

# DSC 40A

## Theoretical Foundations of Data Science I

- \* fill out SET
- \* HW due Wednesday
- \* final logistics will be sent Monday

# Question

Answer at [q.dsc40a.com](http://q.dsc40a.com)

Remember, you can always ask questions at  
[q.dsc40a.com](http://q.dsc40a.com)!

If the direct link doesn't work, click the "Lecture Questions" link in the top right corner of [dsc40a.com](http://dsc40a.com).

# Last Time

- Classification is the problem of predicting a categorical variable.
- The Naïve Bayes algorithm gives a strategy for classifying data according to its features.
- It is naïve because it relies on an assumption of conditional independence of the features, for a given class.

# Naive Bayes Algorithm

- Bayes' Theorem shows how to calculate  $P(\text{class} \mid \text{features})$ .

$$P(\text{class} \mid \text{features}) = \frac{P(\text{features} \mid \text{class}) * P(\text{class})}{P(\text{features})}$$

- Rewrite the numerator, using the naïve assumption of conditional independence of features given the class.
- Estimate each term in the numerator based on the training data.
- Select class based on whichever has the larger numerator.

# Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

$P(\text{ripe} | \text{green-black, firm, Zutano})$   
 conditional independence  
 $P(\text{firm}|\text{ripe}) \cdot P(\text{gb}|\text{ripe}) \cdot P(\text{zutano}|\text{ripe})$   
 $= \frac{1}{7} \cdot \frac{3}{7} \cdot \frac{2}{7}$   
 $P(\text{ripe}|\text{features}) \propto \frac{1 \cdot 3 \cdot 2}{7 \cdot 7 \cdot 7} \cdot \frac{7}{11} = \frac{6}{539}$

$P(\text{unripe} | \text{green-black, firm, Zutano})$   
 $P(\text{unripe}) = \frac{7}{11}$  (prior)

# Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

*common mistake*

$$P(A|B) = 1 - P(\bar{A}|B) \neq 1 - P(A|\bar{B})$$

Assuming conditional independence of features given the class, calculate  $P(\text{firm, green-black, Zutano} | \text{unripe})$ .

- A. 0
- B. 1/4
- C. 3/16
- ☒ D.  $1 - (1/7 * 3/7 * 2/7) \neq 1 - P(\text{firm, gb, zutano} | \text{ripe})$

# Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

$$\rightarrow P(\text{firm}|\text{unripe}) \cdot P(\text{gb}|\text{unripe}) \cdot P(\text{Zutano}|\text{unripe})$$

$$= \frac{3}{4}$$

$$\cdot \frac{2}{4}$$

$$\cdot \frac{2}{4}$$

$$P(\text{unripe}) = \frac{4}{11} \Rightarrow \text{numerator: } \frac{3}{16} \cdot \frac{4}{11} = \frac{3}{44} > \frac{6}{53}$$

# Naïve Bayes Algorithm

- Bayes' Theorem shows how to calculate  $P(\text{class} | \text{features})$ .

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

- Rewrite the numerator, using the naïve assumption of conditional independence of features given the class.
- Estimate each term in the numerator based on the training data.
- Select class based on whichever has the larger numerator.



# Agenda

- How can we use naïve Bayes to classify text?

# Naïve Bayes Classifier for Text Classification

# Bayes' Theorem for Text Classification

Text classification problems include:

- sentiment analysis
  - positive and negative customer reviews
- determining genre
  - news articles, blog posts, etc.
- email foldering
  - promotions tab in Gmail
- **spam filtering**
  - **separating spam from ham (good, non-spam email)**



# Features

Represent an email as a vector or array of features

$$(x_1, x_2, x_3, \dots, x_n)$$

where  $i$  is an index into a dictionary of  $n$  possible words, and

$x_i = 1$  if word  $i$  is present in the email  
 $x_i = 0$  otherwise

} indicator variable

# Features

Called the “bag of words” model:

Ignores location of words within the email, and the frequency of words

look - sook

**Example:**

1 can

2 do

3 each

4 something

5 understand

6 when

“I do understand if something comes  
up and you cannot do this today”

$(\frac{0}{1} \frac{1}{2} \frac{0}{3} \frac{1}{4} \frac{1}{5} \frac{0}{6})$



usually  $n = 10,000$  to  
50,000 words in practice

# Naïve Bayes Spam Classifier

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

To classify an email, we will use Bayes' Theorem to calculate the probability of it belonging to each class:

$$P(\text{spam} | \text{features}) \text{ and } P(\text{ham} | \text{features})$$

Then choose the class according to the **larger** of these two probabilities.

