

Naive Bayes Classifier

A naive Bayes classifier is a simple probabilistic classifier based on applying Bayes' theorem with strong (naive) independence assumptions. It assumes the conditional independence of attribute values given the class:

$$p(v_1, v_2, \dots, v_n | c) = \prod_i p(v_i | c)$$

Naive Bayes formula:

$$p(c | v_1, v_2, \dots, v_n) = p(c) \cdot \prod_i \frac{p(c | v_i)}{p(c)}$$

Classifying a new instance (v_1, v_2, \dots, v_n)

Let's say that our dataset has m classes (c_1, c_2, \dots, c_m) (target variable with m values). The Naive Bayes classifier calculates for each class c_i the conditional probability of class c_i given evidence (v_1, v_2, \dots, v_n)

$$p(c_i | v_1, v_2, \dots, v_n)$$

according to the naive Bayes formula. It classifies the example into the class with the highest probability.

Example

Will the spider catch an ant?

Past experiences of the spider catching ants:

Color	Size	Time	Caught
black	large	day	YES
white	small	night	YES
black	small	day	YES
red	large	night	NO
black	large	night	NO
white	large	night	NO

Ant 1: Color = white, Time = night

$$v_1 = \text{"Color = white"}$$

$$v_2 = \text{"Time = night"}$$

$$c_1 = YES$$

$$c_2 = NO$$

$$p(c_1|v_1, v_2) = (1)$$

$$p(\text{Caught} = YES | \text{Color} = \text{white}, \text{Time} = \text{night}) = (2)$$

$$p(\text{Caught} = YES) * \frac{p(\text{Caught} = YES | \text{Color} = \text{white})}{p(\text{Caught} = YES)} * \frac{p(\text{Caught} = YES | \text{Time} = \text{night})}{p(\text{Caught} = YES)} = (3)$$

$$\frac{1}{2} * \frac{1}{2} * \frac{1}{4} = \frac{1}{4} (4)$$

$$p(c_2|v_1, v_2) = (5)$$

$$p(\text{Caught} = NO | \text{Color} = \text{white}, \text{Time} = \text{night}) = (6)$$

$$p(\text{Caught} = NO) * \frac{p(\text{Caught} = NO | \text{Color} = \text{white})}{p(\text{Caught} = NO)} * \frac{p(\text{Caught} = NO | \text{Time} = \text{night})}{p(\text{Caught} = NO)} = (7)$$

$$\frac{1}{2} * \frac{1}{2} * \frac{3}{4} = \frac{3}{4} (8)$$

The spider will not catch the white ant at night because $p(\text{Caught} = NO | \text{Color} = \text{white}, \text{Time} = \text{night}) > p(\text{Caught} = YES | \text{Color} = \text{white}, \text{Time} = \text{night})$.

Ant 2: Color = black, Size = large, Time = day

$$v_1 = \text{“Color = black”}$$

$$v_2 = \text{“Size = large”}$$

$$v_3 = \text{“Time = day”}$$

$$c_1 = YES$$

$$c_2 = NO$$

$$p(c_1|v_1, v_2, v_3) = \quad (9)$$

$$p(\text{Caught} = YES | \text{Color} = \text{black}, \text{Size} = \text{large}, \text{Time} = \text{day}) = \quad (10)$$

$$p(\text{Caught} = YES) * \frac{p(\text{Caught} = YES | \text{Color} = \text{black})}{p(\text{Caught} = YES)} * \dots \quad (11)$$

$$\dots * \frac{p(\text{Caught} = YES | \text{Size} = \text{large})}{p(\text{Caught} = YES)} * \frac{p(\text{Caught} = YES | \text{Time} = \text{day})}{p(\text{Caught} = YES)} = \quad (12)$$

$$\frac{1}{2} * \frac{2}{3} * \frac{1}{4} * \frac{1}{2} = \frac{2}{3} \quad (13)$$

$$p(c_2|v_1, v_2, v_3) = \quad (14)$$

$$p(\text{Caught} = NO | \text{Color} = \text{black}, \text{Size} = \text{large}, \text{Time} = \text{day}) = \quad (15)$$

$$p(\text{Caught} = NO) * \frac{p(\text{Caught} = NO | \text{Color} = \text{black})}{p(\text{Caught} = NO)} * \dots \quad (16)$$

$$\dots * \frac{p(\text{Caught} = NO | \text{Size} = \text{large})}{p(\text{Caught} = NO)} * \frac{p(\text{Caught} = NO | \text{Time} = \text{day})}{p(\text{Caught} = NO)} * = \quad (17)$$

$$\frac{1}{2} * \frac{1}{3} * \frac{3}{4} * \frac{0}{2} = 0 \quad (18)$$

The spider will catch the large black ant at night because $p(\text{Caught}=\text{YES} | \text{Color} = \text{black}, \text{Size} = \text{large}, \text{Time} = \text{day}) > p(\text{Caught}=\text{NO} | \text{Color} = \text{black}, \text{Size} = \text{large}, \text{Time} = \text{day})$.

To think over:

When calculating probabilities $p(c_1|v_1, v_2)$ and $p(c_2|v_1, v_2)$ for a two class problem using naive Bayes formula, the probabilities sometimes do not sum up to 1: $p(c_1|v_1, v_2) + p(c_2|v_1, v_2) \neq 1$. Why? ¹

¹Petra Kralj Novak, Petra.Kralj@ijs.si, <http://kt.ijs.si/PetraKralj/IPSKnowledgeDiscovery0809.html>