

# 基于合一的Earley句法分析

詹卫东

zwd@pku.edu.cn

冯志伟、孙乐译《自然语言处理综论》电子工业出版社2005年版。第11章。

Daniel Jurafsky & James H. Martin, 2000, Speech and Language Processing, Pearson Education, Inc., Prentice Hall.

# 对Earley算法进行改进

- (1) 根据重写规后所附的合一约束为一个节点生成特征结构
- (2) 对chart中的状态进行改进
- (3) 对Predicator, Scanner, Completer等操作进行改进
- (4) 生成新状态时对特征结构的蕴涵关系进行检查

# 带合一约束的CFG规则

- 将重写规则中所附的合一约束转写为节点的特征结构

$S \rightarrow NP VP$

$\langle NP \text{ HEAD AGREEMENT} \rangle = \langle VP \text{ HEAD AGREEMENT} \rangle$

$\langle S \text{ HEAD} \rangle = \langle VP \text{ HEAD} \rangle$

$$\left[ \begin{array}{ll} S & [ \text{HEAD} \quad \boxed{1} ] \\ NP & [ \text{HEAD} \quad [ \text{AGREEMENT} \quad \boxed{2} ] ] \\ VP & [ \text{HEAD} \quad \boxed{1} [ \text{AGREEMENT} \quad \boxed{2} ] ] \end{array} \right]$$

- 特征结构可以表示为有向无环图(DAG)

# 为Earley分析法中的状态 增加特征结构字段

- 改进后的状态包含4部分：
  - (1) 重写规则，代表分析子树
  - (2) 子树分析的完成状况，用点标记·表示
  - (3) 子树完成部分与输入中词的位置对应关系
  - (4) 特征结构

$$(1) \quad \underbrace{S \rightarrow \cdot NP VP}_{(2)} \quad , \quad \underbrace{[0,0]}_{(3)} \quad , \quad \underbrace{DAG}_{(4)}$$

# Predicator操作

- 每当一个状态被Predicator操作加入状态表时，该规则对应的特征结构也作为状态的一个字段加入。

**Predicator:** 对于状态  $Z \rightarrow \alpha \cdot X \beta \ [j, k] \ \text{DAG}_Z$  其中X是非终结符  
对于语法中每条形如  $X \rightarrow \gamma \ \text{DAG}_X$  的规则，都可以形成一个新状态： $X \rightarrow \cdot \gamma \ [k, k] \ \text{DAG}_X$

# Scanner操作

**Scanner:** 对于状态  $Z \rightarrow \alpha \cdot X \beta \ [j, k] \text{ DAG}_z$  其中  $X$  是终结符  
如果  $X$  与输入字符串中第  $k$  个字符匹配，就将词典中  
 $X$  的特征结构  $\text{DAG}_x$  跟  $\text{DAG}_z$  合一，若成功，则形成  
一个新状态：

