONOMICS, MANAGEMENT AND QUANTITATIVE METHODS

UNIVERSITÀ DEGLI STUDI DI MILANO

DATA SCIENCE AND ECONOMICS

## STATISTICAL LEARNING PROJECT

# ANALYSIS AND PREDICTION MODELS FOR NEW YORK CITY AIRBNB

REPORT

*Author:*

Andrea IERARDI

ACADEMIC YEAR 2019-2020

---

# CONTENTS

<b>List of Figures</b>	<b>4</b>
<b>1 Abstract</b>	<b>5</b>
<b>2 Problem Definition and Algorithm</b>	<b>6</b>
2.1 Two main Goals . . . . .	6
2.1.1 Develop predictive models for price . . . . .	6
2.1.2 Define clusters and groups . . . . .	6
2.2 Algorithms . . . . .	7
2.2.1 Linear Regression . . . . .	7
2.2.2 Decison Trees . . . . .	7
2.2.3 Random Forest . . . . .	7
2.2.4 Ranger Random Forest . . . . .	7
2.2.5 Neural Networks . . . . .	7
2.2.6 K-means . . . . .	7
2.2.7 Principal Component Analysis . . . . .	7
<b>3 Experimental Evaluation</b>	<b>8</b>
3.1 Methodology . . . . .	8
3.1.1 Data Inspection . . . . .	8
3.1.2 Data Cleaning and Pre-processing . . . . .	10
3.2 Results . . . . .	12
3.2.1 Linear Regression Results . . . . .	13
3.2.2 Decison Trees Results . . . . .	15
3.2.3 Random Forest Results . . . . .	16
3.2.4 Ranger Random Forest Results . . . . .	16
3.2.5 Neural Network Results . . . . .	17
3.2.6 K-means Results . . . . .	17
3.2.7 Principal Component Analysis Results . . . . .	17
3.3 Discussion . . . . .	18
3.3.1 Linear Regression Discussion . . . . .	18
3.3.2 Decision Tree Discussion . . . . .	18
3.3.3 Random Forest Discussion . . . . .	18

3.3.4	Ranger Random Forest Discussion . . . . .	18
3.3.5	Neural Network Discussion . . . . .	18
3.3.6	K-Means Discussion . . . . .	18
3.3.7	Principal Component Analysis Discussion . . . . .	19
<b>4</b>	<b>Conclusion</b>	<b>19</b>
4.1	Linear Regression . . . . .	19
4.2	Decision Tree . . . . .	19
4.3	Random Forest . . . . .	19
4.3.1	Ranger Random Forest . . . . .	19
<b>5</b>	<b>Appendix</b>	<b>20</b>

---

## LIST OF FIGURES

1	Price summary . . . . .	9
2	Distribution of all houses in NY colored by prices . . . . .	10
3	Distribution of all houses in NY of price between 15\$ and 500\$ per day . .	10
4	Approximated distribution map of all houses in Manhattan . . . . .	11
5	Summary of the MSE for all the models for each subset of the dataset . . .	12
6	Linear Regression output for the entire dataset . . . . .	13
7	Linear Regression output filtering by Manhattan . . . . .	14
8	Linear Regression output filtering by Manhattan and Entire home/Apartment . . . . .	15

# 1

---

## ABSTRACT

The aim of the project is to analyse, develop prediction models and define data clusters from the "New York City Airbnb Open Data" from a Kaggle competition.

In particular, one part is focused on the development of predictive models to forecast house prices using these Supervised Learning technics:

- Linear Regression
- Decision Tree
- Random Forest
- Ranger Random Forest
- Neural Newtworks

For each of these, a comparison between the Mean Square Error and between all  $R^2$  measure of all methods has been made to highlight which have the best performance. Also, for training all the models, a partition of the dataset in three parts has been applied: entire dataset, filtered by neigborhood group and filtered for neighborhood group and room type. In this way, it is possible to check the perfomance giving less or more features in input.

