

Grammars and Parsing

Article #14

Derivations (Text § 4.2)

We look at a description of a grammar as a way of **deriving** a language and how to design a grammar suitable for a recursive descent parser:

- our objective is a recursive descent parser where the parse tree is not actually generated but is implicit in the derivations
- a recursive descent parser can only work for certain types of grammars with the right kind of productions.

The derivation view is that the sentences of the language are generated by repeated application of the productions.

In general a string α can consist of a mixture of terminals and nonterminals. A production that derives the string α from the nonterminal A has the form

$$A \rightarrow \alpha$$

Since strings are derived from other strings, a series of productions that generates α_n from α_1 is written:

$$\alpha_1 \Rightarrow \alpha_2 \Rightarrow \alpha_3 \Rightarrow \dots \Rightarrow \alpha_n$$

The symbol \Rightarrow means “generates in one step”

To summarize the above equation we can use the symbol \Rightarrow^* meaning “generates in 0 or more steps”. Thus

$$\alpha_1 \Rightarrow^* \alpha_n \quad \mathbf{1.}$$

In general a string α is in **sentential** form if it can be generated by the grammar:

$$S \Rightarrow^* \alpha$$

If a string ω contains only terminals, then it must have been derived from the start symbol by at least one production. We write this as:

$$S \Rightarrow^+ \omega$$

This equation defines the grammar. Such a string is a **sentence** in the language.

Derives in one step. The metasymbol for this is \Rightarrow .

We say a string v **derives** a string w if a production can be used to produce w from v . So

$$v \Rightarrow w$$

if we can write for some strings v and w (where x and z are other strings of terminals and nonterminals)

$$v = xYz, \quad w = xYz$$

where $Y \rightarrow y$ is a production in the grammar that produces the string y from the nonterminal Y .

If the strings x and v are empty, then we have

$$Y \Rightarrow y$$

Derives in one or more steps. The metasymbol for this is \Rightarrow^+

We say a string v derives a string w in one or more steps if w is produced from v after applying one or many productions:

$$v \Rightarrow w_0 \Rightarrow w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n = w$$

summarised as:

$$v \Rightarrow^+ w$$

Derives in zero or more steps The metasymbol for this is \Rightarrow^*

This is \Rightarrow^+ together with no derivation:

$$v \Rightarrow^* w \text{ if } v \Rightarrow^+ w \text{ or } v = w$$

Finally, we arrive at another way of thinking of a language.

Let $G[V_t, V_n, P, S]$ be a grammar (S is the start symbol).

- A string v is called a **sentential form** if v is derivable from S :

$$S \Rightarrow^* v$$

- A string is a **sentence** if v consists only of terminals $v \in V_t^+$
- The language $L(G[V_t, V_n, P, S])$ is therefore the set of sentences produced by the grammar:

$$L(G[V_t, V_n, P, S]) = \{v \mid S \Rightarrow^* v \text{ and } v \text{ in } T\}$$

A language is the set of terminal strings (sentences) that is produced by the grammar. The sentences are a subset of all possible terminal strings.

]

leftmost and rightmost derivations

Consider the following grammar:

$$\begin{array}{lcl} \text{expr} & \rightarrow & \text{expr op expr} \mid (\text{expr}) \mid - \text{expr} \mid \text{id} \\ \text{op} & \rightarrow & + \mid - \mid * \mid / \end{array}$$

This first line is productions for **expr**. The second line is the productions for **op**. **expr** and **op** are nonterminals. **id** and **+**, **-**, ***** and **/** are terminals.

Consider this example from the grammar.

The string $-(x + y)$ is a sentence of the grammar because:

leftmost derivation

$$\begin{array}{l} \text{expr} \Rightarrow -\text{expr} \Rightarrow -(\text{expr}) \Rightarrow -(\mathbf{expr} \text{ op expr}) \Rightarrow -(\mathbf{id} \text{ op expr}) \\ \Rightarrow -(\text{x op expr}) \Rightarrow -(\text{x} + \text{expr}) \Rightarrow -(\text{x} + \mathbf{id}) \Rightarrow -(\text{x} + \text{y}) \end{array}$$

or $\text{expr} \Rightarrow^* -(\text{x} + \text{y})$

This is **leftmost** because the leftmost nonterminal is replaced at each stage.

