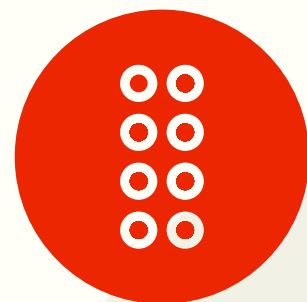


Naive Bayes Classifier



A simple yet powerful probabilistic model that predicts class membership based on Bayes' theorem, assuming features are independent of each other.

19 OCT, 2025



BAYESIANBRAT*



Naive Bayes

Naive Bayes classifier is a simple supervised learning model.

Supervised learning means the model is trained on labeled data (i.e., input features along with their correct output, class).

Naive Bayes is called Naive because it assumes that all features are independent of each other.

Example: If we are predicting whether a patient is sick or healthy based on attributes like "fever", "cough", and "headache", Naive Bayes assumes these symptoms are independent of each other given the class (though in reality, they may be co-related).

Given values of predictor attribute $x = (x_1, x_2, \dots, x_m)$ predict the class of target attribute y .

y is one of the classes $y \in \{\text{improved, stable, degraded}\}$ or $y \in \{0, 1\}$

Classifiers commonly assign to most likely class given a data instance = maximum a posteriori class.

- Given instance data x , assign to class y where $P(Y=y | x=x)$ is maximum.

- Mathematically, we calculate the probability of each class given the features, i.e., $P(Y=y | x=x)$

- Then, we choose the class with the highest probability.

- That is, $\hat{y} = \operatorname{argmax}_y P(Y=y | x=x)$

- Assume that each x_m is categorical, ~~say~~ such as $x_1 \in \{\text{dry, wet, snowy}\}$ and $x_2 \in \{\text{medium, rough, smooth}\}$

- Direct estimate of $P(Y=y | x=x) = P(Y=y | x_1=\text{dry}, x_2=\text{medium})$ can be calculated if predictor count M is small.

Table : Example weather and surface

| Count of delay Weather | surface | delay major | delay minor | delay none |
|---------------------------|---------|----------------|----------------|---------------|
| dry | medium | 3 | 6 | 1 |
| dry | rough | 0 | 3 | 0 |
| dry | smooth | 0 | 2 | 1 |
| snow | medium | 0 | 2 | 0 |
| snow | rough | 0 | 0 | 0 |
| snow | smooth | 0 | 0 | 0 |
| wet | medium | 1 | 0 | 0 |
| wet | rough | 1 | 0 | 0 |
| wet | smooth | 0 | 0 | 0 |
| Grand | Total | 5 | 13 | 2 |

$W = \text{Weather}$ $S = \text{Surface}$

Assign to class with maximum class probability estimates.

$$P(Y = \text{major} \mid W = \text{dry}, S = \text{medium}) = 3/10$$

$$P(Y = \text{minor} \mid W = \text{dry}, S = \text{medium}) = 6/10$$

$$P(Y = \text{none} \mid W = \text{dry}, S = \text{medium}) = 1/10$$

However ...

Suppose that each predictor x_m has 4 possible values and $M = 30$

The count of possible unique indicator values is $4^{30} = 2^{60} = (2^{10})^6 \approx (10^3)^6 = 10^{18}$

Try a different approach: Bayes Theorem

$$P(Y = \text{major} \mid W = \text{dry}, S = \text{medium}) =$$

$$\frac{P(W = \text{dry}, S = \text{medium} \mid Y = \text{major}) P(Y = \text{major})}{P(W = \text{dry}, S = \text{medium})}$$

$$P(Y = y \mid X = x) = \frac{P(X = x \mid Y = y) P(Y = y)}{P(X = x)}$$

$$P(Y = y \mid X = x) = \frac{P(X = x \mid Y = y) P(Y = y)}{P(X = x)}$$

Assign instance to class with maximum class probability estimate computed
Bayes Theorem

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(X = x \mid Y = y) P(Y = y)$$

Why is $P(X = x)$ dropped?

Make a strong assumption: X 's are conditionally independent given $Y = y$ so that

$$P(W = \text{dry}, S = \text{medium} \mid Y = \text{major}) =$$

$$\frac{P(W = \text{dry} \mid Y = \text{major}) \times P(S = \text{medium} \mid Y = \text{major})}{P(Y = \text{major})}$$

$$P(X = x \mid Y = y) = \frac{P(X_1 = x_1 \mid Y = y) \times P(X_2 = x_2 \mid Y = y) \times \dots \times P(X_n = x_n \mid Y = y)}{P(Y = y)}$$

