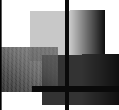


6.863J Natural Language Processing

Lecture 6: The Red Pill or the Blue Pill, Episode 1: part-of-speech tagging



Instructor: Robert C. Berwick
berwick@ai.mit.edu

The Menu Bar



• Administrivia:

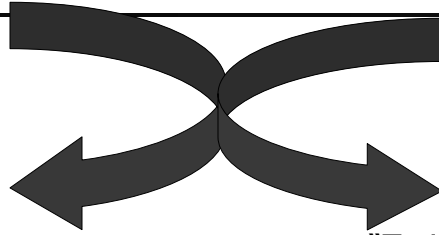
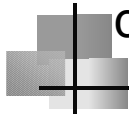
- Schedule alert: Lab1b due today
 - Lab 2a, released today; Lab 2b, this Weds

Agenda:

Red vs. Blue:

- Ngrams as models of language
- Part of speech 'tagging' via statistical models
- Ch. 6 & 8 in Jurafsky

The Great Divide in NLP: the red pill or the blue pill?

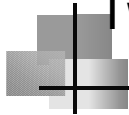


“Knowledge Engineering” approach
Rules built by hand w/
K of Language
“Text understanding”

“Trainable Statistical”
Approach
Rules inferred from lots
of data (“corpora”)
“Information retrieval”

6.863J/9.611J SP04 Lecture 6

Two ways



- Probabilistic model - some constraints on morpheme sequences using prob of one character appearing before/after another
 $\text{prob}(\text{ing} \mid \text{stop})$ vs. $\text{prob}(\text{ly} \mid \text{stop})$
- Generative model – concatenate then fix up joints
 - $\text{stop} + \text{-ing} = \text{stopping}$, $\text{fly} + \text{s} = \text{flies}$
 - Use a cascade of transducers to handle all the fixups

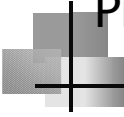
6.863J/9.611J SP04 Lecture 6



The big picture II

- In general: 2 approaches to NLP
- Knowledge Engineering Approach
 - Grammars constructed by hand
 - Domain patterns discovered by human expert via introspection & inspection of 'corpus'
 - Laborious tuning
- Automatically Trainable Systems
 - Use statistical methods when possible
 - Learn rules from annotated (or o.w. processed) corpora

6.863J/9.611J SP04 Lecture 6



Preview of tagging

- What is tagging?
- Input: word sequence:
Police police police
- Output: classification (binning) of words -
Noun Verb Noun *or*
[Help!]

6.863J/9.611J SP04 Lecture 6

Preview of tagging & pills: red pill and blue pill methods

- Method 1: statistical (n-gram)
- Method 2: more symbolic (but still includes some probabilistic training + fixup) – ‘example based’ learning

6.863J/9.611J SP04 Lecture 6

What is part of speech tagging & why?

Input: the lead paint is unsafe

Output: the/Det lead/N paint/N is/V unsafe/Adj

Or: BOS the lyric beauties of Schubert `s Trout Quintet : its elemental rhythms and infectious melodies : make it a source of pure pleasure for almost all music listeners ./

6.863J/9.611J SP04 Lecture 6

Tagging for this..

The/DT lyric/JJ beauties/NNS of/IN
Schubert/NNP 's/POS Trout/NNP Quintet/NNP
--/:
its/PRP\$ elemental/JJ rhythms/NNS
and/CC infectious/JJ melodies/NNS
--/: make/VBP it/PRP
a/DT source/NN of/IN pure/JJ pleasure/NN
for/IN almost/RB all/DT music/NN listeners/NNS ./.

6.863J/9.611J SP04 Lecture 6

Tagging words

- Well defined
- Easy, but not too easy (not AI-complete)
- Data available for machine learning methods
- Evaluation methods straightforward

6.863J/9.611J SP04 Lecture 6

Why should we care?

- The first statistical NLP task
- Been done to death by different methods
- Easy to evaluate (how many tags are correct?)
- Canonical finite-state task
 - Can be done well with methods that look at local context
 - Though should “really” do it by parsing!

6.863J/9.611J SP04 Lecture 6

Why should we care?

- “Simplest” case of recovering surface, underlying form via statistical means
- We are modeling $p(\text{word seq}, \text{tag seq})$
- The tags are hidden, but we see the words
- Is tag sequence T likely given these words?

6.863J/9.611J SP04 Lecture 6

Tagging as n-grams

- Most likely word? Most likely tag t given a word w ? = $P(\text{tag}|\text{word})$ – not quite
 - Task of predicting the next word
 - Woody Allen:
“I have a gub”
- But in general: predict the N^{th} tag from the preceding $n-1$ word (tags) aka N-gram

6.863J/9.611J SP04 Lecture 6

Summary of n-grams

- n-grams define a probability model over sequences
- we have seen examples of sequences of words, but one can also look at characters
- n-grams deal with sparse data by using the Markov assumption

6.863J/9.611J SP04 Lecture 6

Markov models: the 'pure' statistical model...

- 0th order Markov model: $P(w_i)$
- 1st order Markov model: $P(w_i|w_{i-1})$
- 2nd order Markov model: $P(w_i|w_{i-1} w_{i-2})$
- ...
- Where do these probability estimates come from?
- Counts: $P(w_i | w_{i-1}) = \text{count}(w_i, w_{i-1}) / \text{count}(w_{i-1})$
(so-called maximum likelihood estimate - MLE)

6.863J/9.611J SP04 Lecture 6

N-grams

- But...How many possible distinct probabilities will be needed?, i.e. parameter values
- Total number of word tokens in our training data
- Total number of unique words: word types is our vocabulary size

6.863J/9.611J SP04 Lecture 6

n-gram Parameter Sizes – large!

- Let V be the vocabulary, size of V is $|V|$, 3000 distinct types say
- $P(W_i=x)$ how many different values for W_i ?

- $P(W_i = x \mid W_j = y)$, # distinct doubles =

$$3 \times 10^3 \times 3 \times 10^3 = 9 \times 10^6$$

- $P(W_i = x \mid W_k = z, W_j = y)$, how many distinct triples?

$$27 \times 10^9$$

6.863J/9.611J SP04 Lecture 6

Choosing n

Suppose we have a vocabulary (V) = 20,000 words

n	Number of bins
2 (bigrams)	400,000,000
3 (trigrams)	8,000,000,000,000
4 (4-grams)	1.6×10^{17}

6.863J/9.611J SP04 Lecture 6

How far into the past should we go?


- "*long distance* ___"
 - Next word? *Call*?
 - $p(w_n|w...)$
 - Consider special case above
 - Approximation says that
| long distance call|/|distance call| \approx |distance call|/|distance
 - If context 1 word back = bigram
- But even better approx if 2 words back: *long distance* ___
Not always right: long distance runner/long distance call
Further you go: *collect long distance* _____

6.863J/9.611J SP04 Lecture 6

Parameter size vs. corpus size

- Corpus: said the joker to the thief
|V| = 5
- What's the max # of parameters?
- What's observed? (All pairs)
- We observe only |V| many bigrams!
- V had better be large wrt # parameters

6.863J/9.611J SP04 Lecture 6



Reliability vs. discrimination


“large green _____”

tree? mountain? frog? car?

“swallowed the large green _____”

pill? broccoli?

