

Logistic Regression

- is used for classification problem
- it uses sigmoid function to predict the values
- the output of sigmoid function is always between 0 & 1
- mostly preferred for binary classification but it also handle multiclass classification

```
In [10]: 1 import numpy as np
          2 import pandas as pd
          3 import matplotlib.pyplot as plt
          4 import seaborn as sns
          5 import warnings
          6 warnings.filterwarnings('ignore')
```

```
In [34]: 1 df = pd.read_csv(r"C:\Users\Bhupendra\Desktop\DataCenter\Logistic Regression
          2 df.head()
```

```
Out[34]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Outcome
0	6	148	72	35	0	33.6	0.627	1
1	1	85	66	29	0	26.6	0.351	0
2	8	183	64	0	0	23.3	0.672	1
3	1	89	66	23	94	28.1	0.167	0
4	0	137	40	35	168	43.1	2.288	1

```
In [36]: 1 df.Outcome.value_counts()
```

```
Out[36]: 0    500
          1    268
          Name: Outcome, dtype: int64
```

In [37]:

```
1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype  
---  -
 0   Pregnancies           768 non-null   int64  
 1   Glucose               768 non-null   int64  
 2   BloodPressure         768 non-null   int64  
 3   SkinThickness         768 non-null   int64  
 4   Insulin               768 non-null   int64  
 5   BMI                   768 non-null   float64 
 6   DiabetesPedigreeFunction 768 non-null   float64 
 7   Age                   768 non-null   int64  
 8   Outcome               768 non-null   int64  
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

In [39]:

```
1 X = df.drop('Outcome', axis = 1)
2 y = df.Outcome
```

Train Test Split

In [58]:

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size = 0.2, ra
```

Model Building

In [62]:

```
1 from sklearn.linear_model import LogisticRegression
2
3 model = LogisticRegression()
4 model.fit(X_train, y_train)    # model training
```

Out[62]:

```
▼ LogisticRegression
LogisticRegression()
```

Prediction

```
In [67]: 1 y_pred = model.predict(X_test)
        2 y_pred          # predicted results
```

```
Out[67]: array([0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0,
                0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0,
                0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
                1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1,
                0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0,
                0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0],
                dtype=int64)
```

```
In [68]: 1 y_test.values          # actual results
```

```
Out[68]: array([0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0,
                0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1,
                0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,
                0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0,
                1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1,
                0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0,
                0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0],
                dtype=int64)
```

Evaluation

Accuracy score of the model

```
In [69]: 1 model.score(X_test, y_test)
```

```
Out[69]: 0.7857142857142857
```

Using different function to calculate accuracy, precision, recall score

```
In [70]: 1 from sklearn.metrics import accuracy_score, precision_score, recall_score
```

```
In [71]: 1 # print accuracy score
        2
        3 accuracy_score(y_test, y_pred)
```

```
Out[71]: 0.7857142857142857
```

```
In [72]: 1 # print precision score
        2
        3 precision_score(y_test, y_pred)
```

```
Out[72]: 0.7837837837837838
```

```
In [73]: 1 # print recall score
          2
          3 recall_score(y_test, y_pred)
```

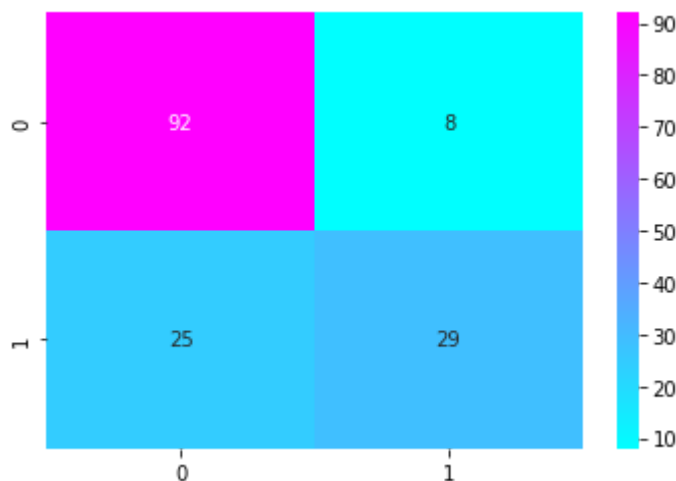
Out[73]: 0.5370370370370371

```
In [74]: 1 from sklearn.metrics import confusion_matrix, classification_report
```

```
In [76]: 1 cm = confusion_matrix(y_test,y_pred)
          2 cm
```

Out[76]: array([[92, 8],
 [25, 29]], dtype=int64)

```
In [82]: 1 # visualize the confusion matrix
          2
          3 sns.heatmap(cm, annot = True, cmap = 'cool');
```



