

## Part-of-Speech Tagging & Parsing

Chloé Kiddon

(slides adapted and/or stolen outright from Andrew McCallum, Christopher Manning, and Julia Hockenmaier)

## Announcements

- We do have a Hadoop cluster!
  - It's offsite. I need to know all groups who want it!
- You all have accounts for MySQL on the cubist machine (cubist.cs.washington.edu)
  - Your folder is /projects/instr/cse454/a-f
- I'll have a better email out this afternoon I hope
- Grading HW1 should be finished by next week.

## Timely warning

- POS tagging and parsing are two large topics in NLP
- Usually covered in 2-4 lectures
- We have an hour and twenty minutes. ☺

## Part-of-speech tagging

- Often want to know what part of speech (POS) or word class (noun,verb,...) should be assigned to words in a piece of text
- **Part-of-speech tagging** assigns POS labels to words

JJ      JJ    NNS   VBP    RB  
 Colorless green ideas sleep furiously.

## Why do we care?

- Parsing (come to later)
- Speech synthesis
  - INsult or inSULT, overFLOW or OVERflow, REAd or reAD
- Information extraction: entities, relations
  - Romeo loves Juliet vs. lost loves found again
- Machine translation

## Penn Treebank Tagset

1. CC	Coordinating conjunction	20. RB	Adverb
2. CD	Cardinal number	21. RBR	Adverb, comparative
3. DT	Determiner	22. RBS	Adverb, superlative
4. EX	Existential there	23. RP	Particle
5. FW	Foreign word	24. SYM	Symbol
6. IN	Preposition or subordinating conjunction	25. TO	to
7. JJ	Adjective	26. UH	Interjection
8. JJR	Adjective, comparative	27. VB	Verb, base form
9. JJS	Adjective, superlative	28. VBD	Verb, past tense
10. LS	List item marker	29. VBG	Verb, gerund or present participle
11. MD	Modal	30. VBN	Verb, past participle
12. NN	Noun, singular or mass	31. VBP	Verb, non-3rd person singular present
13. NNS	Noun, plural	32. VBZ	Verb, 3rd person singular present
14. NP	Proper noun, singular	33. WDT	Wh-determiner
15. NPS	Proper noun, plural	34. WP	Wh-pronoun
16. PDT	Predeterminer	35. WP\$	Possessive wh-pronoun
17. POS	Possessive ending	36. WRB	Wh-adverb
18. PP	Personal pronoun		
19. PP\$	Possessive pronoun		

## Ambiguity

Buffalo buffalo buffalo.

## How many words are ambiguous?

	87-tag Original Brown	45-tag Treebank Brown
<b>Unambiguous (1 tag)</b>	<b>44,019</b>	<b>38,857</b>
<b>Ambiguous (2-7 tags)</b>	<b>5,490</b>	<b>8844</b>
Details:		
2 tags	4,967	6,731
3 tags	411	1,621
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> )

Hockenmaier

## Naïve approach!

- Pick the most common tag for the word

Word	POS listings in Brown		
heat	noun/89	verb/5	
oil	<b>noun/87</b>		
in	<b>prep/20731</b>	noun/1	adv/462
a	<b>det/22943</b>	noun/50	noun-proper/30
large	<b>adj/354</b>	noun/2	adv/5
pot	<b>noun/27</b>		

- 91% success rate!

Andrew McCallum

## We have more information

- We are not just tagging words, we are tagging *sequences* of words

For a sequence of words  $W$ :

$$W = w_1 w_2 w_3 \dots w_n$$

We are looking for a sequence of tags  $T$ :

$$T = t_1 t_2 t_3 \dots t_n$$

where  $P(T|W)$  is maximized

Andrew McCallum

## In an ideal world...

- Find all instances of a sequence in the dataset and pick the most common sequence of tags
  - Count("heat oil in a large pot") = 0 ????
  - Uhh...
- Sparse data problem
- Most sequences will never occur, or will occur too few times for good predictions

## Bayes' Rule

- To find  $P(T|W)$ , use Bayes' Rule:

$$P(T|W) \propto \frac{P(W|T) \times P(T)}{P(W)}$$

- We can maximize  $P(T|W)$  by maximizing  $P(W|T) \times P(T)$

Andrew McCallum

## Finding P(T)

- Generally,

$$P(t_1 t_2 \dots t_n) = P(t_1) \times P(t_2 \dots t_n | t_1)$$

$$P(t_1 t_2 \dots t_n) = P(t_1) \times P(t_2 | t_1) \times P(t_3 \dots t_n | t_1 t_2)$$

$$P(t_1 t_2 \dots t_n) = \prod_i P(t_i | t_1 t_2 \dots t_{i-1})$$

- Usually not feasible to accurately estimate more than tag bigrams (possibly trigrams)

