


N-Gram – Part 2

ICS 482 Natural Language  
Processing



Lecture 8: N-Gram – Part 2

Husni Al-Muhtaseb

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

# ICS 482 Natural Language Processing



Lecture 8: N-Gram – Part 2  
Husni Al-Muhtaseb

# NLP Credits and Acknowledgment

These slides were adapted from  
presentations of the Authors of  
the book

**SPEECH and LANGUAGE PROCESSING:  
An Introduction to Natural Language Processing,  
Computational Linguistics, and Speech Recognition**

and some modifications from  
presentations found in the WEB  
by several scholars including the  
following

# NLP Credits and Acknowledgment



If your name is missing please contact me  
muhtaseb  
At  
Kfupm.  
Edu.  
sa

# NLP Credits and Acknowledgment

Husni Al-Muhtaseb

James Martin

Jim Martin

Dan Jurafsky

Sandiway Fong

Song youngin

Paula Matuszek

Mary-Angela

Papalaskari

Dick Crouch

Tracy Kin

L. Venkata

Subramaniam

Martin Volk

Bruce R. Maxim

Jan Hajič

Srinath Srinivasa

Simeon Ntafos

Paolo Pirjanian

Ricardo Vilalta

Tom Lenaerts

Heshaam Feili

Björn Gambäck

Christian Korthals

Thomas G.

Dietterich

Devika

Subramanian

Duminda

Wijesekera

Lee McCluskey

David J.

Kriegman

Kathleen

McKeown

Michael J. Ciaraldi

David Finkel

Min-Yen Kan

Andreas Geyer-  
Schulz

Franz J. Kurfess

Tim Finin

Nadjet Bouayad

Kathy McCoy

Khurshid Ahmad

Staffan Larsson

Robert Wilensky

Feiyu Xu

Jakub Piskorski

Rohini Srihari

Mark Sanderson

Andrew Elks

Marc Davis

Ray Larson

Jimmy Lin

Marti Hearst

Andrew

McCallum

Nick Kushmerick

Mark Craven

Chia-Hui Chang

Diana Maynard

James Allan

Martha Palmer

julia hirschberg

Elaine Rich

Christof Monz

Bonnie J. Dorr

Nizar Habash

Massimo Poesio

David Goss-

Grubbs

Thomas K Harris

John Hutchins

Alexandros

Potamianos

Mike Rosner

Latifa Al-Sulaiti

Giorgio Satta

Jerry R. Hobbs

Christopher

Manning

Hinrich Schütze

Alexander

Gelbukh

Gina-Anne Levow

Guitao Gao

Qing Ma

# Previous Lectures

---

- Pre-start questionnaire
- Introduction and Phases of an NLP system
- NLP Applications - Chatting with Alice
- Finite State Automata & Regular Expressions & languages
- Deterministic & Non-deterministic FSAs
- Morphology: Inflectional & Derivational
- Parsing and Finite State Transducers
- Stemming & Porter Stemmer
- 20 Minute Quiz
- Statistical NLP – Language Modeling
- N Grams

# Today's Lecture

---

- NGrams
- Bigram
- Smoothing and NGram
  - Add one smoothing
  - Witten-Bell Smoothing

# Simple N-Grams

---

- An **N-gram model** uses the previous N-1 words to predict the next one:
  - $P(w_n | w_{n-1})$
  - We'll be dealing with  $P(\langle \text{word} \rangle | \langle \text{some previous words} \rangle)$
- unigrams:  $P(\text{dog})$
- bigrams:  $P(\text{dog} | \text{big})$
- trigrams:  $P(\text{dog} | \text{the big})$
- quadrigrams:  $P(\text{dog} | \text{the big dopey})$



# Chain Rule

---

conditional probability: 
$$P(A | B) = \frac{P(A \wedge B)}{P(B)}$$

$$P(A \wedge B) = P(A | B)P(B)$$

and

$$P(A \wedge B) = P(B | A)P(A)$$

So: 
$$P(A \wedge B) = P(B | A)P(A)$$

“the dog”: 
$$P(The \wedge dog) = P(dog | the)P(the)$$

“the dog bites”:

$$P(The \wedge dog \wedge bites) = P(The)P(dog | The)P(bites | The \wedge dog)$$

# Chain Rule

---

the probability of a word sequence is the probability of a conjunctive event.

