

Code_EDA

Madhura Jadhav

```
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':  
##  
##   filter, lag
```

```
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union
```

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr   0.3.4  
## v tibble  3.1.4      v stringr 1.4.0  
## v tidyr   1.1.3      v forcats 0.5.1  
## v readr   2.1.2
```

```
## Warning: package 'readr' was built under R version 4.1.2
```

```
## Warning: package 'stringr' was built under R version 4.1.2
```

```
## -- Conflicts ----- tidyverse_conflicts() --  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()    masks stats::lag()
```

```
library(ggplot2)
```

```
library(skimr)
```

```
## Warning: package 'skimr' was built under R version 4.1.3
```

```
library(summarytools)
```

```
## Warning: package 'summarytools' was built under R version 4.1.3
```

```
##  
## Attaching package: 'summarytools'
```

```
## The following object is masked from 'package:tibble':  
##  
## view
```

```
library(lubridate)
```

```
## Warning: package 'lubridate' was built under R version 4.1.2
```

```
##  
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':  
##  
## date, intersect, setdiff, union
```

```
library(stringr)  
library(Hmisc)
```

```
## Warning: package 'Hmisc' was built under R version 4.1.3
```

```
## Loading required package: lattice
```

```
## Loading required package: survival
```

```
## Loading required package: Formula
```

```
##  
## Attaching package: 'Hmisc'
```

```
## The following objects are masked from 'package:summarytools':  
##  
## label, label<-
```

```
## The following objects are masked from 'package:dplyr':  
##  
## src, summarize
```

```
## The following objects are masked from 'package:base':  
##  
## format.pval, units
```

```
library(broom)
library(car)
```

```
## Warning: package 'car' was built under R version 4.1.3
```

```
## Loading required package: carData
```

```
## Warning: package 'carData' was built under R version 4.1.3
```

```
##
## Attaching package: 'car'
```

```
## The following object is masked from 'package:purrr':
##
##      some
```

```
## The following object is masked from 'package:dplyr':
##
##      recode
```

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.1.2
```

```
##
## Attaching package: 'caret'
```

```
## The following object is masked from 'package:survival':
##
##      cluster
```

```
## The following object is masked from 'package:purrr':
##
##      lift
```

```
library(rpart)
library(rpart.plot)
```

```
## Warning: package 'rpart.plot' was built under R version 4.1.3
```

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 4.1.3
```

```
## randomForest 4.7-1
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##  
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':  
##  
##     margin
```

```
## The following object is masked from 'package:dplyr':  
##  
##     combine
```

```
library(xgboost)
```

```
## Warning: package 'xgboost' was built under R version 4.1.3
```

```
##  
## Attaching package: 'xgboost'
```

```
## The following object is masked from 'package:dplyr':  
##  
##     slice
```

```
library(vtreat)
```

```
## Warning: package 'vtreat' was built under R version 4.1.3
```

```
## Loading required package: wrapr
```

```
## Warning: package 'wrapr' was built under R version 4.1.3
```

```
##  
## Attaching package: 'wrapr'
```

```
## The following object is masked from 'package:car':  
##  
##     bc
```

```
## The following object is masked from 'package:summarytools':  
##  
##     view
```

```
## The following objects are masked from 'package:tidyr':  
##  
##    pack, unpack
```

```
## The following object is masked from 'package:tibble':  
##  
##    view
```

```
## The following object is masked from 'package:dplyr':  
##  
##    coalesce
```

```
original_data = read.csv('analysisData.csv')  
dim(original_data)
```

```
## [1] 34404    91
```

Filtering the data to remove out of scope / redundant attributes

```
data=original_data[c('id','host_name', 'host_since', 'host_response_time','host_response_rate', 'host_is_superhost', 'host_total_listings_count', 'neighbourhood_cleansed', 'neighbourhood_group_cleansed', 'zipcode', 'property_type', 'room_type', 'accommodates', 'bathrooms', 'bedrooms', 'beds', 'amenities', 'price', 'cleaning_fee', 'guests_included', 'extra_people', 'minimum_nights', 'maximum_nights', 'availability_30', 'availability_60', 'availability_90', 'availability_365', 'number_of_reviews', 'number_of_reviews_ltm', 'first_review', 'last_review', 'review_scores_rating', 'review_scores_accuracy', 'review_scores_cleanliness', 'review_scores_checkin', 'review_scores_communication', 'review_scores_location', 'review_scores_value', 'instant_bookable', 'calculated_host_listings_count', 'reviews_per_month')]  
dim(data)
```

