

# 第5章 句法分析(II)

# 目录

2

- 5.1 CYK分析算法
- 5.2 PCFG, 概率上下文无关文法

## 5.1 CYK分析算法

# 5.1 CYK分析算法

4

- Cocke-Younger-Kasami (CYK) 算法
  - ▣ 对 Chomsky 文法进行范式化:  
 $A \rightarrow a$  或  $A \rightarrow BC$
- 自下而上的分析方法
- 构造  $(n+1) \times (n+1)$  识别矩阵,  $n$  为输入句子长度。
- 假设输入句子  $x = a_1 a_2 \dots a_n$ ,  $a_i$  为构成句子的单词,  $n = |x|$ 。

# 5.1 CYK分析算法

5

## 识别矩阵的构成

- 方阵对角线以下全部为0
- 主对角线以上的元素由文法G的非终结符构成
- 主对角线上的元素由输入句子的终结符号(单词) 构成

	0	1	2	3	4
0	0	P			
1		他	V	VP	
2			喜欢	V	
3				读	N
4					书

## 6



# 5.1 CYK分析算法

7

## □ 识别矩阵构造步骤

- ▣ (1) 首先构造主对角线，令  $t_{0,0}=0$ ，然后，从  $t_{1,1}$  到  $t_{n,n}$  在主对角线的位置上依次放入输入句子  $x$  的单词  $a_i$ 。
- ▣ (2) 构造  $t_{i,i+1}$  主对角线以上紧靠主对角线的元素，其中， $i = 0, 1, 2, \dots, n-1$ 。对于输入句子  $x = a_1 a_2 \dots a_n$ ，从  $a_1$  开始分析。

# 5.1 CYK分析算法

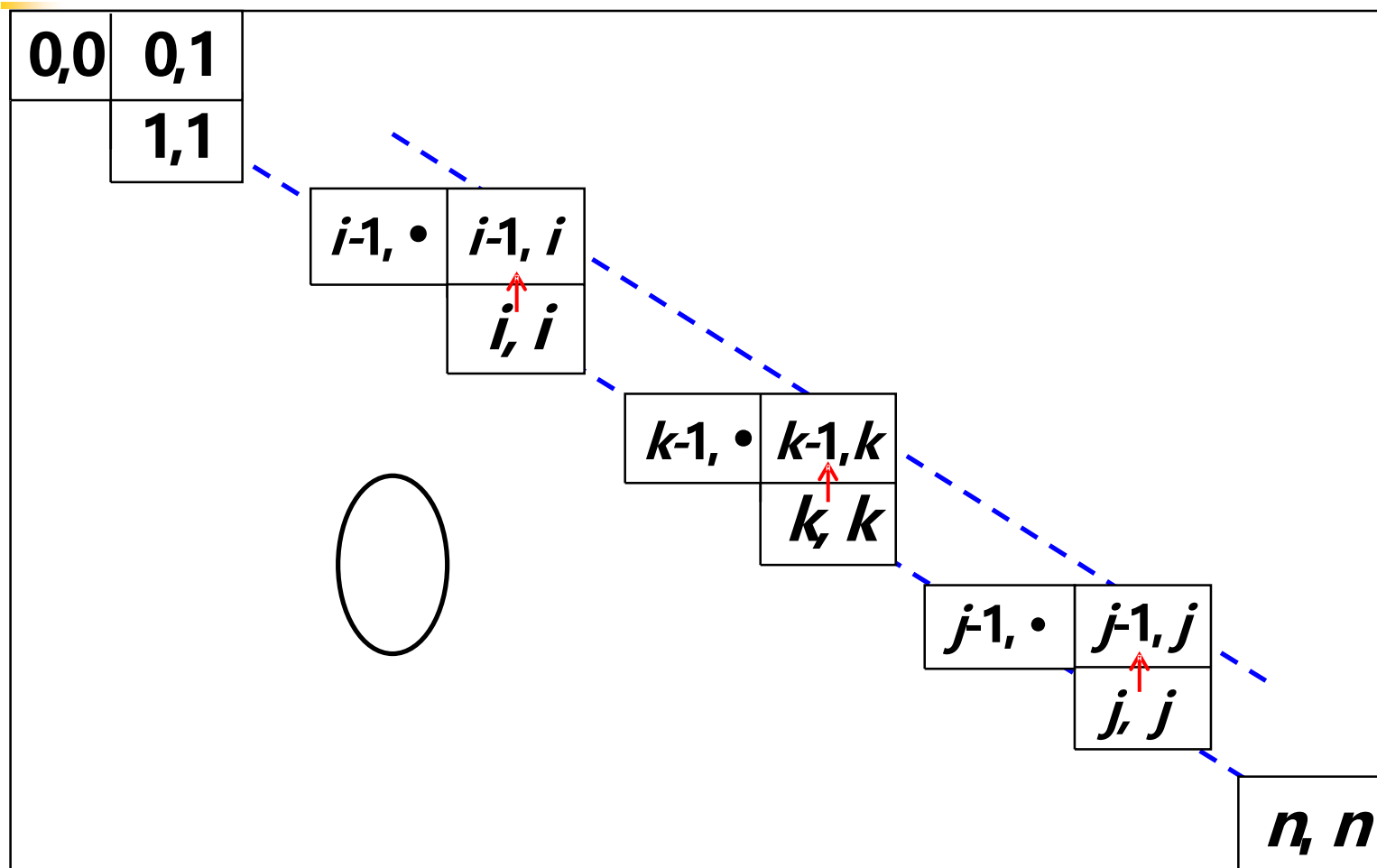
8

- 如果在文法 $G$ 的产生式集合中有一条规则： $A \rightarrow a_1$  则  $t_{0,1}=A$ 。
- 依此类推，如果有  $A \rightarrow a_{i+1}$ ，则  $t_{i,i+1}=A$ 。
- 即，对于主对角线上的每一个终结符 $a_i$ ，所有可能推导出它的非终结符写在它的右边主对角线上方的位置上。



# 5.1 CYK分析算法

9



# 5.1 CYK分析算法

10

- (4)  $P \rightarrow$  他
- (5)  $V \rightarrow$  喜欢
- (6)  $V \rightarrow$  读
- (7)  $N \rightarrow$  书

	0	1	2	3	4
0	0	P			
1		他	V	VP	
2			喜欢	V	
3				读	N
4					书

# 5.1 CYK分析算法

11

- ▣(3) 按平行于主对角线的方向，一层一层地向上填写矩阵的各个元素  $t_{i,j}$ ，其中， $i = 0, 1, \dots, n-d$ ， $j = d+i$ ， $d=2, 3, \dots, n$ 。如果存在一个正整数  $k$ ， $i+1 \leq k \leq j-1$ ，在文法  $G$  的规则集中有产生式  $A \rightarrow BC$ ，并且， $B \in t_{i,k}$ ， $C \in t_{k,j}$ ，那么，将  $A$  写到矩阵  $t_{i,j}$  位置上。

# 5.1 CYK分析算法

12

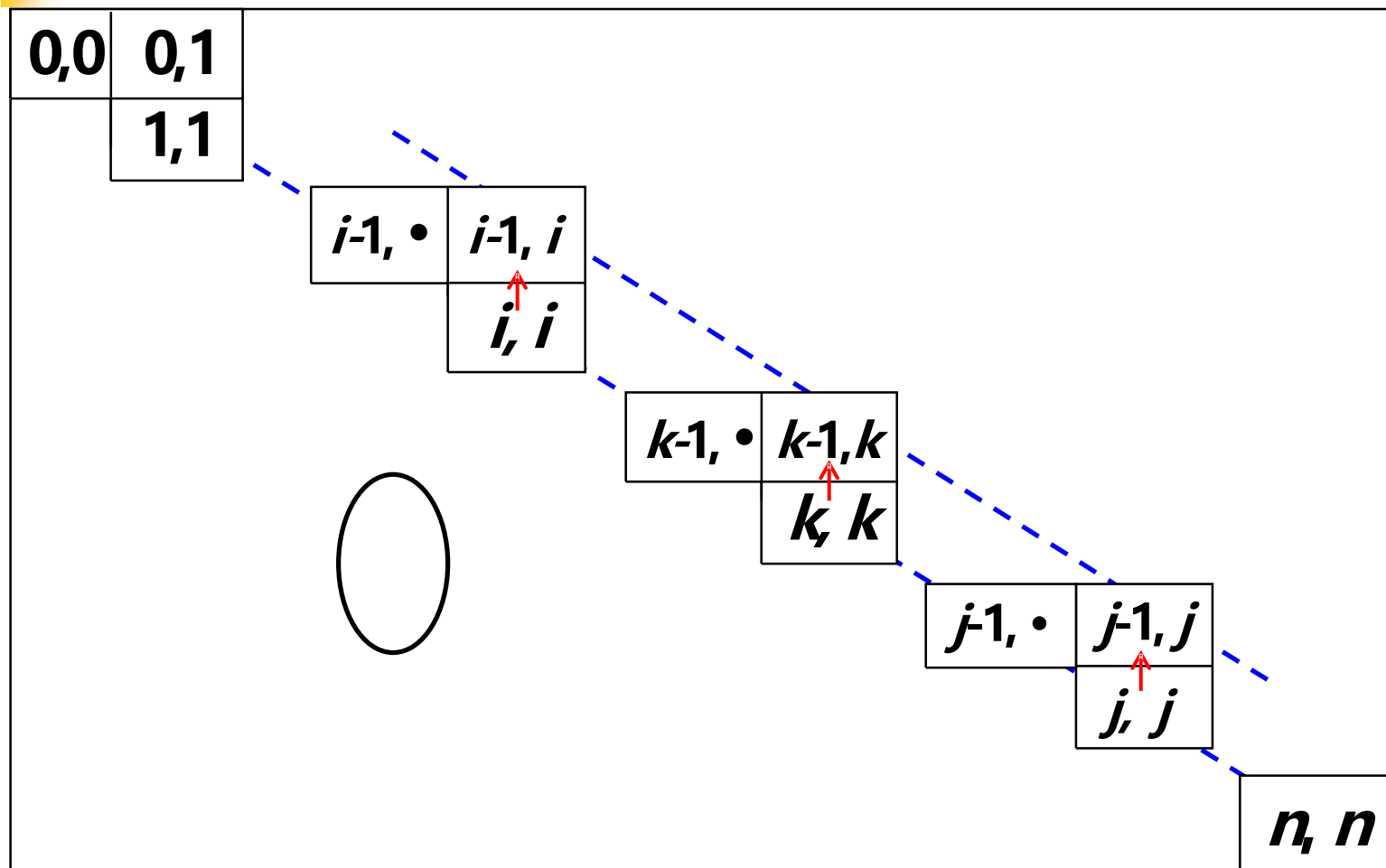
- (2)  $VP \rightarrow V V$
- (4)  $P \rightarrow$  他
- (5)  $V \rightarrow$  喜欢
- (6)  $V \rightarrow$  读
- (7)  $N \rightarrow$  书