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2 | w_1)P(w_3 | w_1^2) \dots P(w_n | w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k | w_1^{k-1}) \end{aligned}$$

Unfortunately, that's really not helpful in general. Why?

# Markov Assumption

---

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

- $P(w_n)$  can be approximated using only  $N-1$  previous words of context
- This lets us collect statistics in practice
- Markov models are the class of probabilistic models that assume that we can predict the probability of some future unit without looking too far into the past
- Order of a Markov model: length of prior context

# Language Models and N-grams

---

- Given a word sequence:  $w_1 w_2 w_3 \dots w_n$
- **Chain rule**
  - $p(w_1 w_2) = p(w_1) p(w_2|w_1)$
  - $p(w_1 w_2 w_3) = p(w_1) p(w_2|w_1) p(w_3|w_1 w_2)$
  - ...
  - $p(w_1 w_2 w_3 \dots w_n) = p(w_1) p(w_2|w_1) p(w_3|w_1 w_2) \dots p(w_n|w_1 \dots w_{n-2} w_{n-1})$
- Note:
  - It's not easy to collect (meaningful) statistics on  $p(w_n|w_{n-1} w_{n-2} \dots w_1)$  for all possible word sequences
- **Bigram approximation**
  - *just look at the previous word only (not all the proceedings words)*
  - Markov Assumption: finite length history
  - 1st order Markov Model
  - $p(w_1 w_2 w_3 \dots w_n) = p(w_1) p(w_2|w_1) p(w_3|w_1 w_2) \dots p(w_n|w_1 \dots w_{n-3} w_{n-2} w_{n-1})$
  - $p(w_1 w_2 w_3 \dots w_n) \approx p(w_1) p(w_2|w_1) p(w_3|w_2) \dots p(w_n|w_{n-1})$
- Note:
  - $p(w_n|w_{n-1})$  is a lot easier to estimate well than  $p(w_n|w_1 \dots w_{n-2} w_{n-1})$

# Language Models and N-grams

---

- Given a word sequence:  $w_1 w_2 w_3 \dots w_n$
- **Chain rule**
  - $p(w_1 w_2) = p(w_1) p(w_2|w_1)$
  - $p(w_1 w_2 w_3) = p(w_1) p(w_2|w_1) p(w_3|w_1 w_2)$
  - ...
  - $p(w_1 w_2 w_3 \dots w_n) = p(w_1) p(w_2|w_1) p(w_3|w_1 w_2) \dots p(w_n|w_1 \dots w_{n-2} w_{n-1})$
- **Trigram approximation**
  - 2nd order Markov Model
  - *just look at the preceding two words only*
  - $p(w_1 w_2 w_3 w_4 \dots w_n) = p(w_1) p(w_2|w_1) p(w_3|w_1 w_2)$   
 $p(w_4|w_1 w_2 w_3) \dots p(w_n|w_1 \dots w_{n-3} w_{n-2} w_{n-1})$
  - $p(w_1 w_2 w_3 \dots w_n) \approx p(w_1) p(w_2|w_1) p(w_3|w_1 w_2) p(w_4|w_2 w_3) \dots p(w_n | w_{n-2} w_{n-1})$
- Note:
  - $p(w_n|w_{n-2} w_{n-1})$  is a lot easier to estimate well than  $p(w_n|w_1 \dots w_{n-2} w_{n-1})$  but harder than  $p(w_n|w_{n-1})$

# Corpora

---

- Corpora are (generally online) collections of text and speech
- e.g.
  - Brown Corpus (1M words)
  - Wall Street Journal and AP News corpora
  - ATIS, Broadcast News (speech)
  - TDT (text and speech)
  - Switchboard, Call Home (speech)
  - TRAINS, FM Radio (speech)

# Sample Word frequency (count) Data

(The Text REtrieval Conference) - (from B. Croft, UMass)

<b>Frequent Word</b>	<b>Number of Occurrences</b>	<b>Percentage of Total</b>
the	7,398,934	5.9
of	3,893,790	3.1
to	3,364,653	2.7
and	3,320,687	2.6
in	2,311,785	1.8
is	1,559,147	1.2
for	1,313,561	1.0
The	1,144,860	0.9
that	1,066,503	0.8
said	1,027,713	0.8

Frequencies from 336,310 documents in the 1GB TREC Volume 3 Corpus  
125,720,891 total word occurrences; 508,209 unique words

# Counting Words in Corpora

---

- Probabilities are based on counting things, so ....
- What should we count?
- Words, word classes, word senses, speech acts ...?
- What is a word?
  - e.g., are **cat** and **cats** the same word?
  - **September** and **Sept**?
  - **zero** and **oh**?
  - Is seventy-two one word or two? AT&T?
- Where do we find the things to count?



