4/15/2023

# CS-351-Artificial Intelligence

2020244 - Mohsin Zia
2020474 - Syed Irtaza Haider

INSTRUCTOR: PROFESSOR ZAHID HALIM

## Clustering Techniques!

*HOW TO RUN?*

1. Open the folder in vsCode.
2. Go to main.py.
3. Click on run python file.

*EXPLAINATION:*

➢ **UTILITY LIBRARIES USED:**

```python
import numpy as np
import matplotlib.pyplot as plt
from numpy import ndarray
import pandas as pd
import math
import networkx as nx
import os
import shutil
```

➢ **TASK 01:**

✓ **FUNCTIONS FOR READ/WRITE FROM/TO A TEXT FILE:**
- Function read_data reads data from a text file and returns a tuple containing (numberOfRows, numberOfColumns, Data).
- Function write_data writes data from matrix to a text file including (numberOfRows, numberOfColumns, Data).

```python
def read_data(file_name: str, seperator=None) -> tuple: ...


def write_data(file_name: str, arrayType: str, data: list) -> ndarray
```

✓ **FUNCTIONS FOR COMPUTING PEARSON'S CORRELATION COEFFICIENT:**

- Function correlationCoefficient computes and returns the correlation value between two rows using Pearson's correlation coefficient formula.
- Function pearsonCorrelationCoefficientMatrix computes and returns the correlation matrix of NxN dimension.

```python
def correlationCoefficient(dataX: list, dataY: list) → float: ⋯



def pearsonCorrelationCoefficientMatrix(rows: int, data: ndarray) → ndarray: ⋯
```

✓ **FUNCTION FOR COMPUTING DISCRETIZED MATRIX:**

- Function discretize computes and returns the discretize matrix from the correlation matrix data. If the value in the correlation matrix is ≥ mean of it's respective column, then it's replaced with 0 else it is kept 1. Here doing vice versa of what we should've made us closer to sample output.

```python
def discretize(totalRows: int, corr_matrix: list) → ndarray: ⋯
```

