

Probabilistic Context Free Grammars

Lecture #14

Introduction to Natural Language Processing
CMPSCI 585, Fall 2004



Andrew McCallum

(including slides from Jason Eisner)

Ambiguity in Parsing

- Time flies like an arrow.
- Fruit flies like a banana.
- I saw the man with the telescope.

How to solve this combinatorial explosion of ambiguity?

1. First try parsing without any weird rules, throwing them in only if needed.
2. Better: every rule has a weight. A tree's weight is total weight of all its rules. Pick the overall "lightest" parse of sentence.
3. Can we pick the weights automatically? We'll get to this later ...

CYK Parser

Input: A string of words, grammar in CNF

Output: yes/no

Data structure: $n \times n$ table

rows labeled 0 to $n-1$, columns 1 to n
cell (i,j) lists constituents spanning i,j

For each i from 1 to n

Add to $(i-1,i)$ all Nonterminals that could produce the word at $(i-1,i)$

CYK Parser

For **width** from 2 to n

For **start** from 0 to n -width

Define **end** to be **start+width**

For **mid** from **start+1** to **end-1**

For every constituent in (**start, mid**)

For every constituent in (**mid,end**)

For all ways of combining them (if any)

Add the resulting constituent to (**start,end**).

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3				
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	
4					N 8

- 1 $S \rightarrow NP VP$
- 6 $S \rightarrow Vst NP$
- 2 $S \rightarrow S PP$
- 1 $VP \rightarrow V NP$
- 2 $VP \rightarrow VP PP$
- 1 $NP \rightarrow Det N$
- 2 $NP \rightarrow NP PP$
- 3 $NP \rightarrow NP NP$
- 0 $PP \rightarrow P NP$

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3				
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10			
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8			
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3 Vst 3	NP 10 S 8 S 13			
0					
1		NP 4 VP 4			
2			P 2 V 5		
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3 Vst 3	NP 10 S 8 S 13			
0					
1		NP 4 VP 4			
2			P 2 V 5		PP 12
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3 Vst 3	NP 10 S 8 S 13			
0					
1		NP 4 VP 4			
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3 Vst 3	NP 10 S 8 S 13			
0					
1		NP 4 VP 4			
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3 Vst 3	NP 10 S 8 S 13			
0					
1		NP 4 VP 4			NP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3 Vst 3	NP 10 S 8 S 13			
0					
1		NP 4 VP 4			NP 18 S 21
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24 S 27 S 22
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24 S 27 S 22 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Follow backpointers ...

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24 S 27 S 22 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24 S 27 S 22 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP



time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

1 S → NP VP
 6 S → Vst NP
 2 S → S PP
 1 VP → V NP
 2 VP → VP PP
 1 NP → Det N
 2 NP → NP PP
 3 NP → NP NP
 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

1 S → NP VP
 6 S → Vst NP
 2 S → S PP
 1 VP → V NP
 2 VP → VP PP
 1 NP → Det N
 2 NP → NP PP
 3 NP → NP NP
 0 PP → P NP

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

1 S → NP VP
 6 S → Vst NP
 2 S → S PP
 1 VP → V NP
 2 VP → VP PP
 1 NP → Det N
 2 NP → NP PP
 3 NP → NP NP
 0 PP → P NP

Which entries do we need?

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

1 S → NP VP
 6 S → Vst NP
 2 S → S PP
 1 VP → V NP
 2 VP → VP PP
 1 NP → Det N
 2 NP → NP PP
 3 NP → NP NP
 0 PP → P NP

Which entries do we need?

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

1 S → NP VP
 6 S → Vst NP
 2 S → S PP
 1 VP → V NP
 2 VP → VP PP
 1 NP → Det N
 2 NP → NP PP
 3 NP → NP NP
 0 PP → P NP

Not worth keeping ...

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

1 S → NP VP
 6 S → Vst NP
 2 S → S PP
 1 VP → V NP
 2 VP → VP PP
 1 NP → Det N
 2 NP → NP PP
 3 NP → NP NP
 0 PP → P NP

... since it just breeds worse options

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Keep only best-in-class!

