

# NAÏVE BAYES CONTINUED

David Kauchak  
CS159 Fall 2020

1

## Admin

Assignment 6b

No class Tuesday

Assignment 7 out Monday

2

## Final project

1. Your project should relate to something involving NLP
2. Your project must include a solid experimental evaluation
3. Your project should be in a group of 2-4. If you'd like to do it solo, please come talk to me.

3

## Final project

date	time	description
11/5	in-class	Project proposal presentation
11/11	11:59pm	Project proposal write-up
11/11	11:59pm	Status report
11/23	11:59pm	Paper draft
11/24	in-class	Presentation
11/25	11:59pm	Final paper and code

Read the final project handout ASAP!

Start forming groups and thinking about what you want to do

4

## Final project ideas

- pick a text classification task
  - evaluate different machine learning methods
  - implement a machine learning method
  - analyze different feature categories
- n-gram language modeling
  - implement and compare other smoothing techniques
  - implement alternative models
- parsing
  - lexicalized PCFG (with smoothing)
  - n-best list generation
  - parse output reranking
  - implement another parsing approach and compare
  - parsing non-traditional domains (e.g. twitter)
- EM
  - try and implement IBM model 2
  - word-level translation models

5

## Final project application areas

- spelling correction
- part of speech tagger
- text chunker
- dialogue generation
- pronoun resolution
- compare word similarity measures (more than the ones we looked at)
- word sense disambiguation
- machine translation
- information retrieval
- information extraction
- question answering
- summarization
- speech recognition

6

## Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

### Probabilistic models

Which model do we use, i.e. how do we calculate  $p(\text{feature}, \text{label})$ ?

How do train the model, i.e. how to we **estimate the probabilities** for the model?

How do we deal with overfitting?

7

## Naïve Bayes assumption

$$p(\text{features}, \text{label}) = p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})$$

$$p(x_i | y, x_1, x_2, \dots, x_{i-1}) = p(x_i | y)$$

Assumes feature  $i$  is independent of the the other features given the label

8

## Generative Story

To classify with a model, we're given an example and we obtain the probability

We can also ask how a given model would **generate** an example

This is the "generative story" for a model

Looking at the generative story can help understand the model

We also can use generative stories to help develop a model

9

## Bernoulli NB generative story

$$p(y) \prod_{j=1}^m p(x_j | y)$$

1. Pick a label according to  $p(y)$ 
  - roll a biased, num\_labels-sided die
2. For each feature:
  - Flip a *biased* coin:
    - if heads, include the feature
    - if tails, don't include the feature

What does this mean for text classification, assuming unigram features?

10

## Bernoulli NB generative story

$$p(y) \prod_{j=1}^m p(w_j | y)$$

1. Pick a label according to  $p(y)$ 
  - roll a biased, num\_labels-sided die
2. For each word in your vocabulary:
  - Flip a *biased* coin:
    - if heads, include the word in the text
    - if tails, don't include the word

11

## Bernoulli NB

### Pros

- Easy to implement
- Fast!
- Can be done on large data sets

### Cons

- Naïve Bayes assumption is generally not true
- Performance isn't as good as other models
- For text classification (and other sparse feature domains) the  $p(x_i=0 | y)$  can be problematic

