

Natural Language Processing

Lecture 8—9/17/2015
Jim Martin

Today

- Finish up smoothing
 - Kneser-Ney example
- HMMs
 - POS tagging example
 - Basic HMM model
 - Decoding
 - Viterbi

9/17/15 Speech and Language Processing - Jurafsky and Martin 2

Absolute Discounting

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

9/17/15 Speech and Language Processing - Jurafsky and Martin 3

Absolute Discounting w/ Interpolation

$$P_{\text{AbsoluteDiscounting}}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) - d}{c(w_{i-1})} + \lambda(w_{i-1}) P(w)$$

discounted bigram Interpolation weight
unigram

4

Kneser-Ney Smoothing

- Better estimate for probabilities of lower-order unigrams!
 - Shannon game: *I can't see without my reading [Fglaciers](#)?*
 - "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 a word to appear as a novel continuation?"
 - For each word, count the number of bigram types it completes

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

9/17/15

Speech and Language Processing - Jurafsky and Martin

6

Kneser-Ney Smoothing

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

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

Kneser-Ney Smoothing

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

λ 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}

8

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

9/17/15

Speech and Language Processing - Jurafsky and Martin

9

BERP

- Let's look at "chinese food". We'll need:
 - Count("chinese food") 82
 - Count("chinese") 158
 - P_continuation("food")
 - Count of bigrams "food" completes 110
 - Count of all bigram types 9421
 - Count of bigrams that "chinese" starts 17

9/17/15 Speech and Language Processing - Jurafsky and Martin 10

BERP

- Let's look at "chinese food". We'll need:
 - Count("chinese food")
 - Count("chinese")
 - P_continuation("food")
 - Count of bigrams "food" completes
 - Count of all bigram types
 - Count of bigrams that "chinese" starts

9/17/15 Speech and Language Processing - Jurafsky and Martin 11

Break

9/17/15 Speech and Language Processing - Jurafsky and Martin 12

Word Classes: Parts of Speech

- 8 (ish) traditional parts of speech
 - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc
 - Called: parts-of-speech, lexical categories, word classes, morphological classes, lexical tags...
 - Lots of debate within linguistics about the number, nature, and universality of these
 - We' ll completely ignore this debate.

9/17/15

Speech and Language Processing - Jurafsky and Martin

13

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	DET
keys	N
on	P
the	DET
table	N

9/17/15

Speech and Language Processing - Jurafsky and Martin

14

Penn TreeBank POS Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	and, but, or	SYM	symbol	+, %, &
CD	cardinal number	one, two, three	TO	"to"	to
DT	determiner	a, the	UH	interjection	ah, oops
EX	existential 'there'	there	VB	verb, base form	eat
FW	foreign word	mea culpa	VBD	verb, past tense	ate
IN	preposition/sub-conj	of, in, by	VBG	verb, gerund	eating
JJ	adjective	yellow	VBN	verb, past participle	eaten
JJR	adj., comparative	bigger	VBP	verb, non-3sg pres	eat
JJS	adj., superlative	widest	VBZ	verb, 3sg pres	eats
LS	list item marker	1, 2, One	WDT	wh-determiner	which, that
MD	modal	can, should	WP	wh-pronoun	what, who
NN	noun, sing. or mass	llama	WPS	possessive wh-	whose
NNS	noun, plural	llamas	WRB	wh-adverb	how, where
NNP	proper noun, singular	IBM	\$	dollar sign	\$
NNPS	proper noun, plural	Carobmas	#	pound sign	#
PDT	predeterminer	all, both	"	left quote	" or "
POS	possessive ending	's	"	right quote	" or "
PP	personal pronoun	I, you, he	(left parenthesis	[, (, {, <
PRP	possessive pronoun	your, one's)	right parenthesis], }, >
RB	adverb	quickly, never	,	comma	,
RBR	adverb, comparative	faster	.	sentence-final punc	! , ?
RBS	adverb, superlative	fastest	:	mid-sentence punc	! ; ... - -
RP	particle	up, off			

9/17/15

Speech and Language Processing - Jurafsky and Martin

15

POS Tagging

- Words often have more than one part of speech: *back*
 - The *back* door = JJ
 - On my *back* = NN
 - Win the voters *back* = RB
 - Promised to *back* the bill = VB
- The POS tagging problem is to determine the POS tag for a particular instance of a word in context
 - Usually a sentence

9/17/15

Speech and Language Processing - Jurafsky and Martin

16

POS Tagging

- Note this is distinct from the task of identifying which sense of a word is being used given a particular part of speech. That's called word sense disambiguation. We'll get to that later.
 - "... *backed* the car into a pole"
 - "... *backed* the wrong candidate"

