# Assignment 1.1

In this assignment we are supposed to analyze two datasets. The first dataset is from the NY Airbnb data and the second one is about international football results.
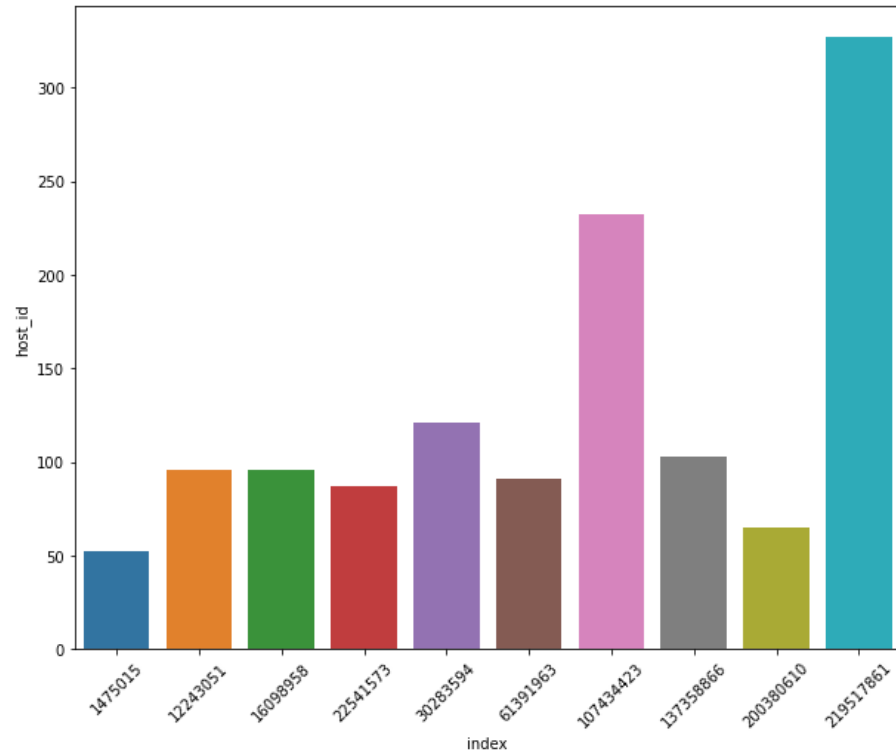
In the first dataset (NY):

We can see that the null data for 'last_review' and 'reviews_per_month' columns are common. And their corresponding 'number_of_reviews' is also zero. So, that means there are no reviews for that instance at all, hence can be filled with zero.

