

Naive Bayes

Naive Bayes is a probabilistic machine learning algorithm used primarily for solving classification problems. It is based on Bayes' theorem, which calculates the probability of an event occurring given evidence. Despite its simplicity, it's surprisingly effective in many real-world applications like text classification, spam detection, sentiment analysis and more.

Bayes Theorem: Bayes Theorem provides a way to update the probability of a hypothesis given new evidence. It is expressed as —

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

likelihood prior probability
 ↙ posterior probability ↗ marginal likelihood

- $P(H|E)$ is the posterior probability: the probability of hypothesis H being true given the evidence E.
- $P(E|H)$ is the likelihood: the probability of the evidence E given that hypothesis H is true.
- $P(H)$ is the prior probability: the initial probability of the hypothesis H before seeing the evidence.
- $P(E)$ is the marginal likelihood: the probability of evidence (E) under all possible hypotheses.

The Naive Assumption: The "naive" part of Naive Bayes assumes that all features (or predictors) are independent of each other, given the class label. This is often unrealistic in real-world scenarios but simplifies the computation, making the model fast and efficient.

Working of Naive Bayes Classifier

Given the class variable ' y_k ' and dependent feature vector x_1, x_2, \dots, x_n :

$$P(y_k|x_1, x_2, x_3, \dots, x_n) = \frac{P(x_1, x_2, x_3, \dots, x_n|y_k)P(y_k)}{P(x_1, x_2, x_3, \dots, x_n)}$$

Using our naive conditional independence assumption

$$P(x_i | y_k, x_1, x_2, \dots, x_n) = P(x_i | y_k)$$

for all such i , the relationship is simplified to

$$P(y_k | x_1, x_2, x_3, \dots, x_n) = P(y_k) \cdot \frac{\prod_{i=1}^n P(x_i | y_k)}{P(x_1, x_2, \dots, x_n)}$$

Since $P(x_1, x_2, \dots, x_n)$ is constant given the input, we can use the following classification rule:

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

$$\hat{y} = \arg \max_{y_k} P(y_k) \prod_{i=1}^n P(x_i | y_k)$$

and we can use Maximum A Posteriori (MAP) estimation to estimate $P(y_k)$ and $P(x_i | y_k)$; the former is then the relative frequency of class y_k in the training set.

$\rightarrow P(y_k)$ is the prior probability of class y_k .

$\rightarrow P(x_i | y_k)$ is the likelihood of feature x_i given class y_k . We calculate this for each class y_k and choose the class with highest posterior probability.

Eg:1 Let's understand the working of Naive Bayes Classifier using an example -

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	Normal	Weak	No
Sunny	Cool	High	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	Normal	Strong	Yes
Overcast	Hot	High	Weak	Yes
Rain	Mild	Normal	Strong	No

In this example, based on the given dataset find the probability a person will be able to play tennis given outlook = sunny, temperature = cool, humidity high, wind = strong.

Solution: We need to find

$$P(Y=\text{Yes} | \text{outlook}=\text{sunny}, \text{temperature}=\text{cool}, \text{humidity}=\text{high}, \text{wind}=\text{strong}) = [P(\text{outlook}=\text{sunny}) * P(\text{temp}=\text{cool}|\text{yes}) * P(\text{hum}=\text{high}|\text{yes}) * P(\text{wind}=\text{strong}|\text{yes})] * p(\text{yes})$$

Similarly,

$$P(Y=\text{No} | \text{out}=\text{sunny}, \text{temp}=\text{cool}, \text{humid}=\text{high}, \text{wind}=\text{strong}) = [P(\text{out}=\text{sunny}|\text{no}) * P(\text{temp}=\text{cool}|\text{No}) * P(\text{humid}=\text{high}|\text{No}) * P(\text{wind}=\text{strong}|\text{No})] * P(\text{No})$$

Step 1: We calculate the prior probabilities, $P(\text{yes})$ & $P(\text{No})$

$$P(\text{yes}) = \frac{9}{14}, \quad P(\text{No}) = \frac{5}{14}$$

Step 2: We calculate the likelihoods for each column one by one

Outlook			Temperature		
	Yes	No		Yes	No
Sunny	2/9	3/5	Hot	2/9	2/5
Overcast	4/9	0	Mild	4/9	2/5
Rain	3/9	2/5	Cool	3/9	1/5

