

Name : Kundan Patil

Email id: patilkundan.1718@gmail.com (<mailto:patilkundan.1718@gmail.com>)

contact: +91 8485010139

linkedIn: www.linkedin.com/in/kundanpatilds (<http://www.linkedin.com/in/kundanpatilds>)

Github: <https://github.com/patilkundan?tab=repositories> (<https://github.com/patilkundan?tab=repositories>)

Q. 7 Write Python Programming for given Dataset "3Classdata.csv" and 2Classdata.csv ¶

Expected Output:

1. Need to use required libraries
2. Apply Classification Algorithm and calculate performance score

In [24]:

```
1 #Let's start with importing required libraries
2 import sklearn
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
10 from sklearn.naive_bayes import GaussianNB
11 from sklearn.svm import SVC
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.ensemble import AdaBoostClassifier
16 from sklearn.model_selection import cross_val_score
17 from sklearn.model_selection import GridSearchCV
18 from sklearn.svm import LinearSVC
19 from scipy import stats
20 import warnings
```

In [99]:

```
1 df = pd.read_csv("2Classdata.csv") # Reading the Data
2 df.head()
```

Out[99]:

	pelvic_incidence	pelvic_tilt_numeric	lumbar_lordosis_angle	sacral_slope	pelvic_radius	degree_spondylolisthesis	class
0	63.027817	22.552586	39.609117	40.475232	98.672917	-0.254400	Abnormal
1	39.056951	10.060991	25.015378	28.995960	114.405425	4.564259	Abnormal
2	68.832021	22.218482	50.092194	46.613539	105.985135	-3.530317	Abnormal
3	69.297008	24.652878	44.311238	44.644130	101.868495	11.211523	Abnormal
4	49.712859	9.652075	28.317406	40.060784	108.168725	7.918501	Abnormal

In [100]:

```
1 df.shape
```

Out[100]:

(310, 7)

In [101]:

```
1 df['class'].unique()
```

Out[101]:

array(['Abnormal', 'Normal'], dtype=object)

In [102]:

```
1 df.info() #df.dtypes also this way can be done
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 310 entries, 0 to 309
Data columns (total 7 columns):
#   Column                Non-Null Count  Dtype
---  -
0   pelvic_incidence       310 non-null    float64
1   pelvic_tilt_numeric     310 non-null    float64
2   lumbar_lordosis_angle   310 non-null    float64
3   sacral_slope            310 non-null    float64
4   pelvic_radius           310 non-null    float64
5   degree_spondylolisthesis 310 non-null    float64
6   class                  310 non-null    object
dtypes: float64(6), object(1)
memory usage: 17.1+ KB
```

In [103]:

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

Out[103]:

```
pelvic_incidence    0
pelvic_tilt_numeric  0
lumbar_lordosis_angle  0
sacral_slope        0
pelvic_radius        0
degree_spondylolisthesis 0
class               0
dtype: int64
```

In [104]:

```
1 duplicate = df[df.duplicated()]
2 print("Duplicate Rows :")
```

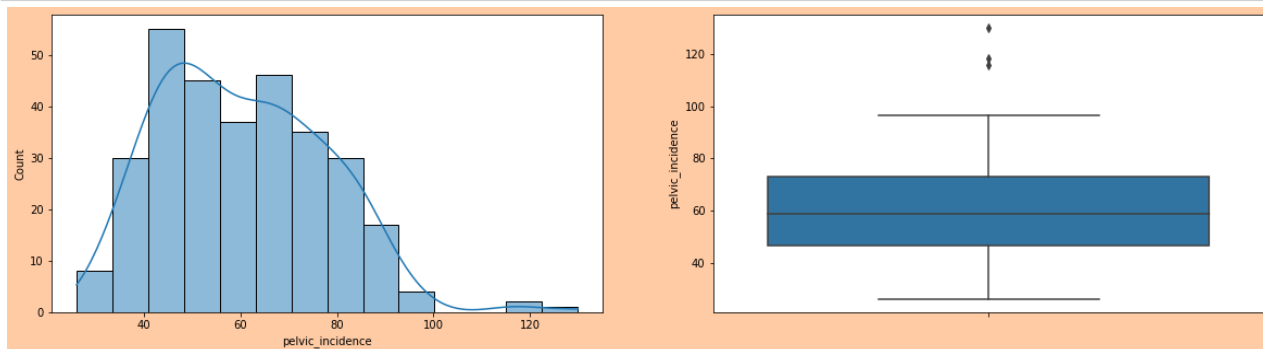
Duplicate Rows :

In [105]:

```
1 def EDA(df,i): # create a function for Continuous variables
2     plt.figure(figsize=(20,5),facecolor='#FFCBA4')
3     plt.subplot(1,2,1)
4     sns.histplot(x=i,data=df,kde=True)
5     plt.subplot(1,2,2)
6     sns.boxplot(y=i,data=df)
7     plt.show()
8     print('skewness of' ,i, 'column---> ',df[i].skew() )
```

In [106]:

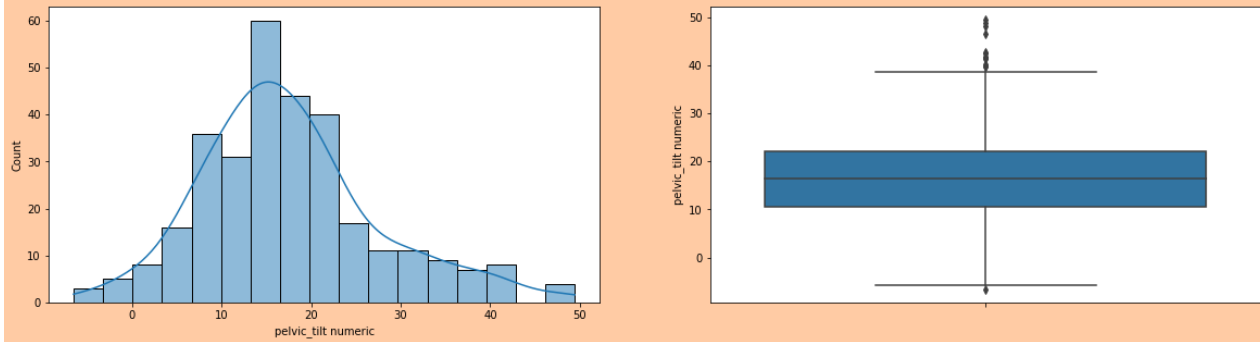
```
1 EDA(df,'pelvic_incidence')
```



skewness of pelvic_incidence column---> 0.5204398948625644

In [107]:

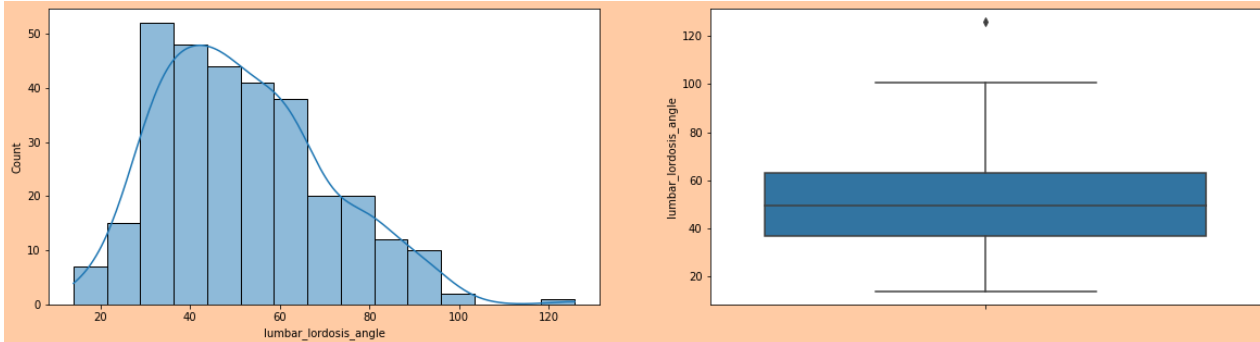
```
1 EDA(df, 'pelvic_tilt_numeric')
```



skewness of pelvic_tilt_numeric column--> 0.6765533590425815

In [108]:

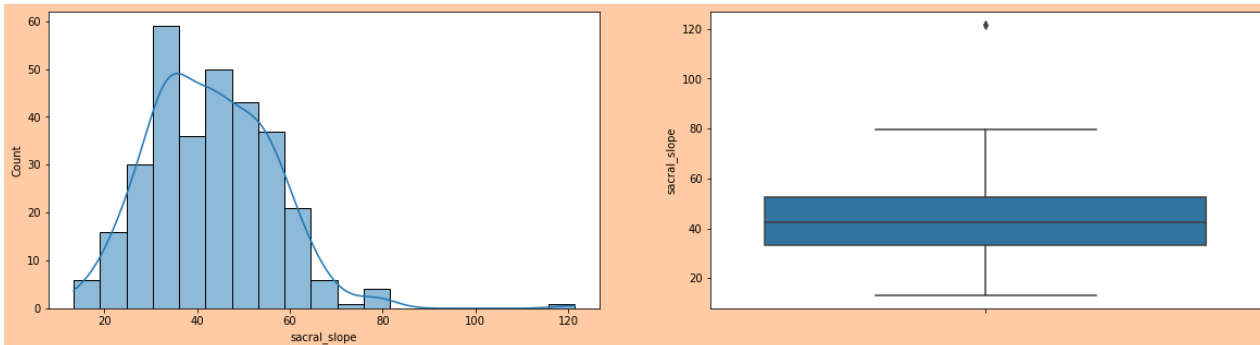
```
1 EDA(df, 'lumbar_lordosis_angle')
```



skewness of lumbar_lordosis_angle column--> 0.5994514775939379

In [109]:

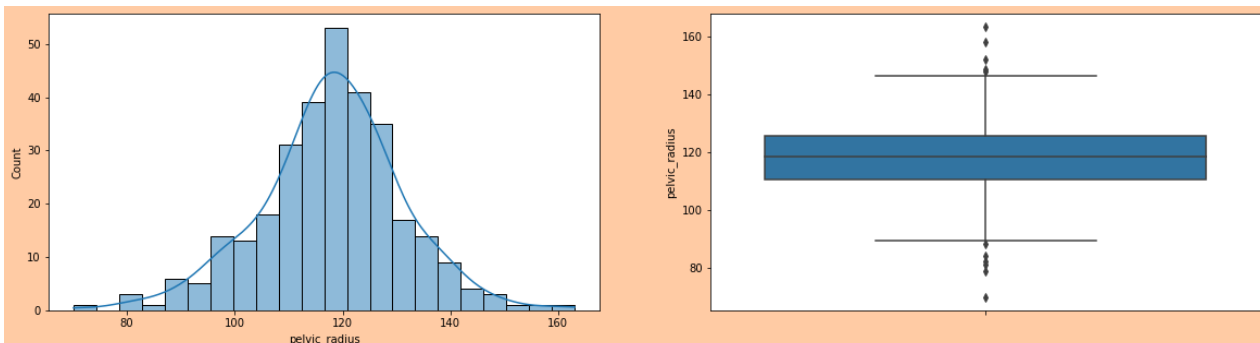
```
1 EDA(df, 'sacral_slope')
```



skewness of sacral_slope column--> 0.7925766941630668

In [110]:

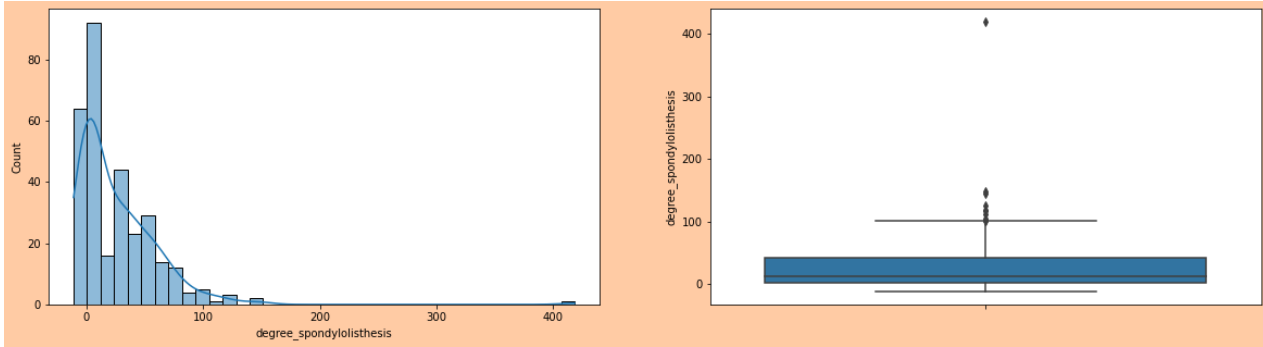
```
1 EDA(df, 'pelvic_radius')
```



skewness of pelvic_radius column--> -0.17683486805355644

In [111]:

```
1 EDA(df, 'degree_spondylolisthesis')
```



skewness of degree_spondylolisthesis column---> 4.317953644012235

Lets check out target variable is imbalanced or not ...!

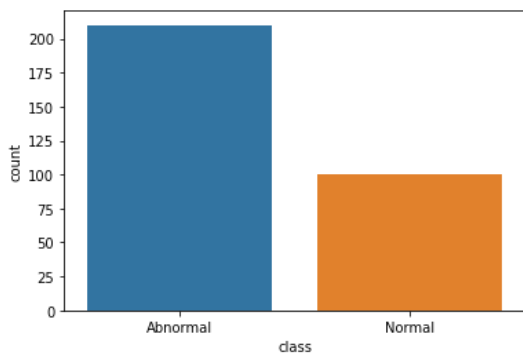
- because this is binary classification problem so we must have balanced data set

In [112]:

```
1 sns.countplot(x='class', data=df)
```

Out[112]:

<AxesSubplot:xlabel='class', ylabel='count'>



In [113]:

```
1 df['class'].unique()
```

Out[113]:

array(['Abnormal', 'Normal'], dtype=object)

Yes data set is balanced

Encoding the target variable

In [114]:

```
1 df['class']=df['class'].replace('Normal', 0)
2 df['class']=df['class'].replace('Abnormal', 1)
```

Outliers treatment

There are several ways to treat outliers when working with data.

- 1) delete the outliers

Z-score method

IQR method

- 2) Transforming the outliers

log or the square root, can help reduce the impact of outliers

In [115]:

```
1 from sklearn.preprocessing import PowerTransformer
2 scaler = PowerTransformer (method='box-cox')
3 df['pelvic_incidence'] = scaler.fit_transform(df[['pelvic_incidence']].values)
```

In [116]:

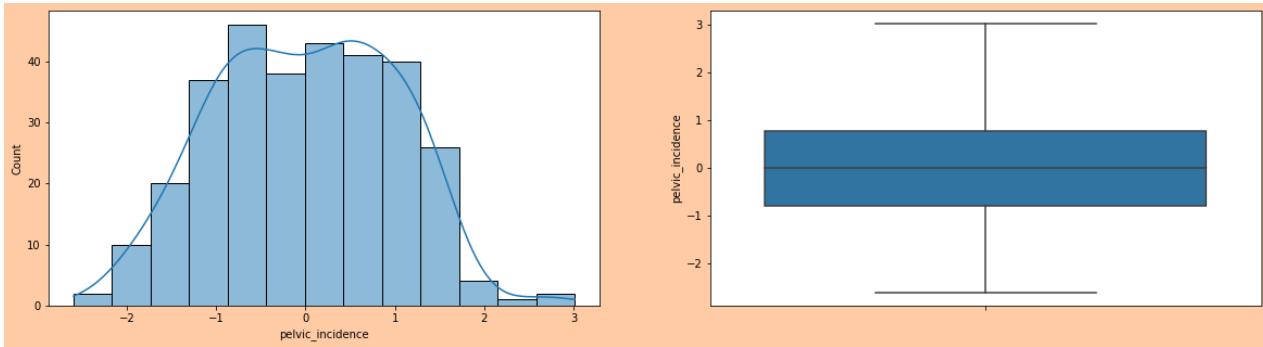
```
1 df['pelvic_incidence'].skew()
```

Out[116]:

-0.011150581982432725

In [117]:

```
1 EDA(df, 'pelvic_incidence')
```



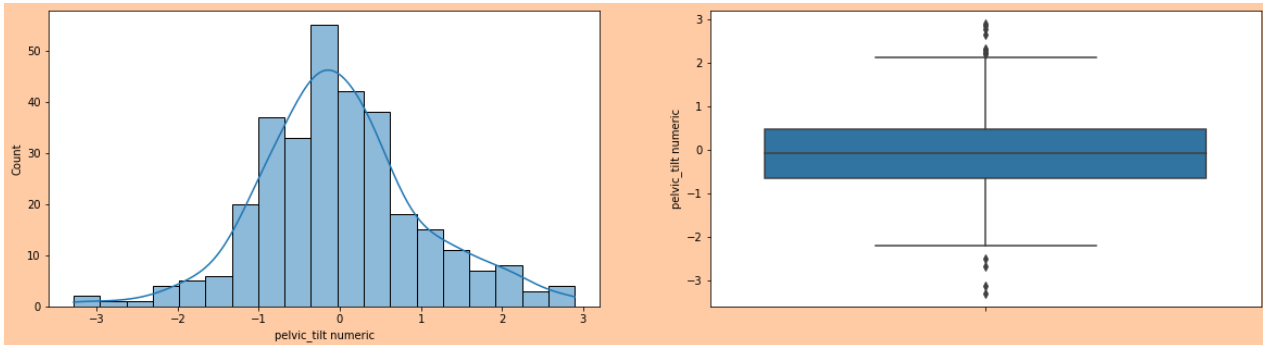
skewness of pelvic_incidence column--> -0.011150581982432725

In [118]:

```
1 from sklearn.preprocessing import PowerTransformer
2 scaler = PowerTransformer (method='yeo-johnson')
3 df['pelvic_tilt numeric'] = scaler.fit_transform(df[['pelvic_tilt numeric']].values)
```