The alternative is to replace the rightmost nonterminal at each stage:

rightmost derivation

$$\begin{array}{l} \text{expr} \Rightarrow -\text{expr} \Rightarrow -(\text{expr}) \Rightarrow -(\text{expr op} \mathbf{expr}) \Rightarrow -(\text{expr op} \mathbf{id}) \\ \Rightarrow -(\text{expr op} \mathbf{y}) \Rightarrow -(\text{expr} + \text{y}) \Rightarrow -(\mathbf{id} + \text{y}) \Rightarrow -(\text{x} + \text{y}) \end{array}$$

leftmost and rightmost derivations generate the same parse tree and every parse tree has associated with it a unique leftmost and unique rightmost derivation.

(A type of parsing algorithm LL(1) that we will meet later takes its name from the fact that it reads the input token string from **Left** (first L) to **Right**, does a **Leftmost** (second L) derivation, and looks one token ahead to determine the production to use)

Left and Right Associativity and Left and Right Recursion (Old Dragon Text § 2.2, New Dragon Text § 4.3.2)

Consider the following strings: $9 - 5 - 2$ and $a = b = c$

1. $9 - 5 - 2$

The $-$ operator is **left associative** since **the 5 associates with the $-$ on its left**.
The normal interpretation of the expression is:

$$(9 - 5) - 2 = 2$$

not

$$9 - (5 - 2) = 6 \quad \text{wrong}$$

2. $a = b = c$

The $=$ operator is right associative since the b associates with the $=$ on its right
The normal interpretation of the expression is:

$$a = (b = c) \quad \text{c is assigned to b, then to a}$$

not

$$(a = b) = c \quad \text{b may be undefined}$$

How can this associativity be represented in grammars?

Consider two grammars that implement this associativity:

1. left associativity for strings of digits with arithmetic operators

list \rightarrow **list** + digit | **list** - digit | digit
digit \rightarrow [0-9]

Left associativity is enforced by evaluating left associating expressions lower down in the parse tree where precedence is higher.

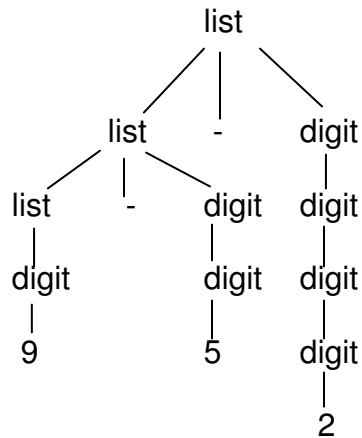
2. right associativity for strings of characters with assignment

right \rightarrow letter = **right** | letter
letter \rightarrow [a-z]

Right associativity is enforced by evaluating right associating expressions lower down in the parse tree where precedence is higher.

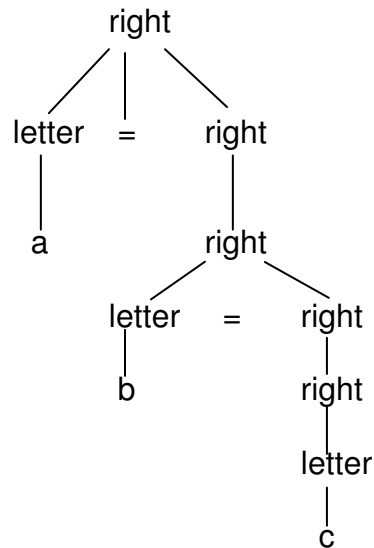
Parse trees generated for these expressions are (using a leftmost derivation and one derivation at each stage)):

1. **9 – 5 – 2**



Note how the left associative tree expands to the left

2. **a = b = c**



Note how the right associative tree grows down to the right.

Left and Right Recursion (Text § 2.4, 4.3)

Left-associative grammars naturally support the associativity of some operators, but there is a serious problem with recursive-descent parsers. Recursive descent parsers implement the productions as function calls. So if the nonterminal on the LHS is also the first symbol on the RHS, the function will enter infinite recursion.

