

NAME : RETHINAGIRI G

ROLL NO : 225229130

COURSE TITLE : PRACTICAL MACHINE LEARNING LAB

LAB6. Predictive Analytics for Hospitals

Step.1

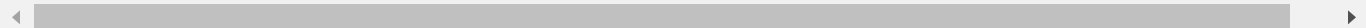
```
In [1]: import pandas as pan
```

```
In [2]: db=pan.read_csv("Diabetes.csv")
```

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

```
Out[3]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outco
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	



```
In [4]: db.size
```

```
Out[4]: 6912
```

```
In [5]: db.shape
```

```
Out[5]: (768, 9)
```

```
In [6]: db.describe()
```

```
Out[6]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunc
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.00
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.47
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.33
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.07
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.24
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.37
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.62
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.42

```
In [7]: db.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 [8]: type(db)
```

```
Out[8]: pandas.core.frame.DataFrame
```

```
In [9]: db.columns
```

```
Out[9]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
               'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
              dtype='object')
```

```
In [10]: db["Glucose"].value_counts
```

```
Out[10]: <bound method IndexOpsMixin.value_counts of 0      148
1         85
2        183
3         89
4        137
...
763      101
764      122
765      121
766      126
767        93
Name: Glucose, Length: 768, dtype: int64>
```

```
In [11]: db.count
```

```
Out[11]: <bound method DataFrame.count of
s  Insulin  BMI  \
0         6   148    72    35    0  33.6
1         1    85    66    29    0  26.6
2         8   183    64     0    0  23.3
3         1    89    66    23   94  28.1
4         0   137   40    35  168  43.1
..      ...   ...   ...   ...   ...   ...
763      10   101    76    48  180  32.9
764         2   122    70    27    0  36.8
765         5   121    72    23  112  26.2
766         1   126    60     0    0  30.1
767         1    93    70    31    0  30.4

DiabetesPedigreeFunction  Age  Outcome
0          0.627    50         1
1          0.351    31         0
2          0.672    32         1
3          0.167    21         0
4          2.288    33         1
..      ...   ...   ...
763        0.171    63         0
764        0.340    27         0
765        0.245    30         0
766        0.349    47         1
767        0.315    23         0

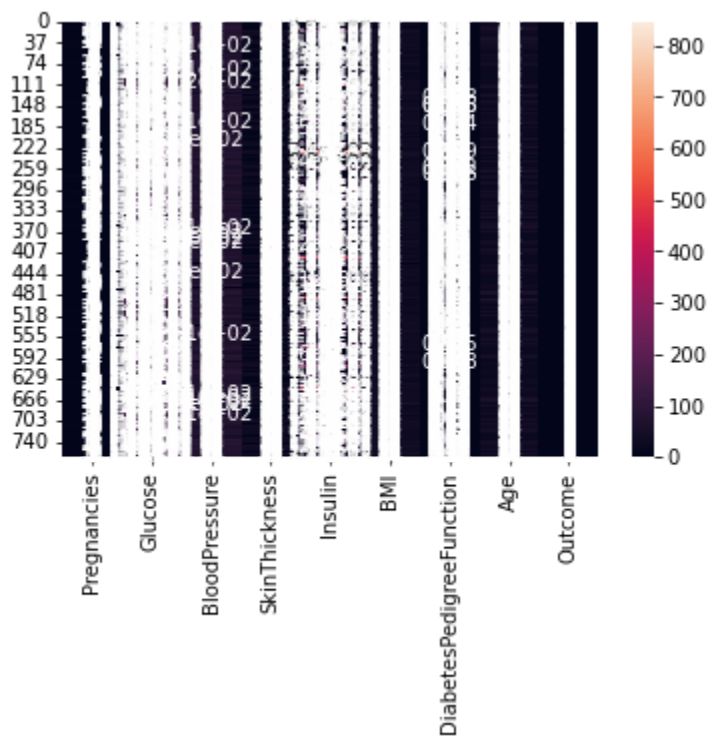
[768 rows x 9 columns]>
```

