

# 一种有效的概率上下文无关文法分析算法

朱胜火 周明 刘昕 黄昌宁

(清华大学计算机科学与技术系智能技术与系统国家重点实验室 北京 100084)

**摘要** 句法分析的研究是自然语言处理领域的一个重要组成部分. 该文提出并实现了一种有效的概率上下文无关文法 SCFG(stochastic context-free grammar)的分析算法. 首先对原有的 GLR 分析表加以改进, 以便能够利用分析过程的控制结构来计算有关的概率; 然后对分析过程中的每个状态增设了下标, 以区分不同的归约路径. 通过上述手段, 成功地引入了状态的前向(Forward)概率和内(Inner)概率. 利用这两个概率可以计算输入句子的所有可能分析树的概率, 用于选择最佳的分析结果. 通过对大规模真实文本进行实验, 结果表明, 这种算法具有较高的分析效率.

**关键词** SCFG(stochastic context-free grammar), 分析算法, GLR 算法, Earley 算法, 自然语言处理.

**中图法分类号** TP301, TP18

概率的上下文无关文法 SCFG(stochastic context-free grammar)或 PCFG(probabilistic context-free grammar)在上下文无关文法 CFG(context-free grammar)的基础上加入了概率模型, 用于选择优先扩展的分析路径, 或者筛选有歧义的分析结果. SCFG 的分析算法通常是由原有的上下文无关文法的分析算法扩展而来的. 例如, Fujisaki 等人对 CYK 算法的扩展<sup>[1]</sup>, Stolcke 对 Earley 算法的扩展<sup>[2]</sup>, Wright 对 GLR 算法的扩展<sup>[3]</sup>等等.

本文提出并实现了一种有效的概率上下文无关文法(SCFG)的分析算法. 第1节详细描述分析算法的设计, 主要包括分析表的构造方法、分析控制结构、状态的前向概率和内概率以及概率化的分析过程, 并给出一个分析实例; 第2节叙述了对大规模真实文本的实验, 并与相关的研究工作进行了比较.

## 1 分析算法的设计

GLR 算法是一种比较有效的分析方法.<sup>[4]</sup>它采用了 LR 分析表、图栈(Graph-Structured Stack)和分析森林(Parse Forest)等措施, 达到了较快的分析速度(GLR 算法的时间复杂度与其他算法相同, 是  $O(n^3)$ , 但它具有较小的常数). 因此, 本研究采用了类似于 GLR 算法的分析控制结构来计算分析树概率.

LR 分析表最初被设计用于确定的上下文无关语言(如程序设计语言)的分析, 在生成算法中采用合并项目集的方法避免可能的冲突. 如果直接使用 LR 分析表来计算项目的前向概率和内概率, 则因为在一个项目集中的不同项目可能具有不同的分析概率, 所以难以计算项目集中每一个项目的前向概率和内概率. 针对这个问题, 我们在生成分析表时对某一个状态的相同符号的移进(或转移)动作产生的项目集不进行合并, 用其核心项目的前向概率和内概率来定义该状态的前向概率和内概率. 而该状态中的其他项目的前向概率和内概率则可通过其核心项目的前向概率和内概率来推定.

此外, 在 GLR 算法的分析过程中, 遇到两个(或多个)动作移进(或转移)到同一个项目集情况时, 只产生一个新的状态. 可是, 经过若干步以后, 该状态的后续状态可能要归约到两个(或多个)不同的状态, 而归约后, 转移产生两个(或多个)状态的概率一般说来是不同的, 所以, 前述遇到两个(或多个)动作移进(或转移)到同一个项目集情况时, 只产生一个新的状态的做法, 不能有效地利用分析控制结构来计算这些状态的概率. 针对这个问题, 我们通过对分析过程中的每个状态增设下标来表示分析的起始位置, 以区分不同的归约路径, 从而可以方便地计算每个归约后转移产生的状态的概率.

\* 本研究得到国家自然科学基金和航天预研基金资助. 作者朱胜火, 1974年生, 硕士, 主要研究领域为自然语言处理. 周明, 1964年生, 博士后, 副研究员, 主要研究领域为机器翻译. 刘昕, 1974年生, 硕士, 主要研究领域为自然语言处理. 黄昌宁, 1937年生, 教授, 博士生导师, 主要研究领域为计算语言学.

