

Modeling Natural Text

David Kauchak

CS458
Fall 2012

Admin

○ Final project

- Paper draft
 - due next Friday by midnight
 - Saturday, I'll e-mail out 1-2 paper drafts for you to read
 - Send me your reviews by Sunday at midnight
 - Monday morning, I'll forward these so you can integrate comments
- Initial code submission
 - Make sure to start integrating your code sooner than later
 - Initial code submission due next Friday

Admin

○ Final project continued

- At the beginning of class on Tuesday and Thursday we'll spend 15 min. discussing where things are at
- Any support from me?
 - let me know sooner than later...

Watson paper discussion

○ First application attempts

- How did the discussion go?
- One more paper discussion next Tuesday...

Modeling natural text

Questions

- what are the key topics in the text?
- what is the sentiment of the text?
- who/what does the article refer to?
- what are the key phrases?
- ...



Phenomena

- synonymy
- sarcasm/hyperbole
- variety of language (slang), misspellings
- coreference (e.g. pronouns like he/she)
- ...

Applications

search engines



- search
- advertising
- corporate databases

language generation



speech recognition

美 清 分
在 加 英
梅 知 思
月 年 爱
安 平 建
安 手 智

machine translation

I think, therefore I am
↓
I am

text simplification

text classification and clustering



SPAM



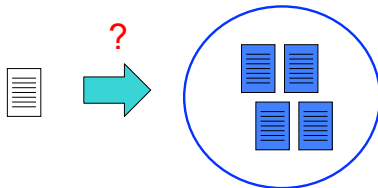
document hierarchies



sentiment analysis

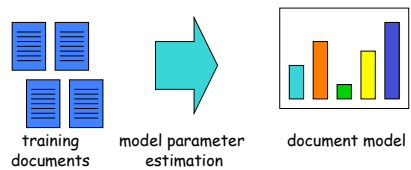
Document modeling: learn a probabilistic model of documents

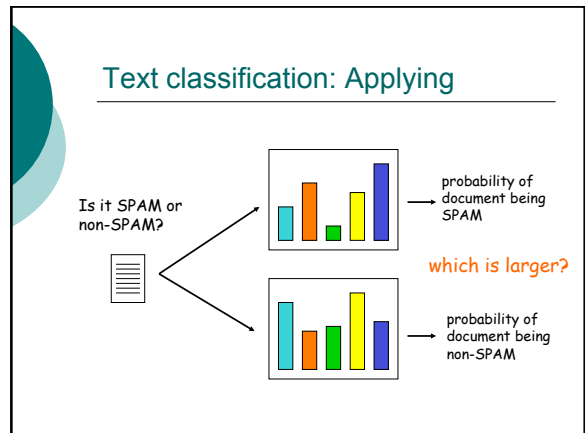
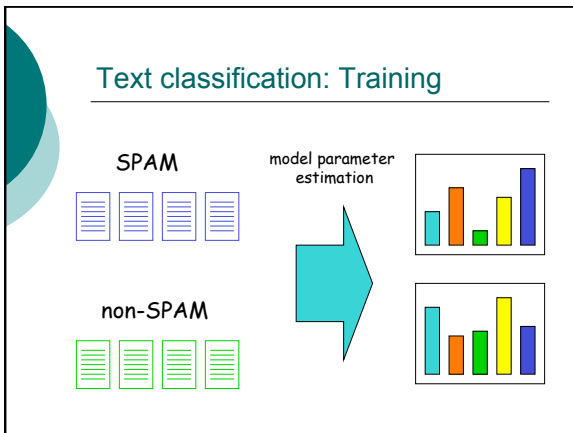
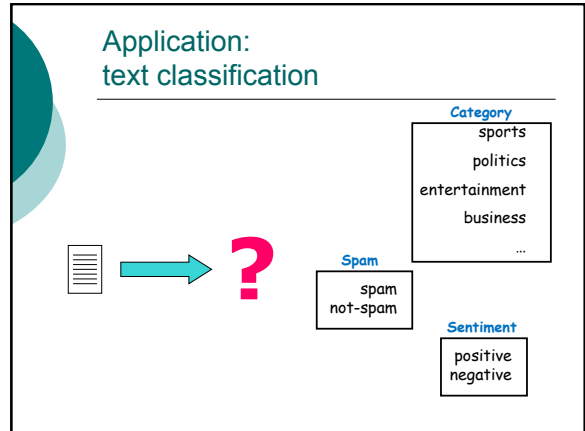
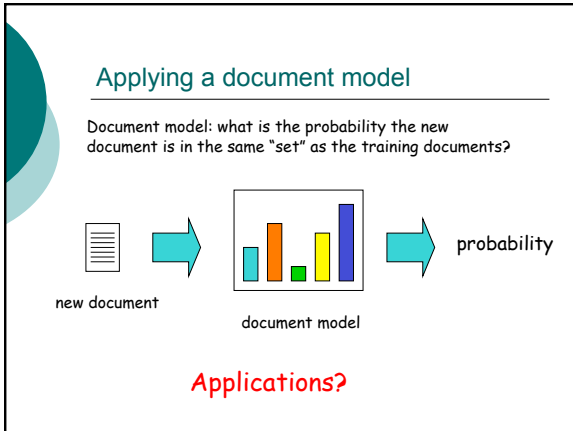
Predict the likelihood that an unseen document belongs to a set of documents



