

一个简单的 语法制导翻译器

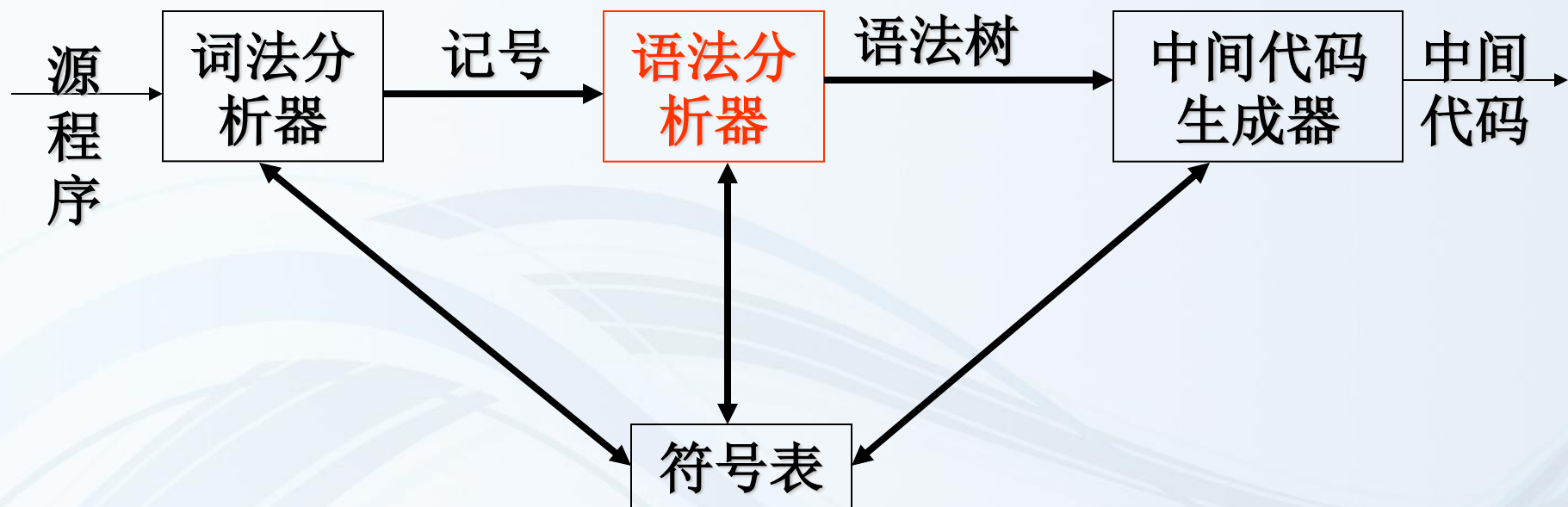
待编译的代码

```
{  
    int i , j ;  
    float a[100], v, x;  
    while ( true ) {  
        do{ i = i + 1 ; }while ( a[i] < v );  
        do{ j = j - 1 ; }while ( a[j] > v );  
        if ( i >= j ) break;  
        x=a[i]; a[i]=a[j]; a[j] = x;  
    }  
}
```

简化的中间代码

```
1:      i=i+1  
2:      t1=a[i]  
3:      if t1<v goto 1  
4:      j=j-1  
5:      t2=a[j]  
6:      if t2>v goto 4  
7:      ifFalse i>=j goto 9  
8:      goto 14  
9:      x=a[i]  
10:     t3=a[j]  
11:     a[i]=t3  
12:     a[j]=x  
13:     goto 1  
14:
```

2.1 引言（编译器前端的模型）



2.2 文法的形式化描述

- ❖ 一个文法包括若干**产生式** (production), 产生式也称**重写规则** (rewriting rule, substitution rule) 或**规则**。

(上下文无关文法) **产生式**写作:

$$U ::= u \quad \text{或} \quad U \rightarrow u$$

其中:

U是一个文法符号, 称为规则的**左部**或**头**

u是有穷文法符号串, 称为规则的**右部**或**体**

- ❖ **例:**

$\langle \text{sentence} \rangle \rightarrow \langle \text{subject} \rangle \langle \text{verb} \rangle$