```
## [1] 34404    41
```

```
#Library(summarytools)  
print(dfSummary(data,style='grid',graph.col = T),method = 'render')
```

Data Frame Summary


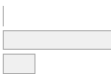



data

Dimensions: 34404 x 41

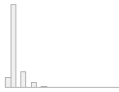


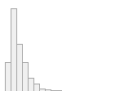
Duplicates: 0

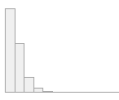

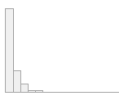



No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
----	----------	----------------	--------------------	-------	-------	---------

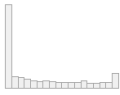
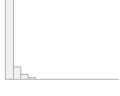
No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
1	id [integer]	<div>Mean (sd) : 445480 (210539.2)</div> <div>min ≤ med ≤ max: 100015 ≤ 437248 ≤ 898796</div> <div>IQR (CV) : 338110.8 (0.5)</div>	34404 distinct values		34404 (100.0%)	0 (0.0%)
2	host_name [character]	<div>1. Michael</div> <div>2. David</div> <div>3. John</div> <div>4. Alex</div> <div>5. Sonder (NYC)</div> <div>6. Sarah</div> <div>7. Maria</div> <div>8. Daniel</div> <div>9. Anna</div> <div>10. Mike</div> <div>[9246 others]</div>	<div>310 (0.9%)</div> <div>280 (0.8%)</div> <div>229 (0.7%)</div> <div>195 (0.6%)</div> <div>185 (0.5%)</div> <div>166 (0.5%)</div> <div>158 (0.5%)</div> <div>156 (0.5%)</div> <div>141 (0.4%)</div> <div>139 (0.4%)</div> <div>32445 (94.3%)</div>		34404 (100.0%)	0 (0.0%)
3	host_since [character]	<div>1. 2018-10-08</div> <div>2. 2017-06-27</div> <div>3. 2016-03-03</div> <div>4. 2015-12-16</div> <div>5. 2014-05-28</div> <div>6. 2017-03-07</div> <div>7. 2014-10-14</div> <div>8. 2018-05-22</div> <div>9. 2013-07-15</div> <div>10. 2017-03-14</div> <div>[3679 others]</div>	<div>201 (0.6%)</div> <div>81 (0.2%)</div> <div>74 (0.2%)</div> <div>69 (0.2%)</div> <div>64 (0.2%)</div> <div>49 (0.1%)</div> <div>48 (0.1%)</div> <div>48 (0.1%)</div> <div>47 (0.1%)</div> <div>46 (0.1%)</div> <div>33677 (97.9%)</div>		34404 (100.0%)	0 (0.0%)
4	host_response_time [character]	<div>1. (Empty string)</div> <div>2. a few days or more</div> <div>3. N/A</div> <div>4. within a day</div> <div>5. within a few hours</div> <div>6. within an hour</div>	<div>21 (0.1%)</div> <div>581 (1.7%)</div> <div>8680 (25.2%)</div> <div>3457 (10.0%)</div> <div>5530 (16.1%)</div> <div>16135 (46.9%)</div>		34404 (100.0%)	0 (0.0%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing																																	
5	host_response_rate [character]	<table><tr><td>1. 100%</td></tr><tr><td>2. N/A</td></tr><tr><td>3. 90%</td></tr><tr><td>4. 80%</td></tr><tr><td>5. 98%</td></tr><tr><td>6. 50%</td></tr><tr><td>7. 97%</td></tr><tr><td>8. 94%</td></tr><tr><td>9. 92%</td></tr><tr><td>10. 93%</td></tr><tr><td>[74 others]</td></tr></table>	1. 100%	2. N/A	3. 90%	4. 80%	5. 98%	6. 50%	7. 97%	8. 94%	9. 92%	10. 93%	[74 others]	<table><tr><td>18090</td><td>(52.6%)</td></tr><tr><td>8680</td><td>(25.2%)</td></tr><tr><td>1203</td><td>(3.5%)</td></tr><tr><td>598</td><td>(1.7%)</td></tr><tr><td>507</td><td>(1.5%)</td></tr><tr><td>322</td><td>(0.9%)</td></tr><tr><td>313</td><td>(0.9%)</td></tr><tr><td>280</td><td>(0.8%)</td></tr><tr><td>278</td><td>(0.8%)</td></tr><tr><td>271</td><td>(0.8%)</td></tr><tr><td>3862</td><td>(11.2%)</td></tr></table>	18090	(52.6%)	8680	(25.2%)	1203	(3.5%)	598	(1.7%)	507	(1.5%)	322	(0.9%)	313	(0.9%)	280	(0.8%)	278	(0.8%)	271	(0.8%)	3862	(11.2%)		34404 (100.0%)	0 (0.0%)
1. 100%																																							
2. N/A																																							
3. 90%																																							
4. 80%																																							
5. 98%																																							
6. 50%																																							
7. 97%																																							
8. 94%																																							
9. 92%																																							
10. 93%																																							
[74 others]																																							
18090	(52.6%)																																						
8680	(25.2%)																																						
1203	(3.5%)																																						
598	(1.7%)																																						
507	(1.5%)																																						
322	(0.9%)																																						
313	(0.9%)																																						
280	(0.8%)																																						
278	(0.8%)																																						
271	(0.8%)																																						
3862	(11.2%)																																						
6	host_is_superhost [character]	<table><tr><td>1. (Empty string)</td></tr><tr><td>2. f</td></tr><tr><td>3. t</td></tr></table>	1. (Empty string)	2. f	3. t	<table><tr><td>21</td><td>(0.1%)</td></tr><tr><td>26654</td><td>(77.5%)</td></tr><tr><td>7729</td><td>(22.5%)</td></tr></table>	21	(0.1%)	26654	(77.5%)	7729	(22.5%)		34404 (100.0%)	0 (0.0%)																								
1. (Empty string)																																							
2. f																																							
3. t																																							
21	(0.1%)																																						
26654	(77.5%)																																						
7729	(22.5%)																																						
7	host_total_listings_count [integer]	<table><tr><td>Mean (sd) : 7.4 (50.5)</td></tr><tr><td>min ≤ med ≤ max:</td></tr><tr><td>0 ≤ 1 ≤ 1080</td></tr><tr><td>IQR (CV) : 1 (6.8)</td></tr></table>	Mean (sd) : 7.4 (50.5)	min ≤ med ≤ max:	0 ≤ 1 ≤ 1080	IQR (CV) : 1 (6.8)	75 distinct values		34383 (99.9%)	21 (0.1%)																													
Mean (sd) : 7.4 (50.5)																																							
min ≤ med ≤ max:																																							
0 ≤ 1 ≤ 1080																																							
IQR (CV) : 1 (6.8)																																							
8	neighbourhood_cleansed [character]	<table><tr><td>1. Williamsburg</td></tr><tr><td>2. Bedford-Stuyvesant</td></tr><tr><td>3. Harlem</td></tr><tr><td>4. Bushwick</td></tr><tr><td>5. Hell's Kitchen</td></tr><tr><td>6. East Village</td></tr><tr><td>7. Upper West Side</td></tr><tr><td>8. Upper East Side</td></tr><tr><td>9. Crown Heights</td></tr><tr><td>10. East Harlem</td></tr><tr><td>[209 others]</td></tr></table>	1. Williamsburg	2. Bedford-Stuyvesant	3. Harlem	4. Bushwick	5. Hell's Kitchen	6. East Village	7. Upper West Side	8. Upper East Side	9. Crown Heights	10. East Harlem	[209 others]	<table><tr><td>2782</td><td>(8.1%)</td></tr><tr><td>2778</td><td>(8.1%)</td></tr><tr><td>1925</td><td>(5.6%)</td></tr><tr><td>1706</td><td>(5.0%)</td></tr><tr><td>1431</td><td>(4.2%)</td></tr><tr><td>1290</td><td>(3.7%)</td></tr><tr><td>1269</td><td>(3.7%)</td></tr><tr><td>1198</td><td>(3.5%)</td></tr><tr><td>1146</td><td>(3.3%)</td></tr><tr><td>837</td><td>(2.4%)</td></tr><tr><td>18042</td><td>(52.4%)</td></tr></table>	2782	(8.1%)	2778	(8.1%)	1925	(5.6%)	1706	(5.0%)	1431	(4.2%)	1290	(3.7%)	1269	(3.7%)	1198	(3.5%)	1146	(3.3%)	837	(2.4%)	18042	(52.4%)		34404 (100.0%)	0 (0.0%)
1. Williamsburg																																							
2. Bedford-Stuyvesant																																							
3. Harlem																																							
4. Bushwick																																							
5. Hell's Kitchen																																							
6. East Village																																							
7. Upper West Side																																							
8. Upper East Side																																							
9. Crown Heights																																							
10. East Harlem																																							
[209 others]																																							
2782	(8.1%)																																						
2778	(8.1%)																																						
1925	(5.6%)																																						
1706	(5.0%)																																						
1431	(4.2%)																																						
1290	(3.7%)																																						
1269	(3.7%)																																						
1198	(3.5%)																																						
1146	(3.3%)																																						
837	(2.4%)																																						
18042	(52.4%)																																						
9	neighbourhood_group_cleansed [character]	<table><tr><td>1. Bronx</td></tr><tr><td>2. Brooklyn</td></tr><tr><td>3. Manhattan</td></tr><tr><td>4. Queens</td></tr><tr><td>5. Staten Island</td></tr></table>	1. Bronx	2. Brooklyn	3. Manhattan	4. Queens	5. Staten Island	<table><tr><td>837</td><td>(2.4%)</td></tr><tr><td>14428</td><td>(41.9%)</td></tr><tr><td>14532</td><td>(42.2%)</td></tr><tr><td>4309</td><td>(12.5%)</td></tr><tr><td>298</td><td>(0.9%)</td></tr></table>	837	(2.4%)	14428	(41.9%)	14532	(42.2%)	4309	(12.5%)	298	(0.9%)		34404 (100.0%)	0 (0.0%)																		
1. Bronx																																							
2. Brooklyn																																							
3. Manhattan																																							
4. Queens																																							
5. Staten Island																																							
837	(2.4%)																																						
14428	(41.9%)																																						
14532	(42.2%)																																						
4309	(12.5%)																																						
298	(0.9%)																																						

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing																																																							
10	zipcode [character]	<table><tr><td>1. 11211</td></tr><tr><td>2. 11221</td></tr><tr><td>3. 11206</td></tr><tr><td>4. 11216</td></tr><tr><td>5. 10002</td></tr><tr><td>6. 10019</td></tr><tr><td>7. 10009</td></tr><tr><td>8. 11238</td></tr><tr><td>9. 11222</td></tr><tr><td>10. 11233</td></tr><tr><td>[191 others]</td></tr></table>	1. 11211	2. 11221	3. 11206	4. 11216	5. 10002	6. 10019	7. 10009	8. 11238	9. 11222	10. 11233	[191 others]	<table><tr><td>1556</td><td>(</td><td>4.5%</td><td>)</td></tr><tr><td>1348</td><td>(</td><td>3.9%</td><td>)</td></tr><tr><td>1073</td><td>(</td><td>3.1%</td><td>)</td></tr><tr><td>1023</td><td>(</td><td>3.0%</td><td>)</td></tr><tr><td>912</td><td>(</td><td>2.7%</td><td>)</td></tr><tr><td>852</td><td>(</td><td>2.5%</td><td>)</td></tr><tr><td>815</td><td>(</td><td>2.4%</td><td>)</td></tr><tr><td>796</td><td>(</td><td>2.3%</td><td>)</td></tr><tr><td>754</td><td>(</td><td>2.2%</td><td>)</td></tr><tr><td>728</td><td>(</td><td>2.1%</td><td>)</td></tr><tr><td>24547</td><td>(</td><td>71.3%</td><td>)</td></tr></table>	1556	(4.5%)	1348	(3.9%)	1073	(3.1%)	1023	(3.0%)	912	(2.7%)	852	(2.5%)	815	(2.4%)	796	(2.3%)	754	(2.2%)	728	(2.1%)	24547	(71.3%)		34404 (100.0%)	0 (0.0%)
1. 11211																																																													
2. 11221																																																													
3. 11206																																																													
4. 11216																																																													
5. 10002																																																													
6. 10019																																																													
7. 10009																																																													
8. 11238																																																													
9. 11222																																																													
10. 11233																																																													
[191 others]																																																													
1556	(4.5%)																																																										
1348	(3.9%)																																																										
1073	(3.1%)																																																										
1023	(3.0%)																																																										
912	(2.7%)																																																										
852	(2.5%)																																																										
815	(2.4%)																																																										
796	(2.3%)																																																										
754	(2.2%)																																																										
728	(2.1%)																																																										
24547	(71.3%)																																																										
11	property_type [character]	<table><tr><td>1. Apartment</td></tr><tr><td>2. House</td></tr><tr><td>3. Townhouse</td></tr><tr><td>4. Condominium</td></tr><tr><td>5. Loft</td></tr><tr><td>6. Guest suite</td></tr><tr><td>7. Serviced apartment</td></tr><tr><td>8. Boutique hotel</td></tr><tr><td>9. Hotel</td></tr><tr><td>10. Other</td></tr><tr><td>[21 others]</td></tr></table>	1. Apartment	2. House	3. Townhouse	4. Condominium	5. Loft	6. Guest suite	7. Serviced apartment	8. Boutique hotel	9. Hotel	10. Other	[21 others]	<table><tr><td>26942</td><td>(</td><td>78.3%</td><td>)</td></tr><tr><td>2935</td><td>(</td><td>8.5%</td><td>)</td></tr><tr><td>1250</td><td>(</td><td>3.6%</td><td>)</td></tr><tr><td>1075</td><td>(</td><td>3.1%</td><td>)</td></tr><tr><td>1048</td><td>(</td><td>3.0%</td><td>)</td></tr><tr><td>311</td><td>(</td><td>0.9%</td><td>)</td></tr><tr><td>267</td><td>(</td><td>0.8%</td><td>)</td></tr><tr><td>112</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>94</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>60</td><td>(</td><td>0.2%</td><td>)</td></tr><tr><td>310</td><td>(</td><td>0.9%</td><td>)</td></tr></table>	26942	(78.3%)	2935	(8.5%)	1250	(3.6%)	1075	(3.1%)	1048	(3.0%)	311	(0.9%)	267	(0.8%)	112	(0.3%)	94	(0.3%)	60	(0.2%)	310	(0.9%)		34404 (100.0%)	0 (0.0%)
1. Apartment																																																													
2. House																																																													
3. Townhouse																																																													
4. Condominium																																																													
5. Loft																																																													
6. Guest suite																																																													
7. Serviced apartment																																																													
8. Boutique hotel																																																													
9. Hotel																																																													
10. Other																																																													
[21 others]																																																													
26942	(78.3%)																																																										
2935	(8.5%)																																																										
1250	(3.6%)																																																										
1075	(3.1%)																																																										
1048	(3.0%)																																																										
311	(0.9%)																																																										
267	(0.8%)																																																										
112	(0.3%)																																																										
94	(0.3%)																																																										
60	(0.2%)																																																										
310	(0.9%)																																																										
12	room_type [character]	<table><tr><td>1. Entire home/apt</td></tr><tr><td>2. Hotel room</td></tr><tr><td>3. Private room</td></tr><tr><td>4. Shared room</td></tr></table>	1. Entire home/apt	2. Hotel room	3. Private room	4. Shared room	<table><tr><td>17859</td><td>(</td><td>51.9%</td><td>)</td></tr><tr><td>3</td><td>(</td><td>0.0%</td><td>)</td></tr><tr><td>15731</td><td>(</td><td>45.7%</td><td>)</td></tr><tr><td>811</td><td>(</td><td>2.4%</td><td>)</td></tr></table>	17859	(51.9%)	3	(0.0%)	15731	(45.7%)	811	(2.4%)		34404 (100.0%)	0 (0.0%)																																			
1. Entire home/apt																																																													
2. Hotel room																																																													
3. Private room																																																													
4. Shared room																																																													
17859	(51.9%)																																																										
3	(0.0%)																																																										
15731	(45.7%)																																																										
811	(2.4%)																																																										
13	accommodates [integer]	<table><tr><td>Mean (sd) : 2.9 (1.9)</td></tr><tr><td>min ≤ med ≤ max:</td></tr><tr><td>1 ≤ 2 ≤ 16</td></tr><tr><td>IQR (CV) : 2 (0.6)</td></tr></table>	Mean (sd) : 2.9 (1.9)	min ≤ med ≤ max:	1 ≤ 2 ≤ 16	IQR (CV) : 2 (0.6)	16 distinct values		34404 (100.0%)	0 (0.0%)																																																			
Mean (sd) : 2.9 (1.9)																																																													
min ≤ med ≤ max:																																																													
1 ≤ 2 ≤ 16																																																													
IQR (CV) : 2 (0.6)																																																													
14	bathrooms [numeric]	<table><tr><td>Mean (sd) : 1.1 (0.4)</td></tr><tr><td>min ≤ med ≤ max:</td></tr><tr><td>0 ≤ 1 ≤ 7</td></tr><tr><td>IQR (CV) : 0 (0.4)</td></tr></table>	Mean (sd) : 1.1 (0.4)	min ≤ med ≤ max:	0 ≤ 1 ≤ 7	IQR (CV) : 0 (0.4)	15 distinct values		34404 (100.0%)	0 (0.0%)																																																			
Mean (sd) : 1.1 (0.4)																																																													
min ≤ med ≤ max:																																																													
0 ≤ 1 ≤ 7																																																													
IQR (CV) : 0 (0.4)																																																													

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
15	bedrooms [integer]	<div>Mean (sd) : 1.2 (0.7)</div> <div>min ≤ med ≤ max:</div> <div>0 ≤ 1 ≤ 11</div> <div>IQR (CV) : 0 (0.6)</div>	12 distinct values		34404 (100.0%)	0 (0.0%)
16	beds [integer]	<div>Mean (sd) : 1.6 (1.1)</div> <div>min ≤ med ≤ max:</div> <div>0 ≤ 1 ≤ 26</div> <div>IQR (CV) : 1 (0.7)</div>	20 distinct values		34360 (99.9%)	44 (0.1%)
17	amenities [character]	<div>1. TV, Cable TV ,Wifi, Air c</div> <div>2. .</div> <div>3. TV, Cable TV ,Wifi, Air c</div> <div>4. TV, Cable TV ,Wifi, Air c</div> <div>5. TV,Wifi, Air conditioning</div> <div>6. TV,Wifi, Air conditioning</div> <div>7. TV, Cable TV ,Wifi, Air c</div> <div>8. TV,Wifi, Air conditioning</div> <div>9. Wifi, Air conditioning ,K</div> <div>10. TV, Cable TV ,Wifi, Air c</div> <div>[32280 others]</div>	<div>129 (0.4%)</div> <div>36 (0.1%)</div> <div>23 (0.1%)</div> <div>18 (0.1%)</div> <div>18 (0.1%)</div> <div>18 (0.1%)</div> <div>16 (0.0%)</div> <div>14 (0.0%)</div> <div>12 (0.0%)</div> <div>10 (0.0%)</div> <div>34110 (99.1%)</div>		34404 (100.0%)	0 (0.0%)
18	price [integer]	<div>Mean (sd) : 135.1 (106)</div> <div>min ≤ med ≤ max:</div> <div>0 ≤ 100 ≤ 999</div> <div>IQR (CV) : 102 (0.8)</div>	517 distinct values		34404 (100.0%)	0 (0.0%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
19	cleaning_fee [integer]	Mean (sd) : 62.2 (50.4) min ≤ med ≤ max: 0 ≤ 50 ≤ 600 IQR (CV) : 65 (0.8)	186 distinct values		29128 (84.7%)	5276 (15.3%)
20	guests_included [integer]	Mean (sd) : 1.6 (1.2) min ≤ med ≤ max: 1 ≤ 1 ≤ 16 IQR (CV) : 1 (0.8)	15 distinct values		34404 (100.0%)	0 (0.0%)
21	extra_people [integer]	Mean (sd) : 16.1 (24.7) min ≤ med ≤ max: 0 ≤ 10 ≤ 300 IQR (CV) : 25 (1.5)	104 distinct values		34404 (100.0%)	0 (0.0%)
22	minimum_nights [integer]	Mean (sd) : 5.8 (17.7) min ≤ med ≤ max: 1 ≤ 2 ≤ 1250 IQR (CV) : 3 (3)	88 distinct values		34404 (100.0%)	0 (0.0%)
23	maximum_nights [integer]	Mean (sd) : 63620.1 (11578280) min ≤ med ≤ max: 1 ≤ 365 ≤ 2147483647 IQR (CV) : 1096 (182)	248 distinct values		34404 (100.0%)	0 (0.0%)
24	availability_30 [integer]	Mean (sd) : 7.3 (9.4) min ≤ med ≤ max: 0 ≤ 3 ≤ 30 IQR (CV) : 12 (1.3)	31 distinct values		34404 (100.0%)	0 (0.0%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
25	availability_60 [integer]	Mean (sd) : 16.6 (19.7) min ≤ med ≤ max: 0 ≤ 8 ≤ 60 IQR (CV) : 30 (1.2)	61 distinct values		34404 (100.0%)	0 (0.0%)
26	availability_90 [integer]	Mean (sd) : 28.2 (30.9) min ≤ med ≤ max: 0 ≤ 15 ≤ 90 IQR (CV) : 54 (1.1)	91 distinct values		34404 (100.0%)	0 (0.0%)
27	availability_365 [integer]	Mean (sd) : 119 (130.5) min ≤ med ≤ max: 0 ≤ 62 ≤ 365 IQR (CV) : 239 (1.1)	366 distinct values		34404 (100.0%)	0 (0.0%)
28	number_of_reviews [integer]	Mean (sd) : 28 (47.3) min ≤ med ≤ max: 0 ≤ 9 ≤ 604 IQR (CV) : 28 (1.7)	376 distinct values		34404 (100.0%)	0 (0.0%)
29	number_of_reviews_ltm [integer]	Mean (sd) : 10.9 (16.7) min ≤ med ≤ max: 0 ≤ 3 ≤ 283 IQR (CV) : 13.2 (1.5)	144 distinct values		34404 (100.0%)	0 (0.0%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing																																																							
30	first_review [character]	<table><tr><td>1. 2019-01-01</td></tr><tr><td>2. 2018-01-01</td></tr><tr><td>3. 2016-01-02</td></tr><tr><td>4. 2019-06-30</td></tr><tr><td>5. 2018-01-02</td></tr><tr><td>6. 2019-01-02</td></tr><tr><td>7. 2019-07-01</td></tr><tr><td>8. 2019-04-21</td></tr><tr><td>9. 2017-01-02</td></tr><tr><td>10. 2017-01-01</td></tr><tr><td>[2967 others]</td></tr></table>	1. 2019-01-01	2. 2018-01-01	3. 2016-01-02	4. 2019-06-30	5. 2018-01-02	6. 2019-01-02	7. 2019-07-01	8. 2019-04-21	9. 2017-01-02	10. 2017-01-01	[2967 others]	<table><tr><td>150</td><td>(</td><td>0.4%</td><td>)</td></tr><tr><td>141</td><td>(</td><td>0.4%</td><td>)</td></tr><tr><td>116</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>112</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>104</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>100</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>95</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>93</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>92</td><td>(</td><td>0.3%</td><td>)</td></tr><tr><td>84</td><td>(</td><td>0.2%</td><td>)</td></tr><tr><td>33317</td><td>(</td><td>96.8%</td><td>)</td></tr></table>	150	(0.4%)	141	(0.4%)	116	(0.3%)	112	(0.3%)	104	(0.3%)	100	(0.3%)	95	(0.3%)	93	(0.3%)	92	(0.3%)	84	(0.2%)	33317	(96.8%)		34404 (100.0%)	0 (0.0%)
1. 2019-01-01																																																													
2. 2018-01-01																																																													
3. 2016-01-02																																																													
4. 2019-06-30																																																													
5. 2018-01-02																																																													
6. 2019-01-02																																																													
7. 2019-07-01																																																													
8. 2019-04-21																																																													
9. 2017-01-02																																																													
10. 2017-01-01																																																													
[2967 others]																																																													
150	(0.4%)																																																										
141	(0.4%)																																																										
116	(0.3%)																																																										
112	(0.3%)																																																										
104	(0.3%)																																																										
100	(0.3%)																																																										
95	(0.3%)																																																										
93	(0.3%)																																																										
92	(0.3%)																																																										
84	(0.2%)																																																										
33317	(96.8%)																																																										
31	last_review [character]	<table><tr><td>1. 2019-07-21</td></tr><tr><td>2. 2019-07-28</td></tr><tr><td>3. 2019-08-04</td></tr><tr><td>4. 2019-07-22</td></tr><tr><td>5. 2019-07-31</td></tr><tr><td>6. 2019-06-30</td></tr><tr><td>7. 2019-07-29</td></tr><tr><td>8. 2019-08-03</td></tr><tr><td>9. 2019-07-30</td></tr><tr><td>10. 2019-07-14</td></tr><tr><td>[1849 others]</td></tr></table>	1. 2019-07-21	2. 2019-07-28	3. 2019-08-04	4. 2019-07-22	5. 2019-07-31	6. 2019-06-30	7. 2019-07-29	8. 2019-08-03	9. 2019-07-30	10. 2019-07-14	[1849 others]	<table><tr><td>928</td><td>(</td><td>2.7%</td><td>)</td></tr><tr><td>822</td><td>(</td><td>2.4%</td><td>)</td></tr><tr><td>760</td><td>(</td><td>2.2%</td><td>)</td></tr><tr><td>633</td><td>(</td><td>1.8%</td><td>)</td></tr><tr><td>567</td><td>(</td><td>1.6%</td><td>)</td></tr><tr><td>521</td><td>(</td><td>1.5%</td><td>)</td></tr><tr><td>510</td><td>(</td><td>1.5%</td><td>)</td></tr><tr><td>465</td><td>(</td><td>1.4%</td><td>)</td></tr><tr><td>458</td><td>(</td><td>1.3%</td><td>)</td></tr><tr><td>452</td><td>(</td><td>1.3%</td><td>)</td></tr><tr><td>28288</td><td>(</td><td>82.2%</td><td>)</td></tr></table>	928	(2.7%)	822	(2.4%)	760	(2.2%)	633	(1.8%)	567	(1.6%)	521	(1.5%)	510	(1.5%)	465	(1.4%)	458	(1.3%)	452	(1.3%)	28288	(82.2%)		34404 (100.0%)	0 (0.0%)
1. 2019-07-21																																																													
2. 2019-07-28																																																													
3. 2019-08-04																																																													
4. 2019-07-22																																																													
5. 2019-07-31																																																													
6. 2019-06-30																																																													
7. 2019-07-29																																																													
8. 2019-08-03																																																													
9. 2019-07-30																																																													
10. 2019-07-14																																																													
[1849 others]																																																													
928	(2.7%)																																																										
822	(2.4%)																																																										
760	(2.2%)																																																										
633	(1.8%)																																																										
567	(1.6%)																																																										
521	(1.5%)																																																										
510	(1.5%)																																																										
465	(1.4%)																																																										
458	(1.3%)																																																										
452	(1.3%)																																																										
28288	(82.2%)																																																										
32	review_scores_rating [integer]	<table><tr><td>Mean (sd) : 93.9 (9.1)</td></tr><tr><td>min ≤ med ≤ max:</td></tr><tr><td>20 ≤ 97 ≤ 100</td></tr><tr><td>IQR (CV) : 8 (0.1)</td></tr></table>	Mean (sd) : 93.9 (9.1)	min ≤ med ≤ max:	20 ≤ 97 ≤ 100	IQR (CV) : 8 (0.1)	55 distinct values		34404 (100.0%)	0 (0.0%)																																																			
Mean (sd) : 93.9 (9.1)																																																													
min ≤ med ≤ max:																																																													
20 ≤ 97 ≤ 100																																																													
IQR (CV) : 8 (0.1)																																																													
33	review_scores_accuracy [integer]	<table><tr><td>Mean (sd) : 9.6 (0.9)</td></tr><tr><td>min ≤ med ≤ max:</td></tr><tr><td>2 ≤ 10 ≤ 10</td></tr><tr><td>IQR (CV) : 1 (0.1)</td></tr></table>	Mean (sd) : 9.6 (0.9)	min ≤ med ≤ max:	2 ≤ 10 ≤ 10	IQR (CV) : 1 (0.1)	<table><tr><td>2 :</td><td>133</td><td>(</td><td>0.4%</td><td>)</td></tr><tr><td>3 :</td><td>1</td><td>(</td><td>0.0%</td><td>)</td></tr><tr><td>4 :</td><td>79</td><td>(</td><td>0.2%</td><td>)</td></tr><tr><td>5 :</td><td>26</td><td>(</td><td>0.1%</td><td>)</td></tr><tr><td>6 :</td><td>321</td><td>(</td><td>0.9%</td><td>)</td></tr><tr><td>7 :</td><td>275</td><td>(</td><td>0.8%</td><td>)</td></tr><tr><td>8 :</td><td>1535</td><td>(</td><td>4.5%</td><td>)</td></tr><tr><td>9 :</td><td>6607</td><td>(</td><td>19.2%</td><td>)</td></tr><tr><td>10 :</td><td>25427</td><td>(</td><td>73.9%</td><td>)</td></tr></table>	2 :	133	(0.4%)	3 :	1	(0.0%)	4 :	79	(0.2%)	5 :	26	(0.1%)	6 :	321	(0.9%)	7 :	275	(0.8%)	8 :	1535	(4.5%)	9 :	6607	(19.2%)	10 :	25427	(73.9%)		34404 (100.0%)	0 (0.0%)						
Mean (sd) : 9.6 (0.9)																																																													
min ≤ med ≤ max:																																																													
2 ≤ 10 ≤ 10																																																													
IQR (CV) : 1 (0.1)																																																													
2 :	133	(0.4%)																																																									
3 :	1	(0.0%)																																																									
4 :	79	(0.2%)																																																									
5 :	26	(0.1%)																																																									
6 :	321	(0.9%)																																																									
7 :	275	(0.8%)																																																									
8 :	1535	(4.5%)																																																									
9 :	6607	(19.2%)																																																									
10 :	25427	(73.9%)																																																									

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
34	review_scores_cleanliness [integer]	<div>Mean (sd) : 9.3 (1.1)</div> <div>min ≤ med ≤ max:</div> <div>2 ≤ 10 ≤ 10</div> <div>IQR (CV) : 1 (0.1)</div>	2 : 168 (0.5%) 3 : 10 (0.0%) 4 : 166 (0.5%) 5 : 93 (0.3%) 6 : 621 (1.8%) 7 : 829 (2.4%) 8 : 3533 (10.3%) 9 : 10028 (29.1%) 10 : 18956 (55.1%)		34404 (100.0%)	0 (0.0%)
35	review_scores_checkin [integer]	<div>Mean (sd) : 9.7 (0.8)</div> <div>min ≤ med ≤ max:</div> <div>2 ≤ 10 ≤ 10</div> <div>IQR (CV) : 0 (0.1)</div>	2 : 102 (0.3%) 3 : 1 (0.0%) 4 : 50 (0.1%) 5 : 18 (0.1%) 6 : 227 (0.7%) 7 : 176 (0.5%) 8 : 929 (2.7%) 9 : 4449 (12.9%) 10 : 28452 (82.7%)		34404 (100.0%)	0 (0.0%)
36	review_scores_communication [integer]	<div>Mean (sd) : 9.7 (0.8)</div> <div>min ≤ med ≤ max:</div> <div>2 ≤ 10 ≤ 10</div> <div>IQR (CV) : 0 (0.1)</div>	2 : 112 (0.3%) 3 : 5 (0.0%) 4 : 49 (0.1%) 5 : 21 (0.1%) 6 : 214 (0.6%) 7 : 192 (0.6%) 8 : 952 (2.8%) 9 : 4047 (11.8%) 10 : 28812 (83.7%)		34404 (100.0%)	0 (0.0%)
37	review_scores_location [integer]	<div>Mean (sd) : 9.6 (0.8)</div> <div>min ≤ med ≤ max:</div> <div>2 ≤ 10 ≤ 10</div> <div>IQR (CV) : 1 (0.1)</div>	2 : 68 (0.2%) 3 : 2 (0.0%) 4 : 40 (0.1%) 5 : 10 (0.0%) 6 : 253 (0.7%) 7 : 206 (0.6%) 8 : 1822 (5.3%) 9 : 8719 (25.3%) 10 : 23284 (67.7%)		34404 (100.0%)	0 (0.0%)
38	review_scores_value [integer]	<div>Mean (sd) : 9.4 (1)</div> <div>min ≤ med ≤ max:</div> <div>2 ≤ 10 ≤ 10</div> <div>IQR (CV) : 1 (0.1)</div>	2 : 111 (0.3%) 3 : 6 (0.0%) 4 : 108 (0.3%) 5 : 55 (0.2%) 6 : 409 (1.2%) 7 : 373 (1.1%) 8 : 2474 (7.2%) 9 : 11229 (32.6%) 10 : 19639 (57.1%)		34404 (100.0%)	0 (0.0%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Valid	Missing
39	instant_bookable [character]	1. f 2. t	21020 (61.1%) 13384 (38.9%)		34404 (100.0%)	0 (0.0%)
40	calculated_host_listings_count [integer]	Mean (sd) : 5.2 (27.5) min ≤ med ≤ max: 1 ≤ 1 ≤ 371 IQR (CV) : 1 (5.3)	52 distinct values		34404 (100.0%)	0 (0.0%)
41	reviews_per_month [numeric]	Mean (sd) : 1.4 (1.7) min ≤ med ≤ max: 0 ≤ 0.8 ≤ 66.6 IQR (CV) : 1.8 (1.2)	924 distinct values		34402 (100.0%)	2 (0.0%)

Generated by summarytools (<https://github.com/dcomtois/summarytools>) 1.0.0 (R (<https://www.r-project.org/>) version 4.1.1)
2022-05-02

Missing Values

Missing Values Columns: cleaning_fee 5276, beds 44, host_total_listings_count 21, reviews_per_month 2, host_response_time 8679, host_response_rate 8679

Will drop rows with missing values <50 for a column

