#### Natural Language Processing

#### Lecture 15—3/5/2015 Martha Palmer

## Today

- Review CKY
- Earley
- Partial parsing
  - Finite-state methods
  - Chunking
    - Sequence labeling methods

## **CKY Algorithm**

**function** CKY-PARSE(*words*, *grammar*) **returns** *table* 

$$\begin{array}{ll} \mbox{for } j \leftarrow \mbox{from 1 to LENGTH}(words) \mbox{ do} & \mbox{Looping over the columns} \\ table[j-1,j] \leftarrow \{A \mid A \rightarrow words[j] \in gram \mbox{ Filling the bottom cell} \\ \mbox{for } i \leftarrow \mbox{from } j - 2 \mbox{ downto 0 do} & \mbox{Filling row i in column j} \\ \mbox{for } k \leftarrow i + 1 \mbox{ to } j - 1 \mbox{ do} & \mbox{Filling row i in column j} \\ \mbox{table}[i,j] \leftarrow table[i,j] \cup & \mbox{Looping over the possible split locations} \\ \mbox{table}[i,j] \leftarrow table[i,j] \cup & \mbox{between i and j.} \\ \mbox{Check the grammar for rules that} \\ \mbox{link the constituents in [i,k] with} \\ \mbox{to se in [k,j]. For each rule} \\ \mbox{found store the LHS of the rule in} & \mbox{c} \in table[i,k], \\ \mbox{C} \in table[k,j] \\ \end{table} \end{array}$$

Book	the	flight	through	Houston	
S, VP, Verb Nominal, Noun	10 21	S,VP,X2	10 41	S,VP,X2	
[[0,1]	Det	NP	[0,4]	NP	
	[1,2]	[1,3]	[1,4]	[1,5] Neminal	
		Noun		Nominal	
		[2,3]	[2,4]	[2,5]	
			Prep	PP	
			[3,4]	[3,5]	
				NP, Proper- Noun	
				[4,5]	

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

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

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

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun		S,VP,X2		
[0,1]	[0,2]	[0,3]	[0,4]	[0,5]
	Det ←	NP [1.3]	[1.4]	NP [1]5]
	[ * ,=]	Nominal, Noun	[2 4]	Nominal
		[[2,0]	Prep	PP
			[[3,4]	[3,5] NP, Proper- Noun
				[4,5]

Book	the	flight	through	Houstor	า
S, VP, Verb; Nominal, Noun	<	S, VP, <b>←</b> X2 <b>←</b>		- S <sub>1</sub> ,VP, X2 	2
[[0,1]	Det	[U,3] NP	[U,4]	NP	
	[1,2]	[1.3]	[1,4]	[1,5]	
		Nominal, Noun		Nominal	
		[2,3]	[2,4]	[2,5]	
			Prep	PP	
			[3,4]	[3,5]	
				NP, Proper- Noun	
				[4,5]	

## Note

- An alternative is to fill a diagonal at a time.
  - That still satisfies our requirement that the component parts of each constituent/cell will already be available when it is filled in.

### **CKY Notes**

- Since it's bottom up, CKY populates the table with a lot of phantom constituents.
  - Segments that by themselves are constituents but cannot really occur in the context in which they are being suggested.
  - To avoid this we can switch to a top-down control strategy
  - Or we can add some kind of filtering that blocks constituents where they can not happen in a final analysis.

## **Earley Parsing**

- Allows arbitrary CFGs
- Top-down control
- Fills a table in a single sweep over the input
  - Table is length N+1; N is number of words
  - Table entries represent
    - Completed constituents and their locations
    - In-progress constituents
    - Predicted constituents

#### **States**

 The table-entries are called states and are represented with dotted-rules.

$$S \rightarrow \cdot VP$$
 A VP is predicted

- $NP \rightarrow Det \cdot Nominal$
- $VP \rightarrow V NP \cdot$

- An NP is in progress
- A VP has been found

### States/Locations

S → • VP [0,0]

 A VP is predicted at the start of the sentence

NP → Det • Nominal
 [1,2]

• VP  $\rightarrow$  V NP • [0,3]

- An NP is in progress; the Det goes from 1 to 2
- A VP has been found starting at 0 and ending at 3

## Earley

- As with most dynamic programming approaches, the answer is found by looking in the table in the right place.
- In this case, there should be an S state in the final column that spans from 0 to N and is complete. That is,

• S  $\rightarrow \alpha \bullet [0,N]$ 

If that's the case you're done.

## Earley

- So sweep through the table from 0 to N...
  - New predicted states are created by starting top-down from S
  - New incomplete states are created by advancing existing states as new constituents are discovered
  - New complete states are created in the same way.