Model should capture text characteristics

Training a document model





Representation and Notation

Standard representation: bag of words

- Fixed vocabulary ~50K words
- Documents represented by a count vector, where each dimension represents the frequency of a word

Clinton said banana repeatedly last week on tv, "banana, banana, banana"

(4, 1, 1, 0, 0, 0, 1, 0, 0, ...)

banana
clinton
said
california
across
tv
wrong
capitol

Representation allows us to generalize across documents

Downside?

Representation and Notation

- Standard representation: bag of words
 - Fixed vocabulary ~50K words
 - Documents represented by a count vector, where each dimension represents the frequency of a word

Clinton said banana repeatedly last week on tv, "banana, banana, banana"

(4, 1, 1, 0, 0, 0, 1, 0, 0, ...)

banana
clinton
said
california
across
tv
wrong
capitol

Representation allows us to generalize across documents

Downside: lose word ordering information

Word burstiness

What is the probability that a political document contains the word "Clinton" *exactly* once?

The Stacy Koon-Lawrence Powell defense! The decisions of Janet Reno and Bill Clinton in this affair are essentially the moral equivalents of Stacy Koon's. ...

$p(\text{"Clinton"}=1|\text{political})= 0.12$

Word burstiness

What is the probability that a political document contains the word "Clinton" *exactly twice*?

The Stacy Koon-Lawrence Powell defense! The decisions of Janet Reno and Bill Clinton in this affair are essentially the moral equivalents of Stacy Koon's. Reno and Clinton have the advantage in that they investigate themselves.

$p(\text{"Clinton"}=2|\text{political})= 0.05$

Word burstiness in models

$$p(\text{"Clinton"}=2|\text{political}) = 0.05$$

Many models incorrectly predict:

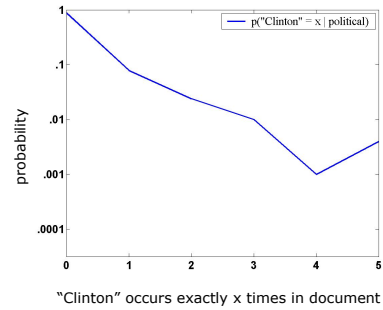
$$p(\text{"Clinton"}=2|\text{political}) \approx p(\text{"Clinton"}=1|\text{political})^2$$

$0.05 \neq \mathbf{0.0144} (0.12^2)$

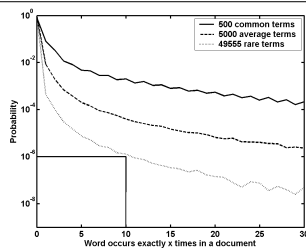
And in general, predict:

$$p(\text{"Clinton"}=i|\text{political}) \approx p(\text{"Clinton"}=1|\text{political})^i$$

$p(\text{"Clinton"} = x | \text{political})$

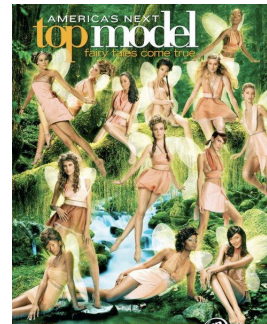


Word count probabilities




common words – 71% of word occurrences and 1% of the vocabulary
 average words – 21% of word occurrences and 10% of the vocabulary
 rare words – 8% of word occurrences and 89% of the vocabulary

The models...