本文通讯联系人: 黄昌宁, 北京 100084, 清华大学计算机科学与技术系智能技术与系统国家重点实验室

本文 1996-11-07 收到原稿, 1997-07-18 收到修改稿

通过上述措施,既保留了原有 GLR 算法高效的特点,又可利用 Stolcke 计算前向概率和内概率的方法来优选分析路径和评价分析结果,从而有效地提高了分析效率。

### 1.1 分析表的构造算法

项目是一个带点的概率化的产生式,形如  $\langle A \rightarrow \alpha \cdot \beta, \rho \rangle$ 。核心项目是  $\langle \rightarrow \cdot S, \rho \rangle$ ,或点不在最左边的项目  $\langle A \rightarrow \alpha \cdot \beta, \rho \rangle$  ( $|\alpha| > 0$ )。项目集<sup>\*</sup>由一个或多个项目组成,每个项目集中只有一个核心项目,每个项目集可以通过计算核心项目的所有预测项目得到,每个核心项目仅属于一个项目集。KERNEL( $I$ )表示项目集  $I$  的核心项目。项目集的动作有移进(Shift)、归约(Reduce)和转移(Goto)。每个项目集对同一个终结符(或非终结符)可以有一个动作,也可以有多个动作。分析表的构造方法如下:

- (1) 起始项目集  $I_0$ ;  $\text{KERNEL}(I_0) = \langle \rightarrow \cdot S, \rho \rangle$ ;
- (2) 对终结符  $b$ , 如果  $\langle A \rightarrow \alpha \cdot b\beta, \rho \rangle \in I_m, \langle A \rightarrow \alpha b \cdot \beta, \rho \rangle \in I_n$ , 则  $\langle \text{shift-to-}n \rangle \in \text{action}[m, b]$ ;
- (3) 对非终结符  $B$ , 如果  $\langle A \rightarrow \alpha \cdot B\beta, \rho \rangle \in I_m, \langle A \rightarrow \alpha B \cdot \beta, \rho \rangle \in I_n$ , 则  $\langle \text{go-to-}n \rangle \in \text{action}[m, B]$ ;
- (4) 如果  $\langle \rightarrow \cdot S, \rho \rangle \in I_m$ , 则  $\langle \text{accept} \rangle \in \text{action}[m, \$]$ ;
- (5) 如果  $B \rightarrow \gamma \cdot, \rho \in I_m$ , 则对每个  $b \in \text{FOLLOW}(B)$ , 有  $\langle \text{reduce-}l \rangle \in \text{action}[m, b]$ 。<sup>\*\*</sup>

### 1.2 分析控制结构

为了引入概率化的分析,这里对 GLR 算法的分析控制结构做了一些修改。

我们将每个状态表示为  $\langle m, i, k \rangle$ , 其中  $m$  是项目集  $I_m$  的标号,  $i$  是下一个输入终结符的下标, 简称输入下标,  $k$  是  $\text{KERNEL}(I_m)$  的起始下标。<sup>\*\*\*</sup>分析过程的初始状态是  $\langle 0, 0, 0 \rangle$ 。

如果  $v = \langle m, i, k \rangle$  是活动状态,

- (1) 如果  $\langle \text{accept} \rangle \in \text{action}[m, x_i]$ , 则得到一个分析结果。
- (2) 如果  $\langle \text{shift-to-}n \rangle \in \text{action}[m, x_i]$ , 且  $\text{KERNEL}(I_n)$  是  $\text{KERNEL}(I_m)$  的下一个项目, 则将  $w = \langle n, i+1, k \rangle$  压入分析栈, 位居  $v$  的上方;
- (3) 如果  $\langle \text{shift-to-}n \rangle \in \text{action}[m, x_i]$ , 且  $\text{KERNEL}(I_n)$  不是  $\text{KERNEL}(I_m)$  的下一个项目, 则将  $w = \langle n, i+1, i \rangle$  压入分析栈, 位居  $v$  的上方;
- (4) 如果  $\langle \text{reduce-}l \rangle \in \text{action}[m, x_i]$ ,  $\text{KERNEL}(I_m) = \langle Y \rightarrow \nu \cdot, \rho \rangle$ , 对在栈中  $v$  之前所有第  $l$  个状态  $v'' = \langle m'', k, k'' \rangle$ ,  $\text{KERNEL}(I_{m''}) = \langle X \rightarrow \lambda \cdot Z\mu, \rho'' \rangle$ ,  $\langle \text{go-to-}n \rangle \in \text{action}[m'', Y]$ ;