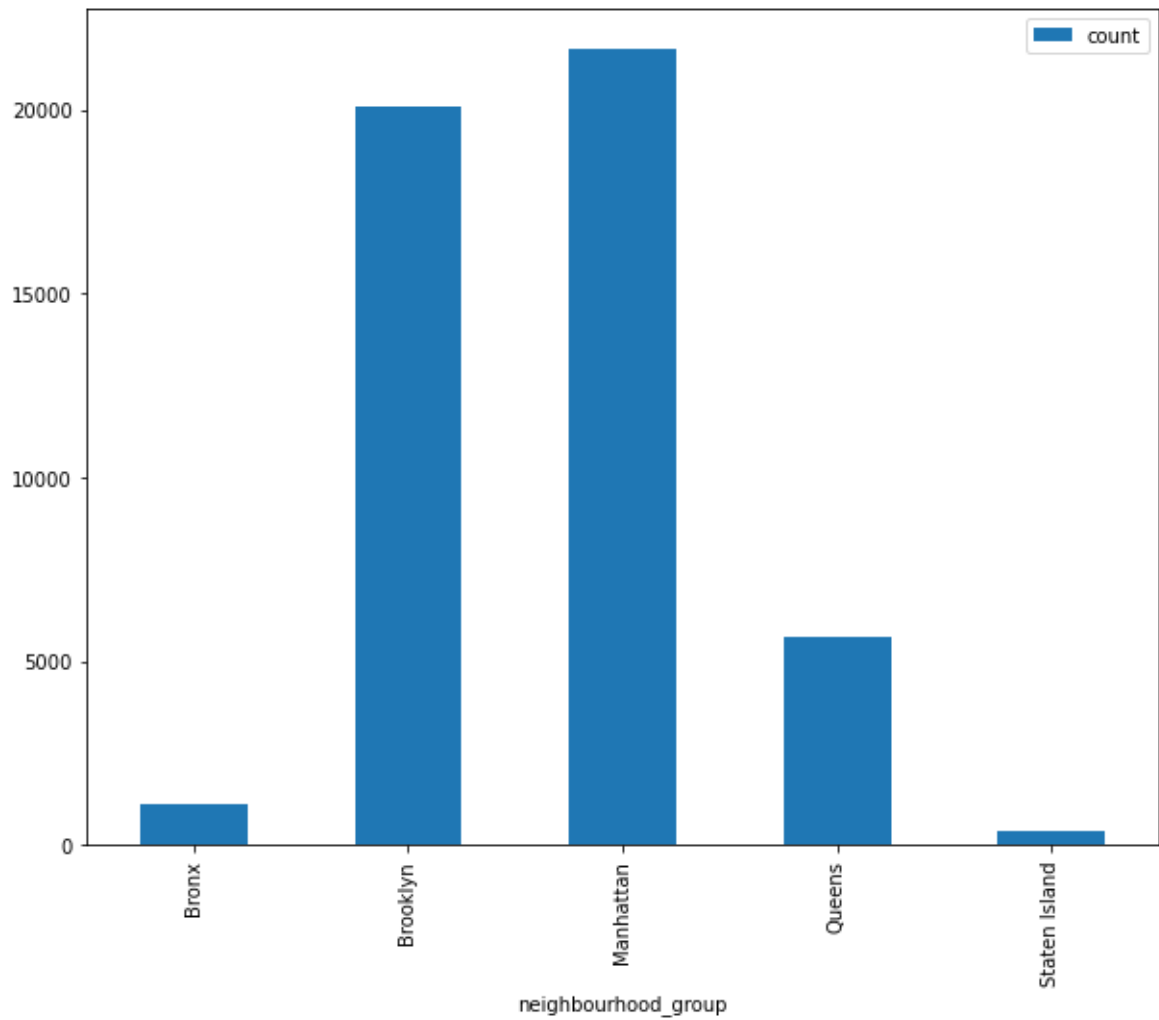We can also see that there are no duplicate instances in the data.

I cleaned some of the columns I didn't find useful, like the id column, and the rows with price of zero.