Multinomial model



20 rolls of a fair, 6-side die - each number is equally probable

(1, 10, 5, 1, 2, 1)


ones twos threes fours fives sixes

(3, 3, 3, 3, 4, 4)

ones twos threes fours fives sixes

Which is more probable?

Multinomial model



20 rolls of a fair, 6-side die - each number is equally probable

(1, 10, 5, 1, 2, 1)


ones twos threes fours fives sixes

(3, 3, 3, 3, 4, 4)

ones twos threes fours fives sixes

How much more probable?

Multinomial model



20 rolls of a fair, 6-side die - each number is equally probable

(1, 10, 5, 1, 2, 1)

ones twos threes fours fives sixes

0.000000764

(3, 3, 3, 3, 4, 4)

ones twos threes fours fives sixes

0.000891

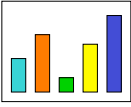
1000 times more likely

Multinomial model for text

Many more "sides" on the die than 6, but the same concept...

(4, 1, 1, 0, 0, 1, 0, 0, ...)

barbara clinton said california across tv wrapping capital



multinomial document model

probability

Generative Story

To apply a model, we're given a document and we obtain the probability

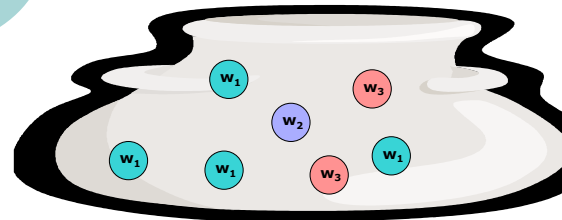
We can also ask how a given model would *generate* a document

This is the "generative story" for a model



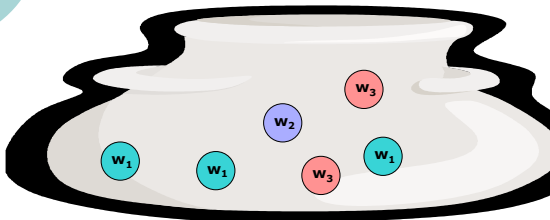
Multinomial Urn: Drawing words from a multinomial

Selected:



Drawing words from a multinomial

Selected: w_1

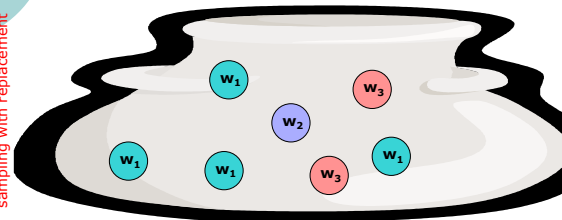


Drawing words from a multinomial

Selected: w_1

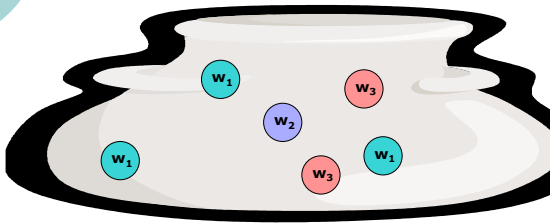
Put a copy of w_1 back

sampling with replacement



Drawing words from a multinomial

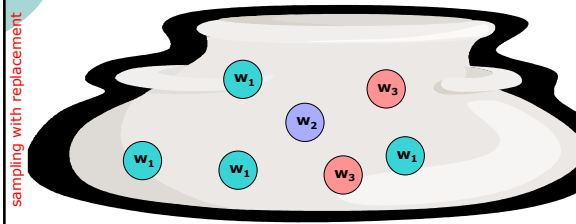
Selected: w_1 w_1



Drawing words from a multinomial

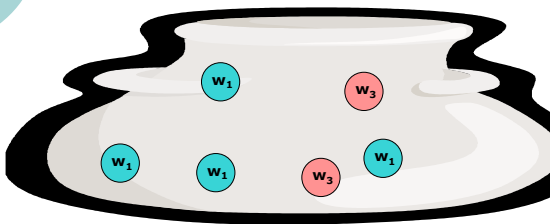
Selected: w_1 w_1

Put a copy of w_1 back



Drawing words from a multinomial

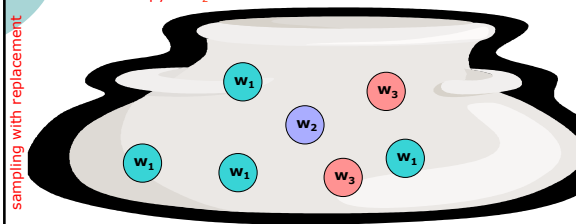
Selected: w_1 w_1 w_2



Drawing words from a multinomial

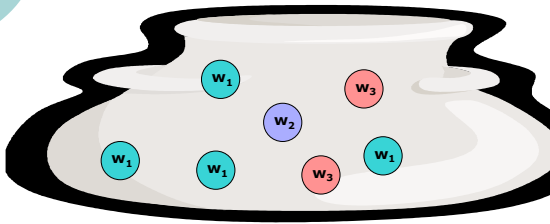
Selected: w_1 w_1 w_2

Put a copy of w_2 back



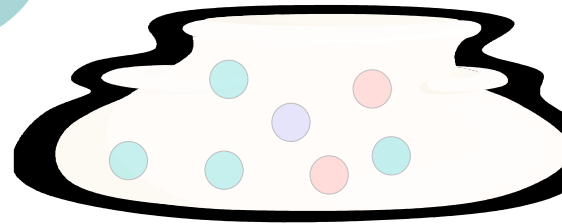
Drawing words from a multinomial

Selected: w_1 w_1 w_2 ...



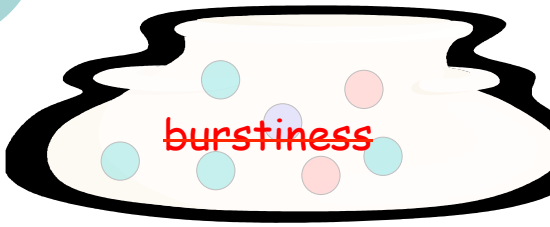
Drawing words from a multinomial

Does the multinomial model capture burstiness?



Drawing words from a multinomial

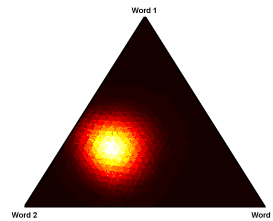
$p(\text{word})$ remains constant, independent of which words have already been drawn (in particular, how many of this particular word have been drawn)



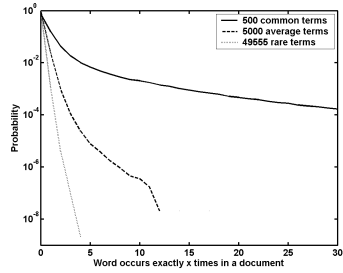
Multinomial probability simplex

Generate documents containing 100 words from a multinomial with just 3 possible words

