## ⌄  Hands-on Activity 6.1 Introduction to Data Analysis and Tools

CPE311 Computational Thinking with Python

**Name:** De los Reyes, Jann Moises Nyll B. **\*Section: \*** CPE22S3

**Performed on :** 03/06/2024

**Submitted on :** 03/06/2024

**Submitted to :** 03/06/2024

## 6.1 Intended Learning Outcome

1. Use pandas and numpy data analysis tools.
2. Demonstrate how to analyze data using numpy and pandas

## 6.2 Resources:

- Personal Computer
- Jupyter Notebook
- Internet Connection

**6.3 Supplementary Activities:**

**Exercise 1**

Run the given code below for exercises 1 and 2, perform the given tasks without using any Python modules

```
1 import random
2 random.seed(0)
3 salaries = [round(random.random()*1000000, -3) for _ in range(100)]
```

Using the data generated above, calculate the following statistics without importing anything from the statistics module in the standard library (https://docs.python.org/3/library/statistics.html) and then confirm your results match up to those that are obtained when using the statistics module (where possible):

- Mean
- Median
- Mode (hint: check out the Counter in the collections module of the standard library at https://docs.python.org/3/library/collections.html#collections.Counter)
- Sample variance
- Sample standard deviation

```
1 #Write a comment per statistical function
2
3 def mean(dataset):
4   # Calculate mean by the sum of all dataset
5   #divided by the number of the data
6   return sum(dataset)/len(dataset)
7
8 #To test
9 mean(salaries)
```

    585690.0

```
1 def median(dataset):
2   #Sort the dataset in increasing order
3   dataset.sort()
4   #Get the number of data
5   length = len(dataset)
6   # if the number of data is odd, Median = [(n + 1)/2] th term
7   if length % 2 == 1:
8     return dataset[int(length/2)]
9   #If the number of data is even,
10   #Median = [ (n/2) th term + (n/2 + 1) th term]/2.
11   else:
12     x = dataset[int((length/2)-1)]
13     y =dataset[int((length/2))]
14     ave= (x+y)/2
15     return ave
16
17 #To test
18 print(median(salaries))
```

    589000.0

```
1 import collections
2
3 def mode(dataset):
4    #create a dictionary-like object that count the occurence of each unique element in the dataset.
5    # the .most_common() method sort the element by frequency in desceding order and return them as a list of tuple.
6    modes = collections.Counter(dataset).most_common()
7
8    #the indexing retrieves the most common element(mode) from the sorted list
9    return modes[0][0]
10
11 #To test
12 mode(salaries)
```

    477000.0

```
1 def sample_variance(dataset):
2    n = len(dataset)
3    mean = sum(dataset) / n
4
5    # Calculate the squared deviations
6    squared_deviations = [(x - mean) ** 2 for x in dataset]
7
8    # Sum up the squared deviations
9    sum_squared_deviations = sum(squared_deviations)
10
11    # Compute sample variance
12    sample_variance = sum_squared_deviations / (n - 1)
13    return sample_variance
14
15 #To test:
16 var = sample_variance(salaries)
17 #Display the sample variance  with  6 decimal places
18 print(f"{var:.6f}")
```

    70664054444.444443

```
1 def sample_std(dataset):
2    n = len(dataset)
3    mean = sum(dataset)/n
4
5    # Calculate the squared deviations
6    squared_deviations = [(x - mean) ** 2 for x in dataset]
7
8    # Sum up the squared deviations
9    sum_squared_deviations = sum(squared_deviations)
10
11    # Compute sample variance
12    sample_std = (sum_squared_deviations / (n - 1))**0.5
13    return sample_std
14
15 #To test:
16 std =sample_std(salaries)
17 #Display the sample standard deviation with  6 decimal places
18 print(f"{std:.6f}")
```

    265827.113825

**Exercise 2**

Using the same data, calculate the following statistics using the functions in the statistics module where appropriate:

- Range
- Coefficient of variation
- Interquartile range
- Quartile coefficient of dispersion

```
1 #Write a comment per statistical function
2
3 import statistics
4
5 def calc_range(dataset):
6    #To check if the data set contains at least two values
7    if len(dataset)< 2:
8      raise ValueError("The data must contain at least two values.")
9
10    #Compute the range by subtracting the maximum and minimum value of the dataset.
11    return max(dataset)- min(dataset)
12
13 range_value = calc_range(salaries)
14 print(f" The range value is {range_value:.2f}")
```

