

Naive Bayes Examples



The Naïve Bayes Algorithm:

A diagram showing the Bayes' Theorem equation with arrows pointing from labels to its components. The equation is $P(c | x) = \frac{P(x | c)P(c)}{P(x)}$. An arrow points from 'Likelihood' to $P(x | c)$. An arrow points from 'Class Prior Probability' to $P(c)$. An arrow points from 'Posterior Probability' to $P(c | x)$. An arrow points from 'Evidence (Normalization)' to $P(x)$.

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

Labels and arrows:

- Likelihood points to $P(x | c)$
- Class Prior Probability points to $P(c)$
- Posterior Probability points to $P(c | x)$
- Evidence (Normalization) points to $P(x)$

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

Bayes Theorem provides a way of calculating the poster probability: $P(c|x)$.

Naïve Bayes classifier assumes the effect of predictor x is **independent** of the values of other predictors. This assumption is called conditional independence.

Naïve Bayes:

Example : SPAM vs. HAM

D1	send us your password	SPAM
D2	send us your review	HAM
D3	review your password	HAM
D4	review us	SPAM
D5	send your password	SPAM
D6	send us your account	SPAM

New Document.. Spam or Ham?

D7	review us	?
----	-----------	---

Naïve Bayes:

Example : SPAM vs. HAM

SPAM	HAM	
2/4	1/2	password
1/4	2/2	review
3/4	1/2	send
3/4	1/2	us
3/4	1/2	your
1/4	0/2	account

$$P(\text{review us} \mid \text{spam}) = P(0,1,0,1,0,0 \mid \text{spam}) = (1-2/4)(1/4)(1-3/4)(3/4)(1-3/4)(1-1/4) = 0.00439$$

$$P(\text{review us} \mid \text{ham}) = P(0,1,0,1,0,0 \mid \text{ham}) = (1-1/2)*1*(1-1/2)*(1/2)*(1-1/2)*(1-0) = 0.0625$$

$$P(\text{ham} \mid \text{review us}) = ((0.625) * 1/3) / (0.0625 * 1/3 + 0.0044 * 2/3) = 0.87$$

Naïve Bayes: How do we handle words with zero probabilities?

- Method 1: Additive smoothing
 - Add a constant δ to the counts of each word

The diagram shows the formula for Laplace smoothing:
$$p(w | d) = \frac{c(w, d) + 1}{|d| + |V|}$$
 The formula is highlighted in yellow. Annotations include: an arrow from "Counts of w in d" pointing to $c(w, d)$; an arrow from "“Add one”, Laplace smoothing" pointing to the $+1$ in the numerator; an arrow from "Vocabulary size" pointing to $|V|$ in the denominator; and an arrow from "Length of d (total counts)" pointing to $|d|$ in the denominator.

For the previous example : account

Previous: $P(\text{account} | \text{ham}) = 0 / 2$

Now : $P(\text{account} | \text{ham}) = 0 + 1 / (2 + 6)$

Naïve Bayes: The Multinomial Approach

	Document	Class
	1 Hamburger NYC Hamburger	A
	2 Hamburger Hamburger Texas	A
	3 Hamburger Cheeseburger	A
	4 Montreal Toronto Hamburger	C
5	Hamburger Hamburger Hamburger Toronto Montreal	??

1) Develop our likelihoods for 'seeing' each word in given the class:

$$P(\text{Hamburger} | A) = (6 + 1) / (8 + 6) =$$

times
'hamburger'
appeared in
'A' Docs

smoothing

of words
in 'A' docs

of unique
words in all
docs

Naïve Bayes: The Multinomial Approach

	Document	Class
	Hamburger NYC 1 Hamburger	A
	Hamburger Hamburger 2 Texas	A
	3 Hamburger Cheeseburger	A
	Montreal Toronto 4 Hamburger	C
5	Hamburger Hamburger Hamburger Toronto Montreal	??

1) Develop our likelihoods for 'seeing' each word in given the class:

