

Natural Language Processing
SoSe 2015



Language Modelling

Dr. Mariana Neves

April 20th, 2015

(based on the slides of Dr. Saeedeh Momtazi)

Outline

- Motivation
- Estimation
- Evaluation
- Smoothing

Outline

- Motivation
- Estimation
- Evaluation
- Smoothing

Language Modelling

- Finding the probability of a sentence or a sequence of words
 - $P(S) = P(w_1, w_2, w_3, \dots, w_n)$

... all of a sudden I notice three guys standing on the sidewalk ...

... on guys all I of notice sidewalk three a sudden standing the ...

Language Modelling

- Finding the probability of a sentence or a sequence of words
 - $P(S) = P(w_1, w_2, w_3, \dots, w_n)$

... all of a sudden I notice three guys standing on the sidewalk ...

... on guys all I of notice sidewalk three a sudden standing the ...

Language Modelling

- Applications
 - Word prediction
 - Speech recognition
 - Handwriting recognition
 - Machine translation
 - Spell checking

Applications

- Word prediction
 - „natural language..“
 - processing
 - management
 - orange
 - ...

Applications

- Speech recognition
 - „Computers can recognize speech.“
 - „Computers can wreck a nice peach.“



Communication tip:

Phonological ambiguities or Give peas a chance!

One of my favourite ways to have fun with communication are phonological ambiguities.

Phonological ambiguities are two or more words which sound the same and have different meanings.



Language can contain ambiguities - and more than one way to compose a set of sounds into words.

So listen to yourself: It is always good to notice a spoken sentence often contains many words which are (sometimes not)

intended to be heard.

English examples:

- there - their
- here - hear
- plane - plain
- Hamburger (Citizens of Hamburg) - hamburger (burger, food)
- sea - see
- Friday - fry day
- weekend - weak end
- ice cream - I scream.
- new direction - nude erection
- new day - nude, eh?
- I don't know! - I don't - no!
- but - butt
- Wait - Weight
- psychotherapist - psycho the rapist
- You're unconscious now... - Your unconscious now...
- Your students... - You're students...
- Two - too - to

German examples:

- Du hast Gewehre. (You have got guns.) - Du hasst Gewehre. (You hate guns.)
- Lehrer (teacher) - leerer (emptier)

http://worldsgreatestsmile.com/html/phonological_ambiguity.html

Applications

- Handwriting recognition



<https://play.google.com/store/apps/details?id=com.metamoji.mazecen>

Applications

- Handwriting recognition
 - „Take the money and run“, Woody Allen: „I have a gub.“ instead of „I have a gun.“



<https://www.youtube.com/watch?v=-UHOgkDbVqc>

Applications

- Machine translation
 - „The cat eats...”
 - „Die Katze frisst...”
 - „Die Katze isst...”

 - Chinese to English:
 - „He briefed **to** reporters on the **chief** contents of the statements”
 - „He briefed reporters on the **chief** contents of the statements”
 - „He briefed **to** reporters on the **main** contents of the statements”
 - „He briefed reporters on the **main** contents of the statements”

Applications

- Spell checking
 - „I want to **adver** this project“
 - „adverb“ (noun)
 - „advert“ (verb)
 - „They are leaving in about fifteen **minuets** to go to her house.“
 - „minutes“
 - „The design **an** construction of the system will take more than a year.“
 - „and“

Outline

- Motivation
- Estimation
- Evaluation
- Smoothing

Language Modeling

- Finding the probability of a sentence or a sequence of words
 - $P(S) = P(w_1, w_2, w_3, \dots, w_n)$
 - „Computers can recognize speech.“
 - $P(\text{Computer, can, recognize, speech})$