## Markov assumption

- Assume that the probability of a tag only depends on the tag that came **directly before** it

$$P(t_i | t_1 t_2 \dots t_{i-1}) = P(t_i | t_{i-1})$$

- Then,

$$P(t_1 t_2 \dots t_n) = P(t_1) \times P(t_2 | t_1) \times \prod_i P(t_i | t_{i-1}) \times \dots \times P(t_n | t_{n-1})$$

- Only need to count tag bigrams.

## Putting it all together

- We can similarly assume

$$P(w_i | t_1 \dots t_n) = P(w_i | t_i)$$

- So:

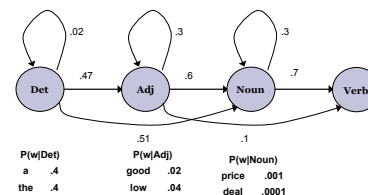
$$P(w_1 \dots w_n | t_1 \dots t_n) = P(w_1 | t_1) \times P(w_2 | t_2) \times \dots \times P(w_n | t_n)$$

- And the final equation becomes:

$$P(W | T) \times P(T) = P(w_1 | t_1) \times P(w_2 | t_2) \times \dots \times P(w_n | t_n) \times P(t_1) \times P(t_2 | t_1) \times P(t_3 | t_2) \times \dots \times P(t_n | t_{n-1})$$

## Process as an HMM

- Start in an initial state  $t_0$  with probability  $\pi(t_0)$
- Move from state  $t_i$  to  $t_j$  with transition probability  $a(t_j | t_i)$
- In state  $t_i$ , emit symbol  $w_k$  with emission probability  $b(w_k | t_i)$



## Three Questions for HMMs

1. **Evaluation** - Given a sequence of words  $\mathbf{W} = w_1 w_2 w_3 \dots w_n$  and an HMM model  $\Theta$ , what is  $P(\mathbf{W}|\Theta)$
2. **Decoding** - Given a sequence of words  $\mathbf{W}$  and an HMM model  $\Theta$ , find the most probable parse  $\mathbf{T} = t_1 t_2 t_3 \dots t_n$
3. **Learning** - Given a tagged (or untagged) dataset, find the HMM  $\Theta$  that maximizes the data

