

Part-of-Speech Tagging

Chapter 8
(8.1-8.4.6)

Outline

- Parts of speech (POS)
- Tagsets
- POS Tagging
 - Rule-based tagging
 - Probabilistic (HMM) tagging

Garden Path Sentences

- The old dog the footsteps of the young

9/18/2018

Speech and Language Processing - Jurafsky and Martin

3

Parts of Speech

- 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.

9/18/2018

Speech and Language Processing - Jurafsky and Martin

4

Parts of Speech

- Traditional parts of speech
 - ~ 8 of them

Parts of Speech	
NOUN Name of a person, place, thing or idea. Examples: Daniel, London, table, hope - Mary uses a blue pen for her notes.	PRONOUN A pronoun is used in place of a noun or noun phrase to avoid repetition. Examples: I, you, it, we, us, them, those - I want her to dance with me.
ADJECTIVE Describes, modifies or gives more information about a noun or pronoun. Examples: cold, happy, young, two, fun - The little girl has a pink hat.	VERB Shows an action or a state of being. Examples: go, speak, eat, live, are, is - I listen to the word and then repeat it.
ADVERB Modifies a verb, an adjective or another adverb. It tells how (often), where, when. Examples: slowly, very, always, well, too - Yesterday, I ate my lunch quickly.	PREPOSITION Shows the relationship of a noun or pronoun to another word. Examples: at, on, in, from, with, about - I left my keys on the table for you.
CONJUNCTION Joins two words, ideas, phrases together and shows how they are connected. Examples: and, or, but, because, yet, so - I was hot and tired but still finished it.	INTERJECTION A word or phrase that expresses a strong emotion. It is a short exclamation. Examples: Ouch! Hey! Oh! Watch out! - Wow! I passed my English exam.

5

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*

POS Tagging

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

WORD

tag

the
koala
put
the
keys
on
the
table

9/18/2018

Speech and Language Processing - Jurafsky and Martin

7

POS Tagging

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

WORD

tag

the
koala
put
the
keys
on
the
table

DET

9/18/2018

Speech and Language Processing - Jurafsky and Martin

8

POS Tagging

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

WORD

tag

the
koala
put
the
keys
on
the
table

DET
N

9/18/2018

Speech and Language Processing - Jurafsky and Martin

9

POS Tagging

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

WORD

tag

the
koala
put
the
keys
on
the
table

DET
N
V

9/18/2018

Speech and Language Processing - Jurafsky and Martin

10

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	
on	
the	
table	

9/18/2018

Speech and Language Processing - Jurafsky and Martin

11

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	
the	
table	

9/18/2018

Speech and Language Processing - Jurafsky and Martin

12

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	
table	

9/18/2018

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	

9/18/2018

Speech and Language Processing - Jurafsky and Martin

14

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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

15

Why is POS Tagging Useful?

- First step of many practical tasks, e.g.
- Speech synthesis (aka text to speech)
 - How to pronounce "lead"?
 - OBject obJECT
 - CONtent conTENT
- Parsing
 - Need to know if a word is an N or V before you can parse
- Information extraction
 - Finding names, relations, etc.
- Language modeling
 - Backoff

9/18/2018

Speech and Language Processing - Jurafsky and Martin

16

Why is POS Tagging Difficult?

- Words often have more than one POS:

back

- The back door = adjective
- On my back =
- Win the voters back =
- Promised to back the bill =

Why is POS Tagging Difficult?

- Words often have more than one POS:

back

- The back door = adjective
- On my back = noun
- Win the voters back =
- Promised to back the bill =

Why is POS Tagging Difficult?

- Words often have more than one POS:

back

- The back door = adjective
- On my back = noun
- Win the voters back = adverb
- Promised to back the bill =

Why is POS Tagging Difficult?

- Words often have more than one POS:

back

- The back door = adjective
- On my back = noun
- Win the voters back = adverb
- Promised to back the bill = verb
- The POS tagging problem is to determine the POS tag for a particular instance of a word.