```
data = data %>% drop_na(reviews_per_month)
data = data %>% drop_na(host_total_listings_count)
data = data %>% drop_na(beds)
```

Cleaning fees is option so absence of it is cleaning fee=0

```
data = data %>% mutate(cleaning_fee = ifelse(is.na(data$cleaning_fee), 0, data$cleaning_fee))
```

We will drop host_response_time 8679, host_response_rate 8679 as we have ~25% missing data

```
data = data[ , -which(names(data) %in% c("host_response_rate", "host_response_time"))]
```

Data Transformations

Converting columns with logical data type to 1 and 0

host_is_superhost, instant_bookable

```
data$host_is_superhost = as.numeric(ifelse(data$host_is_superhost == 't', 1, 0))
data$instant_bookable = as.numeric(ifelse(data$instant_bookable == 't', 1, 0))
```

New Features

host_since can be used to calculate host active duration in months

first_review & last_review can be used to calculate listing active duration in months

we can drop column last_review, first_review, host_since after calculating new features

```
#max(data$last_review) is '2020-01-03'

data$listing_active_duration = round(as.numeric(difftime(data$last_review, data$first_review, units =
"days"))/(365.25/12),2)
data$host_active_duration = round(as.numeric(difftime('2020-01-03', data$host_since, units ="days"))/(3
65.25/12),2)

data = subset(data, select = -c(host_since,first_review,last_review) )
```

Amenities can be extracted to check important individual amenity

```
#Library(stringr)
data$TV = str_extract(data$amenities, "TV")
data$TV[is.na(data$TV)] <- 0
data$TV[data$TV == "TV"] <- 1

data$Elevator = str_extract(data$amenities, "Elevator")
data$Elevator[is.na(data$Elevator)] <- 0
data$Elevator[data$Elevator == "Elevator"] <- 1
```

Combine variable levels

property_type

combine levels with <1% data to 'Others' making total number of levels for property_type as 6

```
data %>% group_by(property_type) %>% summarise(n=n()) %>% arrange(desc(n))
```

```
## # A tibble: 31 x 2
##   property_type      n
##   <chr>          <int>
## 1 Apartment      26894
## 2 House          2924
## 3 Townhouse     1249
## 4 Condominium   1073
## 5 Loft          1048
## 6 Guest suite    311
## 7 Serviced apartment 266
## 8 Boutique hotel  111
## 9 Hotel          94
## 10 Other         60
## # ... with 21 more rows
```

```
data=data %>% mutate(property_type = fct_lump(property_type, prop = 0.01))
```

room_type

updating 'Hotel room' level as 'Private room' as we have just 3 'Hotel room' records

```
data$room_type = fct_collapse(data$room_type, 'Private room' = c('Hotel room', 'Private room'))
```

bathrooms

updating 0.5 bathroom levels and reducing number of levels

```
data$bathrooms[data$bathrooms < 1] <- 1  
data$bathrooms[ data$bathrooms == 1.5] <- 2  
data$bathrooms[ data$bathrooms == 2.5] <- 3  
data$bathrooms[ data$bathrooms == 3.5] <- 4  
data$bathrooms[ data$bathrooms > 4] <- 4
```

bedrooms

reducing number of levels from 12 to 7

```
data$bedrooms[data$bedrooms > 6] <- 6
```

Delete erroneous data

price

delete listings with price=0, assuming it is data error

```
data= data[!(data$price==0),]
```

zipcode

clean dirty zipcode data, treat NAs with mode value

```
data$zipcode=str_replace_all(data$zipcode,"NY ","")  
data$zipcode=str_replace_all(data$zipcode," ","")  
data$zipcode[data$zipcode=="11385-2308"] = '11385'  
data$zipcode[data$zipcode=="11103-3233"] = '11103'  
data$zipcode[data$zipcode=="10003-8623"] = '10003'  
data$zipcode[data$zipcode=="11413-3220"] = '11413'  
data$zipcode[data$zipcode=="10065"] = '10021'  
data$zipcode[data$zipcode=="11249"] = '11211'  
data$zipcode <- gsub('NA','11211',data$zipcode)
```



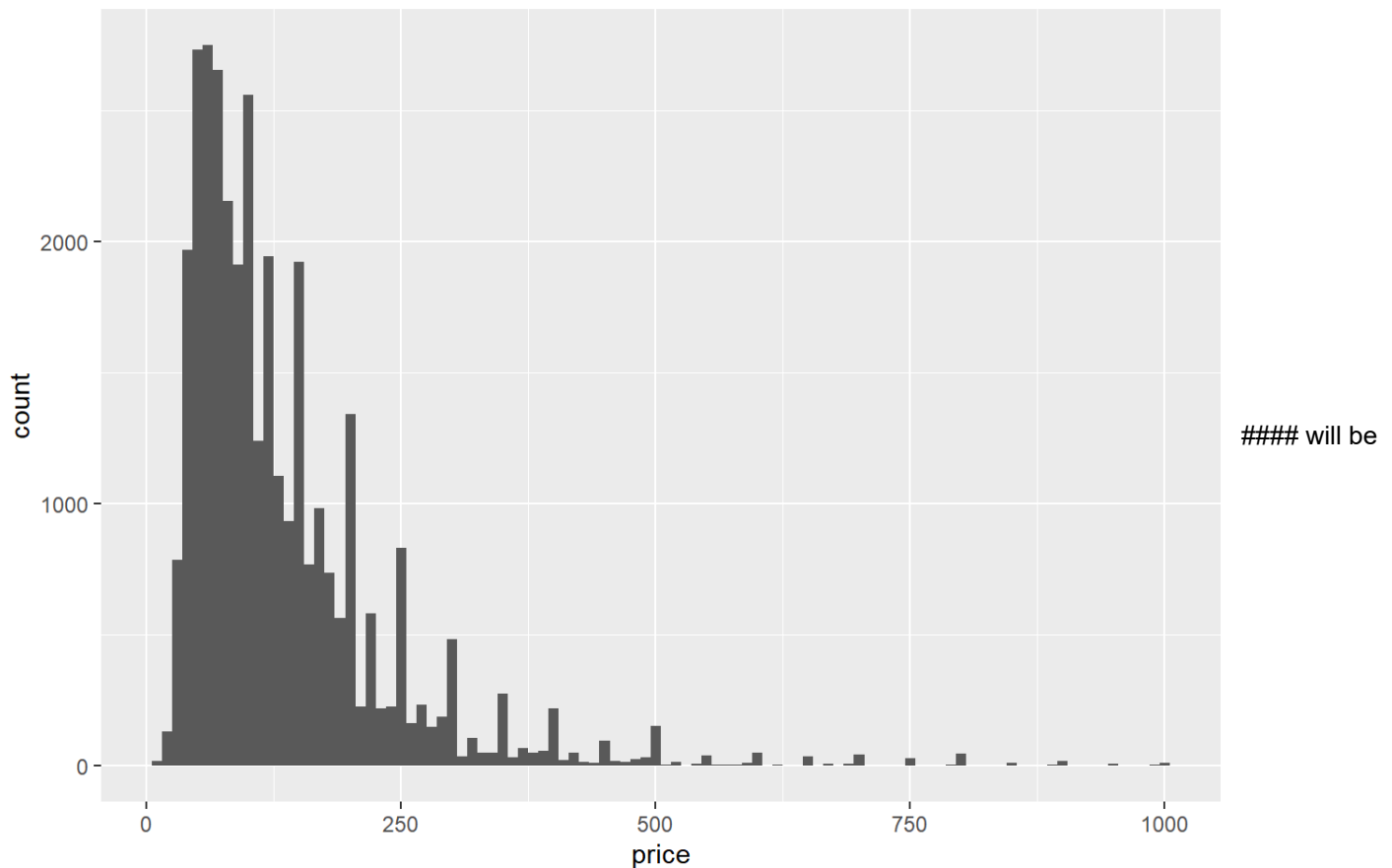
```
# odata %>%
#   filter(neighbourhood_group_cleansed=='Staten Island') %>%
#   count(zipcode, sort = TRUE)

#Brooklyn - 11211
#Bronx - 10469
#Manhattan - 10002
#Queens- 11385
#SA - 10301
```

EDA

##Target Variable - price distribution

```
ggplot(data, aes(x=price)) +
  geom_histogram(binwidth = 10)
```



analyzing numeric variables across price with scatter plots #### will be analyzing categorical and few numeric variables across price with box plots

```
box = geom_boxplot(varwidth=T)
scatter = geom_point()
```

Analyzing Categorical Variable

categorical:

neighbourhood_cleansed,neighbourhood_group_cleansed

neighbourhood_cleansed – explored neighbourhood_cleansed for each borough.. but Too many levels and no so DROP

```
# data %>%  
#   filter(neighbourhood_group_cleansed == 'Brooklyn') %>%  
#   ggplot(aes(x=neighbourhood_cleansed,y=price)) + box
```

```
# data %>%  
#   filter(neighbourhood_group_cleansed == 'Bronx') %>%  
#   ggplot(aes(x=neighbourhood_cleansed,y=price)) + box
```

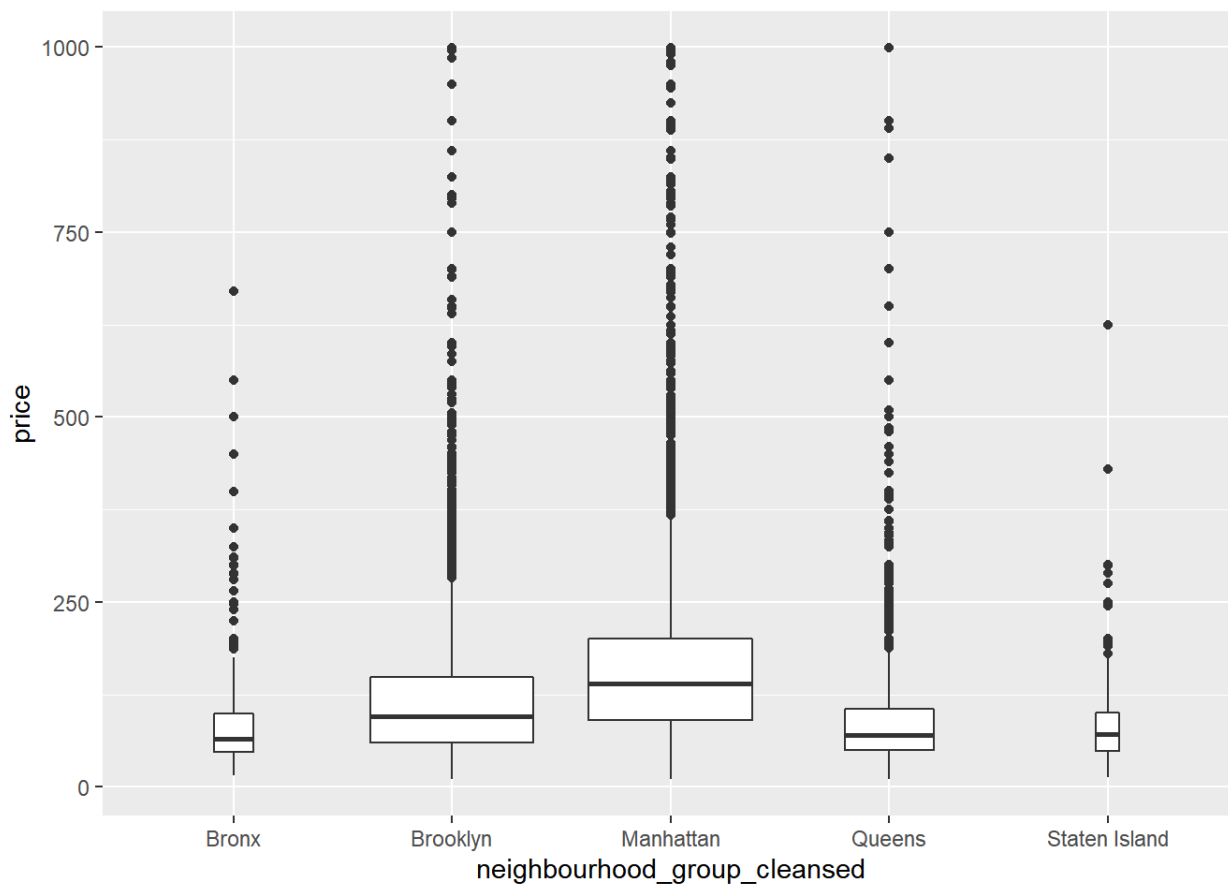
```
# data %>%  
#   filter(neighbourhood_group_cleansed == 'Manhattan') %>%  
#   ggplot(aes(x=neighbourhood_cleansed,y=price)) + box
```

```
# data %>%  
#   filter(neighbourhood_group_cleansed == 'Queens') %>%  
#   ggplot(aes(x=neighbourhood_cleansed,y=price)) + box
```

```
# data %>%  
#   filter(neighbourhood_group_cleansed == 'Staten Island') %>%  
#   ggplot(aes(x=neighbourhood_cleansed,y=price)) + box
```

neighbourhood_group_cleansed

```
ggplot(data, aes(x=neighbourhood_group_cleansed,y=price)) + box
```

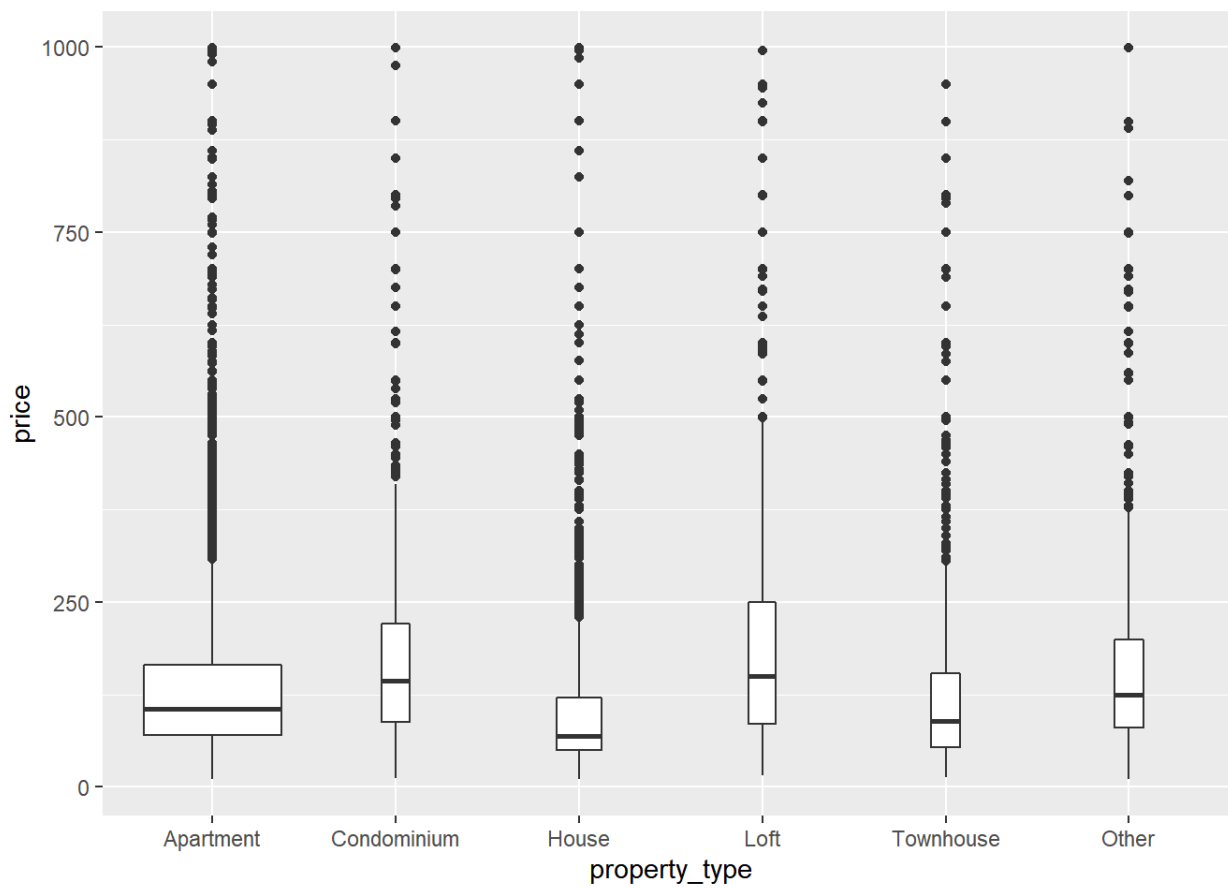


Analyzing Factor Variables

factor: property_type, room_type

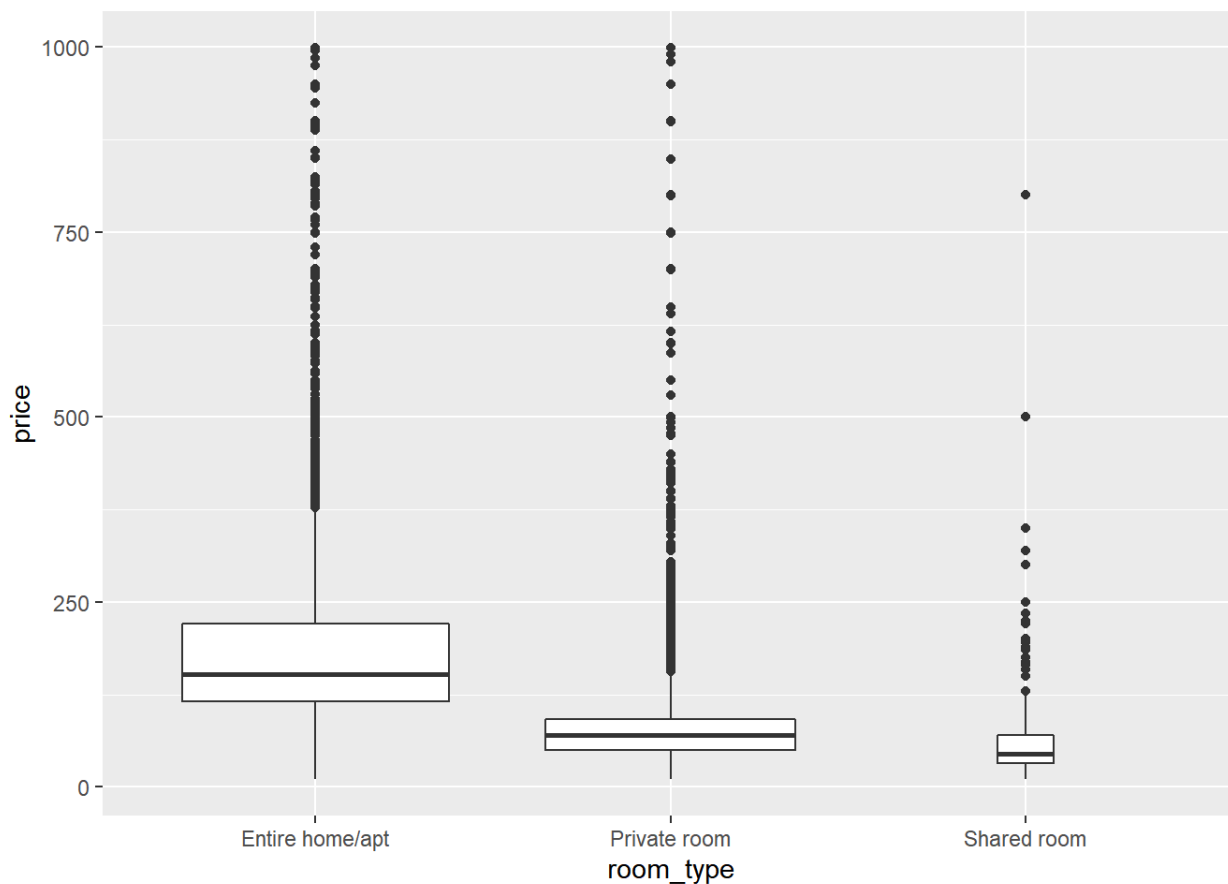
property_type