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
0	Vst 3	S 8			S 22
		S 13			S 27
					NP 24
					S 27
					S 22
					S 27
1		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

inferior stock

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Keep only best-in-class!

(and backpointers so you can recover parse)

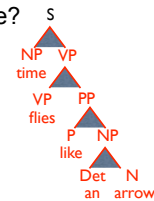
time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10			NP 24
1	Vst 3	S 8			S 22
		NP 4			NP 18
		VP 4			S 21
					VP 18
2			P 2		PP 12
			V 5		VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Probabilistic Trees

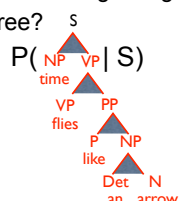
- Instead of lightest weight tree, take highest probability tree
- Given any tree, your assignment generator would have some probability of producing it!
- Just like using n-grams to choose among strings ...
- What is the probability of this tree?



Andrew McCallum, UMass

Probabilistic Trees

- Instead of lightest weight tree, take highest probability tree
- Given any tree, your assignment generator would have some probability of producing it!
- Just like using n-grams to choose among strings ...
- What is the probability of this tree?
- You rolled a lot of independent dice...



Chain rule: One word at a time

$$\begin{aligned}
 & p(\text{time flies like an arrow}) \\
 &= p(\text{time}) \\
 & \quad * p(\text{flies} \mid \text{time}) \\
 & \quad * p(\text{like} \mid \text{time flies}) \\
 & \quad * p(\text{an} \mid \text{time flies like}) \\
 & \quad * p(\text{arrow} \mid \text{time flies like an})
 \end{aligned}$$

Andrew McCallum, UMass

Chain rule + backoff (to get trigram model)

$$\begin{aligned}
 p(\text{time flies like an arrow}) &= p(\text{time}) \\
 &\quad * p(\text{flies} \mid \text{time}) \\
 &\quad * p(\text{like} \mid \text{time flies}) \\
 &\quad * p(\text{an} \mid \text{time flies like}) \\
 &\quad * p(\text{arrow} \mid \text{time flies like an})
 \end{aligned}$$

Andrew McCallum, UMass

Chain rule – written differently

$$\begin{aligned}
 p(\text{time flies like an arrow}) &= p(\text{time}) \\
 &\quad * p(\text{time flies} \mid \text{time}) \\
 &\quad * p(\text{time flies like} \mid \text{time flies}) \\
 &\quad * p(\text{time flies like an} \mid \text{time flies like}) \\
 &\quad * p(\text{time flies like an arrow} \mid \text{time flies like an})
 \end{aligned}$$

$$\text{Proof: } p(x, y \mid x) = p(x \mid x) * p(y \mid x, x) = 1 * p(y \mid x)$$

Andrew McCallum, UMass

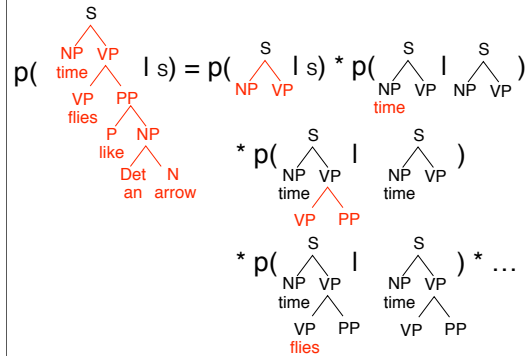
Chain rule + backoff

$$\begin{aligned}
 p(\text{time flies like an arrow}) &= p(\text{time}) \\
 &\quad * p(\text{time flies} \mid \text{time}) \\
 &\quad * p(\text{time flies like} \mid \text{time flies}) \\
 &\quad * p(\text{time flies like an} \mid \text{time flies like}) \\
 &\quad * p(\text{time flies like an arrow} \mid \text{time flies like an})
 \end{aligned}$$

$$\text{Proof: } p(x, y \mid x) = p(x \mid x) * p(y \mid x, x) = 1 * p(y \mid x)$$

