

## ITD 112 Laboratory Exercises #2

Name: Clint Joshua O. Velasquez

**Complete the codes** of the following exercises. Screenshot your output. Submit in pdf format.

**Source Code Link:** <https://github.com/kiyojiii/ITD112>

**Exercise 2.01:** Using NumPy to Compute the Mean, Median, Variance, and Standard Deviation of a Dataset

The image displays two screenshots of a Jupyter Notebook interface, showing the progression of a data analysis task.

**Top Screenshot:**

- Header:** jupyter Exercise 2.01 Last Checkpoint: 13 hours ago (unsaved changes)
- Menu:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help
- Toolbar:** Includes icons for file operations, running cells, and a dropdown menu set to 'Markdown'.
- Section Title:** Exercise 2.01: Using NumPy to Compute the Mean, Median, Variance, and Standard Deviation of a Dataset
- Description:** In this activity, you will consolidate the skills you've acquired in the last exercise and use NumPy to do some very basic mathematical calculations on our normal\_distribution dataset. NumPy has a consistent API, so it should be rather easy to transfer your knowledge of the mean method to median and variance.
- Code Cell 1:**

```
In [2]: # importing the necessary dependencies
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```
- Code Cell 2:**

```
In [4]: #1 Load the "normal_distribution.csv" dataset
dataset = np.genfromtxt('C:/Users/user/Desktop/ITD112/Laboratory_Exercise_2/Datasets/normal_distribution.csv', delimiter=',')
```
- Code Cell 3:**

```
In [5]: #Look at the dataset
dataset
```
- Output:** Out[5]: array([[ 99.14931546, 104.03852715, 107.43534677, 97.85230675, 98.74986914, 98.80833412, 96.81964892, 98.56783189],
 [ 92.02628776, 97.10439252, 99.32066924, 97.24584816, 92.9267508 , 92.65657752, 105.7197853 , 101.23162942],
 [ 95.66253664, 95.17750125, 90.93318132, 110.18889465,

**Bottom Screenshot:**

- Code Cell 4:**

```
In [6]: #2 Look at the first two rows of the dataset
dataset[0:2]
```
- Output:** Out[6]: array([[ 99.14931546, 104.03852715, 107.43534677, 97.85230675, 98.74986914, 98.80833412, 96.81964892, 98.56783189],
 [ 92.02628776, 97.10439252, 99.32066924, 97.24584816, 92.9267508 , 92.65657752, 105.7197853 , 101.23162942]])
- Section Title:** Mean
- Code Cell 5:**

```
In [7]: #3 calculate the mean of the third row
np.mean(dataset[2, :])
```
- Output:** Out[7]: 100.20466135250001
- Code Cell 6:**

```
In [8]: #4 calculate the mean of the last column
np.mean(dataset[:, -1])
```
- Output:** Out[8]: 100.4404927375
- Code Cell 7:**