Humidity			Wind		
	Yes	No		Yes	No
High	3/9	4/5	Weak	6/9	2/5
Normal	6/9	1/5	Strong	3/9	3/5

$$P(Y=\text{yes}) = P(\text{outlook}=\text{sunny}|\text{Yes}) * P(\text{temp}=\text{cool}|\text{yes}) * P(\text{Humid}=\text{high}|\text{yes}) * P(\text{Wind}=\text{strong}|\text{yes}) * P(\text{yes}) \\ = \frac{2}{9} * \frac{3}{8} * \frac{3}{8} * \frac{3}{14} * \frac{9}{14} = 0.0053$$

$$P(Y=\text{no}) = P(\text{out}=\text{sunny}|\text{No}) * P(\text{temp}=\text{cool}|\text{No}) * P(\text{Humid}=\text{high}|\text{No}) * P(\text{Wind}=\text{strong}|\text{No}) * P(\text{No}) \\ = \frac{3}{5} * \frac{1}{5} * \frac{4}{5} * \frac{3}{14} * \frac{5}{14} = \frac{36}{1750} = 0.0205$$

Since $P(\text{No}) > P(\text{Yes})$, therefore we will predict the person will not be able to play tennis.

Eg. 2 Given a dataset of cars being stolen or not, predict a car will get stolen given colour is Red, Type is SUV, Origin is domestic.

Colour	Type	Origin	Stolen
Red	Sports	Domestic	Yes
Red	Sports	Domestic	No
Red	Sports	Domestic	Yes
Yellow	Sports	Domestic	No
Yellow	Sports	Imported	Yes
Yellow	SUV	Imported	No
Yellow	SUV	Imported	Yes
Yellow	SUV	Domestic	No
Red	SUV	Imported	No
Red	Sports	Imported	Yes

Sln:

$$P(\text{Yes} | \text{colour}=\text{Red}, \text{Type}=\text{SUV}, \text{Origin}=\text{Domestic}) = P(\text{Yes}) *$$

$$P(\text{colour}=\text{Red} | \text{Yes}) * P(\text{Type}=\text{SUV} | \text{Yes}) * P(\text{origin}=\text{dom} | \text{Yes})$$

and, $P(\text{No} | \text{col}=\text{Red}, \text{Ty}=\text{SUV}, \text{Origin}=\text{Domestic}) = P(\text{No}) *$

$$P(\text{col}=\text{Red} | \text{No}) * P(\text{Ty}=\text{SUV} | \text{No}) * P(\text{origin}=\text{domestic} | \text{No})$$

probability

Step 1: We first calculate the prior and likelihoods

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

Colour	Yes	No	Type	Yes	No
Red	3/5	2/5	Sports	4/5	2/5
Yellow	2/5	3/5	SUV	1/5	3/5

Origin	Yes	No
Domestic	2/5	3/5
Imported	3/5	2/5

$$\therefore P(\text{Yes} | \text{all conditions}) = \frac{5}{10} * \frac{3}{5} * \frac{1}{2} * \frac{2}{5} = 0.024$$

$$P(\text{No} | \text{all conditions}) = \frac{5}{10} * \frac{2}{5} * \frac{3}{5} * \frac{3}{5} = 0.072$$

$\therefore P(\text{No}) > P(\text{Yes})$, we will say car will not be stolen.

Types of Naive Bayes Classifier : There are several types of Naive Bayes classifier depending on how the likelihood $P(X_i | Y_k)$ gets calculated.

w Gaussian Naive Bayes: These are used when the features are numerical and continuous in nature, and assume a Gaussian (Normal) distribution.

The likelihood is calculated using the probability density function of a normal distribution :

$$P(X_i | Y_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right)$$

where μ_k and σ_k^2 are mean and standard deviation of feature X_i for class Y_k .

Eg:3 Based on the given data, determine the gender of a person having height 6 ft, weight 130 pounds, and foot size 8 inches.

Height	Weight	Foot size	Gender
6.00	180	12	Male
5.92	190	11	Male
5.58	170	12	Male
5.92	165	10	Male
5.00	100	6	Female
5.50	150	8	Female
5.42	130	7	Female
5.75	150	9	Female

Step 1: We find the prior probabilities

$$P(\text{Male}) = 4/8, P(\text{Female}) = 4/8$$

Step 2: For each feature we need to calculate their likelihoods, so we will first calculate mean and variance for all features one by one for each class.

