

Statistical Data Analysis

1. First of all, we have to generate CSV file that is mentioned in code. Lets do it:

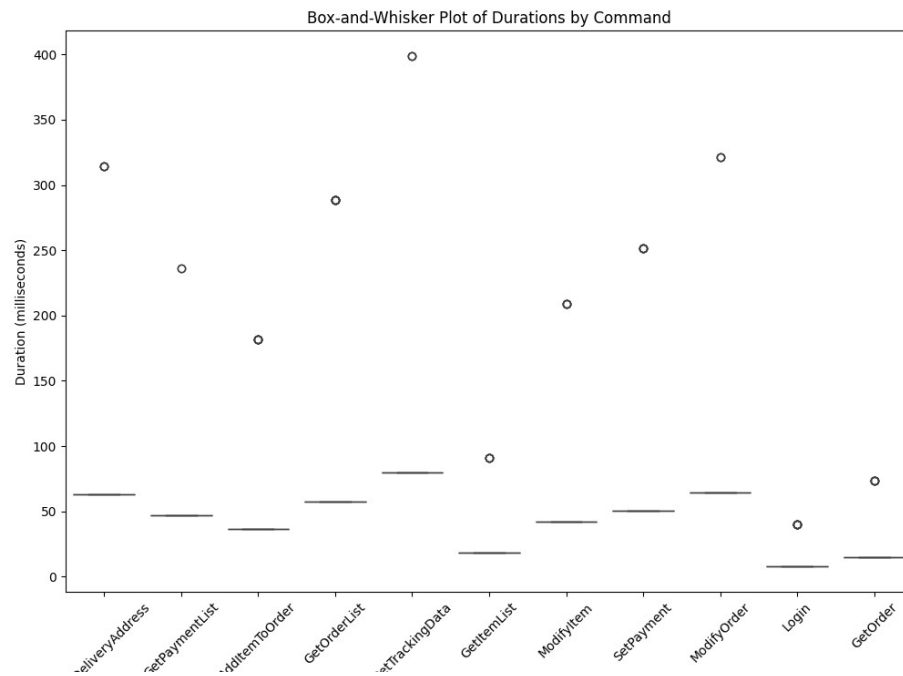
```
61 log_data.append([command,request_date.strftime('%Y-%m-%d %H:%M:%S.%f'),res]
62
63 print('Expected outliers: ',expected_outliers)
64
65 df = pd.DataFrame(log_data[1:], columns=log_data[0])
66
67 file_path = 'assignment.csv'
68 df.to_csv(file_path, index=False)
69
70
71
72
```

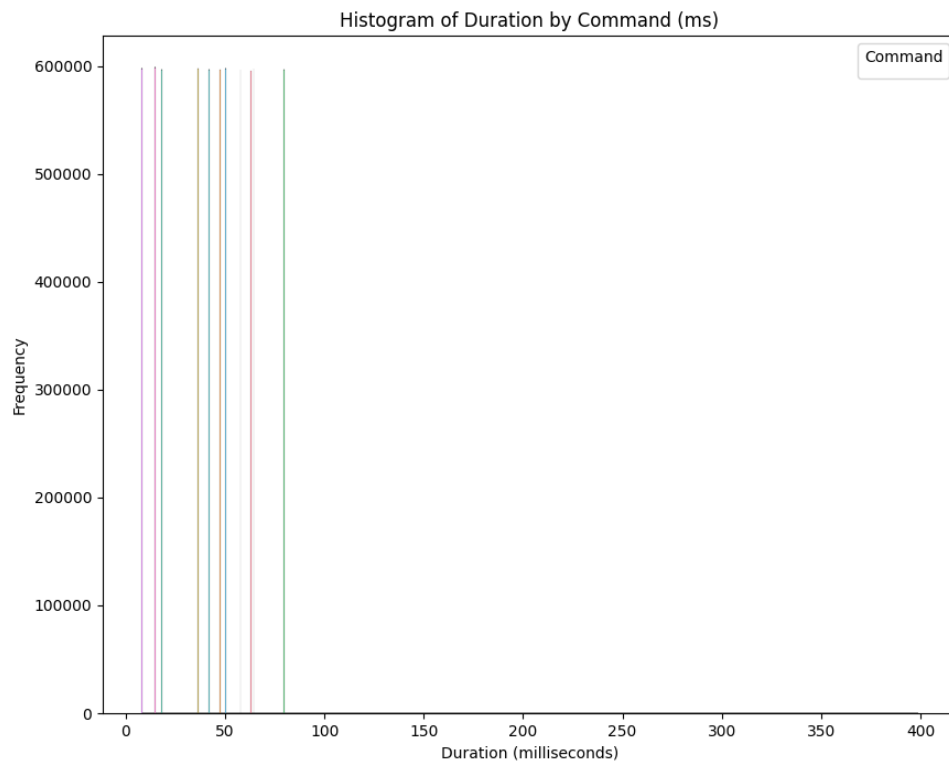
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS Python + - [] [X] ... ^ X

```
PS C:\Users\user\Desktop\POM2> & "C:/Program Files/Python312/python.exe" "c:/Users/user/Desktop/POM2/Python data generator.py"
Please enter your name/surname: Farid Mehdiyev
Please enter your student ID: 13670
Task generation time: 2024-05-19 00:58:22.791999
Expected outliers: {'GetOrderList': 1, 'GetPaymentList': 8, 'ModifyItem': 2, 'ModifyOrder': 1, 'AddItemToOrder': 3, 'GetTrackingData': 5, 'Login': 1, 'SetDeliveryAddress': 1, 'GetItemList': 2, 'GetOrder': 5, 'SetPayment': 2}
PS C:\Users\user\Desktop\POM2> |
```

