

## ➤ Importing Some Required Library and the Data Set :

```
R 4.1.3 · ~/ ~
> library(tidyverse)
> library(glue)
> library(dplyr)
> library(esquisse)
> library(ggthemes)
> library(lubridate)
> library(tinytex)
> library(shiny)
> library(knitr)
>
>
> #Importing Data
> airbnb <- read.csv("C:/Users/manis/Desktop/Industry_Assignment_2/Assignment2_Solution/Airbnb.csv",header = TRUE)
>
> #Viewing the data set
> head(airbnb)
```

id name host\_id host\_name neighbourhood\_group neighbourhood latitude longitude room\_type price minimum\_nights number\_of\_reviews last\_review reviews\_per\_month

1	Clean & quiet apt home by the park	2787	John	Brooklyn	Kensington	40.64749
2	Skylit Midtown Castle	2845	Jennifer	Manhattan	Midtown	40.75362
3	THE VILLAGE OF HARLEM...NEW YORK !	4632	Elisabeth	Manhattan	Harlem	40.80902
4	Cozy Entire Floor of Brownstone	4869	LisaRoxanne	Brooklyn	Clinton Hill	40.68514
5	Entire Apt: Spacious Studio/Loft by central park	7192	Laura	Manhattan	East Harlem	40.79851
6	Large Cozy 1 BR Apartment In Midtown East	7322	Chris	Manhattan	Murray Hill	40.74767

1 -73.97237 Private room 149 1 9 10/19/2018 0.21  
 2 -73.98377 Entire home/apt 225 1 45 5/21/2019 0.38  
 3 -73.94190 Private room 150 3 0 NA  
 4 -73.95976 Entire home/apt 89 1 270 7/5/2019 4.64  
 5 -73.94399 Entire home/apt 80 10 9 11/19/2018 0.10  
 6 -73.97500 Entire home/apt 200 3 74 6/22/2019 0.59

```
R 4.1.3 · ~/ ~
> summary(airbnb)
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood
Min. :	2539	Length:48895	Min. : 2438	Length:48895	Length:48895	Length:48895
1st Qu.:	9471945	Class :character	1st Qu.: 7822033	Class :character	Class :character	Class :character
Median :	19677284	Mode :character	Median : 30793816	Mode :character	Mode :character	Mode :character
Mean :	19017143		Mean : 67620011			
3rd Qu.:	29152178		3rd Qu.:107434423			
Max. :	36487245		Max. :274321313			

	latitude	longitude	room_type	price	minimum_nights	number_of_reviews
Min. :	40.50	Min. :-74.24	Length:48895	Min. : 0.0	Min. : 1.00	Min. : 0.00
1st Qu.:	40.69	1st Qu.:-73.98	Class :character	1st Qu.: 69.0	1st Qu.: 1.00	1st Qu.: 1.00
Median :	40.72	Median :-73.96	Mode :character	Median : 106.0	Median : 3.00	Median : 5.00
Mean :	40.73	Mean :-73.95		Mean : 152.7	Mean : 7.03	Mean : 23.27
3rd Qu.:	40.76	3rd Qu.:-73.94		3rd Qu.: 175.0	3rd Qu.: 5.00	3rd Qu.: 24.00
Max. :	40.91	Max. :-73.71		Max. :10000.0	Max. :1250.00	Max. :629.00

	last_review	reviews_per_month	calculated_host_listings_count	availability_365
Length:	48895	Min. : 0.010	Min. : 1.000	Min. : 0.0
Class :	character	1st Qu.: 0.190	1st Qu.: 1.000	1st Qu.: 0.0
Mode :	character	Median : 0.720	Median : 1.000	Median : 45.0
		Mean : 1.373	Mean : 7.144	Mean :112.8
		3rd Qu.: 2.020	3rd Qu.: 2.000	3rd Qu.:227.0
		Max. :58.500	Max. :327.000	Max. :365.0
		NA's :10052		

```
> glimpse(airbnb)
Rows: 48,895
Columns: 16
$ id
$ name
$ host_id
$ host_name
$ neighbourhood_group
$ neighbourhood
$ latitude
$ longitude
$ room_type
$ price
$ minimum_nights
$ number_of_reviews
$ last_review
$ reviews_per_month
$ calculated_host_listings_count
$ availability_365
```

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Project: (None)

Source

R 4.1.3 · ~/airbnb · max. 1.0000.00 max. 1.200.00 max. 1.025.00
last_review reviews_per_month calculated_host_listings_count availability_365
Length:48895 Min. : 0.010 Min. : 1.000 Min. : 0.0
Class :character 1st Qu.: 0.190 1st Qu.: 1.000 1st Qu.: 0.0
Mode :character Median : 0.720 Median : 1.000 Median : 45.0
Mean : 1.373 Mean : 7.144 Mean :112.8
3rd Qu.: 2.020 3rd Qu.: 2.000 3rd Qu.:227.0
Max. :58.500 Max. :327.000 Max. :365.0
NA's :10052

> glimpse(airbnb)
Rows: 48,895
Columns: 16
$ id
$ name
$ host_id
$ host_name
$ neighbourhood_group
$ neighbourhood
$ latitude
$ longitude
$ room_type
$ price
$ minimum_nights
$ number_of_reviews
$ last_review
$ reviews_per_month
$ calculated_host_listings_count
$ availability_365
>

```

## ➤ Analyzing the Data Set to Perform :

- 1) Pre - Processing**
- 2) Exploratory Data Analysis.**
- 3) Plot Graphs Based On Some Conditions.**

## **1) Pre - Processing :**

### a) Treat missing values :

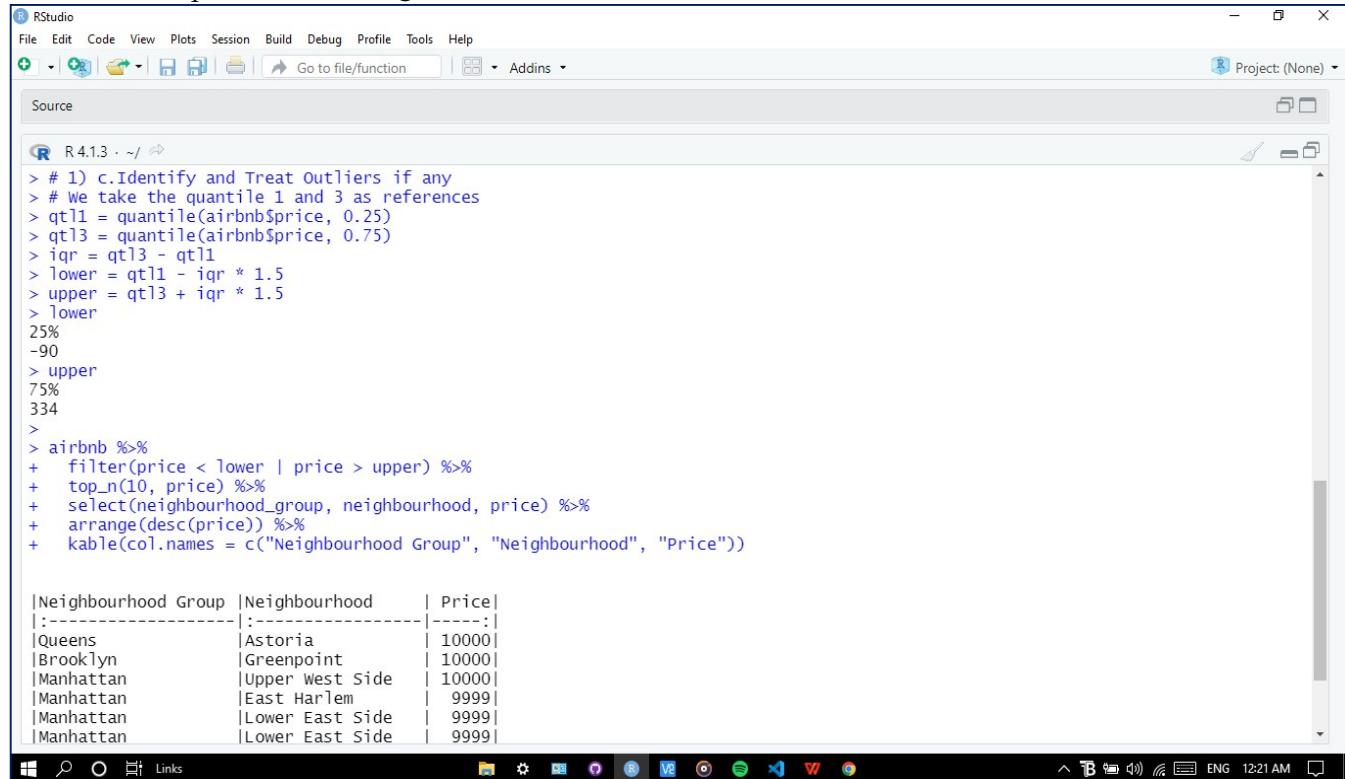
#There are Total 10052 values missing from reviews per month column

