



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Name:	
Roll No:	
Class/Sem:	TE/V
Experiment No.:	7
Title:	Implementation of any one classifier using languages like JAVA/ python.
Date of Performance:	
Date of Submission:	
Marks:	
Sign of Faculty:	



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Aim: To implement Naïve Bayesian classification

Objective: Develop a program to implement Bayesian classification.

Theory:

The Naive Bayes is a classification algorithm that is suitable for binary and multiclass classification. Naïve Bayes performs well in cases of categorical input variables compared to numerical variables. It is useful for making predictions and forecasting data based on historical results.

The naïve Bayesian classifier, or simple Bayesian classifier, works as follows:

- 1) Let D be a training set of tuples and their associated class labels. As usual, each tuple is represented by an n -dimensional attribute vector, $X = (x_1, x_2, \dots, x_n)$, depicting n measurements made on the tuple from n attributes, respectively, A_1, A_2, \dots, A_n .
- 2) Suppose that there are m classes, C_1, C_2, \dots, C_m . Given a tuple, X , the classifier will predict that X belongs to the class having the highest posterior probability, conditioned on X . That is, the naïve Bayesian classifier predicts that tuple X belongs to the class C_i if and only if

$P(C_i|X) > P(C_j|X)$ Thus we maximize $P(C_i|X)$. The class C_i for which $P(C_i|X)$ is maximized is called the maximum posteriori hypothesis.

By Bayes' theorem,

$$P(C_i|X) = P(X|C_i) * P(C_i) / P(X)$$

- 3) As $P(X)$ is constant for all classes, only $P(X|C_i)P(C_i)$ need be maximized. If the class prior probabilities are not known, then it is commonly assumed that the classes are equally likely, that is, $P(C_1) = P(C_2) = \dots = P(C_m)$, and we would therefore maximize $P(X|C_i)$. Otherwise, we maximize $P(X|C_i)P(C_i)$. Note that the class prior probabilities may be estimated by $P(C_i) = |C_i, D| / |D|$, where $|C_i, D|$ is the number of training tuples of class C_i in D .

The equation: Posterior = Prior x (Likelihood over Marginal probability)

There are four parts:

- Posterior probability (updated probability after the evidence is considered)
- Prior probability (the probability before the evidence is considered)
- Likelihood (probability of the evidence, given the belief is true)
- Marginal probability (probability of the evidence, under any circumstance)

Bayes' Rule can answer a variety of probability questions, which help us (and machines) understand the complex world we live in.



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Example:

Car No	Color	Type	Origin	Stolen
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	No
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

$$P(\text{yes}) = 5/10$$

$$P(\text{No}) = 5/10$$

-Color:

$$P(\text{Red/Y}) = 3/5 \quad P(\text{yellow/Y}) = 2/5$$

$$P(\text{Red/N}) = 2/5 \quad P(\text{yellow/N}) = 3/5$$

-Type:

$$P(\text{SUV/Y}) = 1/5 \quad P(\text{Sports/Y}) = 4/5$$

$$P(\text{SUV/N}) = 3/5 \quad P(\text{Sports/N}) = 2/5$$

-Origin:

$$P(\text{Domestic/Y}) = 2/5 \quad P(\text{Imported/Y}) = 3/5$$

$$P(\text{Domestic /N}) = 3/5 \quad P(\text{Imported/N}) = 2/5$$

$$P(x|\text{Yes}).P(\text{Yes}) = 0.024$$

$$P(x|\text{No}).P(\text{No}) = 0.072$$

So, Bayesian Classification Predicts the class "NO"



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Code:

```
import numpy as np
import pandas as pd
# Import necessary modules
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
d1 = pd.read_csv('diabetes_csv.csv')

d1.head()

# Loading data
# Create feature and target arrays
X = d1[d1.columns[:-1]]
y = d1[d1.columns[-1]]
# Split into training and test set
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size = 0.2, random_state=42)
gnb = GaussianNB()
gnb.fit(X_train, y_train)
y_pred = gnb.predict(X_test)

d1

from sklearn import metrics
print("Gaussian Naive Bayes model accuracy(in %):",
metrics.accuracy_score(y_test, y_pred)*100)
# Predict on dataset which model has not seen before
print(gnb.predict(X_test))
```

Output:

[illegible]

CSL503: Data warehousing and Mining Lab