

* Assignment DA2 *

* 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 Classification on given dataset.

* SW packages and H/W apparatus used:

① OS:- 64 bit open source linux.

② programming language python/R.

* Outcomes:

Students will be able to summarize the properties of the dataset, split the dataset into training & test data and apply naive bayes algorithm for classification of application.

Related mathematics:-

Mathematical Model

Let S be system set

$S = \{s; e; x; y; fme; pp, NDP, FC, sc\}$

where dataset is loaded into the dataframe

s = start state

e = end state i.e. Classification of sample from the test dataset.

x = set of inputs.

$x = \{x_1\}$

where

x_1 = Pima Indians diabetes dataset where,

y = set of outputs.

① Splitting of dataset into training & 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 dataset into training & test data.

f_3 :- function to invoke naive Bayes Classifier

DD:- Deterministic dataset.

PIMA Indians diabetes datasets.

NDP:- Non-deterministic data.

NULL values in dataset.

FC:- Failure Case.

Failed to classify the record into correct class.

Theory:-

Naive Bayes Classifiers are a collection of classification algorithms 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 features 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 features. In datasets features like outlook, temperature, humidity and windy are dependent features.

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

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

* Bayes theorem:-

Bayes theorem finds the probability of an event occurring given the probability of another event occurring given the probability of another event that has already occurred. Bayes's theorem is stated mathematically

as following Equation.

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

where $A \& B$ are events and $P(B)$ 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 a posteriori probability of B , i.e. probability of event after evidence is seen.

Now, with regards to our dataset, we can apply Baye's theorem in following way,

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

where y is Class Variable and x is a dependent feature vector (of size n) where:

$$x = (x_1, x_2, x_3, \dots, x_n)$$

*Naive Assumption:-

Now, we can put naive assumption to baye's theorem which is independent among the features. so now, we split evidence into the independent parts,

Now, if any two events A & B are independent.

then

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

* Test Cases:-

For given dataset

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

and accuracy score is 0.71354266

* Conclusion:- In this way, naive bayes classifier is used for pima indians dataset analysis.