```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 sns.set_theme(color_codes = True)
```

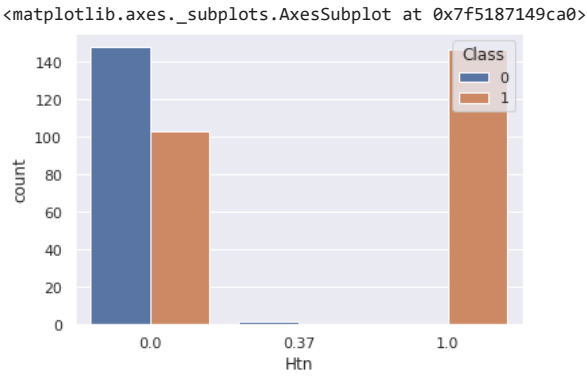## Dataset :

https://www.kaggle.com/datasets/abhia1999/chronic-kidney-disease

```
1 df = pd.read_csv('new_model.csv')
2 df
```

|  | Bp | Sg | Al | Su | Rbc | Bu | Sc | Sod | Pot | Hemo | Wbcc | Rbcc | Htn | Class |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 80.0 | 1.020 | 1.0 | 0.0 | 1.0 | 36.0 | 1.2 | 137.53 | 4.63 | 15.4 | 7800.0 | 5.20 | 1.0 | 1 |
| 1 | 50.0 | 1.020 | 4.0 | 0.0 | 1.0 | 18.0 | 0.8 | 137.53 | 4.63 | 11.3 | 6000.0 | 4.71 | 0.0 | 1 |
| 2 | 80.0 | 1.010 | 2.0 | 3.0 | 1.0 | 53.0 | 1.8 | 137.53 | 4.63 | 9.6 | 7500.0 | 4.71 | 0.0 | 1 |
| 3 | 70.0 | 1.005 | 4.0 | 0.0 | 1.0 | 56.0 | 3.8 | 111.00 | 2.50 | 11.2 | 6700.0 | 3.90 | 1.0 | 1 |
| 4 | 80.0 | 1.010 | 2.0 | 0.0 | 1.0 | 26.0 | 1.4 | 137.53 | 4.63 | 11.6 | 7300.0 | 4.60 | 0.0 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 395 | 80.0 | 1.020 | 0.0 | 0.0 | 1.0 | 49.0 | 0.5 | 150.00 | 4.90 | 15.7 | 6700.0 | 4.90 | 0.0 | 0 |
| 396 | 70.0 | 1.025 | 0.0 | 0.0 | 1.0 | 31.0 | 1.2 | 141.00 | 3.50 | 16.5 | 7800.0 | 6.20 | 0.0 | 0 |
| 397 | 80.0 | 1.020 | 0.0 | 0.0 | 1.0 | 26.0 | 0.6 | 137.00 | 4.40 | 15.8 | 6600.0 | 5.40 | 0.0 | 0 |
| 398 | 60.0 | 1.025 | 0.0 | 0.0 | 1.0 | 50.0 | 1.0 | 135.00 | 4.90 | 14.2 | 7200.0 | 5.90 | 0.0 | 0 |
| 399 | 80.0 | 1.025 | 0.0 | 0.0 | 1.0 | 18.0 | 1.1 | 141.00 | 3.50 | 15.8 | 6800.0 | 6.10 | 0.0 | 0 |

