

✓ [35 pts] Visualizing Data

✓ [5 pts] Load the data + statistics

- load the dataset
- display the first few rows of the data

```
1 import pandas as pd
2 data = pd.read_csv("https://drive.google.com/uc?id=1Wa5ixqnsckyQQhb1Ap4wl37mSgnPI12F&export=download")
3 data.head()
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	lat
0	2539	Clean & quiet apt home by the park	2787	John	Brooklyn	Kensington	4
1	2595	Skyliit Midtown Castle	2845	Jennifer	Manhattan	Midtown	4
2	3647	THE VILLAGE OF HARLEM....NEW YORK !	4632	Elisabeth	Manhattan	Harlem	4
3	3831	Cozy Entire Floor of Brownstone	4869	LisaRoxanne	Brooklyn	Clinton Hill	4
4	5022	Entire Apt: Spacious Studio/Loft by central park	7192	Laura	Manhattan	East Harlem	4

Next steps: [Generate code with data](#) [View recommended plots](#)

- pull up info on the data type for each of the data fields. Will any of these be problematic feeding into your model (you may need to do a little research on this)? Discuss:

Something about the names being non-numeric with thousands of possibilities, none of which are relevant to this analysis, since we are not analyzing something like the most popular baby names.

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48895 entries, 0 to 48894
Data columns (total 16 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   id                    48895 non-null  int64
 1   name                  48879 non-null  object
 2   host_id               48895 non-null  int64
 3   host_name             48874 non-null  object
 4   neighbourhood_group    48895 non-null  object
 5   neighbourhood         48895 non-null  object
 6   latitude              48895 non-null  float64
 7   longitude             48895 non-null  float64
 8   room_type             48895 non-null  object
 9   price                 48895 non-null  int64
10  minimum_nights        48895 non-null  int64
11  number_of_reviews     48895 non-null  int64
12  last_review           38843 non-null  object
13  reviews_per_month     38843 non-null  float64
14  calculated_host_listings_count  48895 non-null  int64
15  availability_365       48895 non-null  int64
dtypes: float64(3), int64(7), object(6)
memory usage: 6.0+ MB
```

All the data fields with data type "object" (name, host_name, neighbourhood, neighborhood_group, room_type, last_review) could be problematic later on. Because their data types are unclear to us. It could possibly mean that the data type is messed up in some fields, or there are some

empty values.

Double-click (or enter) to edit

- drop the following columns: name, host_id, host_name, last_review, and reviews_per_month
- display a summary of the statistics of the loaded data

```
1 data_new = data.drop(columns = ['name', 'host_id', 'host_name', 'last_review', 'reviews_per_month'])
```

```
1 data_new.describe()
```

