

SURUTHI S

225229141

PML LAB 6

```
In [2]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
```

```
In [3]: 1 df = pd.read_csv("diabetes.csv")
```

```
In [4]: 1 df.head()
```

Out[4]:

	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

```
In [5]: 1 df.sample(6)
```

Out[5]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
50	1	103	80	11	82	19.4	0.491	22	0
469	6	154	78	41	140	46.1	0.571	27	0
260	3	191	68	15	130	30.9	0.299	34	0
427	1	181	64	30	180	34.1	0.328	38	1
742	1	109	58	18	116	28.5	0.219	22	0
754	8	154	78	32	0	32.4	0.443	45	1

```
In [6]: 1 df.shape
```

Out[6]: (768, 9)

```
In [7]: 1 df.columns
```

Out[7]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'], dtype='object')

```
In [8]: 1 df.dtypes
```

Out[8]: Pregnancies int64  
Glucose int64  
BloodPressure int64  
SkinThickness int64  
Insulin int64  
BMI float64  
DiabetesPedigreeFunction float64  
Age int64  
Outcome int64  
dtype: object

```
In [9]: 1 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 [10]: 1 df.value_counts()

Out[10]: Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin  BMI  DiabetesPedigreeFunction  Age  Outcome
0          57          60             0             0             0      21.7    0.735             67    0             1
          67          76             0             0             0      45.3    0.194             46    0             1
          103         108             37             0             0      39.2    0.305             65    0             1
          104          74             0             0             0      28.8    0.153             48    0             1
          105          72             29            325            36.9    0.159             28    0             1
          ..
          84          50             23             76            30.4    0.968             21    0             1
          85          65             0             0            39.6    0.930             27    0             1
          87           0             23             0            28.9    0.773             25    0             1
          ..
          ..          58             16             52            32.7    0.166             25    0             1
          17         163          72             41            114            40.9    0.817             47    1             1
Length: 768, dtype: int64
```

```
In [11]: 1 df.isnull().sum()

Out[11]: Pregnancies      0
Glucose      0
BloodPressure  0
SkinThickness  0
Insulin      0
BMI          0
DiabetesPedigreeFunction  0
Age          0
Outcome      0
dtype: int64
```

## STEP 2

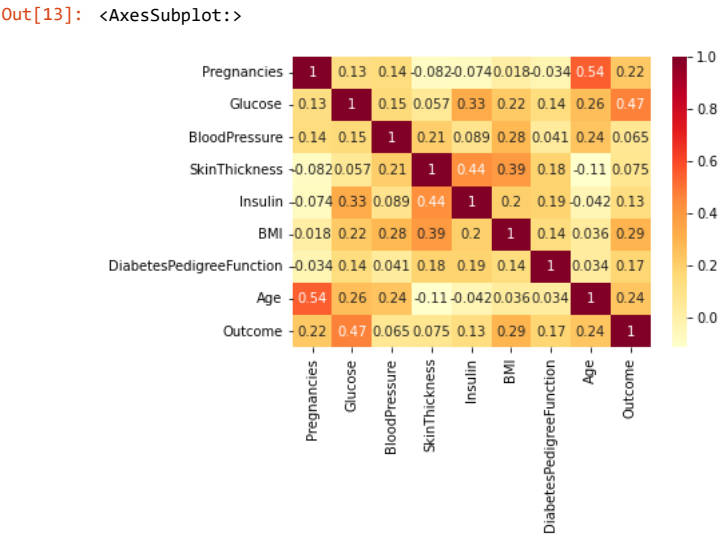
### IDENTIFYING RELATIONSHIPS

```
In [12]: 1 import seaborn as sns
2 df.corr()
```

Out[12]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
Pregnancies	1.000000	0.129459	0.141282	-0.081672	-0.073535	0.017683	-0.033523	0.544341	0.221898
Glucose	0.129459	1.000000	0.152590	0.057328	0.331357	0.221071	0.137337	0.263514	0.466581
BloodPressure	0.141282	0.152590	1.000000	0.207371	0.088933	0.281805	0.041265	0.239528	0.065068
SkinThickness	-0.081672	0.057328	0.207371	1.000000	0.436783	0.392573	0.183928	-0.113970	0.074752
Insulin	-0.073535	0.331357	0.088933	0.436783	1.000000	0.197859	0.185071	-0.042163	0.130548
BMI	0.017683	0.221071	0.281805	0.392573	0.197859	1.000000	0.140647	0.036242	0.292695
DiabetesPedigreeFunction	-0.033523	0.137337	0.041265	0.183928	0.185071	0.140647	1.000000	0.033561	0.173844
Age	0.544341	0.263514	0.239528	-0.113970	-0.042163	0.036242	0.033561	1.000000	0.238356
Outcome	0.221898	0.466581	0.065068	0.074752	0.130548	0.292695	0.173844	0.238356	1.000000

```
In [13]: 1 sns.heatmap(df.corr(),annot=True,cmap='YlOrRd')
```



