

### Group Details:

- Name: Ashley Assous ID: 328956081
- Name: Yonatan Vinograd ID: 211690144
- Name: Aya Talmon ID: 208325712

## > Helper Functions and Imports

[Show code](#)

## ✓ Introduction to Data Science - Lab #2

### Exploratory Data Analysis

## ✓ Case Study: Rental Listings in Jerusalem

In this lab we will practice our exploratory data analysis skills using real data!

We will explore data of rental pricings in Jersualem. The dataset consists of listings published in <https://www.komo.co.il/> during the summer of 2022.

We will use two python packages for visualizing the data: `matplotlib` (and specifically its submodule `pypplot` imported here as `plt`) and `seaborn` (imported as `sns`). Seaborn is a package that "wraps" `matplotlib` and introduces more convenient functions for quickly creating standard visualizations based on dataframes.


Please **briefly** go over this [quick start guide](#) to `matplotlib`, the [first](#) `seaborn` introduction page until the "Multivariate views on complex datasets" section (not included), and the [second](#) introduction page until the "Combining multiple views on the data" section.

## ✓ Loading the dataset

```
#@title Loading the dataset
rent_df = load_df(RENT_ID)[['propertyID', 'neighborhood', 'monthlyRate', 'mefarsem', 'rooms'],
rent_df = rent_df.drop_duplicates(subset='propertyID').reset_index(drop=True)
rent_df_backup_for_exercise = rent_df.copy()
clean_df_area_filtered = None
clean_df = None
```

Let's print a random sample:


```
np.random.seed(2)
rent_df.sample(5)
```



	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	entry
403	3981729	גבעת מרדכי	4500.0	private	3.0	6.0	62.0	10/08/2022
457	3991612	קריית משה	3000.0	private	3.5	3.0	NaN	NaN

And print some summary statistics:

```
rent_df.describe(include='all')
```



	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	
count	6.120000e+02	612	612.000000	612	612.000000	611.000000	29
unique	NaN	54	NaN	2	NaN	NaN	
top	NaN	קריית יובל	NaN	private	NaN	NaN	
freq	NaN	66	NaN	600	NaN	NaN	
mean	3.981582e+06	NaN	4717.393791	NaN	2.927288	1.916530	8
std	6.525543e+04	NaN	2195.215139	NaN	1.007350	1.581006	27
min	2.494041e+06	NaN	0.000000	NaN	1.000000	-2.000000	
25%	3.981694e+06	NaN	3500.000000	NaN	2.000000	1.000000	4
50%	3.987901e+06	NaN	4400.000000	NaN	3.000000	2.000000	6
75%	3.992605e+06	NaN	5800.000000	NaN	3.500000	3.000000	8

The variables we will focus on are:

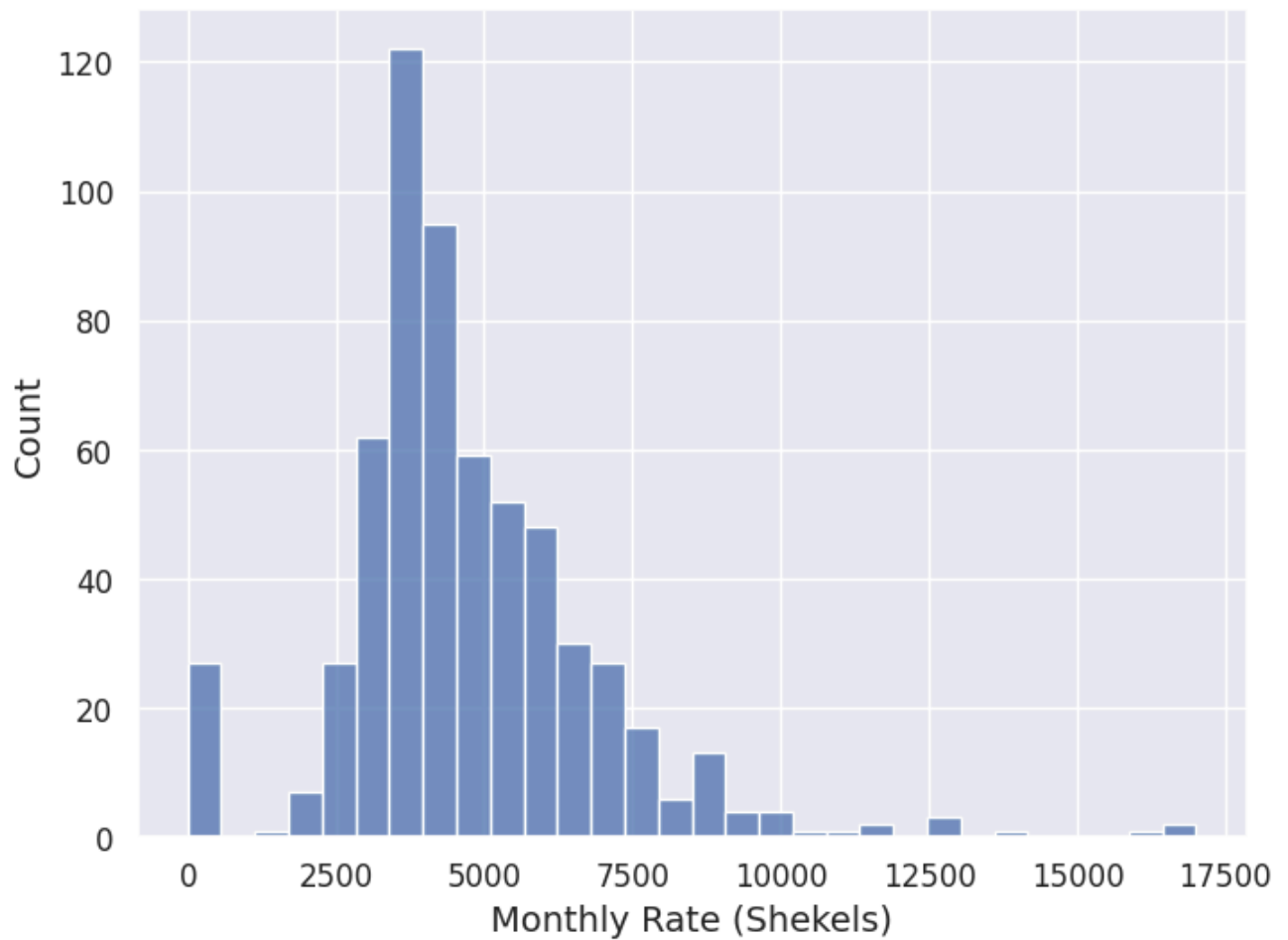
1. neighborhood: The hebrew name of the neighborhood in jerusalem where the listing is located
2. monthlyRate: The monthly rate (שכר דירה) in shekels
3. rooms: The number of rooms in the apartment
4. floor: The floor in which the apartment is located
5. area: The area of the apartment in squared meters
6. numFloors: The total number of floors in the building

## What is the distribution of prices in this dataset?

Q: Plot a histogram with 30 bins of the monthly rates in this dataset:

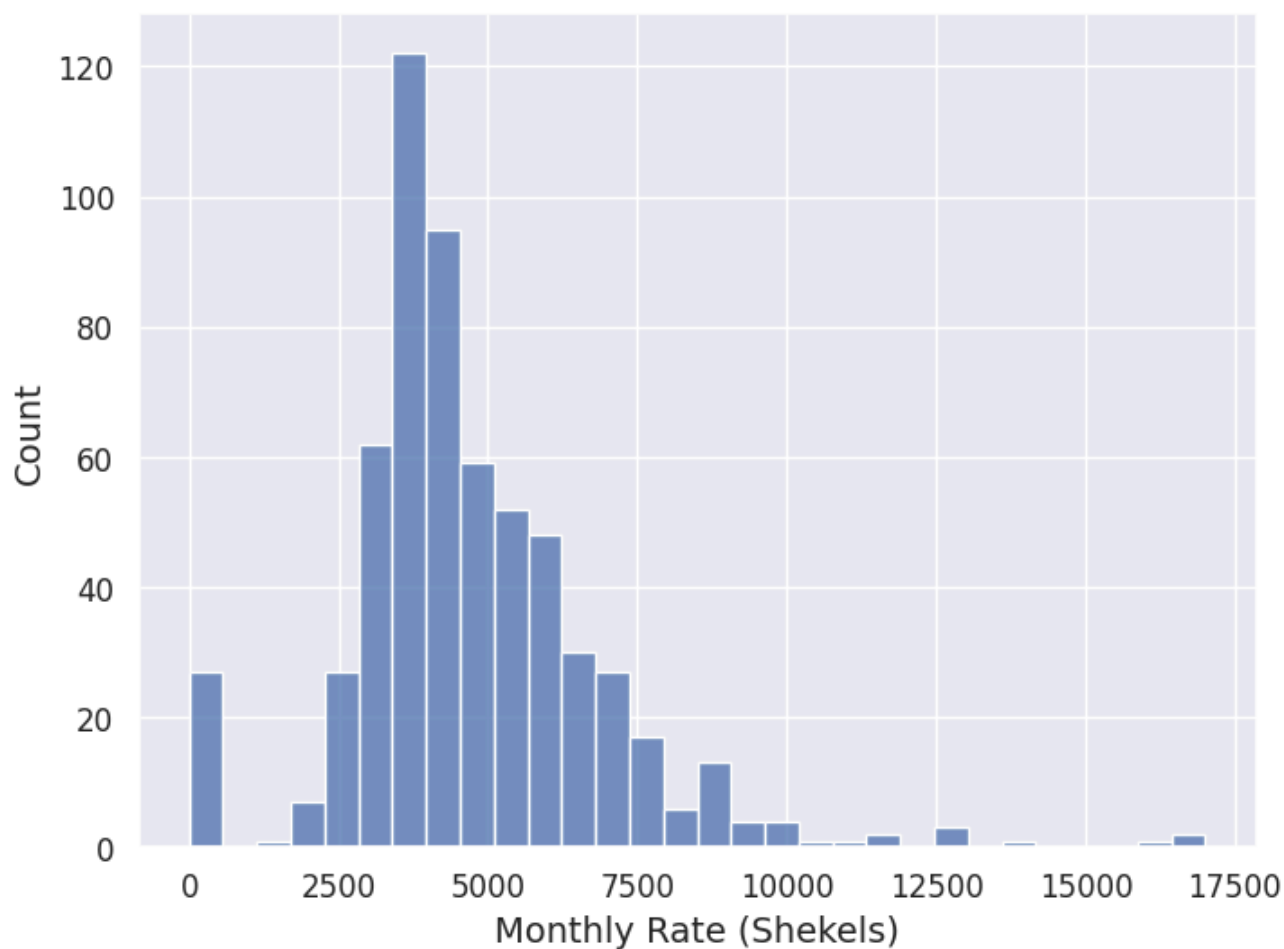
### ✓ Solution 1

```
# @title Solution 1
plt.figure(figsize=(8,6))
sns.histplot(rent_df['monthlyRate'], bins=30)
plt.xlabel("Monthly Rate (Shekels)");
```



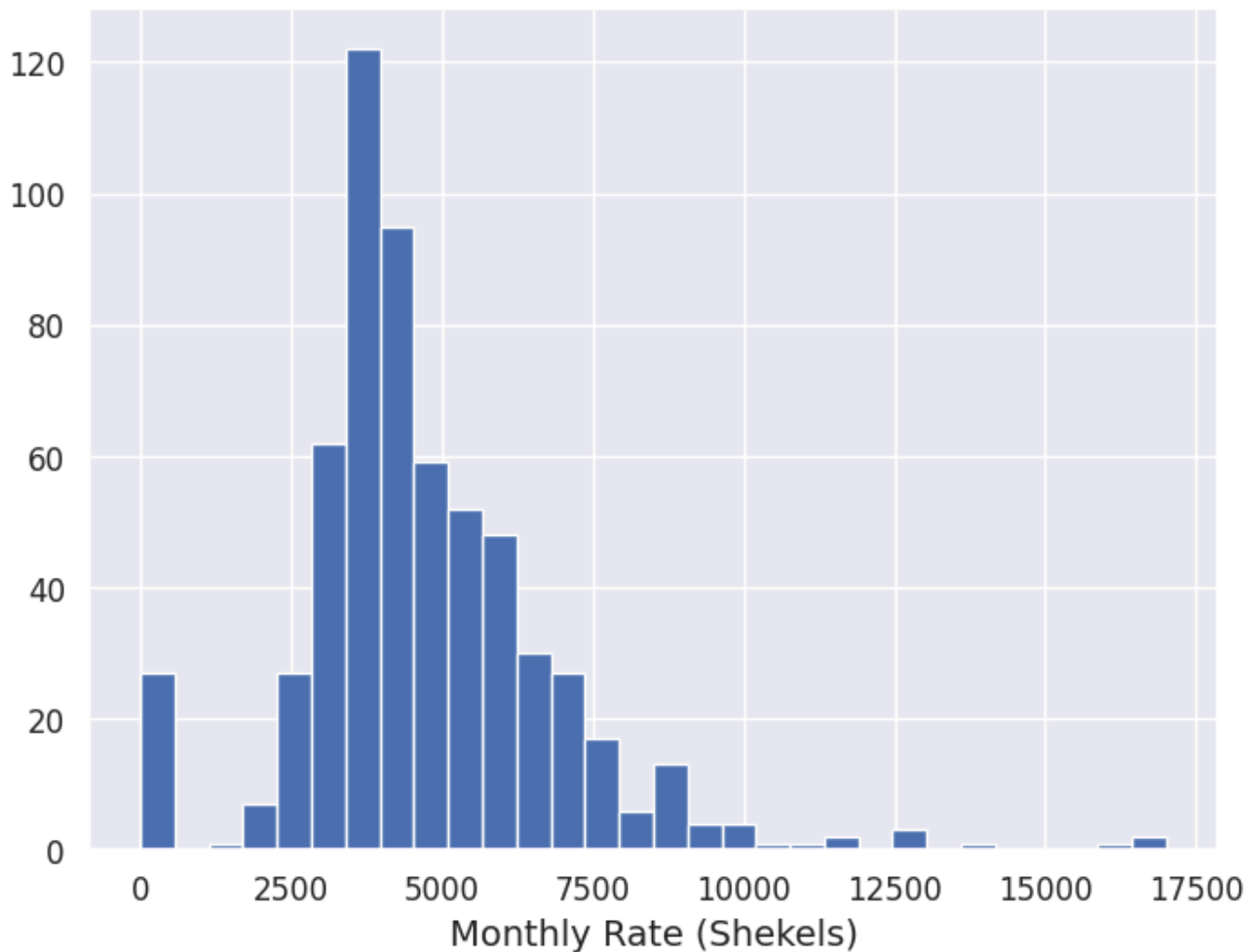
### ✓ Solution 1

```
# @title Solution 1
plt.figure(figsize=(8,6))
sns.histplot(rent_df['monthlyRate'], bins=30)
plt.xlabel("Monthly Rate (Shekels)");
```



## ✓ Solution 2

```
# @title Solution 2
rent_df["monthlyRate"].hist(bins=30, figsize=(8,6))
plt.xlabel("Monthly Rate (Shekels)");
```



We see that the prices distribution peaks around ~3500 Shekels and that it is right skewed, as there are some very expensive apartments. We can also see a peak at zero which makes sense as sometimes listings do not include a price. We would want to filter those out when we analyze prices later on.

