

Neural Networks Assignment – 4

ID:700739769

Name: Anjani Priya Marthati

1. Data Manipulation

- a. Read the provided CSV file 'data.csv'.
- b. <https://drive.google.com/drive/folders/1h8C3mLsso-R-slOLsvoYwPLzy2fJ4lOF?usp=sharing>
- c. Show the basic statistical description about the data.
- d. Check if the data has null values. i. Replace the null values with the mean
- e. Select at least two columns and aggregate the data using: min, max, count, mean.
- f. Filter the dataframe to select the rows with calories values between 500 and 1000.
- g. Filter the dataframe to select the rows with calories values > 500 and pulse < 100.
- h. Create a new "df_modified" dataframe that contains all the columns from df except for "Maxpulse".
- i. Delete the "Maxpulse" column from the main df dataframe
- j. Convert the datatype of Calories column to int datatype. k. Using pandas create a scatter plot for the two columns (Duration and Calories).

```

In [5]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import os
4
5 # Read the CSV file into a Pandas dataframe
6 df = pd.read_csv('C:\\Neural networks\\data.csv')
7 #Show the basic statistical description about the data
8 print("Statistics of Data:\n{} \n".format(df.describe()))
9 # Check for null values
10 print("Number of null Values in data per column: \n{} \n".format(df.isnull().sum()))
11 # Replace null values with the mean
12 df.fillna(df.mean(), inplace=True)
13 #Select at least two columns and aggregate the data using: min, max, count, mean.
14 cols = ['Duration', 'Calories']
15 agg = df[cols].agg(['min', 'max', 'count', 'mean'])
16 print("Aggregate data of two columns (Duration, Calories) : \n {} \n".format(agg))
17 # Filter data with calories between 500 and 1000
18 df_500_1000 = df[(df['Calories'] >= 500) & (df['Calories'] <= 1000)]
19 print("Data with calories between 500 and 1000: \n {} \n".format(df_500_1000))
20 # Filter data with calories > 500 and pulse < 100
21 df_500_pulse = df[(df['Calories'] > 500) & (df['Pulse'] < 100)]
22 print("Data with calories > 500 and pulse < 100: \n {} \n".format(df_500_pulse))
23 # Create new dataframe without "Maxpulse" column
24 df_modified = df.drop('Maxpulse', axis=1)
25 # Delete "Maxpulse" column from the main df dataframe
26 df.drop('Maxpulse', axis=1, inplace=True)
27 # Convert "Calories" column to int datatype
28 df['Calories'] = df['Calories'].astype(int)
29 # Scatter plot for "Duration" and "Calories"
30
31 plt.scatter(df['Duration'], df['Calories'])
32 plt.xlabel('Duration')
33 plt.ylabel('Calories')
34 plt.show()
35
36 df.drop('Maxpulse', axis=1, inplace=True)
37 # Convert "Calories" column to int datatype
38 df['Calories'] = df['Calories'].astype(int)
39 # Scatter plot for "Duration" and "Calories"
40
41 plt.scatter(df['Duration'], df['Calories'])
42 plt.xlabel('Duration')
43 plt.ylabel('Calories')
44 plt.show()

```

```

Statistics of Data:
      Duration      Pulse      Maxpulse      Calories
count  169.000000  169.000000  169.000000  164.000000
mean    63.846154  107.461538  134.047337  375.790244
std     42.299949   14.510259   16.450434   266.379919
min     15.000000   80.000000  100.000000   50.300000
25%     45.000000  100.000000  124.000000  250.925000
50%     60.000000  105.000000  131.000000  318.600000
75%     60.000000  111.000000  141.000000  387.600000
max     300.000000  159.000000  184.000000  1860.400000

```

```

Number of null Values in data per column:
Duration      0
Pulse         0
Maxpulse      0
Calories      5
dtype: int64

```

```

Aggregate data of two columns (Duration, Calories) :
      Duration      Calories
min    15.000000   50.300000
max    300.000000  1860.400000
count  169.000000  169.000000
mean    63.846154  375.790244