```
In [9]: #5 calculate the mean of the intersection of the first 3 rows and first 3 columns
# Extract the intersection of the first 3 rows and first 3 columns
intersection = dataset[4:, :4]
```

Jupyter Exercise 2.01 Last Checkpoint: 13 hours ago (unsaved changes)

Python 3 (ipykernel)

File Edit View Insert Cell Kernel Widgets Help

Out[8]: 100.4404927375

In [9]: #5 calculate the mean of the intersection of the first 3 rows and first 3 columns

# Extract the intersection of the first 3 rows and first 3 columns  
intersection = dataset[4:, :4]  
  
# Calculate the Mean of the intersection  
intersection\_mean = np.mean(intersection)  
  
print(intersection\_mean)  
  
100.80743067

Median

In [10]: #6 calculate the median of the last row

np.median(dataset[-1, :])

Out[10]: 99.18748092

In [11]: #7 calculate the median of the last 3 columns

np.median(dataset[:, -3:])

Out[11]: 99.08416696500001

Jupyter Exercise 2.01 Last Checkpoint: 13 hours ago (unsaved changes)

Python 3 (ipykernel)

File Edit View Insert Cell Kernel Widgets Help

In [12]: #8 calculate the median of each row

np.median(dataset, axis=1)

Out[12]: array([ 98.77910163, 97.17512034, 98.58782879, 100.68449836,  
101.00170737, 97.76908825, 101.85002253, 100.04756697,  
102.24292555, 99.59514997, 100.4955753 , 99.8860714 ,  
 99.00647994, 98.67276177, 102.44376222, 96.61933565,  
104.0968893 , 100.72023043, 98.70877396, 99.75008654,  
104.89344428, 101.00634942, 98.30543801, 99.18748092])

Variance

In [13]: #10 calculate the variance of each column

np.var(dataset, axis=0)

Out[13]: array([23.64757465, 29.78886109, 20.50542011, 26.03204443, 28.38853175,  
19.09960817, 17.67291174, 16.17923204])

In [14]: #11 calculate the variance of the intersection of the last 2 rows and first 2 columns

# Extract the intersection of the last 2 rows and first 2 columns  
intersection = dataset[-3:, :3]  
  
# Calculate the variance of the intersection  
variance\_intersection = np.var(intersection)

Jupyter Exercise 2.01 Last Checkpoint: 13 hours ago (autosaved)

Python 3 (ipykernel)

File Edit View Insert Cell Kernel Widgets Help

19.09960817, 17.67291174, 16.17923204])

In [14]: #11 calculate the variance of the intersection of the last 2 rows and first 2 columns

# Extract the intersection of the last 2 rows and first 2 columns  
intersection = dataset[-3:, :3]  
  
# Calculate the variance of the intersection  
variance\_intersection = np.var(intersection)  
  
print(variance\_intersection)  
  
16.046702056236402

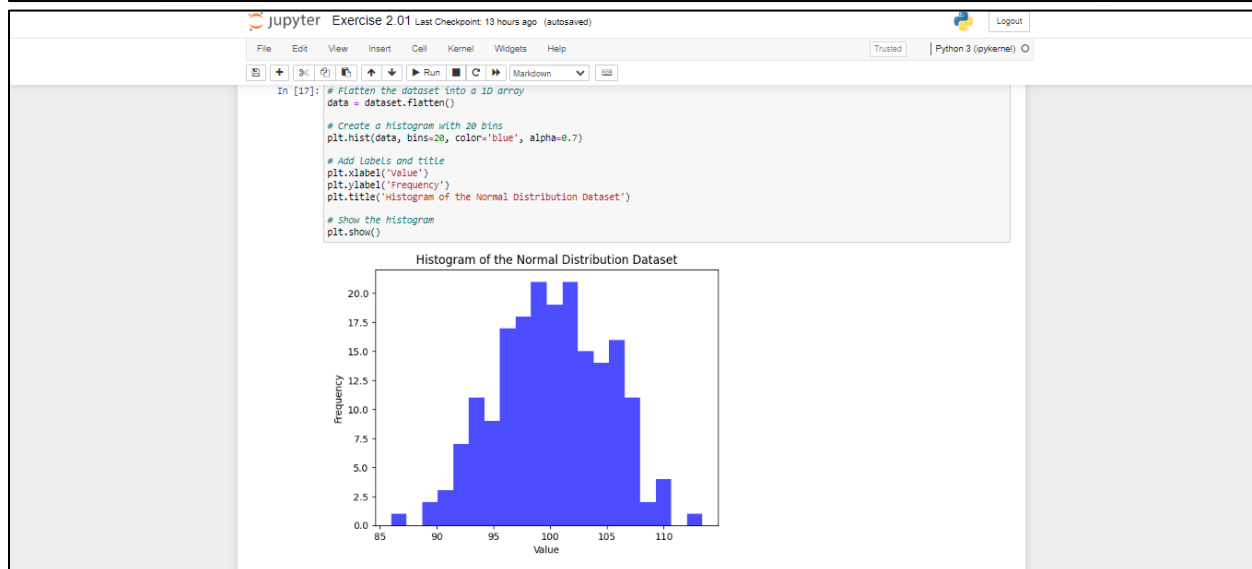
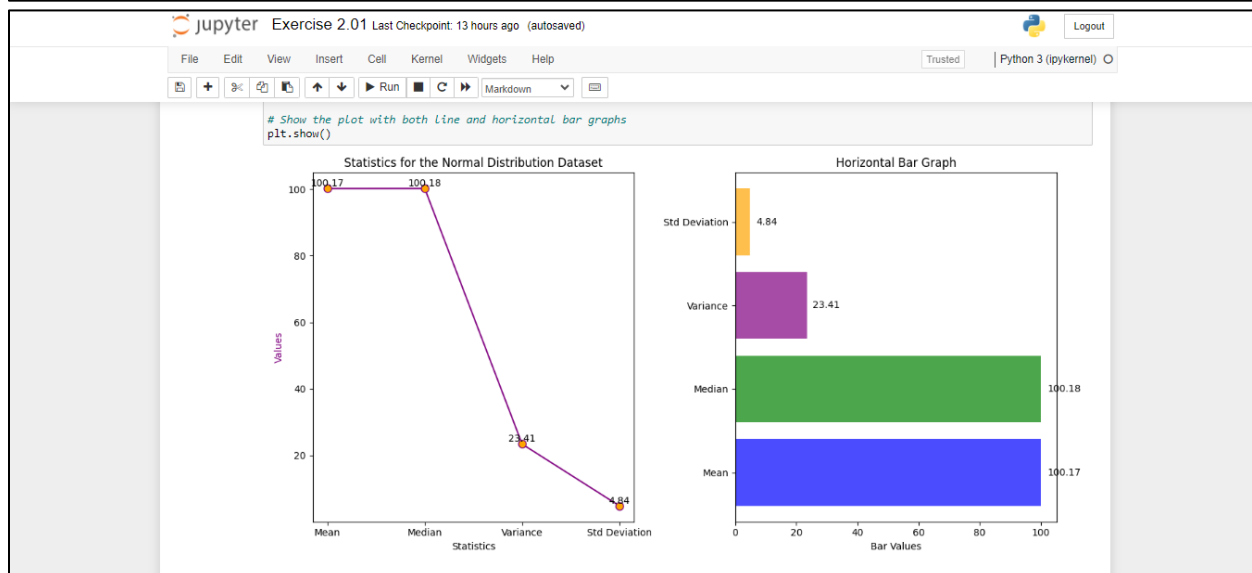
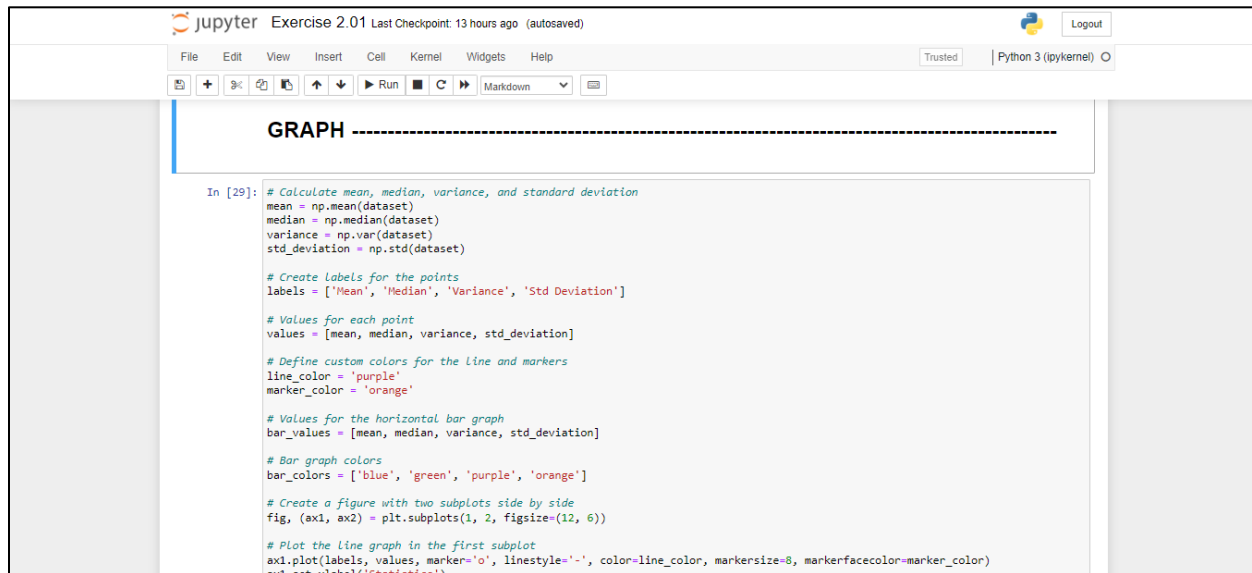
