

Statistical Language Models

Hongning Wang

CS@UVa

Today's lecture

1. How to represent a document?
 - Make it computable
2. How to infer the relationship among documents or identify the structure within a document?
 - Knowledge discovery

What is a statistical LM?

- A model specifying probability distribution over word sequences
 - $p(\textit{“Today is Wednesday”}) \approx 0.001$
 - $p(\textit{“Today Wednesday is”}) \approx 0.0000000000000001$
 - $p(\textit{“The eigenvalue is positive”}) \approx 0.00001$
- It can be regarded as a probabilistic mechanism for “generating” text, thus also called a “generative” model

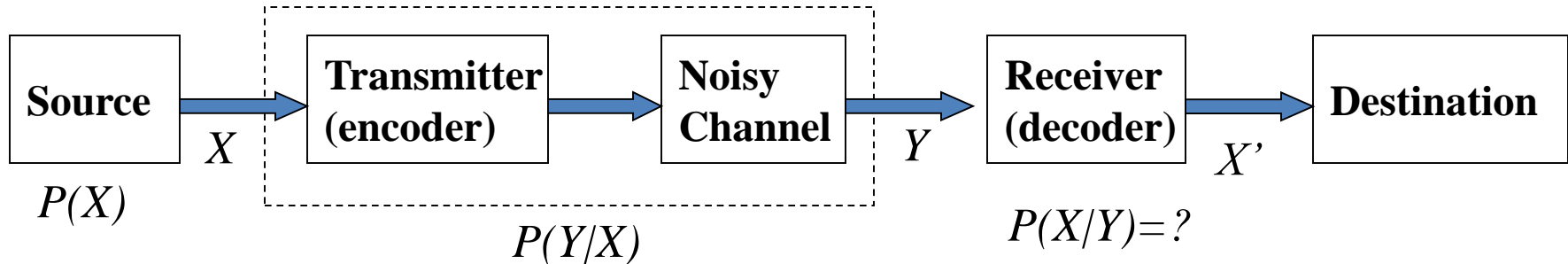
Why is a LM useful?

- Provide a principled way to quantify the uncertainties associated with natural language
- Allow us to answer questions like:
 - Given that we see “*John*” and “*feels*”, how likely will we see “*happy*” as opposed to “*habit*” as the next word?
(speech recognition)
 - Given that we observe “baseball” three times and “game” once in a news article, how likely is it about “sports” v.s. “politics”?
(text categorization)
 - Given that a user is interested in sports news, how likely would the user use “baseball” in a query?
(information retrieval)

Measure the fluency of documents

- How likely this document is generated by a given language model
 - If $p_{text-mining}(d) > p_{health}(d)$, d belongs to text mining related documents
 - If $p_{user_a}(d_1) > p_{user_a}(d_2)$, recommend d_1 to $user_a$
 - A relative concept

Source-Channel framework [Shannon 48]



$$\hat{X} = \arg \max_X p(X | Y) = \arg \max_X p(Y | X) p(X) \quad (\text{Bayes Rule})$$

When X is text, $p(X)$ is a language model

Many Examples:

Speech recognition:	X =Word sequence	Y =Speech signal
Machine translation:	X =English sentence	Y =Chinese sentence
OCR Error Correction:	X =Correct word	Y = Erroneous word
Information Retrieval:	X =Document	Y =Query
Summarization:	X =Summary	Y =Document

Basic concepts of probability

- Random experiment
 - An experiment with uncertain outcome (e.g., tossing a coin, picking a word from text)
- Sample space (S)
 - All possible outcomes of an experiment, e.g., tossing 2 fair coins, $S=\{HH, HT, TH, TT\}$
- Event (E)
 - $E \subseteq S$, E happens iff outcome is in S, e.g., $E=\{HH\}$ (all heads), $E=\{HH, TT\}$ (same face)
 - Impossible event ($\{\}$), certain event (S)
- Probability of event
 - $0 \leq P(E) \leq 1$

Recap: how to assign weights?

- Important!
- Why?
 - Corpus-wise: some terms carry more information about the document content
 - Document-wise: not all terms are equally important
- How?
 - Two basic heuristics
 - TF (Term Frequency) = Within-doc-frequency
 - IDF (Inverse Document Frequency)

Recap: TF-IDF weighting

- Combining TF and IDF
 - Common in doc \rightarrow high tf \rightarrow high weight
 - Rare in collection \rightarrow high idf \rightarrow high weight
 - $w(t, d) = TF(t, d) \times IDF(t)$
- Most well-known document representation schema in IR! (G Salton et al. 1983)



"Salton was perhaps the leading computer scientist working in the field of information retrieval during his time." - wikipedia

[Gerard Salton Award](#)

– highest achievement award in IR

Recap: cosine similarity

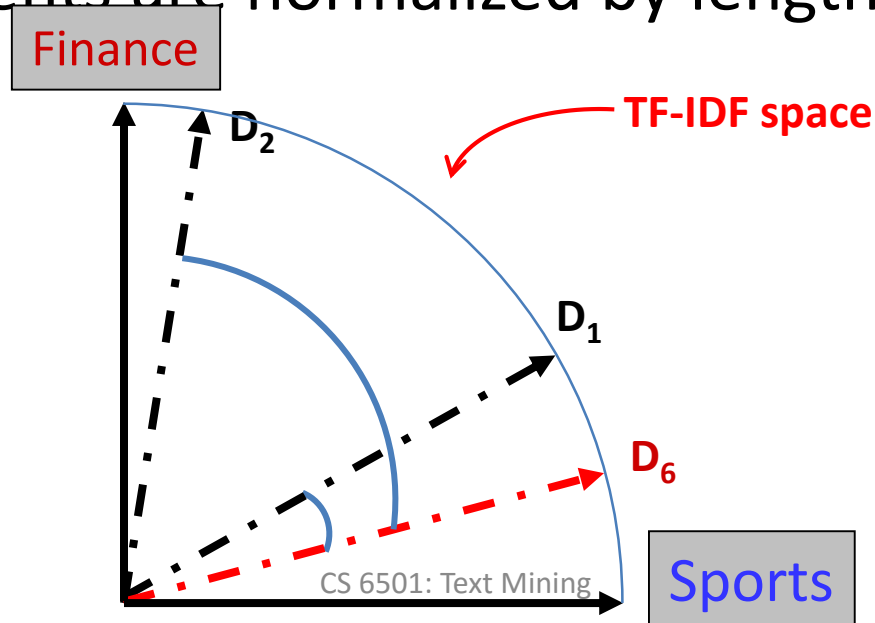
- Angle between two vectors

$$- \text{cosine}(d_i, d_j) = \frac{V_{d_i}^T V_{d_j}}{\|V_{d_i}\|_2 \|V_{d_j}\|_2}$$