The second part is focused on the cluster and data reduction technics using these Unsupervised Learning technics:

- K-means Algorithm
- Clustering for mixed-type data
- Principal Component Analysis

# 2

---

## PROBLEM DEFINITION AND ALGORITHM

### 2.1 TWO MAIN GOALS

#### 2.1.1 DEVELOP PREDICTIVE MODELS FOR PRICE

The first objective is the forecast of the prices given some input information. This could be useful for a lot of scenarios. For example from a AirBnB customer point of view, he/she would like to get the list of houses more in alignment with his/her preference choice; or for a host point of view, where given the position and other information he/she could get information about the possible per day price of his property in New York City.

#### 2.1.2 DEFINE CLUSTERS AND GROUPS

The second objective is the definition of group or partition between the houses with different characteristics. For a user point of view could be useful to have information about available houses similar to those booked in the past. This could be also useful after the booking for a suggestion analysis having the information about last booked houses in New York or houses in similar cities around the world.

## 2.2 ALGORITHMS

### 2.2.1 LINEAR REGRESSION

Linear regression is a linear approach to modeling the relationship between a dependent variable and one or more independent variables.

### 2.2.2 DECISON TREES

In decision analysis, a decision tree can be used to visually and explicitly represent decisions and decision making.

### 2.2.3 RANDOM FOREST

Random forests are an ensemble learning method for classification, regression and other tasks that operate by constructing a multitude of decision trees.

### 2.2.4 RANGER RANDOM FOREST

Ranger is a fast implementation of random forests or recursive partitioning, particularly suited for high dimensional data.

### 2.2.5 NEURAL NETWORKS

Neural networks are a set of algorithms, modeled loosely after the human brain, that are designed to recognize patterns.

### 2.2.6 K-MEANS

K-means is a method of vector quantization, originally from signal processing, that aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

### 2.2.7 PRINCIPAL COMPONENT ANALYSIS

Principal Component Analysis (PCA) is mostly used as a tool in exploratory data analysis and for making predictive models.

# 3

---

## EXPERIMENTAL EVALUATION

### 3.1 METHODOLOGY

#### 3.1.1 DATA INSPECTION

The dataset is part of a Kaggle competition, called the New York City Airbnb Open Data. It contains 48.000 rows per 16 columns. The dataset is structured with these columns:

- **id**
- **name**: name of the listing
- **host\_id**
- **host\_name**
- **neighbourhood\_group**: location
- **neighbourhood**: area
- **latitude**: coordinates
- **longitude**: coordinates
- **room\_type**: space type
- **price**: in dollars
- **minimum\_nights**: amount of nights minimum
- **number\_of\_reviews**: number of reviews
- **last\_review**: latest review
- **reviews\_per\_month**: number of reviews per month
- **calculated\_host\_listings\_count**: amount of listing per host
- **availability\_365**: number of days when listing is available for booking

It has been select only 5 of these variables: price, latitude, longitude, neighbourhood\_group and room\_type. The reason is that it is reasonable to select them to predict the prices and to obtain different clusters. The position and the neighbourhood is important since if a property is positioned near the city center will have a higher price with respect to those situated in the outskirts; also the type of room, since an entire apartment will cost more than a single room.

From the Figure 1 is possible to see the distribution of the price. The minimum price is 0 and the maximum is 10000\$. It is not possible to rent a house for free, so it is possible to filter the price with a price higher than 15\$. It is possible that a luxury house cost a lot per day, but these value can not be consider in the model, instead they are outliers and for this reason it is convenient to filter the price again and take those that have a value lower than 500\$. The reason is that the third quantile has a price of 175\$ which is a far from 10000\$ (Figure 3). It has also been checked the null and missing value in the dataset and has not been found apart from the reviews\_per\_month column. It is not a problem, since this feature has been not taken in account to train the models.

```
price
Min. : 0.0
1st Qu.: 69.0
Median : 106.0
Mean : 152.7
3rd Qu.: 175.0
Max. : 10000.0
```

Figure 1: Price summary

From Figure 2, it is possible to see the distribution of all houses in New York City and the price. This picture is not really informative since it can be noticed that the prices for the most part are in the range 0-500\$ and only a low number of instances have a price greater than 500\$. Deleting the outliers, the Figure 3 is more informative than the one before.