Step.2 [Identify relationships between feature]

```
In [12]: import seaborn as sb
```

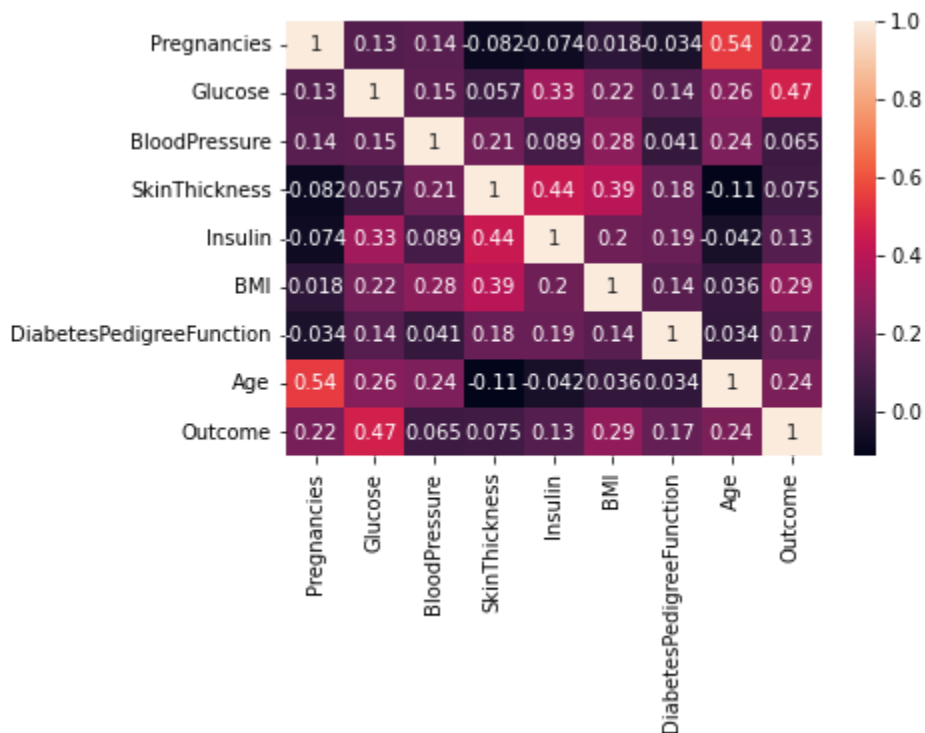
```
In [13]: sb.heatmap(db,annot=True)
```

```
Out[13]: <AxesSubplot:>
```



```
In [14]: c=db.corr()  
sb.heatmap(c,annot=True)
```

```
Out[14]: <AxesSubplot:>
```



Step3. [predicton using one feature]

```
In [15]: X=db[['Age']]
```

```
In [16]: y=db[["Outcome"]]
```

```
In [17]: from sklearn.model_selection import train_test_split as tts
X_train,X_test,y_train,y_test=tts(X,y,test_size=.25,random_state=42)
```

```
In [18]: from sklearn.linear_model import LogisticRegression as LOR
reg=LOR()
reg.fit(X_train,y_train)
```

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
```

Out[18]: LogisticRegression()

```
In [19]: reg.predict(y_test)
```

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\base.py:493: FutureWarning: The feature names should match those that were passed during fit. Starting version 1.2, an error will be raised.
```

```
Feature names unseen at fit time:
```

- Outcome

Feature names seen at fit time, yet now missing:

- Age

```
warnings.warn(message, FutureWarning)
```

[illegible]

```
In [20]: C=reg.coef_
```

In [21]: C

```
Out[21]: array([[0.05221912]])
```

```
In [22]: I=reg.intercept_
```

In [23]:

```
I
```

Out[23]: array([-2.39506398])

In [24]:

```
lrf=C*60+I  
  
from scipy.special import expit
```

In [25]:

```
if expit(lrf)>0.5:  
    print(expit(lrf))  
    print("HE WILL BECOME DIABETIC")  
else:  
    print("No")
```

```
[[0.67657656]]  
HE WILL BECOME DIABETIC
```