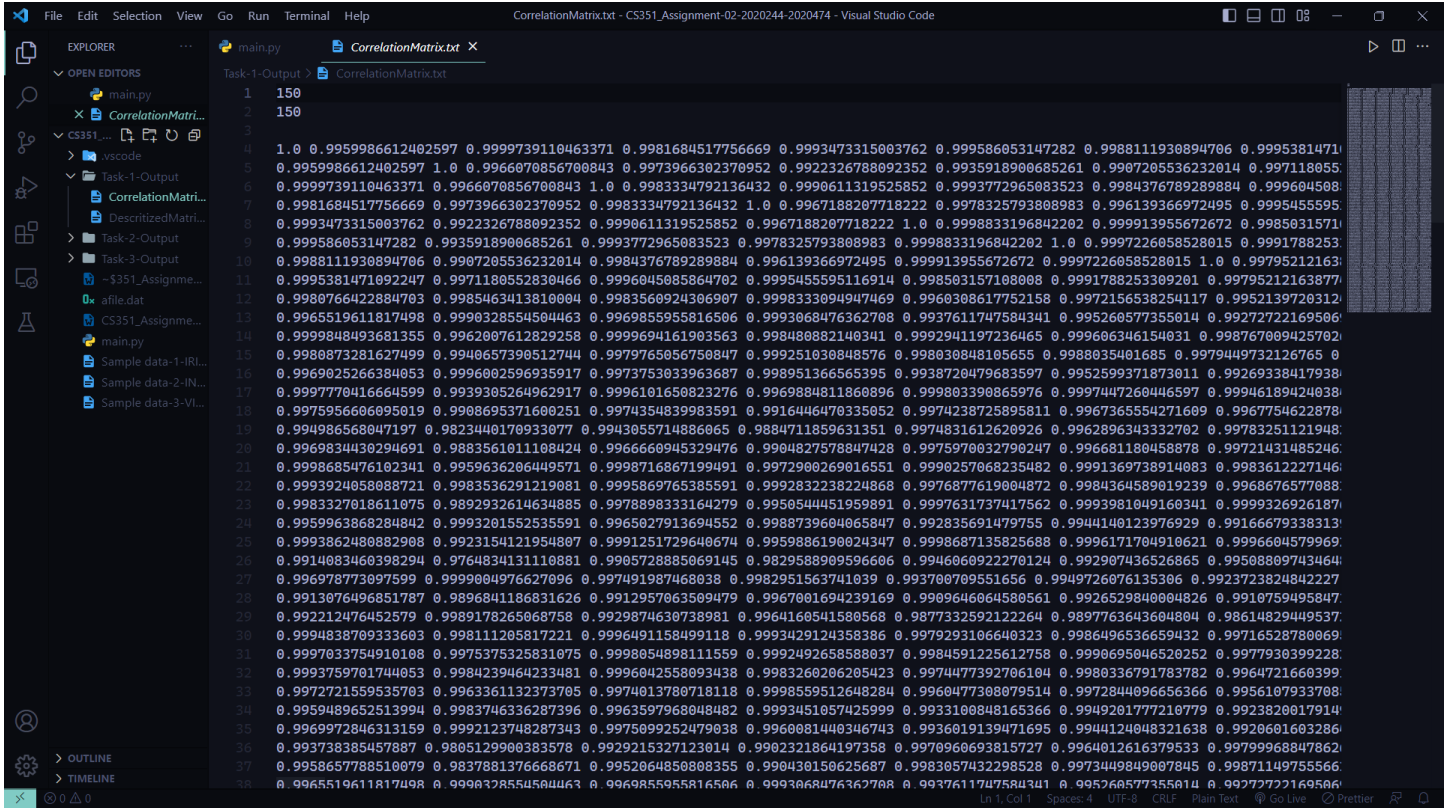✓ **FUNCTION FOR COMPUTING PIXEL MATRIX FOR COLOR CODED IMAGE:**

- Function colourCodedImage computes and returns the Pixel Matrix of the correlation matrix data. It is calculated by dividing each element by maximum value of their respective column and then multiplying each value with 255. Then it's placed in 3-D Matrix [R, G, B]. We compute the value of G here as we are required to print sample image in green shade and then place it like [0    G    0].

