

CSBB 311 : Machine Learning

LAB ASSIGNMENT 1: Analysis Of Data Using Python Library

Submitted By:

Name: KARTIK MITTAL

Roll No: 221210056

Branch: CSE

Semester: 5th

Group: 2

Submitted To: Dr. Preeti Mehta

Department of Computer Science and Engineering

NATIONAL INSTITUTE OF TECHNOLOGY
DELHI



2024

Q1 : - Loading of data file

Code :-

```
1 import pandas as pd
2
3 # Load the Titanic dataset from seaborn's built-in dataset
4 data = pd.read_csv('tested.csv')
5
6 # Display the first few rows of the dataset
7 print("Data Loaded:")
8 print(data.head())
```

Result :-

```
PS C:\Users\HP\Desktop\college\semester 5\Machine Learning> python -u "c:\Users\HP\Desktop\college\semester 5\Machine Learning\loadData.py"
Q
1      893      1      3      Wilkes, Mrs. James (Ellen Needs) ... 363272  7.0000  NaN  S
2      894      0      2      Myles, Mr. Thomas Francis ... 240276  9.6875  NaN  Q
3      895      0      3      Wirz, Mr. Albert ... 315154  8.6625  NaN  S
4      896      1      3  Hirvonen, Mrs. Alexander (Helga E Lindqvist) ... 3101298 12.2875  NaN  S

[5 rows x 12 columns]
PS C:\Users\HP\Desktop\college\semester 5\Machine Learning>
```

Q2 : - Size of features and names of features

Code :-

```
1 import loadData
2 from loadData import data
3
4 # Size of the dataset
5 data_shape = data.shape
6 print(f"Size of the dataset: {data_shape}")
7
8 # Names of the features
9 features = data.columns.tolist()
10 print(f"Names of Features: {features}")
11
```

Result :-

```

python -u "c:\Users\HP\Desktop\college\semester 5\Machine Learning\2_checkFeatures.py"
Data Loaded:
  PassengerId  Survived  Pclass    Name  ...  Ticket    Fare  Cabin  Embarked
0         892         0        3  Kelly, Mr. James  ...  330911    7.8292   NaN      Q
1         893         1        3  Wilkes, Mrs. James (Ellen Needs)  ...  363272    7.0000   NaN      S
2         894         0        2    Myles, Mr. Thomas Francis  ...  240276    9.6875   NaN      Q
3         895         0        3    Wirz, Mr. Albert  ...  315154    8.6625   NaN      S
4         896         1        3  Hirvonen, Mrs. Alexander (Helga E Lindqvist)  ...  3101298   12.2875   NaN      S

[5 rows x 12 columns]
Size of the dataset: (418, 12)
Names of Features: ['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked']

```

Q3 :- Finding the missing entities

Code :-

```

# Finding missing data
missing_data = data.isnull().sum()
print("Missing Entities in Each Feature:")
print(missing_data)

```

Result :-

```

Missing Entities in Each Feature:
PassengerId      0
Survived          0
Pclass           0
Name             0
Sex              0
Age             86
SibSp            0
Parch            0
Ticket           0
Fare             1
Cabin          327
Embarked         0
dtype: int64
PS C:\Users\HP\Desktop\college\semester 5\Machine Learning>

```

Q4 :- Creating of file 1 and file 2 for the missing entities

Code :-

```

1  import pandas as pd
2
3  # Load the Titanic dataset
4  df = pd.read_csv('tested.csv')
5
6  # Identify missing values in critical and non-critical columns
7  missing_critical = df[df[['Age', 'Fare']].isnull().any(axis=1)]
8  missing_non_critical = df[df[['Cabin', 'Ticket']].isnull().any(axis=1)]
9
10 # Save to file1.csv (critical columns)
11 missing_critical.to_csv('file1.csv', index=False)
12
13 # Save to file2.csv (non-critical columns)
14 missing_non_critical.to_csv('file2.csv', index=False)
15

```

Q5 :- Normalization of new files using two approaches

APPROACH 1 :- Using Normalization

Code :-

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3
4  # Load the normalized files
5  file1 = pd.read_csv('file1_normalized.csv')
6  file2 = pd.read_csv('file2_normalized.csv')
7
8  # Plot histograms for numerical columns in file1
9  file1.select_dtypes(include=['float64', 'int64']).hist(figsize=(12, 10), bins=30)
10 plt.suptitle('Histograms for file1')
11 plt.show()
12
13 # Plot histograms for numerical columns in file2
14 file2.select_dtypes(include=['float64', 'int64']).hist(figsize=(12, 10), bins=30)
15 plt.suptitle('Histograms for file2')
16 plt.show()
17

```

APPROACH 2 :- Using Standardized

Code :-