$$Z \rightarrow \alpha X \cdot \beta \ [j, k+1] \text{ New-DAG}$$

否则，不改变当前状态集。

# Completer操作

**Completer** : 对于一个已经“完成”的状态  $Z \rightarrow \gamma \cdot [j, k] \text{ DAG}_z$

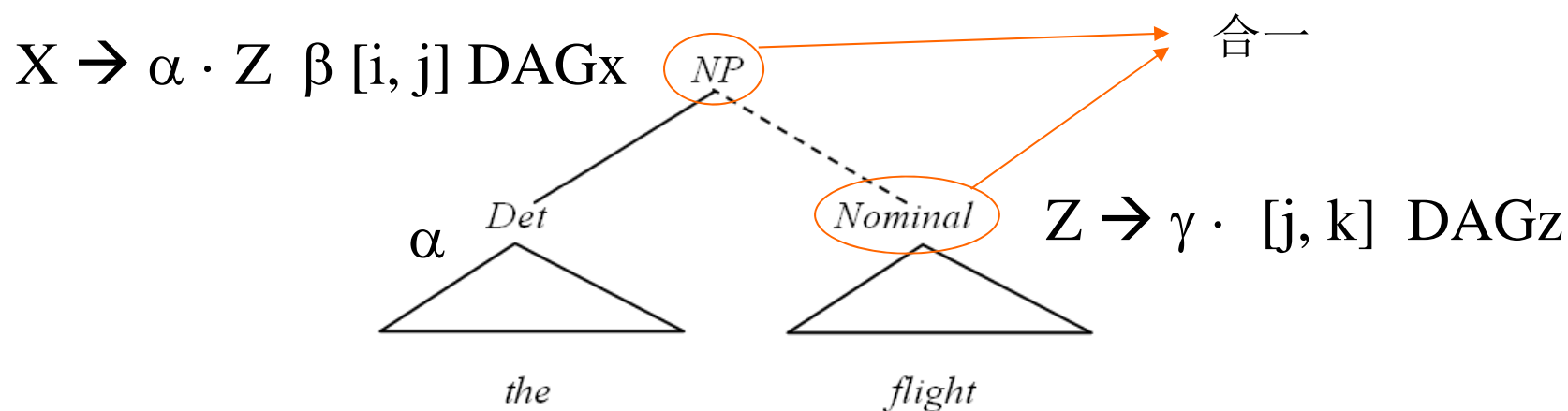
如果已有状态集合中有形如  $X \rightarrow \alpha \cdot Z \beta [i, j] \text{ DAG}_x$  这样的状态，就将  $\text{DAG}_z$  跟  $\text{DAG}_x$  进行合一运算，二者合一的结果（记作 **New-DAG**）若为成功，则形成一个新状态：

$$X \rightarrow \alpha Z \cdot \beta [i, k] \text{ New-DAG}$$

否则，不改变当前状态集。

# Completer操作

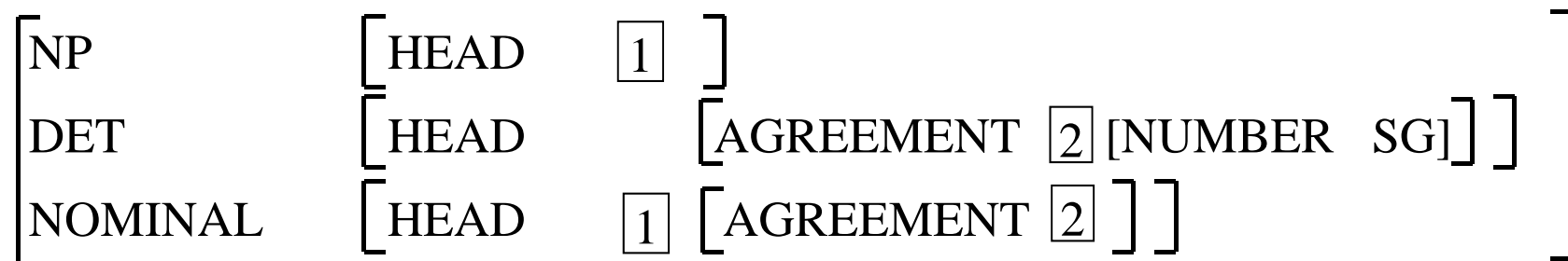
- 每当一个子树( $s_1$ )完成分析时，触发Completer操作
- 检查当前状态集中是否有子树( $s_2$ )等待子树( $s_1$ )来完成分析
- 将 $s_1$ 对应的特征结构与 $s_2$ 对应的特征结构进行合一
- 若合一失败，不产生新的状态
- 若合一成功，则产生新状态，并把合一结果作为新状态的特征结构