To get a better understanding of which hosts have the most dedicated datapoints to, I inspect this plot:
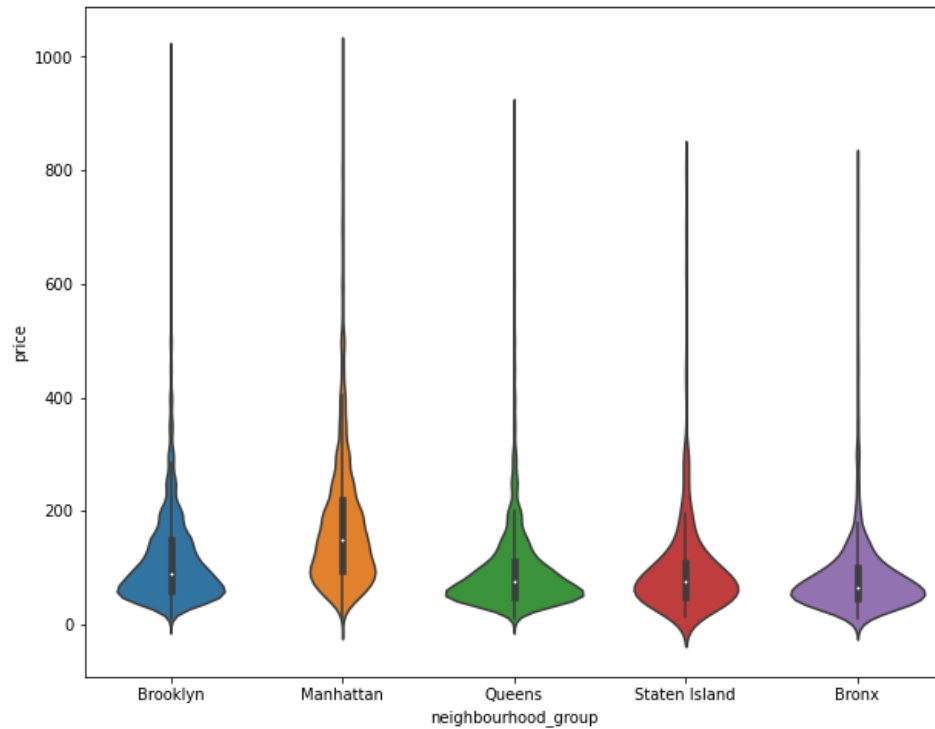
The host with the id no. 219517861 has 327 indexes dedicated to it. The properties of this host are in Manhattan which is a hot spot.
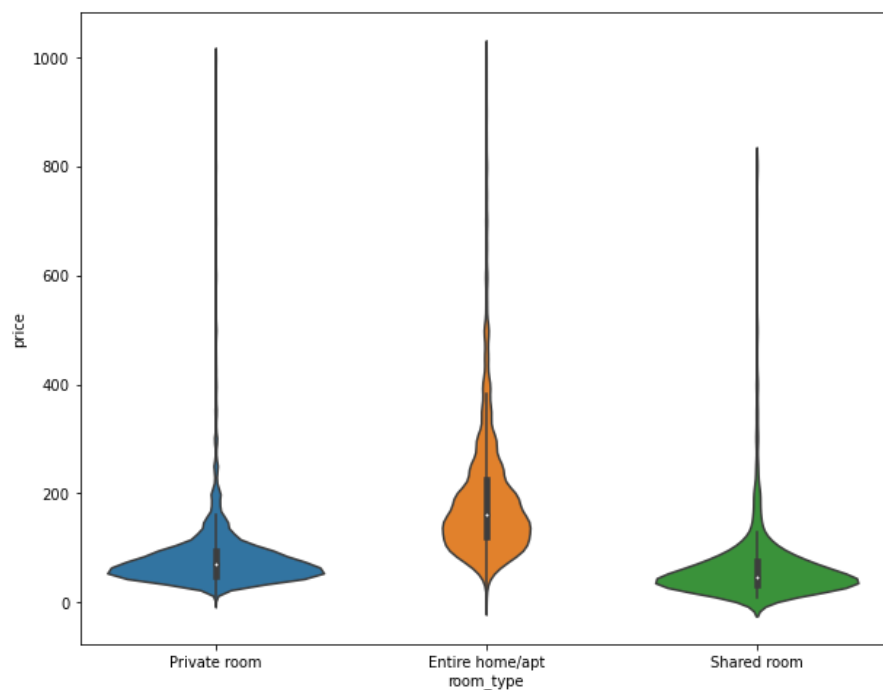
To get a better understanding of the neighborhood data count, lets take a look at this plot.