     The range value is 995000.00

```
1 from statistics import stdev,mean
2
3 def calc_COV(dataset):
4    #Compute the  coefficient variation by calculating
5    # the standard deviation over the mean of the dataset
6    COV = stdev(dataset)/mean(dataset)
7
8    #Convert the COV in percentage
9    P_COV = COV  * 100
10
11   calc_COV = print(f" Coefficient of Variation: {COV}\n COV in Percentage : {P_COV:.2f}")
12   return calc_COV
13
14 #To test:
15 calc_COV(salaries)

     Coefficient of Variation: 0.45386998894439035
     COV in Percentage : 45.39 %
```

```
1 from statistics import quantiles
2 # Computation for interquartile range
3 # 1st step
4 def calc_quartiles(dataset):
5
6    #use the .quantiles() module in which .quantiles(data, n=4 )
7    #Divide data into n continuous intervals with equal probability.
8    #Returns a list of n - 1 cut points separating the intervals.
9    quartiles = quantiles(dataset, n = 4)
10
11   return quartiles
12
13 #To test:
14 calc_quartiles(salaries)
15
16

     [400500.0, 589000.0, 822250.0]
```

```
1 # 2nd Step
2 def calc_IQR(dataset):
3    #Assign a  variable  for the quartile list.
4    Qlist = quantiles(dataset)
5
6    #Compute for Interquartile Range where IQR = Q3 -Q1
7    calc_IQR = Qlist[-1]- Qlist[0]
8    return calc_IQR
9
10 #To test:
11 calc_IQR(salaries)

     421750.0
```

```
1 def calc_QCD(dataset):
2
3    #Assign a variable for the quartile list.
4    Qlist = quantiles(dataset)
5
6    #compute for  Quartile Coefficient of Dispersion where
7    #QCD =  (Q3-Q1)/(Q3+Q1) * 100
8    #use the calc_IQR for the numerator
9    QCD = calc_IQR(dataset)/(Qlist[-1]+ Qlist[0])
10
11   #Convert the QCD in percentage
12   P_QCD = QCD * 100
13
14   #Display the value of QCD
15   calc_QCD = print(f"Quartile Coefficient of Dispersion: {QCD} \nQCD in Percentage: {P_QCD:.6f} %")
16
17   return calc_QCD
18
19 calc_QCD(salaries)
20

     Quartile Coefficient of Dispersion: 0.34491923941934166
     QCD in Percentage: 34.491924 %
```

**Exercise 3: Pandas for Data Analysis**

Load the diabetes.csv file. Convert the diabetes.csv into dataframe Perform the following tasks in the diabetes dataframe:

1. Identify the column names

2. Identify the data types of the data

3. Display the total number of records

4. Display the first 20 records

5. Display the last 20 records

6. Change the Outcome column to Diagnosis

7. Create a new column Classification that display "Diabetes" if the value of outcome is 1 , otherwise "No Diabetes"

8. Create a new dataframe "withDiabetes" that gathers data with diabetes

9. . Create a new dataframe "noDiabetes" thats gathers data with no diabetes

10. Create a new dataframe "Pedia" that gathers data with age 0 to 19

11. Create a new dataframe "Adult" that gathers data with age greater than 19

12. Use numpy to get the average age and glucose value.

13. Use numpy to get the median age and glucose value.

14. Use numpy to get the middle values of glucose and age

15. Use numpy to get the standard deviation of the skinthickness.

```
1 # Indicate which item you're answering with a comment
```

```
1 #Uploading  diabetes.csv file and converting dataframe
2 filepath = '/content/diabetes.csv'
3 import numpy as np
4 import pandas as pd
5 data = pd.read_csv(filepath)
6
7 data
```

|  | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 | 0 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 | 0 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 | 0 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 | 1 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 | 0 |

768 rows × 9 columns

Next steps:  ◉ View recommended plots

```
1 #converting the data as dataframe
2 exercise_df = pd.DataFrame(data)
3 exercise_df
```

|  | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 | 0 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 | 0 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 | 0 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 | 1 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 | 0 |

768 rows × 9 columns

Next steps:  ◉ View recommended plots

```
1 #1. Identify the column names
2 exercise_df = pd.DataFrame(data)
3 exercise_df.columns
```

```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome', 'Classification'],
      dtype='object')
```

```
1 #2. Identify the data types of the data
2 exercise_df = pd.DataFrame(data)
3 exercise_df.dtypes
```

```
Pregnancies                 int64
Glucose                     int64
BloodPressure               int64
SkinThickness               int64
Insulin                     int64
BMI                         float64
DiabetesPedigreeFunction    float64
Age                         int64
Outcome                     int64
dtype: object
```

```
1 #3. Display the total number of records
2 exercise_df = pd.DataFrame(data)
3 len(exercise_df)
```

```
768
```

```
1 #4 Display the first 20 records
2 exercise_df.iloc[:20]
```