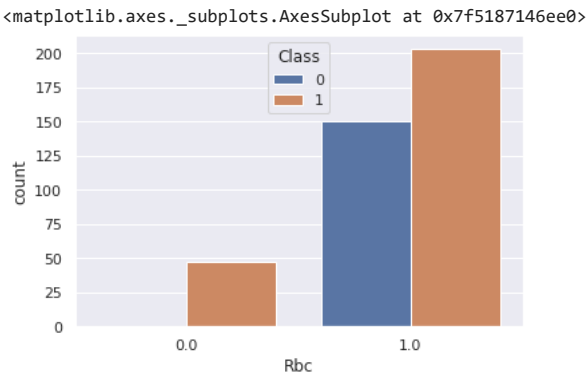400 rows × 14 columns

## Exploratory Data Analysis

```
1 sns.countplot(data=df, x="Htn", hue="Class")
```
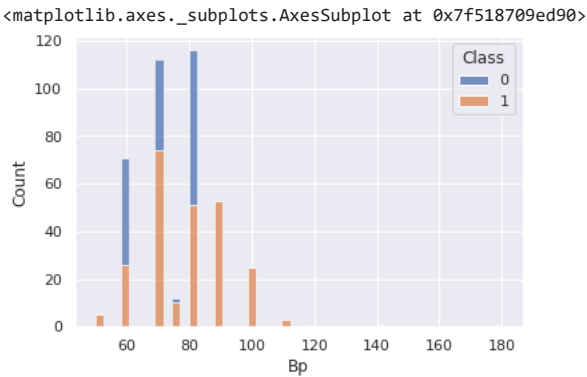
```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5187149ca0>
```



```
1 sns.countplot(data=df, x="Rbc", hue="Class")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5187146ee0>
```



```
1 sns.histplot(data=df, x="Bp", hue="Class", multiple="stack")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f518709ed90>
```

## Data Preprocessing

```
1 df.isnull().sum()
```

```
Bp       0
Sg       0
Al       0
Su       0
Rbc      0
Bu       0
Sc       0
Sod      0
Pot      0
Hemo     0
Wbcc     0
Rbcc     0
Htn      0
Class    0
dtype: int64
```

```
1 #replace 0 value with NaN
2 df_copy = df.copy(deep = True) #deep = True -> Buat salinan indeks dan data dalam dataframe
3 df_copy[['Bp','Sg','Bu','Sc','Sod','Pot','Hemo','Wbcc','Rbcc']] = df_copy[['Bp','Sg','Bu','Sc','Sod','Pot','Hemo','Wbcc','Rbcc']].replace(
4
5 # Showing the Count of NANs
6 print(df_copy.isnull().sum())
```

```
Bp       0
Sg       0
Al       0
Su       0
Rbc      0
Bu       0
Sc       0
Sod      0
Pot      0
Hemo     0
Wbcc     0
Rbcc     0
Htn      0
Class    0
dtype: int64
```

## Check if the class label is balanced or not

```
1 sns.countplot(df['Class'])
2 print(df.Class.value_counts())
```

```
1    250
0    150
Name: Class, dtype: int64
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From
  warnings.warn(
```
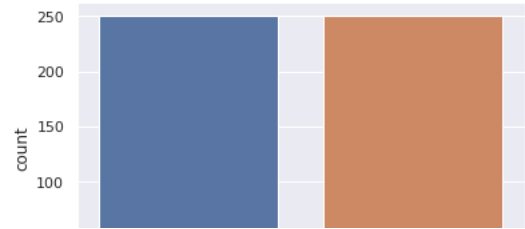


## Do Oversampling Minority Class to Balance the class label

```
 1 from sklearn.utils import resample
 2 #create two different dataframe of majority and minority class
 3 df_majority = df[(df['Class']==1)]
 4 df_minority = df[(df['Class']==0)]
 5 # upsample minority class
 6 df_minority_upsampled = resample(df_minority,
 7                                  n_samples= 250,
 8                                  random_state=0)
 9
10 # Combine majority class with upsampled minority class
11 df2 = pd.concat([df_minority_upsampled, df_majority])
```
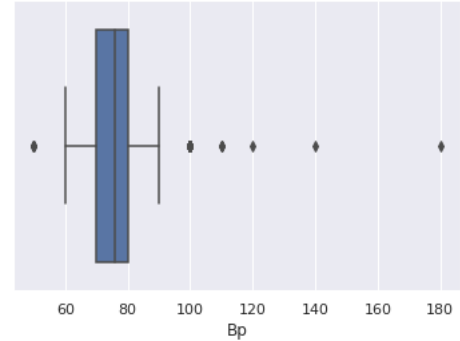
```
1 sns.countplot(df2['Class'])
2 print(df2.Class.value_counts())
```

```
/usr/local/lib/python3.8/dist-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From
  warnings.warn(
0    250
1    250
Name: Class, dtype: int64
```



