# Data Classification Analysis
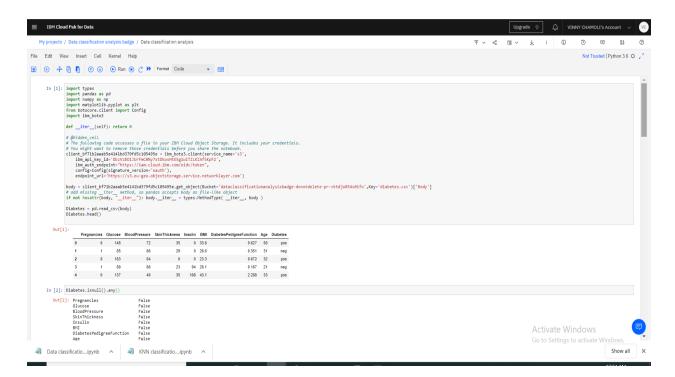
# Diabetes Prediction Using Machine Learning

## *Logistic Regression*

**Logistic Regression is used when the dependent variable (target) is categorical.**
**For example,**
- To predict whether an email is spam (1) or (0)
- Whether the tumor is malignant (1) or not (0)

Logistic regression **is used to describe data and to explain the relationship between one dependent binary variable and one or more nominal, ordinal, interval independent variables. In** Logistic Regression, **we don't directly fit a straight line to our data like in linear regression. Instead, we fit a S shaped curve, called** Sigmoid, **to our observations.**

File  Edit  View  Insert  Cell  Kernel  Help

Not Trusted | Python 3.6 ○

```
In [3]: from sklearn.preprocessing import LabelEncoder
        labelencoder_y = LabelEncoder()
        Diabetes['Diabetes'] = labelencoder_y.fit_transform(Diabetes['Diabetes'])
```

```
In [4]: Diabetes.head()
```

Out[4]:

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Diabetes |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |

```
In [5]: fig = plt.figure(figsize=(18, 14))
        corr = Diabetes.corr()
        c = plt.pcolor(corr)
        plt.yticks(np.arange(0.5, len(corr.index), 1), corr.index)
        plt.xticks(np.arange(0.5, len(corr.columns), 1), corr.columns)
        fig.colorbar(c)
```

Out[5]: <matplotlib.colorbar.Colorbar at 0x7f892551fac8>

---

```
In [7]: from sklearn.model_selection import train_test_split
```

```
In [7]:  from sklearn.model_selection import train_test_split
         X = Diabetes.iloc[:, :-1].values
         y= Diabetes.iloc[:,-1].values
         X_train, X_test, y_train, y_test = train_test_split(X, y,test_size = 0.2,random_state= 0)
```

```
In [8]:  from sklearn.preprocessing import StandardScaler
         sc1 = StandardScaler()
         X_train = sc1.fit_transform(X_train)
         X_test = sc1.transform(X_test)
```

```
In [9]:  from sklearn.linear_model import LogisticRegression
         classifier = LogisticRegression()
         classifier.fit(X_train,y_train)
```
        /opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
          FutureWarning)

```
Out[9]:  LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                  intercept_scaling=1, max_iter=100, multi_class='warn',
                  n_jobs=None, penalty='l2', random_state=None, solver='warn',
                  tol=0.0001, verbose=0, warm_start=False)
```

```
In [10]: y_predict=classifier.predict(X_test)
         y_predict
```

```
Out[10]: array([1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0,
                0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1,
                1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1,
                1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1,
                0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0])
```

```
In [11]: y_predict1 = classifier.predict([[6,48,72,35,0,33.6,0.627,50]])
         y_predict1
```

```
Out[11]: array([1])
```

```
In [12]: from sklearn.metrics import accuracy_score
         ac = accuracy_score(y_test,y_predict)
         ac
```

```
Out[12]: 0.8246753246753247
```

---

```
Out[11]: array([1])
```

```
In [12]: from sklearn.metrics import accuracy_score
         ac = accuracy_score(y_test,y_predict)
         ac
```

```
Out[12]: 0.8246753246753247
```

```
In [13]: from sklearn.metrics import confusion_matrix
         cm = confusion_matrix(y_test,y_predict)
         cm
```

```
Out[13]: array([[98,  9],
                [18, 29]])
```

```
In [14]: import sklearn.metrics as metrics
         fpr, tpr, threashhold =  metrics.roc_curve(y_test,y_predict)
```

```
In [15]: roc_auc=metrics.auc(fpr,tpr)
         roc_auc
```

```
Out[15]: 0.7664545635315172
```

```
In [16]: import matplotlib.pyplot as plt
         plt.plot(fpr,tpr,label='AUC=%0.2f'%roc_auc)
         plt.legend()
         plt.show()
```

# _K Nearest Neighbours Classification Algorithm_

The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement supervised machine learning algorithm that can be used to solve both classification. KNN has no model other than storing the entire dataset, so there is no learning required.