TF-IDF vector

Unit vector

- Documents are normalized by length



Recap: disadvantages of VS model

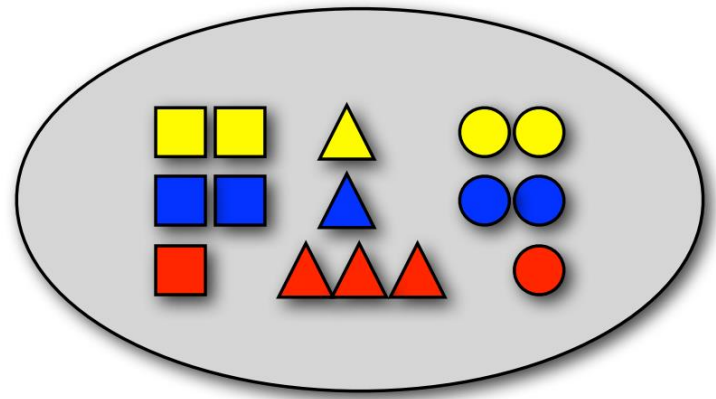
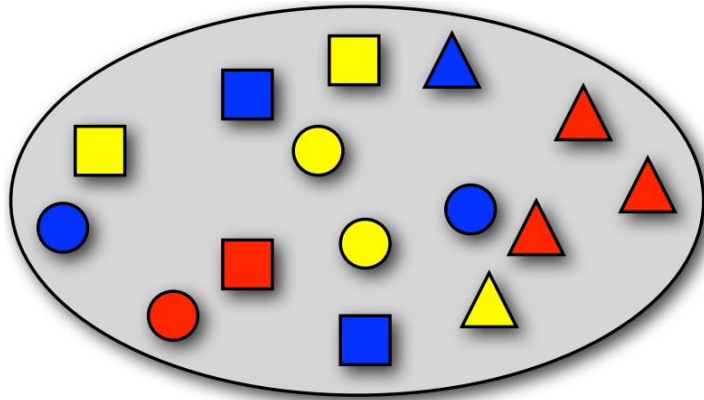
- Assume term independence
- Lack of “predictive adequacy”
 - Arbitrary term weighting
 - Arbitrary similarity measure
- Lots of parameter tuning!

Recap: what is a statistical LM?

- A model specifying probability distribution over word sequences
 - $p(\textit{“Today is Wednesday”}) \approx 0.001$
 - $p(\textit{“Today Wednesday is”}) \approx 0.0000000000000001$
 - $p(\textit{“The eigenvalue is positive”}) \approx 0.00001$
- It can be regarded as a probabilistic mechanism for “generating” text, thus also called a “generative” model

Sampling with replacement

- Pick a random shape, then put it back in the bag



$$P(\text{blue square}) = 2/15$$

$$P(\text{blue}) = 5/15$$

$$P(\text{blue} | \text{square}) = 2/5$$

$$P(\text{red square}) = 1/15$$

$$P(\text{red}) = 5/15$$


$$P(\text{square}) = 5/15$$

$$P(\text{red square or blue triangle}) = 2/15$$

$$P(\text{triangle} | \text{red}) = 3/5$$

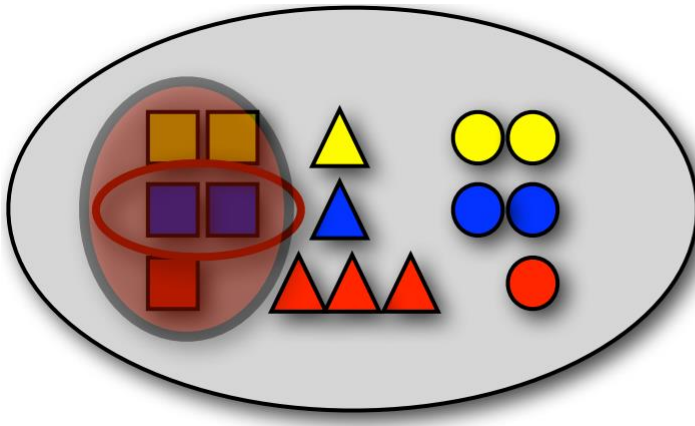
Essential probability concepts

- Probability of events
 - Mutually exclusive events
 - $P(A \cup B) = P(A) + P(B)$
 - General events
 - $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 - Independent events
 - $P(A \cap B) = P(A)P(B)$

 Joint probability, or
simply as $P(A, B)$

Essential probability concepts

- Conditional probability
 - $P(B|A) = P(A, B)/P(A)$
 - Bayes' Rule: $P(B|A) = P(A|B)P(B)/P(A)$
 - For independent events, $P(B|A) = P(B)$



$$P(\text{blue} \mid \blacksquare) = 2/5$$

Language model for text

We need independence assumptions!

- Probability distribution over word sequences

- $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2) \dots p(w_n|w_1, w_2, \dots, w_{n-1})$

- Complexity - $O(V^{n^*})$

- n^* - maximum ~~document~~ length sentence

Chain rule: from conditional probability to joint probability

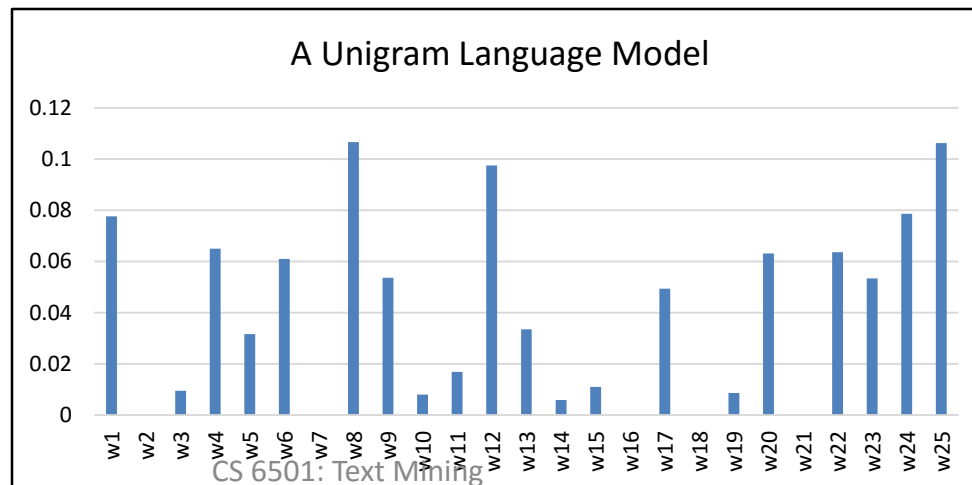
- 475,000 main headwords in Webster's Third New International Dictionary
 - Average English sentence length is 14.3 words
 - A rough estimate: $O(475000^{14})$

How large is this? $\frac{475000^{14}}{8\text{bytes} \times (1024)^4} \approx 3.38e^{66}TB$

Unigram language model

- Generate a piece of text by generating each word independently
 - $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2) \dots p(w_n)$
 - $s. t. \{p(w_i)\}_{i=1}^N, \sum_i p(w_i) = 1, p(w_i) \geq 0$
- Essentially a multinomial distribution over the vocabulary

The simplest and most popular choice!



