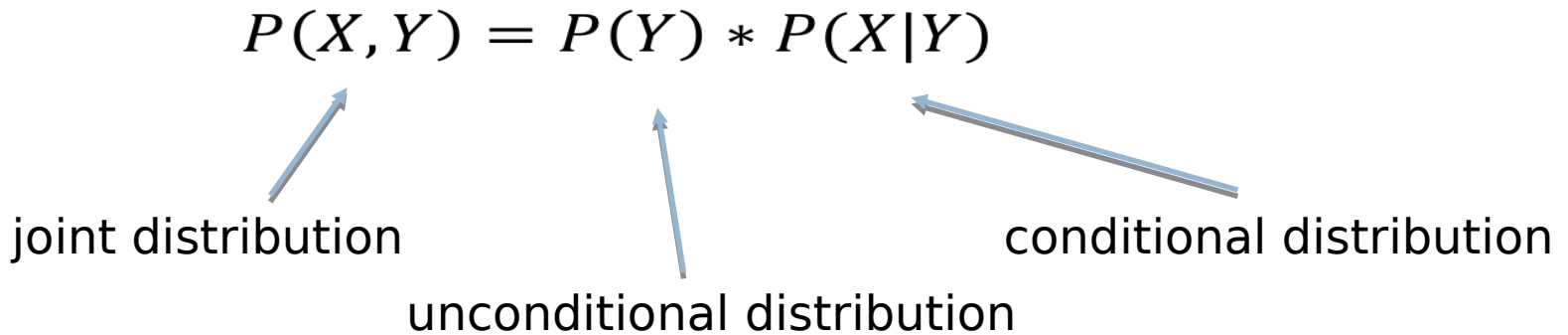


NAÏVE BAYES

Joseph C. Osborn
CS 51A - Spring 2020

Relationship between distributions

$$P(X, Y) = P(Y) * P(X|Y)$$


joint distribution

unconditional distribution

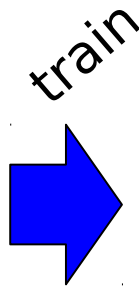
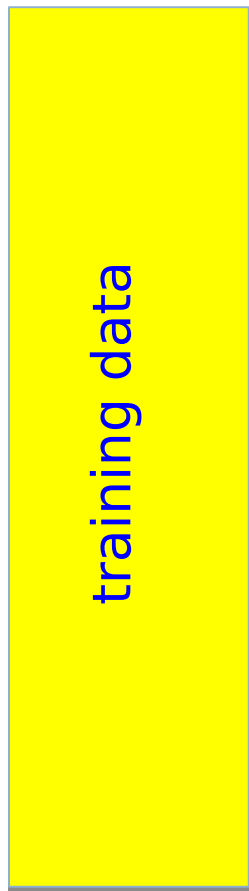
conditional distribution

Can think of it as describing the two events happening in two steps:

The likelihood of X and Y happening:

1. How likely it is that Y happened?
2. Given that Y happened, how likely is it that X happened?

Back to probabilistic modeling



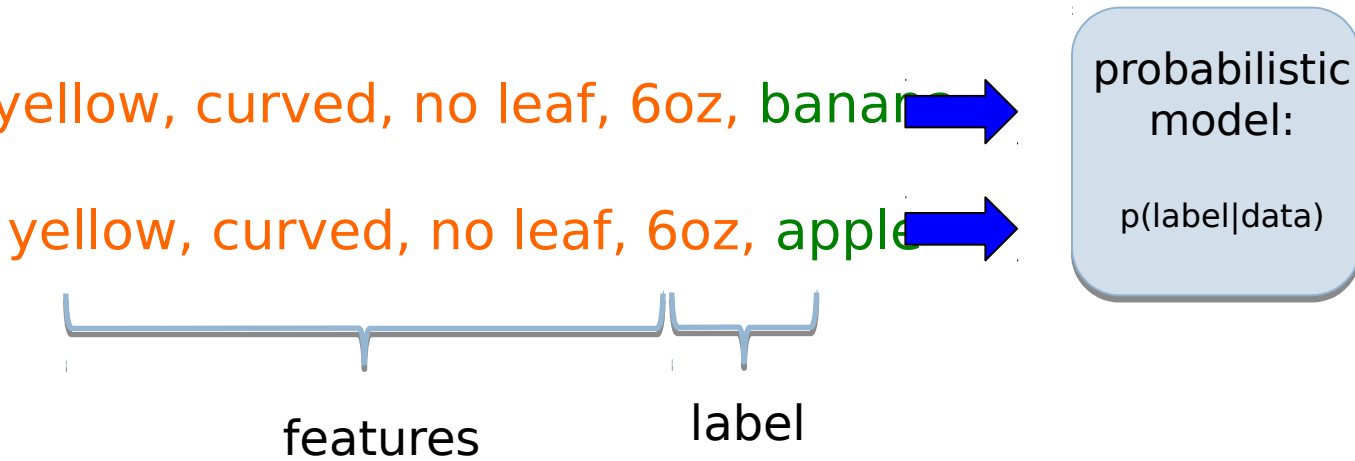
Build a model of the conditional distribution:

$P(\text{label} \mid \text{data})$

How likely is a label given the data

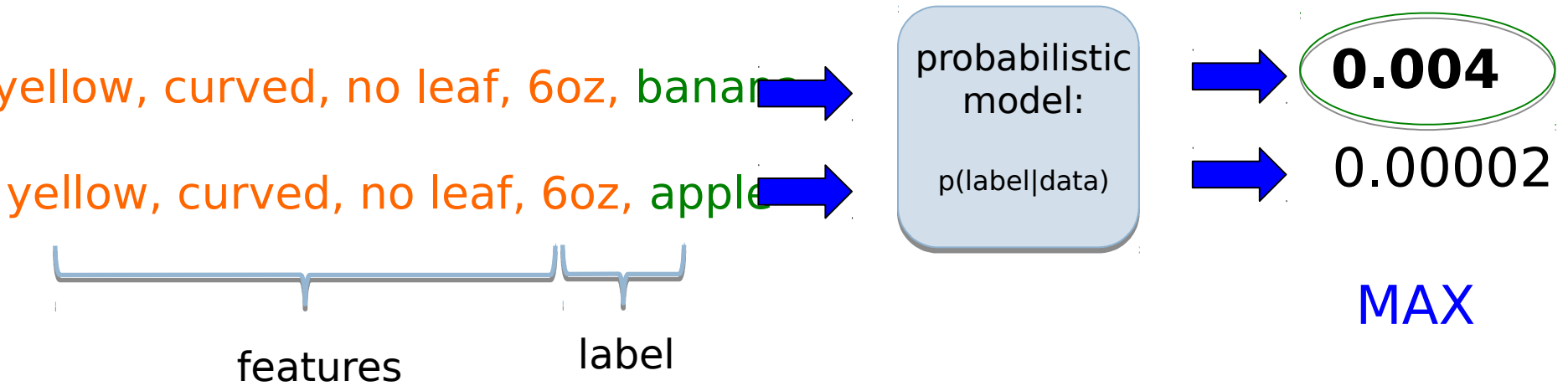
Back to probabilistic models

For each label, calculate the probability of the label given the data



Back to probabilistic models

Pick the label with the highest probability



Naïve Bayes model

Two parallel ways of breaking down the joint distribution

$$P(\text{data}, \text{label}) = P(\text{label}) * P(\text{data}|\text{label})$$

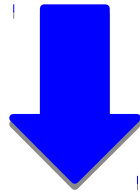
$$P(\text{data}, \text{label}) = P(\text{data}) * P(\text{label}|\text{data})$$

$$P(\text{label}) * P(\text{data}|\text{label}) = P(\text{data}) * P(\text{label}|\text{data})$$

What is $P(\text{label}|\text{data})$?

Naïve Bayes

$$P(\textit{label}) * P(\textit{data}|\textit{label}) = P(\textit{data}) * P(\textit{label}|\textit{data})$$



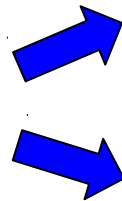
$$P(\textit{label}|\textit{data}) = \frac{P(\textit{label}) * P(\textit{data}|\textit{label})}{P(\textit{data})}$$

(This is called Bayes' rule!)