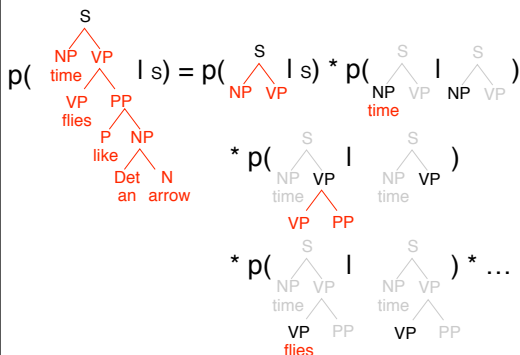
45

Chain rule: One node at a time



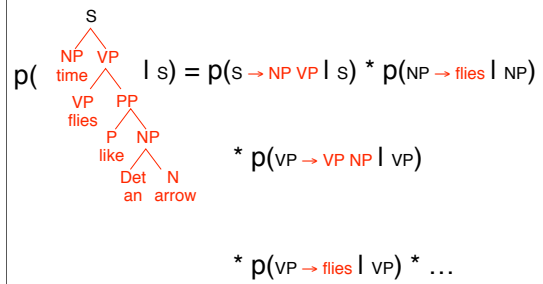
46

Chain rule + backoff



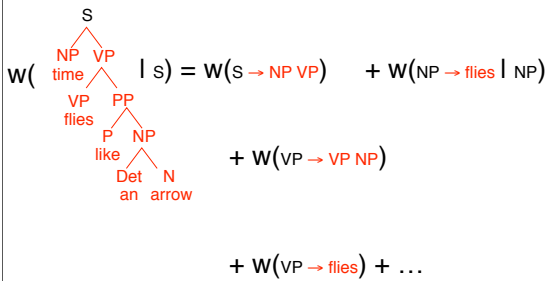
47

Simplified notation



48

Already have a CKY alg for weights ...



Just let $w(x \rightarrow yz) = -\log p(x \rightarrow yz | x)$
Then lightest tree has highest prob ⁴⁹

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10		NP 24
0	Vst 3	S 8		S 22
		S 13		S 27
				NP 24
				S 27
1		NP 4		NP 18
		VP 4		S 21
				VP 18
2			P 2	PP 12
			V 5	VP 16
3			Det 1	NP 10
4				N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Need only best-in-class to get best parse

time 1 flies 2 like 3 an 4 arrow 5

	NP 3	NP 10		NP 24
0	Vst 3	S 8		S 22
		S 13		S 27
				NP 24
				S 27
1		NP 4		NP 18
		VP 4		S 21
				VP 18
2			P 2	PP 12
			V 5	VP 16
3			Det 1	NP 10
4				N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Why probabilities not weights?

- We just saw probabilities are really just a special case of weights ...
- ... *but* we can estimate them from training data by counting and smoothing! Use all of our lovely probability theory machinery!

Probabilistic Context Free Grammars (PCFGs)

A PCFG G consists of the usual parts of a CFG

- A set of terminals, $\{w^k\}, k = 1, \dots, V$
- A set of nonterminals, $\{N^i\}, i = 1, \dots, n$
- A designated start symbol, N^1
- A set of rules, $\{N^i \rightarrow \zeta^j\}$, (where ζ^j is a sequence of terminals and nonterminals)

and

- A corresponding set of probabilities on rules such that:

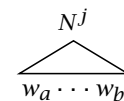
$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$

