

Analysis Of Airbnb Listings

Sai Venkata Manoj Vungarala | Ishal Abhishek Mummidivarapu | Sujith Naarayan Hirendra Babu

Summary

The goal of the project is to analyze the Airbnb listings in six major cities (New York, Los Angeles, San Francisco, Boston, Chicago, and DC) in the United States. Airbnb Inc. is an American company that operates an online marketplace for lodging, primarily homestays for vacation rentals, and tourism activities. It lists more than 6 million rooms and houses in 81,000 cities. The platform has more than 150 million users worldwide with 2 million people staying in Airbnb rentals across the world on any given night. Given the extremely large quantity of listings, this project's dataset will contain listing in only six US cities (New York, Los Angeles, San Francisco, Boston, Chicago, and DC).

The data contains 74111 rows and 29 columns. Some important columns are price, property_type, room_type, bedrooms, bathrooms, city, amenities, etc.. Apart from analysis, this project will also include modeling and will list out the important factors while listing a property on the website. Additionally, the project also provides the basic common amenities required while listing a property by utilizing the concept of tokenization.

A brief description of the project's methods would include:

1. Data Preprocessing: Tidying the data
2. Data Visualization : Visualizing various graphs and making conclusions based on that.
3. Tokenization : Helps in identifying the top 10 common amenities in a listing.
4. Modeling : Applied Linear Regression, Random Forest and SVM which will list out the best response variable for price predictor.

In the results sections, the project will conclude the important factors like amenities, type of the listing, etc that are to be considered while listing an Airbnb property.

Methods

Data Preprocessing:

The dataset was taken from Kaggle ([link](#)) with the title "US Airbnb Listings". Our Airbnb dataset consists of 29 columns that describe the type of the property, price, amenities, number of rooms, number of bathrooms, etc.

The first thing that was performed was to add the state column so that it would be easier to plot maps. Next, the log_price column has been converted to regular price so that would be easier to understand. Finally, tokenization has been performed to amenities column to identify the top ten common amenities provided by a listing.

Data Visualization:

Exploratory Data Analysis (EDA) helps us in analyzing the data sets to visually summarize their characteristics. The various data from the dataset are visualized using various graphs to provide greater insights about the dataset. Below are the visualizations that this project consists of:

- Number of Airbnb's in Different Cities
- Property Type and its Count
- Count of Room Types
- Average price of Property Type
- Top 10 Amenities
- No. of Airbnb's with Cancellation policy in the six Cities
- Instant Booking option
- Average Ratings in different cities
- Number of Airbnb's Vs type of accommodations in various cities
- Geological location of the listings and its Price Range

Tokenization:

The amenities column consists of a list of amenities that a particular listing offers. Tokenization (process of splitting text into tokens) is used to read and identify that top 10 common amenities that most of the listings offered.

Modeling:

1 Linear Regression:

Linear Regression is a machine learning algorithm based on supervised learning. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables they are considering, and the number of independent variables being used.

Linear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a linear relationship between x (input) and y (output). Hence, the name is Linear Regression.

While training the model we are given: x : input training data (univariate – one input variable(parameter)) y : labels to data (supervised learning) When training the model – it fits the best line to predict the value of y for a given value of x . The model gets the best regression fit line by finding the best θ_1 and θ_2 values. θ_1 : intercept θ_2 : coefficient of x Once we find the best θ_1 and θ_2 values, we get the best fit line. So, when we are finally using our model for prediction, it will predict the value of y for the input value of x . How to update θ_1 and θ_2 values to get the best fit line? Cost Function (J): By achieving the best-fit regression line, the model aims to predict y value such that the error difference between predicted value and true value is minimum. So, it is very important to update the θ_1 and θ_2 values, to reach the best value that minimizes the error between predicted y value (pred) and true y value (y).

2 Support Vector Machines (SVM):

A support vector machine (SVM) is a linear classifier that finds a hyperplane in an n -dimensional space, where n is the total number of features. There are many different planes to draw to separate the features, but the goal of SVM is to identify one that has the greatest distance between the vectors of both classes. These hyperplanes serve as judgment boundaries, allowing the data to be classified more precisely. Kernels are used to transform data, and we picked the linear kernel because our data is linearly separable into two classes. The data was vectorized with the help of the TF-IDF algorithm.

3 Random Forest:

The RF classifier is an ensemble method for bagging that involves training several decision trees in parallel with bootstrapping and aggregation. Using different subsets of available features, several individual decision trees are trained in parallel on different subsets of the training dataset. Bootstrapping ensures that each decision tree in the random forest is unique, lowering the RF classifier's overall variance. RF classifier aggregates individual tree decisions for the final decision; as a result, RF classifier has good generalization. In terms of accuracy, the RF classifier tends to outperform most other classification methods without the risk of overfitting. RF classifiers, like DT classifiers, do not require feature scaling. RF classifier, unlike DT classifier, is more resistant to training sample selection and noise in the training dataset. In comparison to the DT classifier, the RF classifier is more difficult to interpret but easier to tune the hyperparameter.

Results

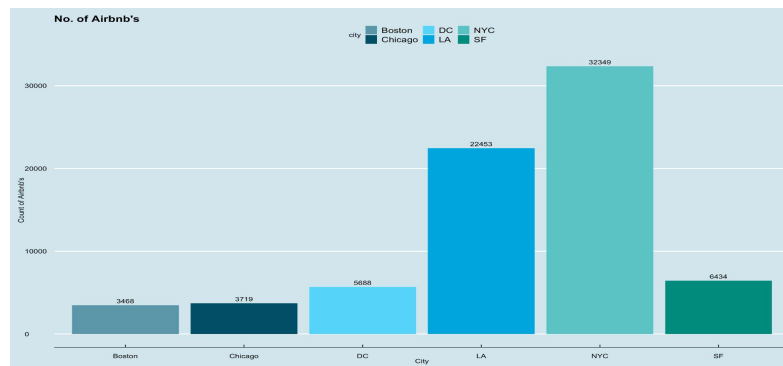


Figure 1. Number of Airbnb's in Different Cities

Figure 1. Shows the graph plot between Count of Airbnb's in the six Cities. Based on the Airbnb listings of major US cities it is observed that New York City (32349) being the financial capital has the maximum number of listings followed by Los Angeles (22453) and the least being Boston (3468). From the graph, we can also infer that the city having many Airbnb's is densely populated.

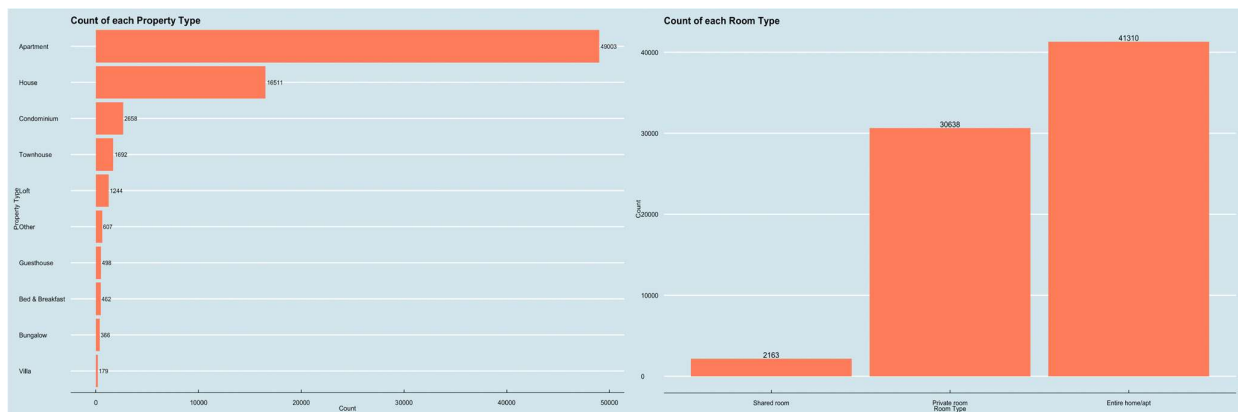


Figure 3a & 3b. Property Type and its Count

From the graph plotted in Figure 3a between Property type vs Count, the count of Apartment as stay is predominant followed by normal Houses with Villa being the least of all.

