

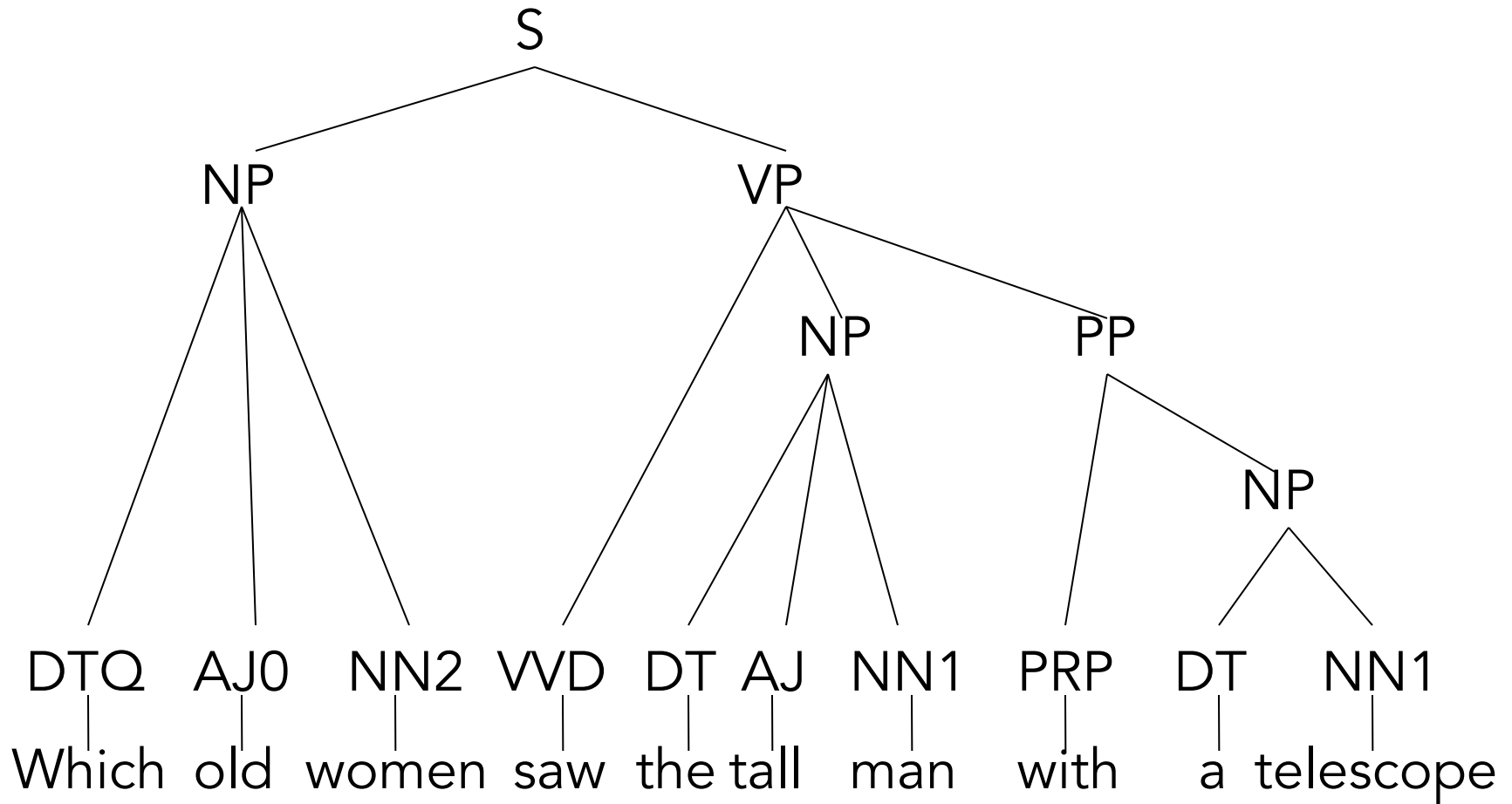
CKY Parsing

CS-585

Natural Language Processing

Derrick Higgins

Phrase Structure Trees



Parsing

- Recognize if a sentence is valid
- Figure out its syntactic structure
- Problem: Find all parse trees licensed by a given grammar, covering exactly the input words.
- Is it possible to parse a sentence deterministically as it is being read (e.g. from left-to-right)?

