



GRAMMAR

- 一般之語言，以數學形式定義如下：
- 定義一： $L(T, N, P, S)$ 為一語言。
 - T：所有終端符號(Terminal symbol)之集合。
 - N：所有非終端符號(Non-terminal symbol)之集合。
 - P：所有造句程序(Production)之集合。
 - S：所有起始符號(Start symbol)之集合。
- 定義二：
 - T^* ：所有屬於T之終端符號所組成字串之集合。
 - ξ ：由起始符號S一造句程序P所產生的字串，屬於 T^* 。

■ 定義三：

若且唯若存在一組 $\lambda_1, \lambda_2, \dots, \lambda_{n-1}$ 對任一 λ_i 可以由 λ_{i-1} 直接產生，則 λ_n 可由 λ_0 產生。

■ 定義四：

(a) $\xi = \alpha \xi' \beta$

(b) $\eta = \alpha \eta' \beta$

(c) 存在一組造句程序 p ， $\xi' ::= \eta'$

若且唯若存在符號串合乎上述三條件時，則可由直接產生，即 $\xi \rightarrow \eta$ 。

- In general, a grammar involves four quantities: *terminals*, *nonterminals*, a *start symbol*, and *productions*.
- The basic symbols of which strings in the language are composed we shall call *terminals*. The word “token” is a synonym for “terminal” when we are talking about programming languages. In the example above, certain keywords, such as **begin** and **else** are terminals, as are punctuation symbols such as ‘;’ and operators such as ‘+’.

- Nonterminals are special symbols that denote sets of strings. The terms “syntactic variable” and “syntactic category” are synonyms for “nonterminal.” In the examples above, the syntactic categories statement, expression, and statement-list are nonterminals; each denotes a set of strings. One nonterminal is selected as the start symbol, and it denotes the language in which we are truly interested. The other nonterminals are used to define other sets of strings, and these help define the language. They also help provide a hierarchical structure for the language at hand.

- The productions (rewriting rules) define the ways in which the syntactic categories may be built up from one another and from the terminals. Each production consists of a nonterminal, followed by a string of nonterminals and terminals.

statement \rightarrow **begin** statement-list **end**

statement-list \rightarrow statement

statement-list \rightarrow statement ; statement-list

(例一)英語簡單語句由下列三條造句程序組成。

P1<sentence> ::= <subject><predicate>.

P2<subject> ::= lions | cats

P3<predicate> ::= cry | fly

P1描述一個語句<sentence>定義為非終端符號主詞<subject>之後，緊跟著另一非終端符號述詞<predicate>，終端符號。句點。P2說明非終端符號主詞<subject>定義為終端符號(terminal symbol) lions或終端符號cats。P3說明非終端符號<predicate>定義為終端符號cry或fly。

- 因此合乎簡單語句<sentence>語法的與具有下列四句是正確，其他為語法錯誤。

語句一

lions cry.

語句二

lions fly.

語句三

cats cry.

語句四

cats fly.

- 以例一說明。 $T = \{\text{lions, cats, cry, fly, .}\}$ 。 $N = \{\langle \text{sentence} \rangle, \langle \text{subject} \rangle, \langle \text{predicate} \rangle\}$ 。 $P = \{P_1, P_2, P_3\}$ 。 $S = \{\langle \text{sentence} \rangle\}$ 。 $L = \text{英語簡單語句}$ 。
由定義四及 P_1 可知 為 $\langle \text{sentence} \rangle$ ， 為 η'
 $\langle \text{subject} \rangle \langle \text{predicate} \rangle$ ，即 $\langle \text{subject} \rangle \langle \text{predicate} \rangle$
可由 $\langle \text{sentence} \rangle$ 直接產生，寫成
 $\langle \text{sentence} \rangle \rightarrow \langle \text{subject} \rangle \langle \text{predicate} \rangle.$