```

1  import pandas as pd
2  from sklearn.preprocessing import StandardScaler
3
4  # Load the files
5  file1 = pd.read_csv('file1.csv')
6  file2 = pd.read_csv('file2.csv')
7
8  # Initialize the StandardScaler
9  scaler = StandardScaler()
10
11 # Normalize the data in both files
12 file1_normalized = file1.copy()
13 file2_normalized = file2.copy()
14
15 # Apply Z-Score Normalization to numerical columns
16 for column in file1_normalized.select_dtypes(include=['float64', 'int64']).columns:
17     file1_normalized[[column]] = scaler.fit_transform(file1_normalized[[column]])
18
19 for column in file2_normalized.select_dtypes(include=['float64', 'int64']).columns:
20     file2_normalized[[column]] = scaler.fit_transform(file2_normalized[[column]])
21
22 # Save the normalized files
23 file1_normalized.to_csv('file1_standardized.csv', index=False)
24 file2_normalized.to_csv('file2_standardized.csv', index=False)

```

Q5 :- Visualization of features using different plots

Using Histograms :-

Code :-

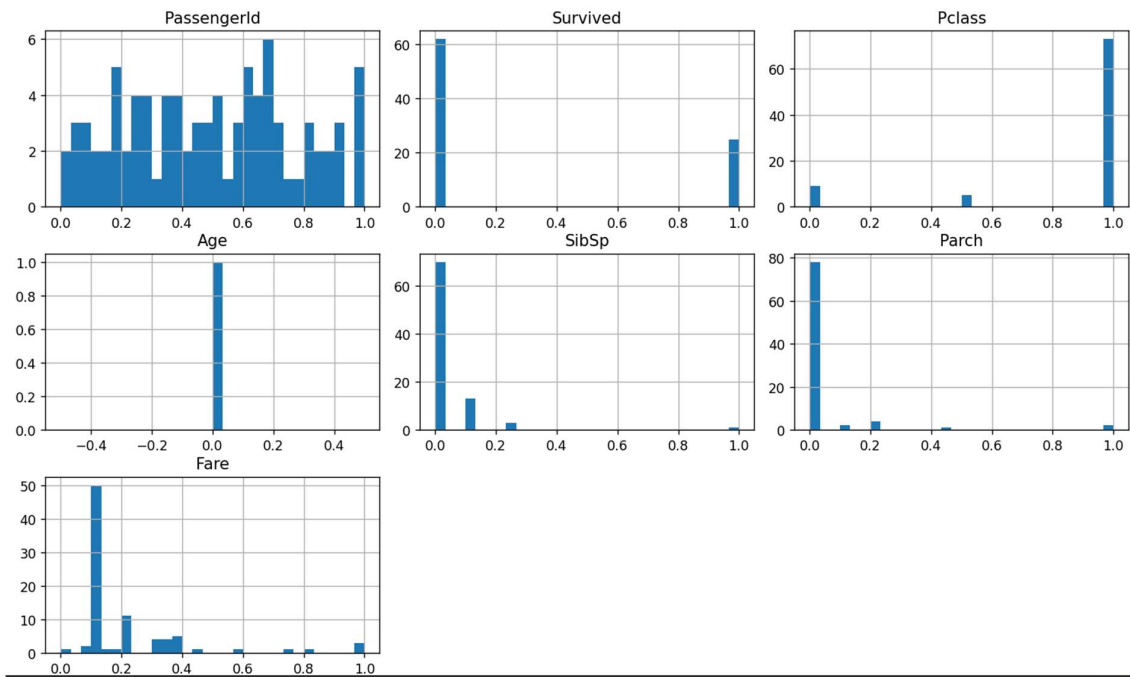
```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3
4  # Load the normalized files
5  file1 = pd.read_csv('file1_normalized.csv')
6  file2 = pd.read_csv('file2_normalized.csv')
7
8  # Plot histograms for numerical columns in file1
9  file1.select_dtypes(include=['float64', 'int64']).hist(figsize=(12, 10), bins=30)
10 plt.suptitle('Histograms for file1')
11 plt.tight_layout(rect=[0, 0.03, 1, 0.95]) # Adjust layout to make room for the title
12 plt.show()
13
14 # Plot histograms for numerical columns in file2
15 file2.select_dtypes(include=['float64', 'int64']).hist(figsize=(12, 10), bins=30)
16 plt.suptitle('Histograms for file2')
17 plt.tight_layout(rect=[0, 0.03, 1, 0.95]) # Adjust layout to make room for the title
18 plt.show()