9/17/15

Speech and Language Processing - Jurafsky and Martin

17

How Hard is POS Tagging? Measuring Ambiguity

	87-tag Original Brown	45-tag Treebank Brown
Unambiguous (1 tag)	44,019	38,857
Ambiguous (2-7 tags)	5,490	8844
Details:		
2 tags	4,967	6,731
3 tags	411	1621
4 tags	91	357
5 tags	17	90
6 tags	2 (<i>well, beat</i>)	32
7 tags	2 (<i>still, down</i>)	6 (<i>well, set, round, open, fit, down</i>)
8 tags		4 (<i>'s, half, back, a</i>)
9 tags		3 (<i>that, more, in</i>)

9/17/15

Speech and Language Processing - Jurafsky and Martin

18

Two Methods for POS Tagging

1. Rule-based tagging
2. Stochastic
 1. Probabilistic sequence models
 - HMM (Hidden Markov Model) tagging
 - MEMMs (Maximum Entropy Markov Models)

9/17/15

Speech and Language Processing - Jurafsky and Martin

19

POS Tagging as Sequence Classification

- We are given a sentence (an “observation” or “sequence of observations”)
 - *Secretariat is expected to race tomorrow*
- What is the best sequence of tags that corresponds to this sequence of observations?
- Probabilistic view
 - Consider all possible sequences of tags
 - Out of this universe of sequences, choose the tag sequence which is most probable given the observation sequence of n words $w_1 \dots w_n$.

9/17/15

Speech and Language Processing - Jurafsky and Martin

20

Getting to HMMs

- We want, out of all sequences of n tags $t_1 \dots t_n$ the single tag sequence such that $P(t_1 \dots t_n | w_1 \dots w_n)$ is highest.

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- Hat $\hat{}$ means “our estimate of the best one”
- $\operatorname{Argmax}_x f(x)$ means “the x such that $f(x)$ is maximized”

9/17/15

Speech and Language Processing - Jurafsky and Martin

21

Getting to HMMs

- This equation should give us the best tag sequence

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- But how to make it operational? How to compute this value?
- Intuition of Bayesian inference:
 - Use Bayes rule to transform this equation into a set of probabilities that are easier to compute (and give the right answer)

9/17/15

Speech and Language Processing - Jurafsky and Martin

22

Using Bayes Rule

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)} \quad \text{Know this.}$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n) P(t_1^n)}{P(w_1^n)}$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n) P(t_1^n)$$

9/17/15

Speech and Language Processing - Jurafsky and Martin

23

Likelihood and Prior



$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \overbrace{P(w_1^n | t_1^n)}^{\text{likelihood}} \overbrace{P(t_1^n)}^{\text{prior}}$$

$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i)$$



$$P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

9/17/15

Speech and Language Processing - Jurafsky and Martin

24

Two Kinds of Probabilities

- Tag transition probabilities $p(t_i|t_{i-1})$
 - Determiners likely to precede adjs and nouns
 - That/DT flight/NN
 - The/DT yellow/JJ hat/NN
 - So we expect $P(NN|DT)$ and $P(JJ|DT)$ to be high
 - Compute $P(NN|DT)$ by counting in a labeled corpus:

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

$$P(NN|DT) = \frac{C(DT, NN)}{C(DT)} = \frac{56,509}{116,454} = .49$$

9/17/15

Speech and Language Processing - Jurafsky and Martin

25

Two Kinds of Probabilities

- Word likelihood probabilities $p(w_i|t_i)$
 - VBZ (3sg Pres Verb) likely to be “is”
 - Compute $P(is|VBZ)$ by counting in a labeled corpus:

$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

$$P(is|VBZ) = \frac{C(VBZ, is)}{C(VBZ)} = \frac{10,073}{21,627} = .47$$

9/17/15

Speech and Language Processing - Jurafsky and Martin

26

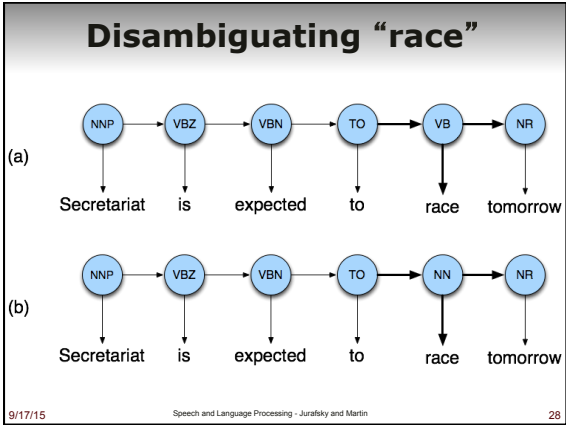
Example: The Verb “race”

