



Stony Brook University

EEO388

FALL 2024

# Assignment 1

*Pete Mills*  
*ID: 115009163*

Professor  
Vibha MANE

September 22, 2024

## Overview

In this assignment we made use of scatter plots and histograms to analyze distributions and relationships, separability, and irregularities in data sets.

### 1 Exercise Set 1A

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from datetime import datetime
6
7 # Create a NumPy array with 1 column and 100 rows, filled with random numbers
8 array_1d = np.random.randn(100)
9
10
11 # Convert the 1D array into a pandas Series
12 series_1d = pd.Series(array_1d)
13 print(series_1d)
14
15
16 # Create a NumPy array with 3 columns and 100 rows, filled with random integers
17 array_3d = np.random.randn(100, 3)
18
19
20 # Convert the 3-column array into a DataFrame with column labels
21 df = pd.DataFrame(array_3d, columns=['X1', 'X2', 'X3'])
22 print(df)
23
24
25 # Use Seaborn to create a pairplot of the DataFrame
26 sns.pairplot(df)
27 #plt.show()
28
29 # Generate current timestamp in the format yyyy-mm-dd-hh-mm-ss
30 timestamp = datetime.now().strftime("%Y-%m-%d-%H-%M-%S")
31
32 # Create the file name with timestamp prepended
33 filename = f"{timestamp}_pairplot.png"
34 plt.savefig(filename)
35
36
```

Figure 1: Source code.

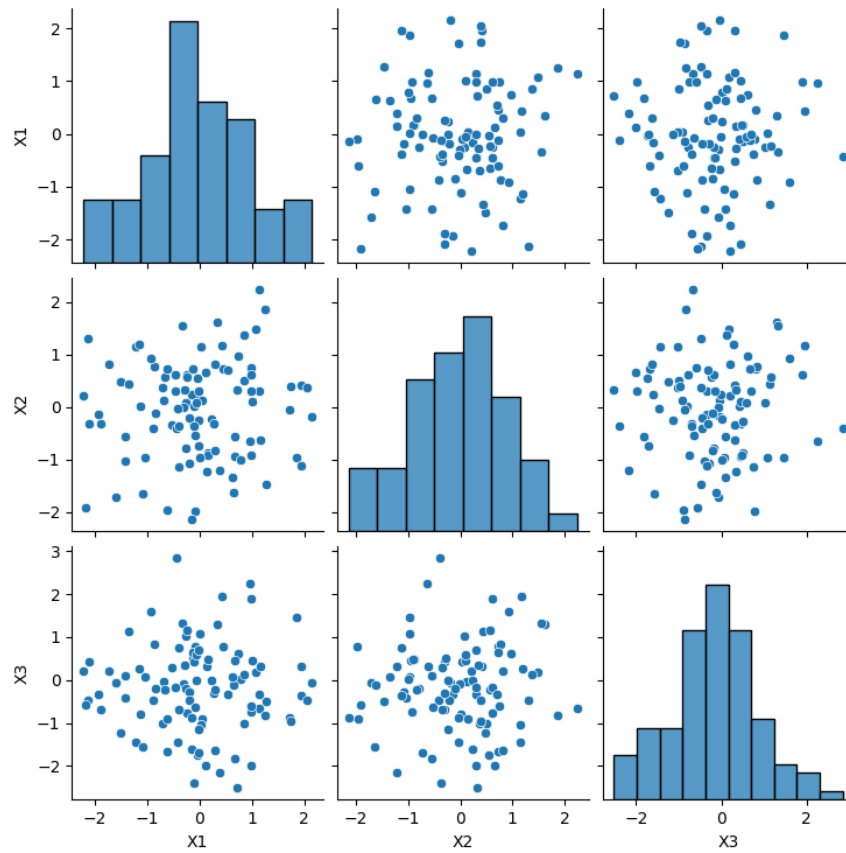


Figure 2: The Pair plot set for these random data.

## 2 Exercise Set 1B

### 2.1 Raisin

The data in this dataset are categorical. Based on the plots there do appear to be some outliers. There are 450 samples in the Kecimen class, and 450 samples in the Bensı class.

By setting hue based on class, we can visualize the two classes of raisins on the same scatter plot. We can then inspect for clustering of data which reveals there is some separability in the data. We can use this to identify characteristics for class sorting.