6.863J/9.611J SP04 Lecture 6



Reliability vs. discrimination

- larger n : more information about the context of the specific instance (greater discrimination)
- smaller n : more instances in training data, better statistical estimates (more reliability)

6.863J/9.611J SP04 Lecture 6



Statistical estimators

Example:

Corpus: five Jane Austen novels

N = 617,091 words

V = 14,585 unique words

Task: predict the next word of the trigram "inferior to ____"

from test data, *Persuasion*:

"[In person, she was] inferior to *both* [sisters.]"

6.863J/9.611J SP04 Lecture 6

Shakespeare in lub...



The unkindest cut of all

- Shakespeare: 884,647 words or *tokens* (Kucera, 1992)
- 29,066 *types* (incl. proper nouns)
- So, # bigrams is $29,066^2 > 844$ million. 1 million word training set doesn't cut it – only 300,000 diff't bigrams appear
- Most entries are zero
- So we can't go very far...

6.863J/9.611J SP04 Lecture 6

Bigram models in practice

- $P(\text{Bush read a book}) =$
 $P(\text{Bush} \mid \text{BOS}) \times$
 $P(\text{read} \mid \text{Bush}) \times$
 $P(\text{a} \mid \text{read}) \times$
 $P(\text{book} \mid \text{a}) \times$
 $P(\text{EOS} \mid \text{book})$

Estimate via counts $P(w_i \mid w_{i-1}) = \text{count}(w_i, w_{i-1}) / \text{count}(w_{i-1})$

On unseen data, $\text{count}(w_i, w_{i-1})$ or worse, $\text{count}(w_{i-1})$ could be zero! What to do?

6.863J/9.611J SP04 Lecture 6

How to Estimate?

- $p(z \mid xy) = ?$
- Suppose our training data includes
... xya ..
... xyd ...
... xyd ...
but never xyz
- Should we conclude
 $p(a \mid xy) = 1/3?$
 $p(d \mid xy) = 2/3?$
 $p(z \mid xy) = 0/3?$
- NO! Absence of xyz might just be bad luck.

6.863J/9.611J SP04 Lecture 6

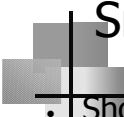


Smoothing

Smoothing deals with events that have been observed zero times

- Smoothing algorithms also tend to improve the accuracy of the model
- Not just unobserved events: what about events observed once?

6.863J/9.611J SP04 Lecture 6



Smoothing the Estimates

- Should we conclude
 - $p(a | xy) = 1/3?$ *reduce this*
 - $p(d | xy) = 2/3?$ *reduce this*
 - $p(z | xy) = 0/3?$ *increase this*
- Discount the positive counts somewhat
- Reallocate that probability to the zeroes
- Especially if the denominator is small ...
 - 1/3 probably too high, 100/300 probably about right
- Especially if numerator is small ...
 - 1/300 probably too high, 100/300 probably about right

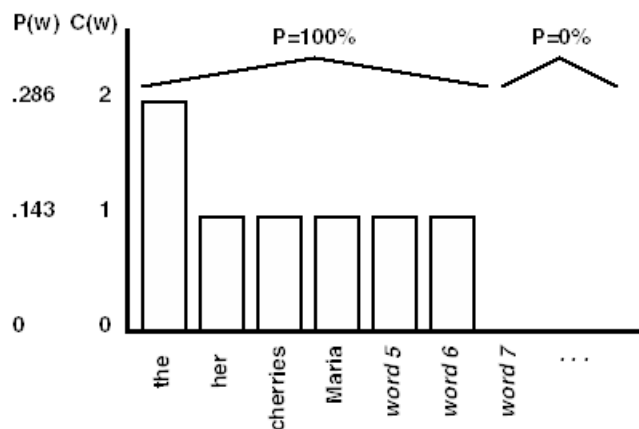
6.863J/9.611J SP04 Lecture 6

Add-one smoothing

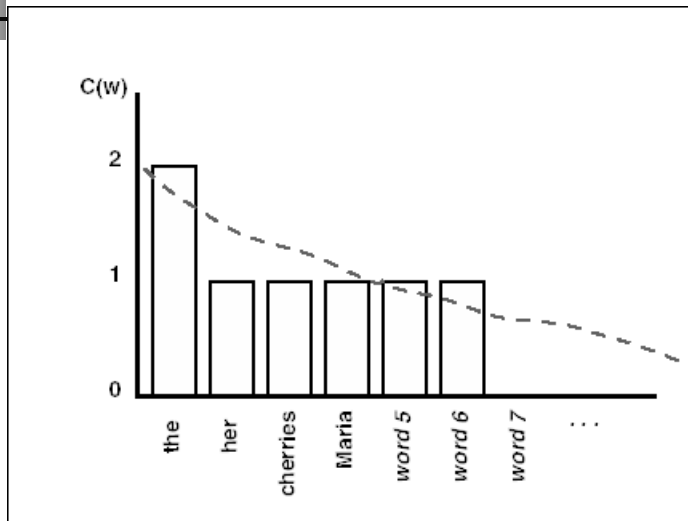
- Let V be the number of words in our vocabulary
- Remember that we observe only V many bigrams
- Assigns count of 1 to unseen bigrams

6.863J/9.611J SP04 Lecture 6

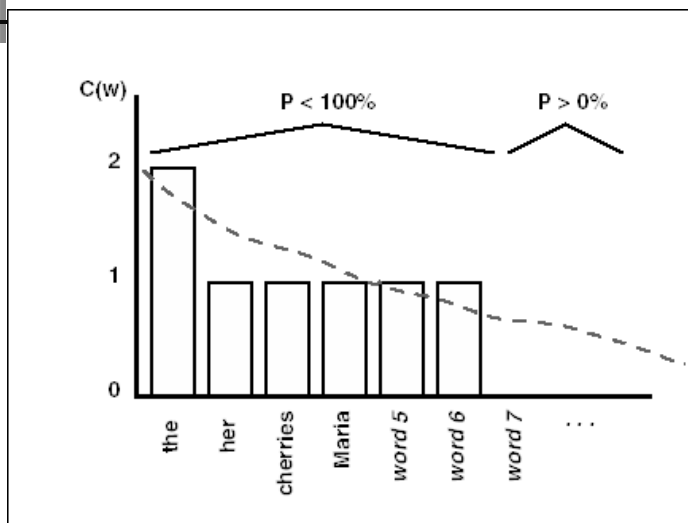
Maximum likelihood estimate



Actual probability distribution



Comparison



Add-One Smoothing

xya	1	1/3	2	2/29
xyb	0	0/3	1	1/29
xyc	0	0/3	1	1/29
xyd	2	2/3	3	3/29
xye	0	0/3	1	1/29
...				
xyz	0	0/3	1	1/29
Total xy	3	3/3	29	29/29

6.863J/9.611J SP04 Lecture 6

Add-one smoothing

6.863J/9.611J SP04 Lecture 6

Example: Bush reads a book

- $P(\text{Bush reads a book})$
- Without smoothing:
- With add-one smoothing (assuming $c(\text{Bush})=1$ but $c(\text{Bush, read}) = 0$)