# Terminology

---

- Sentence: unit of written language
- Utterance: unit of spoken language
- Wordform: the inflected form that appears in the corpus
- Lemma: lexical forms having the same stem, part of speech, and word sense
- Types: number of distinct words in a corpus (vocabulary size)
- Tokens: total number of words

# Training and Testing

---

- Probabilities come from a **training corpus**, which is used to design the model.
  - narrow corpus: probabilities don't generalize
  - general corpus: probabilities don't reflect task or domain
- A separate **test corpus** is used to **evaluate** the model

# Simple N-Grams

---

- An **N-gram model** uses the previous N-1 words to predict the next one:
  - $P(w_n | w_{n-1})$
  - We'll be dealing with  $P(\langle \text{word} \rangle | \langle \text{some prefix} \rangle)$
- unigrams:  $P(\text{dog})$
- bigrams:  $P(\text{dog} | \text{big})$
- trigrams:  $P(\text{dog} | \text{the big})$
- quadrigrams:  $P(\text{dog} | \text{the big red})$

# Using N-Grams

---

- Recall that

- $P(w_n | w_{1..n-1}) \approx P(w_n | w_{n-N+1..n-1})$

- For a bigram grammar

- P(sentence) can be approximated by multiplying all the bigram probabilities in the sequence

- $P(\text{I want to eat Chinese food}) = P(\text{I} | \langle \text{start} \rangle) P(\text{want} | \text{I}) P(\text{to} | \text{want}) P(\text{eat} | \text{to}) P(\text{Chinese} | \text{eat}) P(\text{food} | \text{Chinese}) P(\langle \text{end} \rangle | \text{food})$

# Chain Rule

---

- Recall the definition of conditional probabilities

$$P(A | B) = \frac{P(A \wedge B)}{P(B)}$$

- Rewriting

$$P(A \wedge B) = P(A | B)P(B)$$

- Or...

$$P(\textit{The big}) = P(\textit{big} | \textit{the})P(\textit{the})$$

- Or...

$$P(\textit{The big}) = P(\textit{the})P(\textit{big} | \textit{the})$$

# Example

---

- The big red dog
- $P(\text{The}) * P(\text{big}|\text{the}) * P(\text{red}|\text{the big}) * P(\text{dog}|\text{the big red})$
- Better  $P(\text{The} | \text{<Beginning of sentence>})$   
written as  $P(\text{The} | \text{<S>})$
- Also  $\text{<end>}$  for end of sentence

# General Case

---

- The word sequence from position 1 to n is  $w_1^n$
- So the probability of a sequence is

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2 | w_1)P(w_3 | w_1^2) \dots P(w_n | w_1^{n-1}) \\ &= P(w_1) \prod_{k=2}^n P(w_k | w_1^{k-1}) \end{aligned}$$

# Unfortunately

---

- That doesn't help since its unlikely we'll ever gather the right statistics for the prefixes.



# Markov Assumption

---

- Assume that the entire prefix history isn't necessary.
- In other words, an event doesn't depend on all of its history, just a fixed length near history

# Markov Assumption

---

- So for each component in the product replace each with its approximation (assuming a prefix (Previous words) of N)

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

# N-Grams

## The big red dog

---

- Unigrams:  $P(\text{dog})$
- Bigrams:  $P(\text{dog}|\text{red})$
- Trigrams:  $P(\text{dog}|\text{big red})$
- Four-grams:  $P(\text{dog}|\text{the big red})$

In general, we'll be dealing with  
 $P(\text{Word} | \text{Some fixed prefix})$

Note: prefix is Previous words

---

N-gram models can be trained by  
counting and normalization

Bigram: 
$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

Ngram: 
$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}w_n)}{C(w_{n-N+1}^{n-1})}$$

