

# CSE 312 SECTION 3

## THE NAIVE BAYES CLASSIFIER

LUXI WANG, PEMI NGUYEN,  
MITCHELL ESTBERG AND SHREYA JAYARAMAN  
ALEX TSUN

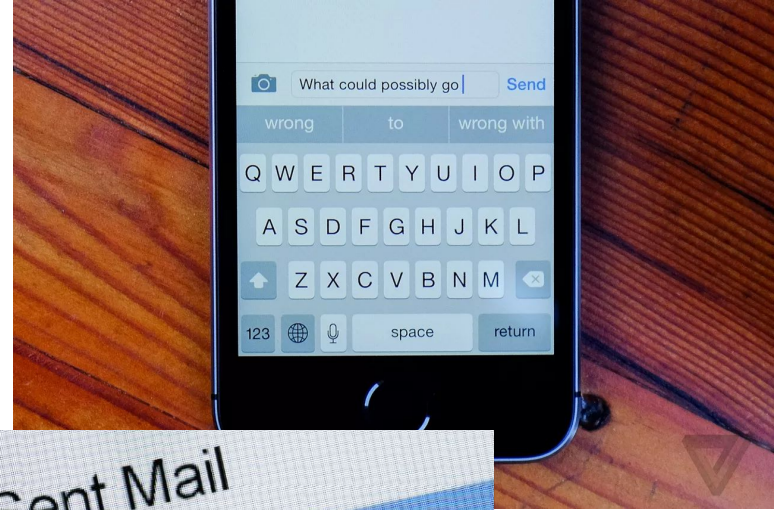
# ANNOUNCEMENTS

- PSET2 DUE THURS, JAN 21, AT 11:59 PM PST
- PSET3 OUT YESTERDAY, JAN 20
  - DUE WED, JAN 27, AT 11:59 PM PDT
- CONCEPT CHECK DUE FRIDAY, JANUARY 22ND AT 9:00 AM PST




