

Chart Parsing: the CYK Algorithm

Informatics 2A: Lecture 18

Shay Cohen

3 November 2015

Deterministic parsing (e.g., LL(1)) aims to address a limited amount of **local ambiguity** – the problem of not being able to decide uniquely which grammar rule to use next in a left-to-right analysis of the input string.

By re-structuring the grammar, the parser can make a unique decision, based on a limited amount of **look-ahead**.

Recursive Descent parsing also demands grammar restructuring, in order to eliminate left-recursive rules that can get it into a hopeless loop.

But grammars for natural human languages should be **revealing**, re-structuring the grammar may destroy this. (Indirectly) left-recursive rules are needed in English.

NP \rightarrow DET N

NP \rightarrow NPR

DET \rightarrow NP 's

These rules generate NPs with possessive modifiers such as:

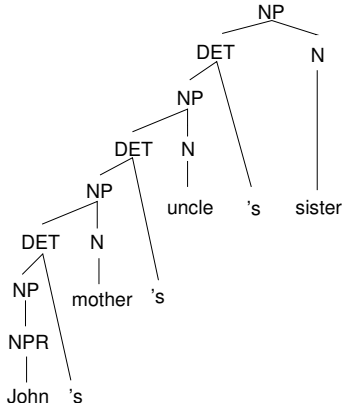
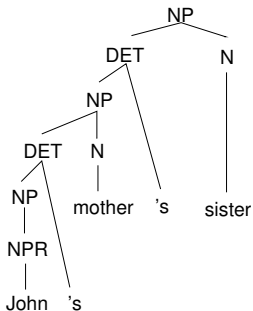
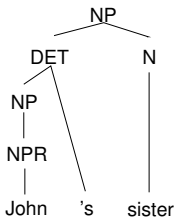
John's sister

John's mother's sister

John's mother's uncle's sister

John's mother's uncle's sister's niece

Left Recursion



We don't want to re-structure our grammar rules just to be able to use a particular approach to parsing. Need an alternative.

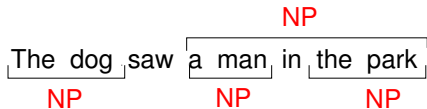
Problems with Parsing as Search

- 1 A **recursive descent parser** (top-down) will do badly if there are many different rules for the same LHS. Hopeless for rewriting parts of speech (preterminals) with words (terminals).
- 2 A **shift-reduce parser** (bottom-up) does a lot of useless work: many phrase structures will be locally possible, but globally impossible. Also inefficient when there is much lexical ambiguity.
- 3 Both strategies do repeated work by **re-analyzing** the same substring many times.

We will see how **chart parsing** solves the re-parsing problem, and also copes well with ambiguity.

Dynamic Programming

With a CFG, a parser should be able to avoid re-analyzing sub-strings because the analysis of any sub-string is **independent** of the rest of the parse.



The parser's exploration of its search space can exploit this independence if the parser uses **dynamic programming**.

Dynamic programming is the basis for all **chart parsing** algorithms.

Parsing as Dynamic Programming

- Given a problem, systematically fill a table of solutions to sub-problems: this is called **memoization**.
- Once solutions to all sub-problems have been accumulated, solve the overall problem by composing them.
- For parsing, the sub-problems are analyses of sub-strings and correspond to **constituents** that have been found.
- Sub-trees are stored in a **chart** (aka **well-formed substring table**), which is a record of all the substructures that have ever been built during the parse.

Solves **re-parsing problem**: sub-trees are looked up, not re-parsed!

Solves **ambiguity problem**: chart implicitly stores all parses!

A **chart** can be depicted as a matrix:

- Rows and columns of the matrix correspond to the start and end positions of a span (ie, starting **right before** the first word, ending **right after** the final one);
- A cell in the matrix corresponds to the sub-string that starts at the row index and ends at the column index.
- It can contain information about the **type** of constituent (or constituents) that span(s) the substring, pointers to its sub-constituents, and/or **predictions** about what constituents might follow the substring.

CYK (Cocke, Younger, Kasami) is an algorithm for recognizing and recording constituents in the chart.

- Assumes that the grammar is in Chomsky Normal Form: rules all have form $A \rightarrow BC$ or $A \rightarrow w$.
- Conversion to CNF can be done automatically.

NP	→	Det Nom	NP	→	Det Nom
Nom	→	N OptAP Nom	Nom	→	<i>book</i> <i>orange</i> AP Nom
OptAP	→	ε OptAdv A	AP	→	<i>heavy</i> <i>orange</i> Adv A
A	→	<i>heavy</i> <i>orange</i>	A	→	<i>heavy</i> <i>orange</i>
Det	→	<i>a</i>	Det	→	<i>a</i>
OptAdv	→	ε <i>very</i>	Adv	→	<i>very</i>
N	→	<i>book</i> <i>orange</i>			

Let's look at a simple example before we explain the general case.