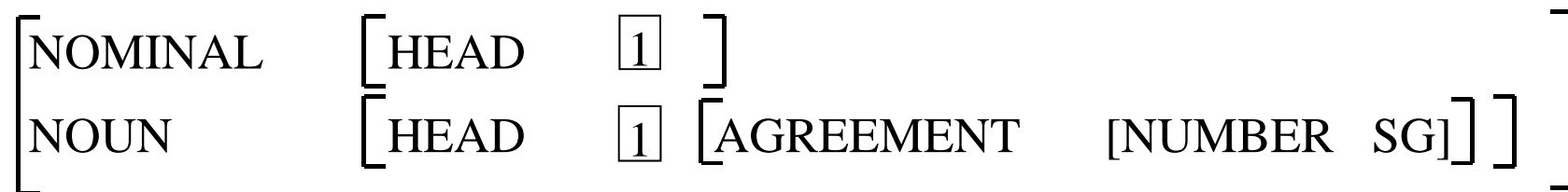


# Completer示例

- $NP \rightarrow Det \cdot Nominal, [0,1], DAG_1$



- $Nominal \rightarrow Noun \cdot, [1,2], DAG_2$



- 将 $DAG_1$ 中的 $Nominal$ 同 $DAG_2$ 中的 $Nominal$ 进行合一

# Completer操作中检查DAG蕴涵关系

- 为了保证不重复分析子树，在Earley算法中，如果新产生的状态和状态集中的某个状态相同，新状态将不被加入状态集
- 在基于合一的Earley算法中，同时还要检查状态所附的特征结构是否具有蕴涵关系，若新状态的特征结构被状态集中的状态的特征结构蕴涵，则不被加入状态集。

例如：

$NP \rightarrow \cdot Det\ NP, [i,i], DAG$

若状态集中的状态对 $Det$ 没有约束，而新产生的状态要求 $Det$ 必须是单数的则新状态不必加入。

# 基于合一的Earley分析算法

设输入字符串长度为 $n$ , 字符间隔可记做 $0,1,2,\dots,n$

(1) 将语法规则中形如  $S \rightarrow \alpha$   $\text{DAG}_s$  的规则形成为状态:

$\langle S \rightarrow \cdot \alpha \ [0, 0] \ \text{DAG}_s \rangle$  加入到状态集合中 (种子状态/seed state)

(2) 对当前分析句子的每个词, 依次进行循环:

对状态集中的每个状态, 依次进行循环:

i) 如果当前状态是[未完成状态], 且点后不是终结符, 则

执行**Predictor**;

ii) 如果当前状态是[未完成状态], 且点后是终结符, 则

执行**Scanner**;

iii) 如果当前状态是[完成状态], 则

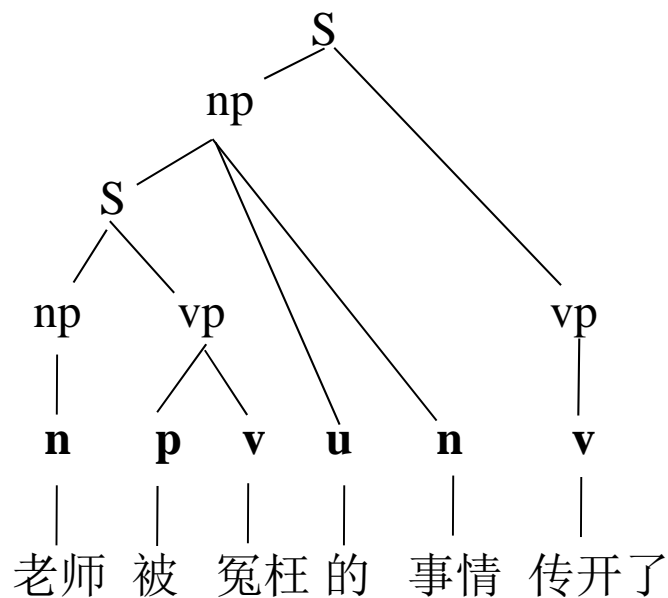
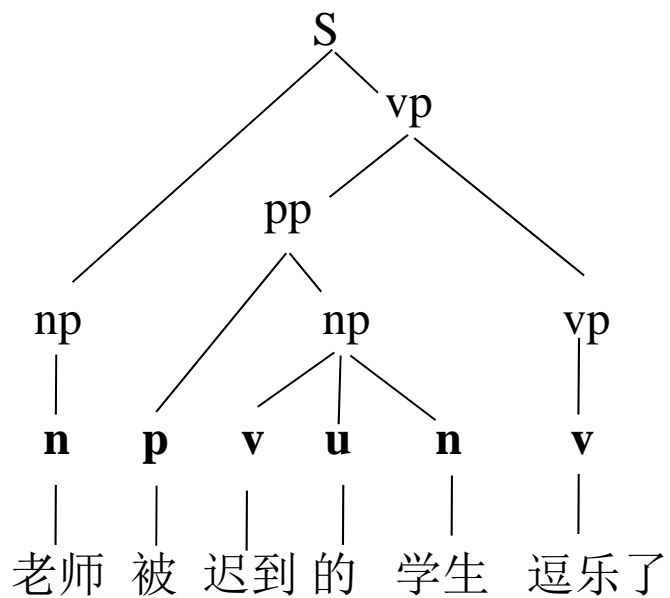
执行**Completer**;