The values of the variance might seem a little bit strange at first.  
You can always go back to the topic that gives you a quick statistical overview to recap what you've learned so far.

Note:  
Just remember, the variance is not the standard deviation.

Try calculation the standard deviation with NumPy to get a more descriptive value when comparing it to our dataset

In [32]: #12 calculate the standard deviation for the dataset

# Calculate the standard deviation for the entire dataset  
std\_deviation = np.std(dataset)  
  
print(std\_deviation)  
  
4.838197554269257



## Exercise 2.02: Forest Fire Size and Temperature Analysis



**Jupyter Exercise 2.02** Last Checkpoint: 11 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted | Python 3 (ipykernel) O

In [18]: `#5 get the smallest area value from our dataset  
smallest_area = filtered_dataset['area'].min()  
print(smallest_area)`

0.09

In [19]: `#6 get the largest area value from our dataset  
largest_area = filtered_dataset['area'].max()  
print(largest_area)`

1099.84

In [20]: `#7 get the standard deviation of values in our dataset  
std_deviation = filtered_dataset['area'].std()  
print(std_deviation)`

86.50163460412126

The largest value is much larger than our mean.  
The standard deviation also is quite large which indicates that the difference between our mean and the "middle value" will be quite high.

Let's look at the last 20 values of our sorted dataset to see if we have more than one very large value.  
Sort the filtered dataset by the `area` column and output the last 20 entries from it.

In [21]: `#8 sorting the filtered dataset and printing the last 20 elements  
sorted_dataset = filtered_dataset.sort_values(by='area', ascending=False)  
last_20_entries = sorted_dataset.tail(20)  
print(last_20_entries)`

**Jupyter Exercise 2.02** Last Checkpoint: 11 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted | Python 3 (ipykernel) O

Sort the filtered dataset by the `area` column and output the last 20 entries from it.

In [21]: `#8 sorting the filtered dataset and printing the last 20 elements  
sorted_dataset = filtered_dataset.sort_values(by='area', ascending=False)  
last_20_entries = sorted_dataset.tail(20)  
print(last_20_entries)`

	Y	month	day	FFHC	DHC	DC	ISI	temp	RH	wind	rain	area
X	2	4	aug	92.2	91.6	503.6	9.6	20.7	70	2.2	0.0	0.75
	1	3	sep	91.2	94.3	744.4	8.4	22.3	48	4.0	0.0	0.72
	1	2	jul	90.0	51.3	296.3	8.7	16.6	53	5.4	0.0	0.71
	8	5	aug	93.1	157.3	666.7	13.5	26.8	25	3.1	0.0	0.68
	4	3	aug	94.2	117.2	581.1	11.0	21.4	44	2.7	0.0	0.68
	8	6	aug	90.1	108.0	529.8	12.5	21.2	51	8.9	0.0	0.61
	1	2	aug	95.5	99.9	513.3	13.2	23.3	31	4.5	0.0	0.55