6.863J/9.611J SP04 Lecture 6

Add-One Smoothing

300 observations instead of 3 – better data, less smoothing

xya	100	100/300	101	101/326
xyb	0	0/300	1	1/326
xyz	0	0/300	1	1/326
xyd	200	200/300	201	201/326
xye	0	0/300	1	1/326
...				
xyz	0	0/300	1	1/326
Total xy	300	300/300	326	326/326

6.863J/9.611J SP04 Lecture 6

Add-One Smoothing

Suppose we're considering 20000 word types, not 26 letters

xya	1	1/3	2	2/29
xyb	0	0/3	1	1/29
xyc	0	0/3	1	1/29
xyd	2	2/3	3	3/29
xye	0	0/3	1	1/29
...				
xyz	0	0/3	1	1/29
Total xy	3	3/3	29	29/29

6.863J/9.611J SP04 Lecture 6

Add-One Smoothing

As we see more word types, smoothed estimates keep falling

see the abacus	1	1/3	2	2/20003
see the abbot	0	0/3	1	1/20003
see the abduct	0	0/3	1	1/20003
see the above	2	2/3	3	3/20003
see the Abram	0	0/3	1	1/20003
...				
see the zygote	0	0/3	1	1/20003
Total	3	3/3	20003	20003/20003

6.863J/9.611J SP04 Lecture 6



Problems...too many mouths to feed

- Suppose we're dealing with a vocab of 20000 words
- As we get more and more training data, we see more and more words that need probability – the probabilities of existing words keep dropping, instead of converging
- This can't be right – eventually they drop too low

6.863J/9.611J SP04 Lecture 6



Good-Turing smoothing

- Add-1 works horribly in practice – adding 1 seems too large
- So...imagine you're sitting at a sushi bar with a conveyor belt
- How likely are you to see a new kind of seafood appear?

6.863J/9.611J SP04 Lecture 6

The sushi bar

10 tuna, 3 unagi, 2 salmon, 1 shrimp, 1 octopus, 1 yellowtail

How likely are you to see another salmon? $< 2/18$

6.863J/9.611J SP04 Lecture 6

Good-Turing smoothing

- How likely are you to see a new type of seafood?
- How many types of seafood (submarines, words) were seen only once? Use this to predict probabilities for unseen events
- Let n_1 be the # of events that occurred once, then the initial est. of this is, $p_0 = n_1/N$
- Let n_2 be the # of events that occurred twice

6.863J/9.611J SP04 Lecture 6

Good-Turing smoothing

- 10 tuna, 3 unagi, 2 salmon, 1 shrimp, 1 octopus, 1 yellowtail
-

- $p_0 = n_1/N = 3/18$
- Now how likely is *octopus*?
- Good-Turing estimate: for any n-gram that occurs r times, we pretend it occurs r^* times,

6.863J/9.611J SP04 Lecture 6

At the sushi bar

- 10 tuna, 3 unagi, 2 salmon, 1 shrimp, 1 octopus, 1 yellowtail
- Octopus occurs 1 time, $r=1$ so we adjust to 1^*
- We need n_1 # of things that occur once = 3
- We need n_2 # of things that occur twice = 1
- Then

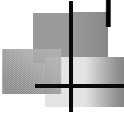
6.863J/9.611J SP04 Lecture 6



Is that the final word?

- No – what happens if something is not seen at all?
- Then you must backoff to bigrams, unigrams...
- Many other new smoothing methods available (see book) – various weighting/discounting schemes (we shall revisit: EM algorithm)
- Are we done? Can we use trigrams now?
- Not quite...

6.863J/9.611J SP04 Lecture 6



Time must have a stop

- Note that in the trigram model the length of the sentence n is variable
- But then, what is the event space for calculating probabilities?
- Suppose our alphabet is just a, b , and the language is all strings over this, eg, $\epsilon, a, b, aa, bb, ab, \dots$
- Assume unigram model, $P(a)=P(b) = 0.5$
- Then $P(aa) = 0.5^2 = 0.25 = P(bb)$, etc...
- But then, adding, $P(a)+P(b)+P(aa)+P(bb) = 1.5$????

What went wrong???

6.863J/9.611J SP04 Lecture 6



Must have a stop

- No probability for $P(\varepsilon)$
- Add special stop symbol:
 - $P(a)=P(b)= 0.25$; $P(\text{stop})= 0.5$
- Now it works out: $P(a \text{ stop}) = P(b \text{ stop}) = 0.25 \times 0.5 = 0.125$; $P(aa \text{ stop}) = 0.25^2 = .03125$, etc.
- Exercise: show the P sum is now 1.

6.863J/9.611J SP04 Lecture 6



OK, back to sushi

- Tagging for the Trout:

6.863J/9.611J SP04 Lecture 6

Tagging for this..

```
The/DT lyric/JJ beauties/NNS of/IN
Schubert/NNP 's/POS Trout/NNP Quintet/NNP
--/:
its/PRP$ elemental/JJ rhythms/NNS
and/CC infectious/JJ melodies/NNS
--/: make/VBP it/PRP
a/DT source/NN of/IN pure/JJ pleasure/NN
for/IN almost/RB all/DT music/NN listeners/NNS ./.
```

6.863J/9.611J SP04 Lecture 6

(Next step: bracketing...)

```
[The/DT lyric/JJ beauties/NNS ]
of/IN
[ Schubert/NNP 's/POS Trout/NNP Quintet/NNP ]
--/:
[ its/PRP$ elemental/JJ rhythms/NNS ]
and/CC [ infectious/JJ melodies/NNS ]
--/: make/VBP [ it/PRP ]
[ a/DT source/NN ] of/IN [ pure/JJ pleasure/NN ]
for/IN almost/RB [ all/DT music/NN listeners/NNS ] ./.
```

6.863J/9.611J SP04 Lecture 6



Tagging methods

- Statistical Tagger T3
- Error-driven Transformation-based Tagger TBT
- Maximum Entropy Tagger MET
- Example-based tagger ET

6.863J/9.611J SP04 Lecture 6



Why tagging? Beyond part of speech

- Syntactic word class
- Word sense
- Attachment
- Shallow parsing
- Sentence boundary detection

6.863J/9.611J SP04 Lecture 6

Two approaches

1. ~~Statistical model~~
2. Deterministic baseline tagger composed with a cascade of fixup transducers

These two approaches are the guts of Lab 2 (lots of others methods: decision trees, ...)

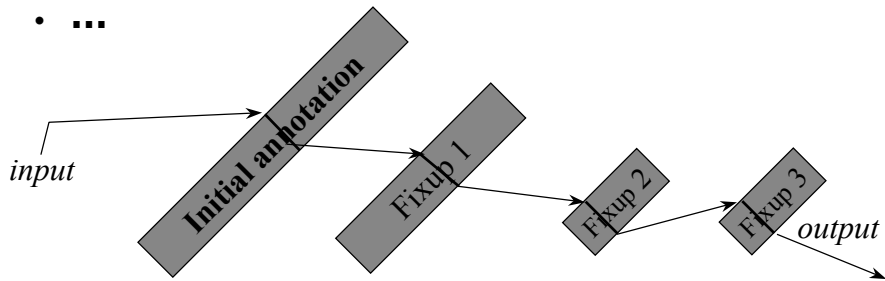
Statistical Model

- $\text{Prob}(\text{Tag sequence, word sequence})$ – based on n-grams: train on marked up text
- We shall see how to do this in detail in a moment

Another FST Paradigm: Successive Fixups

• Like successive markups but *alter*

- Morphology
- Phonology
- Part-of-speech tagging
- ...



