



Green University Of Bangladesh
Department Of Computer Science and Engineering (CSE)
Faculty of Sciences and Engineering
Semester: (Fall, Year: 2023), B.Sc. in CSE (DAY)

LAB REPORT NO - 05
Course Title: Data Mining Lab
Course Code: CSE-436 **Section:** D2

Lab Experiment Name: Linear & Logistic Regression

Student Details

Name		ID
1	Sk. Nahid	201902073

Lab Date : 25/11/2023
Submission Date : 01/12/2023
Course Teacher Name : Rezwanul Haque

Lab Report Status

Mark:.....	Signature:.....
Comments:.....	Date:.....

1 INTRODUCTION

Regression is a method for understanding the relationship between independent variables or features and a dependent variable or outcome. Linear regression analysis is used to predict the value of a variable based on the value of another variable. Logistic regression estimates the probability of an event occurring.

2 OBJECTIVE

This lab report aims to determine logistic & linear regression and their implementation in python.

3 IMPLEMENTATION

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn import preprocessing, svm
6 from sklearn import metrics
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LinearRegression
9 from sklearn.metrics import mean_absolute_error, mean_squared_error
10 from sklearn.linear_model import LogisticRegression
```

Listing 1: Import Library & Dataset

```
1 df = pd.read_csv('/kaggle/input/average-temperature-from-1900-to-2023/
    Average Temperature 1900-2023.csv')
2 df.head()
```

Listing 2: Read the dataset

```
1 X = np.array(df['Year']).reshape((-1, 1))
2 y = np.array(df['Average_Fahrenheit_Temperature']).reshape((-1, 1))
3 df.dropna(inplace = True)
```

Listing 3: Creating X Y

```
1 Linear regression with 20% split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)
3
4 regr = LinearRegression()
5
6 regr.fit(X_train, y_train)
7 print("Linear Regression: ", regr.score(X_test, y_test))
8
9 y_pred = regr.predict(X_test)
```

```

10 plt.scatter(X_test, y_test, color = 'b')
11 plt.plot(X_test, y_pred, color = 'k')
12 plt.show()
13 mae = mean_absolute_error(y_true=y_test, y_pred=y_pred)
14 mse = mean_squared_error(y_true=y_test, y_pred=y_pred) #default=True
15 rmse = mean_squared_error(y_true=y_test, y_pred=y_pred, squared=False)
16 print("MAE:", mae)
17 print("MSE:", mse)
18 print("RMSE:", rmse)

```

Listing 4: Split the dataset in 8:2 and implementing Linear Regression

```

1 Linear regression with 30% split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3)
3
4 regr = LinearRegression()
5
6 regr.fit(X_train, y_train)
7 print("Linear Regression: ", regr.score(X_test, y_test))
8
9 y_pred = regr.predict(X_test)
10 plt.scatter(X_test, y_test, color = 'b')
11 plt.plot(X_test, y_pred, color = 'k')
12 plt.show()
13 #Evaluation Metrics For Regression
14 mae = mean_absolute_error(y_true=y_test, y_pred=y_pred)
15 mse = mean_squared_error(y_true=y_test, y_pred=y_pred) #default=True
16 rmse = mean_squared_error(y_true=y_test, y_pred=y_pred, squared=False)
17 print("MAE:", mae)
18 print("MSE:", mse)
19 print("RMSE:", rmse)

```

Listing 5: Split the dataset in 7:3 and implementing Linear Regression

```

1 col_names = ['AGE', 'SMOKING']
2 df = pd.read_csv("/kaggle/input/lung-cancer/survey_lung_cancer.csv")
3 X = np.array(df['AGE']).reshape((-1, 1))
4 y = np.array(df['SMOKING']).reshape((-1, 1))
5
6 df.dropna(inplace = True)

```

Listing 6: Preparing dataset for logistic regression

```

1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
0.2, random_state=23)
2
3 logreg = LogisticRegression()
4
5 logreg.fit(X_train, y_train)
6 y_pred=logreg.predict(X_test)
7 cnf_matrix = metrics.confusion_matrix(y_test, y_pred)
8 cnf_matrix

```

```

9
10 class_names=[0,1] # name of classes
11 fig, ax = plt.subplots()
12 tick_marks = np.arange(len(class_names))
13 plt.xticks(tick_marks, class_names)
14 plt.yticks(tick_marks, class_names)
15 # create heatmap
16 sns.heatmap(pd.DataFrame(cnf_matrix), annot=True, cmap="YlGnBu",fmt='g')
17 ax.xaxis.set_label_position("top")
18 plt.tight_layout()
19 plt.title('Confusion matrix', y=1.1)
20 plt.ylabel('Actual label')
21 plt.xlabel('Predicted label')
22 print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
23 print("Precision:",metrics.precision_score(y_test, y_pred))
24 print("Recall:",metrics.recall_score(y_test, y_pred))

```

Listing 7: Logistic Regression on 20% testing

