

# Natural Language Processing

Lecture 7—9/15/2015

Jim Martin

---

---

---

---

---

---

---

---

## Today

### More language modeling

- Probabilistic model
  - Independence assumptions
- Smoothing
  - Laplace smoothing (add-1)
  - Kneser-Ney

9/15/15

Speech and Language Processing - Jurafsky and Martin

2

---

---

---

---

---

---

---

---

## Markov Assumption

So for each component in the product replace with the approximation (assuming a prefix of N - 1)

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

Bigram version

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

9/15/15

Speech and Language Processing - Jurafsky and Martin

3

---

---

---

---

---

---

---

---

## Estimating Bigram Probabilities

- The Maximum Likelihood Estimate (MLE)

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

9/15/15

Speech and Language Processing - Jurafsky and Martin

4

---

---

---

---

---

---

---

---

## An Example

- <s> I am Sam </s>
- <s> Sam I am </s>
- <s> I do not like green eggs and ham </s>

$$P(I | \langle s \rangle) = \frac{2}{3} = .67 \quad P(\text{Sam} | \langle s \rangle) = \frac{1}{3} = .33 \quad P(\text{am} | I) = \frac{2}{3} = .67$$

$$P(\langle s \rangle | \text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam} | \text{am}) = \frac{1}{2} = .5 \quad P(\text{do} | I) = \frac{1}{3} = .33$$

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

9/15/15

Speech and Language Processing - Jurafsky and Martin

5

---

---

---

---

---

---

---

---

## Bigram Counts

- Vocabulary size is 1446 |V|
- Out of 9222 sentences
  - Eg. "I want" occurred 827 times

	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

9/15/15

Speech and Language Processing - Jurafsky and Martin

6

---

---

---

---

---

---

---

---

## Bigram Probabilities

- Divide bigram counts by prefix unigram counts to get bigram probabilities.

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

	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

9/15/15

Speech and Language Processing - Jurafsky and Martin

7

---

---

---

---

---

---

---

---

---

---

---

---

## Bigram Estimates of Sentence Probabilities

- $P(\langle s \rangle \text{ I want english food } \langle /s \rangle) =$   
 $P(i | \langle s \rangle) \cdot$   
 $P(\text{want} | i) \cdot$   
 $P(\text{english} | \text{want}) \cdot$   
 $P(\text{food} | \text{english}) \cdot$   
 $P(\langle /s \rangle | \text{food})$   
 $= .000031$

9/15/15

Speech and Language Processing - Jurafsky and Martin

8

---

---

---

---

---

---

---

---

---

---

---

---

## Perplexity

- Perplexity is the probability of a test set (assigned by the language model), as normalized by the number of words:
 
$$PP(W) = P(w_1 w_2 \dots w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}$$
- Chain rule:
 
$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$
- For bigrams:
 
$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$$
- Minimizing perplexity is the same as maximizing probability
  - The best language model is one that best predicts an unseen test set

9/15/15

Speech and Language Processing - Jurafsky and Martin

9

---

---

---

---

---

---

---

---

---

---

---

---

## Lower perplexity means a better model

- Training 38 million words, test 1.5 million words, WSJ

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

9/15/15

Speech and Language Processing - Jurafsky and Martin

10

---

---

---

---

---

---

---

---

## Practical Issues

- Once we start looking at test data, we'll run into words that we haven't seen before.
- Standard solution
  - Given a corpus
    - Create an unknown word token <UNK>
      - Create a fixed lexicon L, of size V
        - From a dictionary or
        - A subset of terms from the training set
      - At text normalization phase, any training word not in L is changed to <UNK>
      - Collect counts for that as for any normal word
    - At test time
      - Use UNK counts for any word not seen in training

9/15/15

Speech and Language Processing - Jurafsky and Martin

11

---

---

---

---

---

---

---

---

## Practical Issues

- Multiplying a bunch of small numbers between 0 and 1 is a really bad idea
  - Multiplication is slow
  - And underflow is likely
- So do everything in log space
  - Avoid underflow
  - Adding is faster than multiplying

$$p_1 \times p_2 \times p_3 \times p_4 = \exp(\log p_1 + \log p_2 + \log p_3 + \log p_4)$$

9/15/15

Speech and Language Processing - Jurafsky and Martin

12

---

---

---

---

---

---

---

---

## Smoothing

- **Back to Shakespeare**
  - Recall that Shakespeare produced 300,000 bigram types out of  $V^2 = 844$  million possible bigrams...
  - So, 99.96% of the possible bigrams were never seen (have zero entries in the table)
  - Does that mean that any sentence that contains one of those bigrams should have a probability of 0?
  - For generation (shannon game) it means we'll never emit those bigrams
  - But for analysis it's problematic because if we run across a new bigram in the future then we have no choice but to assign it a probability of zero.