	0	1	2	3	4
0	0	P			
1		他	V	VP	
2			喜欢	V	
3				读	N
4					书

$$t_{1,2} = V, t_{2,3} = V \rightarrow t_{1,3} = VP$$

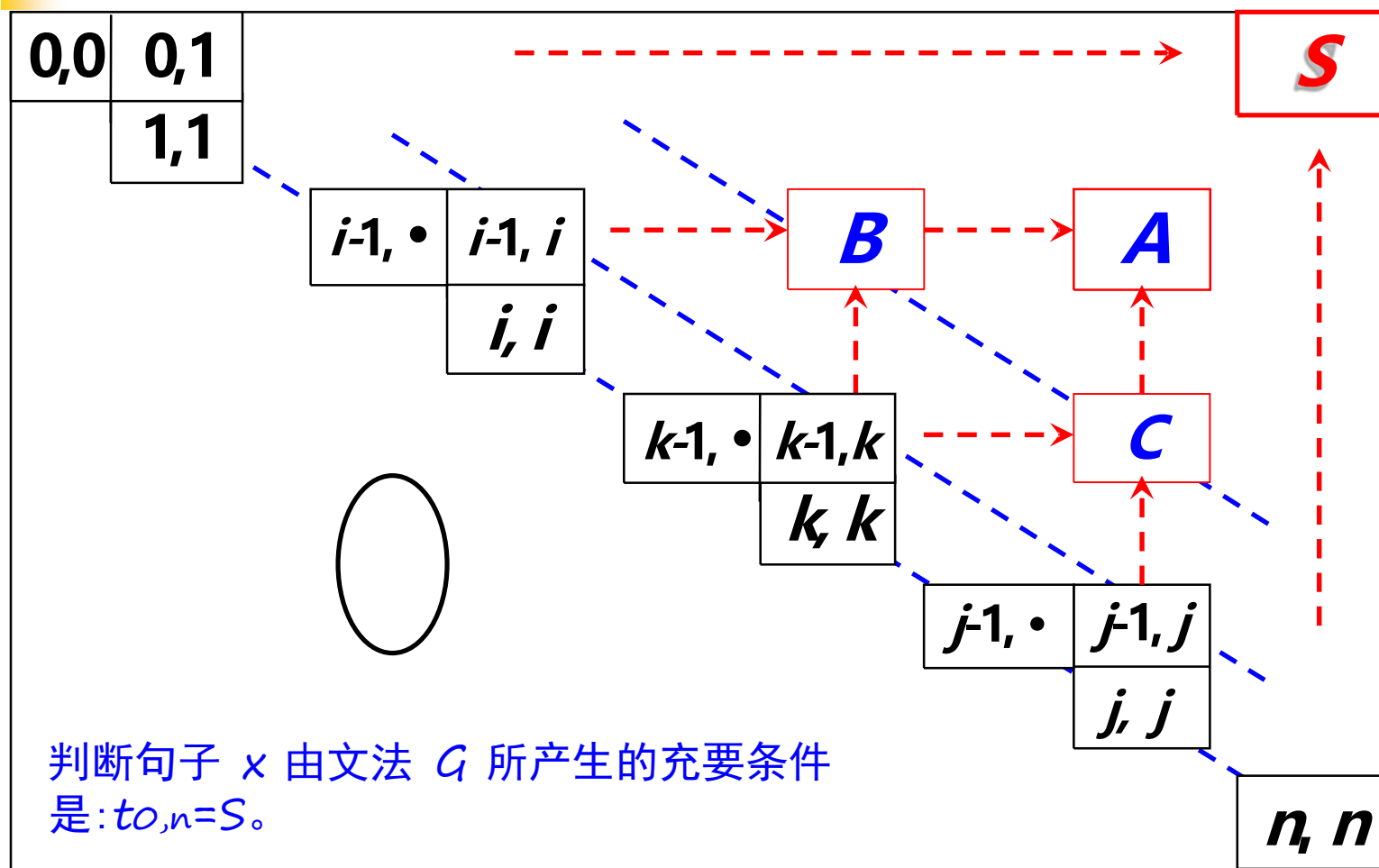
# 5.1 CYK分析算法

13



# 5.1 CYK分析算法

14



# 5.1 CYK分析算法

15

- 例子：给定文法 $G(S)$ ：
  - (1)  $S \rightarrow P VP$
  - (2)  $VP \rightarrow V V$
  - (3)  $VP \rightarrow VP N$
  - (4)  $P \rightarrow \text{他}$
  - (5)  $V \rightarrow \text{喜欢}$
  - (6)  $V \rightarrow \text{读}$
  - (7)  $N \rightarrow \text{书}$
- 请用CYK算法分析句子：他喜欢读书

# 5.1 CYK分析算法

16

□ 汉语分词和词性标注以后：

□ 他/P 喜欢/V 读/V 书/N n=4

□ 构造识别矩阵：

□ 执行分析过程。

(1)  $S \rightarrow P VP$

(2)  $VP \rightarrow V V$

(3)  $VP \rightarrow VP N$

	0	1	2	3	4
0	0	P			
1		他	V	VP	
2			喜欢	V	
3				读	N
4					书



# 5.1 CYK分析算法

17

□ 汉语分词和词性标注以后：

▣ 他/P 喜欢/V 读/V 书/N  $n=4$

□ 构造识别矩阵：

□ 执行分析过程。

(1)  $S \rightarrow P VP$

(2)  $VP \rightarrow V V$

(3)  $VP \rightarrow VP N$

	0	1	2	3	4
0	0	P → P			
1		他	V → VP		
2			喜欢	V	N
3				读	N
4					书

# 5.1 CYK分析算法

18

□ 汉语分词和词性标注以后：

▣ 他/P 喜欢/V 读/V 书/N n=4

□ 构造识别矩阵：

□ 执行分析过程。

(1)  $S \rightarrow P VP$

(2)  $VP \rightarrow V V$

(3)  $VP \rightarrow VP N$

	0	1	2	3	4
0	0	P	P	S	
1		他	V	VP	?
2			喜欢	V	N
3				读	N
4					书

# 5.1 CYK分析算法

19

□ 汉语分词和词性标注以后：

▣ 他/P 喜欢/V 读/V 书/N  $n=4$

□ 构造识别矩阵：

□ 执行分析过程。

(1)  $S \rightarrow P VP$

(2)  $VP \rightarrow V V$

(3)  $VP \rightarrow VP N$

	0	1	2	3	4
0	0	P → P			
1		他	V → VP → VP		
2			喜欢	V → N	
3				读	N
4					书

# 5.1 CYK分析算法

20

□ 汉语分词和词性标注以后：

▣ 他/P 喜欢/V 读/V 书/N  $n=4$

□ 构造识别矩阵：

□ 执行分析过程。

(1)  $S \rightarrow P VP$

(2)  $VP \rightarrow V V$

(3)  $VP \rightarrow VP N$

	0	1	2	3	4
0	0	P → P → P			
1		他	V → VP → VP		
2			喜欢	V	N
3				读	N
4					书

# 5.1 CYK分析算法

21

□ 汉语分词和词性标注以后：

▣ 他/P 喜欢/V 读/V 书/N  $n=4$

□ 构造识别矩阵：

□ 执行分析过程。

(1)  $S \rightarrow P VP$

(2)  $VP \rightarrow V V$

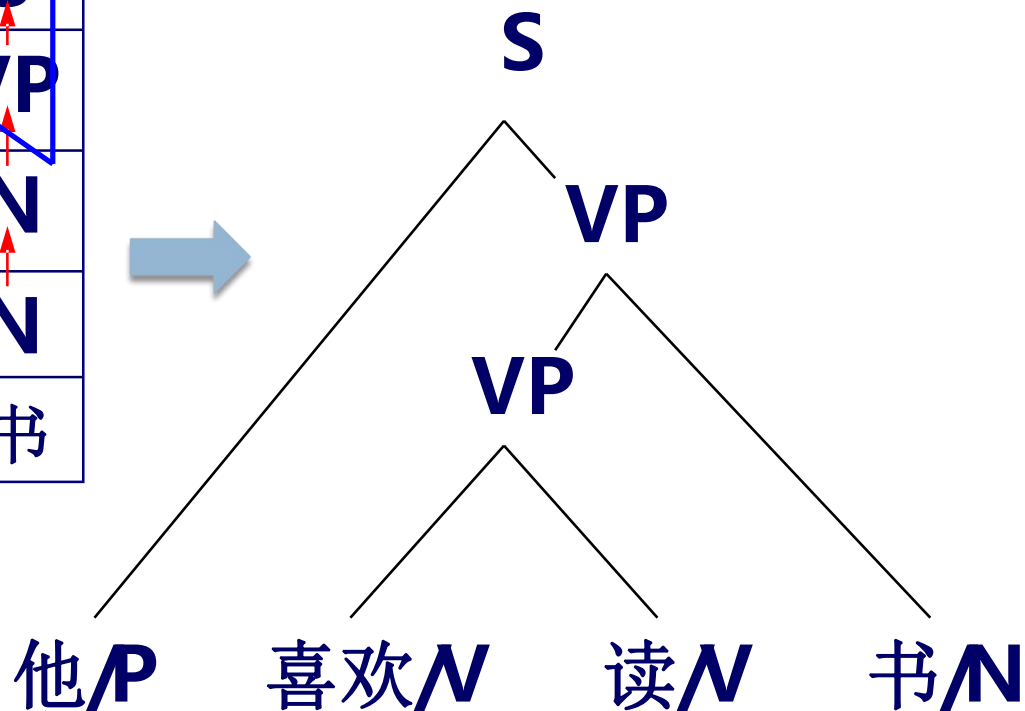
(3)  $VP \rightarrow VP N$

	0	1	2	3	4
0	0	P → P → P → S			
1		他	V → VP → VP		
2			喜欢	V → N	
3				读	N
4					书

# 5.1 CYK分析算法

22

	0	1	2	3	4
0	0	P →	P →	P →	S
1		他	V →	VP →	VP
2			喜欢	V	N
3				读	N
4					书



# 5.1 CYK分析算法

23

## □ CYK算法的评价

### □ 优点

- 简单易行，执行效率高

### □ 弱点

- 必须对文法进行范式化处理
- 无法区分歧义

## 5.2 概率上下文无关文法PCFG



# 5.2 PCFG

25

## □ PCFG规则

形式:  $A \rightarrow \alpha, p$

约束:  $\sum_{\alpha} p(A \rightarrow \alpha) = 1$

例如:  $\left. \begin{array}{l} \text{NP} \rightarrow \text{NN NN}, 0.60 \\ \text{NP} \rightarrow \text{NN CC NN}, 0.40 \end{array} \right\} \sum p = 1$