Q: Print the number of listings that have no monthly rate:

## ✓ Solution

```
# @title Solution
print("Number of apartments without a price: ", rent_df['monthlyRate'].value_counts()[0].
```

➞ Number of apartments without a price: 25

We want to remove those listings, but we don't want to lose these entries, as we might want to know how many and what type of outliers we originally removed. So we create another dataframe that has the listings we removed and the reason for removal.

```
outlier_df = pd.DataFrame(columns=rent_df.columns.to_list()+['reason']) # will save the o

outliers = rent_df[rent_df['monthlyRate'] <= 0].reset_index(drop=True)
outliers['reason'] = "monthlyRate <= 0"
outlier_df = pd.concat([outlier_df, outliers], axis=0, ignore_index=True).drop_duplicates
outlier_df.tail()
```

➞ [Show hidden output](#)

We will now remove those listings and save the result to a new variable `clean_df`:

```
clean_df = rent_df[rent_df['monthlyRate'] > 0].reset_index(drop=True)
```

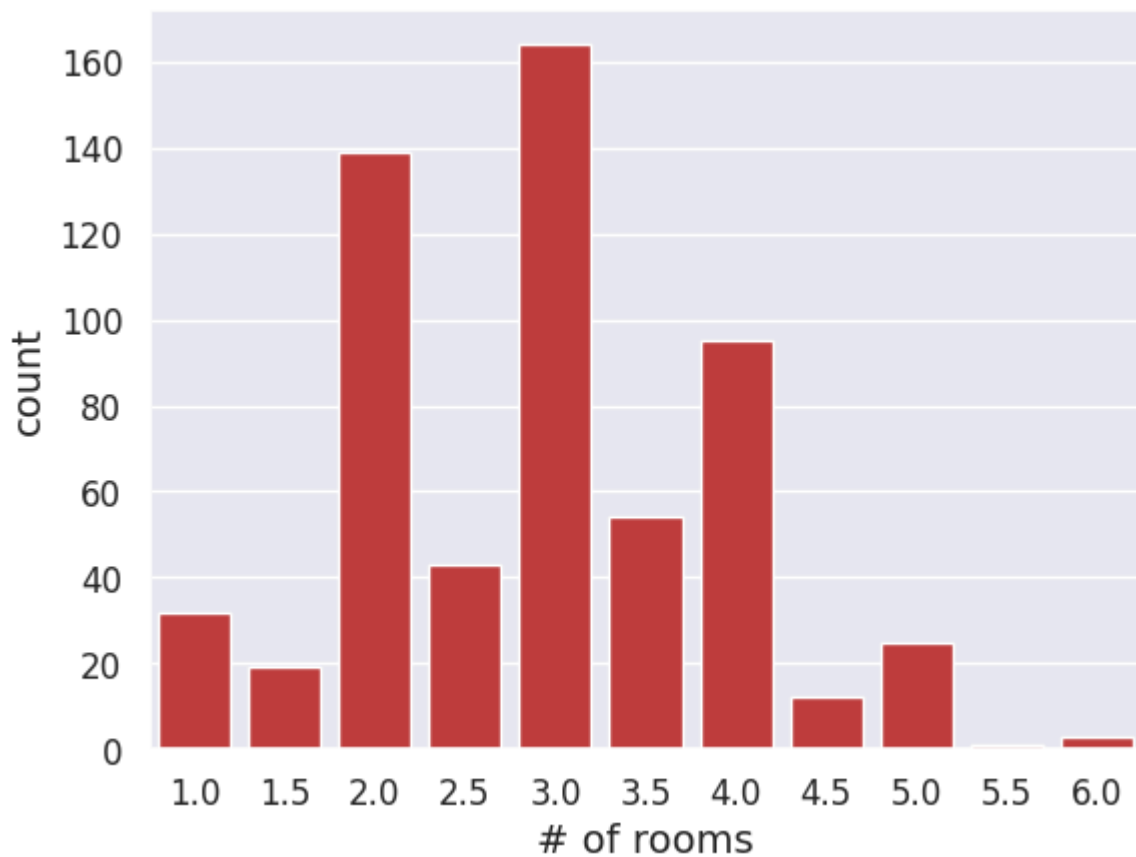
### What is the distribution of the number of rooms?

Q: Use `sns.countplot` to compare the counts of listings with different numbers of rooms. Plot all bars in the same [color](#) of your choice.

Start coding or [generate](#) with AI.

### ✓ Solution

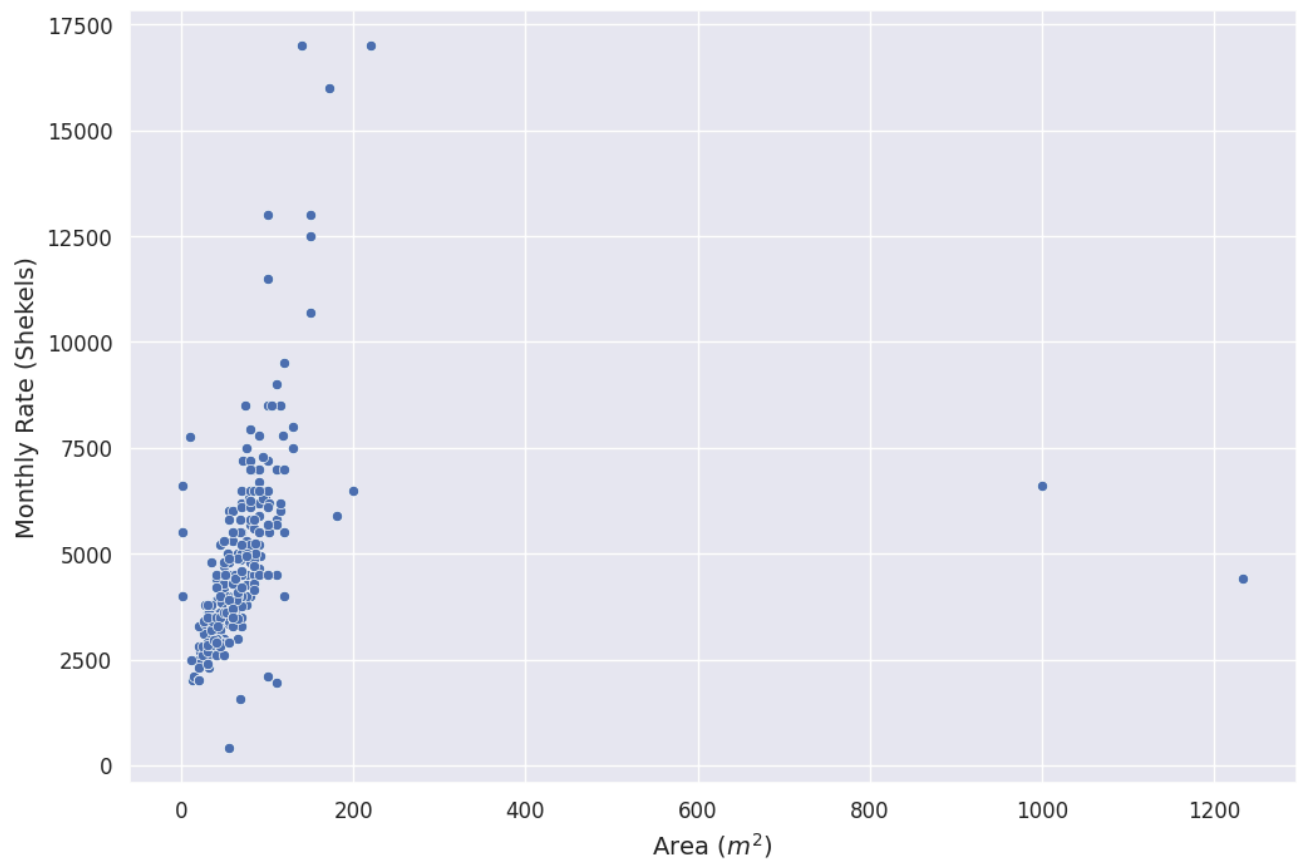
```
# @title Solution
if clean_df is None:
    print("Can't run until 'clean_df' is created!")
else:
    sns.countplot(x='rooms', data=clean_df, color='tab:red')
    plt.xlabel("# of rooms");
```



The distribution peaks at three rooms and we also see that "half rooms" are less common.

### Can we see an association between apartment area and price?

```
if clean_df is None:
    print("Can't run until 'clean_df' is created!")
else:
    plt.figure(figsize=(12,8))
    sns.scatterplot(x='area', y='monthlyRate', data=clean_df)
    plt.ylabel("Monthly Rate (Shekels)")
    plt.xlabel("Area ($m^2$));
```



We see clear outliers here! We know that area is measured in squared meters and it is unlikely that there are any apartments of  $\sim 1000m^2$ .

Let's look at those samples to see if we can understand what happened there:

```
if clean_df is None:
    print("Can't run until 'clean_df' is created!")
else:
    display(clean_df.sort_values('area', ascending=False).head(4))
```





	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	entry
185	3964340	תלפיות	4400.0	private	2.0	2.0	1234.0	10/08/2022
543	3956561	זכרון משה	6600.0	private	3.5	3.0	1000.0	01/07/2022



And inspect the description of one of those listings:

```
if clean_df is None:
    print("Can't run until 'clean_df' is created!")
else:
    display(clean_df.at[543, 'description'])
```



דירה מהממת בלב ירושלים. צמודה לרכבת הקלה- תחנת הדוידקה. 3 חדרים ענקיים ולכל חדר מרפסת גדולה. חלל כניס'ה עם פינת ישיבה. מתאימה מאוד ל- 3 שותפים

Clearly not a 1000 m<sup>2</sup> apartment...

Q: Save a new dataframe named `clean_df_area_filtered` with all listings with area smaller than 800 m<sup>2</sup>. Again, add the removed outliers to the `outliers_df` dataframe.

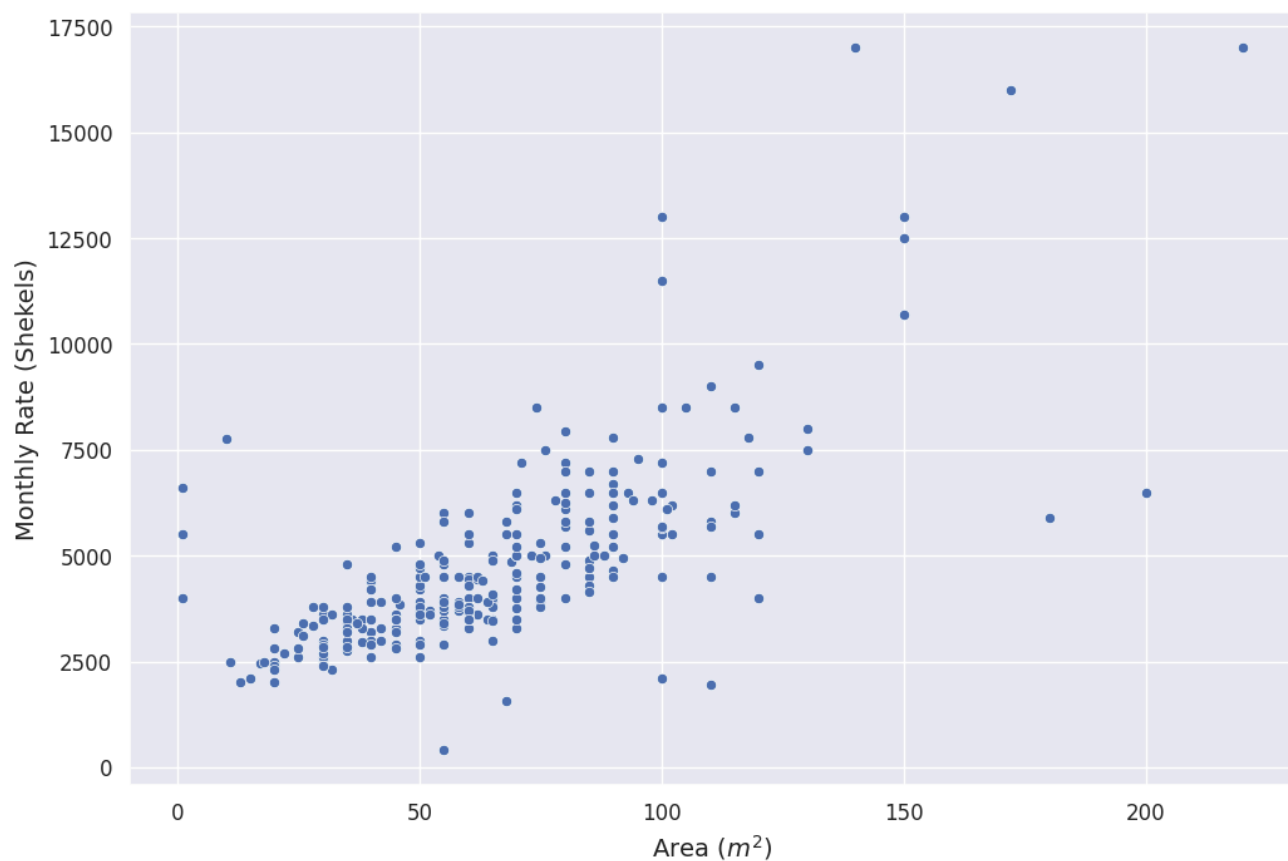
Plot again the scatter of area vs. monthly rate after removing the outliers.