**b)** Remove duplicates if any :

#There are No Duplicated Values in the data set

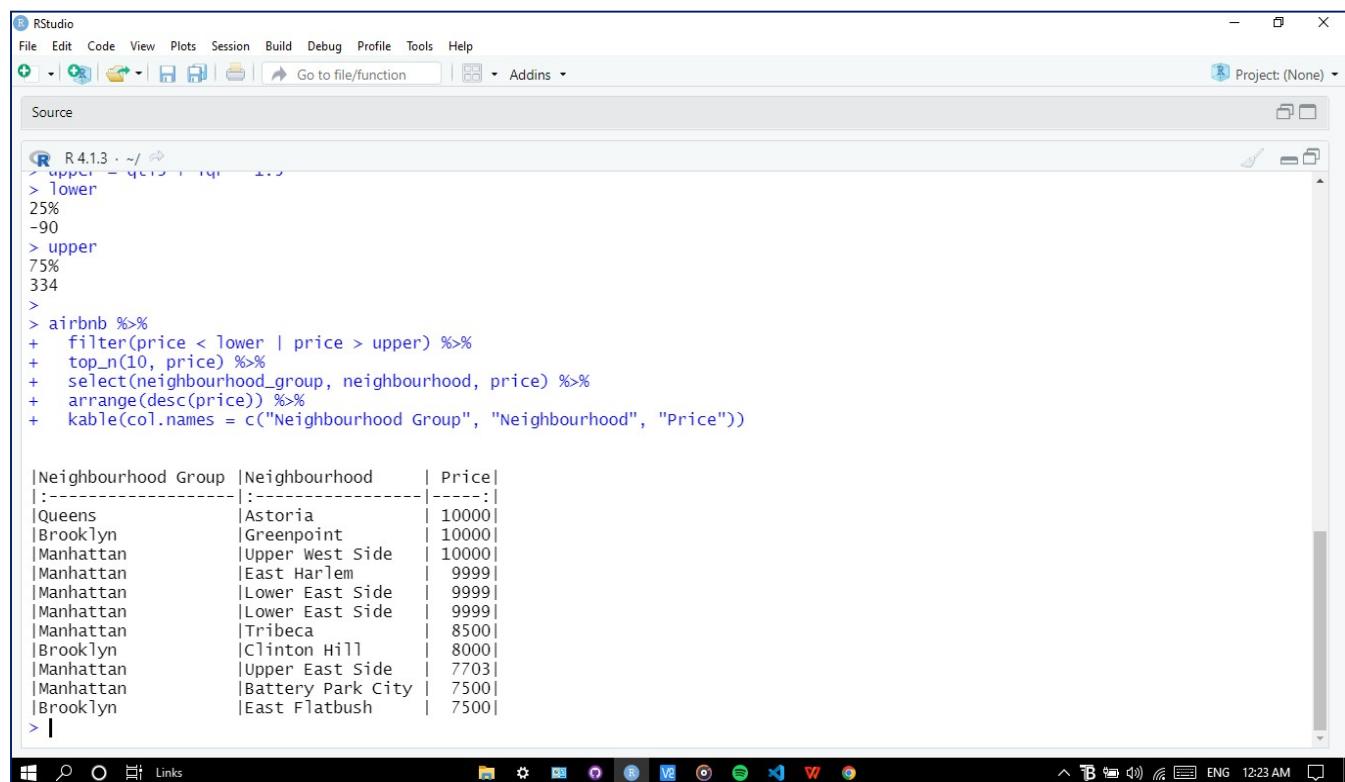
## c) Identify and treat outliers if any :

Before starting our analysis, we also want to check the outlier points in this data set and we take the quantile 1 and 3 as references.



```
R 4.1.3 · ~/ ◀
> # 1) c. Identify and Treat Outliers if any
> # We take the quantile 1 and 3 as references
> qt1l = quantile(airbnb$price, 0.25)
> qt13 = quantile(airbnb$price, 0.75)
> iqr = qt13 - qt1l
> lower = qt1l - iqr * 1.5
> upper = qt13 + iqr * 1.5
> lower
25%
-90
> upper
75%
334
>
> airbnb %>%
+   filter(price < lower | price > upper) %>%
+   top_n(10, price) %>%
+   select(neighbourhood_group, neighbourhood, price) %>%
+   arrange(desc(price)) %>%
+   kable(col.names = c("Neighbourhood Group", "Neighbourhood", "Price"))

|Neighbourhood Group|Neighbourhood|Price|
|:-----|:-----|-----:|
|Queens|Astoria|10000|
|Brooklyn|Greenpoint|10000|
|Manhattan|Upper West Side|10000|
|Manhattan|East Harlem|9999|
|Manhattan|Lower East Side|9999|
|Manhattan|Lower East Side|9999|
```



```
R 4.1.3 · ~/ ◀
> upper = qt13 + iqr * 1.5
> lower
25%
-90
> upper
75%
334
>
> airbnb %>%
+   filter(price < lower | price > upper) %>%
+   top_n(10, price) %>%
+   select(neighbourhood_group, neighbourhood, price) %>%
+   arrange(desc(price)) %>%
+   kable(col.names = c("Neighbourhood Group", "Neighbourhood", "Price"))

|Neighbourhood Group|Neighbourhood|Price|
|:-----|:-----|-----:|
|Queens|Astoria|10000|
|Brooklyn|Greenpoint|10000|
|Manhattan|Upper West Side|10000|
|Manhattan|East Harlem|9999|
|Manhattan|Lower East Side|9999|
|Manhattan|Lower East Side|9999|
|Manhattan|Tribeca|8500|
|Brooklyn|Clinton Hill|8000|
|Manhattan|Upper East Side|7703|
|Manhattan|Battery Park City|7500|
|Brooklyn|East Flatbush|7500|
```

**d) Transform Data based on judgement and explain the methodology used :**

### **Methodology I Used :**

The methodology I have used here is that first I have checked the unique values in neighborhood\_group, room\_type, neighbourhood, host\_name, and minimum\_nights and then I have converted the neighbourhood\_group, neighbourhood and room type to factors to see the types of it with the help of level() function.

In the data set the last review column is not properly defined, so in order to remove the complexity there should be a specific format. I have separated the dates with slashes (/), saved them by Month, Day, and Year , and then replaced NA values with zero (0), then converted them back to integers!

Replacing NA values in review\_per\_month column with zero(0).

Viewing up the summary!

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins ▾
Project: (None) ▾
Source
R 4.1.3 · ~/ ▾
> # 1) d.Transform Data based on judgement and explain the methodology used
>
> #Checking Up the Unique Values in neighborhood_group, room_type column(nominal qualitative)
> #For extracting unique values
> c(unique(airbnb[["neighbourhood_group"]]))
$neighbourhood_group
[1] "Brooklyn"      "Manhattan"     "Queens"        "Staten Island" "Bronx"
> c(unique(airbnb[["room_type"])))
$room_type
[1] "Private room"   "Entire home/apt" "Shared room"
> c(unique(airbnb[["neighbourhood"])))
$neighbourhood
[1] "Kensington"      "Midtown"       "Harlem"        "Clinton Hill"
[5] "East Harlem"    "Murray Hill"  "Bedford-Stuyvesant" "Hell's Kitchen"
[9] "Upper West Side" "Chinatown"    "South Slope"   "West Village"
[13] "Williamsburg"   "Fort Greene"  "Chelsea"       "Crown Heights"
[17] "Park Slope"     "Windsor Terrace" "Inwood"        "East Village"
[21] "Greenpoint"      "Bushwick"    "Flatbush"      "Lower East Side"
[25] "Prospect-Lefferts Gardens" "Long Island City" "Kips Bay"      "SoHo"
[29] "Upper East Side" "Prospect Heights" "Washington Heights" "Woodside"
[33] "Brooklyn Heights" "Carroll Gardens" "Gowanus"      "Flatlands"
[37] "Cobble Hill"    "Flushing"    "Boerum Hill"   "Sunnyside"
[41] "DUMBO"          "St. George"   "Highbridge"    "Financial District"
[45] "Ridgewood"       "Morningside Heights" "Jamaica"      "Middle Village"
[49] "NoHo"           "Ditmars Steinway" "Flatiron District" "Roosevelt Island"
[53] "Greenwich Village" "Little Italy" "East Flatbush" "Tompkinsville"
[57] "Astoria"        "Clason Point" "Eastchester"  "Kingsbridge"
[61] "Two Bridges"    "Queens Village" "Rockaway Beach" "Forest Hills"
```

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins ▾
Project: (None) ▾
Source
R 4.1.3 · ~/ ▾
> c(unique(airbnb[["host_name"]]))
$host_name
[1] "John"          "Jennifer"      "Elisabeth"
[4] "LisaRoxanne"   "Laura"         "Chris"
[7] "Garon"         "Shunichi"     "MaryEllen"
[10] "Ben"           "Lena"          "Kate"
[13] "Laurie"        "Claudio"      "Alina"
[16] "Allen & Irina" "Jane"          "Doti"
[19] "Adam And Charity" "Sing"          "Chaya"
[22] "LiseI"         "Nathalie"     "Gregory"
[25] "Claude & Sophie" "Tommi"        "Shon"
[28] "Dana"          "Ssameer Or Trip" "Teri"
[31] "Andrea"        "Angela"        "Vt"
[34] "Tyrome"        "Harriet"      "Edward"
[37] "Abdu1"         "Yolande"      "Cyn"
[40] "Earl"          "Rana"          "Orestes"
[43] "Adreinne"      "Alexander"   "JT And Tiziana"
[46] "Joya"          "James"         "Jeanne"
[49] "Francesca"    "Joanna"       "Bianca"
[52] "Luiz"          "Ted"          "Cristina"
[55] "Petra"         "D"            "Dimitri"
[58] "Patricia"      "Mark"         "Sara"
[61] "Reka"          "Daniel"       "Casey"
[64] "Robin"         "Anna"         "Enzo"
[67] "Tye And Etienne" "George"      "Josh"
[70] "Victoria"      "Justin"       "Blaise"
[73] "DAVID And RICK" "LulaO"       "Sybilla"
[76] "JoLynn"         "Gaia"          "Ana"
[79] "Maggie"        "Starlee"      "Pas"
[82] "Aurietia"      "Sora & Lymotto" "Erica"
```

RStudio  
File Edit Code View Plots Session Build Debug Profile Tools Help  
+ - Go to file/function Addins Project: (None) ▾

Source R 4.1.3 · ~/ ▾

```
> c(unique(airbnb$minimum_nights))
$minimum_nights
 [1]   1   3  10  45   2   5   4   90   7  14   60   29   30  180   9   31   6   15   8   26   28  200   50   17
[25]  21  11  25  13  35   27  18  20  40   44   65   55  120  365  122  19  240   88  115  150  370  16   80  181
[49] 265 300  59 185 360   56  12  70  39   24   32 1000 110 270  22  75  250   62  23 1250 364  74 198 100
[73] 500  43  91 480  53   99  160   47 999 186 366  68   93  87 183 299 175   98 133 354   42  33  37 225
[97] 400 105 184 153 134  222   58 210 275 182 114  85   36

>
>
> #Converting neighbourhood_group,neighbourhood and room_type into factors
> airbnb$neighbourhood_group <- as.factor(airbnb$neighbourhood_group)
> levels(airbnb$neighbourhood_group)
[1] "Bronx"          "Brooklyn"        "Manhattan"       "Queens"         "Staten Island"
>
> airbnb$room_type <- as.factor(airbnb$room_type)
> levels(airbnb$room_type)
[1] "Entire home/apt" "Private room"    "Shared room"
>
> airbnb$neighbourhood <- as.factor(airbnb$neighbourhood)
> levels(airbnb$neighbourhood)
[1] "Allerton"           "Arden Heights"      "Arrochar"          "Arverne"
[5] "Astoria"            "Bath Beach"         "Battery Park City" "Bay Ridge"
[9] "Bay Terrace"         "Bay Terrace, Staten Island" "Baychester"        "Bayside"
[13] "Bayswater"          "Bedford-Stuyvesant" "Belle Harbor"      "Bellerose"
[17] "Belmont"             "Bensonhurst"        "Bergen Beach"      "Boerum Hill"
[21] "Borough Park"       "Breezy Point"        "Briarwood"          "Brighton Beach"
[25] "Bronxdale"           "Brooklyn Heights"   "Brownsville"        "Bull's Head"
[29] "Bushwick"            "Cambria Heights"    "Canarsie"          "Carroll Gardens"
```

