

CSCI 5832
Natural Language Processing

Jim Martin
Lecture 8

2/7/08 1

Today 2/7

- Finish remaining LM issues
 - ♦ Smoothing
 - ♦ Backoff and Interpolation
- Parts of Speech
- POS Tagging
- HMMs and Viterbi

2/7/08 2

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: $P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$
- Reconstructed counts: $c_i^* = (c_i + 1) \frac{N}{N + V}$

2/7/08 3

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

2/7/08 4

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

2/7/08 5

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

2/7/08 6

Big Changes to Counts

- C(count to) went from 608 to 238!
- P(to|want) from .66 to .26!
- Discount $d = c^d/c$
 - d for "chinese food" = .10!!! A 10x reduction
 - So in general, Laplace is a blunt instrument
 - Could use more fine-grained method (add-k)
- Despite its flaws Laplace (add-k) is however still used to smooth other probabilistic models in NLP, especially
 - For pilot studies
 - in domains where the number of zeros isn't so huge.

2/7/08

7

Better Discounting Methods

- Intuition used by many smoothing algorithms
 - Good-Turing
 - Kneser-Ney
 - Witten-Bell
- Is to use the count of things we've seen once to help estimate the count of things we've never seen

2/7/08

8

Good-Turing

- Imagine you are fishing
 - There are 8 species: carp, perch, whitefish, trout, salmon, eel, catfish, bass
- You have caught
 - 10 carp, 3 perch, 2 whitefish, 1 trout, 1 salmon, 1 eel
 - = 18 fish (tokens)
 - = 6 species (types)
- How likely is it that you'll next see another trout?

2/7/08

9

Good-Turing

- Now how likely is it that next species is new (i.e. catfish or bass)

There were 18 distinct events... 3 of those represent singleton species

3/18

2/7/08

10

Good-Turing

- But that 3/18s isn't represented in our probability mass. Certainly not the one we used for estimating another trout.

2/7/08

11

Good-Turing Intuition

- Notation: N_x is the frequency-of-frequency-x
 - So $N_{10}=1$, $N_1=3$, etc
- To estimate total number of unseen species
 - Use number of species (words) we've seen once
 - $c_0^* = c_1$ $p_0 = N_1/N$
- All other estimates are adjusted (down) to give probabilities for unseen

$$c^* = (c + 1) \frac{N_{c+1}}{N_c}$$

2/7/08

Slide from Josh Goodman

12

Good-Turing Intuition

- Notation: N_x is the frequency-of-frequency-x
 - So $N_{10}=1, N_1=3$, etc
- To estimate total number of unseen species
 - Use number of species (words) we've seen once
 - $c_0^* = c_1 \quad p_0 = N_1/N \quad p_0 = N_1/N = 3/18$
- All other estimates are adjusted (down) to give probabilities for unseen

$$P_{GT}^*(\text{things with frequency zero in training}) = \frac{N_1}{N}$$

$$P(\text{cell}) = c^*(1) = (1+1) 1/3 = 2/3$$

$$c^* = (c + 1) \frac{N_{c+1}}{N_c}$$

2/7/08

Slide from Josh Goodman

13

Bigram frequencies of frequencies and GT re-estimates

AP Newswire			Berkeley Restaurant—		
c (MLE)	N_c	c^* (GT)	c (MLE)	N_c	c^* (GT)
0	74,671,100,000	0.0000270	0	2,081,496	0.002553
1	2,018,046	0.446	1	5315	0.533960
2	449,721	1.26	2	1419	1.357294
3	188,933	2.24	3	642	2.373832
4	105,668	3.24	4	381	4.081365
5	68,379	4.22	5	311	3.781350
6	48,190	5.19	6	196	4.500000

