

Data Science

Unit-VI

Bayes' Classifier

Bayes' Classifier

Naïve Bayes Classifier

According to *Bayes theorem*, we need to calculate the posterior probability

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

Diagram labels for the Naïve Bayes formula:

- $P(x | c)$ is labeled **Likelihood**
- $P(c)$ is labeled **Class Prior Probability**
- $P(c | x)$ is labeled **Posterior Probability**
- $P(x)$ is labeled **Predictor Prior Probability**

Or simply we calculate in expanded form:

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

We need to calculate it for each class and then compare the results to find which gives the higher score.



Thomas Bayes
1702 - 1761

BAYE'S THEOREM

1. Conditional Probability

$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$

where A = Events 1

B = Events 2. (given).

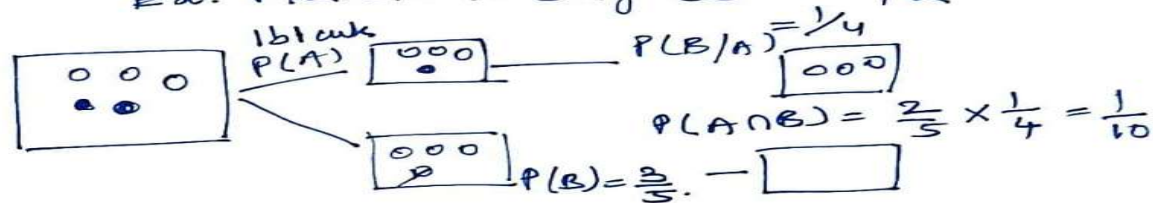
2. Independent Events

Ex: Tossing two coins.

Probability = 0.5, $q = 0.5$

3. Dependent Events

Ex: Marbles in bag (2 Red + 2 black)



$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$

$$\therefore A \cap B = B \cap A$$

$$P(B/A) = \frac{P(B \cap A)}{P(A)}$$

$$\therefore P(A/B) \times P(B) = P(B/A) \times P(A)$$

$$P(A/B) = \frac{P(B/A) \times P(A)}{P(B)}$$

where $P(A/B)$ = Posterior probability

$P(B/A)$ = Likelihood.

$P(A)$ = Prior Probability.

$P(B)$ = marginal probability

Naive Bayes' classifier

Apply Baye's theorem on dataset.

Let Dataset

$$\text{record } x = \{x_1, x_2, x_3, \dots, x_n\} \{y\}$$

$$P(y/x_1, x_2, \dots, x_n) = \frac{P(x_1/y) P(x_2/y) \dots P(x_n/y) * P(y)}{P(x_1) P(x_2) \dots P(x_n)}$$

$$P(y/x_1, x_2, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i/y)}{\underbrace{P(x_1) P(x_2) \dots P(x_n)}_{\leftarrow \text{constant.}}}$$

$$P(y/x_1, x_2, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i/y)$$

$$y = \underbrace{\arg \max_y}_{\text{takes max value.}} P(y) * \prod_{i=1}^n P(x_i/y).$$

Example

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$$P(\text{PlayTennis} = \text{yes}) = 9/14 = .64$$

$$P(\text{PlayTennis} = \text{no}) = 5/14 = .36$$

Outlook	Y	N	Humidity	Y	N
sunny	2/9	3/5	high	3/9	4/5
overcast	4/9	0	normal	6/9	1/5
rain	3/9	2/5			
Temperature			Windy		
hot	2/9	2/5	Strong	3/9	3/5
mild	4/9	2/5	Weak	6/9	2/5
cool	3/9	1/5			

Example

$\langle \text{Outlook} = \text{sunny}, \text{Temperature} = \text{cool}, \text{Humidity} = \text{high}, \text{Wind} = \text{strong} \rangle$

$$\begin{aligned} v_{NB} &= \operatorname{argmax}_{v_j \in \{\text{yes}, \text{no}\}} P(v_j) \prod_i P(a_i | v_j) \\ &= \operatorname{argmax}_{v_j \in \{\text{yes}, \text{no}\}} P(v_j) \cdot P(\text{Outlook} = \text{sunny} | v_j) P(\text{Temperature} = \text{cool} | v_j) \\ &\quad \cdot P(\text{Humidity} = \text{high} | v_j) P(\text{Wind} = \text{strong} | v_j) \end{aligned}$$