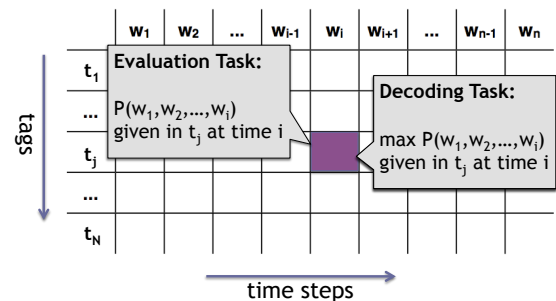
## Three Questions for HMMs

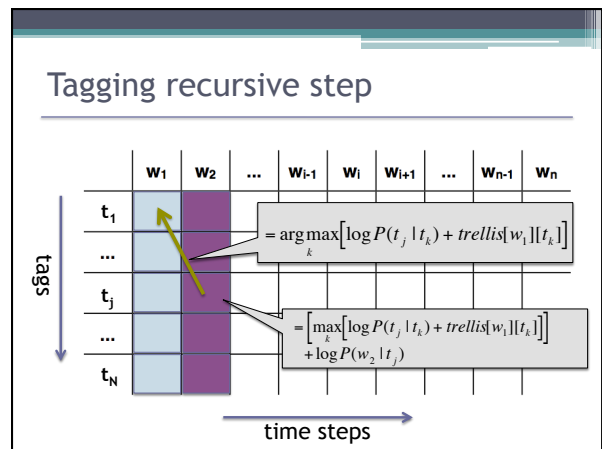
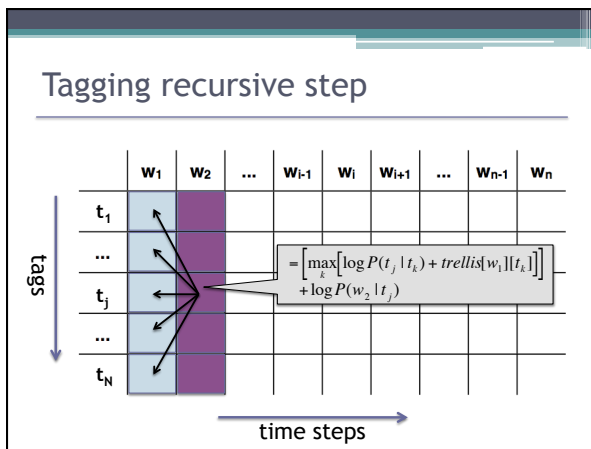
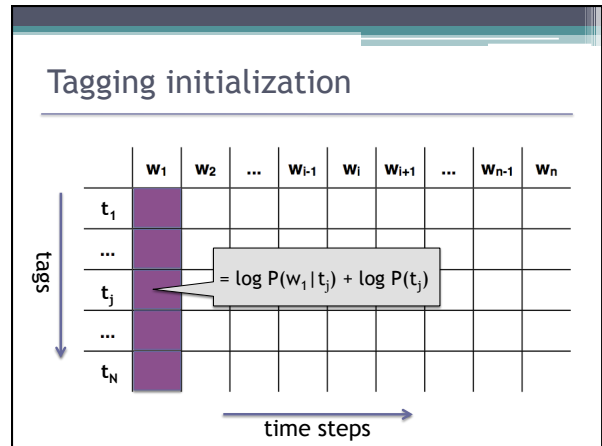
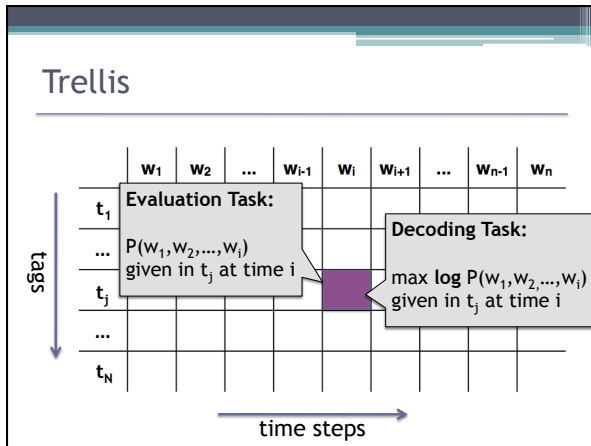
1. **Evaluation** - Given a sequence of words  $\mathbf{W} = w_1 w_2 w_3 \dots w_n$  and an HMM model  $\Theta$ , what is  $P(\mathbf{W}|\Theta)$
2. **Tagging** - Given a sequence of words  $\mathbf{W}$  and an HMM model  $\Theta$ , find the most probable parse  $\mathbf{T} = t_1 t_2 t_3 \dots t_n$
3. **Learning** - Given a tagged (or untagged) dataset, find the HMM  $\Theta$  that maximizes the data

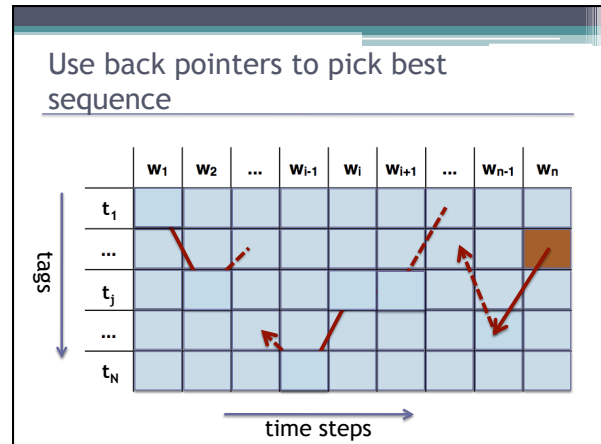
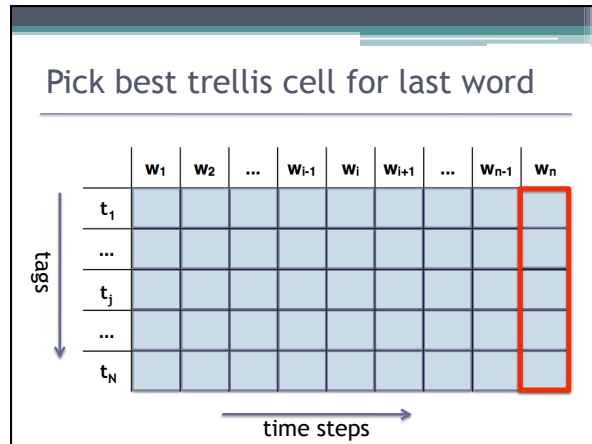
## Tagging