如果  $\text{KERNEL}(I_n)$  是  $\text{KERNEL}(I_m)$  的下一个项目, 将  $w = \langle n, i, k'' \rangle$  压入分析栈, 位居  $v''$  的上方;

如果  $\text{KERNEL}(I_n)$  不是  $\text{KERNEL}(I_m)$  的下一个项目, 将  $w = \langle n, i, k \rangle$  压入分析栈, 位居  $v''$  的上方。

通过上述步骤可以得到分析的图栈, 用类似于 GLR 算法中构造分析森林的方法<sup>[5]</sup>来得到分析森林。

### 1.3 状态的前向概率和内概率

Stolcke 曾经利用计算前缀概率的方法将 Earley 算法成功地扩展成 SCFG 的分析算法。<sup>[2]</sup>这里, 我们借鉴了 Stolcke 计算前缀概率的方法, 同时, 为了适合我们的算法的分析控制结构, 又做了一些调整。为了引入概率化的分析过程, 我们首先给出如下定义。

定义 1. (a) 路径是由状态构成的、由移进和转移连接的状态串。

(b) 如果一条路径中所有移进的终结符连接成串  $x$ , 则称之为生成串  $x$  的受限路径。

(c) 如果一条路径的第 1 个状态的核心项目与最后一个状态的核心项目具有相同产生式和起始下标, 则称它为一条可归约的路径。

(d) 路径的长度是指路径中扫描状态的个数。

(e) 路径  $\mathscr{P}$  的概率  $P(\mathscr{P})$  是在  $\mathscr{P}$  中所有用到的归约产生式的概率之积。

(f) 非终结符生成的子串概率  $P(X \Rightarrow x)$  是开始于  $X$  的生成串  $x$  的所有可归约受限路径的概率之和。当  $X=S$  时,  $P$  就是句子概率。

(g) 前缀概率  $P(S \Rightarrow_l x)$  是开始于初始状态、终止于某个扫描状态且生成串  $x$  的所有可归约受限路径的概率

\* 有时也叫作状态。为了有所区分, 将在分析表中的称为项目集, 在分析过程中的称为状态。

\*\* 这里使用了类似于 SLR(1) 的分析表, 为了简化, 后面的例子仅使用类似于 LR(0) 的分析表。

\*\*\* 这里类似于 Earley 状态。设  $\text{KERNEL}(I_m)$  是  $X \rightarrow \lambda \cdot \mu$ , 则其对应的 Earley 状态是  $i; kX \rightarrow \lambda \cdot \mu$ 。

之和.

通过在每个状态上定义以下两个概率,就可以计算子串概率和前缀概率.

定义 2. (a) 状态  $\langle m, i, k \rangle$  的前向概率  $\alpha(\langle m, i, k \rangle)$  是所有以该状态为结束状态、长度为  $i$  的受限路径的概率之和.

(b) 状态  $\langle m, i, k \rangle$  的内概率  $\gamma(\langle m, i, k \rangle)$  是所有以  $\langle m, k, k \rangle$  为起始状态、以  $\langle m, i, k \rangle$  为结束状态的、长度为  $i-k$  的路径的概率之和.

推论. 当输入串为  $x(|x|=n)$  时,

假定  $S \Rightarrow_L x_0 \dots x_{l-1} X$  是语法的最左推导,那么,非终结符  $X$  生成子串  $x_k \dots x_{l-1}$  的概率是  $P(X \Rightarrow x_k \dots x_{l-1}) = \Sigma \gamma(\langle m, i, k \rangle)$  (对所有满足  $\text{LEFT}(\text{KERNEL}(I_m)) = X$  的  $m$  进行求和 \*);

$P(S \Rightarrow x) = \gamma(\langle m, n, 0 \rangle) = \alpha(\langle m, n, 0 \rangle)$ , 其中  $m$  满足  $\langle \rightarrow S \cdot, p \rangle \in I_m$ ;

$P(S \Rightarrow_L x) = \Sigma \alpha(\langle m, n, k \rangle)$  (对所有满足  $\langle X \rightarrow \lambda x_{n-1} \cdot \mu, p \rangle \in I_m$  的  $m$  进行求和).

由推论可知,通过计算前向概率和内概率可以得到子串和前缀概率.利用前缀概率可以优选分析路径,利用子串概率可以得到输入串的分析树概率.