Test Phase: Predict class for instar

- Given a new instance,

$\mathbf{x}' = (\text{Outlook} = \text{Sunny}, \text{Temperature} = \text{Cool}, \text{Humidity} = \text{High}, \text{Wind} = \text{Strong})$

- Look up tables

$$P(\text{Outlook} = \text{Sunny} | \text{Play} = \text{Yes}) = 2/9$$

$$P(\text{Outlook} = \text{Sunny} | \text{Play} = \text{No}) = 3/5$$

$$P(\text{Temperature} = \text{Cool} | \text{Play} = \text{Yes}) = 3/9$$

$$P(\text{Temperature} = \text{Cool} | \text{Play} = \text{No}) = 1/5$$

$$P(\text{Humidity} = \text{High} | \text{Play} = \text{Yes}) = 3/9$$

$$P(\text{Humidity} = \text{High} | \text{Play} = \text{No}) = 4/5$$

$$P(\text{Wind} = \text{Strong} | \text{Play} = \text{Yes}) = 3/9$$

$$P(\text{Wind} = \text{Strong} | \text{Play} = \text{No}) = 3/5$$

$$P(\text{Play} = \text{Yes}) = 9/14$$

$$P(\text{Play} = \text{No}) = 5/14$$

- MAP rule

$$P(\text{Yes} | \mathbf{x}'): [P(\text{Sunny} | \text{Yes})P(\text{Cool} | \text{Yes})P(\text{High} | \text{Yes})P(\text{Strong} | \text{Yes})]P(\text{Play} = \text{Yes}) = 0.0053$$

$$P(\text{No} | \mathbf{x}'): [P(\text{Sunny} | \text{No})P(\text{Cool} | \text{No})P(\text{High} | \text{No})P(\text{Strong} | \text{No})]P(\text{Play} = \text{No}) = 0.0206$$

Given the fact $P(\text{Yes} | \mathbf{x}') < P(\text{No} | \mathbf{x}')$, we label \mathbf{x}' to be “No”.

Bayesian Classification: Why?

A statistical classifier: performs *probabilistic prediction*, i.e., predicts class membership probabilities

Foundation: Based on Bayes' Theorem.

Performance: A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers

Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct — prior knowledge can be combined with observed data

Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

Bayes' Theorem: Basics

Total probability Theorem:
$$P(B) = \sum_{i=1}^M P(B|A_i)P(A_i)$$

Bayes' Theorem:

$$P(H | \mathbf{X}) = \frac{P(\mathbf{X} | H)P(H)}{P(\mathbf{X})} = P(\mathbf{X} | H) \times P(H) / P(\mathbf{X})$$

Let \mathbf{X} be a data sample ("*evidence*"): class label is unknown

Let H be a *hypothesis* that X belongs to class C

Classification is to determine $P(H|\mathbf{X})$, (i.e., *posteriori probability*): the probability that the hypothesis holds given the observed data sample \mathbf{X}

$P(H)$ (*prior probability*): the initial probability

E.g., \mathbf{X} will buy computer, regardless of age, income, ...

$P(\mathbf{X})$: probability that sample data is observed

$P(\mathbf{X}|H)$ (*likelihood*): the probability of observing the sample \mathbf{X} , given that the hypothesis holds

E.g., Given that \mathbf{X} will buy computer, the prob. that X is 31..40, medium income

Prediction Based on Bayes' Theorem

Given training data \mathbf{X} , *posteriori probability of a hypothesis* H , $P(H|\mathbf{X})$, follows the Bayes' theorem

$$P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$$

Informally, this can be viewed as

posteriori = likelihood x prior/evidence

Predicts \mathbf{X} belongs to C_i iff the probability $P(C_i|\mathbf{X})$ is the highest among all the $P(C_k|\mathbf{X})$ for all the k classes

Practical difficulty: It requires initial knowledge of many probabilities, involving significant computational cost

Classification Is to Derive the Maximum Posteriori

Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n -D attribute vector $\mathbf{X} = (x_1, x_2, \dots, x_n)$

Suppose there are m classes C_1, C_2, \dots, C_m .