PCFG notation

Sentence: sequence of words $w_1 \dots w_m$

w_{ab} : the subsequence $w_a \dots w_b$

N_{ab}^i : nonterminal N^i dominates $w_a \dots w_b$



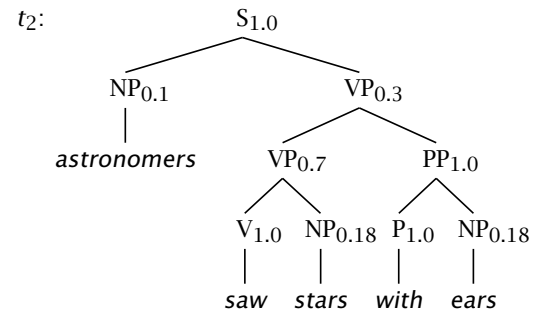
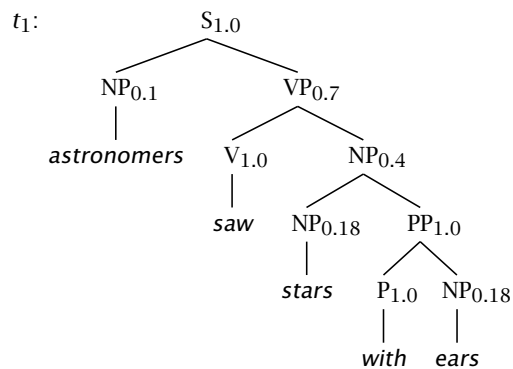
$N^i \xrightarrow{*} \zeta$: Repeated derivation from N^i gives ζ .

PCFG probability of a string

$$\begin{aligned}
 P(w_{1n}) &= \sum_t P(w_{1n}, t) \quad t \text{ a parse of } w_{1n} \\
 &= \sum_{\{t: \text{yield}(t)=w_{1n}\}} P(t)
 \end{aligned}$$

A simple PCFG (in CNF)

$S \rightarrow NP VP$	1.0	$NP \rightarrow NP PP$	0.4
$PP \rightarrow P NP$	1.0	$NP \rightarrow \textit{astronomers}$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow \textit{ears}$	0.18
$VP \rightarrow VP PP$	0.3	$NP \rightarrow \textit{saw}$	0.04
$P \rightarrow \textit{with}$	1.0	$NP \rightarrow \textit{stars}$	0.18
$V \rightarrow \textit{saw}$	1.0	$NP \rightarrow \textit{telescopes}$	0.1



The two parse trees' probabilities and the sentence probability

$$\begin{aligned}
 P(t_1) &= 1.0 \times 0.1 \times 0.7 \times 1.0 \times 0.4 \\
 &\quad \times 0.18 \times 1.0 \times 1.0 \times 0.18 \\
 &= 0.0009072 \\
 P(t_2) &= 1.0 \times 0.1 \times 0.3 \times 0.7 \times 1.0 \\
 &\quad \times 0.18 \times 1.0 \times 1.0 \times 0.18 \\
 &= 0.0006804 \\
 P(w_{15}) &= P(t_1) + P(t_2) = 0.0015876
 \end{aligned}$$

Assumptions of PCFGs

1. Place invariance (like time invariance in HMM):

$$\forall k \quad P(N_{k(k+c)}^j \rightarrow \zeta) \text{ is the same}$$

2. Context-free:

$$P(N_{kl}^j \rightarrow \zeta | \text{words outside } w_k \dots w_l) = P(N_{kl}^j \rightarrow \zeta)$$

3. Ancestor-free:

$$P(N_{kl}^j \rightarrow \zeta | \text{ancestor nodes of } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$$

The sufficient statistics of a PCFG are thus simply counts of how often different local tree configurations occurred (= counts of which grammar rules were applied).

Some features of PCFGs

Reasons to use a PCFG, and some idea of their limitations:

- Partial solution for grammar ambiguity: a PCFG gives some idea of the plausibility of a sentence.
- But, in the simple case, not a very good idea, as independence assumptions are too strong (e.g., not lexicalized).
- Gives a probabilistic language model for English.
- In the simple case, a PCFG is a worse language model for English than a trigram model.
- Better for grammar induction (Gold 1967 vs. Horning 1969)
- Robustness. (Admit everything with low probability.)

Andrew McCallum, UMass

Some features of PCFGs

- A PCFG encodes certain biases, e.g., that smaller trees are normally more probable.
- One can hope to combine the strengths of a PCFG and a trigram model.