6.863J/9.611J SP04 Lecture 6

Evaluation Criteria

- Precision/recall
- Accuracy/ambiguity

6.863J/9.611J SP04 Lecture 6

An exemplar for the divide: "tagging" text

- Input: the lead paint is unsafe
Output: the/Det lead/N paint/N is/V unsafe/Adj
- Can be challenging:
I know that
I know that block
I know that blocks the sun
- new words (OOV= out of vocabulary); words can be whole phrases ("I can't believe it's not butter")

6.863J/9.611J SP04 Lecture 6

What are tags?

- Bridge from words to parsing – but not quite the morphemic details that Kimmo provides (but see next slide)
- Idea is more divide-and-conquer – and depends on task
- "Shallow" analysis for "shallow parsing"

6.863J/9.611J SP04 Lecture 6

More sophisticated – use features

- Word form: $A^+ \rightarrow 2^{(L,C1,C2,\dots,Cn)} \rightarrow T$
 - He always books the violin concert tickets early.
 - books $\rightarrow \{(book-1,Noun,Pl,-,-),(book-2,Verb,Sg,Pres,3)\}$
 - tagging (disambiguation): ... $\rightarrow (Verb,Sg,Pres,3)$
 - ...was pretty good. However, she did not realize...
 - However $\rightarrow \{(however-1,Conj/coord,-,-,-),$
(however- 2,Adv,-,-,-)\}
 - tagging: ... $\rightarrow (Conj/coord,-,-,-)$

6.863J/9.611J SP04 Lecture 6

Two approaches

1. Noisy Channel Model (statistical) – what's that?? (we will have to learn some statistics)
 2. Deterministic baseline tagger composed with a cascade of fixup transducers
- These two approaches will be the guts of Lab 2 (lots of others: decision trees, ...)

6.863J/9.611J SP04 Lecture 6

Example tagsets

- 87 tags - Brown corpus
 - Three most commonly used:
 1. Small: 45 Tags - Penn treebank (Medium size: 61 tags, British national corpus)
 2. Large: 146 tags
- Big question: have we thrown out the right info? Impoverished? How?

6.863J/9.611J SP04 Lecture 6

Brown/Upenn corpus tags

J. text,
p. 297
Fig 8.6
1M words
60K tag
counts

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>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>
PP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, ({ , <</i>
PP\$	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>			

Current (computer/human) performance

Input: the lead paint is unsafe

Output: the/Det lead/N paint/N is/V unsafe/Adj

- How many tags are correct?
 - About 97% currently
 - But baseline is already 90%:
 - Baseline is Homer Simpson algorithm:
 - Tag every word with its most frequent tag (Unigram frequency)
 - Tag unknown words as nouns
- How well do people do?

6.863J/9.611J SP04 Lecture 6

Ok, what should we look at?

correct tags

PN Verb Det Noun Prep Noun Prep Det Noun
Bill directed a cortege of autos through the dunes

PN (Adj) Det Noun Prep Noun Prep Det Noun
Verb Verb Noun Verb

Adj

some possible tags for

Prep

each word (maybe more)

...?

Each unknown tag is constrained by its word and by the tags to its immediate left and right.

But those tags are unknown too ...

6.863J/9.611J SP04 Lecture 6

Ok, what should we look at?

correct tags

PN Verb Det Noun Prep Noun Prep Det Noun
Bill directed a cortege of autos through the dunes

PN ~~Adj~~ Det Noun Prep Noun Prep Det Noun
Verb ~~Verb~~ Noun Verb

Adj

some possible tags for

Prep

each word (maybe more)

...?

Each unknown tag is constrained by its word
and by the tags to its immediate left and right.
But those tags are unknown too ...

6.863J/9.611J SP04 Lecture 6

Ok, what should we look at?

correct tags

PN Verb Det Noun Prep Noun Prep Det Noun
Bill directed a cortege of autos through the dune

PN ~~Adj~~ Det Noun Prep Noun Prep Det Noun
Verb ~~Verb~~ Noun Verb

Adj

some possible tags for

Prep

each word (maybe more)

...?

Each unknown tag is constrained by its word
and by the tags to its immediate left and right.
But those tags are unknown too ...

6.863J/9.611J SP04 Lecture 6

Ok, what *should* we look at?

correct tags

PN Verb Det Noun Prep Noun Prep Det Noun
Bill directed a cortege of autos through the dunes

PN (Adj) Det Noun Prep Noun Prep Det Noun
Verb Verb Noun Verb

Adj

some possible tags for

Prep

each word (maybe more)

...?

Each unknown tag is constrained by its word
and by the tags to its immediate left and right.
But those tags are unknown too ...

6.863J/9.611J SP04 Lecture 6

We can use n-grams for tagging

- Replace 'words' with 'tags'
- Find best maximum likelihood estimates
- Estimates calculated this way:
 - $P(\text{noun}|\text{det}) = p(\text{det}, \text{noun})/p(\text{det})$ replace:
 - $\approx \frac{\text{count}(\text{det at position } i-1 \ \& \ \text{noun at } i)}{\text{count}(\text{det at position } i-1)}$
 - Correction: include frequency of context word
 $\approx \frac{\text{count}(\text{det at position } i-1 \ \& \ \text{noun at } i)}{\text{count}(\text{det at position } i-1) * \text{count}(\text{noun at } i)}$

Find optimal path – highest p , using dynamic programming algorithm, approx. linear in length of tag sequence

6.863J/9.611J SP04 Lecture 6

Example

The guy still saw her
Det NN NN NN PPO
VB VB VBD PP\$
RB

Table 2 from DeRose (1988)

Det=determiner, NN=noun, VB=verb, RB=adverb,
VBD=past-tense-verb, PPO=object pronoun and
PP\$=possessive pronoun

Find the Max likelihood estimate (MLE) path through
this 'trellis'

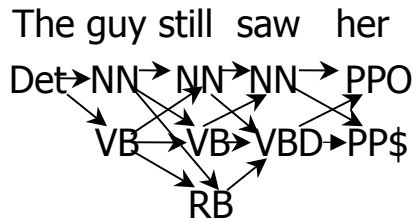
6.863J/9.611J SP04 Lecture 6

Transitional probability estimates from counts

	DT	NN	PPO	PP\$	RB	VB	VBD
DT	0	186	0	8	1	8	9
NN	40	1	3	40	9	66	186
PPO	7	3	16	164	109	16	313
PP\$	176	0	0	5	1	1	2
RB	5	3	16	164	109	16	313
VB	22	694	146	98	9	1	59
VBD	11	584	143	160	2	1	91

6.863J/9.611J SP04 Lecture 6

Tagging search tree (trellis)



Step 1. $c(\text{DT-NN}) = 186$

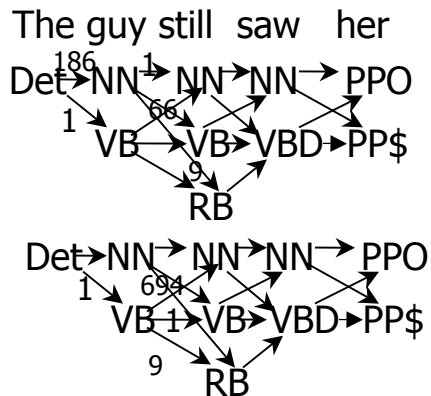
$c(\text{DT-VB}) = 1$

Keep both paths. (Why?)