12

### Another generative story

Randomly draw words from a "bag of words" until document length is reached

13

### Draw words from a fixed distribution

Selected:  $w_1$

14

### Draw words from a fixed distribution

Selected:  $w_1$   
Put a copy of  $w_1$  back

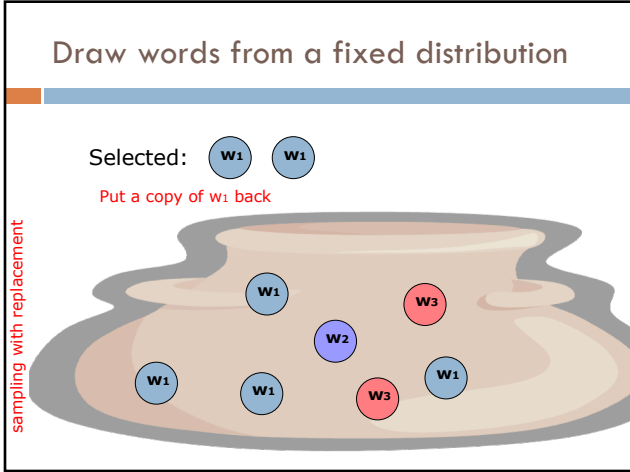
sampling with replacement

15

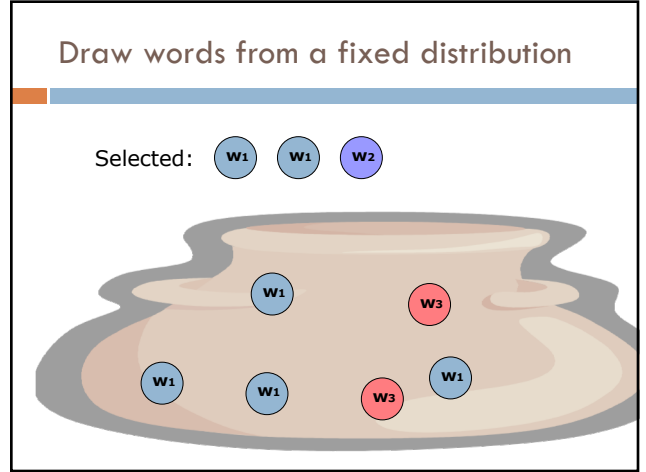
### Draw words from a fixed distribution

Selected:  $w_1$   $w_1$

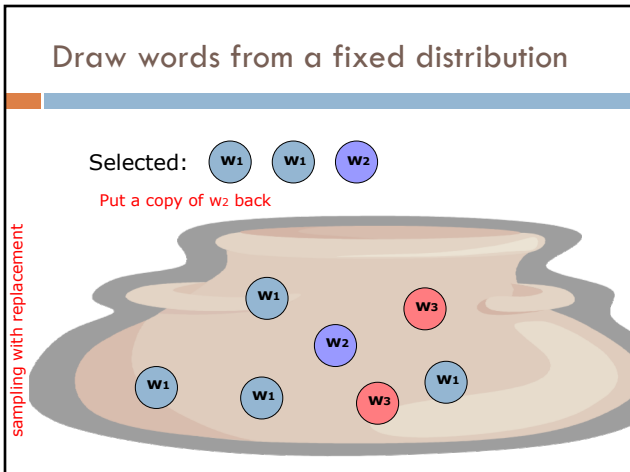
16



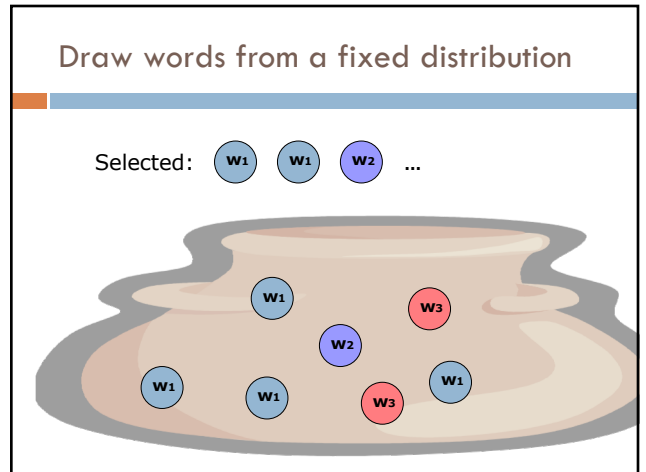
17



18



19



20

Draw words from a fixed distribution

Is this a NB model, i.e. does it assume each individual word occurrence is independent?

21

Draw words from a fixed distribution

Yes! Doesn't matter what words were drawn previously, still the same probability of getting any particular word

22

Draw words from a fixed distribution

Does this model handle multiple word occurrences?

23

Draw words from a fixed distribution

Selected:  $W_1$   $W_1$   $W_2$  ...

