



CS 4240: Compilers

Lecture 22: LL(1) Parsing

Instructor: Vivek Sarkar
(vsarkar@gatech.edu)

April 10, 2019

ANNOUNCEMENTS & REMINDERS

- » Project 3 assigned on April 8th
 - » Due by 11:59pm on Tuesday, April 23rd
 - » Automatic penalty-free extension until 11:59pm on Tuesday, April 30th
 - » 10% of course grade
- » Homework 3 assigned today
 - » Due by 11:59pm on Tuesday, April 23rd
 - » Automatic penalty-free extension until 11:59pm on Friday, April 26th
 - » 5% of course grade
- » FINAL EXAM: Wednesday, May 1, 2:40 pm - 5:30 pm
 - » 30% of course grade

Worksheet # 21

Solution

(From Lecture #21 given on 4/8/2019)

1. Consider the following left recursive grammar.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \epsilon$$

Convert the above grammar to one that is suitable for LL parsing by removing left recursion from it.

• The grammar can be represented as a 4-tuple. $G = (S, NT, T, P)$

S is the start symbol

NT = {**S**, **A**} (Set of non-terminal variables)

T = {**a**, **b**, **c**, **d**} (Set of terminal variables)

P = {**S** \Rightarrow **Aa**, **S** \Rightarrow **b**, **A** \Rightarrow **Ac**, **A** \Rightarrow **Sd**, **A** \Rightarrow ϵ } (Set of productions)

1. Consider the following left recursive grammar.

$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow Ac \mid Sd \mid \epsilon \end{aligned}$$

Convert the above grammar to one that is suitable for LL parsing by removing left recursion from it.

- A grammar is left recursive if $\exists A \in NT$ such that \exists a derivation $A \Rightarrow^+ A\alpha$, for some string $\alpha \in (NT \cup T)^+$
- $S \Rightarrow Aa \Rightarrow Sda$ (indirect left recursion)
- $A \Rightarrow Ac$ (direct left recursion)

Algorithm to remove left recursion

- The algorithm removes left recursion via 2 techniques.

1. **convert indirect left recursion to direct left recursion**

2. **Rewrite direct left recursion as right recursion**

impose an order on the nonterminals, A_1, A_2, \dots, A_n

for $i \leftarrow 1$ to n do;

for $j \leftarrow 1$ to $i - 1$ do;

if \exists a production $A_i \rightarrow A_j \gamma$

*then replace $A_i \rightarrow A_j \gamma$ with one or more
productions that expand A_j*

end;

rewrite the productions to eliminate

any direct left recursion on A_i

end;

Using the algorithm : 1st iteration of outer loop

- Let's impose the order of non-terminals : **S, A**

impose an order on the nonterminals, A_1, A_2, \dots, A_n

for $i \leftarrow 1$ to n do;

for $j \leftarrow 1$ to $i - 1$ do;

if \exists a production $A_i \rightarrow A_j \gamma$

*then replace $A_i \rightarrow A_j \gamma$ with one or more
productions that expand A_j*

end;

*rewrite the productions to eliminate
any direct left recursion on A_i*

end;

$S \rightarrow Aa \mid b$

$A \rightarrow Ac \mid Sd \mid \epsilon$

Nothing is done in the inner loop.

No direct recursion on S exists.

No change is made to the grammar.

Using the algorithm : 2nd iteration of outer loop

- Imposed order of non-terminals : **S, A**

impose an order on the nonterminals, A_1, A_2, \dots, A_n

for $i \leftarrow 1$ to n do;

for $j \leftarrow 1$ to $i - 1$ do;

if \exists a production $A_i \rightarrow A_j \gamma$

*then replace $A_i \rightarrow A_j \gamma$ with one or more
productions that expand A_j*

end;

*rewrite the productions to eliminate
any direct left recursion on A_i*

end;

$S \Rightarrow Aa \mid b$

$A \Rightarrow Ac \mid \mathbf{S}d \mid \varepsilon$

In the inner loop,
production $\mathbf{A} \Rightarrow \mathbf{S}d$ is
replaced with
 $\mathbf{A} \Rightarrow \mathbf{A}ad, \mathbf{A} \Rightarrow \mathbf{b}d$.