Step 2. Pick *max* to each of the tags NN, VB, RB
 need keep only the *max*. Why?

6.863J/9.611J SP04 Lecture 6

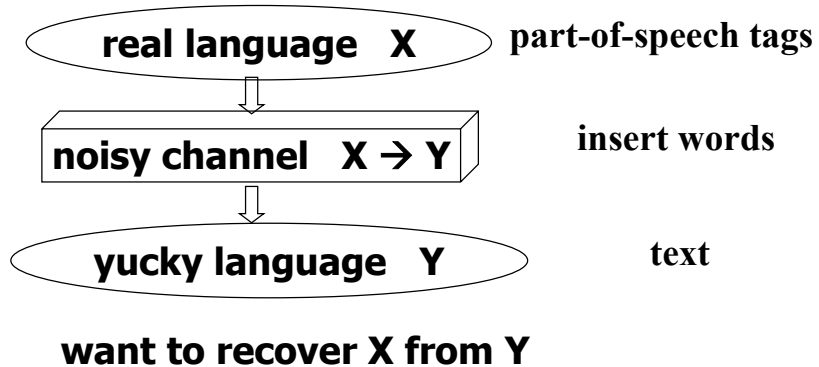
Trellis search



6.863J/9.611J SP04 Lecture 6

Finite-state approaches

- Noisy Channel Muddle (statistical)



6.863J/9.611J SP04 Lecture 6

So far, then...

- n-gram models are a.k.a. Markov models/chains/processes.
- They are a model of how a sequence of observations comes into existence.
- The model is a probabilistic walk on a FSA.
- $Pr(a/b)$ = probability of entering state a , given that we're currently in state b .

6.863J/9.611J SP04 Lecture 6

How well does this work for tagging?

- 90% accuracy (for unigram) pushed up to 96%
- So what?
- How *good* is this? Evaluation!

6.863J/9.611J SP04 Lecture 6

Evaluation of systems

- The principal measures for information extraction tasks are recall and precision.
- Recall is the number of answers the system got right divided by the number of possible right answers
 - It measures how complete or comprehensive the system is in its extraction of relevant information
- Precision is the number of answers the system got right divided by the number of answers the system gave
 - It measures the system's correctness or accuracy
 - Example: there are 100 possible answers and the system gives 80 answers and gets 60 of them right, its recall is 60% and its precision is 75%.

6.863J/9.611J SP04 Lecture 6

A better measure - Kappa

- Takes baseline & complexity of task into account – if 99% of tags are Nouns, getting 99% correct no great shakes
- Suppose no “Gold Standard” to compare against?
- $P(A)$ = proportion of times hypothesis *agrees* with standard (% correct)
- $P(E)$ = proportion of times hypothesis and standard would be *expected* to agree by chance (computed from some other knowledge, or actual data)

6.863J/9.611J SP04 Lecture 6

Kappa [p. 315 J&M text]

- Note K ranges between 0 (no agreement, except by chance; to complete agreement, 1)
- Can be used even if no ‘Gold standard’ that everyone agrees on
- $K > 0.8$ is good

6.863J/9.611J SP04 Lecture 6



Kappa

- A = actual agreement; E = expected agreement
- consistency is more important than “truth”
- methods for raising consistency
 - style guides (often have useful insights into data)
 - group by task, not chronologically, etc.
 - annotator acclimatization

6.863J/9.611J SP04 Lecture 6



Statistical Tagging Methods

- Simple bigram – ok, done
- Combine bigram and unigram
- OUR GOAL: maximize $P(T,w)$ where T =a tag sequence (guessed); and w = the observed word sequence – note this is a joint probability
- So, why not use our formula for joint probabilities...

6.863J/9.611J SP04 Lecture 6



Our first try...

- $P(T,w) = P(T) P(w | T)$

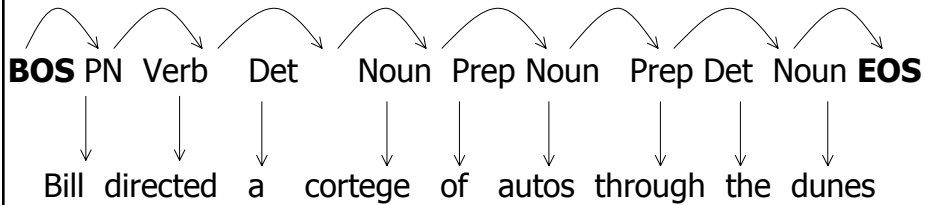


Our plan to find $P(T,w)$

- Find best $P(T)$ – probability of a tag sequence
- How? A: use bigrams
- Find best $P(w|T)$ --????? How?

Punchline – recovering (words, tags)

tags $X \rightarrow$



words $Y \rightarrow$

GOAL: Find tag sequence X that maximizes probability product

6.863J/9.611J SP04 Lecture 6

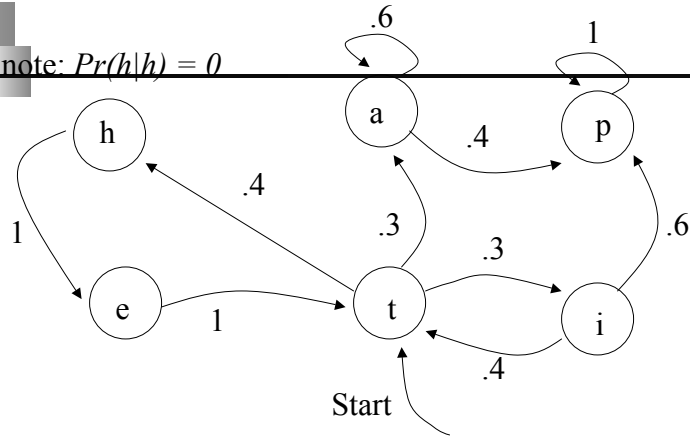
Break apart in 2 stages

- If we were just predicting tags, we could just use bigrams
- We can model this as a Markov process, in particular, an fsa with probabilities on the arcs...