We'll look at simple PCFGs first. They have certain inadequacies, but we'll see that most of the state-of-the-art probabilistic parsers are fundamentally PCFG models, just with various enrichments to the grammar

Andrew McCallum, UMass

A slightly different task

- Been asking: What is probability of generating a given *tree*?
 - To pick tree with highest prob: useful in parsing.
- But could also ask: What is probability of generating a given *string* with the generator?
 - To pick string with highest prob: useful in speech recognition, as substitute for an n-gram model.
 - ("Put the file in the folder" vs. "Put the file and the folder")
 - To get prob of generating string, must add up probabilities of all trees for the string ...

63

Could just add up the parse probabilities

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 8 S 13			NP 24 S 22 S 27 NP 24 S 27 S 22 S 27
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 12 VP 16
3				Det 1	NP 10
4					N 8

oops, back to finding exponentially many parses

- 1 S → NP VP
- 6 S → Vst NP
- 2 S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Any more efficient way?

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 2 ⁻⁸ S 2 ⁻¹³			NP 24 S 22 S 27 NP 24 S 27 S 2 ⁻²² S 2 ⁻²⁷
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 2 ⁻¹² VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2⁻² S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Add as we go ... (the "inside algorithm")

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 2 ^{-8+2⁻¹³}			NP 24 S 22 S 27 NP 24 S 27 S 2 ⁻²² S 2 ⁻²⁷
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 2 ⁻¹² VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2⁻² S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Add as we go ... (the "inside algorithm")

time 1 flies 2 like 3 an 4 arrow 5

0	NP 3 Vst 3	NP 10 S 2 ⁻⁸ +2 ⁻¹³			NP 2 ⁻²² +2 ⁻²⁷ S 2 ⁻²² +2 ⁻²⁷ +2 ⁻²⁷ +2 ⁻²² +2 ⁻²⁷
1		NP 4 VP 4			NP 18 S 21 VP 18
2			P 2 V 5		PP 2 ⁻¹² VP 16
3				Det 1	NP 10
4					N 8

- 1 S → NP VP
- 6 S → Vst NP
- 2⁻² S → S PP
- 1 VP → V NP
- 2 VP → VP PP
- 1 NP → Det N
- 2 NP → NP PP
- 3 NP → NP NP
- 0 PP → P NP

Inside and Outside Probabilities

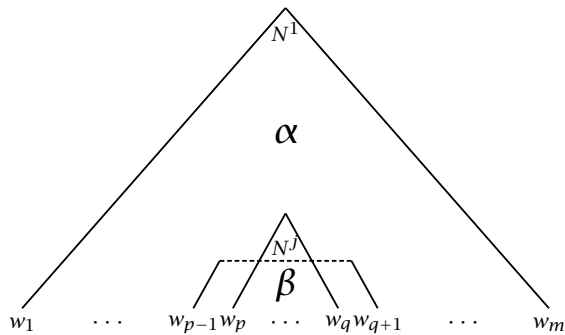
Probability of all possible rule re-writes for generating words inside position p to q, given that non-terminal j exactly spans p to q.

$$\text{Inside} = \beta_j(p, q) = P(w_{pq} | N_{pq}^j, G)$$

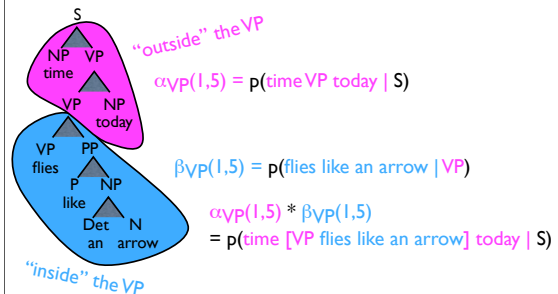
Probability of all possible rule re-writes for generating words *outside* position p to q, and that non-terminal j exactly spans p to q.

$$\text{Outside} = \alpha_j(p, q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$$