## Earley

- More specifically...
  - 1. Predict all the states you can upfront
  - 2. Read a word
    - 1. Extend states based on matches
    - 2. Generate new predictions
    - 3. Go to step 2
  - 3. When you' re out of words, look at the chart to see if you have a winner

## **Core Earley Code**

function EARLEY-PARSE(words, grammar) returns chart

```
ENQUEUE((\gamma \rightarrow \bullet S, [0,0]), chart[0])
for i \leftarrow from 0 to LENGTH(words) do
 for each state in chart[i] do
   if INCOMPLETE?(state) and
            NEXT-CAT(state) is not a part of speech then
      PREDICTOR(state)
   elseif INCOMPLETE?(state) and
            NEXT-CAT(state) is a part of speech then
       SCANNER(state)
   else
      COMPLETER(state)
 end
end
return(chart)
```

## **Earley Code**

procedure PREDICTOR( $(A \rightarrow \alpha \bullet B \beta, [i, j])$ ) for each  $(B \rightarrow \gamma)$  in GRAMMAR-RULES-FOR(B, grammar) do ENQUEUE( $(B \rightarrow \bullet \gamma, [j, j])$ , chart[j]) end

**procedure** SCANNER( $(A \rightarrow \alpha \bullet B \beta, [i, j])$ ) **if** B  $\subset$  PARTS-OF-SPEECH(*word[j]*) **then** ENQUEUE( $(B \rightarrow word[j], [j, j+1]), chart[j+1]$ )

procedure COMPLETER( $(B \rightarrow \gamma \bullet, [j,k])$ ) for each  $(A \rightarrow \alpha \bullet B \beta, [i,j])$  in chart[j] do ENQUEUE( $(A \rightarrow \alpha B \bullet \beta, [i,k]), chart[k]$ ) end

- Book that flight
- We should find... an S from 0 to 3 that is a completed state...

# Chart[0]

S0	$\gamma \rightarrow \bullet S$	[0,0]	Dummy start state
S1	$S \rightarrow \bullet NP VP$	[0,0]	Predictor
S2	$S \rightarrow \bullet Aux NP VP$	[0,0]	Predictor
<b>S</b> 3	$S \rightarrow \bullet VP$	[0,0]	Predictor
S4	$NP \rightarrow \bullet Pronoun$	[0,0]	Predictor
S5	$NP \rightarrow \bullet Proper-Noun$	[0,0]	Predictor
S6	$NP \rightarrow \bullet Det Nominal$	[0,0]	Predictor
S7	$VP \rightarrow \bullet Verb$	[0,0]	Predictor
S8	$VP \rightarrow \bullet Verb NP$	[0,0]	Predictor
S9	$VP \rightarrow \bullet Verb NP PP$	[0,0]	Predictor
S10	$VP \rightarrow \bullet Verb PP$	[0,0]	Predictor
S11	$VP \rightarrow \bullet VP PP$	[0,0]	Predictor

#### Note that given a grammar, these entries are the same for all inputs; they can be pre-loaded.

## Chart[1]

S12	$Verb \rightarrow book \bullet$	[0,1]	Scanner
S13	$VP \rightarrow Verb \bullet$	[0,1]	Completer
S14	$VP \rightarrow Verb \bullet NP$	[0,1]	Completer
S15	$VP \rightarrow Verb \bullet NP PP$	[0,1]	Completer
S16	$VP \rightarrow Verb \bullet PP$	[0,1]	Completer
S17	$S \rightarrow VP \bullet$	[0,1]	Completer
S18	$VP \rightarrow VP \bullet PP$	[0,1]	Completer
S19	$NP \rightarrow \bullet Pronoun$	[1,1]	Predictor
S20	$NP \rightarrow \bullet Proper-Noun$	[1,1]	Predictor
S21	$NP \rightarrow \bullet Det Nominal$	[1,1]	Predictor
S22	$PP \rightarrow \bullet Prep NP$	[1,1]	Predictor

# Charts[2] and [3]

S23	$Det \rightarrow that \bullet$	[1,2]
S24	$NP \rightarrow Det \bullet Nominal$	[1,2]
S25	$Nominal \rightarrow \bullet Noun$	[2,2]
S26	$Nominal \rightarrow \bullet Nominal Noun$	[2,2]
S27	Nominal $\rightarrow \bullet$ Nominal PP	[2,2]
S28	Noun $\rightarrow$ flight $\bullet$	[2,3]
S29	Nominal $\rightarrow$ Noun $\bullet$	[2,3]
S30	$NP \rightarrow Det Nominal \bullet$	[1,3]
S31	$Nominal \rightarrow Nominal \bullet Noun$	[2,3]
S32	$Nominal \rightarrow Nominal \bullet PP$	[2,3]
S33	$VP \rightarrow Verb NP \bullet$	[0,3]
S34	$VP \rightarrow Verb NP \bullet PP$	[0,3]
S35	$PP \rightarrow \bullet Prep NP$	[3,3]
S36	$S \rightarrow VP \bullet$	[0,3]
S37	$VP \rightarrow VP \bullet PP$	[0,3]

Scanner

Completer

Predictor

Predictor

Predictor

Scanner

Completer

Completer

Completer

Completer

Completer

Completer

Completer

Completer

Predictor

## Efficiency

- For such a simple example, there seems to be a lot of useless stuff in there.
- Why?

• It's predicting things that aren't consistent with the input

•That's the flipside to the CKY problem.

#### Details

 As with CKY that isn't a parser until we add the backpointers so that each state knows where it came from.