6.863J/9.611J SP04 Lecture 6

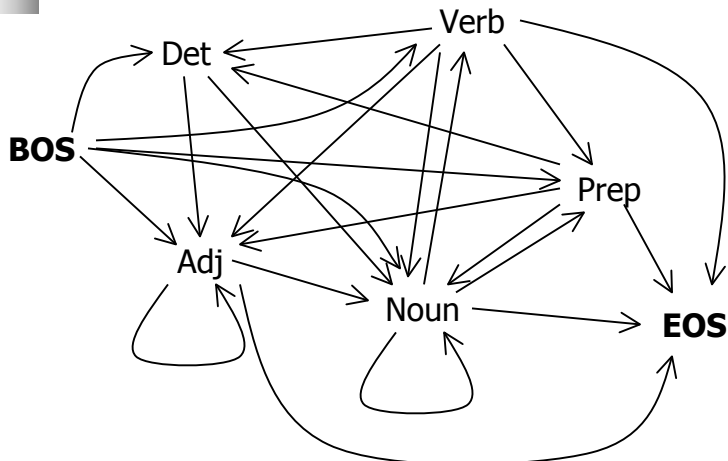
Markov chain...pr of letter sequences

note: $Pr(h|h) = 0$



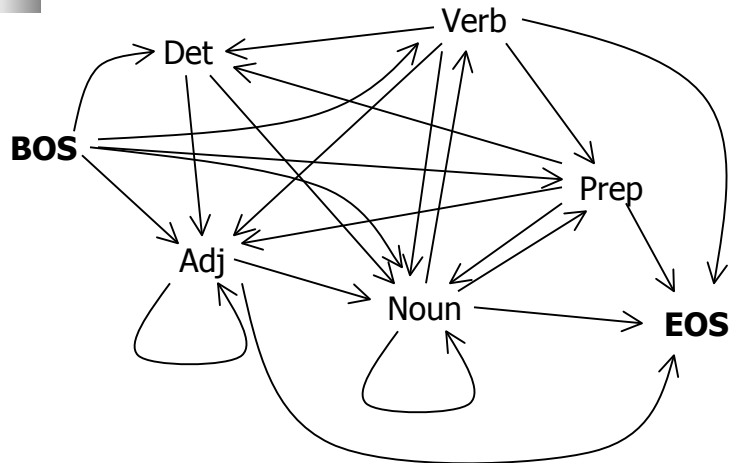
6.863J/9.611J SP04 Lecture 6

First-order Markov (bigram tag) model as fsa



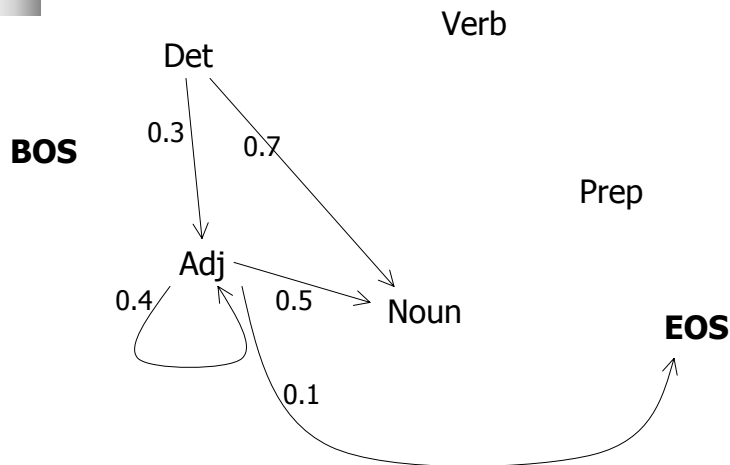
6.863J/9.611J SP04 Lecture 6

First-order Markov (bigram tag) model as fsa



6.863J/9.611J SP04 Lecture 6

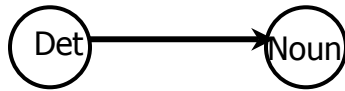
Add in transition probs from training data - sum to 1



6.863J/9.611J SP04 Lecture 6

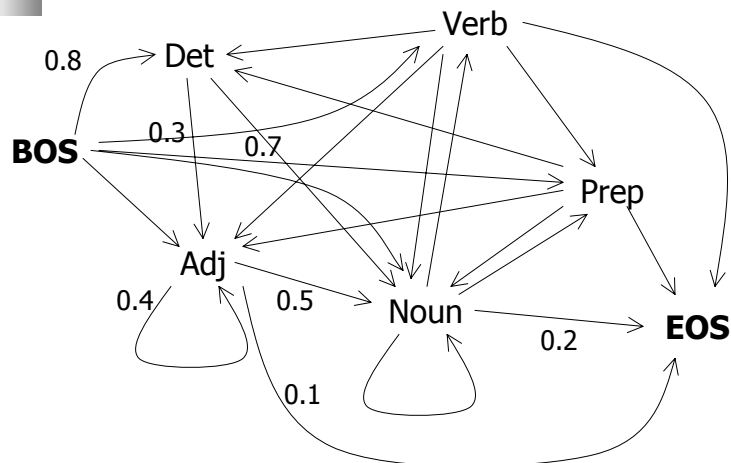
Same as bigram...estimate the same way

$$P(\text{Noun}|\text{Det})=0.7 \equiv$$



6.863J/9.611J SP04 Lecture 6

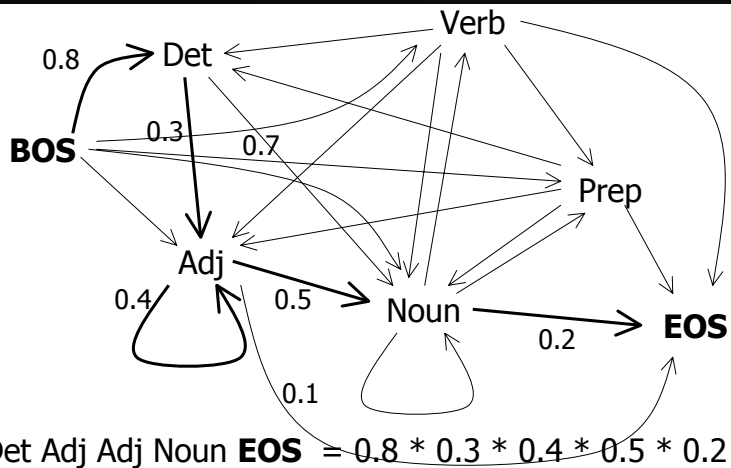
Add in start & etc.



6.863J/9.611J SP04 Lecture 6

Markov Model – bigram tag sequence

$p(\text{tag sequence})$



6.863J/9.611J SP04 Lecture 6

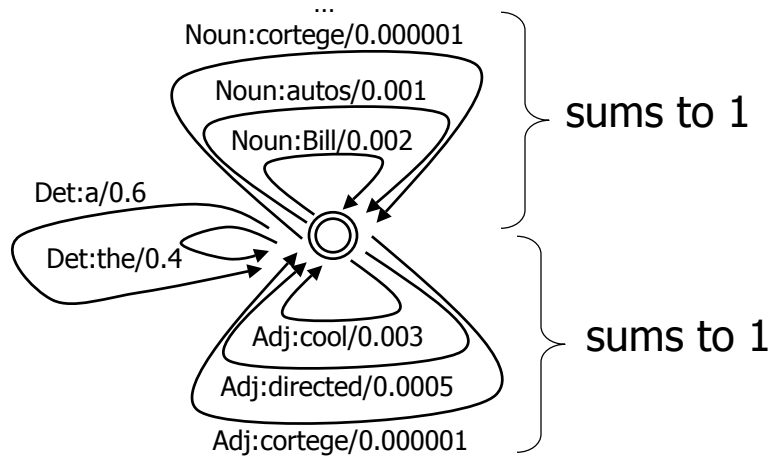
So what?

