Lab 5:

Step 1: import pandas as pd
           import csv

```
dlab = pd.read_csv("diabetes.csv")
dlab.
dlab.head()
dlab.shape.
df = pd.read_csv("diabetes.csv")
df.
f = df.columns
         b df.info()
            df.dtypes
            df.dtypes.value_counts()
```

Step 2:
```
X = df.drop("Outcome", axis=1)
Y = df[["Outcome"]]
X
Y
from sklearn.model_selection import
        stratified_shuffle_split.
```

## Lab5. Diabetes Classification using Logistic Regression

### Objectives
In this lab, you will classify using Logistic Regression model, whether a person will become diabetic or not using PIMA diabetes dataset available in UCI. You will also predict a new person who is not in the dataset will become diabetic of not based on his details.

### Learning Outcomes
After completing this lab, you will be able to
- Understand data and build baseline LogisticRegression model
- Create Heatmap with confusion matrix values
- Apply scaling using StandardScaler and MinMaxScaler and rebuild LoR model
- Print classification metrics scores and plot ROC curve
- Compare the performance of LoR model with LogisticRegressionCV with L1 and L2 regularization

### Business Use Case
You are a data scientist. A leading hospital in your city has approached you with the medical details of their patients related to their diabetes information. The hospital has given you a file (diabetes.csv) that contains details of 768 patients and each patient is described with these 9 features. Here, *Outcome* is the *dependent variable* and *all other 8 variables are independent variables.*

- Pregnancies: Number of times pregnant
- Glucose: Plasma glucose concentration over 2 hours in an oral glucose tolerance test
- BloodPressure: Diastolic blood pressure (mm Hg)
- SkinThickness: Triceps skin fold thickness (mm)
- Insulin: 2-Hour serum insulin (mu U/ml)
- BMI: Body mass index (weight in kg/(height in m)2)
- DiabetesPedigreeFunction: Diabetes pedigree function (a function which scores likelihood of diabetes based on family history)
- Age: Age (years)
- Outcome: Class variable (0 if non-diabetic, 1 if diabetic)

The hospital wants you to build a model so that the model will assess when a new patient visits them he will become diabetic or not.

**Step1. [Understand Data].** Using Pandas, import **"diabetes.csv"** file and print properties such as head, shape, columns, dtype, info and value_counts.

**Step2. [Build Logistic Regression Model]**
- Prepare X matrix (8 feature columns) and y vector (ie., Outcome column)
- Split dataset with stratified shuffle split for training and testing as X_train, X_test, y_train, y_test (use **25%** test size).
- Create LogisticRegression model, fit on training set and predict on test set

**Step3. [Predict on a new sample]**
- Will this person become diabetic?. His details are given below.
- new_person = [[6, 200, 90, 10, 25, 23.3, 0.672, 42]]

**Step3. [Compute Classification Metrics]**
- Compute and print Accuracy, Precision, Recall and AUC scores

**Step4. [Understand Correlation]**
- Create confusion matrix between y_test and y_pred and plot confusion matrix values in a Heatmap. Explain the meaning of the 4 numbers you get.

```
stratified-shuffle-split()
    shuf = stratified shuffle split (n-split=4, test-size=0.35,
                                        random-state=0)

        shuf.get-n-splits(X,y)

        import-warnings
        warnings.filterwarnings("ignore")

    for train, test in shuf.split(x,y):
        X-train, X-test = X.iloc[train], X.loc[test]
        y-train, y-test = y.loc[train], y.iloc[test]

    from sklearn.linear.model.import logistic Regression.

        logmodel = Logistic Regression()
        logmodel.fit(X-train, y-train)
            y-predic = logmodel.predict(y-test)
            y-predic
        logmodel.score(y-train, y-train)

    step 3:-
            new-person = [[6,200,90,10,25,33.3,0.172,400]]
            print(logmodel.predict(new-person))

        precision:-
            = from sklearn.metrices import precision-score
                print(precision-score,(y-test, y-predic))

        Recall:-
            = from sklearn.metrices import recall-score
                print(recall-score(y-test, y-predic))
```

### Step5. [Normalization using MinmaxScaler and rebuild LoR]
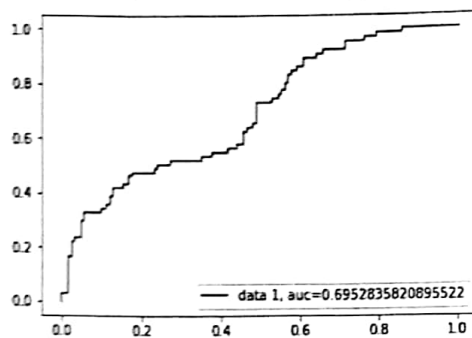- Now, normalize your X_train and X_test values using MinmaxScaler
- Create a new LogisticRegression model, fit on normalized training set and predict on the normalized test set
- Compute and print Accuracy, Precision, Recall and AUC scores

### Step6. [Normalization using StandardScaler and rebuild LoR]
- Repeat Step5 with StandardScaler
- Among the 3 models, which model gives better classification scores?
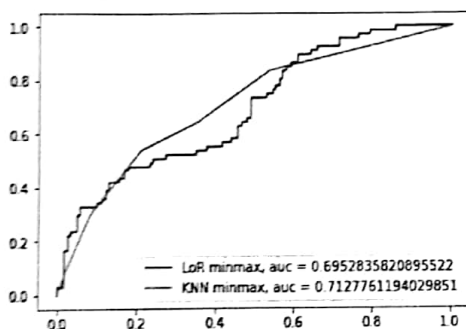
### Step7. [Plot ROC curve]
- Plot ROC curve as shown below. You can use the MinmaxScaler scaled values of X_test for computing predict_proba() score.