$S \Rightarrow Aa \mid b$

$A \Rightarrow Ac \mid Aad \mid bd \mid \varepsilon$

Using the algorithm : 2nd iteration of outer loop

- Imposed order of non-terminals : $A_1 = S, A_2 = A$

impose an order on the nonterminals, A_1, A_2, \dots, A_n

for $i \leftarrow 1$ to n do;

for $j \leftarrow 1$ to $i - 1$ do;

if \exists a production $A_i \rightarrow A_j \gamma$

then replace $A_i \rightarrow A_j \gamma$ with one or more productions that expand A_j

end;

rewrite the productions to eliminate any direct left recursion on A_i

end;

$S \Rightarrow Aa \mid b$

$A \Rightarrow Ac \mid Aad \mid$

$bd \mid \varepsilon$

**We have to
rewrite direct left
recursion as
right recursion.**

Removing direct left recursion in the grammar

```
1 A_productions = [ A  $\Rightarrow$  Ac , A  $\Rightarrow$  Aad, A  $\Rightarrow$  bd, A  $\Rightarrow$   $\epsilon$  ]
2 Aprime_productions = []
3 for production in A_productions:
4     if production result starts with A:
5         move A in the production result to end of result
6         popped = A_productions.pop(production)
7         replace every A in popped to A'
8         A_prime_productions.push(popped)
9     else:
10        append A' to the production result
11 Aprime_productions.push(A'  $\Rightarrow$   $\epsilon$ )
```

$S \Rightarrow Aa \mid b$

$A \Rightarrow Ac \mid Aad \mid$

$bd \mid \epsilon$

- **Resulting grammar (right-recursive grammar)**

$S \Rightarrow Aa \mid b$

$A \Rightarrow bdA' \mid A'$

$A' \Rightarrow cA' \mid adA' \mid \epsilon$

Page 27 of slides from lecture21,
covers on how to remove direct
left recursion in a grammar.

Roadmap (Where are we?)

We set out to study parsing

- » Specifying syntax
 - Context-free grammars ✓
 - Ambiguity ✓
- » Top-down parsers
 - Algorithm & its problem with left recursion ✓
 - Left-recursion removal ✓
- » Predictive top-down parsing
 - The LL(1) condition
 - Simple recursive descent parsers
 - Table-driven LL(1) parsers

Predictive Parsing

Basic idea

Given $A \rightarrow \alpha \mid \beta$, the parser should be able to choose between α & β

FIRST sets

For some rhs $\alpha \in G$, define $\text{FIRST}(\alpha)$ as the set of tokens that appear as the first symbol in some string that derives from α

That is, $\underline{x} \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* \underline{x} \gamma$, for some γ

The LL(1) Property

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar, we would like

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$$

This would allow the parser to make a correct choice with a lookahead of exactly one symbol !



This is almost correct
See the next slide

Predictive Parsing

What about ϵ -productions?

⇒ They complicate the definition of LL(1)

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ and $\epsilon \in \text{FIRST}(\alpha)$, then we need to ensure that $\text{FIRST}(\beta)$ is disjoint from $\text{FOLLOW}(A)$, too, where

$\text{FOLLOW}(A)$ = the set of terminal symbols that can immediately follow A in a sentential form

Define $\text{FIRST}^+(A \rightarrow \alpha)$ as

- » $\text{FIRST}(\alpha) \cup \text{FOLLOW}(A)$, if $\epsilon \in \text{FIRST}(\alpha)$
- » $\text{FIRST}(\alpha)$, otherwise

Then, a grammar is LL(1) iff $A \rightarrow \alpha$ and $A \rightarrow \beta$ implies

$$\text{FIRST}^+(A \rightarrow \alpha) \cap \text{FIRST}^+(A \rightarrow \beta) = \emptyset$$