$\langle \text{sentence} \rangle \rightarrow \langle \text{subject} \rangle \langle \text{verb} \rangle \langle \text{complement} \rangle$

文法的形式化描述 (II)

❖ 巴科斯-瑙尔范式 (Backus-Naur Form, BNF)

用**尖括号**包围非终结符来标记非终结符

终结符带有**下划线** (如果很明确, 也常省略)

符号 **::=** 表示 “推导出”

符号 **|** 表示 “还能推导出”

$$\begin{aligned} \langle \text{sentence} \rangle ::= & \langle \text{subject} \rangle \langle \text{verb} \rangle \\ & | \langle \text{subject} \rangle \langle \text{verb} \rangle \langle \text{complement} \rangle \end{aligned}$$

↓ 常常写成

$$\begin{aligned} \textit{sentence} \rightarrow & \textit{subject} \ \textit{verb} \\ & | \textit{subject} \ \textit{verb} \ \textit{complement} \end{aligned}$$

文法例1

<句子> ::= <主语> <谓语> <宾语>

<主语> ::= <代词> | <名词>

<谓语> ::= <动词>

<宾语> ::= <代词> | <名词>

<代词> ::= **I** | **you** | **he** | **she** | **it** | **we** | **they**

<名词> ::= **cat** | **dog** | **book** | **pen** | **computer**

<动词> ::= **eat** | **read** | **write** | **use** | **see**

最后一条相当于**5**条规则

产生式的使用

<句子> ⇒ <主语> <谓语> <宾语>

⇒ <代词> <谓语> <宾语>

⇒ I <谓语> <宾语>

⇒ I <动词> <宾语>

⇒ I use <宾语>

⇒ I use <名词>

⇒ I use computer

文法例2

<URL> ::= <协议>://<域名>

<协议> ::= http | https | ftp

<域名> ::= <域名段>.<域名段>.<顶级域名>

<域名段> ::= <字母数字串>

<顶级域名> ::= com | cn | org | net | edu

**<字母数字串> ::= <字母> | <数字> | <字母> <字母数字串>
| <数字> <字母数字串>**

<字母> ::= a | b | ... | z | A | B | ... | Z

<数字> ::= 0 | 1 | ... | 9

推导 *Derivation*

❖ 从开始符出发，不断将某个非终结符替换为该非终结符的某个产生式的右部。

❖ 推导的符号： \Rightarrow

❖ 语言：从开始符出发，利用推导能得到的所有终结符号串的集合。

推导 (例)

$list \Rightarrow list + digit \Rightarrow list - digit + digit$

$\Rightarrow digit - digit + digit \Rightarrow 9 - digit + digit$

$\Rightarrow 9 - 5 + digit \Rightarrow 9 - 5 + 2$

$list \rightarrow list + digit$

$list \rightarrow list - digit$

$list \rightarrow digit$

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

上下文无关文法Context-Free Grammar

❖ 也称2型文法，是四元组：

➤ 终结符 (terminal) 集 (常记作 T 或 V_T)

✓ 文法所定义语言的基本符号的集合

✓ 仅仅出现在产生式右部 (体)

➤ 非终结符 (non-terminal) 集 (常记作 N , N_T 或 V_N)

✓ 在某个产生式的左部出现过的文法符号

➤ 产生式 (production) 集

➤ 开始符 (start symbol)

✓ 两种说明方式：

(1) 明确指出, 如 $G[A]$

(2) 没有明确指出时, 第1条产生式规则左部的文法符号

正则文法Regular Grammar

❖ 也称3型文法，文法的任何产生式均为 $A \rightarrow aB$ 或 $A \rightarrow a$ 的形式，其中A和B是非终结符，a是终结符或空串