# An example

---

- `<s> I am Sam <\s>`
- `<s> Sam I am <\s>`
- `<s> I do not like green eggs and meet <\s>`

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1} w_n)}{C(w_{n-N+1}^{n-1})}$$

$$P(I | \langle s \rangle) = \frac{2}{3} = 0.67$$

$$P(\text{Sam} | \langle s \rangle) = \frac{1}{3} = 0.33$$

$$P(\text{am} | I) = \frac{2}{3} = 0.67$$

$$P(\langle s \rangle | \text{Sam}) = \frac{1}{2} = 0.5$$

$$P(\langle s \rangle | \text{am}) = \frac{1}{2} = 0.5$$

$$P(\text{Sam} | \text{am}) = \frac{1}{2} = 0.5$$

$$P(\text{do} | I) = \frac{1}{1} = 1.0$$

# BERP Bigram Counts

BErkeley Restaurant Project (speech)

---

	I	Want	To	Eat	Chinese	Food	lunch
I	8	1087	0	13	0	0	0
Want	3	0	786	0	6	8	6
To	3	0	10	860	3	0	12
Eat	0	0	2	0	19	2	52
Chinese	2	0	0	0	0	120	1
Food	19	0	17	0	0	0	0
Lunch	4	0	0	0	0	1	0

# BERP Bigram Probabilities

---

- Normalization: divide each row's counts by appropriate unigram counts

I	Want	To	Eat	Chinese	Food	Lunch
3437	1215	3256	938	213	1506	459

- Computing the probability of **I I**
  - $C(I|I)/C(\text{all } I)$
  - $p = 8 / 3437 = .0023$
- A bigram grammar is an  $N \times N$  matrix of probabilities, where  $N$  is the vocabulary size

# A Bigram Grammar Fragment from BERP

---

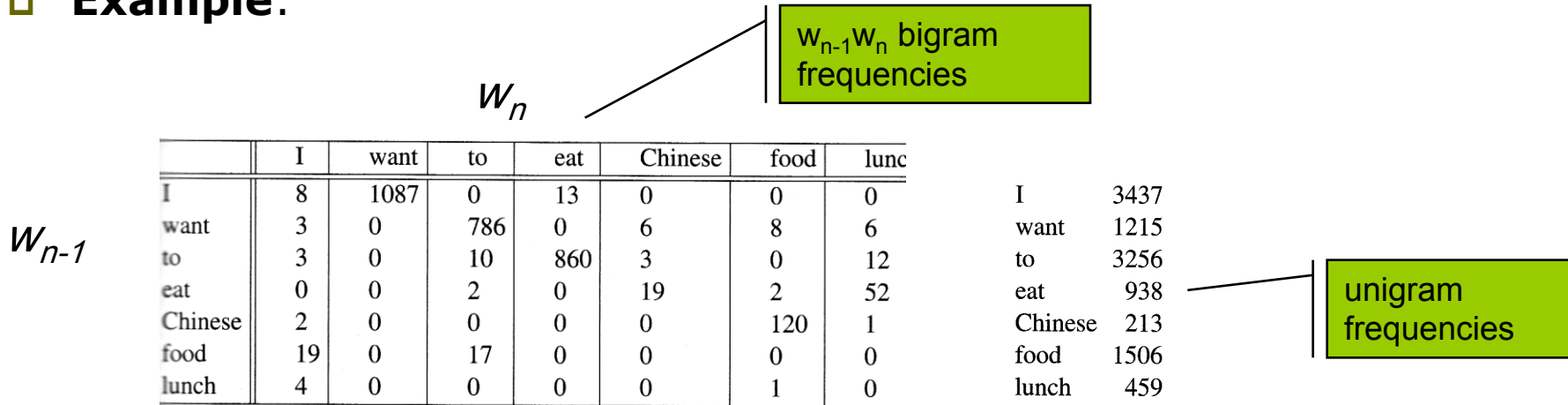
Eat on	.16	Eat Thai	.03
Eat some	.06	Eat breakfast	.03
Eat lunch	.06	Eat in	.02
Eat dinner	.05	Eat Chinese	.02
Eat at	.04	Eat Mexican	.02
Eat a	.04	Eat tomorrow	.01
Eat Indian	.04	Eat dessert	.007
Eat today	.03	Eat British	.001