First part is done, I will upload the “assignment.csv” file to OneDrive

2. Visualize the distribution and range of values for each command type (histogram and box and whisker diagram):





Commands are in order, as a

**“GetDeliveryAddress,
GetPaymentList,
AddItemtoOrde,
GetTrackingData,
GetItemList,
ModifyItem,
SetPayment,
ModifyOrder,
Login and GetOrder”**

```

1  from matplotlib.pylab import norm
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6
7  # Reading the provided assignement csv file
8  log_file = pd.read_csv('assignment.csv')
9
10 # Converting dates to datetime format
11 log_file['Request date'] = pd.to_datetime(log_file['Request date'])
12 log_file['Response date'] = pd.to_datetime(log_file['Response date'])
13
14 # Calculating the duration in milliseconds
15 log_file['Duration (msec)'] = (log_file['Response date'] - log_file['Request date']).dt.total_milliseconds
16
17
18 # Plotting the histogram
19 plt.figure(figsize=(10, 8))
20 sns.histplot(data=log_file, x='Duration (msec)', hue='Command', multiple='stack', kde=False)
21
22 #Title and Labels added to Histogram
23 plt.title('Histogram of Duration by Command (ms)')
24 plt.xlabel('Duration (milliseconds)')
25 plt.ylabel('Frequency')
26 plt.legend(title='Command')
27 plt.show()
28
29
30 # Plot the box-and-whisker plot
31 plt.figure(figsize=(12, 8))
32 sns.boxplot(x='Command', y='Duration (msec)', data=log_file)
33 plt.title('Box-and-Whisker Plot of Durations by Command')
34 plt.xlabel('Command')
35 plt.ylabel('Duration (milliseconds)')
36
37 # For readability, i rotated x-axis
38 plt.xticks(rotation=45)
39 plt.show()

```

Provided Python code generating both a histogram and box-and-whisker plot for visualize the distribution of durations by commands from a dataset that is provided as a “assignment.csv” file. Duration in milliseconds is calculated by the finding difference between “Response date” and “Request date” columns and that is converted the results to the millisecond.