Grammar Rules in CNF

NP	→	Det	Nom
Nom	→	<i>book</i>	<i>orange</i> AP Nom
AP	→	<i>heavy</i>	<i>orange</i> Adv A
A	→	<i>heavy</i>	<i>orange</i>
Det	→	<i>a</i>	
Adv	→	<i>very</i>	

(N.B. Converting to CNF sometimes breeds duplication!)

Now let's parse: *a very heavy orange book*

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

	1 <i>a</i>	2 <i>very</i>	3 <i>heavy</i>	4 <i>orange</i>	5 <i>book</i>
0 <i>a</i>					
1 <i>very</i>					
2 <i>heavy</i>					
3 <i>orange</i>					
4 <i>book</i>					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1 <i>a</i>	2 <i>very</i>	3 <i>heavy</i>	4 <i>orange</i>	5 <i>book</i>
0	a	Det				
1	very					
2	heavy					
3	orange					
4	book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1 <i>a</i>	2 <i>very</i>	3 <i>heavy</i>	4 <i>orange</i>	5 <i>book</i>
0	a	Det				
1	very		Adv			
2	heavy					
3	orange					
4	book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1 <i>a</i>	2 <i>very</i>	3 <i>heavy</i>	4 <i>orange</i>	5 <i>book</i>
0	a	Det				
1	very		Adv			
2	heavy			A,AP		
3	orange					
4	book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

	1 <i>a</i>	2 <i>very</i>	3 <i>heavy</i>	4 <i>orange</i>	5 <i>book</i>
0 a	Det				
1 very		Adv	AP		
2 heavy			A,AP		
3 orange					
4 book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

	1	2	3	4	5
	<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det			
1	very		Adv	AP	
2	heavy		A,AP		
3	orange			Nom,A,AP	
4	book				

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1	2	3	4	5
		<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det				
1	very		Adv	AP		
2	heavy			A,AP	Nom	
3	orange				Nom,A,AP	
4	book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1	2	3	4	5
		<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det				
1	very		Adv	AP	Nom	
2	heavy			A,AP	Nom	
3	orange				Nom,A,AP	
4	book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1	2	3	4	5
		<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det			NP	
1	very		Adv	AP	Nom	
2	heavy			A,AP	Nom	
3	orange				Nom,A,AP	
4	book					

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

	1	2	3	4	5
	<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a Det			NP	
1	very	Adv	AP	Nom	
2	heavy		A,AP	Nom	
3	orange			Nom,A,AP	
4	book				Nom

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1	2	3	4	5
		<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det			NP	
1	very		Adv	AP	Nom	
2	heavy			A,AP	Nom	
3	orange				Nom,A,AP	Nom
4	book					Nom

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1	2	3	4	5
		<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det			NP	
1	very		Adv	AP	Nom	
2	heavy			A,AP	Nom	Nom
3	orange				Nom,A,AP	Nom
4	book					Nom

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

	1	2	3	4	5
	<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a Det			NP	
1	very	Adv	AP	Nom	Nom
2	heavy		A,AP	Nom	Nom
3	orange			Nom,A,AP	Nom
4	book				Nom

Filling out the CYK chart

0 a 1 very 2 heavy 3 orange 4 book 5

		1	2	3	4	5
		<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a	Det			NP	NP
1	very		Adv	AP	Nom	Nom
2	heavy			A,AP	Nom	Nom
3	orange				Nom,A,AP	Nom
4	book					Nom

CYK: The general algorithm

function CKY-Parse(*words*, *grammar*) **returns** *table* **for**

j ← **from** 1 **to** LENGTH(*words*) **do**

table[*j* - 1, *j*] ← {*A* | *A* → *words*[*j*] ∈ *grammar*}

for *i* ← **from** *j* - 2 **downto** 0 **do**

for *k* ← *i* + 1 **to** *j* - 1 **do**

table[*i*, *j*] ← *table*[*i*, *j*] ∪

 {*A* | *A* → *BC* ∈ *grammar*,

B ∈ *table*[*i*, *k*]

C ∈ *table*[*k*, *j*]}

CYK: The general algorithm

function CKY-Parse(*words*, *grammar*) **returns** *table* **for**

j ← **from** 1 **to** LENGTH(*words*) **do**

loop over the columns

$table[j-1, j] \leftarrow \{A \mid A \rightarrow words[j] \in grammar\}$ fill bottom cell

for *i* ← **from** *j* - 2 **downto** 0 **do**

fill row *i* in column *j*

for *k* ← *i* + 1 **to** *j* - 1 **do**

loop over split locations

$table[i, j] \leftarrow table[i, j] \cup$

between *i* and *j*

$\{A \mid A \rightarrow BC \in grammar,$
 $B \in table[i, k]$
 $C \in table[k, j]\}$

Check the grammar for rules that link the constituent in $[i, k]$ with those in $[k, j]$. For each rule found store LHS in cell $[i, j]$.