- We want more!
- We cannot observe the tag sequence –
- But we can estimate $P(\text{words} \mid \text{tags})$
- Also use an fsa – just unigrams

6.863J/9.611J SP04 Lecture 6

Unigram replacement model

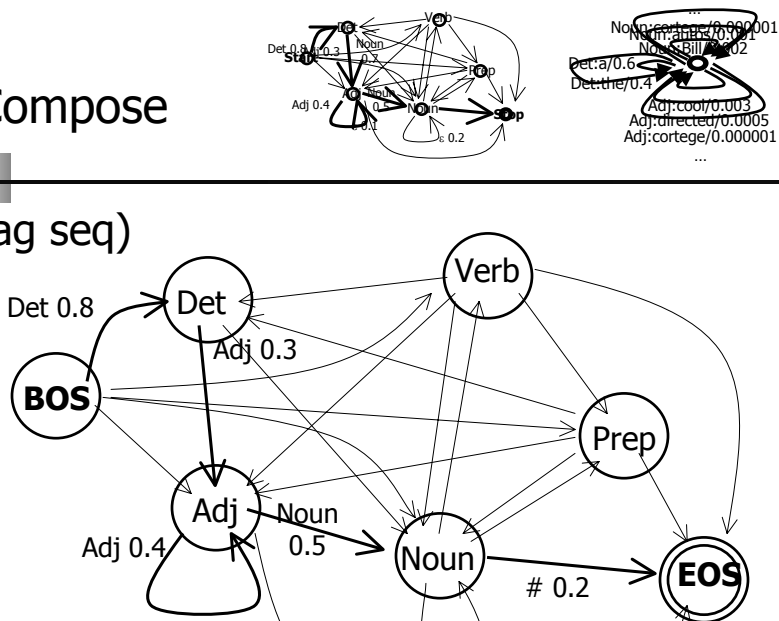
$P(\text{word} | \text{tag})$



6.863/9.611J SP04 Lecture 6

Compose

$p(\text{tag seq})$



6.863/9.611J SP04 Lecture 6

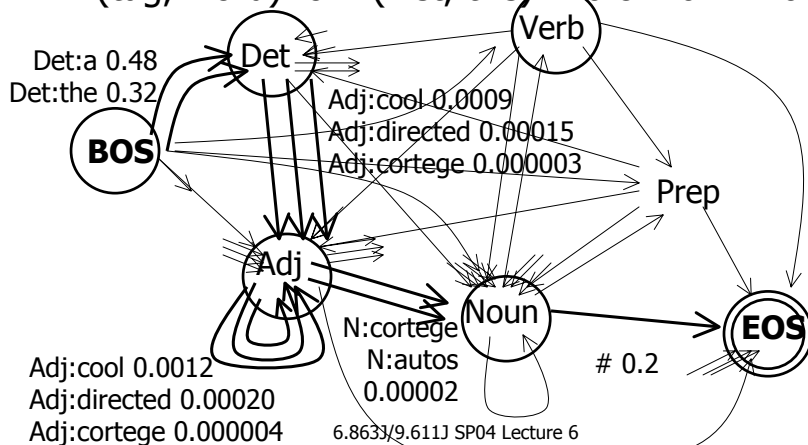
Now we compose (multiply) the nets

- Compose $P(\text{tag sequence})$ with $P(\text{word}|\text{tag})$
- Result: $P(\text{tag}, \text{word})$

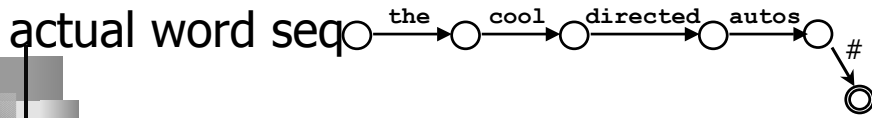
6.863J/9.611J SP04 Lecture 6

Compose

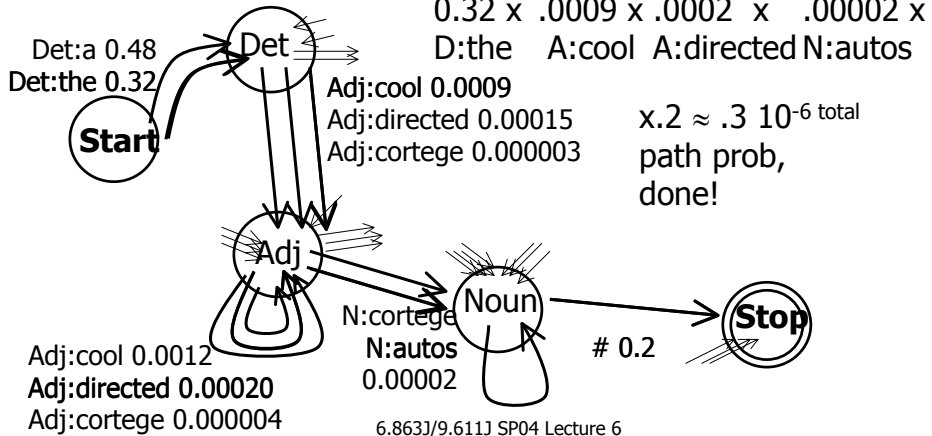
$P(\text{tag seq}) * P(\text{word seq} | \text{tag seq})$ e.g.,
 $P(\text{tag}, \text{word})$ for $P(\text{Det}, \text{the}) = 0.8 \times 0.4 = 0.32$



Compose with



$$p(\text{word seq, tag seq}) = p(\text{tag seq}) * p(\text{word seq} | \text{tag seq})$$



Well, we are almost done!

- The Pr of a sequence is just found by multiplying through as we go from start to stop
- Given the actual words in the sentence, trace through and find the highest value Pr – this will give the most likely tag sequence, word sequence combination
- (What have we wrought?)

This is an Hidden Markov model for tagging

- Each hidden tag state produces a word in the sentence
- Each word is
 - Uncorrelated with all the other words and their tags
 - Probabilistic depending on the N previous tags only

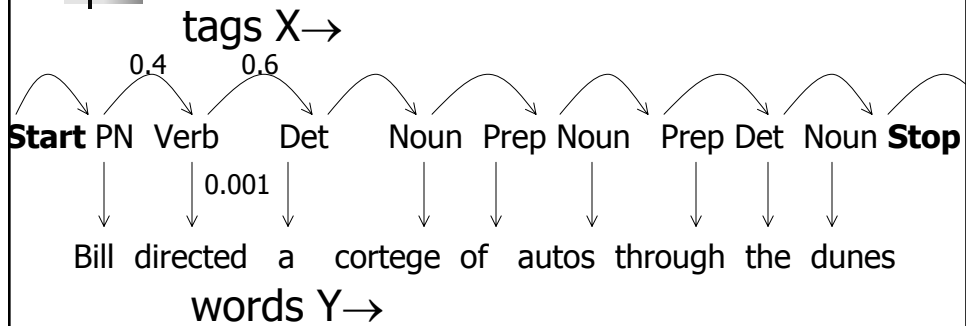
6.863J/9.611J SP04 Lecture 6

The statistical view, in short:

- We are modeling $p(\text{word seq}, \text{tag seq})$
- The tags are *hidden*, but we see the words
- Q: What is the most likely tag sequence?
- Use a finite-state automaton, that can emit the observed words
- FSA has limited memory
- Note that given words, in general, there could be more than 1 underlying state sequence corresponding to the words

6.863J/9.611J SP04 Lecture 6

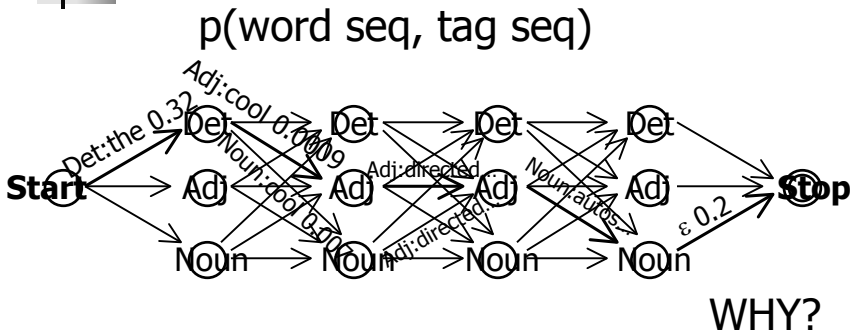
Punchline – ok, where do the pr numbers come from?