### Step8. [Comparison with KNN classifier].
Create a KNN classifier with default values, fit on the scaled X using MinmaxScaler, predict and print classification metric scores.

### Step9. [Update ROC curve]
- Update your ROC curve, this time, with one more curve of KNN classifier, as shown below.



### Step10. [Regularization]
- In order to reduce overfitting of your data, you will use *LogisticRegressionCV* model with L1 and L2 regularization parameters. Create both models using the following statements
    - `model1 = LogisticRegressionCV(Cs=10, cv=4, penalty='l1', solver='liblinear')`
    - `model2 = LogisticRegressionCV(Cs=10, cv=4, penalty='l2')`
- Perform fit using MinmaxScaler scaled values and predict

### Step11. [Update ROC curve]
- Update your ROC curve, this time, with two more curves, as shown below

## Accuracy:

```
from sklearn.metrices. import accuracy-score

log-acscore = accuracy-score (y-test, y-predic)

log-acscore.
```

### AUC Scores

```
from sklearn.metrices import roc_auc-score
print (roc-auc-score (y-test, y-predic))
```

## Step: 4

```
from sklearn.metrices import confussion-matrix

Confu-matrix = confussion-matrix (y-test, y-porodic)

Confu-matrix

Confu-accu-score = accuracy-score (y-test, y-perodic)

Confu-accu-score.

import seaborn as sns.

sns. heatmap (confusion-matrix (y-test, y-periodic)
              /len(y)), cmap="Y·1GnBU" anot=True)
```

## Step: 5

```
from sklearn. preprocessing import minmax scalar.

scales = Minmax scaler ()

X-trained-min = scaler. fit-transform (x-train)

X-test-min = scaler-transform (x-test)

X-trained-min.shape

X-tested-min.shape
```
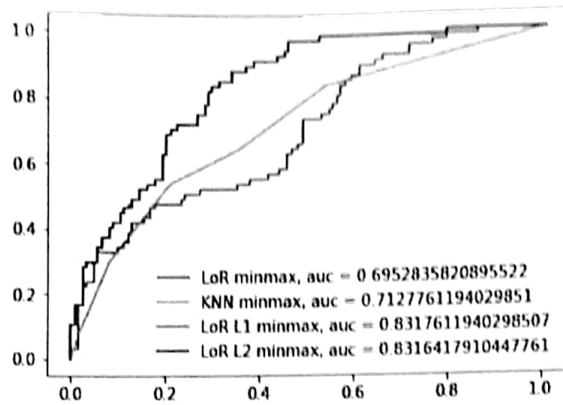
Out of these 4 models, which model performs the best?. How?. Why?.

from sklearn.Linear-model import Logistic Regression.

logmodel I = logistic Regression ()

logmodel I.fit (x-trained-min, y train

y-predict = logmodel 1.predict(x-tested-min)

y-predict.

logmodel 1.score(x-trained.min, y-train)

Precision:

from sklearn.metrices Import precision-score
print (precision-score (y-test, y-predict)

Recall

from sklearn.metrices Import recall-score

print (recall-score (y-test, y-predict)

Accuracy.

from sklearn.metrices. import accuracy-score

min-accuscore = accuracy-score(y-test, y-predict)
min-accscore

Auc Score:

```
from sklearn. metrices import roc_auc_score

log_auc_sc = roc_auc_score (y-test, y-predict)

log_auc1 = (LOR.minman, Auc =, 'log_accsare)

log_auc1.
```

step:6

```
from sklearn. proprocessing import standard scaler
scaler =. Standard scaler()
X_trained, stand = scaler. fit_transform
                                    (x-train)
X_tested, stand = scaler. transform (x-test)
X_trained _stand.shape
X_tested _stand.shape
from sklearn. linear_model import logistic Regression
logmodel 2 = logistic Regression ()
log model 2 fit (Xtrained _stand, y-train)
y_predict_stand = logmodel 2. predict (x_tested
                                            -stand)
y_predict_stand
```

```
logmodel 2. score (x_trained. stand, y-train)
precision Score
from sklearn. metrics import precision_score.
print ( precision _score (y-test, y-predict_stand))
```

Recall score
= =

from. sklearn.metrices import recall_score

Print (recall_score (y_test, y_predict_stand, ))

Accuracy:
=
from sklearn.metrices.import accuracy_Store

Stand_accscore = accuracy_score (y_test, y_predict_stand)

stand_accscore.

Auc Score
= =

from sklearn.metrices import, rou_auc_score

stand_aucscore = roc_auc_score (y_test, y_predict_stand)

stand_aucs = (Auc= stand_aucscore)

stand aur 3.

Print ('logistic Regression mode) ', log_acscore)

Print ('Minmax Scaler:', min_accscore

Print ('standard sealer:', stand_accscore)

Step 8:
=
from sklearn.metrices import. roc_curve

Pred_prb3= log mode 1, predict_proba
(x_tested_min)

```python
fprb3, tprb3, threshold = roc_curve
    (y_test, pred_prb3[:,1], pos_label=1)
import matplotlib.pyplot as plt

plt.style.use('seaborn')
plt.annotate(xy=[0.7,0], s=stand_auc3)
plt.plot(fprb3, tprb3, linestyle=':', color='black',
        label='Logistic Regression')
plt.title('Roc curve')
plt.xlabel('False positive Rate')
plt.ylabel('True positive Rate')
```

step: 8:

```python
from sklearn.neighbors import KNeighbours
                              classifier
log model 3 = KNeighbors classifier(n_neighbors=3)
log model 3.fit(x_trained_min, y_train)
knn_y_pred = log model 3.predict(x_tested_min)
knn_y_pred
```