Step.4 [Prediction using many features]

In [26]:

```
X1=db[['Glucose','BMI','Age']]
```

In [27]:

```
y1=db[["Outcome"]]
```

In [28]:

```
X1_train,X1_test,y1_train,y1_test=tts(X1,y1,test_size=.25,random_state=42)
```

In [29]:

```
reg1=LOR()
```

In [30]:

```
reg1.fit(X1_train,y1_train)
```

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().  
y = column_or_1d(y, warn=True)
```

Out[30]: LogisticRegression()

In [31]:

```
print(reg1.coef_)
```

```
[[0.03326879 0.09717039 0.04404934]]
```

In [32]:

```
reg1.intercept_
```

Out[32]: array([-9.47396587])

```
In [33]: m=reg1.predict([[150,30,40]])
```

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but LogisticRegression was fitted with feature names
  warnings.warn(
```

```
In [34]: expit(m)
```

```
Out[34]: array([0.73105858])
```

```
In [35]: pr=reg1.predict_proba([[150,30,40]])
```

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but LogisticRegression was fitted with feature names
  warnings.warn(
```

```
In [36]: expit(pr)
```

```
Out[36]: array([[0.61106495, 0.63371999]])
```

Step.5 [BuildLoR model with all features]

```
In [37]: X2=db.drop('Outcome',axis=1)
y2=db[['Outcome']]
```

```
In [38]: reg2=LOR()
```

```
In [39]: X2_train,X2_test,y2_train,y2_test=tts(X2,y2,
                                                test_size=.25,random_state=42)
```

```
In [40]: reg2.fit(X2_train,y2_train)
```

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
```

```
C:\Users\user\anaconda3\lib\site-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> (<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 (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
n_iter_i = _check_optimize_result(
```

```
Out[40]: LogisticRegression()
```

```
In [41]: y2_pred=reg2.predict(X2_test)
```

```
In [42]: from sklearn.metrics import *
```

```
In [43]: AUC_value=roc_auc_score(y2_test,y2_pred)
AUC_value
```

```
Out[43]: 0.7122658183103571
```

Step.6 [Forward Selection Procedure]

```
In [44]: from sklearn.feature_selection import f_regression
```

```
In [45]: selected_features = []
best_score = 0
```

```
In [46]: for feature in db.columns:
    X_subset = db[selected_features + [feature]]
    model = LOR().fit(X_subset, y)
    f_score = f_regression(X_subset, y)[0][-1]
    if f_score > best_score:
        selected_features.append(feature)
        best_score = f_score
```

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
In [47]: print(selected_features)
```

```
['Pregnancies', 'Glucose', 'Outcome']
```



```
In [48]: selected_features
```

```
Out[48]: ['Pregnancies', 'Glucose', 'Outcome']
```

```
In [49]: def get_auc(var,tar,df):
          fx = df[var]
          fy = df[tar]
          reg4=LOR()
          reg4.fit(fx,fy)
          pred=reg4.predict_proba(fx)[:,-1]
          auc_val = roc_auc_score(y,pred)
          return auc_val
get_auc(['Glucose','BMI'], ['Outcome'],db)
```

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
y = column_or_1d(y, warn=True)

```
Out[49]: 0.8109328358208956
```

```
In [50]: get_auc(['Pregnancies','BloodPressure','SkinThickness'],
                 ['Outcome'],db)
```

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
y = column_or_1d(y, warn=True)

```
Out[50]: 0.6444962686567164
```

```
In [51]: def best_next(current,cand,tar,df):
          best_auc=-1
          best_var=None
          for i in cand:
              auc_v = get_auc(current+[i],tar,df)
              if auc_v>=best_auc:
                  best_auc=auc_v
                  best_var=i
          return best_var
```

```
In [52]: current=['Insulin','BMI','DiabetesPedigreeFunction','Age']
          cand=['Pregnancies','Glucose','BloodPressure','SkinThickness']
          tar=['Outcome']
          next_var = best_next(current,cand,tar,db)
          next_var
```

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
y = column_or_1d(y, warn=True)