2/7/08

14

GT smoothed bigram probs

	i	want	to	eat	chinese	food	lunch	spend
i	0.0014	0.326	0.00248	0.00355	0.000205	0.0017	0.00073	0.000489
want	0.00134	0.00152	0.656	0.000483	0.00455	0.00455	0.00384	0.000483
to	0.000512	0.00152	0.00165	0.284	0.000512	0.0017	0.00175	0.0873
eat	0.00101	0.00152	0.00166	0.00189	0.0214	0.00166	0.0563	0.000585
chinese	0.00283	0.00152	0.00248	0.00189	0.000205	0.519	0.00283	0.000585
food	0.0137	0.00152	0.0137	0.00189	0.000409	0.00366	0.00073	0.000585
lunch	0.00363	0.00152	0.00248	0.00189	0.000205	0.00131	0.00073	0.000585
spend	0.00161	0.00152	0.00161	0.00189	0.000205	0.0017	0.00073	0.000585

2/7/08

15

Backoff and Interpolation

- Another really useful source of knowledge
- If we are estimating:
 - trigram $p(z|xy)$
 - but $c(xyz)$ is zero
- Use info from:
 - Bigram $p(z|y)$
- Or even:
 - Unigram $p(z)$
- How to combine the trigram/bigram/unigram info?

2/7/08

16

Backoff versus interpolation

- **Backoff:** use trigram if you have it, otherwise bigram, otherwise unigram
- **Interpolation:** mix all three

2/7/08

17

Interpolation

- Simple interpolation

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

$$\sum_i \lambda_i = 1$$

- Lambdas conditional on context:

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

2/7/08

18

How to set the lambdas?

- Use a **held-out** corpus
- Choose lambdas which maximize the probability of some held-out data
 - ♦ I.e. fix the N-gram probabilities
 - ♦ Then search for lambda values
 - ♦ That when plugged into previous equation
 - ♦ Give largest probability for held-out set
 - ♦ Can use EM to do this search

2/7/08

19

Practical Issues

- We do everything in log space
 - ♦ Avoid underflow
 - ♦ (also adding is faster than multiplying)

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

2/7/08

20

Language Modeling Toolkits

- SRILM
- CMU-Cambridge LM Toolkit

2/7/08

21

Google N-Gram Release

All Our N-gram are Belong to You

By Peter Norvig - 8/03/2006 11:26:00 AM

Posted by Alex Franz and Thorsten Brants, Google Machine Translation Team

Here at Google Research we have been using word n-gram models for a variety of R&D projects, such as statistical machine translation, speech recognition, spelling correction, entity detection, information extraction, and others. While such models have usually been estimated from training

to share this enormous dataset with everyone. We processed 1,024,908,267,229 words of running text and are publishing the counts for all 1,176,470,663 five-word sequences that appear at least 40 times. There are 13,568,391 unique words, after discarding words that appear less than 200 times.

2/7/08

22

Google N-Gram Release

- serve as the incoming 92
- serve as the incubator 99
- serve as the independent 794
- serve as the index 223
- serve as the indication 72
- serve as the indicator 120
- serve as the indicators 45
- serve as the indispensable 111
- serve as the indispensible 40
- serve as the individual 234

2/7/08

23

LM Summary

- Probability
 - ♦ Basic probability
 - ♦ Conditional probability
 - ♦ Bayes Rule
- Language Modeling (N-grams)
 - ♦ N-gram Intro
 - ♦ The Chain Rule
 - Perplexity
 - ♦ Smoothing:
 - Add-1
 - Good-Turing

2/7/08

24

Break

- Moving quiz to Thursday (2/14)
- Readings
 - ♦ Chapter 2: All
 - ♦ Chapter 3:
 - Skip 3.4.1 and 3.12
 - ♦ Chapter 4
 - Skip 4.7, 4.9, 4.10 and 4.11
 - ♦ Chapter 5
 - Read 5.1 through 5.5