(3) 如果最后得到形如  $\langle S \rightarrow \alpha \cdot \ [0, n] \ \text{DAG}_s \rangle$  这样的状态, 那么输入字符串被接受为合法的句子, 否则分析失败

# 基于合一的Earley分析法示例

- ①  $S \rightarrow np\ vp$
- ②  $vp \rightarrow pp\ vp$
- ③  $pp \rightarrow p\ np$
- ④  $np \rightarrow v\ u\ n$
- ⑤  $np \rightarrow n$
- ⑥  $vp \rightarrow v$

- ⑦  $vp \rightarrow p\ v$
- ⑧  $np \rightarrow S\ u\ n$



# 基于合一的Earley分析法示例

- (1) 迟到 [配价数:1] {施事:[语义类:人]}
- (2) 逗乐 [配价数:2] {施事:[语义类:人], 受事:[语义类:人]}
- (3) 冤枉 [配价数:2] {施事:[语义类:人], 受事:[语义类:人]}
- (4) 传开 [配价数:2] {施事:[语义类:人], 受事:[语义类:事]}

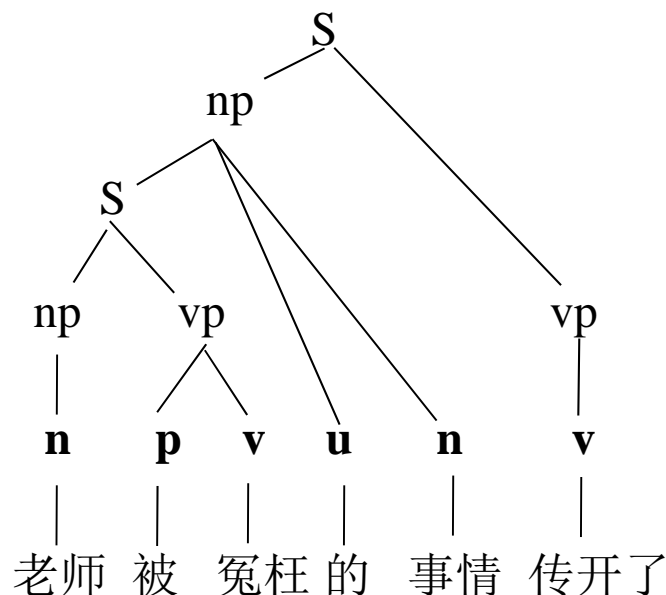
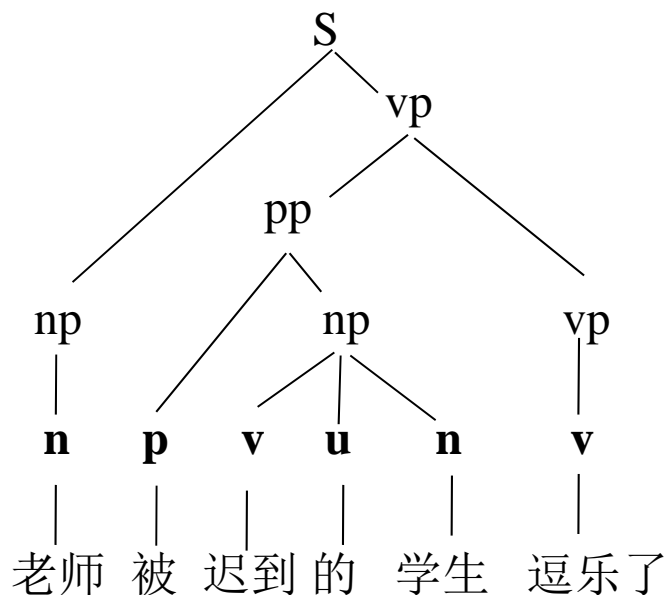
dag1: \$.语态=vp.语态, IF vp.语态=被动 THEN vp.受事=np

dag2: IF pp.lex=被 THEN vp.配价数>1, vp.施事=pp.宾语, \$.语态=被动

dag3: v.施事=n

dag4: IF p.lex=被 THEN v.配价数>1, \$.语态=被动

- ①  $S \rightarrow np\ vp :: dag1$
- ②  $vp \rightarrow pp\ vp :: dag2$
- ③  $pp \rightarrow p\ np$
- ④  $np \rightarrow v\ u\ n :: dag3$
- ⑤  $np \rightarrow n$
- ⑥  $vp \rightarrow v$
- ⑦  $vp \rightarrow p\ v :: dag4$
- ⑧  $np \rightarrow S\ u\ n$