➤ 也称右线性文法

➤ 3型文法也常定义为左线性文法：文法的任何产生式均为 $A \rightarrow Ba$ 或 $A \rightarrow a$ 的形式

➤ 3型文法是特殊的2型文法

上下文无关文法（例）

$$S \rightarrow (S)$$

$$S \rightarrow a$$

也写作： $S \rightarrow (S) \mid a$

终结符集 $T=\{(\,a,\,)\}$ ，非终结符集 $N=\{S\}$ ，开始符是 S

本文法描述了语言 $\{(^n a)^n \mid n \geq 0\}$

上下文无关文法生成的语言称为上下文无关语言

串

❖ 字母表：元素的有穷非空集合，常记作 Σ 。

➤ 字母表中的元素称为**符号**，故字母表又名**符号集**。

❖ 串（符号串）：符号构成的有穷**序列**

➤ 空符号串（空串）：不含任何符号的符号串，用 ε 表示

❖ 串的计算：

➤ 串的**长度**：组成该串的符号个数，例： $x=abca$, 则 $|x|=4$

➤ 串的联结：设 x 和 y 表示两个串，则 xy 是它们的联结，表示将串 y 联结到 x 之后。

✓ $x\varepsilon=\varepsilon x=x$

✓ $x=abca, y=bb$, 则 $xy=abcabb$

✓ 串的幂运算： $x^0=\varepsilon$, $x^2=xx$, $x^2y^2=xxyy$

闭包

❖ 符号串集合的乘积：设A和B为两个符号串集合，则它们的乘积 $AB = \{xy \mid x \in A, y \in B\}$ 。

➤ $A = \{a, ab\}$, $B = \{b, bb\}$, 则 $AB = \{ab, abb, abbb\}$

➤ $\{\epsilon\}A = A\{\epsilon\} = A$

➤ 符号串集合的幂运算： $A = \{a, b\}$, 则 $A^0 = \{\epsilon\}$, $A^2 = \{aa, ab, ba, bb\}$

❖ 符号串集合的闭包：符号串集合的各次方幂之并集

➤ 正闭包：不考虑0次幂

➤ 串的闭包和正闭包： $a^* = \{\epsilon, a, aa, aaa, \dots\}$ $a^+ = \{a, aa, aaa, \dots\}$

➤ $A = \{a, b, c\}$, 则 $A^* = \{\epsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, \dots\}$

$A^+ = ?$

例子

已知: $L = \{A, B, C, D\}$ $M = \{1, 2, 3\}$

$L \cup M = \{A, B, C, D, 1, 2, 3\}$

$LM = \{A1, A2, A3, B1, B2, B3, C1, C2, C3, D1, D2, D3\}$

$L^2 = \{AA, AB, AC, AD, BA, BB, BC, BD, CA, \dots DD\}$

$L^* = \{ \text{All possible strings of } L \text{ plus } \epsilon \}$

$L^+ = L^* - \{\epsilon\}$

$L(L \cup M) = ??$

$L(L \cup M)^* = ??$

句型 and 句子

❖ 推导长度：执行**一步推导**的次数

➤ 一步推导 \Rightarrow

➤ 大于等于**1**步的推导 $\stackrel{+}{\Rightarrow}$

➤ 大于等于**0**步的推导 $\stackrel{*}{\Rightarrow}$

❖ 对于上下文无关文法 $G[S]$:

➤ 如果 $S \stackrel{*}{\Rightarrow} \alpha$, 则 α 为文法 G 的**句型**

➤ 如果 $S \stackrel{\pm}{\Rightarrow} W$ ($W \in T^*$), 则 W 为文法 G 的**句子**。

✓ 一个文法的句子集合称为**语言**。

✓ 文法 G 对应的语言记为 **$L(G)$**

例：以下文法会产生多少个不同的句子

(1) $A \rightarrow BB$ $B \rightarrow CC$ $C \rightarrow 1 \mid 2$

(2) $A \rightarrow BB$ $B \rightarrow CC$ $C \rightarrow 1 \mid 2 \mid \varepsilon$

例：以下文法会产生多少个不同的句型

(1) $A \rightarrow BB$ $B \rightarrow CC$ $C \rightarrow 1 \mid 2$

$$0 + 1 + 1 + 2 * 3^2 + 3^4 = 101$$

(2) $A \rightarrow BB$ $B \rightarrow CC$ $C \rightarrow 1 \mid 2 \mid \varepsilon$

$$1 + 5 + 4^2 + (3^2 + 3^2 + 3^3) + 3^4 = 148$$

文法和语言

❖ 文法和语言的关系

- 给定文法，就**唯一**地确定语言
- 给定语言，能给出文法，但文法**不唯一**

❖ 等价文法

- 设 G_1 和 G_2 是两个文法，若 $L(G_1)=L(G_2)$ ，则称 G_1 与 G_2 为等价文法。

文法和语言 (例)