Computer Languages vs. Natural Languages

A parser for a computer language

- yields a unique tree for each string
- must be deterministic
- is allowed (basically) unrestricted memory

A natural language parser

- must allow for more than one parse
- should predict which parse will most likely be offered as the first choice by a native speaker
- human short-term memory is restricted

Properties of NL Parsing

- Highly ambiguous
- Different theoretical notions of *grammar*
- All solutions should be found (in principle)
- Analysis problem more complex
- Solutions based on saving partial parses

Context-Free Grammar

- Start symbol S
- Set of non-terminal symbols $\{NP, VP, \dots\}$
- Set of terminal symbols (words)
- Set of production rules, of the form

$$NT \rightarrow a \ b \ c \ \dots$$

where NT is a non-terminal and $a \ b \ c$ comprise a sequence of 1 or more terminals and non-terminals

Example

S	→	NP VP	Name	→	joe
NP	→	Name	N	→	ice
NP	→	N	N	→	drinks
NP	→	NP PP	N	→	water
VP	→	V NP	V	→	drinks
VP	→	V	P	→	with
VP	→	VP PP			
PP	→	P NP			

“joe drinks water with ice”

Parser Properties

Soundness: A parser is sound if every parse returned is valid in the grammar.

Completeness: A parser is complete if for every grammar and sentence it returns all valid parses for that sentence.

- Soundness is key...
- ...but completeness may be difficult or even undesirable, e.g. for highly ambiguous grammars...

Top-Down Parsing

- Start with list [S] of goals to achieve
- Iteratively rewrite the goal list by expanding rules until it matches the sentence
- Choices at each step:
 1. Which rule to use if several apply to a given nonterminal
 2. Which order to address the subgoals (depth-first, breadth-first, etc.)

Example

joe drinks water with ice

[S]

[NP, VP]

[Name, VP]

[joe, VP]

[joe, V, NP]

[joe, drinks, NP]

[joe, drinks, Name]

...

[joe, drinks, NP]

[joe, drinks, N]

[joe, drinks, ice]

...

[joe, drinks, NP]...

[joe, drinks, NP, PP]...

[joe, drinks, water, P, N]

[joe, drinks, water, with, ice]

Issues with Top-Down Parsing

- LR depth-first parser won't terminate if grammar has left-recursive rules
 - E.g., $NP \rightarrow NP PP$ will produce
$$NP \rightarrow NP PP \rightarrow NP PP PP \rightarrow NP PP PP PP \rightarrow \dots$$
 - Even combinations may be left-recursive:
$$NPos \rightarrow NP \text{ “ ’ ” “ s ” } \quad (\text{john's book})$$
$$NP \rightarrow NPos NP$$
- Many rules with same LHS may lead to a lot of backtracking...

Bottom-Up Parsing

- Goal list initialized as list of terminals in the string to be parsed
- If sequence of goals matches RHS of a rule, replace it with the LHS of the rule
- Parsing complete when producing S
- Choices:
 1. RHS of multiple rules may match
 2. Order of subgoals (depth-first, breadth-first)

NB: Inefficient when grammar has lexical ambiguity

Shift-Reduce Parsing

DoneGoals	String	Operation
	Nancy ate waffles today	shift
NP	ate waffles today	shift
NP V	waffles today	shift
NP V NP	today	reduce
NP VP	today	shift
NP VP AV		reduce
NP VP		reduce
S		

Effects of ambiguity

- Produces inefficiency due to backtracking
 - Multiple matching LHSs in top-down
 - Lexical ambiguity (multiple RHSs) in bottom-up
- How to produce all (most? many?) parses efficiently?
 - There may be an exponential number of them

Chart Parsing

- Remember intermediate results
- Explore all possible solutions in parallel

Sentence: $w_1 w_2 w_3 w_4 \dots$

Chart: Array whose entries show the set of categories that could generate words from n to $n+m$