Conditional Probability

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

$$P(A, B) = P(A) \cdot P(B|A)$$

$$P(A, B, C, D) = P(A) \cdot P(B|A) \cdot P(C|A, B) \cdot P(D|A, B, C)$$

Conditional Probability

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

$P(A) = 0.266$ or 26.6%



$P(B) = 0.332$ or 33.2%



$P(A \cap B) = 0.033$ or 3.3%

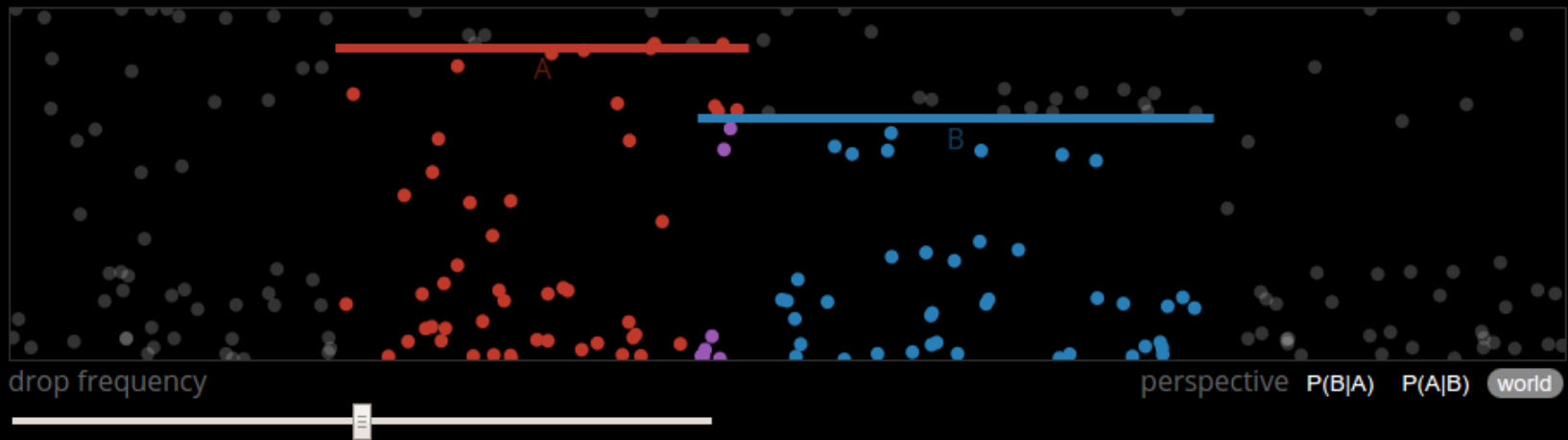


$P(B|A) = 0.124$ or 12.4%

If we have a ball and we know it hit the red shelf, there's a 12.4% chance it also hit the blue shelf.

$P(A|B) = 0.099$ or 9.9%

If we have a ball and we know it hit the blue shelf, there's a 9.9% chance it also hit the red shelf.



<http://setosa.io/conditional/>

Conditional Probability

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

$$P(S) = \prod_{i=1}^n P(w_i|w_1, w_2, \dots, w_{i-1})$$

$$P(\text{Computer, can, recognize, speech}) = P(\text{Computer}) \cdot P(\text{can}|\text{Computer}) \cdot P(\text{recognize}|\text{Computer can}) \cdot P(\text{speech}|\text{Computer can recognize})$$

Corpus

- Probabilities are based on counting things
- A computer-readable collection of text or speech
 - The Brown Corpus
 - A million-word collection of samples
 - 500 written texts from different genres (newspaper, fiction, non-fiction, academic, ...)
 - Assembled at Brown University in 1963-1964
- Can also be used for evaluation and comparison purposes

Corpus

1	<p>Confidence in the pound is widely expected to take another sharp dive if trade figures for September, due for release tomorrow, fail to show a substantial improvement from July and August 's near-record deficits.</p>
2	<p>Chancellor of the Exchequer Nigel Lawson 's restated commitment to a firm monetary policy has helped to prevent a freefall in sterling over the past week.</p>
3	<p>But analysts reckon underlying support for sterling has been eroded by the chancellor 's failure to announce any new policy measures in his Mansion House speech last Thursday.</p>
4	<p>This has increased the risk of the government being forced to increase base rates to 16% from their current 15% level to defend the pound, economists and foreign exchange market analysts say.</p>
5	<p>"The risks for sterling of a bad trade figure are very heavily on the down side," said Chris Dillow, senior U.K. economist at Nomura Research Institute.</p>
6	<p>"If there is another bad trade number, there could be an awful lot of pressure," noted Simon Briscoe, U.K. economist for Midland Montagu, a unit of Midland Bank PLC.</p>

<http://weaver.nlplab.org/~brat/demo/latest/#/not-editable/CoNLL-00-Chunking/train.txt-doc-1>

Corpus

- Text Corpora
 - Corpus of Contemporary American English
 - The British National Corpus
 - The International Corpus of English
 - The Google N-gram Corpus (<https://books.google.com/ngrams>)
 - WBI repository (biomedical domain) (<http://corpora.informatik.hu-berlin.de/>)

Word occurrence

- A language consists of a set of „V“ words (Vocabulary)
- A text is a sequence of the words from the vocabulary
- A word can occur several times in a text
 - Word Token: each occurrence of words in text
 - Word Type: each unique occurrence of words in the text