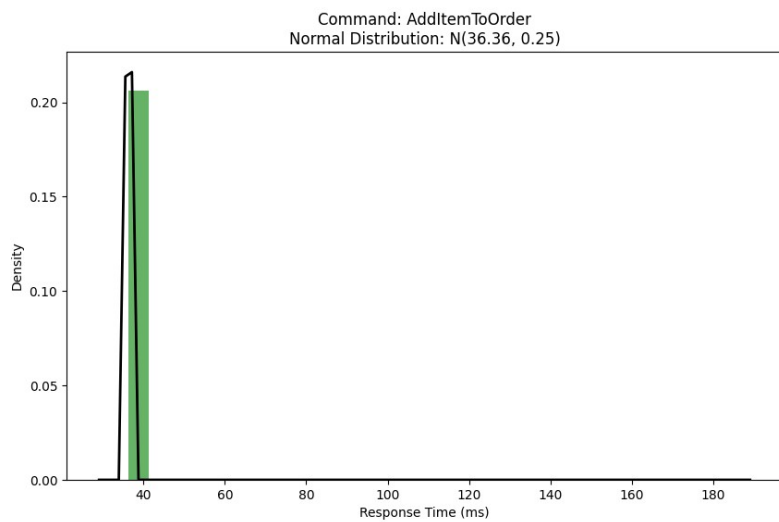
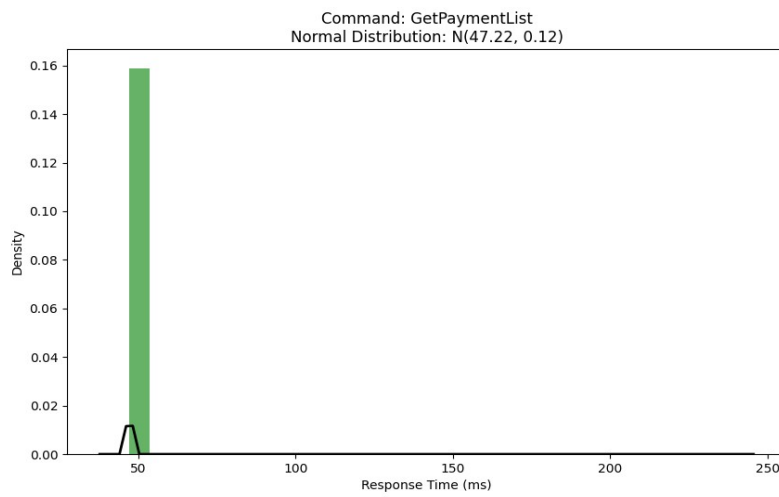
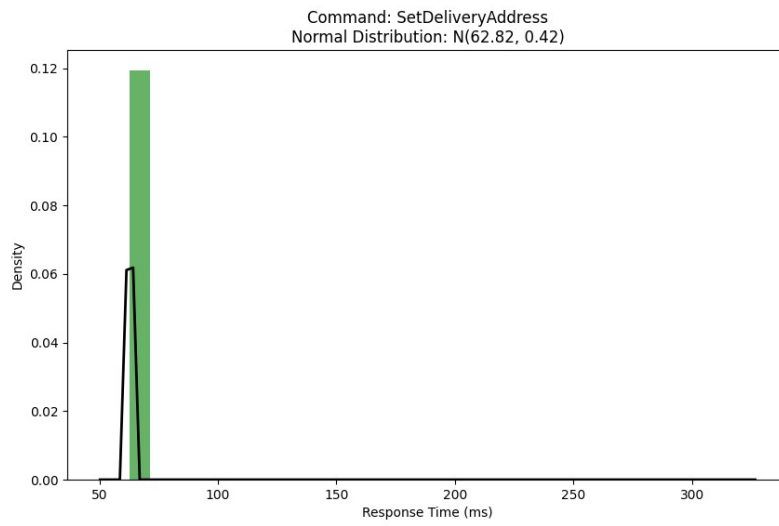
Histogram - for visualization, I have used Matplotlib and Seaborn Python libraries and they are very useful. Histogram is created using Seaborn’s ‘histplot’ function. Customization of the histogram is x and y axis labels for enhancing readability.