In [119]:

```
1 EDA(df, 'pelvic_tilt numeric')
```



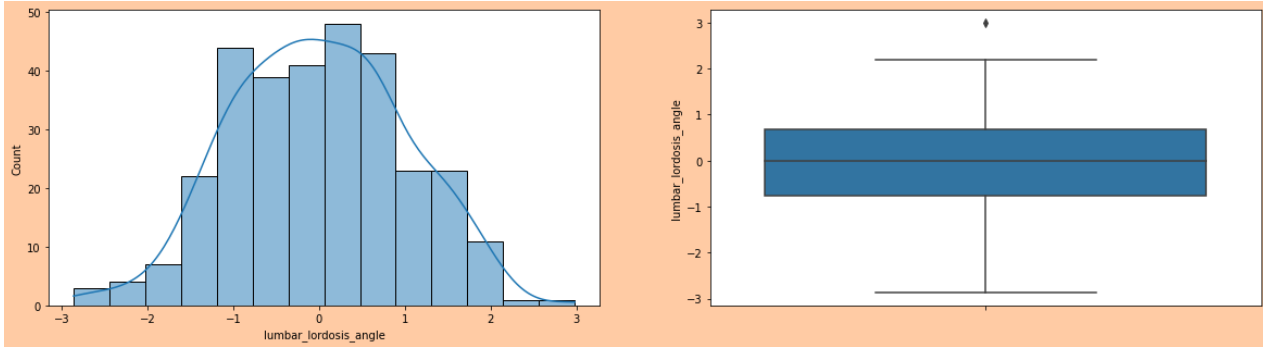
skewness of pelvic_tilt numeric column--> 0.2273319485942726

In [120]:

```
1 from sklearn.preprocessing import PowerTransformer
2 scaler = PowerTransformer (method='box-cox')
3 df['lumbar_lordosis_angle'] = scaler.fit_transform(df[['lumbar_lordosis_angle']].values)
```

In [121]:

```
1 EDA(df, 'lumbar_lordosis_angle')
```



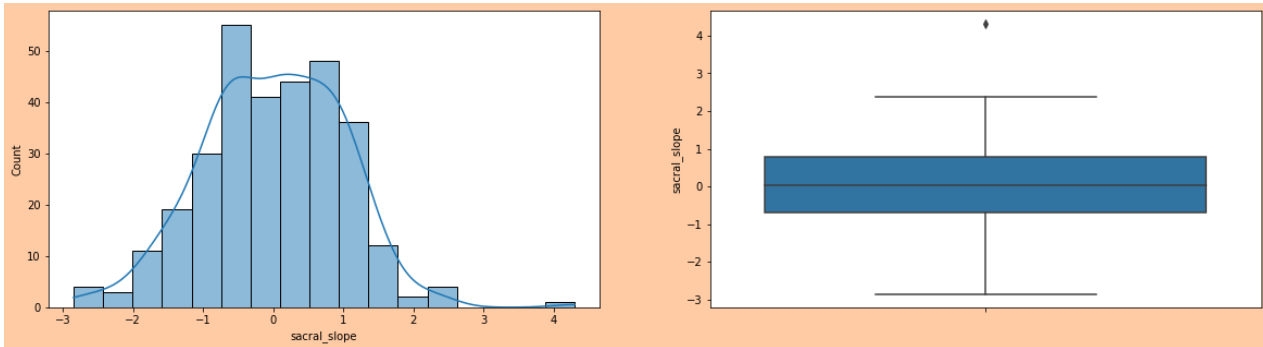
skewness of lumbar_lordosis_angle column---> -0.010697570815072556

In [122]:

```
1 from sklearn.preprocessing import PowerTransformer
2 scaler = PowerTransformer (method='box-cox')
3 df['sacral_slope'] = scaler.fit_transform(df[['sacral_slope']].values)
```

In [123]:

```
1 EDA(df, 'sacral_slope')
```



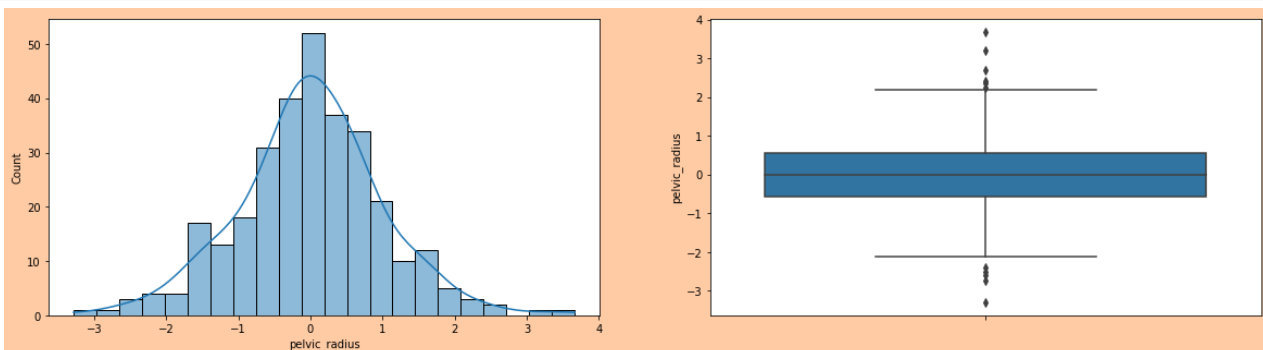
skewness of sacral_slope column---> 0.020529522618659736

In [124]:

```
1 from sklearn.preprocessing import PowerTransformer
2 scaler = PowerTransformer (method='box-cox')
3 df['pelvic_radius'] = scaler.fit_transform(df[['pelvic_radius']].values)
```