# Naïve Bayes Spam Classifier

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

**Observe:** the formulas for  $P(\text{spam} | \text{features})$  and  $P(\text{ham} | \text{features})$  have the same denominator,  $P(\text{features})$ .

We can find the larger of the two probabilities by just comparing numerators.

$P(\text{features} | \text{spam}) * P(\text{spam})$  vs.  $P(\text{features} | \text{ham}) * P(\text{ham})$

# Naive Bayes Spam Classifier

$$P(\overset{B}{\text{class}}|\overset{A}{\text{features}}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

$$P(A|B) + P(\bar{A}|B) = \frac{P(A \cap B)}{P(B)} + \frac{P(\bar{A} \cap B)}{P(B)} = \frac{P(S \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1$$

To use Bayes' Theorem, need to determine four quantities:

$$P(A|B) + P(\bar{A}|B) \neq 1 \text{ (not necessary)}$$

i.  $P(\text{features} | \text{spam})$

ii.  $P(\text{features} | \text{ham})$

iii.  $P(\text{spam})$

iv.  $P(\text{ham})$

$$P(\bar{B}) + P(B) = 1$$

Which of these probabilities should add to 1?

~~A.~~ i, ii

☒ B. iii, iv

~~C.~~ both A and B

~~D.~~ neither A nor B



# Estimating Parameters with Training Data

parameter:

$P(\text{spam})$

estimate:

$\frac{\# \text{ spam emails in training}}{\text{size of training set}}$

$\Rightarrow$  counting

parameter:

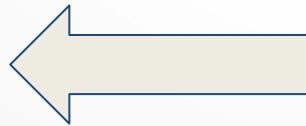
$P(\text{ham})$

estimate:

$\frac{\# \text{ ham emails in training set}}{\text{size of training set}}$

$P(\text{features} \mid \text{spam})$

$P(\text{features} \mid \text{ham})$



harder to estimate

# Assumption of Conditional Independence

To estimate  $P(\text{features} \mid \text{spam})$  and  $P(\text{features} \mid \text{ham})$ , we assume that the probability of a word appearing in an email of a given class is not affected by other words in the email.

$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots \mid \text{spam}) = \leftarrow \text{assumed equal}$

$$P(x_1=0 \mid \text{spam}) * P(x_2=1 \mid \text{spam}) * P(x_3=1 \mid \text{spam}) * \dots$$

*Is this a reasonable assumption? No, this is naïve. Even given ham/spam we have word co-occurrences with different probabilities.*

# Estimating Parameters with Training Data

$$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots | \text{spam}) = \text{← assumed equal}$$

$$P(x_1=0 | \text{spam}) * P(x_2=1 | \text{spam}) * P(x_3=1 | \text{spam}) * \dots$$

parameter:

$$P(x_1=0 | \text{spam})$$

counting

word 1 doesn't appear in spam

estimate:

# spam emails in training set not containing the first word in the dictionary

# spam emails in training set

# Estimating Parameters with Training Data

$$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots | \text{spam}) = \leftarrow \text{assumed equal}$$

$$P(x_1=0 | \text{spam}) * P(x_2=1 | \text{spam}) * P(x_3=1 | \text{spam}) * \dots$$

parameter:

$$P(x_2=1 | \text{spam})$$

prob. of word 2 appearing in  
spam email

estimate:


$$\frac{\text{\# spam emails in training set not containing the first word in the dictionary}}{\text{\# spam emails in training set}}$$

# Estimating Parameters with Training Data

$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots | \text{spam}) = \leftarrow \text{assumed equal}$

$$P(x_1=0 | \text{spam}) * P(x_2=1 | \text{spam}) * P(x_3=1 | \text{spam}) * \dots$$

parameter:

$P(x_3=1 | \text{spam})$   word 3 appeared in spam

estimate:

$$\frac{\text{\# spam emails in training set not containing the first word in the dictionary}}{\text{\# spam emails in training set}}$$

Can term-by-term estimate  $P(\text{features} | \text{class})$

# Naïve Bayes Spam Classifier: Recap

Bayes' Theorem shows how to calculate  $P(\text{spam} \mid \text{features})$  and  $P(\text{ham} \mid \text{features})$ .

$$P(\text{class} \mid \text{features}) = \frac{\overset{\text{conditional indep.}}{P(\text{features} \mid \text{class})} * \overset{\text{prior}}{P(\text{class})}}{\underbrace{P(\text{features})}}$$

Rewrite the numerator, using the naive assumption of conditional independence of words given the class.