Start coding or [generate](#) with AI.

## ✓ Solution

```
# @title Solution
if clean_df is None:
    print("Can't run until 'clean_df' is created!")
elif outlier_df is None:
    print("Can't run until 'outlier_df' is created!")
else:
    # save outliers
    outliers = clean_df[clean_df['area'] >= 800].reset_index(drop=True)
    outliers['reason'] = "'area' >= 800"
    outlier_df = pd.concat([outlier_df, outliers], axis=0, ignore_index=True).drop_duplicates()

    # remove the outliers from the dataset
    clean_df_area_filtered = clean_df[clean_df['area'] < 800].reset_index(drop=True)
    plt.figure(figsize=(12,8))
    sns.scatterplot(x='area', y='monthlyRate', data=clean_df_area_filtered)
    plt.xlabel("Area ($m^2$)")
    plt.ylabel("Monthly Rate (Shekels)");
```



Again, we see some strange behavior of apartments with almost zero area but with a high monthly rate. Let's check them out:

We start with all apartments with an area between 0 to 25  $m^2$ :

```
# Show all apartments with area between 0 and 25
clean_df_area_filtered[clean_df_area_filtered['area'].between(0,25)]
```



	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	entry
0	3994505	קריית יובל	2000.0	private	1.0	2.0	13.0	10/08/2022
1	3981298	רחביה	2450.0	private	1.0	1.0	17.0	10/08/2022
3	3993997	בית וגן	2100.0	private	1.0	0.0	15.0	10/08/2022
5	3993552	הר נוף	2000.0	private	1.0	0.0	20.0	10/08/2022
6	3972039	גבעת שאול	2700.0	private	1.0	0.0	22.0	10/08/2022
7	3988096	המושבה הגרמנית	2500.0	private	1.0	0.0	18.0	10/08/2022
8	3992809	נחלאות	3200.0	private	1.0	2.0	25.0	10/08/2022
10	3983516	הגבעה הצרפתית	2000.0	private	1.0	2.0	20.0	10/08/2022



Some make sense and others do not. Let's focus on the expensive ones (between 5,000 and 10,000 shekels):

```
# Show all apartments with area between 0 and 25 that also have a price between 5000 and
if clean_df_area_filtered is None:
    print("Can't run until 'clean_df_area_filtered' is created!")
else:
    display(clean_df_area_filtered[clean_df_area_filtered['area'].between(0,25) & clean_df_
```



	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	entry	
	197	3984483	ארנונה	6600.0	private	4.0	2.0	1.0	01/09/2022

Those are clearly wrong too... Besides that the relationship between the area and the price seems linear. Let's remove these outliers too:

```
#remove the outliers
if clean_df_area_filtered is None:
    print("Can't run until 'clean_df_area_filtered' is created!")
elif outlier_df is None:
    print("Can't run until 'outlier_df' is created!")
else:
    non_outliers = clean_df_area_filtered['area'] > 10 # get non outliers series of true/false

    # save outliers
    outliers = clean_df_area_filtered[~non_outliers].reset_index(drop=True) # get the outliers
    outliers['reason'] = "'area' <= 10"
    outlier_df = pd.concat([outlier_df, outliers], axis=0, ignore_index=True).drop_duplicates()

    # remove them
    clean_df_area_filtered = clean_df_area_filtered[non_outliers].reset_index(drop=True)
```

### ✓ Can we see a different pattern for top floor apartments?

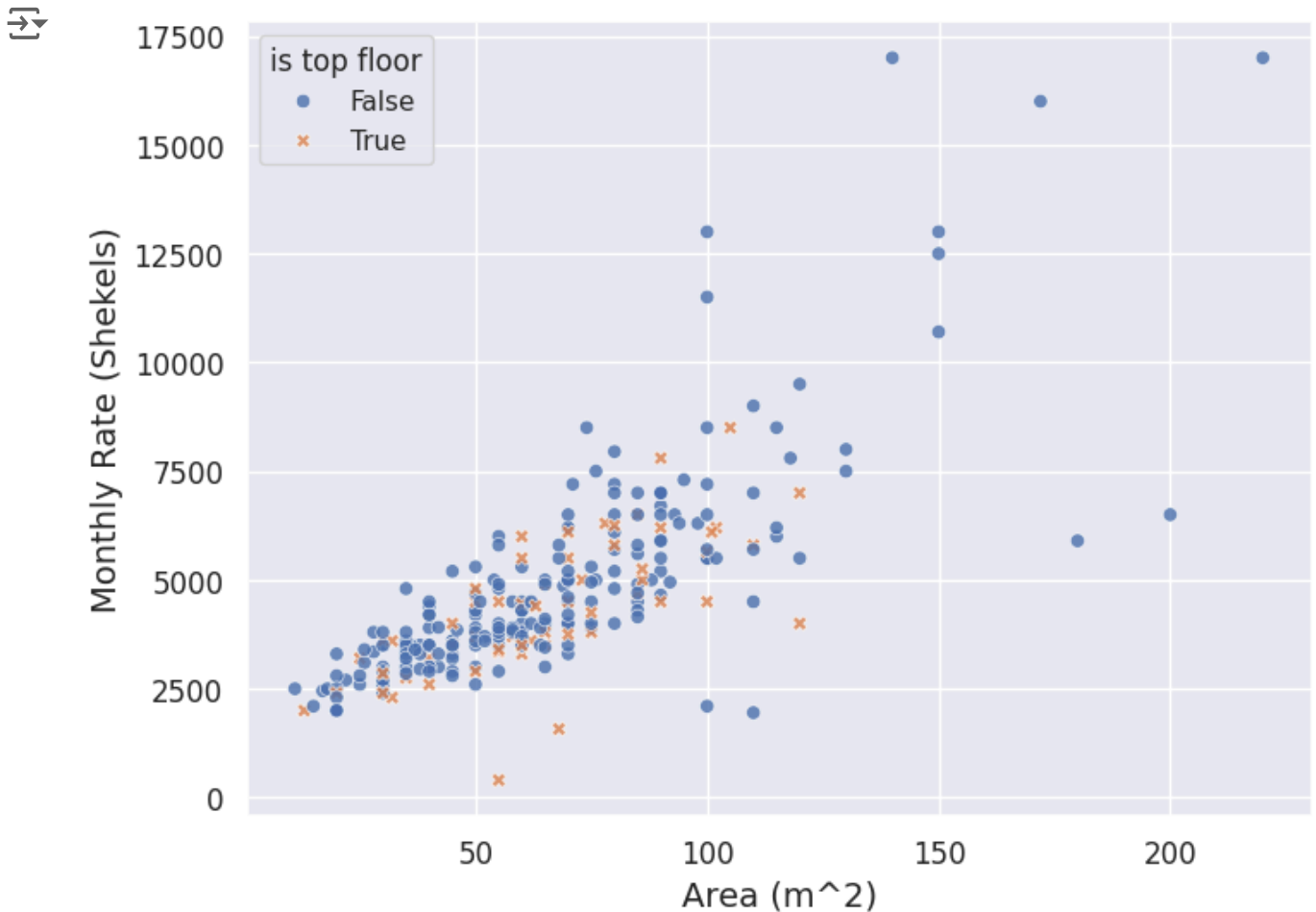
Q: Plot again a scatter of area vs. monthly rate. This time distinguish (by color / marker style or both) between apartments that are in the top floor and the rest of the apartments. (To do that you should create a new column in `clean_df_area_filtered` called `is_top_floor` and set it to 1 if the apartment is in the top floor and 0 otherwise.)

Start coding or [generate](#) with AI.

### ✓ Solution

```
# @title Solution
```

```
if clean_df_area_filtered is None:
    print("Can't run until 'clean_df_area_filtered' is created!")
else:
    clean_df_area_filtered['is top floor'] = clean_df_area_filtered['floor'] == clean_df_ar
    plt.figure(figsize=(8,6))
    sns.scatterplot(x='area', y='monthlyRate', data=clean_df_area_filtered, alpha=0.8, hue=
    plt.xlabel("Area (m^2)")
    plt.ylabel("Monthly Rate (Shekels)");
```



We can take a deeper look on the apartments with the very high monthly rate (to see if those are outliers or not):

```
if clean_df_area_filtered is None:
    print("Can't run until 'clean_df_area_filtered' is created!")
else:
    display(clean_df_area_filtered[clean_df_area_filtered['monthlyRate'] > 11000])
```



	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	entry
198	3956418	רחביה	13000.0	agent	4.0	1.0	100.0	NaN
236	3985051	טלביה	17000.0	private	4.0	4.0	140.0	10/08/2022

We can see some representation of the more expensive neighborhoods of Jerusalem here..  
More on the neighborhoods later on!

### Is there also a relation between the number of rooms and the listing price?

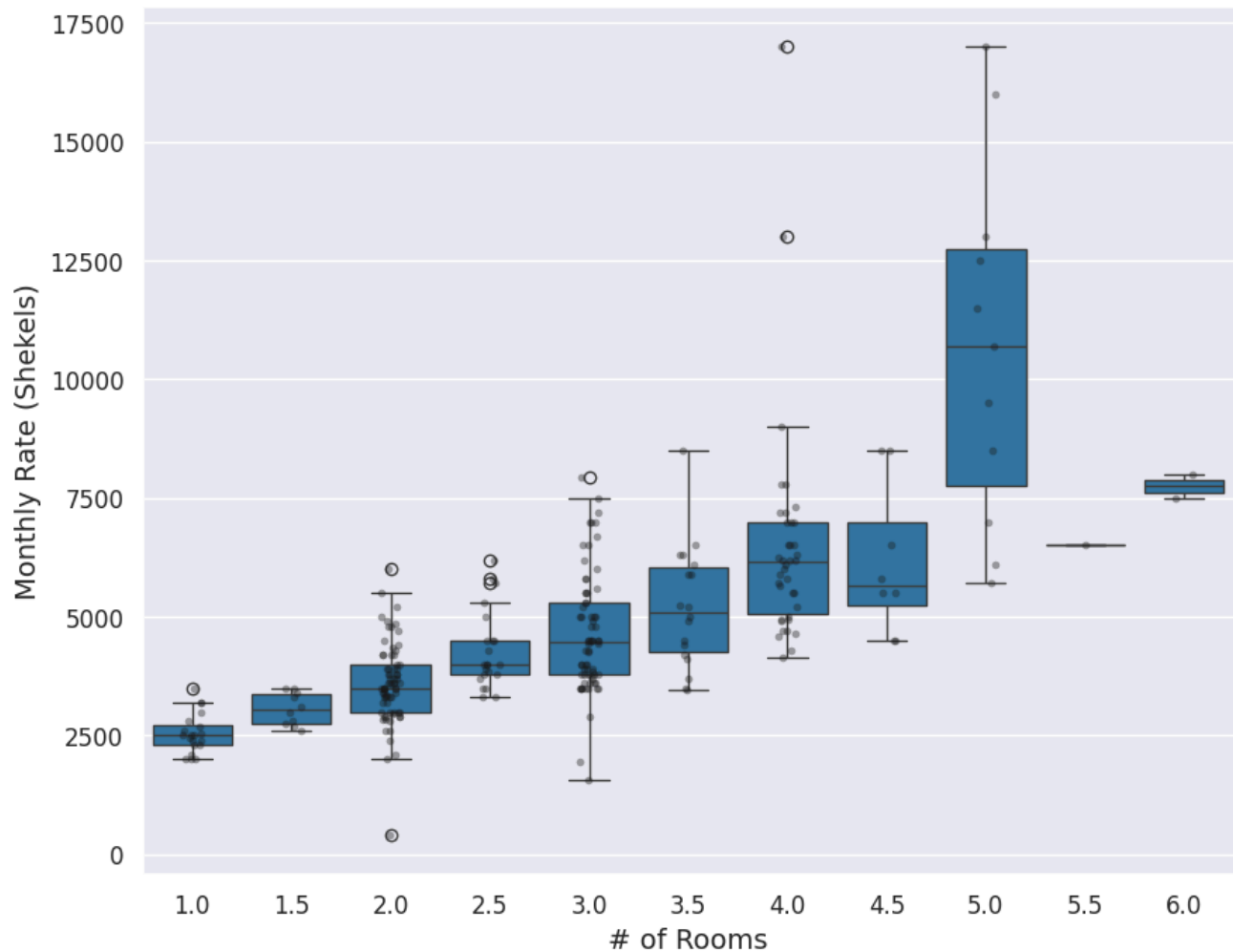
Q: Create a visualization that compares the distribution of prices for different number of rooms. Your visualization should provide information about central tendency (mean/median/mode) and some information about the distribution of individual values around it (standard deviation/interquartile range) for each number of rooms. Also, show the real prices of the listings per number of rooms.