The screenshot shows an RStudio interface with the following details:

- Header:** RStudio, File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Toolbar:** Contains icons for file operations like Open, Save, and Print, along with Go to file/function and Addins dropdown.
- Project:** Project (None) is selected.
- Source Editor:** The code is written in R, starting with `R 4.1.3` and `#~`.
- Code Content:**

```
> #Converting neighbourhood_group,neighbourhood and room_type into factors
> airbnb$neighbourhood_group <- as.factor(airbnb$neighbourhood_group)
> levels(airbnb$neighbourhood_group)
[1] "Bronx"           "Brooklyn"        "Manhattan"       "Queens"          "Staten Island"
>
> airbnb$room_type <- as.factor(airbnb$room_type)
> levels(airbnb$room_type)
[1] "Entire home/apt" "Private room"    "Shared room"
>
> airbnb$neighbourhood <- as.factor(airbnb$neighbourhood)
> levels(airbnb$neighbourhood)
[1] "Allerton"         "Arden Heights"   "Arrochar"        "Arverne"
[5] "Astoria"          "Bath Beach"      "Battery Park City" "Bay Ridge"
[9] "Bay Terrace"       "Bay Terrace, Staten Island" "Baychester"      "Bayside"
[13] "Bayswater"         "Bedford-Stuyvesant" "Belle Harbor"    "Bellerose"
[17] "Belmont"          "Bensonhurst"     "Bergen Beach"   "Boerum Hill"
[21] "Borough Park"    "Breezy Point"     "Briarwood"       "Brighton Beach"
[25] "Bronxdale"         "Brooklyn Heights" "Brownsville"     "Bull's Head"
[29] "Bushwick"         "Cambria Heights"  "Canarsie"        "Carroll Gardens"
[33] "Castle Hill"       "Castleton Corners" "Chelsea"         "Chinatown"
[37] "City Island"       "Civic Center"     "Claremont Village" "Clason Point"
[41] "Clifton"          "Clinton Hill"    "Co-op City"     "Cobble Hill"
[45] "College Point"    "Columbia St"    "Concord"         "Concourse"
[49] "Concourse Village" "Coney Island"    "Corona"          "Crown Heights"
[53] "Cypress Hills"    "Ditmars Steinway" "Dongan Hills"   "Douglaston"
[57] "Downtown Brooklyn" "DUMBO"           "Dyker Heights"  "East Elmhurst"
[61] "East Flatbush"     "East Harlem"     "East Morrisania" "East New York"
[65] "East Village"      "Eastchester"     "Edenwald"        "Edgemere"
[69] "Elmhurst"          "Eltingville"     "Emerson Hill"   "Far Rockaway"
```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins Project: (None)

Source

```
R 4.1.3 · ~/ ↘
>
> #Separating Dates
> airbnb <- tidyverse::separate(airbnb, last_review, c("Month", "Day", "Year"), sep = "/")
Warning message:
Expected 3 pieces. Missing pieces filled with `NA` in 10052 rows [3, 20, 27, 37, 39, 194, 205, 261, 266, 268, 277, 346, 350, 391,
426, 433, 438, 487, 546, 586, ...].
>
> #Replacing NA with 0
> airbnb$Year[is.na(airbnb$Year) == TRUE] = 0
> airbnb$Month[is.na(airbnb$Month) == TRUE] = 0
> airbnb$Day[is.na(airbnb$Day) == TRUE] = 0
>
> #Datatype Conversion
> airbnb$Month <- as.integer(airbnb$Month)
> airbnb$Year <- as.integer(airbnb$Year)
> airbnb$Day <- as.integer(airbnb$Day)
> View(airbnb)
>
>
> #Replacing NA in review_per_month with 0
> airbnb$reviews_per_month[is.na(airbnb$reviews_per_month) == TRUE] = 0
>
> #Checking for any NA value remaining
> sapply(airbnb, function(x) sum(is.na(x)))
      id             name          host_id        host_name
      0              0              0              0
neighbourhood_group neighbourhood      latitude      longitude
      0              0              0              0
      room_type       price      minimum_nights number_of_reviews
      0              0              0              0
```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins Project: (None)

Source

```
R 4.1.3 · ~/ ↘
> #Replacing NA in review_per_month with 0
> airbnb$reviews_per_month[is.na(airbnb$reviews_per_month) == TRUE] = 0
>
> #Checking for any NA value remaining
> sapply(airbnb, function(x) sum(is.na(x)))
      id             name          host_id        host_name
      0              0              0              0
neighbourhood_group neighbourhood      latitude      longitude
      0              0              0              0
      room_type       price      minimum_nights number_of_reviews
      0              0              0              0
      Month          Day           Year reviews_per_month
      0              0           0              0
calculated_host_listings_count availability_365
      0              0
```

> glimpse(airbnb)

Rows: 48,895  
Columns: 18

	<code>\$ id</code>	<code>\$ name</code>	<code>\$ host_id</code>	<code>\$ host_name</code>	<code>\$ neighbourhood_group</code>	<code>\$ neighbourhood</code>	<code>\$ latitude</code>	<code>\$ longitude</code>	<code>\$ room_type</code>	<code>\$ price</code>	<code>\$ minimum_nights</code>	<code>\$ calculated_host_listings_count</code>	<code>\$ availability_365</code>
	<code>&lt;int&gt;</code> 2539, 2595, 3647, 3831, 5022, 5099, 5121, 5178, 5203, 5238, 5295, 5441, 5803, 6021, 6~	<code>&lt;chr&gt;</code> "Clean & quiet apt home by the park", "Skylit Midtown Castle", "THE VILLAGE OF HARLEM~	<code>&lt;int&gt;</code> 2787, 2845, 4632, 4869, 7192, 7322, 7356, 8967, 7490, 7549, 7702, 7989, 9744, 11528, ~	<code>&lt;chr&gt;</code> "John", "Jennifer", "Elisabeth", "LisaRoxanne", "Laura", "Chris", "Garon", "Shunichi"~	<code>&lt;fct&gt;</code> Brooklyn, Manhattan, Manhattan, Brooklyn, Manhattan, Brooklyn, Manhattan, ~	<code>&lt;fct&gt;</code> "Kensington", "Midtown", "Harlem", "Clinton Hill", "East Harlem", "Murray Hill", "Bed~	<code>&lt;dbl&gt;</code> 40.64749, 40.75362, 40.80902, 40.68514, 40.74767, 40.68688, 40.76489, 40.80~	<code>&lt;dbl&gt;</code> -73.97237, -73.98377, -73.94190, -73.95976, -73.94399, -73.97500, -73.95596, -73.9849~	<code>&lt;fct&gt;</code> Private room, Entire home/apt, Private room, Entire home/apt, Entire home/apt, Entire~	<code>&lt;int&gt;</code> 149, 225, 150, 89, 80, 200, 60, 79, 150, 135, 85, 89, 85, 120, 140, 215, 140, 99, ~	<code>&lt;int&gt;</code> 1, 1, 3, 1, 10, 3, 45, 2, 2, 1, 5, 2, 4, 2, 90, 2, 2, 1, 3, 7, 3, 2, 1, 2, 2, 1, 4, 1~	<code>&lt;int&gt;</code> 9, 45, 0, 270, 9, 74, 49, 430, 118, 160, 53, 188, 167, 113, 27, 148, 198, 260, 53, 0, ~	

R 4.1.3 · ~/

```
> glimpse(airbnb)
Rows: 48,895
Columns: 18
$ id
$ name
$ host_id
$ host_name
$ neighbourhood_group
$ neighbourhood
$ latitude
$ longitude
$ room_type
$ price
$ minimum_nights
$ number_of_reviews
$ Month
$ Day
$ Year
$ reviews_per_month
$ calculated_host_listings_count
$ availability_365
>
> summary(airbnb)
```

	id	name	host_id	host_name	neighbourhood_group
Min.	: 2539	Length:48895	Min. : 2438	Length:48895	Bronx : 1091
1st Qu.:	9471945	Class :character	1st Qu.: 7822033	Class :character	Brooklyn : 20104
Median :	19677284	Mode :character	Median : 30793816	Mode :character	Manhattan : 21661
Mean :	19017143		Mean : 67620011		Queens : 5666
3rd Qu.:	29152178		3rd Qu.:107434423		Staten Island: 373
Max.:	36487245		Max. :274321313		