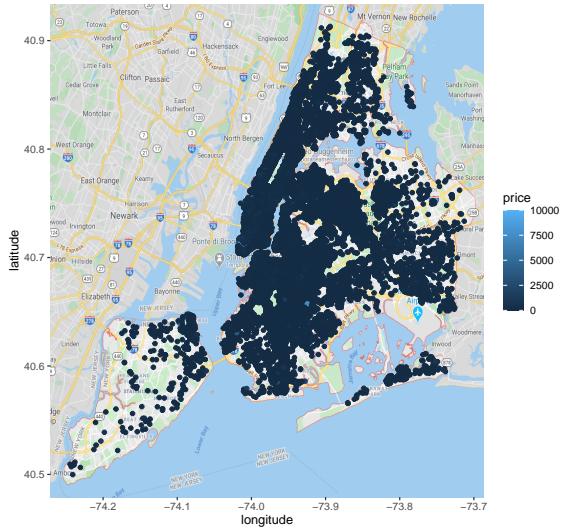


Figure 2: Distribution of all houses in NY colored by prices

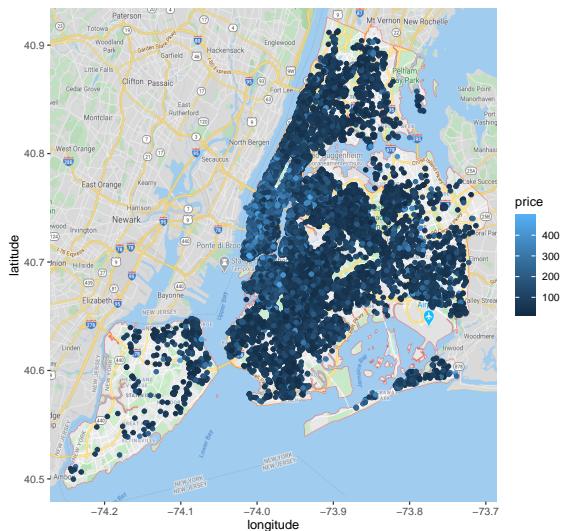


Figure 3: Distribution of all houses in NY of price between 15\$ and 500\$ per day

### 3.1.2 DATA CLEANING AND PRE-PROCESSING

The dataset has approximately 48.000 rows for each column, so an important part of the work is related to the pre-processing due to the large amount of data. For this reason, some variable has been rescaled to let the models to learn faster and better and to perform a better prediction. Latitude and longitude have not been rescaled for forecast price model, while for clustering methods to have a consistent distance calculation. Categorical variables have also been rescaled assigning a numerical value to each category, resulting as unordered factors.

The selected categorical variables are: neighbourhood\_group and room type.

The selected numerical variables are: latitude, longitude and price.

Considering that every neighbourhood group has his characteristic, also taking in account the room type that could impact largely the price. It is possible to define of different subset of the original set and try different methods of forecast. This is due to the fact that a customer should choose which of the different neighbourhood and room type is interested in. Models have been run to different scenarios (Figure 4): users interested in all New York City houses and all type of room, users interested only in a single neighbourhood and users interested in a single neighbourhood and a single room type.

Including different scenario could be potentially computationally expensive. In reality, the training proceeded without any problem.<sup>1</sup> Computational problem may emerge for the case of hyperparametrisation tuning, even using multiprocessing and multithreading technics. For semplicity, in this project no model tuning have been made.

Also for clustering have been filter every scenario, but not for the case of Hierarchical Clustering because it need to generate a readable dendrogram using the mean of all neighbourhood or room type price, Principal Component Analysis since avaible data are really small for some specific neighbourhood like Staten Island and specific room type.