While analyzing for each Airbnb listing based on the general category, we were able to obtain the following plot (Figure 3b) where we could say that people prefer the entire home or apartment compared to the shared rooms.

The average price of different properties is listed in Figure 4. The property type Hostel is the one that is the cheapest and the property type Timeshare is the costliest.

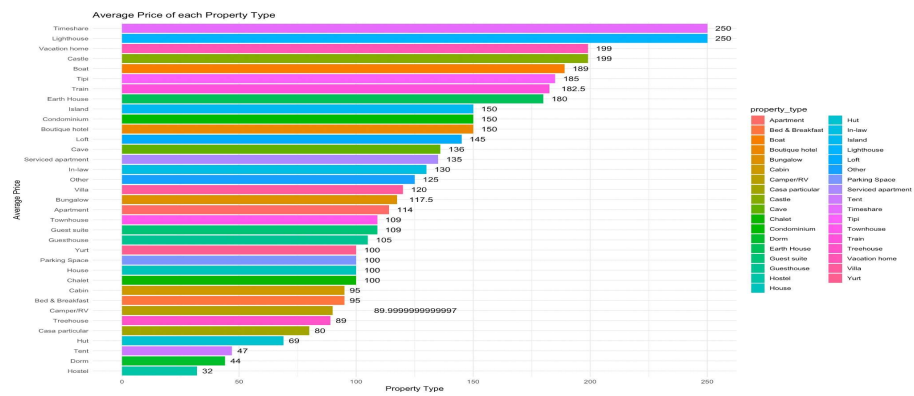


Figure 4. Average price of Property Type

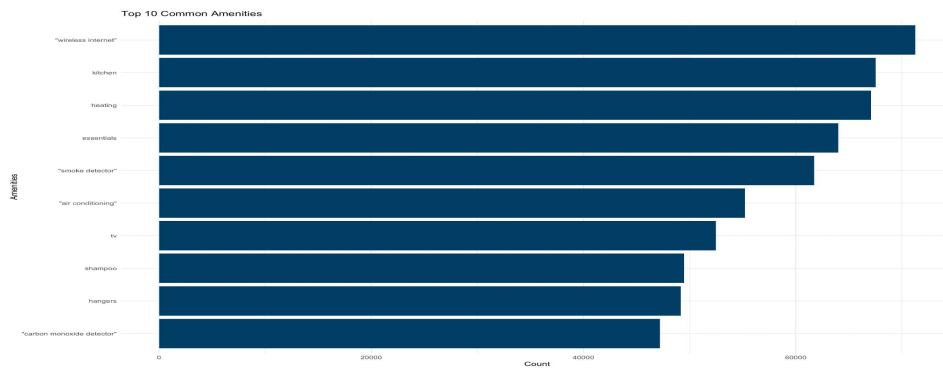


Figure 5. Top 10 Amenities

Wireless Internet (Wifi) is the common amenities across all the properties followed by Kitchen and Heaters. Almost 44000+ properties have all the above-mentioned amenities. (Figure 5)

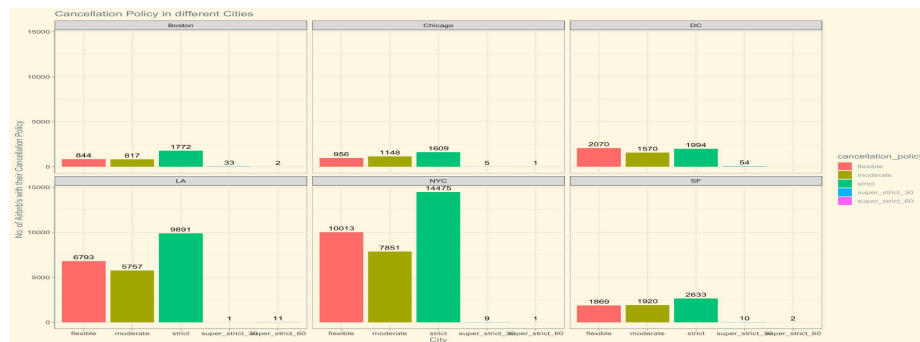


Figure 6. No. of Airbnb's with Cancellation policy in six Cities

The group of Bar graphs show the different cancellation policies followed in different properties in different cities. In Boston, cancellation policy seems not flexible, that is they have low flexibility and NYC Airbnb's having more properties with good flexibility for cancellation policy. The rest of the city listings can be observed above. (Figure 6)

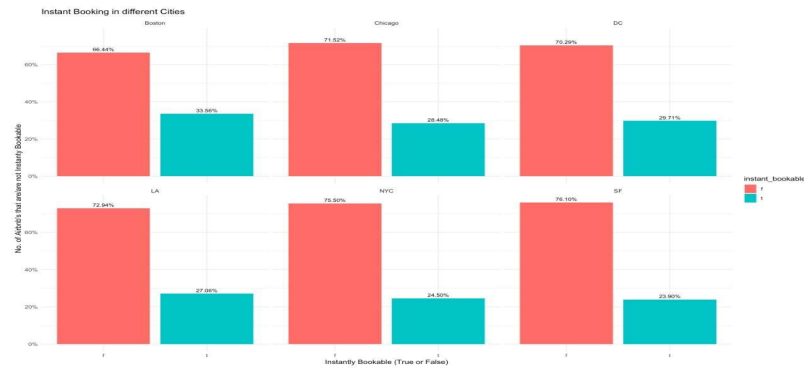


Figure 7. Instant Booking Option

Cities	Instant booking	Not Instant Booking
Boston	2304	1164
Chicago	2660	1059
DC	3998	1690
LA	16377	6076
NYC	24425	7924
San Francisco	4896	1638

As observed from the previous graphs, densely populated cities such as New York and Los Angeles have the premium option of instant booking and the ratio is equally managed between instant and non-instant booking subject to their high volume of listings. This can be clearly understood by observing Figure 7.