```
ggplot(data, aes(x=property_type,y=price)) + box
```



room_type

```
ggplot(data, aes(x=room_type,y=price)) + box
```

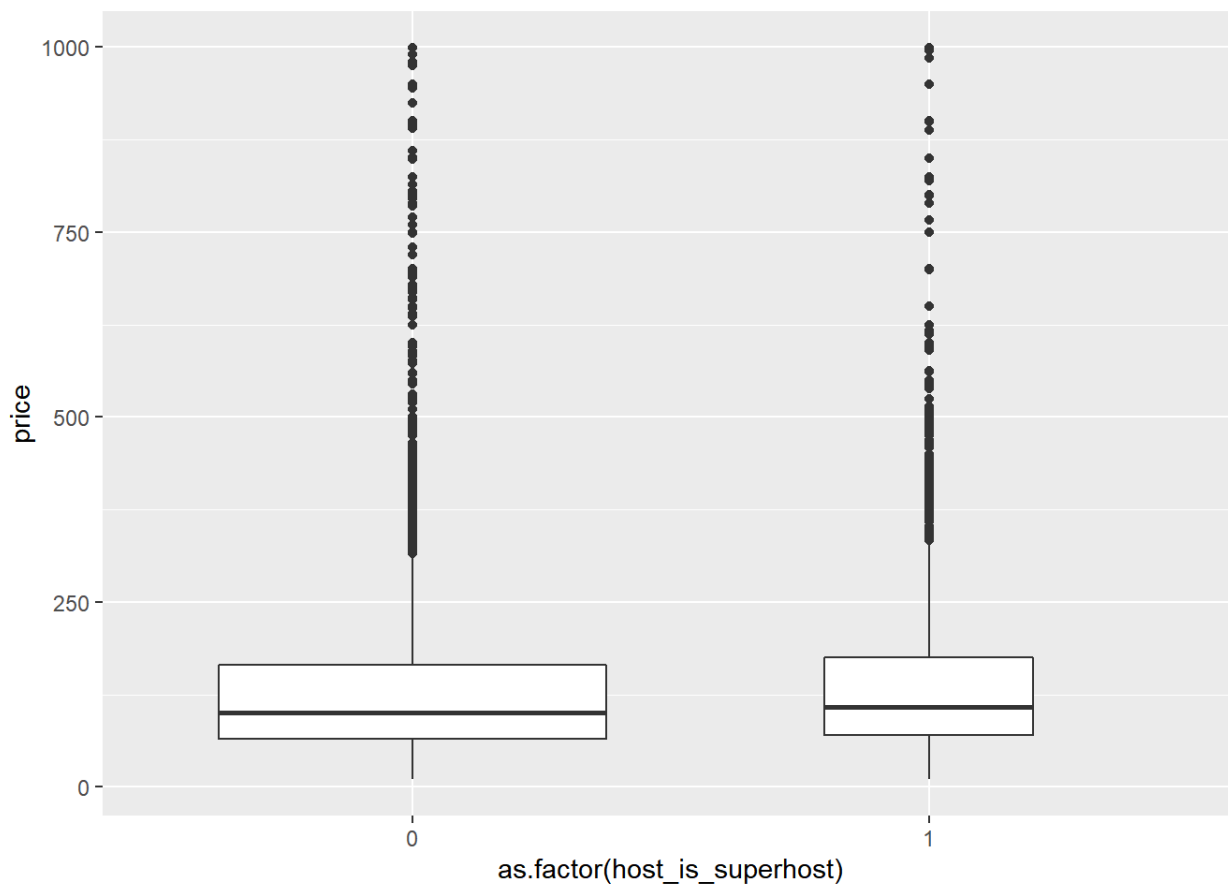


Analyzing Numeric Factor Variables

numeric factor: host_is_superhost, accommodates, bathrooms, bedrooms, beds, availability_30, review_scores_rating, review_scores_accuracy, review_scores_cleanliness, review_scores_checkin, review_scores_communication, review_scores_location, review_scores_value, instant_bookable

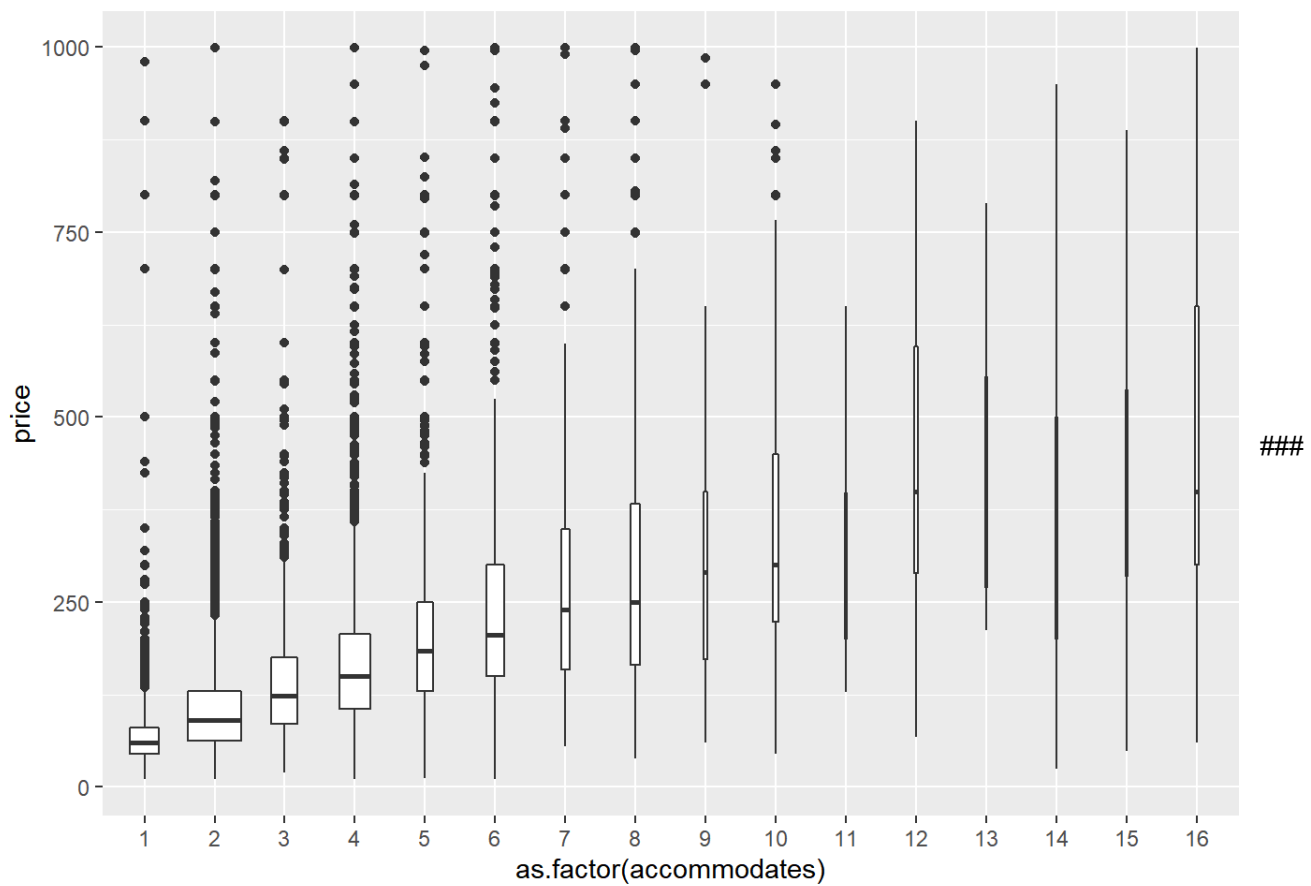
host_is_superhost

```
ggplot(data, aes(x=as.factor(host_is_superhost),y=price)) + box
```



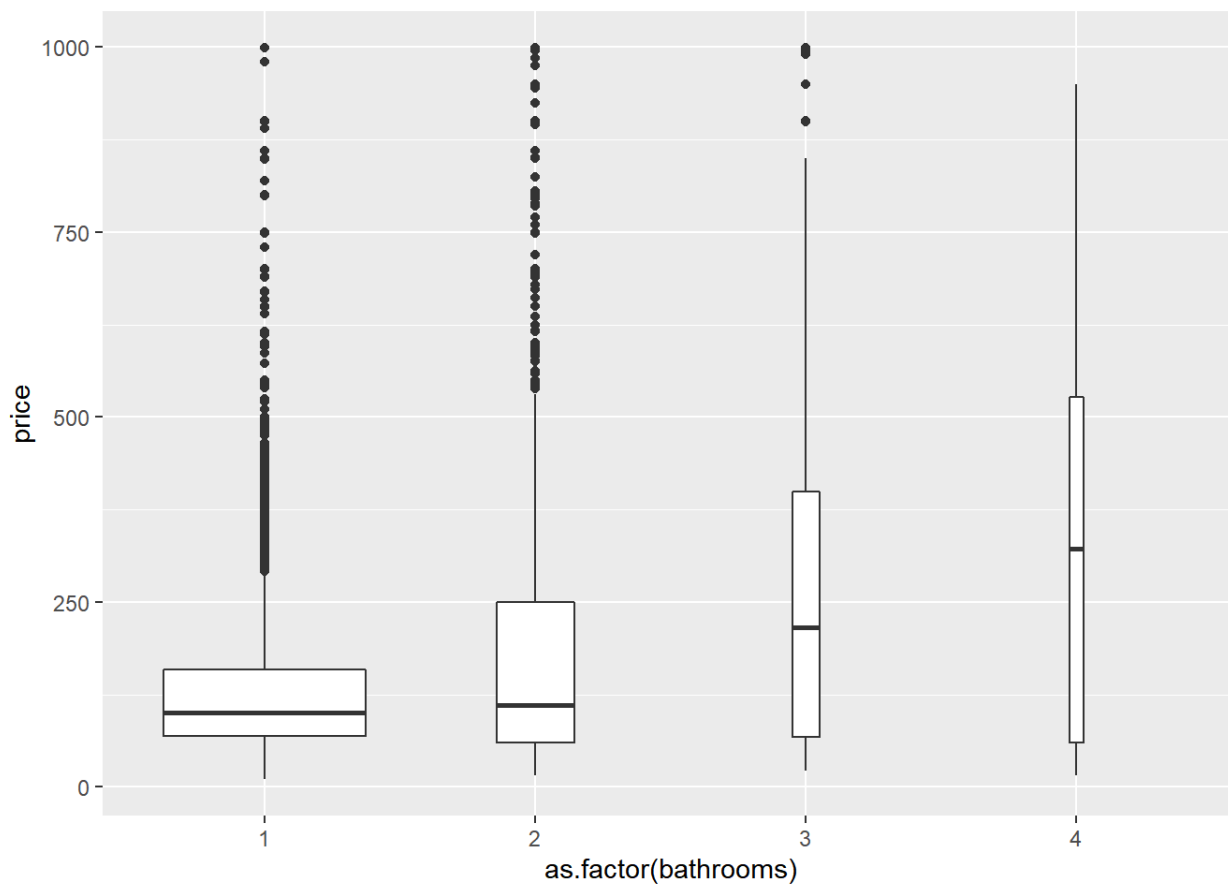
accommodates – might convert 9 onwards as 9

```
ggplot(data, aes(x=as.factor(accommodates),y=price)) + box
```



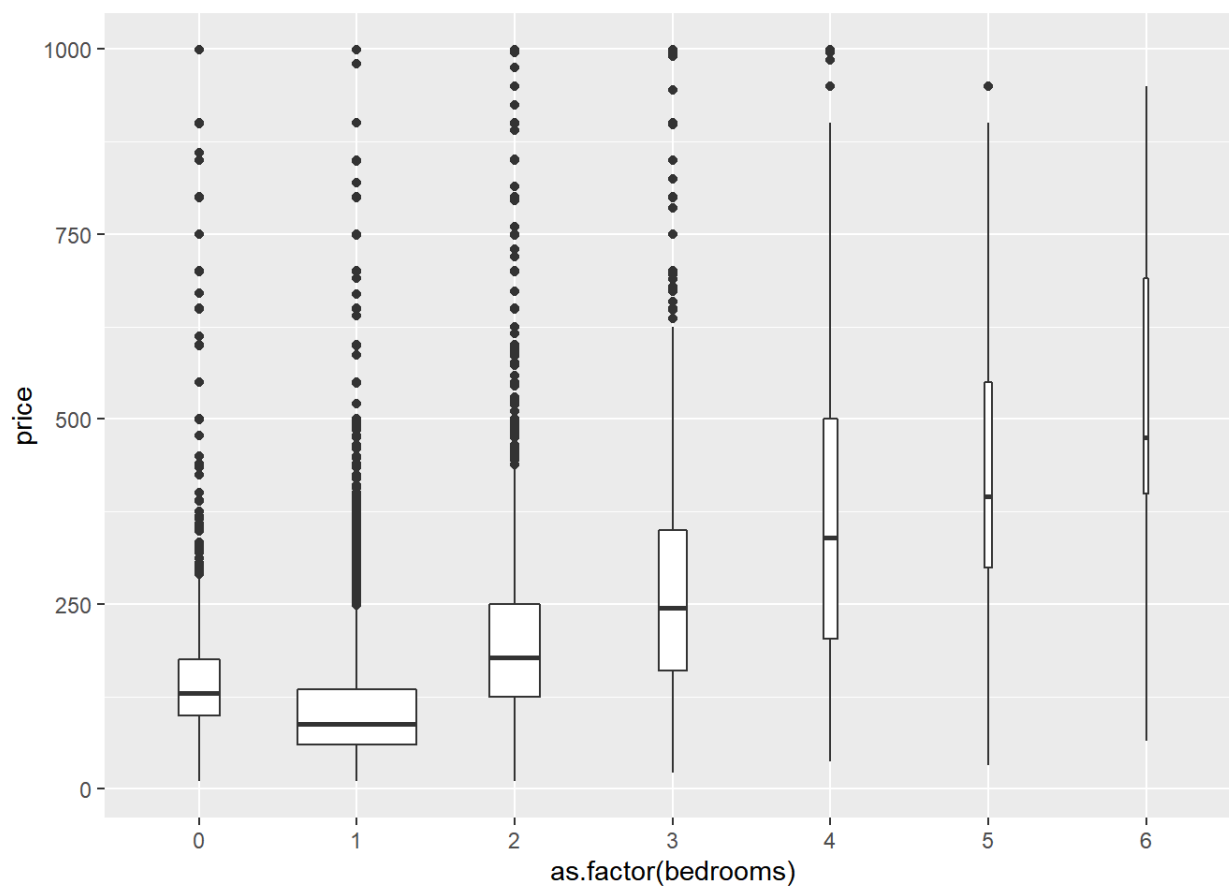
bathrooms- we have converted

```
ggplot(data, aes(x=as.factor(bathrooms),y=price)) + box
```



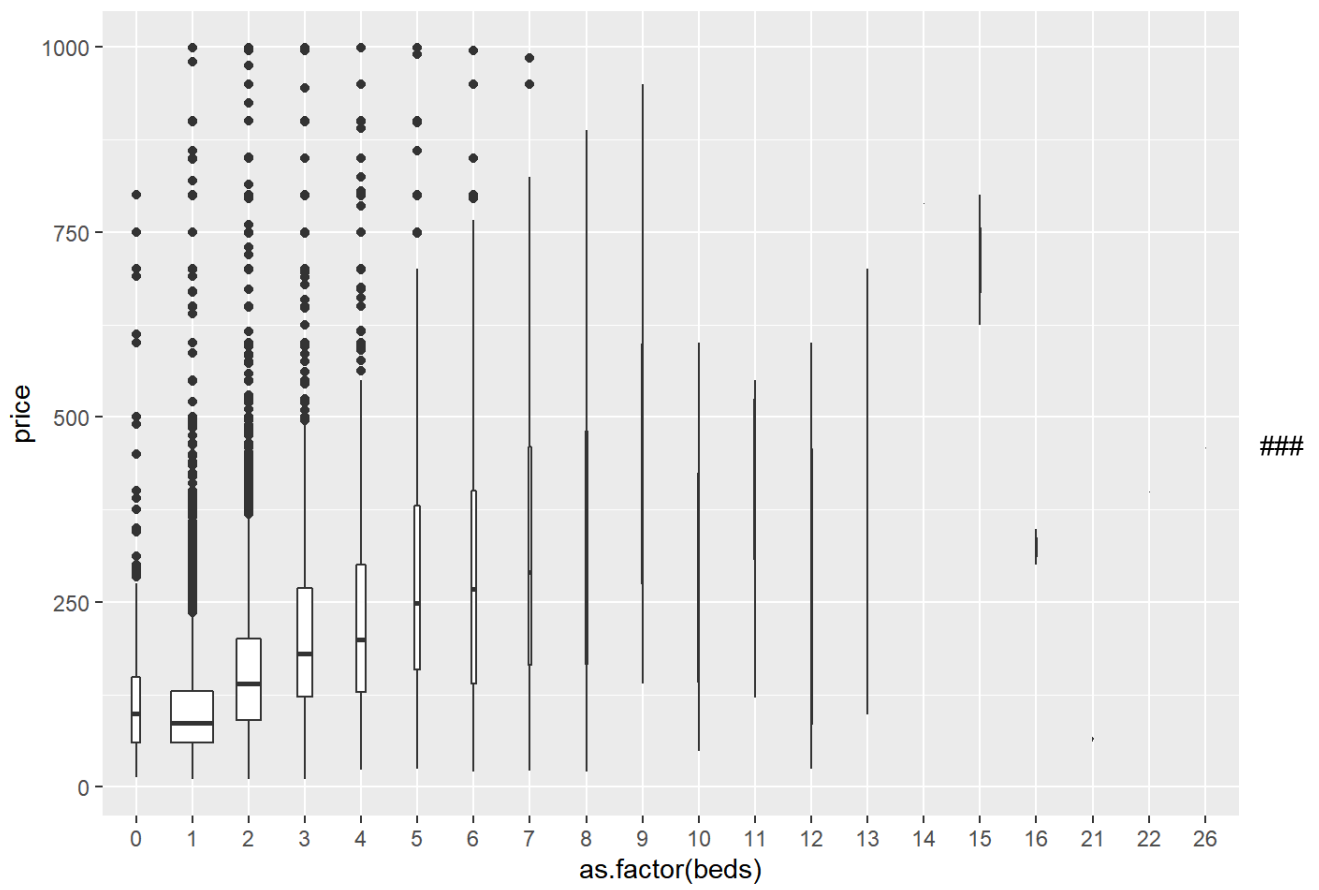
bedrooms