R 4.1.3 · ~/

```
> summary(airbnb)
```

	id	name	host_id	host_name	neighbourhood_group
Min.	: 2539	Length:48895	Min. : 2438	Length:48895	Bronx : 1091
1st Qu.:	9471945	Class :character	1st Qu.: 7822033	Class :character	Brooklyn : 20104
Median :	19677284	Mode :character	Median : 30793816	Mode :character	Manhattan : 21661
Mean :	19017143		Mean : 67620011		Queens : 5666
3rd Qu.:	29152178		3rd Qu.:107434423		Staten Island: 373
Max.:	36487245		Max. :274321313		

	neighbourhood	latitude	longitude	room_type	price	minimum_nights
Williamsburg	: 3920	Min. :40.50	Min. :-74.24	Entire home/apt:25409	Min. : 0.0	Min. : 1.00
Bedford-Stuyvesant	: 3714	1st Qu.:40.69	1st Qu.:-73.98	Private room :22326	1st Qu.: 69.0	1st Qu.: 1.00
Harlem	: 2658	Median :40.72	Median :-73.96	Shared room : 1160	Median : 106.0	Median : 3.00
Bushwick	: 2465	Mean :40.73	Mean :-73.95		Mean : 152.7	Mean : 7.03
Upper West Side	: 1971	3rd Qu.:40.76	3rd Qu.:-73.94		3rd Qu.: 175.0	3rd Qu.: 5.00
Hell's Kitchen	: 1958	Max. :40.91	Max. :-73.71		Max. :10000.0	Max. :1250.00
(Other)	: 32209					

	number_of_reviews	Month	Day	Year	reviews_per_month	calculated_host_listings_count
Min.:	0.00	Min. : 1.000	Min. : 0.00	Min. : 0	Min. : 0.000	Min. : 1.000
1st Qu.:	1.00	1st Qu.: 5.000	1st Qu.: 1.00	1st Qu.:2016	1st Qu.: 0.040	1st Qu.: 1.000
Median :	5.00	Median : 6.000	Median :11.00	Median :2019	Median : 0.370	Median : 1.000
Mean :	23.27	Mean : 6.174	Mean :12.54	Mean :1603	Mean : 1.091	Mean : 7.144
3rd Qu.:	24.00	3rd Qu.: 7.000	3rd Qu.:23.00	3rd Qu.:2019	3rd Qu.: 1.580	3rd Qu.: 2.000
Max.:	629.00	Max. :12.000	Max. :31.00	Max. :2019	Max. :58.500	Max. :327.000
NA's:	10052					

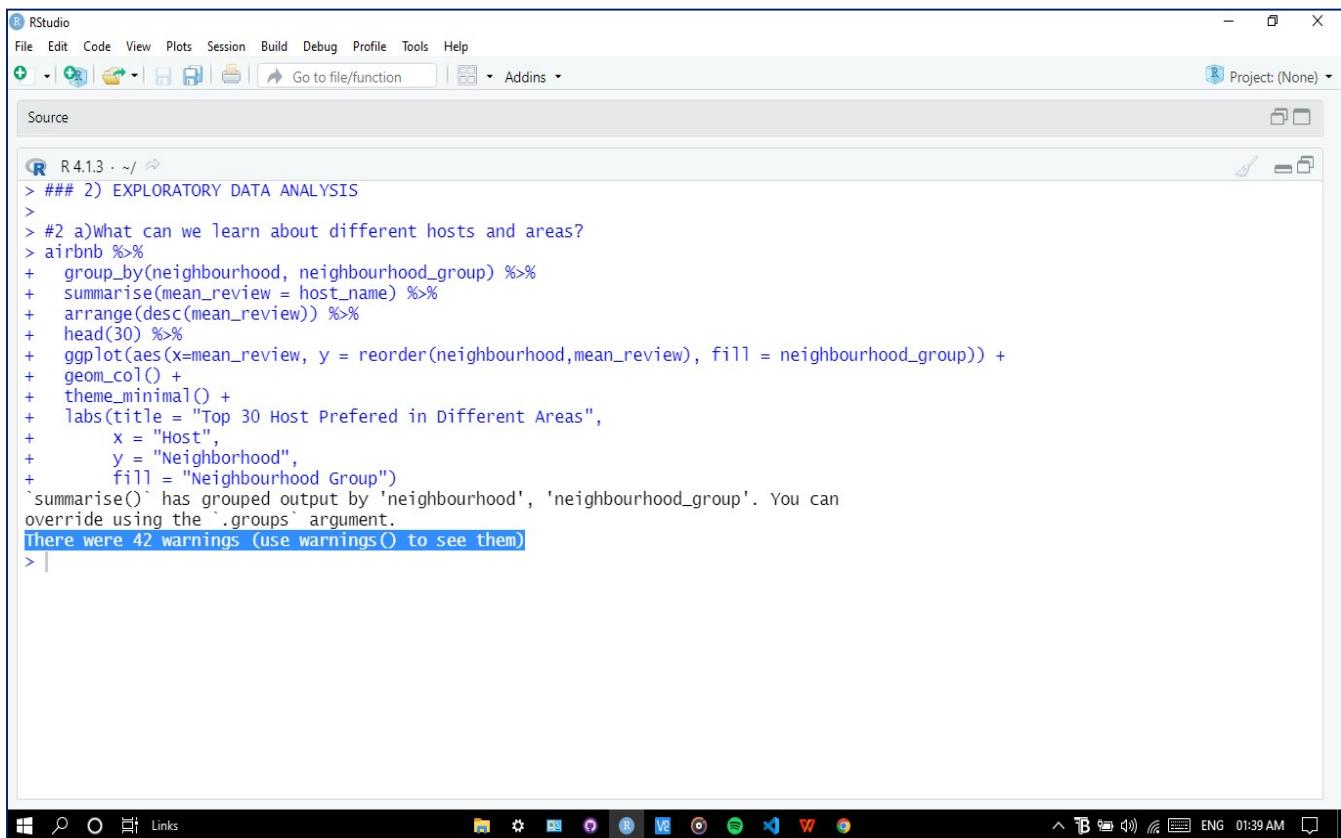
  

	availability_365
Min.:	0.0
1st Qu.:	0.0
Median :	45.0
Mean :	112.8

## 2) Exploratory Data Analysis :

### a) What can we learn about different hosts and areas?

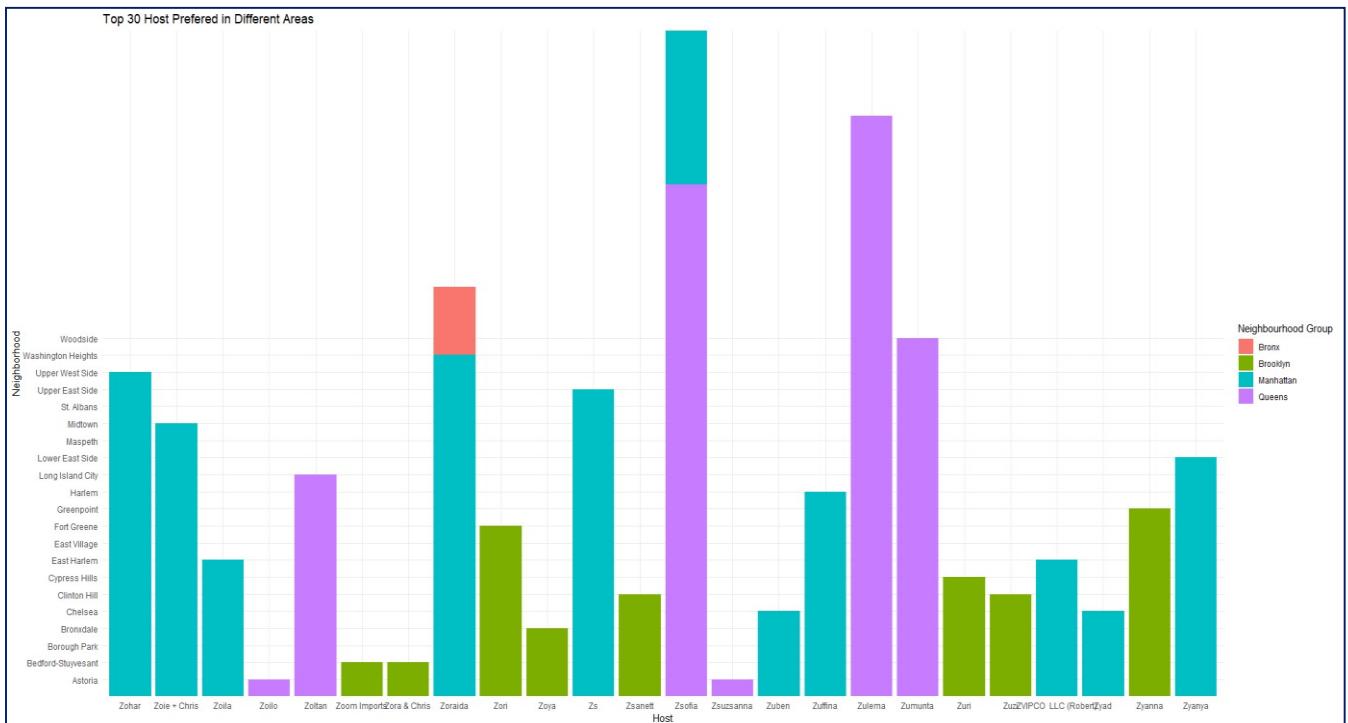
- The graph represents the different areas under different hosts.
- Hosts are stated on the x axis and different areas are shown on the y axis.
- Different colors in the graph represents the different neighborhood groups.
- As per the graph we could say that zsofia host has the more number of area which are queens and manhattan neighborhood groups.
- The Least number of area is under zsuzsanna has the least number of area which is of queens neighborhood groups.



The screenshot shows the RStudio interface with the following R code in the Source pane:

```
R 4.1.3 · ~/ 
> ##### 2) EXPLORATORY DATA ANALYSIS
>
> #2 a)What can we learn about different hosts and areas?
> airbnb %>%
+   group_by(neighbourhood, neighbourhood_group) %>%
+   summarise(mean_review = host_name) %>%
+   arrange(desc(mean_review)) %>%
+   head(30) %>%
+   ggplot(aes(x=mean_review, y = reorder(neighbourhood,mean_review), fill = neighbourhood_group)) +
+   geom_col() +
+   theme_minimal() +
+   labs(title = "Top 30 Host Preferred in Different Areas",
+       x = "Host",
+       y = "Neighborhood",
+       fill = "Neighbourhood Group")
`summarise()` has grouped output by 'neighbourhood', 'neighbourhood_group'. You can
override using the `.`groups` argument.
There were 42 warnings (use warnings() to see them)
> |
```

The RStudio interface includes a menu bar (File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help), a toolbar with various icons, and a status bar at the bottom showing system information like battery level, signal strength, and the time (01:39 AM).