```python
def colourCodedImage(totalRows: int, corr_matrix: list) → ndarray: ⋯
```
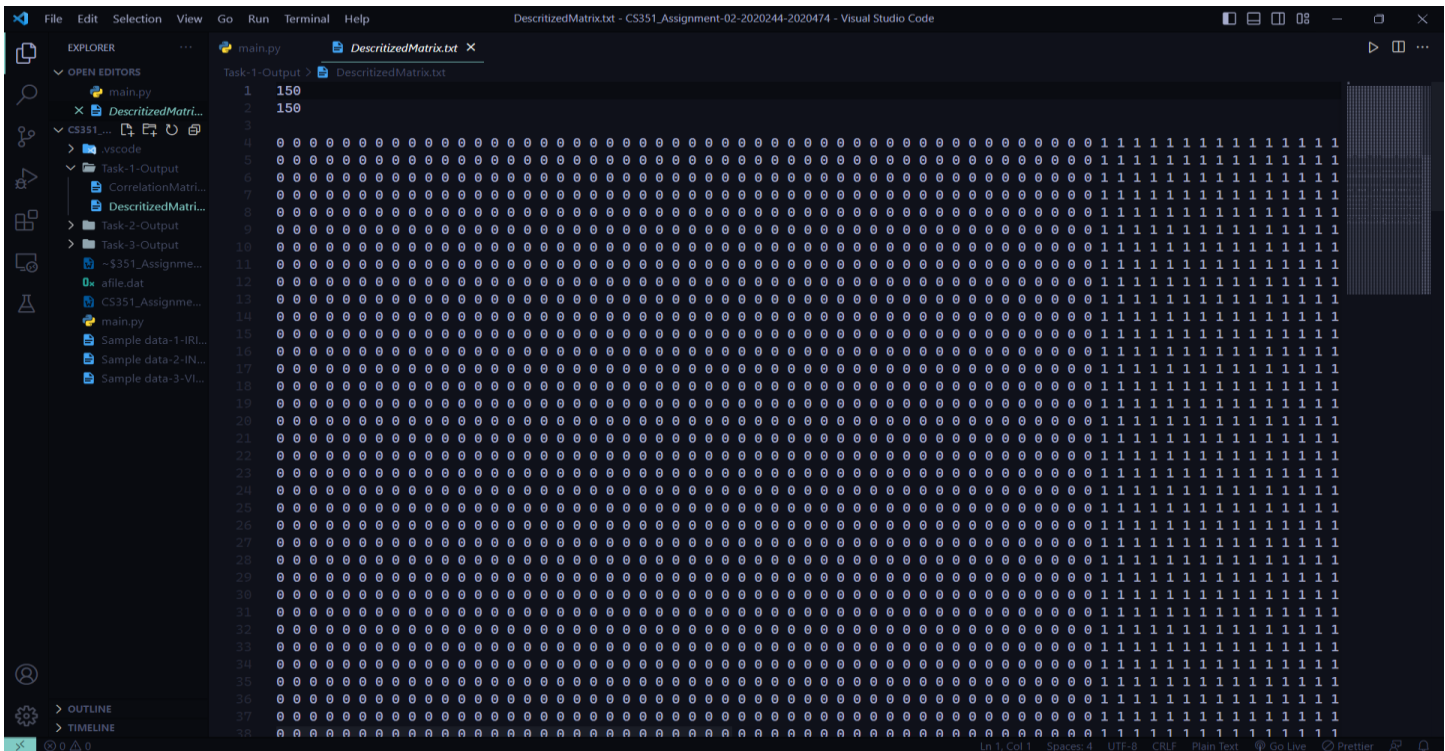
✓ **USING MATPLOTLIB FOR VISUALIZATION:**

```python
fig, axes = plt.subplots(nrows=1, ncols=2)
axes[0].set_title('Discretized Matrix')
axes[0].imshow(discretizedMatrix, cmap='gray')
axes[1].set_title('Color Coded Correlation Matrix')
axes[1].imshow(pixelMatrix.astype(int))
plt.show()
```
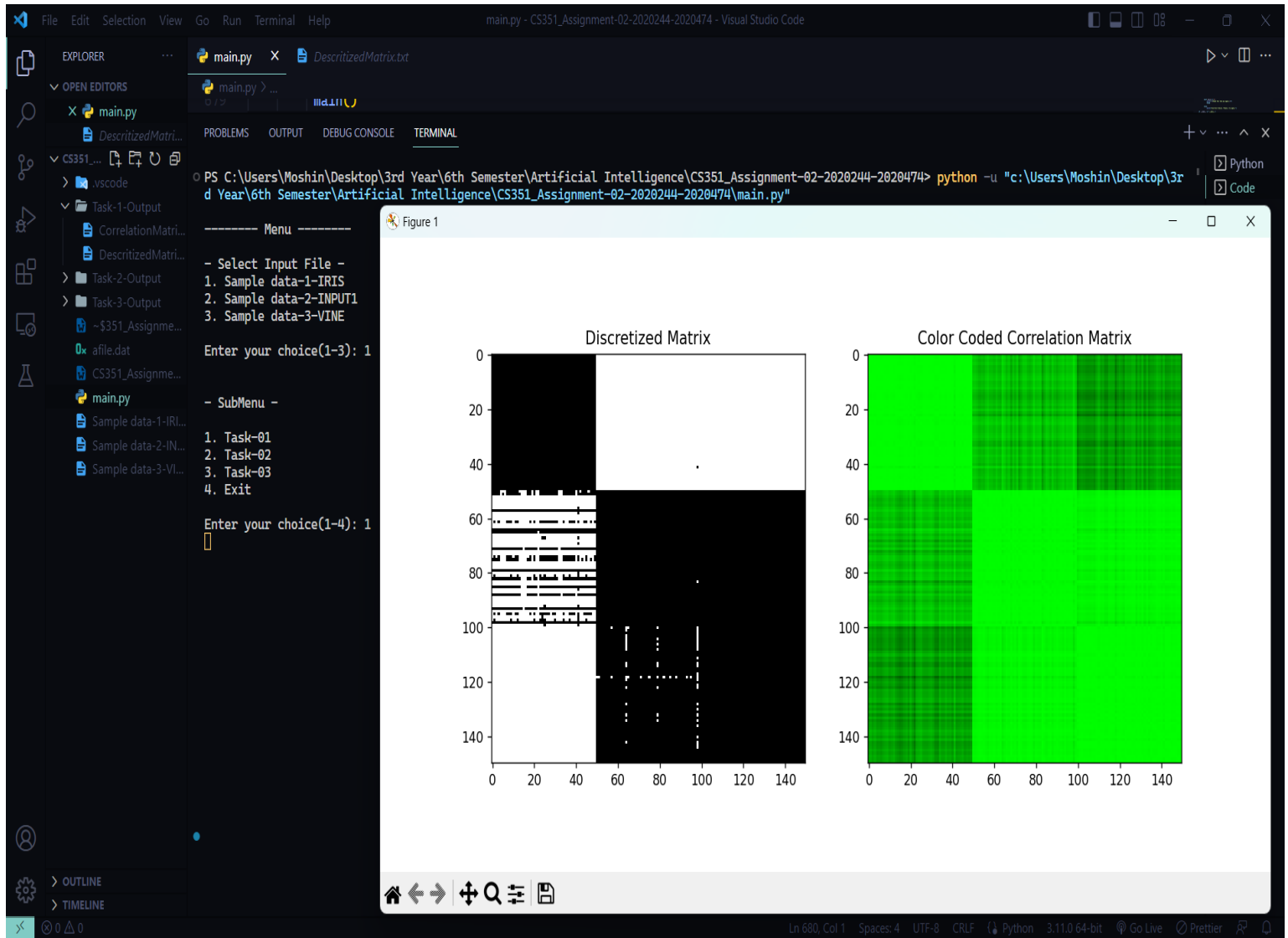
✓ **OUTPUT SCREENSHOTS:**