POS Tagging

- Input: Plays well with others
- Ambiguity: NNS/VBZ UH/JJ/NN/RB IN NNS
- Output: Plays/VBZ well/RB with/IN others, NNS

Penn
Treebank
POS tags

POS tagging performance

- How many tags are correct? (Tag accuracy)
 - About 97% currently
 - But baseline is already 90%
 - Baseline is performance of stupidest possible method
 - Tag every word with its most frequent tag
 - Tag unknown words as nouns
 - Partly easy because
 - Many words are unambiguous
 - You get points for them (*the*, *a*, etc.) and for punctuation marks!

Deciding on the correct part of speech can be difficult even for people

- Mrs/NNP Shaefer/NNP never/RB got/VBD around/RP to/TO joining/VBG
- All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB around/IN the/DT corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ around/RB 250/CD

How difficult is POS tagging?

- About 11% of the word types in the Brown corpus are ambiguous with regard to part of speech
- But they tend to be very common words.
E.g., *that*
 - I know *that* he is honest = IN
 - Yes, *that* play was nice = DT
 - You can't go *that* far = RB
- 40% of the word tokens are ambiguous

Open vs. Closed Classes

- **Closed class:** *why?*
 - Determiners: a, an, the
 - 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:** *why?*
 - English has 4: Nouns, Verbs, Adjectives, Adverbs
 - Many languages have these 4, but not all!

9/18/2018

Speech and Language Processing - Jurafsky and Martin

25

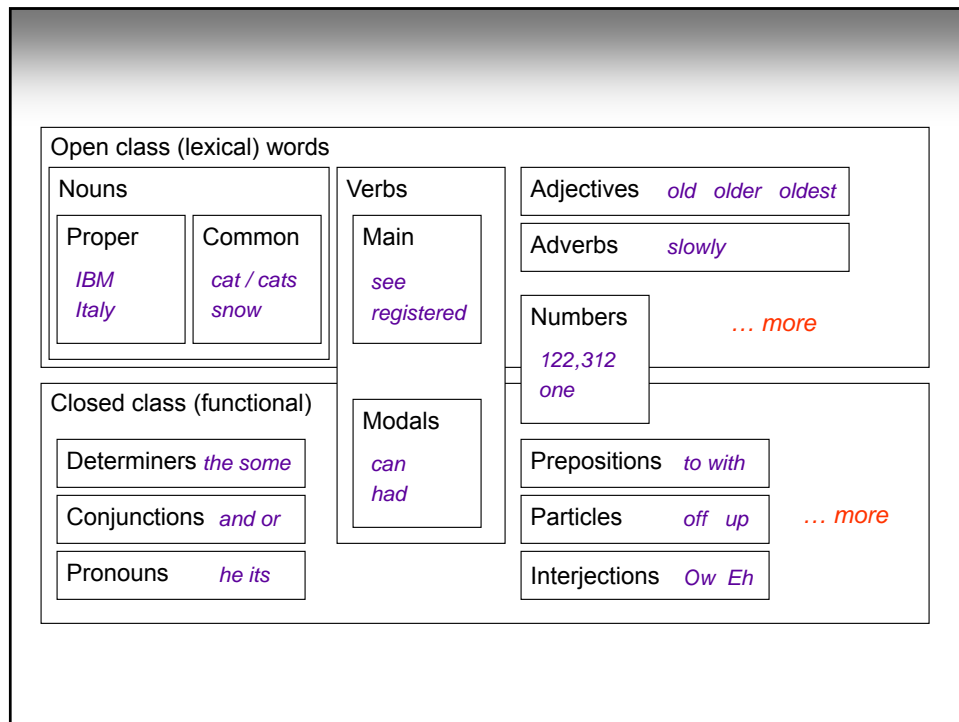
Open vs. Closed Classes

- **Closed class:** a small fixed membership
 - Determiners: a, an, the
 - 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!