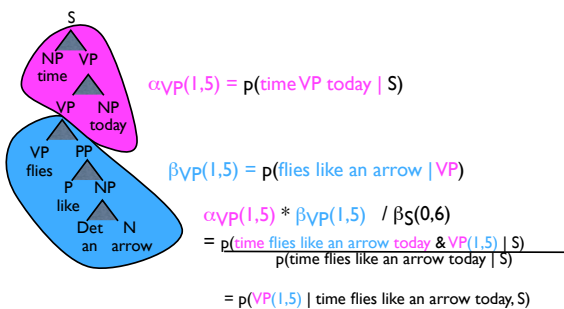
Inside and Outside Probabilities



Inside & Outside Probabilities



Inside & Outside Probabilities



So $\alpha_{VP}(1,5) * \beta_{VP}(1,5) / \beta_S(0,6)$
 is probability that there is a VP here,
 given all of the observed data (words)

Probability of a string

Inside probability

$$P(w_{1m} | G) = P(N^1 \Rightarrow w_{1m} | G)$$

$$= P(w_{1m}, N_{1m}^1, G) = \beta_1(1, m)$$

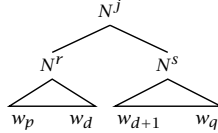
Base case: We want to find $\beta_j(k, k)$ (the probability of a rule $N^j \rightarrow w_k$):

$$\beta_j(k, k) = P(w_k | N_{kk}^j, G)$$

$$= P(N^j \rightarrow w_k | G)$$

Probability of a string

Induction: We want to find $\beta_j(p, q)$, for $p < q$. As this is the inductive step using a Chomsky Normal Form grammar, the first rule must be of the form $N^j \rightarrow N^r N^s$, so we can proceed by induction, dividing the string in two in various places and summing the result:



These inside probabilities can be calculated bottom up.

For all j ,

$$\begin{aligned}
 \beta_j(p, q) &= P(w_p q | N^j_{pq}, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(w_p d | N^r_{pd}, w_{(d+1)q} | N^s_{(d+1)q}, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(N^r_{pd}, N^s_{(d+1)q} | N^j_{pq}, G) \\
 &\quad P(w_p d | N^r_{pd}, N^s_{(d+1)q}, G) \\
 &\quad P(w_{(d+1)q} | N^r_{pd}, N^s_{(d+1)q}, w_p d, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(N^r_{pd}, N^s_{(d+1)q} | N^j_{pq}, G) \\
 &\quad P(w_p d | N^r_{pd}, G) P(w_{(d+1)q} | N^s_{(d+1)q}, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)
 \end{aligned}$$

Inside probabilities as CYK

1	2	3	4	5
$\beta_{NP} = 0.1$		$\beta_S = 0.0126$		$\beta_S = 0.0015876$
	$\beta_{NP} = 0.04$ $\beta_V = 1.0$	$\beta_{VP} = 0.126$		$\beta_{VP} = 0.015876$
		$\beta_{NP} = 0.18$		$\beta_{NP} = 0.01296$
			$\beta_P = 1.0$	$\beta_{PP} = 0.18$
				$\beta_{NP} = 0.18$
astronomers	saw	stars	with	ears

Outside probabilities

Probability of a string: For any k , $1 \leq k \leq m$,

$$\begin{aligned}
 P(w_1 m | G) &= \sum_j P(w_1(k-1), w_k, w_{(k+1)m} | N^j_{kk}, G) \\
 &= \sum_j P(w_1(k-1), N^j_{kk}, w_{(k+1)m} | G) \\
 &\quad \times P(w_k | w_1(k-1), N^j_{kk}, w_{(k+1)m}, G) \\
 &= \sum_j \alpha_j(k, k) P(N^j \rightarrow w_k)
 \end{aligned}$$

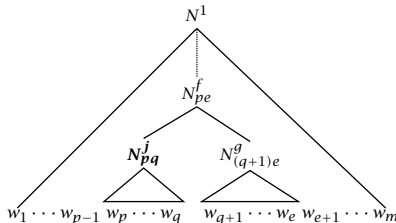
Inductive (DP) calculation: One calculates the outside probabilities top down (after determining the inside probabilities).