- *Correlation Matrix:*



- *Discretized Matrix:*

- *Bitmap Visualization:*

➢ **TASK 02:**

✓ **FUNCTION FOR PERMUTATION/SHUFFLING:**
- Function permutateMatrix provides constant shuffle over a period. It shuffles the rows of inputData. A random seed (50) is used to ensure that results are reproducible.

```
> def permutateMatrix(data: list) → list: …
```

✓ **FUNCTION FOR CALCULATING SIGNATURE VALUE AND SORTING VIA SIGNATURE:**
- Function signatureCalculation computes and returns sorted input data via their signature values. We passed the permutated data as input to this function, but I guess it really doesn't matter here if we pass the non-permutated data, in the end it's getting sorted via it's signature values.

```
> def signatureCalculation(number_of_rows: int, number_of_cols: int, shuffledData: list) → list: …
```

✓ **USING MATPLOTLIB FOR VISUALIZATION:**

```python
# Plotting Discritized Matrix before and after permutation.
fig, axes = plt.subplots(nrows=1, ncols=2)
axes[0].set_title('Before Permutation')
axes[0].imshow(discretizedMatrixBeforePermutation, cmap='gray')
axes[1].set_title('After Permutation')
axes[1].imshow(discretizedMatrixAfterPermutation, cmap='gray')
plt.show()
```

```python
# Plotting Pixel Matrix after SignatureWiseSorting And Rearrangement.
plt.suptitle('After Signature-Wise-Sorting')
plt.imshow(pixelMatrix.astype(int))
plt.show()
```

✓ **OUTPUT SCREENSHOTS:**

- *Input after permutation:*



- *Input after sorting via signature:*

- *Discretized Matrix after permutation:*



- *Correlation Matrix after sorting via signature.*

- *Bitmap Visualization:*

> ➢ **TASK 03:**

> ✓ **FUNCTION FOR SETUP OF PERMUTATED CORRELATION MATRIX:**
>   - Function setPermutatedMatrix takes permutated data set's correlation matrix and based on a threshold, it updates that matrix. It iterates through each value in the matrix and if the value is less than the threshold it makes it 0, else it is **not** updated.

```
> def setPermutatedMatrix(number_of_rows: int, number_of_cols: int, shuffledData: list) → list: …
```

> ✓ **FUNCTIONS FOR COMPUTING NODE WEIGHTS OF THE UPDATED MATRIX WITH RESPECT TO THRESHOLD:**
>   - Function getNodeWeights computes and returns the nodeWeights matrix by calculating node weights of each row. It is calculated as the sum of all columns in a respective row.

```
> def getNodeWeights(number_of_rows: int, number_of_cols: int, weightedMatrix: list) → list: …
```