Classification is to derive the maximum posteriori, i.e., the maximal $P(C_i | \mathbf{X})$

This can be derived from Bayes' theorem

Since $P(\mathbf{X})$ is constant for all classes, only

needs to be maximized
$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

$$P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i)P(C_i)$$

Naïve Bayes Classifier

A simplified assumption: attributes are conditionally independent (i.e., no dependence relation between attributes):

$$P(\mathbf{X} | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

This greatly reduces the computation cost: Only counts the class distribution

If A_k is categorical, $P(x_k | C_i)$ is the # of tuples in C_i having value x_k for A_k divided by $|C_{i,D}|$ (# of tuples of C_i in D)

If A_k is continuous-valued, $P(x_k | C_i)$ is usually computed based on Gaussian distribution with a mean μ and standard deviation σ

and $P(x_k | C_i)$ is

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

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

Naïve Bayes Classifier: Training Dataset

Given the training data in the table below (*Buy Computer* data), predict the class of the following new example using Naïve Bayes classification: age \leq 30, income=medium, student=yes, credit-rating=fair

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data to be classified:

X = (age \leq 30,

Income = medium,

Student = yes

Credit_rating = Fair)

RID	age	income	student	credit_rating	Class: buys_computer
1	≤ 30	high	no	fair	no
2	≤ 30	high	no	excellent	no
3	31 ... 40	high	no	fair	yes
4	> 40	medium	no	fair	yes
5	> 40	low	yes	fair	yes
6	> 40	low	yes	excellent	no
7	31 ... 40	low	yes	excellent	yes
8	≤ 30	medium	no	fair	no
9	≤ 30	low	yes	fair	yes
10	> 40	medium	yes	fair	yes
11	≤ 30	medium	yes	excellent	yes
12	31 ... 40	medium	no	excellent	yes
13	31 ... 40	high	yes	fair	yes
14	> 40	medium	no	excellent	no

Naïve Bayes Classifier: An Example

Solution:

$E = \text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair}$

E_1 is $\text{age} \leq 30$, E_2 is $\text{income} = \text{medium}$, E_3 is $\text{student} = \text{yes}$, E_4 is $\text{credit-rating} = \text{fair}$

We need to compute $P(\text{yes}|E)$ and $P(\text{no}|E)$ and compare them.

$$P(\text{yes} | E) = \frac{P(E_1 | \text{yes}) P(E_2 | \text{yes}) P(E_3 | \text{yes}) P(E_4 | \text{yes}) P(\text{yes})}{P(E)}$$

$$P(\text{yes}) = 9/14 = 0.643$$

$$P(\text{no}) = 5/14 = 0.357$$

$$P(E_1 | \text{yes}) = 2/9 = 0.222$$

$$P(E_1 | \text{no}) = 3/5 = 0.6$$

$$P(E_2 | \text{yes}) = 4/9 = 0.444$$

$$P(E_2 | \text{no}) = 2/5 = 0.4$$

$$P(E_3 | \text{yes}) = 6/9 = 0.667$$

$$P(E_3 | \text{no}) = 1/5 = 0.2$$

$$P(E_4 | \text{yes}) = 6/9 = 0.667$$

$$P(E_4 | \text{no}) = 2/5 = 0.4$$

$$P(\text{yes} | E) = \frac{0.222 \cdot 0.444 \cdot 0.667 \cdot 0.668 \cdot 0.443}{P(E)} = \frac{0.028}{P(E)} \quad P(\text{no} | E) = \frac{0.6 \cdot 0.4 \cdot 0.2 \cdot 0.4 \cdot 0.357}{P(E)} = \frac{0.007}{P(E)}$$

Hence, the Naïve Bayes classifier predicts $\text{buys_computer} = \text{yes}$ for the new example.

