

```

In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.model_selection import train_test_split

```

```
In [2]: df = pd.read_csv('D:\\24 - Machine_Learning\\download files\\diabetes.csv')
```

Out[2]:	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
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
...
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0

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

Out[3]:	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
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	23.1	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 [4]: df.describe()
```

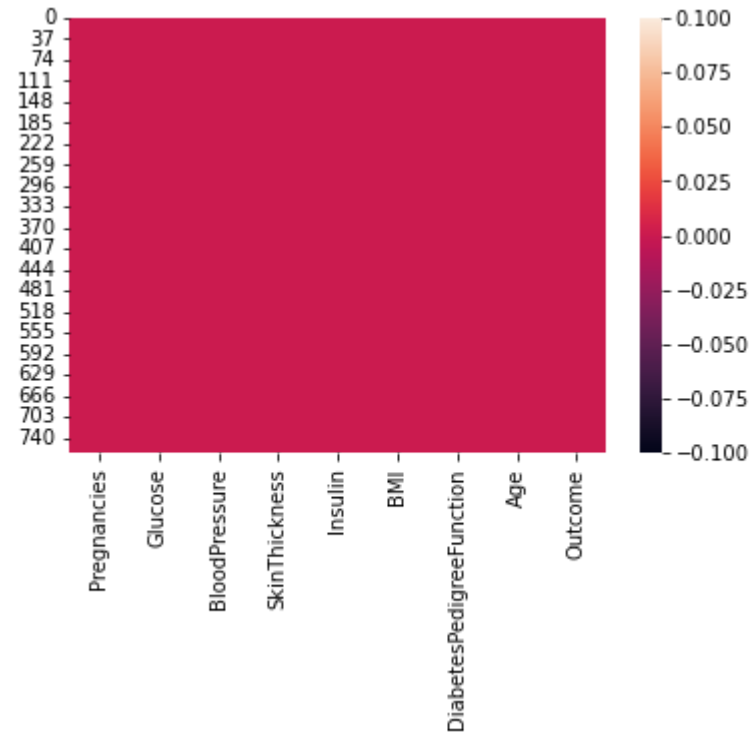
Out [4]:	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.848502	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

```
In [5]: df.isnull()
```

[illegible]

```
In [6]: sns.heatmap(df.isnull())
```

```
Out[6]: <AxesSubplot:>
```



```
In [9]: Y = df.Outcome
```

```
Out[9]:
```

0	1
1	0
2	1
3	0
4	1
...	...
763	0
764	0
765	0
766	1
767	0

Name: Outcome, Length: 768, dtype: int64

```
In [8]: X = df.drop(['Outcome'],axis=1)
```

	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 x 8 columns

```
in [10]: # FOR UNSCALING SKIP THIS STEP
```

```
from sklearn import preprocessing
# Get column names first
names = X.columns
# Creating scaler object
scaler = preprocessing.StandardScaler()
# fit data in the scaler object
scaled_df = scaler.fit_transform(X)
X = pd.DataFrame(scaled_df, columns=names)
```

```
In [11]: X
```

Age	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	0.639947	0.848324	1.049641	0.907270	-0.692891	0.204013	0.468492	1.425995
1	-0.844885	-1.123396	-0.160546	0.539092	-0.692891	-0.684422	-0.360561	-0.190672
2	1.233880	1.943724	-0.263941	-1.288212	-0.692891	1.103255	0.604397	-0.105584
3	-0.844885	-0.998208	-0.160546	0.154533	0.123302	-0.490443	-0.920763	-1.041549
4	-1.141852	0.504055	0.150487	0.907270	0.765836	1.409746	0.584909	-0.020496
...
763	1.827813	-0.622642	0.356432	1.722735	0.870031	0.115169	-0.908682	5.523136
764	-0.547919	0.034598	0.046245	0.405445	-0.692891	0.610154	-0.388282	-0.530102
765	0.342981	0.003301	1.049641	0.154533	0.279594	0.735190	-0.685193	-0.275760
766	-0.844885	0.159787	-0.407332	-1.288212	-0.692891	-0.240205	-0.371101	1.170732
767	-0.844885	-0.873019	0.046245	0.656358	-0.692891	-0.202129	-0.473785	-0.873019