Word occurrence

- Example:
 - „This is a sample text from a book that is read every day.“

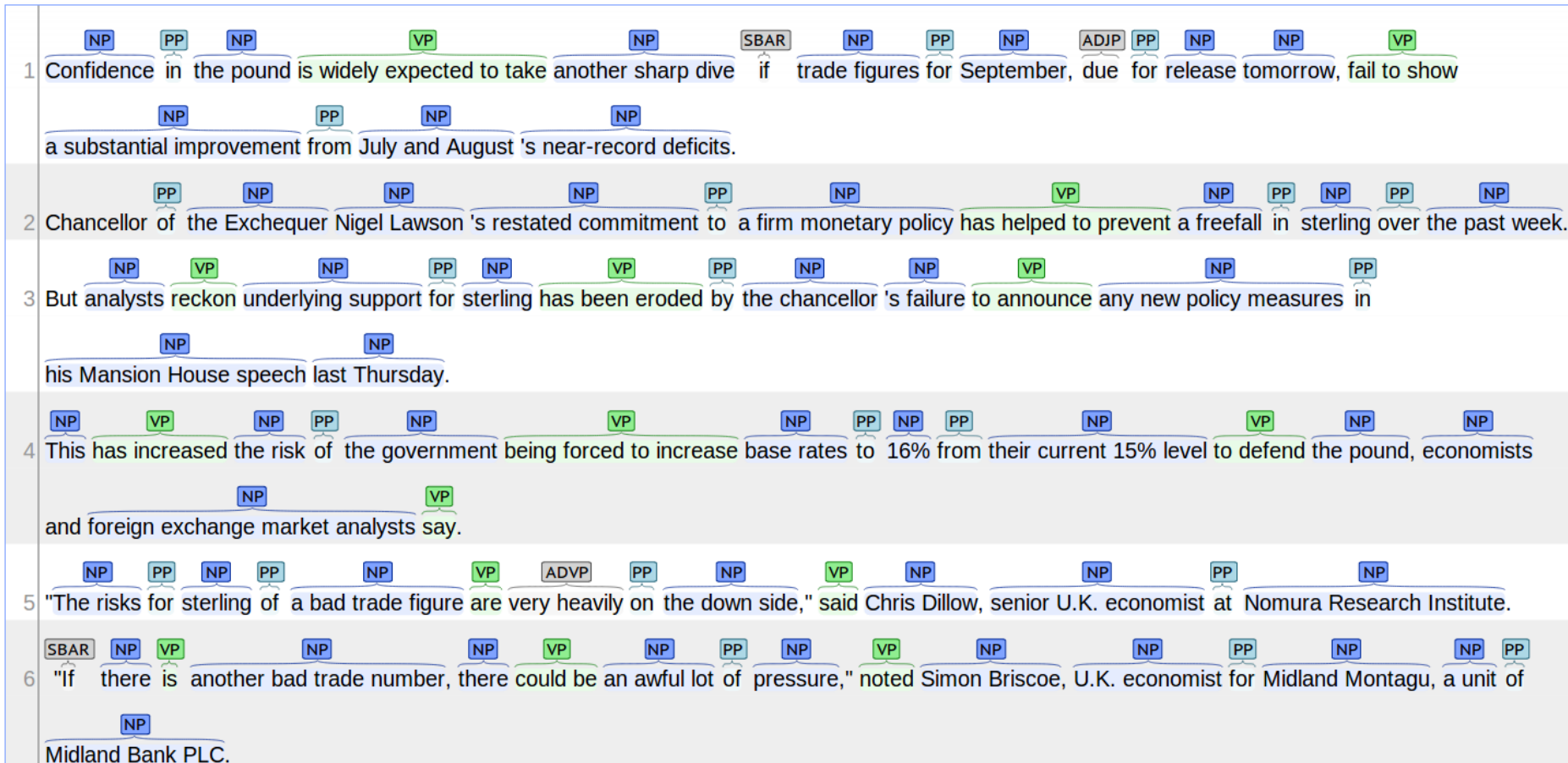
Word occurrence

- Example:
 - „This is a sample text from a book that is read every day.“
 - # Word Tokens: 13
 - # Word Types: 11

Counting

- The Brown corpus
 - 1,015,945 word tokens
 - 47,218 word types
- Google N-Gram corpus
 - 1,024,908,267,229 word tokens
 - 13,588,391 word types
- Why so many word types?
 - Large English dictionaries have around 500k word types

Corpora include numbers, misspellings, names, acronyms, etc.