### ✓ Solution

```
# @title Solution
if clean_df_area_filtered is None:
    print("Can't run until 'clean_df_area_filtered' is created!")
else:
    plt.figure(figsize=(10,8))
    sns.boxplot(x='rooms', y='monthlyRate', data=clean_df_area_filtered, color='tab:blue')
    sns.stripplot(x='rooms', y='monthlyRate', alpha=0.4 ,size=4,color='k',data=clean_df_are
    plt.xlabel("# of Rooms")
    plt.ylabel("Monthly Rate (Shekels)");

# Or:
# plt.figure(figsize=(10,8))
# sns.barplot(x='rooms', y='monthlyRate', data=clean_df_area_filtered, color='tab:blue'
# # Can also use mean but median is more informative in this case as prices are skewed.
# sns.stripplot(x='rooms', y='monthlyRate', alpha=0.4 ,color='k',data=clean_df_area_fil
# plt.xlabel("# of Rooms")
# plt.ylabel("Monthly Rate (Shekels)");

#Violin plot completely fails for very small subsets:
# plt.figure(figsize=(10,8))
# sns.violinplot(x='rooms', y='monthlyRate', data=clean_df_area_filtered, color='tab:bl
# plt.xlabel("# of Rooms")
# plt.ylabel("Monthly Rate (Shekels)");
```



Now that we finished pre-processing the data, we can see the state of our outliers VS the data that remains:

```
if outlier_df is None:
    print("Can't run until 'outlier_df' is created!")
else:
    # describe the outlier data
    display(outlier_df.groupby('reason').describe())
    print(f"Proportion removed: {100*len(outlier_df) / (len(outlier_df)+len(clean_df_area_f
```





reason	monthlyRate								rooms	
	count	mean	std	min	25%	50%	75%	max	count	n
'area' <= 10	5.0	5870.0	1399.821417	4000.0	5500.0	5500.0	6600.0	7750.0	5.0	5
'area' >= 800	2.0	5500.0	1555.634919	4400.0	4950.0	5500.0	6050.0	6600.0	2.0	2
monthlyRate <= 0	25.0	0.0	0.000000	0.0	0.0	0.0	0.0	0.0	25.0	25

3 rows × 40 columns

## ✓ Submission Exercises

### ✓ Part 1: Diving deeper into rental prices

We will create a copy of the dataset and work on that. We want to make sure that we do not modify the original dataset.

#### ✓ Part 1 - Create a DataFrame

```
# @title Part 1 - Create a DataFrame
part1_df = rent_df_backup_for_exercise.copy()
```

Let's go back to the distribution of monthly rental prices in the dataset. Are there interesting trends in the distribution that we missed in the visualizations before?

**Use only `part1_df` for the coding questions in this part**

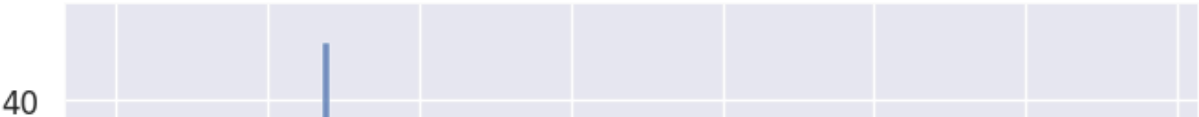
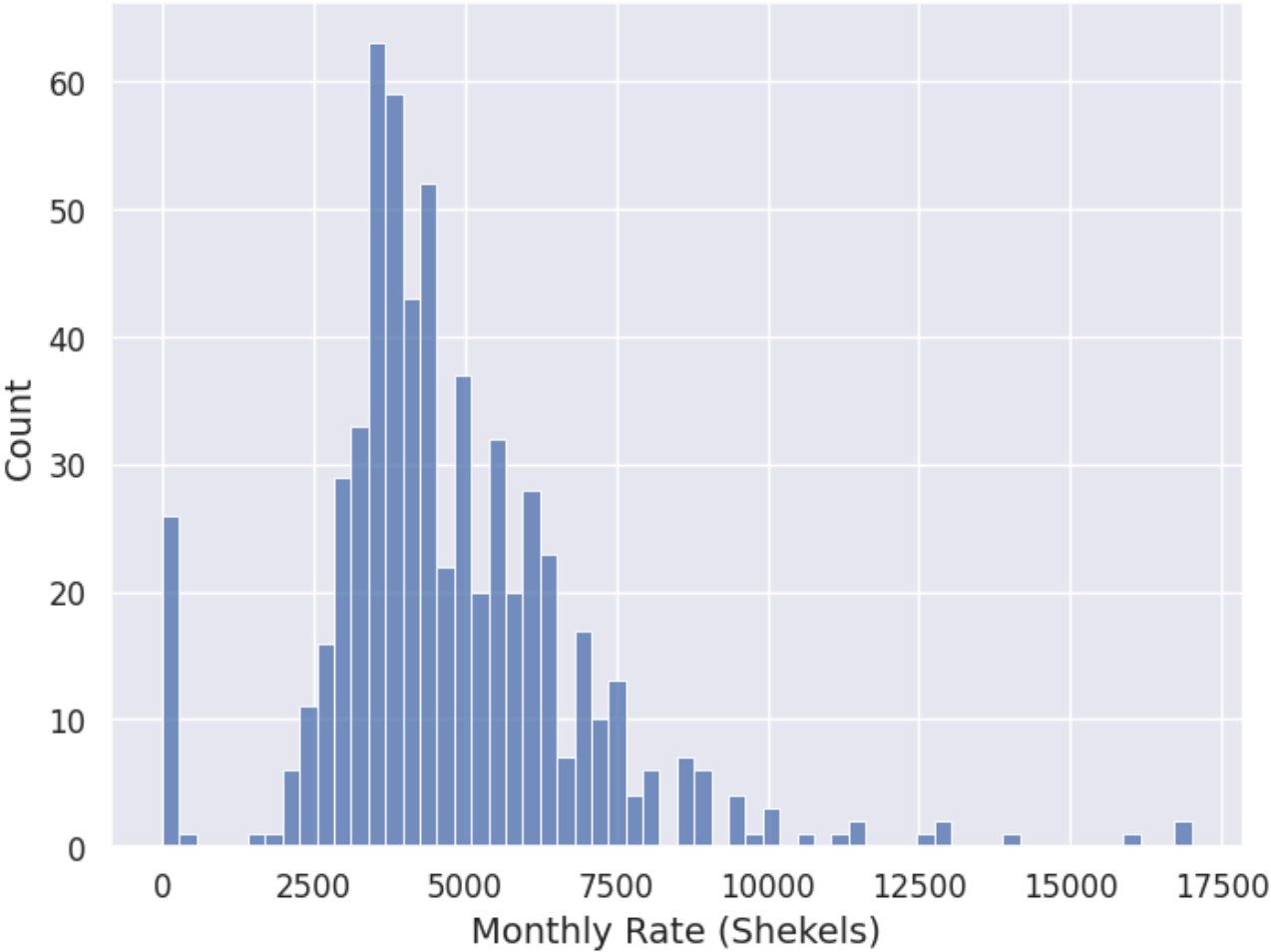
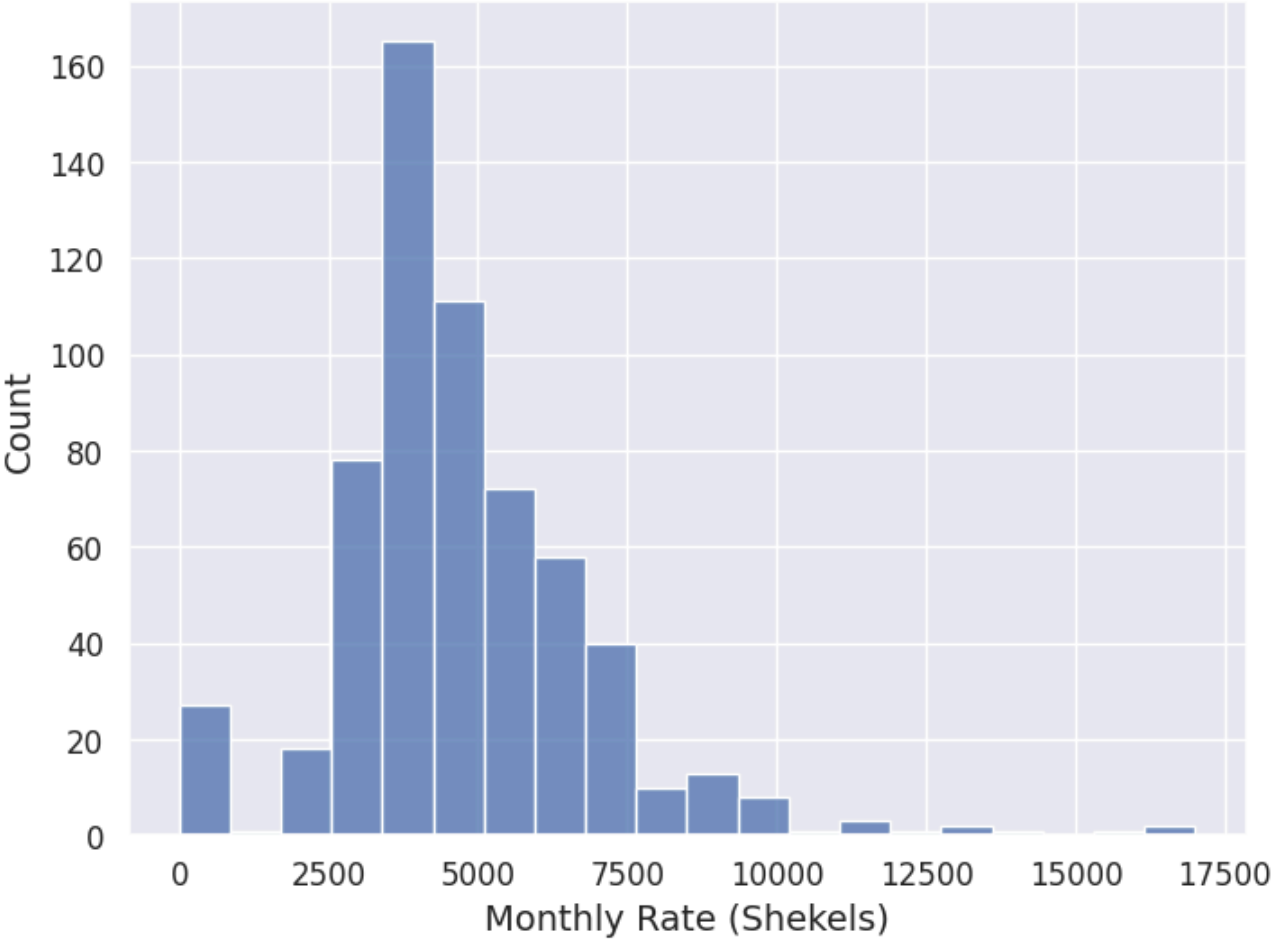
#### ✓ Question 1