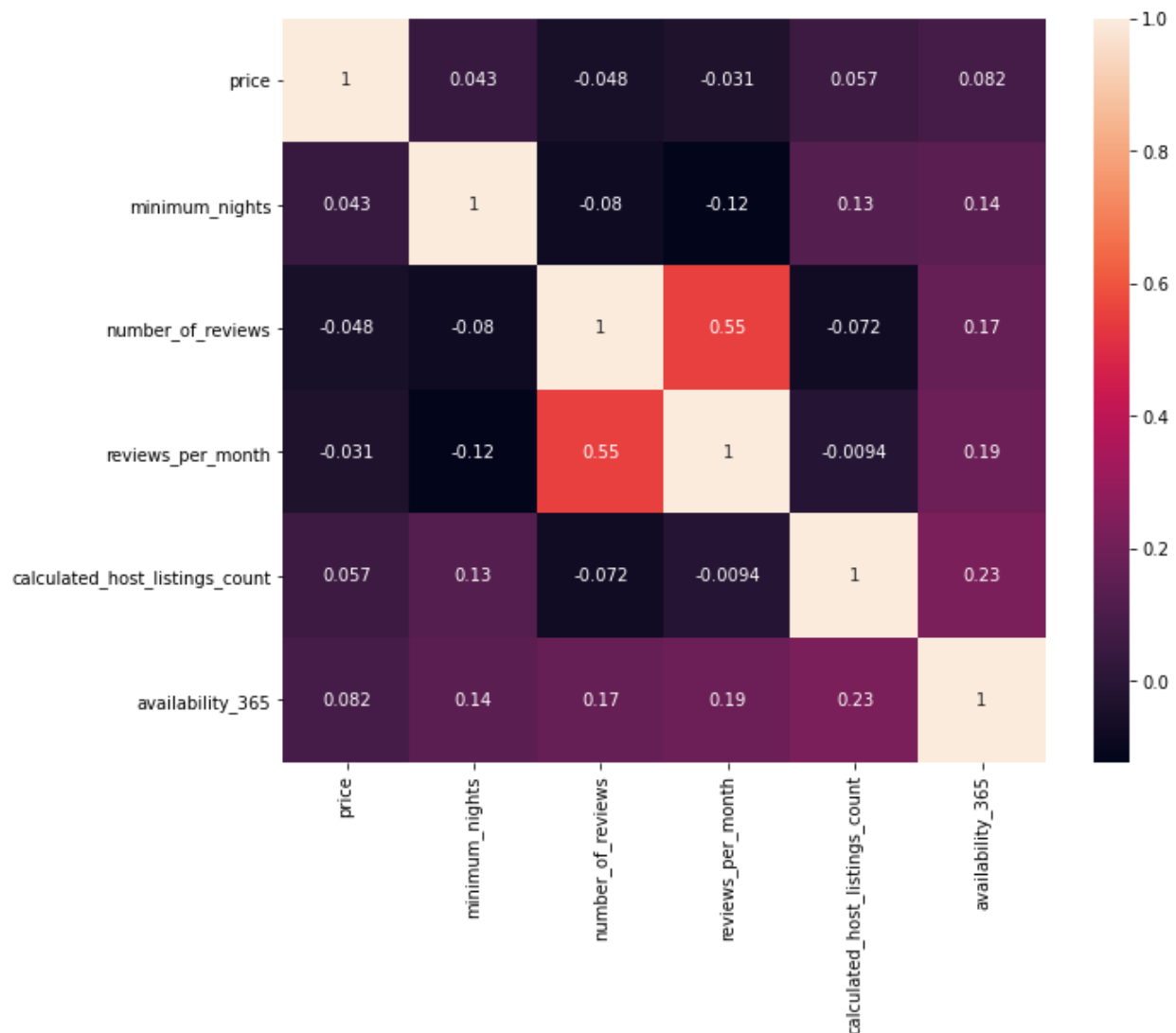
As can be seen in the next plot Clearly Manhattan has higher prices followed by Brooklyn.
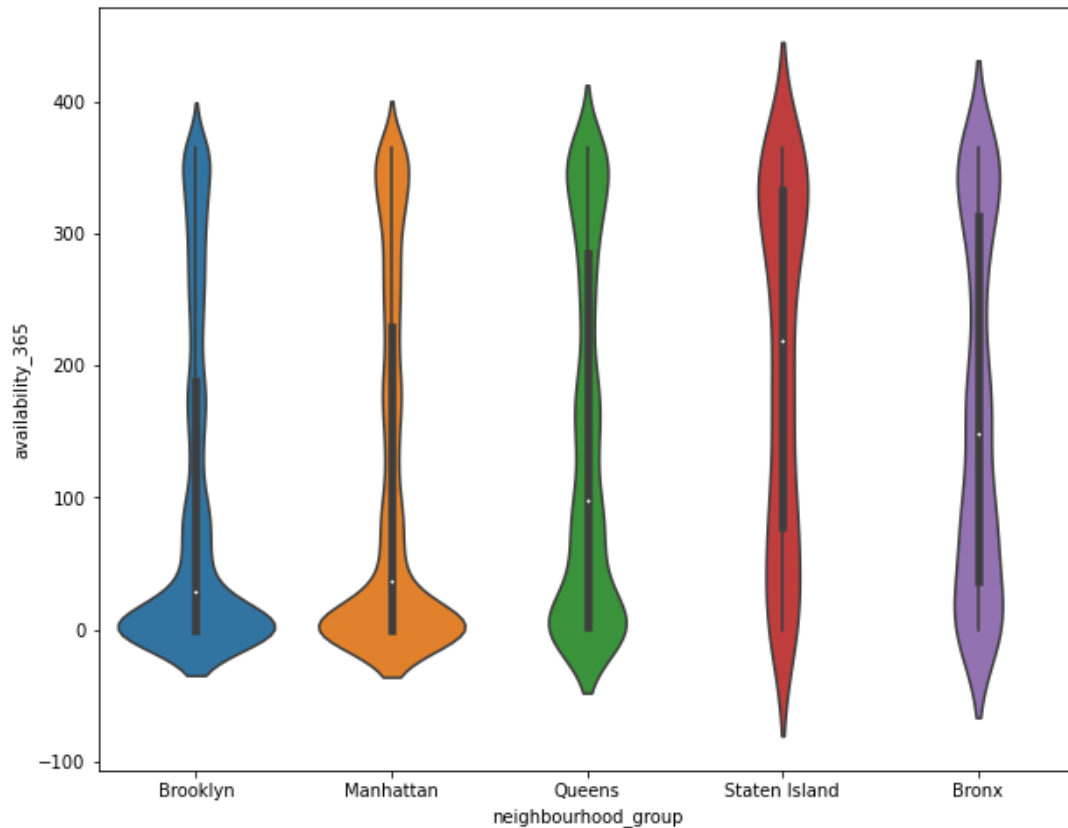
To inspect the room type and its correlation with the price, if we inspect this plot we can get a better understanding that prices of Entire home/apt is more and shared rooms have less cost, which is expected.