```

1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Load the data
6 df_raisin = pd.read_csv('../data-sets/Raisin_Dataset.csv')
7
8 # Scatter plot matrix
9 sns.pairplot(df_raisin, hue='Class', markers=["o", "s"])
10 plt.show()
11
12 # Histograms
13 df_raisin.hist(bins=30, figsize=(10, 7))
14 plt.tight_layout()
15 plt.show()
16
17 # Check for class distribution
18 print(df_raisin['Class'].value_counts())
19

```

Figure 3: Source code.

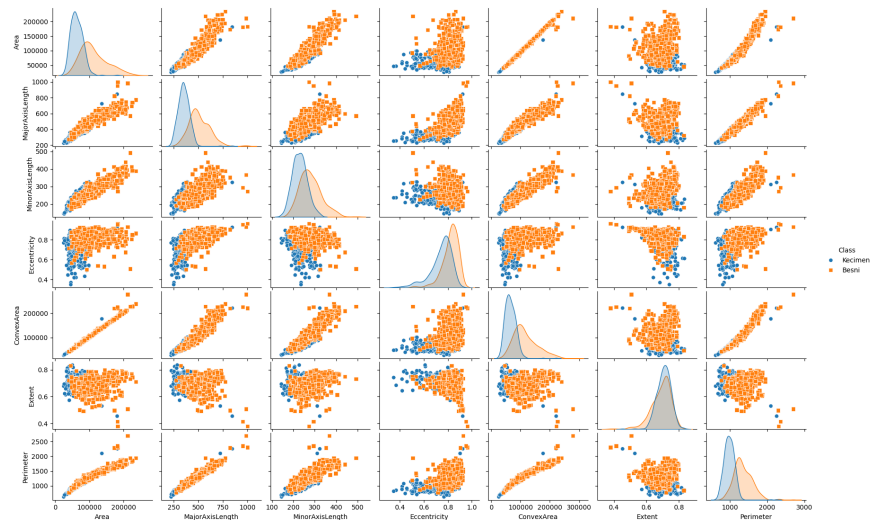


Figure 4: Scatter plots.

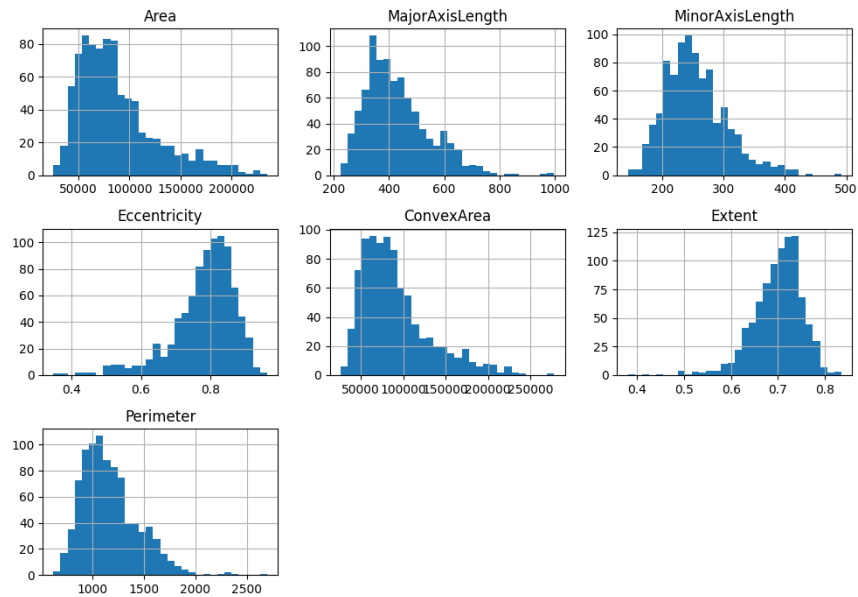


Figure 5: Histograms.

## 2.2 Deep Space

The data in this dataset are numerical. There does not appear to be any outliers.