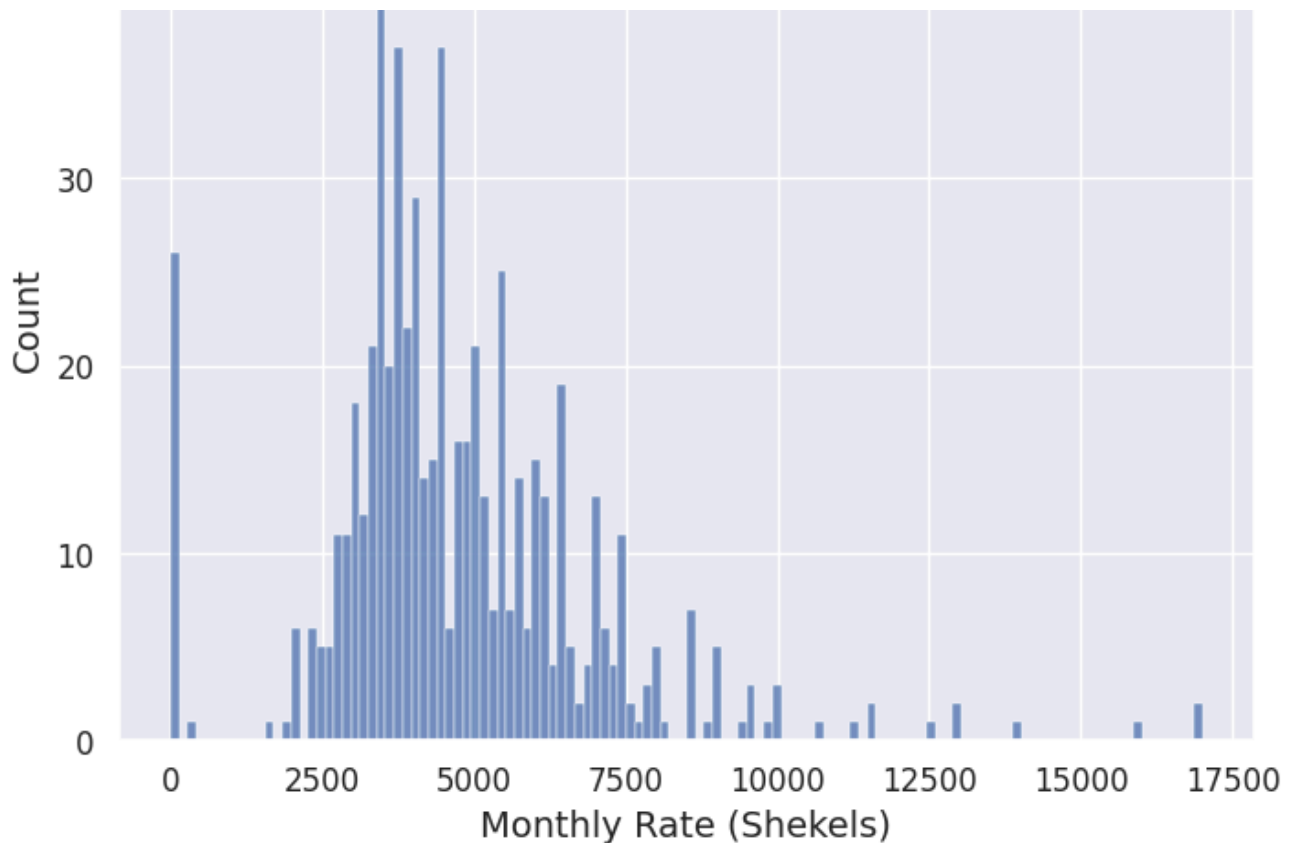
Plot 3 different histograms of the monthly prices with 20, 60 and 120 bins respectively, each in a different axis/figure.

```
# Part 1 - Question 1
plt.figure(figsize=(8,6))
sns.histplot(part1_df['monthlyRate'], bins=20)
plt.xlabel("Monthly Rate (Shekels)");

plt.figure(figsize=(8,6))
sns.histplot(part1_df['monthlyRate'], bins=60)
plt.xlabel("Monthly Rate (Shekels)");

plt.figure(figsize=(8,6))
sns.histplot(part1_df['monthlyRate'], bins=120)
plt.xlabel("Monthly Rate (Shekels)");
```





## ✓ Question 2

For 60 and 120 bins, you can see a repeating pattern of "peaks" and "vallies" in the distribution (mostly in the range between 500 and 7000). Is this pattern due to people rounding the rental prices? Please create a visualization that answers this question. Describe in words how the graph shows what the answer is (Hint: you can use the '%' operator to compute the remainder of dividing values in a pandas Series by a scalar number).

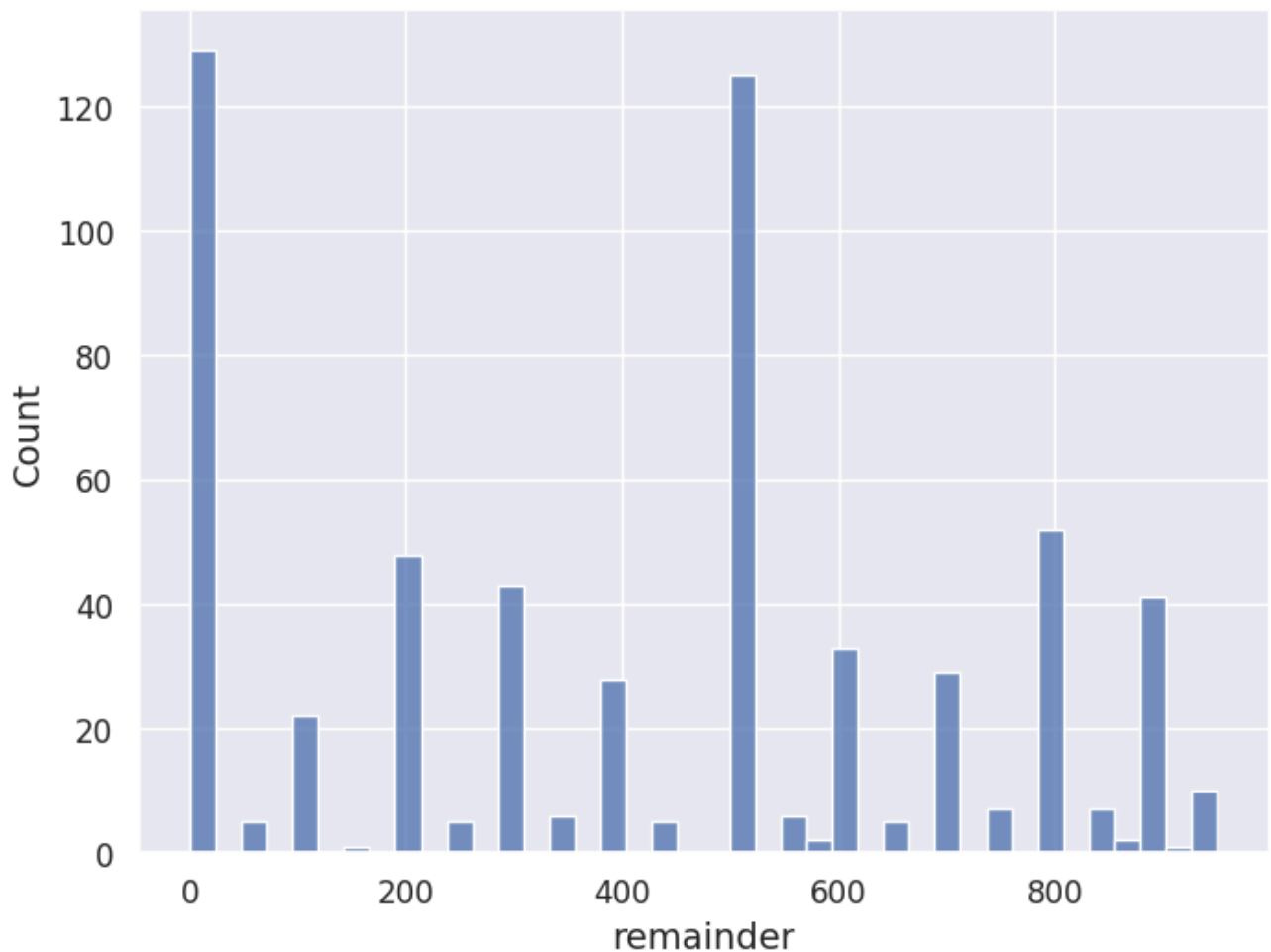
- **extra hint:** please open this cell only after discussing with the course staff the best solution you could come up with

[Show code](#)

```
# Part 1 - Question 2
# Your code goes here:
part1_df['remainder'] = part1_df['monthlyRate'] % 1000
np.random.seed(2)
part1_df.sample(5)

plt.figure(figsize=(8,6))
sns.histplot(part1_df['remainder'], bins=40)
```

↗ <Axes: xlabel='remainder', ylabel='Count'>



### Part 1 Question 2 - textual Answer:

Write your answer here: ארצה להציג אולי היסטוגרמה של האם המחירים מתחלקים באלף ולראות האם ההיסטוגרמה של המתחלקים יותר גבוהה. ניתן לראות שככל שהשארית מחלוקה ב-1000 עגולה יותר - יש יותר דירות

הרבה דירות כאשר השארית היא 0 וגם 500

### ✓ Question 3

We expect to see a "drop" in prices frequency near the 5000 Shekels mark due to tax considerations (See [here](#) for an explanation). Create a histogram visualization of the data with the smallest possible bins such that every bin will include exactly one multiplication of 500 (Hint: read the `bins` parameter documentation and what types it accepts). Explain why does this choice of bin size ensures that we will not see rounding effects. Do you see a "drop" around 5000 Shekels? Are there other "drops"?

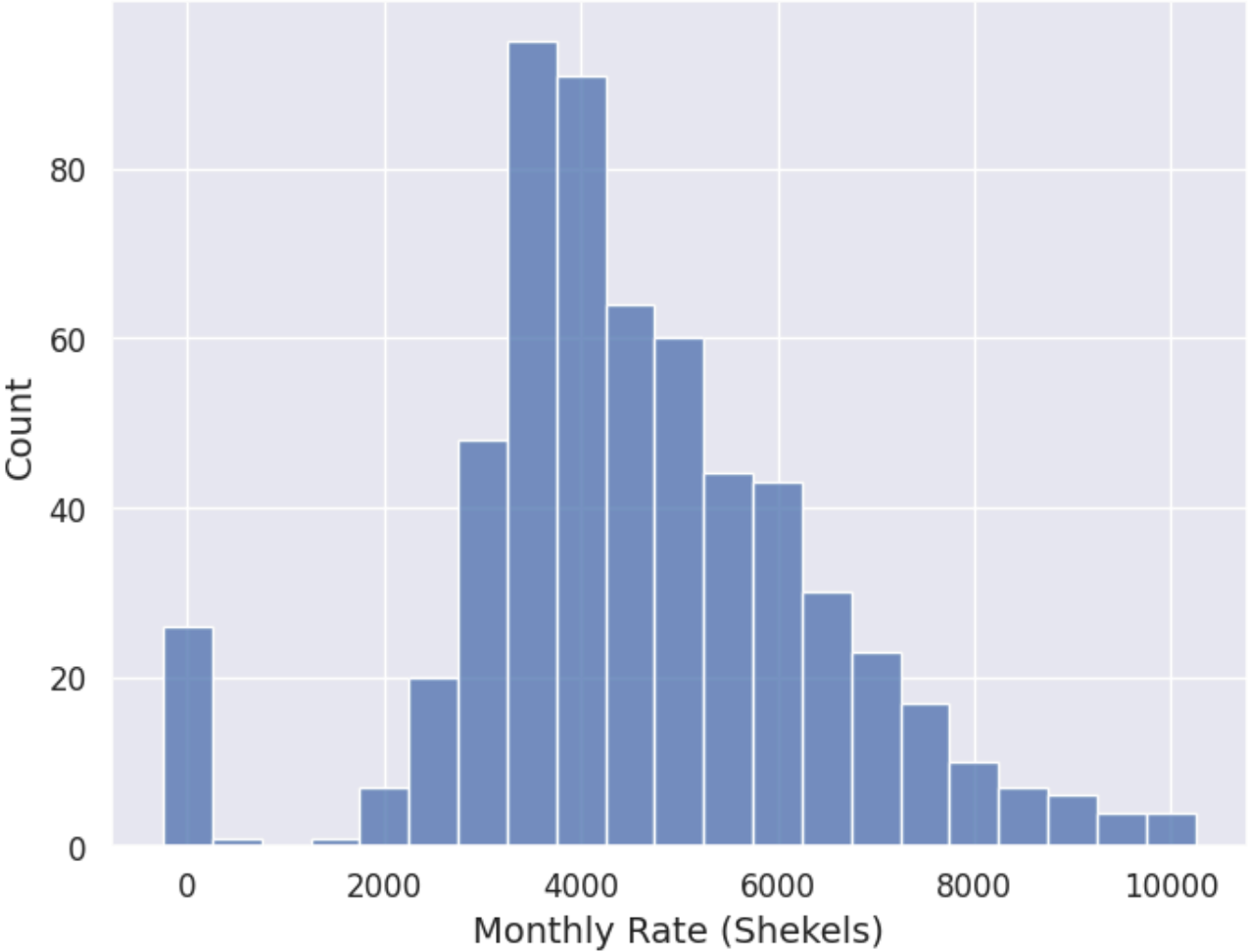
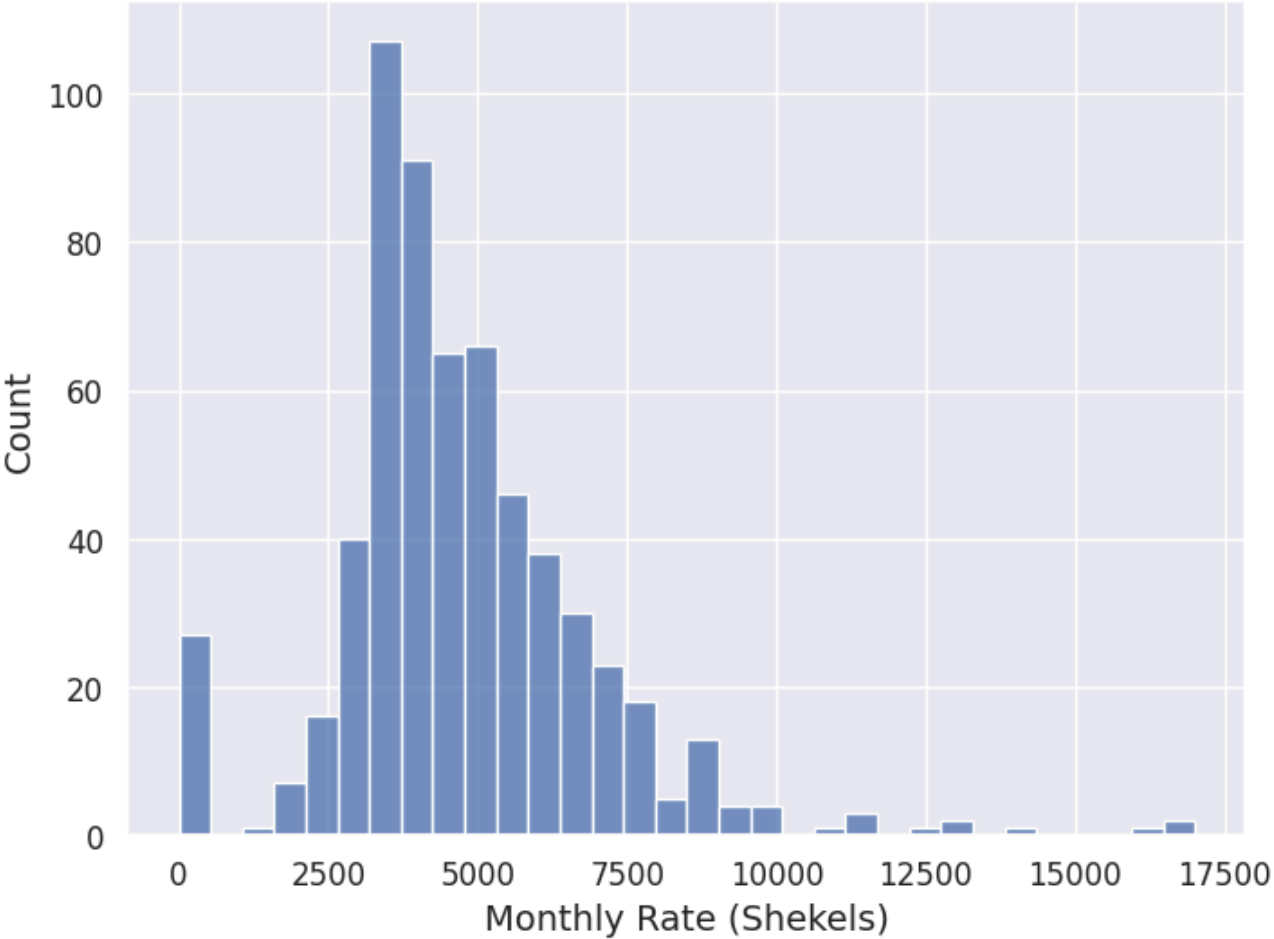
```
# Part 1 - Question 3
# Your code goes here:
```

```
plt.figure(figsize=(8,6))
sns.histplot(part1_df['monthlyRate'], bins='fd')
plt.xlabel("Monthly Rate (Shekels)");
```

```
bin_centers = np.arange(0, 10000 + 1, 500) # Adjust range according to your data
```

```
# Calculate bin edges from bin centers
bin_edges = np.append(bin_centers - 250, bin_centers[-1] + 250)
```

```
plt.figure(figsize=(8,6))
sns.histplot(part1_df['monthlyRate'], bins=bin_edges)
plt.xlabel("Monthly Rate (Shekels)");
```