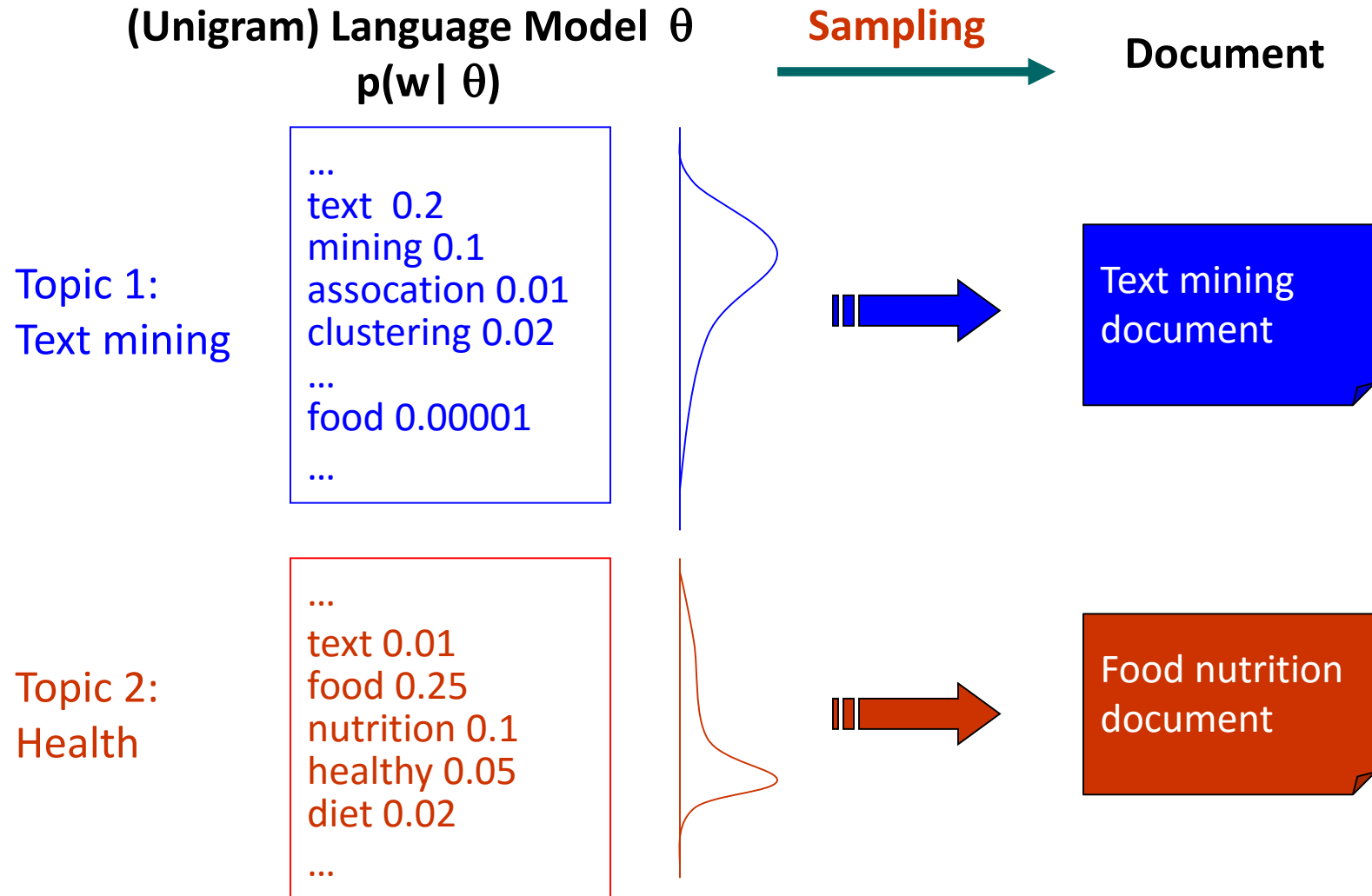
More sophisticated LMs

- N-gram language models
 - Conditioned only on the past $n-1$ words
 - E.g., bigram: $p(w_1 \dots w_n) = p(w_1)p(w_2|w_1)p(w_3|w_2) \dots p(w_n|w_{n-1})$
- Remote-dependence language models (e.g., Maximum Entropy model)
- Structured language models (e.g., probabilistic context-free grammar)

Why just unigram models?

- Difficulty in moving toward more complex models
 - They involve more parameters, so need more data to estimate
 - They increase the computational complexity significantly, both in time and space
- Capturing word order or structure may not add so much value for “topical inference”
- But, using more sophisticated models can still be expected to improve performance ...

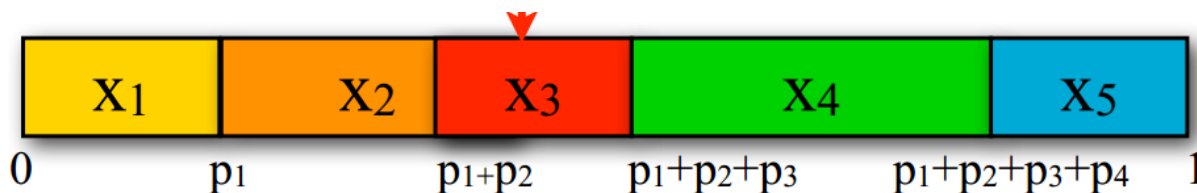
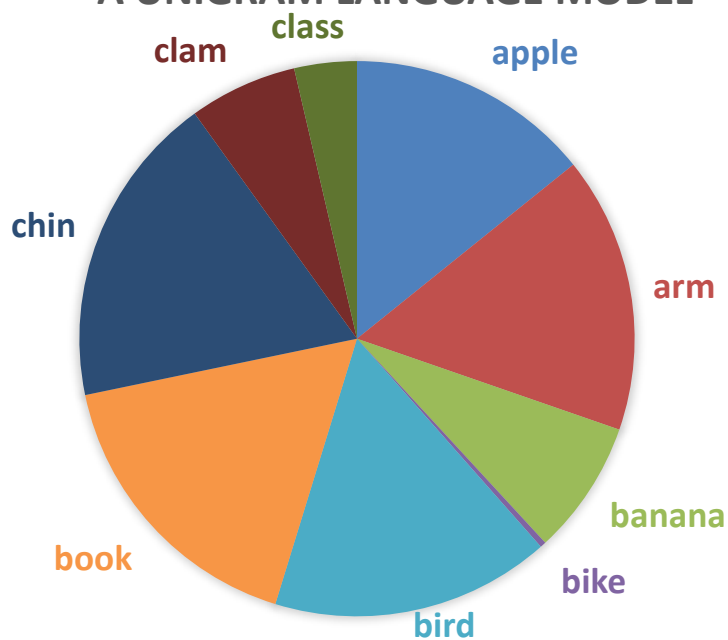
Generative view of text documents



How to generate text from an N-gram language model?