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| 5 | 5 | 116 | 74 | 0 | 0 | 25.6 | 0.201 | 30 | 0 |
| 6 | 3 | 78 | 50 | 32 | 88 | 31.0 | 0.248 | 26 | 1 |
| 7 | 10 | 115 | 0 | 0 | 0 | 35.3 | 0.134 | 29 | 0 |
| 8 | 2 | 197 | 70 | 45 | 543 | 30.5 | 0.158 | 53 | 1 |
| 9 | 8 | 125 | 96 | 0 | 0 | 0.0 | 0.232 | 54 | 1 |
| 10 | 4 | 110 | 92 | 0 | 0 | 37.6 | 0.191 | 30 | 0 |
| 11 | 10 | 168 | 74 | 0 | 0 | 38.0 | 0.537 | 34 | 1 |
| 12 | 10 | 139 | 80 | 0 | 0 | 27.1 | 1.441 | 57 | 0 |
| 13 | 1 | 189 | 60 | 23 | 846 | 30.1 | 0.398 | 59 | 1 |
| 14 | 5 | 166 | 72 | 19 | 175 | 25.8 | 0.587 | 51 | 1 |
| 15 | 7 | 100 | 0 | 0 | 0 | 30.0 | 0.484 | 32 | 1 |
| 16 | 0 | 118 | 84 | 47 | 230 | 45.8 | 0.551 | 31 | 1 |
| 17 | 7 | 107 | 74 | 0 | 0 | 29.6 | 0.254 | 31 | 1 |
| 18 | 1 | 103 | 30 | 38 | 83 | 43.3 | 0.183 | 33 | 0 |
| 19 | 1 | 115 | 70 | 30 | 96 | 34.6 | 0.529 | 32 | 1 |

```
1 #5.Display the last 20 records
2 exercise_df.iloc[-20:]
```

|     | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|-----|-------------|---------|---------------|---------------|---------|------|--------------------------|-----|---------|
| 748 | 3 | 187 | 70 | 22 | 200 | 36.4 | 0.408 | 36 | 1 |
| 749 | 6 | 162 | 62 | 0 | 0 | 24.3 | 0.178 | 50 | 1 |
| 750 | 4 | 136 | 70 | 0 | 0 | 31.2 | 1.182 | 22 | 1 |
| 751 | 1 | 121 | 78 | 39 | 74 | 39.0 | 0.261 | 28 | 0 |
| 752 | 3 | 108 | 62 | 24 | 0 | 26.0 | 0.223 | 25 | 0 |
| 753 | 0 | 181 | 88 | 44 | 510 | 43.3 | 0.222 | 26 | 1 |
| 754 | 8 | 154 | 78 | 32 | 0 | 32.4 | 0.443 | 45 | 1 |
| 755 | 1 | 128 | 88 | 39 | 110 | 36.5 | 1.057 | 37 | 1 |
| 756 | 7 | 137 | 90 | 41 | 0 | 32.0 | 0.391 | 39 | 0 |
| 757 | 0 | 123 | 72 | 0 | 0 | 36.3 | 0.258 | 52 | 1 |
| 758 | 1 | 106 | 76 | 0 | 0 | 37.5 | 0.197 | 26 | 0 |
| 759 | 6 | 190 | 92 | 0 | 0 | 35.5 | 0.278 | 66 | 1 |
| 760 | 2 | 88 | 58 | 26 | 16 | 28.4 | 0.766 | 22 | 0 |
| 761 | 9 | 170 | 74 | 31 | 0 | 44.0 | 0.403 | 43 | 1 |
| 762 | 9 | 89 | 62 | 0 | 0 | 22.5 | 0.142 | 33 | 0 |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 | 0 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 | 0 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 | 0 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 | 1 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 | 0 |

```
1 #6. Change the Outcome column to Diagnosis
2 exercise_df.rename(columns={"Outcome":"Diagnosis"})
```

|     | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Diagnosis |
|-----|-------------|---------|---------------|---------------|---------|------|--------------------------|-----|-----------|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 | 0 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 | 0 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 | 0 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 | 1 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 | 0 |

768 rows × 9 columns

```
1 #7. Create a new column Classification that display "Diabetes" if the value of outcome is 1 , otherwise "No Diabetes"
2 exercise_df['Classification'] = exercise_df['Outcome'].apply(lambda x: 'Diabetes' if x ==1 else 'No Diabetes')
3 exercise_df
```

|     | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|-----|-------------|---------|---------------|---------------|---------|------|--------------------------|-----|---------|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 | 0 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 | 0 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 | 0 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 | 1 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 | 0 |