In [125]:

```
1 EDA(df, 'pelvic_radius')
```



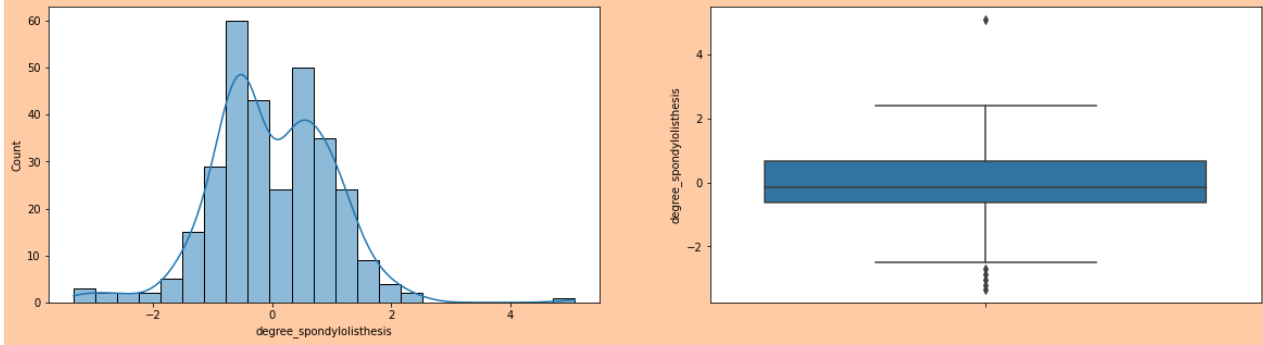
skewness of pelvic_radius column---> 0.04470701353253126

In [128]:

```
1 from sklearn.preprocessing import PowerTransformer
2 scaler = PowerTransformer (method='yeo-johnson')
3 df['degree_spondylolisthesis'] = scaler.fit_transform(df[['degree_spondylolisthesis']].values)
```

In [129]:

```
1 EDA(df, 'degree_spondylolisthesis')
```



skewness of degree_spondylolisthesis column---> 0.10770868217363619

Relationship exploration: Categorical Vs Continuous -- Box Plots

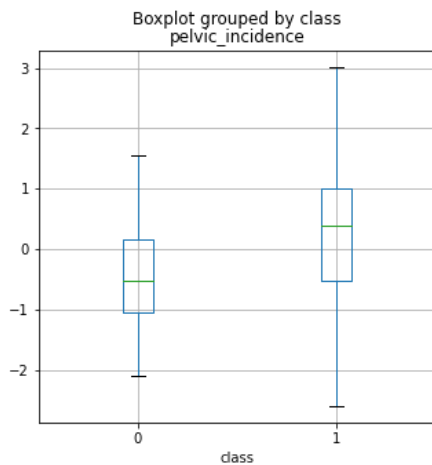
When variable is Continuous and the label/target variable is Categorical we analyze the relation using Boxplots and measure the strength of relation using Anova test

In [141]:

```
1 df.boxplot(column='pelvic_incidence', by='class', figsize=(5,5), vert=True)
```

Out[141]:

<AxesSubplot:title={'center': 'pelvic_incidence'}, xlabel='class'>

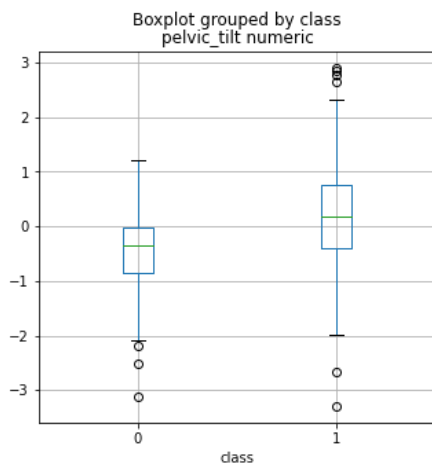


In [142]:

```
1 df.boxplot(column='pelvic_tilt numeric', by='class', figsize=(5,5), vert=True)
```

Out[142]:

<AxesSubplot:title={'center': 'pelvic_tilt numeric'}, xlabel='class'>

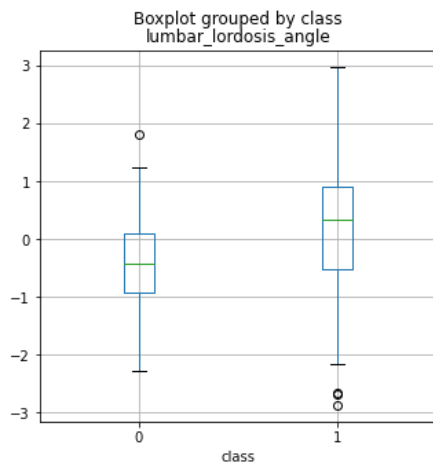


In [143]:

```
1 df.boxplot(column='lumbar_lordosis_angle', by='class', figsize=(5,5), vert=True)
```

Out[143]:

<AxesSubplot:title={'center':'lumbar_lordosis_angle'}, xlabel='class'>

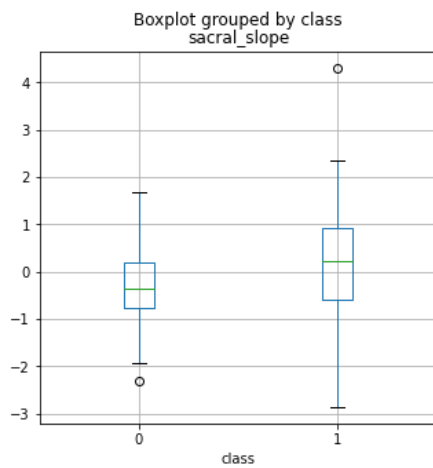


In [144]:

```
1 df.boxplot(column='sacral_slope', by='class', figsize=(5,5), vert=True)
```

Out[144]:

<AxesSubplot:title={'center':'sacral_slope'}, xlabel='class'>

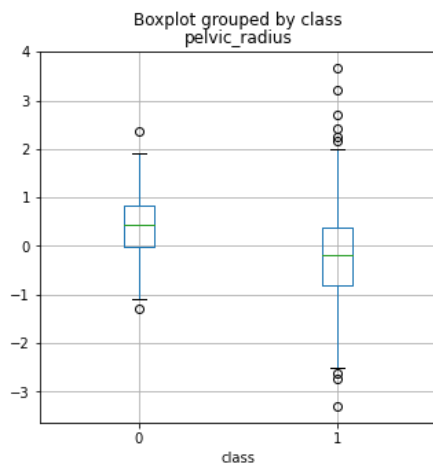


In [147]:

```
1 df.boxplot(column='pelvic_radius', by='class', figsize=(5,5), vert=True)
```