## Check the Outlier using Boxplot

```
1 sns.boxplot(x=df2["Bp"])
```
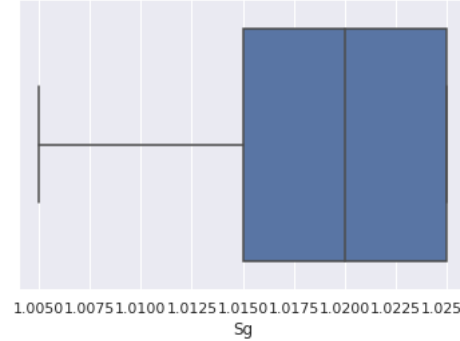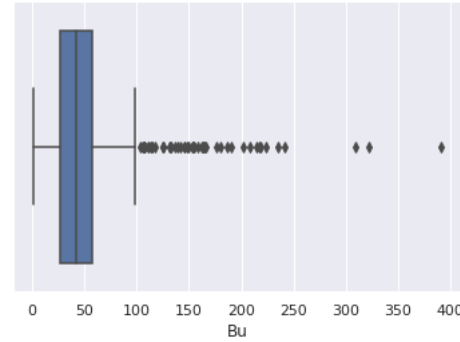
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186ed9970>



```
1 sns.boxplot(x=df2["Sg"])
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f5186e83d90>



```
1 sns.boxplot(x=df2["Bu"])
```
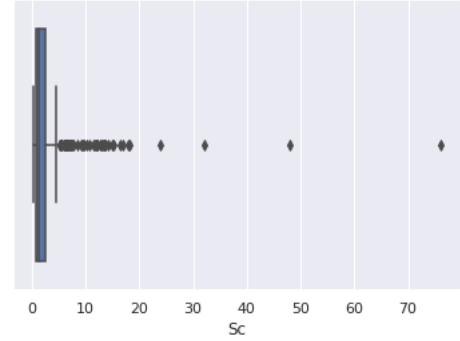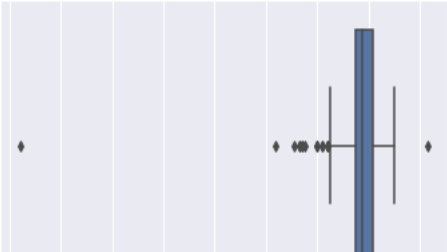
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186de2be0>