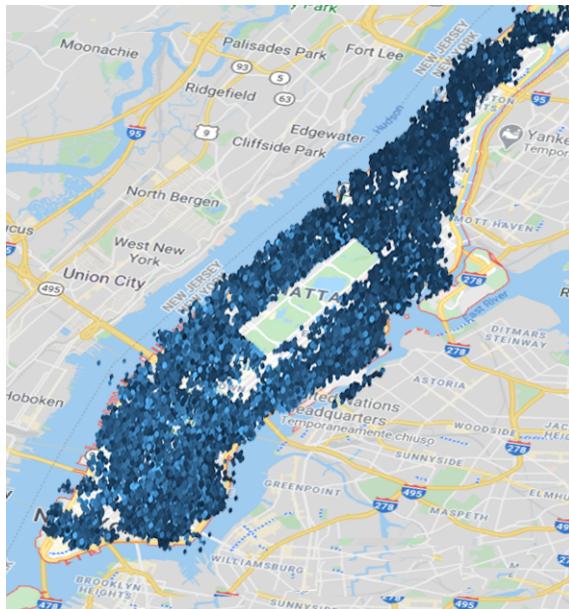


Figure 4: Approximated distribution map of all houses in Manhattan

---

<sup>1</sup>The models have been trained on a machine with quad-core 3.5 GHz processor, 8 Gb RAM and 4 Gb dedicated GPU.

## 3.2 RESULTS

Since the number of models trained are very high, the study result will be presented for different macro groups that depend on the different subset of the dataset: entire dataset, filtering by neighborhood and filtering by neighborhood and room type. This is due to the fact that neighborhood in general have similar result (like Manhattan has similar shape of Brooklyn and Staten Island has a similar shape of Bronx).

===== DA RIVEDERE ===

Subset	Linear.Regression	Decision.Tree	Random.Forest	Ranger.Random.Forest	Neural.Networks
<b>Brooklyn</b>	0.4351730	0.4395665	0.4175302	0.4141696	0.5295534
<b>Manhattan</b>	0.7936143	0.7964149	0.7483505	0.7474648	0.8801335
<b>Queens</b>	0.3721530	0.3699898	0.3514677	0.3525291	0.3807765
<b>Staten Island</b>	0.3467308	0.3471879	0.3415798	0.3492562	0.6710544
<b>Bronx</b>	0.2741482	0.2744778	0.2774524	0.2699379	0.3388045
<b>Brooklyn/Private room</b>	0.1858138	0.1806116	0.1936970	0.1925123	0.2464772
<b>Brooklyn/Entire home/apt</b>	0.7262980	0.7179003	0.7702742	0.7709408	0.7863187
<b>Brooklyn/Shared room</b>	0.1157458	0.1320524	0.0881610	0.0855227	0.1180578
<b>Manhattan/Private room</b>	0.4445611	0.3753113	0.3778408	0.3774914	0.5059889
<b>Manhattan/Entire home/apt</b>	0.9361941	0.9585893	0.9728082	0.9773220	1.0688116
<b>Manhattan/Shared room</b>	0.5616574	0.5452982	0.5875275	0.5895540	1.1310954
<b>Queens/Private room</b>	0.1616290	0.1551588	0.1516938	0.1525349	0.2684766
<b>Queens/Entire home/apt</b>	0.6846303	0.6591288	0.6550463	0.6523420	0.6980945
<b>Queens/Shared room</b>	0.0709897	0.2989787	0.2102031	0.2043754	0.0560464
<b>Staten Island/Private room</b>	0.1241196	0.1802672	0.1355774	0.1375752	0.1395840
<b>Staten Island/Entire home/apt</b>	0.5703345	0.7737648	0.7185084	0.7119659	1.2433252
<b>Staten Island/Shared room</b>	0.0518804	0.0960227	0.0921335	0.0865780	0.2093138
<b>Bronx/Private room</b>	0.2146868	0.2159434	0.2161088	0.2180725	0.5847212
<b>Bronx/Entire home/apt</b>	0.5627967	0.7622148	0.6682157	0.6587408	0.7178106
<b>Bronx/Shared room</b>	0.2071295	0.1939584	0.1759492	0.1750096	0.1756585
<b>All</b>	0.6036113	0.6005651	0.5672255	0.5436218	0.6334017

Figure 5: Summary of the MSE for all the models for each subset of the dataset

### 3.2.1 LINEAR REGRESSION RESULTS

**Linear Regression on the entire dataset** The results on the entire dataset are acceptable. All variables are significant for the model and the  $R^2$  value is 40% (Figure 7). The Mean Square Error (Figure 5) has a value of 0.6 which is acceptable. All the neighbourhood groups have a positive effect on the price apart for the 4<sup>th</sup> one. This could be due to the fact that Staten Island does not have expensive houses compared to the other group.