Formally:

$$chart(m, n) = \{ A \mid A \rightarrow^* w_n \dots w_{n+m} \}$$

Example

The₀ man₁ drinks₂ water₃ with₄ ice₅

n (constituent start index)

	0	1	2	3	4	5
0	Det	N, NP	V, VP, NP	N, NP	P	N, NP
1	NP	S	VP	{ }	PP	
2	S	S	{ }	NP		
3	S	{ }	VP			
4	{ }	S				
5	S					

m (constituent length -1)

Chomsky Normal Form

- Constraint on form of the grammar:
 - Each RHS is either 2 non-terminals or a terminal
- All CFGs can be written in CNF

S	→	NP	VP	Det	→	the
NP	→	NP	PP	NP	→	joe
NP	→	Det	NP	NP	→	ice
VP	→	V	NP	NP	→	drinks
VP	→	VP	PP	NP	→	water
PP	→	P	NP	V	→	drinks
				VP	→	drinks
				P	→	with

Cocke-Younger-Kasami (CYK)

Assume “Chomsky Normal Form” grammar

```
for n := 0 to  $N_w-1$  do:
```

```
  chart[0, n] := {X |  $X \rightarrow \text{word}_n$ }
```

```
for m := 1 to  $N_w-1$  do:
```

```
  for n := 0 to  $N_w-m-1$  do:
```

```
    chart[m, n] := {}
```

```
    for k := n+1 to n+m do
```

```
      for every rule  $A \rightarrow B C$  do
```

```
        if  $B \in \text{chart}[k-n-1, n]$  and  $C \in \text{chart}[n+m-k, k]$  then
```

```
          chart[m, n] := chart[m, n]  $\cup$  {A}
```

```
if  $S \in \text{chart}[N_w-1, 0]$  then accept else reject
```

CYK Example (in CNF)

S	→	NP	VP	NP	→	joe
NP	→	NP	PP	NP	→	ice
VP	→	V	NP	NP	→	drinks
VP	→	VP	PP	NP	→	water
PP	→	P	NP	V	→	drinks
				VP	→	drinks
				P	→	with

“joe drinks water with ice”

Cocke-Younger-Kasami (CYK)

```
for n := 0 to  $N_w - 1$  do:  
  chart[0, n] := {X |  $X \rightarrow \text{word}_n$ }
```

```
for m := 1 to  $N_w - 1$  do:  
  for n := 0 to  $N_w - m - 1$  do:
```

```
    chart[m, n] := {}
```

```
    for k := 1 to m do:  
      Initialize chart with terminal symbols
```

```
      for every rule  $A \rightarrow B C$  do
```

```
        if  $B \in \text{chart}[k - n - 1, n]$  and  $C \in \text{chart}[n + m - k, k]$  then
```

```
          chart[m, n] := chart[m, n]  $\cup$  {A}
```

```
if  $S \in \text{chart}[N_w - 1, 0]$  then accept else reject
```

Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1					
2					
3					
4					

$N_w = 5$

Cocke-Younger-Kasami (CYK)

Look for increasingly longer phrases

for $n := 0$ to $N_w - 1$ do
chart $[0, n] := \{X \mid X \text{ is a word}\}$ Start from the left-hand edge and stop if the constituent would run past the end of the sentence

for $m := 1$ to $N_w - 1$ do:

for $n := 0$ to $N_w - m - 1$ do:

chart $[m, n] := \{\}$

for $k := n + 1$ to $n + m$ do

for every rule $A \rightarrow B C$ do

if $B \in \text{chart}[k - n - 1, n]$ and $C \in \text{chart}[n + m - k, k]$ then

chart $[m, n] := \text{chart}[m, n] \cup \{A\}$

Consider all ways you could divide the text span into two parts, and look for a rule that matches

if $S \in \text{chart}[N_w - 1, 0]$ then accept else reject

Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S				
2					
3					
4					

$$N_w = 5$$

$$k = n+1$$

S → NP VP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP			
2					
3					
4					