## STEP 3 -PREDICTION USING ONE FEATURE

```
In [14]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import classification_report
```

```
In [15]: 1 df.Outcome.value_counts()
```

```
Out[15]: 0    500
1     268
Name: Outcome, dtype: int64
```

```
In [16]: 1 # stratified shuffle split
2
3 from sklearn.model_selection import StratifiedShuffleSplit
4
5 sss = StratifiedShuffleSplit(n_splits=4,test_size = 0.30,random_state = 42)
6
7
```

```
In [17]: 1 features = df[['Age']]
2 labels = df[['Outcome']]
3
4 for train_index,test_index in sss.split(features,labels):
5     x_train,x_test = features.iloc[train_index],features.iloc[test_index]
6     y_train,y_test = labels.iloc[train_index],labels.iloc[test_index]
```

```
In [18]: 1 x_train.shape
```

Out[18]: (537, 1)

```
In [19]: 1 x_test.shape
```

Out[19]: (231, 1)

```
In [20]: 1 logreg = LogisticRegression()
2 logreg.fit(x_train,y_train)
3 y_pred = logreg.predict(x_test)
4 print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.67	0.93	0.78	150
1	0.52	0.15	0.23	81
accuracy			0.65	231
macro avg	0.60	0.54	0.50	231
weighted avg	0.62	0.65	0.59	231

C:\Users\SURUTHI S\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)

```
In [21]: 1 from sklearn.metrics import accuracy_score
2
3 accuracy_score(y_test,y_pred)
```

Out[21]: 0.6536796536796536

```
In [22]: 1 logreg.coef_
```

Out[22]: array([[0.03700533]])

```
In [23]: 1 logreg.intercept_
```

Out[23]: array([-1.87818016])

```
In [24]: 1 # age = 60 , outcome = ?
2 age = 60
3 lrf = logreg.coef_*age + logreg.intercept_
```

```
In [25]: 1 lrf
```

Out[25]: array([[0.34213957]])

```
In [26]: 1 from scipy.special import expit
2
3 expit(lrf)
4
5 #Expit (a.k.a. Logistic sigmoid) ufunc for ndarrays.
6
7 # The expit function, also known as the logistic sigmoid function,
8 # is defined as expit(x) = 1/(1+exp(-x)). It is the inverse of the Logit function.
```

Out[26]: array([[0.58471016]])

the output is greater than 0.5 yes the person is diabetic

## STEP 4 PREDICTION USING MANY FEATURES

```
In [27]: 1 feat = df[['Glucose','BMI','Age']]
2
3
4 for train_index,test_index in sss.split(feat,labels):
5     x_train_ , x_test_ = feat.iloc[train_index],feat.iloc[test_index]
6     y_train_ , y_test_ = labels.iloc[train_index],labels.iloc[test_index]
```

```
In [28]: 1 x_train_.shape
```

Out[28]: (537, 3)

```
In [29]: 1 x_test_.shape

Out[29]: (231, 3)
```

```
In [30]: 1 logreg_ = LogisticRegression()
2
3 logreg_.fit(x_train_,y_train_)
4 y_pred_ = logreg_.predict(x_test_)
5
6 accuracy_score(y_test_,y_pred_)

C:\Users\SURUTHI S\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]: 0.7792207792207793
```

```
In [31]: 1 glucose = 150
2 bmi = 30
3 age = 40
4
5 w=[glucose,bmi,age]
```

```
In [32]: 1 logreg_.coef_

Out[32]: array([[0.03884178, 0.08275889, 0.02420389]])
```

```
In [33]: 1 logreg_.intercept_

Out[33]: array([-8.99044306])
```

```
In [34]: 1 lrf_ = np.dot(logreg_.coef_,w) + logreg_.intercept_
```

```
In [35]: 1 lrf_

Out[35]: array([0.28674638])
```

```
In [36]: 1 from scipy.special import expit
2
3 expit(lrf_)

Out[36]: array([0.57119941])
```

the output is greater than 0.5 yes the person is diabetic

```
In [37]: 1 logreg_.predict_proba([w])
2
3 # the outcome as 1 class having large probability

C:\Users\SURUTHI S\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(

Out[37]: array([[0.42880059, 0.57119941]])
```

## STEP 5 - BUILD LOR WITH ALL FEATURES