9/15/15

Speech and Language Processing - Jurafsky and Martin

13

---

---

---

---

---

---

---

---

## Zero Counts

- Some of those zeros are really zeros...
  - Things that really aren't ever going to happen
- **On the other hand, some of them are just rare events.**
  - If the training corpus had been a little bigger they would have had a count
    - What would that count be in all likelihood?
- Zipf's Law (long tail phenomenon):
  - A small number of events occur with high frequency
  - A large number of events occur with low frequency
  - You can quickly collect statistics on the high frequency events
  - You might have to wait an arbitrarily long time to get valid statistics on low frequency events
- **Result:**
  - Our estimates are sparse! We have no counts at all for the vast number of things we want to estimate!
- **Answer:**
  - Estimate the likelihood of unseen (zero count) N-grams!

9/15/15

Speech and Language Processing - Jurafsky and Martin

14

---

---

---

---

---

---

---

---

## Laplace Smoothing

- Also called Add-One smoothing
- Just add one to all the counts!
- Very simple



- MLE estimate:

$$P(w_i) = \frac{c_i}{N}$$

- Laplace estimate:

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$$

9/15/15

Speech and Language Processing - Jurafsky and Martin

15

---

---

---

---

---

---

---

---

## Bigram Counts

	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

9/15/15

Speech and Language Processing - Jurafsky and Martin

16

---

---

---

---

---

---

---

---

---

---

---

---

## Laplace-Smoothed Bigram Counts

	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

9/15/15

Speech and Language Processing - Jurafsky and Martin

17

---

---

---

---

---

---

---

---

---

---

---

---

## Laplace-Smoothed Bigram Probabilities

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

	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

9/15/15

Speech and Language Processing - Jurafsky and Martin

18

---

---

---

---

---

---

---

---

---

---

---

---

## Reconstituted Counts

$$c^*(w_{n-1}w_n) = \frac{[C(w_{n-1}w_n) + 1] \times C(w_{n-1})}{C(w_{n-1}) + V}$$

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

9/15/15

Speech and Language Processing - Jurafsky and Martin

19

---

---

---

---

---

---

---

---

---

---

---

---

## Reconstituted Counts (2)

	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

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

9/15/15

Speech and Language Processing - Jurafsky and Martin

20

---

---

---

---

---

---

---

---

---

---

---

---

## Big Change to the Counts!

- $C(\text{want to})$  went from 608 to 238!
- $P(\text{to}|\text{want})$  from .66 to .26!
- Discount  $d = c^*/c$ 
  - $d$  for "chinese food" = .10!!! A 10x reduction
  - So in general, Laplace is a blunt instrument
  - Could use more fine-grained method (add-k)
- Because of this Laplace smoothing not often used for language models, as we have much better methods
- Despite its flaws Laplace (add-1) is still used to smooth other probabilistic models in NLP and IR, especially
  - For pilot studies
  - In document classification
  - In domains where the number of zeros isn't so huge.

9/15/15

Speech and Language Processing - Jurafsky and Martin

21

---

---

---

---

---

---

---

---

---

---

---

---

## Better Smoothing

- Two key ideas
  - Use less context if the counts are missing for longer contexts
  - Use the count of things we've seen once to help estimate the count of things we've never seen

9/15/15

Speech and Language Processing - Jurafsky and Martin

22

---

---

---

---

---

---

---

---

## Backoff and Interpolation

- Sometimes it helps to use less context
  - Condition on less context for contexts you haven't learned much about
- Backoff
  - use trigram if you have good evidence,
  - otherwise bigram, otherwise unigram
- Interpolation
  - mix unigram, bigram, trigram
- Interpolation works better

---

---

---

---

---

---

---

---

## Linear Interpolation

- Simple interpolation

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1 P(w_n|w_{n-2}w_{n-1}) + \lambda_2 P(w_n|w_{n-1}) + \lambda_3 P(w_n) \quad \sum_i \lambda_i = 1$$

- Lambdas conditional on context

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1(w_{n-2}^{n-1}) P(w_n|w_{n-2}w_{n-1}) + \lambda_2(w_{n-2}^{n-1}) P(w_n|w_{n-1}) + \lambda_3(w_{n-2}^{n-1}) P(w_n)$$

---

---

---

---

---

---

---

---



## How to set the lambdas?

- Use a held-out corpus



- Choose  $\lambda$ s to maximize the probability of held-out data:
  - Fix the N-gram probabilities (using the training data)
  - Then search for  $\lambda$ s that give largest probability to held-out set.