Each conditional probability $P(X_m = x_m | Y=y)$ is easy to estimate from data as

$$P(W = \text{dry} | Y = \text{major}) \approx 3/5$$

$$P(W = \text{wet} | Y = \text{major}) \approx 2/5$$

$$P(W = \text{snow} | Y = \text{major}) = ?$$

$$P(W = \text{dry} | Y = \text{minor}) = 11/13$$

$$P(W = \text{wet} | Y = \text{minor}) = 0/13$$

$$P(W = \text{snow} | Y = \text{minor}) = 2/13$$

$$P(W = \text{dry} | Y = \text{none}) = 2/2$$

$$P(W = \text{wet} | Y = \text{none}) = ?$$

$$P(W = \text{snow} | Y = \text{none}) = ?$$

$P(Y=y)$ estimator is the proportion of data instances with $Y=y$

$$P(X_1 = x_1 | Y=y) \dots P(X_M = x_M | Y=y) \underline{P(Y=y)}$$

Underline a reminder to include the a priori probability that $Y=y$

$$P(Y = \text{major} | W = \text{wet}, S = \text{medium}) \propto \\ P(W = \text{wet} | Y = \text{major}) \times P(S = \text{medium} | Y = \text{major}) \\ \times P(Y = \text{major})$$

$$= \frac{2}{5} \times \frac{4}{5} \times \frac{5}{20}$$

$$= \frac{40}{500}$$

$$P(Y = \text{minor} | W = \text{wet}, S = \text{medium}) \propto \\ P(W = \text{wet} | Y = \text{minor}) \times P(S = \text{medium} | Y = \text{minor}) \times P(Y = \text{minor})$$

$$= \frac{0}{13} \times \frac{8}{13} \times \frac{13}{20}$$

$$P(Y = \text{none} | W = \text{wet}, S = \text{medium}) \propto P(W = \text{wet} | Y = \text{none}) \times P(S = \text{medium} | Y = \text{none}) \times P(Y = \text{none}) \\ = 0$$

Assign the class $\hat{y} = \text{major}$

→ The problem here is, when we have a single predictor that does the whole probability zero, puts a huge dent in the calculation.

→ There's a need of some smoothing technique

Applying law of total probability and, using denominator.

$$P(X = x) = \sum_{y'} P(X = x | Y = y') P(Y = y')$$

$$P(Y = \text{major} | X) = \frac{0.08}{P(X = x)} = \frac{0.08}{\sum_{y'} P(X = x | Y = y') P(Y = y')} \\ = \frac{0.08}{0.08} = 1$$

$$P(Y = \text{minor} | X) = 0 \\ P(Y = \text{none} | X) = 0$$

Laplace Smoothing:

→ What happens if the estimate of $P(X_m = x_m | Y = y) = 0$? More likely with many predictors.

→ Software commonly adjusts with Laplace smoothing. Suppose that x has V unique values x_1, x_2, \dots, x_V .

→ Laplace smoothing adds α to each count $N(X = x_i)$ to remove zero counts.

→ In simpler terms:

If one of these probabilities = 0,
We then multiply

$$P(X_1 | Y) \cdot P(X_2 | Y) \dots P(Y)$$

The whole product becomes zero.

That class can never be chosen even if the class is actually possible.

This is called the zero frequency problem.

To fix this, we pretend, we have seen at each possible category at least once.

We add small constant α (often 1)
to each count

Then we re-normalize

Suppose feature x has V possible values (x_1, \dots, x_V)

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

With Laplace smoothing:

$$P(x = x_v | y = y) = \frac{N(x = x_v, y = y) + \alpha}{N(y = y) + \alpha V}$$

Numerator: We add α

Denominator: We add αV (because we added α to each V possible values).

Example:

For class = "none"

Without smoothing: $P(W = \text{wet} | Y = \text{none}) = \frac{0}{2}$

With Laplace ($\alpha = 1$, and because it has 3 categories: dry, wet, snow)

$$P(W = w_{set} | Y = none) = \frac{0+1}{2+3}$$

$$= \frac{1}{5} = 0.2$$

Smoothing prevents classes from being completely wiped out.

Especially important when you have many features and categories, since some combinations will be missing in training data.

Example:

1. Without smoothing (raw counts)

Imagine a categorical feature x with 3 possible values

a, b, c

Suppose for class $Y = y$ in training data we saw

$$N(x = a, Y = y) = 2$$

$$N(x = b, Y = y) = 0$$

$$N(x = c, Y = y) = 1$$

$$\text{Total } N(Y = y) = 3$$

$$P(a|y) = \frac{2}{3} = 0.667$$

$$P(b|y) = \frac{0}{3} = 0$$

$$P(c|y) = \frac{1}{3} = 0.333$$

2. Doing a Laplace smoothing:

$$P(X = x_v | Y = y) =$$

$$\frac{N(X = x_v, Y = y) + \alpha}{N(Y = y) + \alpha v}$$

$$\text{Here } v = 3, \quad N(Y = y) = 3 \quad \alpha = 1$$

$$a = \frac{2+1}{3+3} = \frac{3}{6} = 0.5$$

$$b = \frac{0+1}{3+3} = \frac{1}{6} = 0.167$$

$$c = \frac{1+1}{3+3} = \frac{2}{6} = 0.333$$

Not all probabilities are nonzero and still sum to 1.