### 1.4 概率化的分析过程

可将每个状态  $\langle m, i, k \rangle$  扩展为  $\langle m, i, k, \alpha, \gamma \rangle$ ,  $\alpha = \alpha(\langle m, i, k \rangle)$ ,  $\gamma = \gamma(\langle m, i, k \rangle)$ . 分析过程的初始状态是  $\langle 0, 0, 0, 1, 1 \rangle$ .

如果  $v = \langle m, i, k, \alpha, \gamma \rangle$  是活动状态,

(1) 如果  $\langle \text{accept} \rangle \in \text{action}[m, x_i]$ , 则得到一个分析结果;

(2) 如果  $\langle \text{shift-to-}n \rangle \in \text{action}[m, x_i]$ , 且  $\text{KERNEL}(I_n)$  是  $\text{KERNEL}(I_m)$  的下一个项目, 则将  $w = \langle n, i+1, k, \alpha', \gamma' \rangle$  压入分析栈, 位居  $v$  的上方. 其中

$$\alpha' = \alpha \quad \gamma' = \gamma$$

(3) 如果  $\langle \text{shift-to-}n \rangle \in \text{action}[m, x_i]$ , 且  $\text{KERNEL}(I_n)$  不是  $\text{KERNEL}(I_m)$  的下一个项目, 则将  $w = \langle n, i+1, i, \alpha', \gamma' \rangle$  压入分析栈, 位居  $v$  的上方. 其中

$$\begin{aligned} \alpha' + &= \alpha R(Z \Rightarrow_L Y) P(Y \Rightarrow v) \\ \gamma' &= P(Y \Rightarrow v) \end{aligned}$$

(4) 如果  $\langle \text{reduce-}l \rangle \in \text{action}[m, x_i]$ ,  $\text{KERNEL}(I_m) = \langle Y \rightarrow v \cdot, p \rangle$ , 对在栈中  $v$  之前所有第  $l$  个状态  $v'' = \langle m'', k, k'', \alpha'', \gamma'' \rangle$ ,  $\text{KERNEL}(I_{m''}) = \langle X \rightarrow \lambda \cdot Z \mu, p'' \rangle$ ,  $\langle \text{go-to-}n \rangle \in \text{action}[m'', Y]$ :

如果  $\text{KERNEL}(I_n)$  是  $\text{KERNEL}(I_{m''})$  的下一个项目, 将  $w = \langle n, i, k'', \alpha', \gamma' \rangle$  压入分析栈, 位居  $v''$  的上方. 其中

$$\begin{aligned} \alpha' + &= \alpha'' \gamma R(Z \Rightarrow Y) \\ \gamma' + &= \gamma'' \gamma R(Z \Rightarrow Y) \end{aligned}$$

如果  $\text{KERNEL}(I_n)$  不是  $\text{KERNEL}(I_{m''})$  的下一个项目, 将  $w = \langle n, i, k, \alpha', \gamma' \rangle$  压入分析栈, 位居  $v''$  的上方. 其中

$$\begin{aligned} \alpha' + &= \rho \alpha'' \gamma R(Z \Rightarrow Y) R(Z \Rightarrow_L W) \\ \gamma' + &= \rho \gamma R(Z \Rightarrow Y) \end{aligned}$$

$W = \text{LEFT}(\text{KERNEL}(I_n))$ ,  $Z$  是  $\text{LEFT}(\text{KERNEL}(I_{m''}))$ ,  $Y \rightarrow v$  是  $\text{KERNEL}(I_n)$  的产生式.

从分析森林中选取概率最大的分析树作为分析结果. 具体算法见文献[6].

### 1.5 一个例子

利用图 1 的文法  $G$  分析句子  $I$  saw a man in the park with a scope. 可以生成分析表, 如图 2 所示.

注意, 在分析表中每个项目集对同一个终结符(或非终结符)可以有 1 个动作, 也可以有多个动作. 例如, 项目集 0 中对非终结符  $S$  有两个动作:  $g1$  和  $g5$ .