---

## Part 1 Question 3 - textual Answer:

*Write your answer here:*

לרוב המשכירים מעגלים סביב 500 ולכן אם נייצג ככה את הנתונים, אנחנו מבטלים את העיגול כי אנחנו לכאורה מייצגים את הממוצעים סביב העיגול אנחנו רואים שההיסטוגרמה באמת הוחלקה ושיש ירידה לא מאוד גדולה סביב 500 גם אחרי 4000 יש ירידה משמעותית

---

### ✓ Part 2: Size or number of rooms?

### ✓ Part 2 - Create a DataFrame for Part 2

```
# @title Part 2 - Create a DataFrame for Part 2
```

```
# Create the dataframe and remove the outliers we found in the intro part:
part2_df = rent_df_backup_for_exercise.copy()
part2_df = part2_df[part2_df['monthlyRate'] > 0].reset_index(drop=True);
part2_df = part2_df[part2_df['area'] < 800].reset_index(drop=True)
part2_df = part2_df[part2_df['area'] > 10].reset_index(drop=True)
```

We saw that both the number of rooms and the area of an apartment are strongly associated with the monthly rate. We now want to check if those are just two perspectives of the same relation (how big is the apartment) or is there something more to it. We will use the cleaned dataframe for this exercise.

**Use only `part2_df` for the coding questions in this part**

### ✓ Question 1

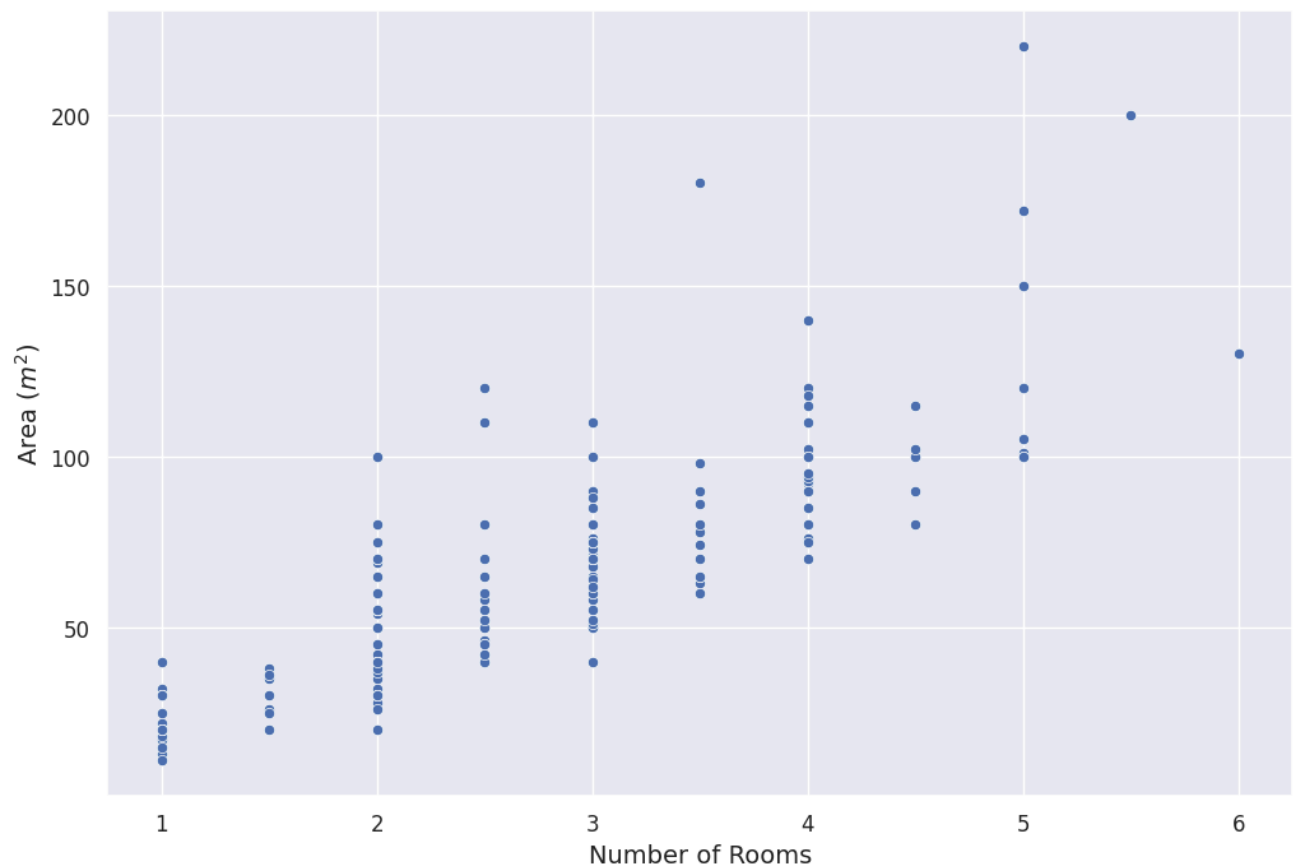
Generate a visualization to show that there is a strong association between the number of rooms and the area of the apartment. Explain your choice of plot type and your conclusion from the graph.



```
# Part 2 - Question 1  
# Your code goes here:
```

```
np.random.seed(2)  
part2_df.sample(5)
```

```
plt.figure(figsize=(12,8))  
sns.scatterplot(x='rooms', y='area', data=part2_df)  
plt.ylabel("Area ($m^2$)")  
plt.xlabel("Number of Rooms");
```



Part 2 Question 1 - textual Answer:

*Write your answer here:*

ניתן לראות קשר חיובי די משמעותי בין מספר החדרים לבין גודל הדירה.

בחרנו את הגרף הזה כי הוא מראה קשר בין שני משתנים

כמו שבדקנו בשיעור את הקשר בין גיל המשתתפים לתוצאה שלהם (בסבב הראשון)

## ✓ Question 2

Add a new column to the dataframe named "averageRoomSize" with the average room size in the given listing.

```
# Part 2 - Question 2
# Your code goes here:
```

```
part2_df['averageRoomSize'] = part2_df['area'] / part2_df['rooms']
```

```
np.random.seed(2)
part2_df.sample(5)
```



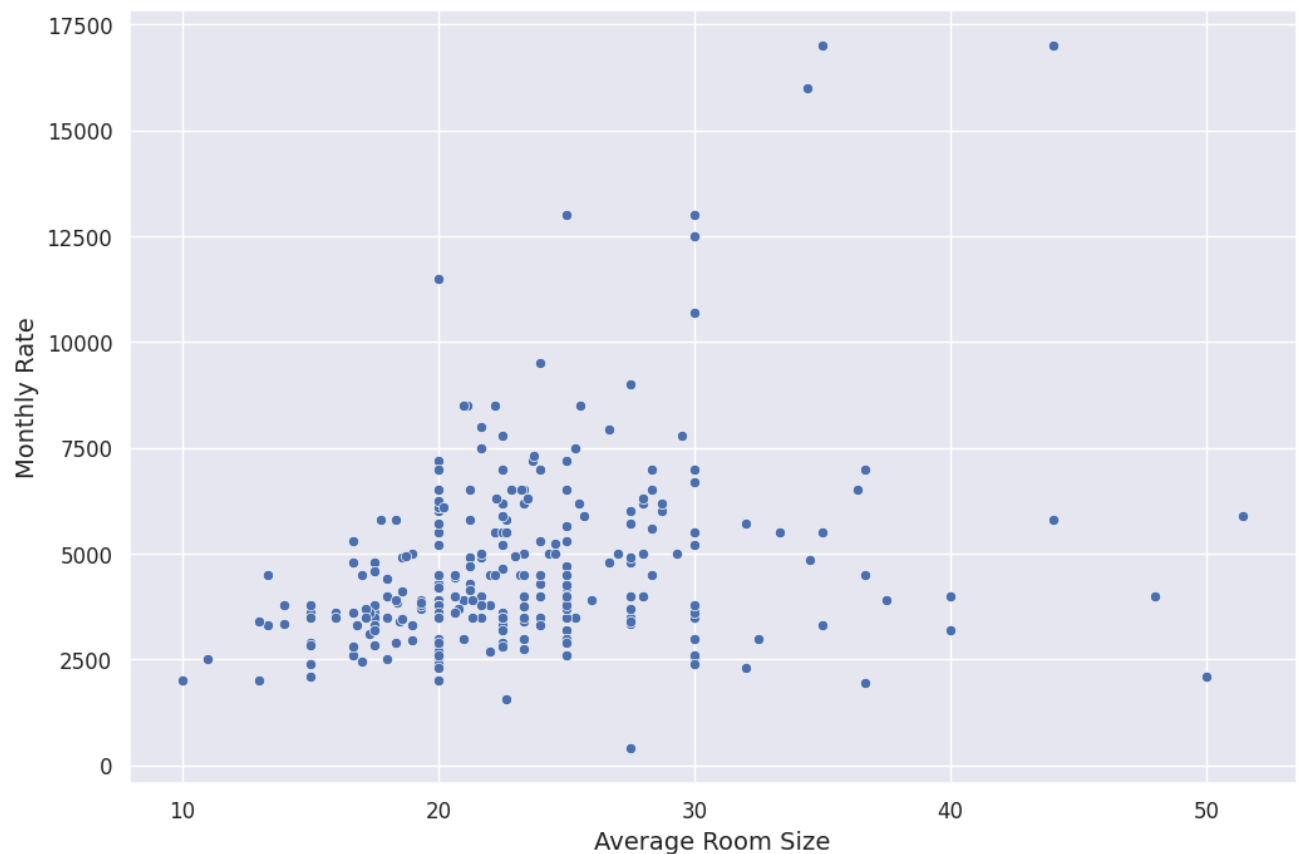
	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	entry	di
66	3981542	בקעה	4200.0	private	2.0	1.0	40.0	NaN	0
7	3988096	המושבה הגרמנית	2500.0	private	1.0	0.0	18.0	10/08/2022	ת

## ✓ Question 3

Create a plot of the relation between the average room size and the monthly rate.

```
# Part 2 - Question 3
# Your code goes here:
```

```
plt.figure(figsize=(12,8))
sns.scatterplot(x='averageRoomSize', y='monthlyRate', data=part2_df)
plt.ylabel("Monthly Rate")
plt.xlabel("Average Room Size");
```



#### ✓ Question 4 - **bonus**

We can see that the variance of the monthly rate increases with the average room size.