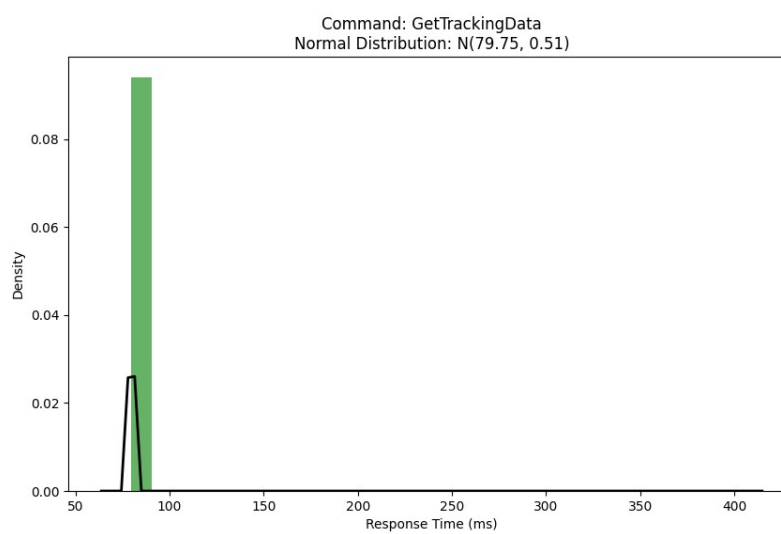
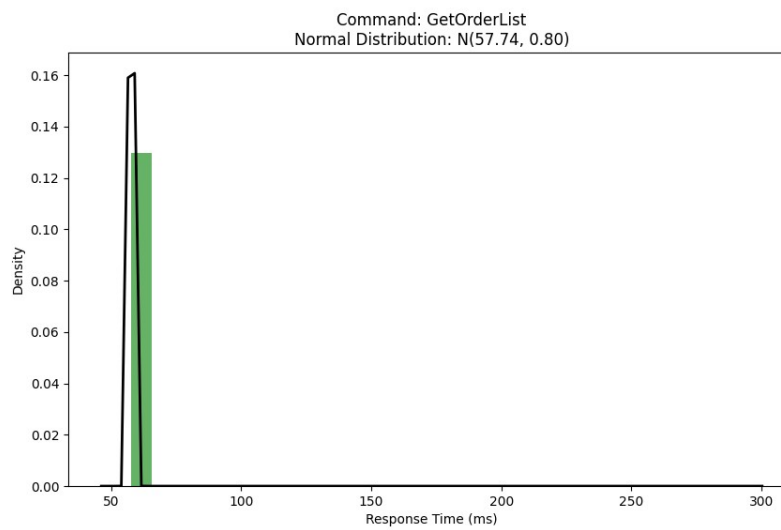
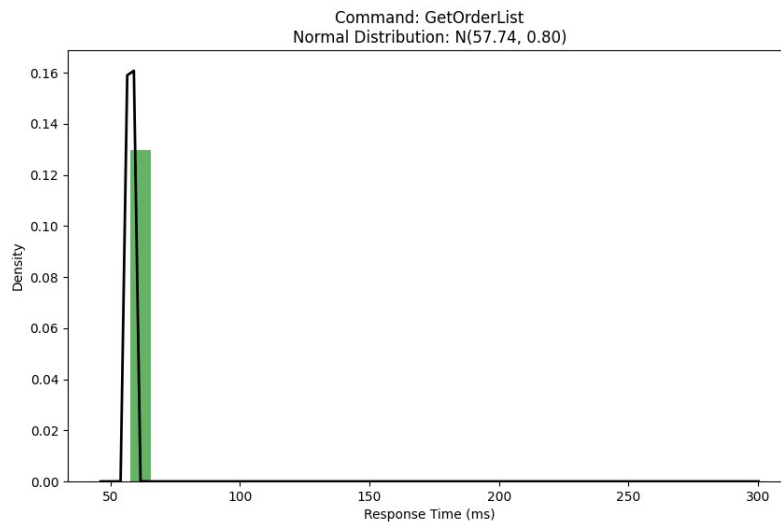
Box-and-Whisker Plot – it is the same as the histogram as you can see that – same data reading and processing log data. I have used ‘boxplot’ function for box-and-whisker plot and customization is same as histogram.