$$N_w = 5$$

$$k = n+1$$

VP \rightarrow V NP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }		
2					
3					
4					

$$N_w = 5$$

$$k = n+1$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2					
3					
4					

$$N_w = 5$$

$$k = n+1$$

PP → P NP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S				
3					
4					

$$N_w = 5$$

$$k = n+1$$

S → NP VP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

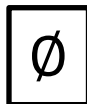
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S				
3					
4					

$$N_w = 5$$

$$k = n+2$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

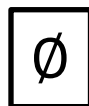
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }			
3					
4					

$$N_w = 5$$

$$k = n+1$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

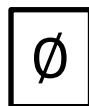
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }			
3					
4					

$$N_w = 5$$

$$k = n+2$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3					
4					

$$N_w = 5$$

$$k = n+1$$

NP → NP PP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

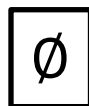
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3					
4					

$$N_w = 5$$

$$k = n+2$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

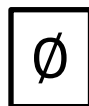
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }				
4					

$$N_w = 5$$

$$k = n+1$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

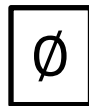
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }				
4					

$$N_w = 5$$

$$k = n+2$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

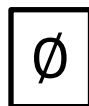
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }				
4					

$$N_w = 5$$

$$k = n+3$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4					

$$N_w = 5$$

$$k = n+1$$

VP \rightarrow V NP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4					

$$N_w = 5$$

$$k = n+2$$

VP → VP PP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

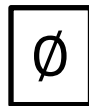
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4					

$$N_w = 5$$

$$k = n+3$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4	S				

$$N_w = 5$$

$$k = n+1$$

S → NP VP

Example

Joe₀ drinks₁ water₂ with₃ ice₄

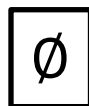
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4	S				

$$N_w = 5$$

$$k = n+2$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

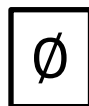
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4	S				

$$N_w = 5$$

$$k = n+3$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

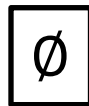
n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4	S				

$$N_w = 5$$

$$k = n+4$$



Example

Joe₀ drinks₁ water₂ with₃ ice₄

n (constituent start index)

m (constituent length -1)

	0	1	2	3	4
0	NP	V, VP, NP	NP	P	NP
1	S	VP	{ }	PP	
2	S	{ }	NP		
3	{ }	VP			
4	S				

$N_w = 5$

Generalizing to non-CNF grammars

- Well-Formed Substring Table (WFST)
 - Same as CYK table, but allows RHS of rules to have many non-terminals
- Complexity: $O(n^{k+1})$ where k is largest number of non-terminals in a rule

WFST Example

S → NP VP

NP → NP PP

VP → V_{trans} NP

VP → V_{ditrans} NP NP

VP → VP PP

PP → P NP

NP → joe

NP → mary

NP → ice

NP → drinks

NP → water

V_{trans} → drinks

V_{ditrans} → gives

P → with

“joe gives mary water”

Active Charts

- WFST stores complete analyses....
...but this requires redoing partial results:

$VP \rightarrow \mathbf{V}_{ditrans} \mathbf{NP} PP_{to}$

$VP \rightarrow \mathbf{V}_{ditrans} \mathbf{NP} NP$

Idea: Store partial results as well!

- Active chart contains:
 - Passive items: complete results
 - Active items: partial results



E

OF TECHNOLOGY

Transforming Lives. Inventing the Future. www.iit.edu

Dotted Rules

- Partial results (in chart) must consider how much has been recognized so far:

Dotted rule	For a VP we still need...
$VP \rightarrow \blacklozenge Vdit \ NP \ PP_{to}$	$Vdit$, NP , and PP_{to}
$VP \rightarrow Vdit \ \blacklozenge NP \ PP_{to}$	NP and PP_{to}
$VP \rightarrow Vdit \ NP \ \blacklozenge PP_{to}$	PP_{to}
$VP \rightarrow Vdit \ NP \ PP_{to} \ \blacklozenge$	nothing