```
In [85]: 1 cr = classification_report(y_test, y_pred)
          2
          3 print(cr)
```

	precision	recall	f1-score	support
0	0.79	0.92	0.85	100
1	0.78	0.54	0.64	54
accuracy			0.79	154
macro avg	0.79	0.73	0.74	154
weighted avg	0.79	0.79	0.77	154

read more about micro, macro and weighted everage from here

<https://towardsdatascience.com/micro-macro-weighted-averages-of-f1-score-clearly-explained-b603420b292f> (<https://towardsdatascience.com/micro-macro-weighted-averages-of-f1-score->

[clearly-explained-b603420b292f](#))

ROC curve

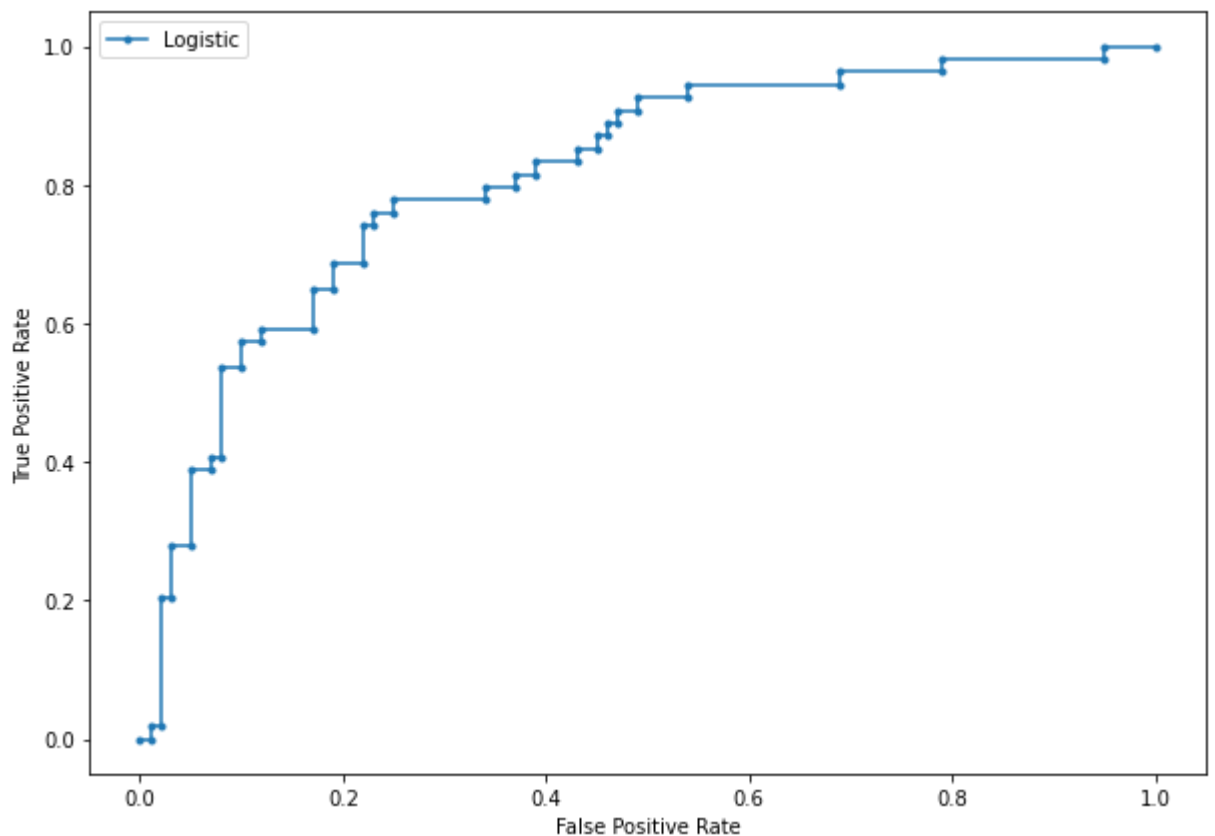
- which stands for receiver operating characteristics shows the tradeoff between **True Positive Rate** vs **False Positive Rate**

In [102]:

```
1 from sklearn.metrics import roc_curve, roc_auc_score
```

```
In [111]: 1 plt.figure(figsize = (10,7))
2 lr_probs = model.predict_proba(X_test)
3
4 # keep probabilities for the positive outcome only
5 lr_probs = lr_probs[:, 1]
6
7 # calculate scores
8 lr_auc = roc_auc_score(y_test, lr_probs)
9
10 # summarize scores
11 print('Logistic: ROC AUC=%.3f' % (lr_auc))
12
13
14 # calculate roc curves
15 lr_fpr, lr_tpr, _ = roc_curve(y_test, lr_probs)
16
17 # plot the roc curve for the model
18 plt.plot(lr_fpr, lr_tpr, marker='.', label='Logistic')
19
20
21 # axis labels
22 plt.xlabel('False Positive Rate')
23 plt.ylabel('True Positive Rate')
24 plt.legend()
25 plt.show()
```

Logistic: ROC AUC=0.815



In []:

1