Male : Mean (Height) = $\frac{6+5.92+5.58+5.92}{4} = 5.855$

Variance (Height) = $\frac{\sum (x_i - \bar{x})^2}{n-1}$ = 0.035055

Female : Mean (Height) = $\frac{5+5.5+5.42+5.75}{4} = 5.4175$

Variance (Height) = $\frac{\sum (x_i - \bar{x})^2}{n-1}$ = $(5-5.4175)^2 + (5.5-5.4175)^2 + (5.42-5.4175)^2 + (5.75-5.4175)^2$
 $= 0.097225$

Similarly for Weight —

Male : Mean (Weight) = 176.25

Variance (Weight) = 122.92

Female : Mean (Weight) = 132.5

Variance (Weight) = 558.33

Similarly for foot size —

Male : Mean (size) = 11.25

Variance (Size) = 0.92

Female : Mean (size) = 7.5

Variance (Size) = 1.67

Now, we have to calculate

P(Gender=Male | height=6 ft, weight=130 pounds, foot size=8 inches)

and P(Gender=Female | height=6 ft, weight=130 pounds, foot size=8 inches)

$P(\text{Gender}=\text{Male}) = P(M) * P(H|M) * P(W|M) * P(S|M)$

$P(\text{Gender}=\text{Female}) = P(F) * P(H|F) * P(W|F) * P(S|F)$

Now, we need to calculate the likelihood for each numerical feature using the probability density function for a gaussian / normal distribution :

$$\therefore P(H|M) = \frac{1}{\sqrt{2\pi \cdot 0.035}} * e^{-\frac{(6-5.855)^2}{2 \cdot 0.035}} = 1.5789$$

$P(W|M) = 5.9881 e^{-6}, P(S|M) = 1.3112 e^{-3}$

$P(H|F) = 2.2346 e^{-1}, P(W|F) = 1.6789 e^{-2}, P(S|F) = 2.8669 e^{-4}$

$\therefore P(M) = 0.5 * 1.5789 * 5.9881 e^{-6} * 1.3112 e^{-3} = 6.1984 e^{-9}$

$P(F) = 0.5 * 2.2346 e^{-1} * 1.6789 e^{-2} * 2.8669 e^{-4} = 5.377 e^{-9}$

Since $P(F) > P(M)$, so the new example will be classified as Female.

(b) Multinomial Naive Bayes: It is typically used for discrete data, such as word counts in text classification. The likelihood $P(X_i | C_k)$ is proportional to the frequency of feature X_i within class C_k .

(c) Bernoulli Naive Bayes: It is suitable for binary / boolean features or distribution which follows bernoulli's distribution. The likelihood $P(X_i | C_k)$ represents the probability of feature X_i being 1 (present) or 0 (absent) in class C_k .

* Laplace Smoothing: Laplace smoothing (additive smoothing) is a technique used to handle the zero probability problem. It involves adding a small constant (usually 1) to the count of each feature during likelihood estimation to ensure that no probability is ever zero.

For Eg. if we are estimating the likelihood of a word in text classification :

$$P(\text{word} | C_k) = \frac{\text{count of (word in class } C_k)}{\text{total words in class } C_k + V} + 1$$

where V is the total number of unique words (vocab size).

Eg: Given the dataset as shown below, predict the output value for new instance where outlook=overcast, temperature=60, humidity=62 and windy=false.

Outlook	Temp	Humidity	Wind	Play Tennis
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	70	96	False	Yes
Rainy	68	80	False	Yes
Rainy	65	70	True	No
Overcast	64	65	True	Yes
Sunny	72	95	False	No
Sunny	69	70	False	Yes
Rainy	75	80	False	Yes
Sunny	75	70	True	Yes
Overcast	72	90	True	Yes
Overcast	81	75	False	Yes
Rainy	71	91	True	No

Soln: We need to calculate the following

$$P(\text{Yes} | \text{outlook} = \text{overcast}, \text{temp} = 60, \text{humid} = 62, \text{wind} = \text{False})$$

$$\& P(\text{No} | \text{out} = \text{overcast}, \text{temp} = 60, \text{humid} = 62, \text{wind} = \text{False}).$$

Step 1: We calculate the prior probabilities,

$$P(\text{Yes}) = \frac{9}{14}, P(\text{No}) = \frac{5}{14}$$

Step 2: Now, we calculate the likelihood for each feature.