9/18/2018

Speech and Language Processing - Jurafsky and Martin

26



Open Class Words

■ Nouns

- Proper nouns (Pittsburgh, Patti Beeson)
 - 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)

Closed Class Words

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, ...*

9/18/2018

Speech and Language Processing - Jurafsky and Martin

29

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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

30

POS Tagging

Choosing a Tagset

- There are so many parts of speech, potential distinctions we can draw
- To do POS tagging, we 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 "Penn TreeBank tagset", 45 tags
- Even more fine-grained tagsets exist

9/18/2018

Speech and Language Processing - Jurafsky and Martin

31

Penn TreeBank POS Tagset

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>	VCN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</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>	WP\$	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>
PRP\$	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>			

9/18/2018

Speech and Language Processing - Jurafsky and Martin

32

Using the Penn Tagset

- The/? grand/? jury/? commmented/? on/?
a/? number/? of/? other/? topics/? ./?

9/18/2018

Speech and Language Processing - Jurafsky and Martin

33

Using the Penn Tagset

- The/DT grand/JJ jury/NN commented/VBD
on/IN a/DT number/NN of/IN other/JJ
topics/NNS ./.

9/18/2018

Speech and Language Processing - Jurafsky and Martin

34

Recall POS Tagging Difficulty

- 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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

35

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>)

9/18/2018

Speech and Language Processing - Jurafsky and Martin

36

Tagging Whole Sentences with POS is Hard too

- Ambiguous POS contexts
 - E.g., Time flies like an arrow.
- Possible POS assignments
 - Time/[V,N] flies/[V,N] like/[V,Prep] an/Det arrow/N
 - Time/N flies/V like/Prep an/Det arrow/N
 - Time/V flies/N like/Prep an/Det arrow/N
 - Time/N flies/N like/V an/Det arrow/N
 -

37

How Do We Disambiguate POS?

- Many words have only one POS tag (e.g. is, Mary, smallest)
- Others have a single *most likely* tag (e.g. Dog is less used as a V)
- Tags also tend to *co-occur* regularly with other tags (e.g. Det, N)
- In addition to conditional probabilities of words $P(w_1|w_{n-1})$, we can look at POS likelihoods $P(t_1|t_{n-1})$ to disambiguate sentences and to assess sentence likelihoods

38

More and Better Features → Feature-based tagger

- Can do surprisingly well just looking at a word by itself:
 - Word the: the → DT
 - Lowercased word Importantly: importantly → RB
 - Prefixes unfathomable: un- → JJ
 - Suffixes Importantly: -ly → RB
 - Capitalization Meridian: CAP → NNP
 - Word shapes 35-year: d-x → JJ

Overview: POS Tagging Accuracies

- Rough accuracies:

- Most freq tag: ~90% / ~50%
- Trigram HMM: ~95% / ~55%
- Maxent $P(t|w)$: 93.7% / 82.6%
- Upper bound: ~98% (human)

Most errors
on unknown
words

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.

9/18/2018

Speech and Language Processing - Jurafsky and Martin

41

Start With a Dictionary

- she:
- promised:
- to
- back:
- the:
- bill:

9/18/2018

Speech and Language Processing - Jurafsky and Martin

42

Start With a Dictionary

- she: PRP
- promised: VBN,VBD
- to TO
- back: VB, JJ, RB, NN
- the: DT
- bill: NN, VB

9/18/2018

Speech and Language Processing - Jurafsky and Martin

43

Assign Every Possible Tag

			NN			
			RB			
	VBN		JJ		VB	
PRP	VBD		TO	VB	DT	NN
She	promised	to	back	the	bill	

9/18/2018

Speech and Language Processing - Jurafsky and Martin

44

Write Rules to Eliminate Tags

Eliminate VBN if VBD is an option when
VBN|VBD follows “<start> PRP”

			NN			
			RB			
	VBN		JJ		VB	
PRP	VBD		TO	VB	DT	NN
She	promised		to	back	the	bill

9/18/2018

Speech and Language Processing - Jurafsky and Martin

45

POS tag sequences

- Some tag sequences are more likely occur than others
- POS Ngram view
<https://books.google.com/ngrams/graph?content=ADJ+NOUN%2CADV+NOUN%2C+ADV+VERB>

Existing methods often model POS tagging as a sequence tagging problem

46

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 $w_1 \dots w_n$.