$\left. \begin{array}{l} \text{CD} \rightarrow \text{QP}, 0.99 \\ \text{CD} \rightarrow \text{LST}, 0.01 \end{array} \right\} \sum p = 1$

## 5.2 PCFG

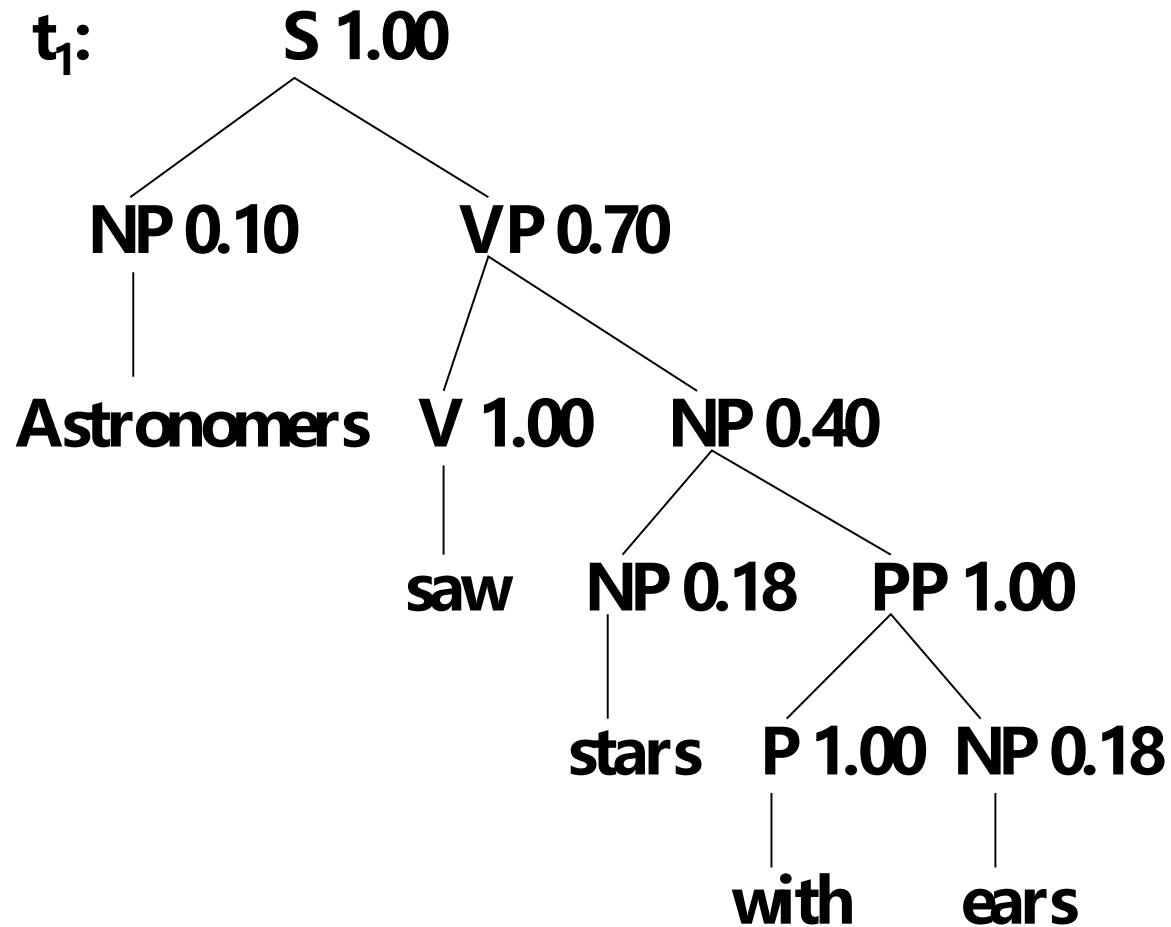
26

◆例-1:  $S \rightarrow NP VP, 1.00$        $NP \rightarrow NP PP, 0.40$   
          $NP \rightarrow \text{astronomers}, 0.10$   
          $NP \rightarrow \text{ears}, 0.18$        $NP \rightarrow \text{saw}, 0.04$   
          $NP \rightarrow \text{stars}, 0.18$        $NP \rightarrow \text{telescopes}, 0.1$   
          $PP \rightarrow P NP, 1.00$        $P \rightarrow \text{with}, 1.00$   
          $VP \rightarrow V NP, 0.70$        $VP \rightarrow VP PP, 0.30$   
          $V \rightarrow \text{saw}, 1.00$