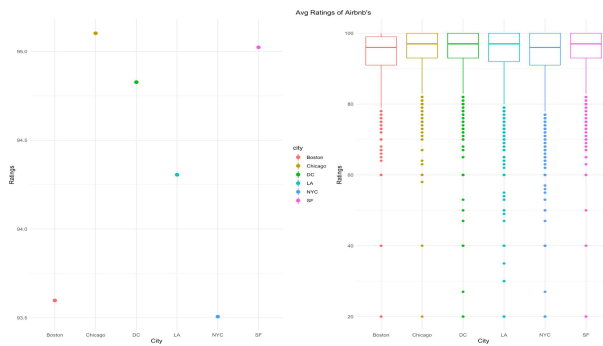


Figure 8. Average Ratings in different cities

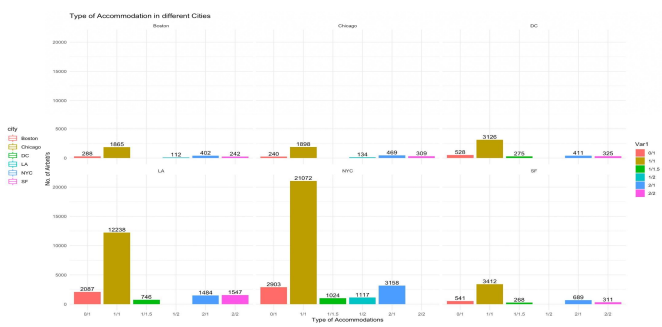


Figure 9. Number of Airbnb's Vs type of accommodations

Chicago and San Francisco properties have the most positive reviews of above 95 ratings, whereas New York city is the least rated one with 93.5. For much more clarity we have also plotted the box plot for each individual city. (Figure 8)

Among all the cities, 1 bedroom with 1 bathroom is the most common type of accommodation/listing. (Figure 9)

The population density of various kinds of properties across the Map have been mapped and placed in Figure 10a and Figure 10b. It gives a brief idea about the various properties in the cities. Finding an exact price is not extremely important when surveying an Airbnb listing. A price range can be just as telling. So, we discretized the log price into six different categories that is from 0 to 1500 and plotted the density map for each city for which we got the following result.

The listing costs are largely in line with the location scores. Highly rated locations also tend to be the most expensive ones since the highly rated location would also tend to be costly.

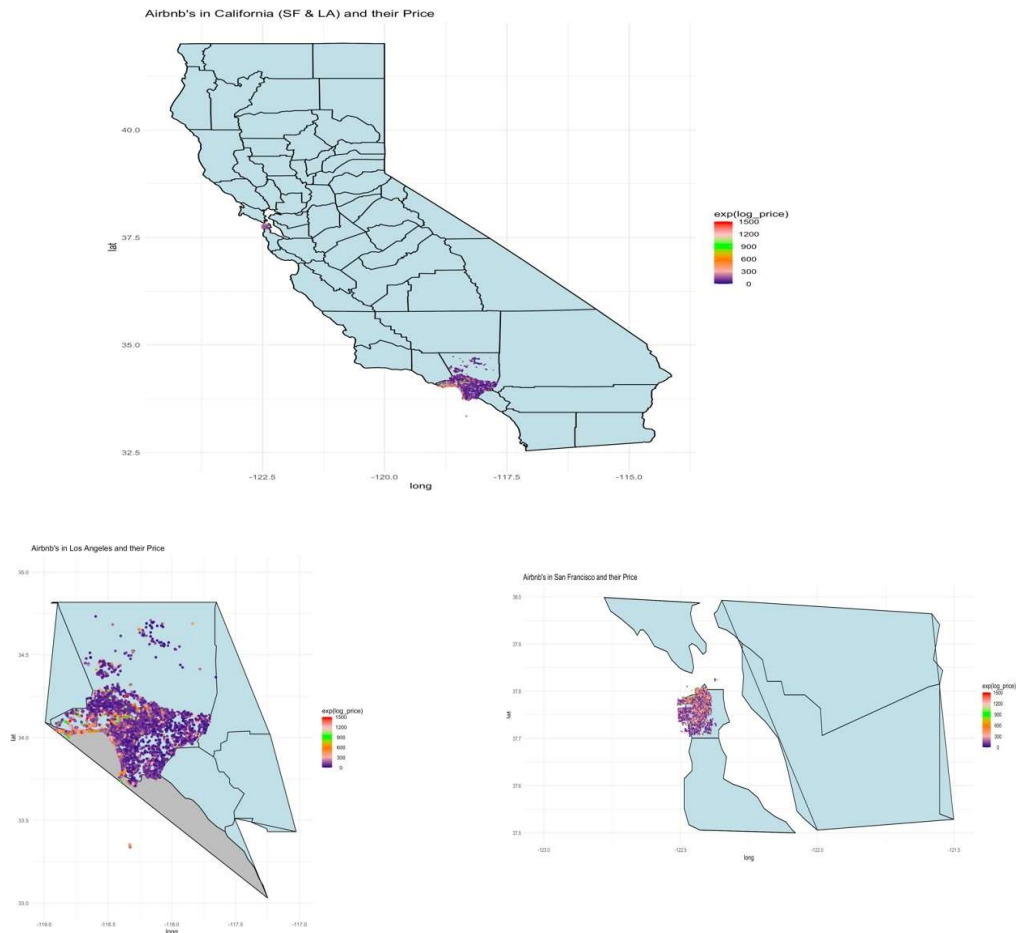


Figure 10. California (San Francisco and Los Angeles)

Above 3 are the graphs for California state along with its cities San Francisco and Los Angeles. We could observe that the price of a listing is expensive if it is closer to the bay area.

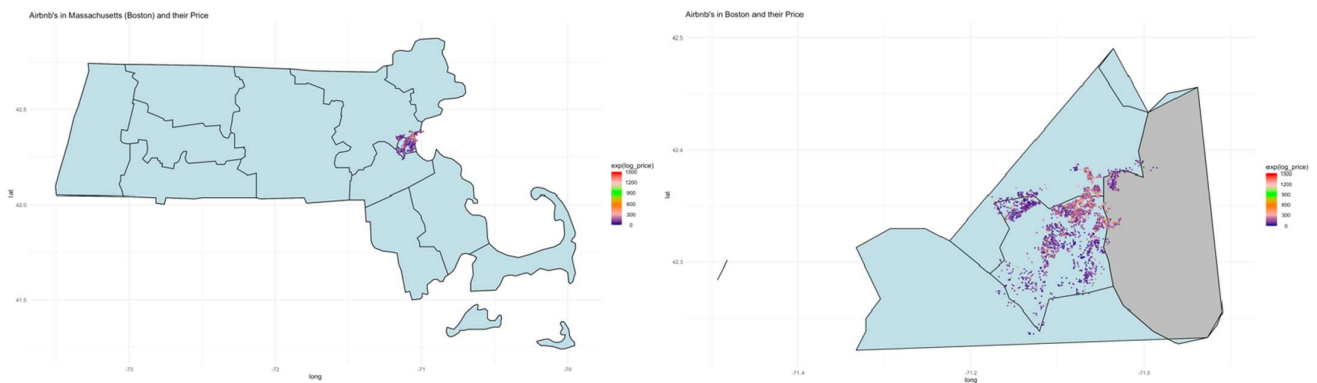


Figure 11. Massachusetts (Boston)

Figure 11 displays the locations of the listings in Boston, Massachusetts. The same observations can be made here that closer the location to the ocean, Charles River and Boston downtown, higher is its price.

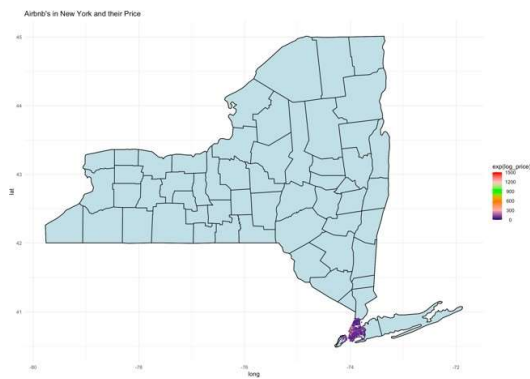


Figure 12. New York

In New York Properties along the Hudson River are the expensive when compared to other listings.