The tags are not observable & they are states of some fsa
We estimate transition probabilities between states
We also have 'emission' pr's from states
En tout: a Hidden Markov Model (HMM)

But...how do we find this 'best' path???

Unroll the fsa - All paths together form 'trellis'

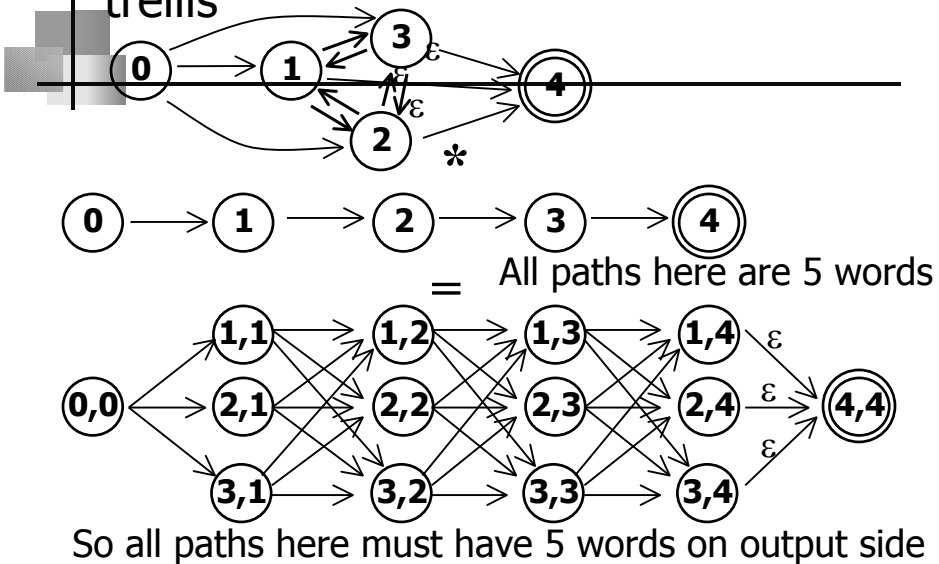


The best path:

BOS Det Adj Adj Noun **EOS** = $0.32 * 0.0009 \dots$
 the cool directed autos

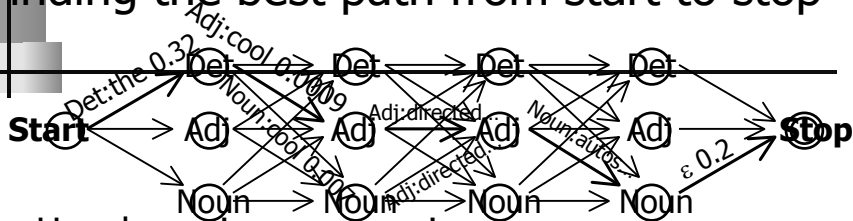
6.863J/9.611J SP04 Lecture 6

Cross-product construction forms trellis



6.863J/9.611J SP04 Lecture 6

Finding the best path from start to stop



- Use dynamic programming
- What is best path from Start to each node?
 - Work from left to right
 - Each node stores its best path from Start (as probability plus one backpointer)
- Special acyclic case of Dijkstra's shortest-path algorithm
- Faster if some arcs/states are absent

6.863J/9.611J SP04 Lecture 6

Method: Viterbi algorithm

- For each path reaching state s at step (word) t , we compute a path probability. We call the max of these viterbi(s,t)
- [Base step] Compute $\text{viterbi}(0,0)=1$
- [Induction step] Compute $\text{viterbi}(s',t+1)$, assuming we know $\text{viterbi}(s,t)$ for all s

6.863J/9.611J SP04 Lecture 6

Viterbi recursion

$$\text{path-prob}(s'|s,t) = \text{viterbi}(s,t) * a[s,s']$$

probability of path to s' through s max path score * transition p
 for state s at time t $s \rightarrow s'$

$$\text{viterbi}(s',t+1) = \max_{s \text{ in STATES}} \text{path-prob}(s' | s,t)$$

Method...

- This is *almost* correct...but again, we need to factor in the *unigram* prob of a state s' given an observed surface word w
- So the correct formula for the path prob is:
 $\text{path-prob}(s'|s,t) = \text{viterbi}(s,t) * a[s,s'] * b_{s'}(o_t)$

bigram unigram

Or as in your text...p. 179

function VITERBI(*observations* of len T , *state-graph*) **returns** *best-path*

$num\text{-}states \leftarrow \text{NUM-OF-STATES}(state\text{-}graph)$

Create a path probability matrix $viterbi[num\text{-}states+2, T+2]$

$viterbi[0,0] \leftarrow 1.0$

for each time step t **from** 0 **to** T **do**

for each state s **from** 0 **to** $num\text{-}states$ **do**

for each transition s' from s specified by *state-graph*

$new\text{-}score \leftarrow viterbi[s, t] * a[s, s'] * b_{s'}(o_t)$

if $((viterbi[s', t+1] = 0) \parallel (new\text{-}score > viterbi[s', t+1]))$

then

$viterbi[s', t+1] \leftarrow new\text{-}score$

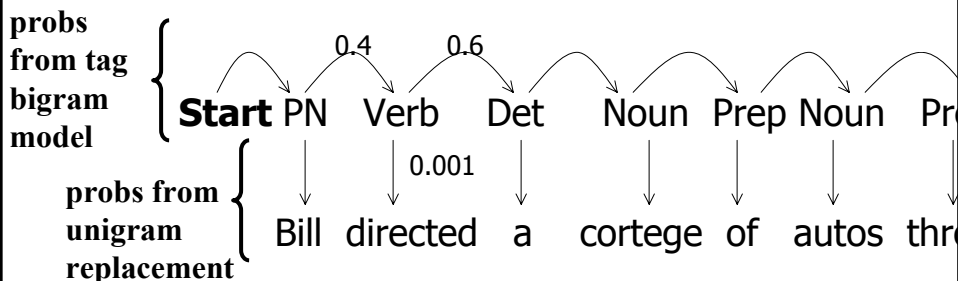
$back\text{-}pointer[s', t+1] \leftarrow s$

Backtrace from highest probability state in the final column of $viterbi[]$ and return path

6.863J/9.611J SP04 Lecture 6

Summary

- We are modeling $p(\text{word seq}, \text{tag seq})$
- The tags are hidden, but we see the words
- Is tag sequence X likely with these words?
- Model is a "Hidden Markov Model":



- Find X that maximizes probability product

6.863J/9.611J SP04 Lecture 6