9/18/2018

Speech and Language Processing - Jurafsky and Martin

47

How do you predict the tags?

- Two types of information are useful
 - Relations between words and tags
 - Relations between tags and tags
 - DT NN, DT JJ NN...

48

Getting to HMMs (Hidden Markov Models)

- 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”

9/18/2018

Speech and Language Processing - Jurafsky and Martin

49

Getting to HMMs

- 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 this equation into a set of other probabilities that are easier to compute

9/18/2018

Speech and Language Processing - Jurafsky and Martin

50

Using Bayes Rule

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

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

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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

51

Statistical POS tagging

- What is the most likely sequence of tags for the given sequence of words w

$$\begin{aligned} \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}|\mathbf{w}) &= \operatorname{argmax}_{\mathbf{t}} \frac{P(\mathbf{t}, \mathbf{w})}{P(\mathbf{w})} \\ &= \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}, \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t})P(\mathbf{w}|\mathbf{t}) \end{aligned}$$

$P(\text{DT JJ NN} \mid \text{a smart dog}) =$

52

Statistical POS tagging

- What is the most likely sequence of tags for the given sequence of words w

$$\begin{aligned}\operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}|\mathbf{w}) &= \operatorname{argmax}_{\mathbf{t}} \frac{P(\mathbf{t}, \mathbf{w})}{P(\mathbf{w})} \\ &= \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}, \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t})P(\mathbf{w}|\mathbf{t})\end{aligned}$$

$$\begin{aligned}P(\text{DT JJ NN} | \text{a smart dog}) &= P(\text{DD JJ NN a smart dog}) / P(\text{a smart dog}) \\ &\propto P(\text{DD JJ NN a smart dog}) \\ &= P(\text{DD JJ NN}) P(\text{a smart dog} | \text{DD JJ NN})\end{aligned}$$

53

Likelihood and Prior



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

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



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

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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

54

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
 - But $P(DT|JJ)$ to be:
 - Compute $P(NN|DT)$ by counting in a labeled corpus:

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

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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

55

Two Kinds of Probabilities

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

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

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

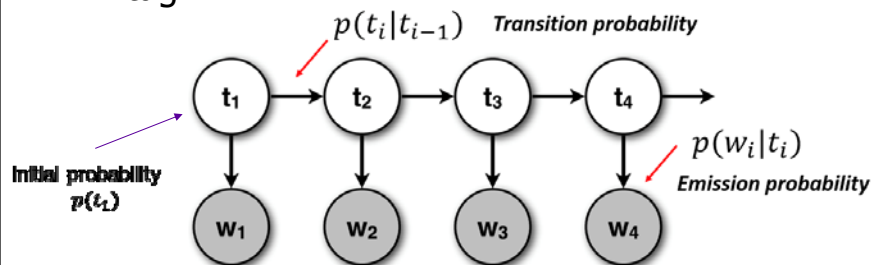
9/18/2018

Speech and Language Processing - Jurafsky and Martin

56

Put them together

- Two independent assumptions
 - Approximate $P(t)$ by a bi(or N)-gram model
 - Assume each word depends only on its POS tag



57

Table representation

Transition Matrix A

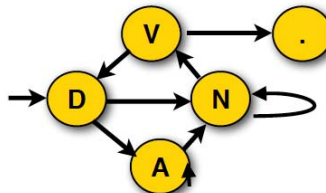
	D	N	V	A	.
D		0.8		0.2	
N		0.7	0.3		
V	0.6				0.4
A		0.8		0.2	
.					