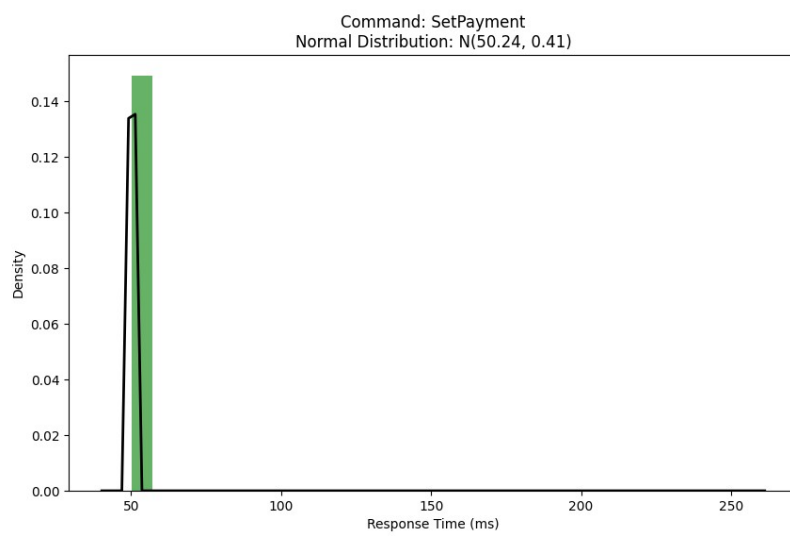
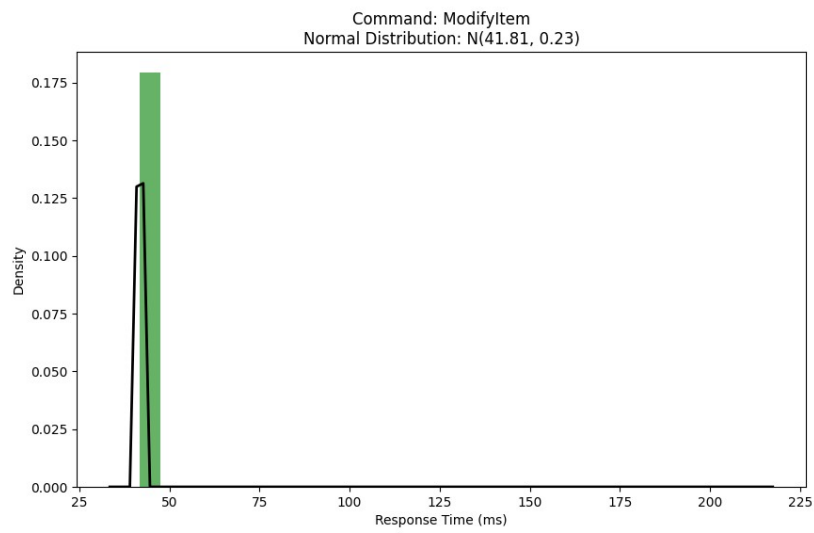
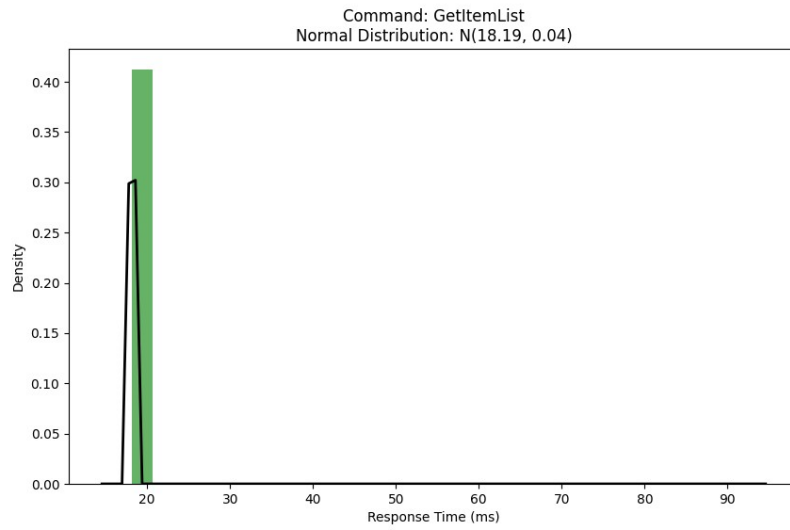
3. Show the distribution of each command as Normal distribution: $N(m, \text{variation})$.

