# Task 21

Data Manipulation with Pandas and NumPy

Name: Basel Amr Barakat
Email: baselamr52@gmail.com

# Requirement One

Creating and Manipulating Pandas DataFrames

# 1.1 Create Dataframe from lists
# 1.2 Create Dataframe from Dictionary

# 1.3 Create Dataframe from NumPy array
# 1.4 Data Selection and Filtering

## 3.3 Create DataFrame from NumPy Array

```python
1 # Creating a DataFrame from a NumPy array
2 array = np.random.rand(5, 3)                    # ---> Build a random numpy array 5x3 (5 rows and 3 columns)
3 df_array = pd.DataFrame(array, columns=['Column1', 'Column2', 'Column3'])
4 print(df_array)
```

```
    Column1   Column2   Column3
0  0.956032  0.669054  0.638814
1  0.190019  0.122922  0.665823
2  0.916984  0.022632  0.550137
3  0.016407  0.101865  0.612821
4  0.996306  0.224719  0.618217
```

## 3.4 Data Selection and Filtering

```python
[3]   1 # Selecting a single column
      2 print("Filtering the names only of the dataframe")
      3 print(df_dict['Name'])
      4 # Filtering using loc methode
      5 Names = df_dict.loc[:,"Name"]
      6 print(Names)
      7 # Filtering Using iloc methode
      8 Names = df_dict.iloc[:,0]
      9 print(Names)
     10 # Filtering using condition
     11
     12 # Filtering rows based on conditions
     13 filtered_df = df_dict[df_dict['Age'] > 25]
     14 print("Filtering the rows of age >25")
     15 print(filtered_df)
```

```
Filtering the names only of the dataframe
0       Basel
1        Omar
2     Mohamed
3     Abanoub
4         Aya
Name: Name, dtype: object
0       Basel
1        Omar
2     Mohamed
3     Abanoub
4         Aya
Name: Name, dtype: object
0       Basel
1        Omar
2     Mohamed
3     Abanoub
4         Aya
Name: Name, dtype: object
Filtering the rows of age >25
      Name  Age    City
1     Omar   27  Assuit
3  Abanoub   26   Cairo
```

# 2.1 Handling Missing Data

```python
1  # Create a DataFrame with missing values
2  data = {
3      'A': [1, 2, np.nan, 4,5,2,1,2,4,10],
4      'B': [5, np.nan, 7, 8,np.nan,1,3,2,2,1],
5      'C': [9, 10, 11, np.nan,np.nan,2,1,3,5,np.nan],
6      'D':[1,2,3,4,5,6,7,8,9,10],
7      'E':[11,22,31,41,51,61,71,81,91,101],
8      'F':([1,2,3,4,5,6,7,8,9,10]),
9      'G':[11,22,np.nan,41,51,np.nan,71,81,91,101]
10 }
11
12 df_missing = pd.DataFrame(data)
13 print("Original DataFrame with missing values:")
14 print(df_missing)
15
16 # Fill missing values
17 df_filled = df_missing.fillna(0,inplace=False)
18 print("After filling missing values:")
19 print(df_filled)
20
21 # Drop missing values
22 df_dropped = df_missing.dropna(axis=0,inplace=False)
23 print("After dropping rows with missing values:")
24 print(df_dropped)
25
26 df_dropped = df_missing.dropna(axis=1,inplace=False)
27 print("After dropping columns with missing values:")
28 print(df_dropped)
```

```
Original DataFrame with missing values:
      A    B     C   D    E   F      G
0   1.0  5.0   9.0   1   11   1   11.0
1   2.0  NaN  10.0   2   22   2   22.0
2   NaN  7.0  11.0   3   31   3    NaN
3   4.0  8.0   NaN   4   41   4   41.0
4   5.0  NaN   NaN   5   51   5   51.0
5   2.0  1.0   2.0   6   61   6    NaN
6   1.0  3.0   1.0   7   71   7   71.0
7   2.0  2.0   3.0   8   81   8   81.0
8   4.0  2.0   5.0   9   91   9   91.0
9  10.0  1.0   NaN  10  101  10  101.0
After filling missing values:
      A    B     C   D    E   F      G
0   1.0  5.0   9.0   1   11   1   11.0
1   2.0  0.0  10.0   2   22   2   22.0
2   0.0  7.0  11.0   3   31   3    0.0
3   4.0  8.0   0.0   4   41   4   41.0
4   5.0  0.0   0.0   5   51   5   51.0
5   2.0  1.0   2.0   6   61   6    0.0
6   1.0  3.0   1.0   7   71   7   71.0
7   2.0  2.0   3.0   8   81   8   81.0
8   4.0  2.0   5.0   9   91   9   91.0
9  10.0  1.0   0.0  10  101  10  101.0
After dropping rows with missing values:
     A    B    C  D   E  F      G
0  1.0  5.0  9.0  1  11  1   11.0
6  1.0  3.0  1.0  7  71  7   71.0
7  2.0  2.0  3.0  8  81  8   81.0
8  4.0  2.0  5.0  9  91  9   91.0
After dropping columns with missing values:
    D    E   F
0   1   11   1
1   2   22   2
2   3   31   3
3   4   41   4
4   5   51   5
5   6   61   6
6   7   71   7
7   8   81   8
8   9   91   9
9  10  101  10
```