Emission Matrix B

	the	man	ball	throws	sees	red	blue	.
D	1.0							
N		0.7	0.3					
V				0.6	0.4			
A						0.8	0.2	
.								1

Initial state vector π

	D	N	V	A	.
π	1.0				

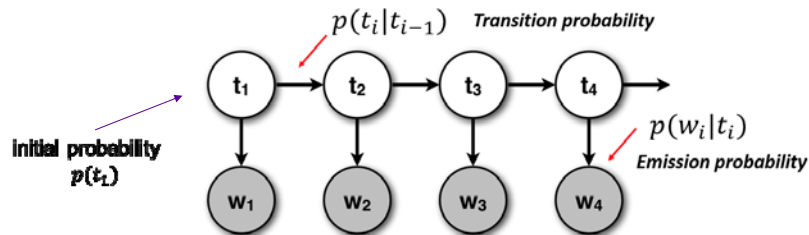


Let $\lambda = \{A, B, \pi\}$ represents all parameters

58

Prediction in generative model

- **Inference:** What is the most likely sequence of tags for the given sequence of words **w**



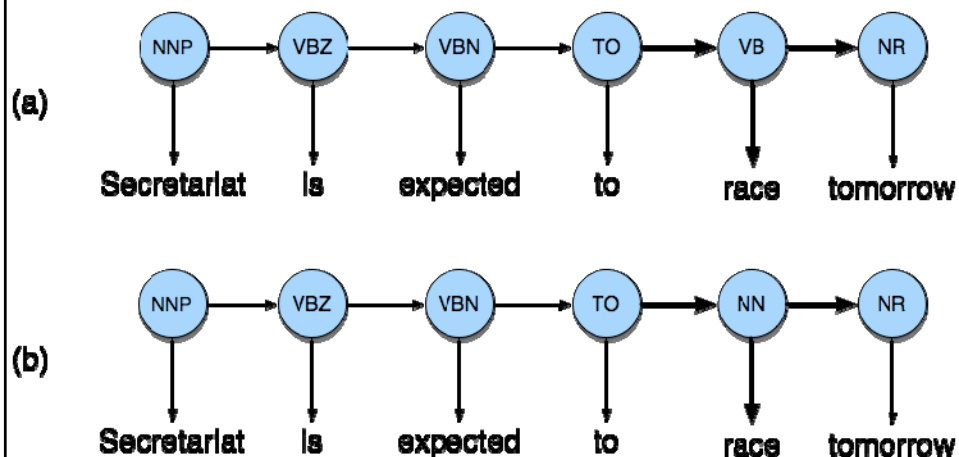
- What are the latent states that most likely generate the sequence of word **w**

59

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?

Disambiguating "race"



9/18/2018

Speech and Language Processing - Jurafsky and Martin

61

Example

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

9/18/2018

Speech and Language Processing - Jurafsky and Martin

62

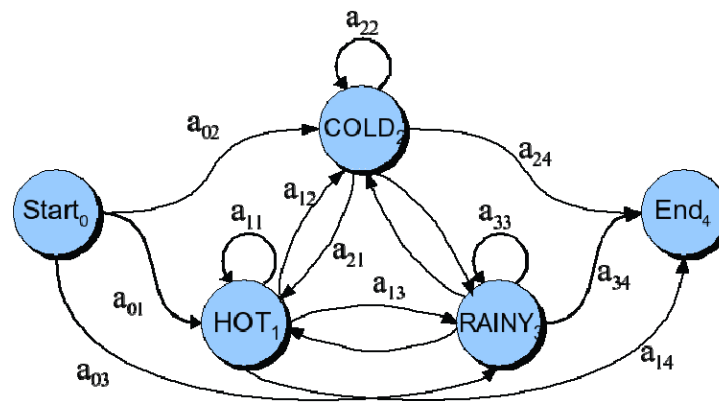
Hidden Markov Models

- What we've described with these two kinds of probabilities is a Hidden Markov Model (HMM)

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 WFSA 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