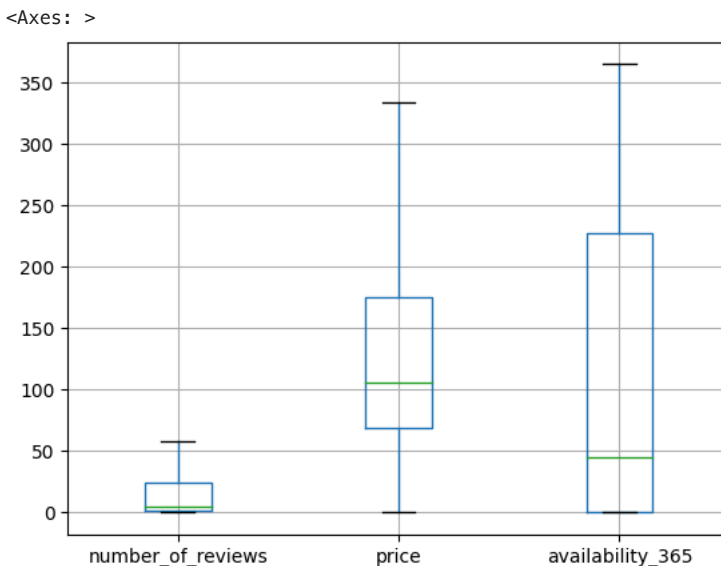


	id	latitude	longitude	price	minimum_nights	number_o
count	4.889500e+04	48895.000000	48895.000000	48895.000000	48895.000000	48
mean	1.901714e+07	40.728949	-73.952170	152.720687	7.029962	
std	1.098311e+07	0.054530	0.046157	240.154170	20.510550	
min	2.539000e+03	40.499790	-74.244420	0.000000	1.000000	
25%	9.471945e+06	40.690100	-73.983070	69.000000	1.000000	
50%	1.967728e+07	40.723070	-73.955680	106.000000	3.000000	
75%	2.915218e+07	40.763115	-73.936275	175.000000	5.000000	
max	3.648724e+07	40.913060	-73.712990	10000.000000	1250.000000	

✓ [5 pts] Boxplot 3 features of your choice

- plot [boxplots](#) for 3 features of your choice

```
1 import matplotlib
2 import matplotlib.pyplot as plt
3 data_new.boxplot(column = ['number_of_reviews', 'price', 'availability_365'], showfliers = False)
```



- describe what you expected to see with these features and what you actually observed

This question should be open-ended, such as some observations, e.g. The range varies a lot, plus, long tail distribution, etc.

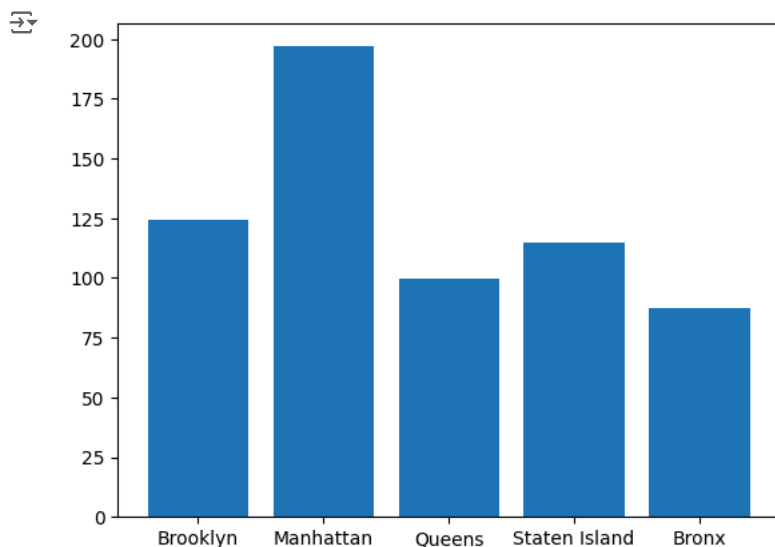
Looks like a lot of the other tickers don't have the best variability in comparison to availability and price.

✓ [10 pts] Plot average price of a listing per neighbourhood_group

```

1 grouped = data_new.groupby('neighbourhood_group')
2
3 temp = dict()
4
5 for ng in data_new['neighbourhood_group'].unique():
6     temp[ng] = data_new[data_new['neighbourhood_group'] == ng]['price'].mean()
7
8 plt.bar(temp.keys(), temp.values())
9 plt.show()

```



- describe what you expected to see with these features and what you actually observed