分析过程大致如下. 起始状态是  $\langle 0, 0, 0, 1, 1 \rangle$ , 第 1 个输入的词是  $I$ , 其词性是  $-n$  (名词), 所以输入符号是  $-n$ . 根据分析表, 动作是  $s7$ , 通过 1.4 节中分析过程步骤(3)得到新的状态是  $\langle 7, 0, 1, 0.5, 0.333\ 333 \rangle$ . 然后, 查分析表知动作是  $r3$ , 则根据第 3 条产生式进行规约, 通过 1.4 节中分析过程步骤(4)得到状态  $\langle 3, 0, 1, 0.333\ 333, 0.2 \rangle$  和  $\langle 10, 0, 1, 0.166\ 667, 0.111\ 111 \rangle$ . 以此类推得到句子的分析结果. 图 3 表示了分析的前 24 个状态, 图 4 表示了分析的全部状态

\*  $\text{LEFT}(P)$  表示项目或产生式  $P$  的左边, 如  $\text{LEFT}(X \rightarrow \lambda \cdot \mu) = X$ .

\*\*  $+$  表示对右边的表达式求和.

(为清楚起见,省略了图 3 中一些没有扩展的状态)。

- 1. S: NP VP 0.6
- 2. S: S PP 0.4
- 3. NP: <sub>-n</sub> 0.333
- 4. NP: <sub>-d</sub> <sub>-n</sub> 0.333
- 5. NP: NP PP 0.333
- 6. PP: <sub>-p</sub> NP 1
- 7. VP: <sub>-v</sub> NP 1

	<sub>-n</sub>	<sub>-d</sub>	<sub>-p</sub>	<sub>-v</sub>	\$	S	NP	PP	VP
0	s7	s8				g1/g5	g3/g10		
1					s2				
2			a						
3				s14					g4
4			r1						
5			s12					g6	
6			r2						
7			r3						
8	s9								
9			r4						
10			s12					g11	
11			r5						
12	s7	s8					g10/g13		
13			r6						
14	s7	s8					g15/g10		
15			r7						

图1 文法G

图2 文法G的分析表

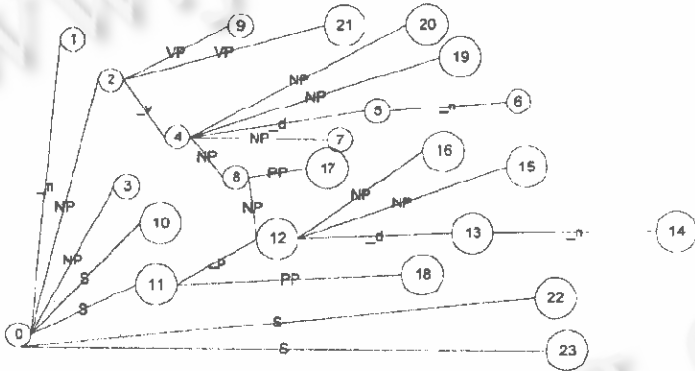


图3 状态转移图(1)

