

Logistic Regression in Machine Learning

<https://www.geeksforgeeks.org/understanding-logistic-regression/>

Logistic regression is a supervised machine learning algorithm mainly used for classification tasks where the goal is to predict the probability that an instance of belonging to a given class or not. It is a kind of statistical algorithm, which analyze the relationship between a set of independent variables and the dependent binary variables. It is a powerful tool for decision-making. For example email spam or not.

```
1 import pandas as pd
2 import pylab as pl
3 import numpy as np
4 import scipy.optimize as opt
5 import statsmodels.api as sm
6 from sklearn import preprocessing
7
8 import matplotlib.pyplot as plt
9 import matplotlib.mlab as mlab
10 import seaborn as sn
11
12 from google.colab import files
13 uploaded = files.upload()
```

```
[ ] 1 # dataset
2 disease_df = pd.read_csv("framingham.csv")
3 disease_df.drop(['education'], inplace = True, axis = 1)
4 disease_df.rename(columns = {'male': 'Sex_male'}, inplace = True)
5
6 # removing NaN / NULL values
7 disease_df.dropna(axis = 0, inplace = True)
8 print(disease_df.head(), disease_df.shape)
9 print(disease_df.TenYearCHD.value_counts())
```

```
1 # counting no. of patients affected with CHD
2 plt.figure(figsize=(7, 5))
3 sn.countplot(x='TenYearCHD', data=disease_df,
4             palette="BuGn_r")
5 plt.show()
```

```
1 laste = disease_df['TenYearCHD'].plot()
2 plt.show(laste)
```

```
1 X = np.asarray(disease_df[['age', 'Sex_male', 'cigsPerDay',
2                             'totChol', 'sysBP', 'glucose']])
3 y = np.asarray(disease_df['TenYearCHD'])
4
5 # normalization of the dataset
6 X = preprocessing.StandardScaler().fit(X).transform(X)
7
8 # Train-and-Test -Split
9 from sklearn.model_selection import train_test_split
10 X_train, X_test, y_train, y_test = train_test_split(
11     X, y, test_size = 0.3, random_state = 4)
12
13
14 print ('Train set:', X_train.shape, y_train.shape)
15 print ('Test set:', X_test.shape, y_test.shape)
```

```
1 from sklearn.linear_model import LogisticRegression
2
3 logreg = LogisticRegression()
4 logreg.fit(X_train, y_train)
5 y_pred = logreg.predict(X_test)
6
7 # Evaluation and accuracy
8 from sklearn.metrics import jaccard_similarity_score
9
10 print('')
11 print('Accuracy of the model in jaccard similarity score is = ',
12       jaccard_similarity_score(y_test, y_pred))
```

```
1 # This code is contributed by @amartajisce
2 from sklearn.ensemble import RandomForestClassifier
3
4 rf = RandomForestClassifier()
5 rf.fit(X_train, y_train)
6
7 score = rf.score(x_test,y_test)*100
8 print('Accuracy of the model is = ', score)
```

```
1 # Confusion matrix
2 from sklearn.metrics import confusion_matrix, classification_report
3
4 cm = confusion_matrix(y_test, y_pred)
5 conf_matrix = pd.DataFrame(data = cm,
6                             columns = ['Predicted:0', 'Predicted:1'],
7                             index = ['Actual:0', 'Actual:1'])
8
9 plt.figure(figsize = (8, 5))
10 sn.heatmap(conf_matrix, annot = True, fmt = 'd', cmap = "Greens")
11
12 plt.show()
13
14 print('The details for confusion matrix is =')
15 print (classification_report(y_test, y_pred))
16
```

Sr.No	Linear Regression	Logistic Regression
1	Linear regression is used to predict the continuous dependent variable using a given set of independent variables.	Logistic regression is used to predict the categorical dependent variable using a given set of independent variables.
2	Linear regression is used for solving Regression problem.	It is used for solving classification problems.
3	In this we predict the value of continuous variables	In this we predict values of categorical variables
4	In this we find best fit line.	In this we find S-Curve .
5	Least square estimation method is used for estimation of accuracy.	Maximum likelihood estimation method is used for Estimation of accuracy.

6	The output must be continuous value,such as price,age,etc.	Output is must be categorical value such as 0 or 1, Yes or no, etc.
7	It required linear relationship between dependent and independent variables.	It not required linear relationship.
8	There may be collinearity between the independent variables.	There should not be collinearity between independent variable.