> ✓ **FUNCTION FOR MAKING CLUSTER OF THE ROW IN CORRELATION MATRIX WITH MAXIMUM NODE WEIGHTS:**
>   - Function getIndexOfMaxWeight computes and returns the index of the node with maximum weightage. Once it is visited, matrix of node weights at that particular index is placed with -1 so that that node is excluded from our to-cluster list.
>   - Function makeCluster makes the graph of element with maximum node weight, it connects the edges to the elements in graphList that have survived the threshold test and then it connects to them and forms a cluster. Adding nodes and edges in between is done in this function.

```
> def getIndexOfMaxWeight(number_of_rows: int, nodeWeights: list) → int: …



> def makeClusters(index: int, graphList: list): …
```

> ✓ **FUNCTION FOR VISUALIZATION OF WEIGHTED GRAPH IN FORM OF CLUSTERS:**
>   - Function visualizeWeightedGraph takes input a weighted graph and index with maximum weight, it then used matplotlib to create a figure. Further it uses networkx library to draw the nodes and edges between source and destination. It can also add weights in between edges but it doesn't look good when we are visualizing it. It saves the figure in png format.

```
> def visualizeWeightedGraph(G, index: int): …
```

✓ **CREATING A DIRECTORY FOR OUR CLUSTER IMAGES:**

```
605     # removing if the directory is already present
606     shutil.rmtree('./Task-3-Output/GraphClusterVisualization',
607                   ignore_errors=True)
608
609     # creating a directory for saving png file of clusters
610     os.mkdir('./Task-3-Output/GraphClusterVisualization')
```

✓ **OUTPUT SCREENSHOTS:**

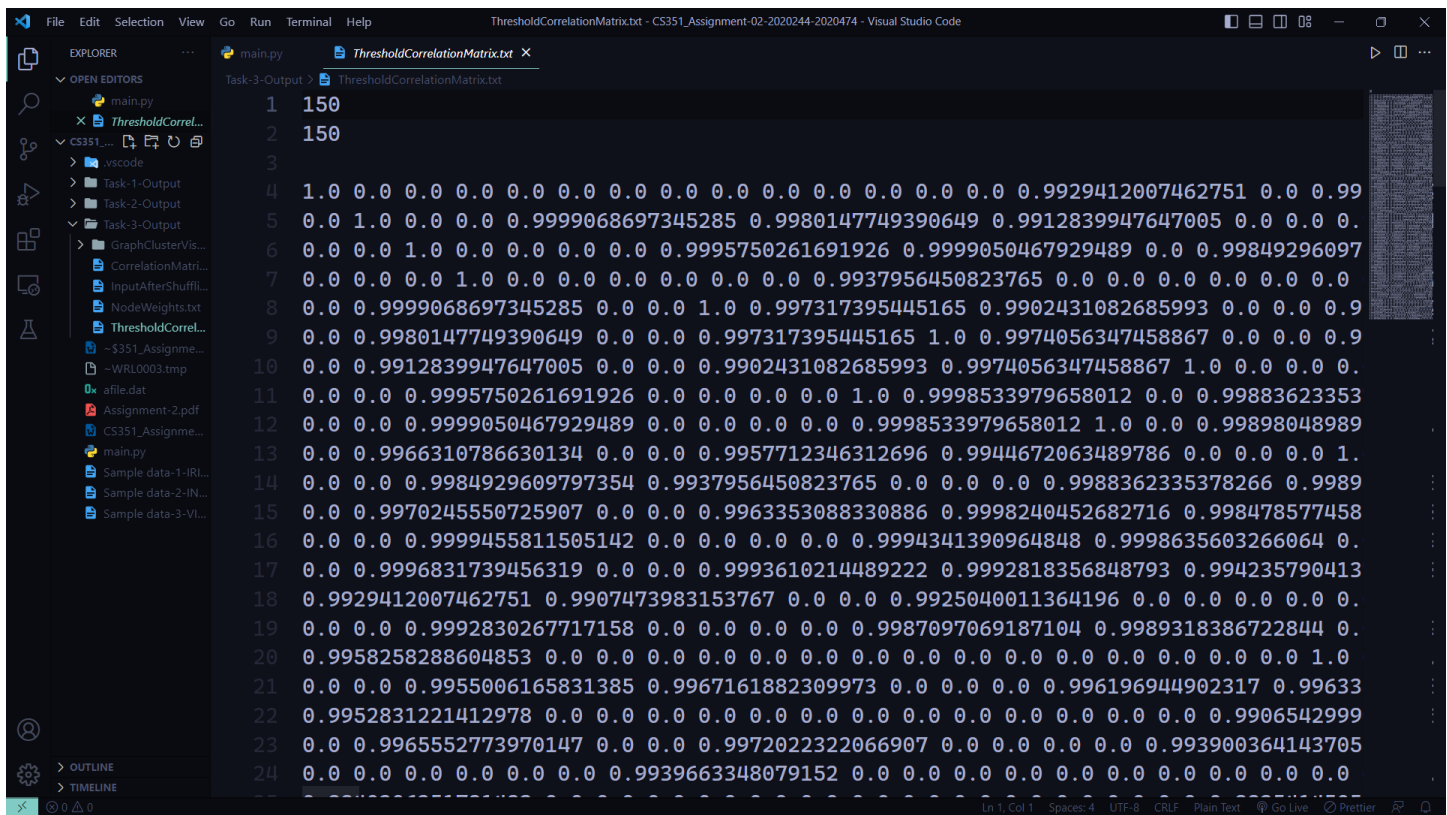- *Node Weights Matrix:*
    - The threshold is set to 0.99.