```
Out[52]: 'Pregnancies'
```

```
In [53]: tar = ['Outcome']
current=[]
cand=['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesP
max_num = 7
num_it = min(max_num, len(cand))
for i in range(0, num_it):
    next_var = best_next(current, cand, tar, db)
    current += [next_var]
    cand.remove(next_var)
    print("variable addd in step "+str(i+1)+" is "+ next_var + " .")
```

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

variable addd in step 1 is Pregnancies .

variable addd in step 2 is Glucose .

variable addd in step 3 is BloodPressure .

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

variable addd in step 4 is SkinThickness .

variable addd in step 5 is Insulin .

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

variable addd in step 6 is BMI .

variable addd in step 7 is DiabetesPedigreeFunction .

```
C:\Users\user\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:814: Conver  
genceWarning: 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> (<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 (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
n_iter_i = _check_optimize_result(
```

```
In [54]: current
```

```
Out[54]: ['Pregnancies',  
          'Glucose',  
          'BloodPressure',  
          'SkinThickness',  
          'Insulin',  
          'BMI',  
          'DiabetesPedigreeFunction']
```

```
In [ ]:
```

STEP 7 [PLOT LINE GRAPH OF AUC VALUES SELECT CUT-OFF]

```
In [55]: X2_train,X2_test,y2_train,y2_test = tts(X2,y,stratify=y,test_size=.5,random_state=42)  
         prediction=reg2.predict_proba(X2_test)  
         train = pan.concat([X2_train,y2_train],axis =1)  
         test = pan.concat([X2_test,y2_test],axis =1)
```

```
In [56]: def auc_train_test (variables,target, train, test):
X_train = train[variables]
X_test = test[variables]
Y_train =train[target]
Y_test = test[target]
Lor=LOR()
Lor.fit(X_train,Y_train)
prediction_train = Lor.predict_proba(X_train)[: ,1]
prediction_test = Lor.predict_proba(X_test)[: ,1]
auc_train = roc_auc_score(Y_train, prediction_train)
auc_test = roc_auc_score(Y_train,prediction_test)
return (auc_train,auc_test)
auc_values_train=[]
auc_values_test=[]
variable_evaluate=[]
for v in X2.columns:
    variable_evaluate.append(v)
    auc_train,auc_test = auc_train_test(variable_evaluate,['Outcome'],train,test)
    auc_values_train.append(auc_train)
    auc_values_test.append(auc_test)
```

```

C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
C:\Users\user\anaconda3\lib\site-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:

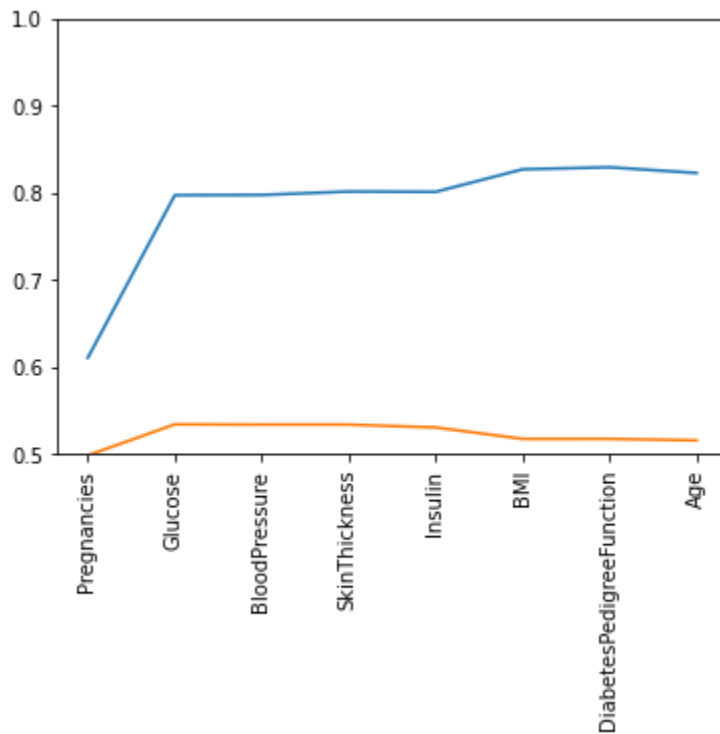
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
n_iter_i = _check_optimize_result(
```