```
In [38]: 1 allFeatures = df.drop('Outcome',axis=1)
2
```

```
In [39]: 1 sss = StratifiedShuffleSplit(n_splits = 4 ,test_size = 0.30,random_state =42)
2
3 for train_index,test_index in sss.split(allFeatures,labels):
4     x_train_all ,x_test_all = allFeatures.iloc[train_index],allFeatures.iloc[test_index]
5     y_train_all , y_test_all = labels.iloc[train_index],labels.iloc[test_index]
6
```

```
In [40]: 1 x_train_all.shape

Out[40]: (537, 8)
```

```
In [41]: 1 x_test_all.shape

Out[41]: (231, 8)
```

```
In [42]: 1 logreg_all = LogisticRegression()
2         logreg_all.fit(x_train_all,y_train_all)
3         y_pred_all = logreg_all.predict(x_test_all)
4
5         accuracy_score(y_test_all,y_pred_all)
```

C:\Users\SURUTHI S\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\SURUTHI S\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) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

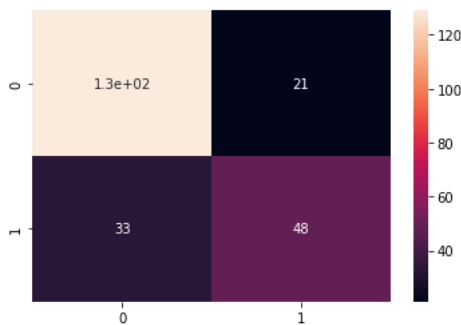
```
n_iter_i = _check_optimize_result(
```

Out[42]: 0.7662337662337663

```
In [43]: 1 from sklearn.metrics import confusion_matrix
2         cm = confusion_matrix(y_test_all,y_pred_all)
```

```
In [44]: 1 sns.heatmap(cm,annot=True)
```

Out[44]: <AxesSubplot:>



```
In [45]: 1 from sklearn.metrics import roc_auc_score
2
3         y_pred_proba = logreg_all.predict_proba(x_test_all)[: ,1]
4         auc = roc_auc_score(y_test_all,y_pred_proba)
5         auc
```

Out[45]: 0.7995884773662552

<https://quantifyinghealth.com/stepwise-selection/> (<https://quantifyinghealth.com/stepwise-selection/>)

## STEP 6

```
In [46]: 1 import warnings
2         warnings.filterwarnings(action='ignore')
```

```
In [47]: 1 import mlxtend
```

```
In [48]: 1 from mlxtend.feature_selection import SequentialFeatureSelector as SFS
```

```
In [49]: 1 sfs = SFS(logreg_all,
2             k_features=4,
3             scoring='accuracy')
```

```
In [50]: 1 sfs.fit(x_train_all,y_train_all)
```

Out[50]: SequentialFeatureSelector(estimator=LogisticRegression(), k\_features=(4, 4), scoring='accuracy')

```
In [51]: 1 sfs.k_feature_names_
```

Out[51]: ('Glucose', 'BloodPressure', 'BMI', 'DiabetesPedigreeFunction')

```
In [52]: 1 def get_next(i,model):
2         sfs = SFS(model,k_features=i,scoring='accuracy')
3         sfs.fit(x_train_all,y_train_all)
4         return sfs.k_feature_names_
5
```

```
In [53]: 1 columns = x_train_all.shape[1]
```

```
In [54]: 1 FSP = get_next(columns-1,logreg_all)
```

```

In [55]: 1 FSP

Out[55]: ('Glucose',
          'BloodPressure',
          'SkinThickness',
          'Insulin',
          'BMI',
          'DiabetesPedigreeFunction',
          'Age')

In [56]: 1 a = FSP[0][1]
          2 a

Out[56]: '1'

In [57]: 1 def get_auc(cols):
          2     lr = LogisticRegression()
          3     x_train = x_train_all[cols]
          4     y_train = y_train_all
          5
          6     x_test = x_test_all[cols]
          7     y_test = y_test_all
          8
          9     lr.fit(x_train,y_train)
         10     ypred = lr.predict(x_test)
         11
         12     y_pred_proba = lr.predict_proba(x_test)[: ,1]
         13
         14     auc = roc_auc_score(y_test_all,y_pred_proba)
         15
         16     return auc
         17
         18

In [58]: 1 get_auc(['BloodPressure'])

Out[58]: 0.6496296296296296

In [59]: 1 get_auc(['BMI'])

Out[59]: 0.7565020576131688

1 http://rasbt.github.io/mlxtend/
2
3 https://www.codespeedy.com/sequential-forward-selection-with-python-and-scikit-learn/

In [ ]: 1

```

## step:7 [plot Line graph of AUC values and select cut-off]