Also the latitude and longitude have a slight negative effect on the price. Entire apartment have a slight positive effect while the shared house negative. This could be due to the fact that an entire house will cost more than a shared room, so the price will increase for this type of room.

```
[1] "===== all ====="
call:
lm(formula = price ~ ., data = trains[[sub]])
Residuals:
    Min      1Q  Median      3Q     Max 
-2.1739 -0.4434 -0.1419  0.2269  4.9515 
Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -1.916e+02  1.420e+01 -13.497 < 2e-16 ***
neighbourhood_group2 5.023e-01  1.618e-02  31.041 < 2e-16 ***
neighbourhood_group3 2.216e-01  1.978e-02  11.202 < 2e-16 ***
neighbourhood_group4 -9.535e-01  5.882e-02 -16.211 < 2e-16 ***
neighbourhood_group5 2.687e-01  3.914e-02   6.865 6.81e-12 ***
latitude        -1.726e+00  1.393e-01 -12.385 < 2e-16 ***
longitude       -3.532e+00  1.595e-01 -22.145 < 2e-16 *** 
room_type2      9.996e-01  9.622e-03 103.887 < 2e-16 *** 
room_type3     -2.379e-01  3.068e-02  -7.754 9.18e-15 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
Residual standard error: 0.784 on 28680 degrees of freedom
Multiple R-squared:  0.3914, Adjusted R-squared:  0.3912 
F-statistic: 2306 on 8 and 28680 DF, p-value: < 2.2e-16
```

Figure 6: Linear Regression output for the entire dataset

### Linear Regression for specific Neighbourhood group

Models give different results for the specific neighbourhood group. For the first three (Brooklyn, Manhattan and Queens)  $R^2$  are over approximately 30%. For the MSE, the value have different ranges based on the distribution of prices in the singular neighbourhood. To be noticed, is the fact that MSE of Brooklyn is significantly lower with respect to Manhattan, given they have a similar number of houses and have almost the same price distribution. The only difference between the two is the fact that Manhattan has more entire apartment type of room. All variables have a significant effect on the price apart for the latitude in Manhattan and longitude for Queens with a 0.1.