<http://weaver.nlpplab.org/~brat/demo/latest/#/not-editable/CoNLL-00-Chunking/train.txt-doc-1>

Word frequency

Rank	Word	Count	Freq(%)
1	The	69970	6.8872
2	of	36410	3.5839
3	and	28854	2.8401
4	to	26154	2.5744
5	a	23363	2.2996
6	in	21345	2.1010
7	that	10594	1.0428
8	is	10102	0.9943
9	was	9815	0.9661
10	He	9542	0.9392
11	for	9489	0.9340
12	it	8760	0.8623
13	with	7290	0.7176
14	as	7251	0.7137
15	his	6996	0.6886
16	on	6742	0.6636
17	be	6376	0.6276
18	at	5377	0.5293
19	by	5307	0.5224
20	I	5180	0.5099

Zipf's Law

- The frequency of any word is inversely proportional to its rank in the frequency table
- Given a corpus of natural language utterances, the most frequent word will occur approximately
 - twice as often as the second most frequent word,
 - three times as often as the third most frequent word,
 - ...
- Rank of a word times its frequency is approximately a constant
 - $\text{Rank} \cdot \text{Freq} \approx c$
 - $c \approx 0.1$ for English

Word frequency

Rank	Word	Count	Freq(%)	Freq x Rank
1	The	69970	6.8872	0.06887
2	of	36410	3.5839	0.07167
3	and	28854	2.8401	0.08520
4	to	26154	2.5744	0.10297
5	a	23363	2.2996	0.11498
6	in	21345	2.1010	0.12606
7	that	10594	1.0428	0.07299
8	is	10102	0.9943	0.07954
9	was	9815	0.9661	0.08694
10	He	9542	0.9392	0.09392
11	for	9489	0.9340	0.10274
12	it	8760	0.8623	0.10347
13	with	7290	0.7176	0.09328
14	as	7251	0.7137	0.09991
15	his	6996	0.6886	0.10329
16	on	6742	0.6636	0.10617
17	be	6376	0.6276	0.10669
18	at	5377	0.5293	0.09527
19	by	5307	0.5224	0.09925
20	I	5180	0.5099	0.10198

Freq · Rank ≈ c

Word frequency

- Zipf's Law is not very accurate for very frequent and very infrequent words

Rank	Word	Count	Freq(%)	Freq x Rank
1	The	69970	6.8872	0.06887
2	of	36410	3.5839	0.07167
3	and	28854	2.8401	0.08520
4	to	26154	2.5744	0.10297
5	a	23363	2.2996	0.11498

Word frequency

- Zipf's Law is not very accurate for very frequent and very infrequent words

Rank	Word	Count	Freq(%)	Freq x Rank
1000	current	104	0.0102	0.10200
1001	spent	104	0.0102	0.10210
1002	eight	104	0.0102	0.10220
1003	covered	104	0.0102	0.10230
1004	Negro	104	0.0102	0.10240
1005	role	104	0.0102	0.10251
1006	played	104	0.0102	0.10261
1007	l'd	104	0.0102	0.10271
1008	date	103	0.0101	0.10180
1009	council	103	0.0101	0.10190
1010	race	103	0.0101	0.10201

Maximum Likelihood Estimation

- $P(\text{speech}|\text{Computer can recognize})$

$$P(\text{speech}|\text{Computer can recognize}) = \frac{\#(\text{Computer can recognize speech})}{\#(\text{Computer can recognize})}$$

- Too many phrases
- Limited text for estimating probabilities
- Simplification assumption

Markov assumption

$$P(S) = \prod_{i=1}^n P(w_i | w_1, w_2, \dots, w_{i-1})$$



$$P(S) = \prod_{i=1}^n P(w_i | w_{i-1})$$

Markov assumption

$$P(\text{Computer, can, recognize, speech}) = P(\text{Computer}) \cdot P(\text{can}|\text{Computer}) \cdot P(\text{recognize}|\text{Computer can}) \cdot P(\text{speech}|\text{Computer can recognize})$$



$$P(\text{Computer, can, recognize, speech}) = P(\text{Computer}) \cdot P(\text{can}|\text{Computer}) \cdot P(\text{recognize}|\text{can}) \cdot P(\text{speech}|\text{recognize})$$

$$P(\text{speech}|\text{recognize}) = \frac{\#(\text{recognize speech})}{\#(\text{recognize})}$$

N-gram model

- Unigram $P(S) = \prod_{i=1}^n P(w_i)$
- Bigram $P(S) = \prod_{i=1}^n P(w_i | w_{i-1})$
- Trigram $P(S) = \prod_{i=1}^n P(w_i | w_{i-1}, w_{i-2})$
- N-gram $P(S) = \prod_{i=1}^n P(w_i | w_1, w_2, \dots, w_{i-1})$

