

## Lecture 3

# 上下文无关文法

徐 辉

xuh@fudan.edu.cn

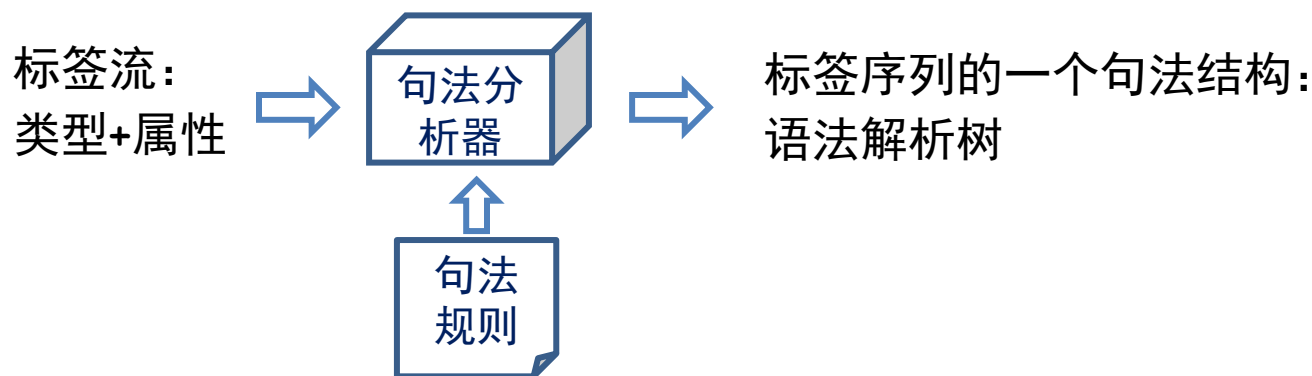


# 主要内容

- 一、上下文无关文法
- 二、语言分析问题

# 句法解析：问题定义

- 给定一个句子和句法规则，找到可生成该句子的一个句法推导
- 通过词法分析已经将句子转换为了标签流
- 句法规则（句法）定义了：
  - 什么是句法分析器可接受的标签序列
  - 及其推导方式
- 语法 = 词法 + 句法
  - 下文将句法统称为语法，但省略了词法分析步骤



# 基本概念

- 一门语言 (language) 是多个句子 (sentences) 的集合。
- 句子 (sentence) 是由终结符 (terminal symbols) 组成的序列 (sequence)。
- 字符串 (string) 是包含终结符和非终结符的序列。
  - 非终结符:  $X$ 、 $Y$ 、 $Z$
  - 终结符 (标签):  $\langle \text{BINOP} \rangle$ 、 $\langle \text{NUM} \rangle$ 
    - 如考虑词法解析: 字符
  - 字符串符号:  $\alpha$ 、 $\beta$ 、 $\gamma$
- 语法 (grammar) 包括一个开始符号  $S$  和多条推导规则 (productions)
  - $S \rightarrow \beta$
  - ...

# 语法推导

- 语法 $G$ 的语言 $L(G)$ 是该语法可推导的所有句子的集合。
- 问题：下列语法是否可推导出句子 $1 + 2 \times 3$ ？

## 语法规则

[1]  $E \rightarrow E + E$   
[2]  $E \rightarrow E - E$   
[3]  $E \rightarrow E \times E$   
[4]  $E \rightarrow E / E$   
[5]  $E \rightarrow \langle \text{NUM} \rangle$

## 推导

[1]  $E \rightarrow E + E$   
[5]  $E \rightarrow 1 + E$   
[3]  $E \rightarrow 1 + E \times E$   
[5]  $E \rightarrow 1 + 2 \times E$   
[5]  $E \rightarrow 1 + 2 \times 3$

# 上线文无关语法和BNF范式

- 上下文无关语法（CFG: Context-Free Grammar）是一个四元组 $(T, NT, S, P)$ 
  - T: 终结符
  - NT: 非终结符
  - S: 起始符号
  - P: 产生式规则集合 $X \rightarrow \gamma$ ,
    - $X$  是非终结符
    - $\gamma$  是可能包含终结符和非终结符的字符串
- BNF范式（Backus-Naur form）：经典CFG语法表示形式
  - $\langle \text{symbol} \rangle ::= \_\_\text{expression}\_\_\_$