For the last two groups (Staten Island and Bronx) all variables does not have significance in the response variable apart for the Entire Apartment dummy which have a positive effect. This is due to the fact that entire properties will cost more than shared rooms.

```

[1] "***** 1 *****"
Call:
lm(formula = price ~ ., data = trains[[sub]])

Residuals:
    Min     1Q   Median     3Q    Max 
-1.8239 -0.3502 -0.1259  0.1721  5.3480 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -555.16875  20.65023 -26.884 < 2e-16 ***
latitude      5.11535   0.20932  24.577 < 2e-16 *** 
longitude    -0.66795  -0.21202 -30.261 < 2e-16 *** 
room_type2    0.96345   0.01137  84.701 < 2e-16 *** 
room_type3   -0.17115   0.03990  -4.289  1.8e-05 *** 
...
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.6774 on 14889 degrees of freedom
Multiple R-squared:  0.3801, Adjusted R-squared:  0.38 
F-statistic: 2283 on 4 and 14889 DF, p-value: < 2.2e-16

[1] "***** 2 *****"
Call:
lm(formula = price ~ ., data = trains[[sub]])

Residuals:
    Min     1Q   Median     3Q    Max 
-2.3475 -0.5308 -0.1869  0.2803  4.7768 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -895.75437  58.11443 -15.414 < 2e-16 *** 
latitude     -0.10719   0.35247  -0.304   0.761  
longitude    -12.16468   0.61365 -19.823 < 2e-16 *** 
room_type2    1.02538   0.01494  68.621 < 2e-16 *** 
room_type3   -0.25643   0.04737  -5.413 6.29e-08 *** 
...
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.877 on 15653 degrees of freedom
Multiple R-squared:  0.3391, Adjusted R-squared:  0.339 
F-statistic: 2008 on 4 and 15653 DF, p-value: < 2.2e-16

[1] "***** 3 *****"
Call:
lm(formula = price ~ ., data = trains[[sub]])

Residuals:
    Min     1Q   Median     3Q    Max 
-1.3833 -0.3036 -0.1202  0.1360  4.9322 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept)  3.74775  12.41542   0.302   0.7628  
latitude    -0.77393   0.27981  -2.766  0.0057 **  
longitude   -0.36613   0.19934  -1.837  0.0663  
room_type2   0.81079   0.01957  41.426 < 2e-16 *** 
room_type3   -0.23506   0.05230  -4.494 7.17e-06 *** 
...
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.6063 on 4219 degrees of freedom
Multiple R-squared:  0.3065, Adjusted R-squared:  0.3058 
F-statistic: 466.1 on 4 and 4219 DF, p-value: < 2.2e-16

[1] "***** 4 *****"
Call:
lm(formula = price ~ ., data = trains[[sub]])

Residuals:
    Min     1Q   Median     3Q    Max 
-0.9669 -0.2906 -0.1359  0.1124  4.9772 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 62.09595  75.12511   0.827   0.409  
latitude    -0.46076   0.89539  -0.515   0.607  
longitude   0.59638   0.72879   0.818   0.413  
room_type2   0.67380   0.04732  14.238 < 2e-16 *** 
room_type3   -0.15156   0.09529  -1.591   0.112  
...
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.6266 on 806 degrees of freedom
Multiple R-squared:  0.2199, Adjusted R-squared:  0.216 
F-statistic: 56.79 on 4 and 806 DF, p-value: < 2.2e-16

[1] "***** 5 *****"
Call:
lm(formula = price ~ ., data = trains[[sub]])

Residuals:
    Min     1Q   Median     3Q    Max 
-0.9669 -0.2906 -0.1359  0.1124  4.9772 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 30.68887 144.852858  -0.212   0.832  
latitude    0.734021 1.531167   0.486   0.628  
longitude  -0.001178 1.331628  -0.001   0.999  
room_type2  0.739498 0.078599  9.408 < 2e-16 *** 
room_type3  0.159594 0.292213   0.546   0.585  
...
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 0.6387 on 269 degrees of freedom
Multiple R-squared:  0.2493, Adjusted R-squared:  0.2382 
F-statistic: 22.34 on 4 and 269 DF, p-value: 6.126e-16

```

Figure 7: Linear Regression output filtering by Manhattan

## Linear Regression selecting the Neighboorhood group and the type of room

```

[1] "===== n2-r1 ====="
Call:
lm(formula = price ~ ., data = trains[[sub]])

Residuals:
    Min      1Q  Median      3Q     Max 
-1.2007 -0.3766 -0.1682  0.1281  4.7297 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -673.3536   81.3847 -8.274 <2e-16 ***
latitude     -0.5113    0.4675 -1.094   0.274    
longitude    -9.3809   0.8670 -10.820 <2e-16 ***  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6845 on 5252 degrees of freedom
Multiple R-squared:  0.1039, Adjusted R-squared:  0.1036 
F-statistic: 304.5 on 2 and 5252 DF,  p-value: < 2.2e-16

```

Figure 8: Linear Regression output filtering by Manhattan and Entire home/Apartment

Linear Regression	MSE	R <sup>2</sup>
Entire dataset	0.59	.40
Filter: Mahnattan	0.76	0.355
Filter: Manhattan & Apt	0.44	0.12

### 3.2.2 DECISON TREES RESULTS

For the Decision tree the model give different results based on the variables used.

#### Decison tree without filters

#### Decison tree selecting the Neighboorhood group

#### Decison tree selecting the Neighboorhood group and room type

Decision Tree	MSE	R <sup>2</sup>
Entire dataset	0.59	..
Filter: Mahnattan	0.77	..
Filter: Manhattan & Apt	0.40	..