Out[147]:

<AxesSubplot:title={'center':'pelvic_radius'}, xlabel='class'>

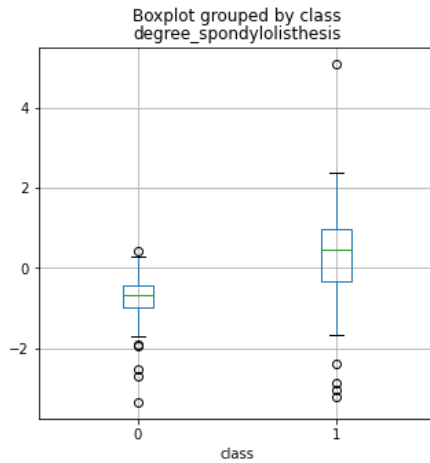


In [140]:

```
1 df.boxplot(column='degree_spondylolisthesis', by='class', figsize=(5,5), vert=True)
```

Out[140]:

```
<AxesSubplot:title={'center':'degree_spondylolisthesis'}, xlabel='class'>
```



Box-Plots interpretation What should you look for in these box plots?

These plots gives an idea about the data distribution of continuous predictor in the Y-axis for each of the category in the X-Axis.

If the distribution looks similar for each category(Boxes are in the same line), that means the the continuous variable has NO effect on the target variable. Hence, the variables are not correlated to each other.

On the other hand if the distribution is different for each category(the boxes are not in same line!). It hints that these variables might be correlated with Rating.

In this data, all three categorical predictors looks correlated with the Target variable.

We confirm this by looking at the results of ANOVA test below

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

Checking the relation with target variable

- relation can be check by using visualization techniques
- statistical tools (Corr, Anova Test, Chi2)

Here we are using anova test because we have categorical vs continuous for this type of scenario ANOVA performance well

Assumption(H0): There is NO relation between the given variables (i.e. The average(mean) values of the numeric variable is same for all the groups in the categorical target variable) ANOVA Test result: Probability of H0 being true

In [149]:

```
1 # get all the continuous columns name
2 num_cols = list(df.select_dtypes(exclude='object').columns)
3 print(f'Continuous columns: {num_cols}')
```

Continuous columns: ['pelvic_incidence', 'pelvic_tilt numeric', 'lumbar_lordosis_angle', 'sacral_slope', 'pelvic_radius', 'degree_spondylolisthesis', 'class']

In [150]:

```
1 categorical_col = ['pelvic_incidence', 'pelvic_tilt numeric', 'lumbar_lordosis_angle', 'sacral_slope', 'pelvic_radius', 'degree_sp
```

In [151]:

```
1 # Defining a function to find the statistical relationship with all the categorical variables
2 def FunctionAnova(inpData, TargetVariable, categorical_col):
3     from scipy.stats import f_oneway
4
5     # Creating an empty List of final selected predictors
6     SelectedPredictors=[]
7
8     print('##### ANOVA Results ##### \n')
9     for predictor in categorical_col:
10         CategoryGroupLists=inpData.groupby(TargetVariable)[predictor].apply(list)
11         AnovaResults = f_oneway(*CategoryGroupLists)
12
13         # If the ANOVA P-Value is <0.05, that means we reject H0
14         if (AnovaResults[1] < 0.05):
15             print(TargetVariable, 'is correlated with', predictor, '| P-Value:', AnovaResults[1])
16             SelectedPredictors.append(predictor)
17         else:
18             print(TargetVariable, 'is NOT correlated with', predictor, '| P-Value:', AnovaResults[1])
19
20     return(SelectedPredictors)
```

In [152]:

```
1 FunctionAnova(inpData=df,
2               TargetVariable='class',
3               categorical_col=categorical_col)
```

ANOVA Results

```
class is correlated with pelvic_incidence | P-Value: 2.97444222457454e-10
class is correlated with pelvic_tilt numeric | P-Value: 7.3276004791601866e-09
class is correlated with lumbar_lordosis_angle | P-Value: 1.3034613241263046e-07
class is correlated with sacral_slope | P-Value: 0.000665309055218086
class is correlated with pelvic_radius | P-Value: 4.166406479179887e-08
class is correlated with degree_spondylolisthesis | P-Value: 1.5372258010635345e-23
```

Out[152]:

```
['pelvic_incidence',
 'pelvic_tilt numeric',
 'lumbar_lordosis_angle',
 'sacral_slope',
 'pelvic_radius',
 'degree_spondylolisthesis']
```

The results of ANOVA confirm that our visual analysis using box plots above.

All continuous variables are correlated with the Target variable.

Now Its time to check multicollinearity in our data set

In [153]:

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 #Finding variance inflation factor in each scaled column i.e X_scaled.shape [1] (1/(1-R2))
3 vif = pd.DataFrame()
4 vif["vif"] = [variance_inflation_factor (df, i) for i in range (df.shape[1])]
5 vif["Features"] = df.columns
6 #Let's check the values
```

In [154]:

1 vif

Out[154]:

	vif	Features
0	106.893700	pelvic_incidence
1	35.838951	pelvic_tilt numeric
2	2.660915	lumbar_lordosis_angle
3	62.474009	sacral_slope
4	1.212631	pelvic_radius
5	1.979380	degree_spondylolisthesis
6	1.143276	class

Note,

1) std value for vif is 5

2) if vif value is > 5 then there possibilities multicollinearity problem

3) std value for vif can be different by project / it depend on data set or project

Conclusion : That means there is multicollinearity problem exist in our data set .

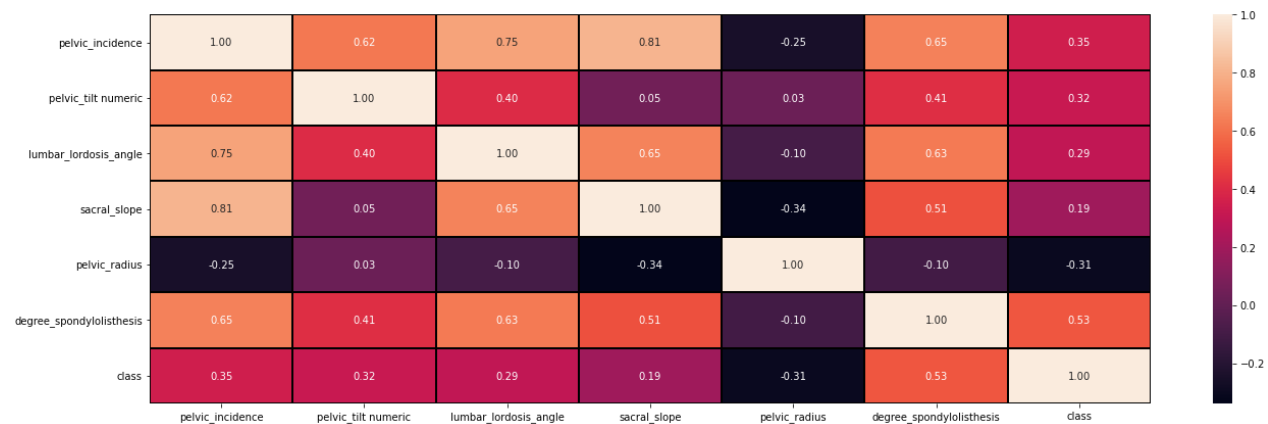
We need to drop correlated function

In [156]:

```
1 plt.figure(figsize=(22,7)) # plotting the heat map
2 sns.heatmap(df.corr(),annot=True,linewidths=0.1,linecolor="black",fmt="0.2f")
```

Out[156]:

<AxesSubplot:>

**with the following function we can select highly correlated features**

In [157]:

```
1 # with the following function we can select highly correlated features
2 # it will remove the first feature that is correlated with anything other feature
3
4 def correlation(dataset, threshold):
5     col_corr = set() # Set of all the names of correlated columns
6     corr_matrix = dataset.corr()
7     for i in range(len(corr_matrix.columns)):
8         for j in range(i):
9             if abs(corr_matrix.iloc[i, j]) > threshold: # we are interested in absolute coeff value
10                 colname = corr_matrix.columns[i] # getting the name of column
11                 col_corr.add(colname)
12     return col_corr
```

In [159]:

```
1 corr_features = correlation(df, 0.7)
2 len(set(corr_features))
```

Out[159]:

2

In [160]:

```
1 corr_features # 80% highly correlated features name
```

Out[160]:

```
{'lumbar_lordosis_angle', 'sacral_slope'}
```

In [161]:

```
1 df.drop (columns = ['sacral_slope'],inplace=True,)
```

In [162]:

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 #Finding variance inflation factor in each scaled column i.e X_scaled.shape [1] (1/(1-R2))
3 vif = pd.DataFrame()
4 vif["vif"]= [variance_inflation_factor (df, i) for i in range (df.shape[1])]
5 vif["Features"] = df. columns
6 #Let's check the values
```

In [163]:

```
1 vif
```

Out[163]:

	vif	Features
0	4.024323	pelvic_incidence
1	1.818716	pelvic_tilt numeric
2	2.634955	lumbar_lordosis_angle
3	1.212520	pelvic_radius
4	1.979290	degree_spondylolisthesis
5	1.140573	class

Now its fine

Now we have finalized the final predictors for ML so we are moving towards model buliding

We are not applying scaling technique as we have already used log transformation

Seprate the features and label

In [167]:

```
1 f = df.drop (columns = ['class'])
2 l = df[ 'class']
```

In [168]:

```
1 log_reg = LogisticRegression ()
```

In [169]:

```
1 maxacc=0
2 maxrn=0
3
4 for i in range(1,100):
5     x_train,x_test,y_train,y_test=train_test_split(f,l,test_size=.30,random_state=i)
6     log_reg.fit(x_train,y_train)
7     pred=log_reg.predict(x_test)
8     score=accuracy_score(pred,y_test)
9     if score>maxacc:
10         maxacc=score
11         maxrn=i
12 print('accuracy_score:-',maxacc,'Random state:-',maxrn)
```

```
accuracy_score:- 0.8817204301075269 Random state:- 55
```

In [170]:

```
1 x_train,x_test,y_train,y_test= train_test_split(f , l, test_size= 0.25, random_state = 55)
```

In [171]:

```
1 log_reg.fit(x_train, y_train)
```

Out[171]:

```
LogisticRegression()
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [172]:

```
1 from sklearn.metrics import classification_report
```

In [173]:

```
1 y_pred = log_reg.predict(x_test)
2 y_pred
```

Out[173]:

```
array([1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1,
       1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0,
       0, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0])
```

In [174]:

```
1 #Model Accuracy
2 accuracy = accuracy_score(y_test,y_pred)
3 accuracy
```

Out[174]:

```
0.8846153846153846
```

In [175]:

```
1 # Confusion Matrix
2 conf_mat = confusion_matrix(y_test,y_pred)
3 conf_mat
```

Out[175]:

```
array([[22,  6],
       [ 3, 47]], dtype=int64)
```

In [176]:

```
1 print (classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.88	0.79	0.83	28
1	0.89	0.94	0.91	50
accuracy			0.88	78
macro avg	0.88	0.86	0.87	78
weighted avg	0.88	0.88	0.88	78

In [180]:

```
1 cross_val_score (log_reg, f, l, cv=7)
```

Out[180]:

```
array([0.55555556, 0.68888889, 0.77272727, 0.90909091, 0.84090909,
       0.84090909, 0.86363636])
```