```

```

Data with calories between 500 and 1000:
  Duration  Pulse  Maxpulse  Calories
51         80    123       146    643.1
62        160    109       135    853.0
65        180     90       130    800.4
66        150    105       135    873.4
67        150    107       130    816.0
72         90    100       127    700.0
73        150     97       127    953.2
75         90     98       125    563.2
78        120    100       130    500.4
83        120    100       130    500.0
90        180    101       127    600.1
99         90     93       124    604.1
101        90     90       110    500.0
102        90     90       100    500.0
103        90     90       100    500.4
106        180     90       120    800.3
108         90     90       120    500.3

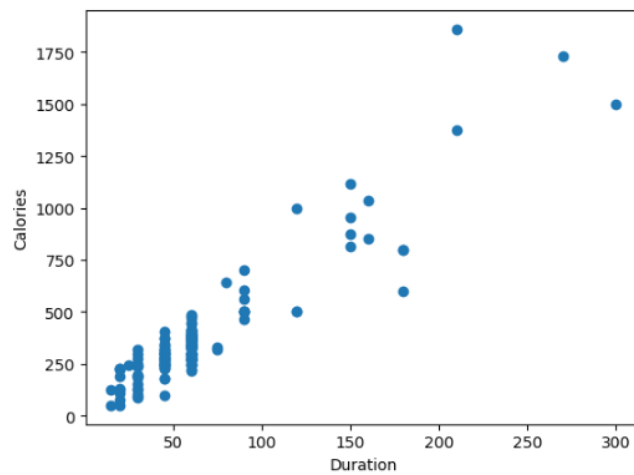
Data with calories > 500 and pulse < 100:
  Duration  Pulse  Maxpulse  Calories
65        180     90       130    800.4
70        150     97       129   1115.0
73        150     97       127    953.2
75         90     98       125    563.2
99         90     93       124    604.1
103        90     90       100    500.4
106        180     90       120    800.3
108         90     90       120    500.3

```

```

75         90     98       125    563.2
99         90     93       124    604.1
103        90     90       100    500.4
106        180     90       120    800.3
108         90     90       120    500.3

```



2. Linear Regression

- Import the given "Salary_Data.csv"
- Split the data in train_test partitions, such that 1/3 of the data is reserved as test subset.
- Train and predict the model.
- Calculate the mean_squared error

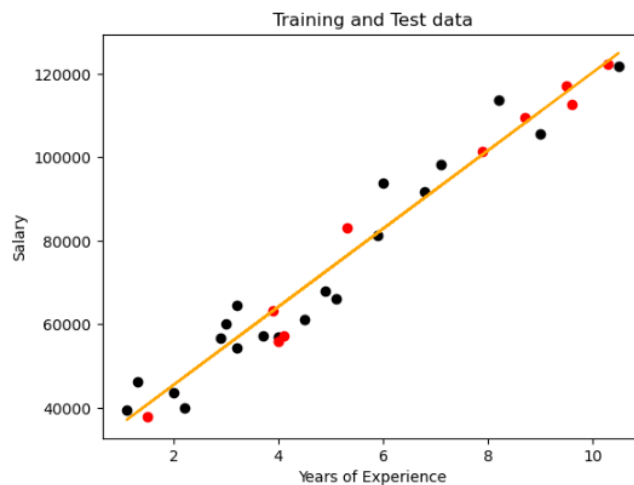
e) Visualize both train and test data using scatter plot.

```
In [9]: 1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6 import matplotlib.pyplot as plt
7
8 # Import the data
9 df = pd.read_csv("C:\\Neural networks\\Salary_Data (2).csv")
10 # Split the data in train_test partitions, such that 1/3 of the data is reserved as test subset
11 X = df[['YearsExperience']]
12 y = df[['Salary']]
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=1/3, random_state=0)
14 # Train and predict the model
15 reg = LinearRegression()
16 reg.fit(X_train, y_train)
17 y_pred = reg.predict(X_test)
18
19 # Calculate the mean squared error
20 mse = mean_squared_error(y_test, y_pred)
21 print("Mean Squared Error: ", mse)
22
23 # Visualize the train and test data using scatter plot
24 plt.scatter(X_train, y_train, color='black')
25 plt.scatter(X_test, y_test, color='red')
26 plt.plot(X_train, reg.predict(X_train), color='orange')
27 plt.xlabel('Years of Experience')
28 plt.ylabel('Salary')
29 plt.title('Training and Test data')
30 plt.show()
```

Mean Squared Error: 21026037.329511296

```
20 plt.ylabel('Salary')
29 plt.title('Training and Test data')
30 plt.show()
```

Mean Squared Error: 21026037.329511296



In []: 1

Video Link:

https://vimeo.com/manage/videos/908491809/e83b72fb19?studio_recording=true&record_session_id=ca777c5e-9963-46ab-942b-73c9fa1c3d31

Github Link:

<https://github.com/Priyamarthati/Assignment4>