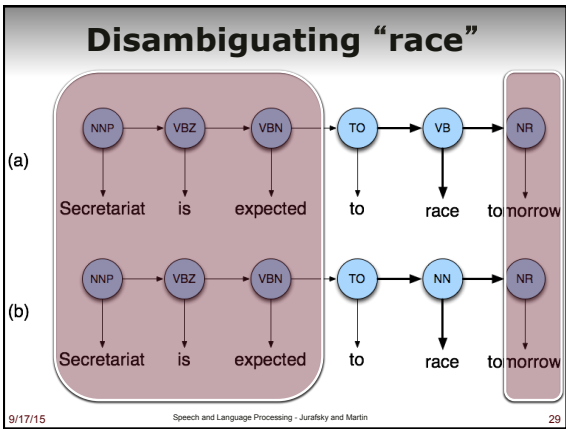
- Secretariat/NNP is/VBZ expected/VBN to/TO **race**/VB tomorrow/NR
- People/NNS continue/VB to/TO inquire/VB the/DT reason/NN for/IN the/DT **race**/NN for/IN outer/JJ space/NN
- How do we pick the right tag?

9/17/15

Speech and Language Processing - Jurafsky and Martin

27





Example

- $P(NN|TO) = .00047$
- $P(VB|TO) = .83$
- $P(race|NN) = .00057$
- $P(race|VB) = .00012$
- $P(NR|VB) = .0027$
- $P(NR|NN) = .0012$

(a)

(b)

- $P(VB|TO)P(NR|VB)P(race|VB) = .00000027$
- $P(NN|TO)P(NR|NN)P(race|NN) = .0000000032$

▪ So we (correctly) choose the verb tag for "race"

9/17/15 Speech and Language Processing - Jurafsky and Martin 30

Hidden Markov Models

- What we've just described is called a Hidden Markov Model (HMM)
- This is a kind of *generative* model.
 - There is a **hidden** underlying generator of observable events
 - The hidden generator can be modeled as a network of states and transitions
 - We want to infer the underlying state sequence given the observed event sequence

9/17/15

Speech and Language Processing - Jurafsky and Martin

31

Hidden Markov Models


- States $Q = q_1, q_2, \dots, q_N$
- Observations $O = o_1, o_2, \dots, o_N$
 - Each observation is a symbol from a vocabulary $V = \{v_1, v_2, \dots, v_V\}$
- Transition probabilities
 - Transition probability matrix $A = \{a_{ij}\}$
$$a_{ij} = P(q_t = j \mid q_{t-1} = i) \quad 1 \leq i, j \leq N$$
- Observation likelihoods
 - Output probability matrix $B = \{b_i(k)\}$
$$b_i(k) = P(X_t = o_k \mid q_t = i)$$
- Special initial probability vector π
$$\pi_i = P(q_1 = i) \quad 1 \leq i \leq N$$

9/17/15

Speech and Language Processing - Jurafsky and Martin

32

HMMs for Ice Cream

- You are a climatologist in the year 2799 studying global warming
- You can't find any records of the weather in Baltimore for summer of 2007
- But you find Jason Eisner's diary which lists how many ice-creams Jason ate every day that summer 
- Your job: figure out how hot it was each day

9/17/15

Speech and Language Processing - Jurafsky and Martin

33

Eisner Task

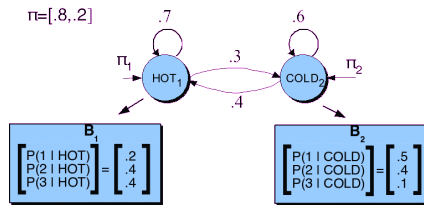
- Given
 - Ice Cream Observation Sequence:
1,2,3,2,2,2,3...
- Produce:
 - Hidden Weather Sequence:
H,C,H,H,H,C, C...

9/17/15

Speech and Language Processing - Jurafsky and Martin

34

HMM for Ice Cream



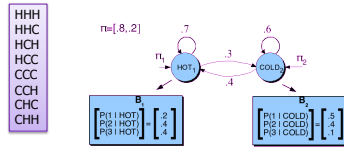
9/17/15

Speech and Language Processing - Jurafsky and Martin

35

Ice Cream HMM

- Let's just do 131 as the sequence
 - How many underlying state (hot/cold) sequences are there?



- How do you pick the right one?

Argmax $P(\text{sequence} | 131)$

9/17/15

Speech and Language Processing - Jurafsky and Martin

36

Question

- If there are 30 or so tags in the Penn set
- And the average sentence is around 20 words...
- How many tag sequences do we have to enumerate to argmax over in the worst case scenario?

30^{20}