Predictive Parsing

Given a grammar that has the LL(1) property

- » Can write a simple routine to recognize each lhs
- » Code is both simple & fast

Consider $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with

$$\text{FIRST}^+(A \rightarrow \beta_i) \cap \text{FIRST}^+(A \rightarrow \beta_j) = \emptyset \text{ if } i \neq j$$

```
/* find an A */  
if (current_word  $\in$  FIRST( $A \rightarrow \beta_1$ ))  
    find a  $\beta_1$  and return true  
else if (current_word  $\in$  FIRST( $A \rightarrow \beta_2$ ))  
    find a  $\beta_2$  and return true  
else if (current_word  $\in$  FIRST( $A \rightarrow \beta_3$ ))  
    find a  $\beta_3$  and return true  
else  
    report an error and return false
```

Grammars with the LL(1) property are called predictive grammars because the parser can "predict" the correct expansion at each point in the parse.

Parsers that capitalize on the LL(1) property are called predictive parsers.

One kind of predictive parser is the recursive descent parser.

Recursive Descent Parsing

Recall the expression grammar, after transformation

1	<i>Goal</i>	\rightarrow	<i>Expr</i>
2	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>
3	<i>Expr'</i>	\rightarrow	<i>+ Term Expr'</i>
4		$ $	<i>- Term Expr'</i>
5		$ $	ϵ
6	<i>Term</i>	\rightarrow	<i>Factor Term'</i>
7	<i>Term'</i>	\rightarrow	<i>* Factor Term'</i>
8		$ $	<i>/ Factor Term'</i>
9		$ $	ϵ
10	<i>Factor</i>	\rightarrow	<u>number</u>
11		$ $	<u>id</u>
12		$ $	<u>(Expr)</u>

This produces a parser with six mutually recursive routines:

- *Goal*
- *Expr*
- *EPrime*
- *Term*
- *TPrime*
- *Factor*

Each recognizes one NT or T

The term descent refers to the direction in which the parse tree is built.

Recursive Descent Parsing (Procedural)

A couple of routines from the expression parser

Goal()

```
token ← next_token();  
if (Expr() = true & token = EOF)  
    then next compilation step;  
else  
    report syntax error;  
    return false;
```

Expr()

```
if (Term() = false)  
    then return false;  
else return Eprime();
```

looking for Number, Identifier, or
"(", found token instead, or failed
to find Expr or ")" after "("

Factor()

```
if (token = Number) then  
    token ← next_token();  
    return true;  
else if (token = Identifier) then  
    token ← next_token();  
    return true;  
else if (token = Lparen)  
    token ← next_token();  
    if (Expr() = true & token = Rparen) then  
        token ← next_token();  
        return true;  
    // fall out of if statement  
    report syntax error;  
    return false;
```

EPrime, **Term**, & **TPrime** follow the same
basic lines

Recursive Descent (Summary)

1. Massage grammar to have LL(1) condition
 - a. Remove left recursion
 - b. Build FIRST (and FOLLOW) sets
 - c. Left factor it, as needed
2. Define a procedure for each non-terminal
 - a. Implement a case for each right-hand side
 - b. Call procedures as needed for non-terminals
3. Add extra code, as needed
 - a. Perform context-sensitive checking
 - b. Build an IR to record the code

Can we automate this process?

Roadmap (Where are we?)