- 同理

$\langle \text{subject} \rangle \rightarrow \text{lions}$

$\langle \text{predicate} \rangle \rightarrow \text{cry}$

故

*

$\langle \text{sentence} \rangle \rightarrow \text{lions cry}.$

- 因 $\langle \text{sentence} \rangle \xrightarrow{*} \text{lions cry}$. 依據定義二知其為英語簡單語句中之一個語句，也就是說 lions cry . 語法正確。由例一可以看出非終端符號 $\langle \text{subject} \rangle$ 和 $\langle \text{predicate} \rangle$ 僅能存在於轉換過程，當轉換完成時所有的符號均為終端符號。依語法規則造出語句，所以稱這些語法規則為造句程序 (production)，習慣上以 p 表示。

語言表示法

- 每一種語言(language)均有一套語法(syntax)。合乎語法為正確。描述語法最常用的為EBNF法(Extended Backus Naur Form)。BNF因描述程式語言ALGOL 60而著名，後來擴展到其他語言的描述，例如英文語句(sentence)可以表示為：

$\langle \text{sentence} \rangle ::= \langle \text{subject} \rangle \langle \text{predicate} \rangle$

- 像BNF用來描述另一種語言，稱為類語言(metalanguage)，其使用的符號為類符號(metasymbol)。BNF的類符號有三種：

- (1) ::= 定義為。
- (2) | 或。
- (3) <> 非終端符號(non-terminal symbol)。

- 後來為了描述方便起見，增加數種：

- (4) { } 重複發生0次，1次，2次，...，多次
- (5) [] 選擇。方括號內之結構可以選用，也可以不選用。
- (6) \ \ 在背斜線(back slash)中之結構必須選一種。

以上六種類符號構成了EBNF類語言。

語句分析

- 辨認語句是依輸入的語句，導出產生該語句的造句程序。一般稱此工作為剖析(parsing)，其工作通常複雜而困難，有時甚至不可能做到。而用那一類型的造句程序來定義語言是影響辨認工作是否困難的主要原因。常用的為由上而下(top-down)的剖析方法，其法由起始符號找出產生整個語句的一系列造句程序。例如要剖析lions cry. 是否合乎英語簡單語句，從起始符號<subject>開始，由 P_2 得知<subject>可以直接產生lions，由 P_3 的起始符號<predicate>直接產生cry，故知lions cry. 為英語簡單語句。