- Sample from a discrete distribution $p(X)$

A UNIGRAM LANGUAGE MODEL



Generating text from language models

$$P(\text{of}) = 3/66$$

$$P(\text{Alice}) = 2/66$$

$$P(\text{was}) = 2/66$$

$$P(\text{to}) = 2/66$$

$$P(\text{her}) = 2/66$$

$$P(\text{sister}) = 2/66$$

$$P(,) = 4/66$$

$$P(') = 4/66$$

Under a unigram language model:



Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

Generating text from language models

$$P(\text{of}) = 3/66$$

$$P(\text{Alice}) = 2/66$$

$$P(\text{was}) = 2/66$$

$$P(\text{to}) = 2/66$$

$$P(\text{her}) = 2/66$$

$$P(\text{sister}) = 2/66$$

$$P(,) = 4/66$$

$$P(') = 4/66$$

Under a unigram language model:



The same likelihood!

beginning by, very Alice but was and?
reading no tired of to into sitting
sister the, bank, and thought of without
her nothing: having conversations Alice
once do or on she it get the book her had
peeped was conversation it pictures or
sister in, 'what is the use had twice of
a book' 'pictures or' to

N-gram language models will help

Generated from language models of New York Times

- Unigram
 - Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a q acquire to six executives.
- Bigram
 - Last December through the way to preserve the Hudson corporation N.B.E.C. Taylor would seem to complete the major central planners one point five percent of U.S.E. has already told M.X. corporation of living on information such as more frequently fishing to keep her.
- Trigram
 - They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions.

Turing test: generating Shakespeare

A	<ul style="list-style-type: none"> • To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have • Every enter now severally so, let • Hill he late speaks; or! a more to leg less first you enter • Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like
B	<ul style="list-style-type: none"> • What means, sir. I confess she? then all sorts, he is trim, captain. • Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow. • What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first • El SClgen - An Automatic CS Paper Generator com- mand of fear not a liberal largess given away, Falstaff! Exeunt
C	<ul style="list-style-type: none"> • Sweet prince, Falstaff shall die. Harry of Monmouth's grave. • This shall forbid it should be branded, if renown made it empty. • Indeed the duke; and had a very good friend. • Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.
D	<ul style="list-style-type: none"> • King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in; • Will you not tell me who I am? • It cannot be but so. • Indeed the short and the long. Marry, 'tis a noble Lepidus.

Estimation of language models

Unigram Language Model θ

$$p(w | \theta) = ?$$

...
text ?
mining ?
association ?
database ?
...
query ?
...

Estimation



Document

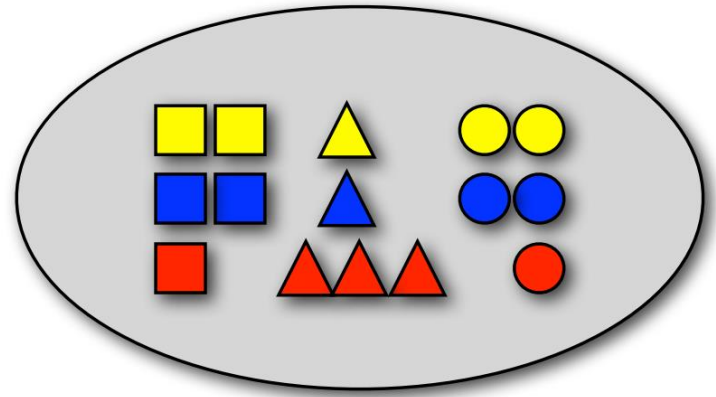
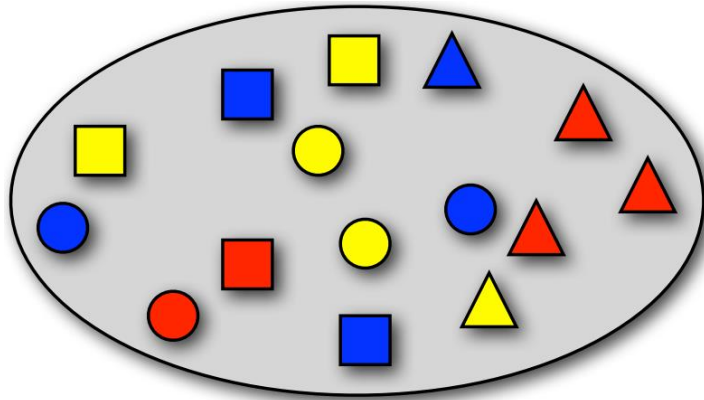
text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1



A “text mining” paper
(total #words=100)

Sampling with replacement

- Pick a random shape, then put it back in the bag



$$P(\text{blue square}) = 2/15$$

$$P(\text{blue}) = 5/15$$

$$P(\text{blue} | \text{square}) = 2/5$$

$$P(\text{red square}) = 1/15$$

$$P(\text{red}) = 5/15$$

$$P(\text{square}) = 5/15$$

$$P(\text{red or blue triangle}) = 2/15$$

$$P(\text{triangle} | \text{red}) = 3/5$$

Recap: language model for text

We need independence assumptions!

- Probability distribution over word sequences

- $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2) \dots p(w_n|w_1, w_2, \dots, w_{n-1})$

- Complexity - $O(V^{n^*})$

- n^* - maximum ~~document~~ length sentence

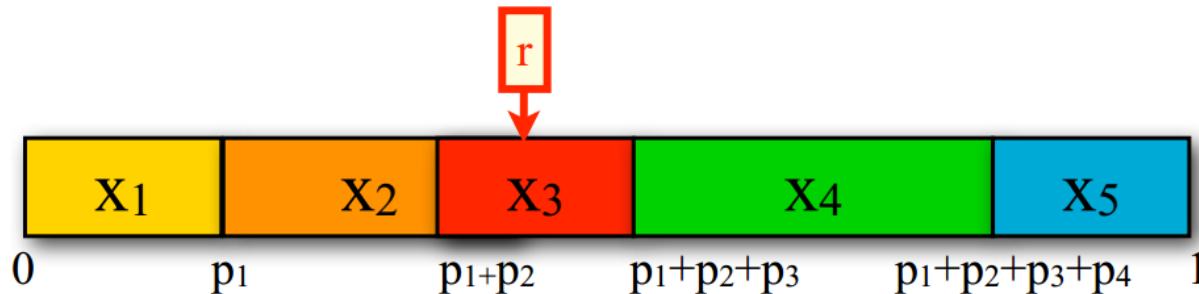
Chain rule: from conditional probability to joint probability

- 475,000 main headwords in Webster's Third New International Dictionary
 - Average English sentence length is 14.3 words
 - A rough estimate: $O(475000^{14})$

How large is this? $\frac{475000^{14}}{8\text{bytes} \times (1024)^4} \approx 3.38e^{66}TB$

Recap: how to generate text from an N-gram language model?