24

## NB generative story

<h3 style="text-align: center; color: blue;">Bernoulli NB</h3> <ol style="list-style-type: none"> <li>1. Pick a label according to <math>p(y)</math> <ul style="list-style-type: none"> <li>- roll a biased, num_labels-sided die</li> </ul> </li> <li>2. For each word in your vocabulary:           <ul style="list-style-type: none"> <li>- Flip a biased coin:               <ul style="list-style-type: none"> <li>- if heads, include the word in the text</li> <li>- if tails, don't include the word</li> </ul> </li> </ul> </li> </ol>	<h3 style="text-align: center; color: blue;">Multinomial NB</h3> <ol style="list-style-type: none"> <li>1. Pick a label according to <math>p(y)</math> <ul style="list-style-type: none"> <li>- roll a biased, num_labels-sided die</li> </ul> </li> <li>2. Keep drawing words from <math>p(\text{words}   y)</math> until text length has been reached.</li> </ol>
---	---

25

## Probabilities

<h3 style="text-align: center; color: blue;">Bernoulli NB</h3> <ol style="list-style-type: none"> <li>1. Pick a label according to <math>p(y)</math> <ul style="list-style-type: none"> <li>- roll a biased, num_labels-sided die</li> </ul> </li> <li>2. For each word in your vocabulary:           <ul style="list-style-type: none"> <li>- Flip a biased coin:               <ul style="list-style-type: none"> <li>- if heads, include the word in the text</li> <li>- if tails, don't include the word</li> </ul> </li> </ul> </li> </ol> $p(y) \prod_{j=1}^m p(x_j   y)$ <p>(1, 1, 1, 0, 0, 1, 0, 0, ...)</p>	<h3 style="text-align: center; color: blue;">Multinomial NB</h3> <ol style="list-style-type: none"> <li>1. Pick a label according to <math>p(y)</math> <ul style="list-style-type: none"> <li>- roll a biased, num_labels-sided die</li> </ul> </li> <li>2. Keep drawing words from <math>p(\text{words}   y)</math> until document length has been reached</li> </ol> <div style="text-align: center; color: red; font-size: 2em;">?</div> <p>(4, 1, 2, 0, 0, 7, 0, 0, ...)</p>
--	--

26

## A digression: rolling dice

What's the probability of getting a 3 for a single roll of this dice?

1/6

27


## A digression: rolling dice

What is the probability distribution over possible single rolls?

1/6	1/6	1/6	1/6	1/6	1/6
1	2	3	4	5	6

28

### A digression: rolling dice




What if I told you 1 was twice as likely as the others?

$\frac{2}{7}$	$\frac{1}{7}$	$\frac{1}{7}$	$\frac{1}{7}$	$\frac{1}{7}$	$\frac{1}{7}$
1	2	3	4	5	6

29

### A digression: rolling dice




What if I rolled 400 times and got the following number?

1: 100  
2: 50  
3: 50  
4: 100  
5: 50  
6: 50

$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$
1	2	3	4	5	6

30

### A digression: rolling dice




1. What is the probability of rolling a 1 and a 5 (in any order)?
2. Two 1s and a 5 (in any order)?
3. Five 1s and two 5s (in any order)?

$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$
1	2	3	4	5	6

31

### A digression: rolling dice



1. What is the probability of rolling a 1 and a 5 (in any order)?  
 $(\frac{1}{4} * \frac{1}{8}) * 2 = \frac{1}{16}$   
prob. of those two rolls      number of ways that can happen (1,5 and 5,1)
1. Two 1s and a 5 (in any order)?  
 $((\frac{1}{4})^2 * \frac{1}{8}) * 3 = \frac{3}{128}$
2. Five 1s and two 5s (in any order)?  
 $((\frac{1}{4})^5 * (\frac{1}{8})^2) * 21 = \frac{21}{524,288} = 0.00004$       **General formula?**

$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$
1	2	3	4	5	6

32



### Multinomial distribution

Multinomial distribution: independent draws over  $m$  possible categories

If we have frequency counts  $x_1, x_2, \dots, x_m$  over each of the categories, the probability is:

$$p(x_1, x_2, \dots, x_m | \theta_1, \theta_2, \dots, \theta_m) = \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m \theta_j^{x_j}$$

number of different ways to get those counts
probability of particular counts

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$\theta_5$	$\theta_6$	...
1	2	3	4	5	6	...

33

### Multinomial distribution

$$p(x_1, x_2, \dots, x_m | \theta_1, \theta_2, \dots, \theta_m) = \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m \theta_j^{x_j}$$

What are  $\theta_j$ ?

Are there any constraints on the values that they can take?

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$\theta_5$	$\theta_6$	...
1	2	3	4	5	6	...

34

### Multinomial distribution

$$p(x_1, x_2, \dots, x_m | \theta_1, \theta_2, \dots, \theta_m) = \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m \theta_j^{x_j}$$

$\theta_j$ : probability of rolling "j"

$\theta_j \geq 0$

$\sum_{j=1}^m \theta_j = 1$

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$\theta_5$	$\theta_6$	...
1	2	3	4	5	6	...

