# Neural Network Deep Learning

Assignment – 4

Name: Kishor Kumar Andekar

Student ID: 700744713

Github Link : https://github.com/kishorreyansh/Neural-Network-Deep-Learning/tree/main/Assignment-4

1. Data Manipulation

a. Read the provided CSV file 'data.csv'.

b. https://drive.google.com/drive/folders/1h8C3mLsso-R-sIOLsvoYwPLzy2fJ4IOF?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 the below code snippet, we are doing data manipulation using Pandas:

Reading the provided CSV file 'data.csv' using the read_csv function, assigning a variable df_dataframe to it, and showing the statistical description of the dataframe using the description() function. Replacing the null values in all the columns with the mean. Selecting two columns (duration and pulse) and aggregating the data using the.agg() function and using separate variables to filter the dataframe to select the rows with calorie values between 500 and 1000 and to select the rows with calorie values > 500 and pulse <100. Next, create a new "df_modified" dataframe that contains all the columns from the dataframe except for "Maxpulse" and delete the "Maxpulse" column from the main df_dataframe dataframe. Converting the datatype of the Calories column to an int datatype using the.astype() function and finally using Pandas to create a scatter plot for the two columns (Duration and Calories).

```
datamanipulation.py U ×     linearregression.py U

Assignments > Assignment 4 >  datamanipulation.py > ...
   1    # 1. Data Manipulation
   2        # a. Read the provided CSV file 'data.csv'.
   3        # b. https://drive.google.com/drive/folders/1h8C3mLsso-R-sIOLsvoYwPLzy2fJ4IOF?usp=sharing
   4        # c. Show the basic statistical description about the data.
   5        # d. Check if the data has null values.
   6        # i. Replace the null values with the mean
   7        # e. Select at least two columns and aggregate the data using: min, max, count, mean.
   8        # f. Filter the dataframe to select the rows with calories values between 500 and 1000.
   9        # g. Filter the dataframe to select the rows with calories values > 500 and pulse < 100.
  10        # h. Create a new "df_modified" dataframe that contains all the columns from df except for "Maxpulse".
  11        # i. Delete the "Maxpulse" column from the main df dataframe
  12        # j. Convert the datatype of Calories column to int datatype.
  13        # k. Using pandas create a scatter plot for the two columns (Duration and Calories).
  14
  15    import numpy as np
  16    import matplotlib.pyplot as plt
  17    import pandas as pd   # pip install pandas
  18
  19    # Importing the datasets from data.csv
  20    df_dataframe = pd.read_csv('data.csv')
  21
  22    # Show the basic statistical description about the data.
  23    statistical_description = df_dataframe.describe()
  24    print("********** STATISTICAL DESCRIPTION **********")
  25    print(statistical_description)
  26    print(" ")
  27
```

```python
28    # Applying Only on variables with NaN values
29    # Check if the data has null values. Replace the null values with the mean
30    for i in df_dataframe.columns[df_dataframe.isnull().any(axis=0)]:
31        df_dataframe[i].fillna(df_dataframe[i].mean(),inplace=True)
32    print(df_dataframe)
33    print(" ")
34
35    # Select at least two columns and aggregate the data using: min, max, count, mean
36    selected_colums = ['Duration', 'Pulse']
37    aggregrated_data = df_dataframe[selected_colums].agg(['min', 'max','count','mean'])
38    print(" ****** AGGREGRATED DATA ********** ")
39    print(aggregrated_data)
40
41    # Filter the dataframe to select the rows with calories values between 500 and 1000
42    filtered_df1 = df_dataframe[(df_dataframe['Calories'] >= 500) & (df_dataframe['Calories'] <= 1000)]
43
44    # Filter the dataframe to select the rows with calories values > 500 and pulse < 100
45    filtered_df2 = df_dataframe[(df_dataframe['Calories'] > 500) & (df_dataframe['Pulse'] < 100)]
46
47    # Create a new "df_modified" dataframe that contains all the columns from df except for "Maxpulse".
48    df_modified = df_dataframe.drop(columns='Maxpulse')
49
50    # Delete the "Maxpulse" column from the main df dataframe
51    df_dataframe.drop(columns='Maxpulse',inplace=True)
52
```

```python
53    # Convert the datatype of Calories column to int datatype.
54    df_dataframe['Calories'] = df_dataframe['Calories'].astype(int)
55    print(" DATAFRAME Calories Column INT")
56    print(df_dataframe.info())
57    print(" ******* DATA MODIFIED DATA FRAME ******* ")
58    print(df_modified)
59
60    # Using pandas create a scatter plot for the two columns (Duration and Calories)
61    plt.scatter(df_dataframe['Duration'], df_dataframe['Calories'])
62    plt.xlabel('Duration')
63    plt.ylabel('Calories')
64    plt.title('DURATION AND CALORIES GRAPH')
65    plt.show()
```