以剖析樹(parsing tree)表示如下：

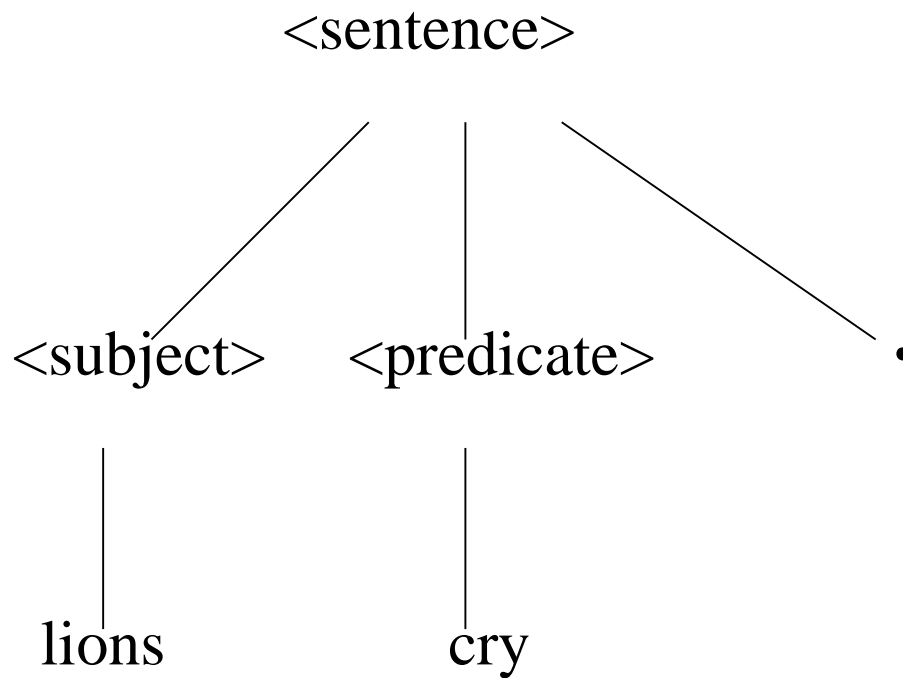


圖2.1 剖析樹。說明lions cry . 之剖析情形

- (例二) S語言以EBNF定義如下，以小寫字母表示者為終端符號。

$P_1 \quad \langle S \rangle ::= x \langle A \rangle$

$P_2 \quad \langle A \rangle ::= z \mid y \langle A \rangle$

請以由上而下的剖析方法剖析xyz是否合乎S語言。

(解)

$\langle S \rangle$	xyz	
$x \langle A \rangle$	xyz	依造句程序 P_1
$\langle A \rangle$	yz	
$y \langle A \rangle$	yz	依造句程序 P_2
$\langle A \rangle$	z	
z		依造句程序 P_2