	Outlook	Yes	No	Wind	Yes	No
Sunny	2/9	3/5	True	3/9	3/5	
Overcast	4/9	0/5	False	6/9	2/5	
Rainy	3/9	2/5				

From these we can get,

$$P(\text{Sunny} | \text{Yes}) = \frac{2}{9}, P(\text{Sunny} | \text{No}) = \frac{3}{5}, P(\text{Rainy} | \text{Y}) = \frac{3}{9}, P(\text{Rainy} | \text{N}) = \frac{2}{5}$$

$$P(\text{W}_T | \text{Y}) = \frac{3}{9}, P(\text{W}_F | \text{Y}) = \frac{6}{9}, P(\text{W}_T | \text{N}) = \frac{3}{5}, P(\text{W}_F | \text{N}) = \frac{2}{5}$$

$$\text{But } P(\text{Overcast} | \text{Yes}) = 0 \& P(\text{Overcast} | \text{No}) = 0.$$

Now, when $P(\text{Overcast} | \text{No})$ becomes 0, whole

$P(\text{No} | \text{any condition}) = 0$, to overcome this limitation

we make use of the Laplace Smoothing and add a constant

value 1.

$$\therefore P(\text{Overcast} | \text{No}) = \frac{0+1}{5+3} = \frac{1}{8}$$

$$P(\text{Overcast} | \text{Yes}) = \frac{4+1}{9+3} = \frac{5}{12} \rightarrow \text{since there are 3 unique categories in outlook}$$

Now, we need to calculate the prior probabilities for the numerical columns as well, so we first find the mean and standard dev. variance of both the features.

$$\text{Mean}(\text{Temp} | \text{Yes}) = 73, \text{Variance}(\text{Temp} | \text{Y}) = 6.2$$

$$\text{Mean}(\text{Temp} | \text{No}) = 74.6, \text{Variance}(\text{Temp} | \text{N}) = 8.0$$

$$\text{Mean}(\text{Hum} | \text{Yes}) = 79.1, \text{Variance}(\text{Hum} | \text{Y}) = 10.2$$

$$\text{Mean}(\text{Hum} | \text{No}) = 86.2, \text{Variance}(\text{Hum} | \text{N}) = 9.7$$

Step 3: We now need to calculate -

$$p(\text{temp} = 60 \mid \text{Yes}) = \frac{1}{6.2\sqrt{2\pi}} e^{-\frac{(60-73)^2}{2*(6.2)^2}} = 0.071$$

$$p(\text{temp} = 60 \mid \text{No}) = \frac{1}{8\sqrt{2\pi}} e^{-\frac{(60-74.6)^2}{2*8^2}} = 0.0094$$

$$p(\text{humid} = 62 \mid \text{Yes}) = \frac{1}{10.2\sqrt{2\pi}} e^{-\frac{(62-79.1)^2}{2*(10.2)^2}} = 0.0096$$

$$p(\text{humid} = 62 \mid \text{No}) = \frac{1}{9.7\sqrt{2\pi}} e^{-\frac{(62-86.2)^2}{2*(9.7)^2}} = 0.0018$$

Step 4: $P(\text{Yes} \mid E) = P(\text{out=overcast} \mid \text{Yes}) * P(\text{temp}=60 \mid \text{Yes}) * P(\text{humid}=62 \mid \text{Yes})$
 $* P(\text{wind=false} \mid \text{Yes}) * P(\text{Yes})$

$$= \frac{5}{12} * 0.071 * 0.0096 * \frac{6}{9} * \frac{9}{14} = 1.22 * 10^{-5}$$

$$P(\text{No} \mid E) = P(\text{out=overcast} \mid \text{No}) * P(\text{temp}=60 \mid \text{No}) * P(\text{humid}=62 \mid \text{No})$$
 $* P(\text{wind=false} \mid \text{No}) * P(\text{No})$
 $= \frac{1}{8} * 0.0094 * 0.0018 * \frac{2}{5} * \frac{5}{14} = 3.02 * 10^{-7}$

Since $P(\text{Yes} \mid E) > P(\text{No} \mid E)$, so we classify the new instance as Yes.

Advantages of Naive Bayes :

1. It is simple to implement and computationally efficient, especially for large datasets.
2. Training and prediction are fast due to the independence assumption and closed-form solutions for likelihoods.
3. Despite the naive assumption, it often performs well even on small amount of data.

Disadvantages of Naive Bayes :

1. The assumption that features are independent given the class label is rarely true in real-world data. Violations of this assumption can degrade performance.
2. If a certain feature X_i never appears in the training data for a particular class, then $P(X_i \mid k)$ will be zero, making entire product zero. This is prevented by using Laplace Smoothing.
3. For continuous data, the Gaussian assumption might not always hold.