# 回顾：括号匹配问题

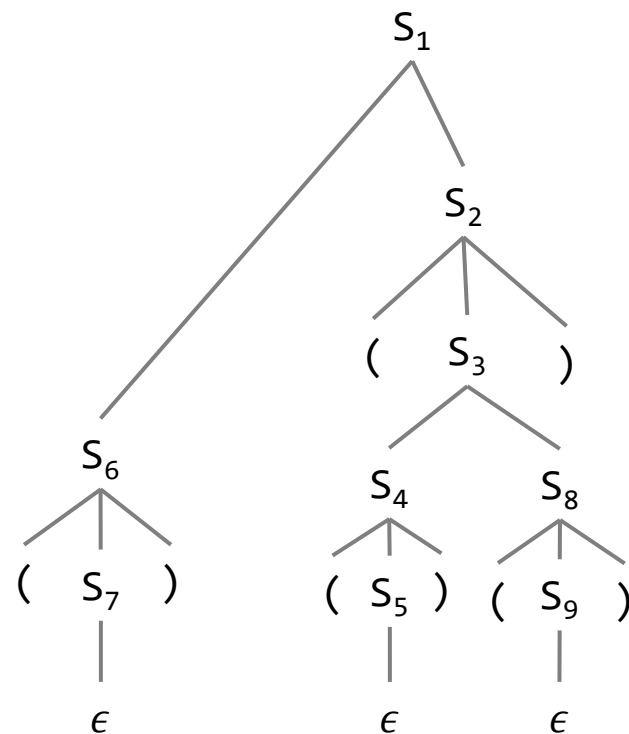
- 用CFG语法设计一套括号匹配规则
- 验证： $()(())()$ 是该语法的一个推导吗？

语法规则

[1]	$S \rightarrow \epsilon$
[2]	$  (S)$
[3]	$  SS$

推导

[3]  $S \rightarrow SS$   
[2]  $S \rightarrow S(S)$   
[3]  $S \rightarrow S(SS)$   
[2]  $S \rightarrow S(S(S))$   
[1]  $S \rightarrow S(S())$   
[2]  $S \rightarrow S((S)())$   
[1]  $S \rightarrow S(())$   
[2]  $S \rightarrow (S)(())$   
[1]  $S \rightarrow ()(())$



语法解析树

# 写出计算器的CFG文法

```
[1] E → E <ADD> E
[2]   | E <SUB> E
[3]   | E <MUL> E
[4]   | E <DIV> E
[5]   | E <POW> E
[6]   | <LPAR> E <RPAR>
[7]   | NUM
[8] NUM → <UNUM>
[9]   | <SUB> <UNUM>
```



# 二义性问题 (ambiguity)

- $L(G)$ 中的某个句子存在一个以上的最左（或最右）推导
- 语法解析树不同

算式:

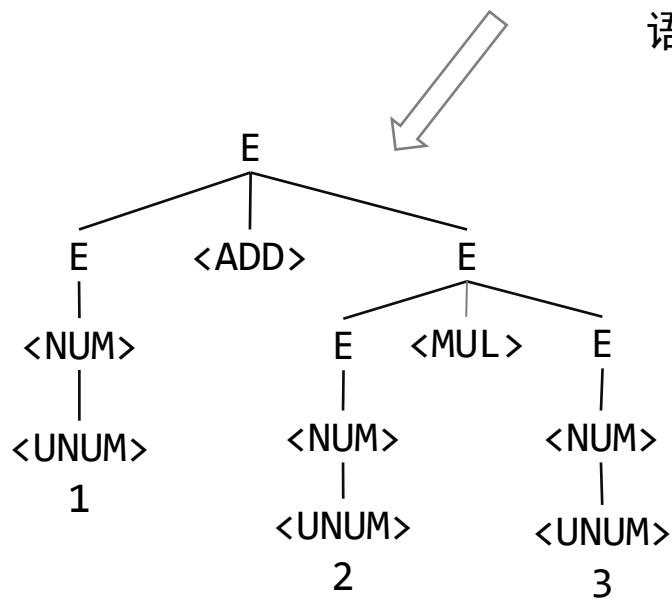
1 + 2 \* 3



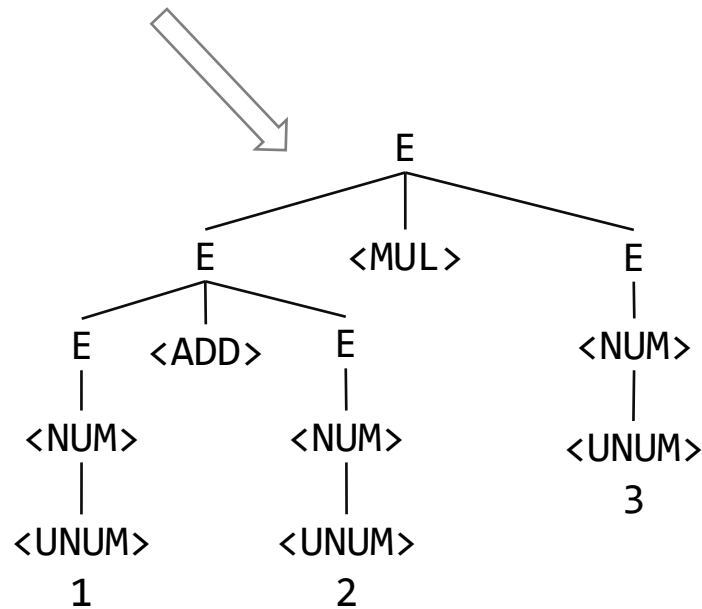
词法解析

标签流: <NUM> <ADD> <UNUM> <MUL> <UNUM>

语法解析



语法解析树1



语法解析树2





A programmer's wife asks him to go to the grocery. She says "Get a gallon of milk. If they have eggs, get 12."

The programmer returns with 12 gallons of milk.

# 消除二义性

- 将运算符特性加入到语法规则中：
  - 优先级： $\wedge > \times / \div > + / -$
  - 结合性： $\times / \div > + / -$  左结合， $\wedge$  右结合

```
[1] E → E <ADD> E
[2]   | E <SUB> E
[3]   | E <MUL> E
[4]   | E <DIV> E
[5]   | E <EXP> E
[6]   | <LPAR> E <RPAR>
[7]   | NUM
[8] NUM → <UNUM>
[9]     | <SUB> <UNUM>
```



```
[1] E → E OP1 E1
[2]   | E1
[3] E1 → E1 OP2 E2
[4]   | E2
[5] E2 → E3 OP3 E2
[6]   | E3
[7] E3 → NUM
[8]   | <LPAR> E <RPAR>
[9] NUM → <UNUM>
[10]    | <SUB> <UNUM>
[11] OP1 → <ADD>
[12]    | <SUB>
[13] OP2 → <MUL>
[14]    | <DIV>
[15] OP3 → <POW>
```

# 练习：语法设计

- 为下列语言设计语法规则：
  - 1) 所有0和1组成的字符串，每一个0后面紧跟着若干个1
  - 2) 所有0和1组成的字符串，0和1的个数相同
  - 3) 所有0和1组成的字符串，0和1的个数不相同

# 练习：语法设计

- 为描述正则语言的正则表达式语法设计一种CFG
  - 支持字符 [A-Za-z0-9]
  - 支持连接、或|、闭包\*运算
  - 支持()
- 检查语法是否有二义性？

## 二、语言分析问题

---

# 概念回顾

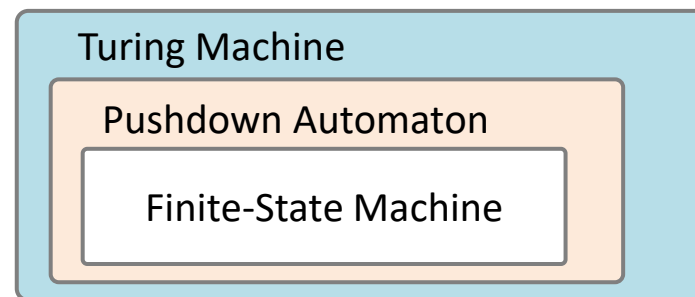
- 一门语言 (language) 是多个句子 (sentences) 的集合。
- 句子 (sentence) 是由终结符 (terminal symbols) 组成的序列 (sequence)。
- 字符串 (string) 是包含终结符和非终结符的序列。
  - 非终结符:  $X$ 、 $Y$ 、 $Z$
  - 终结符 (标签):  $\langle \text{BINOP} \rangle$ 、 $\langle \text{NUM} \rangle$
  - 字符串符号:  $\alpha$ 、 $\beta$ 、 $\gamma$
- 语法 (grammar) 包括一个开始符号  $S$  和多条推导规则 (productions)
  - $S \rightarrow \beta$
  - ...