- Need to find the **most likely** tag sequence given a sequence of words
  - maximizes  $P(\mathbf{W}|\mathbf{T}) * P(\mathbf{T})$  and thus  $P(\mathbf{T}|\mathbf{W})$
- Use Viterbi!

## Trellis







### Learning a POS-tagging HMM

- Estimate the parameters in the model using counts
 
$$P(t_i | t_{i-1}) \rightarrow$$

$$P(w_i | t_i) \rightarrow$$
- With smoothing, this model can get 95-96% correct tagging

### Problem with supervised learning

- Requires a large hand-labeled corpus
  - Doesn't scale to new languages
  - Expensive to produce
  - Doesn't scale to new domains
- Instead, apply unsupervised learning with Expectation Maximization (EM)
  - Expectation step: calculate probability of all sequences using set of parameters
  - Maximization step: re-estimate parameters using results from E-step

## Lots of other techniques!

- Trigram models (more common)
- Text normalization
- Error-based transformation learning (“Brill learning”)
  - Rule-based system
    - Calculate initial states: proper noun detection, tagged corpus
    - Acquire transformation rules
      - Change VB to NN when prev word was adjective
      - The long race finally ended
- Minimally supervised learning
  - Unlabeled data but have a dictionary

## Seems like POS-tagging is solved

- Penn Treebank POS-tagging accuracy = human ceiling
  - Human agreement 97%
- In other languages, not so much

## So now we are HMM Masters

- We can use HMMs to...
  - Tag words in a sentence with their parts of speech
  - Extract entities and other information from a sentence
- Can we use them to determine syntactic categories?

## Syntax

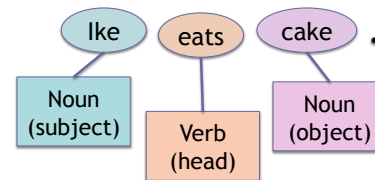
- Refers to the study of the way words are arranged together, and the relationship between them.
- Prescriptive vs. Descriptive
- Goal of syntax is to model the knowledge of that people unconsciously have about the grammar of their native language
- Parsing extracts the syntax from a sentence



## Parsing applications

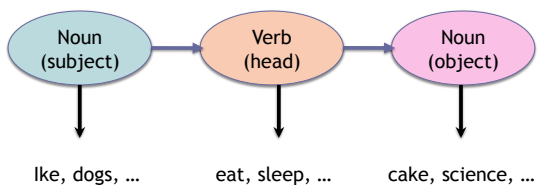
- High-precision Question-Answering systems
- Named Entity Recognition (NER) and information extraction
- Opinion extraction in product reviews
- Improved interaction during computer applications/games

## Basic English sentence structure



Hockenmaier

## Can we build an HMM?



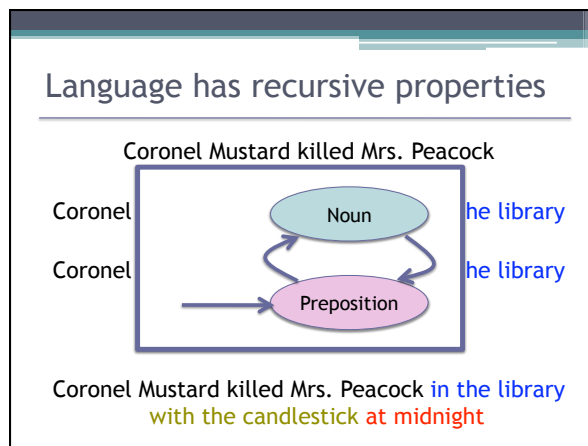
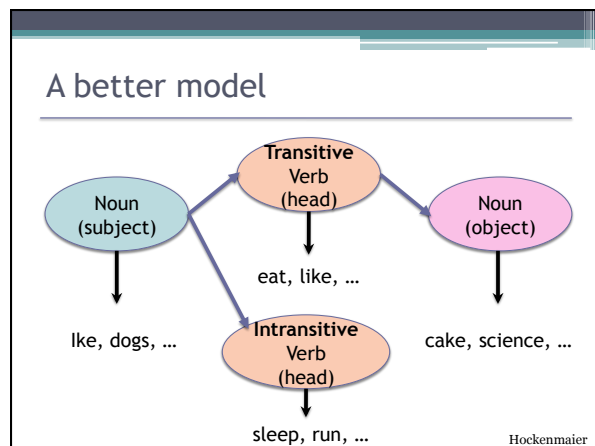
Hockenmaier