## b) Which hosts are the busiest and why?

- Considering how many reviews it has, we can say the host is busiest as the number of reviews increases, therefore we can say the host is busiest as it is trying their best to manage the stuff.
- Here Maya, has 2273 reviews, so we can conclude that the host is busiest because she is giving her 100% to manage items.
- So, Maya is the most busiest host here!

```

> #2 b) Which hosts are the busiest and why?
> #We can check the host is busiest ,according to number of reviews
> airbnb |> group_by(host_name,neighbourhood,neighbourhood_group) |>
+   tally(number_of_reviews) |>
+   arrange(desc(n)) |>
+   head(10)
# A tibble: 10 x 4
# Groups:   host_name, neighbourhood [10]
  host_name      neighbourhood neighbourhood_group     n
  <chr>          <fct>           <fct>            <int>
1 Maya           East Elmhurst    Queens             2273
2 Brooklyn& Breakfast -Len- Prospect Heights Brooklyn          2205
3 Danielle        East Elmhurst    Queens             2017
4 Yasu & Akiko   Hell's Kitchen   Manhattan         1971
5 Brady          Bedford-Stuyvesant Brooklyn          1818
6 jj             Harlem           Manhattan         1798
7 Michael        Bedford-Stuyvesant Brooklyn          1667
8 John           Williamsburg    Brooklyn          1551
9 John           East Village    Manhattan         1484
10 Randy          Bedford-Stuyvesant Brooklyn          1379
>

```

### c) Is there any noticeable difference in traffic among different areas and what could be the reason for it?

- The graph represents the Top 50 Neighbourhood according to the Number of reviews
- Different colors in the graph represents the different neighborhood groups.
- Yes! There is a lot of difference in traffic among different areas as we can see in the bar plot there is huge traffic in QUEENS EAST ELMHURST as compared to others neighbourhood group.
- From bar plot we can clearly see that the number of reviews of QUEEN EAST ELMHURST is highest ,so this means the number of visits are more in QUEENS EAST ELMHURST so we could conclude that there is more traffic in QUEENS EAST ELMHURST neighbourhood group.
- So, according to number of reviews we can analyse traffic among different areas

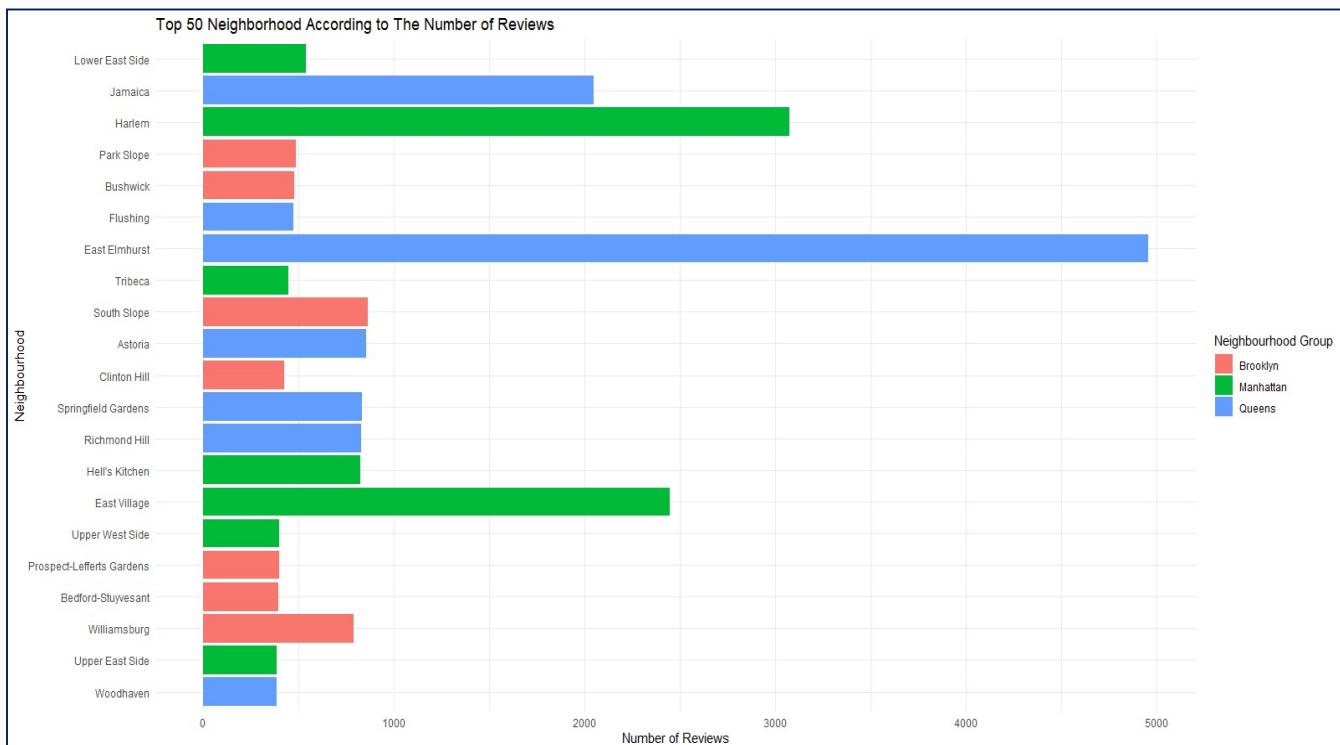
RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

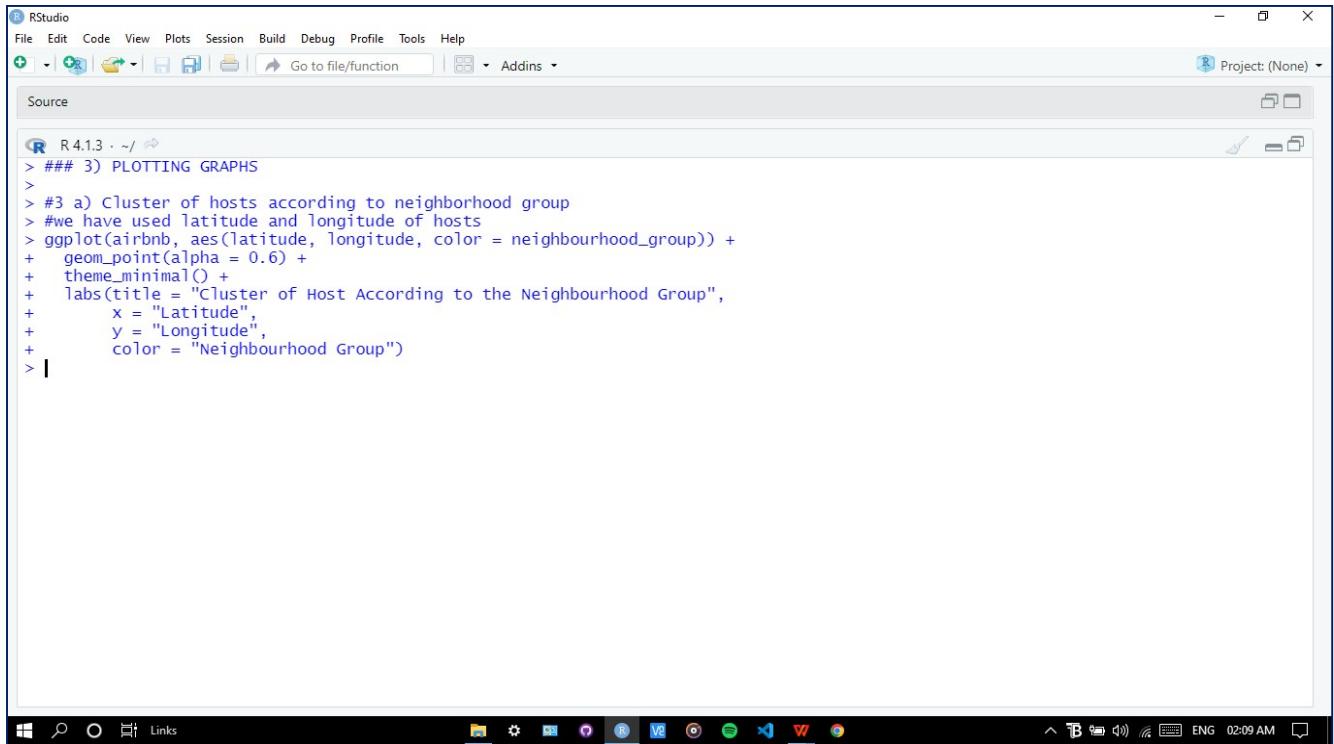
Source

```
R 4.1.3 · ~/ ◊
> # 2) Is there any noticeable difference in traffic among different areas and what could be the reason for it?
> airbnb %>%
+   group_by(neighbourhood, neighbourhood_group) %>%
+   summarise(review = number_of_reviews) %>%
+   arrange(desc(review)) %>%
+   head(50) %>%
+   ggpplot(aes(x=review, y = reorder(neighbourhood,review), fill = neighbourhood_group)) +
+   geom_col() +
+   theme_minimal() +
+   labs(title = "Top 50 Neighborhood According to The Number of Reviews",
+        x = "Number of Reviews",
+        y = "Neighbourhood",
+        fill = "Neighbourhood Group")
`summarise()` has grouped output by 'neighbourhood', 'neighbourhood_group'. You can override using the `groups` argument.
> |
```



