

Lab program 4

Naïve Boyer

Code:

```
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import GaussianNB
from sklearn import metrics

df = pd.read_csv("pima_indian.csv")

feature_col_names = ['num_preg', 'glucose_conc', 'diastolic_bp', 'thickness', 'insulin', 'bmi',
'diab_pred', 'age']

predicted_class_names = ['diabetes']

X = df[feature_col_names].values # these are factors for the prediction
y = df[predicted_class_names].values # this is what we want to predict

#splitting the dataset into train and test data

xtrain,xtest,ytrain,ytest=train_test_split(X,y,test_size=0.33)

print ('\n the total number of Training Data :',ytrain.shape)
print ('\n the total number of Test Data :',ytest.shape)

# Training Naive Bayes (NB) classifier on training data.

clf = GaussianNB().fit(xtrain,ytrain.ravel())
predicted = clf.predict(xtest)
predictTestData= clf.predict([[6,148,72,35,0,33.6,0.627,50]])

#printing Confusion matrix, accuracy, Precision and Recall
```

```

print('\n Confusion matrix')
print(metrics.confusion_matrix(ytest,predicted))

print('\n Accuracy of the classifier is',metrics.accuracy_score(ytest,predicted))

print('\n The value of Precision', metrics.precision_score(ytest,predicted))

print('\n The value of Recall', metrics.recall_score(ytest,predicted))

print("Predicted Value for individual Test Data:", predictTestData)

```

Output:

```

the total number of Training Data : (514, 1)

the total number of Test Data : (254, 1)

Confusion matrix
[[145  20]
 [ 34  55]]

Accuracy of the classifier is 0.7874015748031497

The value of Precision 0.7333333333333333

The value of Recall 0.6179775280898876
Predicted Value for individual Test Data: [1]

```

Code2:

```

import csv
import random
import math

def loadcsv(filename):
    lines = csv.reader(open(filename, "r"));
    dataset = list(lines)

```

```
for i in range(len(dataset)):
    dataset[i] = [float(x) for x in dataset[i]]
return dataset
```

```
def splitdataset(dataset, splitratio):
    trainsize = int(len(dataset) * splitratio);
    trainset = []
    copy = list(dataset);
    while len(trainset) < trainsize:
        index = random.randrange(len(copy));
        trainset.append(copy.pop(index))
    return [trainset, copy]
```

```
def separatebyclass(dataset):
    separated = {}
    for i in range(len(dataset)):
        vector = dataset[i]
        if (vector[-1] not in separated):
            separated[vector[-1]] = []
        separated[vector[-1]].append(vector)
    return separated
```

```
def mean(numbers):
    return sum(numbers)/float(len(numbers))
```

```
def stdev(numbers):
    avg = mean(numbers)
    variance = sum([pow(x-avg,2) for x in numbers])/float(len(numbers)-1)
    return math.sqrt(variance)
```

```
def summarize(dataset):
```

```

summaries = [(mean(attribute), stdev(attribute)) for attribute in zip(*dataset)];
del summaries[-1] #excluding labels +ve or -ve
return summaries

```

```

def summarizebyclass(dataset):
    separated = separatebyclass(dataset);
    summaries = {}
    for classvalue, instances in separated.items():
        summaries[classvalue] = summarize(instances)
    return summaries

```

```

def calculateprobability(x, mean, stdev):
    exponent = math.exp(-(math.pow(x-mean,2)/(2*math.pow(stdev,2))))
    return (1 / (math.sqrt(2*math.pi) * stdev)) * exponent

```

```

def calculateclassprobabilities(summaries, inputvector):
    probabilities = {}
    for classvalue, classsummaries in summaries.items():
        probabilities[classvalue] = 1
        for i in range(len(classsummaries)):
            mean, stdev = classsummaries[i]
            x = inputvector[i] #testvector's first attribute
            probabilities[classvalue] *= calculateprobability(x, mean, stdev);#use normal dist
    return probabilities

```

```

def predict(summaries, inputvector):
    probabilities = calculateclassprobabilities(summaries, inputvector)
    bestLabel, bestProb = None, -1
    for classvalue, probability in probabilities.items():
        if bestLabel is None or probability > bestProb:
            bestProb = probability
            bestLabel = classvalue

```

```
return bestLabel
```

```
def getpredictions(summaries, testset):
```

```
    predictions = []
```

```
    for i in range(len(testset)):
```

```
        result = predict(summaries, testset[i])
```

```
        predictions.append(result)
```

```
    return predictions
```

```
def getaccuracy(testset, predictions):
```

```
    correct = 0
```

```
    for i in range(len(testset)):
```

```
        if testset[i][-1] == predictions[i]:
```

```
            correct += 1
```

```
    return (correct/float(len(testset))) * 100.0
```

```
def main():
```

```
    filename = 'naivedata.csv'
```

```
    splitratio = 0.67
```

```
    dataset = loadcsv(filename);
```

```
    trainingset, testset = splitdataset(dataset, splitratio)
```

```
    print('Split {0} rows into train={1} and test={2} rows'.format(len(dataset), len(trainingset),  
len(testset)))
```

```
    summaries = summarizebyclass(trainingset);
```

```
    predictions = getpredictions(summaries, testset)
```

```
    accuracy = getaccuracy(testset, predictions)
```

```
    print('Accuracy of the classifier is : {0}%'.format(accuracy))
```

```
main()
```

output:

```
Split 768 rows into train=514 and test=254 rows  
Accuracy of the classifier is : 74.01574803149606%
```