	8	6	aug	92.1	207.0	672.6	8.2	26.8	35	1.3	0.0	0.54
	2	2	aug	92.1	152.6	658.2	14.3	21.8	56	3.1	0.0	0.52
	7	4	sep	88.2	55.2	732.3	11.6	15.2	64	3.1	0.0	0.52
	2	5	sep	90.9	126.5	686.5	7.0	21.9	39	1.8	0.0	0.47
	6	5	aug	91.0	166.9	752.6	7.1	18.2	62	5.4	0.0	0.43
	1	4	sep	91.0	129.5	692.6	7.0	21.7	38	2.2	0.0	0.43
	7	4	sep	89.6	84.1	714.3	5.7	17.1	53	5.4	0.0	0.41
	9	9	jul	85.8	48.3	313.4	3.9	18.0	42	2.7	0.0	0.36
	1	3	sep	91.1	91.3	738.1	7.2	19.1	46	2.2	0.0	0.33
	8	5	aug	93.1	157.3	666.7	13.5	24.0	36	3.1	0.0	0.24
	6	5	aug	93.1	157.3	666.7	13.5	22.1	37	3.6	0.0	0.21
	6	5	aug	94.3	131.7	607.1	22.7	19.4	55	4.0	0.0	0.17
	5	4	aug	91.8	175.1	700.7	13.8	25.7	39	5.4	0.0	0.09

As we can see here, only 11 out of the 270 rows contain values that are larger than 100.  
After 20 values we are close to the area value of 60.

Let's imagine our dataset contained only 1 or 2 values that were much higher than the other ones, e.g. an area size value of 10254.91. Simply by observing the dataset, this feels like there might have been an error on adding this to the dataset.  
In a smaller dataset, the mean value would get heavily distorted by this one entry. A more stable value to use in such a case is the median value of the dataset.

**Jupyter Exercise 2.02** Last Checkpoint: 11 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted | Python 3 (ipykernel) O

Use the median value for the `area` column.

In [22]: `#9 calculate the median value for the area column  
median_area = filtered_dataset['area'].median()  
print(median_area)`

6.37

**Note:**  
Remember that the median is not the same as then mean of your dataset. While the median is simple the "value in the middle", the mean is much more prone to distortion by outliers.

**Finding the month with the most forest fires**

In this second task we want to quickly see which months have the most forest fires and whether or not the temperature has a direct connection to it.

Get a list of month values that are present in our dataset.

In [23]: `#10 get a list of month values from the dataset  
unique_months = dataset['month'].unique()  
print(unique_months)`

['mar' 'oct' 'aug' 'sep' 'apr' 'jun' 'jul' 'feb' 'jan' 'dec' 'may' 'nov']

In addition to the unique values we also want use the shape element of our dataset to determine how many rows it has.

Filter the dataset for only rows that contain the month `mar` and print the number of rows using `shape`.

In [62]: `#11 get the number of forest fires for the month of march  
fires_in_march = dataset[dataset['month'] == 'mar'].shape[0]  
print(fires_in_march)`

Jupyter Exercise 2.02 Last Checkpoint: 11 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