Estimate each term in the numerator based on the training data. *by counting*

Select class based on whichever has the larger numerator.

# Modifications and Extensions

- features are pairs (or longer sequences) of words rather than individual words
  - better captures dependencies between words
  - less naïve
  - much bigger feature space
    - $n$  words  $\rightarrow n^2$  pairs of words
- features are the number of occurrences of each word
  - captures low-frequency vs. high-frequency words
- smoothing
  - better handling of previously unseen words

# Smoothing

$$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots | \text{spam}) = \frac{0}{\# \text{spam}} = 0$$
$$P(x_1=0 | \text{spam}) * P(x_2=1 | \text{spam}) * P(x_3=1 | \text{spam}) * \dots$$

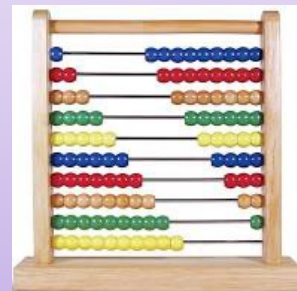
## Dictionary

1. a
2. aardvark
3. abacus
4. abandon
5. abate

...

Suppose you are classifying an email containing the word “abacus,” which does not appear in any emails in your training data. For this new email’s features, what is  **$P(\text{features} | \text{spam})$**  according to a naive Bayes classifier?

- A.  $P(\text{features} | \text{spam}) = \text{undefined}$
- ☒ B.  $P(\text{features} | \text{spam}) = 0$
- C.  $P(\text{features} | \text{spam}) = 1/n$
- D.  $P(\text{features} | \text{spam}) = 1$





# Smoothing

$$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots | \text{ham}) = \underbrace{\frac{0}{\# \text{ham}}}_{\text{red}} P(x_1=0 | \text{ham}) * P(x_2=1 | \text{ham}) * P(x_3=1 | \text{ham}) * \dots$$

## Dictionary

1. a
2. aardvark
3. abacus
4. abandon
5. abate

...

Suppose you are classifying an email containing the word “abacus,” which does not appear in any emails in your training data. For this new email’s features, what is  **$P(\text{features} | \text{ham})$**  according to a naïve Bayes classifier?

- A.  $P(\text{features} | \text{ham}) = \text{undefined}$
- ☒ B.  $P(\text{features} | \text{ham}) = 0$
- C.  $P(\text{features} | \text{ham}) = 1/n$
- D.  $P(\text{features} | \text{ham}) = 1$

# Smoothing

If any  $P(\text{feat}_i | \text{class}) = 0 \Rightarrow P(\text{features} | \text{class}) = 0$

$P(\text{features} | \text{spam}) = 0$  and  $P(\text{features} | \text{ham}) = 0$

**Tiebreaker:** randomly select one of the classes?

# Smoothing

$P(\text{features} \mid \text{spam}) = 0$  and  $P(\text{features} \mid \text{ham}) = 0$

**Tiebreaker:** randomly select one of the classes?

**Better solution:** make sure probabilities can't be zero

**Key idea:** just because you've never seen something happen doesn't mean it's impossible  $\Rightarrow$  make it a small prob event

# Estimating Parameters with Training Data

#spam  $\rightarrow$  #spam + 1

#ham  $\rightarrow$  #ham + 1

smoothing  
prob. of  
class

Without  
Smoothing

With  
Smoothing

parameter:

P(spam)

estimate:

$$\frac{\text{\#spam}}{\text{\#spam} + \text{\#ham}}$$

$$\frac{\text{\#spam} + 1}{\text{\#spam} + 1 + \text{\#ham} + 1}$$

+

+

+

parameter:

P(ham)

estimate:

$$\frac{\text{\#ham}}{\text{\#spam} + \text{\#ham}}$$

$$\frac{\text{\#ham} + 1}{\text{\#spam} + 1 + \text{\#ham} + 1}$$

//

//

//

1

1

1

$$\frac{\text{\#spam}}{\text{\#spam} + \text{\#ham}} + \frac{\text{\#ham}}{\text{\#spam} + \text{\#ham}} = 1$$

$$\frac{\text{\#ham} + 1}{\text{\#spam} + 1 + \text{\#ham} + 1} + \frac{\text{\#spam} + 1}{\text{\#spam} + 1 + \text{\#ham} + 1} = 1$$