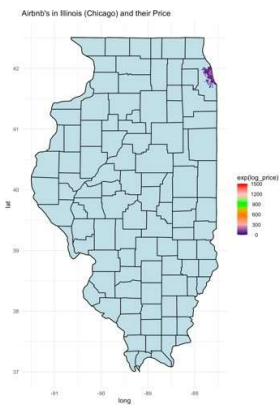


Figure 13. Illinois (Chicago)

From Figure 13 and 14 it would be clear that listings price changes with the change in its location. Prime located listings tend to have higher price when compared to others.

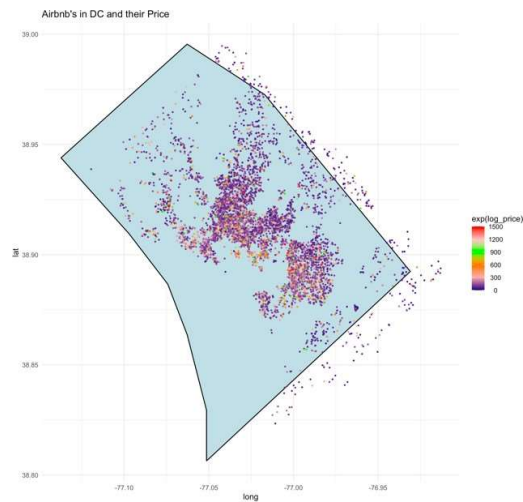


Figure 14. DC

Model	Train RMSE	Test RMSE
Linear Regression	0.43	0.44
Random Forest	0.40	0.42
SVM	0.41	0.42

Above are the results for the 3 models that have been implemented. Based on the results, the important factors that are to be considered while booking a property are:

- Room type
- Ratings
- Bedrooms
- Bathrooms

The most listed type of listings is "Apartments" and 1B1B type are more.

Properties Near Oceans, Rivers and Downtowns have high prices.

Also, the 10 common amenities in a listing are:

- Wireless Internet
- Kitchen
- Heating
- Essentials
- Smoke Detectors
- Air Conditioning
- TV
- Shampoo
- Hangers
- Carbon Monoxide Detectors

Discussions

From the results one can understand what people are looking in a listing before booking it. They would be considering room type, ratings, bedrooms, accommodations, and bathrooms. Moreover, from the analysis it is evident that price of the property may increase based on the location of it within the city. For example, in Los Angeles if a property is close to the bay area, the cost would be more when compared to a listing which is centrally located. This makes sense as most people would prefer the beach view and so is the demand for it.

This project will help people who want to post their listing on Airbnb. They will know what the customers are looking in a property, the amenities to be provided, importance of location and pricing.

Further, this project can be improved by gathering even more data across different cities and adding few other parameters like property area, its connectivity, etc. This would help in better analyzing the price.

Statement of Contributions

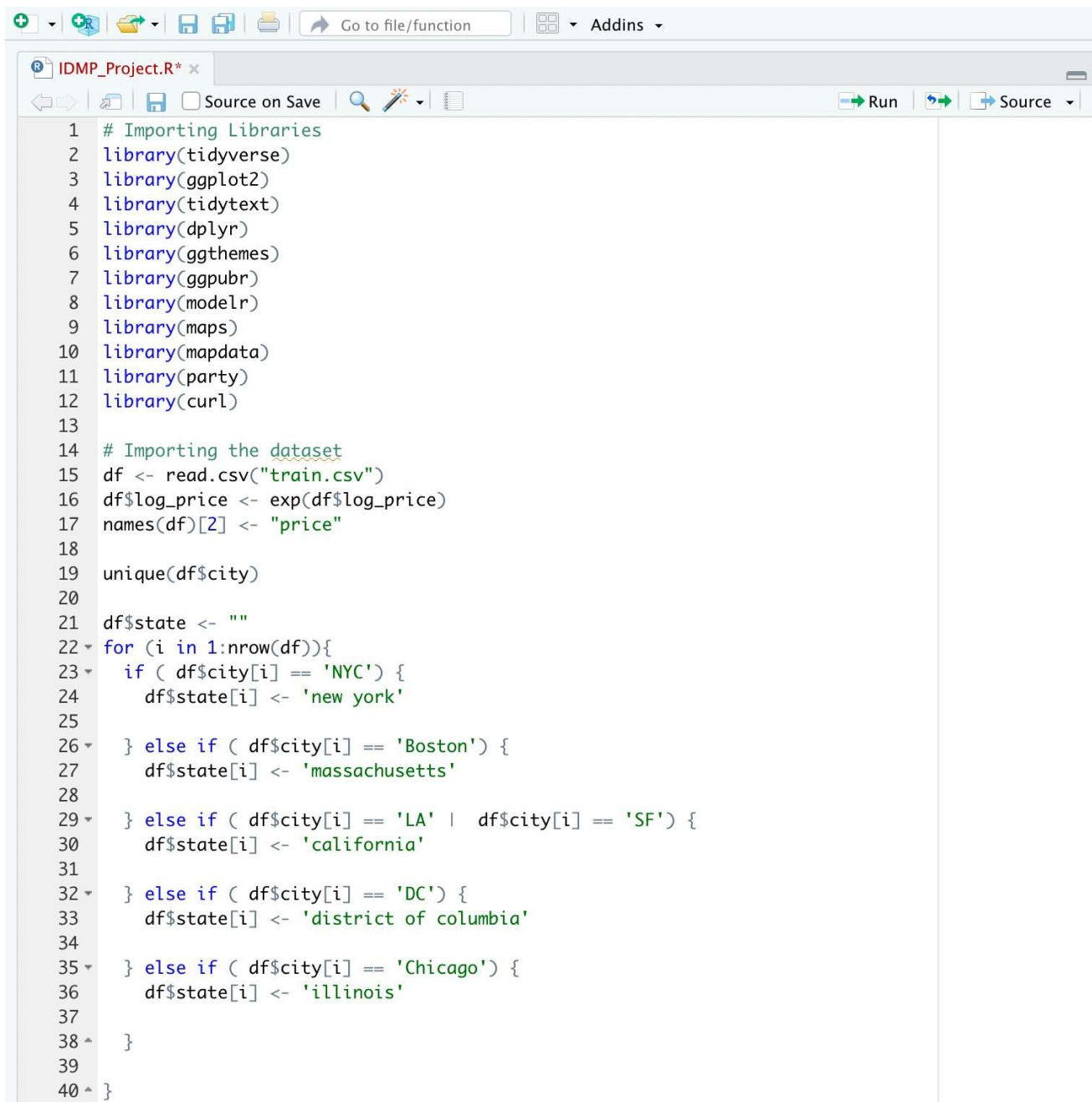
- Sai Venkata Manoj Vungarala: Data Collection, Data Transformation, Tokenization, Data Visualization, Exploratory Data Analysis, Step wise Modeling and Random Forest Regressor Model.
- Ishal Abhishek Mummidivarapu: General Trends analysis, Data Visualization, Linear Regression and SVM Modeling.
- Sujith Naarayan Hirendra Babu: Preliminary Analysis, Data Collection, Cleaning and Preprocessing, Data Transformation, Data Visualization, Project Content and Information Collection and Preparation of Presentation.

For the report, everyone worked on the portion that they did the analysis and generated visualizations.