# AGENDA

- WHAT IS MACHINE LEARNING?
- FEATURIZING EMAILS
- NAIVE BAYES

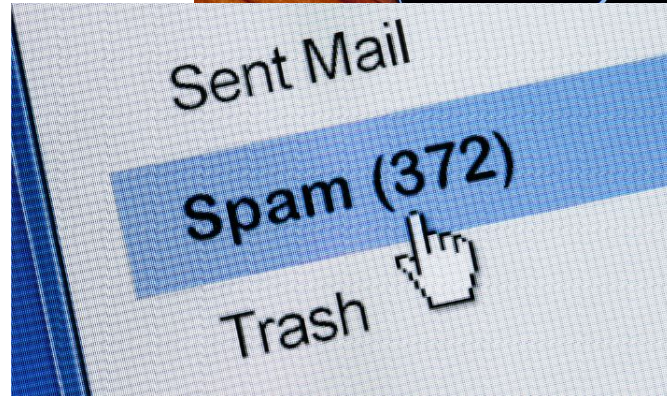
# MACHINE LEARNING IN THE REAL WORLD



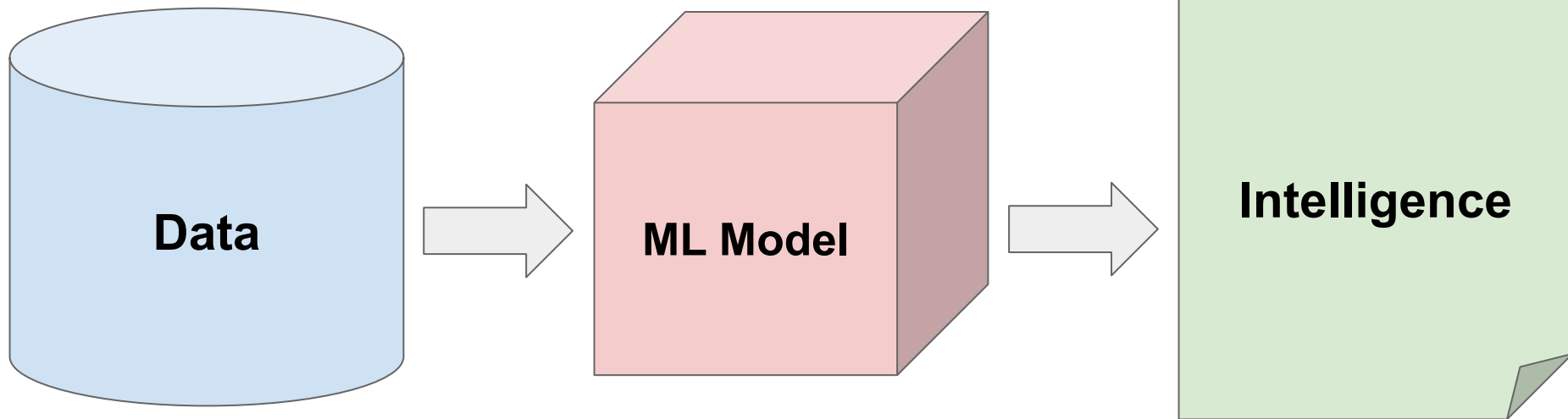
Jobs you may be interested in

 <b>Engineering Manager - Data Infrastructure</b> Twilio Inc. — San Francisco, CA, US <a href="#">View Job</a>	 <b>Chief Architect</b> Appthority — San Francisco Bay Area <a href="#">View Job</a>	 <b>Sr. Engineering Manager</b> Comcast Silicon Valley Innovation Center — San <a href="#">View Job</a>
--	--	---

[See more jobs >](#)

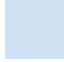






# ML PIPELINE



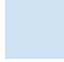




From **Wikipedia**: “Machine learning is the study of computer algorithms that improve automatically through experience.”

# YOU ARE A MACHINE!

Number	Shape	“Label”
3		12
5		15
-2		-8
7		21
-4		???

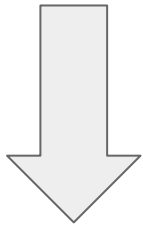
Given examples with correct “labels”, make predictions!

# YOU ARE A MACHINE!

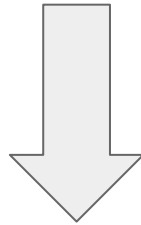
Number	Shape	“Label”
3		12
5		15
-2		-8
7		21
-4		-16

Given examples with correct “labels”, make predictions!

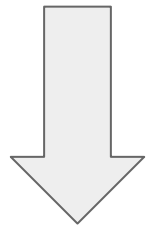
# REGRESSION: IDEA



\$ 340,135



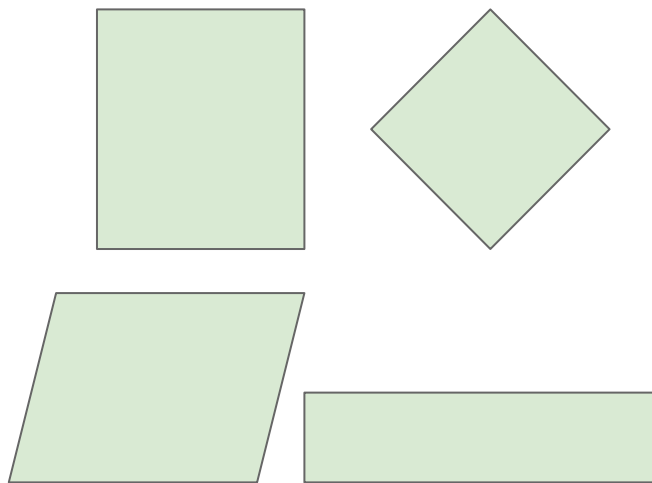
\$801,353



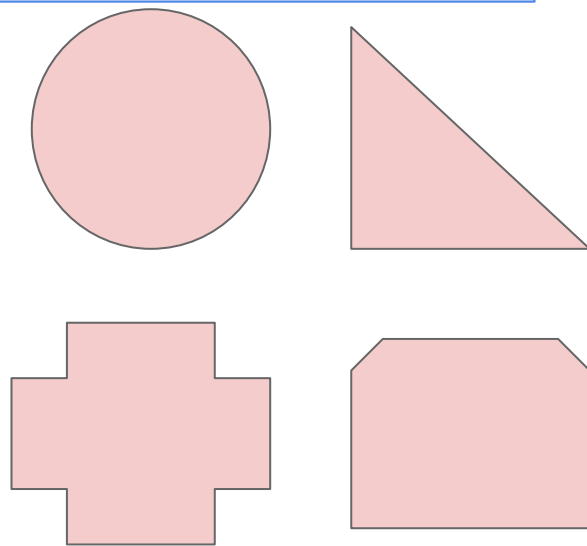
??????



