Lecture 6: General and Bottom-Up Parsing

Last modified: Wed Sep 19 10:34:38 2012

Project #1 Notes

- Project involves generating an AST for Python dialect.
- Our tools provide extended BNF (BNF + regular-expression notations like '*', '+', and '?') both for context-free and lexical definitions.
- Tools also provide largely automatic AST building:
 - Tokens double as AST operators.
 - By default, each rule computes the list of all trees built by its right-hand side.
 - The 'a' notation allows you to build a tree designating the operator.
 - Or, in an action, you can use ' $\$^(...)$ ' to build an AST node, and '\$*' to denote the list of children's ASTs.
- We've also provided methods to print nodes.

Project #1 Notes (II)

In my solution, a majority of grammar rules look like this:

```
attributeref: primary "."! identifier
               \{ \$\$ = \$^(ATTRIBUTEREF, \$*); \}
```

and all the printing, etc. is taken care of.

- Dummy tokens like ATTRIBUTEREF are first defined with
 - %token ATTRIBUTEREF "@attributeref"
- In a few cases, I can just write

```
expr1 : expr1 "or"^ expr1
```

and the action is generated automatically.

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 $(\alpha, \beta, ...)$.
- 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.

For example,

- $A: \alpha$ might describe the production e: e '+' t,
- $B \Rightarrow \alpha A \gamma \Rightarrow \alpha \beta \gamma$ 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 substring derivable from A.
- That is, parse(A, $c_1c_2 \dots c_n$) = k, where

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

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
```

- Assume that the grammar contains one production for the start symbol: p: γ \dashv .
- We'll say that a call to parse returns a value if some set of choices for productions (the blue step) would return a value (just like NFA).
- Then if parse(p, S) returns a value, S must be in the language.

Consider parsing S="ID*ID→" with a grammar from last time:

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

Consider parsing S="ID*ID→" with a grammar from last time:

```
p : e^{\gamma} \rightarrow A failing path through the program:
  e:t
                    parse(p, S):
    | e '/' t
                       Choose p : e '⊢':
    l 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
                                  return 1 (and add to k_1)
variable k in the
                          Check S[2] == S[k_1+1] == '-|': GIVE UP (S[2] == '*')
call to parse that
is nested i deep."
Outermost k is
```

 k_1 .

Consider parsing S="ID*ID→" with a grammar from last time:

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

 k_1 +=1; return 4

S.

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).

Earley's Algorithm: I

- 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 <= 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_i."""
```

Or diagrammatically, parse returns an integer j such that:

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

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").

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

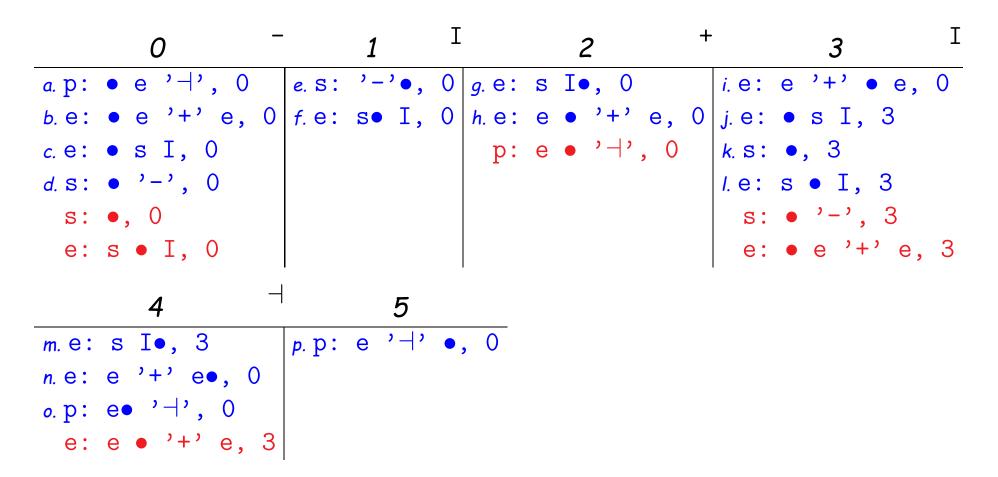
- I + I ⊢

Chart. Headings are values of k and c_{k+1} (raised symbols).

	0	_	1	Ι	2	+	3	I
_							e: e '+'	
b. e: ●6	e '+' e,	0 f.	e: s•I,	0 h. e	: e •'+'	e, 0 j.	e: ●s I,	3
c. e: ●s	s I, O					k.	s: •, 3 e: s •I,	
d. S∶ •	·-·, 0					1.	e: s •I,	3
	4	\dashv	5					
m.e: s	I•, 3	p.]	p: e '⊢	, • , 0				
n. e: e	'+' e•,	0						
o. p: e	•'⊢', 0							

Example, completed

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



Adding Semantic Actions

- 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 characters $c_{s+1} \cdots c_j$.
- This value is actually computed during calls of the form parse(A: α' •, s, k) (i.e., where the β part is empty).
- ullet Assume that we have attached these values to the nonterminals in lpha, so that they are available when computing the value for A.

Ambiguity

- Ambiguity only important here when computing semantic actions.
- Rather than being satisfied with a single path through the chart, we look at all paths.
- And we attach the set of possible results of parse(Y: $\bullet \kappa$, s, k) to the nonterminal Y in the algorithm.