N-grams

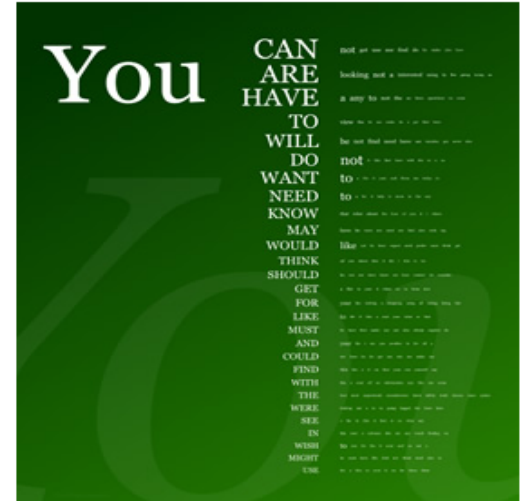
1. (*unigram*) Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives
2. (*bigram*) Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her
3. (*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

Google N-gram corpus



(<https://books.google.com/ngrams>)

WebTrigrams



(<http://www.chrisharrison.net/index.php/Visualizations/WebTrigrams>)

Maximum Likelihood Estimation

- `<s> I saw the boy </s>`
- `<s> the man is working </s>`
- `<s> I walked in the street </s>`

- Vocabulary:
 - $V = \{I, \text{saw}, \text{the}, \text{boy}, \text{man}, \text{is}, \text{working}, \text{walked}, \text{in}, \text{street}\}$

 - walked boy working
 - The boy is working
 - street saw the man

Maximum Likelihood

- <s> I saw the boy </s>
- <s> the man is working </s>
- <s> I walked in the street </s>

boy	I	in	is	man	saw	street	the	walked	working
1	2	1	1	1	1	1	3	1	1

	boy	I	in	is	man	saw	street	the	walked	working
boy	0	0	0	0	0	0	0	0	0	0
I	0	0	0	0	0	1	0	0	1	0
in	0	0	0	0	0	0	0	1	0	0
is	0	0	0	0	0	0	0	0	0	1
man	0	0	0	1	0	0	0	0	0	0
saw	0	0	0	0	0	0	0	1	0	0
street	0	0	0	0	0	0	0	0	0	0
the	1	0	0	0	1	0	1	0	0	0
walked	0	0	1	0	0	0	0	0	0	0
working	0	0	0	0	0	0	0	0	0	0

Maximum Likelihood

- $\langle s \rangle$ I saw the man $\langle /s \rangle$

boy	I	in	is	man	saw	street	the	walked	working	$\langle s \rangle$
1	2	1	1	1	1	1	3	1	1	3

	boy	I	in	is	man	saw	street	the	walked	working
boy	0	0	0	0	0	0	0	0	0	0
I	0	0	0	0	0	1	0	0	1	0
in	0	0	0	0	0	0	0	1	0	0
is	0	0	0	0	0	0	0	0	0	1
man	0	0	0	1	0	0	0	0	0	0
saw	0	0	0	0	0	0	0	1	0	0
street	0	0	0	0	0	0	0	0	0	0
the	1	0	0	0	1	0	1	0	0	0
walked	0	0	1	0	0	0	0	0	0	0
working	0	0	0	0	0	0	0	0	0	0

$\langle s \rangle$

2

1

$$P(S) = P(I|\langle s \rangle) \cdot P(\text{saw}|I) \cdot P(\text{the}|\text{saw}) \cdot P(\text{man}|\text{the})$$

$$P(S) = \frac{\#(\langle s \rangle I)}{\#(\langle s \rangle)} \cdot \frac{\#(I \text{ saw})}{\#(I)} \cdot \frac{\#(\text{saw the})}{\#(\text{saw})} \cdot \frac{\#(\text{the man})}{\#(\text{the})} \quad P(S) = \frac{2}{3} \cdot \frac{1}{2} \cdot \frac{1}{1} \cdot \frac{1}{3}$$

Unkown words

- `<s>` I saw the woman `</s>`
- Closed vocabulary: test set can only contain words from this lexicon
- Open vocabulary: test set can contain unknown words
- Out of vocabulary (OOV) words:
 - Choose a vocabulary
 - Convert unknown (OOV) words to `<UNK>` word token
 - Estimate probabilities for `<UNK>`
- Alternatively,
 - Replace the first occurrence of every word type in the training data by `<UNK>`

Outline

- Motivation
- Estimation
- **Evaluation**
- Smoothing

Branching Factor

- Branching factor is the number of possible words that can be used in each position of a text
 - Maximum branching factor for each language is V
 - A good language model should be able to
 - minimize this number
 - give a higher probability to the words that occur in real texts

