

Classification

Lesson Objectives

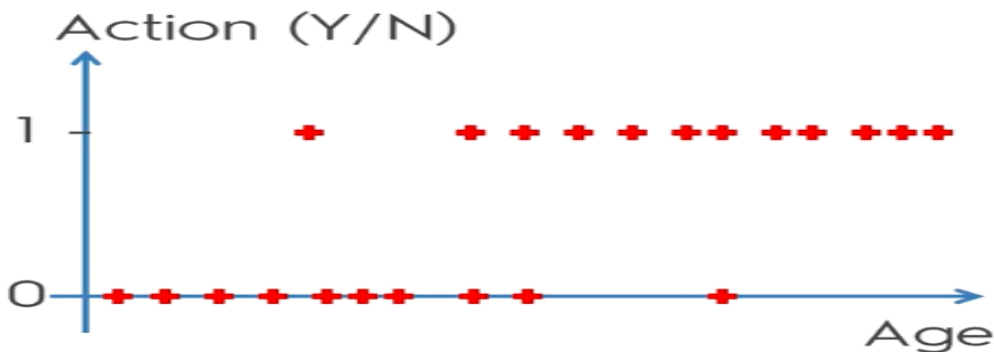


On completion of this lesson, you will be able to understand :

- The key-steps of how to create a Logistic-Regression Model for Classification :
- Train the Model
- Test the Model



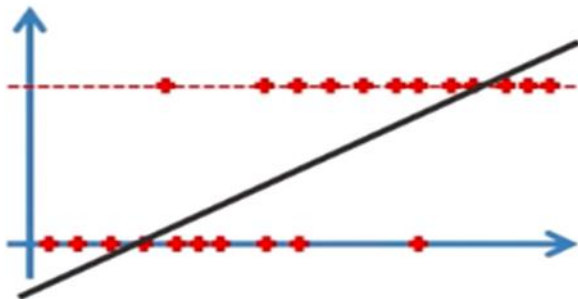
- In classification you predict a category of an Observation
- Classification models are a type of Supervised-Learning
- Classification models include Logistic Regression, SVM, K-NN, Kernel SVM and Random Forests
- Example would be to predict whether a customer would buy or NOT buy a particular product



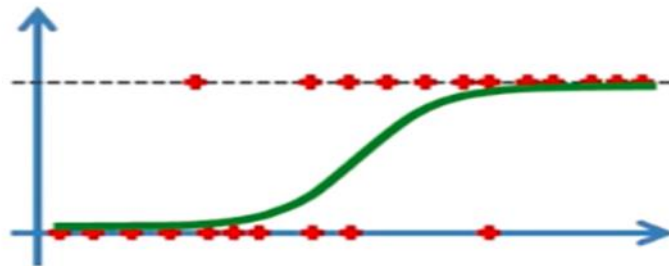


Logistic Linear Regression

- Why Logistic Regression and why not Linear Regression?