- *Threshold Correlation Matrix:*



- *Directory of png's overview/Clusters saved:*

- *Most dense and least dense clusters formed:*



o *Most dense cluster:*



o *Least dense cluster:*

> ➢ **COMPARISON OF TASK 2 AND 3:**

>> ✓ **TASK-2:**
>> - The primary goal of task 2 is to demonstrate the dissimilarity in bitmap by carrying out permutations on the data and presenting the bitmap of the discretized iris data before and after the permutation process. Subsequently, the permutated data undergoes sorting based on its signature values, and a color-coded image of the correlation data is generated and displayed.

>> ✓ **TASK-3:**
>> - The primary objective of task 3 is to identify clusters within the correlation matrix subsequent to the permutation of data. To accomplish this, a threshold value is established, which in our case is 0.99, and the correlation matrix is updated accordingly. Following this, we pinpoint the node with the highest weight and construct a cluster including its surviving neighbors from the threshold process. This process is repeated iteratively until clusters are generated for each row in our graph list. Ultimately, this technique enables us to effectively determine the clusters present within the correlation matrix.

> ➢ **WORK DISTRIBUTION AMONG GROUP MEMBERS:**

Throughout the completion of this assignment, it is evident that we have made substantial contributions towards its successful completion. From the initiation of the project to the final submission, we both have actively participated in the task allocation, research, and implementation phases. Additionally, the collaborative nature of the project enabled both of us to offer valuable insights and suggestions, leading to the development of a comprehensive and well-structured final submission. We both have played an integral role throughout the assignment, and combined efforts have resulted in a successful outcome.

> ➢ **INSIGHTS:**

During the completion of this assignment, we have gained several valuable insights. Firstly, we discovered that the Pearson correlation matrix, although correlation is always symmetric process, didn't produce symmetric bitmap for discretized data. While initially considering that we may have made an error in our formula, using the predefined pd.DataFrame.corr() function produced results that were consistent with our output. Despite this, doubts remained about the accuracy of the sample image produced by us. It should also be noted that discretized data is calculated by taking the mean of the column and updating each value in that column. I think that it's the reason why the bitmap is not looking symmetric. Another observation made was that sorting via signature sorts both the input data and shuffled data to the same matrix, rendering the shuffling process is unnecessary here. Finally, we recognized the immense power of the Python language. Tasks that would have taken a significant amount of time in other programming languages such as C++, such as bitmap visualization (I first tried visualization using QT after completing these tasks in C++ and it was becoming too complex 😠 ), were easily accomplished using Matplotlib in just two lines of code in Python.

*THE END .*