```

In [32]: X2_train,X2_test,y2_train,y2_test=train_test_split(x_train_all,y_train_all,stratify=y_train_all,test_size=0.5,random_state=42)

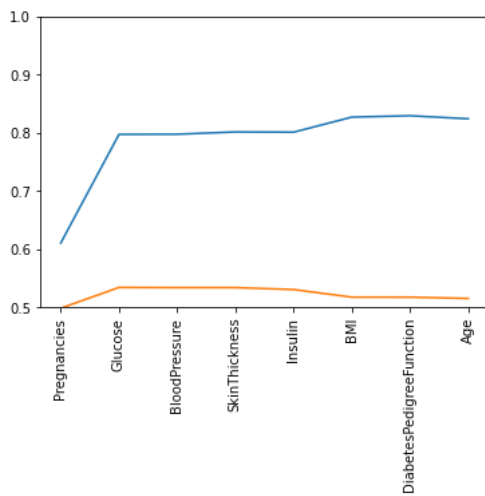
In [33]: 1 prediction = lr2.predict_proba(x_test_all)

In [34]: 1 train = pd.concat([x_train_all, y_train_all], axis =1)
          2 test = pd.concat([x_test_all, y_test_all], axis =1)
          3 def auc_train_test (variables, target, train, test):
          4     X_train = train[variables]
          5     X_test = test[variables]
          6     Y_train = train[target]
          7     Y_test = test[target]
          8     lr3 = LogisticRegression()
          9     lr3.fit(X_train, Y_train)
         10     predictions_train = lr3.predict_proba(X_train)[: ,1]
         11     predictions_test = lr3.predict_proba(X_test)[: ,1]
         12     auc_train = roc_auc_score(Y_train, predictions_train)
         13     auc_test = roc_auc_score(Y_train, predictions_test)
         14     return (auc_train, auc_test)
         15 auc_values_train =[]
         16 auc_values_test =[]
         17 variables_evaluate=[]
         18
         19 for v in X2.columns:
         20     variables_evaluate.append(v)
         21     auc_train, auc_test = auc_train_test(variables_evaluate, ['Outcome'],train,test)
         22     auc_values_train.append(auc_train)
         23     auc_values_test.append(auc_test)

```

...

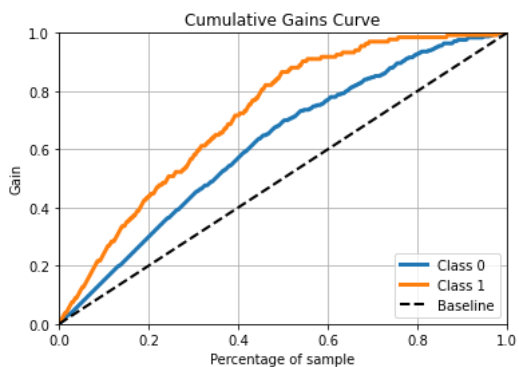
```
In [35]: 1 import matplotlib.pyplot as plt
2 import numpy as np
3 x = np.array(range(0,len(auc_values_train)))
4 my_train = np.array(auc_values_train)
5 my_test = np.array(auc_values_test)
6 plt.xticks(x, X2.columns, rotation =90)
7
8 plt.plot(x, my_train)
9 plt.plot(x, my_test)
10 plt.ylim(0.5, 1)
11 plt.show()
```



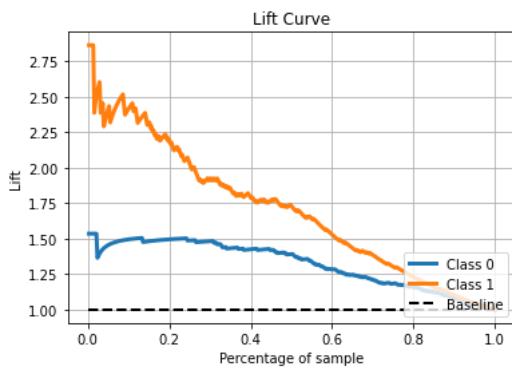
## step:8 [Draw cumulative gain chart and lift chart]

```
In [36]: 1 !pip install scikit-plot
2 from scikitplot.estimators import plot_feature_importances
3 from scikitplot.metrics import plot_confusion_matrix, plot_roc
```

```
In [37]: 1 import scikitplot as skplt
2 skplt.metrics.plot_cumulative_gain(y2_test, prediction)
3 plt.show()
4 plt.figure(figsize=(7,7))
5 skplt.metrics.plot_lift_curve(y2_test, prediction)
6 plt.show()
```



<Figure size 504x504 with 0 Axes>



```
In [ ]: 1
```