35

### Back to words...

Why the digression?

$$p(x_1, x_2, \dots, x_m | \theta_1, \theta_2, \dots, \theta_m) = \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m \theta_j^{x_j}$$

Drawing words from a bag is the same as rolling a die!

number of sides = number of words in the vocabulary

36

## Back to words...

Why the digression?

$$p(x_1, x_2, \dots, x_m | \theta_1, \theta_2, \dots, \theta_m) = \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m \theta_j^{x_j}$$

$$p(\text{features}, \text{label}) = p(y) \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m (\theta_y)^{x_j}$$

$\theta_y$  for class  $y$


37

## Basic steps for probabilistic modeling

Model each class as a multinomial:

$$p(\text{features}, \text{label}) = p(y) \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m (\theta_y)^{x_j}$$

Step 2: figure out how to estimate the probabilities for the model



How do we train the model, i.e. estimate  $\theta$ , for each class?

38

## A digression: rolling dice


What if I rolled 400 times and got the following number?


1: 100  
2: 50  
3: 50  
4: 100  
5: 50  
6: 50

1/4	1/8	1/8	1/4	1/8	1/8
1	2	3	4	5	6

39

## Training a multinomial


label1: 


label2: 

1/4	1/8	1/8	1/4	1/8	1/8
1	2	3	4	5	6

40

### Training a multinomial



label: 

For each label, y:

w1: 100 times  
w2: 50 times  
w3: 10 times  
w4: ...


$$\theta_j = \frac{\text{count}(w_j, y)}{\sum_{k=1}^m \text{count}(w_k, y)}$$

=  $\frac{\text{number of times word } w_j \text{ occurs in label } y \text{ docs}}{\text{total number of words in label } y \text{ docs}}$

1/4	1/8	1/8	1/4	1/8	1/8
1	2	3	4	5	6

41

### Classifying with a multinomial

 (10, 2, 6, 0, 0, 1, 0, 0, ...)

$w_1$   $w_2$   $w_3$   $w_4$   $w_5$   $w_6$   $w_7$   $w_8$


$p(y=1) \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m (\theta_j)^{x_j}$       $p(y=2) \frac{n!}{\prod_{j=1}^m x_j!} \prod_{j=1}^m (\theta_j)^{x_j}$

Any way I can make this simpler?

pick largest

42

### Classifying with a multinomial

 (10, 2, 6, 0, 0, 1, 0, 0, ...)

$w_1$   $w_2$   $w_3$   $w_4$   $w_5$   $w_6$   $w_7$   $w_8$

$p(y=1) \prod_{j=1}^m (\theta_j)^{x_j}$       $p(y=2) \prod_{j=1}^m (\theta_j)^{x_j}$

pick largest

$\frac{n!}{\prod_{j=1}^m x_j!}$  is a constant!

43

### Multinomial finalized

**Training:**

- Calculate p(label)
- For each label, calculate  $\theta$ s

$$\theta_j = \frac{\text{count}(w_j, y)}{\sum_{k=1}^m \text{count}(w_k, y)}$$

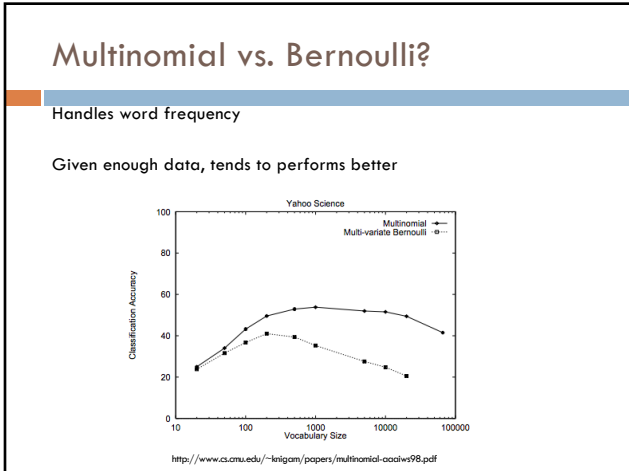
**Classification:**

- Get word counts
- For each label you had in training, calculate:

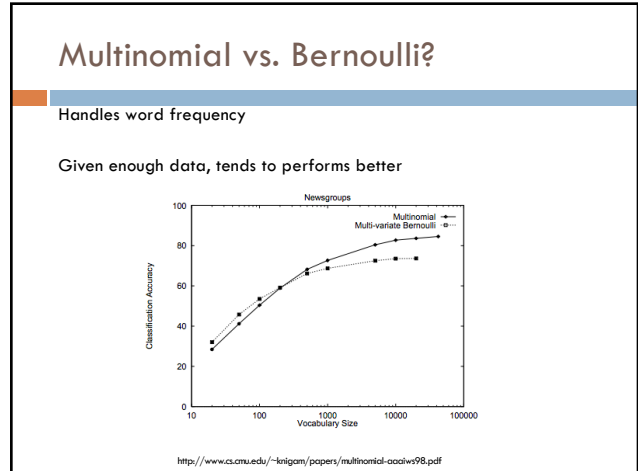
$$p(y) \prod_{j=1}^m \theta_j^{x_j}$$

and pick the largest

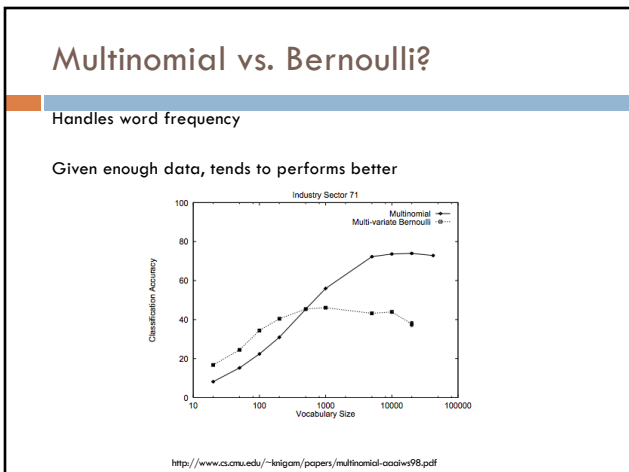
44



45



46



47

### Maximum likelihood estimation

Intuitive

Sets the probabilities so as to maximize the probability of the training data

**Problems?**

- Overfitting!
- Amount of data
  - particularly problematic for rare events
- Is our training data representative

48

### Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

**Probabilistic models**

Which model do we use, i.e. how do we calculate  $p(\text{feature}, \text{label})$ ?

How do train the model, i.e. how to we we **estimate the probabilities** for the model?

How do we deal with overfitting?

49

### Unseen events

training data

→

positive

negative

banana: 2

banana: 0

$$\theta_j = \frac{\text{count}(w_j, y)}{\sum_{k=1}^m \text{count}(w_k, y)}$$

What will  $\theta_{\text{banana}}$  be for the negative class?

50

### Unseen events

training data

→

positive

negative

banana: 2

banana: 0

$$\theta_j = \frac{\text{count}(w_j, y)}{\sum_{k=1}^m \text{count}(w_k, y)}$$

What will  $\theta_{\text{banana}}$  be for the negative class?

O! Is this a problem?

51

### Unseen events

training data

→

positive

negative

banana: 2

banana: 0

p("I ate a bad banana", negative) = ?

52

### Unseen events

training data → positive banana: 2  
negative banana: 0

$p(\text{"I ate a bad banana"}, \text{negative}) = 0$   
 $p(\text{"... banana ..."}, \text{negative}) = 0$

Solution?

53

### Add lambda smoothing

training data → positive banana: 2  
negative banana: 0

$$\theta_j = \frac{\text{count}(w_j, y)}{\sum_{k=1}^m \text{count}(w_k, y)}$$

$$\theta_j = \frac{\text{count}(w_j, y) + \lambda}{\lambda m + \sum_{k=1}^m \text{count}(w_k, y)}$$

for each label, pretend like we've seen each feature/word occur in  $\lambda$  additional examples

54

### Different than...

training data → positive banana: 0  
negative banana: 0

How is this problem different?

55

### Different than...

training data → positive banana: 0  
negative banana: 0

$p(\text{"I ate a bad banana"}, \text{positive})$  →  $p(\text{"I ate a bad"}, \text{positive})$   
 $p(\text{"I ate a bad banana"}, \text{negative})$  →  $p(\text{"I ate a bad"}, \text{negative})$

Out of vocabulary. Many ways to solve... for our implementation, we'll just ignore them.

56