# 2.2 Merging Dataframes

```
[27]    1 # Create two DataFrames for merging
        2 data1 = {'ID': [1, 2, 3], 'Name': ['Basel', 'Amr', 'Barakat']}
        3 data2 = {'ID': [1, 2, 4], 'Score': [85, 90, 95]}
        4 df1 = pd.DataFrame(data1)
        5 df2 = pd.DataFrame(data2)
        6 print("Before Merging the two dataframes")
        7 print(df1)
        8 print(df2)
        9 # Merging on ID
       10 merged_df = pd.merge(df1, df2, on='ID', how='inner')
       11 print("Merged DataFrame:")
       12 print(merged_df)
       13 print("We found that the two dataframes are merged together depending on col `ID` and that the ID 3 is not in the second dataframe")
       14 print("We found that the two dataframes are merged together depending on col `ID` and that the ID 4 is not in the first dataframe")
```

```
Before Merging the two dataframes
   ID     Name
0   1    Basel
1   2      Amr
2   3  Barakat
   ID  Score
0   1     85
1   2     90
2   4     95
Merged DataFrame:
   ID   Name  Score
0   1  Basel     85
1   2    Amr     90
We found that the two dataframes are merged together depending on col `ID` and that the ID 3 is not in the second dataframe
We found that the two dataframes are merged together depending on col `ID` and that the ID 4 is not in the first dataframe
```

# Requirement Three

Combining Pandas and NumPy for Analysis

# 3.1 Aggregations and statistical Analysis

## 5.1 Aggregations and Statistical Analysis

```python
1  # Creating a Sample DataFrame for Analysis
2  data_analysis = {
3      'Product': ['A', 'B', 'C', 'D'],
4      'Sales_Price': [1000, 2500, 30000, 4000],
5      'Cost': [50, 80, 120, 150]
6  }
7  df_analysis = pd.DataFrame(data_analysis)
8  print("\nAnalysis DataFrame:\n", df_analysis)
9
10 # Using NumPy for Mathematical Operations
11 df_analysis['Profit'] = np.array(df_analysis['Sales_Price']) - np.array(df_analysis['Cost'])
12 print("\nDataFrame with Profit Calculated:\n", df_analysis)
13
14 # Aggregation and Statistical Analysis
15 mean_sales = np.mean(df_analysis['Sales_Price'])
16 print("\nMean Sales:", mean_sales)
17
18 total_cost = np.sum(df_analysis['Cost'])
19 print("Total Cost:", total_cost)
20
21 max_profit = np.max(df_analysis['Profit'])
22 print("Maximum Profit:", max_profit)
23
24 min_profit = np.min(df_analysis['Profit'])
25 print("Minimum Profit:", min_profit)
26
27 median_profit = np.median(df_analysis['Profit'])
28 print("Median Profit:", median_profit)
29
30 std_profit = np.std(df_analysis['Profit'])
31 print("Standard Deviation of Profit:", std_profit)
32
33 df_analysis.describe()
34 print("We have noticed that the std equation of NumPy is different from the equation of describe() ")
```

```
Analysis DataFrame:
   Product  Sales  Cost
0        A   1000    50
1        B   2500    80
2        C  30000   120
3        D   4000   150

DataFrame with Profit Calculated:
   Product  Sales  Cost  Profit
0        A   1000    50     950
1        B   2500    80    2420
2        C  30000   120   29880
3        D   4000   150    3850

Mean Sales: 9375.0
Total Cost: 400
Maximum Profit: 29880
Minimum Profit: 950
Median Profit: 3135.0
Standard Deviation of Profit: 11940.407237611287
```