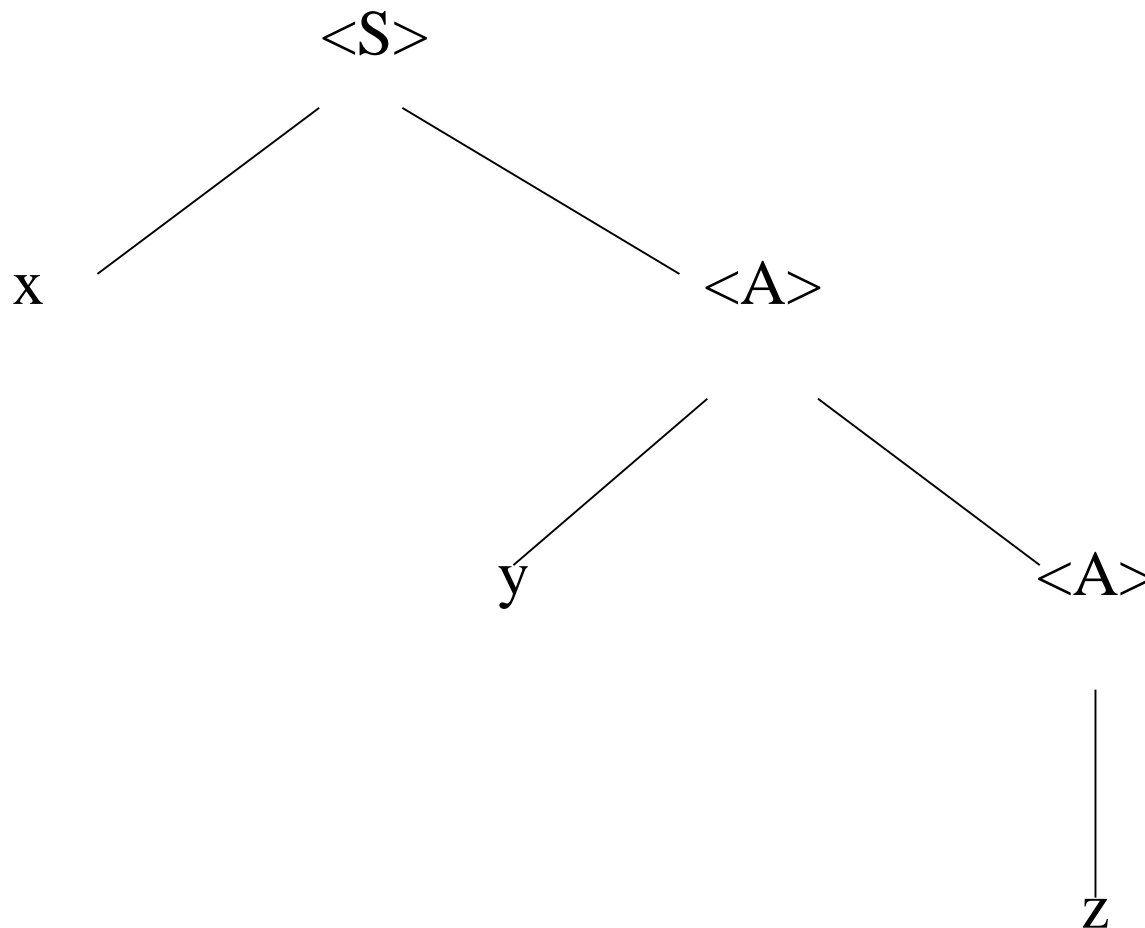


圖2-2 由上而下剖析xyz是否合乎S語言

- (例三) 有一T語言，以EBNF表示如下。

$P_1 \quad \langle T \rangle ::= \langle A \rangle \mid \langle B \rangle$

$P_2 \quad \langle A \rangle ::= x \langle A \rangle \mid y$

$P_3 \quad \langle B \rangle ::= x \langle B \rangle \mid z$

請以由上而下方法剖析xxz使否合乎T語言。

(解)

$\langle T \rangle$	xxz	$\langle T \rangle$	xxz
$\langle A \rangle$	xxz	$\langle B \rangle$	xxz
$x\langle A \rangle$	xxz	$x\langle B \rangle$	xxz
$\langle A \rangle$	xz	$\langle B \rangle$	xz
$x\langle A \rangle$	xz	$x\langle B \rangle$	xz
$\langle A \rangle$	z	$\langle B \rangle$	z
?		z	

- 在剖析過程中， $\langle T \rangle$ 之轉換，若選 $\langle A \rangle$ ，發現xxz不合乎T語言，若選 $\langle B \rangle$ ，發現xxz合乎T語言。問題在造句程序 P_2 及 P_3 之起始符號均為x，以致無法判斷鑰匙用那一造句程序來剖析，這種情形應該避免。

- 規定甲：

若有一造句程序

$$\langle A \rangle ::= \xi_1 \mid \xi_2 \mid \dots \mid \xi_n$$

則所有能從產生語句的起始符號之交集必須是空集合。

即 $start(\xi_i) \cap start(\xi_j) = \phi$

但 $i, j = 1, 2, \dots, n$

且 $i \neq j$

- 集合 $start(\xi)$ 是由所有能從產生語句的起始符號所組成，可由下列兩個規則得到。

- 規則一：若起始符號就是終端符號，則為該終端符號之集合。即

$$start(x\xi) = \{x\}$$

- 規則二：若起始符號為非終端符號，其造句程序為

$$\langle A \rangle ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

則

$$start(\langle A \rangle) = start(\alpha_1) \cup start(\alpha_2) \cup \dots \cup start(\alpha_n)$$

- 在例三中，終端符號 x 同時屬於 $\text{start}(\langle A \rangle)$ 和 $\text{start}(\langle B \rangle)$ ，因此其交集不是空集合，違反規定甲。例三的語法若改為

$$\langle S \rangle ::= \langle C \rangle \mid x \langle S \rangle$$
$$\langle C \rangle ::= y \mid z$$

則符合規定甲。但還有一些問題必須克服，請看例四。

■ (例四) 有一語言U以EBNF表示

$$\langle U \rangle ::= \langle A \rangle x$$

$$\langle A \rangle ::= x \mid \varepsilon$$

請剖析x是否為U語言之語句。

(解)

$\langle U \rangle$	x	$\langle U \rangle$	x
$\langle A \rangle x$	x	$\langle A \rangle x$	x
x x	x	εx	x
x	?	x	x

- 剖析得兩種不同的結果，問題在於 $\langle A \rangle$ 選x或 ε 。稱為空符號串問題，這種現象僅有在非終端符號能產生空符號串時才發生，為了避免這種問題的發生，必須另有規定。