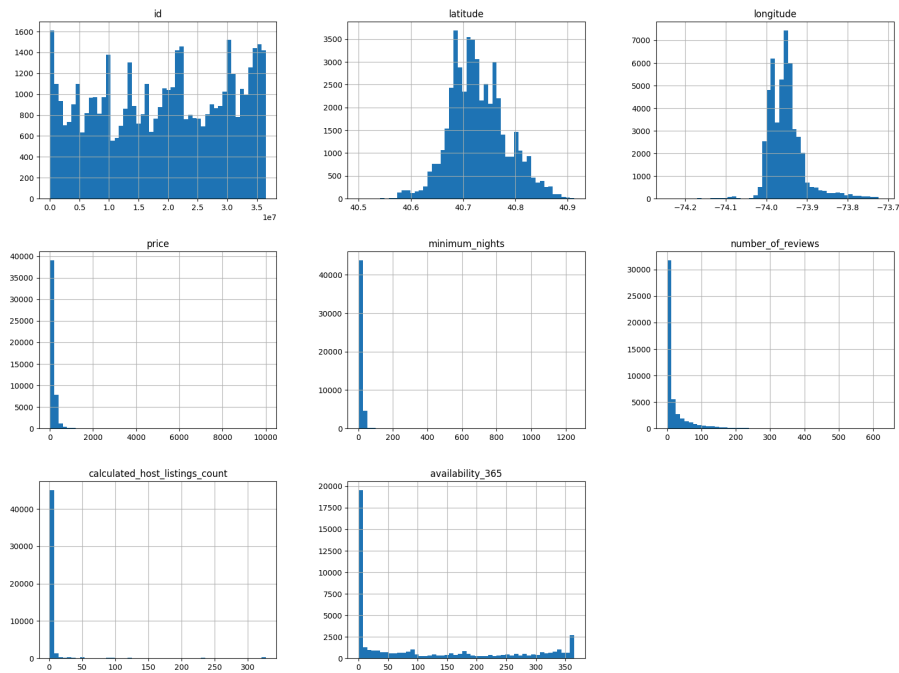
This question should be open-ended, they can talk about things such as "expecting the rooms with higher price to be have fewer reviews / more days available, but this is not the case". Or "location affects the price" and indeed it does.

- So we can see different neighborhoods have dramatically different pricepoints, but how does the price breakdown by range. To see let's do a histogram of price by neighborhood to get a better sense of the distribution.

```

1 import numpy as np
2
3 # for neighborhood in set(data_new["neighbourhood_group"]):
4 #     data_tmp = data_new[data_new["neighbourhood_group"] == neighborhood]
5 #     data_tmp = data_tmp[~data_tmp["price"].isna()]
6 #     n_chunks = 10
7 #     bins = list()
8 #     labels = list()
9 #     granularity = 50
10 #     for i in range(n_chunks):
11 #         bins.append(i * granularity)
12 #         labels.append(i)
13 #     bins.append(np.inf)
14 #     data_tmp["price_cut"] = pd.cut(data_tmp["price"],
15 #                                     bins=bins,
16 #                                     labels=labels)
17
18 #     data_tmp["price_cut"].hist()
19 #     plt.title(neighborhood)
20 #     plt.show()
21
22
23 data_new.hist(bins=50, figsize=(20,15))
24
25 plt.show()
26
27