Output:



```
datamanipulation.py U ×    linearregression.py U

Assignments > Assignment 4 >  datamanipulation.py > ...
53    # Convert the datatype of Calories column to int datatype.
54    df_dataframe['Calories'] = df_dataframe['Calories'].astype(int)
55    print(" DATAFRAME Calories Column INT")
56    print(df_dataframe.info())
57    print(" ****** DATA MODIFIED DATA FRAME ******* ")
58    print(df_modified)
59
60    # Using pandas create a scatter plot for the two columns (Duration and Calories)
61    plt.scatter(df_dataframe['Duration'], df_dataframe['Calories'])
62    plt.xlabel('Duration')
63    plt.ylabel('Calories')
64    plt.title('DURATION AND CALORIES GRAPH')

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS

PS D:\UCM\Kishor\Neural Network Deep Learning\Assignments\Assignment 4> python .\datamanipulation.py
D:\UCM\Kishor\Neural Network Deep Learning\Assignments\Assignment 4\datamanipulation.py:17: DeprecationWarning:
Pyarrow will become a required dependency of pandas in the next major release of pandas (pandas 3.0),
(to allow more performant data types, such as the Arrow string type, and better interoperability with other libraries)
but was not found to be installed on your system.
If this would cause problems for you,
please provide us feedback at https://github.com/pandas-dev/pandas/issues/54466

  import pandas as pd  # pip install pandas
********** STATISTICAL DESCRIPTION **********
          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

D:\UCM\Kishor\Neural Network Deep Learning\Assignments\Assignment 4\datamanipulation.py:31: FutureWarning: A value is trying to be set on a copy of a DataFrame or Se
ries through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perfo
rm the operation inplace on the original object.
```



```
datamanipulation.py U ×    linearregression.py U

Assignments > Assignment 4 >  datamanipulation.py > ...
53    # Convert the datatype of Calories column to int datatype.
54    df_dataframe['Calories'] = df_dataframe['Calories'].astype(int)
55    print(" DATAFRAME Calories Column INT")
56    print(df_dataframe.info())
57    print(" ****** DATA MODIFIED DATA FRAME ******* ")

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS

  df_dataframe[i].fillna(df_dataframe[i].mean(),inplace=True)
     Duration  Pulse  Maxpulse  Calories
0          60    110       130     409.1
1          60    117       145     479.0
2          60    103       135     340.0
3          45    109       175     282.4
4          45    117       148     406.0
..        ...    ...       ...       ...
164        60    105       140     290.8
165        60    110       145     300.0
166        60    115       145     310.2
167        75    120       150     320.4
168        75    125       150     330.4

[169 rows x 4 columns]

 ****** AGGREGRATED DATA **********
          Duration       Pulse
min      15.000000   80.000000
max     300.000000  159.000000
count   169.000000  169.000000
mean     63.846154  107.461538
 DATAFRAME Calories Column INT
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 169 entries, 0 to 168
Data columns (total 3 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   Duration  169 non-null    int64
 1   Pulse     169 non-null    int64
 2   Calories  169 non-null    int32
dtypes: int32(1), int64(2)
memory usage: 3.4 KB
```
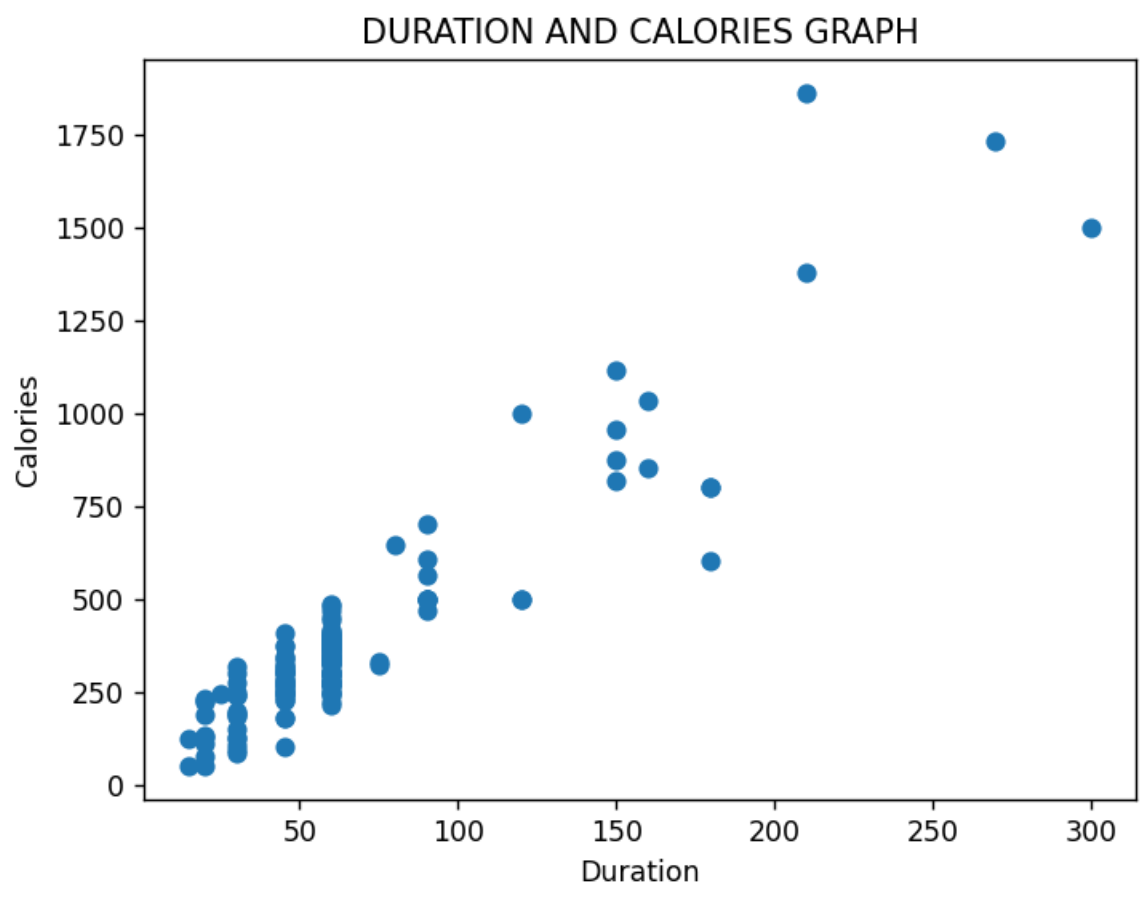
2. Linear Regression
   a) Import the given "Salary_Data.csv"
   b) Split the data in train_test partitions, such that 1/3 of the data is reserved as test subset.
   c) Train and predict the model.
   d) Calculate the mean_squared error
   e) Visualize both train and test data using scatter plot