We set out to study parsing

- » Specifying syntax
 - Context-free grammars ✓
 - Ambiguity ✓
- » Top-down parsers
 - Algorithm & its problem with left recursion ✓
 - Left-recursion removal ✓
- » Predictive top-down parsing ✓
 - The LL(1) condition ✓
 - Simple recursive descent parsers ✓
 - First and Follow sets
 - Table-driven LL(1) parsers

FIRST and FOLLOW Sets

FIRST(α)

For some $\alpha \in (T \cup NT)^*$, define **FIRST(α)** as the set of tokens that appear as the first symbol in some string that derives from α

That is, $\underline{x} \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* \underline{x} \gamma$, for some γ

FOLLOW(A)

For some $A \in NT$, define **FOLLOW(A)** as the set of symbols that can occur immediately after A in a valid sentential form

$\text{FOLLOW}(S) = \{\text{EOF}\}$, where S is the start symbol

To build **FOLLOW** sets, we need **FIRST** sets ...

Computing FIRST Sets

```
for each  $x \in T$ ,  $\text{FIRST}(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $\text{FIRST}(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin;
       $\text{rhs} \leftarrow \text{FIRST}(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in \text{FIRST}(B_i)$  do
         $\text{rhs} \leftarrow \text{rhs} \cup (\text{FIRST}(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in \text{FIRST}(B_k)$ 
      then  $\text{rhs} \leftarrow \text{rhs} \cup \{\epsilon\}$ 
     $\text{FIRST}(A) \leftarrow \text{FIRST}(A) \cup \text{rhs}$ 
  end // for loop
end // while loop
```

Outer loop is monotone
increasing for FIRST
sets

$\rightarrow |T \cup NT \cup \epsilon|$ is
bounded, so it terminates

Inner loop is bounded
by the length of the
productions in the
grammar

Computing FOLLOW Sets

for each $A \in NT$, $FOLLOW(A) \leftarrow \emptyset$

$FOLLOW(S) \leftarrow \{EOF\}$

while (FOLLOW sets are still changing)

for each $p \in P$, of the form $A \rightarrow B_1 B_2 \dots B_k$

TRAILER $\leftarrow FOLLOW(A)$

for $i \leftarrow k$ down to 1

if $B_i \in NT$ then // domain check

$FOLLOW(B_i) \leftarrow FOLLOW(B_i) \cup TRAILER$

if $\epsilon \in FIRST(B_i)$ // add right context

then $TRAILER \leftarrow TRAILER \cup (FIRST(B_i) - \{\epsilon\})$

else $TRAILER \leftarrow FIRST(B_i)$ // no $\epsilon \Rightarrow$ no right context

else $TRAILER \leftarrow \{B_i\}$ // $B_i \in T \Rightarrow$ only 1 symbol

Example of FIRST and FOLLOW sets

1	<i>Goal</i>	\rightarrow	<i>Expr</i>
2	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>
3	<i>Expr'</i>	\rightarrow	<i>+ Term Expr'</i>
4			<i>- Term Expr'</i>
5			ϵ
6	<i>Term</i>	\rightarrow	<i>Factor Term'</i>
7	<i>Term'</i>	\rightarrow	<i>* Factor Term'</i>
8			<i>/ Factor Term'</i>
9			ϵ
10	<i>Factor</i>	\rightarrow	<u>number</u>
11			<u>id</u>
12			<u>(Expr)</u>

$\text{FIRST}^+(A \rightarrow \beta)$ is identical to $\text{FIRST}(\beta)$ except for productions 5 and 9

$\text{FIRST}^+(\text{Expr}' \rightarrow \epsilon)$ is $\{\epsilon,), \text{eof}\}$

$\text{FIRST}^+(\text{Term}' \rightarrow \epsilon)$ is $\{\epsilon, +, -,), \text{eof}\}$

Symbol	FIRST	FOLLOW
<u>num</u>	<u>num</u>	
<u>name</u>	<u>name</u>	
+	+	
-	-	
*	*	
/	/	
((
))	
<u>eof</u>	<u>eof</u>	
ϵ	ϵ	
<i>Goal</i>	<u>(, name, num</u>	eof
<i>Expr</i>	<u>(, name, num</u>), eof
<i>Expr'</i>	+, -, ϵ), eof
<i>Term</i>	<u>(, name, num</u>	+, -,), eof
<i>Term'</i>	*, /, ϵ	+, -,), eof
<i>Factor</i>	<u>(, name, num</u>	+, -, *, /,), eof