```

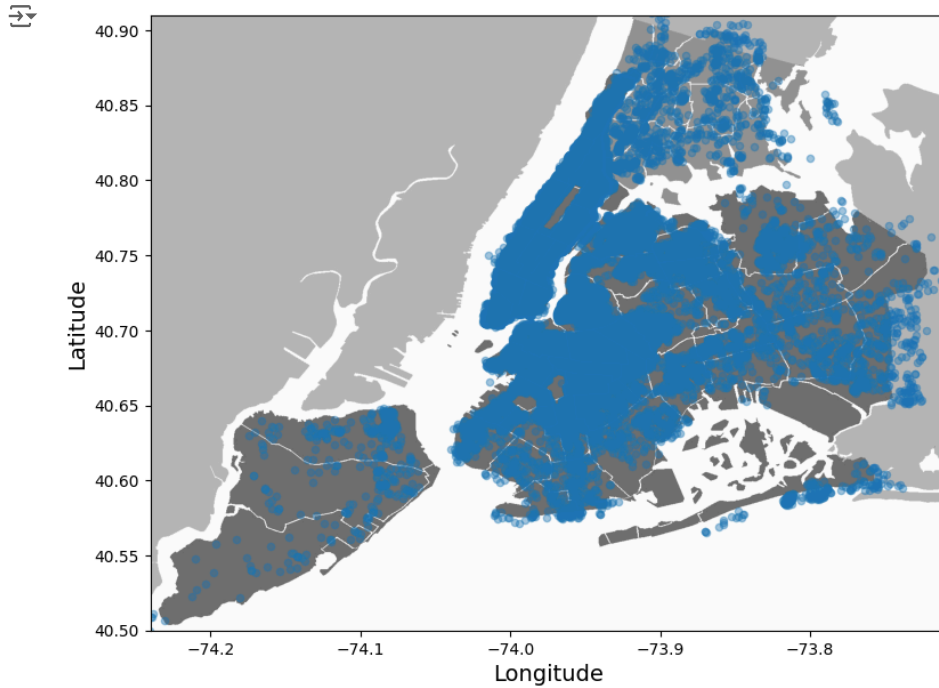


- ✓ [5 pts] Plot map of airbnbs throughout New York (if it gets too crowded take a subset of the data, and try to make it look nice if you can :)).

```

1 images_path = os.path.join('.', "images")
2 os.makedirs(images_path, exist_ok=True)
3 filename = "newyork.png"
4
5 import matplotlib.image as mpimg
6 newyork_img=mpimg.imread(os.path.join(images_path, filename))
7 ax = data.plot(kind="scatter", x="longitude", y="latitude", figsize=(10,7), alpha=0.4)
8 plt.imshow(newyork_img, extent=[-74.24, -73.71, 40.50, 40.91])
9 plt.ylabel("Latitude", fontsize=14)
10 plt.xlabel("Longitude", fontsize=14)
11
12 plt.show()
13
14 # This one was done for me it seems?
15

```

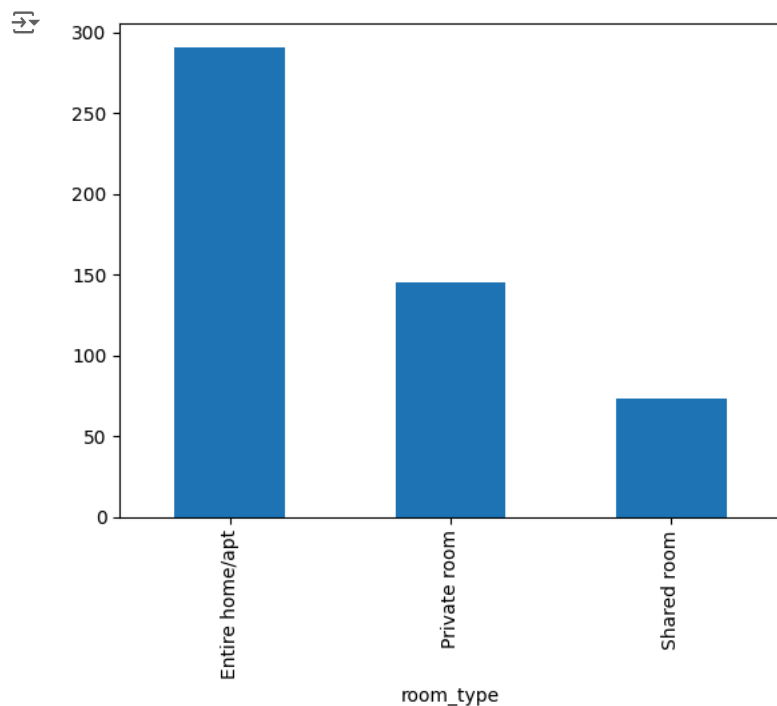


✓ [5 pts] Plot average price of room types who have availability greater than 180 days and neighbourhood_group is Manhattan

```

1 data_temp = data_new[data_new["neighbourhood_group"] == "Manhattan"]
2 data_temp_180 = data_temp[data_temp["availability_365"] > 180].groupby('room_type').mean(numeric_only=True)
3 data_temp_180["price"].plot(kind="bar")
4 plt.show()