Branching Factor

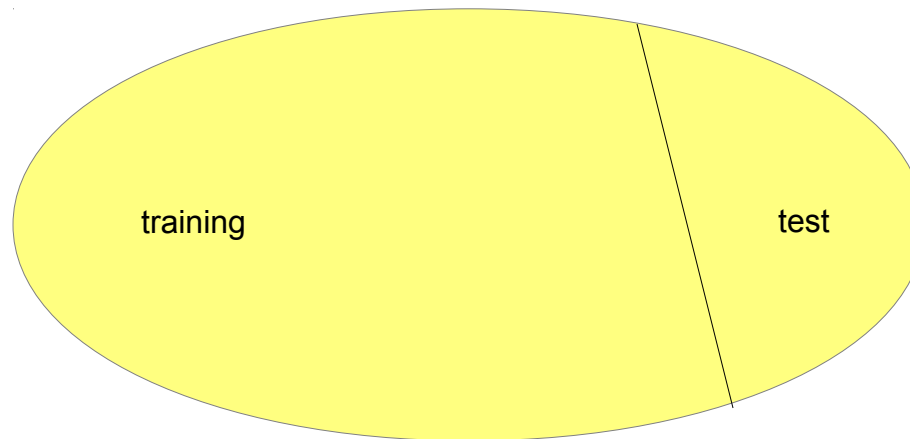
- John eats an ...
 - computer, book, apple, banana, umbrella, orange, desk

Branching Factor

- John eats an ...
 - computer, book, apple, banana, umbrella, orange, desk

Evaluation

- Dividing the corpus to two parts
- Building a language model from the training set
 - Word frequencies, etc..
- Estimating the probability of the test set
- Calculate the average branching factor of the test set



Perplexity

- Goal: giving higher probability to frequent texts
 - minimizing the perplexity of the frequent texts

$$P(S) = P(w_1, w_2, \dots, w_n)$$

$$\text{Perplexity}(S) = P(w_1, w_2, \dots, w_n)^{-\frac{1}{n}} = \sqrt[n]{\frac{1}{P(w_1, w_2, \dots, w_n)}}$$

$$\text{Perplexity}(S) = \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_1, w_2, \dots, w_{i-1})}}$$

Perplexity

- Maximum branching factor for each language is $|V|$

$$\text{Perplexity}(S) = \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_1, w_2, \dots, w_{i-1})}}$$

- Example: predicting next characters instead of next words:
 - $|V| = 26$, five next characters:

$$\text{Perplexity}(S) = \left(\left(\frac{1}{26} \right)^5 \right)^{-\frac{1}{5}} = 26$$

Perplexity

- Wall Street Journal (19,979 word vocabulary)
 - Training set: 38 million word tokens
 - Test set: 1.5 million words
- Perplexity:
 - Unigram: 962
 - Bigram: 170
 - Trigram: 109

Outline

- Motivation
- Estimation
- Evaluation
- **Smoothing**

Maximum Likelihood

- `<s> I saw the boy </s>`
- `<s> the man is working </s>`
- `<s> I walked in the street </s>`
- `<s> I saw the man </s>`

$$P(S) = P(I) \cdot P(\text{saw}|I) \cdot P(\text{the}|\text{saw}) \cdot P(\text{man}|\text{the})$$

$$P(S) = \frac{\#(I)}{\#(<s>)} \cdot \frac{\#(I \text{ saw})}{\#(I)} \cdot \frac{\#(\text{saw the})}{\#(\text{saw})} \cdot \frac{\#(\text{the man})}{\#(\text{the})}$$

$$P(S) = \frac{2}{3} \cdot \frac{1}{2} \cdot \frac{1}{1} \cdot \frac{1}{3}$$

Zero probability

- $\langle s \rangle$ I saw the man in the street $\langle /s \rangle$

boy	I	in	is	man	saw	street	the	walked	working
1	2	1	1	1	1	1	3	1	1

	boy	I	in	is	man	saw	street	the	walked	working
boy	0	0	0	0	0	0	0	0	0	0
I	0	0	0	0	0	1	0	0	1	0
in	0	0	0	0	0	0	0	1	0	0
is	0	0	0	0	0	0	0	0	0	1
man	0	0	0	1	0	0	0	0	0	0
saw	0	0	0	0	0	0	0	1	0	0
street	0	0	0	0	0	0	0	0	0	0
the	1	0	0	0	1	0	1	0	0	0
walked	0	0	1	0	0	0	0	0	0	0
working	0	0	0	0	0	0	0	0	0	0

$$P(S) = P(I) \cdot P(\text{saw}|I) \cdot P(\text{the}|\text{saw}) \cdot P(\text{man}|\text{the}) \cdot P(\text{in}|\text{man}) \cdot P(\text{the}|\text{in}) \cdot P(\text{street}|\text{the})$$

$$P(S) = \frac{\#(I)}{\#(\langle s \rangle)} \cdot \frac{\#(I \text{ saw})}{\#(I)} \cdot \frac{\#(\text{saw the})}{\#(\text{saw})} \cdot \frac{\#(\text{the man})}{\#(\text{the})} \cdot \frac{\#(\text{man in})}{\#(\text{man})} \cdot \frac{\#(\text{in the})}{\#(\text{in})} \cdot \frac{\#(\text{the street})}{\#(\text{the})}$$

$$P(S) = \frac{2}{3} \cdot \frac{1}{2} \cdot \frac{1}{1} \cdot \frac{1}{3} \cdot \frac{0}{1} \cdot \frac{1}{1} \cdot \frac{1}{3}$$