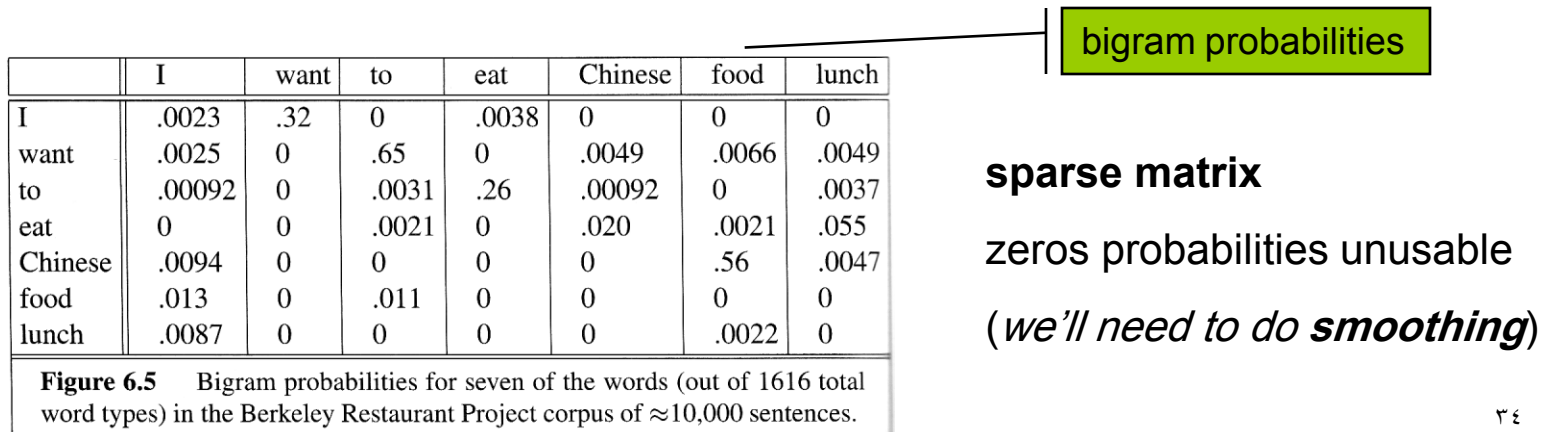
<start> I	.25	Want some	.04
<start> I'd	.06	Want Thai	.01
<start> Tell	.04	To eat	.26
<start> I'm	.02	To have	.14
I want	.32	To spend	.09
I would	.29	To be	.02
I don't	.08	British food	.60
I have	.04	British restaurant	.15
Want to	.65	British cuisine	.01
Want a	.05	British lunch	.01

# Language Models and N-grams

## □ Example:



**Figure 6.4** Bigram counts for seven of the words (out of 1616 total word types) in the Berkeley Restaurant Project corpus of  $\approx 10,000$  sentences.



# Example

---

- $P(\text{I want to eat British food}) =$   
 $P(\text{I} | \langle \text{start} \rangle) P(\text{want} | \text{I}) P(\text{to} | \text{want})$   
 $P(\text{eat} | \text{to}) P(\text{British} | \text{eat}) P(\text{food} | \text{British})$   
 $= .25 * .32 * .65 * .26 * .001 * .60 =$   
**0.0000081** *(different from textbook)*
- vs.  $\text{I want to eat Chinese food} = .00015$

# Note on Example

---

- Probabilities seem to capture “syntactic” facts, “world knowledge”
  - **eat** is often followed by a NP
  - British food is not too popular

# What do we learn about the language?

---

- What's being captured with ...
  - $P(\text{want} \mid I) = .32$
  - $P(\text{to} \mid \text{want}) = .65$
  - $P(\text{eat} \mid \text{to}) = .26$
  - $P(\text{food} \mid \text{Chinese}) = .56$
  - $P(\text{lunch} \mid \text{eat}) = .055$

# Some Observations

---

- $P(I \mid I)$
- $P(\text{want} \mid I)$
- $P(I \mid \text{food})$
- I I I want
- I want I want to
- The food I want is

---

□ What about

- $P(I \mid I) = .0023$  I I I I want
- $P(I \mid \text{want}) = .0025$  I want I want
- $P(I \mid \text{food}) = .013$  the kind of food I want is ...