References

- [1] Data from Kaggle (<https://www.kaggle.com/rudymizrahi/airbnb-listings-in-major-us-cities-deloitte-ml>)
- [2] Text-Mining (<https://www.analyticsvidhya.com/blog/2020/05/what-is-tokenization-nlp/>)
- [3] Modeling (<https://www.sciencedirect.com/topics/engineering/random-forest>)

Appendix



The image shows a screenshot of the RStudio IDE interface. The top toolbar includes icons for file operations (new, open, save, print) and a search bar labeled 'Go to file/function'. Below the toolbar, the file explorer shows a project named 'IDMP_Project.R*'. The main editor pane displays R code with line numbers from 1 to 40. The code imports various libraries, reads a CSV file, and processes the data by adding state information based on city names.

```
1 # Importing Libraries
2 library(tidyverse)
3 library(ggplot2)
4 library(tidytext)
5 library(dplyr)
6 library(ggthemes)
7 library(ggpubr)
8 library(modelr)
9 library(maps)
10 library(mapdata)
11 library(party)
12 library(curl)
13
14 # Importing the dataset
15 df <- read.csv("train.csv")
16 df$log_price <- exp(df$log_price)
17 names(df)[2] <- "price"
18
19 unique(df$city)
20
21 df$state <- ""
22 for (i in 1:nrow(df)){
23   if ( df$city[i] == 'NYC') {
24     df$state[i] <- 'new york'
25   } else if ( df$city[i] == 'Boston') {
26     df$state[i] <- 'massachusetts'
27   } else if ( df$city[i] == 'LA' | df$city[i] == 'SF') {
28     df$state[i] <- 'california'
29   } else if ( df$city[i] == 'DC') {
30     df$state[i] <- 'district of columbia'
31   } else if ( df$city[i] == 'Chicago') {
32     df$state[i] <- 'illinois'
33   }
34 }
35
36
37
38
39
40 }
```

Of IDMP Project.R

Source on Save

@

4 Run

^4

Source

```

42 # The data
43 head(df)
44
45 # Count of each Property Type |
46 df %>% count(property_type, sort=TRUE) %>%
47   top_n(10) %>%
48   ggplot(aes(x=reorder(property_type, n), y=n, label=n)) +
49     geom_col(fill="#F18464") +
50     coord_flip(j+
51     geom_text(size=3.5, hjust=-0.1) +
52     labs(x='Property Type', y='Count', title="Count of each Property Type") *
53     theme_economist
54     scale_fill_economist()
55
56 # Plotting the count of different Room Types
57
58 df %>% count(room_type, sort=TRUE) %>%
59   top_n(10) %>%
60   ggplot(aes(x=reorder(room_type, n), y=n, label=n)) +
61     geom_col(fill="#F18464") +
62     geom_text(size=4.5, vjust=-0.2) +
63     labs(x='Room Type', y='Count', title="Count of each Room Type") +
64     theme_economist
65     scale_fill_economist()
66
67 # Plotting the average price of each property type
68 df %>%
69   group_by(property_type) %>%
70   summarise(MedPrice=median(price, na.rm=TRUE)) %>%
71   ggplot(aes(x=reorder(property_type, MedPrice), y=MedPrice, label=MedPrice, fill=property_type)) +
72     geom_col() +
73     coord_flip(j+
74     geom_text(size=4.0, hjust=-0.5) +
75     labs(x="Average Price", y="Property Type", title="Average Price of each Property Type") +
76     theme_minimal()
77
78
79 # Plotting the most common Amenities
80 df$amenities <- substr(df$amenities, start=2, stop=nchar(df$amenities)-1)
81 df_tz <- unnest_to_kens df, Amenities, amenities, to_kens = 'regex', pattern=","
82
83 df_tz %>%
84   count(Amenities, sort=TRUE) %>%
85   top_n(10) %>%
86   ggplot(aes(x=reorder(Amenities, n), y=n)) +
87     geom_col(fill="#023E62")

```

IDMP Project.R* x



```

| | @ Q Source on Save | @ - • | Q I-6 Run | I'6 | Source
87 coord_ft | pC$ +
88 t absCx="Amenities " , y="Count" ,
89 title="Top 10 Common Amenities"} +
90 theme minimol(J
91
92 # Cancellation Policy in the six Cities
93 ggplot(df, aes(x=cancellation_policy, fill=cancellation_policy)) +
94   geom_bar(J+
95   facet_wrap(~cityJ+
96   geom_text(stat='count', aes(label=..count..), vjust=-0.5)+

97   tabsCx = 'City' , y = "No. of Airbnb's with the Cancellation Policy" ,
98   title = "Cancellation Policy in different Cities") +
99   theme_solarizedC}

100
101 # Instant Booking in the six Cities
102 df %>% group_by(city) %>%
103   count(instant_bookable = factor(instant_bookable)) %>%
104   mutate(pct = prop.table(n)) %>%
105   ggplot(aes(x=instant_bookable, y=pct, fill=instant_bookable, label=scales::percent(pct))) +
106   geom_bar(position='dodge', stat='identity') +
107   facet_wrap(~cityJ+
108   geom_text(position=position_dodge(width=.9),
109             vjust=-0.5,
110             size=3) +
111   scale_y_continuous(labels=scales::percent)+
112   tabsCx = 'Instantly Bookable (True or False)' , y = "No. of Airbnb's that are or are not Instantly B
113   title = "Instant Booking in different Cities" } +
114   theme_minimal{}
115
116
117 # No. of Airbnb's
118 ggplot(df, aes(x=city, fill=city)) +
119   geom_bar(J+
120   geom_text(stat='count', aes(label=..count..), vjust=-0.5)+
121   labs(x="City", y="Count of Airbnb's", title="No. of Airbnb's")+
122   theme_economist(J +
123   scale_fill_economist()

125 # Avg Ratings of Airbnb's
126
127 plot1 <- ggplot(df, aes(x=city, y=review_scores_rating, color=city)) +
128   stat_summary(fun="mean") +
129   tabsCx="City" , y="Ratings" +
130   theme_minimal(J

```

IDMP Project.R* x



```

132 plot2 <- ggplot(df,aes(x=city,y=review_scores_nating,color=city))+
133   geom_boxplot()+
134   tabsCx="City",y="Rotings")+
135   theme_minimal()
136
137 Fig <- ggarrangeCptot1,ptot2}
138 annotate_figure{fig, top = text_grob("c")}
139
140 # Bedrooms/Bathrooms
141
142 ddf<- df %>% mutate(concated_column = paste(bedrooms, bathrooms, sep = '/')
143 ddf<-os.data.frame(table(ddf#concated_column,ddf#cityJJ %>%
144   group_by{Vor2} B>S
145   statistic_max{order_by = Freq, n = SQ B>B
146   group_by{Vor2} B>S statistic_maxCorder_by = Freq, n = 5)
147
148 ggplot(ddf, aesCx = Var1, y = Freq, facet=Var1)+
149   geom_bar(stat = "identity") +
150   geom_text(aes(label = Freq), vjust = -0.3) +
151   facet_wrapC-Vor2)+
152   labs(x = 'Type of Accommodations', y = "No. of Airbnb's",
153        title = "Type of Accommodation in different Cities") +
154   theme_minimal()
155
156 # Maps
157
158 mop_df <- mop_data("stote")
159 mop_data <- left_join(mop_df, df, by = c("region"="stote"))
160
161 states <- map_data("state")
162 counties <- map_data("county")
163
164 # Colifonnio Mop
165 co_df <- subset(states, region == "colifonnio")
166 co_county <- subset(counties, region == "colifonnio")
167 city_co <- map_data("stote",region='colifonnio')
168 df_ca <- filter(df,df$stote == 'colifonnio')
169
170 co_base <- ggplot(data = co_df, mapping = aes(x = long, y = lat, group = group)) +
171   coord_fixed(1.3) +
172   geom_polygon(color = "black", fill = "gray")
173
174 ca_base <- ggplot(data = ca_df, mapping = aes(x = long, y = lat, group = group)) +
175   coord_fixed(1.3) +
176   geom_polygon(color = "black", fill = "gray")