Naïve Bayes

$$P(\text{label}|\text{data}) = \frac{P(\text{label}) * P(\text{data}|\text{label})}{P(\text{data})}$$

probabilistic
model:
 $p(\text{label}|\text{data})$



$$\frac{P(\text{positive}) * P(\text{data}|\text{positive})}{P(\text{data})}$$

$$\frac{P(\text{negative}) * P(\text{data}|\text{negative})}{P(\text{data})}$$

MAX

One observation

$$\frac{P(\textit{positive}) * P(\textit{data}|\textit{positive})}{P(\textit{data})}$$

MAX

$$\frac{P(\textit{negative}) * P(\textit{data}|\textit{negative})}{P(\textit{data})}$$

For picking the largest $P(\textit{data})$ doesn't matter!

One observation

$$P(\textit{positive}) * P(\textit{data}|\textit{positive})$$

$$P(\textit{negative}) * P(\textit{data}|\textit{negative})$$

MAX

For picking the largest $P(\textit{data})$ doesn't matter!

A simplifying assumption (for this class)

$$\frac{P(\textit{positive}) * P(\textit{data}|\textit{positive})}{P(\textit{negative}) * P(\textit{data}|\textit{negative})} \text{MAX}$$

If we assume $P(\textit{positive}) = P(\textit{negative})$ then:

$$\frac{P(\textit{data}|\textit{positive})}{P(\textit{data}|\textit{negative})} \text{MAX}$$

$P(\text{data}|\text{label})$

$$P(\text{data}|\text{label}) = P(f_1, f_2, \dots, f_n|\text{label})$$

$$\begin{aligned} &\approx P(f_1|\text{label}) * \\ &* P(f_2|\text{label}) * \\ &\quad \dots \\ &P(f_n|\text{label}) \end{aligned}$$

This is generally not true!

However..., it makes our life easier.

This is why the model is called **Naïve** Bayes

Naïve Bayes

$$P(f_1|\textit{positive}) * P(f_2|\textit{positive}) * \dots * P(f_n|\textit{positive})$$

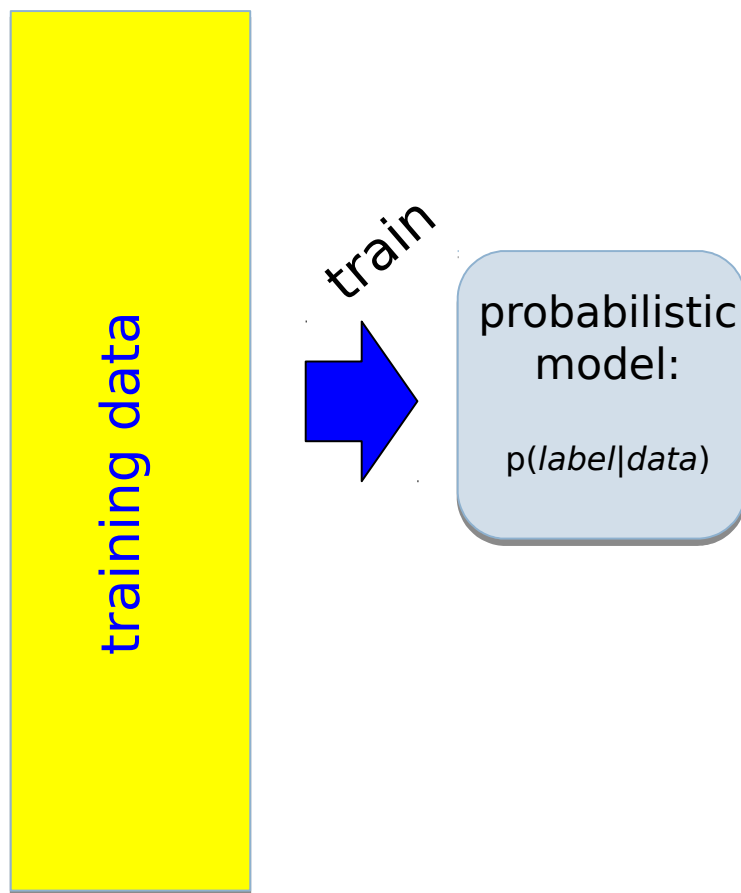
MAX

$$P(f_1|\textit{negative}) * P(f_2|\textit{negative}) * \dots * P(f_n|\textit{negative})$$



Where do these come from?

Training Naïve Bayes



An aside: P(heads)

What is the P(heads) on a fair coin?