- Sample from a discrete distribution $p(X)$
 - Assume n outcomes in the event space X
 1. Divide the interval $[0,1]$ into n intervals according to the probabilities of the outcomes
 2. Generate a random number r between 0 and 1
 3. Return x_i where r falls into



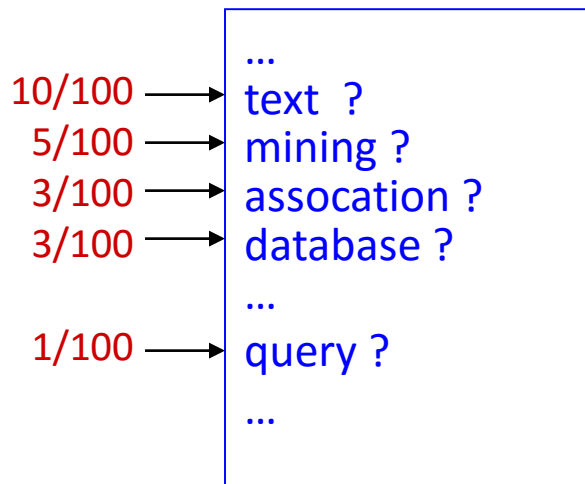
Recap: estimation of language models

- Maximum likelihood estimation

~~Unigram~~ Language Model θ
 $p(w | \theta) = ?$

Estimation

Document



A blue box containing the text of a document. The text is: text 10, mining 5, association 3, database 3, algorithm 2, query 1, efficient 1. The box has a folded bottom-right corner.

text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

A “text mining” paper
(total #words=100)

Parameter estimation

- General setting:
 - Given a (hypothesized & probabilistic) model that governs the random experiment
 - The model gives a probability of any data $p(X|\theta)$ that depends on the parameter θ
 - Now, given actual sample data $X=\{x_1,\dots,x_n\}$, what can we say about the value of θ ?
- Intuitively, take our best guess of θ -- “best” means “best explaining/fitting the data”
- Generally an optimization problem

Maximum likelihood vs. Bayesian

- Maximum likelihood estimation
 - “Best” means “data likelihood reaches maximum”

$$\hat{\theta} = \operatorname{argmax}_{\theta} P(X|\theta)$$

- Issue: small sample size

ML: Frequentist's point of view

- Bayesian estimation

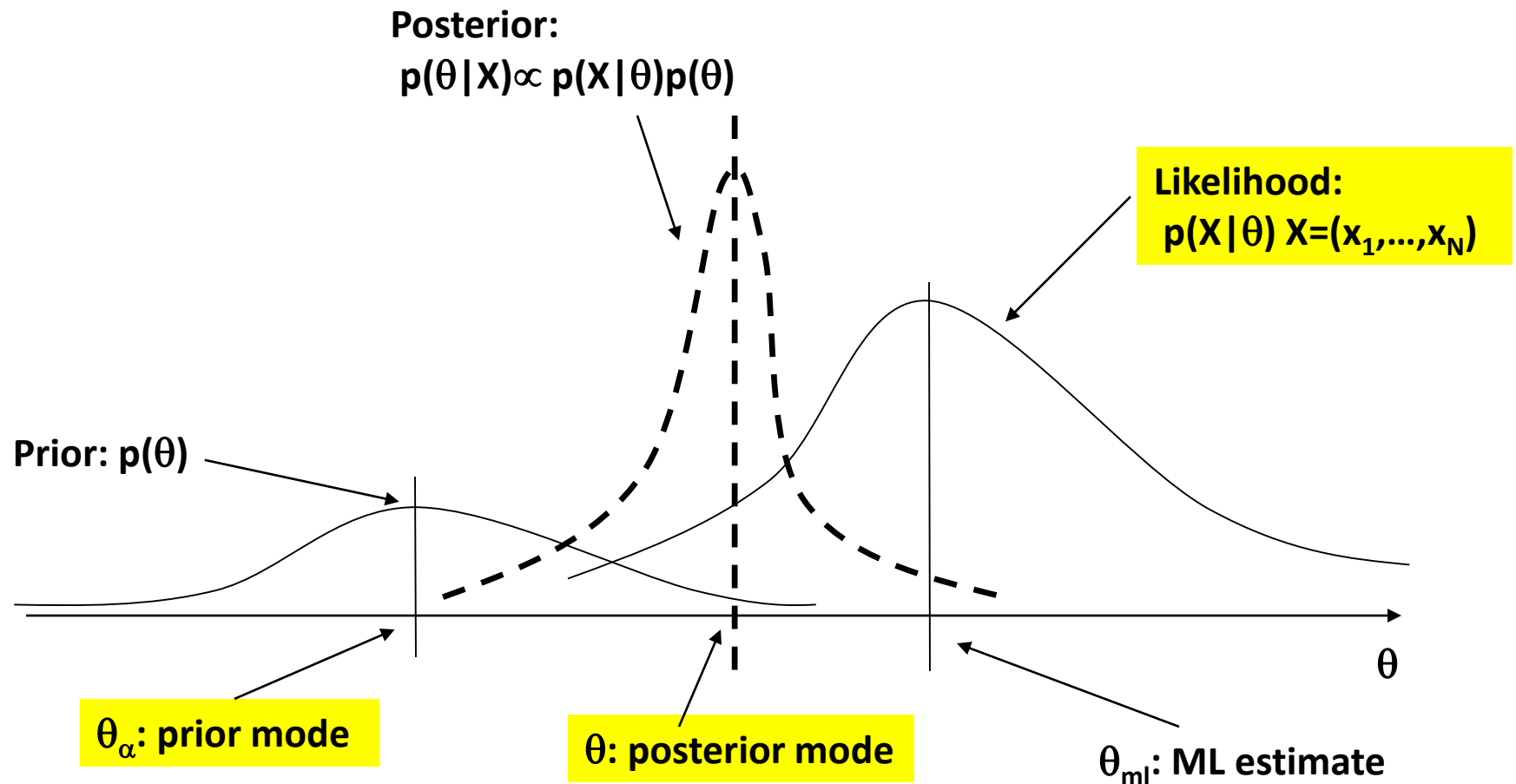
- “Best” means being consistent with our “prior” knowledge and explaining data well

$$\hat{\theta} = \operatorname{argmax}_{\theta} P(\theta|X) = \operatorname{argmax}_{\theta} P(X|\theta)P(\theta)$$

- A.k.a, Maximum a Posterior estimation
 - Issue: how to define prior?

