

Experiment No-6

AIM:

A program to implement the naive Bayesian classifier for a sample training data set stored as a CSV file. Compute the accuracy of the classifier, Considering few test data sets.

Description:

Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

It is probabilistic classifier, which means it predicts on the basis of the probability of an object.

If it is not a single algorithm, but a family of algorithms where all of them share a common principle.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$P(h|D)$ is the probability of hypothesis h given data D . This is called the posterior probability.

$P(D|h)$ is probability of data d given hypothesis h is true.

$P(h)$ prior Probability of h .

$P(D)$ prior probability of D .

Source Code

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

Variance = Var((1/len(x) * sum((x - avg)**2)) for x in numbers)) /

float(len(numbers)-1)

return math.sqrt(Variance)

def unnormalize(dataSet):

Summaries = [Crosses(attribute), stdDev(attribute)]

for attribute in zip(*dataSet):

del Summaries[-1]

return summaries

def unnormalize_ByClass(dataSet):

Separated = separateByClass(dataSet)

summaries = {}

for classValue, instances in Separated.items():

summaries[classValue] = unnormalize(instances)

return summaries

def calculateClassProbabilities(summaries, inputVector):

probabilities = {}

for classValue, classSummaries in summaries.items():

itertools:

probabilities[classValue] = 1

for i in range(len(classSummaries)):

mean, stdDev = classSummaries[i]

x = inputVector[i]

probabilities[classValue] *= calculateProbability(x, mean, stdDev)

(x, mean, stdDev) 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 = 'data.csv'  
    splitRatio = 0.61  
    dataset = loadCSV(filename)  
    trainingSet, testSet = splitDataset(dataset, splitRatio)  
    print('split {} rows into train = {} and test = {}'.format(  
        len(dataset), len(trainingSet),  
        len(testSet)))
```

✓

```
    summaries = summarizeByClass(trainingSet)  
    predictions = getPredictions(summaries, testSet)  
    accuracy = getAccuracy(testSet, predictions)  
    print('Accuracy : {} %'.format(accuracy))  
main()
```

Output

Split 306 rows into train = 205 and test = 101 rows.

Accuracy : 70.2772277228%

Churn prediction result showing 70% accuracy

Decision Tree

Churn prediction using decision tree

Churn prediction using decision tree

Decision Tree

Churn prediction using decision tree

Decision Tree

Churn prediction using decision tree

Experiment - 7

AIM:

A naïve Bayesian classifier model to perform, built-in java classes/API can be used to write the program. Calculate the accuracy, precision, and recall for your data set.

Description:

Load the training dataset from the CSV file into a list of dictionaries, where each dictionary represents single row and attribute names (columns).

Determine the class variable for each instance in the training data set and add it as a new Key Value pair.

Compute the prior probabilities for each class variable by counting the number of instances (in the training) dataset that belong to each class variable and dividing by the total number of instances.

For each attribute training data set, compute the conditional probabilities for each attribute value given each class variable by counting the number of instances in the training data set.

Load the test data sets from csv files into lists of dictionaries.

For each instance dataset, posterior probability for each class variable given the attribute values in the instance, using the Naïve Bayesian.

Source code:

```
import pandas as pd
msg = pd.read_csv('document.csv', names=['message', 'label'])
print("Total Instances of Dataset : ", msg.shape[0])
msg['labelnum'] = msg['label'].map({'pos': 1, 'neg': 0})
X = msg['message']
Y = msg['labelnum']

from sklearn.model_selection import train_test_split
Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, Y)
from sklearn.feature_extraction.text import CountVectorizer
CountV = CountVectorizer()
Xtrain_dm = CountV.fit_transform(Xtrain)
Xtest_dm = CountV.transform(Xtest)
df = pd.DataFrame(Xtrain_dm.toarray()), columns=CountV.get_features().names_out()
print(df[0:5])

from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB()
clf.fit(Xtrain_dm, Ytrain)
pred = clf.predict(Xtest_dm)

for doc, p in zip(Xtest, pred):
    P = 'pos' if p == 1 else 'neg'
    print("%s → %s (%d, %s)" % (doc, p))

from sklearn.metrics import accuracy_score, confusion_matrix, precision_score, recall_score
```