0.5

What if you didn't know that, but had a coin to experiment with?

$$P(\text{heads}) = \frac{\text{number of times heads came up}}{\text{total number of coin tosses}}$$

Try it out...



P(feature|label)

$$P(\text{heads}) = \frac{\text{number of times heads came up}}{\text{total number of coin tosses}}$$

Can we do the same thing here? What is the probability of a feature given positive, i.e. the probability of a feature occurring in in the positive label?

$$P(\text{feature}|\text{positive}) = ?$$

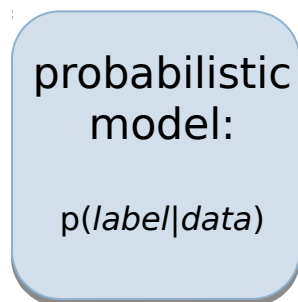
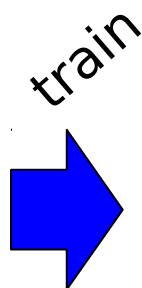
P(feature|label)

$$P(\text{heads}) = \frac{\text{number of times heads came up}}{\text{total number of coin tosses}}$$

Can we do the same thing here? What is the probability of a feature given positive, i.e. the probability of a feature occurring in in the positive label?

$$P(\text{feature}|\text{positive}) = \frac{\text{number of positive examples with that feature}}{\text{total number of positive examples}}$$

Training Naïve Bayes



1. Count how many examples have each label
2. For all examples with a particular label, count how many times each feature occurs
3. Calculate the conditional probabilities of each feature for all labels:

$$P(\text{feature}|\text{label}) = \frac{\text{number of ``label'' examples with that feature}}{\text{total number of examples with that label}}$$

Classifying with Naïve Bayes

For each label, calculate the product of $p(\text{feature}|\text{label})$ for each label

yellow, curved, no leaf, 6oz

$P(\text{yellow}|\text{banana}) * \dots * P(6\text{oz}|\text{banana})$

$P(\text{yellow}|\text{apple}) * \dots * P(6\text{oz}|\text{apple})$

MAX

Naïve Bayes Text Classification

Positive

I loved it
I loved that movie
I hated that I loved it

Negative

I hated it
I hated that movie
I loved that I hated it

Given examples of text in different categories, learn to predict the category of new examples

Sentiment classification: given positive/negative examples of text (sentences), learn to predict whether new text is positive/negative

Text classification training

Positive

I loved it
I loved that movie
I hated that | loved it

Negative

I hated it
I hated that movie
I loved that | hated it

We'll assume words just occur once in any given sentence

Text classification training

Positive

I loved it
I loved that movie
I hated that loved it

Negative

I hated it
I hated that movie
I loved that hated it

We'll assume words just occur once in any given sentence

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

Negative

I hated it
I hated that movie
I loved that hated it

For each word and each label, learn:

$$p(\text{word} \mid \text{label})$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

Negative

I hated it
I hated that movie
I loved that hated it

$P(I \mid \text{positive}) = ?$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$$P(I \mid \text{positive}) = 3/3 = 1.0$$

Negative

I hated it
I hated that movie
I loved that hated it

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

Negative

I hated it
I hated that movie
I loved that hated it

$$P(I \mid \text{positive}) = 1.0$$
$$P(\text{loved} \mid \text{positive}) = ?$$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

Negative

I hated it
I hated that movie
I loved that hated it

$$P(I \mid \text{positive}) = 1.0$$
$$P(\text{loved} \mid \text{positive}) = 3/3$$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

Negative

I hated it
I hated that movie
I loved that hated it

$$\begin{aligned}P(I \mid \text{positive}) &= 1.0 \\P(\text{loved} \mid \text{positive}) &= 3/3 \\P(\text{hated} \mid \text{positive}) &= ?\end{aligned}$$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$$\begin{aligned}P(I \mid \text{positive}) &= 1.0 \\P(\text{loved} \mid \text{positive}) &= 3/3 \\P(\text{hated} \mid \text{positive}) &= 1/3\end{aligned}$$

...

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Negative

I hated it
I hated that movie
I loved that hated it

$$P(I \mid \text{negative}) = ?$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$$\begin{aligned}P(I \mid \text{positive}) &= 1.0 \\P(\text{loved} \mid \text{positive}) &= 3/3 \\P(\text{hated} \mid \text{positive}) &= 1/3\end{aligned}$$

Negative

I hated it
I hated that movie
I loved that hated it

$$P(I \mid \text{negative}) = 1.0$$

...