```

1  import pandas as pd
2  import seaborn as sns
3  import matplotlib.pyplot as plt
4
5  # Load the data
6  df_deepspace = pd.read_csv('../data-sets/DeepSpaceData.csv')
7
8  # Scatter plot matrix without class labels
9  sns.pairplot(df_deepspace)
10 plt.show()
11
12
13 # Histograms
14 df_deepspace.hist(bins=30, figsize=(10, 7))
15 plt.tight_layout()
16 plt.show()

```

Figure 6: Source code.

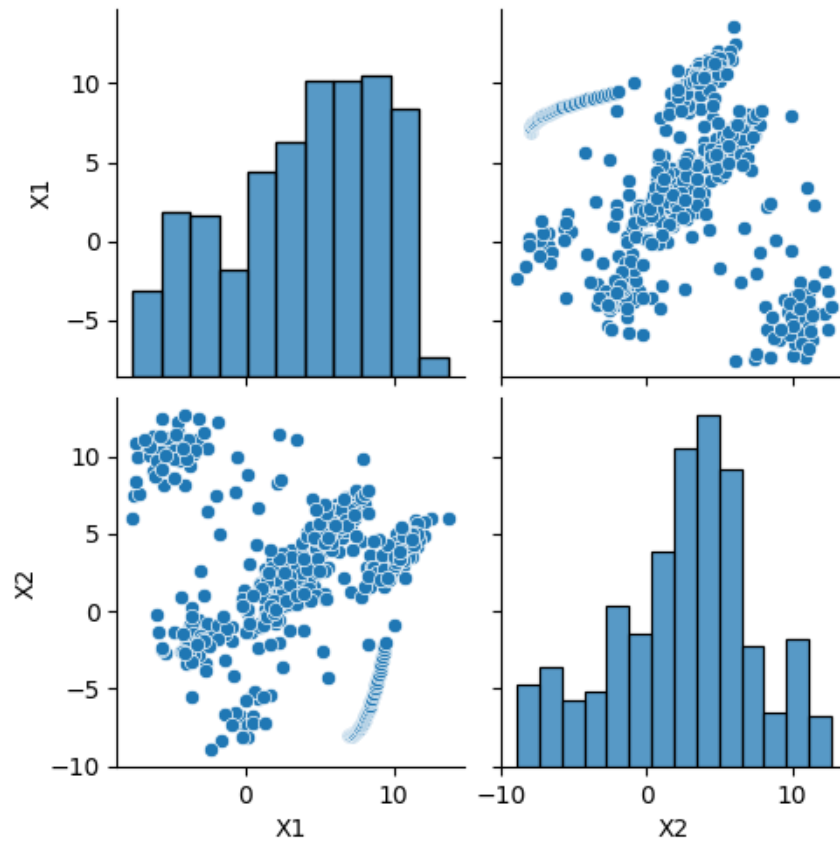


Figure 7: Scatter Plots.

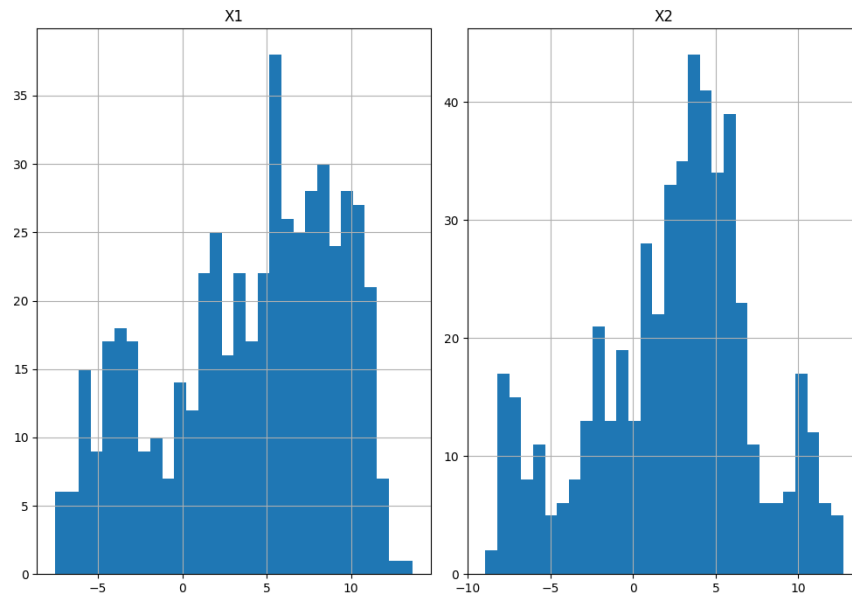


Figure 8: Histograms.