A succinct representation of CKY

We have a Boolean table called *Chart*, such that $\text{Chart}[A, i, j]$ is true if there is a sub-phrase according to the grammar that dominates words i through words j

Build this chart recursively, similarly to the Viterbi algorithm:

For $j > i + 1$:

$$\text{Chart}[A, i, j] = \bigvee_{k=i+1}^{j-1} \bigvee_{A \rightarrow BC} \text{Chart}[B, i, k] \wedge \text{Chart}[C, k, j]$$

Seed the chart, for $i + 1 = j$:

$\text{Chart}[A, i, i + 1] = \text{True}$ if there exists a rule $A \rightarrow w_{i+1}$ where w_{i+1} is the $(i + 1)$ th word in the string

From CYK Recognizer to CYK Parser

- So far, we just have a chart **recognizer**, a way of determining whether a string belongs to the given language.
- Changing this to a **parser** requires recording which existing constituents were combined to make each new constituent.
- This requires another field to record the one or more ways in which a constituent spanning (i,j) can be made from constituents spanning (i,k) and (k,j) . (More clearly displayed in **graph** representation, see next lecture.)
- In any case, for a fixed grammar, the CYK algorithm runs in time $O(n^3)$ on an input string of n tokens.
- The algorithm identifies **all possible parses**.

CYK-style parse charts

Even without converting a grammar to CNF, we can draw *CYK-style* parse charts:

	1	2	3	4	5
	<i>a</i>	<i>very</i>	<i>heavy</i>	<i>orange</i>	<i>book</i>
0	a Det			NP	NP
1	very	OptAdv	OptAP	Nom	Nom
2	heavy		A,OptAP	Nom	Nom
3	orange			N,Nom,A,AP	Nom
4	book				N,Nom

(We haven't attempted to show ϵ -phrases here. Could in principle use cells below the main diagonal for this ...)

However, CYK-style parsing will have run-time worse than $O(n^3)$ if e.g. the grammar has rules $A \rightarrow BCD$.

Grammar Rules in CNF

$S \rightarrow NP VP$	$Nominal \rightarrow book flight money$
$S \rightarrow X1 VP$	$Nominal \rightarrow Nominal noun$
$X1 \rightarrow Aux VP$	$Nominal \rightarrow Nominal PP$
$S \rightarrow book include prefer$	$VP \rightarrow book include prefer$
$S \rightarrow Verb NP$	$VPVerb \rightarrow NP$
$S \rightarrow X2$	$VP \rightarrow X2 PP$
$S \rightarrow Verb PP$	$X2 \rightarrow Verb NP$
$S \rightarrow VP PP$	$VP \rightarrow Verb NP$
$NP \rightarrow TWA Houston$	$VP \rightarrow VP PP$
$NP \rightarrow Det Nominal$	$PP \rightarrow Preposition NP$
$Verb \rightarrow book include prefer$	$Noun \rightarrow book flight money$

Let's parse *Book the flight through Houston!*

Grammar Rules in CNF

$S \rightarrow NP VP$

$S \rightarrow X1 VP$

$X1 \rightarrow Aux VP$

$S \rightarrow book|include|prefer$

$S \rightarrow Verb NP$

$S \rightarrow X2$

$S \rightarrow Verb PP$

$S \rightarrow VP PP$

$NP \rightarrow TWA|Houston$

$NP \rightarrow Det Nominal$

$Verb \rightarrow book|include|prefer$

$Nominal \rightarrow book|flight|money$

$Nominal \rightarrow Nominal noun$

$Nominal \rightarrow Nominal PP$

$VP \rightarrow book|include|prefer$

$VPVerb \rightarrow NP$

$VP \rightarrow X2 PP$

$X2 \rightarrow Verb NP$

$VP \rightarrow Verb NP$

$VP \rightarrow VP PP$

$PP \rightarrow Preposition NP$

$Noun \rightarrow book|flight|money$

Let's parse *Book the flight through Houston!*

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]				

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]				
	Det [1, 2]			

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]				
	Det [1, 2]			
		Nominal, Noun [2, 3]		