```
ggplot(data, aes(x=as.factor.bedrooms),y=price)) + box
```



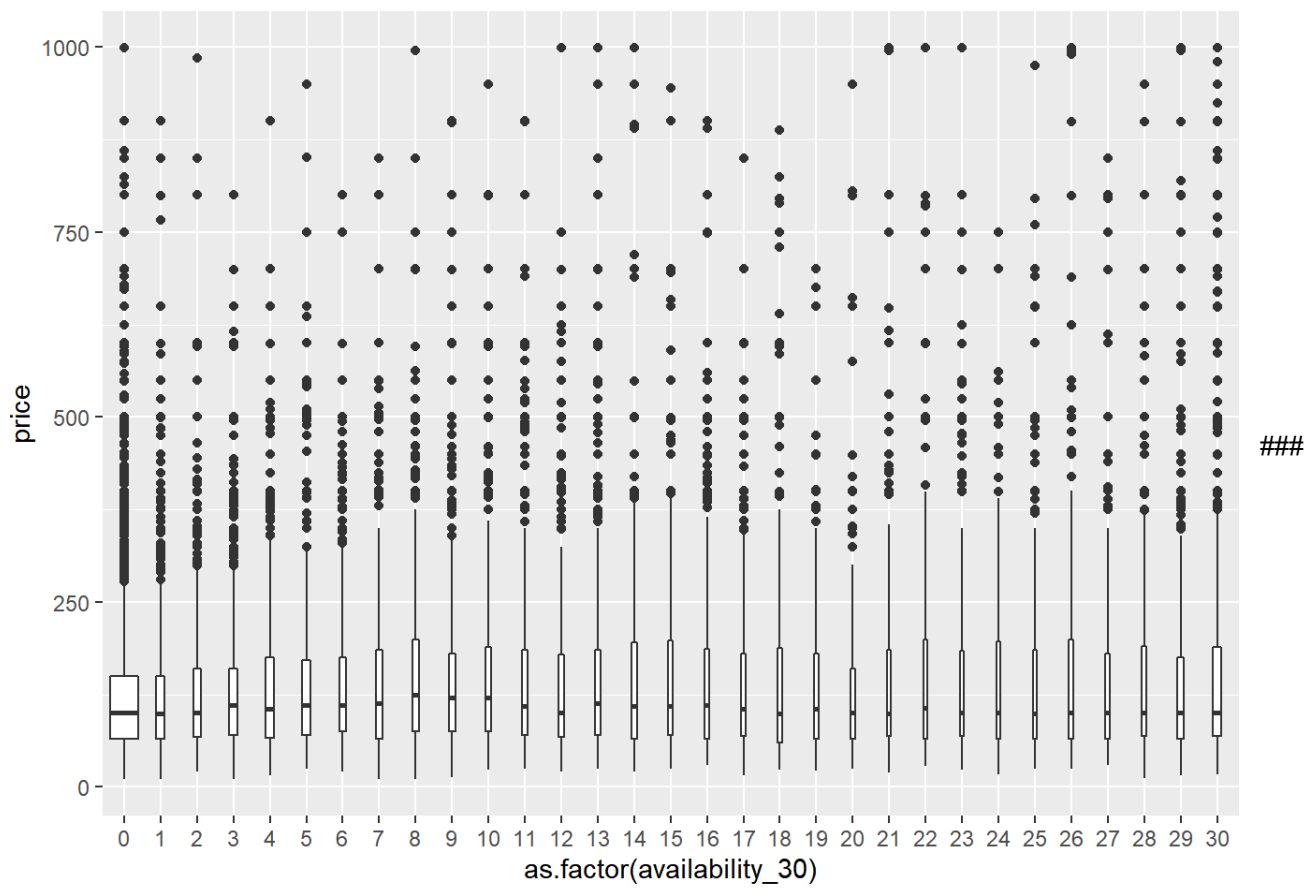
###beds – might convert 8 onwards as 8

```
ggplot(data, aes(x=as.factor(beds),y=price)) + box
```

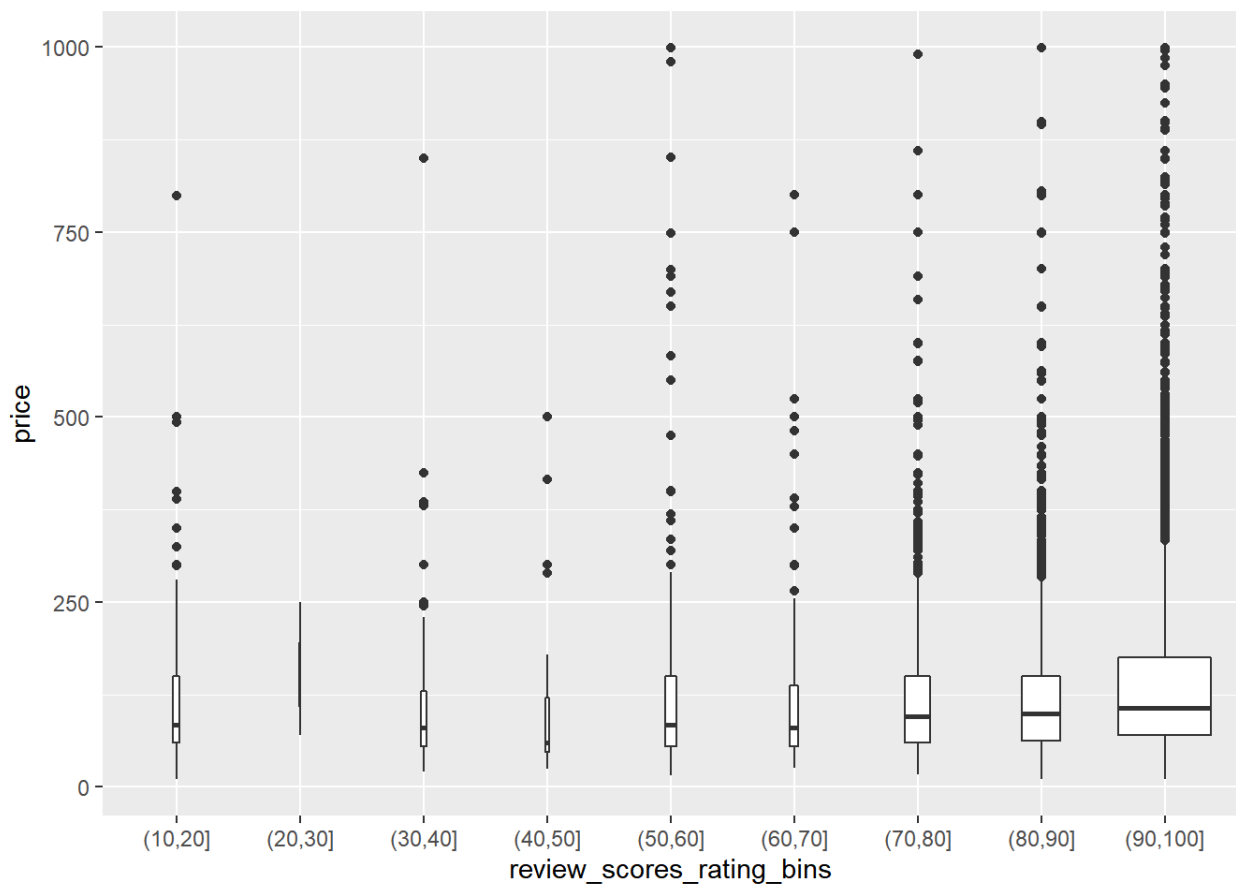
availability_30

```
ggplot(data, aes(x=as.factor(availability_30),y=price)) + box
```



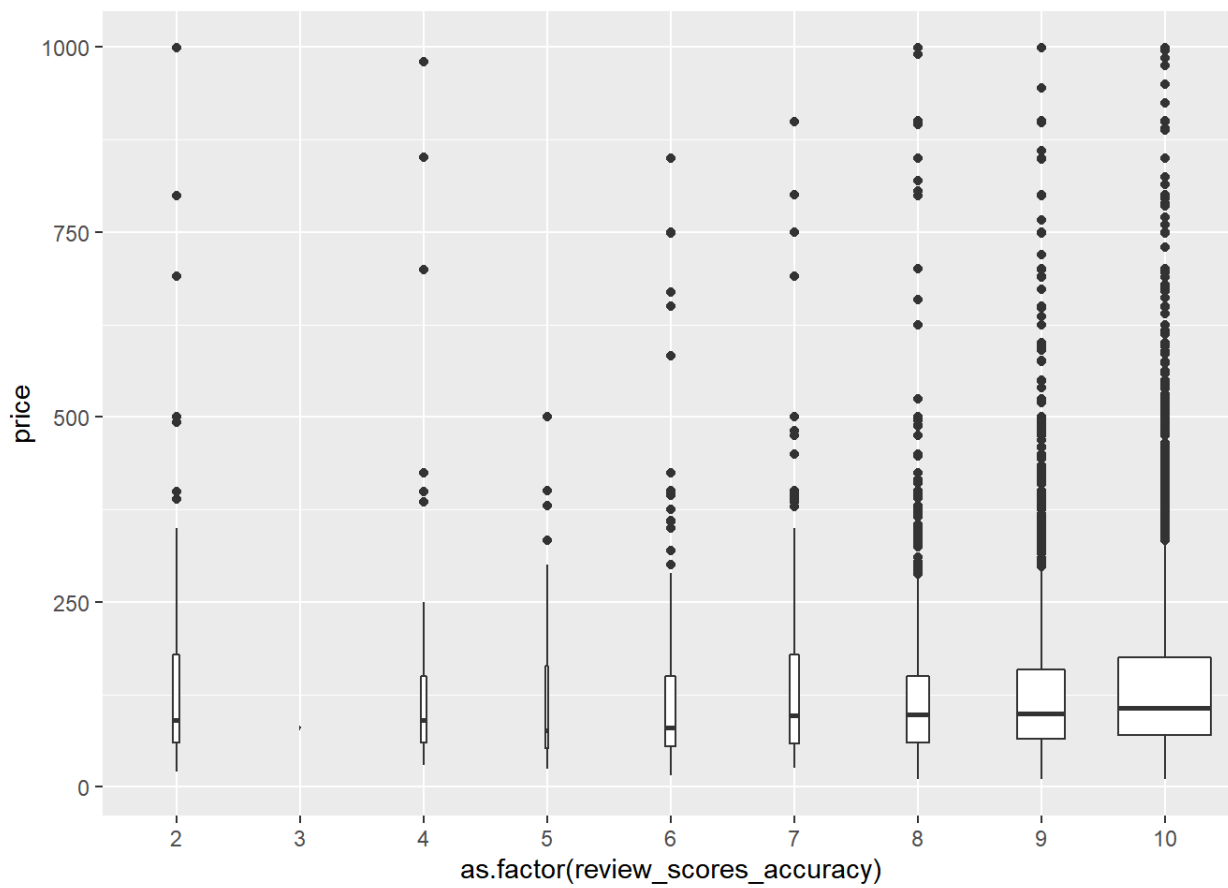
review_scores_rating

```
data%>%
  mutate(review_scores_rating_bins = cut(review_scores_rating, breaks = c(0,10,20,30,40,50,60,70,80,90,
    100))) %>%
  ggplot(aes(review_scores_rating_bins,price)) + box
```



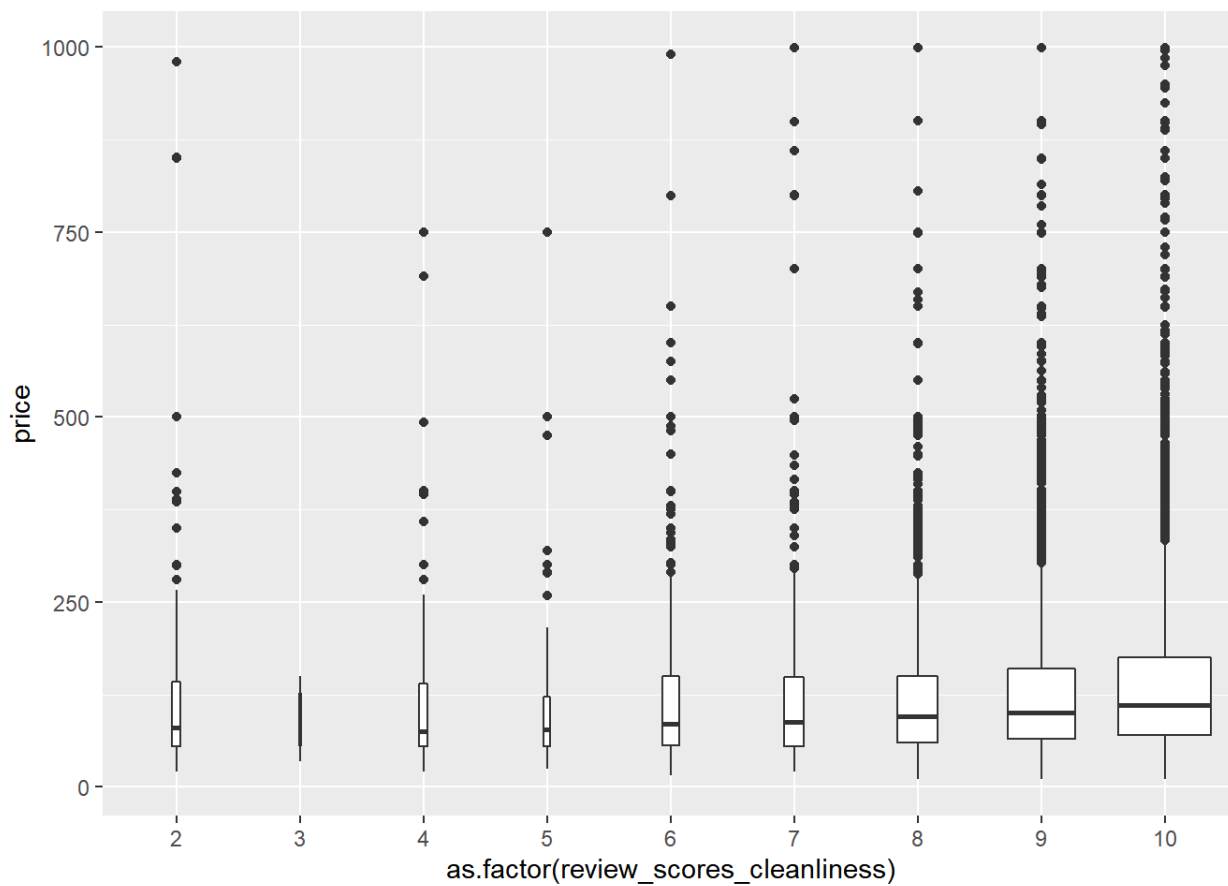
review_scores_accuracy

```
ggplot(data, aes(x=as.factor(review_scores_accuracy),y=price)) + box
```



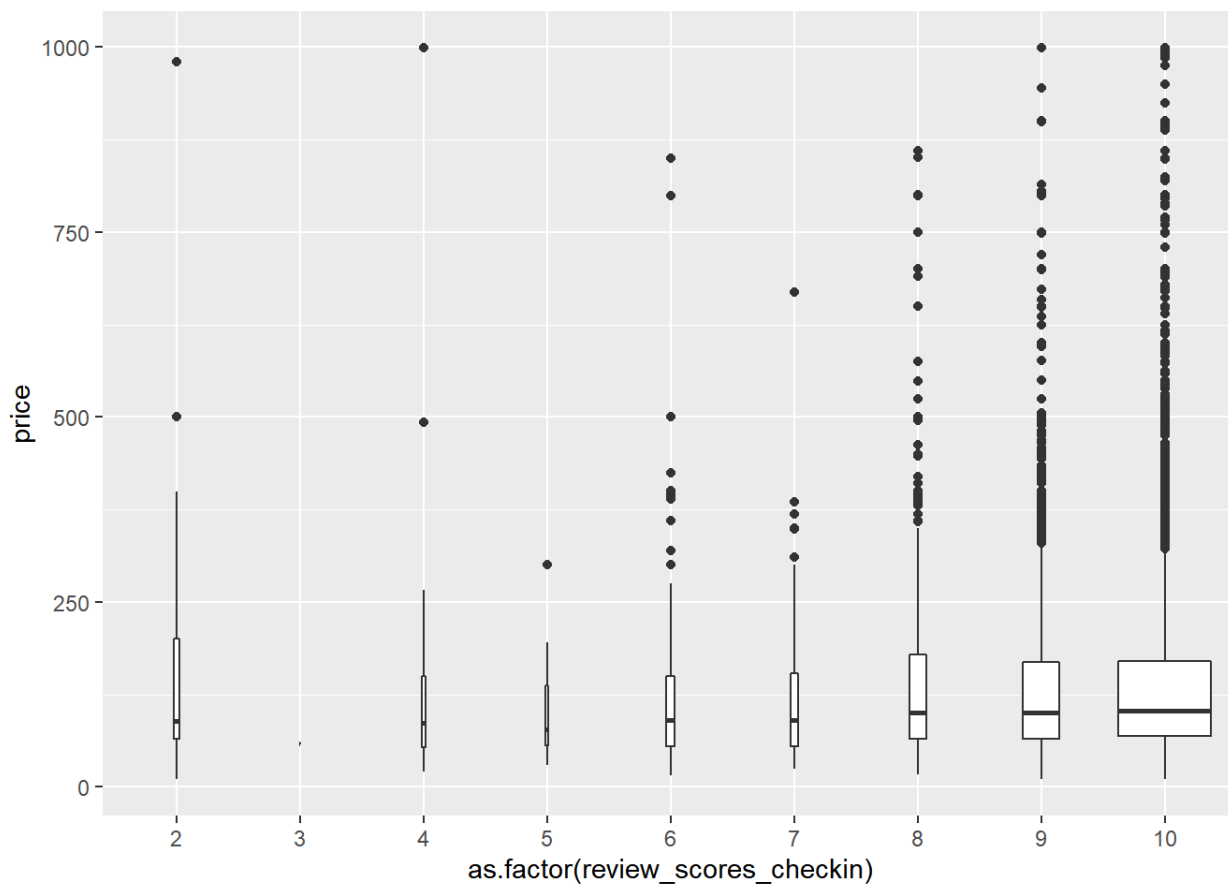
review_scores_cleanliness