$$P(\text{Hamburger} \mid A) = (5 + 1) / (8 + 6) = 3/7$$

$$P(\text{Toronto} \mid A) = (0 + 1) / (8 + 6) = 1/14$$

$$P(\text{Montreal} \mid A) = (0 + 1) / (8 + 6) = 1/14$$

$$P(\text{Hamburger} \mid C) = (1 + 1) / (3 + 6) = 2/9$$

$$P(\text{Toronto} \mid C) = (1 + 1) / (3 + 6) = 2/9$$

$$P(\text{Montreal} \mid C) = (1 + 1) / (3 + 6) = 2/9$$

$$P(A \mid d5) \sim (3/4) * ((1/14)^2) (3/7)^3 = 0.0003$$

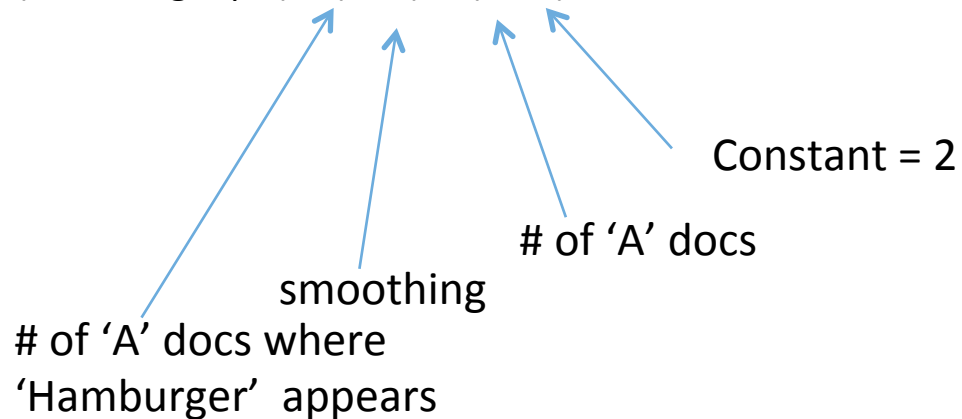
$$P(C \mid d5) \sim (1/4) * (2/9)^3 * (2/9)(2/9) = 0.0001$$

Naïve Bayes: The Bernouli Approach

	Document	Class
	1 Hamburger NYC Hamburger	A
	2 Hamburger Hamburger Texas	A
	3 Hamburger Cheeseburger	A
	4 Montreal Toronto Hamburger	C
5	Hamburger Hamburger Hamburger Toronto Montreal	??

1) Develop our likelihoods for 'seeing' each word in given the class:

$$P(\text{Hamburger} | A) = (3+1) / (3+ 2) = 4/5$$



Naïve Bayes: The Bernouli Approach

	Document	Class
	1 Hamburger NYC Hamburger	A
	2 Hamburger Hamburger Texas	A
	3 Hamburger Cheeseburger	A
	4 Montreal Toronto Hamburger	C
5	Hamburger Hamburger Hamburger Toronto Montreal	??

1) Develop our likelihoods for 'seeing' each word in given the class:

$$P(\text{Hamburger} \mid A) = (3+1)/(3+2) = 4/5$$

$$P(\text{NYC} \mid A) = P(\text{Texas} \mid A) = P(\text{Cheeseburger} \mid A) = (1+1)/5 = 2/5$$

$$P(\text{Montreal} \mid A) = P(\text{Toronto} \mid A) = (0+1)/5$$

$$P(\text{Hamburger} \mid C) = P(\text{Toronto} \mid C) = P(\text{Montreal} \mid C) = (1+1)/(1+2) = 2/3$$

$$P(\text{NYC} \mid C) = P(\text{Texas} \mid C) = P(\text{Cheeseburger} \mid C) = 1/3$$

$$P(A \mid d_5) \sim (3/4) * (4/5) (1/5)^2 (1-2/5)^3 = .005$$

$$P(C \mid d_5) \sim (1/4) * (2/3)^3 (1-1/3)^3 = .022$$

Naïve Bayes

Bernoulli : models the fraction of documents of class C that contain the word 'w' (ignores number of occurrences)

Vs.

Multinomial: models the fraction of *positions* in documents of class C that contain the word 'w' (keeps track of number of occurrences)

But why does Naïve Bayes work so well-
(Considering that it is Naïve) ?

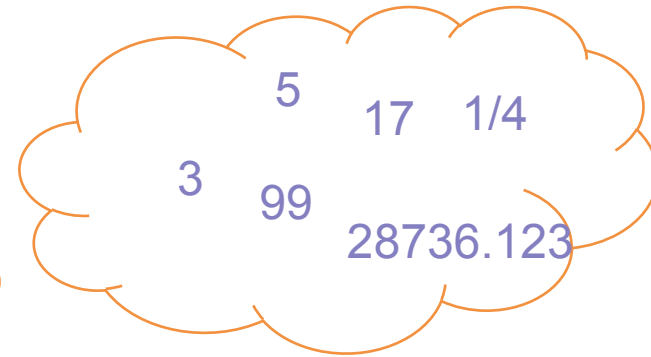
NB chooses among possible classes to find the class with the highest associated probability.

Naiveté doesn't hurt, because correctness is based on classification, not prediction

Advantages of Naïve Bayes:

- Simple & Fast. Just doing a bunch of counts!
- Will converge quickly. Requires less training data
- Can handle sparse matrices
- Can handle multiple classes well

How about numeric features?



Naïve Bayes: The Gaussian Approach

$$p(h_x|c) = \frac{1}{\sqrt{2\pi} \sigma_{h,c}^2} \exp - \frac{1}{2} \left(\frac{(h_x - \mu_{h,c})^2}{\sigma_{h,c}^2} \right)$$

$$p(w_x|c) = \frac{1}{\sqrt{2\pi} \sigma_{w,c}^2} \exp - \frac{1}{2} \left(\frac{(w_x - \mu_{w,c})^2}{\sigma_{w,c}^2} \right)$$

$$p(h_x|a) = \frac{1}{\sqrt{2\pi} \sigma_{h,a}^2} \exp - \frac{1}{2} \left(\frac{(h_x - \mu_{h,a})^2}{\sigma_{h,a}^2} \right)$$

$$p(w_x|a) = \frac{1}{\sqrt{2\pi} \sigma_{w,a}^2} \exp - \frac{1}{2} \left(\frac{(w_x - \mu_{w,a})^2}{\sigma_{w,a}^2} \right)$$

$$P(x|a) = p(h_x|a) p(w_x|a)$$

$$P(a|x) \sim P(x|a) * P(a)$$