Outside probabilities

Base Case:

$$\begin{aligned}
 \alpha_1(1, m) &= 1 \\
 \alpha_j(1, m) &= 0, \text{ for } j \neq 1
 \end{aligned}$$

Inductive Case: it's either a left or right branch - we will sum over both possibilities and calculate using outside and inside probabilities



Outside probabilities, Inductive case

$$\begin{aligned}
 \alpha_j(p, q) &= \left[\sum_{f,g} \sum_{e=q+1}^m P(w_1(p-1), w_{(q+1)m}, N^f_{pe}, N^j_{pq}, N^g_{(q+1)e}) \right. \\
 &\quad \left. + \left[\sum_{f,g} \sum_{e=1}^{p-1} P(w_1(p-1), w_{(q+1)m}, N^f_{pe}, N^g_{e(p-1)}, N^j_{pq}) \right] \right] \\
 &= \left[\sum_{f,g} \sum_{e=q+1}^m P(w_1(p-1), w_{(e+1)m}, N^f_{pe}) P(N^j_{pq}, N^g_{(q+1)e} | N^f_{pe}) \right. \\
 &\quad \times P(w_{(q+1)e} | N^g_{(q+1)e}) \left. + \left[\sum_{f,g} \sum_{e=1}^{p-1} P(w_1(e-1), w_{(q+1)m}, N^f_{eg}) \right. \right. \\
 &\quad \times P(N^g_{e(p-1)}, N^j_{pq} | N^f_{eg}) P(w_{(e+1)e} | N^g_{e(p-1)}) \left. \left. \right] \right] \\
 &= \left[\sum_{f,g} \sum_{e=q+1}^m \alpha_f(p, e) P(N^f \rightarrow N^j N^g) \beta_g(q+1, e) \right] \\
 &\quad + \left[\sum_{f,g} \sum_{e=1}^{p-1} \alpha_f(e, q) P(N^f \rightarrow N^j N^g) \beta_g(e, p-1) \right]
 \end{aligned}$$

Probability that a rule is used

As with a HMM, we can form a product of the inside and outside probabilities. This time:

$$\begin{aligned} & \alpha_j(p, q)\beta_j(p, q) \\ &= P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m}|G)P(w_{pq}|N_{pq}^j, G) \\ &= P(w_{1m}, N_{pq}^j|G) \end{aligned}$$

$$P(N_{pq}^j|w_{1m}, G) = \frac{P(N_{pq}^j|w_{1m}, G)}{P(w_{1m}|G)} = \frac{\alpha_j(p, q)\beta_j(p, q)}{\beta_1(1, m)}$$

This is an “expected count” for the number of times this rule occurred.

Overall probability of a node existing

As with a HMM, we can form a product of the inside and outside probabilities. This time:

$$\begin{aligned} & \alpha_j(p, q)\beta_j(p, q) \\ &= P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m}|G)P(w_{pq}|N_{pq}^j, G) \\ &= P(w_{1m}, N_{pq}^j|G) \end{aligned}$$

Therefore,

$$p(w_{1m}, N_{pq}|G) = \sum_j \alpha_j(p, q)\beta_j(p, q)$$

Just in the cases of the root node and the preterminals, we know there will always be some such constituent.

Learning PCFGs (1)

- We would like to calculate how often each rule is used:

$$\hat{P}(N^j \rightarrow \zeta) = \frac{C(N^j \rightarrow \zeta)}{\sum_y C(N^j \rightarrow y)}$$

- If we have labeled data, we count and find out
- Relative frequency again gives maximum likelihood probability estimates
- This is the motivation for building *Trebanks* of hand-parsed sentences

Learning PCFGs (2) Inside-Outside

- Otherwise we work iteratively from expectations of current model.
- We construct an EM training algorithm, as for HMMs
- For each sentence, at each iteration, we work out expectation of how often each rule is used using inside and outside probabilities
- We assume sentences are independent and sum expectations over parses of each
- We re-estimate rules based on these ‘counts’