6	$S \rightarrow np\ vp \cdot$ dag1-(2)	$vp \rightarrow pp\ vp \cdot$ dag2-(2)				$vp \rightarrow v \cdot$	归约
5		$vp \rightarrow pp \cdot vp$ $pp \rightarrow p\ np \cdot$	$np \rightarrow v\ u\ n \cdot$ dag3-(1)			$vp \rightarrow \cdot v$ $vp \rightarrow \cdot pp\ vp$ $vp \rightarrow \cdot p\ v$ $pp \rightarrow \cdot p\ np$	预测 归约 扫描
4			$np \rightarrow v\ u \cdot n$				扫描
3		<del><math>vp \rightarrow p\ v \cdot</math> dag4-(1)</del>	$np \rightarrow v \cdot u\ n$ $vp \rightarrow v \cdot$				扫描
2		$pp \rightarrow p \cdot np$ $vp \rightarrow p \cdot v$	$S \rightarrow \cdot np\ vp$ $np \rightarrow \cdot S\ u\ n$ $np \rightarrow \cdot v\ u\ n$ $np \rightarrow \cdot n$				预测 扫描
1	$S \rightarrow np \cdot vp$ $np \rightarrow n \cdot$	$pp \rightarrow \cdot p\ np$ $vp \rightarrow \cdot v$ $vp \rightarrow \cdot pp\ vp$ $vp \rightarrow \cdot p\ v$					预测 归约 扫描
0	$np \rightarrow \cdot S\ u\ n$ $np \rightarrow \cdot v\ u\ n$ $np \rightarrow \cdot n$ $S \rightarrow \cdot np\ vp$						预测 种子
	0	1	2	3	4	5	6

n  
老师

p  
被

v  
迟到

u  
的

n  
学生

v  
逗乐了

6	$S \rightarrow np\ vp \cdot$ dag1-(4)					$vp \rightarrow v \cdot$	归约
5	$S \rightarrow np \cdot vp$ $np \rightarrow S\ u\ n \cdot$		<del><math>np \rightarrow v\ u\ n</math> dag3-(3)</del>			$vp \rightarrow \cdot v$ $vp \rightarrow \cdot pp\ vp$ $vp \rightarrow \cdot p\ v$ $pp \rightarrow \cdot p\ np$	预测 归约 扫描
4	$np \rightarrow S\ u \cdot n$		$np \rightarrow v\ u \cdot n$				扫描
3	$np \rightarrow S \cdot u\ n$ $S \rightarrow np\ vp \cdot$ dag1-(3)	$vp \rightarrow p\ v \cdot$ dag4 - (1)	$np \rightarrow v \cdot u\ n$ $vp \rightarrow v \cdot$				归约 扫描
2		$pp \rightarrow p \cdot np$ $vp \rightarrow p \cdot v$	$S \rightarrow \cdot np\ vp$ $np \rightarrow \cdot S\ u\ n$ $np \rightarrow \cdot v\ u\ n$ $np \rightarrow \cdot n$				预测 扫描
1	$S \rightarrow np \cdot vp$ $np \rightarrow n \cdot$	$pp \rightarrow \cdot p\ np$ $vp \rightarrow \cdot v$ $vp \rightarrow \cdot pp\ vp$ $vp \rightarrow \cdot p\ v$					预测 归约 扫描
0	$np \rightarrow \cdot S\ u\ n$ $np \rightarrow \cdot v\ u\ n$ $np \rightarrow \cdot n$ $S \rightarrow \cdot np\ vp$						预测 种子
	0	1	2	3	4	5	6

n

p

v

u

n

v

老师

被

冤枉

的

事情

传开了