```

1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
  0.3, random_state=23)
2
3 logreg = LogisticRegression()
4
5 logreg.fit(X_train,y_train)
6 y_pred=logreg.predict(X_test)
7 cnf_matrix = metrics.confusion_matrix(y_test, y_pred)
8 cnf_matrix
9
10 class_names=[0,1] # name of classes
11 fig, ax = plt.subplots()
12 tick_marks = np.arange(len(class_names))
13 plt.xticks(tick_marks, class_names)
14 plt.yticks(tick_marks, class_names)
15 # create heatmap
16 sns.heatmap(pd.DataFrame(cnf_matrix), annot=True, cmap="YlGnBu",fmt='g')
17 ax.xaxis.set_label_position("top")
18 plt.tight_layout()
19 plt.title('Confusion matrix', y=1.1)
20 plt.ylabel('Actual label')
21 plt.xlabel('Predicted label')
22 print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
23 print("Precision:",metrics.precision_score(y_test, y_pred))
24 print("Recall:",metrics.recall_score(y_test, y_pred))

```

Listing 8: Logistic Regression on 30% testing

4 OUTPUT

	Year	Average_Fahrenheit_Temperature
0	1900	53.9
1	1901	53.5
2	1902	52.1
3	1903	50.6
4	1904	51.8

Figure 1: Dataset details

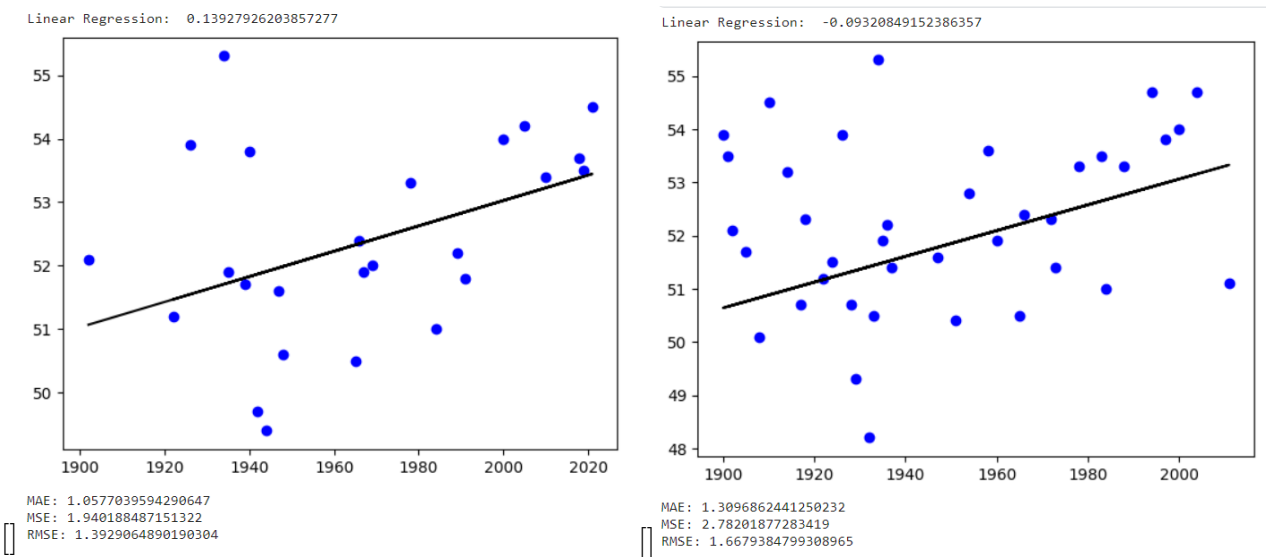


Figure 2: Linear Regression (a) 20% Dataset (b) 30% Dataset

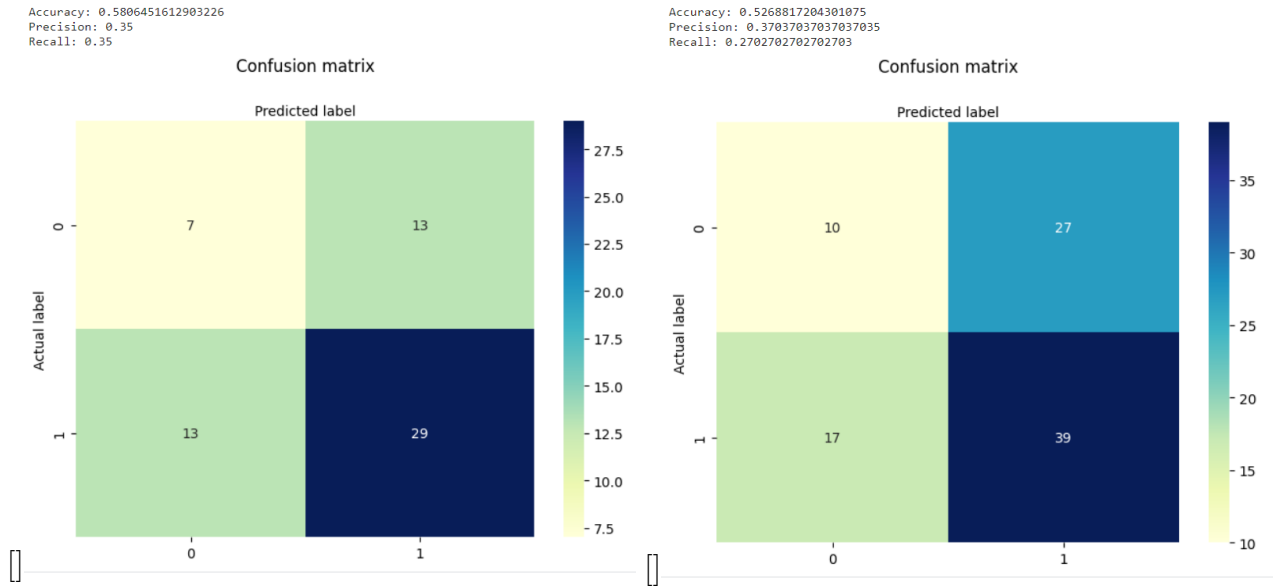


Figure 3: Logistic Regression (a) 20% Dataset (b) 30% Dataset

5 DISCUSSION & ANALYSIS

In this study, logistic regression was applied to a dataset with categorical outcomes, while linear regression was used for a dataset with continuous numerical outcomes. This approach allowed us to address distinct aspects of our research question, leveraging the strengths of each regression technique. The use of different datasets was intentional and aligned with the nature of the variables and the specific objectives of each analysis.