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]				
	Det [1, 2]			
		Nominal, Noun [2, 3]		
			Prep [3, 4]	

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]				
	Det [1, 2]			
		Nominal, Noun [2, 3]		
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]			
	Det [1, 2]			
		Nominal, Noun [2, 3]		
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]			
	Det [1, 2]	NP [1, 3]		
		Nominal, Noun [2, 3]		
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]		
	Det [1, 2]	NP [1, 3]		
		Nominal, Noun [2, 3]		
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]		
	Det [1, 2]	NP [1, 3]		
		Nominal, Noun [2, 3]	[2, 4]	
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]		
	Det [1, 2]	NP [1, 3]	[1, 4]	
		Nominal, Noun [2, 3]	[2, 4]	
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	
	Det [1, 2]	NP [1, 3]	[1, 4]	
		Nominal, Noun [2, 3]	[2, 4]	
			Prep [3, 4]	
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	
	Det [1, 2]	NP [1, 3]	[1, 4]	
		Nominal, Noun [2, 3]	[2, 4]	
			Prep [3, 4]	PP [3, 5]
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	
	Det [1, 2]	NP [1, 3]	[1, 4]	
		Nominal, Noun [2, 3]	[2, 4]	Nominal [2, 5]
			Prep [3, 4]	PP [3, 5]
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	
	Det [1, 2]	NP [1, 3]	[1, 4]	NP [1, 5]
		Nominal, Noun [2, 3]	[2, 4]	Nominal [2, 5]
			Prep [3, 4]	PP [3, 5]
				NP, Proper- Noun [4, 5]

Second example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	S ₁ , VP, X2, S ₂ , VP, S ₃ [0, 5]
	Det [1, 2]	NP [1, 3]	[1, 4]	NP [1, 5]
		Nominal, Noun [2, 3]	[2, 4]	Nominal [2, 5]
			Prep [3, 4]	PP [3, 5]
				NP, Proper- Noun [4, 5]

Visualizing the Chart

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	S ₁ , VP, X2, S ₂ , VP, S ₃ [0, 5]
	Det [1, 2]	NP [1, 3]	[1, 4]	NP [1, 5]
		Nominal, Noun [2, 3]	[2, 4]	Nominal [2, 5]
			Prep ← [3, 4]	PP [3, 5]
				NP, Proper- Noun [4, 5]

Visualizing the Chart

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0, 1]	[0, 2]	S, VP, X2 [0, 3]	[0, 4]	S ₁ , VP, X2, S ₂ , VP, S ₃ [0, 5]
	Det [1, 2]	NP [1, 3]	[1, 4]	NP [1, 5]
		Nominal, Noun [2, 3]	[2, 4]	Nominal [2, 5]
			Prep [3, 4]	PP [3, 5]
				NP, Proper- Noun [4, 5]

A Tribute to CKY (part 1/3)

*You, my CKY algorithm,
dictate every parser's rhythm,
if Cocke, Younger and Kasami hadn't bothered,
all of our parsing dreams would have been shattered.*

*You are so simple, yet so powerful,
and with the proper semiring and time,
you will be truthful,
to return the best parse - anything less would be a crime.*

*With dynamic programming or memoization,
you are one of a kind,
I really don't need to mention,
if it weren't for you, all syntax trees would be behind.*

*Failed attempts have been made to show there are better,
for example, by using matrix multiplication,
all of these impractical algorithms didn't matter
you came out stronger, insisting on just using summation.*

*All parsing algorithms to you hail,
at least those with backbones which are context-free,
you will never become stale,
as long as we need to have a syntax tree.*

*It doesn't matter that the C is always in front,
or that the K and Y can swap,
you are still on the same hunt,
maximizing and summing, nonstop.*

A Tribute to CKY (part 3/3)

*Every Informatics student knows you intimately,
they have seen your variants dozens of times,
you have earned that respect legitimately,
and you will follow them through their primes.*

*CKY, going backward and forward,
inside and out,
it is so straightforward -
You are the best, there is no doubt.*

- Parsing as search is inefficient (typically exponential time).
- Alternative: use dynamic programming and memoize sub-analysis in a chart to avoid duplicate work.
- The chart can be visualized as as a matrix.
- The CYK algorithm builds a chart in $O(n^3)$ time. The basic version gives just a recognizer, but it can be made into a parser if more info is recorded in the chart.

Reading: J&M (2nd ed), Chapter. 13, Sections 13.3–13.4
NLTK Book, Chapter. 8 (*Analyzing Sentence Structure*), Section 8.4

Next lecture: the Earley parser or dynamic programming for top-down parsing