```
print ('Accuracy Metrics : In')
print ('Accuracy % ', accuracy_score(ytest,pred))
print ('Recall : ', recall_score(ytest, pred))
print ('precision : ', precision_score(ytest, pred))
print ('Confusion Matrix : In', Confusion_matrix(ytest, pred))
```

Output

Total Instances of Dataset: 18

tomorrow view we went what will with work

0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0
3	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0

[5 rows x 24 columns]

I can't deal with this → neg

This is an awesome place → pos

What a great holiday → negj

I am sick and tired of this place → pos

I am tired of this stuff → pos

Accuracy Metrics:

Accuracy: 100

Recall: 100%

Precision: 1.0

Confusion Matrix

[[20]]
10 3T

Experiment-8

AIM:

A program to construct a Bayesian network considering medical data. Use this model to demonstrate the diagnosis of heart patients. Using stand Heart Disease Data set.

Description:

Bayesian network is directed acyclic graph in which each edge corresponds to a conditional dependency, and each node corresponds to a unique random variable.

The directed acyclic graph is a set of random variables represented by nodes.

Dataset:

The goal is to calculate the posterior conditional probability distribution of each of each of possible unobserved cause given the observed evidence $P[\text{Cause} | \text{evidence}]$

	Age	Gender	Family history	diet	Lifestyle	Cholesterol	heart disease
2	0	0	1	1	3	0	0
3	0	1	1	1	3	0	0
4	1	0	0	0	2	1	0
5	4	0	1	1	3	2	0
6	3	1	1	0	0	2	0
7	2	0	1	1	1	0	1
8	4	0	1	0	2	0	0
9	0	0	1	1	3	0	0
10	3	1	1	0	0	2	0
11	1	1	0	0	0	2	1
12	4	1	0	1	2	0	1
13	4	0	1	1	3	2	0
14	2	1	0	0	0	0	0
15	2	0	1	1	1	0	1
16	3	1	1	0	0	1	0
17	0	0	1	0	0	2	1
18	1	1	0	1	2	1	1

Source Code:

```
import pandas as pd
data = pd.read_csv("heartdisease.csv")
heart_disease = pd.DataFrame(data)
print(heart_disease)

from pgmpy.models import BayesianModel
model = BayesianModel([
    ('Age', 'Lifestyle'), ('Gender', 'Lifestyle'), ('Family', 'HeartDisease'),
    ('Age', 'Cholesterol'), ('Gender', 'Cholesterol'), ('Lifestyle', 'Cholesterol'),
    ('Cholesterol', 'HeartDisease'), ('Cholesterol', 'Lifestyle')
])

from pgmpy.estimators import MaximumLikelihoodEstimator
model.fit(heart_disease, estimator=MaximumLikelihoodEstimator)

from pgmpy.inference import VariableElimination
HeartDisease_infer = VariableElimination(model)

print("For age Enter {SuperseniorCitizen:0, SeniorCitizen:1,
MiddleAged:2, Youth:3, Teen:4}")
print("For Gender Enter {Male:0, Female:1}")
print("For Family History Enter {Yes:1, No:0}")
print("For diet Enter {High:0, Medium:1}")
print("For lifestyle Enter {Affiliate:0, Active:1, Moderate:2,
Sedentary:3}")

q = HeartDisease_infer.query(variables=['HeartDisease'],
                               evidence={'Age': int(input("Enter age:")),
                                         'Gender': int(input("Enter Gender:")),
                                         'Family': int(input("Enter family history:")),
                                         'Cholesterol': int(input("Enter Cholesterol:")),
                                         'Lifestyle': int(input("Enter Lifestyle:"))
                                         })
```

Experiment - 9

AIM:

Apply EM algorithm to cluster a set of data stored in a .csv file. Use the same data set for clustering using K-means algorithm. Compare the results of these two algorithms.

Description:

EM algorithm is a iterative approach that cycles between two modes. The first mode attempts to estimate the missing or latent variables called the estimation-step or E-step; the second mode attempts to optimize the parameters of the model to best explain the data, called maximization step.

Initially, a set of initial values are considered. A set of incomplete observed data is given to the system with assumption that observed data.

use complete data generated in the preceding "expectation" - step in order to update the values of parameters. It is basically used to update the hypothesis.

Now, in the next step, it is checked whether the values are converging or not, if yes, then stop otherwise repeat step-2 and step-3 until the convergence occurs.

Source Code:

```
from sklearn.cluster import KMeans  
from sklearn.preprocessing import StandardScaler  
from sklearn.mixture import GaussianMixture  
from sklearn.datasets import load_iris  
import sklearn.metrics as sm  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
  
dataset = load_iris()  
# print(dataset)  
x = pd.DataFrame(dataset.data)  
x.columns = ['Sepal-length', 'Sepal-width', 'Petal-length',  
             'Petal-width']  
y = pd.DataFrame(dataset.target)  
y.columns = ['Targets']  
# print(x)  
plt.figure(figsize=(14, 7))  
colormap = np.array(['red', 'lime', 'black'])  
  
# Real plot  
plt.subplot(1, 3, 1)  
plt.scatter(x.Petal-length, x.Petal-width, c=colormap  
           [y.Targets], s=40)  
plt.title('Real')  
  
# K-plot  
plt.subplot(1, 3, 2)  
model = KMeans(n_clusters=3)  
model.fit(x)  
predY = np.choose(model.labels_, [0, 1, 2]).astype(np.int64)  
plt.scatter(x.Petal-length, x.Petal-width, c=colormap  
           [predY], s=40)  
plt.title('K-Means')
```

GMM PLOT

scaler = preprocessing, StandardScaler()

scaler • fit(x)

Xsa = scaler • transform(x)

Xs = pd • DataFrame(Xsa, columns = x • columns)

gmm = GaussianMixture(n_components = 3)

gmm • fit(xs)

y_Cluster_gmm = gmm • predict(xs)

plt • subplot(1, 3, 3)

plt • Scatter(x • petal_Length, x • petal_Width, c = colormap

[y_Cluster_gmm], s = 400)

plt • title("GMM Classification")

Output



Experiment - 10

AIM:

A program to implement K-Nearest Neighbour algorithm to the Iris data set.

Description:

We are going to write a function, which will find the distance between two given 2-D points in the X-Y plane. import numpy, to take help of numpy arrays for storing the coordinates. Finding the distance between two points will help in finding the nearest neighbor of input point, from existing point, classify the point to insert into set of existing points and then insert accordingly.

Build a function find_nearest_neighbour() to find nearest given point. It will take in point we wish to insert, set of existing points and K-helps with indices as parameters to the function, visualize the situation by plotting the X-Y plane filled with points with the help of matplotlib.

Source code:

```
from sklearn.datasets import load_iris  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.model_selection import train_test_split  
import numpy as np
```

```

dataset = load_iris
# print(dataset)
x_train, x_test, y_train, y_test = train_test_split(dataset
[["data"]], dataset[["target"]], random_state=0)

kn = KNeighborsClassifier(n_neighbors=1)
kn.fit(x_train, y_train)

for i in range(len(x_test)):
    x = x_test[i]
    x_new = np.array([x])
    prediction = kn.predict(x_new)
    print(f"TARGET = {y_test[i]}, dataset[\"target_names\"]"
[y-test[i]], "PREDICTED = ", prediction, dataset
[["target_names"]][prediction])