```
In [12]: X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state=41)
```

```
In [13]: LR = LogisticRegression()  
LR.fit(X_train, Y_train)
```

```
Out[13]: LogisticRegression()
```

```
In [14]: Y_pred = LR.predict(X_test)
Y_pred
```

[illegible]

```
In [15]: Y_test
```

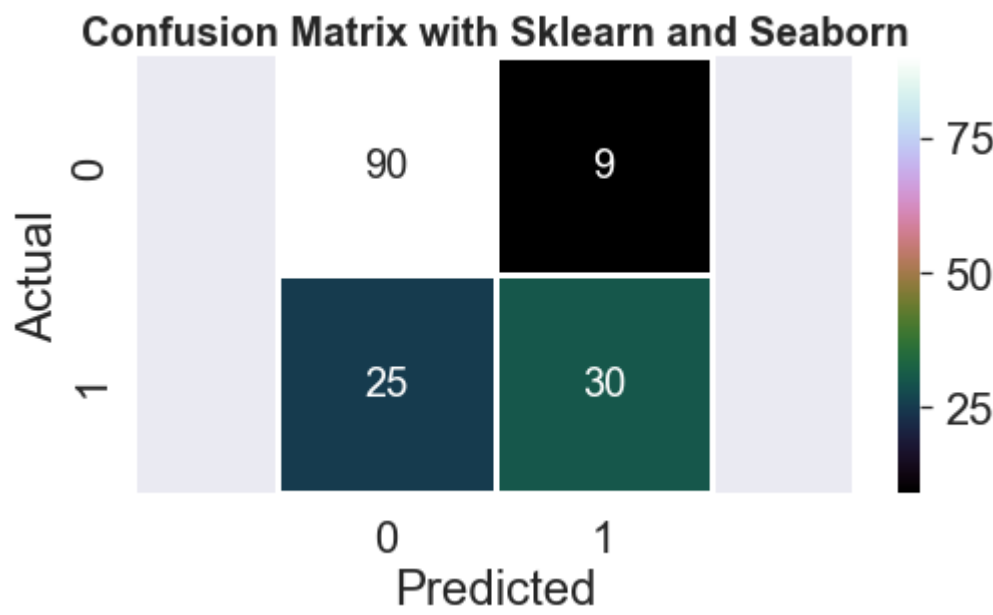
```
Out[15]:
679      0
345      0
486      0
20       0
457      0
..
660      0
627      0
537      0
544      0
518      0
Name: Outcome, Length: 154, dtype: int64
```

```
In [16]: LR.score(X_test, Y_test)
```

```
Out[16]: 0.7792207792207793
```

```
In [18]: cf = confusion_matrix(Y_test, Y_pred)
         cf
```

```
[19]: with plt.style.context('seaborn'):
plt.figure(figsize=(8,4))
sns.set(font_scale=2)
sns.heatmap(cf, annot = True, square=True, annot_kws={'size':20}, linewidth=3, cmap='cubehelix')
plt.xlabel('Predicted');
plt.ylabel('Actual');
plt.axis('equal');
plt.title('Confusion Matrix with Sklearn and Seaborn', fontweights='bold', fontsize=20)
plt.show()
```



```
In [20]: # CLASSIFICATION REPORT
print(classification_report(Y_test, Y_pred))
```

	precision	recall	f1-score	support
0	0.78	0.91	0.84	99
1	0.77	0.55	0.64	55
accuracy			0.78	154
macro avg	0.78	0.73	0.74	154
weighted avg	0.78	0.78	0.77	154