I investigate the correlation between 'price', 'minimum_nights', 'number_of_reviews', 'reviews_per_month', 'calculated_host_listings_count', 'availability_365 with the correlation heatmap plot. reviews_per_month and number_of_reviews have high correlation, which was expected



The availability 365 days were inspected for different neighborhoods and as a result, Staten Island has more availability followed by Bronx in a year. This is expected as Manhattan and Brooklyn have more population and has tourist spots.

To analyze the second dataset, I started by adding two columns of winner and loser to the dataset to initialize the results of each match.

The number of most won matches are by these 11 first teams between 300 teams in this dataset:

```
Brazil         635
England        583
Germany        561
Argentina      529
Sweden         508
South Korea    458
Mexico         445
Hungary        443
Italy          435
France         428
Spain          412
```

To get a better view of which teams won how many matches, let's take a look at the winners plot:

To make the interpretability of dates easier, I added a year scale timeline that divides dates of games into old, middle, modern matches.

And also I converted the date column for easier use

I added a Goal column and to understand the number of goals overtime, we can see this plot. This could also be interpreted into the increasing number of footbal matches over time



For better undertanding, lets look at a better plot of goals in decades

Another way to look at the data we have is to inspect home team scores:
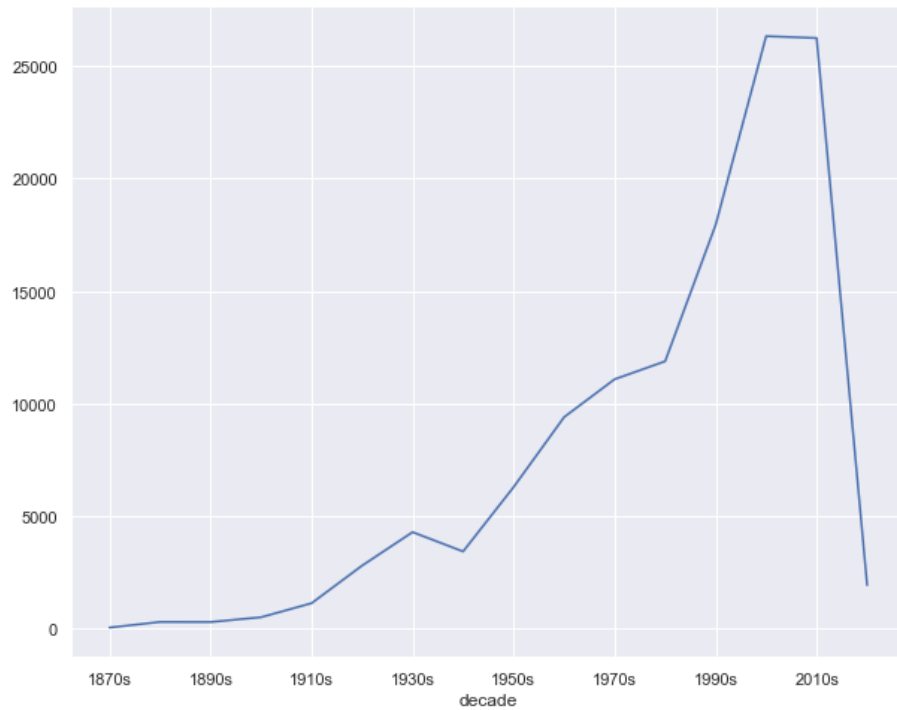
To see the number of goals in each tonament, we can take a look at this plot:



We can see that the number of goals in friendly tournamets are the highest overal.