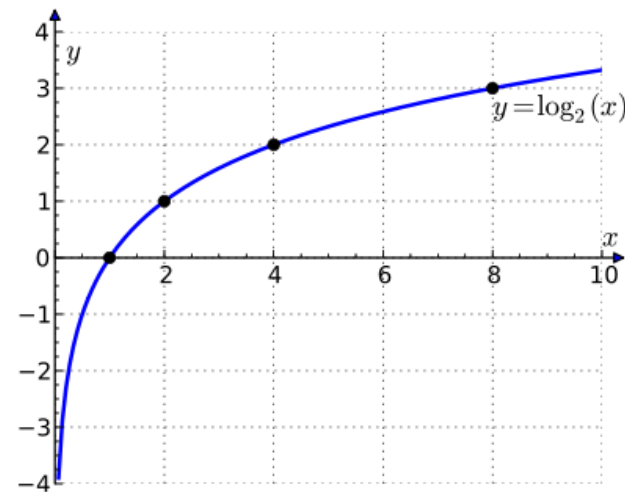
MAP: Bayesian's point of view

Illustration of Bayesian estimation



Maximum likelihood

- Data: a collection of words, w_1, w_2, \dots
- Model: multinomial distribution $p(W|\theta)$
- Maximum likelihood estimator: $\hat{\theta} = \operatorname{argmax}_{\theta} p(W|\theta)$



$$p(W|\theta) = \binom{N}{c(w_1), \dots, c(w_N)} \prod_{i=1}^N \theta_i^{c(w_i)} \propto \prod_{i=1}^N \theta_i^{c(w_i)} \Rightarrow \log p(W|\theta) = \sum_{i=1}^N c(w_i) \log \theta_i$$

$$\Rightarrow L(W, \theta) = \sum_{i=1}^N c(w_i) \log \theta_i + \lambda \left(\sum_{i=1}^N \theta_i - 1 \right)$$

$$\Rightarrow \frac{\partial L}{\partial \theta_i} = \frac{c(w_i)}{\theta_i} + \lambda \rightarrow \theta_i = -\frac{c(w_i)}{\lambda}$$

$$\Rightarrow \text{Since } \sum_{i=1}^N \theta_i = 1 \text{ we have } \lambda = -\sum_{i=1}^N c(w_i)$$

$$\Rightarrow \theta_i = \frac{c(w_i)}{\sum_{i=1}^N c(w_i)}$$

Using Lagrange multiplier approach, we'll tune θ_i to maximize $L(W, \theta)$

Set partial derivatives to zero

Requirement from probability


ML estimate

Maximum likelihood estimation

- For N-gram language models

$$- p(w_i | w_{i-1}, \dots, w_{i-n+1}) = \frac{c(w_i, w_{i-1}, \dots, w_{i-n+1})}{c(w_{i-1}, \dots, w_{i-n+1})}$$

- $c(\emptyset) = N$

 *Length of document or total
number of words in a corpus*

Problem with MLE

- Unseen events

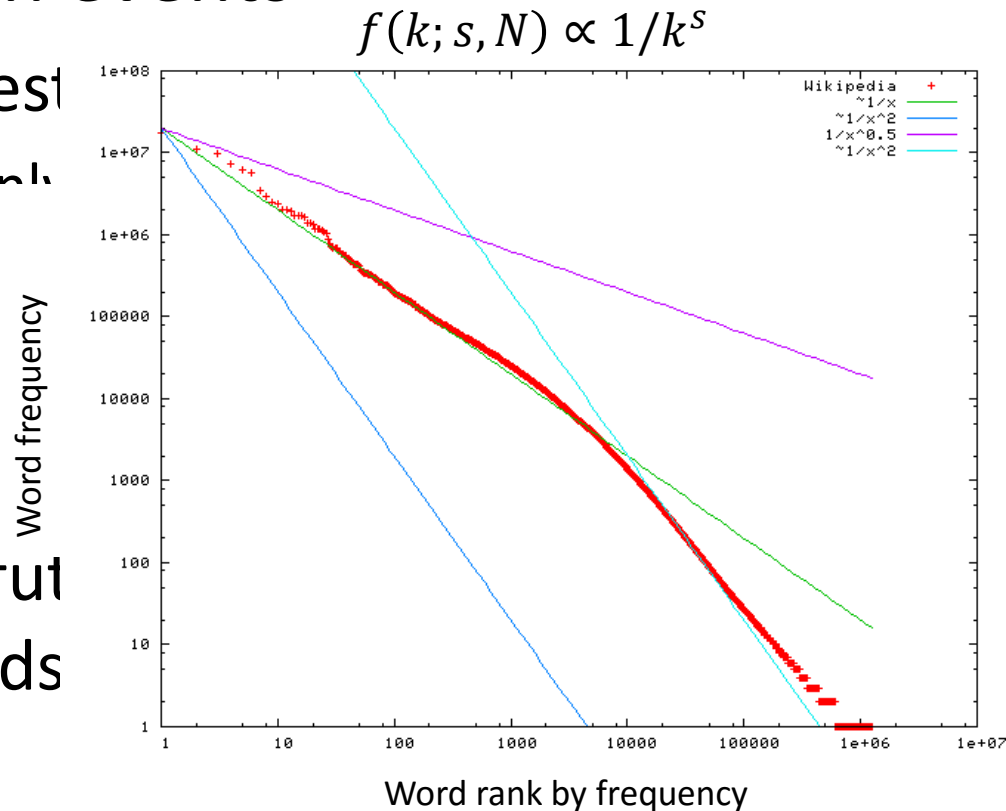
– We estimate

- $O \sim 1/k$

- $O \sim 1/k^2$

➤ This is the frequency of words in the training set

➤ No test words



A plot of word frequency in Wikipedia (Nov 27, 2006)

tokens, but:

red

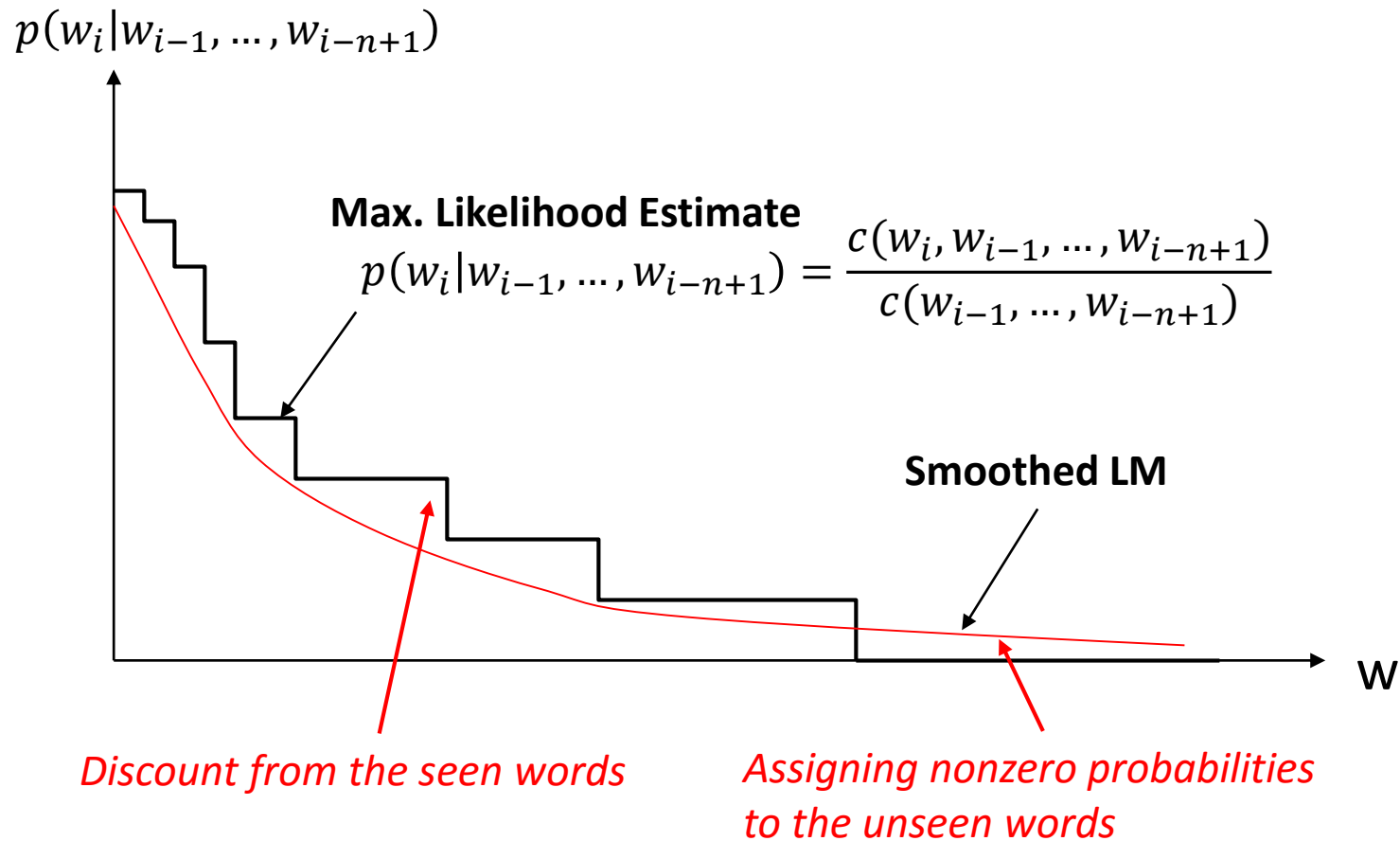
es not occur
ty!

ise unseen

Smoothing

- If we want to assign non-zero probabilities to unseen words
 - Unseen words = new words, new N-grams
 - Discount the probabilities of observed words
- General procedure
 1. Reserve some probability mass of words seen in a document/corpus
 2. Re-allocate it to unseen words

Illustration of N-gram language model smoothing



Smoothing methods

- Additive smoothing
 - Add a constant δ to the counts of each word
 - Unigram language model as an example

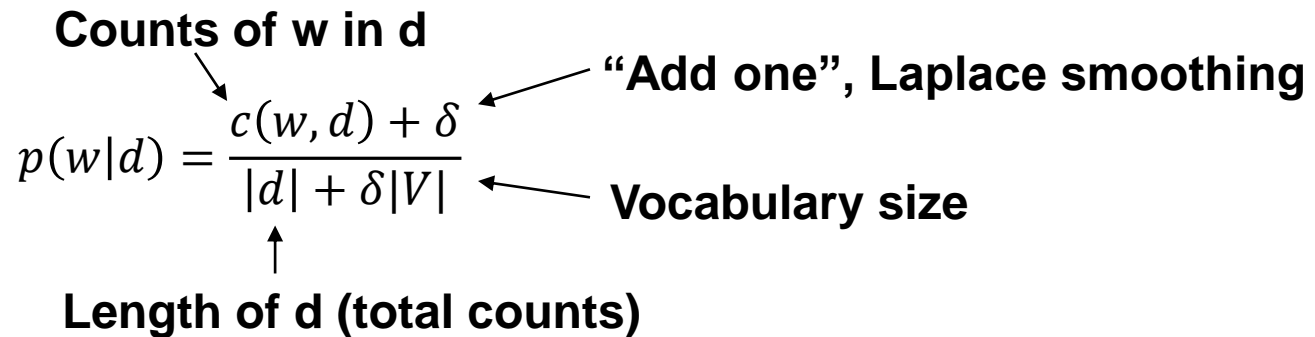
Counts of w in d

“Add one”, Laplace smoothing

$$p(w|d) = \frac{c(w, d) + \delta}{|d| + \delta|V|}$$

Vocabulary size

Length of d (total counts)



- Problems?
 - Hint: all words are equally important?