```
1 sns.boxplot(x=df2["Sc"])
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f5186dc1dc0>
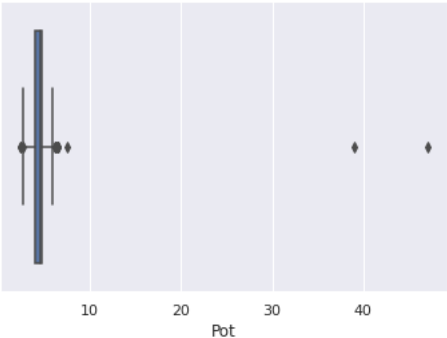


```
1 sns.boxplot(x=df2["Sod"])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186d0c790>
```



```
1 sns.boxplot(x=df2["Pot"])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186e8aa60>
```
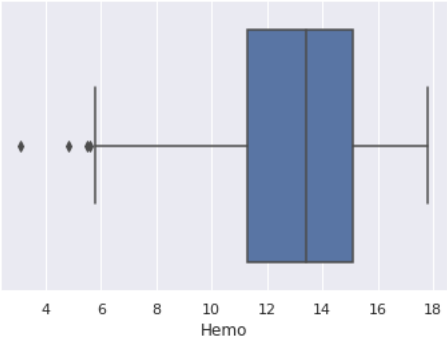


```
1 sns.boxplot(x=df2["Hemo"])
```
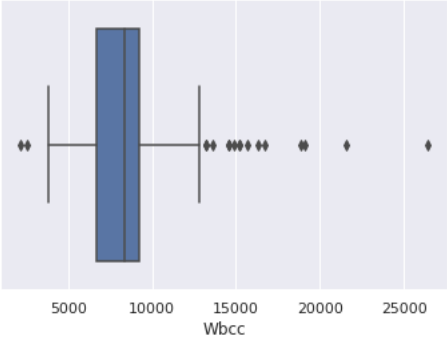
```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186cc4220>
```
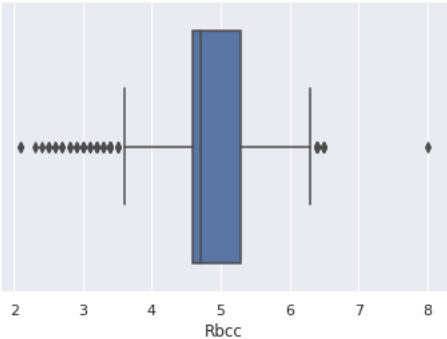


```
1 sns.boxplot(x=df2["Wbcc"])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186c775b0>
```



```
1 sns.boxplot(x=df2["Rbcc"])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186be32b0>
```
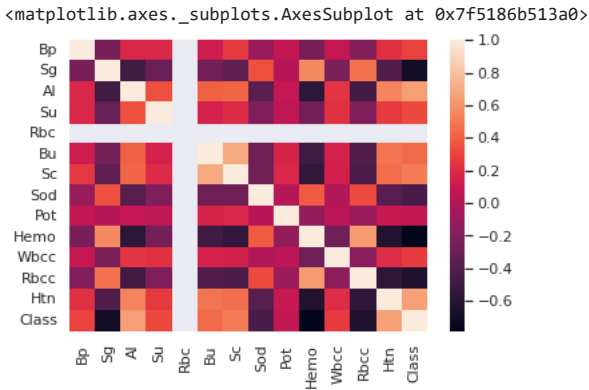


## Remove Outlier using Z-Score

```
1 import scipy.stats as stats
2 z = np.abs(stats.zscore(df2))
3 data_clean = df2[(z<3).all(axis = 1)]
4 data_clean.shape
```
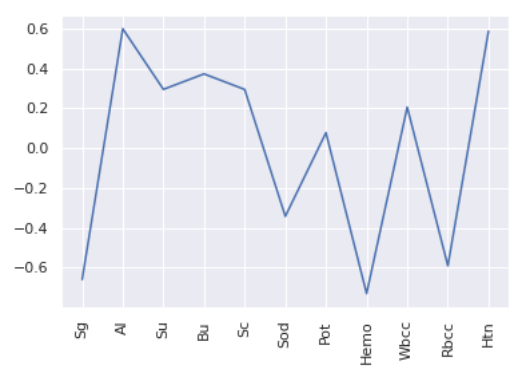
```
(420, 14)
```

# Heatmap Data Correlation

```
1 sns.heatmap(data_clean.corr(), fmt='.2g')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f5186b513a0>
```



```
1 #Rbc attribute is irrlevant, so we have to remove it
2 data_clean2 = df.drop(columns=['Rbc'])
```

```
1 corr = data_clean2[data_clean2.columns[1:]].corr()['Class'][:-1]
2 plt.plot(corr)
3 plt.xticks(rotation=90)
4 plt.show()
```



# Machine Learning Model Building

```
1 X = data_clean2.drop('Class', axis=1)
2 y = data_clean2['Class']
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import accuracy_score
3 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2,random_state=0)
```