# To avoid underflow use Logs

---

- ❑ You don't really do all those multiplies. The numbers are too small and lead to underflows
- ❑ Convert the probabilities to logs and then do additions.
- ❑ To get the real probability (if you need it) go back to the antilog.



# Generation

---

- Choose N-Grams according to their probabilities and string them together

# BERP

---

- I want  
    want to  
        to eat  
            eat Chinese  
                Chinese food  
                    food .

# Some Useful Observations

---

- A small number of events occur with high frequency
  - You can collect reliable statistics on these events with relatively small samples
- A large number of events occur with small frequency
  - You might have to wait a long time to gather statistics on the low frequency events

# Some Useful Observations

---

- Some zeroes are really zeroes
  - Meaning that they represent events that can't or shouldn't occur
- On the other hand, some zeroes aren't really zeroes
  - They represent low frequency events that simply didn't occur in the corpus

# Problem

---

- Let's assume we're using N-grams
- How can we assign a probability to a sequence where one of the component n-grams has a value of zero
- Assume all the words are known and have been seen
  - Go to a lower order n-gram
  - Back off from bigrams to unigrams
  - Replace the zero with something else

# Add-One

---

- Make the zero counts 1.
- Justification: They're just events you haven't seen yet. If you had seen them you would only have seen them once. so make the count equal to 1.

# Add-one: Example

unsmoothed bigram counts:  $\underbrace{\hspace{15em}}_{2^{nd} \text{ word}}$

}	<i>1<sup>st</sup> word</i>		<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total (N)</b>
	<i>I</i>	8	1087	0	13	0	0	0	0		3437
	<i>want</i>	3	0	786	0	6	8	6		1215	
	<i>to</i>	3	0	10	860	3	0	12		3256	
	<i>eat</i>	0	0	2	0	19	2	52		938	
	<i>Chinese</i>	2	0	0	0	0	120	1		213	
	<i>food</i>	19	0	17	0	0	0	0		1506	
	<i>lunch</i>	4	0	0	0	0	1	0		459	
...											

unsmoothed normalized bigram probabilities:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total</b>
<i>I</i>	.0023 (8/3437)	.32	0	.0038 (13/3437)	0	0	0		1
<i>want</i>	.0025	0	.65	0	.0049	.0066	.0049		1
<i>to</i>	.00092	0	.0031	.26	.00092	0	.0037		1
<i>eat</i>	0	0	.0021	0	.020	.0021	.055		1
<i>Chinese</i>	.0094	0	0	0	0	.56	.0047		1
<i>food</i>	.013	0	.011	0	0	0	0		1
<i>lunch</i>	.0087	0	0	0	0	.0022	0		1
...									

# Add-one: Example (con't)

add-one smoothed bigram counts:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total (N+V)</b>
<i>I</i>	8- 9	<del>1087</del> 1088	1	14	1	1	1		3437 5053
<i>want</i>	3 4	1	787	1	7	9	7		2831
<i>to</i>	4	1	11	861	4	1	13		4872
<i>eat</i>	1	1	23	1	20	3	53		2554
<i>Chinese</i>	3	1	1	1	1	121	2		1829
<i>food</i>	20	1	18	1	1	1	1		3122
<i>lunch</i>	5	1	1	1	1	2	1		2075

add-one normalized bigram probabilities:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total</b>
<i>I</i>	.0018 (9/5053)	.22	.0002	.0028 (14/5053)	.0002	.0002	.0002		1
<i>want</i>	.0014	.00035	.28	.00035	.0025	.0032	.0025		1
<i>to</i>	.00082	.00021	.0023	.18	.00082	.00021	.0027		1
<i>eat</i>	.00039	.00039	.0012	.00039	.0078	.0012	.021		1
<i>Chinese</i>	.0016	.00055	.00055	.00055	.00055	.066	.0011		1
<i>food</i>	.0064	.00032	.0058	.00032	.00032	.00032	.00032		1
<i>lunch</i>	.0024	.00048	.00048	.00048	.00048	.0022	.00048		1



# The example again

unsmoothed bigram counts:

$V = 1616$  word types

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total (N)</b>
<i>I</i>	8	1087	0	13	0	0	0		3437
<i>want</i>	3	0	786	0	6	8	6		1215
<i>to</i>	3	0	10	860	3	0	12		3256
<i>eat</i>	0	0	2	0	19	2	52		938
<i>Chinese</i>	2	0	0	0	0	120	1		213
<i>food</i>	19	0	17	0	0	0	0		1506
<i>lunch</i>	4	0	0	0	0	1	0		459

$V = 1616$

Smoothed  $P(I \text{ eat})$

$= (C(I \text{ eat}) + 1) / (\text{nb bigrams starting with "I"} + \text{nb of possible bigrams starting with "I"})$

$= (13 + 1) / (3437 + 1616)$

$= 0.0028$

# Smoothing and N-grams

I	3437
want	1215
to	3256
eat	938
Chinese	213
food	1506
lunch	459

## □ Add-One Smoothing

- add 1 to all frequency counts

## □ Bigram

- $p(w_n|w_{n-1}) = (C(w_{n-1}w_n)+1)/(C(w_{n-1})+V)$

- $(C(w_{n-1}w_n)+1) * C(w_{n-1}) / (C(w_{n-1})+V)$

## □ Frequencies

	I	want	to	eat	Chinese	food	lunch
I	8	1087	0	13	0	0	0
want	3	0	786	0	6	8	6
to	3	0	10	860	3	0	12
eat	0	0	2	0	19	2	52
Chinese	2	0	0	0	0	120	1
food	19	0	17	0	0	0	0
lunch	4	0	0	0	0	1	0

	I	want	to	eat	Chinese	food	lunch
I	6.12	740.05	0.68	9.52	0.68	0.68	0.68
want	1.72	0.43	337.76	0.43	3.00	3.86	3.00
to	2.67	0.67	7.35	575.41	2.67	0.67	8.69
eat	0.37	0.37	1.10	0.37	7.35	1.10	19.47
Chinese	0.35	0.12	0.12	0.12	0.12	14.09	0.23
food	9.65	0.48	8.68	0.48	0.48	0.48	0.48
lunch	1.11	0.22	0.22	0.22	0.22	0.44	0.22

### Remarks:

add-one causes large changes in some frequencies due to relative size of  $V(1616)$

*want to*: 786  $\Rightarrow$  338

$$= (786 + 1) * 1215 / (1215 + 1616)$$

$$(c_i+1) \frac{N}{N+V}$$

# Problem with add-one smoothing

- bigrams starting with *Chinese* are boosted by a factor of 8 ! (1829 / 213)

unsmoothed bigram counts:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total (N)</b>
1st word	<i>I</i>	8	1087	0	13	0	0	0	<b>3437</b>
	<i>want</i>	3	0	786	0	6	8	6	<b>1215</b>
	<i>to</i>	3	0	10	860	3	0	12	<b>3256</b>
	<i>eat</i>	0	0	2	0	19	2	52	<b>938</b>
	<i>Chinese</i>	2	0	0	0	0	120	1	<b><u>213</u></b>
	<i>food</i>	19	0	17	0	0	0	0	<b>1506</b>
	<i>lunch</i>	4	0	0	0	0	1	0	<b>459</b>

add-one smoothed bigram counts:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<b>Total (N+V)</b>
1st word	<i>I</i>	9	1088	1	14	1	1	1	<b>5053</b>
	<i>want</i>	4	1	787	1	7	9	7	<b>2831</b>
	<i>to</i>	4	1	11	861	4	1	13	<b>4872</b>
	<i>eat</i>	1	1	23	1	20	3	53	<b>2554</b>
	<i>Chinese</i>	3	1	1	1	1	121	2	<b><u>1829</u></b>
	<i>food</i>	20	1	18	1	1	1	1	<b>3122</b>
	<i>lunch</i>	5	1	1	1	1	2	1	<b>2075</b>

# Problem with add-one smoothing (con't)

- Data from the AP from (Church and Gale, 1991)
  - Corpus of 22,000,000 bigrams
  - Vocabulary of 273,266 words (i.e. 74,674,306,756 possible bigrams)
  - 74,671,100,000 bigrams were unseen
  - And each unseen bigram was given a frequency of 0.000295

<i>Freq. from training data</i>	$f_{MLE}$	$f_{empirical}$	$f_{add-one}$	<i>Add-one smoothed freq.</i>
	0	0.000027	<b>0.000295</b>	too high
<i>Freq. from held-out data</i>	1	0.448	0.000274	
	2	1.25	0.000411	too low
	3	2.24	0.000548	
	4	3.23	0.000685	
	5	4.21	0.000822	

- Total probability mass given to unseen bigrams =  
(74,671,100,000 × 0.000295) / 22,000,000 ~ **99.96** !!!!