给定句子 S: *Astronomers saw stars with ears.*

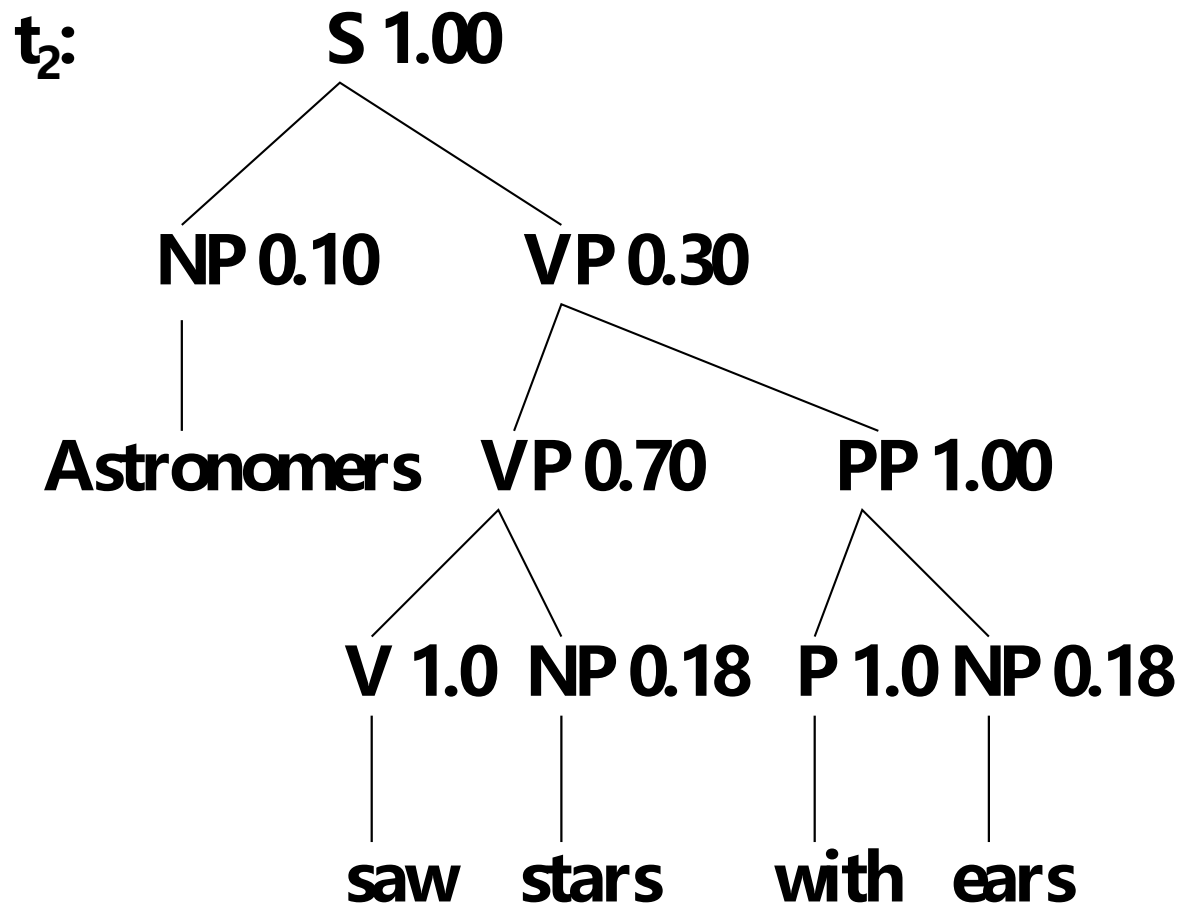
# 5.2 PCFG

27



# 5.2 PCFG

28



# 5.2 PCFG

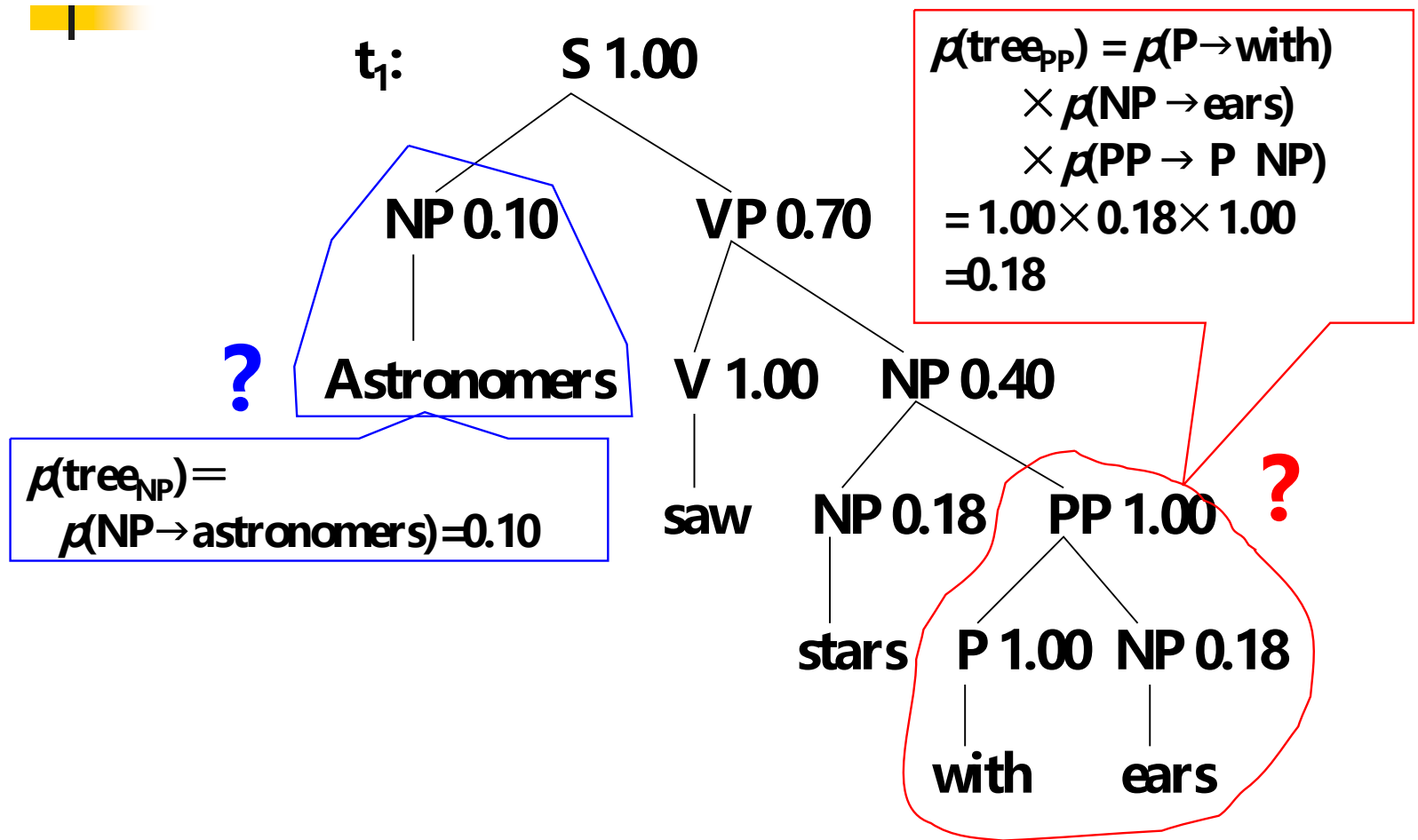
29

## □ 计算分析树概率的基本假设

- **位置不变性**: 子树的概率与其管辖的词在整个句子中所处的位置无关, 即对于任意的  $k$   
 $p(A_{k(k+C)} \rightarrow w)$  一样。
- **上下文无关性**: 子树的概率与子树管辖范围以外的词无关, 即  $p(A_{kl} \rightarrow w / \text{任何超出 } k \sim l \text{ 范围的上下文}) = p(A_{kl} \rightarrow w)$ 。
- **祖先无关性**: 子树的概率与推导出该子树的祖先结点无关, 即  $p(A_{kl} \rightarrow w / \text{任何除 } A \text{ 以外的祖先结点}) = p(A_{kl} \rightarrow w)$ 。

