

PROBLEM STATEMENT :

Implement Naive Bayes to predict the work type for a person with following parameters: age: 30, Qualification: MTech, Experience: 8

Following table provides the details of the available data:

Work Type	Age	Qualification	Experience
Consultancy	30	Ph.D.	9
Service	21	MTech.	1
Research	26	MTech.	2
Service	28	BTech.	10
Consultancy	40	MTech.	14
Research	35	Ph.D.	10
Research	27	BTech.	6
Service	32	MTech.	9
Consultancy	45	Btech.	17
Research	36	Ph.D.	7

OBJECTIVE :

- To apply algorithmic strategies while solving problems
- To develop time and space efficient algorithms
- To study algorithmic examples in distributed, concurrent and parallel environments

THEORY :

The Naive Bayesian classifier is based on Bayes' theorem with independence assumptions between predictors. A Naive Bayesian model is easy to build, with no complicated iterative parameter estimation which makes it particularly useful for very large datasets. Despite its simplicity, the Naive Bayesian classifier often does surprisingly well and is widely used because it often outperforms more sophisticated classification methods.

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

$$P(C_i | X) > P(C_j | X) \text{ for } 1 \leq j \leq m, j \neq i.$$

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,

$$P(C_i | X) = [P(X | C_i)P(C_i)] / P(X)$$

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 $P(C_i) = |C_i|/|D|$, where $|C_i|$ is the number of training tuples of class C_i in D .

Given data sets with many attributes, it would be extremely computationally expensive to compute $P(X|C_i)$. In order to reduce computation in evaluating $P(X|C_i)$, the naive assumption of class conditional independence is made. This presumes that the values of the attributes are conditionally independent of one another, given the class label of the tuple (i.e., that there are no dependence relationships among the attributes). Thus,

$$\begin{aligned} P(X|C_i) &= \prod_{k=1}^n P(x_k|C_i) \\ &= P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i). \end{aligned}$$

We can easily estimate the probabilities $P(x_1|C_i)$, $P(x_2|C_i)$, \dots , $P(x_n|C_i)$ from the training tuples. Recall that here x_k refers to the value of attribute A_k for tuple X . For each attribute, we look at whether the attribute is categorical or continuous-valued. For instance, to compute $P(X|C_i)$, we consider the following:

(a) If A_k is categorical, then $P(x_k|C_i)$ is the number of tuples of class C_i in D having the value x_k for A_k , divided by $|C_i|$, the number of tuples of class C_i in D .

(b) If A_k is continuous-valued, then we need to do a bit more work, but the calculation is pretty straightforward. A continuous-valued attribute is typically assumed to have a Gaussian distribution with a mean μ and standard deviation σ , defined by

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

so that,

$$P(x_k|C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i}).$$

These equations may appear daunting, but hold on! We need to compute μ_{C_i} and σ_{C_i} , which are the mean (i.e., average) and standard deviation, respectively, of the values of attribute A_k for training tuples of class C_i . We then plug these two quantities into Equation (6.13), together with x_k , in order to estimate $P(x_k|C_i)$.

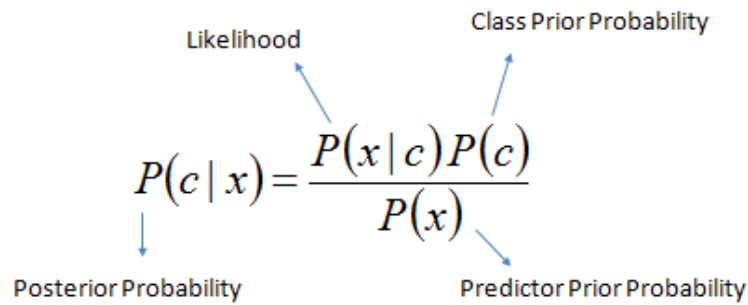
In order to predict the class label of X , $P(X|C_i)P(C_i)$ is evaluated for each class C_i . The classifier predicts that the class label of tuple X is the class C_i if and only if

$$P(X|C_i)P(C_i) > P(X|C_j)P(C_j) \text{ for } 1 \leq j \leq m, j \neq i.$$

In other words, the predicted class label is the class C_i for which $P(X|C_i)P(C_i)$ is the maximum.