KNN works by finding the distances between a query and all the examples in the data, selecting the specified number examples (K) closest to the query, then votes for the most frequent label (in the case of classification)

A case is classified by a majority vote of its neighbors, with the case being assigned to the class most common amongst its K nearest neighbors measured by a distance function. If K = 5, then the case is simply assigned to the class which has highest majority. It is a memory based algorithm. In this test time is greater than training time and more over K should be an odd number as to classify into either of the categories k should be an odd number. Feature Scaling is must for KNN as the algorithm uses distance calculation for its approach.

File   Edit   View   Insert   Cell   Kernel   Help      Trusted | Python 3.6 ○

Format   Code

```
In [10]: from sklearn.preprocessing import LabelEncoder
         labelencoder_y = LabelEncoder()
         Diabetes['Diabetes'] = labelencoder_y.fit_transform(Diabetes['Diabetes'])
```

```
In [11]: Diabetes.head()
```

Out[11]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Diabetes |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |

```
In [12]: fig = plt.figure(figsize=(18, 14))
         corr = Diabetes.corr()
         c = plt.pcolor(corr)
         plt.yticks(np.arange(0.5, len(corr.index), 1), corr.index)
         plt.xticks(np.arange(0.5, len(corr.columns), 1), corr.columns)
         fig.colorbar(c)
```

Out[12]: <matplotlib.colorbar.Colorbar at 0x7fa09963f5c0>

Format    Code

```python
In [15]: from sklearn.model_selection import train_test_split
         X = Diabetes.iloc[:, :-1].values
         Y= Diabetes.iloc[:,-1].values
         X_train, X_test, Y_train, Y_test = train_test_split(X, Y,test_size = 0.2,random_state= 0)
```

```python
In [16]: from sklearn.preprocessing import StandardScaler
         from sklearn.pipeline import Pipeline
```

```python
In [34]: from sklearn.pipeline import Pipeline
         from sklearn.neighbors import KNeighborsClassifier
         pipe = Pipeline(
             [('rescale', StandardScaler()),
              ('classifier',KNeighborsClassifier(n_neighbors=5, metric = 'euclidean', p=2))
             ])
         pipe.fit(X_train,Y_train)
```

```
Out[34]: Pipeline(memory=None,
              steps=[('rescale', StandardScaler(copy=True, with_mean=True, with_std=True)), ('classifier', KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='euclidean',
                     metric_params=None, n_jobs=None, n_neighbors=5, p=2,
                     weights='uniform'))])
```

```python
In [35]: Y_pred = pipe.predict(X_test)
         Y_pred
```

```
Out[35]: array([1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0,
                0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
                1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 1,
                1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
                1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1,
                0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0])
```

```python
In [36]: y_predict1 = pipe.predict([[1,85,66,29,0,26.6,0.351,31]])
         y_predict1
```

```
Out[36]: array([0])
```

```python
In [37]: from sklearn.metrics import accuracy_score
         ac = accuracy_score(Y_test,Y_pred)
         ac
```

```
Out[37]: 0.7987012987012987
```

---

Format    Code

```
Out[36]: array([0])
```

```python
In [37]: from sklearn.metrics import accuracy_score
         ac = accuracy_score(Y_test,Y_pred)
         ac
```

```
Out[37]: 0.7987012987012987
```

```python
In [38]: from sklearn.metrics import confusion_matrix
         cm = confusion_matrix(Y_test,Y_pred)
         cm
```

```
Out[38]: array([[93, 14],
                [17, 30]])
```

```python
In [39]: import sklearn.metrics as metrics
         fpr, tpr, threashhold =  metrics.roc_curve(Y_test,Y_pred)
```

```python
In [40]: roc_auc=metrics.auc(fpr,tpr)
         roc_auc
```

```
Out[40]: 0.7537283754225492
```

```python
In [41]: import matplotlib.pyplot as plt
         plt.plot(fpr,tpr,label='AUC=%0.2f'%roc_auc)
         plt.legend()
         plt.show()
```