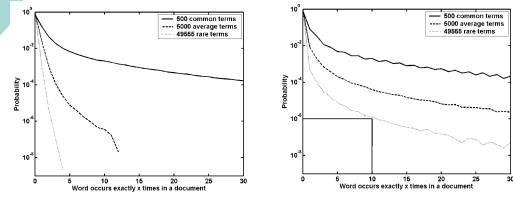
word 1 word 2 word 3
{0.31, 0.44, 0.25}



Multinomial word count probabilities



Multinomial does not model burstiness of average and rare words



Better model of burstiness: DCM

Dirichlet Compound Multinomial

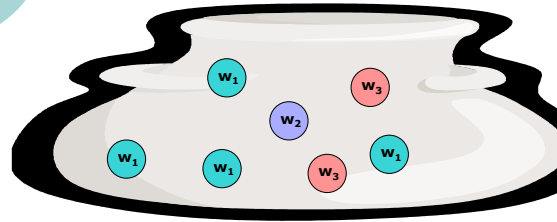
Polya Urn process

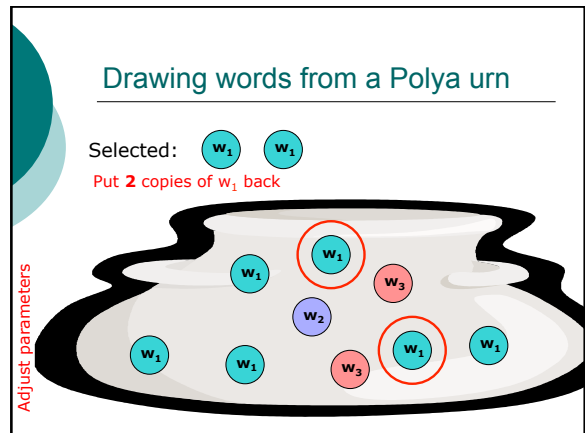
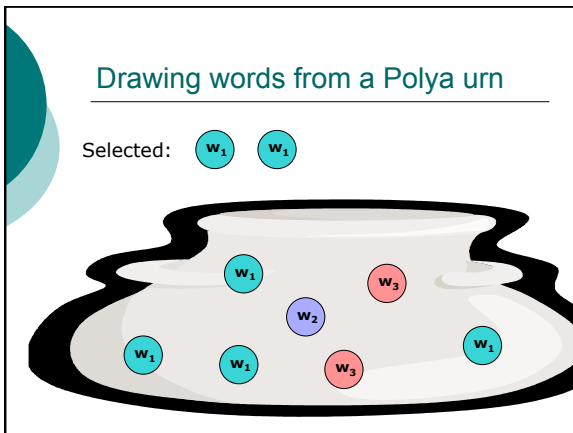
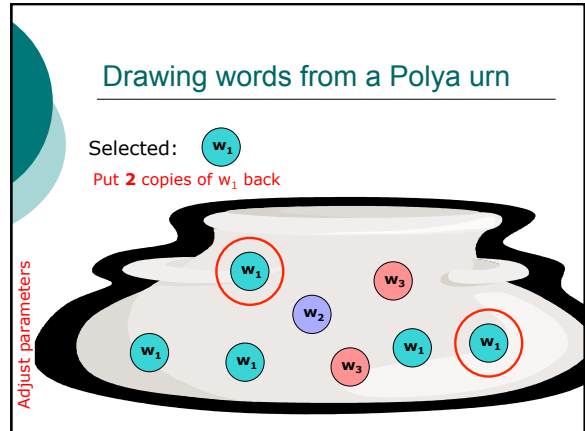
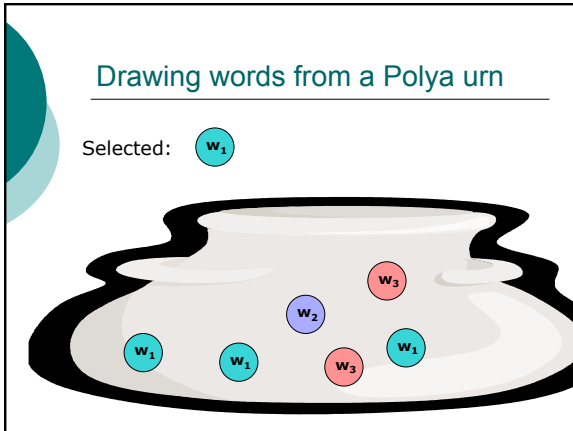
- **KEY:** Urn distribution changes based on previous words drawn
- Generative story:
 - Repeat until document length hit
 - Randomly draw a word from urn – call it w_i
 - Put **2** copies of w_i back in urn