# 语言分析问题分类：按难度

## Chomsky Hierarchy

类型	文法名称	自动机模型	生成式形式	语言示例
0 型	递归枚举	图灵机	无限制	
1 型	上下文敏感	Linear bounded TM	左侧可以多个符号 $\alpha S \rightarrow \beta$	$a^n b^n c^n$
2 型	上下文无关	下推自动机	左侧仅一个符号 $S \rightarrow \beta$	$a^n b^n$
3 型	正则	有穷自动机	右侧全部为终结符 $S \rightarrow \langle a \rangle \langle b \rangle$	$a^n$



# 正则语言 VS 上下文无关语言

- 正则语言也可以用CFG规则形式表示：
  - $X \rightarrow \gamma$
  - $\gamma \rightarrow \gamma_1$
  - ...
- 特点：右侧的非终结符均可替换为终结符

[1]	$S \rightarrow A B$
[2]	$A \rightarrow (0?1)^*$
[3]	$B \rightarrow (1?0)^*$

 $\Rightarrow S \rightarrow (0?1)^*|(1?0)^*$

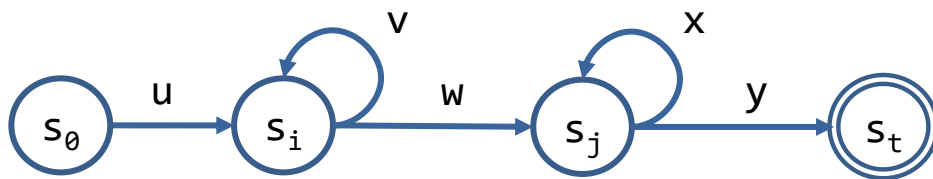
# 非CFG语言：上下文敏感语法

- $L = \{a^n b^n c^n, n > 0\}$  不是CFG语言
- 如何定义？
  - 上下文敏感语法规则形式：  $aS \rightarrow \beta$

[1]	$S \rightarrow aBC$
[2]	$\quad \mid aSBC$
[3]	$CB \rightarrow BC$
[4]	$aB \rightarrow ab$
[5]	$bB \rightarrow bb$
[6]	$bC \rightarrow bc$
[7]	$cC \rightarrow cc$

# 非CFG语言的泵引理

- CFG语言的泵引理（必要条件）：
  - 任意长度超过 $p$ （泵长）的句子可以被拆分为 $uvwxy$ ,
  - 子句 $v$ 和 $x$ 被重复任意次后得到的新句子（如 $uvvwxxy$ ）仍属于该语言。
- 正则属于CFG:  $uv^n w \epsilon^n \epsilon$



# 练习：下列语言是否为正则语言？

- 集合表示

1)  $L = \{a^n b^n | n \leq 100\}$

2)  $L = \{a^n | n \geq 1\}$

3)  $L = \{a^{2^n} | n \geq 1\}$

4)  $L = \{a^p | p \text{ is prime}\}$

- Regex/CFG语法表示

1)  $S \rightarrow (0? 1)^*$

2)  $S \rightarrow aT | \epsilon, T \rightarrow Sb$

3)  $S \rightarrow 0S1S | 1S0S | \epsilon$

# 思考

- 1) 用正则表达式可以定义所有的正则语言吗？
- 2) 有穷自动机可以解析任意正则表达式吗？
- 3) 用CFG可以定义任意正则语言吗？
- 4) 用CFG可以定义任意上下文无关语言吗？
- 5) 用下推自动机可以解析任意正则表达式吗？
- 6) 用下推自动机可以解析任意CFG吗？
- 7) 用通用图灵机可以解析任意CFG吗？
- 8) 用通用图灵机可以解析任意程序吗？