---

---

---

---

---

---

---

---

## Types, Tokens and Fish

- Much of what's coming up was first studied by biologists who are often faced with 2 related problems
  - Determining how many species occupy a particular area (types)
  - And determining how many individuals of a given species are living in a given area (tokens)

9/15/15

Speech and Language Processing - Jurafsky and Martin

26

---

---

---

---

---

---

---

---

## One Fish Two Fish

- Imagine you are fishing
  - There are 8 species: carp, perch, whitefish, trout, salmon, eel, catfish, bass
    - Not sure where this fishing hole is...
- You have caught up to now
  - 10 carp, 3 perch, 2 whitefish, 1 trout, 1 salmon, 1 eel = 18 fish
- How likely is it that the next fish to be caught is an eel?
- How likely is it that the next fish caught will be a member of newly seen species?
- Now how likely is it that the next fish caught will be an eel?

9/15/15

Slide adapted from Josh Goodman  
Speech and Language Processing - Jurafsky and Martin

27

---

---

---

---

---

---

---

---

## Fish Lesson

- We need to steal part of the observed probability mass to give it to the as yet unseen N-Grams. Questions are:
  - How much to steal
  - How to redistribute it

9/15/15

Speech and Language Processing - Jurafsky and Martin

28

---

---

---

---

---

---

---

---

## Absolute Discounting

- Just subtract a fixed amount from all the observed counts (call that  $d$ ).
- Redistribute it proportionally based on observed data

9/15/15

Speech and Language Processing - Jurafsky and Martin

29

---

---

---

---

---

---

---

---

## Absolute Discounting w/ Interpolation

$$P_{\text{AbsoluteDiscounting}}(w_i | w_{i-1}) = \frac{\overset{\text{discounted bigram}}{c(w_{i-1}, w_i)} - d}{c(w_{i-1})} + \overset{\text{Interpolation weight}}{\lambda(w_{i-1})} \underset{\text{unigram}}{P(w)}$$

30

---

---

---

---

---

---

---

---

## Kneser-Ney Smoothing

- Better estimate for probabilities of lower-order unigrams!
  - Shannon game: *I can't see without my reading [Figliaccino](#)?*
  - "Francisco" is more common than "glasses"
  - ... but "Francisco" frequently follows "San"
- So  $P(w)$  isn't what we want

---

---

---

---

---

---

---

---

## Kneser-Ney Smoothing

- $P_{\text{continuation}}(w)$ : "How likely is  $w$  to appear as a novel continuation?"
  - For each word, count the number of bigram *types* it completes

$$P_{\text{CONTINUATION}}(w) \propto |\{w_{i-1} : c(w_{i-1}, w) > 0\}|$$

9/15/15

Speech and Language Processing - Jurafsky and Martin

32

---

---

---

---

---

---

---

---

## Kneser-Ney Smoothing

- Normalize by the total number of word bigram types to get a true probability

$$P_{\text{CONTINUATION}}(w) = \frac{|\{w_{i-1} : c(w_{i-1}, w) > 0\}|}{|\{(w_{j-1}, w_j) : c(w_{j-1}, w_j) > 0\}|}$$

---

---

---

---

---

---

---

---

## Kneser-Ney Smoothing

$$P_{KN}(w_i | w_{i-1}) = \frac{\max(c(w_{i-1}, w_i) - d, 0)}{c(w_{i-1})} + \lambda(w_{i-1})P_{CONTINUATION}(w_i)$$

$\lambda$  is a normalizing constant; the probability mass we've discounted

$$\lambda(w_{i-1}) = \frac{d}{c(w_{i-1})} |\{w : c(w_{i-1}, w) > 0\}|$$

the normalized discount

The number of word types that can follow  $w_{i-1}$   
= # of word types we discounted  
= # of times we applied normalized discount

34

---

---

---

---

---

---

---

---