```
ggplot(data, aes(x=as.factor(review_scores_cleanliness),y=price)) + box
```



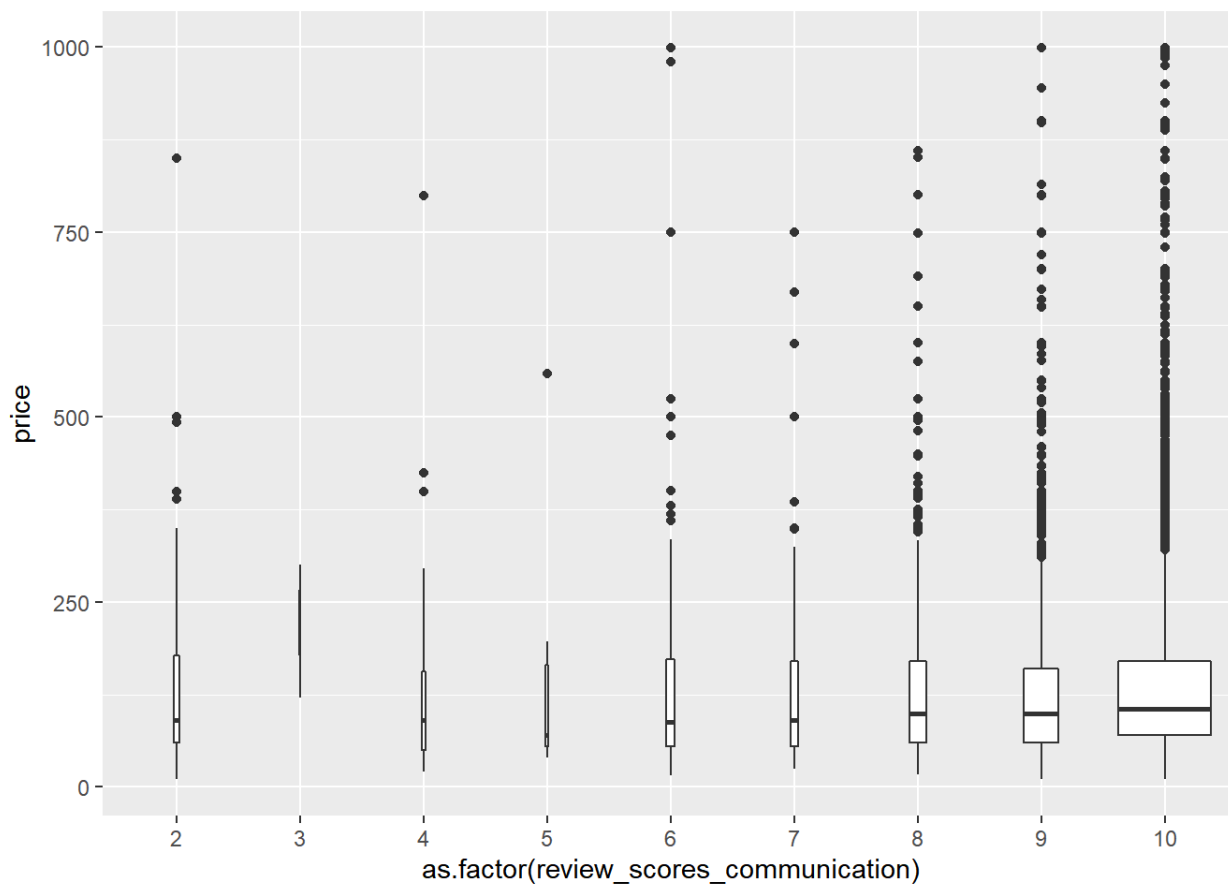
review_scores_checkin

```
ggplot(data, aes(x=as.factor(review_scores_checkin),y=price)) + box
```



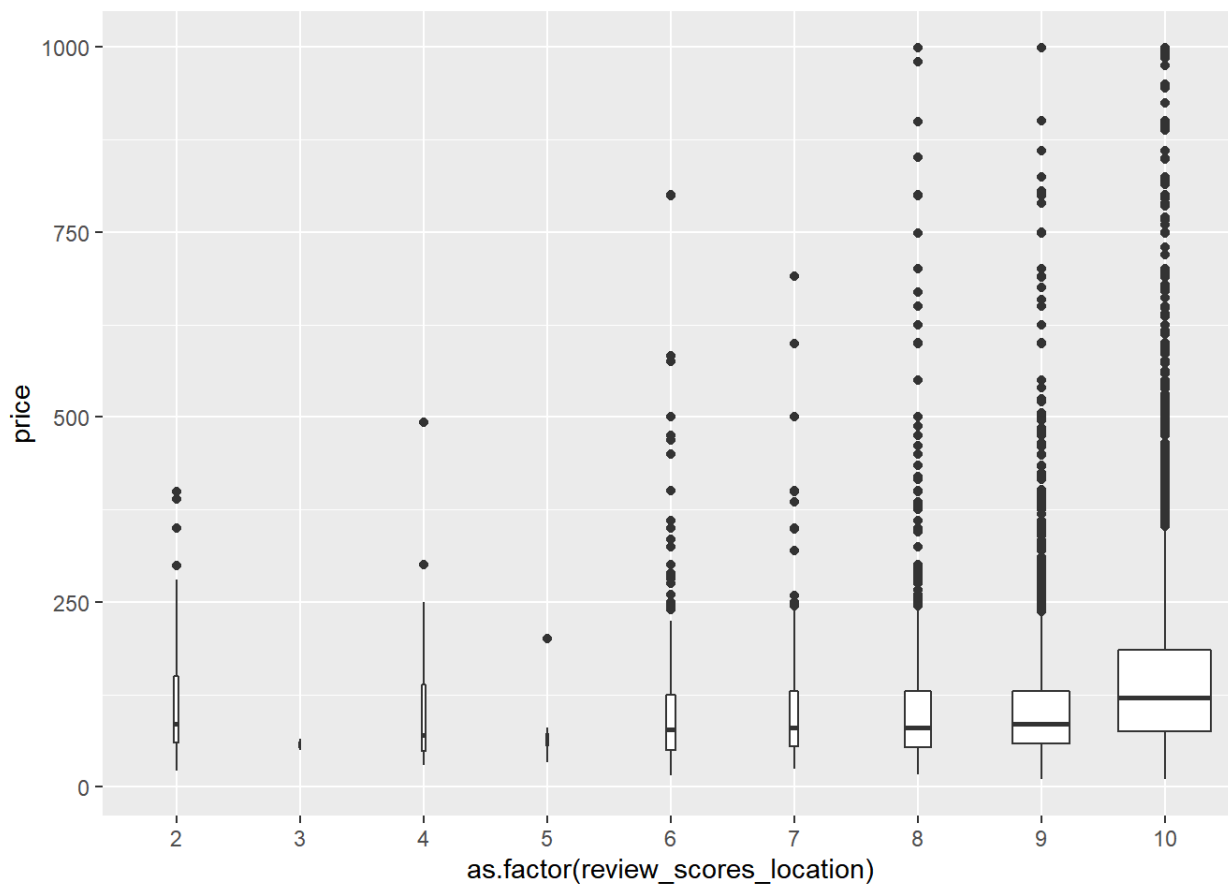
review_scores_communication

```
ggplot(data, aes(x=as.factor(review_scores_communication),y=price)) + box
```



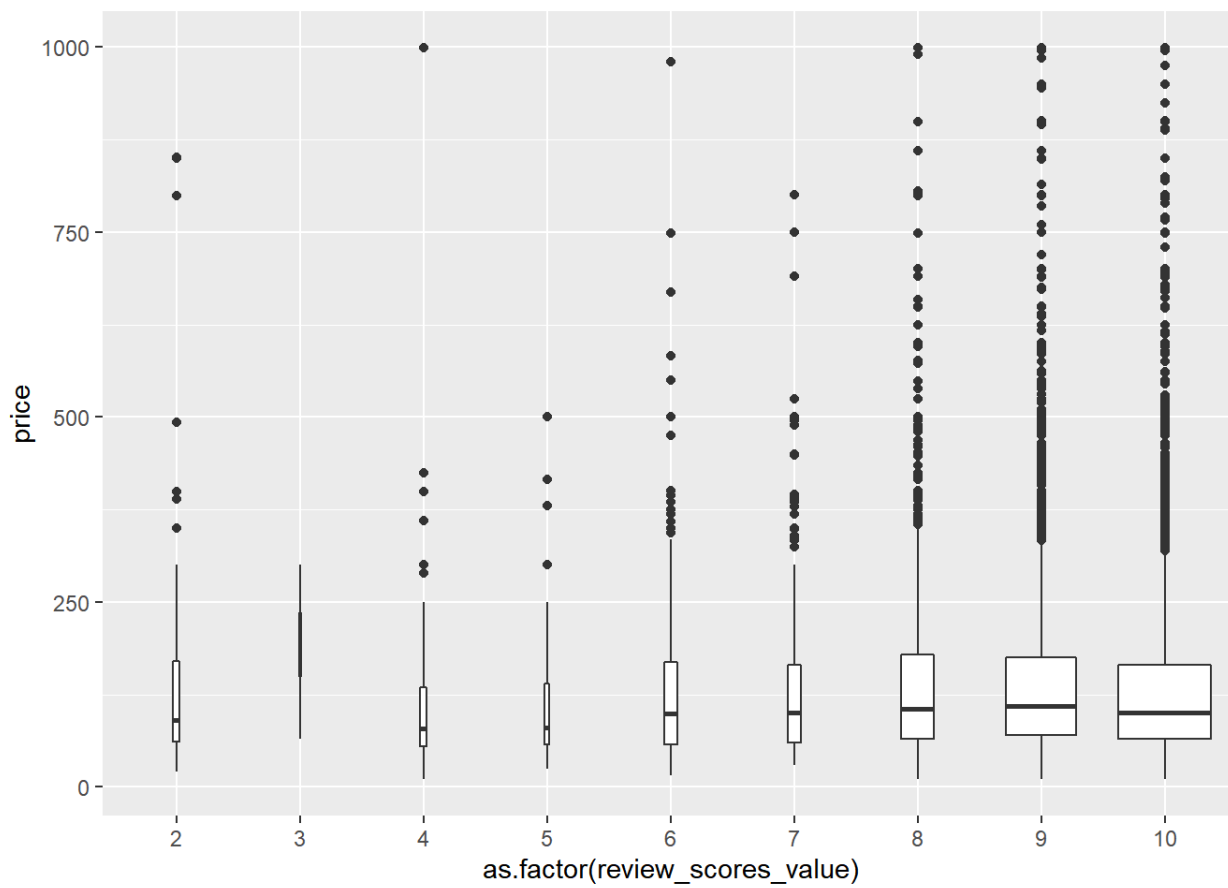
review_scores_location

```
ggplot(data, aes(x=as.factor(review_scores_location),y=price)) + box
```



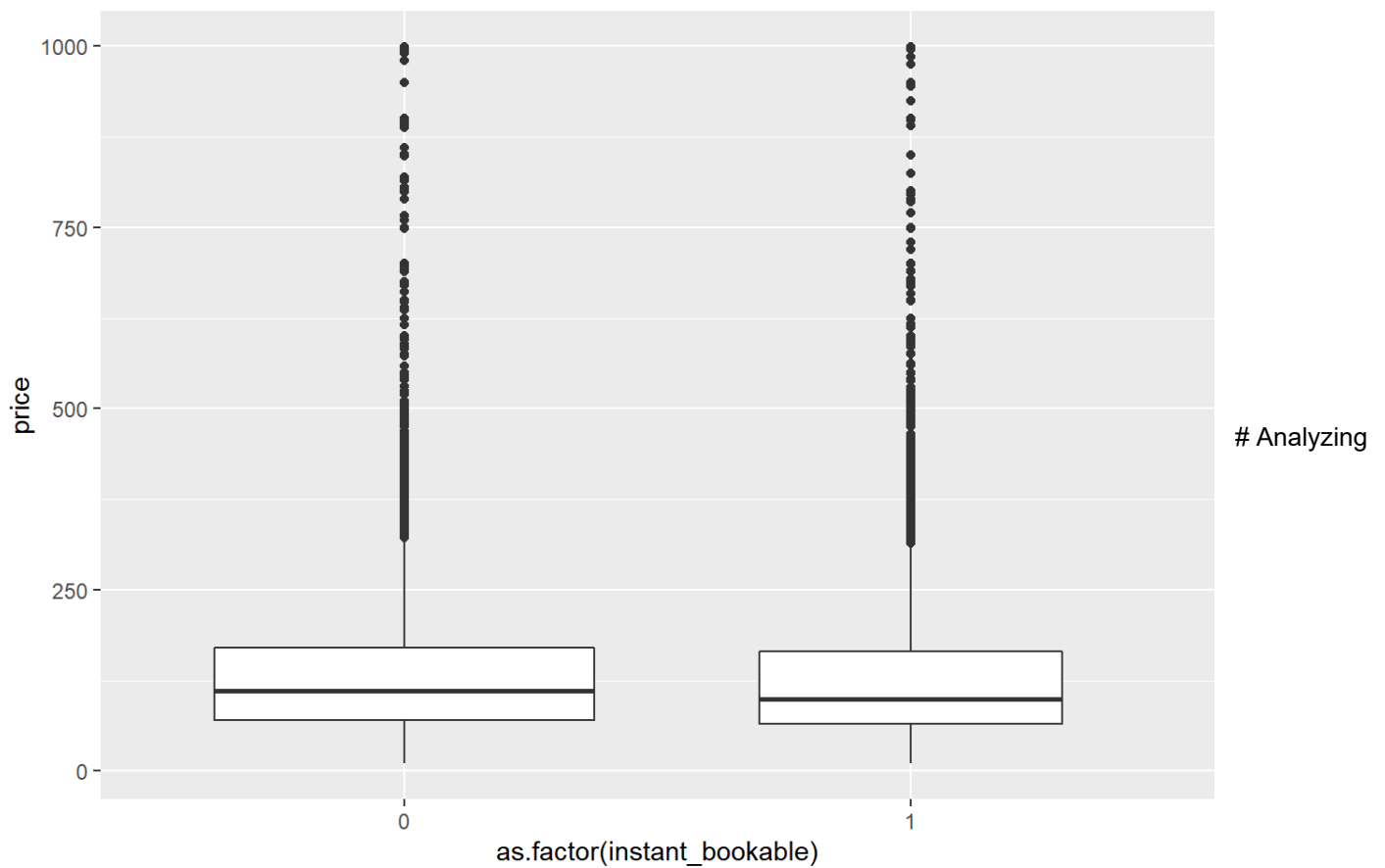
review_scores_value

```
ggplot(data, aes(x=as.factor(review_scores_value),y=price)) + box
```

instant_bookable

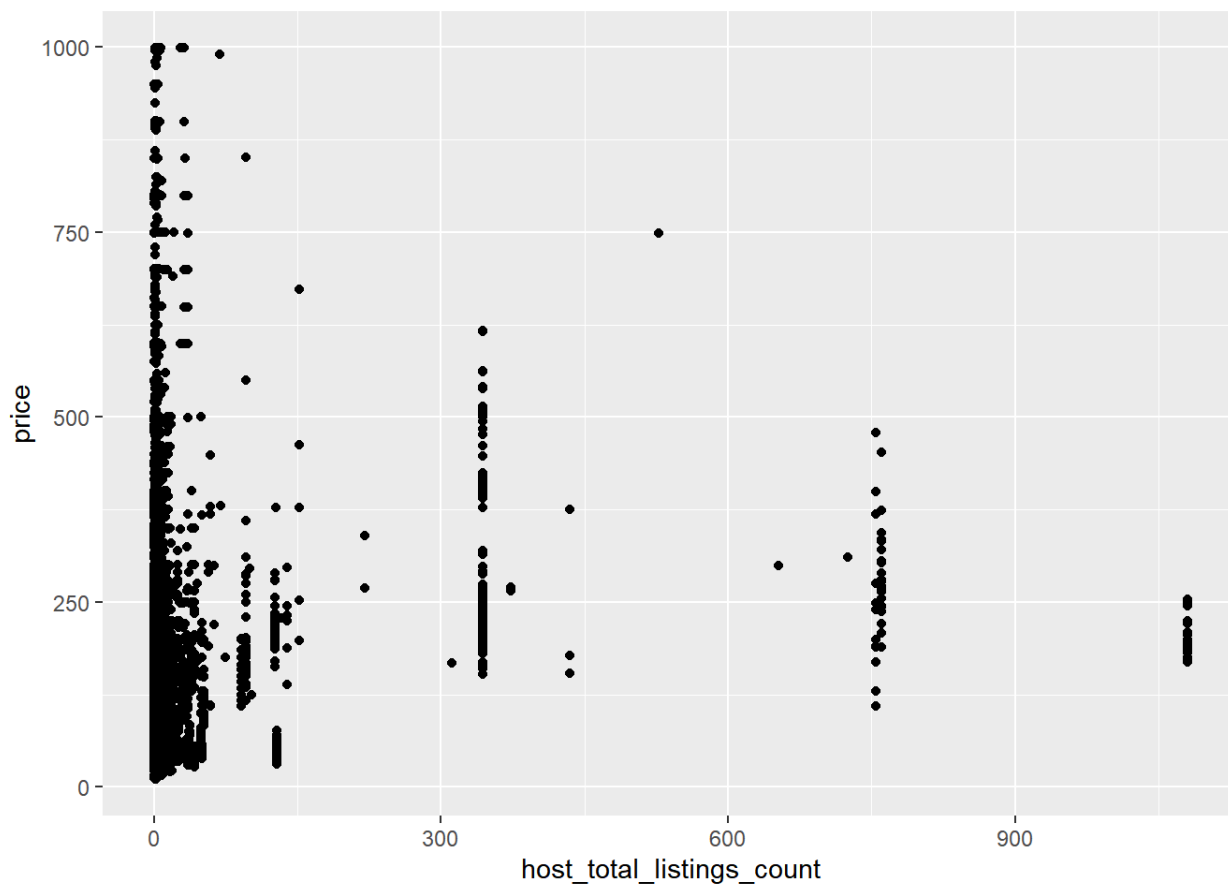
```
ggplot(data, aes(x=as.factor(instant_bookable),y=price)) + box
```



Numeric Variables ## numeric: host_total_listings_count, cleaning_fee, minimum_nights, maximum_nights, number_of_reviews, number_of_reviews_ltm, calculated_host_listings_count, reviews_per_month, listing_active_duration, host_active_duration

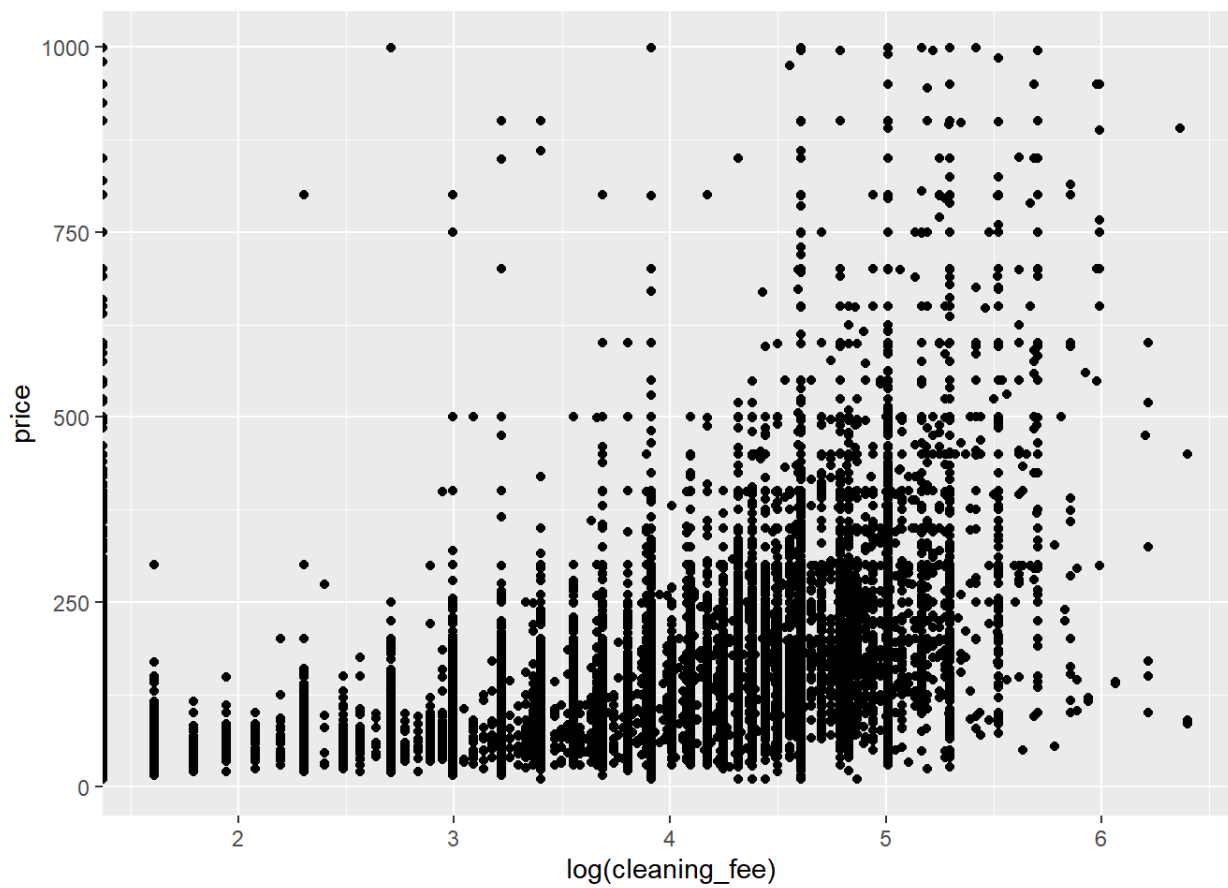
host_total_listings_count