文法确定语言
注意指数的取值范围

(1) 文法 $A \rightarrow aAb \mid ab$ 对应的语言为:

$$L(G[A]) = \{a^n b^n \mid n \geq 1\}$$

(2) 文法 $S: S \rightarrow xSx \mid xS \mid y$ 所识别的语言为:

$$L(G[S]) = \{x^m y x^n \mid m \geq n \geq 0\}$$

(3) 文法 $G: S \rightarrow xxS \mid y$ 所识别的语言为:

$$L(G[S]) = \{x^{2^n} y \mid n \geq 0\}$$

也可以写成 $L(G[S]) = (xx)^* y$

文法和语言 (例2)

多个文法可以
产生同一语言

(1) 语言 $L = \{ab^n a \mid n \geq 1\}$ 对应的文法可以为:

$$A \rightarrow aBa \quad B \rightarrow b \mid bB \quad \text{或}$$

$$A \rightarrow aB \quad B \rightarrow ba \mid bB \quad \text{或} \dots$$

(2) 语言 $L = \{a^{2n}b \mid n \geq 1\}$ 对应的文法可以为:

$$A \rightarrow Bb \quad B \rightarrow aa \mid aaB \quad \text{或}$$

$$A \rightarrow aab \mid aaA \quad \text{或} \dots$$

(3) 语言 $L = \{a^m b^n c^k \mid m=n \text{ 或 } n=k\}$ 对应的文法可以为:

$$S \rightarrow AB \mid DE \quad A \rightarrow aAb \mid \varepsilon \quad B \rightarrow cB \mid \varepsilon$$