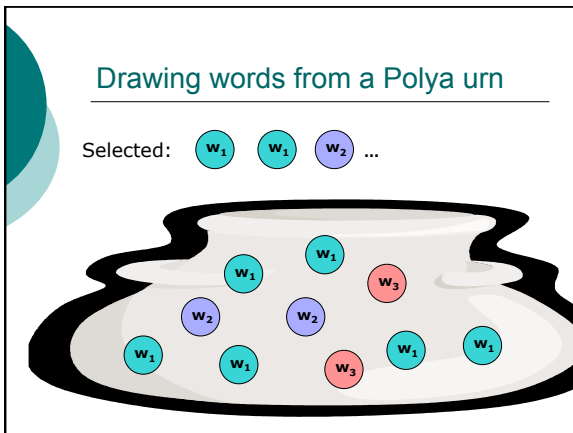
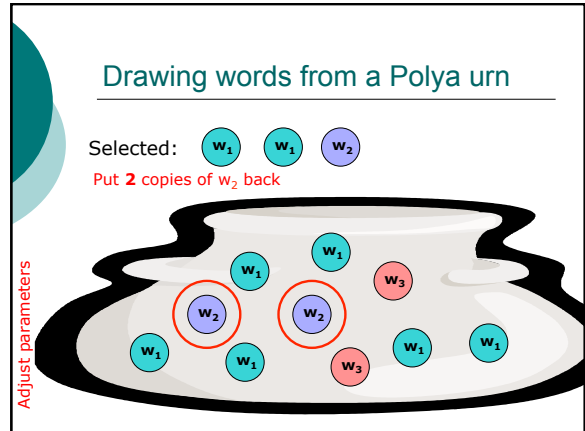
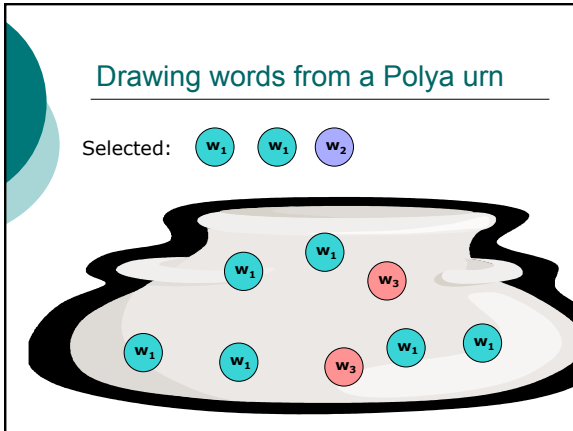


Drawing words from a Polya urn

Selected:







Polya urn

★ Words already drawn are more likely to be seen again

Results in the DCM distribution

We can modulate burstiness by increasing/ decreasing the number of words in the urn while keeping distribution the same

Controlling burstiness

Same distribution of words

Which is more bursty?

more bursty less bursty

Burstiness with DCM

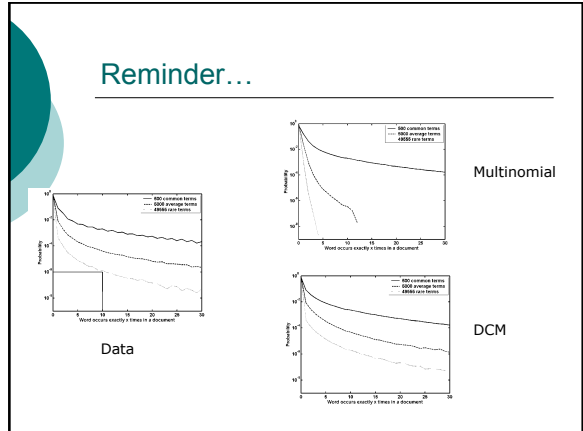
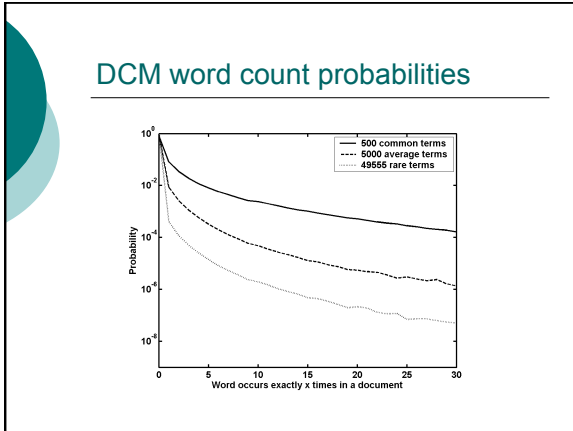
Multinomial

DCM

Down scaled
{.31, .44, .25}

Medium scaled
{.93, 1.32, .75}

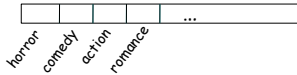
Up scaled
{2.81, 3.94, 2.25}



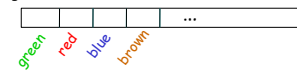
Modeling burstiness in other applications

Which model would be better: multinomial, DCM, other?