# CLASSIFICATION: IDEA

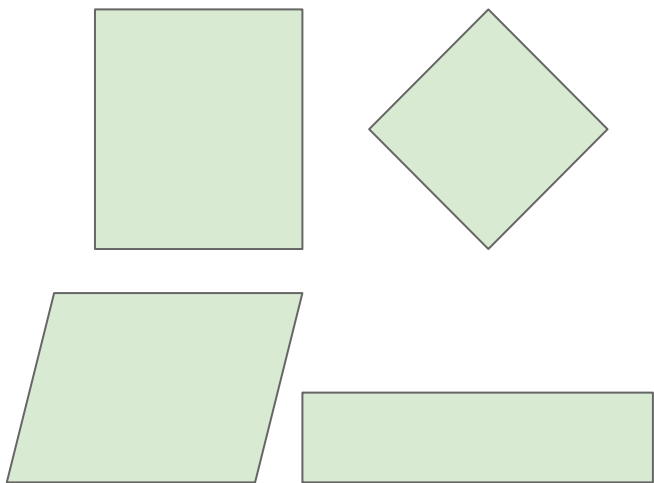


“Green” class

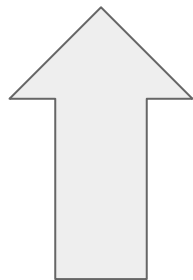


“Red” class

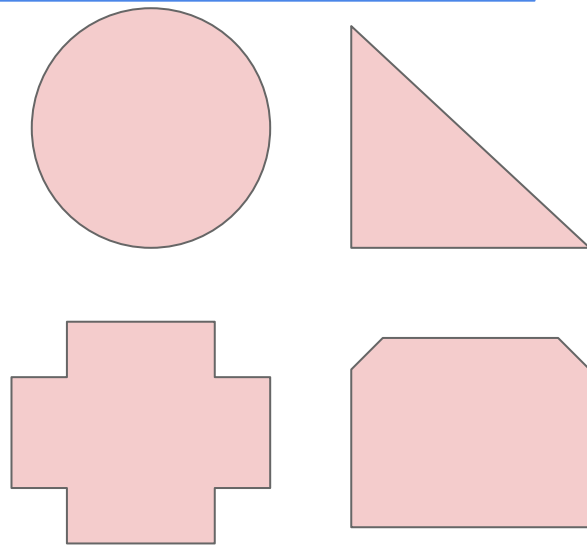
# CLASSIFICATION: IDEA



“Green” class



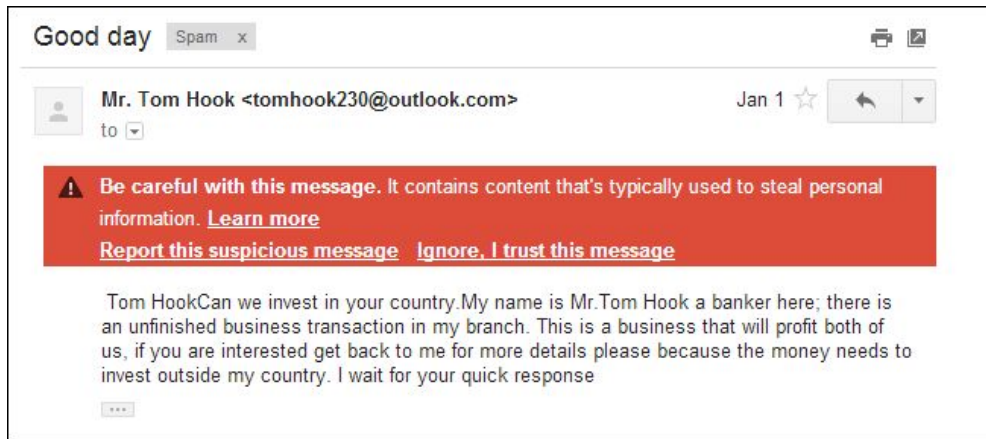
Is this new shape  
supposed to be  
“green” or “red”?



“Red” class

# SPAM FILTER

- In real life, you may have seen a lot of spam emails like this.
- Building a good spam filter helps protect users from potential scams, unnecessary advertising, or malware links.



# EVALUATING PERFORMANCE

## Training Set

## Test Set

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

Email	Label
You buy viagra!	Spam
You need viagra sir.	Spam
I hope you are healthy.	Ham
...	...
...	...

We “**train**” our spam filter on the training set, and **evaluate** performance using a test set (data that is unseen by the spam filter initially). This gives an unbiased estimate of performance.

# SPAM FILTER TASK

## Training Set

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham



**Predict** whether this email is spam or ham:

You buy Viagra!