```

```

178 p<-c(0,300,600,900,1200,1500)
179
180 co_bose
181   geom_polygon(doto = co_county, fill='#C7DEE5', color  "black") +
182   geom_polygon(color  "black", fill = NA) +
183   geom_potnt(doto=df_co,oes(x=longitude,y=lotitude,color=prtce),tnherтт.oes = FALSE,stze=0.1) *
184   labs(title = "Atrbnb's in Colifornia (SF & LA) ond their Price")*
185   scale_color_gradientn(limtts  c(0,1500),
186                         colours=c("novybl ue", "#EEB1AA", "darkorange1", 'green', 'pi nk', "red",
187                                   breaks=p, labels=format{p}})+
188   theme_mtnimal{g
189
190 co_bose
191   geom_polygon(doto = co_county, fill='#C7DEE5', color  "black") +
192   geom_polygon(color  "black", fill = NA) +
193   geom_potnt(doto=df_co,oes(x=longitude,y=lotitude,color=prtce),tnherтт.oes = FALSE,stze=0.5) *
194   xli m(-123,-121.5) > yli in(37.5,38)*
195   scale_color_g rradientn(l imtl s  c(0,1500),
196                           col ours=c("novybl ue", "#EEB1AA", "dar korange1", 'gre en', 'pi nk', "red",
197                                   b rea ks=p, l abel s=format p} }) +
198   labs{titl e  "Atrbnb's I n San Franct sco and the I r Prt ce " }+
199   theme_mlni mat
200
201 co_bose
202   geom_polygon(doto = co_county, fill='#C7DEE5', color  "black") +
203   geom_polygon(color = "black", fill = NA) +
204   geom_potnt(doto=df_co,oes(x=longitude,y=lotitude,color=prtce),tnherтт.oes = FALSE,stze=0.5) *
285   xli m(-119,-117) + yli in(33,35g+
286   scale_color_g rradientn(l imtl s  c(0,1500),
287                           col ours=c("novybl ue", "#EEB1AA", "dar korange1", 'gre en', 'pi nk', "red",
288                                   b rea ks=p, l abel s=format p} }) +
289   labs{titl e  "Atrbnb's I n Los Ang eles and I hel r Prt ce " +
218   theme_mlni mat
211
212 A Boston MQp
213 mo_df <- subset(stotes, regton == "mossochusetts")
214 mo_county < subset(counties, region == "mossochusetts")
215 city_mo< mop_doto("stOte",region='mossochusetts')j
216 df_mo <- ftilter(df,df$stote == 'mossochusetts']
217
218 ma_base < ggplot data  ma_df, mapping= aes{x  long, y  lat, group  group} } +
219   coord_ftxed{1.3} +
220   geom_pol ygon colo r  "bl ack", fi TI  "gray " }
221
222 ma base < oaoilotLdata  ma df. fflODDino  aesix  l ono. v  l at. aroua  arouo  +

```



```

218 ma_bose <- ggplot(dota = ma_df, mapping = aes(x = long, y = tot, group = group)) +
219   coord_fixed(C1.3) +
220   geom_polygon(color = "black", fill = "gray")
221
222 mo_bose <- ggplot(dota = mo_df, mapping = aes(x = long, y = lot, group = group)) +
223   coord_fixed(C1.3) +
224   geom_polygon(color = "black", fill = "gray")
225
226 ma_bose +
227   geom_polygon(data = ma_county, fill = 'DC7DEE', color = "black") +
228   geom_polygon(color = "black", fill = NA) +
229   geom_point(dota = df_mo, aes(x = longitude, y = latitude, color = price), inherit.aes = FALSE, size = 0.1) +
230   labs(title = "Airbnb's in Massachusetts Boston and their Price") +
231   scale_color_gradientn(colors = c(0, 1500),
232     colors = c("navyblue", "#EEB1AA", "dorkorange1", "green", "pink", "red"),
233     breaks = p, labels = format(p)) +
234   theme_minimal()
235
236 ma_bose +
237   geom_polygon(data = ma_county, fill = 'DC7DEE', color = "black") +
238   geom_polygon(color = "black", fill = NA) +
239   geom_point(dota = df_mo, aes(x = longitude, y = latitude, color = price), inherit.aes = FALSE, size = 0.5) +
240   xlim(-71.5, -70.5) + ylim(42.2, 42.5) +
241   labs(title = "Airbnb's in Boston and their Price") +
242   scale_color_gradientn(colors = c(0, 1500),
243     colors = c("navyblue", "#EEB1AA", "dorkorange1", "green", "pink", "red"),
244     breaks = p, labels = format(p)) +
245   theme_minimal()
246
247 # New York Map
248 ny_df <- subset(states, region == "new york")
249 ny_county <- subset(counties, region == "new york")
250 city_ny <- map_data("state", region = "new york")
251 df_ny <- filter(df, df$state == "new york")
252
253 ny_bose <- ggplot(dota = ny_df, mapping = aes(x = long, y = lot, group = group)) +
254   coord_fixed(C1.3) +
255   geom_polygon(color = "black", fill = "gray")
256
257 ny_bose <- ggplot(dota = ny_df, mapping = aes(x = long, y = lot, group = group)) +
258   coord_fixed(C1.3) +
259   geom_polygon(color = "black", fill = "gray")
260
261 ny_bose +
262   geom_polygon(data = ny_county, fill = 'DC7DEE', color = "black") +

```



```

260
261 ny_bose +
262   geom_polygon(data = ny_county, VIII= 'DC 7DEE5' , color = "black" ) +
263   geom_polygon(color = "black" , VIII = NA)+
264   geom_point(data=df_ny, aes(x=longitude, y=latitude, color=price), inherit.aes = FALSE , size=0.1)+
265   scale_color_gradientn(colors = c(0,1500),
266     colors=c("navyblue", "#EEB1AA", "darkorange1", 'green', 'pink', "red"),
267     breaks=p , labels=format(p))+
268   theme_minimal()
269
270
271 ny_bose +
272   geom_polygon(data = ny_county, VIII= 'DC 7DEE5' , color = "black" ) +
273   geom_polygon(color = "black" , VIII = NA)+
274   geom_point(data=df_ny, aes(x=longitude, y=latitude, color=price), inherit.aes = FALSE , size=0.5)+
275   xlim(-74.3, -73.5) + ylim(40.5, 41)+
276   scale_color_gradientn(colors = c(0,1500),
277     colors=c("navyblue", "#EEB1AA", "darkorange1", 'green', 'pink', "red"),
278     breaks=p , labels=format(p))+
279   theme_minimal()
280
281
282 # Chicago Map
283 chi_df <- subset(Cstates, region == "Illinois")
284 chi_county <- subset(counties, region == "Illinois")
285 city_chi <- map_dof("state", region= "Illinois")
286 df_chi <- filter(df, state == "Illinois")
287
288 chi_bose <- ggplot(df_chi, mapping = aes(x = long, y = lat, group = groupJ) +
289   coord_fixed(1.3) +
290   geom_polygon(color = "black", fill = "gray"))
291
292 chi_bose <- ggplot(df_chi, mapping = aes(x = long, y = lat, group = groupJ) +
293   coord_fixed(1.3) +
294   geom_polygon(color = "black", fill = "gray"))
295
296 chi_base +
297   geom_polygon(df_chi, fill='DC7DEE5', color = "black")+
298   geom_polygon(color = "black" , VIII = NA)+
299   geom_point(data=df_chi, aes(x=longitude, y=latitude, color=price), inherit.aes = FALSE , size=0.1)+
300   scale_color_gradientn(colors = c(0,1500),
301     colors=c("navyblue", "#EEB1AA", "darkorange1", 'green', 'pink', "red"),
302     breaks=p , labels=format(p))+
303   theme_minimal()
304