Example of FIRST⁺ sets

1	<i>Goal</i>	→	<i>Expr</i>
2	<i>Expr</i>	→	<i>Term Expr'</i>
3	<i>Expr'</i>	→	<i>+ Term Expr'</i>
4			<i>- Term Expr'</i>
5			ϵ
6	<i>Term</i>	→	<i>Factor Term'</i>
7	<i>Term'</i>	→	<i>* Factor Term'</i>
8			<i>/ Factor Term'</i>
9			ϵ
10	<i>Factor</i>	→	<u>number</u>
11			<u>id</u>
12			<u>(Expr)</u>

Prod'n	FIRST ⁺
1	(, <u>name</u> , <u>num</u>
2	(, <u>name</u> , <u>num</u>
3	+
4	-
5	ϵ , <u>,</u> , eof
6	(, <u>name</u> , <u>num</u>
7	*
8	/
9	ϵ , +, -, <u>,</u> , eof
10	number
11	name
12	(

Building Top-down Parsers

Strategy

- » Encode knowledge in a table
- » Use a standard "skeleton" parser to interpret the table

Example

- » The non-terminal *Factor* has 3 expansions
 - (*Expr*) or Identifier or Number
- » Table might look like:

1	<i>Goal</i>	→	<i>Expr</i>
2	<i>Expr</i>	→	<i>Term Expr'</i>
3	<i>Expr'</i>	→	+ <i>Term Expr'</i>
4			- <i>Term Expr'</i>
5			ϵ
6	<i>Term</i>	→	<i>Factor Term'</i>
7	<i>Term'</i>	→	* <i>Factor Term'</i>
8			/ <i>Factor Term'</i>
9			ϵ
10	<i>Factor</i>	→	<u>number</u>
11			<u>id</u>
12			(<i>Expr</i>)

» Table might look like:

	+	-	*	/	Num	Id	()	EOF
<i>Factor</i>	—	—	—	—	10	11	12	—	—

Error on `+`

Reduce by rule 10 on Num

Building Top Down Parsers

Building the complete table

- » Need an interpreter for the table (skeleton parser)
- » Need a row for every NT & a column for every T

LL(1) Skeleton Parser

```
word ← NextWord()           // Initial conditions, including
push EOF onto Stack          // a stack to track local goals
push the start symbol, S, onto Stack
TOS ← top of Stack

loop forever
  if TOS = EOF and word = EOF then
    break & report success // exit on success
  else if TOS is a terminal then
    if TOS matches word then
      pop Stack              // recognized TOS
      word ← NextWord()
    else report error looking for TOS // error exit
  else                      // TOS is a non-terminal
    if TABLE[TOS,word] is  $A \rightarrow B_1 B_2 \dots B_k$  then
      pop Stack              // get rid of A
      push  $B_k, B_{k-1}, \dots, B_1$  // in that order
    else break & report error expanding TOS

TOS ← top of Stack
```

Building Top Down Parsers

Building the complete table

- Need a row for every NT & a column for every T
- Need a table-driven interpreter for the table
- Need an algorithm to build the table

Filling in $TABLE[X,y]$, $X \in NT$, $y \in T$

1. entry is the rule $X \rightarrow \beta$, if $y \in FIRST^+(X \rightarrow \beta)$
 - entry is **error** if rule 1 does not define

If any entry has more than one rule, G is not LL(1)

This is the LL(1) table construction algorithm

Example of LL(1) Parsing Table (entry = production #)

	+	-	*	/	Num	Id	()	EOF
Goal	—	—	—	—	1	1	1	—	—
Expr	—	—	—	—	2	2	2	—	—
Expr'	3	4	—	—	—	—	—	5	5
Term	—	—	—	—	6	6	6	—	—
Term'	9	9	7	8	—	—	—	9	9
Factor	—	—	—	—	10	11	12	—	—