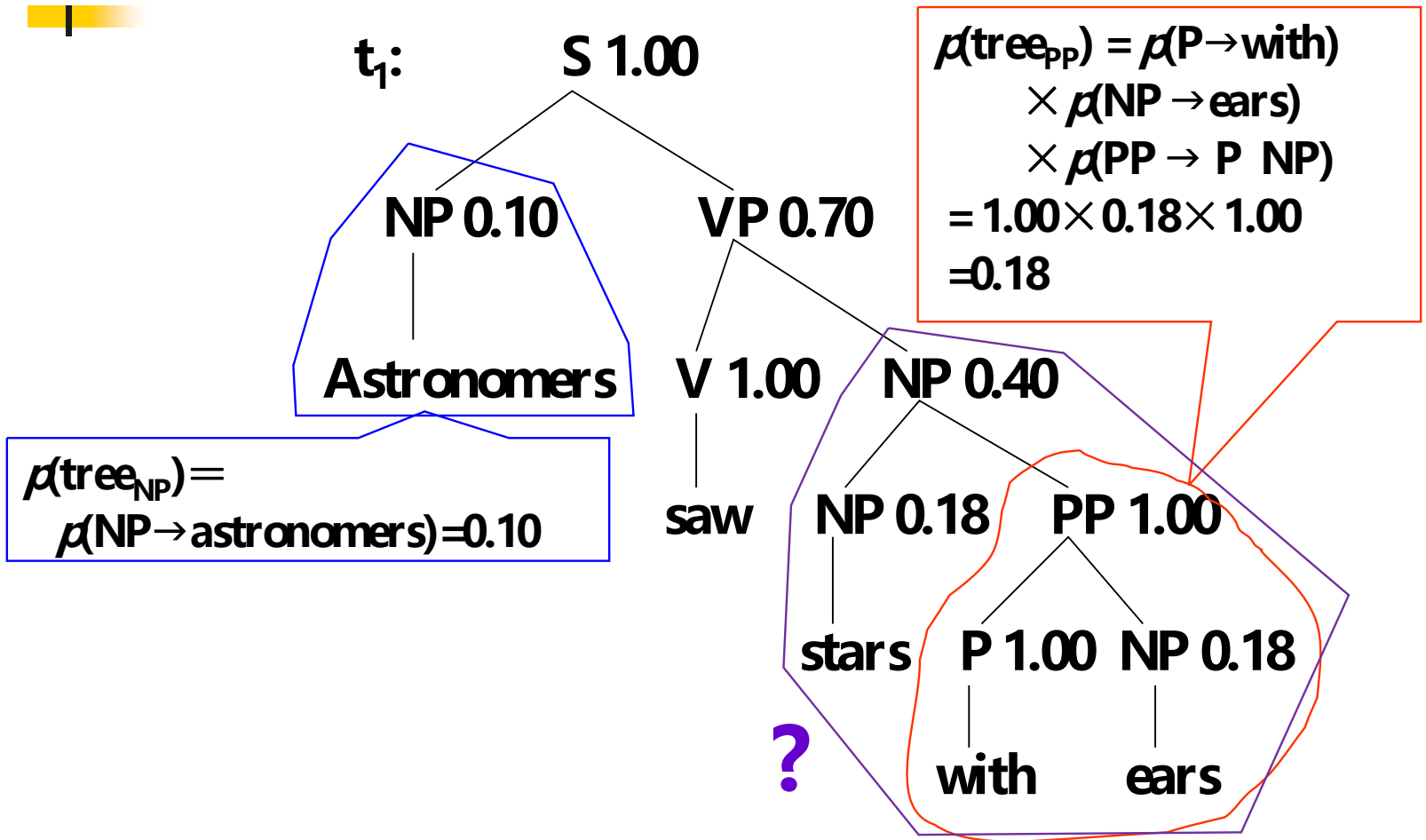
# 5.2 PCFG

30



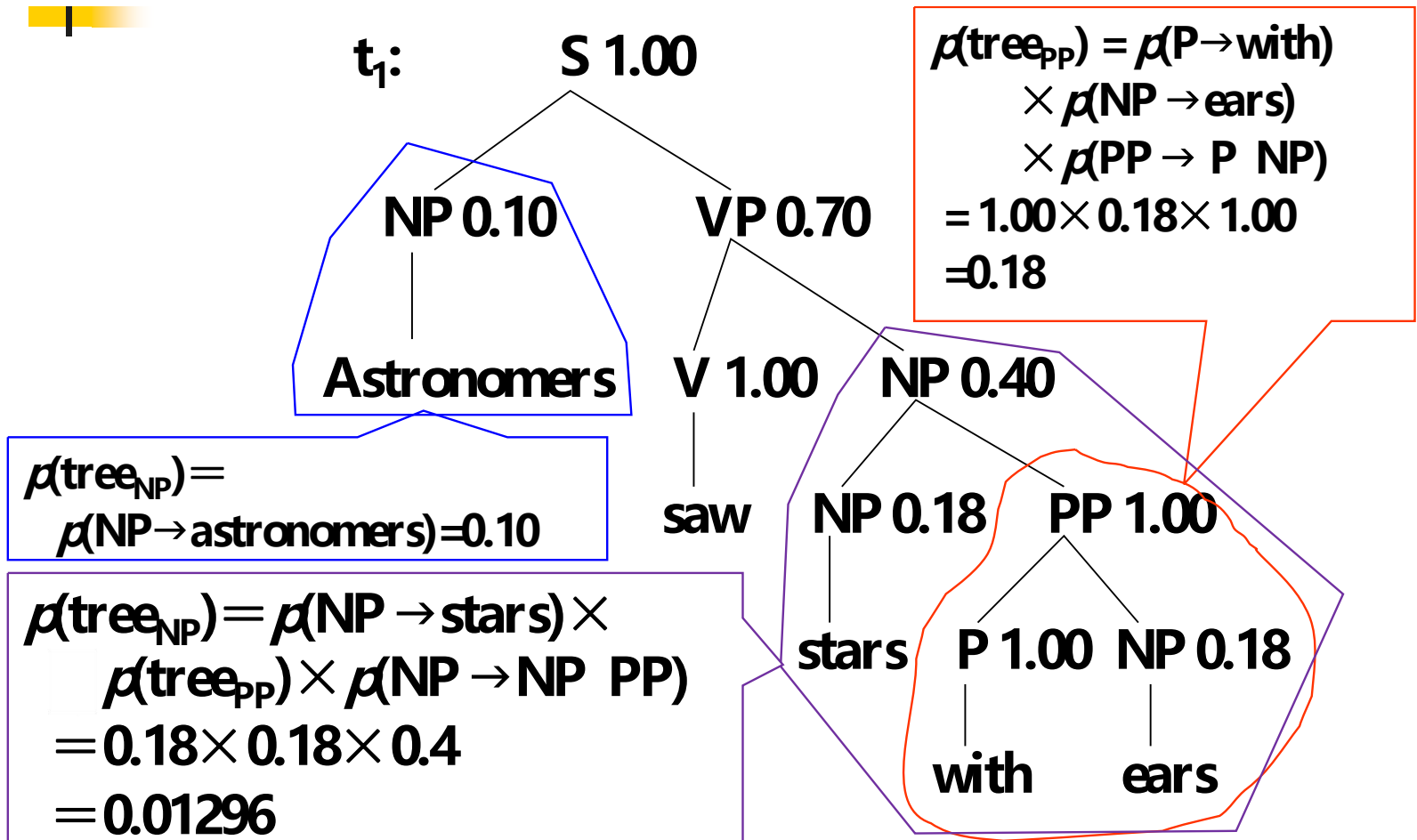
# 5.2 PCFG

31



# 5.2 PCFG

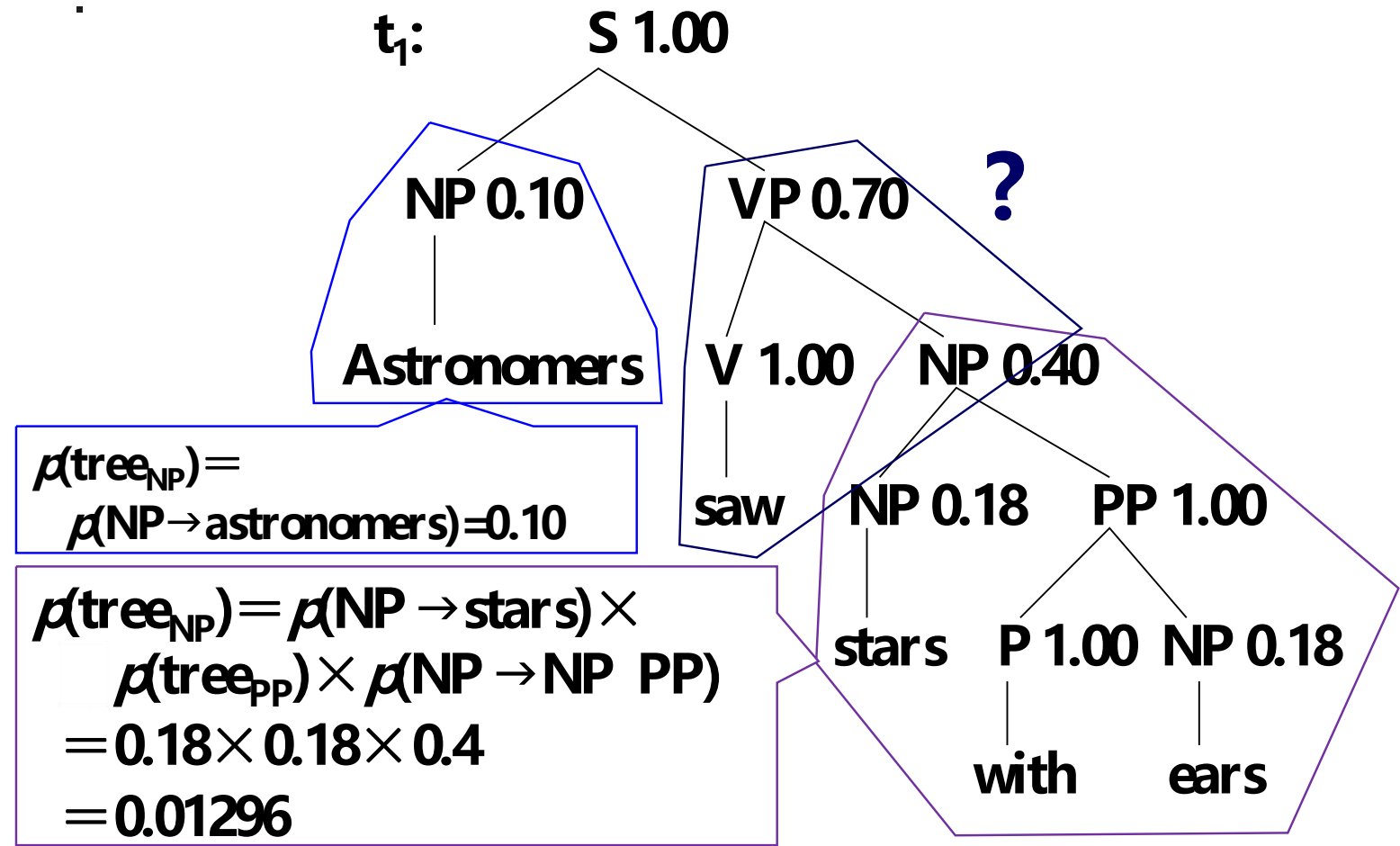
32





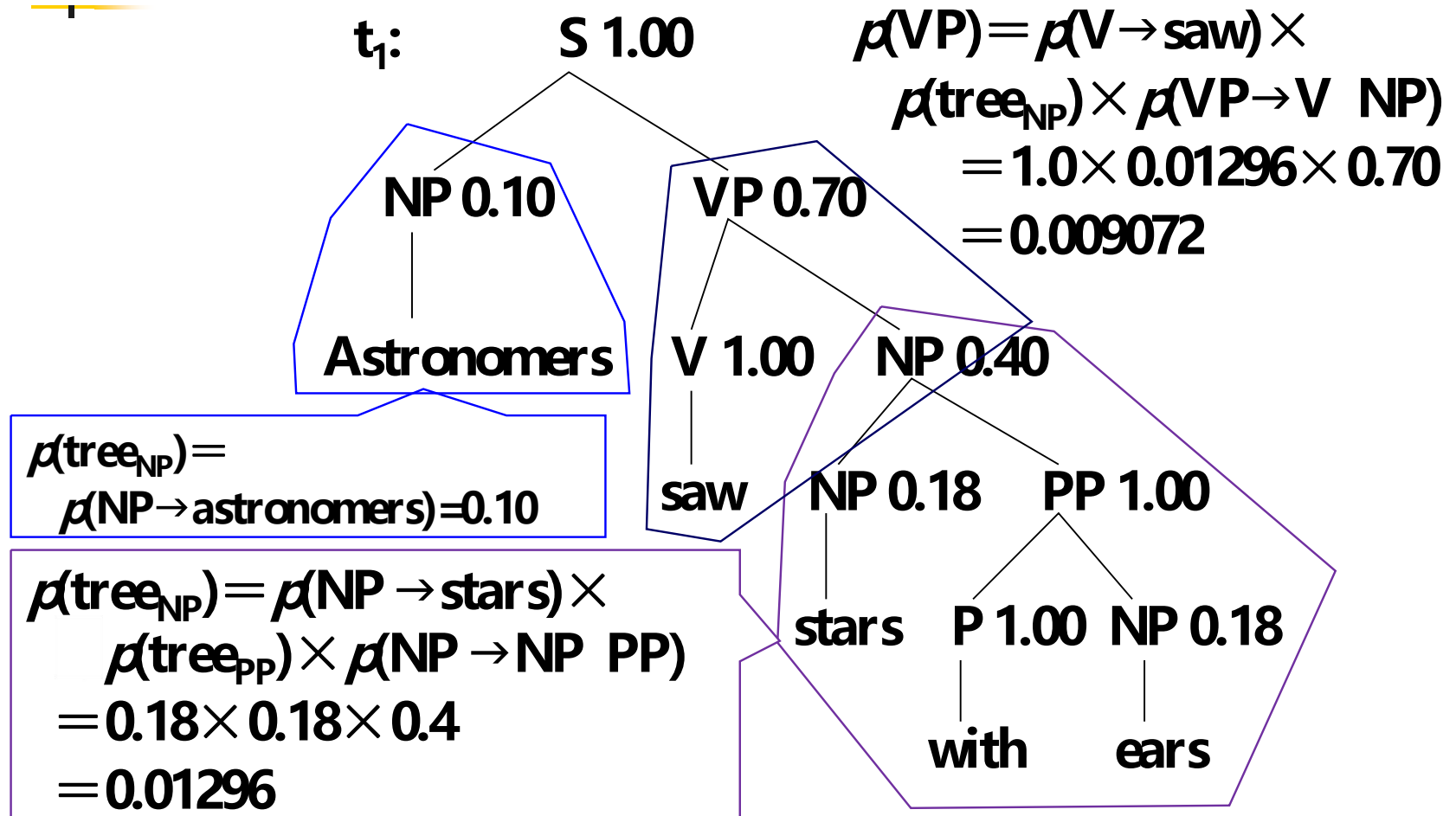
# 5.2 PCFG

33



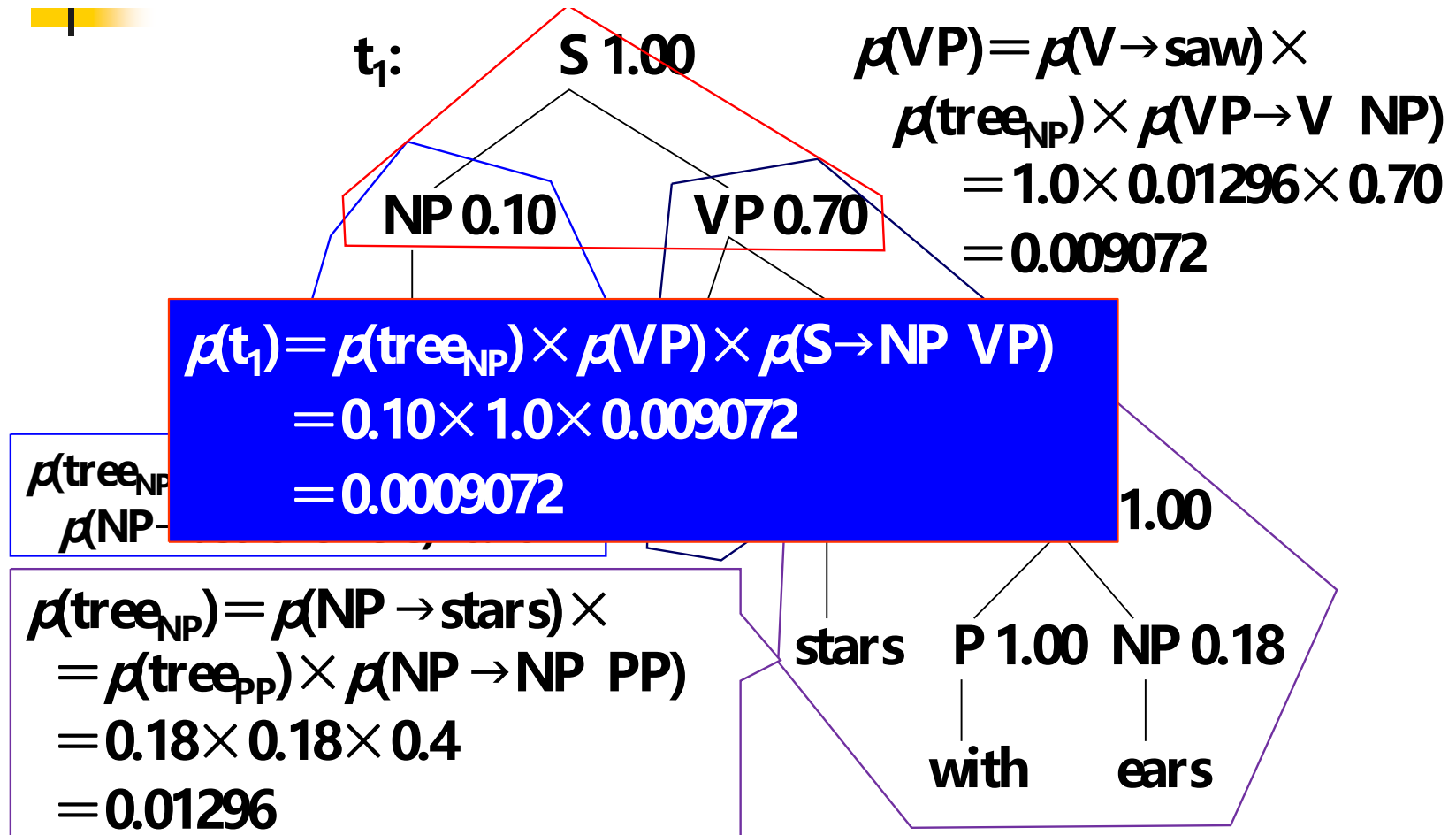
# 5.2 PCFG

34



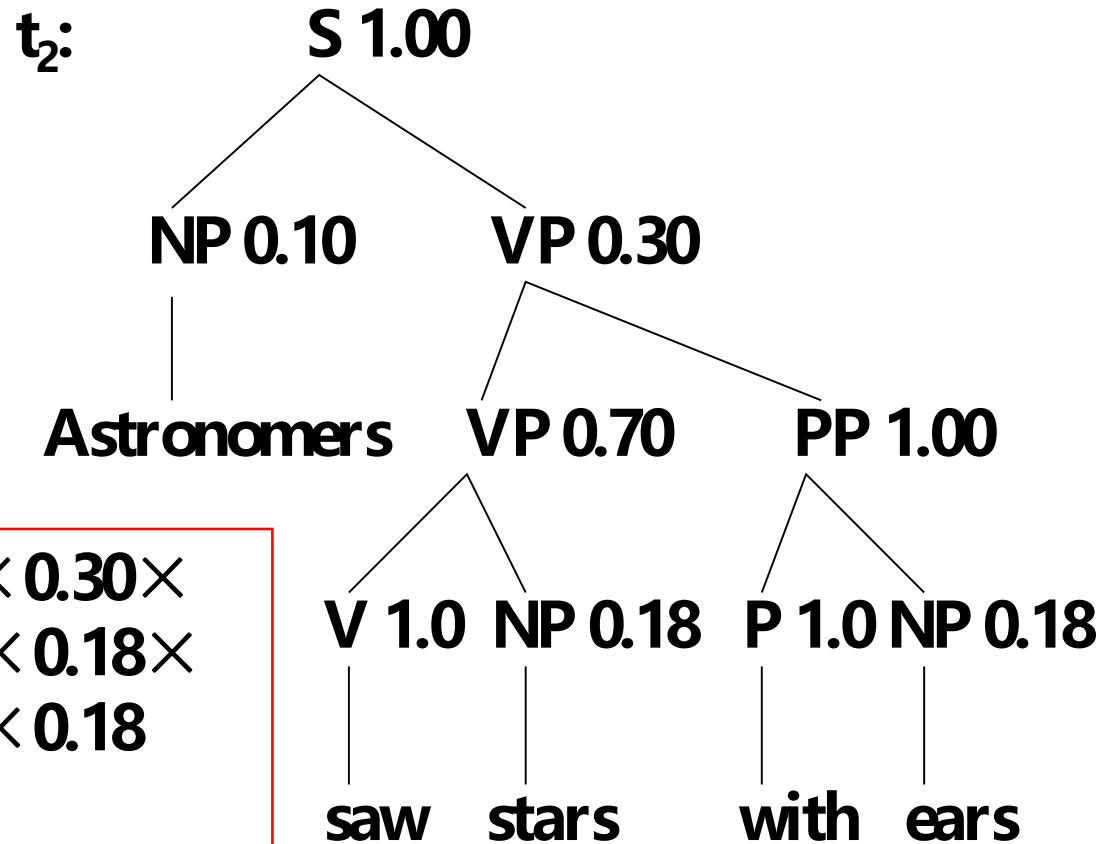
# 5.2 PCFG

35



# 5.2 PCFG

36



$$\begin{aligned} p(t_2) &= 1.00 \times 0.10 \times 0.30 \times \\ &\quad 0.70 \times 1.00 \times 0.18 \times \\ &\quad 1.00 \times 1.00 \times 0.18 \\ &= 0.0006804 \end{aligned}$$

## 5.2 PCFG

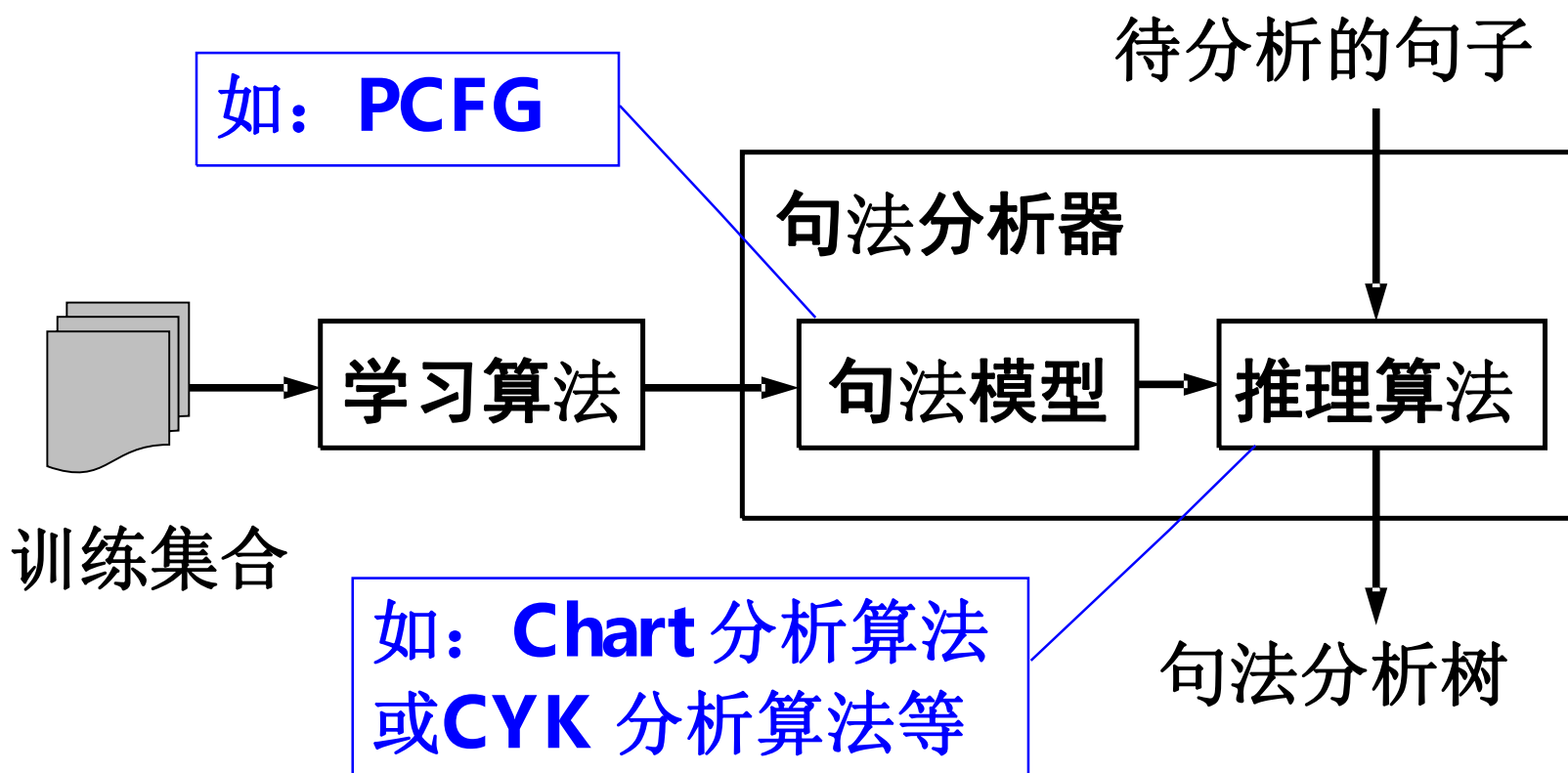
37

- 对于给定的句子S, 两棵句法分析树的概率不等,  $P(t_1) > P(t_2)$ , 因此, 可以得出结论:
- 分析结果 $t_1$ 正确的可能性大于 $t_2$ 。

## 5.2 PCFG的分析实例

38

统计句法分析器实现的一般方法：



## 5.2 PCFG的分析实例

39