2/7/08

25

Outline

- Probability
- Part of speech tagging
 - ♦ Parts of speech
 - ♦ Tag sets
 - ♦ Rule-based tagging
 - ♦ Statistical tagging
 - Simple most-frequent-tag baseline
 - ♦ Important Ideas
 - Training sets and test sets
 - Unknown words
 - Error analysis
 - ♦ HMM tagging

2/7/08

26

Part of Speech tagging

- Part of speech tagging
 - ♦ Parts of speech
 - ♦ What's POS tagging good for anyhow?
 - ♦ Tag sets
 - ♦ Rule-based tagging
 - ♦ Statistical tagging
 - Simple most-frequent-tag baseline
 - ♦ Important Ideas
 - Training sets and test sets
 - Unknown words
 - ♦ HMM tagging

2/7/08

27

Parts of Speech

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

2/7/08

28

POS examples

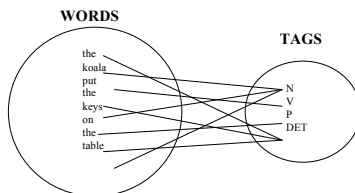
- N noun *chair, bandwidth, pacing*
- V verb *study, debate, munch*
- ADJ adjective *purple, tall, ridiculous*
- ADV adverb *unfortunately, slowly*
- P preposition *of, by, to*
- PRO pronoun *I, me, mine*
- DET determiner *the, a, that, those*

2/7/08

29

POS Tagging: Definition

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



2/7/08

30

POS Tagging example

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

2/7/08

31

What is POS tagging good for?

- First step of a vast number of practical tasks
- Speech synthesis
 - How to pronounce "lead"?
 - INsult inSULT
 - OBject obJECT
 - OVerflow overFLOW
 - DIScount disCOUNT
 - CONtent conTENT
- Parsing
 - Need to know if a word is an N or V before you can parse
- Information extraction
 - Finding names, relations, etc.
- Machine Translation

2/7/08

32

Open and Closed Classes

- Closed class: a relatively fixed membership
 - ♦ Prepositions: of, in, by, ...
 - ♦ Auxiliaries: may, can, will had, been, ...
 - ♦ Pronouns: I, you, she, mine, his, them, ...
 - ♦ Usually function words (short common words which play a role in grammar)
- Open class: new ones can be created all the time
 - ♦ English has 4: Nouns, Verbs, Adjectives, Adverbs
 - ♦ Many languages have these 4, but not all!

2/7/08

33

Open class words

- **Nouns**
 - Proper nouns (Boulder, Granby, Eli Manning)
 - English capitalizes these.
 - Common nouns (the rest).
 - Count nouns and mass nouns
 - Count: have plurals, get counted: goat/goats, one goat, two goats
 - Mass: don't get counted (snow, salt, communism) (*two snows)
- **Adverbs: tend to modify things**
 - Unfortunately, John walked home extremely slowly yesterday
 - Directional/locative adverbs (here, home, downhill)
 - Degree adverbs (extremely, very, somewhat)
 - Manner adverbs (slowly, slinkily, delicately)
- **Verbs:**
 - In English, have morphological affixes (eat/eats/eaten)

2/7/08

34

Closed Class Words

- **Idiosyncratic**
- **Examples:**
 - prepositions: on, under, over, ...
 - particles: up, down, on, off, ...
 - determiners: a, an, the, ...
 - pronouns: she, who, I, ...
 - conjunctions: and, but, or, ...
 - auxiliary verbs: can, may, should, ...
 - numerals: one, two, three, third, ...