$$D \rightarrow aD \mid \varepsilon \quad E \rightarrow bEc \mid \varepsilon$$

语法分析树

❖ 以图形方式描述推导过程

❖ 语法分析树的构成：

➤ 根节点是开始符

➤ 叶节点是终结符 (*token*) 和 ϵ

➤ 内部节点（非叶结点）是非终结符

➤ 如果应用了规则 $A \rightarrow x_1x_2\dots x_n$, 则 A 是内部结点;
 x_1, x_2, \dots, x_n 是子结点

使用分析树来描述推导过程

$list \Rightarrow list + digit$

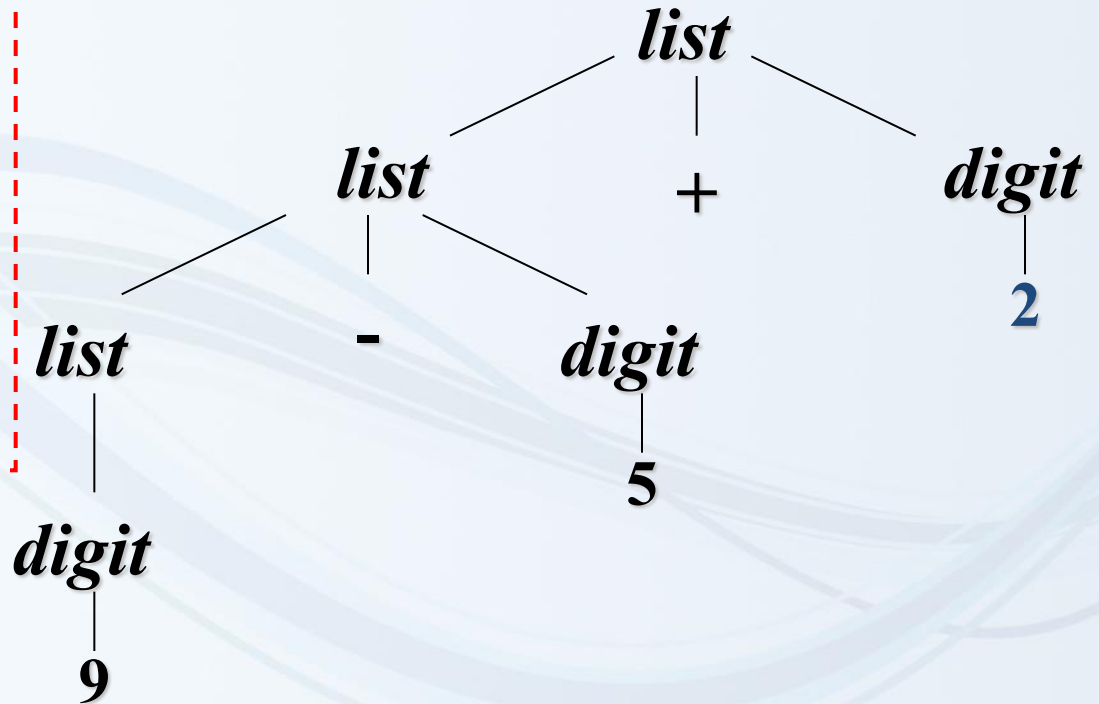
$\Rightarrow list - digit + digit$

$\Rightarrow digit - digit + digit$

$\Rightarrow 9 - digit + digit$

$\Rightarrow 9 - 5 + digit$

$\Rightarrow 9 - 5 + 2$



综合例

考虑文法： $S \rightarrow SS+ \mid SS^* \mid a$

1) 试推导出串 $aa+a^*$

$S \Rightarrow SS^* \Rightarrow SS+S^* \Rightarrow aS+S^* \Rightarrow aa+S^* \Rightarrow aa+a^*$

2) 试为这个串构造一棵语法分析树

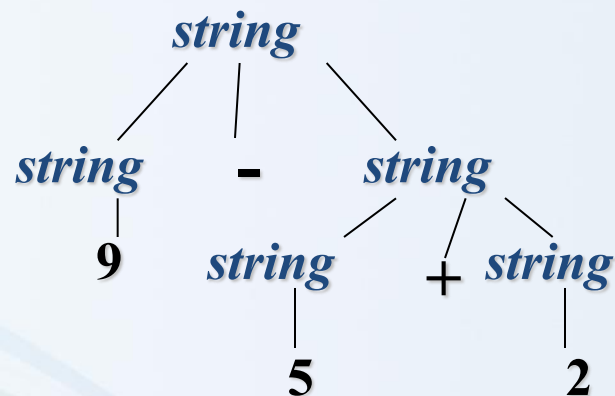
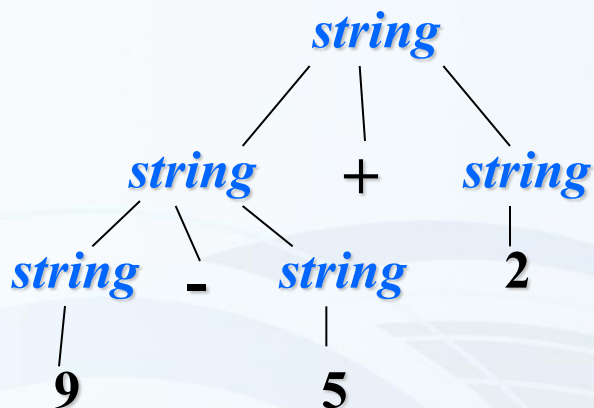
3) 该文法生成的语言是什么？

后缀表达式

二义性 *Ambiguous*

- ❖ 基于一个文法，若有**多棵分析树**生成**同一个终结符号串**，则此文法具有二义性。
- ❖ 本课程只讨论**文法二义**，不涉及**语义二义**。
 - 语义二义例：Jack said Tom left *his* assignment at home.

