# ICP – 4

## Naga Phaneendra Kumara Gupta Mogili

## 700757977

GitHub link: https://github.com/nagaphaneendra2001/Deep_Learning_Neural_Networks.git

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).

**Source Code:**

```
import pandas as pandas
import pandas as pd

# Reading CSV file
data_value = pd.read_csv("data.csv")

# Statistical description
data_value.describe()

#Checking for null values
null_values = data_value.isnull().sum()
print(null_values)
```

```python
#replacing nullvalues with the mean
data_value.fillna(data_value.mean(), inplace=True)
print(data_value)

#Selecting two columns
data_value = data_value[["Duration", "Pulse" ]]

#Aggregating the data
agg_dict = {"Duration": ["max", "min", "count", "mean"],
        "Pulse": ["max", "min", "count", "mean"]}
agg_data_value = data_value.agg(agg_dict)
print(agg_data_value)

# Filtering the dataframe to select the rows with calories values between 500
and 1000.
data_value = pd.read_csv("data.csv")
Calories_filter = (data_value["Calories"] >= 500) & (data_value["Calories"] <=
1000)
filtered_data_value = data_value[Calories_filter]
print(filtered_data_value)

#Filtering the dataframe to select the rows with calories values > 500 and
pulse < 100.
data_value = pd.read_csv("data.csv")
Calories_filter = (data_value["Calories"] > 500) & (data_value["Pulse"] < 100)
filtered_data_value = data_value[Calories_filter]
print(filtered_data_value)

#Creating a new "df_modified" dataframe that contains all the columns from
df except for "Maxpulse" and deleting that maxpulse column from the main df
fataframe
df_modified = data_value.drop(columns=["Maxpulse"])
print(df_modified)

# Converting the datatype of calories column to int datatype
data_value['Calories'] = data_value['Calories'].fillna(0).astype(int)
print(data_value)

# Plotting the output
```
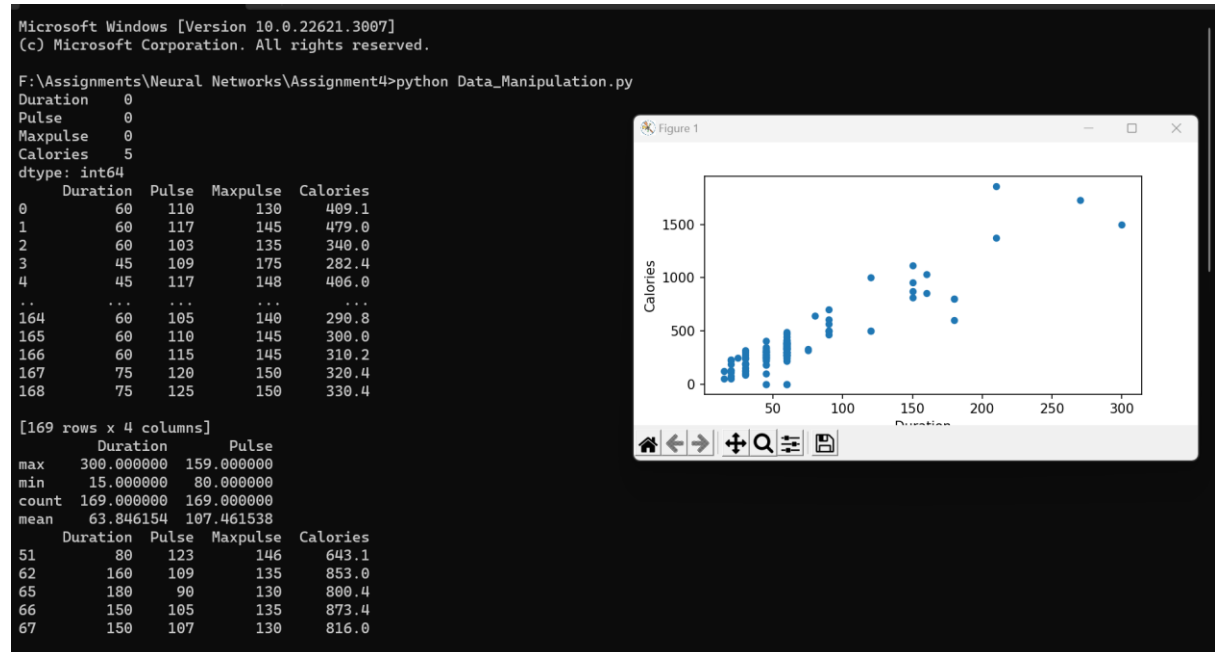
```
import matplotlib.pyplot as plt

data_value.plot(kind='scatter', x='Duration', y='Calories', figsize=(6,3))
plt.show()
```

**Output:**

```
108          90      90       120      500.3
     Duration  Pulse  Calories
0          60     110     409.1
1          60     117     479.0
2          60     103     340.0
3          45     109     282.4
4          45     117     406.0
..        ...    ...       ...
164        60     105     290.8
165        60     110     300.0
166        60     115     310.2
167        75     120     320.4
168        75     125     330.4

[169 rows x 3 columns]
     Duration  Pulse  Maxpulse  Calories
0          60     110       130       409
1          60     117       145       479
2          60     103       135       340
3          45     109       175       282
4          45     117       148       406
..        ...    ...       ...       ...
164        60     105       140       290
165        60     110       145       300
166        60     115       145       310
167        75     120       150       320
168        75     125       150       330

[169 rows x 4 columns]
```

**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.

**Source Code:**

# Simple Linear Regression

# Importing the libraries

```python
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import sklearn as sk


#Importing the datasets
dataset = pd.read_csv("Salary_Data.csv")


X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, 1]


# Splitting the dataset into the Training set and Test set


from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=1/3, random_state=0)


# Fitting Simple Linear Regression to the training set


from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)


# Predicting the Test set result


y_pred = regressor.predict(X_test)


from sklearn.metrics import mean_squared_error
meansquareerror = mean_squared_error(y_true=y_test,y_pred=y_pred)
```

```
print("Mean Square Error:",meansquareerror)


# Visualizing the training set results


viz_train = plt

viz_train.scatter(X_train, y_train, color='green')

viz_train.plot(X_train, regressor.predict(X_train), color='black')

viz_train.title('Training set')

viz_train.xlabel('Years Experience')

viz_train.ylabel('Salary')

viz_train.show()


viz_test = plt

viz_test.scatter(X_test, y_test, color='green')

viz_test.plot(X_train, regressor.predict(X_train), color='black')

viz_test.title('Test set')

viz_test.xlabel('Years Experience')

viz_test.ylabel('Salary')

viz_test.show()
```

**Output:**

Training set



Test set