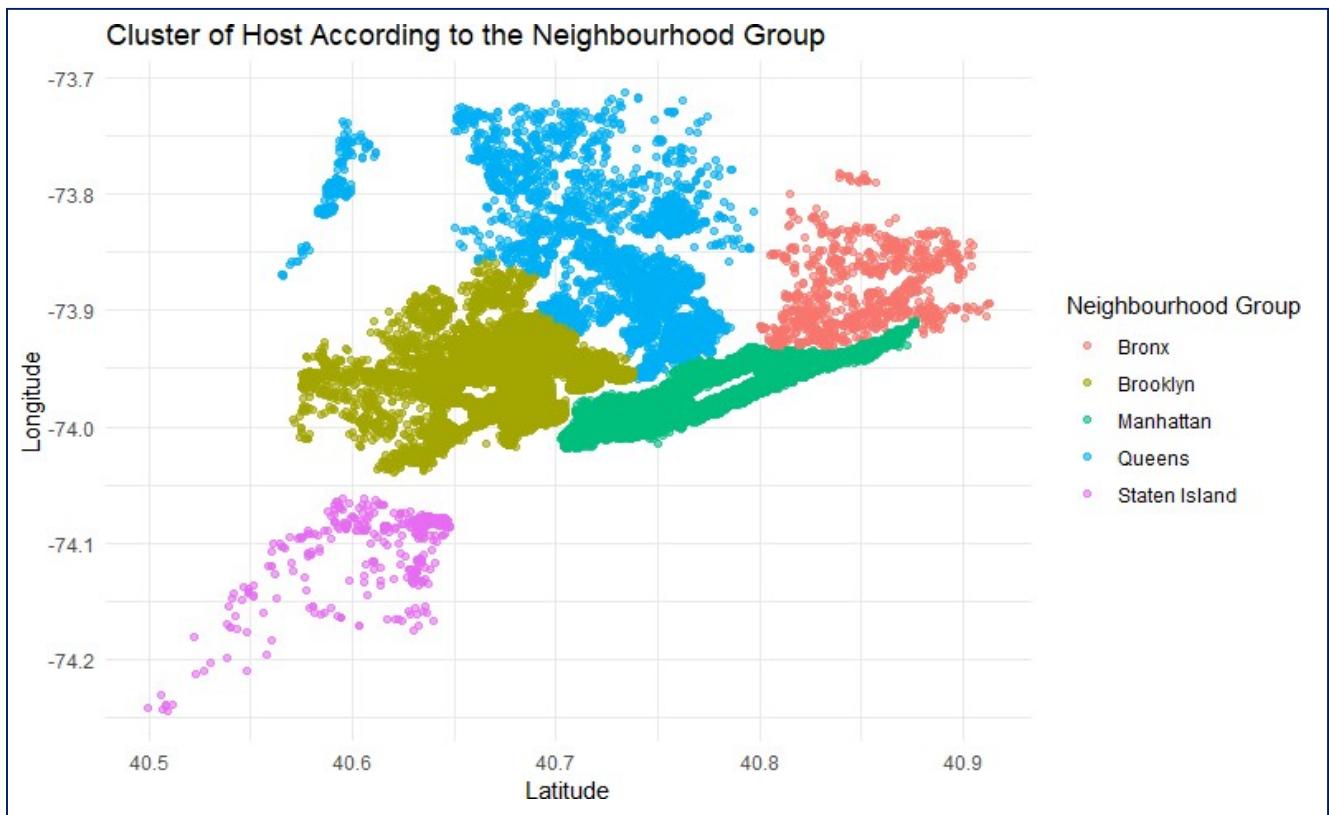
### 3) Plot Graphs Based On Some Conditions

#### a) Cluster of hosts according to neighborhood group :



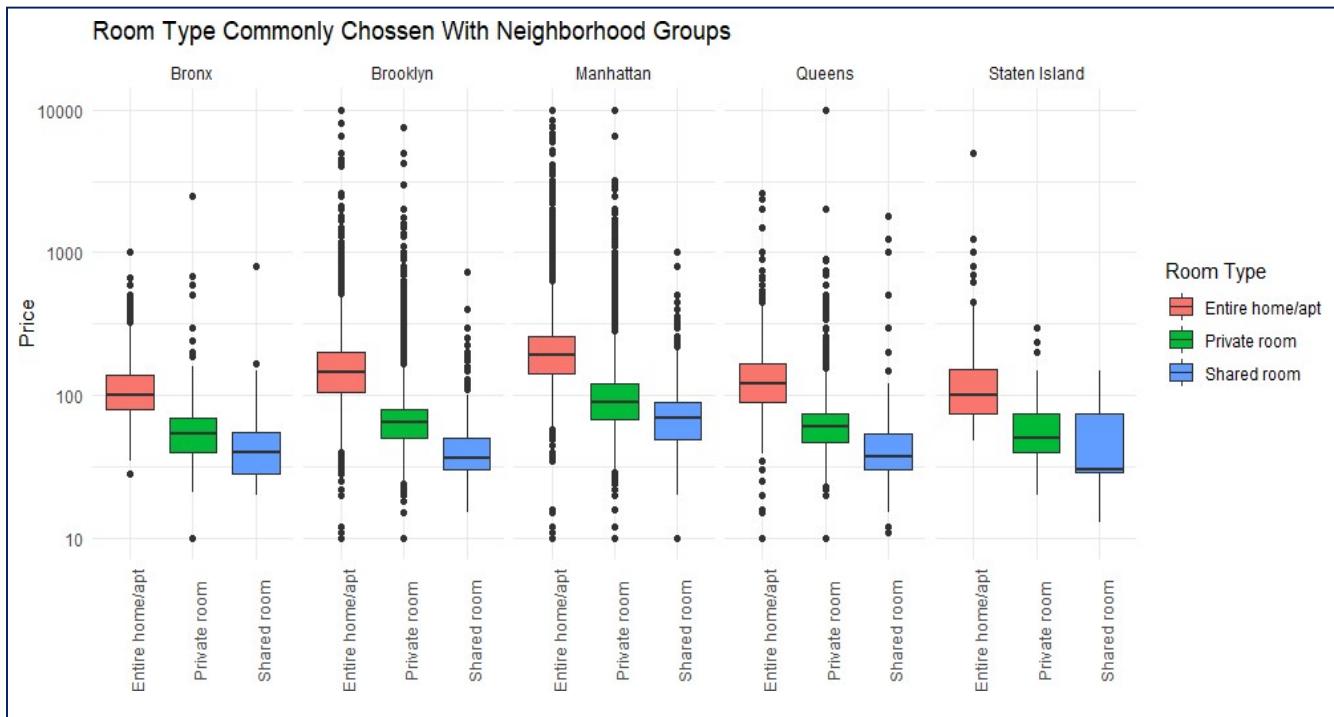
The screenshot shows the RStudio interface with the following code in the Source editor:

```
R 4.1.3 · ~/ ↗
> ### 3) PLOTTING GRAPHS
>
> #3 a) Cluster of hosts according to neighborhood group
> #we have used latitude and longitude of hosts
> ggplot(airbnb, aes(latitude, longitude, color = neighbourhood_group)) +
+   geom_point(alpha = 0.6) +
+   theme_minimal() +
+   labs(title = "Cluster of Host According to the Neighbourhood Group",
+       x = "Latitude",
+       y = "Longitude",
+       color = "Neighbourhood Group")
> |
```



**b) Plot a graph showing room\_type chosen commonly with neighborhood group :**

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Source
R 4.1.3 · ~/ ↘
> #3 b)Plot a graph showing room_type chosen commonly with neighborhood group.
> ggplot(airbnb, aes(x = room_type, y = price, fill = room_type)) + scale_y_log10() +
+   geom_boxplot() +
+   theme_minimal() +
+   labs (x="", y= "Price") +
+   facet_wrap(~neighbourhood_group) +
+   facet_grid(.~ neighbourhood_group) +
+   theme(axis.text.x = element_text(angle = 90), legend.position = "right") +
+   labs(title = "Room Type Commonly Chossen with Neighborhood Groups",
+       fill = "Room Type")
Warning messages:
1: Transformation introduced infinite values in continuous y-axis
2: Removed 11 rows containing non-finite values (stat_boxplot).
> |
```



### c) Plot average price with neighborhood :

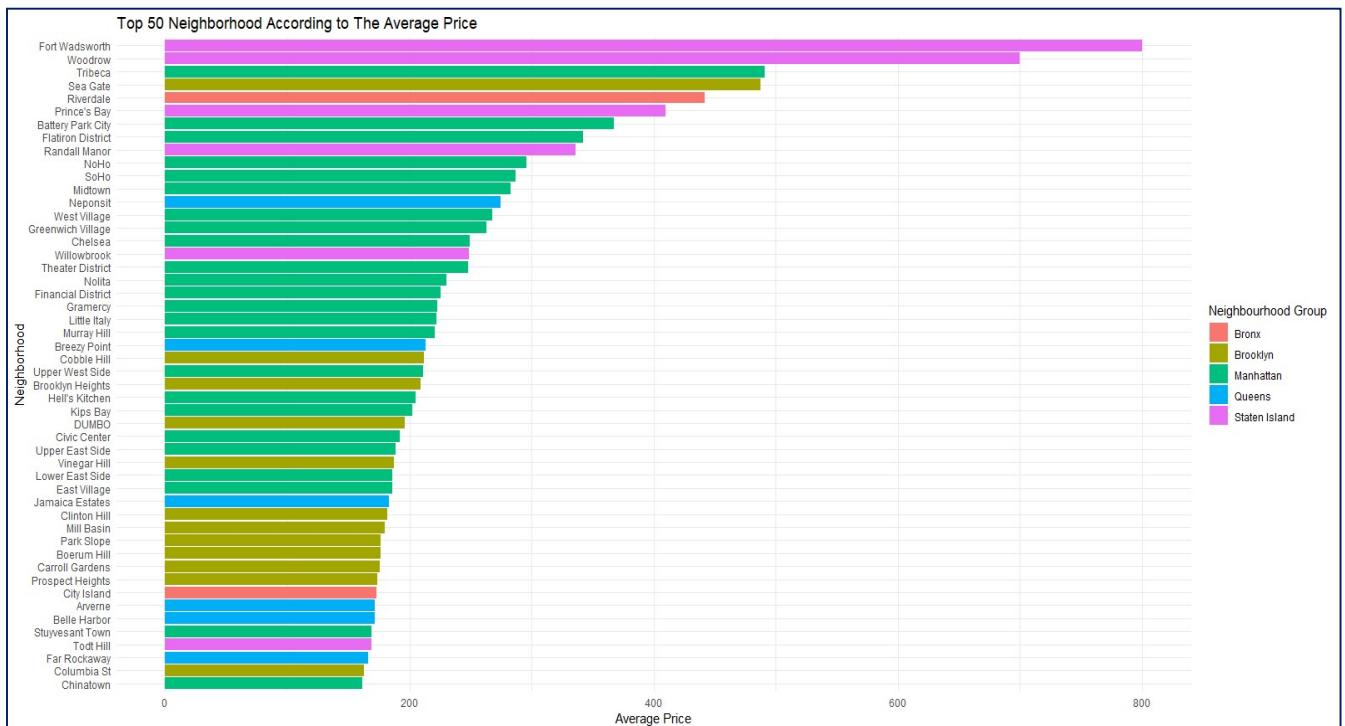
RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins Project: (None)

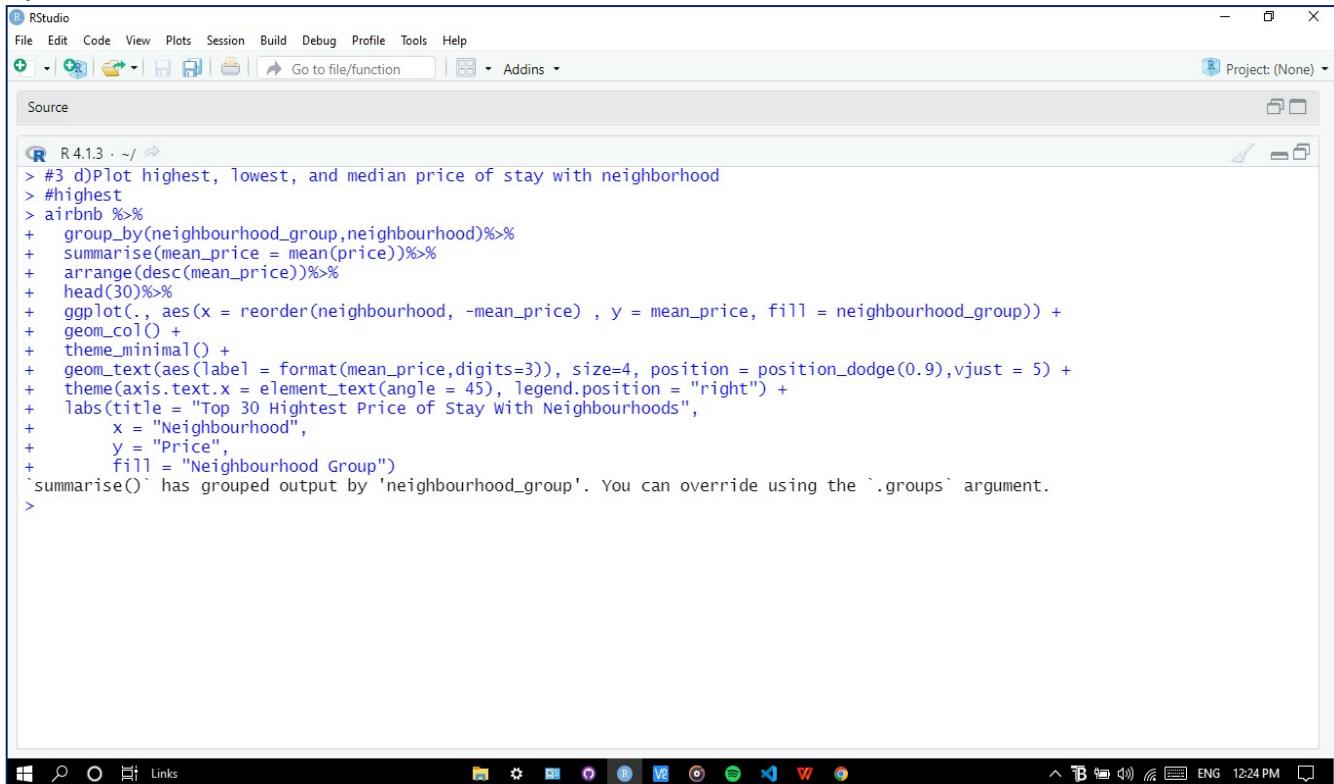
Source