给定如下 **PCFG G(S)**:

非终结符集合:  $N = \{S, NP, VP, PP, DT, Vi, Vt, NN, IN\}$

终结符集合:  $= \{\text{sleeps, saw, man, woman, dog, telescope, the, with, in}\}$

规则集:	(1) $S \rightarrow NP VP$	1.0	(9) $Vt \rightarrow \text{saw}$	1.0
	(2) $VP \rightarrow Vi$	0.3	(10) $NN \rightarrow \text{boy}$	0.1
	(3) $VP \rightarrow Vt NP$	0.4	(11) $NN \rightarrow \text{girl}$	0.1
	(4) $VP \rightarrow VP PP$	0.3	(12) $NN \rightarrow \text{telescope}$	0.3
	(5) $NP \rightarrow DT NN$	0.8	(13) $NN \rightarrow \text{dog}$	0.5
	(6) $NP \rightarrow NP PP$	0.2	(14) $DT \rightarrow \text{the}$	0.5
	(7) $PP \rightarrow IN NP$	1.0	(15) $DT \rightarrow \text{a}$	0.5
	(8) $Vi \rightarrow \text{sleeps}$	1.0	(16) $IN \rightarrow \text{with}$	0.6
			(17) $IN \rightarrow \text{in}$	0.4