# Smoothing and N-grams

---

## □ **Witten-Bell Smoothing**

- equate zero frequency items with frequency 1 items
- use frequency of things seen once to estimate frequency of things we haven't seen yet
- *smaller impact than Add-One*

## □ **Unigram**

- a zero frequency word (unigram) is “an event that hasn't happened yet”
- count the number of words (T) we've observed in the corpus (Number of types)
- $p(w) = T / (Z * (N + T))$ 
  - w is a word with zero frequency
  - Z = number of zero frequency words
  - N = size of corpus

# Distributing

---

- The amount to be distributed is

$$\frac{T}{N + T}$$

- The number of events with count zero

$$Z$$

- So distributing evenly gets us

$$\frac{1}{Z} \frac{T}{N + T}$$

# Smoothing and N-grams

---

## □ Bigram

- $p(w_n|w_{n-1}) = C(w_{n-1}w_n)/C(w_{n-1})$  (original)
- $p(w_n|w_{n-1}) = T(w_{n-1})/(Z(w_{n-1})*(T(w_{n-1})+N))$   
for zero bigrams (after Witten-Bell)
  - $T(w_{n-1})$  = number of bigrams beginning with  $w_{n-1}$
  - $Z(w_{n-1})$  = number of unseen bigrams beginning with  $w_{n-1}$
  - $Z(w_{n-1})$  = total number of possible bigrams beginning with  $w_{n-1}$  minus the ones we've seen
  - $Z(w_{n-1}) = V - T(w_{n-1})$
- $T(w_{n-1})/ Z(w_{n-1}) * C(w_{n-1})/(C(w_{n-1})+ T(w_{n-1}))$ 
  - estimated zero bigram frequency
- $p(w_n|w_{n-1}) = C(w_{n-1}w_n)/(C(w_{n-1})+T(w_{n-1}))$ 
  - for non-zero bigrams (after Witten-Bell)

# Smoothing and N-grams

## □ Witten-Bell Smoothing

- use frequency (count) of things seen once to estimate frequency (count) of things we haven't seen yet

## □ Bigram

- $T(w_{n-1}) / Z(w_{n-1}) * C(w_{n-1}) / (C(w_{n-1}) + T(w_{n-1}))$  estimated zero bigram frequency (count)
  - $T(w_{n-1})$  = number of bigrams beginning with  $w_{n-1}$
  - $Z(w_{n-1})$  = number of unseen bigrams beginning with  $w_{n-1}$

	I	want	to	eat	Chinese	food	lunch
I	8	1087	0	13	0	0	0
want	3	0	786	0	6	8	6
to	3	0	10	860	3	0	12
eat	0	0	2	0	19	2	52
Chinese	2	0	0	0	0	120	1
food	19	0	17	0	0	0	0
lunch	4	0	0	0	0	1	0

**Remark:**  
*smaller changes*

	I	want	to	eat	Chinese	food	lunch
I	7.785	1057.763	0.061	12.650	0.061	0.061	0.061
want	2.823	0.046	739.729	0.046	5.647	7.529	5.647
to	2.885	0.084	9.616	826.982	2.885	0.084	11.539
eat	0.073	0.073	1.766	0.073	16.782	1.766	45.928
Chinese	1.828	0.011	0.011	0.011	0.011	109.700	0.914
food	18.019	0.051	16.122	0.051	0.051	0.051	0.051
lunch	3.643	0.026	0.026	0.026	0.026	0.911	0.026



# Distributing Among the Zeros

---

- If a bigram “ $w_x w_i$ ” has a zero count

Number of bigram types starting with  $w_x$

$$P(w_i | w_x) = \frac{1}{Z(w_x)} \frac{T(w_x)}{N(w_x) + T(w_x)}$$

Number of bigrams starting with  $w_x$  that were not seen

Actual frequency (count) of bigrams beginning with  $w_x$

Thank you

---

السلام عليكم ورحمة الله