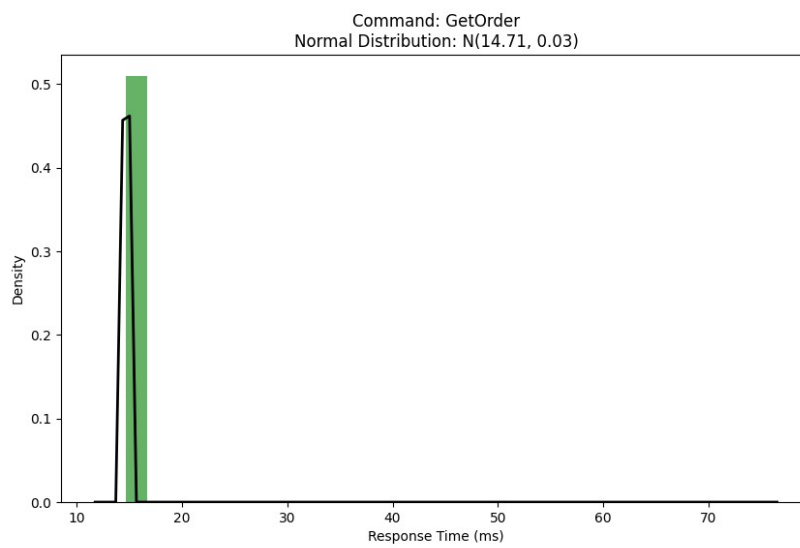
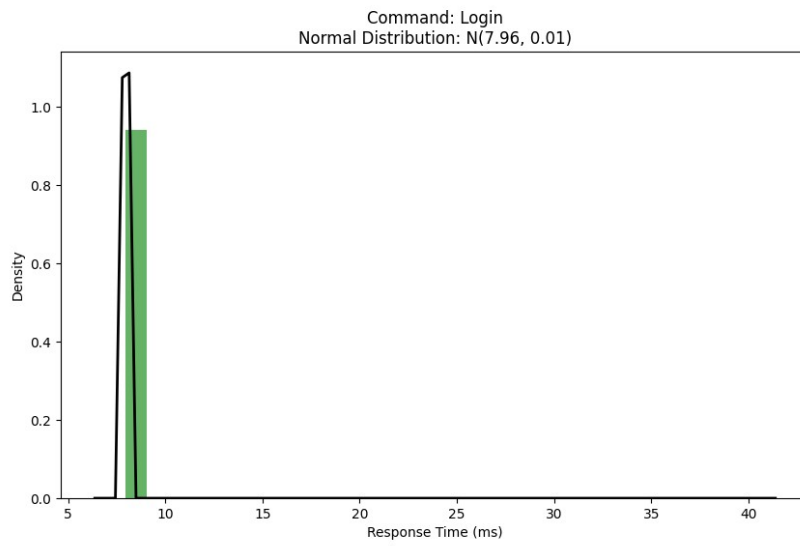
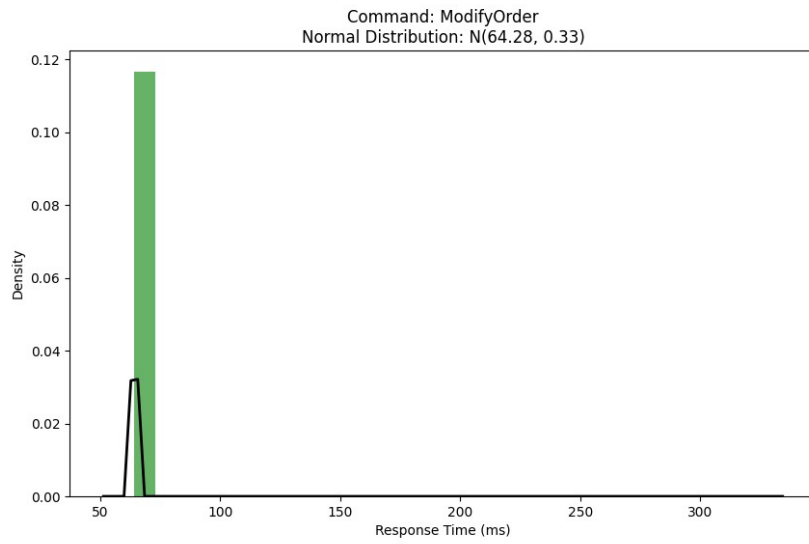
As we can observe that, code is used Pandas to read the CSV file into a DataFrame name '**df**'. For the statistical analysis by the calculating mean and variance of response times for each unique command in the dataset and likewise it iterates unique commands and subsets the *DataFrame* for calculating statistics for each command by separately. '**mean()**' and '**var()**' functions are used for finding mean and variance. For the visualization part I have used SciPy's statistical functions.

```
1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import numpy as np
4  import scipy.stats as stats
5
6  # Loading the generated CSV file that is provided
7  df = pd.read_csv('assignment.csv')
8
9  # Converting datetime strings to datetime objects
10 df['Request date'] = pd.to_datetime(df['Request date'])
11 df['Response date'] = pd.to_datetime(df['Response date'])
12
13 # Calculating response times in milliseconds
14 df['Response time'] = (df['Response date'] - df['Request date']).dt.total_seconds() * 1000
15
16 # Initializing a dictionary to store mean and variance for each command
17 command_stats = {}
18
19 # Calculating mean and variance for each command
20 for command in df['Command'].unique():
21     response_times = df[df['Command'] == command]['Response time']
22     mean = response_times.mean()
23     variance = response_times.var()
24     command_stats[command] = (mean, variance)
25
26 # Plot the distribution of response times
27 plt.figure(figsize=(10, 6))
28 plt.hist(response_times, bins=30, density=True, alpha=0.6, color='g')
29
30 # Plot the normal distribution with calculated mean and variance
31 xmin, xmax = plt.xlim()
32 x = np.linspace(xmin, xmax, 100)
33 p = stats.norm.pdf(x, mean, np.sqrt(variance))
34 plt.plot(x, p, 'k', linewidth=2)
35 plt.title(f'Command: {command}\nNormal Distribution: N({mean:.2f}, {variance:.2f})')
36 plt.xlabel('Response Time (ms)')
37 plt.ylabel('Density')
38 plt.show()
39
```