768 rows × 10 columns

```
1 #8. Create a new dataframe "withDiabetes" that gathers data with diabetes
2 withDiabetes = exercise_df[exercise_df['Classification']=='Diabetes']
3 withDiabetes
```

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome | |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 | |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 | |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 | |
| 6 | 3 | 78 | 50 | 32 | 88 | 31.0 | 0.248 | 26 | 1 | |
| 8 | 2 | 197 | 70 | 45 | 543 | 30.5 | 0.158 | 53 | 1 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 755 | 1 | 128 | 88 | 39 | 110 | 36.5 | 1.057 | 37 | 1 | |
| 757 | 0 | 123 | 72 | 0 | 0 | 36.3 | 0.258 | 52 | 1 | |
| 759 | 6 | 190 | 92 | 0 | 0 | 35.5 | 0.278 | 66 | 1 | |
| 761 | 9 | 170 | 74 | 31 | 0 | 44.0 | 0.403 | 43 | 1 | |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 | 1 | |

268 rows × 10 columns

```
1 #9. Create a new dataframe "noDiabetes" thats gathers data with no diabetes
2 noDiabetes = exercise_df[exercise_df['Classification']=='No Diabetes']
3 noDiabetes
```

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 | |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 | |
| 5 | 5 | 116 | 74 | 0 | 0 | 25.6 | 0.201 | 30 | 0 | |
| 7 | 10 | 115 | 0 | 0 | 0 | 35.3 | 0.134 | 29 | 0 | |
| 10 | 4 | 110 | 92 | 0 | 0 | 37.6 | 0.191 | 30 | 0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 762 | 9 | 89 | 62 | 0 | 0 | 22.5 | 0.142 | 33 | 0 | |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 | 0 | |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 | 0 | |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 | 0 | |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 | 0 | |

500 rows × 10 columns

```
1 #10. Create a new dataframe "Pedia" that gathers data with age 0 to 19
2 Pedia = exercise_df[exercise_df['Age'] <= 19]
3 Pedia
```

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome | Clas |
|---|---|---|---|---|---|---|---|---|---|---|

```
1 #11. Create a new dataframe "Adult" that gathers data with age greater than 19
2 Adult = exercise_df[exercise_df['Age'] > 19]
3 Adult
```

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| | | | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |

Next steps: ◉ View recommended plots

```
 1 #12. Use numpy to get the average age and glucose value.
 2 #Extraction of ages  and glucose value in dataset
 3 ages= np.array(exercise_df['Age'])
 4 glucose_value =np.array(exercise_df['Glucose'])
 5
 6 #computing the average value of age and glucose using numpy
 7 average_age = np.average(ages)
 8 average_glucose =np.average(glucose_value)
 9
10 #Display the average value of age as a whole number and
11 #glucose average value into 2 decimal places
12 print(f"Average age:{average_age:.0f} years old")
13 print(f"Average Glucose value: {average_glucose:.2f}")
```

    Average age:33 years old
    Average Glucose value: 120.89

```
 1 #13.Use numpy to get the median age and glucose value
 2 # Using the extraction above we can get the median
 3
 4 #Compute the median of age and value using numpy
 5 median_age = np.median(ages)
 6 median_glucose = np.median(glucose_value)
 7
 8 #Display the median
 9 print(f"Median age:{median_age} years old")
10 print(f"Median Glucose value: {median_glucose:.2f}")
```

    Median age:29.0 years old
    Median Glucose value: 117.00

```
 1 #14.Use numpy to get the middle values of glucose and age
 2 #Using the extraction above we can get the middle value
 3 #middle value is the value of a given dataset that is not arranged in order.
 4 index = int(len(ages)/2)
 5 midage = ages[index]
 6 midglucose = glucose_value[index]
 7 print("Middle Age: " + str(midage) + " years old")
 8 print("Middle Glucose: " + str(midglucose))
 9
```

    Middle Age: 25 years old
    Middle Glucose: 125

```
 1 #15.Use numpy to get the standard deviation of the skinthickness.
 2 skinthickness = np.array(exercise_df['SkinThickness'])
 3
 4 #compute for the std
 5 std_skinthickness = np.std(skinthickness)
 6
 7 #Display the STD with 6 decimal places
 8 print(f"Standard Deviation of SkinThickness: {std_skinthickness:.6f}")
```

    Standard Deviation of SkinThickness: 15.941829

## 6.4 Conclusion

In this activity, I learn how to make a import a csv file and manipulate data using numpy and pandas. I also implement some data analysis functions without using statistics module. The things I loved while doing the activity is the preparation,manipulation and analyzing the data. Learning numpy and pandas is a valuable skill for anyone who wants to work with data analysis, machine learning, or scientific computing.