

Lab: 6

```
import pandas as pd
data = pd.read_csv('diabetes.csv')
data.head()

data.shape
data.columns
data.info
data.value_counts()
```

```
import seaborn as sns
import matplotlib.pyplot as plt
plt.figure(figsize=(10,10))
sns.heatmap(data.head(20), cmap='Reds', annot=True,
            linewidth=.5)
```

```
X = data[['Age']]
Y = data['outcome']
```

```
from sklearn.linear_model import LogisticRegression
```

```
|rm| = LogisticRegression()
```

```
|rm|.fit(X, Y)
```

```
|rm|.coef_
```

```
|rm|.intercept_
```

```
years-old = [[60]]
```

```
|rm|.predict(years-old)
```

```
rb = |rm|.coef_ * 60 + |rm|.intercept_
```

```
from scipy.special import expit
```

## Lab6. Predictive Analytics for Hospitals

### Objectives

In this lab, you will continue to work on Diabetes data from Lab4 and apply Predictive Analytics principles to recommend the course of action to be taken by the hospital.

### Learning Outcomes

After completing this lab, you will be able to

- Perform prediction on single and multiple independent variables for the target
- Apply Forward Selection procedure to find the best features based on AUC scores
- Draw a line plot of AUC scores for variables and select them using cut-off
- Plot Gain curves and Life curves and interpret results

### Business Use Case

As a Data Scientist, you are requested by the hospital to build a predictive model for them. The hospital will apply your predictive model to perform prediction on new patients whether they will become diabetic or not. Further, based on the prediction, the hospital will recommend action plan such as developing a diet plan, physical activities and others.

#### Step1. [Import dataset]

- Using Pandas, import "diabetes.csv" file and print properties such as head, shape, columns, dtype, info and value\_counts.

#### Step2. [Identify relationships between feature]

- Create a Heatmap for the dataset and understand the data

#### Step3. [Prediction using one feature]

*Will older people become diabetic?*

- Create LogisticRegression model, train with "Age" as X and "Outcome" feature as y.
- Print model parameter values: coef\_ and intercept\_
- **Query:** A person is 60 years old. Will he be diabetic?
- Use model parameters and find function value. Your code will be as below.

```
lrf = logreg.coef_ * 60 + logreg.intercept_  
from scipy.special import expit  
expit(lrf)
```

If your output > 0.5, YES he will become diabetic. Otherwise, NO, he will not be diabetic.

#### Step4. [Prediction using many features]

*Will Glucose, BMI and Age values make someone diabetic?*

- Select the three features 'Glucose', 'BMI' and 'Age' from your dataset, call it as X
- Create a new LogisticRegression model, train with X and 'Outcome' as y.
- Query: For a person, glucose=150, bmi=30, age=40. Will he be diabetic?
- Find the value of expit() as before. Output will be: 0.5208271643241003

Note: You can also verify expit() output value as below  
`logreg.predict_proba([[150, 30, 40]])`

#### Step5. [Build LoR model with all features]

- Create LoR model, train it with X\_train and y\_train values
- Now, compute and print its AUC value
- Can we get this AUC value with limited set of good features?. Yes, we are going to find with 'Forward Selection Procedure'.

```
if expit (lrf) > 0.5:
```

```
    print ("YES, he will become diabetic")  
else:
```

```
    print ("No, he will not be diabetic")
```

```
X1 = data[["Age", "BMI", "Glucose"]]
```

```
lm2 = LogisticRegression()
```

```
lm2.fit(X1, y)
```

```
lm2.predict([[150, 30, 40]])
```

```
lm2.predict_proba([[150, 30, 40]])
```

```
import warnings.
```

```
warnings.filterwarnings("ignore")
```

```
X3 = data.drop("Outcome", axis=1)
```

```
lm3 = LogisticRegression()
```

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, Y_train, Y_test = train_test_split
```

