

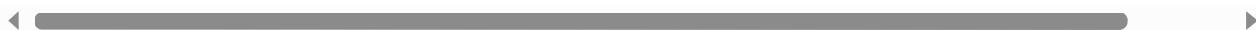
```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
```

```
In [2]: df=pd.read_csv("health.csv")
df
```

```
Out[2]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	6	148	72	35	0	33.6	0.627	50
1	1	85	66	29	0	26.6	0.351	31
2	8	183	64	0	0	23.3	0.672	32
3	1	89	66	23	94	28.1	0.167	21
4	0	137	40	35	168	43.1	2.288	33
...	...	...	...	...	...	...	...	...
763	10	101	76	48	180	32.9	0.171	63
764	2	122	70	27	0	36.8	0.340	27
765	5	121	72	23	112	26.2	0.245	30
766	1	126	60	0	0	30.1	0.349	47
767	1	93	70	31	0	30.4	0.315	23

768 rows × 9 columns



```
In [3]: df.head()
```

```
Out[3]:
```

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



## Data Cleaning and Data Preprocessing

```
In [4]: 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 [5]: `df.describe()`

Out[5]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
<b>count</b>	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
<b>mean</b>	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	33.240659	33.780160	1.016029
<b>std</b>	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	33.998794	5.419054	0.958842
<b>min</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
<b>25%</b>	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	31.000000	24.000000	0.000000
<b>50%</b>	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	32.000000	30.000000	0.000000
<b>75%</b>	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	33.000000	36.000000	1.000000
<b>max</b>	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	67.100000	41.000000	1.000000

In [6]: `df.columns`

Out[6]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'], dtype='object')

In [7]: `feature_matrix = df.iloc[:,0:8]`  
`target_vector = df.iloc[:, -1]`

In [8]: `fs = StandardScaler().fit_transform(feature_matrix)`  
`logr = LogisticRegression()`  
`logr.fit(fs, target_vector)`

Out[8]: LogisticRegression()

In [9]: `observation = [[1, 2, 3, 4, 5, 6, 7, 8]]`  
`prediction = logr.predict(observation)`  
`print(prediction)`

```
[1]
```

```
In [10]: logr.classes_
```

```
Out[10]: array([0, 1], dtype=int64)
```

```
In [11]: logr.predict_proba(observation)
```

```
Out[11]: array([[2.92369487e-04, 9.99707631e-01]])
```

```
In [12]: x = df.iloc[:,0:8]
y = df.iloc[:, -1]
```

```
In [13]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.70)
```

```
In [14]: from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier()
rfc.fit(x_train,y_train)
```

```
Out[14]: RandomForestClassifier()
```

```
In [15]: parameters = {'max_depth':[1,2,3,4,5], 'min_samples_leaf':[5,10,15,20,25],
                        'n_estimators': [10,20,30,40,50]
                        }
```

```
In [16]: from sklearn.model_selection import GridSearchCV
grid_search = GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

```
Out[16]: GridSearchCV(cv=2, estimator=RandomForestClassifier(),
                      param_grid={'max_depth': [1, 2, 3, 4, 5],
                                   'min_samples_leaf': [5, 10, 15, 20, 25],
                                   'n_estimators': [10, 20, 30, 40, 50]},
                      scoring='accuracy')
```

```
In [17]: grid_search.best_score_
```

```
Out[17]: 0.7783804028186206
```

```
In [18]: rfc_best = grid_search.best_estimator_
```

```
In [19]: from sklearn.tree import plot_tree
plt.figure(figsize=(89,40))
plot_tree(rfc_best.estimators_[5], feature_names=x.columns, class_names=['Yes', 'No'],
```

```

Out[19]: [Text(2865.1153846153848, 1993.2, 'Glucose <= 127.5\ngini = 0.459\nsamples = 343\nvalue = [345, 192]\nnclass = Yes'),
  Text(1528.0615384615385, 1630.8000000000002, 'BloodPressure <= 71.0\ngini = 0.31\nsamples = 220\nvalue = [282, 67]\nnclass = Yes'),
  Text(764.0307692307692, 1268.4, 'BMI <= 26.3\ngini = 0.203\nsamples = 127\nvalue = [185, 24]\nnclass = Yes'),
  Text(382.0153846153846, 906.0, 'gini = 0.027\nsamples = 43\nvalue = [72, 1]\nnclass = Yes'),
  Text(1146.0461538461539, 906.0, 'BMI <= 34.85\ngini = 0.281\nsamples = 84\nvalue = [113, 23]\nnclass = Yes'),
  Text(764.0307692307692, 543.5999999999999, 'Age <= 25.5\ngini = 0.317\nsamples = 59\nvalue = [77, 19]\nnclass = Yes'),
  Text(382.0153846153846, 181.19999999999982, 'gini = 0.085\nsamples = 28\nvalue = [43, 2]\nnclass = Yes'),
  Text(1146.0461538461539, 181.19999999999982, 'gini = 0.444\nsamples = 31\nvalue = [34, 17]\nnclass = Yes'),
  Text(1528.0615384615385, 543.5999999999999, 'gini = 0.18\nsamples = 25\nvalue = [36, 4]\nnclass = Yes'),
  Text(2292.0923076923077, 1268.4, 'Insulin <= 11.0\ngini = 0.426\nsamples = 93\nvalue = [97, 43]\nnclass = Yes'),
  Text(1910.076923076923, 906.0, 'gini = 0.486\nsamples = 50\nvalue = [45, 32]\nnclass = Yes'),
  Text(2674.1076923076926, 906.0, 'gini = 0.288\nsamples = 43\nvalue = [52, 11]\nnclass = Yes'),
  Text(4202.169230769231, 1630.8000000000002, 'Insulin <= 149.0\ngini = 0.446\nsamples = 123\nvalue = [63, 125]\nnclass = No'),
  Text(3820.153846153846, 1268.4, 'DiabetesPedigreeFunction <= 0.343\ngini = 0.477\nsamples = 82\nvalue = [49, 76]\nnclass = No'),
  Text(3438.1384615384613, 906.0, 'gini = 0.5\nsamples = 43\nvalue = [31, 33]\nnclass = No'),
  Text(4202.169230769231, 906.0, 'gini = 0.416\nsamples = 39\nvalue = [18, 43]\nnclass = No'),
  Text(4584.184615384615, 1268.4, 'gini = 0.346\nsamples = 41\nvalue = [14, 49]\nnclass = No')]

```