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$$\begin{aligned}P(I \mid \text{positive}) &= 1.0 \\P(\text{loved} \mid \text{positive}) &= 3/3 \\P(\text{hated} \mid \text{positive}) &= 1/3\end{aligned}$$

Negative

I hated it
I hated that movie
I loved that hated it

$$\begin{aligned}P(I \mid \text{negative}) &= 1.0 \\P(\text{movie} \mid \text{negative}) &= ?\end{aligned}$$

...

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$P(I \mid \text{positive}) = 1.0$
 $P(\text{loved} \mid \text{positive}) = 3/3$
 $P(\text{hated} \mid \text{positive}) = 1/3$

Negative

I hated it
I hated that movie
I loved that hated it

$P(I \mid \text{negative}) = 1.0$
 $P(\text{movie} \mid \text{negative}) = 1/3$
...

...

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Classifying

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 1.0$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

Notice that each of these is its own probability distribution

P(loved| positive)

$$P(\text{loved} \mid \text{positive}) = 2/3$$

$$P(\text{no loved} \mid \text{positive}) = 1/3$$

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

How would we classify: “I hated movie”?

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

$$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) = 1.0 * 1/3 * 1/3 = 1/9$$

$$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) = 1.0 * 1.0 * 1/3 = 1/3$$

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

How would we classify: “I hated the movie”?

Trained model

$$P(I \mid \text{positive}) = 1.0$$

$$P(\text{loved} \mid \text{positive}) = 2/3$$

$$p(\text{it} \mid \text{positive}) = 2/3$$

$$p(\text{that} \mid \text{positive}) = 2/3$$

$$p(\text{movie} \mid \text{positive}) = 1/3$$

$$P(\text{hated} \mid \text{positive}) = 1/3$$

$$P(I \mid \text{negative}) = 1.0$$

$$p(\text{hated} \mid \text{negative}) = 1.0$$

$$p(\text{that} \mid \text{negative}) = 2/3$$

$$P(\text{movie} \mid \text{negative}) = 1/3$$

$$p(\text{it} \mid \text{negative}) = 2/3$$

$$p(\text{loved} \mid \text{negative}) = 1/3$$

$$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{the} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) =$$

$$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{the} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) =$$

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

$$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{the} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) =$$
$$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{the} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) =$$

What are these?

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{the} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) =$

$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{the} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) =$

0! Is this a problem?

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{the} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) =$

$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{the} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) =$

Yes. They make the entire product go to 0!

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{the} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) =$

$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{the} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) =$

Our solution: assume any unseen word has a small, fixed probability, e.g. in this example 1/10

Trained model

$P(I \mid \text{positive}) = 1.0$	$P(I \mid \text{negative}) = 1.0$
$P(\text{loved} \mid \text{positive}) = 2/3$	$p(\text{hated} \mid \text{negative}) = 1.0$
$p(\text{it} \mid \text{positive}) = 2/3$	$p(\text{that} \mid \text{negative}) = 2/3$
$p(\text{that} \mid \text{positive}) = 2/3$	$P(\text{movie} \mid \text{negative}) = 1/3$
$p(\text{movie} \mid \text{positive}) = 1/3$	$p(\text{it} \mid \text{negative}) = 2/3$
$P(\text{hated} \mid \text{positive}) = 1/3$	$p(\text{loved} \mid \text{negative}) = 1/3$

$$P(I \mid \text{positive}) * P(\text{hated} \mid \text{positive}) * P(\text{the} \mid \text{positive}) * P(\text{movie} \mid \text{positive}) = 1/90$$

$$P(I \mid \text{negative}) * P(\text{hated} \mid \text{negative}) * P(\text{the} \mid \text{negative}) * P(\text{movie} \mid \text{negative}) = 1/30$$

Our solution: assume any unseen word has a small, fixed probability, e.g. in this example 1/10

Full disclaimer

I've fudged a few things on the Naïve Bayes model for simplicity

Our approach is very close, but it takes a few liberties that aren't technically correct, but it will work just fine

If you're curious, I'd be happy to talk to you offline