

04 - Testing and Manipulating Grammars

Dr. Robert Lowe

Division of Mathematics and Computer Science
Maryville College

Outline

- 1 Grammars and Recursion
- 2 LL(1) Grammars

Sample Grammar G

For this discussion, we will be using the following grammar (found on page 39 of your textbook):

$$S \rightarrow E$$

$$E \rightarrow T \mid E + T$$

$$T \rightarrow F \mid T * F$$

$$F \rightarrow U \mid (E)$$

$$U \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

General Top Down Parsing

- 1 Look at the next symbol of input. This is the **target** symbol.
- 2 Expand the next non-terminal in the sentence.
- 3 If the target symbol does not match, backtrack and select a different non-terminal.
- 4 Keep repeating the process until there are either no non-terminal candidates or until there are no non-terminals left in the sentence.

The Backtracking Problem

- Even with the best of luck, a backtracking parser would be exponential in runtime!
- A left-recursive grammar could lead to an infinite number of candidates.
- Recall that the sample grammar in the textbook has the rule $E \rightarrow E + T$
- Consider the following expansion for the grammar from the textbook:

$$\begin{array}{ll} E + T & 1 + 2 * 3 \\ E + T + T & 1 + 2 * 3 \\ E + T + T + T & 1 + 2 * 3 \\ E + T + T + T + T & 1 + 2 * 3 \\ \dots & \end{array}$$

Ordering Recursion

- Left recursion causes problems in candidate expansions.
- Perhaps we could organize a grammar to mitigate the expansion problem.
- If we move left recursive choices to the end, maybe this would fix it!
- What if we took the grammar G and imposed the order $\langle T, F, U \rangle$ on expansions?

T	$1 + 2 * 3$
F	$1 + 2 * 3$
U	$1 + 2 * 3$
1	$1 + 2 * 3$
λ	$+2 * 3$

Mismatch! Backtrack!

Ordering Recursion (Continued)

$T * F$	$1 + 2 * 3$
$F * F$	$1 + 2 * 3$
$U * F$	$1 + 2 * 3$
$1 * F$	$1 + 2 * 3$
Mismatch! Backtrack!	
$T * F * F$	$1 + 2 * 3$
...	
And there's the loop again...	

Deterministic Parsing

- Backtracking is parsing by “brute force”.
- Backtracking essentially explores every possible production, searching for a match.
- Generally, we want parse times to be proportional to the size of the input, not exponential.
- Undoing parsing is difficult!
- We need some way to determine what production we must have based on the symbols being examined.

LL(k) Grammars

- Instead of guessing and checking, we maintain a buffer of terminals.
- If a grammar is decidable using k terminals, we call this a k -lookahead grammar.
- We can further classify the grammar by its scanning order and which production it expands first.
- An LL(k) grammar is a grammar that is scanned from left to right and expands the left most derivation.
- RL(k) scans input from right to left, expanding left-most derivations.
- LR(k) scans from left to right, expanding left-most derivations.
- All of the above have a look-ahead buffer of k terminals.
- We are really interested in LL(1) grammars.

Defining LL(1) Grammars

- Suppose we have a target expansion $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$
- We must be able to select α_j by looking at the next symbol.
- For each production, we must have a disjoint **director set** $D(A \rightarrow \alpha_j)$.
- For lookup buffer s , $A \rightarrow \alpha_j$ iff $s \in D(A \rightarrow \alpha_j)$.
- We can also have a set of symbols which immediately identify as an error if they are encountered.

Calculating Director Sets

- If $\alpha_j \xRightarrow{*} t\gamma$ for some terminal t
- Then $t \in D(A \rightarrow \alpha_j)$
- Because $A \xRightarrow{+} t\gamma$ is a valid derivation.
- Let \ll be an operator over $(N \cup T)$ such that
$$\beta \ll \alpha \iff \exists \alpha \rightarrow \beta$$
- The reflexive transitive closure \ll^* is therefore the “Can Start” relation
- The **start set** is $START(\alpha) = \beta : \beta \ll^* \alpha$
- Considering $\alpha_j = \beta_1\beta_2 \dots \beta_r$ then
$$t \in START(\beta_1) \implies t \in D(A \rightarrow \alpha_j)$$

Start Sets of G

$$\text{START}(U) = \{U, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

$$\text{START}(F) = \{\{F, () \cup \text{START}(U)\}$$

$$\text{START}(T) = \{\{T\} \cup \text{START}(F)\}$$

$$\text{START}(E) = \{\{E\} \cup \text{START}(T)\}$$

$$\text{START}(S) = \{\{S\} \cup \text{START}(E)\}$$

- Is G an LL(1) grammar?
- NO! In fact, no grammar containing left-recursive rules is LL(1)!
- $D(A \rightarrow A\gamma) \subseteq \text{START}(A)$

The First Function

- Are the start set symbols the only ones in $D(A \rightarrow \alpha_i)$?
- Extend the function START to FIRST which operates on whole strings $\beta_1\beta_2 \dots \beta_r$ over $(N \cup T)^*$ and finds terminals which can start the string.
- This function is defined recursively (where $\gamma \in (N \cup T)$ and $\delta \in (N \cup T)^*$):

$$\text{FIRST}(\lambda) = \emptyset$$

$$\text{FIRST}(\gamma\delta) = \text{terminals of START}(\gamma) \cup \text{FIRST}(\delta) \quad \text{if } \gamma \xRightarrow{*} \lambda$$

$$\text{FIRST}(\gamma\delta) = \text{terminals of START}(\gamma) \quad \text{o.w.}$$

The EMPTY Property

- We need to find if $\gamma \xRightarrow{*} \lambda$ exists.
- If it does we say $\text{EMPTY}(\gamma)$ is true.
- the EMPTY property can be defined as follows:
 - 1 If $\gamma \in T$ then $\text{EMPTY}(\gamma) = \text{false}$
 - 2 If $\gamma \in N$ then
 - 1 If $\exists \gamma \rightarrow \lambda$ then $\text{EMPTY}(\gamma) = \text{true}$
 - 2 If $\exists \gamma \rightarrow \delta_1 \dots \delta_k$ where $\forall 1 \leq i \leq k \text{ EMPTY}(\delta_i) = \text{true}$ then $\text{EMPTY}(\gamma) = \text{true}$
 - 3 For all other γ , $\text{EMPTY}(\gamma) = \text{false}$
- Let's calculate the FIRST for the productions in G .

The FOLLOW Function

- Suppose we have a terminal t in our look-ahead buffer.
- When $\alpha_j \xRightarrow{*} A$, production A is the correct choice for the parser if t can follow A .
- We calculate FOLLOW like this:
 - 1 First, calculate FINISH (the set of all terminals that can end the production)
 - 2 Next, we add the $START(\beta_i)$ for all β_i that can follow A
- Let's do this for G !
- What are the complete director sets for G ?