输入句子: **the boy saw the dog with a telescope**

# 5.2 PCFG的分析实例

40

the boy saw the dog with a telescope

DT 0.5 [0,1]	[0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	[1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		[2,3]	[2,4]	[2,5]	[2,6]	[2,7]	[2,8]
			[3,4]	[3,5]	[3,6]	[3,7]	[3,8]
				[4,5]	[4,6]	[4,7]	[4,8]
					[5,6]	[5,7]	[5,8]
						[6,7]	[6,8]
							[7,8]

DT → the, 0.5

NN → boy, 0.1

NP → DT NN, 0.8

第1步



# 5.2 PCFG的分析实例

41

the boy saw the dog with a telescope

DT 0.5 [0,1]	NP 0.004 [0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	NN 0.1 [1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		[2,3]	[2,4]	[2,5]	[2,6]	[2,7]	[2,8]
			[3,4]	[3,5]	[3,6]	[3,7]	[3,8]
				[4,5]	[4,6]	[4,7]	[4,8]
					[5,6]	[5,7]	[5,8]
						[6,7]	[6,8]
							[7,8]

DT → the, 0.5

NN → boy, 0.1

NP → DT NN, 0.8

第2步

# 5.2 PCFG的分析实例

42

the boy saw the dog with a telescope

DT 0.5 [0,1]	NP 0.004 [0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	NN 0.1 [1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		Vt 1.0 [2,3]	[2,4]	[2,5]	[2,6]	[2,7]	[2,8]
			[3,4]	[3,5]	[3,6]	[3,7]	[3,8]
				[4,5]	[4,6]	[4,7]	[4,8]
					[5,6]	[5,7]	[5,8]
						[6,7]	[6,8]
							[7,8]

DT→the, 0.5

NN→boy, 0.1

NP→DT NN, 0.8

Vt→saw, 1.0

第3步

# 5.2 PCFG的分析实例

43

the boy saw the dog with a telescope

<b>DT 0.5</b> [0,1]	<b>NP 0.004</b> [0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	<b>NN 0.1</b> [1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		<b>Vt 1.0</b> [2,3]	[2,4]	[2,5]	[2,6]	[2,7]	[2,8]
			<b>DT 0.5</b> [3,4]	[3,5]	[3,6]	[3,7]	[3,8]
				[4,5]	[4,6]	[4,7]	[4,8]
					[5,6]	[5,7]	[5,8]
						[6,7]	[6,8]
							[7,8]

DT → the, 0.5

NN → boy, 0.1

NP → DT NN, 0.8

Vt → saw, 1.0

第4步

# 5.2 PCFG的分析实例

44

the boy saw the dog with a telescope

DT 0.5 [0,1]	NP 0.004 [0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	NN 0.1 [1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		Vt 1.0 [2,3]	[2,4]	VP 0.08 [2,5]	[2,6]	[2,7]	[2,8]
			DT 0.5 [3,4]	NP 0.2 [3,5]	[3,6]	[3,7]	[3,8]
				NN 0.5 [4,5]	[4,6]	[4,7]	[4,8]
					[5,6]	[5,7]	[5,8]
						[6,7]	[6,8]
							[7,8]

第5步

DT → the, 0.5

NN → boy, 0.1

NP → DT NN, 0.8

Vt → saw, 1.0

NN → dog, 0.5

VP → Vt Np, 0.4

## 5.2 PCFG的分析实例

45

**the boy saw the dog with a telescope**

DT→the, 0.5

NN → boy, 0.1

NP→DT NN, 0.8

Vt→saw, 1.0

NN  $\rightarrow$  dog, 0.5

VP → Vt Np, 0.4

IN→with, 0.6

<b>DT 0.5</b> [0,1]	<b>NP 0.004</b> [0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	<b>NN 0.1</b> [1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		<b>Vt 1.0</b> [2,3]	[2,4]	<b>VP 0.08</b> [2,5]	[2,6]	[2,7]	[2,8]
			<b>DT 0.5</b> [3,4]	<b>NN 0.2</b> [3,5]	[3,6]	[3,7]	[3,8]
				<b>NN 0.5</b> [4,5]	[4,6]	[4,7]	[4,8]
					<b>IN 0.6</b> [5,6]	[5,7]	[5,8]
						[6,7]	[6,8]
							[7,8]

## 第6步

# 5.2 PCFG的分析实例

46

the boy saw the dog with a telescope

<b>DT 0.5</b> [0,1]	<b>NP 0.004</b> [0,2]	[0,3]	[0,4]	[0,5]	[0,6]	[0,7]	[0,8]
	<b>NN 0.1</b> [1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]
		<b>Vt 1.0</b> [2,3]	[2,4]	<b>VP 0.08</b> [2,5]	[2,6]	[2,7]	[2,8]
			<b>DT 0.5</b> [3,4]	<b>NP 0.2</b> [3,5]	[3,6]	[3,7]	[3,8]
				<b>NN 0.5</b> [4,5]	[4,6]	[4,7]	[4,8]
					<b>IN 0.6</b> [5,6]	[5,7]	[5,8]
						<b>DT 0.5</b> [6,7]	[6,8]
							[7,8]