```
In [62]: #11 get the number of forest fires for the month of march
fires_in_march = dataset[dataset['month'] == 'mar'].shape[0]
print(fires_in_march)
54
```

The last step to fulfill the task is to iterate over all months, filtering our dataset for the rows containing the given month and calculating the mean temperature.

- Iterate over the months from the unique list we created
- Filter our dataset for the rows containing the given month
- Get the number of rows from `shape`
- Get the mean temperature for the given month
- Print a statement with the number of fires, mean temperature and the month

```
In [78]: # 12-15
# iterate over the months list
# get number of forest fires for each month
# get mean temperature for each month
# print out number of fires and mean temperature

# Define a custom order for months
custom_month_order = ['jan', 'feb', 'mar', 'apr', 'may', 'jun', 'jul', 'aug', 'sep', 'oct', 'nov', 'dec']

# Get a list of month values from the dataset
unique_months = dataset['month'].unique()

# Iterate over the custom order of months
for month in custom_month_order:
    if month in unique_months:
        month_data = dataset[dataset['month'] == month]
        num_fires = month_data.shape[0]
        mean_temp = month_data['temp'].mean()
        print(f'Month: {month}, Number of Fires: {num_fires}, Mean Temperature: {mean_temp:.2f}')

Month: jan, Number of Fires: 2, Mean Temperature: 5.25
```

Jupyter Exercise 2.02 Last Checkpoint: 11 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

```
num_fires = month_data.shape[0]
mean_temp = month_data['temp'].mean()
print(f'Month: {month}, Number of Fires: {num_fires}, Mean Temperature: {mean_temp:.2f}')
```

Month: jan, Number of Fires: 2, Mean Temperature: 5.25  
Month: feb, Number of Fires: 20, Mean Temperature: 9.63  
Month: mar, Number of Fires: 54, Mean Temperature: 13.88  
Month: apr, Number of Fires: 9, Mean Temperature: 12.04  
Month: may, Number of Fires: 2, Mean Temperature: 14.65  
Month: jun, Number of Fires: 17, Mean Temperature: 20.49  
Month: jul, Number of Fires: 32, Mean Temperature: 22.11  
Month: aug, Number of Fires: 184, Mean Temperature: 21.63  
Month: sep, Number of Fires: 172, Mean Temperature: 19.61  
Month: oct, Number of Fires: 15, Mean Temperature: 17.09  
Month: nov, Number of Fires: 1, Mean Temperature: 11.80  
Month: dec, Number of Fires: 9, Mean Temperature: 4.52

**GRAPH -----**

```
In [27]: # Get a list of month values from the dataset
unique_months = dataset['month'].unique()

# Initialize lists to store data for the graph
num_fires_per_month = []
mean_temp_per_month = []

# Iterate over the months list, get the number of forest fires and the mean temperature, and store the data
for month in unique_months:
    month_data = dataset[dataset['month'] == month]
    num_fires = month_data.shape[0]
    mean_temp = month_data['temp'].mean()
    num_fires_per_month.append(num_fires)
    mean_temp_per_month.append(mean_temp)
```