```

```

30b chi_base +
307   geom_polygon(doto = cht_county, fill='SC7DEE5', color  "block"]
308   geom_polygon(color  "block", fill = NA)+
309   geom_point(doto=df_chi, aes(x=longitude, y=latitude, color=price), inherit.aes = FALSE, size=0.5)+
310   xlim( 86.9, 88)   ylim(41.6, 42.2)
311   labs(title="Chicago")+
312   scale_color_gradientn(limits = c(0,1500),
313     colours=c("navyblue", "#EEB1AA", "darkorange1", 'green', 'pink', "red",
314     breaks=p, labels=format(p))+
315   labs(title = "Attribute's in Chicago and their Price")+
316   theme_minimal()
317
318 # DC Map
319 dc_df <- subset(states, region  "district of
320 dc_county < subset(counties, region  "district of columbia")
321 city_dc <- map_dota("state", region='district of columbia')
322 df_dc <- filter(df, df$state == 'district
323
324 dc_base <- ggplot(doto = dc_df, mapping = aes(x  long, y  lat, group = group)) +
325   coord_fixed(1.3) +
326   geom_polygon(color  "black", fill  "gray"
327
328 dc_base <- ggplot(doto = dc_df, mapping = aes(x  long, y  lat, group = group)) +
329   coord_fixed(1.3) +
330   geom_polygon(color  "black", fill  "gray")
331
332 dc_base *
333   geom_polygon(doto = dc_county, fill='#C7DEE5', color  "block")
334   geom_polygon(color  "black", fill = NA)+
335   geom_point(doto=df_dc, aes(x=longitude, y=latitude, color=price), inherit.aes = FALSE, size=0.5)+
336   scale_color_gradientn(limits = c(0,1500),
337     colours=c("navyblue", "#EEB1AA", "darkorange1", 'green', 'pink', "red",
338     breaks=p, labels=format(p))+
339   labs(title = "Attribute's in DC and their Price")+
340   theme_minimal()
341
342 # Model
343 set.seed(3)
344 df_part <- resample_partition(df,
345   p=c(train=0.6,
346     val=0.2,
347     test=0.2)

```


Of IDMP Project.R'

4 Run ^4

Source

```

SourceonSave @ •
348 step1 <- functon(response, predictors, condtdotes, portition]
349 (
350   rhs <- poste0(poste0(predtctors, collopse="+"), " ", condidotes)
351   formulas <- lopply(poste0(response, "-", rhs), os.formulo]
352   rmses <- sopply(formulos,
353     function(fm) rmse(lm(fm, doto=portttion$troinj,
354       doto=portittonBvolttd))
355   nomes(rmses] < condtdotes
356   otr(rmses, < rmses#which.min(rmses]]
357   rmses
358- 1
359 model <- NULL
360
361 # Step-1
362 preds <- " 1 "
363 conds <- c("property_type", "room_type", "Occommodotes", "bathrooms",
364   "ci ty", "be drooms", "review_s core s_ratl ng")
365 s1 <- step1{"log_p r l ce", pred s, conds, dE_pa rt
366
367 model <- c(model, att r s1, "best "
368   s1
369
370 # Step-2
371 preds < "room_type"
372 conds <- c("property_type", "Occommodotes", "bothrooms",
373   "city", "bedrooms", "review_scores_rotting"]
374 s1 <- step1{"log_p r l ce", pred s, conds, dE_pa rt
375
376 model <- c(model, att r s1, "best "
377   s1
378
379 # Step-3
380 preds< c("room_type", "review_scores_rotting"]
381 conds <- c("property_type", "Occommodotes", "bothrooms",
382   "city", "bedrooms")
383 s1 <- step1{"log_p r l ce", pred s, conds, dE_pa rt
384
385 model <- c(model, att r s1, "best "
386   s1
387
388 # Step —^
389 preds< c("room_type", "review_scores_rotting", "bedrooms"]
390 conds <- c("property_type", "Occommodotes", "bothrooms",
391
392 s1 <- step1{"log_p r l ce", preds, conds, d f_part

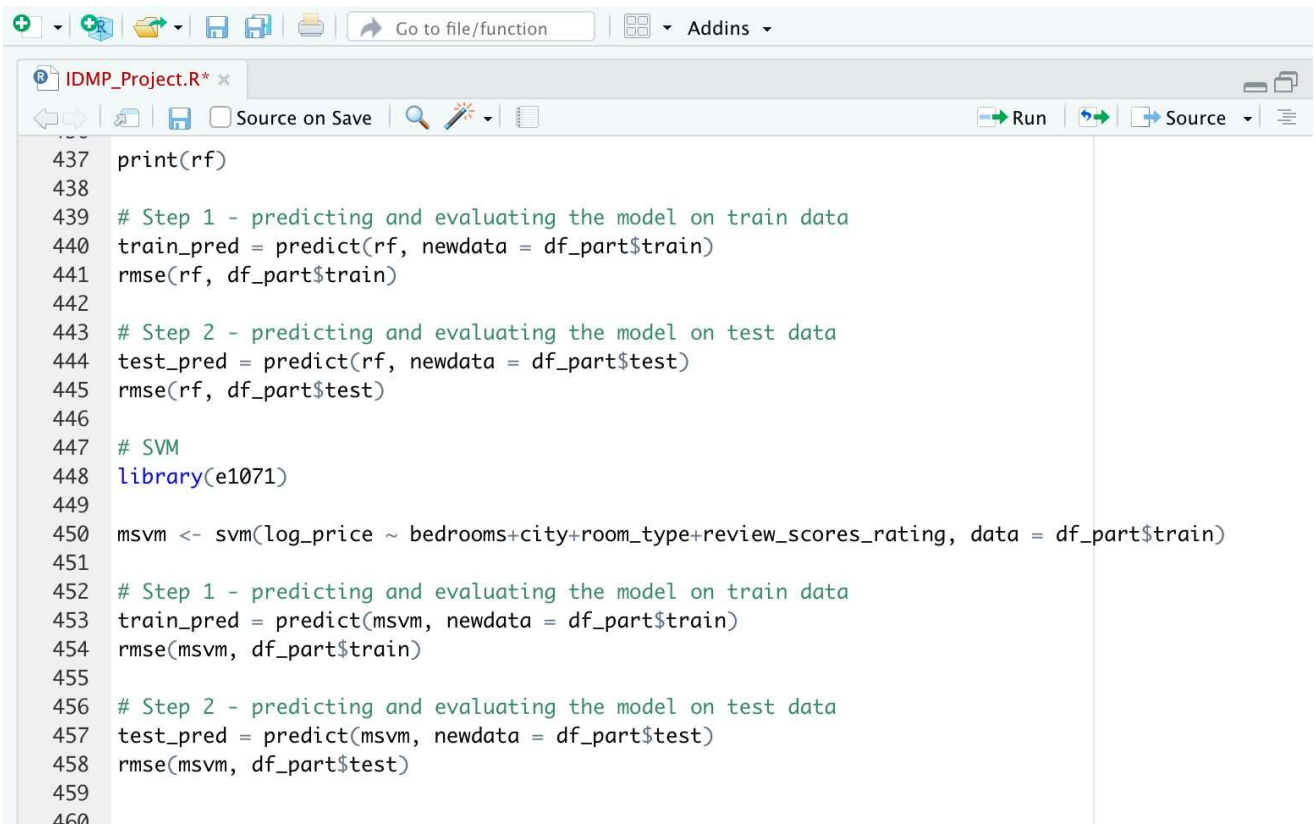
```



```

394 model <- c(model, attr s1, "best"
395 s1
396
397 # Step 4
398 preds <- c("room_type", "review_scores_rating", "bedrooms", "city")
399 conds <- c("property_type", "Oocommodotes", "bothrooms")
400 sl <- step1("log_prtce", preds, conds, df_port]
401
402 model <- c(model, attr s1, "best"
403 s1
404
405 # Step 5
406 preds<- c("room_type", "review_scores_rating", "bedrooms", "city", "oocommodotes")
407 conds <- c("property_type", "bothrooms")
408 sl <- step1("log_prtce", preds, conds, df_port]
409
410 model <- c(model, attr s1, "best"
411 s1
412
413 step_model <- ti bble{i ndex-seq_al ong{ model } ,
414                       variable-factor{names{model, l evels-names{ model } }},
415                       RMSE=model}
416
417 ggplot step_model, aes{y-RMSE}} +
418   geom_pot nt aes{x=var l able}
419   geom_li ne{ aes{x-index} } +
420   l abs{titl e-" Stepwt se model set ection "} +
421   theme_minimal()
422
423 fit_model <- lm(dg prtde - room_type * review_scores_rottnng * bedrooms + ctty ,
424
425
426 ftt5 <- lm{log_pr l ce - roorr_type + review_scores_rat l ng + bedrooms + ctty ,
427          doto = df_portBtrotn]
428 rmse(fit5,df_port5troinj
429 rmse(fit5,df_port5valid)
430 rmse(fit5,df_port5test)
431
432 # Random Forest
433
434 rf <- rondonForest( og_price - room_type review_scores_rottnng bedrooms ctty,
435                    doto = df_port$trotn, mtry=3, importonce = TRUE, no.octton = no.omits
436
437 prtnr(rf)
438

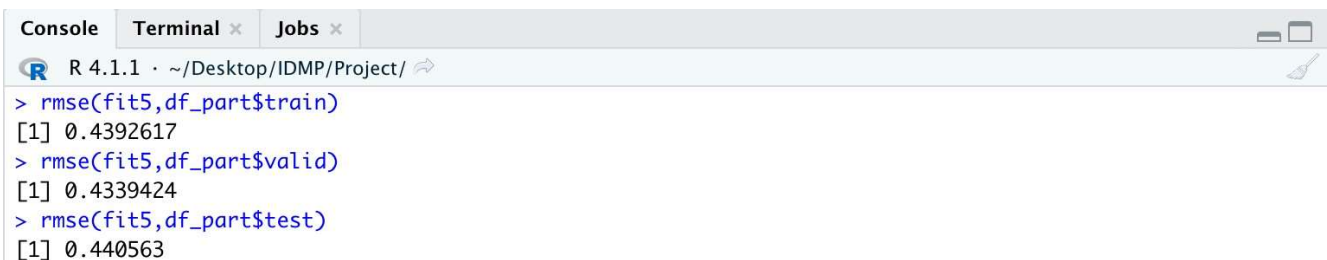
```