# Random Forest

```
1 from sklearn.ensemble import RandomForestClassifier
2 rfc = RandomForestClassifier(random_state=0)
3 rfc.fit(X_train, y_train)
```

```
RandomForestClassifier(random_state=0)
```
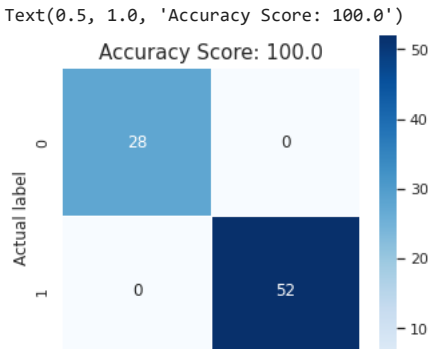
```
1 y_pred = rfc.predict(X_test)
2 print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

```
Accuracy Score : 100.0 %
```

```
1 from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
2 print('F-1 Score : ',(f1_score(y_test, y_pred)))
3 print('Precision Score : ',(precision_score(y_test, y_pred)))
4 print('Recall Score : ',(recall_score(y_test, y_pred)))
```

```
F-1 Score :  1.0
Precision Score :  1.0
Recall Score :  1.0
```

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 cm = confusion_matrix(y_test, y_pred)
3 plt.figure(figsize=(5,5))
4 sns.heatmap(data=cm,linewidths=.5, annot=True,square = True,  cmap = 'Blues')
5 plt.ylabel('Actual label')
6 plt.xlabel('Predicted label')
7 all_sample_title = 'Accuracy Score: {0}'.format(rfc.score(X_test, y_test)*100)
8 plt.title(all_sample_title, size = 15)
```

```
Text(0.5, 1.0, 'Accuracy Score: 100.0')
```



## KNearest Neighbor

```
1 from sklearn.neighbors import KNeighborsClassifier
2 knn = KNeighborsClassifier()
3 knn.fit(X_train, y_train)
```

```
    KNeighborsClassifier()
```
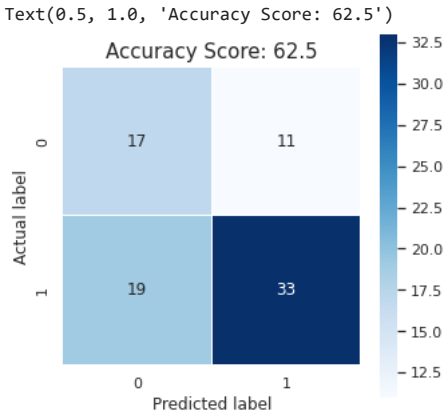
```
1 y_pred = knn.predict(X_test)
2 print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

```
    Accuracy Score : 62.5 %
```

```
1 from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
2 print('F-1 Score : ',(f1_score(y_test, y_pred)))
3 print('Precision Score : ',(precision_score(y_test, y_pred)))
4 print('Recall Score : ',(recall_score(y_test, y_pred)))
```

```
    F-1 Score :  0.6875
    Precision Score :  0.75
    Recall Score :  0.6346153846153846
```

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 cm = confusion_matrix(y_test, y_pred)
3 plt.figure(figsize=(5,5))
4 sns.heatmap(data=cm,linewidths=.5, annot=True,square = True,  cmap = 'Blues')
5 plt.ylabel('Actual label')
6 plt.xlabel('Predicted label')
7 all_sample_title = 'Accuracy Score: {0}'.format(knn.score(X_test, y_test)*100)
8 plt.title(all_sample_title, size = 15)
```

```
    Text(0.5, 1.0, 'Accuracy Score: 62.5')
```