```



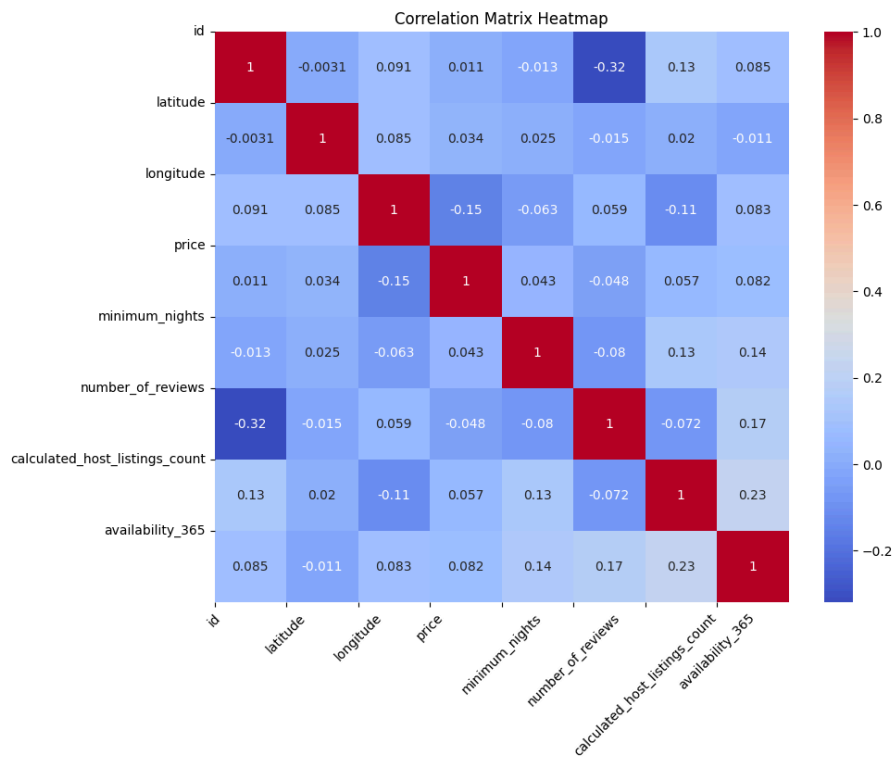
✓ [5 pts] Plot correlation matrix

- which features have positive correlation?
- which features have negative correlation?

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns # Optional, for better visualization
3
4 # Assuming data_new is your DataFrame and numeric_only=True filters numeric columns
5 corr_matrix = data_new.corr(numeric_only=True)
6
7 # Sort correlations related to "price"
8 sorted_correlations = corr_matrix["price"].sort_values(ascending=False)
9
10 # Plotting the correlation matrix heatmap
11 plt.figure(figsize=(10, 8)) # Adjust size as needed
12 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm') # Annotate cells with correlation values
13
14 # Customize labels and ticks
15 plt.xticks(range(len(corr_matrix.columns)), corr_matrix.columns, rotation=45)
16 plt.yticks(range(len(corr_matrix.columns)), corr_matrix.columns)
17
18 plt.title('Correlation Matrix Heatmap')
19 plt.show()