### 3.2.3 RANDOM FOREST RESULTS

There were no problem running Random Forest regression for the price, givin the default parameters. For the parameter tuning there were no possibilities for the entire dataset, due to the large number of data points. For the filtered dataset, instead, were possible to tune the mtry, number of maximum nodes and number of trees.

#### **Random Forest without filters**

There were no possibility to run a forecast of the price

#### **Random Forest selecting the Neighboorhood group**

**Random Forest selecting the Neighboorhood group and room type** Using the tuning of the parameter the results are slightly better. The model start with a 23% explained variance to a value of 25%.

Random Forest	MSE	%Var
Entire dataset	0.56	43.9
Filter: Mahnattan	0.70	38.9
Filter: Manhattan & Apt	0.38	21.6

### 3.2.4 RANGER RANDOM FOREST RESULTS

Ranger Random Forest is known to be computationally light with respect to the classic Random Forest. In fact, for the tuning part there were not problem in running it for the entire dataset.

#### **Ranger without filters**

Ranger outputs for the entire dataset are consistent. We have a  $R^2$  of 0.47 and OOB error of.

#### **Ranger selecting the Neighboorhood group**

The dataset filtered by Neighboorhood ouputs a value of 0.4  $R^2$  and a OOB error of.

**Ranger selecting the Neighboorhood group and room type** The dataset filtered by neighborhood and room type gives as result a  $R^2$  of 0.25.

<b>Ranger RF</b>	MSE	$\%Var$
<b>Entire dataset</b>	0.53	46.6
<b>Filter: Mahnattan</b>	0.71	38.8
<b>Filter: Manhattan &amp; Apt</b>	0.38	21.7

### 3.2.5 NEURAL NETWORK RESULTS

<b>Neural Networks</b>	MSE	Loss
<b>Entire dataset</b>	0.61	0.53
<b>Filter: Mahnattan</b>	0.84	0.64
<b>Filter: Manhattan &amp; Apt</b>	0.48	0.44

### 3.2.6 K-MEANS RESULTS

### 3.2.7 PRINCIPAL COMPONENT ANALYSIS RESULTS

### 3.3 DISCUSSION

#### 3.3.1 LINEAR REGRESSION DISCUSSION

Linear regression model gives interesting results for the non-filtered dataset and also for the filtered by neighbourhood group. All variable results rejected by Null hypothesis, so the model depends on all the selected variables. Latitude and longitude are correlated with the target, room type is strongly positive correlated with the price and neighbourhood group does not seem to have a great contribution in the prediction of the price.

#### 3.3.2 DECISION TREE DISCUSSION

The performance with respect to the other models are not the best, but acceptable. The prediction results are not also very precised for the filtered neighbourhood and room type. Also, the plots of the predicted value are not so consistent since the values are divided in category which correspond to the leaves that are not so strong with respect to the other model predictions.

#### 3.3.3 RANDOM FOREST DISCUSSION

Random forest outputs consistent results and performs better than linear regression and decision tree. Parameters tuning does not give big improvement in performance and also are computationally expensive.

#### 3.3.4 RANGER RANDOM FOREST DISCUSSION

Results of Ranger are the best with respect to the previous models. Also the tuning part was fast and computationally cheaper than the classic Random Forest model but does not give great improvements in performance.

#### 3.3.5 NEURAL NETWORK DISCUSSION

#### 3.3.6 K-MEANS DISCUSSION

# 4

---

## CONCLUSION

From the result the method with the most higher accuracy is the Random Forest method... while the worst are ....

Moreover, Random Forest method is also the worst in term of computation time for the tuning part since it takes for a configuration with 4 core, more or less 1 hour to tune the parameters.

### 4.1 LINEAR REGRESSION

### 4.2 DECISION TREE

Decision tree are one of the most used model in the Machine Learning world since are very familiar to human users and can be easily plotted.

### 4.3 RANDOM FOREST

Random Forest is an ensemble method which use a combination of decision tree to get the prediction.

#### 4.3.1 RANGER RANDOM FOREST

Ranger Random Forest is a computationally light model which results are very close the classical Random Forest.

# 5

---

## APPENDIX

The code is in the file "project-code.Rmd" and is attached to this file.