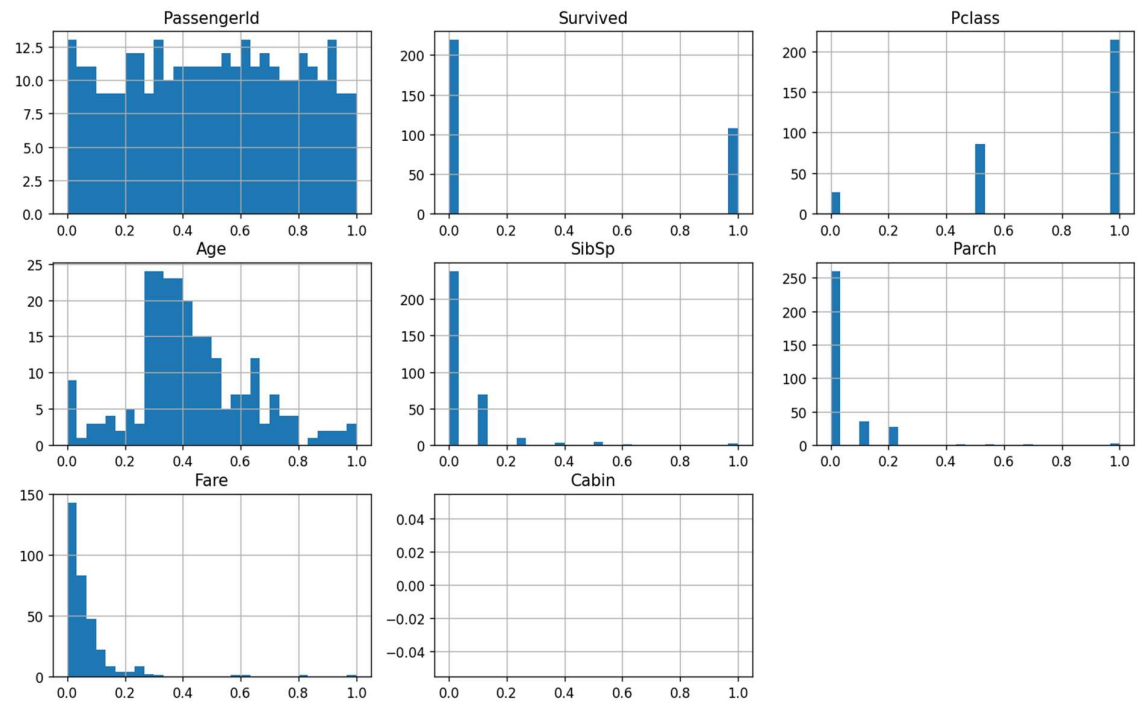
```

Result :-

Histograms for file1



Histograms for file2



Using Scatters :-

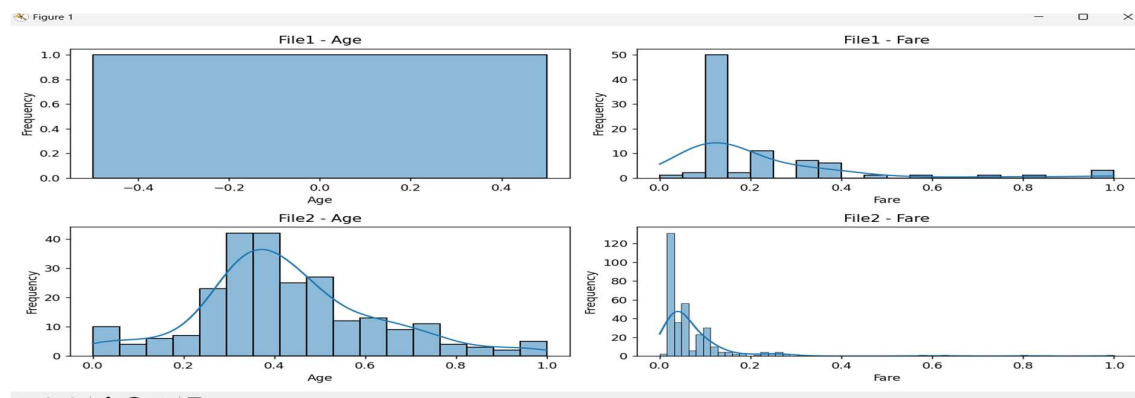
Code :-