Markov Chain for Weather

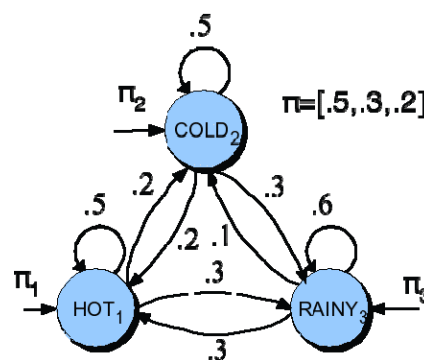


9/18/2018

Speech and Language Processing - Jurafsky and Martin

65

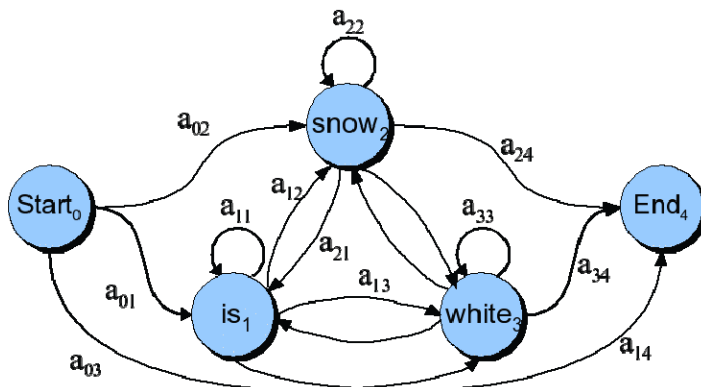
Weather continued



9/18/2018

66

Markov Chain for Words



9/18/2018

Speech and Language Processing - Jurafsky and Martin

67

Markov Chain: "First-order observable Markov Model"

- A set of states
 - $Q = q_1, q_2, \dots, q_N$; 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

$$P(q_i | q_1 \dots q_{i-1}) = P(q_i | q_{i-1})$$

9/18/2018

Speech and Language Processing - Jurafsky and Martin

68

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
- $P(3,3,3,3) =$

9/18/2018

Speech and Language Processing - Jurafsky and Martin

69

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
- $P(3,3,3,3) =$
 - $\pi_3 a_{33} a_{33} a_{33} = 0.2 \times (0.6)^3 = 0.0432$

9/18/2018

Speech and Language Processing - Jurafsky and Martin

70

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 Pittsburgh for summer of 2018
- But you find a diary
- Which lists how many ice-creams someone ate every date that summer
- Our job: figure out how hot it was

9/18/2018

Speech and Language Processing - Jurafsky and Martin

71

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.**

9/18/2018

Speech and Language Processing - Jurafsky and Martin

72

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 $A = \{a_{ij}\}$
 $a_{ij} = P(q_t = j | q_{t-1} = i) \quad 1 \leq i, j \leq N$
- **Observation likelihoods**
 - Output probability matrix $B = \{b_i(k)\}$
 $b_i(k) = P(X_t = o_k | q_t = i)$
 $\pi_i = P(q_1 = i) \quad 1 \leq i \leq N$
- **Special initial probability vector** π

9/18/2018

Speech and Language Processing - Jurafsky and Martin

73

Task

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

9/18/2018

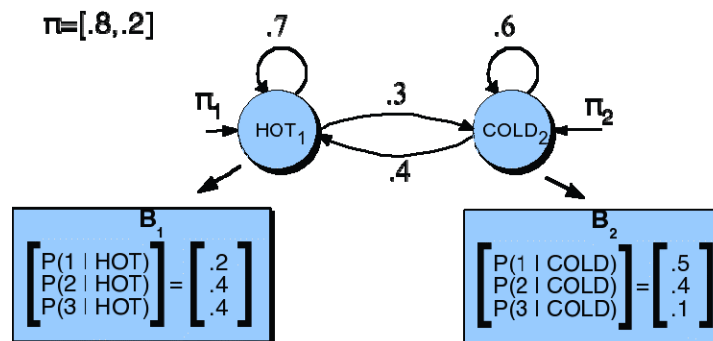
Speech and Language Processing - Jurafsky and Martin

