

Assignment : DA

Title:-

Naive Bayes algorithm for classification on pima Indians dataset.

Problem Statement:-

- Download pima Indians diabetes dataset use Naive Bayes algorithm for classification
- Load the data from CSV and split it into training and test datasets
- Summarise the properties in the training dataset so that we can calculate and make predictions.
- classify the samples from test dataset and a summarized training dataset.

Objective:

- learning Naive Bayes Algorithm.
- learn to use Naive Bayes Algorithm for classification on given dataset.

Software and hardware apparatus used

- 1) OS : 64 bit open source linux
- 2) programming lang: python / R

Outcomes :-

Students will be able to summarise the properties of the dataset, split the dataset into training and test data and apply

Naive Bayes algorithm from classification of application.

Related Mathematics:-

Mathematical model:-

let S be the system set

$$S = \{ s; e; ; x; y; fme; DD; NDD; FC; sc \}$$

where Dataset is loaded into the dataframe

s = start state

e = end state i.e. classification of samples from the test dataset

x = set of inputs.

$$x = \{ x_i \}$$

where

x_1 = Pima Indians Diabetes dataset

where,

y = set of outputs.

- 1) Spitting the dataset into training and test datasets
- 2) Naive bayes classifier.

fme is the set of main functions

$$fme = \{ f_1, f_2, f_3 \}$$

where,

f_1 = function to load dataset into dataframe

f_2 = function to split data set into train & test data

f_3 = function to invoke naive bayes classifier

DD = Deterministic data

PIMA indians diabetes dataset

NDD = Non-Deterministic data

FC = failure case

Failed to classify the record into correct class.

Theory:-

Naive Bayes classifier are a collection of classification algorithms ~~are~~ based on Bayes Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle i.e. every pair of feature being classified is independent of each other.

The dataset is divided into two parts namely, feature matrix and response vector.

Feature matrix contains all the vectors (rows) of dataset in which each vector consists of value of dependent feature. In datasets features like outlook, temperature, humidity and windy are dependent features.

Response vector matrix contains the value of class variable (prediction and output) for each row of feature matrix.

The fundamental Naive Bayes assumption is that each feature makes an independent & equal contribution to outcome.

• Bayes Theorem.

Bayes Theorem finds the probability of an event occurring given probability of another event occurring given probability of another event ~~already~~ already occurred Bayes ~~theorem~~.

theorem is stated mathematically as following equation.

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

Where, A and B are events and $P(B) \neq 0$.

Basically, we are trying to find probability of event A, given the event B is true. Event B is true and also termed as evidence $P(A)$ is the prior of A (The prior probability i.e. probability of event before evidence is seen). The evidence is an attribute value of an unknown instance (here, it is event B). $P(A|B)$ is posteriori probability of B i.e. probability of event after evidence is seen.

Now, with regards to our datasets we can apply Bayes algorithm theorem in following way

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)}$$

where y is class variable and x is a dependent feature vector (of size n) where:
 $x \in (x_1, x_2, x_3, \dots, x_n)$

• Naive assumption

Now, we put naive assumption to bayes algorithm which is independence among the feature so now we split evidence into independent parts.

Now, if any two events A and B are independent then,

$$P(A, B) = P(A) \cdot P(B)$$

Test case:-

For given dataset

Confusion matrix is

$$\begin{bmatrix} 96 & 29 \\ 26 & 41 \end{bmatrix}$$

and accuracy score is 0.713541666.

Conclusion:-

In this way naive bayes classifier is used for pima indians dataset analysis.