## AdaBoost

```
1 from sklearn.ensemble import AdaBoostClassifier
2 ada = AdaBoostClassifier(random_state=0)
3 ada.fit(X_train, y_train)
```

```
    AdaBoostClassifier(random_state=0)
```
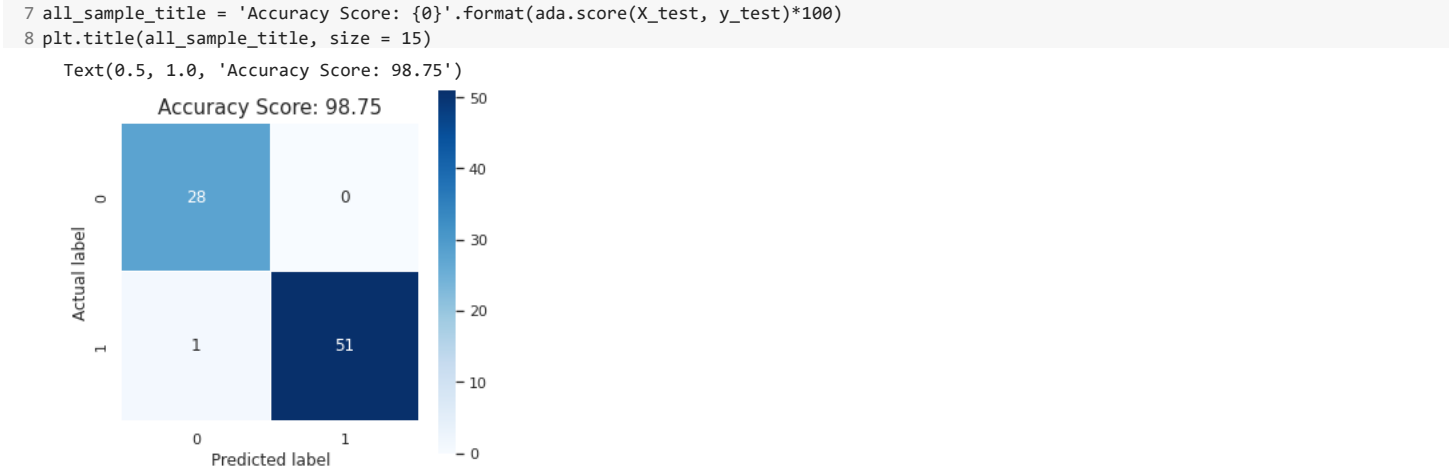
```
1 y_pred = ada.predict(X_test)
2 print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

```
    Accuracy Score : 98.75 %
```

```
1 from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
2 print('F-1 Score : ',(f1_score(y_test, y_pred)))
3 print('Precision Score : ',(precision_score(y_test, y_pred)))
4 print('Recall Score : ',(recall_score(y_test, y_pred)))
```

```
    F-1 Score :  0.9902912621359222
    Precision Score :  1.0
    Recall Score :  0.9807692307692307
```

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 cm = confusion_matrix(y_test, y_pred)
3 plt.figure(figsize=(5,5))
4 sns.heatmap(data=cm,linewidths=.5, annot=True,square = True,  cmap = 'Blues')
5 plt.ylabel('Actual label')
6 plt.xlabel('Predicted label')
```

```
7 all_sample_title = 'Accuracy Score: {0}'.format(ada.score(X_test, y_test)*100)
8 plt.title(all_sample_title, size = 15)
```

```
Text(0.5, 1.0, 'Accuracy Score: 98.75')
```