```
(X3, Y, train_size=0.8, test_size=0.2)
```

```
lm3.fit(X_train, Y_train)
```

```
y_pred = lm3.predict(X_test)
```

```
from sklearn.metrics import roc_auc_score.
```

```
print("AUC", roc_auc_score(Y_test, y_pred))
```

```
type(data.columns)
```

```
def get_auc(row, tan, df):
```

**Step6. [Forward Selection Procedure]**

Now, you have to find and select a good set of features with the best AUC score. The algorithm is

**Forward stepwise variable selection procedure**

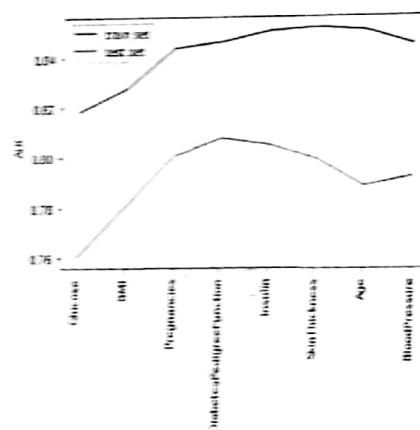
1. Empty set
2. Find best variable  $v$
3. Find best variable  $w$  in combination with  $v$
4. Find best variable  $x$  in combination with  $v, w$
5. ...
6. (Until all variables are added or until predefined number of variables is added)

**Implementation Steps of the forward stepwise procedure**

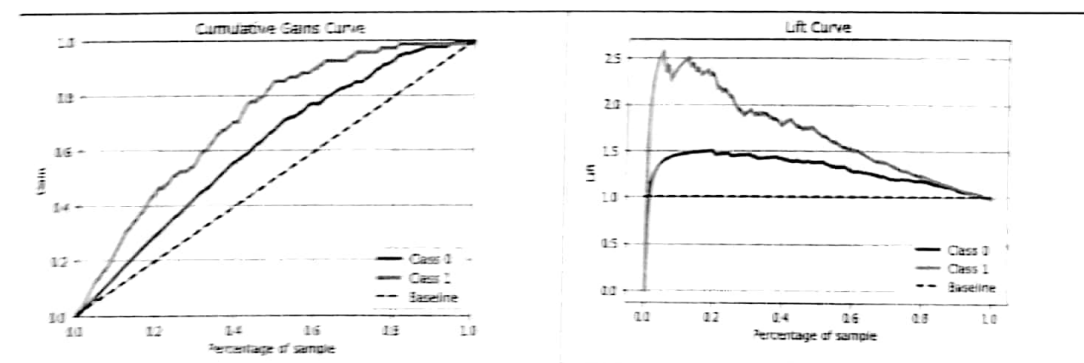
1. Define a function `get_auc()` that calculates AUC given a certain set of variables
2. Define a function `best_next()` that returns next best variable in combination with current variables
3. Loop until desired number of variables

**Step7. [Plot Line graph of AUC values and select cut-off]**

- Split your dataset equally for training and testing
- Plot AUC values for each variable. The graph will appear as below.



The cut-off line can be drawn at the 4<sup>th</sup> variable. It gives the best AUC score around 0.81. Therefore, it will be sufficient to use only the first 4 features for training and testing. Otherwise, AUC is going down gradually.

**Step8. [Draw Cumulative Gain Chart and Lift Chart]**

Interpret these curves.

```
fx = df[var]
```

```
fy = df[tar]
```

```
logseg = LogisticRegression()
```

```
from sklearn.metrics import roc_auc_score
```

```
Print ["Log AUC", roc_auc_score(y_test, y_pred)]
```

```
type(data.columns)
```

```
def get_auc(var, tar, df):
```

```
    fx = df[var]
```

```
    fy = df[tar]
```

```
    logseg = LogisticRegression()
```

```
    logseg.fit(fx, fy)
```

```
    pred = logseg.predict_proba(fx)[0,1]
```

```
    auc_val = roc_auc_score(fy, pred)
```

```
    return auc_val
```

```
get_auc(["BMI", "Glucose"], ["Outcome"], data)
```

```
get_auc(["pregnancies", "Blood Pressure", "Skin Thickness"],  
        ["Outcome"], data)
```

```
def next_best(current_card, var, df):
```

```
    best_auc = -1
```

```
    best_var = None
```

```
    for i in var:
```

```
        auc_v = get_auc([current_card + [i], tar], df)
```

```
        if auc_v > best_auc:
```

```
            best_var = i
```

NOTES

```

best-var = i
return best = var
tar = ["Outcome"]
current = ['Insulin', 'BMI', 'Diabetes pedigree function', 'Age']
card = ['Pregnancies', 'Blood pressure', 'Skin Thickness']
next-var = next-best(current, card, tar, data)
print(next-var)
tar = ["Outcome"]
current = []
card = ['Pregnancies', 'Glucose', 'Blood pressure', 'Skin', 'Thickness', 'Insulin', 'BMI', 'Diabetes pedigree function', 'Age']
max-num = 5
num-it = min(max-num, len(card))
for i in range(0, num-it):
    next-var = next-best(current, card, tar, data)
    current = current + [next-var]
    card.remove(next-var)
    print('variable added in step' + str(i+1) + ' is ' + next-var)
    print(current)
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.5, stratify=y)

```

pred 2 = dm 3. predict = proba (y-test)

train = pd.concat [X-train, y-train, axis=1]

test = pd.concat ([X-test, y-test], axis=1)

def auc\_train\_test (variables, target, train, test):

X\_train = train [variables]

X\_test = train [variables]

Y\_train = train [target]

Y\_test = test [target]

logreg = LogisticRegression()

logreg.fit (X\_train, Y\_train).

predictions\_train = logreg.predict\_proba (X\_train)  
[:, 1]

predictions\_test = logreg.predict\_proba (X\_test)[:, 1]

auc\_train = roc\_auc\_score (X\_train, predictions\_train)

auc\_test = roc\_auc\_score (Y\_test, predictions\_test)

return (auc\_train, auc\_test)

import matplotlib.pyplot as plt

import numpy as np

X = np.array(range (0, len(auc\_values\_train)))

my\_train = np.array(auc\_values\_train)

my\_test = np.array(auc\_values\_test)

plt.xticks (X, X3.columns, rotation = 90)

plt.plot (X, my\_train)

plt.plot (X, my\_test)

plt.ylim (0.6, 1.0) plt.show()

NOTES

```
Import - sklearnplot as skplt.  
skplt.metrics.plot_cumulative_gain(y_test, pred_2)  
plt.show()  
plt.figure(figsize=(7,7))  
skplt.metrics.plot_lift_curve(y_test, pred_2)  
plt.show()
```