|        | Sales        | Cost        | Profit       |
|--------|--------------|-------------|--------------|
| count  | 4.000000     | 4.000000    | 4.000000     |
| mean   | 9375.000000  | 100.000000  | 9275.000000  |
| std    | 13804.437692 | 43.969687   | 13787.594666 |
| min    | 1000.000000  | 50.000000   | 950.000000   |
| 25%    | 2125.000000  | 72.500000   | 2052.500000  |
| 50%    | 3250.000000  | 100.000000  | 3135.000000  |
| 75%    | 10500.000000 | 127.500000  | 10357.500000 |
| max    | 30000.000000 | 150.000000  | 29880.000000 |

# 3.2 Data Transformation

```python
# Apply NumPy transformations
data_analysis = {
    'Sales': [1000, 2500, 30000, 4000],
    'Cost': [50, 80, 120, 150]
}
df_analysis = pd.DataFrame(data_analysis)
df_transformed_log = df_analysis.apply(np.log)
print("Log Transformed DataFrame:")
print(df_transformed_log)

df_transformed_sqrt = df_analysis.apply(np.sqrt)
print("Square Root Transformed DataFrame:")
print(df_transformed_sqrt)

df_tansformed_exp = df_analysis.apply(np.exp)
print("Exponential Transformed DataFrame:")
print(df_tansformed_exp)

df_transformed_sin = df_analysis.apply(np.sin)
print("Sin Transformed DataFrame:")
print(df_transformed_sin)

df_transformed_cos = df_analysis.apply(np.cos)
print("Cos Transformed DataFrame:")
print(df_transformed_cos)

df_transformed_tan = df_analysis.apply(np.tan)
print("Tan Transformed DataFrame:")
print(df_transformed_tan)
```

```
Log Transformed DataFrame:
        Sales      Cost
0    6.907755  3.912023
1    7.824046  4.382027
2   10.308953  4.787492
3    8.294050  5.010635
Square Root Transformed DataFrame:
        Sales       Cost
0    31.622777   7.071068
1    50.000000   8.944272
2   173.205081  10.954451
3    63.245553  12.247449
Exponential Transformed DataFrame:
    Sales          Cost
0     inf  5.184706e+21
1     inf  5.540622e+34
2     inf  1.304181e+52
3     inf  1.393710e+65
Sin Transformed DataFrame:
        Sales      Cost
0    0.826880 -0.262375
1   -0.650128 -0.993889
2   -0.802665  0.580611
3   -0.683504 -0.714876
Cos Transformed DataFrame:
        Sales      Cost
0    0.562379  0.964966
1    0.759825 -0.110387
2   -0.596430  0.814181
3   -0.729947  0.699251
Tan Transformed DataFrame:
        Sales      Cost
0    1.470324 -0.271901
1   -0.855628  9.003655
2    1.345784  0.713123
3    0.936375 -1.022346
```

# 3.1 Aggregations and statistical Analysis

```python
[52]    1 # Plot data
        2 import matplotlib.pyplot as plt
        3 data = {'A': [1, 2, 3, 4], 'B': [5, 6, 7, 8]}
        4 df_analysis = pd.DataFrame(data)
        5
        6 plt.subplot(2, 2, 1)
        7 plt.plot(df_analysis['A'], df_analysis['B'])
        8 plt.title('Line Plot')
        9
       10 plt.subplot(2, 2, 2)
       11 plt.scatter(df_analysis['A'], df_analysis['B'])
       12 plt.title('Scatter Plot')
       13
       14 plt.subplot(2, 2, 3)
       15 plt.bar(df_analysis['A'], df_analysis['B'])
       16 plt.title('Bar Plot')
       17
       18 plt.subplot(2, 2, 4)
       19 plt.hist(df_analysis['A'])
       20 plt.title('Histogram')
       21
       22 plt.tight_layout()
       23 plt.show()
```