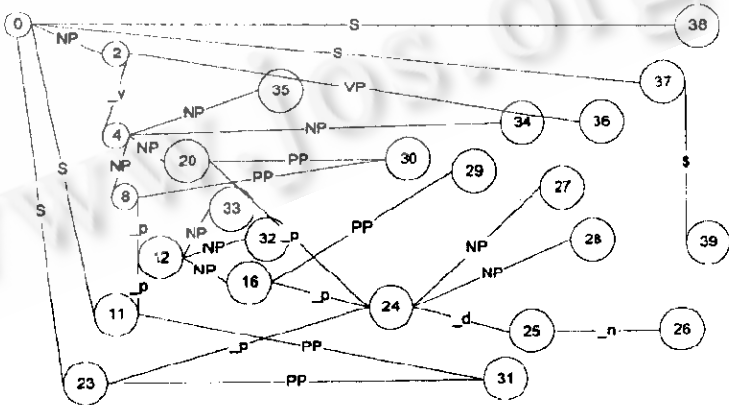


图4 状态转移图(2)

```

state 0: <0,0,0,1,1>
state 1: <7,0,1,0.5,0.333 333>
state 2: <3,0,1,0.333 333,0.2>
state 3: <10,0,1,0.166 667,0.111 111>
state 4: <14,1,2,0,1>
state 5: <8,2,3,0,0.333 333>
state 6: <9,2,4,0,0.333 333>
state 7: <15,1,4,0,0.333 333>
state 8: <10,2,4,0,0.111 111>
state 9: <4,0,4,0.111 111,0.066 666 7>
state 10: <1,0,4,0.066 666 7,0.066 666 7>
state 11: <5,0,4,0.044 444 4,0.026 666 7>
state 12: <12,4,5,0,1>
state 13: <8,5,6,0,0.333 333>
state 14: <9,5,7,0,0.333 333>
state 15: <13,4,7,0,0.333 333>
state 16: <10,5,7,0,0.111 111>
state 17: <11,2,7,0,0.037 037>
state 18: <6,0,7,0.014 814 8,0.008 888 89>
state 19: <15,1,7,0,0.037 037>
state 20: <10,2,7,0,0.012 345 7>
state 21: <4,0,7,0.012 345 7,0.007 407 41>
state 22: <1,0,7,0.016 296 3,0.016 296 3>
state 23: <5,0,7,0.010 864 2,0.006 518 52>
state 24: <12,7,8,0,1>
state 25: <8,8,9,0,0.333 333>
state 26: <9,8,10,0,0.333 333>
state 27: <13,7,10,0,0.333 333>
state 28: <10,8,10,0,0.111 111>
state 29: <11,5,10,0,0.037 037>
state 30: <11,2,10,0,0.008 230 45>
state 31: <6,0,10,0.005 267 49,0.003 160 49>
state 32: <13,4,10,0,0.037 037>
state 33: <10,5,10,0,0.012 345 7>
state 34: <15,1,10,0,0.016 460 9>
state 35: <10,2,10,0,0.005 486 97>
state 36: <4,0,10,0.005 486 97,0.003 292 18>
state 37: <1,0,10,0.009 613 17,0.009 613 17>
state 38: <5,0,10,0.006 408 78,0.003 845 27>
state 39: <2,0,11,0.009 613 17,0.009 613 17>

```

## 2 讨论

### 2.1 相关工作

对 GLR 算法的概率化扩展的其他工作有文献[3]的方法. 该方法利用 LR 转移概率来近似 SCFG 概率, 但是, 由于 LR 转移概率与 SCFG 转移概率并不完全等价, 所以, 基于 LR 转移概率的分析模型并不完全等效于基于 SCFG 的概率分析模型. 并且 LR 分析表的状态等同于项目(带点的产生式)的集合, 任何项目的集合都是一个潜在的状态, 所以在最坏情况下, 状态的个数与项目个数成指数关系.

本研究实现了一个真正的 SCFG 的概率分析算法, 而且, 它的分析表的项目集(状态)与核心项目一一对应, 所以项目集的个数仅与项目个数成线性关系.

本研究是在分析算法上的改进, 原则上与被分析的具体语言无关. 但由于语料库的限制, 本文仅对英语句法分析进行实验. 国内外某些作者已经应用 SCFG 分析汉语, 详见文献[7,8].

### 2.2 分析的时间

为了验证本算法的有效性, 我们分别利用本算法与 Stockle 扩展的 Earley 算法进行对比实验. 实验采用了 Susanne 英语树库. [9]Susanne 树库是一个面向自然语言处理的语法标注的语料库, 它包含 64 个文件, 共 15 万词. 我们从中获取了 723 条带概率的产生式作为实验用的语法. [6]在 Pentium/120 的微机上对所有的 5 671 个句子进行了分析, 其中选取 50 个句子的分析时间, 如图 5 所示. 从图中可以看出, 本算法比 Stockle 扩展的 Earley 算法在分析时间上有明显的提高.

### 2.3 剪枝阈值的设定

为了提高分析速度, 需要在分析过程中对分析图栈进行剪枝, 然而剪枝影响到分析的通过率和正确率. 这里分析的通过率定义为可以得到整个输入句子分析树的百分率; 在能得到分析树的句子中, 分析的正确率定义为在概率最大的分析树中正确分析的短语比率; 如果剪枝阈值为  $r$ , 则表示在具有相同输入下标的所有状态中, 前向概率小于最大值的  $\epsilon^{r\%}$  的状态时不进行移进动作.

我们对 Susanne 语料库 NO1 中的 113 个句子进行分析, 平均句长为 16.5 词, 其正确率、通过率与剪枝阈值的关系如图 6 所示. 从图中可以发现, 通过率随剪枝阈值的增大而提高, 当剪枝阈值达到 10 时, 它基本保持不变, 随剪枝阈值的增大, 正确率稍微有所减小, 在剪枝阈值达到 10 时, 它也基本保持不变. 所以, 在实际分析时将剪枝阈值设为 10, 对 Susanne 英语树库中所有的 5 671 个句子进行了分析, 通过率为 75.3%, 正确率为 73.7%.