```
R 4.1.3 - ~/ ~
> #3 c)Plot average price with neighborhood.
> airbnb %>%
+ group_by(neighbourhood, neighbourhood_group) %>%
+ summarise(mean_price = mean(price)) %>%
+ arrange(desc(mean_price)) %>%
+ head(50) %>%
+ ggplot(aes(x=mean_price, y = reorder(neighbourhood,mean_price)), fill = neighbourhood_group) +
+ geom_col() +
+ theme_minimal() +
+ labs(title = "Top 50 Neighborhood According to The Average Price",
+     x = "Average Price",
+     y = "Neighborhood",
+     fill = "Neighbourhood Group")
`summarise()` has grouped output by 'neighbourhood'. You can override using the `groups` argument.
> |
```

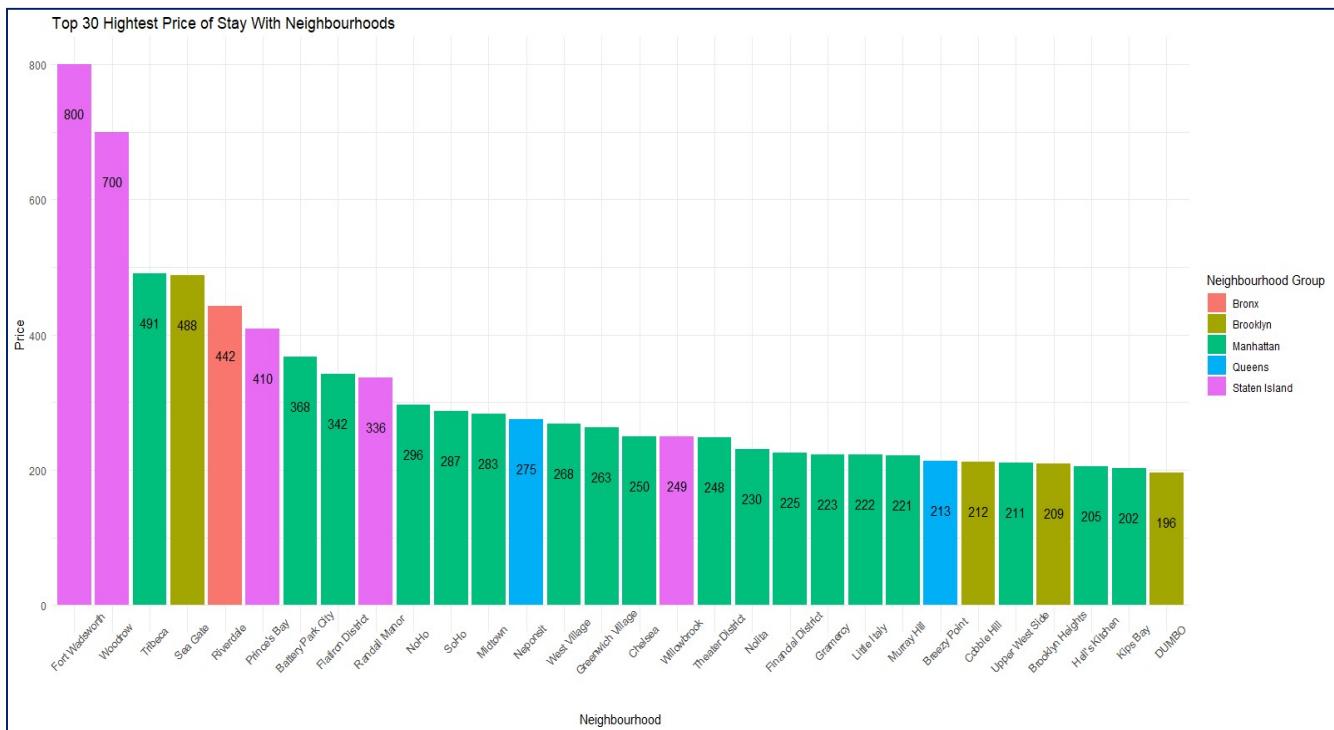


## d) Plot Highest, Lowest, and Median price of stay with neighborhood :

### 1) HIGHEST PRICE OF STAY WITH NEIGHBORHOOD :