2/7/08

35

Prepositions from CELEX

of	540,085	through	14,964	worth	1,563	pace	12
in	331,235	after	13,670	toward	1,390	nigh	9
for	142,421	between	13,275	plus	750	re	4
to	125,691	under	9,525	till	686	mid	3
with	124,965	per	6,515	amongst	525	o'er	2
on	109,129	among	5,090	via	351	but	0
at	100,169	within	5,030	amid	222	ere	0
by	77,794	towards	4,700	underneath	164	less	0
from	74,843	above	3,056	versus	113	midst	0
about	38,428	near	2,026	amidst	67	o'	0
than	20,210	off	1,695	sans	20	thru	0
over	18,071	past	1,575	circa	14	vice	0

2/7/08

36

English particles

aboard	aside	besides	forward(s)	opposite	through
about	astray	between	home	out	throughout
above	away	beyond	in	outside	together
across	back	by	inside	over	under
ahead	before	close	instead	overhead	underneath
alongside	behind	down	near	past	up
apart	below	east, etc.	off	round	within
around	beneath	eastward(s), etc.	on	since	without

2/7/08

37

Conjunctions

and	514,946	yet	5,040	considering	174	forasmuch as	0
that	134,773	since	4,843	lest	131	however	0
but	96,889	where	3,952	albeit	104	immediately	0
or	76,563	nor	3,078	providing	96	in as far as	0
as	54,608	once	2,826	whereupon	85	in so far as	0
if	53,917	unless	2,205	seeing	63	inasmuch as	0
when	37,975	why	1,333	directly	26	insomuch as	0
because	23,626	now	1,290	ere	12	insomuch that	0
so	12,933	neither	1,120	notwithstanding	3	like	0
before	10,720	whenever	913	according as	0	neither nor	0
though	10,329	whereas	867	as if	0	now that	0
than	9,511	except	864	as long as	0	only	0
while	8,144	till	686	as though	0	provided that	0
after	7,042	provided	594	both and	0	providing that	0
whether	5,978	whilst	351	but that	0	seeing as	0
for	5,935	suppose	281	but then	0	seeing as how	0
although	5,424	cos	188	but then again	0	seeing that	0
until	5,072	supposing	185	either or	0	without	0

2/7/08

38

POS tagging: Choosing a tagset

- There are so many parts of speech, potential distinctions we can draw
- To do POS tagging, need to choose a standard set of tags to work with
- Could pick very coarse tagsets
 - N, V, Adj, Adv.
- More commonly used set is finer grained, the "UPenn TreeBank tagset", 45 tags
 - PRP\$, WRB, WP\$, VBG
- Even more fine-grained tagsets exist

2/7/08

39

Penn TreeBank POS Tag set

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>%, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>widest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WPS	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(" or ")</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(' or ')</i>
PRP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([[{ { <</i>
PRPS	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>]] } } ></i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>! ?</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>! ; ... - -)</i>
RP	Particle	<i>up, off</i>			

2/7/08

40

Using the UPenn tagset

- The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.
- Prepositions and subordinating conjunctions marked IN ("although/IN I/PRP..")
- Except the preposition/complementizer "to" is just marked "TO".

2/7/08

41

POS Tagging

- Words often have more than one POS: *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.

These examples from Dekang Lin

2/7/08

42

How hard is POS tagging? Measuring ambiguity

	Original 87-tag corpus	Treebank 45-tag corpus
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>)

2/7/08

43

2 methods for POS tagging

1. Rule-based tagging
 - ♦ (ENGTWOL)
2. Stochastic (=Probabilistic) tagging
 - ♦ HMM (Hidden Markov Model) tagging

2/7/08

44

Rule-based tagging

- Start with a dictionary
- Assign all possible tags to words from the dictionary
- Write rules by hand to selectively remove tags
- Leaving the correct tag for each word.