一个二义的文法

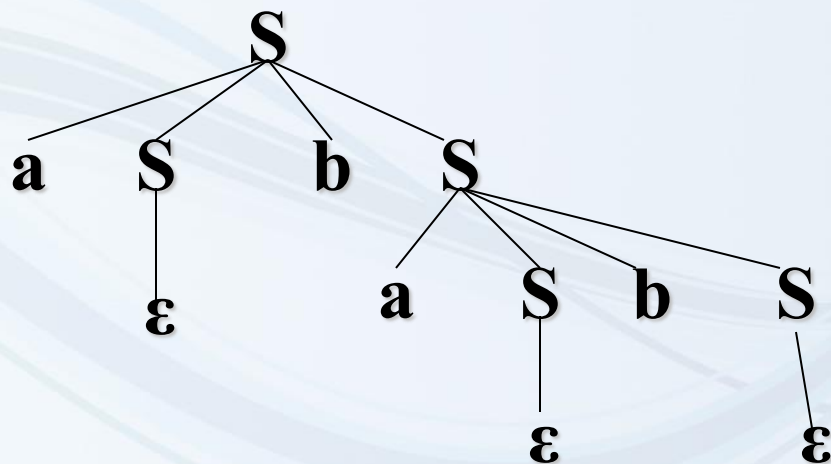
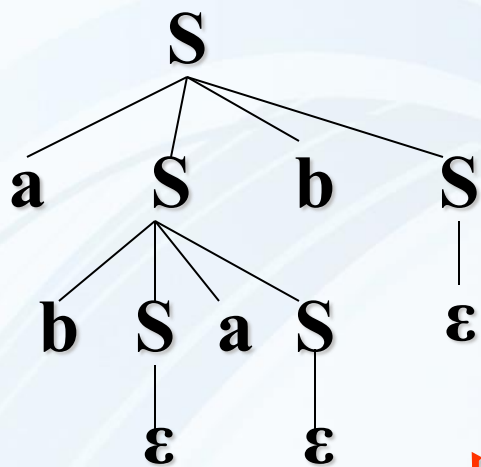


$string \rightarrow string + string \mid string - string \mid 0 \mid 1 \mid \dots \mid 9$

证明以下文法二义。

$$S \rightarrow aSbS \mid bSaS \mid \epsilon$$

证：对于串abab存在两棵不同的分析树：



故本文法二义

Which of the following grammars are ambiguous?

(A) $S \rightarrow SS \mid a \mid b$

(B) $E \rightarrow E + E \mid id$

(C) $S \rightarrow Sa \mid Sb$

(D) $E \rightarrow E' \mid E' + E$

$E' \rightarrow -E' \mid id \mid (E)$

可能基于**优先级**和**结合性**来构建无二义的文法

- ❖ 四则运算存在两个优先级，所以可引入两个非终结符 **expr**和**term**来对应不同的抽象层次。
 - 抽象层次**较低**的非终结符对应**较高**优先级
- ❖ 四则运算均为**左结合**，因此每条规则中，**更为抽象的非终结符应位于左边**

expr → **expr + term | expr - term | term**

term → **term * factor | term / factor | factor**

factor → **digit | (expr)**

digit → **0 | 1 | 2 | ... | 9**

最左推导与最右推导

- ❖ **最左推导** (*Leftmost derivation*) : 总是替换当前句型中**最左边的**非终结符
- ❖ **最右推导** (*rightmost derivation*) : 总是替换当前句型中**最右边的**非终结符

$$\begin{array}{l} \hline E \rightarrow E+T \\ E \rightarrow T \\ T \rightarrow id \\ \hline \end{array}$$

Derivations for $id + id$:

$$\begin{array}{l|l} E \Rightarrow E+T & E \Rightarrow E+T \\ \Rightarrow T+T & \Rightarrow E+id \\ \Rightarrow id+T & \Rightarrow T+id \\ \Rightarrow id+id & \Rightarrow id+id \\ \text{LEFTMOST} & \text{RIGHTMOST} \end{array}$$

最左推导与最右推导（2）

❖ 最左（右）推导与分析树一一对应

- 可以通过说明对于**某一个串存在两个不同的最左推导**来证明文法二义
- 可以通过说明对于**某一个串存在两个不同的最右推导**来证明文法二义

证明以下文法二义：（第二种证法）

$$S \rightarrow aSbS \mid bSaS \mid \varepsilon$$

证：对于串 **abab** 存在两个不同的最左推导：

$$S \Rightarrow aSbS \Rightarrow abSaSbS \Rightarrow abaSbS \Rightarrow ababS \Rightarrow abab$$

$$S \Rightarrow aSbS \Rightarrow abS \Rightarrow abaSbS \Rightarrow ababS \Rightarrow abab$$

故本文法二义