

Grammar Analysis

- Automated construction of syntax analysers
 - simple checks on consistency of a grammar
 - algorithms for constructing lookahead sets (used in recursive descent parsing).

Themes

1. Construction of compilers
2. Use of tools (lex and yacc) to construct compilers — "compiler-compilers"
3. Construction of compiler-compilers

A grammar

$$\text{Expression} ::= \text{Expression Operator primaryExpression} \\ | \text{primaryExpression}$$
$$\text{primaryExpression} ::= \text{Operator primaryExpression}$$

- Check that all symbols (terminal and nonterminal) can be derived from the sentence symbol.
- Check that all nonterminal symbols can generate a terminal string.

Example $G = (\{A, B, C, D, S\}, \{a, b\}, P, S)$ where P is

$$S ::= AB \mid C$$
$$A ::= a \mid aA$$
$$B ::= bA$$
$$C ::= bCa$$
$$D ::= aD \mid b$$

Nomenclature

A nonterminal symbol

α, β, γ strings of nonterminal and/or terminal symbols

X terminal or nonterminal symbol

u, v, w strings of terminal symbols.

directly generates

A grammar G defines a binary relation \rightarrow (read directly generates) on strings. This relation is the smallest relation such that, for all A, α, β, γ

$$\alpha A \beta \rightarrow \alpha \gamma \beta \equiv A ::= \gamma \text{ is a production in } G$$

generates

\rightarrow^* is the reflexive, transitive closure of \rightarrow

So

$$\alpha \rightarrow^* \beta$$

$$\equiv \alpha = \beta \vee \exists (\gamma :: \alpha \rightarrow \gamma \wedge \gamma \rightarrow^* \beta)$$

Definitions

$$\text{nonempty. } X \equiv \exists (u :: X \rightarrow^* u)$$

$$\text{derivable. } X \equiv \exists (u, v :: S \rightarrow^* u X v)$$

(ne.X abbreviates nonempty.X)

$\left. \begin{array}{l} \text{ne.a} = \text{true} \\ \text{ne.b} = \text{true} \end{array} \right\} \text{terminal symbols}$

$\text{ne.S} \Leftarrow (\text{ne.A} \wedge \text{ne.B}) \vee \text{ne.C} \quad S ::= AB|C$

$\text{ne.A} \Leftarrow \text{ne.a} \vee (\text{ne.a} \wedge \text{ne.A}) \quad A ::= a|aA$

$\text{ne.B} \Leftarrow \text{ne.b} \wedge \text{ne.A} \quad B ::= bA$

$\text{ne.C} \Leftarrow \text{ne.b} \wedge \text{ne.C} \wedge \text{ne.a} \quad C ::= bCa$

$\text{ne.D} \Leftarrow (\text{ne.a} \wedge \text{ne.D}) \vee \text{ne.b} \quad D ::= aD|b$

Note:

$\text{ne.C} \Leftarrow \text{ne.b} \wedge \text{ne.C} \wedge \text{ne.a}$

is satisfied by both

$\text{ne.C} = \text{false}$

and $\text{ne.C} = \text{true}$

The solution $\text{ne.C} = \text{false}$ is the least solution (ordering false and true by $\text{false} < \text{true}$).

Joining the Equations

Least solution is found by successive approximation:
update lhs by substituting best-known values in rhs.

```
for each nonterminal A do ne.A := false;  
for each terminal symbol t do ne.t := true;  
{ Invariant:  $\forall(X:: ne.X \Rightarrow \exists(u:: X \rightarrow^* u))$   
  Variant: #nonterminals - n  
  where  $n = \#(\text{nonterminals } A \text{ s.t. } ne.A)$  }  
repeat  
  update lhs values  
until no change in  $\#(\text{nonterminals in } A \text{ s.t. } ne.A)$ 
```

update lhs values:

```
for each production p do  
  ne.(lhs.p)  
  := ne.(lhs.p)  $\vee \forall(i: 0 \leq i < \#(\text{rhs.p})$   
    : ne.((rhs.p)i) )
```

Nomenclature

lhs.p	left hand side of production p
rhs.p	right hand side — " —
$\#(\text{rhs.p})$	length of rhs.p
$(\text{rhs.p})_i$	i th symbol (counting from 0) on rhs of p.

Example

$S ::= AB|C$ $A ::= a|aA$

$B ::= bA$ $C ::= bCa$

$D ::= aD|b$

Initially: $ne.a = true, ne.b = true$
otherwise, $ne.X = false$

1st iteration: $ne.A := true, ne.D := true$

2nd iteration: $ne.B := true$

3rd iteration: $ne.S := true$

4th iteration: no change

(Note: no. of iterations depends on ordering of productions)

Correctness

Termination: $\#(\text{nonterminals } A \text{ s.t. } ne.A)$
never decreases. (Once $ne.A = true$ it remains true)

Conditional correctness:

The invariant property is

$$ne.X \Rightarrow \exists(u :: X \rightarrow^* u)$$

(where u ranges over terminal strings)

It is easy to see that this is initially true
and is maintained by the loop body.

Claim is that on termination: $ne.X = \exists(u :: X \rightarrow^* u)$.

Sketch proof of termination claim.

Introduce count k of the number of times the loop has been executed.

Strengthen invariant to:

$$ne.X = \exists(u :: X \rightarrow^{\leq k} u)$$

Suppose loop terminates after $k+1$ iterations.