2/7/08

45

Stage 1 of ENGTWOL Tagging

- First Stage: Run words through FST morphological analyzer to get all parts of speech.
- Example: *Pavlov had shown that salivation ...*

Pavlov	PAVLOV N NOM SG PROPER
had	HAVE V PAST VFIN SVO HAVE PCP2 SVO
shown	SHOW PCP2 SVOO SVO SV
that	ADV PRON DEM SG DET CENTRAL DEM SG
salivation	CS N NOM SG

2/7/08

49

Stage 2 of ENGTWOL Tagging

- Second Stage: Apply NEGATIVE constraints.
- Example: Adverbial "that" rule
 - Eliminates all readings of "that" except the one in
 - "It isn't *that* odd"

Given input: "that"

if

(+1 A/ADV/QUANT) ;if next word is adj/adv/quantifier

(+2 SENT-LIM) ;following which is E-O-S

(NOT -1 SVOC/A) ; and the previous word is not a

verb like "consider" which

allows adjective complements

; in "I consider that odd"

Then eliminate non-ADV tags

Else eliminate ADV

2/7/08

50

Hidden Markov Model Tagging

- Using an HMM to do POS tagging
- Is a special case of Bayesian inference
 - ♦ Foundational work in computational linguistics
 - ♦ Bledsoe 1959: OCR
 - ♦ Mosteller and Wallace 1964: authorship identification
- It is also related to the "noisy channel" model that's the basis for ASR, OCR and MT

2/7/08

51

POS tagging as a sequence classification task

- We are given a sentence (an “observation” or “sequence of observations”)
 - ♦ Secretariat is expected to race tomorrow
- What is the best sequence of tags which 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 $w_1 \dots w_n$.

2/7/08

52

Getting to HMM

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


2/7/08

53

Getting to HMM

- This equation is guaranteed to 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 classification: 
 - ♦ Use Bayes rule to transform into a set of other probabilities that are easier to compute

2/7/08

54

Using Bayes Rule

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

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

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

2/7/08

55

Likelihood and Prior

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

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

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



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

2/7/08

56

Two Kinds of probabilities (1)

- 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 $P(\text{NN}|\text{DT})$ and $P(\text{JJ}|\text{DT})$ to be high
 - But $P(\text{DT}|\text{JJ})$ to be:
 - ♦ Compute $P(\text{NN}|\text{DT})$ by counting in a labeled corpus:

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

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

2/7/08

57

Two kinds of probabilities (2)

- Word likelihood probabilities $p(w_i|t_i)$
 - ♦ VBZ (3sg Pres verb) likely to be “is”
 - ♦ Compute $P(\text{is}|\text{VBZ})$ by counting in a labeled c:

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

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

2/7/08

58

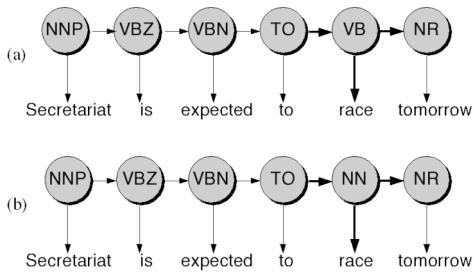
An Example: the verb “race”

- 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?

2/7/08

59

Disambiguating “race”



2/7/08

60

Example

- $P(NN|TO) = .00047$
- $P(VB|TO) = .83$
- $P(\text{race}|NN) = .00057$
- $P(\text{race}|VB) = .00012$
- $P(NR|VB) = .0027$
- $P(NR|NN) = .0012$
- $P(VB|TO)P(NR|VB)P(\text{race}|VB) = .00000027$
- $P(NN|TO)P(NR|NN)P(\text{race}|NN) = .0000000032$
- So we (correctly) choose the verb reading,

2/7/08

61

Hidden Markov Models

- What we've described with these two kinds of probabilities is a Hidden Markov Model
- Let's just spend a bit of time tying this into the model
- First some definitions.