# EMAILS AS WORD COLLECTIONS

Email	Set of Words in the Email
<p>SUBJECT: Top Secret Business Venture</p> <p>Dear Sir. First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret...</p>	<p>{top, secret, business, venture, dear, sir, first, I, must, solicit, your, confidence, in, this, transaction, is, by, virtue, of, its, nature, as, being, utterly, confidential, and}</p>

For simplicity, we will

- Ignore Duplicate Words
- Ignore Punctuation
- Ignore Casing

# EMAILS AS WORD COLLECTIONS

Email	Set of Words in the Email
<p>SUBJECT: Top Secret Business Venture</p> <p>Dear Sir. First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret...</p>	<p>{top, secret, business, venture, dear, sir, first, I, must, solicit, your, confidence, in, this, transaction, is, by, virtue, of, its, nature, as, being, utterly, confidential, and}</p>
<p>Hello hello hello there.</p>	<p>{hello, there}</p>

For simplicity, we will

- Ignore Duplicate Words
- Ignore Punctuation
- Ignore Casing

# EMAILS AS WORD COLLECTIONS

Email	Set of Words in the Email
<p>SUBJECT: Top Secret Business Venture</p> <p>Dear Sir. First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret...</p>	{top, secret, business, venture, dear, sir, first, I, must, solicit, your, confidence, in, this, transaction, is, by, virtue, of, its, nature, as, being, utterly, confidential, and}
Hello hello hello there.	{hello, there}
You buy Viagra!	{you, buy, viagra}

For simplicity, we will

- Ignore Duplicate Words
- Ignore Punctuation
- Ignore Casing



# OUR APPROACH

Compute and Compare:

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"})$$

$$\mathbb{P}(\text{ham} \mid \text{"You buy Viagra!"})$$

Then predict whichever is larger! Can we get away with just computing one of them?

# OUR APPROACH

Compute and Compare:

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"})$$

$$\mathbb{P}(\text{ham} \mid \text{"You buy Viagra!"})$$

Then predict whichever is larger! Can we get away with just computing one of them?

Equivalently, note that these add to 1, so we can just compute

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"})$$

and if it is greater than 0.5, then we predict **spam**.

Otherwise, we predict **ham**.

Note: We resolve the tie in favor of **ham**.

# NAIVE BAYES CLASSIFIER - THE BAYES PART

Bayes Theorem: 
$$\mathbb{P}(A | B) = \frac{\mathbb{P}(B | A) \mathbb{P}(A)}{\mathbb{P}(B)}$$

Apply it to our example:

$$\mathbb{P}(\text{spam} | \text{"You buy Viagra!"}) = \frac{\mathbb{P}(\text{"You buy Viagra!"} | \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"You buy Viagra!"})}$$



# NAIVE BAYES CLASSIFIER - WHAT WE CALCULATE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"}) = \frac{\mathbb{P}(\text{"You buy Viagra!"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"You buy Viagra!"})}$$

# NAIVE BAYES CLASSIFIER - WHAT WE CALCULATE

$$\begin{aligned} \mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"}) &= \frac{\mathbb{P}(\text{"You buy Viagra!"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"You buy Viagra!"})} \\ &= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})} \quad [\text{LTP}] \end{aligned}$$

# NAIVE BAYES CLASSIFIER - WHAT WE CALCULATE

$$\begin{aligned}\mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"}) &= \frac{\mathbb{P}(\text{"You buy Viagra!"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"You buy Viagra!"})} \\ &= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})} \quad [\text{LTP}]\end{aligned}$$

$$\mathbb{P}(\text{spam}) = \frac{\text{total spam emails (in training set)}}{\text{total emails (in training set)}}$$

$$\mathbb{P}(\text{ham}) = \frac{\text{total ham emails (in training set)}}{\text{total emails (in training set)}}$$

(our approximation for these probabilities,  
based on the training set)

# NAIVE BAYES CLASSIFIER - THE NAIVE PART

It is somewhat unlikely that we have the email "You buy Viagra!" in our training data. (In this case we don't!)

# NAIVE BAYES CLASSIFIER - THE NAIVE PART

It is somewhat unlikely that we have the email "You buy Viagra!" in our training data. (In this case we don't!)

We naively assume that words are conditionally independent from each other, given the label (In reality, they aren't):



# NAIVE BAYES CLASSIFIER - THE NAIVE PART

It is somewhat unlikely that we have the email "You buy Viagra!" in our training data. (In this case we don't!)

We naively assume that words are conditionally independent from each other, given the label (In reality, they aren't):

$$\begin{aligned} & \mathbb{P}(\{ \text{"you"}, \text{"buy"}, \text{"viagra"} \} \mid \text{spam}) \\ & \approx \mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \end{aligned}$$

# NAIVE BAYES CLASSIFIER - THE NAIVE PART

It is somewhat unlikely that we have the email "You buy Viagra!" in our training data. (In this case we don't!)