■ 規定乙：

對於任一非終端符號，且具有產生空符號串之能力時，則所有 $\langle A \rangle$ 之起始符號集合，與所有 $\langle A \rangle$ 之跟隨符號集合之交集必須為空集合。即

$$start(\langle A \rangle) \cap follow(\langle A \rangle) = \phi$$

關於集合 $follow(\langle A \rangle)$ 可依下述方法求得。對於任一造句程序

$$\langle B \rangle ::= \xi \langle A \rangle \eta$$

若集合 $S_i = start(\eta_i)$ ，則集合 $follow(\langle A \rangle)$ 為所有 S_i 之聯集，即

$$follow(\langle A \rangle) = start(\eta_1) \cup start(\eta_2) \cup \dots \cup start(\eta_n)$$

- 若存在一 η_i 會產生空符號串時，則集合 $\text{follow}()$ 必包含於集合 $\text{follow}(<A>)$ 之內。即

$$\text{follow}() \subset \text{follow}(<A>)$$

- 在例四之中， $<A>$ 違反規定乙，因為

$$\text{start}(<A>) = \text{follow}(<A>) = \{x\}$$

- 若一語言有表示重複符號串之需求時，通常使用遞迴 (recursive) 結構表示。例如

$$\langle A \rangle ::= \langle B \rangle \mid \langle A \rangle \langle B \rangle$$

此造句程序將產生如下之符號串

$$\langle B \rangle, \langle B \rangle \langle B \rangle, \langle B \rangle \langle B \rangle \langle B \rangle, \dots$$

但此造句程序違反規定甲，因為

$$\begin{aligned} & start(\langle B \rangle) \cap start(\langle A \rangle \langle B \rangle) \\ & = start(\langle B \rangle) \neq \phi \end{aligned}$$

若修改為

$$\langle A \rangle ::= \varepsilon \mid \langle A \rangle \langle B \rangle$$

將產生如下之符號串

$$\varepsilon, \langle B \rangle, \langle B \rangle \langle B \rangle, \langle B \rangle \langle B \rangle \langle B \rangle, \dots$$

卻又違反規定乙，因為

$$\begin{aligned} & start(\langle A \rangle) \cap follow(\langle A \rangle) \\ &= \{\varepsilon, start(\langle B \rangle)\} \cap start(\langle B \rangle) \\ &= start(\langle B \rangle) \neq \phi \end{aligned}$$

因此可知由於規定甲和規定乙的限制，無法使用左遞迴結構來描述語言。有二種方法可以克服困難，其一為改用右遞迴結構，如

$$\langle A \rangle ::= \varepsilon \mid \langle B \rangle \langle A \rangle$$

- 其二為使用EBNF中的 $\{\langle B \rangle\}$ 表示重複結構。
- 除了語法外，語意semantic也是很重要的，如例一中的 lions fly .語法正確，但語意卻有問題。

【例二】多重選擇題

The following grammar $N = \{ B, \Sigma \}$ $T = \{ a, b \}$

$$P = \{ \Sigma \rightarrow a B$$

$$B \rightarrow a B B$$

$$B \rightarrow b \}$$

Generates strings (a)abab (b)aabb (c)aaabbb (d)aababb.

【解】(b)(c)(d)

【說明】(b)(c)(d)的導出過程如下：

$$\Sigma \rightarrow a B \rightarrow a \underline{a B B} \rightarrow a a \underline{b b}$$

$$\Sigma \rightarrow a B \rightarrow a \underline{a B B} \rightarrow a a \underline{a B B} B \rightarrow aaab \underline{b b}$$

$$\Sigma \rightarrow a B \rightarrow a \underline{a B B} \rightarrow a a \underline{b a B B} \rightarrow a a b a \underline{b b}$$

- 【例三】 Suppose that a grammar has the following productions:

$$S \rightarrow a B c$$
$$B \rightarrow b X b$$
$$B \rightarrow b X$$
$$X \rightarrow a$$
$$X \rightarrow a b$$

Is the string “a b a b c” in the language generated by the grammar?

