

Naive Bayes

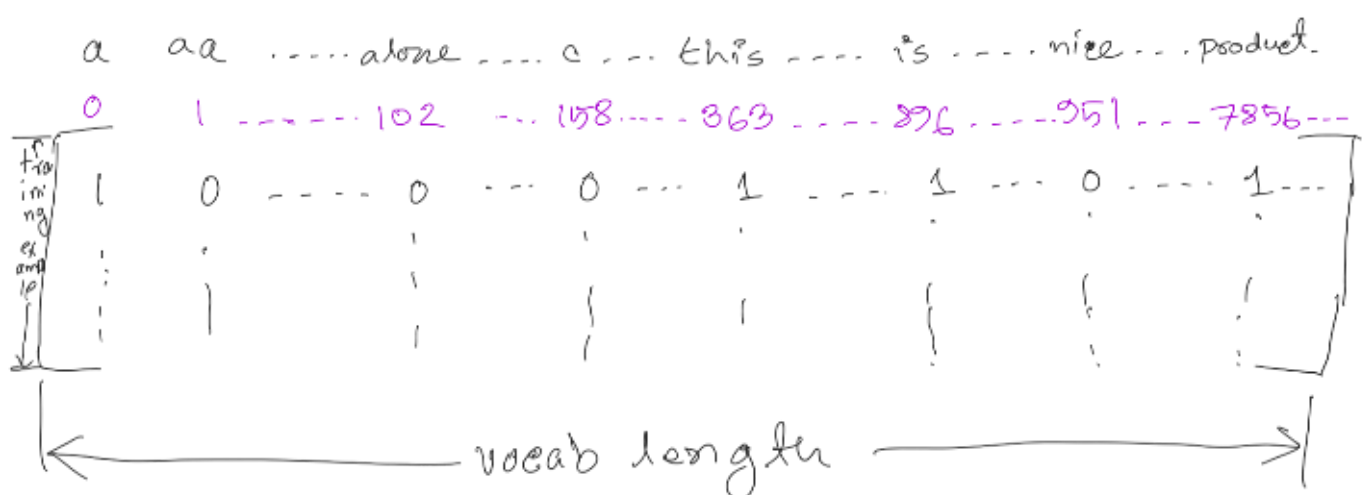
In GDA input values are continuous but in Naive-Bayes input values are discrete. Let's say we want to build a review classifier where our model will predict if a given review is positive(1) or negative(0). First we will convert our input into vector. For that we need to make a vocabulary which is nothing but a number representation of words.

This is a high quality product \rightarrow positive
 $[363 \quad 896 \quad 0 \quad 3392 \quad 1001 \quad 7856] \quad [1]$

In our vocabulary word "This" is represented by number 363

We will convert this vector as one hot encoding.

In one hot encoding a vector will have the same length as the length of vocabulary. If a particular word is present that position will have value 1. And rest of them will have value zero.



Having chosen our feature vector, we now want to build a discriminative model. So, we have to model $p(x|y)$. But if we have, say, a vocabulary of 50000 words, then $x \in \{0, 1\}^{50000}$ (x is a 50000-dimensional vector of 0's and 1's), and if we were to model x explicitly with a multinomial distribution over the 250000 possible outcomes, then we'd end up with a

(250000-1)-dimensional parameter vector. This is clearly too many parameters. I

To reduce this huge parameters we will make a very strong assumption that is " x_i 's are conditionally independent given y ". If we wouldn't mention conditionally that would mean all the input features are independent with each other. But the real assumption is input features are independent only for a particular output. Because of this assumption this algorithm is called Naive Bayes.

Bayes rule is:-

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

To predict $P(y|x)$ we have to find each of the value :-

$P(x|y)$

$$P(x_1, \dots, x_{50,000} | y)$$

$$= P(x_1|y) P(x_2|y, x_1) P(x_3|y, x_1, x_2) \dots P(x_{50,000}|y, x_1, \dots, x_{49,999})$$

$$\Rightarrow P(x_1|y) P(x_2|y) \dots P(x_{50,000}|y)$$

Because of the naive assumption.

$$= \prod_{i=1}^n P(x_i|y) \quad \left] \rightarrow \text{We need to calculate this for each of the output label.} \right.$$

For positive review and word "product" :-

$$n - P(\text{product} | 1, \dots, 1) \quad \# \text{ occurrence of "product" in all reviews}$$

$$Y_{\text{product}}(y=1) = 1 \text{ (if } y=1 \text{)} = \frac{\text{review}}{\# \text{ numbers of positive review}}$$

Generally,

$$\phi_j | y=k = \frac{\sum_{i=1}^m 1 \{x_j^{(i)} = 1 \text{ and } y=k\}}{\sum_{i=1}^m 1 \{y^{(i)} = k\}} \quad \left. \begin{array}{l} m = \text{number} \\ \text{of training} \\ \text{example.} \end{array} \right\}$$

$P(y)$

this is the class probability. That means probability of a class occurrence without seeing the input features. This is also called the prior probability.

Prior probability for class k is.

$$P(y=k) = \frac{\text{Number of occurrences of class 'k'}}{\text{total training example.}}$$

$$= \frac{\sum_{i=1}^m 1 \{y^{(i)} = k\}}{m}$$

$P(x)$

$$P(x) = \sum_{i=1}^K \prod_{j=1}^n P(x_j | y=k)$$

Making prediction:-

We will calculate the prediction for each of the output class. for class ' k '

$$P(y=k | x) = \frac{P(x | y=k) P(y=k)}{P(x)}$$

whichever class which have the highest priority will consider that as output label.

We can apply naive bayes for continuous value input by discretize it.

Living area	< 400	$400 - 800$	$800 - 1200$	$1200 - 1600$	> 1600
x_i	1	2	3	4	5

Laplace smoothing :-

Let say we want to classify the review

"This product is waste of money".

Our training example didn't have any word "waste". So,

$$P(x_{\text{waste}} | y = k) = 0 \quad \text{for all value of 'k'}$$

$$P(x | y = k) = P(x_0 | y = k) P(x_1 | y = k) \dots \underbrace{P(x_{\text{waste}} | y = k)}_{\downarrow 0} \dots$$

$$\therefore P(y | x) = \frac{\overset{=0}{\overbrace{P(x | y)}^{\rightarrow 0}} P(y)}{\underbrace{P(x)}_{\rightarrow 0}}$$

$$= \frac{0}{0}$$

= undefined.

Laplace smoothing is used to solve this problem.

$$P(y=k) = \frac{\sum_{i=1}^m 1\{y^{(i)}=k\} + 1}{n + n_{\text{number of output class}}}$$

$$P(x_j | y=k) = \frac{\sum_{i=1}^m 1\{y^{(i)}=k \text{ and } x_j^{(i)}=1\} + 1}{\sum_{i=1}^m 1\{y^{(i)}=k\} + n}$$

Even if training example doesn't see any word it won't set it's probability to zero. Instead it will make that probability very small number.