Zero probability

- `<s> I saw the boy </s>`
- `<s> the man is working </s>`
- `<s> I walked in the street </s>`

- No „man in“ in our corpus

Smoothing

- Giving a small probability to all as unseen n-grams
 - Laplace smoothing
 - Add one to all counts (Add-one)

	boy	I	in	is	man	saw	street	the	walked	working
boy	0	0	0	0	0	0	0	0	0	0
I	0	0	0	0	0	1	0	0	1	0
in	0	0	0	0	0	0	0	1	0	0
is	0	0	0	0	0	0	0	0	0	1
man	0	0	0	1	0	0	0	0	0	0
saw	0	0	0	0	0	0	0	1	0	0
street	0	0	0	0	0	0	0	0	0	0
the	1	0	0	0	1	0	1	0	0	0
walked	0	0	1	0	0	0	0	0	0	0
working	0	0	0	0	0	0	0	0	0	0

Smoothing

- Giving a small probability to all unseen n-grams
 - Laplace smoothing
 - Add one to all counts (Add-one)

	boy	I	in	is	man	saw	street	the	walked	working
boy	1	1	1	1	1	1	1	1	1	1
I	1	1	1	1	1	2	1	1	2	1
in	1	1	1	1	1	1	1	2	1	1
is	1	1	1	1	1	1	1	1	1	2
man	1	1	1	2	1	1	1	1	1	1
saw	1	1	1	1	1	1	1	2	1	1
street	1	1	1	1	1	1	1	1	1	1
the	2	1	1	1	2	1	2	1	1	1
walked	1	1	2	1	1	1	1	1	1	1
working	1	1	1	1	1	1	1	1	1	1

$$P(w_i|w_{i-1}) = \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} \longrightarrow P(w_i|w_{i-1}) = \frac{\#(w_{i-1}, w_i) + 1}{\#(w_{i-1}) + V}$$

Smoothing

- Giving a small probability to all unseen n-grams
 - Interpolation and Back-off Smoothing
 - Use a background probability

$$P(w_i|w_{i-1}) = \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})}$$

Back-off

$$P(w_i|w_{i-1}) = \begin{cases} \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} & \text{if } \#(w_{i-1}, w_i) > 0 \\ P_{BG} & \text{otherwise} \end{cases}$$

Smoothing

- Giving a small probability to all as unseen n-grams
 - Interpolation and Back-off Smoothing
 - Use a background probability

$$P(w_i|w_{i-1}) = \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})}$$

Interpolation

$$P(w_i|w_{i-1}) = \lambda_1 \cdot \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} + \lambda_2 \cdot P_{BG}$$

$$\sum \lambda = 1$$

Parameter
tuning

Background
probability

Background probability

- Lower levels of n-gram can be used as background probability
 - Trigram » Bigram
 - Bigram » Unigram
 - Unigram » Zerogram $(\frac{1}{V})$

Back-off

$$P(w_i|w_{i-1}) = \begin{cases} \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} & \text{if } \#(w_{i-1}, w_i) > 0 \\ \alpha(w_i)P(w_i) & \text{otherwise} \end{cases}$$

$$P(w_i) = \begin{cases} \frac{\#(w_i)}{N} & \text{if } \#(w_i) > 0 \\ \alpha(w_i)\frac{1}{V} & \text{otherwise} \end{cases}$$

Background probability

- Lower levels of n-gram can be used as background probability
 - Trigram » Bigram
 - Bigram » Unigram
 - Unigram » Zerogram ($\frac{1}{V}$)

Interpolation

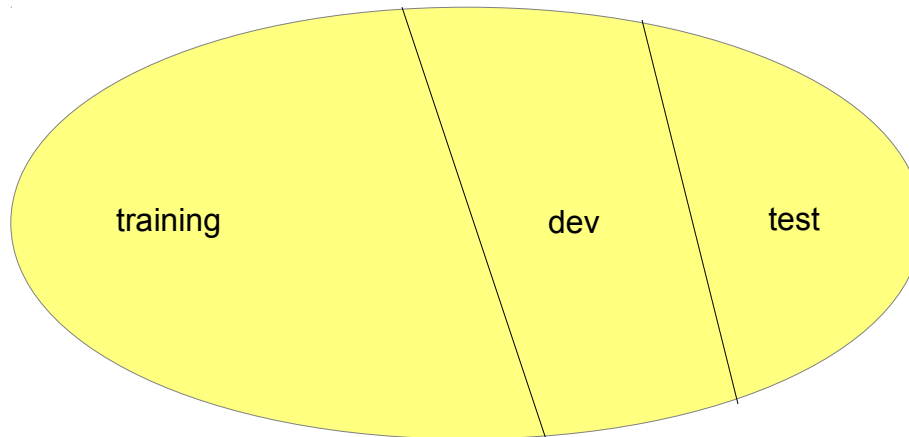
$$P(w_i|w_{i-1}) = \lambda_1 \cdot \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} + \lambda_2 \cdot P(w_i)$$