```
R 4.1.3 · ~/ ◀
> #3 d)Plot highest, lowest, and median price of stay with neighborhood
> #highest
> airbnb %>%
+   group_by(neighbourhood_group,neighbourhood)%>
+   summarise(mean_price = mean(price))%>%
+   arrange(desc(mean_price))%>%
+   head(30)%>%
+   ggplot(., aes(x = reorder(neighbourhood, -mean_price) , y = mean_price, fill = neighbourhood_group)) +
+   geom_col() +
+   theme_minimal() +
+   geom_text(aes(label = format(mean_price,digits=3)), size=4, position = position_dodge(0.9),vjust = 5) +
+   theme(axis.text.x = element_text(angle = 45), legend.position = "right") +
+   labs(title = "Top 30 Highest Price of Stay With Neighbourhoods",
+       x = "Neighbourhood",
+       y = "Price",
+       fill = "Neighbourhood Group")
`summarise()` has grouped output by 'neighbourhood_group'. You can override using the `groups` argument.
>
```

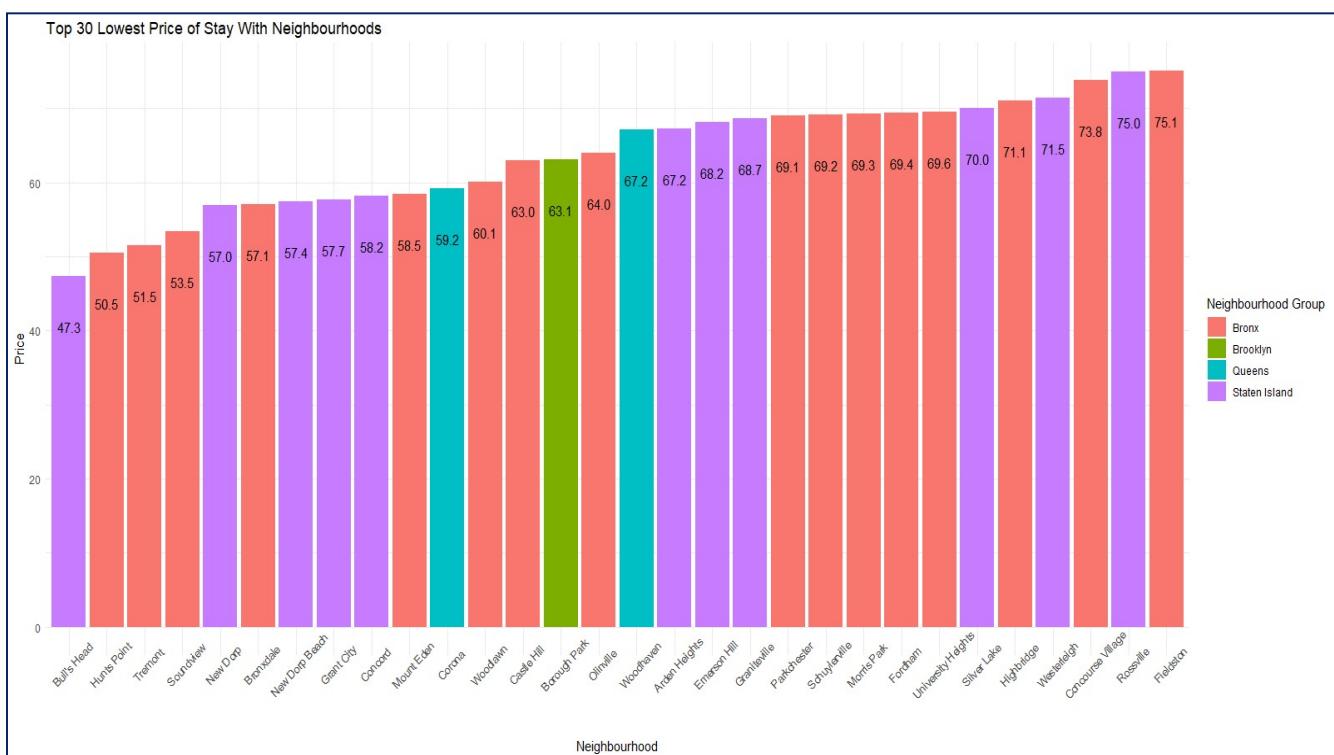


## 2) LOWEST PRICE OF STAY WITH NEIGHBORHOOD :

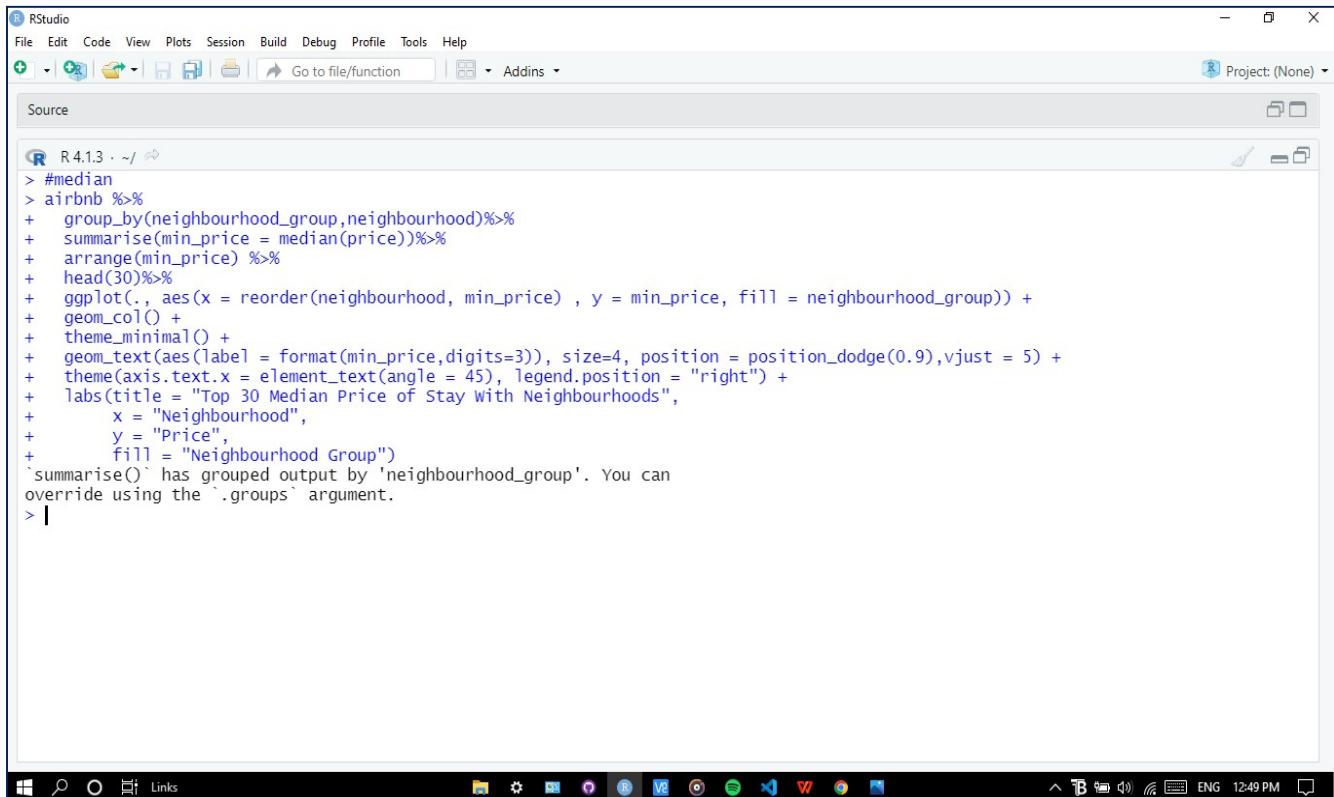
The screenshot shows the RStudio interface with the following R code in the Source pane:

```
> #lowest
> airbnb %>%
+   group_by(neighbourhood_group,neighbourhood)%>%
+   summarise(mean_price = mean(price))%>%
+   arrange(mean_price) %>%
+   head(30)%>%
+   ggplot(., aes(x = reorder(neighbourhood, mean_price) , y = mean_price, fill = neighbourhood_group)) +
+   geom_col() +
+   theme_minimal() +
+   geom_text(aes(label = format(mean_price,digits=3)), size=4, position = position_dodge(0.9),vjust = 5) +
+   theme(axis.text.x = element_text(angle = 45), legend.position = "right") +
+   labs(title = "Top 30 Lowest Price of Stay With Neighbourhoods",
+       x = "Neighbourhood",
+       y = "Price",
+       fill = "Neighbourhood Group")
`summarise()` has grouped output by 'neighbourhood_group'. You can override using the `groups` argument.
```

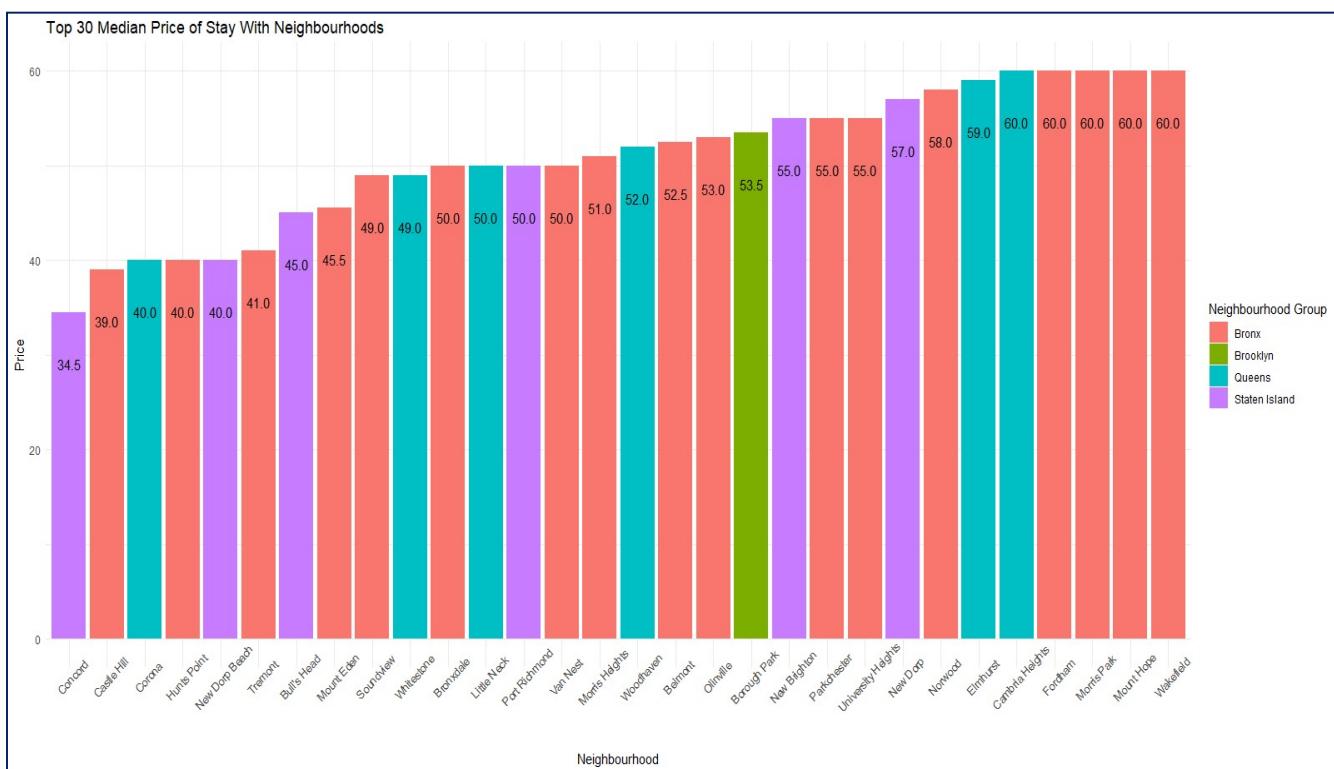
The warning message at the bottom of the code indicates that `summarise()` has grouped output by 'neighbourhood\_group' and suggests overriding using the `groups` argument.



### 3) MEDIAN PRICE OF STAY WITH NEIGHBORHOOD :



```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
+ Go to file/function Addins Project: (None)
Source
R 4.1.3 · ~/ ...
> #median
> airbnb %>%
+ group_by(neighbourhood_group,neighbourhood)%>%
+ summarise(min_price = median(price))%>%
+ arrange(min_price) %>%
+ head(30)%>%
+ ggplot(., aes(x = reorder(neighbourhood, min_price) , y = min_price, fill = neighbourhood_group)) +
+ geom_col() +
+ theme_minimal() +
+ geom_text(aes(label = format(min_price,digits=3)), size=4, position = position_dodge(0.9),vjust = 5) +
+ theme(axis.text.x = element_text(angle = 45), legend.position = "right") +
+ labs(title = "Top 30 Median Price of Stay With Neighbourhoods",
+     x = "Neighbourhood",
+     y = "Price",
+     fill = "Neighbourhood Group")
`summarise()` has grouped output by 'neighbourhood_group'. You can
override using the `groups` argument.
> |
```



---

---