```
ggplot(data, aes(x=host_total_listings_count,y=price)) + scatter
```



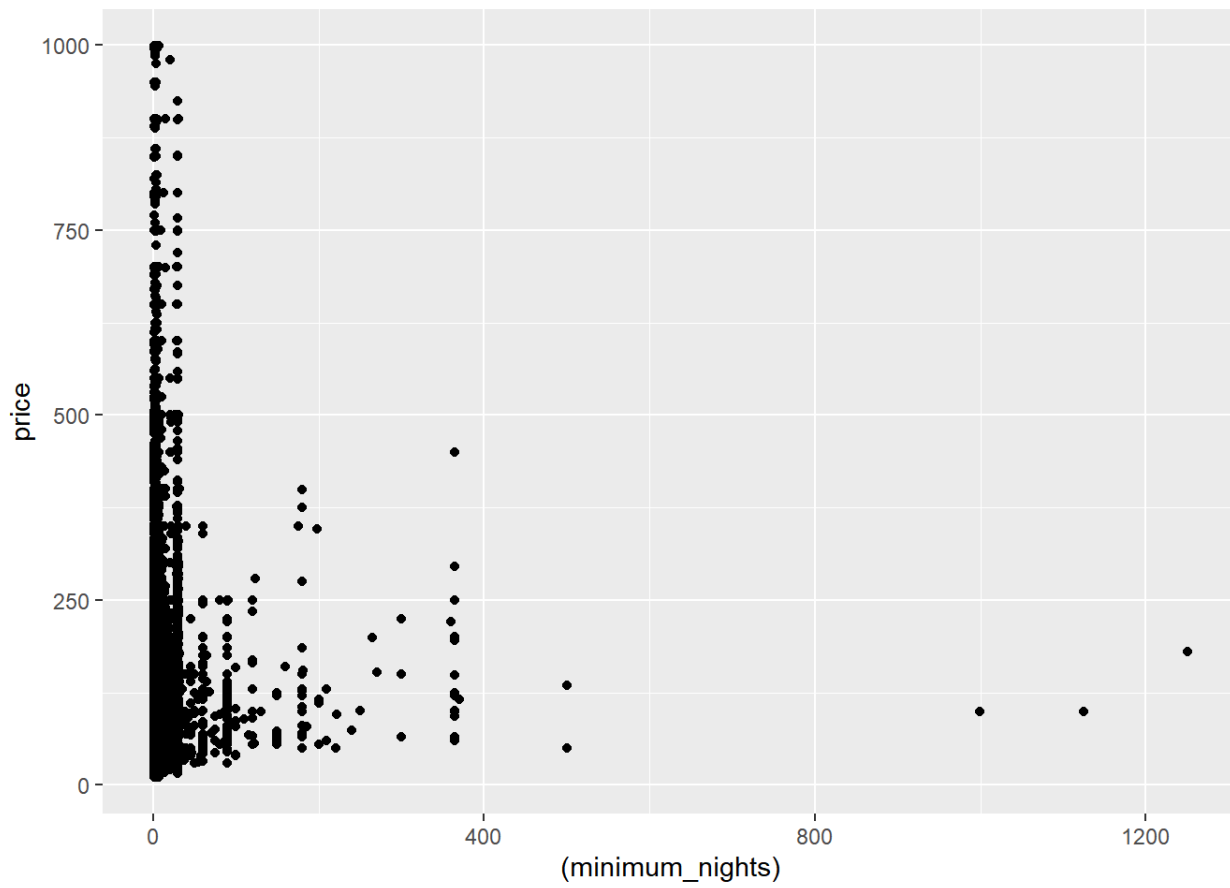
cleaning_fee

```
ggplot(data, aes(x=log(cleaning_fee),y=price)) + scatter #+ geom_smooth(method = "lm", se=FALSE, color = "black", aes(group=1))
```



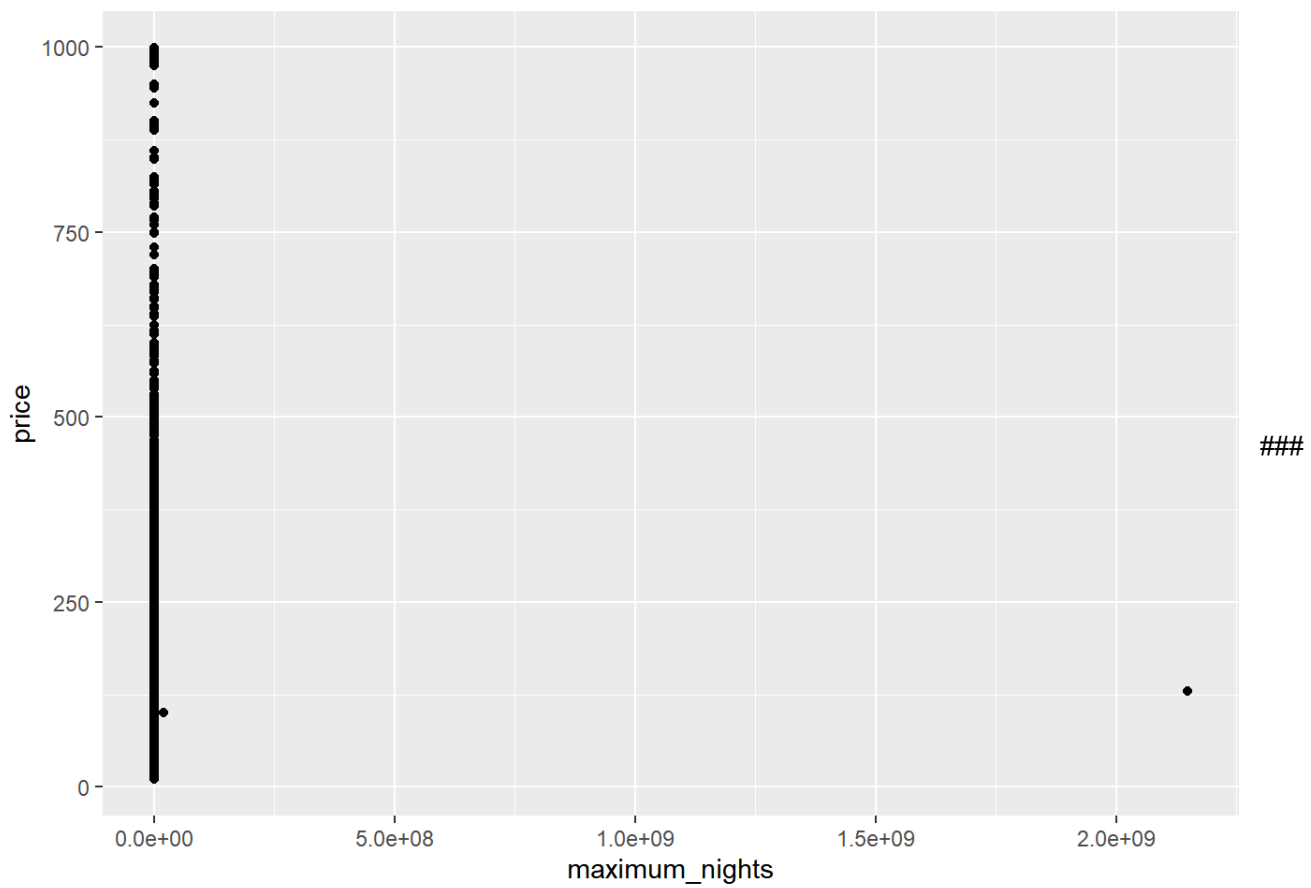
minimum_nights

```
ggplot(data, aes(x=(minimum_nights),y=price)) + scatter
```



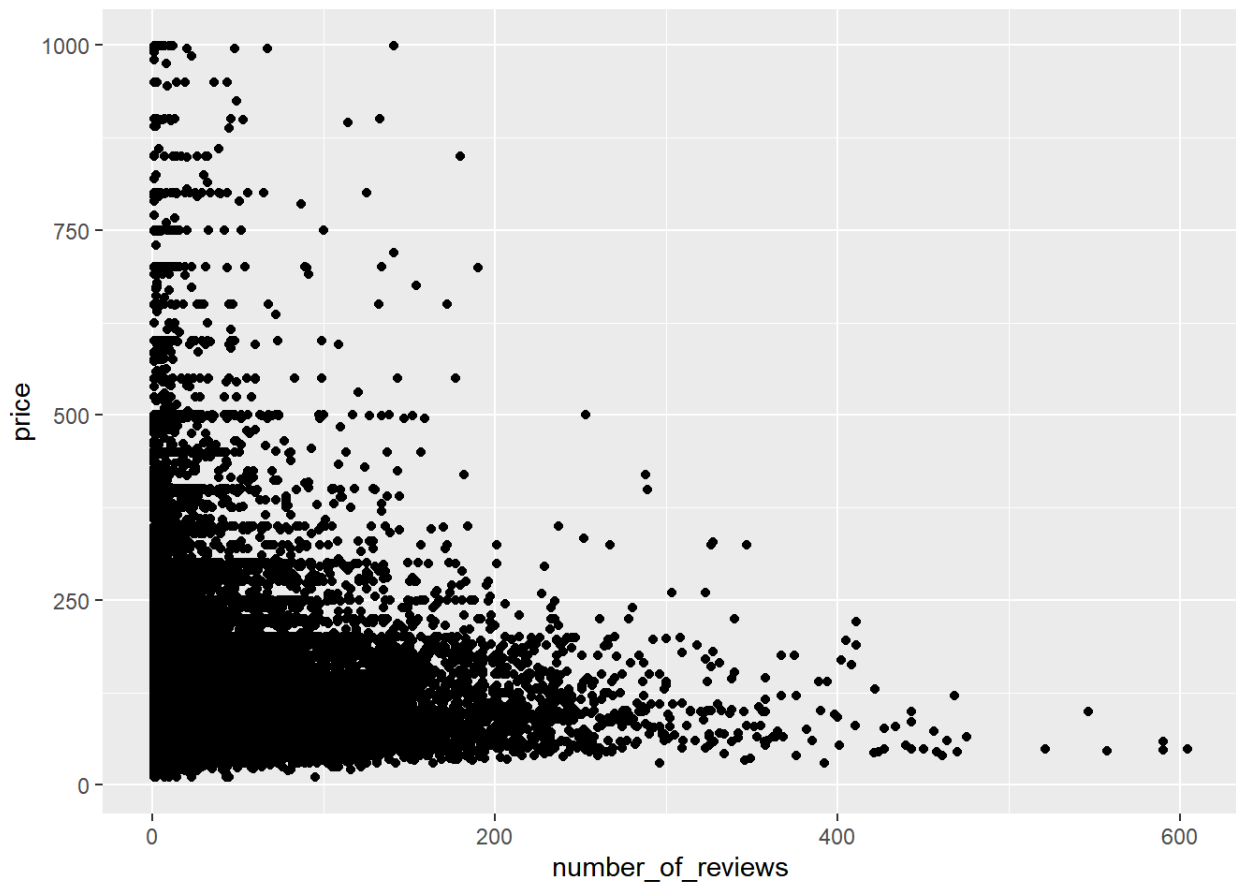
maximum_nights

```
ggplot(data, aes(x=maximum_nights,y=price)) + scatter
```



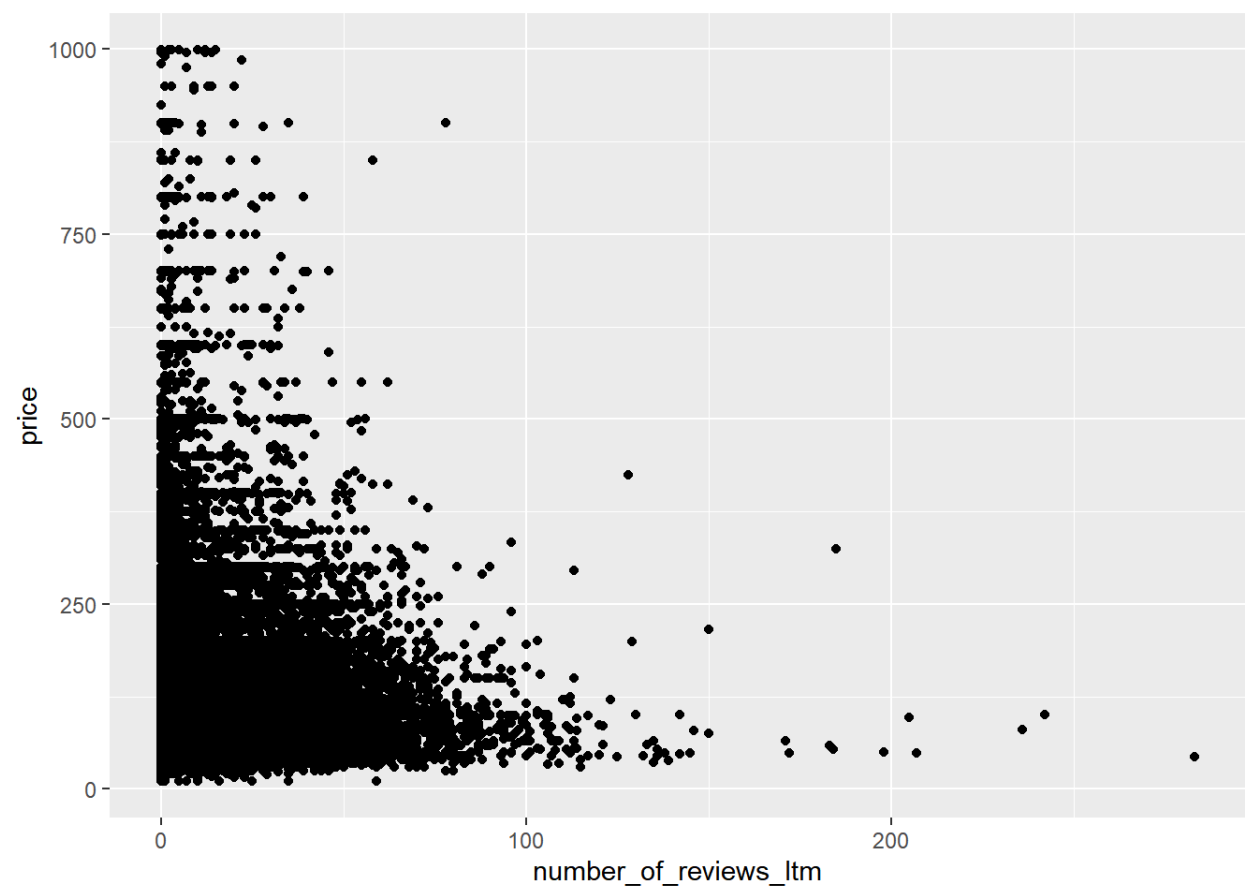
number_of_reviews

```
ggplot(data, aes(x=number_of_reviews,y=price)) + scatter
```



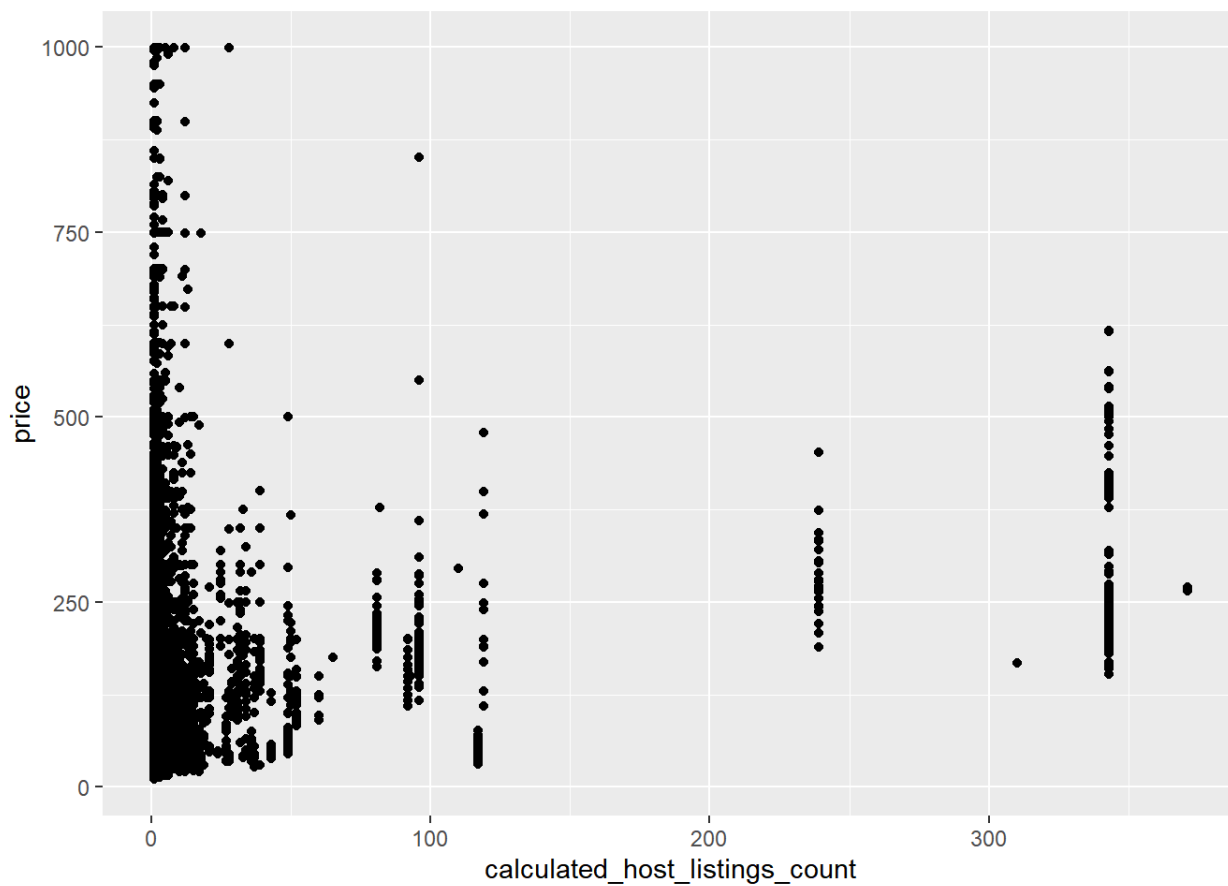
number_of_reviews_ltm