- User movie watching data



- Bags of M&Ms



- Daily Flight delays



A look at the code... multinomial model

Training

```
for i = 1:length(vectors)
    thetas(i,:) = log(sum(vector,1) + ones(1,size(vector,2))) -
        log(sum(sum(vector)) + size(vector,2));
end
```

Applying model

```
for i = 1:length(vectors)
    probs = thetas(:,idx) * vectors{i}(:,idx);
    [temp, decisions{i}] = max(probs);
end
```

DCM model

$$\begin{aligned}
 p(\mathbf{x} | \alpha) &= \int_{\theta} \frac{|\mathbf{x}|!}{\prod_{w=1}^W x_w!} \left(\prod_{w=1}^W \theta_w^{x_w} \right) \frac{\Gamma(\sum_{w=1}^W \alpha_w)}{\prod_{w=1}^W \Gamma(\alpha_w)} \prod_{w=1}^W \theta_w^{\alpha_w - 1} d\theta \\
 &= \frac{|\mathbf{x}|!}{\prod_{w=1}^W x_w!} \frac{\Gamma(\sum_{w=1}^W \alpha_w)}{\prod_{w=1}^W \Gamma(\alpha_w)} \int_{\theta} \prod_{w=1}^W \theta_w^{\alpha_w + x_w - 1} d\theta \\
 &= \frac{|\mathbf{x}|!}{\prod_{w=1}^W x_w!} \frac{\Gamma(\sum_{w=1}^W \alpha_w)}{\prod_{w=1}^W \Gamma(\alpha_w)} \frac{\Gamma(x_w + \alpha_w)}{\Gamma(\alpha_w)}
 \end{aligned}$$

Experiments

How can we test different models quantitatively?



Experiments

Modeling one class: document modeling

Modeling alternative classes: classification



Two standard data sets

Industry sector (web pages)

- More classes
- Less documents per class
- Longer documents

20 newsgroups (newsgroup posts)

- Fewer classes
- More documents per class
- Shorter documents

Modeling a single class: the fruit bowl

Mon Tue Wed Th Fri Sat Sun



Student 1



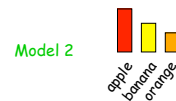
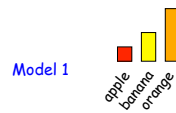
Student 2



Goal: predict what the fruit mix will be for the following Monday (assign probabilities to options)

Modeling a single class/group

How well does a model predict unseen data?



Monday

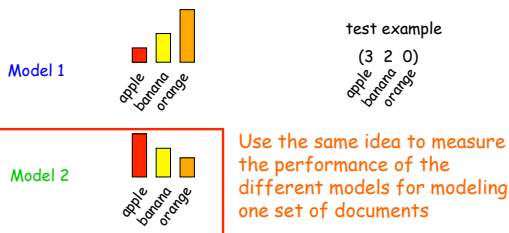
(3 2 0)

Which model is better?

How would you quantify how much better?

Modeling evaluation: perplexity

Perplexity is the average of the negative log of the model probabilities on test data



Perplexity results

20 newsgroups data set

Multinomial **92.1**
DCM **58.7**

Lower is better

- ideally the model would have a perplexity of 0!

Significant increase in modeling performance!

Classification results

Precision = number correct / number of documents

	Industry	20 Newsgroups
Multinomial	0.600	0.853
DCM	0.806	0.890

(results are on par with state of the art discriminative approaches!)

Next steps in text modeling

- Modeling textual phenomena like burstiness in text is important
- Better grounded models like DCM **ALSO** perform better in applications (e.g. classification)

Better models

text substitutability
relax bag of words constraint
(model co-occurrence)

Applications of models

multi-class data modeling
(e.g. clustering)
text similarity

hierarchical models

handling short phrases
(tweets, search queries)

language generation applications
(speech recognition,
translation, summarization)