We **naively assume that words are conditionally independent from each other, given the label (In reality, they aren't):**

$$\begin{aligned} & \mathbb{P}(\{ \text{"you"}, \text{"buy"}, \text{"viagra"} \} \mid \text{spam}) \\ & \approx \mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \end{aligned}$$

Then we estimate for example that

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{\text{number of spam emails containing "you" (in training set)}}{\text{number of spam emails (in training set)}}$$

# NAIVE BAYES CLASSIFIER - THE NAIVE PART

Consider for example the following two emails:

“!!!Lunch free for You!!!!!!” *Spam*

“You free for lunch?” *Ham*

# NAIVE BAYES CLASSIFIER - THE NAIVE PART

Consider for example the following two emails:

“!!!Lunch free for You!!!!!” *Spam*

“You free for lunch?” *Ham*

One shortfalling of our model is that it will make the same prediction for these since they have the same set of words!

# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$
$$= \frac{\mathbb{P}(\text{"you"} \mid \text{spam})\mathbb{P}(\text{"buy"} \mid \text{spam})\mathbb{P}(\text{"viagra"} \mid \text{spam})\mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam})\mathbb{P}(\text{"buy"} \mid \text{spam})\mathbb{P}(\text{"viagra"} \mid \text{spam})\mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham})\mathbb{P}(\text{"buy"} \mid \text{ham})\mathbb{P}(\text{"viagra"} \mid \text{ham})\mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) =$$

$$\mathbb{P}(\text{ham}) =$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) =$$

$$\mathbb{P}(\text{"you"} \mid \text{ham}) =$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) =$$

$$\mathbb{P}(\text{"buy"} \mid \text{ham}) =$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) =$$

$$\mathbb{P}(\text{"viagra"} \mid \text{ham}) =$$



# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham}) \mathbb{P}(\text{"buy"} \mid \text{ham}) \mathbb{P}(\text{"viagra"} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"you"} \mid \text{ham}) = \frac{1}{2}$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) = \text{ed} \quad \mathbb{P}(\text{"buy"} \mid \text{ham}) =$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) = \quad \mathbb{P}(\text{"viagra"} \mid \text{ham}) = \text{ed}$$



# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$\begin{aligned} &= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})} \\ &= \frac{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham}) \mathbb{P}(\text{"buy"} \mid \text{ham}) \mathbb{P}(\text{"viagra"} \mid \text{ham}) \mathbb{P}(\text{ham})} \end{aligned}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"you"} \mid \text{ham}) = \frac{1}{2}$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"buy"} \mid \text{ham}) = 0$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) = 1 \quad \mathbb{P}(\text{"viagra"} \mid \text{ham}) = \frac{1}{2}$$



# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$
$$= \frac{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham}) \mathbb{P}(\text{"buy"} \mid \text{ham}) \mathbb{P}(\text{"viagra"} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"you"} \mid \text{ham}) = \frac{1}{2}$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"buy"} \mid \text{ham}) = 0$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) = 1 \quad \mathbb{P}(\text{"viagra"} \mid \text{ham}) = \frac{1}{2}$$





# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham}) \mathbb{P}(\text{"buy"} \mid \text{ham}) \mathbb{P}(\text{"viagra"} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"you"} \mid \text{ham}) = \frac{1}{2}$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"buy"} \mid \text{ham}) = 0$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) = 1 \quad \mathbb{P}(\text{"viagra"} \mid \text{ham}) = \frac{1}{2}$$



# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham}) \mathbb{P}(\text{"buy"} \mid \text{ham}) \mathbb{P}(\text{"viagra"} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \mathbf{1}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"you"} \mid \text{ham}) = \frac{1}{2}$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"buy"} \mid \text{ham}) = 0$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) = 1 \quad \mathbb{P}(\text{"viagra"} \mid \text{ham}) = \frac{1}{2}$$



# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra"})$$

$$= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"you"} \mid \text{spam}) \mathbb{P}(\text{"buy"} \mid \text{spam}) \mathbb{P}(\text{"viagra"} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{"you"} \mid \text{ham}) \mathbb{P}(\text{"buy"} \mid \text{ham}) \mathbb{P}(\text{"viagra"} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

**= 1 (Marked as spam since no ham email contained "buy")**

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{"you"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"you"} \mid \text{ham}) = \frac{1}{2}$$

$$\mathbb{P}(\text{"buy"} \mid \text{spam}) = \frac{1}{3} \quad \mathbb{P}(\text{"buy"} \mid \text{ham}) = 0$$

$$\mathbb{P}(\text{"viagra"} \mid \text{spam}) = 1 \quad \mathbb{P}(\text{"viagra"} \mid \text{ham}) = \frac{1}{2}$$



# WHAT HAPPENS IF WE GOT A 0?

$P(\text{ham} \mid \text{"You buy Viagra!"}) = 0$  since  $P(\text{"buy"} \mid \text{ham}) = 0$ , since no ham email in our training data contained the word **'buy'**.

But does that mean we will never encounter a ham email with word **'buy'**?



