Video link: https://drive.google.com/file/d/1ROnWeKZzW1d2vbceApyy9kmipLMf1GNn/view?usp=sharing

Github link: https://github.com/naraanjali/Assignment-4

```
min        50.300000    15.000000
max      1860.400000   300.000000
count     169.000000   169.000000
mean      375.790244    63.846154
```

Scatter Plot: Duration vs Calories

```python
#Import Libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import matplotlib.pyplot as plt

#  Load the Dataset"
data = pd.read_csv("Salary_Data.csv")
print("Frist few rows of the data frame")
print(data.head())
# b) Split the data into train_test partitions
X = data[['YearsExperience']]  # Assuming the independent variable is in the 'YearsExperience' column
y = data['Salary']             # Assuming the dependent variable is in the 'Salary' column

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)

#Train and predict the model
model = LinearRegression()
model.fit(X_train, y_train)
y_train_pred = model.predict(X_train)
y_test_pred = model.predict(X_test)

#Calculate the mean_squared error
mse_train = mean_squared_error(y_train, y_train_pred)
mse_test = mean_squared_error(y_test, y_test_pred)

print(f"Mean Squared Error (Train): {mse_train}")
```

```python
    y_train_pred = model.predict(X_train)
    y_test_pred = model.predict(X_test)

    #Calculate the mean_squared error
    mse_train = mean_squared_error(y_train, y_train_pred)
    mse_test = mean_squared_error(y_test, y_test_pred)

    print(f"Mean Squared Error (Train): {mse_train}")
    print(f"Mean Squared Error (Test): {mse_test}")

    #Visualize both train and test data using scatter plot
    plt.scatter(X_train, y_train, color='blue', label='Train Data')
    plt.scatter(X_test, y_test, color='red', label='Test Data')
    plt.plot(X_train, y_train_pred, color='green', linewidth=2, label='Regression Line')
    plt.xlabel('Years of Experience')
    plt.ylabel('Salary')
    plt.title('Salary Prediction Model')
    plt.legend()
    plt.show()
```

[2] ✓ 1.0s                                                                          Python

```
Frist few rows of the data frame
   YearsExperience   Salary
0              1.1  39343.0
1              1.3  46205.0
2              1.5  37731.0
3              2.0  43525.0
4              2.2  39891.0
Mean Squared Error (Train): 29793161.082422983
Mean Squared Error (Test): 35301898.887134895
```



Salary Prediction Model