```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4
5  # Load the normalized files
6  file1_normalized = pd.read_csv('file1_normalized.csv')
7  file2_normalized = pd.read_csv('file2_normalized.csv')
8
9  # Replace 'Feature1' and 'Feature2' with actual column names
10 features = ['Age', 'Fare'] # Example feature names
11
12 # Plot histograms for file1
13 plt.figure(figsize=(12, 6))
14 for i, feature in enumerate(features):
15     plt.subplot(2, len(features), i + 1)
16     sns.histplot(file1_normalized[feature], kde=True)
17     plt.title(f'File1 - {feature}')
18     plt.xlabel(feature)
19     plt.ylabel('Frequency')
20
21 # Plot histograms for file2
22 for i, feature in enumerate(features):
23     plt.subplot(2, len(features), len(features) + i + 1)
24     sns.histplot(file2_normalized[feature], kde=True)
25     plt.title(f'File2 - {feature}')
26
27     plt.xlabel(feature)
28     plt.ylabel('Frequency')
29
30 plt.tight_layout()
31 plt.show()
32
33 # Plot scatter plots for pairs of features
34 plt.figure(figsize=(12, 6))
35 for i in range(len(features)):
36     for j in range(i + 1, len(features)):
37         plt.subplot(len(features) - 1, len(features) - 1, (i * (len(features) - 1)) + j)
38         plt.scatter(file1_normalized[features[i]], file1_normalized[features[j]], alpha=0.5)
39         plt.title(f'File1: {features[i]} vs {features[j]}')
40         plt.xlabel(features[i])
41         plt.ylabel(features[j])
42
43         plt.subplot(len(features) - 1, len(features) - 1, (i * (len(features) - 1)) + j + len(features) - 1)
44         plt.scatter(file2_normalized[features[i]], file2_normalized[features[j]], alpha=0.5)
45         plt.title(f'File2: {features[i]} vs {features[j]}')
46         plt.xlabel(features[i])
47         plt.ylabel(features[j])
48
49 plt.tight_layout()
50 plt.show()

```

Result :-



Using Mean Errors :-

```
1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4  import numpy as np
5
6  # Load the normalized data
7  file1_normalized = pd.read_csv('file1_normalized.csv')
8  file2_normalized = pd.read_csv('file2_normalized.csv')
9
10 # Select 2-3 features to visualize
11 features = ['PassengerId', 'Survived', 'Pclass'] # Replace with actual feature names
12
13 # Set up the plotting environment
14 fig, axes = plt.subplots(nrows=3, ncols=2, figsize=(15, 12))
15
16 # Plot histograms
17 for i, feature in enumerate(features):
18     sns.histplot(file1_normalized[feature], kde=True, ax=axes[i, 0], color='skyblue')
19     axes[i, 0].set_title(f'Histogram of {feature} (File 1)')
20
21     sns.histplot(file2_normalized[feature], kde=True, ax=axes[i, 1], color='orange')
22     axes[i, 1].set_title(f'Histogram of {feature} (File 2)')
23
24 # Create a new figure for scatter plots
25 fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(18, 5))
26
27 # Scatter plots between features
28 for i, feature in enumerate(features):
29     sns.scatterplot(x=file1_normalized[feature], y=file2_normalized[feature], ax=axes[i])
30     axes[i].set_title(f'Scatter Plot of {feature} (File 1 vs File 2)')
31     axes[i].set_xlabel(f'{feature} (File 1)')
32     axes[i].set_ylabel(f'{feature} (File 2)')
33
34 # Create a new figure for mean error plots
35 fig, ax = plt.subplots(figsize=(8, 6))
36
37 # Calculate the mean and standard deviation for each feature in both files
38 means_file1 = file1_normalized[features].mean()
39 stds_file1 = file1_normalized[features].std()
40 means_file2 = file2_normalized[features].mean()
41 stds_file2 = file2_normalized[features].std()
42
43 # Plot mean and error bars
44 ax.errorbar(features, means_file1, yerr=stds_file1, fmt='o', capsize=5, label='File 1', color='skyblue')
45 ax.errorbar(features, means_file2, yerr=stds_file2, fmt='o', capsize=5, label='File 2', color='orange')
46 ax.set_title('Mean and Error Bars for Selected Features')
47 ax.set_xlabel('Features')
48 ax.set_ylabel('Mean  $\pm$  Standard Deviation')
49 ax.legend()
50
51 # Adjust the layout
52 plt.tight_layout()
53
54 # Show all plots
55 plt.show()
```

Result :-