What about the ham:  
"I'll buy sunflowers"

# LAPLACE SMOOTHING

Pretend in spam emails (training set):

- We saw one extra spam email **with** word  $w_i$
- We saw one extra spam email **without** word  $w_i$



# LAPLACE SMOOTHING

Pretend in spam emails (training set):

- We saw one extra spam email **with** word  $w_i$
- We saw one extra spam email **without** word  $w_i$

$$\mathbb{P}(w_i \mid \text{spam}) = \frac{|\text{total spam emails (training set) containing } w_i| + 1}{|\text{total spam emails (training set)}| + 2}$$



# LAPLACE SMOOTHING



Pretend in spam emails (training set):

- We saw one extra spam email **with** word  $w_i$
- We saw one extra spam email **without** word  $w_i$

Same for ham emails.

$$\mathbb{P}(w_i \mid \text{spam}) = \frac{|\text{total spam emails (training set) containing } w_i| + 1}{|\text{total spam emails (training set)}| + 2}$$

$$\mathbb{P}(w_i \mid \text{ham}) = \frac{|\text{total ham emails (training set) containing } w_i| + 1}{|\text{total ham emails (training set)}| + 2}$$



# LAPLACE SMOOTHING



Pretend in spam emails (training set):

- We saw one extra spam email **with** word  $w_i$
- We saw one extra spam email **without** word  $w_i$

Same for ham emails.

$$\mathbb{P}(w_i \mid \text{spam}) = \frac{|\text{total spam emails (training set) containing } w_i| + 1}{|\text{total spam emails (training set)}| + 2}$$

$$\mathbb{P}(w_i \mid \text{ham}) = \frac{|\text{total ham emails (training set) containing } w_i| + 1}{|\text{total ham emails (training set)}| + 2}$$

$$\mathbb{P}(\text{“buy”} \mid \text{ham}) = \frac{0 + 1}{2 + 2} = \frac{1}{4}$$



# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{“You buy Viagra”})$$

$$= \frac{\mathbb{P}(\{\text{“you”}, \text{“buy”}, \text{“viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{“you”}, \text{“buy”}, \text{“viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{“you”}, \text{“buy”}, \text{“viagra”}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$
$$= \frac{\mathbb{P}(\text{“you”} \mid \text{spam})\mathbb{P}(\text{“buy”} \mid \text{spam})\mathbb{P}(\text{“viagra”} \mid \text{spam})\mathbb{P}(\text{spam})}{\mathbb{P}(\text{“you”} \mid \text{spam})\mathbb{P}(\text{“buy”} \mid \text{spam})\mathbb{P}(\text{“viagra”} \mid \text{spam})\mathbb{P}(\text{spam}) + \mathbb{P}(\text{“you”} \mid \text{ham})\mathbb{P}(\text{“buy”} \mid \text{ham})\mathbb{P}(\text{“viagra”} \mid \text{ham})\mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{spam}) =$$

$$\mathbb{P}(\text{“you”} \mid \text{ham}) =$$

$$\mathbb{P}(\text{“buy”} \mid \text{spam}) =$$

$$\mathbb{P}(\text{“buy”} \mid \text{ham}) = \frac{0 + 1}{2 + 2} = \frac{1}{4}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{spam}) =$$

$$\mathbb{P}(\text{“viagra”} \mid \text{ham}) =$$

# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{“You buy Viagra”})$$

$$= \frac{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$
$$= \frac{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{“you”} \mid \text{ham}) \mathbb{P}(\text{“buy”} \mid \text{ham}) \mathbb{P}(\text{“viagra”} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

$$\mathbb{P}(\text{“buy”} \mid \text{spam}) = \text{ed}$$

$$\mathbb{P}(\text{“buy”} \mid \text{ham}) = \frac{0 + 1}{2 + 2} = \frac{1}{4}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{spam}) =$$

$$\mathbb{P}(\text{“viagra”} \mid \text{ham}) = \text{ed}$$

# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{“You buy Viagra”})$$

$$= \frac{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$
$$= \frac{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{“you”} \mid \text{ham}) \mathbb{P}(\text{“buy”} \mid \text{ham}) \mathbb{P}(\text{“viagra”} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5} \quad \mathbb{P}(\text{“you”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

$$\mathbb{P}(\text{“buy”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5} \quad \mathbb{P}(\text{“buy”} \mid \text{ham}) = \frac{0 + 1}{2 + 2} = \frac{1}{4}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{spam}) = \frac{3 + 1}{3 + 2} = \frac{4}{5} \quad \mathbb{P}(\text{“viagra”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{“You buy Viagra”})$$

$$= \frac{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{“you”} \mid \text{ham}) \mathbb{P}(\text{“buy”} \mid \text{ham}) \mathbb{P}(\text{“viagra”} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\frac{2}{5} \cdot \frac{2}{5} \cdot \frac{4}{5} \cdot \frac{3}{5}}{\frac{2}{5} \cdot \frac{2}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} + \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{2} \cdot \frac{2}{5}} \approx 0.7544$$

Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

$$\mathbb{P}(\text{“buy”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5}$$

$$\mathbb{P}(\text{“buy”} \mid \text{ham}) = \frac{0 + 1}{2 + 2} = \frac{1}{4}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{spam}) = \frac{3 + 1}{3 + 2} = \frac{4}{5}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

# EXAMPLE

$$\mathbb{P}(\text{spam} \mid \text{“You buy Viagra”})$$

$$= \frac{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{“you”, “buy”, “viagra”}\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{“you”} \mid \text{spam}) \mathbb{P}(\text{“buy”} \mid \text{spam}) \mathbb{P}(\text{“viagra”} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\text{“you”} \mid \text{ham}) \mathbb{P}(\text{“buy”} \mid \text{ham}) \mathbb{P}(\text{“viagra”} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$= \frac{\frac{2}{5} \cdot \frac{2}{5} \cdot \frac{4}{5} \cdot \frac{3}{5}}{\frac{2}{5} \cdot \frac{2}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} + \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{2} \cdot \frac{2}{5}} \approx 0.7544$$



Email	Label
Buy Viagra!	Spam
You good?	Ham
Viagra help you.	Spam
Good Viagra help.	Spam
I need Viagra for my health condition.	Ham

$$\mathbb{P}(\text{spam}) = \frac{3}{5}$$

$$\mathbb{P}(\text{ham}) = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5}$$

$$\mathbb{P}(\text{“you”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

$$\mathbb{P}(\text{“buy”} \mid \text{spam}) = \frac{1 + 1}{3 + 2} = \frac{2}{5}$$

$$\mathbb{P}(\text{“buy”} \mid \text{ham}) = \frac{0 + 1}{2 + 2} = \frac{1}{4}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{spam}) = \frac{3 + 1}{3 + 2} = \frac{4}{5}$$

$$\mathbb{P}(\text{“viagra”} \mid \text{ham}) = \frac{1 + 1}{2 + 2} = \frac{1}{2}$$

# UNDERFLOW PREVENTION

- Multiplication of many probabilities, each of which will be between 0 and 1, can result in floating-point underflow. The product will be too small and will result in arithmetic underflow.

# UNDERFLOW PREVENTION

- Multiplication of many probabilities, each of which will be between 0 and 1, can result in floating-point underflow. The product will be too small and will result in arithmetic underflow.
- Reminder: Log property:

$$\log(xy) = \log(x) + \log(y)$$

# UNDERFLOW PREVENTION

- Multiplication of many probabilities, each of which will be between 0 and 1, can result in floating-point underflow. The product will be too small and will result in arithmetic underflow.
- Reminder: Log property:

$$\log(xy) = \log(x) + \log(y)$$

- Summing logs of probabilities is better than multiplying probabilities

$$\begin{aligned}\log\left(\prod_{i=1}^n p_i\right) &= \log(p_1 p_2 \dots p_n) = \log(p_1) + \log(p_2) + \dots + \log(p_n) \\ &= \sum_{i=1}^n \log(p_i)\end{aligned}$$



# APPLYING UNDERFLOW PREVENTION

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$\mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

We will output **spam** iff:

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) > \mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\})$$

# APPLYING UNDERFLOW PREVENTION

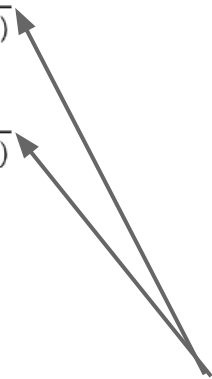
$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$\mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

We will output **spam** iff:

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) > \mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\})$$

$$\iff \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) > \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})$$



Denominators are equal and cancel when comparing

# APPLYING UNDERFLOW PREVENTION

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$\mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

We will output **spam** iff:

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) > \mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\})$$

$$\iff \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) > \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})$$

$$\iff \mathbb{P}(w_1 \mid \text{spam}) \mathbb{P}(w_2 \mid \text{spam}) \cdots \mathbb{P}(w_n \mid \text{spam}) \mathbb{P}(\text{spam}) > \mathbb{P}(w_1 \mid \text{ham}) \mathbb{P}(w_2 \mid \text{ham}) \cdots \mathbb{P}(w_n \mid \text{ham}) \mathbb{P}(\text{ham})$$

# APPLYING UNDERFLOW PREVENTION

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

$$\mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\}) \approx \frac{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}{\mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})}$$