In the below code snippet, we are training data using linear regression:

Importing the given "Salary_Data.csv" into a variable called df_dataframe. Splitting the data using the train_test_split() function, such that 1/3 of the data is reserved as a test subset. Train the model using the LinearRegression() function, predict the values using the predict() function, calculate the mean squared error of the predicted set, and use Pandas to visualize both the train and test data using a scatter plot.

```python
# 2. Linear Regression
    # a) Import the given "Salary_Data.csv"
    # b) Split the data in train_test partitions, such that 1/3 of the data is reserved as test subset.
    # c) Train and predict the model.
    # d) Calculate the mean_squared error
    # e) Visualize both train and test data using scatter plot

import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# pip install -U scikit-learn
from sklearn.model_selection import train_test_split
# Importing a library
from sklearn.linear_model import LinearRegression

# Mean Squared Error using scikit-learn
from sklearn.metrics import mean_squared_error

# Importing the datasets from data.csv
df_dataframe = pd.read_csv('Salary_Data.csv')

X = df_dataframe[['YearsExperience']]
Y = df_dataframe[['Salary']]

print(" ")
X_Train, X_Test, Y_Train, Y_Test = train_test_split(X,Y,test_size=1/3, random_state = 0)
```

```python
29   # Fitting Simple Linear Regression to the training set
30   # Creating an instance of a model i.e., LinearRegression
31   regressor = LinearRegression()
32
33   # Training the Model
34   regressor.fit(X_Train,Y_Train)
35
36   # Predicting the Test set result ⌨
37   Y_Pred = regressor.predict(X_Test)
38
39   # Calculate the mean_squared error or Mean Squared Deviation
40   mse = mean_squared_error(Y_Test,Y_Pred)
41   print("MEAN SQUARED ERROR: ",mse)
42
43   # Visualize both train and test data using scatter plot
44   plt.scatter(X_Train, Y_Train, label="Training Data")
45   plt.scatter(X_Test, Y_Test, label="Test Data")
46   plt.plot(X_Test, Y_Pred, label="Linear Regression")
47   plt.xlabel("Years of Experience")
48   plt.ylabel("Salary")
49   plt.legend()
50   plt.title("Linear Regression: Salary vs Years of Experience")
51   plt.show()
52
53
```

Output:

```
PS D:\UCM\Kishor\Neural Network Deep Learning\Assignments\Assignment 4> python .\linearregression.py
D:\UCM\Kishor\Neural Network Deep Learning\Assignments\Assignment 4\linearregression.py:10: DeprecationWarning:
Pyarrow will become a required dependency of pandas in the next major release of pandas (pandas 3.0),
(to allow more performant data types, such as the Arrow string type, and better interoperability with other libraries)
but was not found to be installed on your system.
If this would cause problems for you,
please provide us feedback at https://github.com/pandas-dev/pandas/issues/54466

  import pandas as pd

MEAN SQUARED ERROR:  21026037.329511296
```