```
ggplot(data, aes(x=number_of_reviews_ltm,y=price)) + scatter
```



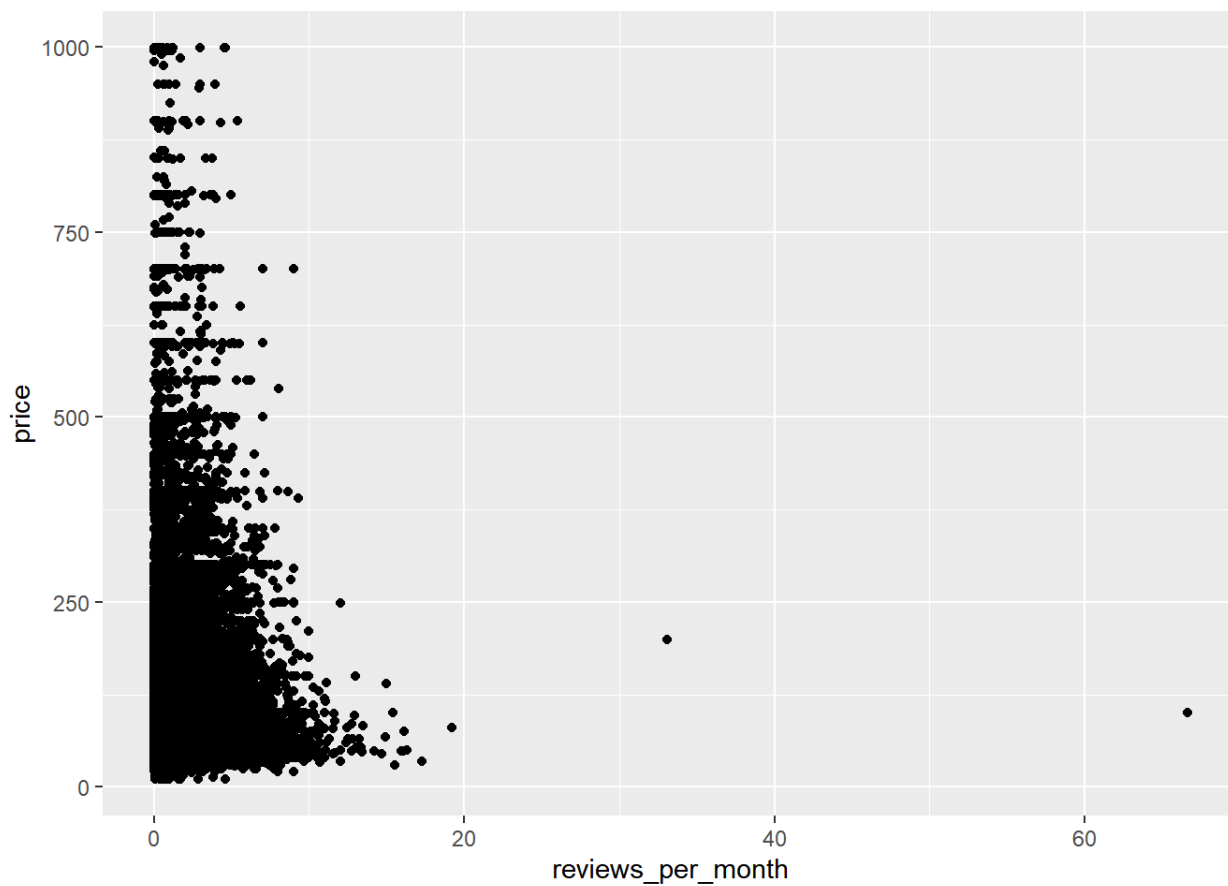
calculated_host_listings_count

```
ggplot(data, aes(x=calculated_host_listings_count,y=price)) + scatter
```



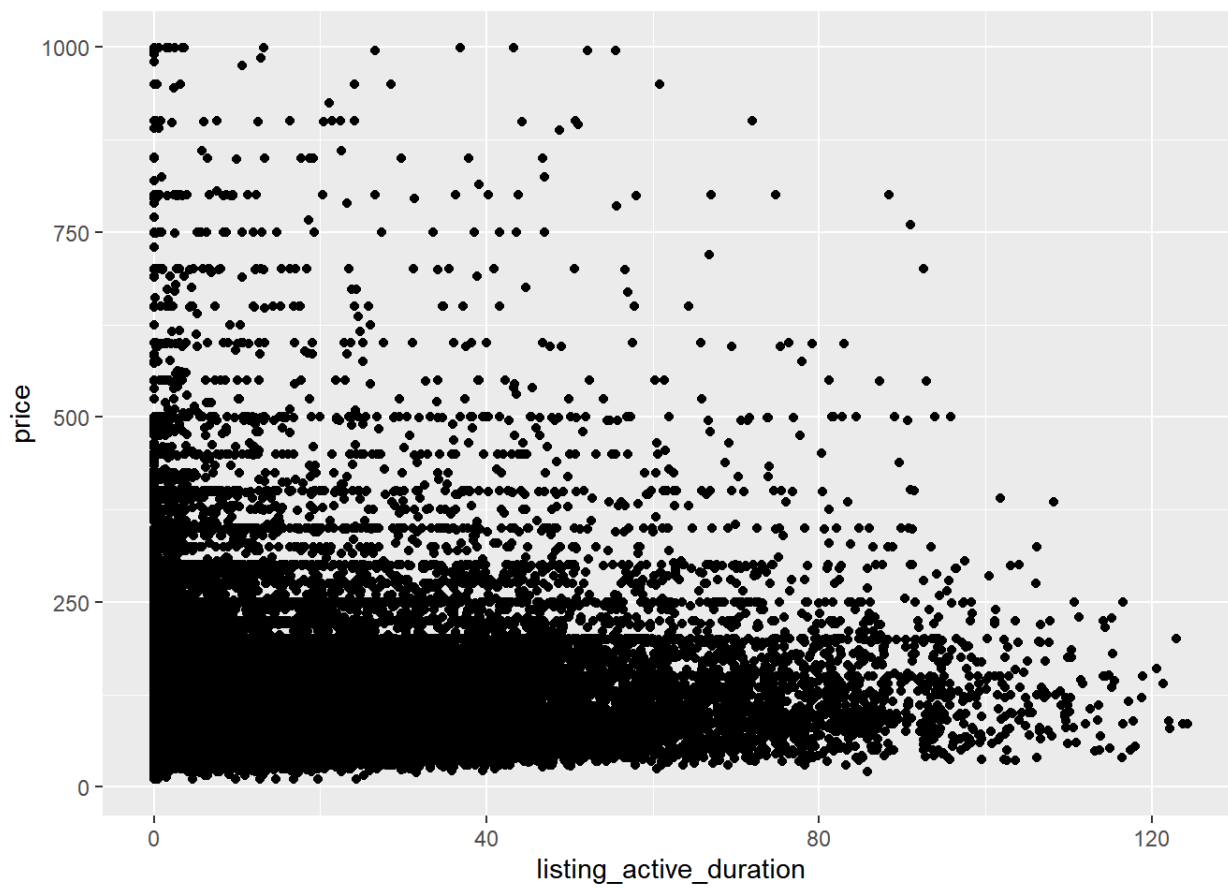
reviews_per_month

```
ggplot(data, aes(x=reviews_per_month,y=price)) + scatter
```

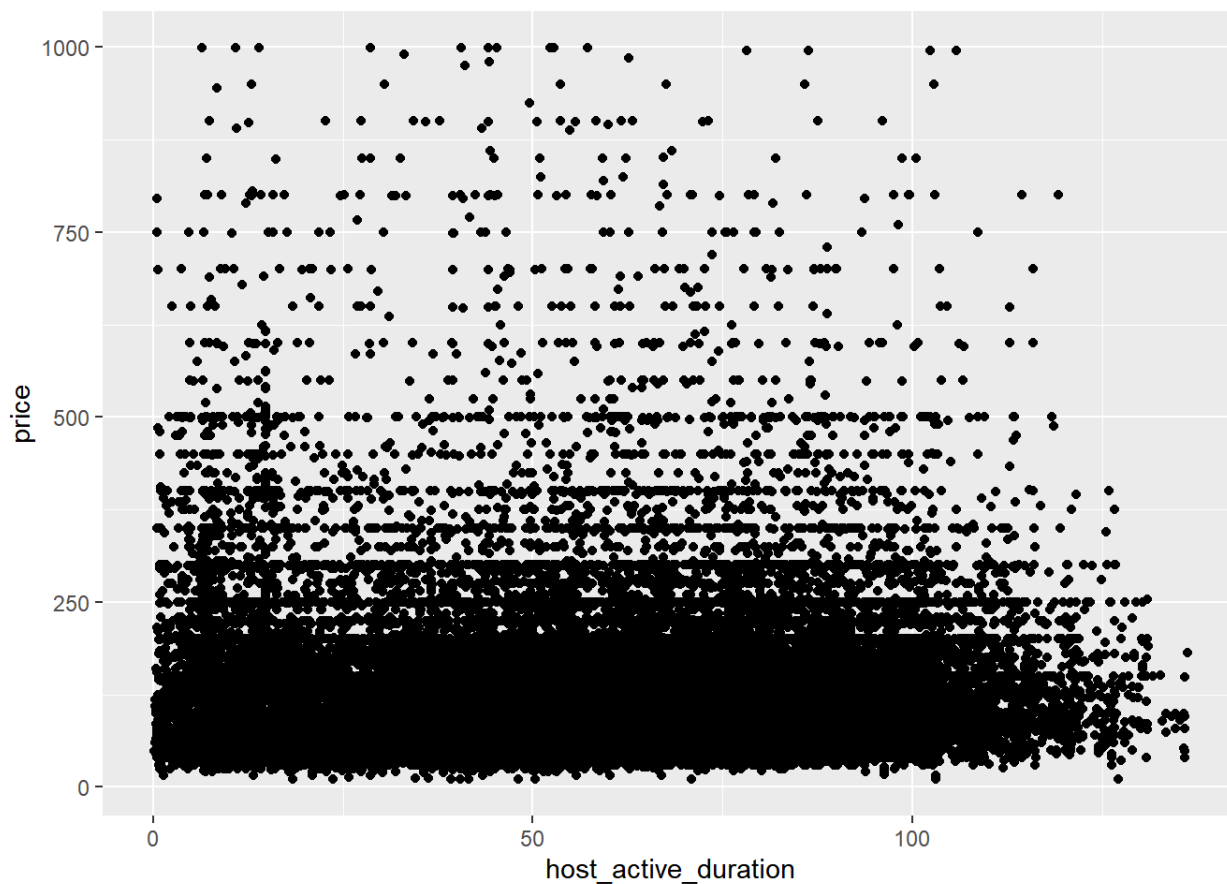
listing_active_duration

```
ggplot(data, aes(x=listing_active_duration,y=price)) + scatter
```



host_active_duration

```
ggplot(data, aes(x=host_active_duration,y=price)) + scatter
```



Correlations

```
#Library(corrplot)
#numeric variables
numericVars <- which(sapply(data, is.numeric))
data_numVar <- data[, numericVars]

#correlations of all numeric variables
cor_numVar <- cor(data_numVar, use="pairwise.complete.obs")

#sorting on decreasing correlations with price
cor_sorted <- as.matrix(sort(cor_numVar[, 'price'], decreasing = TRUE))
print(cor_sorted)
```

```
##                                [,1]
## price                        1.0000000000
## accommodates                 0.5599244814
## cleaning_fee                 0.5204372469
## beds                        0.4304395718
## bedrooms                    0.4298263449
## guests_included              0.3792144020
## bathrooms                   0.2556035195
## review_scores_location       0.1321612753
## extra_people                 0.1028865853
## calculated_host_listings_count 0.0988410625
## host_total_listings_count    0.0966659491
## availability_30              0.0960087133
## availability_365             0.0872138255
## review_scores_cleanliness    0.0751878409
## availability_60              0.0679700281
## review_scores_rating         0.0625075131
## availability_90              0.0497193478
## host_active_duration         0.0279624715
## host_is_superhost            0.0268494262
## review_scores_accuracy       0.0206434046
## listing_active_duration      0.0199215914
## review_scores_communication  0.0148947113
## review_scores_checkin        0.0089322672
## minimum_nights               0.0068358449
## id                           0.0025195714
## maximum_nights               -0.0002593775
## review_scores_value          -0.0131239935
## instant_bookable             -0.0189846554
## reviews_per_month            -0.0357953760
## number_of_reviews            -0.0396449764
## number_of_reviews_ltm        -0.0481008297
```

#Multivariate Filter

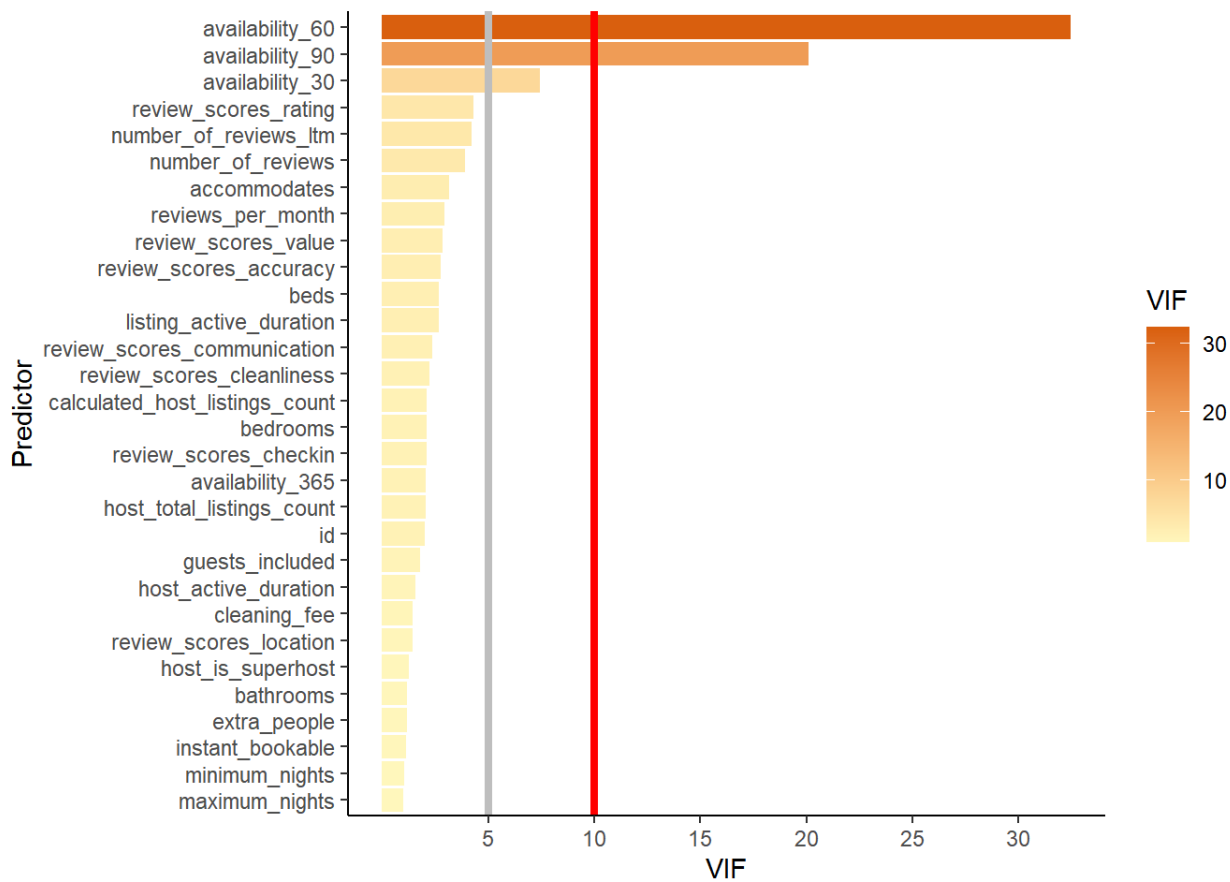
```
model = lm(price~.,data_numVar)
#library(broom)
summary(model) %>%
  tidy()
```

```
## # A tibble: 31 x 5
##   term                estimate std.error statistic  p.value
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)      -96.9      7.05     -13.7  6.92e-43
## 2 id              -0.0000110 0.00000287  -3.84  1.23e- 4
## 3 host_is_superhost  1.12      1.14      0.984  3.25e- 1
## 4 host_total_listings_count 0.0651    0.0121     5.40  6.65e- 8
## 5 accommodates     21.5      0.407     52.8    0
## 6 bathrooms        17.1      1.06     16.1  8.72e-58
## 7 bedrooms         11.2      0.865     12.9  5.65e-38
## 8 beds             -3.99      0.624     -6.38  1.74e-10
## 9 cleaning_fee      0.574      0.00997    57.6    0
## 10 guests_included  2.67      0.478     5.58  2.41e- 8
## # ... with 21 more rows
```

```
#Library(car)
vif(model)
```

```
##          id          host_is_superhost
##      2.010195          1.254679
## host_total_listings_count      accommodates
##      2.049800          3.167454
##      bathrooms          bedrooms
##      1.195990          2.104522
##      beds          cleaning_fee
##      2.669752          1.458843
##      guests_included      extra_people
##      1.776342          1.155236
##      minimum_nights      maximum_nights
##      1.052069          1.000830
##      availability_30      availability_60
##      7.458766          32.455358
##      availability_90      availability_365
##      20.117021          2.052922
##      number_of_reviews      number_of_reviews_ltm
##      3.914190          4.224140
##      review_scores_rating      review_scores_accuracy
##      4.289257          2.772670
##      review_scores_cleanliness      review_scores_checkin
##      2.224826          2.094720
##      review_scores_communication      review_scores_location
##      2.374831          1.449576
##      review_scores_value          instant_bookable
##      2.871729          1.144593
## calculated_host_listings_count      reviews_per_month
##      2.115325          2.959689
##      listing_active_duration      host_active_duration
##      2.662872          1.594275
```

```
data.frame(Predictor = names(vif(model)), VIF = vif(model)) %>%
  ggplot(aes(x=VIF, y = reorder(Predictor, VIF), fill=VIF))+
  geom_col()+
  geom_vline(xintercept=5, color = 'gray', size = 1.5)+
  geom_vline(xintercept = 10, color = 'red', size = 1.5)+
  scale_fill_gradient(low = '#fff7bc', high = '#d95f0e')+
  scale_y_discrete(name = "Predictor")+
  scale_x_continuous(breaks = seq(5,30,5))+
  theme_classic()
```



##Write transformed data to csv file

```
write.csv(data, "clean_full_analysis_data.csv",row.names = F)
```

##Write selective features transformed data to csv file

```
data_trimmed = data[,c('price','zipcode','room_type','bedrooms','accommodates','neighbourhood_group_cleaned','availability_30','host_is_superhost','review_scores_rating','review_scores_location','TV','Elevator','cleaning_fee','property_type','minimum_nights','bathrooms')]

write.csv(data_trimmed, "clean_trimmed_analysis_data.csv",row.names = F)
```

Modelling