## Words take arguments

I eat cake. ☺  
 I sleep cake. ☹  
 I give you cake. ☺  
 I give cake. Hmm...  
 I eat you cake???

- Subcategorization
  - **Intransitive verbs:** take only a subject
  - **Transitive verbs:** take a subject and an object
  - **Ditransitive verbs:** take a subject, object, and indirect object
- Selectional preferences
  - The object of *eat* should be edible

Hockenmaier



### HMMs can't generate hierarchical structure

Coronel Mustard killed Mrs. Peacock **in the library**  
**with the candlestick at midnight.**

- Does Mustard have the candlestick?
- Or is the candlestick just sitting in the library?
- Memoryless
  - Can't make long range decisions about attachments
- Need a better model

### Words work in groups

- **Constituents** - words or groupings of words that function as single units
  - Noun phrases (NPs)
    - The computer science class
    - Peter, Paul, and Mary
    - PAC10 Schools, such as UW,
    - He
    - The reason I was late

## Words work in groups

- **Constituents** - words or groupings of words that function as single units

- **Noun phrases (NPs)**

- *The computer science class listened*

**NPs can appear before a verb.**

- PAC10 Schools, such as UW, dominate ...
- He juggled ...
- The reason I was late was ...
- \*the listened
- \*such sing
- \*late was

## Many different constituents

- |   |   |
|---|---|
| 1. S - simple declarative clause  | mark the head of the NP. Corresponds very roughly to N-bar level but used quite differently.                                  |
| 2. SBAR - Clause introduced by a (possibly empty) subordinating conjunction                       | 15. PP - Prepositional Phrase.  |
| 3. SBARQ - Direct question introduced by a wh-word or a wh-phrase                                 | 16. PRN - Parenthetical.  |
| 4. SINV - Inverted declarative sentence   | 17. PRT - Particle.   |
| 5. SQ - Inverted yes/no question, or main clause of a wh-question                                 | 18. QP - Quantifier Phrase (i.e. complex measure/amount phrase); used within NP.  |
| 6. ADJP - Adjective Phrase.   | 19. RRC - Reduced Relative Clause.  |
| 7. ADVP - Adverb Phrase.  | 20. UCP - Unlike Coordinated Phrase.  |
| 8. CONJP - Conjunction Phrase.  | 21. VP - Verb Phrase.   |
| 9. FRAG - Fragment.   | 22. WHADJP - Wh-adjective Phrase.   |
| 10. INTJ - Interjection.  | 23. WHAVP - Wh-adverb Phrase.   |
| 11. LST - List marker.  | 24. WHNP - Wh-noun Phrase.  |
| 12. NAC - Not a Constituent; used to show the scope of certain pronominal modifiers within an NP. | 25. WHPP - Wh-prepositional Phrase.   |
| 13. NP - Noun Phrase.   | 26. X - Unknown, uncertain, or unbracketable. X is often used for bracketing typos and in bracketing the...the-constructions. |
| 14. NX - Used within certain complex NPs to   |   |

## Many different constituents

1. CC	Coordinating conjunction	20. RB	Adverb
2. CD	Cardinal number	21. RBR	Adverb, comparative
3. DT	Determiner	22. RBS	Adverb, superlative
4. EX	Existential there	23. RP	Particle
5. FW	Foreign word	24. SYM	Symbol
6. IN	Preposition or subordinating conjunction	25. TO	to
7. JJ	Adjective	26. UH	Interjection
8. JJR	Adjective, comparative	27. VB	Verb, base form
9. JJS	Adjective, superlative	28. VBD	Verb, past tense
10. LS	List item marker	29. VBG	Verb, gerund or present participle
11. MD	Modal	30. VBN	Verb, past participle
12. NN	Noun, singular or mass	31. VBP	Verb, non-3rd person singular present
13. NNS	Noun, plural	32. VBZ	Verb, 3rd person singular present
14. NP	Proper noun, singular	33. WDT	Wh-determiner
15. NPS	Proper noun, plural	34. WP	Wh-pronoun
16. PDT	Predeterminer	35. WPS	Possessive wh-pronoun
17. POS	Possessive ending	36. WRB	Wh-adverb
18. PP	Personal pronoun		
19. PPS	Possessive pronoun		