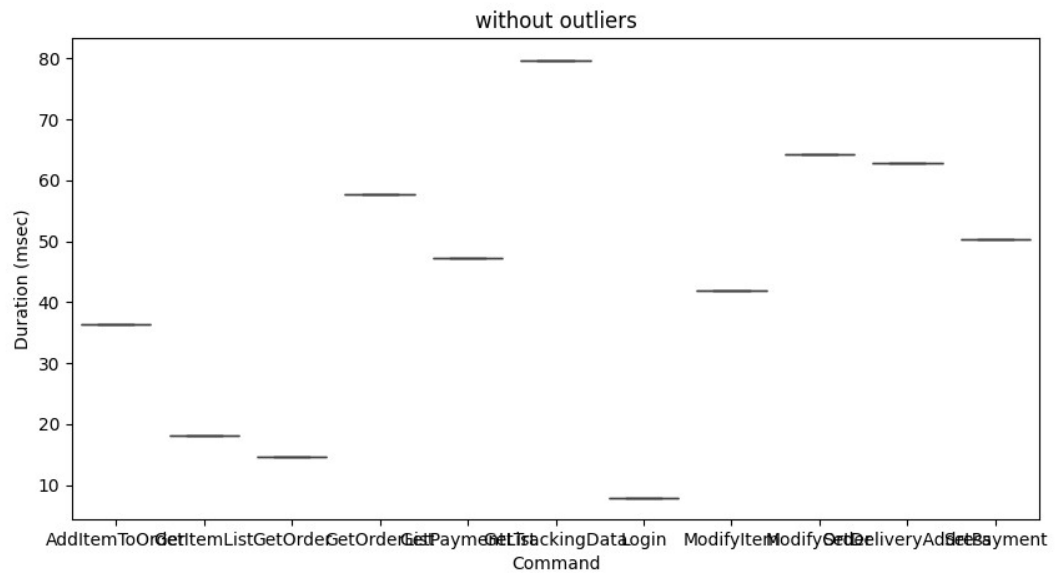




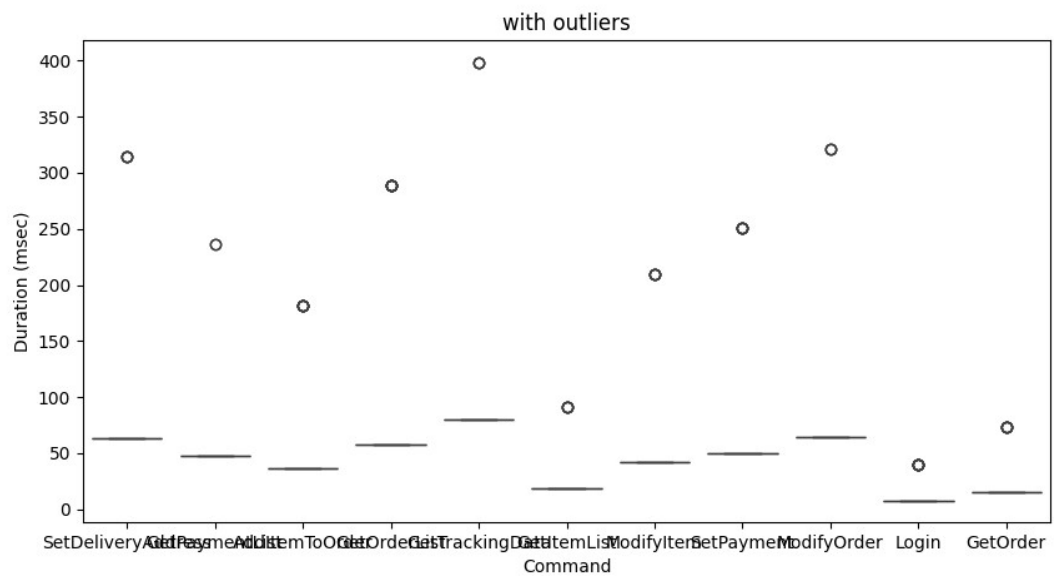


4. Show the outliers for each command and then perform item 2 after excluding the outliers.

Without outliers



With outliers



```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Reading the provided assignment CSV file
6 log_file = pd.read_csv('assignment.csv')
7
8 # Converting dates to datetime format
9 log_file['Request date'] = pd.to_datetime(log_file['Request date'])
10 log_file['Response date'] = pd.to_datetime(log_file['Response date'])
11
12 # Calculating the duration in milliseconds
13 log_file['Duration (msec)'] = (log_file['Response date'] - log_file['Request date']).dt.total_seconds() * 1000
14
15 def filter_outliers(df):
16
17     Quartile1 = df['Duration (msec)'].quantile(0.25) # Calculating first quartile as a Q1
18     Quartile3 = df['Duration (msec)'].quantile(0.75) # Same thing that i mentioned for Q3
19
20     InterQuartileRange = Quartile3 - Quartile1 # Calculating the InterQuartile Range (IQR)
21     LB = Quartile1 - 1.5 * InterQuartileRange
22     upper_bound = Quartile3 + 1.5 * InterQuartileRange
23
24
25     return df[(df['Duration (msec)'] >= LB)
26              & (df['Duration (msec)'] <= upper_bound)]
27
28 # Visualizing data with outliers
29
30 plt.figure(figsize=(10, 5))
31 sns.boxplot(x='Command',
32            y='Duration (msec)',
33            data=log_file)
34 plt.title('with outliers')
35 plt.show()
36
37 # Filtering out the outliers
38
39 clean_data = log_file.groupby('Command').apply(lambda x: filter_outliers(x)).reset_index(drop=True)
40
41 # Visualizing data without outliers
42
43 plt.figure(figsize=(10, 5))
44 sns.boxplot(x='Command'
45            , y='Duration (msec)'
46            , data=clean_data)
47
48 plt.title('without outliers')
49 plt.show()
50

```

As the previous files, first of all we have to reading csv file from 'assignment.csv' file into a dataframe named '**log_file**'. Converting dates I have used Pandas '**to_datetime**' functionality. A function named '**filter_outliers**' is defined to filter outliers from the date and function takes a dataframe as a '**df**' input.

For calculating **Quartile 1 (Q1)** and **Quartile 3 (Q3)** are calculated using the quantile function. The **Interquartile Range (IQR)** is calculated as the *difference* between *Q3 and Q1*.