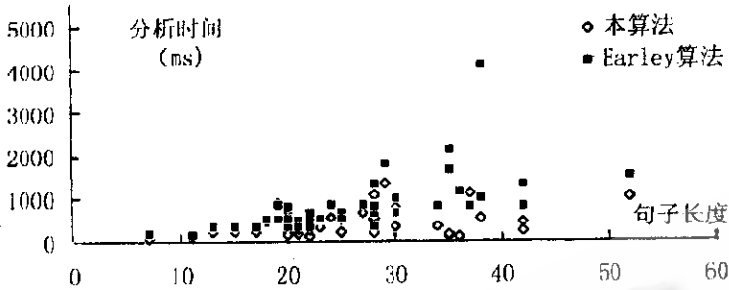


图5 句子长度与分析时间的关系

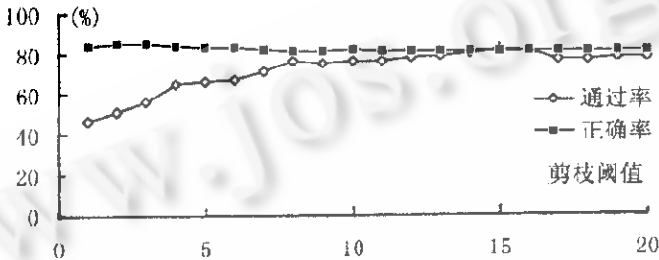


图6 剪枝阈值与正确率、通过率的关系

## 参考文献

- 1 Fujisaki T, Jelinek F, Cocke J *et al.* A probabilistic parsing method for sentence disambiguation. In: Tomita M ed. *Current Issues in Parsing Technology*, International Workshop on Parsing Technologies. Pittsburgh, Boston; Kluwer Academic Publishers, 1991. 139~152
- 2 Stolcke A. An efficient probabilistic context-free parsing algorithm that computes prefix probabilities. *Computational Linguistics*, 1995, 21(2):165~201
- 3 Wright J H. LR parsing of probabilistic grammars with input uncertainty for speech recognition. *Computer Speech and Language*, 1990, (4):297~323
- 4 Shann P. Experiments with GLR and chart parsing. In: Tomita M ed. *Generalized LR Parsing*, International Workshop on Parsing Technologies. Pittsburgh, Boston; Kluwer Academic Publishers, 1991. 17~34
- 5 Tomita M. *Efficient Parsing for Natural Language*. Boston; Kluwer Academic Publishers, 1986
- 6 朱胜火. 从双语对齐文本中获取翻译模板[硕士论文]. 清华大学, 1997  
(Zhu Sheng-huo. Learning translation patterns from bilingual aligned text: [M.S. Thesis]. Tsinghua University, 1997)
- 7 张民. 基于弱限制随机上下文相关文法的汉语树库构造方法研究[博士论文]. 哈尔滨工业大学, 1997  
(Zhang Min. Research on algorithms of Chinese treebank construction based on weakly restricted stochastic context-sensitive grammars [Ph. D. Dissertation]. Harbin; Harbin Institute of Technology, 1997)
- 8 Lee L-S, Chien L-F, Lin L-J *et al.* An efficient natural language processing system specially designed for the Chinese language. *Computational Linguistics*, 1991, 17(4):347~375
- 9 Sampson G. *English for the Computer*. Oxford; Oxford University Press, 1995

## An Efficient Stochastic Context-Free Parsing Algorithm

ZHU Sheng-huo ZHOU Ming LIU Xin HUANG Chang-ning

(State Key Laboratory of Intelligent Technology and Systems Department of Computer Science and Technology  
Tsinghua University Beijing 100084)

**Abstract** The research of parsing is important in the field of natural language processing. An efficient stochastic context-free parsing algorithm is described in this paper. In order to implement a stochastic context-free parser, the authors rebuild a GLR-algorithm-like parsing table so that derivation probabilities can be computed efficiently by making use of the parsing control structure, and add indices to each state in parsing period as identifiers of different parsing paths. Based on the techniques above, the forward and inner probabilities of states are introduced in this paper. With these two probabilities, the probabilities of all parsing trees of the input sentence can be computed to select an optimal parsing result. The experiment shows that the proposed algorithm is efficient.

**Key words** Stochastic context-free grammar, parsing algorithm, GLR algorithm, Earley algorithm, natural language processing.