# Estimating Parameters with Training Data

parameter:

*smoothing conditional prob*

$$P(x_i=1 \mid \text{spam})$$

estimate:

$$\frac{\text{\#spam containing word } i}{\text{\#spam containing word } i + \text{\#spam not containing word } i}$$

Without Smoothing

*if word never appeared*

$$\frac{(\text{\#spam containing word } i) + 1}{(\text{\#spam containing word } i) + 1 + (\text{\#spam not containing word } i) + 1}$$

With Smoothing

$$\approx \frac{\text{\#spam containing word } i + 1}{\text{\#spam} + 2}$$

Similarly for other parameters  $P(x_i=0 \mid \text{spam})$ ,  $P(x_i=1 \mid \text{ham})$ ,  $P(x_i=0 \mid \text{ham})$ .

# Smoothing

$$\frac{\# \text{ham with abacus} + 1}{\# \text{ham with abacus} + 1 + \# \text{ham without abacus} + 1} = \frac{1}{\# \text{ham} + 2}$$

$$P(x_1=0 \text{ and } x_2=1 \text{ and } x_3=1 \text{ and } \dots | \text{ham}) = P(x_1=0 | \text{ham}) * P(x_2=1 | \text{ham}) * P(x_3=1 | \text{ham}) * \dots$$

(low probability)

## Dictionary

1. a
2. aardvark
3. abacus
4. abandon
5. abate

...

Suppose you are classifying an email containing the word “abacus,” which does not appear in any emails in your training data. What is  **$P(x_3 = 1 | \text{ham})$**  according to a naïve Bayes classifier **with smoothing**?

- A.  $P(x_3 = 1 | \text{ham}) = 0$
- B.  $P(x_3 = 1 | \text{ham}) = 1/2$
- C.  $P(x_3 = 1 | \text{ham}) = 1/(\text{total } \# \text{ham} + 1)$
- ☒ D.  $P(x_3 = 1 | \text{ham}) = 1/(\text{total } \# \text{ham} + 2)$  conditioned on ham
- E.  $P(x_3 = 1 | \text{ham}) = 1/(\text{total } \# \text{ham} + \text{total } \# \text{spam} + 2)$

# Smoothing

Suppose your training set includes only six emails, all of which are ham. For a new email, how will we estimate  $P(\text{ham})$

**without smoothing?**      **with smoothing?**

A.  $P(\text{ham}) = 0$

$P(\text{ham}) = 1/8$

B.  $P(\text{ham}) = 0$

$P(\text{ham}) = 6/7$

C.  $P(\text{ham}) = 1$

$P(\text{ham}) = 1/8$

D.  $P(\text{ham}) = 1$

$P(\text{ham}) = 6/7$

E.  $P(\text{ham}) = 1$

$P(\text{ham}) = 7/8$

Consider  
 $\# \text{ham} = 6k$

$$\frac{6k+1}{6k+2} \approx \frac{6k}{6k} \approx 1 > \frac{7}{8}$$

almost  $p=1$

why?  
we have more training data!

all my training were ham

$$P(\text{ham}) = \frac{\# \text{ham} + 1}{\# \text{ham} + 1 + \# \text{spam} + 1} = \frac{6+1}{6+2} = \frac{7}{8}$$

# Smoothing

Suppose your training set includes only six emails, all of which are ham. For a new email, how will we estimate  $P(\text{ham})$

**without smoothing?      with smoothing?**

A.  $P(\text{ham}) = 0$

$P(\text{ham}) = 1/8$

B.  $P(\text{ham}) = 0$

$P(\text{ham}) = 6/7$

C.  $P(\text{ham}) = 1$

$P(\text{ham}) = 1/8$

D.  $P(\text{ham}) = 1$

$P(\text{ham}) = 6/7$

E.  $P(\text{ham}) = 1$

$P(\text{ham}) = 7/8$

Is it still true that  $P(\text{ham}) + P(\text{spam}) = 1$  with smoothing?

$$P(\text{spam}) = \frac{\# \text{spam} + 1}{\# \text{spam} + 1 + \# \text{ham} + 1} = \frac{1}{\# \text{ham} + 2} = \frac{1}{8} \quad \frac{1}{8} + \frac{7}{8} = 1 \quad \checkmark$$



# Summary

- The Naive Bayes algorithm is useful for text classification.
- The bag of words model treats each word in a large dictionary as a feature.
- Smoothing is one modification that allows for better predictions when there are words that have never been seen before.