Lecture 6: General and Bottom-Up Parsing

Last modified: Wed Oct 7 21:40:56 2020 CS164: Lecture #6 1

Parsing So Far (corrected slide)

- Have seen that recursive-descent parsing is a simple way to convert a grammar to a program that parses source using the grammar.
- However, the need to predict which production to take before seeing all the source tokens requires workarounds, as we've seen.
- For example, must perform *left-recursion removal*, as in

```
Original
                                LL(1) rewrite
expr ::= expr "+" term
                              expr ::= term "+" expr
       expr "-" term
                                       term "-" expr
         term ;
                                       term ;
```

• ... and *left factoring*, as in

```
Original
                                    Left-factored
                                  expr ::= term expr_tail
expr ::= term "+" expr
         term "-" expr
                                  expr_tail ::= "+" expr
                                                 "-" expr
         term ;
                                                 \epsilon ;
```

And in this last example, must adjust semantics as well.

Going Bottom Up

- So let's see what happens when we put off the decision about what production to use until after we've examined the text to be produced; this entails processing the children of a node in the parse tree before deciding on the production for that node.
- That is, we determine the parse tree from the bottom up.
- So rather than parsing e ::= e '+' t by
 - Expand the top e of the parse tree.
 - Parse the e on the left of the '+'.
 - Scan the '+'.
 - Parse the t.
- We instead do this as
 - Parse the e on the left of the '+'.
 - Scan the '+'
 - Parse the t
 - Create the top e of the parse tree.
- This order corresponds to a reverse-rightmost derivation.

A Little Notation

Here and in lectures to follow, we'll often have to refer to general productions or derivations. In these, we'll use various alphabets to mean various things:

- Capital roman letters are nonterminals (A, B, ...).
- Lower-case roman letters are terminals (or tokens, characters, etc.)
- Lower-case greek letters are sequences of zero or more terminal and nonterminal symbols, such as appear in sentential forms or on the right sides of productions (α, β, \ldots) .
- Subscripts on lower-case greek letters indicate individual symbols within them, so $\alpha = \alpha_1 \alpha_n \dots \alpha_n$ and each α_i is a single terminal or nonterminal.

So $A := \alpha$ might describe the production e := e' + t,

... and $B \Rightarrow \alpha A \gamma \Rightarrow \alpha \beta \gamma$ might describe the derivation steps $e \Rightarrow e' + t \Rightarrow e' + ID$ (α is e '+'; A is t; B is e; and γ is empty.)

Fixing Recursive Descent

- First, let's define an impractical but simple implementation of a topdown parsing routine.
- For nonterminal A and string $S=c_1c_2\ldots c_n$, we'll define parse(A, S) to return the length of a valid prefix of 5 that is derivable from A.
- That is, parse(A, $c_1c_2 \dots c_n$) = k, where $A \stackrel{*}{\Longrightarrow} c_1c_2 \cdots c_k$:

$$\underbrace{c_1c_2\dots c_k}_{A\stackrel{*}{\Longrightarrow}}c_{k+1}c_{k+2}\dots c_n$$

- ullet Of course, it's possible that A could produce other prefixes of S, and we'll have to figure out which to take.
- As a result we should expect that this parse function is non-deterministic (I did say "impractical.")

Abstract body of parse(A,S)

Can formulate top-down parsing analogously to NFAs.

```
parse (A, S):
    """Assuming A is a nonterminal and S = c_1c_2\dots c_n is a string, return
       integer k such that A can derive the prefix string c_1 \dots c_k of S."""
   Choose production 'A: \alpha_1\alpha_2\cdots\alpha_m' for A (nondeterministically)
   k = 0
   for x in \alpha_1, \alpha_2, \cdots, \alpha_m:
        if x is a terminal:
            if x == c_{k+1}:
                 k += 1
            else:
                 GIVE UP
        else:
            k += parse (x, c_{k+1} \cdots c_n)
   return k
```

- Let the start symbol be p with exactly one production: $p ::= \gamma \dashv$.
- We'll say that a call to parse returns a value if some set of choices for productions (by the blue step) would not give up (just like NFA).
- Then if parse(p, S) returns a value, S must be in the language.

Consider parsing $S = "ID*ID\dashv"$ with a grammar from last time:

```
p ::= e '⊢'
e ::= t
 | e '/' t
  | e '*' t
t ::= ID
```

Consider parsing $S = "ID*ID\dashv"$ with a grammar from last time:

```
p := e^{\gamma} - A failing path through the program:
  e ::= t
                        parse(p, S):
    l e '/' t
                           Choose p ::= e '\dashv ':
       | e '*' t
                              parse(e, S):
  t ::= ID
                                  Choose e ::= t:
                                     parse(t, S):
                                         choose t ::= ID:
                                            check S[1] == ID; OK, so k_3 += 1;
                                            return 1 (= k_3; added to k_2)
k_i means "the vari-
                                     return 1 (and add to k_1)
able k in the call to
                              Check S[2] == S[k_1+1] == '-|': NO: GIVE UP
                                                             (S[2] == '*')
parse that is nested
i deep." Outermost k
```

is k_1 .