## Attachment ambiguities

- Teacher Strikes Idle Kids
- Squad Helps Dog Bite Victim
- Complaints About NBA Referees Getting Ugly
- Soviet Virgin Lands Short of Goal Again
- Milk Drinkers are Turning to Powder

## Attachment ambiguities

- The key parsing decision: How do we ‘attach’ various kinds of constituents - PPs, adverbial or participial phrases, coordinations, etc.
- Prepositional phrase attachment
  - *I saw the man with the telescope.*
- What does with a *telescope* modify?
  - The verb *saw*?
  - The noun *man*?
- Very hard problem. AI Complete.

## Parsing

- We want to run a grammar backwards to find possible structures for a sentence
- Parsing can be viewed as a **search problem**
- Parsing is a **hidden data problem**

## Context-free grammars (CFGs)

- Specifies a set of tree structures that capture constituency and ordering in language
  - A noun phrase can come before a verb phrase

•  $S \rightarrow NP VP$



## Phrase structure grammars = Context-free grammars

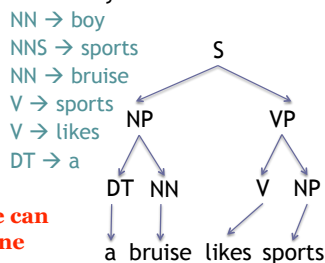
- $G = (T, N, S, R)$ 
  - $T$  is the set of terminals (i.e. words)
  - $N$  is the set of non-terminals
    - Usually separate the set  $P$  of preterminals (POS tags) from the rest of the non-terminals
  - $S$  is the start symbol
  - $R$  is the set of rules/productions of the form  $X \rightarrow \gamma$  where  $X$  is a nonterminal and  $\gamma$  is a sequence of terminals and nonterminals (possibly empty)
- A grammar  $G$  generates a language  $L$

Manning

## A phrase structure grammar

- By convention, S is the start symbol

- $S \rightarrow NP VP$
- $NP \rightarrow DT NN$
- $NP \rightarrow NNS$
- $VP \rightarrow V NP$
- $VP \rightarrow V$
- ...



**But since a sentence can have more than one parse...**

## Probabilistic context-free grammars (PCFGs)

- $G = (T, N, S, R, P)$ 
  - T is the set of terminals (i.e. words)
  - N is the set of non-terminals
    - Usually separate the set P of preterminals (POS tags) from the rest of the non-terminals
  - S is the start symbol
  - R is the set of rules/productions of the form  $X \rightarrow \gamma$  where X is a nonterminal and  $\gamma$  is a sequence of terminals and nonterminals (possibly empty)
  - P(R) gives the probability of each rule

$$\forall X \in N, \sum_{X \rightarrow \gamma \in R} P(X \rightarrow \gamma) = 1$$

- A grammar G generates a language L

Manning

## How to parse

- Top-down:** Start at the top of the tree with an S node, and work your way down to the words.



- Bottom-up:** Look for small pieces that you know how to assemble, and work your way up to larger pieces.



## Given a sentence S...

- We want to find the most likely parse  $\tau$

$$\begin{aligned} \arg \max_{\tau} P(\tau | S) &= \arg \max_{\tau} \frac{P(\tau, S)}{P(S)} \\ &= \arg \max_{\tau} P(\tau, S) \\ &= \arg \max_{\tau} P(\tau) \quad \text{If } S = \text{yield}(\tau) \end{aligned}$$

- How are we supposed to find  $P(\tau)$ ?
- Infinitely many trees in the language!

### Finding P( $\tau$ )

- Define probability distributions over the rules in the grammar
- Context free!