In [185]:

```
1 cross_val_score (log_reg, f, l, cv=7).mean()
```

Out[185]:

```
0.7816738816738816
```

Lets use Ensemble approach

Bagging and boosting

In [186]:

```
1 from sklearn.ensemble import BaggingClassifier
```

In [187]:

```
1 bag_knn = BaggingClassifier(
2     KNeighborsClassifier (n_neighbors=5), # value of k = 5
3     n_estimators=9, max_samples=0.7, # n_estimators means 9 model you have build
4     bootstrap=True, random_state=3, oob_score=True
5 )
```

In [188]:

```
1 bag_knn.fit (x_train, y_train)
2 bag_knn.score (x_test, y_test)
```

Out[188]:

0.8333333333333334

In [189]:

```
1 from sklearn.ensemble import RandomForestClassifier
```

In [190]:

```
1 # Write one function and call as many as times to check accuracy_score of different models
2 def metric_score (clf, x_train,x_test,y_train,y_test, train=True):
3     if train:
4         y_pred = clf.predict (x_train)
5         print("\n=====Train Result=====")
6         print (f"Accuracy Score: {accuracy_score(y_train, y_pred) * 100:.2f}%")
7     elif train==False:
8         pred = clf.predict(x_test)
9         print("\n=====Test Result=====")
10        print (f"Accuracy Score: {accuracy_score(y_test, pred)* 100:.2f}%")
11        print ('\n \n Test Classification Report \n', classification_report(y_test, pred, digits=2)) ## Model confidence/accuracy
```

In [191]:

```
1 # Initiate Decision Tree Classifier with new parameters and train
2 random_clf = RandomForestClassifier()
3 # Train the model
4 random_clf.fit(x_train,y_train)
```

Out[191]:

RandomForestClassifier()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [192]:

```
1 # Call the function and pass dataset to check train and test score
2 metric_score (random_clf,x_train, x_test, y_train,y_test, train=True) # This is for training socre
3 # This is for testing score
4 metric_score(random_clf,x_train,x_test,y_train,y_test,train=False)
```

```
=====Train Result=====
Accuracy Score: 100.00%
```

```
=====Test Result=====
Accuracy Score: 92.31%
```

Test Classification Report				
	precision	recall	f1-score	support
0	0.87	0.93	0.90	28
1	0.96	0.92	0.94	50
accuracy			0.92	78
macro avg	0.91	0.92	0.92	78
weighted avg	0.93	0.92	0.92	78

this is the one the most problem with random forest maximum time it's tends to overfit becuse it is rule base algorithm

Lets Use Boosting

Within the boosting we have few algorithms we are going to work on that

1) Adaptive boosting 2) Gradient boosting3)XGB

In [197]:

```
1 from sklearn.ensemble import AdaBoostClassifier
```

In [199]:

```
1 ada = AdaBoostClassifier()
```

In [200]:

```
1 x_train,x_test,y_train,y_test= train_test_split(f , 1, test_size= 0.25, random_state = 50)
```

In [201]:

```
1 ada. fit (x_train,y_train)
```

Out[201]:

AdaBoostClassifier()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [202]:

```
1 ada.fit (x_train, y_train)
2 ada.score (x_test, y_test)
```

Out[202]:

0.8076923076923077

In [203]:

```
1 #gradient boosting classifier
2 from sklearn.ensemble import GradientBoostingClassifier # GradientBoostingRegressor If we have regression problem
3 from sklearn.metrics import classification_report , accuracy_score
```

In [208]:

```
1 # initiate GradientBoostingClassifier
2 gbdt_clf = GradientBoostingClassifier()
3
```

In [207]:

```
1 maxacc=0
2 maxrn=0
3
4 for i in range(1,100):
5     x_train,x_test,y_train,y_test=train_test_split(f,1,test_size=.30,random_state=i)
6     gbdt_clf.fit(x_train,y_train)
7     pred=gbdt_clf.predict(x_test)
8     score=accuracy_score(pred,y_test)
9     if score>maxacc:
10         maxacc=score
11         maxrn=i
12 print('accuracy_score:-',maxacc,'Random state:-',maxrn)
```

accuracy_score:- 0.8924731182795699 Random state:- 37

In [209]:

```
1 x_train,x_test,y_train,y_test= train_test_split(f , 1, test_size= 0.25, random_state = 37)
```

In [210]:

```
1 gbdt_clf.fit(x_train, y_train)
```

Out[210]:

GradientBoostingClassifier()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [211]:

```

1 def mertric_Score (clf,XX_train,XX_test,y_train,y_test,train= True):
2     if train:
3         y_pred=clf.predict(XX_train)
4         print('=== Training Score ===')
5         print(f"Accuracy score : {accuracy_score(y_train,y_pred)*100 : 2f} %")
6
7     elif train==False:
8         pred = clf.predict(XX_test)
9         print('=== Testing Score ===')
10        print(f"Accuracy Score : {accuracy_score(y_test,pred)*100 : 2f}%")
11
12        print ('\n \n Classification Report \n' , classification_report(y_test,pred,digits=2))

```

In [217]:

```

1 # call the function
2 mertric_Score(gbdt_clf,x_train,x_test,y_train,y_test,train=True)
3 mertric_Score(gbdt_clf,x_train,x_test,y_train,y_test,train=False)

```

```

=== Training Score ===
Accuracy score : 100.000000 %
=== Testing Score ===
Accuracy Score : 87.179487%

```

Classification Report				
	precision	recall	f1-score	support
0	0.84	0.70	0.76	23
1	0.88	0.95	0.91	55
accuracy			0.87	78
macro avg	0.86	0.82	0.84	78
weighted avg	0.87	0.87	0.87	78

In [218]:

```
1 cross_val_score (gbdt_clf, f, l, cv=5)
```

Out[218]:

```
array([0.56451613, 0.72580645, 0.87096774, 0.90322581, 0.85483871])
```

In [219]:

```
1 cross_val_score (gbdt_clf, f, l, cv=5).mean()
```

Out[219]:

```
0.7870967741935485
```

That is How we have Build the multile models on given data set and performance score has been shown with respect to all model s

Scope for the Improvement

1) we can improv the model accuracy by tunning the hyperparameter of models

2) later on we can plot AUC ROC curve to select the best model (generalized model)

Thank you....!

In []:

```
1
```