Consider parsing $S = "ID*ID\dashv"$ with a grammar from last time:

```
p := e^{\gamma} \rightarrow A successful path through the program:
  e ::= t
                        parse(p, S):
                           Choose p ::= e '⊢':
    | e '/' t
                            parse(e, S):
       l e '*' t
                                  Choose e ::= e '*' t:
  t: := TD
                                     parse(e, S):
                                          choose e ::= t:
                                            parse(t, S):
                                                choose t := TD:
                                                   check S[1] == ID; OK, return 1
k_i means "the vari-
                                             return 1 (so k_2 += 1)
able k in the call to
                                      check S[k_2] == '*'; OK, k_2 += 1
parse that is nested
                                      parse(t, S_3): # S_3 == "ID \dashv"
i deep." Outermost k
                                          choose t ::= ID:
is k_1. Likewise for S_i.
                                             check S_3[k_3+1] == S_3[1] == ID; OK
                                            k_3+=1; return 1 (so k_2 += 1)
                                         return 3
                              Check S[k_1+1] == S[4] == '-1': OK
```

 k_1 +=1; return 4

Making a Deterministic Algorithm

- If we had an infinite supply of processors, could just spawn new ones at each "Choose" line.
- Some would give up, some loop forever, but on correct programs, at least one processor would get through.
- To do this for real (say with one processor), need to keep track of all possibilities systematically.
- This is the idea behind Earley's algorithm:
 - Handles any context-free grammar.
 - Finds all parses of any string.
 - Can recognize or reject strings in $O(N^3)$ time for ambiguous grammars, $O(N^2)$ time for "nondeterministic grammars", or O(N) time for deterministic grammars (such as accepted by Bison or CUP).

Earley's Algorithm: I

- ullet First, reformulate to use recursion instead of looping. Assume the string $S=c_1\cdots c_n$ is fixed.
- Redefine parse:

```
parse (A: \alpha \bullet \beta, s, k):

"""Assumes A: \alpha \beta is a production in the grammar,

0 \le s \le k \le n, and \alpha can produce the string c_{s+1} \cdots c_k.

Returns integer j such that \beta can produce c_{k+1} \cdots c_j."""
```

• Or diagrammatically, parse returns an integer j such that:

$$c_1 \cdots c_s \underbrace{c_{s+1} \cdots c_k}_{\alpha \stackrel{*}{\Longrightarrow}} \underbrace{c_{k+1} \cdots c_j}_{\beta \stackrel{*}{\Longrightarrow}} c_{j+1} \cdots c_n$$

ullet So if B is the start symbol, and

can return n, it means that S is recognized to be in the language.

Earley's Algorithm: II

```
parse (A ::= \alpha \bullet \beta, s, k):
    """Assumes A ::= \alpha\beta is a production in the grammar,
       0 <= s <= k <= n, and \alpha can produce the string c_{s+1}\cdots c_k.
       Returns integer j such that \beta can produce c_{k+1} \cdots c_j."""
    if \beta is empty:
       return k
   Assume \beta has the form x\delta
    if x is a terminal:
       if x == c_{k+1}:
             return parse(A ::= \alpha x \bullet \delta, s, k+1)
       else:
             GIVE UP
   else:
       Choose production 'x ::= \kappa' for x (nondeterministically)
       j = parse(x ::= \bullet \kappa, k, k)
       return parse (A ::= \alpha x \bullet \delta, s, j)
```

- Now do all possible choices that result in such a way as to avoid redundant work ("nondeterministic memoization").
- That is, if parse is called with the same three arguments as a previous call, just use the result(s) of the previous call.

Chart Parsing

- Idea is to build up a table (known as a chart) of all calls to parse that have been made
- Only one entry in chart for each distinct triple of arguments (A ::= $\alpha \bullet \beta$, s, k).
- ullet We'll organize table in columns numbered by the k parameter, so that column k represents all calls that are looking at c_{k+1} in the input.
- Each column contains entries with the other two parameters:

 $[A ::= \alpha \bullet \beta, S],$ which are called *items*.

• The columns, therefore, are item sets.