The screenshot shows the RStudio editor interface with a file named 'IDMP_Project.R'. The code is as follows:

```
437 print(rf)
438
439 # Step 1 - predicting and evaluating the model on train data
440 train_pred = predict(rf, newdata = df_part$train)
441 rmse(rf, df_part$train)
442
443 # Step 2 - predicting and evaluating the model on test data
444 test_pred = predict(rf, newdata = df_part$test)
445 rmse(rf, df_part$test)
446
447 # SVM
448 library(e1071)
449
450 msvm <- svm(log_price ~ bedrooms+city+room_type+review_scores_rating, data = df_part$train)
451
452 # Step 1 - predicting and evaluating the model on train data
453 train_pred = predict(msvm, newdata = df_part$train)
454 rmse(msvm, df_part$train)
455
456 # Step 2 - predicting and evaluating the model on test data
457 test_pred = predict(msvm, newdata = df_part$test)
458 rmse(msvm, df_part$test)
459
460
```

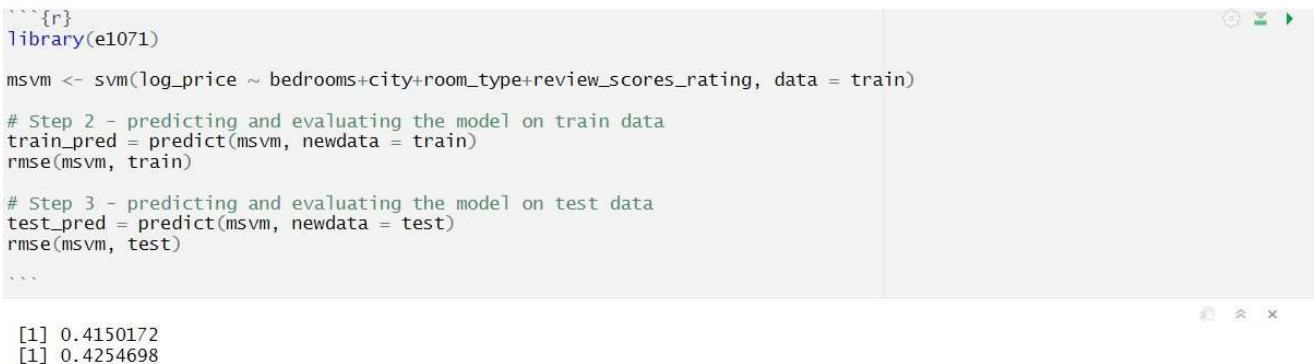
For Linear Regression



The screenshot shows the RStudio Console with the following output:

```
R 4.1.1 · ~/Desktop/IDMP/Project/
> rmse(fit5,df_part$train)
[1] 0.4392617
> rmse(fit5,df_part$valid)
[1] 0.4339424
> rmse(fit5,df_part$test)
[1] 0.440563
```

For SVM



The screenshot shows the RStudio editor interface with the following code:

```
{r}
library(e1071)

msvm <- svm(log_price ~ bedrooms+city+room_type+review_scores_rating, data = train)

# Step 2 - predicting and evaluating the model on train data
train_pred = predict(msvm, newdata = train)
rmse(msvm, train)

# Step 3 - predicting and evaluating the model on test data
test_pred = predict(msvm, newdata = test)
rmse(msvm, test)
```

Below the code, the console output is shown:

```
[1] 0.4150172
[1] 0.4254698
```

For Random Forest

```
##{r}
# Load the library
library(randomForest)

train

rf <- randomForest(log_price ~ bedrooms+city+room_type+review_scores_rating, data = train, mtry=3, importance = TRUE, na.action = na.omit)

# print(rf)

# Step 2 - predicting and evaluating the model on train data
train_pred = predict(rf, newdata = train)
rmse(rf, train)

# Step 3 - predicting and evaluating the model on test data
test_pred = predict(rf, newdata = test)
rmse(rf, test)

##
```

data.frame
25578 x 10

R Console

```
[1] 0.4051719
[1] 0.4266298
```