## ▾ Logistic Regression

```
1 from sklearn.linear_model import LogisticRegression
2 lr = LogisticRegression(random_state = 0)
3 lr.fit(X_train, y_train)
```

```
/usr/local/lib/python3.8/dist-packages/sklearn/linear_model/_logistic.py:814: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
LogisticRegression(random_state=0)
```
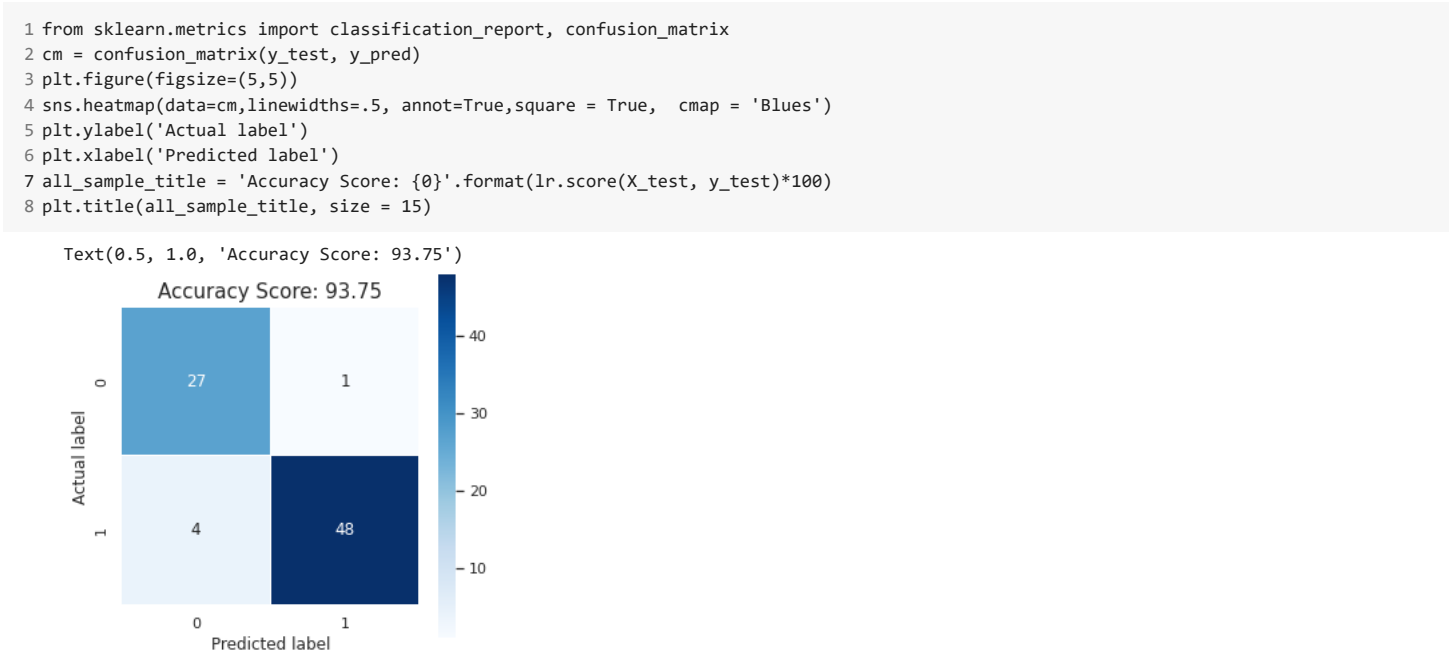
```
1 y_pred = lr.predict(X_test)
2 print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

```
Accuracy Score : 93.75 %
```

```
1 from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
2 print('F-1 Score : ',(f1_score(y_test, y_pred)))
3 print('Precision Score : ',(precision_score(y_test, y_pred)))
4 print('Recall Score : ',(recall_score(y_test, y_pred)))
```

```
F-1 Score :  0.9504950495049506
Precision Score :  0.9795918367346939
Recall Score :  0.9230769230769231
```

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 cm = confusion_matrix(y_test, y_pred)
3 plt.figure(figsize=(5,5))
4 sns.heatmap(data=cm,linewidths=.5, annot=True,square = True,  cmap = 'Blues')
5 plt.ylabel('Actual label')
6 plt.xlabel('Predicted label')
7 all_sample_title = 'Accuracy Score: {0}'.format(lr.score(X_test, y_test)*100)
8 plt.title(all_sample_title, size = 15)
```

```
Text(0.5, 1.0, 'Accuracy Score: 93.75')
```

✓ 0s     completed at 1:05 PM                                                    ● ✕