S	→ NP VP	0.8
S	→ S conj S	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP conj NP	0.2
VP	→ Verb	0.4
VP	→ Verb NP	0.3
VP	→ Verb NP NP	0.1
VP	→ VP PP	0.2
PP	→ P NP	1.0

Hockenmaier

### Finding P( $\tau$ )

- The probability of a tree is the product of the probability of the rules that created it

S	→ NP VP	0.8
S	→ S conj S	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP conj NP	0.2
VP	→ Verb	0.4
VP	→ Verb NP	0.3
VP	→ Verb NP NP	0.1
VP	→ VP PP	0.2
PP	→ P NP	1.0

$$P(\tau) = 0.8 \times 0.3 \times 0.2 \times 1.0 \times 0.2^3 = 0.00384$$

Hockenmaier

### Parsing - Cocke-Kasami-Younger (CKY)

- Like Viterbi but for trees
- Guaranteed to find the most likely parse

For each nonterminal: max probability of the subtree it encompasses

This is the tree yield

### Chomsky Normal Form

- All rules are of the form  $X \rightarrow YZ$  or  $X \rightarrow w$ .
- $n$ -ary rules introduce new nonterminals ( $n > 2$ )
  - $VP \rightarrow V NP PP$  becomes:
    - $VP \rightarrow V @VP-V$  and
    - $@VP-V \rightarrow NP PP$

## CKY Example

Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

John	eats	pie	with	cream					
N	NP	S	S		S	John			
		0.8*0.2*0.4	0.8*0.2*0.08		0.2*0.0024*0.8				
	V	VP	VP		VP	eats			
		0.4	0.3*0.2		max(0.008*0.2, 0.06*0.2*0.2)				
			N	NP	NP	pie			
				0.2	0.2*0.2*0.2				
					P	with			
						PP			
						1*0.2			
							N	NP	cream
								0.2	

S	→	NP VP	0.8
S	→	S conj S	0.2
NP	→	Noun	0.2
NP	→	Det Noun	0.4
NP	→	NP PP	0.2
NP	→	NP conj NP	0.2
VP	→	Verb	0.4
VP	→	Verb NP	0.3
VP	→	Verb NP NP	0.1
VP	→	VP PP	0.2
PP	→	P NP	1.0

Hockenmaier

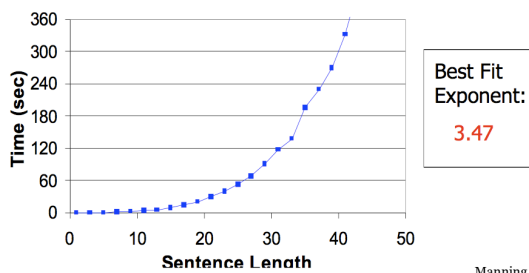
## Estimating $P(X \rightarrow \alpha)$

- Supervised
    - Relative frequency estimation
    - Count what is seen in a treebank corpus
- $$P(X \rightarrow \alpha) = \frac{C(X \rightarrow \alpha)}{C(X)}$$
- Unsupervised
    - Expected relative frequency estimation
    - Use Inside-Outside Algorithm (EM variant)

$$P(X \rightarrow \alpha) = \frac{E[C(X \rightarrow \alpha)]}{E[C(X)]}$$

## How well do PCFGs perform?

- Runtime - supercubic!



## How well do PCFGs perform?

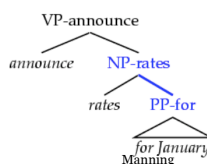
- + Robust to variations in language
- Strong independence assumptions
- ? WSJ parsing accuracy: about 73% LP/LR F1
- Lack of lexicalization
  - A PCFG uses the actual words only to determine the probability of parts-of-speech (the preterminals)
    - I like to eat cake with white frosting.
    - I like to eat cake with a spork.

## Lexicalization

- Lexical heads are important for certain classes of ambiguities (e.g., PP attachment):



- Lexicalizing grammar creates a much larger grammar.
  - Sophisticated smoothing needed
  - Smarter parsing algorithms needed
  - More DATA needed



## Huge area of research

- Coarse-to-fine parsing
  - Parse with a simpler grammar
  - Refine with a more complex one
- Dependency parsing
  - A sentence is parsed by relating each word to other words in the sentence which depend on it.
- Discriminative parsing
  - Given training examples, learn a function that classifies a sentence with its parse tree
- and more!

## The good news!

- Part of speech taggers and sentence parsers are freely available!
- So why did we sit through this lecture?
  - Maybe you'll be interested in this area
  - Useful ideas to be applied elsewhere
    - Write a parser to parse web tables
    - PCFGs for information extraction
  - Like to know how things work

## It's over!

- Thanks!