Grammar

Input String

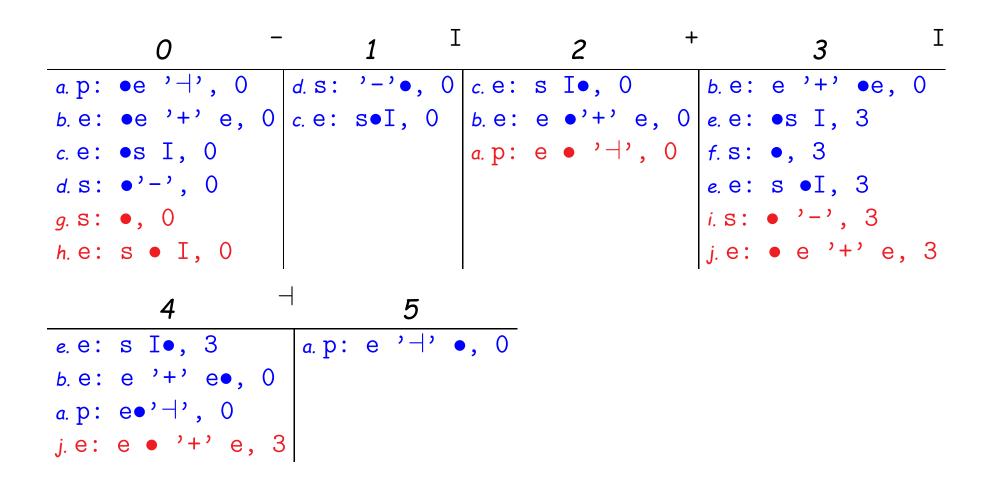
- T + T ⊢

Chart. Headings are values of k and c_{k+1} (raised symbols). Item labels (a-f) trace the "ancestry" of each item. (Have shortened ': :=' to ':' for compactness.)

0	- 1	[[] 2	+	3	Ι
a.p: •e '⊢', 0	d. s: '-'•, 0	c.e: s I•, 0	<i>b.</i> e:	e '+' •e,	0
b.e: •e '+' e, C	c.e: s•I, 0	b. e: e •'+' e,	0 <i>e.</i> e:	•s I, 3	
c.e: •s I, 0				● , 3	
d. s: •'-', 0			e. e:	s •I, 3	
4	⊣ 5				
e.e: s I•, 3	a. p: e ', -, •	, 0			
b.e: e '+' e•, C)				
a.p: e•'⊢', 0					

Example, completed

 Last slide showed only those items that survive and get used. Algorithm actually computes dead ends as well (in red).



Ambiguous Example

Grammar

Input String

 $I + I + I \dashv$

Chart. Only useful items shown.

	0]	-	1	+		2]	-	3	•	+
a. p:	•e '⊢	', 0	c. e:	I •, O		ь. e:	e '+'•	e, 0	d. e:	Ι •,	2	
b. e:	•e '+	'e, 0	b. e:	e •'+'	e, 0	d. e:	● I, 2		b. e:	e '+	, e •,	0
c. e: •	•I, O					e. e :	•e '+'	e, 2	e. e:	e •'	+, e,	2
									b. e:	e •'	+, e,	0
	4]	-	5	_		6					
b. e:	4 e '+'	•e, 0	f. e:	5 I •, 4	_	<i>a.</i> p:	6 e ⊣•,	0				
		•				<i>a</i> . p:	6 e ⊣•,	0				
e. e: (•e, 2	b. e:	I •, 4	• • , 0	a. p:	6 e ⊣•,	0				

Adding Semantic Actions

- Using syntax-directed translation to get semantic values is pretty much like recursive descent.
- The call parse (A: $\alpha \bullet \beta$, s, k) can return, in addition to j, the semantic value of the A that matches symbols $c_{s+1} \cdots c_j$.
- The value is computed during calls of the form parse (A: α' •, s, k) (i.e., where the β part is empty). For terminal symbols, value is provided by the lexer.

Adding Semantic Actions (II)

- ullet On a chart, when we see an item A: \alphaullet , s in column k, it tells us to
 - Perform the semantic action corresponding to the production A $::= \alpha$, getting a semantic value v for the left-hand side A.
 - For each item B: $\beta \bullet A\gamma$, t in column s of the chart, when adding the item B: $\beta A \bullet \gamma$, t to column k, also attach value v to that instance of A in the new item.
 - For all items derived from B: $\beta \bullet A\gamma$, t as its dot is shifted, also attach v to the same instance of A.

This step is what provides the values of nonterminals needed to compute v values (in Bison notation: \$1, \$2, etc.; in CUP notation, labels such as e1 and e2 in the rule e:=e:e1'+'e:e2).

Example with Semantic Values

```
## Compar Compar Compar Comparison of the image of the im
```

Chart. Only useful items shown. Semantic values are subscripts; red items show where they are computed.

Handling Ambiguity in Semantics (Sketch)

- Ambiguity really only important here when computing semantic actions.
- Rather than being satisfied with a single path through the chart, we look at all paths.
- The call parse (A: $\alpha \bullet \beta$, s, k) can return a set of semantic values.
- Accordingly, we attach sets of semantic values to nonterminals.

Last modified: Wed Oct 7 21:40:56 2020