To get into more detail about how each team did in some of these torunaments:

Based on this plot, Brazil had the most number of goals overal In the FIFA.

## Meanwile in friendly tournamets



Germany has the most number of goals in friendly tournamets while brazil does poorly. Which could be interpreted, since Brazil was on top of the list in other plots, that FIFA games results might play a more significant role for teams.

# References

1)Videos and codes that Mr. Sharifi shared

2)Kaggle.com

3)Stackoverflow.com

4)Medium.com

5)Pandas' documents

6)Matplotlib documents

# Appendix

```python
1.   # -*- coding: utf-8 -*-
2.   """
3.   @author: Sayeh
4.   Spyder editor
5.   """
6.
7.   #%%dataset and package import
8.   import pandas as pd
9.   import numpy as np
10.  import matplotlib.pyplot as plt
11.  import seaborn as sns
12.  from scipy.stats import normaltest
13.  import plotly.express as px
14.
15.  df1 = pd.read_csv(r'C:\Users\Asus\Desktop\Data Mining\01_assignment\1.1\NY
     airbnb\AB_NYC_2019.csv')
16.  df2 = pd.read_csv(r'C:\Users\Asus\Desktop\Data Mining\01_assignment\1.1\International
     footbal results\results.csv')
17.
18.  #%% understanding and cleaning and processing the dataset 1
19.
20.  #to get a better understanding of how our data looks we look at the first 10 rows
21.  df1.shape
22.  df1.head(10)
23.
24.  #lets see some info about this dataset and its features
25.  df1.info()
26.
27.  df1.describe()
28.
29.  df1 = df1.dropna(subset=['name'])
30.  df1 = df1.drop(columns=['host_name'])
31.  #df1 = df1.drop(columns=['id'])
32.  #numeric data info
33.  df1._get_numeric_data().describe()
34.
35.  #null managing
36.  df1.isna().sum()
37.
38.  #no data duplications
39.  df1.duplicated().sum()
40.
41.  #the review null data
42.  null_data = df1[df1['reviews_per_month'].isnull()]
43.  null_data.head()
44.  null_data.shape
45.  null_data['number_of_reviews'].value_counts()
46.  df1 = df1.drop(columns=['last_review'])
47.
48.  #cleaning prices = zero
```

```python
49.  df1 = df1[df1['price']>0]
50.  df1.head()
51.  df1['price'].describe()
52.
53.  #to get a better understanding of data in each column
54.  for cols in df1.columns:
55.      #if df1[cols].dtype == 'object' or df1[cols].dtype == 'bool':
56.          print('column : ',cols)
57.          print(df1[cols].value_counts().head())
58.
59.  # to better understand the host ids data
60.  hosts = df1['host_id'].value_counts().reset_index()[:10]
61.  plt.rcParams['figure.figsize'] = (10,8)
62.  ax = sns.barplot(x='index',y='host_id',data=hosts)
63.  ax.set_xticklabels(ax.get_xticklabels(), rotation=45)
64.  plt.show()
65.
66.  #neighbourhood counts
67.  df1.groupby('neighbourhood_group')['id'].agg(['count']).plot(kind="bar")
68.
69.  #neighbourhood vs the prices
70.  sample_data = df1[df1['price']<1000]
71.  sns.violinplot(x='neighbourhood_group', y='price', data=sample_data)
72.  plt.show()
73.
74.  #Room type inspection vs price
75.  df1['room_type'].value_counts()
76.
77.  sns.violinplot(x='room_type', y='price', data=sample_data)
78.  plt.show()
79.
80.  room = df1.groupby('room_type')['price']
81.  private = room.get_group('Private room')
82.  entire = room.get_group('Entire home/apt')
83.  shared = room.get_group('Shared room')
84.
85.  #correlation of price and
86.  sns.heatmap(df1[['price','minimum_nights','number_of_reviews','reviews_per_month','calculate
     d_host_listings_count','availability_365']].corr(), annot=True)
87.  plt.show()
88.  for i in [private,entire,shared]:
89.      print(normaltest(i))
90.
91.  #availability365 of neighborhoods
92.  sns.violinplot(x='neighbourhood_group',y='availability_365', data=df1)
93.
94.  #End of inspections and investigatingg codes of dataframe 1
95.
96.  #%%lets beggin analyzing the second dataset
97.
98.  df2.shape
99.  df2.head(10)
100.
101.  #lets see some info about this dataset and its features
102.  df2.info()
103.
104.  df2.describe()
105.
106.  #introducing a winner column to dataser
107.  def winner(row):
108.      if row['home_score'] > row['away_score']: return row['home_team']
109.      elif row['home_score'] < row['away_score']: return row['away_team']
110.      else: return 'DRAW'
111.
112.  df2['winner'] = df2.apply(lambda row: winner(row), axis=1)
```

```
113.  df2.head()
114.
115.  #introducing a loser column to our dataset
116.  def loser(row):
117.      if row['home_score'] < row['away_score']: return row['home_team']
118.      elif row['home_score'] > row['away_score']: return row['away_team']
119.      else: return 'DRAW'
120.
121.  df2['loser'] = df2.apply(lambda row: loser(row), axis=1)
122.  df2.head()
123.
124.  #lets see what the winners look like
125.  winners = pd.value_counts(df2.winner)
126.  winners = winners.drop('DRAW')
127.  winners.head(11)
128.
129.  #a plot for the winners
130.  fig, ax = plt.subplots(figsize=(25, 100))
131.  sns.set(font_scale=1)
132.  sns.barplot(y = winners.index.tolist(), x = winners.tolist())
133.
134.  #making a yearscale of matches from old to new
135.  df2 = df2.dropna(axis=0)
136.  for i in range(len(df2)):
137.      if df2.iat[i, 0].find('-') != -1:
138.          df2.iat[i, 0] = int(df2.iat[i, 0][:str(df2.iat[i,0]).find('-')])
139.  df2['Year Scale'] = df2['date'].apply(lambda x:(('Old Match(<1976)' , 'Middling
      Match(1976<..<2000)')[x > 1976],'Modern Match(>2000)')[x > 2000])
140.
141.  #changing dates a bit, adding years and decades column
142.  def give_date(x):
143.      t = pd.to_datetime(x['date'])
144.      x['year'] = t.year
145.      x['month'] = t.month
146.      x['day_of_week'] = t.dayofweek
147.      return x
148.
149.  df2 = df2.apply(give_date, axis = 1)
150.
151.  decades = 10 * (df2['year'] // 10)
152.  decades = decades.astype('str') + 's'
153.  df2['decade'] = decades
154.  df2.head()
155.  #
156.  fig = px.scatter(df2, x='home_team', y='away_team',color='tournament',height=1700)
157.  fig.show()
158.
159.  #
160.  x = df2.groupby('home_team')[['home_score','away_score']].sum()
161.  x['total'] = x['home_score'] + x['away_score']
162.  x = x.sort_values('total', ascending = False)
163.  x = x[:10]
164.  x.plot(kind = 'barh')
165.
166.  #goals in years
167.  df2['Goals'] = df2['home_score']+df2['away_score']
168.  x = df2.groupby('decade')['Goals'].sum()
169.  print(x)
170.  x.plot()
171.
172.  #tournamets and goals
173.  x = df2.pivot_table('Goals', index = 'tournament', aggfunc = 'sum')
174.  x = x.sort_values('Goals', ascending  = False)
175.  x = x[:10]
176.  print(x)
```

```
177.    x.plot(kind = 'barh')
178.
179.    #FIFA cup and team performance
180.    x = df2.pivot_table('Goals', index = 'home_team', columns = 'tournament', aggfunc = 'sum',
        fill_value = 0, margins = True, margins_name = 'Total')
181.    x = x.sort_values('Total', ascending = False)
182.    x = x[:20]
183.    x['FIFA World Cup'][1:].plot(kind = 'barh', title = 'FIFA World Cup')
184.
185.    #Friendly matches and teams
186.    x['Friendly'][1:].plot(kind = 'barh', title = 'Friendly')
187.
```