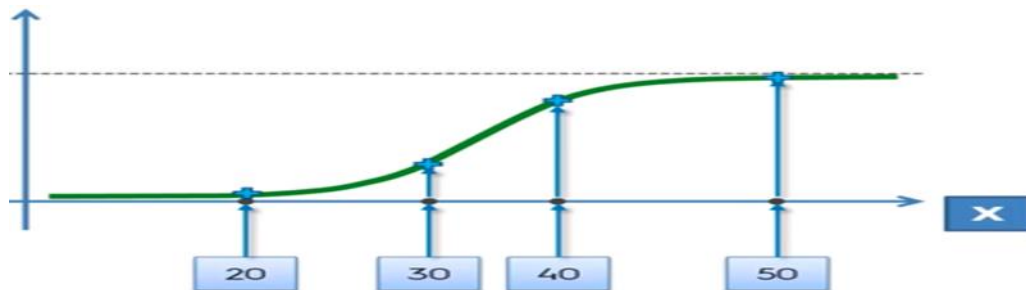


Linear Regression



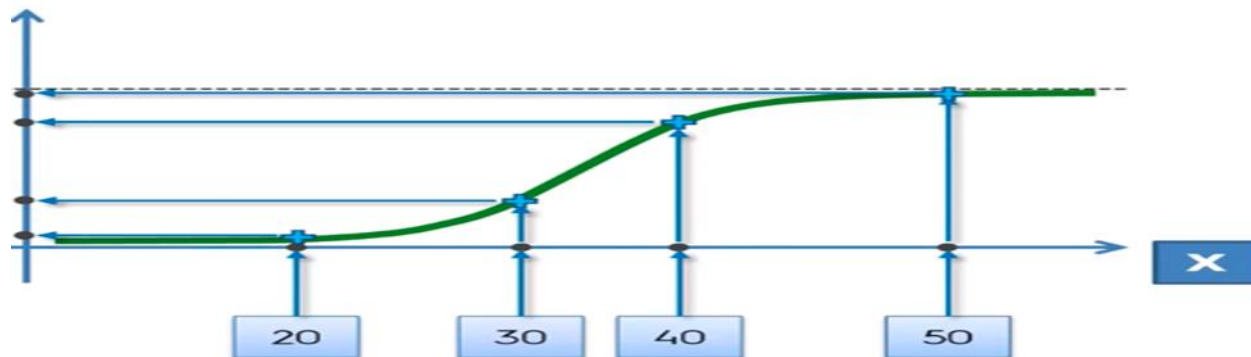
Logistic Regression

- Using Logistic Regression curve

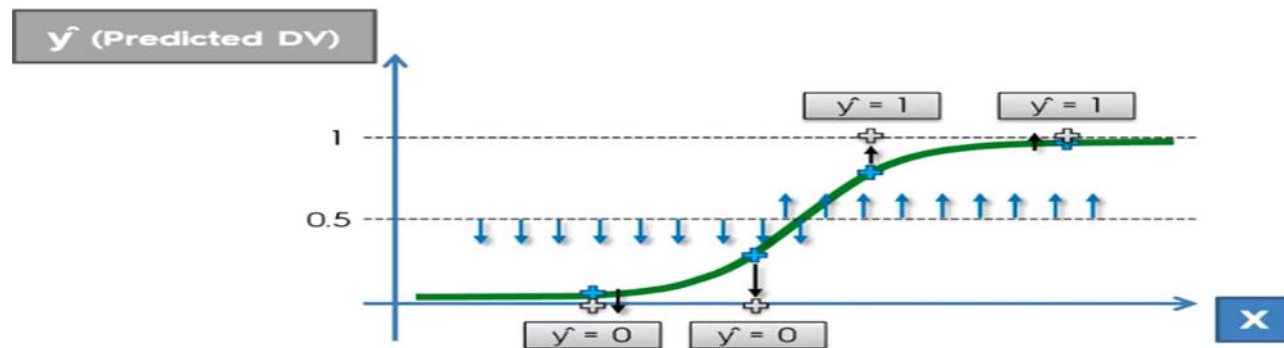




- Plotting the values on the Y-axis



- How to predict





- Change the directory to the path where our data-set lies "C:\BigData\MachineLearning\Machine Learning A-Z Template Folder\Part 3 - Classification\Section 14 - Logistic Regression"
- In a new file in Spyder, paste the template-file code
- Import the basic libraries

```
#Importing the libraries  
import numpy as np  
import matplotlib.pyplot as plt  
import pandas as pd
```



- Import the Dataset

```
# Importing the dataset
```

```
dataset = pd.read_csv('Social_Network_Ads.csv')
```

- Create the matrix of Independent Variables(IVs)

```
x = dataset.iloc[:,[2,3]].values
```

- Create the vector of Dependent Variable(DV)

```
y = dataset.iloc[:,4].values
```

- Check the values of "x" and "y"



Splitting the Data into Training-set and Test-set

- Import the `train_test_split` function from the library `sklearn.model_selection` :
 # Splitting the Data into Training-set and Test-set
 from `sklearn.cross_validation` import `train_test_split`
 or
 from `sklearn.model_selection` import `train_test_split`
- Create 4 variables : `x_train`, `x_test`, `y_train`, `y_test` as follows :
 `x_train, x_test, y_train, y_test = train_test_split(x,y, test_size = 0.25, random_state= 0)`
- Our training-sets : `x_train`(as training data from the matrix of features of independent variables)and `y_train`(as training data from the vector of dependent variables associated with `x_train`)
- Our test-sets : `x_test`(as test-data from the matrix of features of independent variables) and `y_test`(as test-data from the vector of dependent variables associated with `x_test`)



Feature Scaling

- To prevent any IV from the matrix of IVs from dominating the entire ANN architecture, we need to implement feature-scaling

```
#Feature scaling the Data
```

```
from sklearn.preprocessing import StandardScaler
```

```
sc = StandardScaler()
```

```
x_train = sc.fit_transform(x_train)
```

```
x_test = sc.transform(x_test)
```



- Importing the Logistic Regression Class
`from sklearn.linear_model import LogisticRegression`
- Create an object for the Logistic-Regression Model
`classifier = LogisticRegression()`
- Fit the Logistic-Regression Object to the Training-set
`classifier.fit(x_train, y_train)`



Testing the Logistic Regression Model

- Test the model on the Test-dataset
#Predictions on the Test-set
`y_pred = classifier.predict(x_test)`
- Check and Compare the values of the variables `y_pred` and `y_test`
- Using the Confusion-Matrix function to check the accuracy of the predictions

N=200	Predicted FALSE	Predicted TRUE
Actual FALSE	100	12
Actual TRUE	13	75
Percentage of Accuracy $(100+75)/200$		0.875



- Import the function to create the Confusion Matrix
Making the Confusion Matrix
`from sklearn.metrics import confusion_matrix`
- Create an object to implement the Confusion Matrix for our data
`cm = confusion_matrix(y_test, y_pred)`
- Check the value of the "cm" variable

	0	1
0	65	3
1	8	24

N=100	Predicted False	Predicted True
Actual False	65	3
Actual True	8	24



Thank You