74

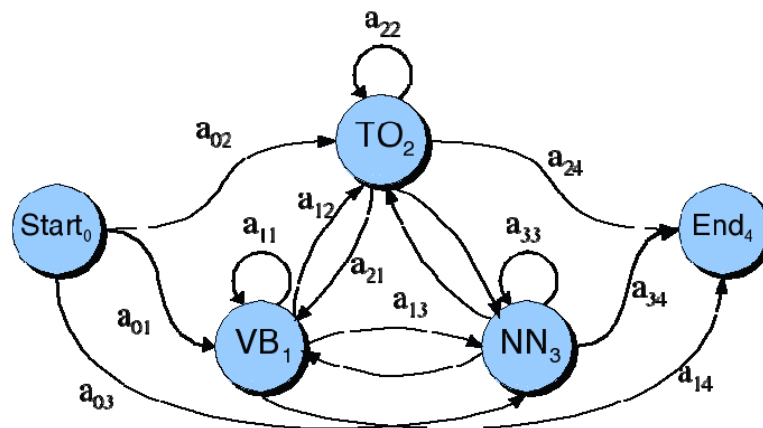
Weather/Ice Cream HMM

- Hidden States: {Hot,Cold}
- Transition probabilities (A Matrix) between H and C
- Observations: {1,2,3} # of ice creams eaten per day

HMM for Ice Cream



Back to POS Tagging: Transition Probabilities

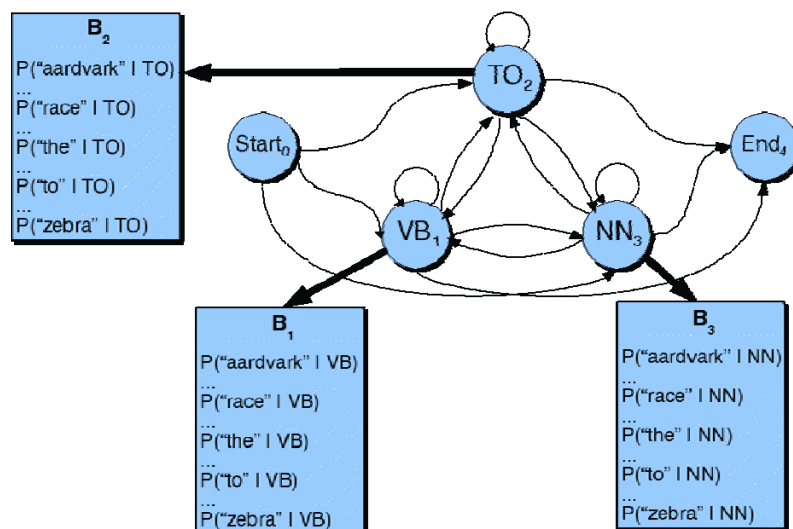


9/18/2018

Speech and Language Processing - Jurafsky and Martin

77

Observation Likelihoods



9/18/2018

Speech and Language Processing - Jurafsky and Martin

78

What can HMMs Do?

- **Likelihood:** Given an HMM λ and an observation sequence O , determine the likelihood $P(O, \lambda)$: *language modeling*
- **Decoding:** Given an observation sequence O and an HMM λ , discover the *best* hidden state sequence Q : Given seq of ice creams, what was the most likely weather on those days? (*tagging*)
- **Learning:** Given an observation sequence O and the set of states in the HMM, learn the HMM *parameters*

9/18/2018

79

Decoding

- Ok, now we have a complete model that can give us what we need. Recall that we need to get

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

- We could just enumerate all paths given the input and use the model to assign probabilities to each.
 - Not a good idea.
 - In practice: Viterbi Algorithm (dynamic programming)

9/18/2018

Speech and Language Processing - Jurafsky and Martin

80

Viterbi Algorithm

- Intuition: since state transition out of a state only depend on the current state (and not previous states), we can record for each state the optimal path
- We record
 - Cheapest cost to state at step
 - Backtrace for that state to best predecessor