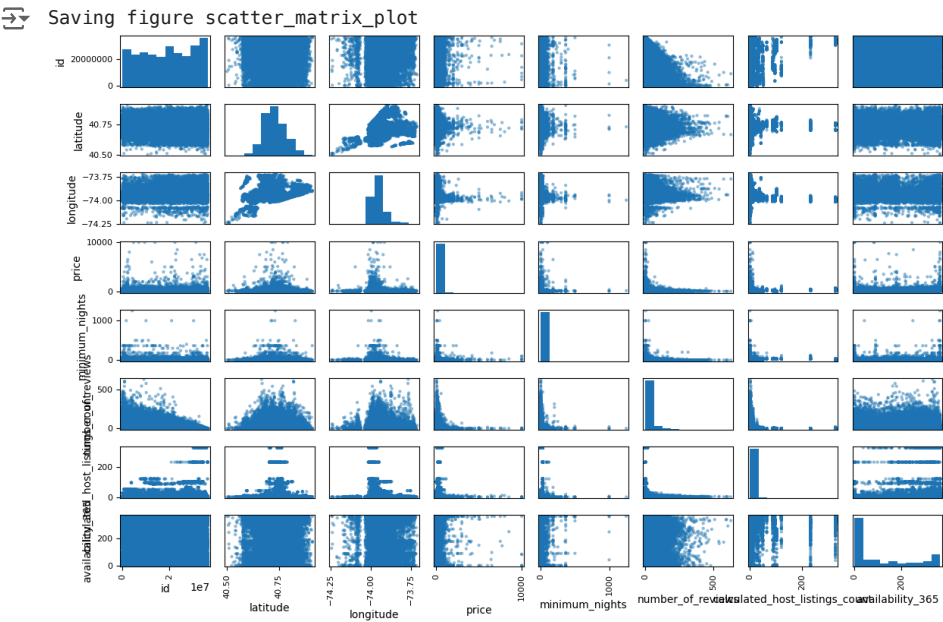
```



```

1 from pandas.plotting import scatter_matrix
2
3 scatter_matrix(data_new, figsize=(12, 8))
4 save_fig("scatter_matrix_plot")

```



(open-ended question, any findings should work)

▾ [30 pts] Prepare the Data

▾ [5 pts] Augment the dataframe with two other features which you think would be useful

Hint: Do Not Use Label (price) to Augment Features

```
1 new_data = data.drop(columns = ['name', 'host_id', 'host_name', 'last_review'],axis = 1)
2 new_data.head()
```

	id	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	reviews_
0	2539	Brooklyn	Kensington	40.64749	-73.97237	Private room	149	1	9	
1	2595	Manhattan	Midtown	40.75362	-73.98377	Entire home/apt	225	1	45	
2	3647	Manhattan	Harlem	40.80902	-73.94190	Private room	150	3	0	
3	3684	Manhattan	Chelsea	40.75514	-73.99570	Entire home/apt	225	1	45	

Next steps:

Generate code with new_data

View recommended plots


```
1 data_new["average_num_review_per_listing"] = data_new["number_of_reviews"] / data_new["calculated_host_listings_count"]
2 data_new['distance_from_center'] = np.sqrt(data_new['latitude']**2 + data_new['longitude']**2)
```

✓ [5 pts] Impute any missing feature with a method of your choice, and briefly discuss why you chose this imputation method

```
1 data_new.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48895 entries, 0 to 48894
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   id                    48895 non-null  int64
 1   neighbourhood_group    48895 non-null  object
 2   neighbourhood          48895 non-null  object
 3   latitude              48895 non-null  float64
 4   longitude             48895 non-null  float64
 5   room_type             48895 non-null  object
 6   price                 48895 non-null  int64
 7   minimum_nights        48895 non-null  int64
 8   number_of_reviews      48895 non-null  int64
 9   calculated_host_listings_count  48895 non-null  int64
10  availability_365       48895 non-null  int64
11  average_num_review_per_listing  48895 non-null  float64
12  distance_from_center   48895 non-null  float64
dtypes: float64(4), int64(6), object(3)
memory usage: 4.8+ MB
```

Had the most NA values in the dataset. Most likely useful if imputed.

✓ [15 pts] Data Preparation

Either using sklearn mixins, or anyway to create ndarray from scratch, is fine.

As long as a feature matrix pf shape (N, D) and a label vector of shape (N,) is prepared

```
1 import numpy as np
2
3 features = ["neighbourhood", "neighbourhood_group", "room_type"]
4
5 from sklearn import preprocessing
6 le = preprocessing.LabelEncoder()
7 for i in features:
8     le.fit(data_new[i])
9     data_new[i]=le.transform(data_new[i])
10
11 data_prepared = data_new.drop("price", axis=1).to_numpy()
12
13 print(data_prepared.shape)
14 print(data_new['price'].shape)
```

```
(48895, 12)
(48895,)
```

✓ [5 pts] Set aside 20% of the data as test test (80% train, 20% test).

```
1 from sklearn.model_selection import train_test_split
2 data_target = data_new["price"]
3 train, test, target, target_test = train_test_split(data_prepared, data_target, test_size=0.2, random_state=0)
```

✓ [15 pts] Fit a model of your choice

The task is to predict the price, you could refer to the housing example on how to train and evaluate your model using MSE. Provide both test and train set MSE values.

```
1 from sklearn.linear_model import LinearRegression
2
3 lin_reg = LinearRegression()
4 lin_reg.fit(train, target)
5
6 data = test
7 labels = target_test
```