Add one smoothing for bigrams

Original:

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Smoothed:

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

After smoothing

- Giving too much to the unseen events

Original:

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Smoothed:

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.00062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

Refine the idea of smoothing

- Should all unseen words get equal probabilities?
- We can use a reference model to discriminate unseen words

$$p(w | d) = \begin{cases} p_{seen}(w | d) & \text{if } w \text{ is seen in } d \\ \alpha_d p(w | REF) & \text{otherwise} \end{cases}$$

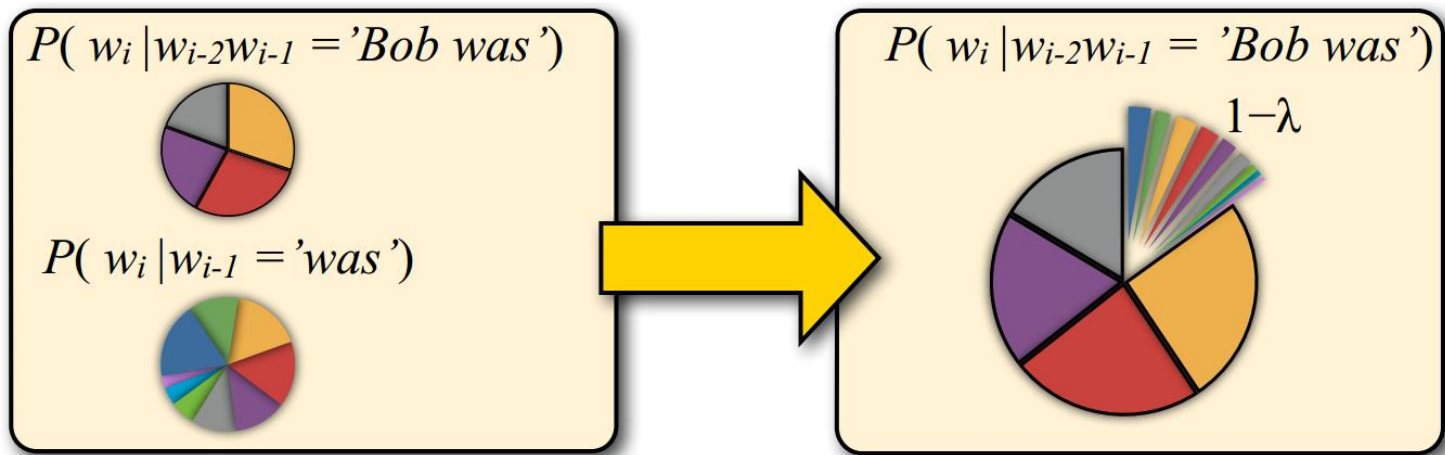
Discounted ML estimate

Reference language model

$$\alpha_d = \frac{1 - \sum_{w \text{ is seen}} p_{seen}(w | d)}{\sum_{w \text{ is unseen}} p(w | REF)}$$

Smoothing methods

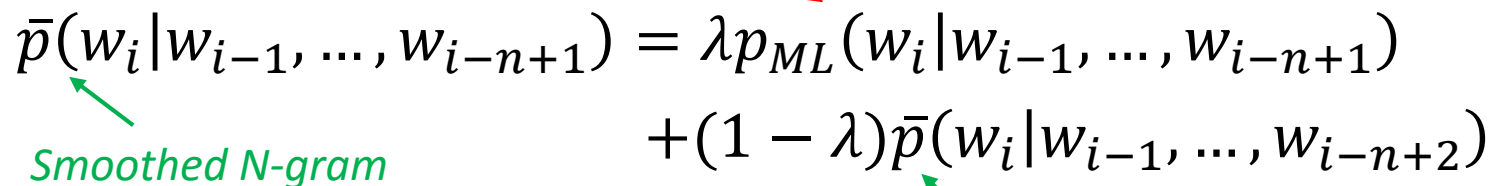
- Linear interpolation
 - Use $(N - 1)$ -gram probabilities to smooth N -gram probabilities
 - We never see the trigram “Bob was reading”, but we do see the bigram “was reading”



Smoothing methods

- Linear interpolation
 - Use (N – 1)-gram probabilities to smooth N-gram probabilities

$$\bar{p}(w_i | w_{i-1}, \dots, w_{i-n+1}) = \lambda p_{ML}(w_i | w_{i-1}, \dots, w_{i-n+1}) + (1 - \lambda) \bar{p}(w_i | w_{i-1}, \dots, w_{i-n+2})$$

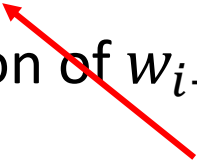


- Further generalize it

$$\bar{p}(w_i | w_{i-1}, \dots, w_{i-n+1}) = \lambda_1 p_{ML}(w_i | w_{i-1}, \dots, w_{i-n+1}) + \lambda_2 \bar{p}(w_i | w_{i-1}, \dots, w_{i-n+2}) + \dots + \lambda_n \bar{p}(w_i)$$

$$s. t. \sum_{i=1}^n \lambda_i = 1$$

Smoothing methods

- Linear interpolation
 - Estimating λs
 - Using a hold-out data set to find the optimal λs
 - An evaluation metric is needed to define “optimality”
 - Define λ as a function of $w_{i-1}, \dots, w_{i-n+2}$
-  We will come back to this later

Smoothing methods

- Absolute discounting
 - Subtract a constant δ from each nonzero n-gram count and then interpolate

$$\bar{p}(w_i | w_{i-1}, \dots, w_{i-n+1}) = \frac{\max(c(w_i, w_{i-1}, \dots, w_{i-n+1}) - \delta, 0)}{c(w_{i-1}, \dots, w_{i-n+1})}$$

Smoothed N-gram

$$+ \lambda \bar{p}(w_i | w_{i-1}, \dots, w_{i-n+2})$$

- If S seen word types occur after $w_{i-1}, \dots, w_{i-n+1}$ in the training data, this reserves the probability mass

$$p(u) = \frac{\delta S}{c(w_{i-1}, \dots, w_{i-n+1})}$$

to be reallocated according to

$$\bar{p}(w_i | w_{i-1}, \dots, w_{i-n+2})$$

- $\lambda = \frac{\delta S}{c(w_{i-1}, \dots, w_{i-n+1})}$

Language model evaluation

- Train the models on the same training set
 - Parameter tuning can be done by holding off some training set for validation
- Test the models on an unseen test set
 - This data set must be disjoint from training data
- Language model A is better than model B
 - If A assigns higher probability to the test data than B

Perplexity

- Standard evaluation metric for language models
 - A function of the probability that a language model assigns to a data set
 - Rooted in the notion of cross-entropy in information theory

Perplexity

- The inverse of the likelihood of the test set as assigned by the language model, normalized by the number of words

$$PP(w_1, \dots, w_N) = \sqrt[N]{\frac{1}{\prod_{i=1}^N p(w_i | w_{i-1}, \dots, w_{i-n+1})}}$$

↖ N-gram language model

An experiment

- Models
 - Unigram, Bigram, Trigram models (with proper smoothing)
- Training data
 - 38M words of WSJ text (vocabulary: 20K types)
- Test data
 - 1.5M words of WSJ text
- Results

	Unigram	Bigram	Trigram
Perplexity	962	170	109

What you should know

- N-gram language models
- How to generate text documents from a language model
- How to estimate a language model
- General idea and different ways of smoothing
- Language model evaluation

Today's reading

- Introduction to information retrieval
 - Chapter 12: Language models for information retrieval
- Speech and Language Processing
 - Chapter 4: N-Grams