Suggest what might be the reason for the increase in the variance and create a visualization to support or refute your suggestion.

```
# Part 2 - Question 4
# Your code goes here:
#
#
```

## Part 2 Question 4 - textual Answer:

*Write your answer here:*

כנראה שכאשר גודל החדר הממוצע גדול יותר - או שהדירה ממש גדולה וקצת חדרים או הפוך

וכאשר גודל החדר הממוצע קטן כנראה שיש רק קצת חדרים וגודל קטן

אפשר לבדוק את זה עם היסטוגרמות - אבדוק את זה אחר כך

---

## ✓ Part 3: Neighborhoods

### ✓ Part 3 - Function Definitions and DataFrame Creation

```
# @title Part 3 - Function Definitions and DataFrame Creation
```

```
def reverse_string(a):
    return a[::-1]
```

```
socialrank_df = load_df(SOCIORANK_ID)
```

```
neighborhood_ranks = {k: v for k,v in zip(socialrank_df['neighborhood'], socialrank_df['s
```

```
def get_neighborhood_rank(neighborhood):
    if neighborhood in neighborhood_ranks:
        return neighborhood_ranks[neighborhood]
    else:
        return None
```

```
# Create the dataframe and remove the outliers we found in the intro part:
```

```
part3_df = rent_df_backup_for_exercise.copy()
part3_df = part3_df[part3_df['monthlyRate'] > 0].reset_index(drop=True);
part3_df = part3_df[part3_df['area'] < 800].reset_index(drop=True)
part3_df = part3_df[part3_df['area'] > 10].reset_index(drop=True)
part3_df["neighborhood_flipped"] = part3_df["neighborhood"].apply(reverse_string) # makin
```

We now want to focus on the differences between different neighborhoods in Jerusalem.

**Use only part3\_df for the coding questions in this part**

\*Use the "neighborhood\_flipped" column for visualizations as seaborn will flip the order of letters in hebrew.

## ✓ Question 1

Print the number of unique neighborhoods that appear in the dataset.

```
# Part 3 - Question 1
# Your code goes here:
```

```
unique_strings = set(part3_df["neighborhood"])

# Get the number of unique strings
num_unique_strings = len(unique_strings)

# Print the number of unique strings
print(f"Number of unique strings in the dataset: {num_unique_strings}")
```

➞ Number of unique strings in the dataset: 46

## ✓ Question 2

Visualize the number of listings per neighborhood in a way that will allow you to easily identify those with the highest count.

```
# Part 3 - Question 2
# Your code goes here:
```

```
listings_per_neighborhood = part3_df['neighborhood_flipped'].value_counts()
plt.figure(figsize=(10, 6))
sns.barplot(x=listings_per_neighborhood.index, y=listings_per_neighborhood.values, palette='magma')
plt.xlabel('Neighborhood')
plt.ylabel('Number of Listings')
plt.title('Number of Listings per Neighborhood')
plt.xticks(rotation=45)
plt.show()
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.

[illegible]

Print the number of neighborhoods with less than 5 listings and the fraction of their total number of listings out of the total number of listings. Also print the fraction of listings from the 8 most frequent neighborhoods out of the total number of listings.

```
# Part 3 - Question 3
# Your code goes here:
neighborhoods_less_than_5 = listings_per_neighborhood[listings_per_neighborhood < 5]
num_neighborhoods_less_than_5 = len(neighborhoods_less_than_5)
print(num_neighborhoods_less_than_5)
fraction_less_than_5 = neighborhoods_less_than_5.sum() / listings_per_neighborhood.sum()

# Top 8 most frequent neighborhoods
top_8_neighborhoods = listings_per_neighborhood.head(8)
fraction_top_8 = top_8_neighborhoods.sum() / listings_per_neighborhood.sum()

print(f"Number of neighborhoods with less than 5 listings: {num_neighborhoods_less_than_5}")
print(f"Fraction of total listings from neighborhoods with less than 5 listings: {fraction_less_than_5}")
print(f"Fraction of total listings from the 8 most frequent neighborhoods: {fraction_top_8}")
```



28

```
Number of neighborhoods with less than 5 listings: 28
Fraction of total listings from neighborhoods with less than 5 listings: 0.23
Fraction of total listings from the 8 most frequent neighborhoods: 0.52
```

Those types of distributions where there are many categories that appear only a few times but together take a large portion of the distribution are called heavy-tailed (or long-tailed) distributions. This is a real issue in many data science applications, since even if we have a large dataset there are still some sub-populations or sub-categories that are not well represented.

## ✓ Question 4

Create a new filtered dataframe with listings from only the 8 most frequent neighborhoods.

```
# Part 3 - Question 4
# Your code goes here:
# Part 3 - Question 3
# Your code goes here:

top_8_neighborhoods = listings_per_neighborhood.head(8).index

filtered_df = part3_df[part3_df['neighborhood_flipped'].isin(top_8_neighborhoods)]

# Display the filtered DataFrame
print(filtered_df)
```



	propertyID	neighborhood	monthlyRate	mefarsem	rooms	floor	area	\
0	3994505	2000.0	קריית יובל	private	1.0	2.0	13.0	
3	3993997	2100.0	בית וגן	private	1.0	0.0	15.0	
4	3994399	2300.0	פסגת זאב	private	1.0	1.0	32.0	
11	3986231	2600.0	קריית יובל	private	1.0	1.0	30.0	
14	3992479	2400.0	קריית יובל	private	1.0	1.0	20.0	
..	...	...	...	...	...	...	...	...

262	3988577	6500.0	פסגת זאב	private	5.5	1.0	200.0
264	3993965	7000.0	בית וגן	private	5.0	3.0	120.0
266	3974914	17000.0	תלפיות	private	5.0	3.0	220.0
270	3981999	13000.0	תלפיות	private	5.0	4.0	150.0
271	3994215	2900.0	קריית יובל	private	2.0	0.0	40.0

	entry		description	numFloors	\
0	10/08/2022	2.0	יחידת דיור להשכרה ברחוב הראשי של קריית יובל, ה...		
3	10/08/2022	3.0	דירת חדר, כ-15 מ"ר, במיקום מרכזי אך שקט, משופצ...		
4	10/08/2022	1.0	בס"ד בפסגת זאב מזרח דירת חדר גדולה משופצת ויפ...		
11	10/08/2022	2.0	הדירה שטופת שמש, מגיעה מרוהטת- מיטה, ארון בגד...		
14	10/08/2022	1.0	דירת חדר חמודה עם גינה קטנה משותפת, מתאים ליחיד...		
..	...		...		...
262	10/08/2022	3.0	דירה בת 5.5 חדרים . בקומה התחתונה סלון , מטבח ...		
264	10/08/2022	3.0	דירת 5 חדרים ובעלת 5 מרפסות קטנות, עברה צביעה ...		
266	10/08/2022	4.0	דירת 5 חדרים ענקית ומהממת, בבנין בוטיק ויחודי ...		
270	10/08/2022	5.0	דירת 5 חדרים חדשה! בדירה יש מרפסת מרפסת שירו...		
271	10/08/2022		NaN	4.0	

	neighborhood_flipped
0	לבוי תיירק
3	גן תיב
4	באז תגספ
11	לבוי תיירק
14	לבוי תיירק
..	...
262	באז תגספ
264	גן תיב
266	תויפלת
270	תויפלת
271	לבוי תיירק

[143 rows x 11 columns]

## ✓ Question 5

Plot a graph to check whether there are different distributions of monthly rates in the eight neighborhoods. Explain your choice for the visualization and your conclusions. Note: Make sure that the neighborhoods are ordered in the plot based on their tendency for higher or lower monthly rates.

Hint: Which is a better descriptor of the central tendency of monthly rates when the distributions are skewed?



```
# Part 3 - Question 5
# Your code goes here:
```

```
# Ordering neighborhoods by median monthly rates
ordered_neighborhoods = filtered_df.groupby('neighborhood_flipped')['monthlyRate'].median

plt.figure(figsize=(12, 8))
sns.violinplot(x='neighborhood_flipped', y='monthlyRate', order=ordered_neighborhoods, da
plt.xlabel("neighborhood")
plt.ylabel("monthlyRate");

listings_per_neighborhood = part3_df['neighborhood_flipped'].value_counts()

# Identifying the 8 most frequent neighborhoods
top_8_neighborhoods = listings_per_neighborhood.head(8).index

# Filtering the DataFrame to include only listings from the 8 most frequent neighborhoods
filtered_df = part3_df[part3_df['neighborhood_flipped'].isin(top_8_neighborhoods)]

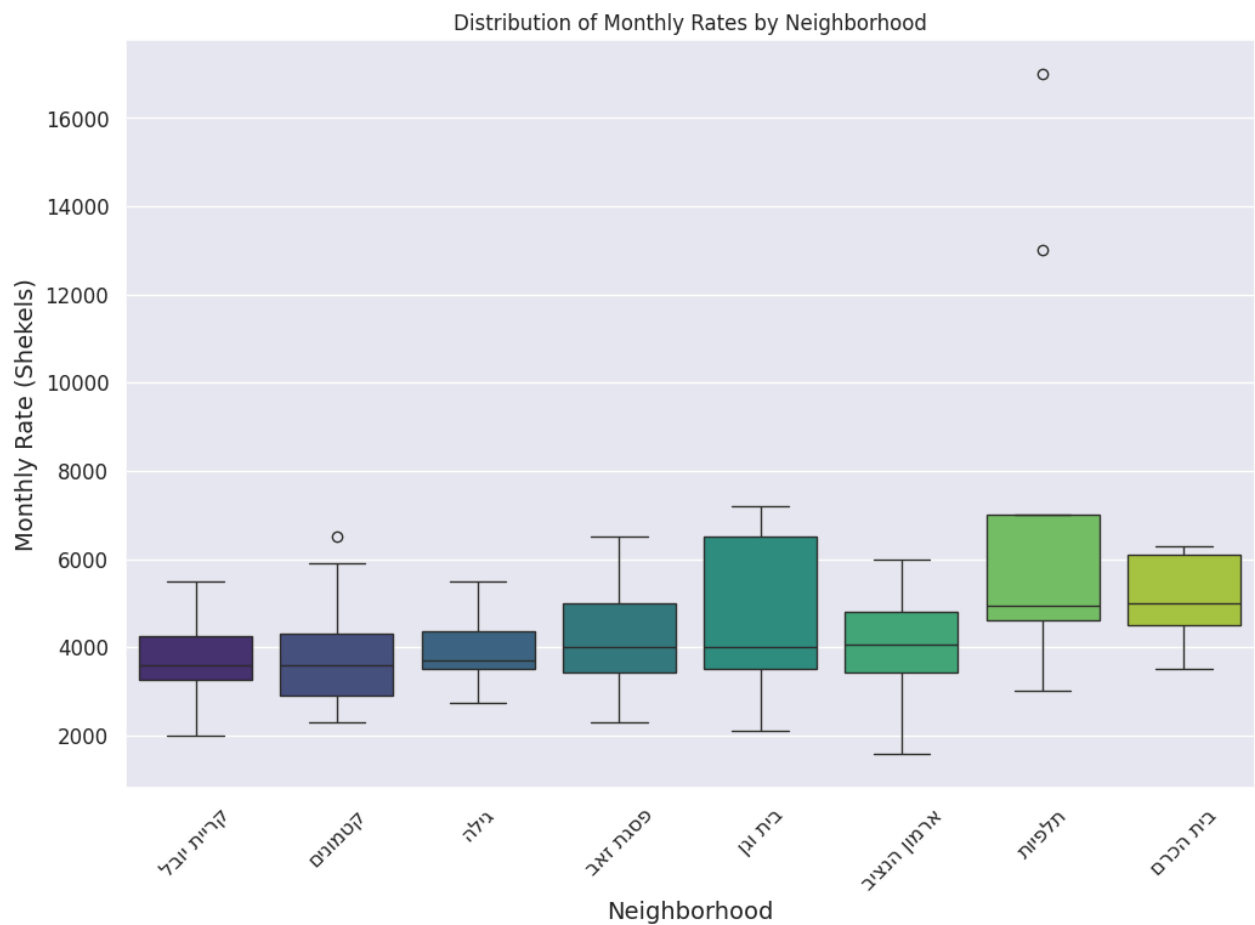
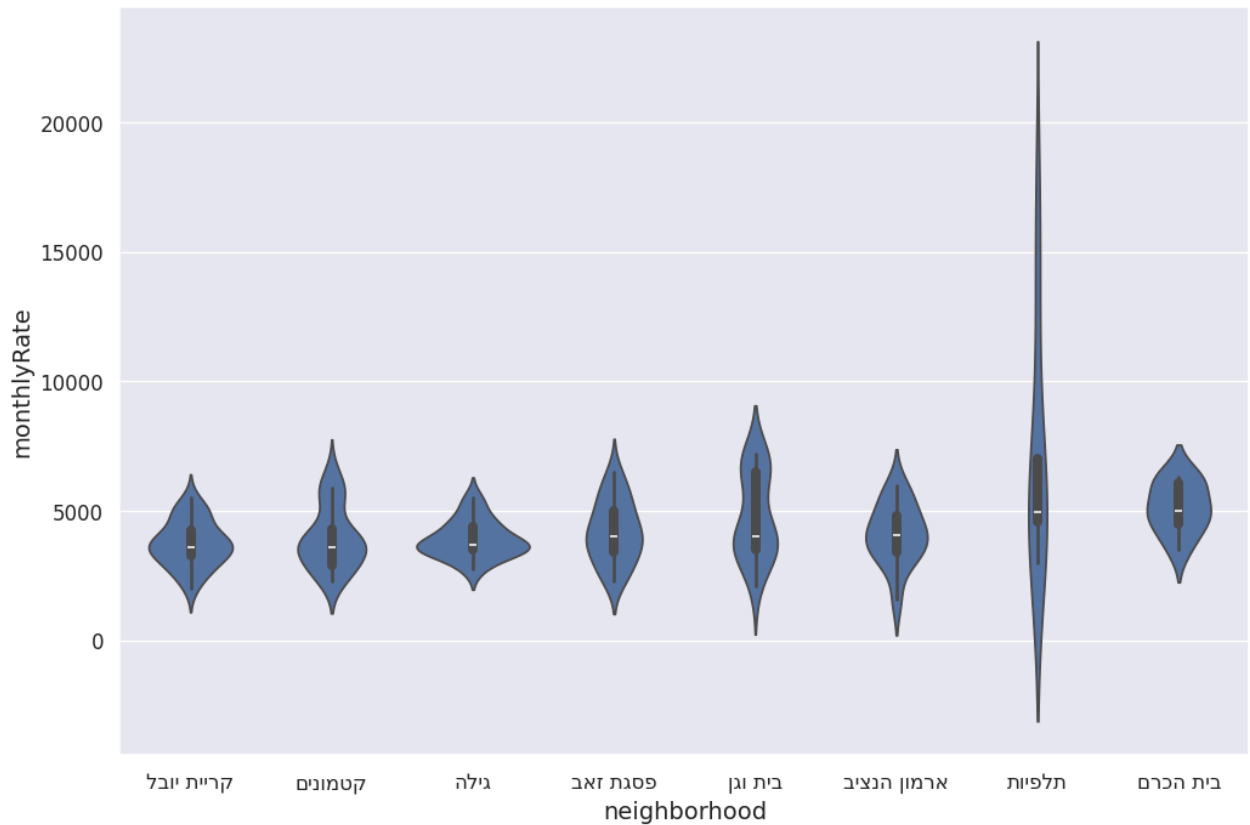
# Plotting the box plot
plt.figure(figsize=(12, 8))
sns.boxplot(data=filtered_df, x='neighborhood_flipped', y='monthlyRate', order=ordered_ne
plt.xlabel("Neighborhood")
plt.ylabel("Monthly Rate (Shekels)")
plt.title("Distribution of Monthly Rates by Neighborhood")
plt.xticks(rotation=45)
plt.show()
```