9/18/2018

Speech and Language Processing - Jurafsky and Martin

81

Viterbi Summary

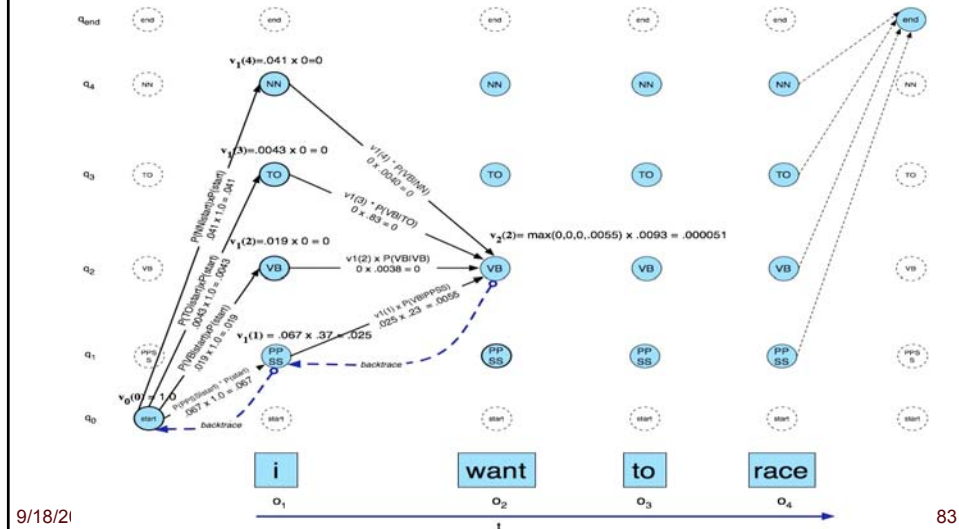
- Create an array
 - With columns corresponding to inputs
 - Rows corresponding to possible states
- Sweep through the array in one pass filling the columns left to right using our transition probs and observations probs
- Dynamic programming key is that we need only store the MAX prob path to each cell (not all paths).

9/18/2018

Speech and Language Processing - Jurafsky and Martin

82

Viterbi Example



Another Viterbi Example

- Analyzing "Fish sleep"
 - Done in class

Evaluation

- So once you have your POS tagger running how do you evaluate it?
 - Overall error rate with respect to a gold-standard test set.
 - Error rates on particular tags
 - Error rates on particular words
 - Tag confusions...
- Need a baseline – just the most frequent tag is 90% accurate!

9/18/2018

Speech and Language Processing - Jurafsky and Martin

85

Error Analysis

- Look at a confusion matrix

	IN	JJ	NN	NNP	RB	VBD	VBN
IN	—	.2			.7		
JJ	.2	—	3.3	2.1	1.7	.2	2.7
NN		8.7	—				.2
NNP	.2	3.3	4.1	—	.2		
RB	2.2	2.0	.5		—		
VBD		.3	.5			—	4.4
VBN		2.8				2.6	—

- See what errors are causing problems
 - Noun (NN) vs ProperNoun (NNP) vs Adj (JJ)
 - Preterite (VBD) vs Participle (VBN) vs Adjective (JJ)

9/18/2018

Speech and Language Processing - Jurafsky and Martin

86

Evaluation

- The result is compared with a manually coded "Gold Standard"
 - Typically accuracy reaches 96-97%
 - This may be compared with result for a baseline tagger (one that uses no context).
- Important: 100% is impossible even for human annotators.

9/18/2018

Speech and Language Processing - Jurafsky and Martin

87

More Complex Issues

- Tag indeterminacy: when 'truth' isn't clear
Caribbean cooking, child seat
- Tagging multipart words
wouldn't --> would/MD n't/RB
- How to handle unknown words
 - Assume all tags equally likely
 - Assume same tag distribution as all other singletons in corpus
 - Use morphology, word length,....

Other Tagging Tasks

- Noun Phrase (NP) Chunking
- [the student] said [the exam] is hard
- Three tags
 - B = beginning of NP
 - I = continuing in NP
 - O = other word
- Tagging result
 - The/B student/I said/O the/B exam/I is/O hard/O

9/18/2018

Speech and Language Processing - Jurafsky and Martin

89

Summary

- Parts of speech
- Tagsets
- Part of speech tagging
- Rule-Based, HMM Tagging

9/18/2018

Speech and Language Processing - Jurafsky and Martin

90