In general productions look like:

left recursion $A \rightarrow A \alpha \mid \beta$ 1.

right recursion $A \rightarrow \alpha A \mid \beta$ 2.

where α and β are sequences of terminals and nonterminals that do not begin with A .

(In general a grammar is left-recursive if it has a nonterminal such that there is a derivation:

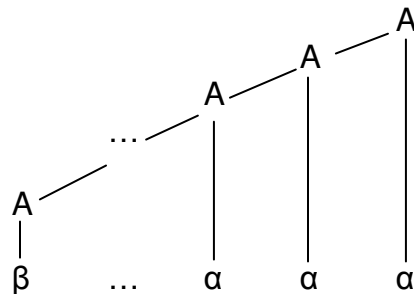
$$A \Rightarrow^+ A \alpha \quad)$$

The strings generated by these productions are of the form:

1.

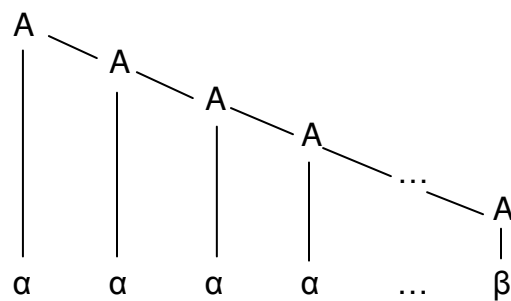
$$A \Rightarrow A \alpha \Rightarrow A \alpha \alpha \Rightarrow A \alpha \alpha \alpha \Rightarrow \dots \beta \alpha \alpha \dots \alpha$$

Parse tree:



2.

$$A \Rightarrow \alpha A \Rightarrow \alpha \alpha A \Rightarrow \alpha \alpha \alpha A \Rightarrow \dots \alpha \alpha \dots \alpha \beta$$



Note that for 1. a recursive-descent parser that attempts to apply the production as a function call would result in the evaluation of the first symbol on RHS which is the same function call, and the beginning of infinite recursion.

Hence in such a parser, we need to remove left recursion but implement the code so as to preserve left associativity.

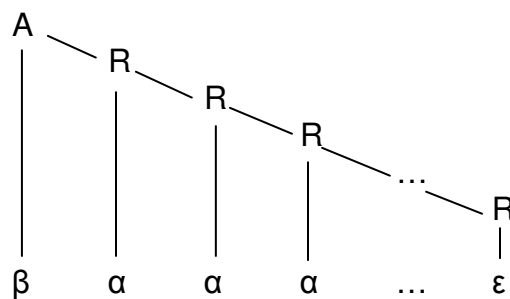
Eliminating Left Recursion (Text § 2.4, 4.3)

We can convert the left recursion (that follows naturally from the grammar) into a right recursion with the following algorithm that can be applied to any production and does not change the set of derivable strings.

original: $A \rightarrow A \alpha \mid \beta$

converted: $A \rightarrow \beta R$
 $R \rightarrow \alpha R \mid \epsilon$

Note that we now have right-recursion and an ϵ -production, but the strings generated are the same as the left recursive case:



Example.

Take the left-recursive part of the grammar above

$\text{list} \rightarrow \text{list} + \text{digit} \mid \text{list} - \text{digit} \mid \text{digit}$

as a set of separate productions:

$\text{list} \rightarrow \text{list} + \text{digit} \mid \text{digit}$
 $\text{list} \rightarrow \text{list} - \text{digit} \mid \text{digit}$

For a single alternation

list \rightarrow list + digit | digit
 with left recursion removed
 list \rightarrow digit list'
 list' \rightarrow + digit list' | ϵ

For both alternations

list \rightarrow list + digit | digit
 list \rightarrow list – digit | digit
 with left recursion removed
 list \rightarrow digit list'
 list' \rightarrow + digit list' | – digit list' | ϵ

The general rule is:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid A\alpha_3 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

Note: these are examples of **immediate** left recursion.
 The recursion could be several steps later:

$S \rightarrow A a \mid b$
 $A \rightarrow A c \mid S d \mid \epsilon$

Left Factoring a Grammar (Text § 2.4, 4.3)

The grammar

relational expression \rightarrow primary expression > primary expression
 | primary expression < primary expression

is not suitable for predictive parsing. It is not possible to make the decision which of the alternatives to choose having one token lookahead because all production alternatives have the same left-most terminal. It is possible to rework the grammar so that the decision can be deferred until the parser can make the right choice. The corresponding grammar transformation is called **left factoring**.