```
<ipython-input-121-3a8270c31932>:25: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.

```
sns.boxplot(data=filtered_df, x='neighborhood_flipped', y='monthlyRate', order=o
```



---

### Part 3 Question 5 - textual Answer:

*Write your answer here:*

בחרנו לראות את ההתפלגויות של כל שכונה באמצעות גרף כינור וראינו משהו מוזר בתלפיות, בדקנו גרף קופסאות וראינו ערכים קיצוניים שמסבירים את זה

---

### ✓ Question 6

Now that we compared the different distributions of monthly rates between neighborhoods, we can check whether we can explain some of the differences using our common-sense and the data we already have. For example, perhaps different neighborhoods have different distributions of apartment sizes?

Think of a new variable that will allow you to check the relationship between neighborhoods and prices fairly, factoring different apartment sizes out of the equation. Save this measure into the dataframe and create a new visualization to answer the question.

```
# Part 3 - Question 6
# Your code goes here:
filtered_df['price_per_mr'] = filtered_df['monthlyRate'] / filtered_df['area']

# Ordering neighborhoods by median monthly rates
ordered_neighborhoods = filtered_df.groupby('neighborhood_flipped')['price_per_mr'].median()

plt.figure(figsize=(12, 8))
sns.violinplot(x='neighborhood_flipped', y='price_per_mr', order=ordered_neighborhoods, d
plt.xlabel("neighborhood")
plt.ylabel("price_per_mr");

# Identifying the 8 most frequent neighborhoods
top_8_neighborhoods = listings_per_neighborhood.head(8).index

# Filtering the DataFrame to include only listings from the 8 most frequent neighborhoods
filtered_df = part3_df[part3_df['neighborhood_flipped'].isin(top_8_neighborhoods)]

listings_per_neighborhood = part3_df['neighborhood_flipped'].value_counts()

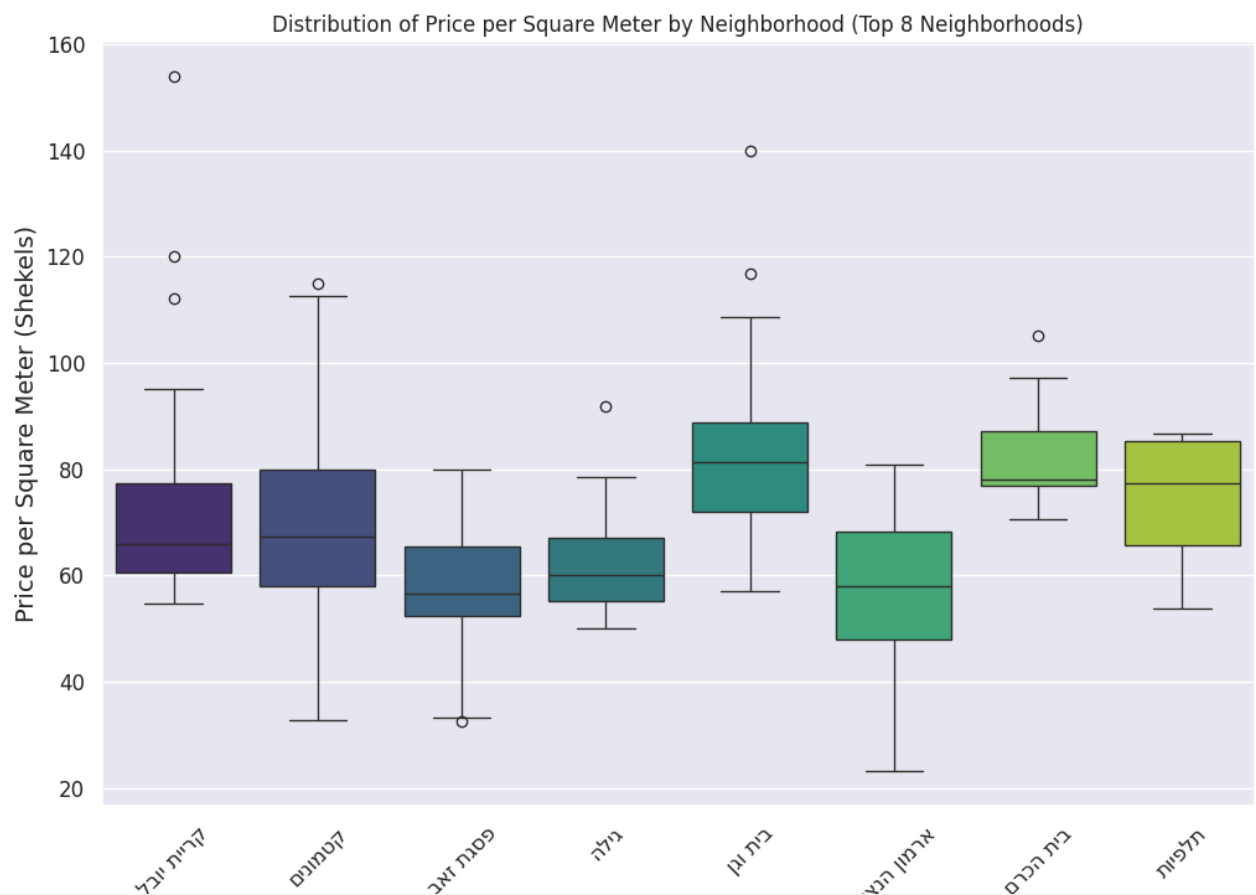
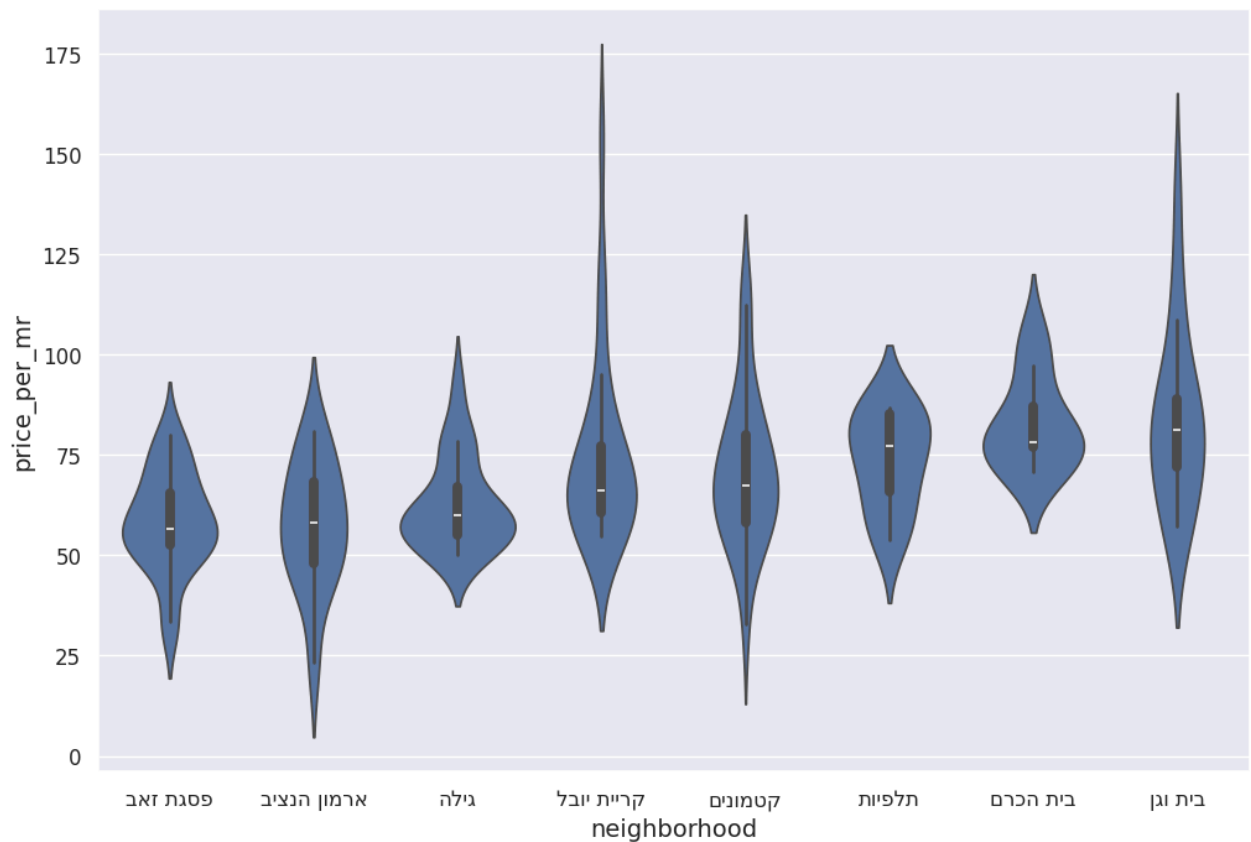
# Plot a box plot for these top 8 neighborhoods
plt.figure(figsize=(12, 8))
sns.boxplot(data=filtered_df_top8, x='neighborhood_flipped', y='price_per_mr', order=top_
plt.xlabel("Neighborhood")
plt.ylabel("Price per Square Meter (Shekels)")
plt.title("Distribution of Price per Square Meter by Neighborhood (Top 8 Neighborhoods)")
plt.xticks(rotation=45)
plt.show()
```



```
filtered_df['price_per_mr'] = filtered_df['monthlyRate'] / filtered_df['area']
<ipython-input-123-3a7256658a79>:24: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.

```
sns.boxplot(data=filtered_df_top8, x='neighborhood_flipped', y='price_per_mr', o
```



---

✓ Part 3 Question 6 - textual Answer:

*Write your answer here:*

נרמלנו את המחירים לפי גודלי הדירות ועכשיו ההתפלגויות נקיות יותר ואיזה כיף לנו

---

Given the conclusions from the previous steps, we may think that the apartment's neighborhood gives us additional information about the expected monthly rate. But the sample size for most neighborhoods is rather small. So let's examine another way to utilize the location information. Luckily, we also have data about the socio-economic rank of most neighborhoods (between 1 and 10).

✓ Question 7 - **bonus**

Use again the full dataset (without filtering by neighborhood).

Create an aggregated dataframe where every record represents a neighborhood, with columns for:

1. neighborhood name
2. flipped neighborhood name
3. The number of listings in a neighborhood
4. The median monthly rate for listings in this neighborhood.

Add a column with the neighborhood socio-economic rank to the dataframe (you can use the provided `get_neighborhood_rank` function that takes as an input a neighborhood name and returns its socio-economic rank.) Use this dataframe to visualize the association between socio-economic rank and pricing for all neighborhoods with at least 5 listings. What is your conclusion?

```
# Part 3 - Question 7
# Your code goes here:
#
#
```

---

Part 3 Question 7 - textual Answer:

*Write your answer here:*

---

✓ Part 4: Are private houses more expensive than apartments?

✓ Part 4 - Create a DataFrame and remove outliers for Part 4

```
# @title Part 4 - Create a DataFrame and remove outliers for Part 4
part4_df = rent_df_backup_for_exercise.copy()
part4_df = part4_df[part4_df['monthlyRate'] > 0].reset_index(drop=True);
part4_df = part4_df[part4_df['area'] < 800].reset_index(drop=True)
part4_df = part4_df[part4_df['area'] > 10].reset_index(drop=True)
```

Finally, we want to check if listings in private houses tend to be more expensive than apartments in a building.

**Use only `part4_df` for the coding questions in this part**

✓ **Question 1**