print(kn.score(x_test, y_test))

```

Output

```

TARGET = 2 virginica PREDICTED = [2] ["virginica"]
TARGET = 1 versicolor PREDICTED = [0] ["versicolor"]
TARGET = 0 setosa PREDICTED = [0] ["setosa"]
TARGET = 2 virginica PREDICTED = [2] ["virginica"]
TARGET = 0 setosa PREDICTED = [0] ["setosa"]
TARGET = 2 virginica PREDICTED = [2] ["virginica"]
TARGET = 1 versicolor PREDICTED = [0] ["versicolor"]
TARGET = 1 versicolor PREDICTED = [1] ["versicolor"]
TARGET = 2 virginica PREDICTED = [2] ["virginica"]
TARGET = 0 setosa PREDICTED = [0] ["setosa"]
TARGET = 1 versicolor PREDICTED = [1] ["versicolor"]
TARGET = 1 versicolor PREDICTED = [0] ["versicolor"]
TARGET = 0 setosa PREDICTED = [0] ["setosa"]
TARGET = 2 virginica PREDICTED = [2] ["virginica"]
TARGET = 1 versicolor PREDICTED = [0] ["versicolor"]
TARGET = 1 versicolor PREDICTED = [1] ["versicolor"]

```

0.9736842105263158

Experiment - II

AIM:

Implement the non-parametric locally weighted Regression algorithm in order to fit datapoints. Select dataset and draw graph.

Description:

Regression is a technique from statistics that is used to predict values of a desired target quantity when the target quantity is continuous.

X is a independent variable

Y is a dependent variable

Loess/ Loimess Regression

Loess regression is a nonparametric technique that uses local weighted regression to fit a smooth curve through points of scatter plot.

Loimess Algorithm.

Locally weighted regression is a very powerful nonparametric model used in statistical learning.

Given dataset X, y , we attempt to find a model parameter β to minimize residual sum of weighted squared errors.

value of model term β using:

$$\hat{\beta}(x_0) = (X^T W X)^{-1} X^T W y$$

prediction = $x_0 * \beta$

Source Code:

```
import numpy as np
from bokeh.plotting import figure, show, output_notebook
from bokeh.layouts import gridplot
from bokeh.io import push_notebook

def local_regression(x0, x, y, tau): # add bias term
    x0 = np.r_[1, x0]
    X = np.c_[np.ones(len(x)), x]
    XW = X.T * radial_Kernel(x0, x, tau) # X transpose * W
    beta = np.linalg.pinv(XW @ X) @ XW @ y # @ Matrix multiplication or Dot product
    # predict value
    return x0 @ beta # @ Matrix multiplication or Dot for prediction

def radial_Kernel(x0, x, tau):
    return np.exp(np.sum((x - x0) ** 2, axis=1) / (-2 * tau * tau))

# weight or Radial Kernel Basis Function
n = 100
# generate dataset
X = np.linspace(-3, 3, num=n)
print("The Data set (10 samples) X: {} x [1:10]".format(X))
Y = np.log(np.abs(X ** 2 - 1) + 0.5)
print("The Fitting curve Dataset (10 samples) Y: {} x [1:10]".format(Y))
# jitter X
x_t = np.random.normal(scale=1, size=n)
print("Normalised (10 samples) X: {} x [1:10]".format(x_t))
```

```

domain = np.linspace(-3, 3, num=300)
print("x0 Domain vspace (10 samples):", domain[1:10])
def plot_lmr(tau):
    # prediction through regression
    prediction = [local_regression(x0, x1, y, tau) for x0 in
                  domain]
    plot = figure(plot_width=400, plot_height=400)
    plot.title.text = f"tau = {log_tau} % tau"
    plot.scatter(x, y, alpha=.3)
    plot.line(domain, prediction, line_width=2, color='red')
    return plot
    show(gridplot([[plot_lmr(10), plot_lmr(1)],

    [plot_lmr(0.01), plot_lmr(0.001)]]))

```

Outputs:

$\tau = 10$:
 $y = 0.5 + 0.5 \sin(\pi x)$
 $\tau = 1$:
 $y = 0.5 + 0.5 \sin(\pi x) + 0.05 \sin(10\pi x)$
 $\tau = 0.01$:
 $y = 0.5 + 0.5 \sin(\pi x) + 0.05 \sin(10\pi x) + 0.005 \sin(100\pi x)$
 $\tau = 0.001$:
 $y = 0.5 + 0.5 \sin(\pi x) + 0.05 \sin(10\pi x) + 0.005 \sin(100\pi x) + 0.0005 \sin(1000\pi x)$