Terminologies involved in Logistic Regression:

Here are some common terms involved in logistic regression:

- **Independent variables:** The input characteristics or predictor factors applied to the dependent variable's predictions.
- **Dependent variable:** The target variable in a logistic regression model, which we are trying to predict.
- **Logistic function:** The formula used to represent how the independent and dependent variables relate to one another. The logistic function transforms the input variables into a probability value between 0 and 1, which represents the likelihood of the dependent variable being 1 or 0.
- **Odds:** It is the ratio of something occurring to something not occurring. it is different from probability as the probability is the ratio of something occurring to everything that could possibly occur.
- **Log-odds:** The log-odds, also known as the logit function, is the natural logarithm of the odds. In logistic regression, the log odds of the dependent variable are modeled as a linear combination of the independent variables and the intercept.

- **Coefficient:** The logistic regression model's estimated parameters, show how the independent and dependent variables relate to one another.
- **Intercept:** A constant term in the logistic regression model, which represents the log odds when all independent variables are equal to zero.
- **Maximum likelihood estimation:** The method used to estimate the coefficients of the logistic regression model, which maximizes the likelihood of observing the data given the model.

The logistic regression model transforms the linear regression function continuous value output into categorical value output using a sigmoid function, which maps any real-valued set of independent variables input into a value between 0 and 1. This function is known as the logistic function.

Let the independent input features be

$$X = \begin{bmatrix} x_{11} & \dots & x_{1m} \\ x_{21} & \dots & x_{2m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nm} \end{bmatrix}$$

and the dependent variable is Y having only binary value i.e. 0 or 1.

$$Y = \begin{cases} 0 & \text{if Class 1} \\ 1 & \text{if Class 2} \end{cases}$$

then apply the multi-linear function to the input variables X

$$z = (\sum_{i=1}^n w_i x_i) + b$$

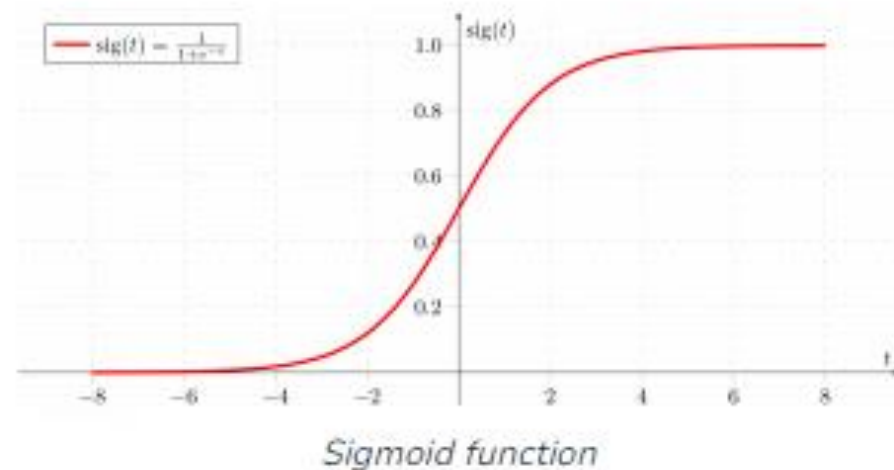
Here x_i is the i th observation of X, $w_i = [w_1, w_2, w_3, \dots, w_m]$ is the weights or Coefficient, and b is the bias term also known as intercept. simply this can be represented as the dot product of weight and bias.

$$z = w \cdot X + b$$

Sigmoid Function

Now we use the sigmoid function where the input will be z and we find the probability between 0 and 1. i.e predicted y .

$$\sigma(z) = \frac{1}{1+e^{-z}}$$



As shown above, the figure sigmoid function converts the continuous variable data into the probability i.e. between 0 and 1.

$$\sigma(z)$$

- tends towards 1 as

$$z \rightarrow \infty$$

$$\sigma(z)$$

- tends towards 0 as

$$z \rightarrow -\infty$$

$$\sigma(z)$$

- is always bounded between 0 and 1

where the probability of being a class can be measured as:

$$P(y = 1) = \sigma(z)$$

$$P(y = 0) = 1 - \sigma(z)$$

Logistic Regression Equation

The odd is the ratio of something occurring to something not occurring. it is different from probability as the probability is the ratio of something occurring to everything that could possibly occur. so odd will be

$$\frac{p(x)}{1-p(x)} = e^z$$

Applying natural log on odd. then log odd will be

$$\begin{aligned}\log \left[\frac{p(x)}{1-p(x)} \right] &= z \\ \log \left[\frac{p(x)}{1-p(x)} \right] &= w \cdot X + b\end{aligned}$$

then the final logistic regression equation will be:

$$p(X; b, w) = \frac{e^{w \cdot X + b}}{1 + e^{w \cdot X + b}} = \frac{1}{1 + e^{-w \cdot X + b}}$$