After k iterations: $\forall(X :: ne.X = \exists(u :: X \rightarrow^{\leq k} u))$

" $k+1$ " : $\forall(X :: ne.X = \exists(u :: X \rightarrow^{\leq k+1} u))$

So, since termination condition is "no change".

$$\forall(X :: \exists(u :: X \rightarrow^{\leq k} u) = \exists(u :: X \rightarrow^{\leq k+1} u))$$

Prove hence: $\forall(X :: \exists(u :: X \rightarrow^{\leq k} u) = \exists(u :: X \rightarrow^* u))$

Automatic Construction of Recursive Descent Parsers

- Basic structure of a recursive descent parser is identical to the structure of the underlying grammar.
- Only non-trivial task is the construction of the LOOKAHEAD sets needed to choose among productions.

nullable

$$\text{nullable}(\alpha) \equiv \alpha \rightarrow^* \varepsilon$$

Equations

$$\text{nullable}(\varepsilon) \equiv \text{true}$$

$$\text{nullable}(A) \equiv \exists(p: \text{lhs}.p = A : \text{nullable}(\text{rhs}.p))$$

$$\text{nullable}(\varepsilon) \equiv \text{true}$$

$$\text{nullable}(X\alpha) \equiv \text{nullable}(X) \wedge \text{nullable}(\alpha)$$

FIRST (starters — Watt and Brown)

$$\text{FIRST}(\alpha) = \{ t \mid \exists(u :: \alpha \rightarrow^* tu) \}$$

Equations

$$\text{FIRST}(\varepsilon) = \{ \varepsilon \}$$

$$\text{FIRST}(A) = \cup(p: \text{lhs}.p = A : \text{FIRST}(\text{rhs}.p))$$

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

$$\text{FIRST}(X\alpha) = \text{FIRST}(X) \cup (\text{if } X \rightarrow^* \varepsilon \\ \text{then } \text{FIRST}(\alpha) \text{ else } \emptyset)$$

FOLLOW

$$\text{FOLLOW}(A) = \{t \mid \exists(u, v :: Z \rightarrow^* uAtv)\}$$

Equations

An rhs occurrence of A , written $B \rightarrow \alpha A \cdot \beta$, consists of a production of the form $B \rightarrow \alpha A \beta$ together with information on the position of A on the rhs (indicated by the dot).

$$\text{FOLLOW}(S) = \{\$ \} \cup \bigcup (\text{rhs occurrences } B \rightarrow \alpha S \cdot \beta \\ :: \text{FIRST}(\beta \cdot \text{FOLLOW}(B)))$$

$$\text{FOLLOW}(A) = \bigcup (\text{rhs occurrences } B \rightarrow \alpha A \cdot \beta \\ :: \text{FIRST}(\beta \cdot \text{FOLLOW}(B)))$$

LOOKAHEAD

$$\text{LOOKAHEAD}(A ::= \alpha)$$

$$= \{t \mid \exists(u, \beta : Z \rightarrow^* uA\beta : t \in \text{FIRST}(\alpha\beta))\}$$

Equations

$$\text{LOOKAHEAD}(A ::= \alpha)$$

$$= \text{FIRST}(\alpha \cdot \text{FOLLOW}(A))$$

Solving the Equations

In each case, use successive approximation

initialise approximate solution ;

repeat

update lhs of equation

using best known approximation to rhs

until no change in approximations

A nonterminal

t terminal

X, Y terminal or non-terminal

p production

lhs.p lefthand side of production p

rhs.p righthand side of production p

u, v string of terminals

α, β string of terminals and/or non-terminals

Assume grammar has a production $Z \rightarrow S\$$

where S is the start symbol. "\$" marks

the right end of the input. Z is not included

in the calculation of FOLLOW sets.

Nullable

for each X do nullable. X := false ;

{ Invariant: $\forall(X :: \text{nullable}.X \Rightarrow X \rightarrow^* \varepsilon)$

Variant: $\#(X :: \neg \text{nullable}.X)$ }

repeat

for each p do

begin let rhs.p = X_1, X_2, \dots, X_n , lhs.p = A ;

nullable.A :=

nullable.A $\vee \forall(i: 1 \leq i \leq n: \text{nullable}.X_i)$

end

until no change occurs in nullable

FIRST

Initialisation

for each t do FIRST. t := $\{t\}$;

foreach A do FIRST. A := \emptyset

Loop body

for each p do

begin let lhs.p = A , rhs.p = X_1, \dots, X_n ;

for $i := 1$ to n do

begin FIRST. A := FIRST. $A \cup$ FIRST. X_i ;

if $\neg \text{nullable}.X_i$ then break

end

end

FOLLOW

Initialisation: for each X do FOLLOW.X := \emptyset ;
FOLLOW.S := { $\$$ } /* end-of-file marker */

Loop body: for each p do
begin let $A = \text{lhs.p}$, $X_1 \dots X_n = \text{rhs.p}$;
 $f := \text{FOLLOW.A}$; $k := n$;
{Invariant: $f = \text{FIRST}(X_{k+1} \dots X_n \cdot \text{FOLLOW.A})$ }
while $k \geq 1$ do
begin FOLLOW.X_k := FOLLOW.X_k \cup f ;
if nullable.X_k then $f := f \cup \text{FIRST.X}_k$
else $f := \text{FIRST.X}_k$;
 $k := k - 1$
end
end