2/7/08

62

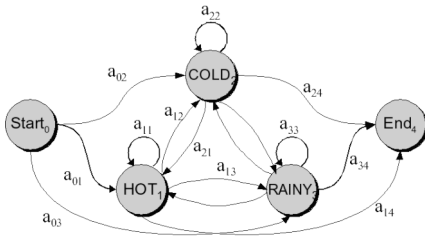
Definitions

- A weighted finite-state automaton adds probabilities to the arcs
 - The sum of the probabilities leaving any arc must sum to one
- A Markov chain is a special case of a WFST in which the input sequence uniquely determines which states the automaton will go through
- Markov chains can't represent inherently ambiguous problems
 - Useful for assigning probabilities to unambiguous sequences

2/7/08

63

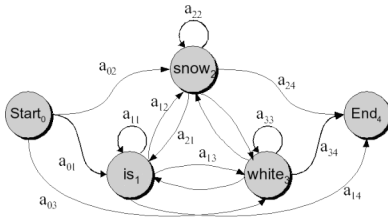
Markov chain for weather



2/7/08

64

Markov chain for words



2/7/08

65

Markov chain = “First-order observable Markov Model”

- A set of states
 - ♦ $Q = q_1, q_2, \dots, q_{N_t}$ the state at time t is q_t
- Transition probabilities:
 - ♦ a set of probabilities $A = a_{01}a_{02} \dots a_{n1} \dots a_{nn}$.
 - ♦ Each a_{ij} represents the probability of transitioning from state i to state j
 - ♦ The set of these is the transition probability matrix A
- Current state only depends on previous state



2/7/08

66

Markov chain for weather

- What is the probability of 4 consecutive rainy days?
- Sequence is rainy-rainy-rainy-rainy
- I.e., state sequence is 3-3-3-3
- $P(3,3,3,3) =$
 - ♦ $\pi_1 a_{11} a_{11} a_{11} a_{11} = 0.2 \times (0.6)^3 = 0.0432$

2/7/08

67

HMM 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, MA for summer of 2007
- But you find Jason Eisner's diary
- Which lists how many ice-creams Jason ate every date that summer
- Our job: figure out how hot it was

2/7/08

68

Hidden Markov Model

- For Markov chains, the output symbols are the same as the states.
 - ♦ See **hot** weather: we're in state **hot**
- But in part-of-speech tagging (and other things)
 - ♦ The output symbols are **words**
 - ♦ But the hidden states are **part-of-speech tags**
- So we need an extension!
- A Hidden Markov Model is an extension of a Markov chain in which the input symbols are not the same as the states.
- This means we don't know which state we are in.

2/7/08

69

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 $A = \{a_{ij}\}$
- Observation likelihoods
 - ♦ Output probability matrix $B = \{b_i(k)\}$
- Special initial probability vector π

2/7/08

70

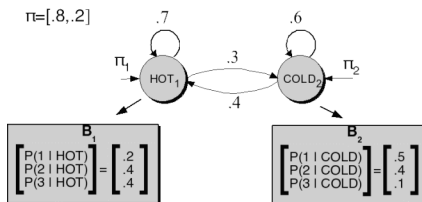
Eisner task

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

2/7/08

71

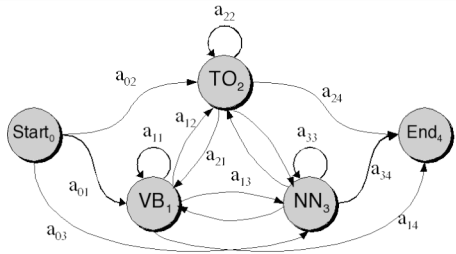
HMM for ice cream



2/7/08

72

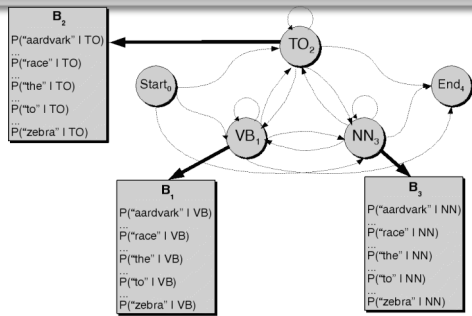
Transitions between the hidden states of HMM, showing A probs



2/7/08

73

B observation likelihoods for POS HMM



2/7/08

74