```
In [57]: import numpy as num
import matplotlib.pyplot as mat
x = num.array(range(0, len(auc_values_train)))
my_train = num.array(auc_values_train)
my_test = num.array(auc_values_test)
mat.xticks(x, X2.columns, rotation=90)
mat.plot(x, my_train)
mat.plot(x, my_test)
mat.ylim(0.5, 1)
mat.show()
```

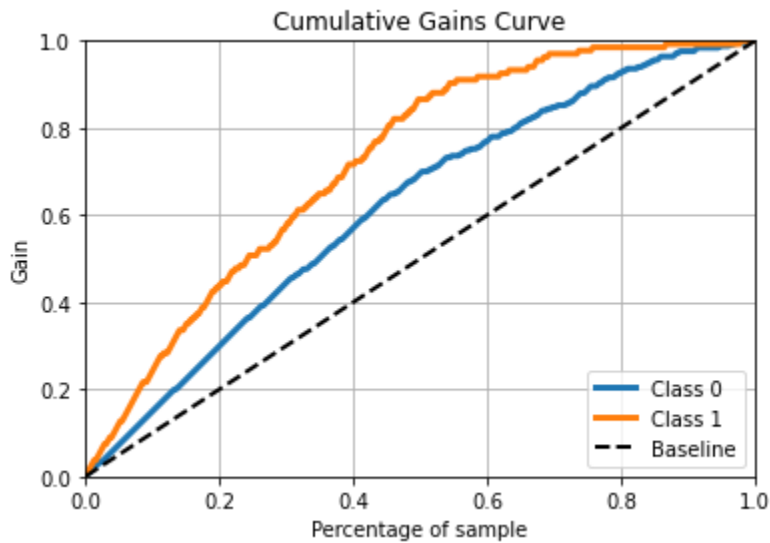


STEP 8 : [DRAW CUMILATIVE GAIN CHART AND LIFT CHART]

```
In [58]: !pip install scikit-plot
```

```
Requirement already satisfied: scikit-plot in c:\users\user\anaconda3\lib\site-packages (0.3.7)
Requirement already satisfied: joblib>=0.10 in c:\users\user\anaconda3\lib\site-packages (from scikit-plot) (1.1.0)
Requirement already satisfied: matplotlib>=1.4.0 in c:\users\user\anaconda3\lib\site-packages (from scikit-plot) (3.5.1)
Requirement already satisfied: scikit-learn>=0.18 in c:\users\user\anaconda3\lib\site-packages (from scikit-plot) (1.0.2)
Requirement already satisfied: scipy>=0.9 in c:\users\user\anaconda3\lib\site-packages (from scikit-plot) (1.7.3)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (3.0.4)
Requirement already satisfied: packaging>=20.0 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (21.3)
Requirement already satisfied: numpy>=1.17 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (1.21.5)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (2.8.2)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (4.25.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (1.3.2)
Requirement already satisfied: cyclor>=0.10 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (0.11.0)
Requirement already satisfied: pillow>=6.2.0 in c:\users\user\anaconda3\lib\site-packages (from matplotlib>=1.4.0->scikit-plot) (9.0.1)
Requirement already satisfied: six>=1.5 in c:\users\user\anaconda3\lib\site-packages (from python-dateutil>=2.7->matplotlib>=1.4.0->scikit-plot) (1.16.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\user\anaconda3\lib\site-packages (from scikit-learn>=0.18->scikit-plot) (2.2.0)
```

```
In [59]: import scikitplot as skplt
skplt.metrics.plot_cumulative_gain(y2_test,prediction)
mat.show()
mat.figure(figsize=(7,7))
skplt.metrics.plot_lift_curve(y2_test, prediction)
mat.show()
```



<Figure size 504x504 with 0 Axes>