$$P(w_i) = \lambda_1 \cdot \frac{\#(w_i)}{N} + \lambda_2 \cdot \frac{1}{V}$$

$$P(w_i|w_{i-1}) = \lambda_1 \cdot \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} + \lambda_2 \cdot \frac{\#(w_i)}{N} + \lambda_3 \cdot \frac{1}{V}$$

Parameter tuning

- Held-out dataset (development set)
- 80% (training), 10% (dev-set), 10% (test)
- Minimize the perplexity of the held-out dataset



Advanced smoothing

$$P(w_i | w_{i-1}) = \frac{\#(w_{i-1}, w_i) + 1}{\#(w_{i-1}) + V}$$

$$P(w_i | w_{i-1}) = \frac{\#(w_{i-1}, w_i) + k}{\#(w_{i-1}) + kV} \quad (\text{add-k, add-}\delta \text{ smoothing})$$

Advanced smoothing

- Absolute discounting
 - Good estimates for high counts, discount won't affect them much
 - Lower counts are anyway not trustworthy

$$P(w_i|w_{i-1}) = \begin{cases} \frac{\#(w_{i-1}, w_i)}{\#(w_{i-1})} & \text{if } \#(w_{i-1}, w_i) > 0 \\ P_{BG} & \text{otherwise} \end{cases}$$

$$P(w_i|w_{i-1}) = \begin{cases} \frac{\#(w_{i-1}, w_i) - \delta}{\#(w_{i-1})} & \text{if } \#(w_{i-1}, w_i) > 0 \\ \alpha(w_i) \cdot P_{BG}(w_i) & \text{otherwise} \end{cases}$$

Advanced smoothing

- Estimation based on the lower-order n-gram
 - I cannot see without my reading ...
 - Going unigram : „Francisco“, „glasses“, ...
- Observations
 - „Francisco“ is more common than „glasses“
 - But „Francisco“ always follows „San“
 - „Francisco“ is not a novel continuation for a text

Advanced smoothing

- Solution
 - Instead of $P(w)$: How likely is „w“ to appear in a text?
 - $P_{\text{continuation}}(w)$: How likely is „w“ to appear as a novel continuation?
 - Count the number of words types after which „w“ appears

$$P_{\text{continuation}}(w) \propto |w_{i-1} : \#(w_{i-1}, w_i) > 0|$$

Advanced smoothing

- How many times does „w“ appear as a novel continuation

$$P_{continuation}(w) \propto |w_{i-1} : \#(w_{i-1}, w_i) > 0|$$

- Normalized by the total number of bigram types

$$P_{continuation}(w) = \frac{|w_{i-1} : \#(w_{i-1}, w_i) > 0|}{|(w_{j-1}, w_j) : \#(w_{j-1}, w_j) > 0|}$$

- Alternatively: normalized by the number of words preceding all words

$$P_{continuation}(w) = \frac{|w_{i-1} : \#(w_{i-1}, w_i) > 0|}{\sum_{w'} |(w'_{i-1}) : \#(w'_{i-1}, w'_i) > 0|}$$

Advanced smoothing

- Kneser-Ney discounting

$$P(w_i|w_{i-1}) = \frac{\max(\#(w_{i-1}, w_i) - \delta, 0)}{\#(w_{i-1})} + \alpha \cdot P_{BG}$$

$$P(w_i|w_{i-1}) = \frac{\max(\#(w_{i-1}, w_i) - \delta, 0)}{\#(w_{i-1})} + \alpha \cdot P_{continuation}$$

Class-based n-grams

- Compute estimation for the bigram „to Shanghai“
- Training data: „to London“, „to Beijing“, „to Denver“
- Classes: CITY_NAME, AIRLINE, DAY_OF_WEEK, MONTH, etc.

$$P(w_i|w_{i-1}) \approx P(c_i|c_{i-1}) \times P(w_i|c_{i-1})$$

Further reading

- Chapter 4
- <http://www.cs.columbia.edu/~mcollins/lm-spring2013.pdf>