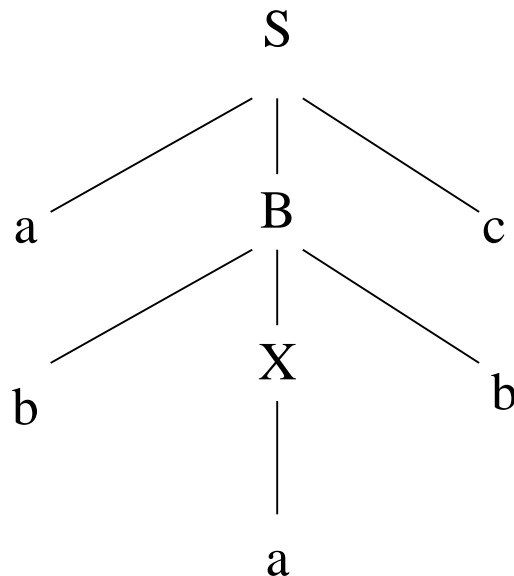
Is the grammar ambiguous? Justify your answer.

Describe the set of all strings in the language which is generated by the grammar.

【解】

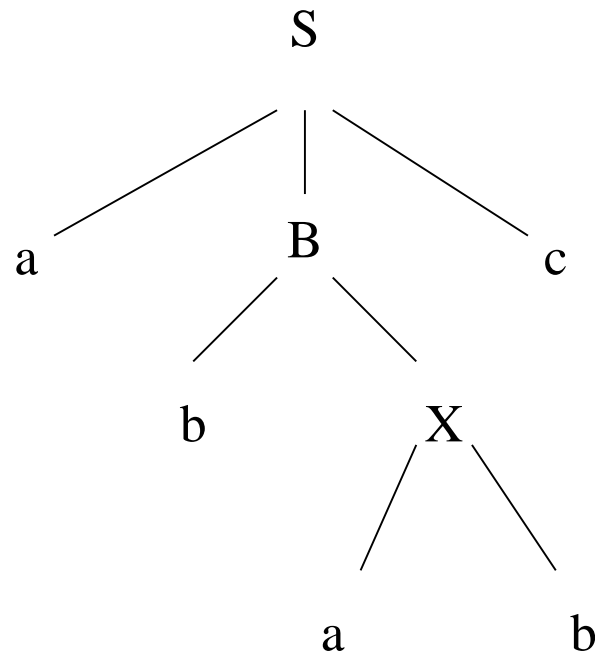
(一)Yes, the string “a b a b c” is in the language generated by the grammar.

【說明】字串 “a b a b c”能從上述文法導出，以剖析樹(parse tree)敘述導出情形如下：



(二) Yes, the grammar is ambiguous.

本例之文法對於某些字串的展開有一種以上的導出方式，所以此文法是語意不明確(ambiguous)之文法；例如：本例（一）中字串 “a b a b c” 有另一種導出方式，以剖析樹敘述導出情形如下：



(三) the set of all strings in the language which is generated by the above grammar = {string of the form $a b a b^n c$ for $0 \leq n \leq 2$ }

【說明】依上述文法的各種情形分別展開之，產生的字串如右： $a b a b c$ 、 $a b a b b c$ 及 $a b a c$ ，這些字串的通式是 $a b a b^n c$ (其中 $0 \leq n \leq 2$)。

- 【例四】 Is the following grammar an ambiguous grammar?

[Yes or no only]

$$E \rightarrow E + E \mid E * E \mid (E) \mid - E \mid \langle id \rangle$$

Where E represents an abbreviation for expression, and $\langle id \rangle$ represents an identifier.

【解】 Yes, the above grammar is an ambiguous grammar.

【說明】 本例的文法對於某些指述(statement)的展開具有一種以上的導出方式，所以此文法是語意不明確之文法 (ambiguous grammar)；例如： $\langle \text{id} \rangle * \langle \text{id} \rangle + \langle \text{id} \rangle$ 利用上述文法展開具有一種以上的導出方式，分別以剖析樹 (parse tree) 表示如下：