We will output **spam** iff:

$$\mathbb{P}(\text{spam} \mid \{w_1, w_2, \dots, w_n\}) > \mathbb{P}(\text{ham} \mid \{w_1, w_2, \dots, w_n\})$$

$$\iff \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{spam}) \mathbb{P}(\text{spam}) > \mathbb{P}(\{w_1, w_2, \dots, w_n\} \mid \text{ham}) \mathbb{P}(\text{ham})$$

$$\iff \mathbb{P}(w_1 \mid \text{spam}) \mathbb{P}(w_2 \mid \text{spam}) \cdots \mathbb{P}(w_n \mid \text{spam}) \mathbb{P}(\text{spam}) > \mathbb{P}(w_1 \mid \text{ham}) \mathbb{P}(w_2 \mid \text{ham}) \cdots \mathbb{P}(w_n \mid \text{ham}) \mathbb{P}(\text{ham})$$

Taking the log of two sides:

$$\iff \log(\mathbb{P}(\text{spam})) + \sum_{i=1}^n \log(\mathbb{P}(w_i \mid \text{spam})) > \log(\mathbb{P}(\text{ham})) + \sum_{i=1}^n \log(\mathbb{P}(w_i \mid \text{ham}))$$

# CRITICAL THOUGHT

Before deploying any ML system, we should think critically if this is a system we should deploy or any potential downsides. What are the possible questions we might ask about the system?

Here are just SOME examples:

- What biases might it encode?
- Does it have the risk of disproportionately impacting a particular group?
- How do we know if this model is a good model?
- What is the cost of a mistake? Does the email get deleted forever?
- Any many more!

# FURTHER READING: CRITICAL PERSPECTIVE OF ML

What does it mean for a model to be “fair”? There are many (sometimes conflicting) definitions!

- Learn more [here](#). Great application of probability!

Learning from human language can be quite difficult when many common datasets contain bias.

- Recent paper from members of the Allen School: [here](#)

More example of unintended consequences of using models:

- ‘Weapons of Math Destruction’ by Cathy O’Neil

# SUMMARY: NAIVE BAYES ALGORITHM STEPS

## 1. TRAINING

1.1. Compute the proportion of emails in the **training set** that is spam or ham:

$$\mathbb{P}(\text{spam}) = \frac{\text{total spam emails (in training set)}}{\text{total emails (in training set)}}$$

$$\mathbb{P}(\text{ham}) = \frac{\text{total ham emails (in training set)}}{\text{total emails (in training set)}}$$

1.2. Iterate over the **training set**, for each unique word **x**, count:

- How many **spam emails** in the training set contain **x**
- How many **ham emails** in the training set contain **x**

# SUMMARY: NAIVE BAYES ALGORITHM STEPS

## 1. TRAINING

1.1. Compute the proportion of emails in the **training set** that is spam or ham:

$$\mathbb{P}(\text{spam}) = \frac{\text{total spam emails (in training set)}}{\text{total emails (in training set)}}$$

$$\mathbb{P}(\text{ham}) = \frac{\text{total ham emails (in training set)}}{\text{total emails (in training set)}}$$

1.2. Iterate over the **training set**, for each unique word **x**, count:

- How many **spam emails** in the training set contain **x**
- How many **ham emails** in the training set contain **x**

## 2. TESTING

Iterate over the **test set**, for each unlabelled email **D**:

- Create a set **S** of **n** unique words appearing in **D**:  $\{w_1, w_2, \dots, w_n\}$
- For each word  $w_i$  in set **S**, calculate:

$$\mathbb{P}(w_i | \text{spam}) = \frac{|\text{total spam emails (training set) containing } w_i| + 1}{|\text{total spam emails (training set)}| + 2}$$

$$\mathbb{P}(w_i | \text{ham}) = \frac{|\text{total ham emails (training set) containing } w_i| + 1}{|\text{total ham emails (training set)}| + 2}$$

- Note: If word  $w_i$  doesn't appear in the training set, we still calculate the above probabilities, with:

$$|\text{total spam emails (training set) containing } w_i| = |\text{total ham emails (training set) containing } w_i| = 0$$

- If  $\log(\mathbb{P}(\text{spam})) + \sum_{i=1}^n \log(\mathbb{P}(w_i | \text{spam})) > \log(\mathbb{P}(\text{ham})) + \sum_{i=1}^n \log(\mathbb{P}(w_i | \text{ham}))$

Predict email **D** as **spam**

Otherwise, predict email **D** as **ham**





QUESTIONS?  
COMMENTS?  
CONCERNS?

LUXI WANG, PEMI NGUYEN, AND SHREYA JAYARAMAN  
ALEX TSUN