age	income	student	credit_rating	com
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Avoiding the Zero-Probability Problem

Naïve Bayesian prediction requires each conditional prob. be **non-zero**. Otherwise, the predicted prob. will be zero

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i)$$

Ex. Suppose a dataset with 1000 tuples, income=low (0), income=medium (990), and income = high (10)

Use **Laplacian correction** (or Laplacian estimator)

Adding 1 to each case

Prob(income = low) = 1/1003

Prob(income = medium) = 991/1003

Prob(income = high) = 11/1003

The “corrected” prob. estimates are close to their “uncorrected” counterparts

Advantages & Disadvantages

➤ Advantages of Naïve Bayes Classifier:

1. Naïve Bayes is one of the fast and easy ML algorithms to predict a class of datasets.
2. It can be used for Binary as well as Multi-class Classifications.
3. It performs well in Multi-class predictions as compared to the other Algorithms.
4. It is the most popular choice for text classification problems.

➤ Disadvantages of Naïve Bayes Classifier:

1. Naive Bayes assumes that all features are independent or unrelated, so it cannot learn the relationship between features.

Baye's Classifier-Example

Confident	Studied	Sick	Result
Yes	No	No	Fail
Yes	No	Yes	Pass
No	Yes	Yes	Fail
No	Yes	No	Pass
Yes	Yes	Yes	Pass

Data set for classification

Find out whether the object with attribute **Confident** = Yes, **Sick** = No will Fail or Pass using Bayesian classification.

Solution:

The data tuples are described by the attributes *Confident*, *Studied* and *Sick*.
The class label attribute, *Result*, has two distinct values (namely, {*Pass*, *Fail*}).

Let, C1 correspond to the class *Result* = *Pass* and
C2 correspond to *Result* = *Fail*.

The tuple we wish to classify is
 $X = (\text{Confident} = \text{Yes}, \text{Sick} = \text{No})$

Formula:

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

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

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

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

Step 1: (Compute prior probability)

We need to maximize $P(X|C_i)P(C_i)$, for $i = 1, 2$. $P(C_i)$, the prior probability of each class, can be computed based on the training tuples:

$$P(\text{Result} = \text{Pass}) = 3/5 = 0.6$$

$$P(\text{Result} = \text{Fail}) = 2/5 = 0.4$$

Step 2: (Compute likelihood probability)

To compute $P(X|C_i)$, for $i = 1, 2$, we compute the following conditional probabilities:

$$P(\text{Confident} = \text{Yes} \mid \text{Result} = \text{Pass}) = 2/3 = 0.6667$$

$$P(\text{Confident} = \text{Yes} \mid \text{Result} = \text{Fail}) = 1/2 = 0.5$$

$$P(\text{Sick} = \text{No} \mid \text{Result} = \text{Pass}) = 1/3 = 0.3333$$

$$P(\text{Sick} = \text{No} \mid \text{Result} = \text{Fail}) = 1/2 = 0.5$$

Step 3: (Compute posterior probability)

$$P(X | \text{Result} = \text{Pass})$$

$$= P(\text{Confident} = \text{Yes} | \text{Result} = \text{Pass}) \times P(\text{Sick} = \text{No} | \text{Result} = \text{Pass})$$

$$= 0.6667 * 0.3333$$

$$= 0.2222$$

$$P(X | \text{Result} = \text{Fail})$$

$$= P(\text{Confident} = \text{Yes} | \text{Result} = \text{Fail}) \times P(\text{Sick} = \text{No} | \text{Result} = \text{Fail})$$

$$= 0.5 * 0.5$$

$$= 0.25$$

Step 4: (predict the class for X)

To find the class, C_i , that maximizes $P(X | C_i)P(C_i)$, we compute

$$P(X | \text{Result} = \text{Pass})P(\text{Result} = \text{Pass}) = 0.2222 \times 0.6 = 0.1333$$

$$P(X | \text{Result} = \text{Fail})P(\text{Result} = \text{Fail}) = 0.25 \times 0.4 = 0.1$$

Therefore, the naive Bayesian classifier predicts *Result = Pass* for tuple X.