ALGORITHM :

Bayes theorem provides a way of calculating the posterior probability, $P(c|x)$, from $P(c)$, $P(x)$, and $P(x|c)$. Naive Bayes classifier assume that the effect of the value of a predictor (x) on a given class (c) is independent of the values of other predictors. This assumption is called class conditional independence.



$$P(c | X) = P(x_1 | c) \times P(x_2 | c) \times \dots \times P(x_n | c) \times P(c)$$

- $P(c|x)$ is the

posterior probability of *class (target)* given *predictor (attribute)*.

- $P(c)$ is the prior probability of *class*.
- $P(x|c)$ is the likelihood which is the probability of *predictor* given *class*.
- $P(x)$ is the prior probability of *predictor*

INPUT :

Enter age , Qualification and Experience

EXPECTED OUTPUT :

Predicted worktype

MATHMODEL

D : Set of tuples

- Each Tuple is an 'n' dimensional attribute vector
- $X : (x_1, x_2, x_3, \dots, x_n)$

Let there be 'm' Classes : $C_1, C_2, C_3, \dots, C_m$

Naïve Bayes classifier predicts X belongs to Class C_i iff

- $P(C_i/X) > P(C_j/X)$ for $1 \leq j \leq m, j \neq i$

Maximum Posteriori Hypothesis

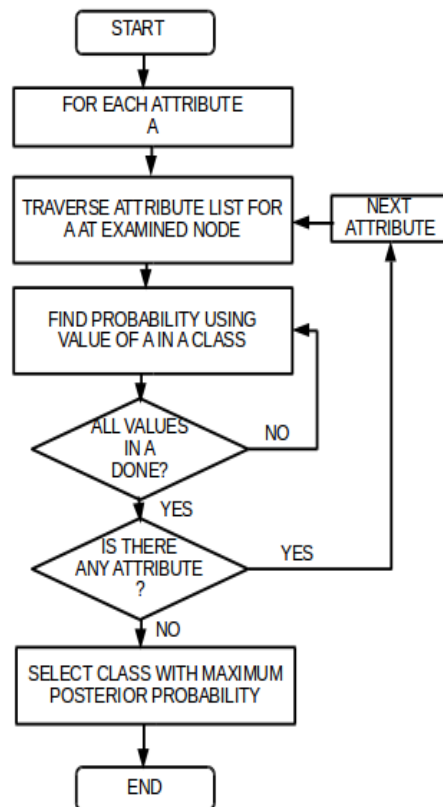
- $P(C_i/X) = P(X/C_i) P(C_i) / P(X)$
- Maximize $P(X/C_i) P(C_i)$ as $P(X)$ is constant

With many attributes, it is computationally expensive to evaluate $P(X/C_i)$.

Naïve Assumption of "class conditional independence"

- $P(X / C_i) = \prod_{k=1}^n P(x_k / C_i)$
- $P(X/C_i) = P(x_1/C_i) * P(x_2/C_i) * \dots * P(x_n/C_i)$

FLOWCHART :



TEST CASES :

TEST CASE	INPUT	EXPECTED OUTPUT	OUTPUT ACHIEVED	REMARKS
1.	ENTER AGE :30 ENTER QUALIFICATION :M.Tech ENTER EXPERIENCE :8	Research	Research	Correct
2.	ENTER AGE :40 ENTER QUALIFICATION :M.Tech ENTER EXPERIENCE :15	Consultancy	Consultancy	Correct
3.	ENTER AGE :40 ENTER QUALIFICATION :P.hD ENTER EXPERIENCE :8	Research	Research	Correct

SPACE AND TIME COMPLEXITIES :

In Baye's classification we use three nested loops which iterate over the entire length of the dataset. Hence time complexity is $O(N^3)$.

CONCLUSION :

Hence we have successfully implemented the naive bayes algorithm to predict the work type of the person.

OUTCOMES ACHIEVED :

COURSE OUTCOME	ACHIEVED(√)
Problem solving abilities for smart devices.	
Problem solving abilities for gamifications.	
Problem solving abilities of pervasiveness,embedded security and NLP.	
To solve problems for multicore or distributed,concurrent/Parallel environments	√