The general algorithm for left-factoring a grammar is

Given a grammar

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$$

if α is not ϵ the left factored grammar is

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

This transformation must be applied repeatedly until no alternatives for a nonterminal have a common prefix.

In many cases the grammar may be transformed without applying the general rule. For example the grammar above can be rewritten as

relational expression \rightarrow
primary expression relational operator primary expression

relational operator $\rightarrow < \mid >$

Building the FIRST and FOLLOW sets for a Grammar (Text § 4.4)

The construction of a predictive parser requires two set of tokens to be built.

The **FIRST** set is a set of tokens that appear at the left-most position after zero or more derivations are applied to a grammar production. The formal definition is

$$\text{FIRST}(A) = \{ a \mid A \Rightarrow^* a\alpha \}$$

The following rules for computing the **FIRST** set can be derived from the formal definition.

BFTS1. If X is a terminal, then $\text{FIRST}(X) = \{X\}$

BFTS2. If X is a nonterminal and $X \rightarrow \epsilon$ is a production, then add ϵ to $\text{FIRST}(X)$.

BFTS3. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production then
 $\text{FIRST}(X) = \{\text{FIRST}(Y_1)\}$ if Y_1 does not derive ϵ ;
 $\text{FIRST}(X) = \{\text{FIRST}(Y_1), \text{FIRST}(Y_2)\}$ if Y_1 derives ϵ ;
 $\text{FIRST}(X) = \{\text{FIRST}(Y_1), \text{FIRST}(Y_2), \text{FIRST}(Y_3)\}$ if Y_1 and Y_2 derive ϵ ;
 $\text{FIRST}(X) = \{\text{FIRST}(Y_1), \text{FIRST}(Y_2), \text{FIRST}(Y_3), \dots, \text{FIRST}(Y_i)\}$ if ϵ is in all $\text{FIRST}(Y_1), \text{FIRST}(Y_2), \text{FIRST}(Y_3), \dots, \text{FIRST}(Y_{i-1})$;
 $\text{FIRST}(X) = \{\text{FIRST}(Y_1), \text{FIRST}(Y_2), \text{FIRST}(Y_3), \dots, \text{FIRST}(Y_k), \epsilon\}$ if ϵ is in all $\text{FIRST}(Y_1), \text{FIRST}(Y_2), \text{FIRST}(Y_3), \dots, \text{FIRST}(Y_k)$;

BFTS4. If X is a nonterminal and $X \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_m$ is a production then
 $\text{FIRST}(X) = \{\text{FIRST}(\alpha_1), \text{FIRST}(\alpha_2), \dots, \text{FIRST}(\alpha_m)\}$.

The **FOLLOW** set is a set of tokens that appear on the right side of a nonterminal after zero or more derivations a grammar production. The formal definition is

$$\text{FOLLOW}(B) = \{ a \mid A \Rightarrow \alpha B a \gamma \}$$

The following rules for computing the **FOLLOW** set can be derived from the formal definition.

BFWS0. The **FOLLOW** set can not contain ϵ .

BFWS1. If **S** is the start symbol (nonterminal) for a grammar place **\$** in the **FOLLOW(S)**, where **\$** is the input end-marker (for example, end of file)

BFWS2. If there is production $A \rightarrow \alpha B a \mid \alpha B b$, then **FOLLOW(B) = {a,b}**

BFWS3. If there is production $A \rightarrow \alpha B \beta$ and β does not derive ϵ , then add **FIRST(β)** to the **FOLLOW(B)** set.

BFWS4. If there is production $A \rightarrow \alpha B$, then add **FOLLOW(A)** to the **FOLLOW(B)** set

BFWS5. If there is production $A \rightarrow \alpha B \beta$ and β derives ϵ , then add **FOLLOW(A)** to the **FOLLOW(B)** set.