第7步

DT → the, 0.5

NN → boy, 0.1

NP → DT NN, 0.8

Vt → saw, 1.0

NN → dog, 0.5

VP → Vt Np, 0.4

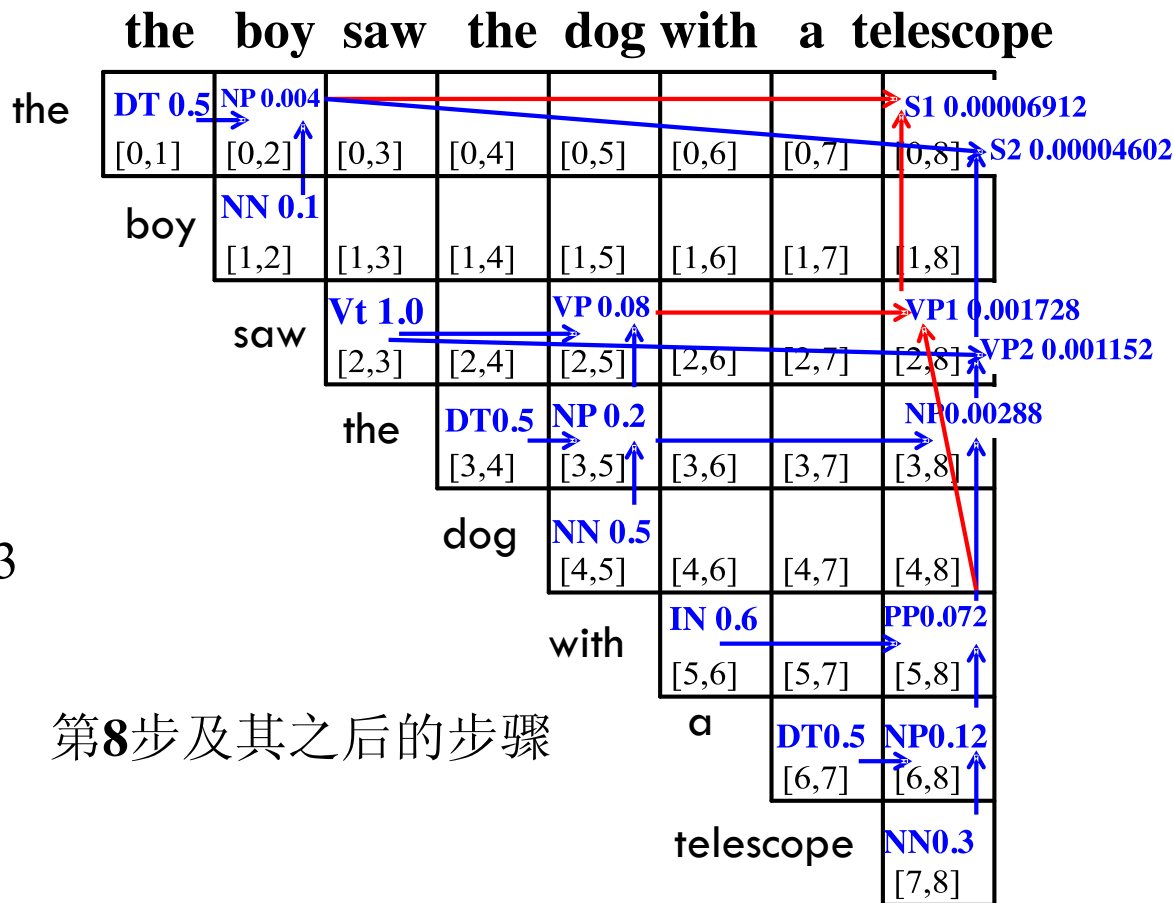
IN → with, 0.6

DT → a, 0.5

# 5.2 PCFG的分析实例

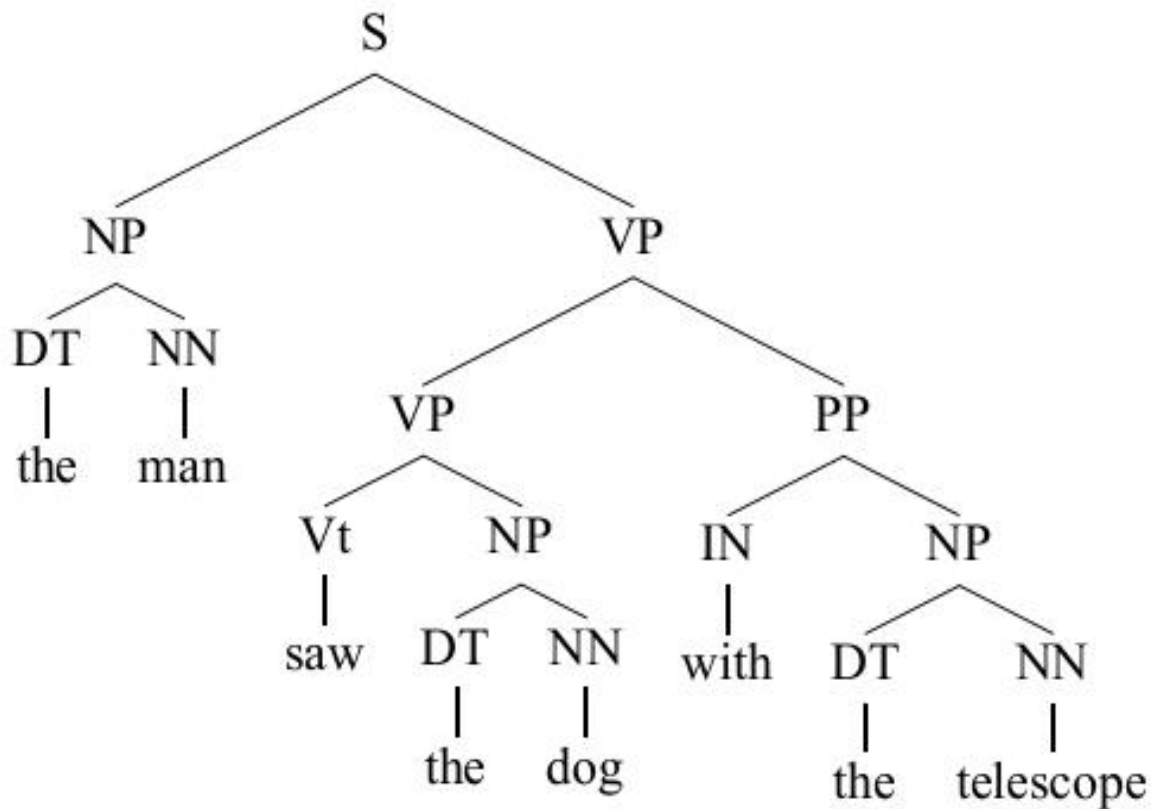
47

DT → the, 0.5  
 NN → boy, 0.1  
 NP → DT NN, 0.8  
 Vt → saw, 1.0  
 NN → dog, 0.5  
 VP → Vt Np, 0.4  
 IN → with, 0.6  
 DT → a, 0.5  
 NN → telescope, 0.3  
 PP → IN NP, 1.0  
 NP → NN PP, 0.2  
 VP → VP PP, 0.3  
 .....



## 5.2 PCFG的分析实例

48

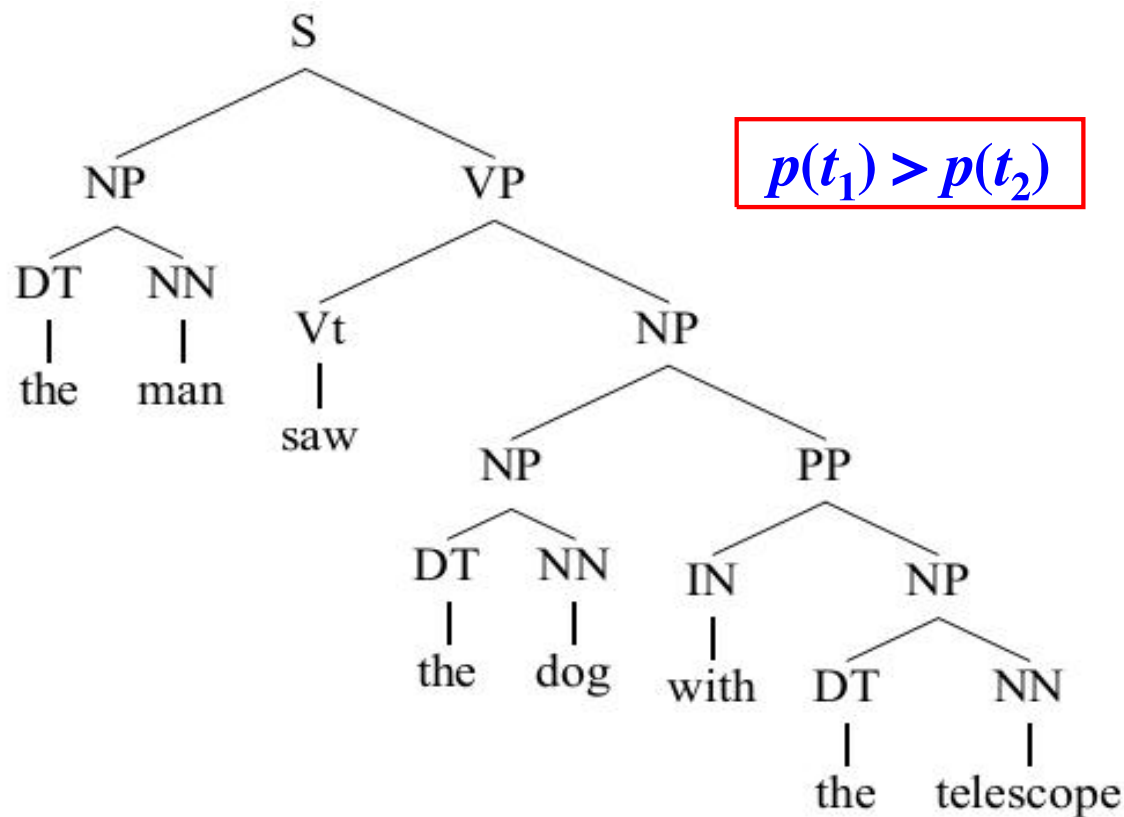


句法树 $t_1$



# 5.2 PCFG的分析实例

49



句法树 $t_2$

# 5.2 PCFG的分析实例

50

## □ PCFG的评价

### □ 优点:

- 可利用概率减少分析过程的搜索空间;
- 可利用概率对概率较小的子树剪枝, 加快分析效率;
- 可以定量地比较两个语法的性能。

### □ 弱点:

- 分析树的概率计算条件非常苛刻, 甚至不够合理。

# 5.2 基于PCFG分析方法的改进

51

- 代表性论文
- **Collins' Parser**
  - ▣ Michael Collins. 2003. Head-driven statistical models for natural language parsing. Computational Linguistics, 29(4): 589-637
- **Charniak Parser**
  - ▣ Eugene Charniak. 2000. A maximum-entropy-inspired parser. In Proceedings of NAACL, pp.132-139
  - ▣ Eugene Charniak and Mark Johnson. 2005. Coarse-to-fine n-best parsing and MaxEnt discriminative reranking. In Proceedings of ACL, pp.173-180
- **Bikel Parser**
  - ▣ Daniel M. Bikel. 2004. Intricacies of Collins' parsing model. Computational Linguistics, 30(4): 479-511

# 5.2 基于PCFG分析方法的改进

52

## □ 代表性论文

## □ Stanford Parser

- ▣ Dan Klein and Christopher D. Manning. 2003. Accurate unlexicalized parsing. In Proceedings of ACL, pp. 423-430

## □ Berkeley Parser

- ▣ Slav Petrov and Dan Klein. 2007. Improved Inference for Unlexicalized Parsing. Proc. NAACL-HLT, pp. 404–411
- ▣ Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning Accurate, Compact, and Interpretable Tree Annotation. Proc. the 21st COLING and 44th ACL, pp.433 - 440

# 5.2 基于PCFG分析方法的改进

53

## □ 开源句法分析器

### □ Collins Parser

- <http://people.csail.mit.edu/mcollins/code.html>

### □ Bikel Parser

- <http://www.cis.upenn.edu/~dbikel/software.html#stat-parser>

### □ Charniak Parser

- <http://www.cs.brown.edu/people/ec/#software>

### □ Oboe Parser (可执行程序)

- <http://www.openpr.org.cn/index.php/NLP-Toolkit-for-Natural-Language-Processing/>

# 5.2 基于PCFG分析方法的改进

54

- 开源句法分析器
- Berkeley Parser
  - ▣ <http://nlp.cs.berkeley.edu/Main.html#Parsing>
- Stanford Parser
  - ▣ <http://nlp.stanford.edu/downloads/lex-parser.shtml>

## 5.2 PCFG的三个问题

55

- 给定句子  $W=w_1w_2...w_n$  和 PCFG  $G$ , 如何快速计算  $p(W|G)$ ?
- 给定句子  $W=w_1w_2...w_n$  和 PCFG  $G$ , 如何快速选择最佳句法结构树?
- 给定句子  $W=w_1w_2...w_n$  和 PCFG  $G$ , 如何调节  $G$  的参数, 使得  $p(W|G)$  最大?

## 5.2 PCFG的三个问题

56

- 内向算法或外向算法解决第一个问题
- 基本思想：利用动态规划算法计算由非终结符 $A$ 推导出的某个字符串片段 $w_i w_{i+1} \dots w_j$ 的概率 $a_{ij}(A)$ 。语句 $W = w_1 w_2 \dots w_n$ 的概率即为文法 $G(S)$ 中 $S$ 推导出的字符串的概率 $a_{1n}(S)$ 。



## 5.2 PCFG的三个问题

57

- Viterbi算法解决第二个问题
- Viterbi变量是由非终结符A推导出语句W中子字符串 $w_i w_{i+1} \dots w_j$ 的最大概率。
- 变量 $\psi$ 用于记忆字符串 $w_1 w_2 \dots w_n$ 的Viterbi语法分析结果。

## 5.2 PCFG的三个问题

58

□ 内外向算法解决第三个问题

□ 基本思路：

- 如果有大量已标注语法结构的训练语料，则可直接通过计算每个语法规则的使用次数，用最大似然估计方法计算 PCFG 规则的概率参数，即：

$$\hat{p}(N^j \rightarrow \zeta) = \frac{C(N^j \rightarrow \zeta)}{\sum_{\gamma} C(N^j \rightarrow \gamma)}$$

## 5.2 PCFG的三个问题

59

- 多数情况下，没有可利用的标注语料，只好借助EM (Expectation Maximization)迭代算法估计PCFG的概率参数。
- 初始时随机地给参数赋值，得到语法 $G_0$ ，依据 $G_0$ 和训练语料，得到语法规则使用次数的期望值，以期望次数运用于最大似然估计，得到语法参数新的估计值，由此得到新的语法 $G_1$ ，由 $G_1$ 再次得到语法规则的使用次数的期望值，然后又可以重新估计语法参数。循环这个过程，语法参数将收敛于最大似然估计值。