

上下文无关文法和下推自动机

上下文无关文法

上下文无关文法可以表示大多数程序设计语言的语法，尤其是程序中配对或者嵌套的结构，应用领域有程序语言的设计、编译器的实现、XML的格式类型定义DTD

形式定义

上下文无关文法(CFG), G 是一个四元组

$$G = (V, T, P, S)$$

- 字母定义

V : 变元的有穷集

T : 终结符的有穷集, $V \cap T = \emptyset$

P : 产生式的有穷集: 左部 \rightarrow 体/右部 $\in (V \cup T)^*$

S : 初始符号

- 简化表达

多个 A 的产生式 $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$ 简写为 $A \rightarrow \alpha_1 | \alpha_2 | \alpha_3 | \dots | \alpha_n$

变元 A 的全体产生式称为 A 产生式

- 归约和派生

归约: 字符串到文法变元的分析过程, 自底向上, 从产生式的体向头

派生: 文法变元到字符串的分析过程, 自顶向下, 从产生式的头向体

$\alpha A \beta \xRightarrow{G} \alpha \gamma \beta$: 当 $A \rightarrow \gamma \in P$ 时, 称为 $\alpha A \beta$ 可派生出 $\alpha \gamma \beta$, 也称 $\alpha \gamma \beta$ 可归约为 $\alpha A \beta$

$\alpha A \beta \Rightarrow \alpha \gamma \beta$: G 在已知情况下, 可以省略

$\alpha A \beta \xRightarrow{*}_G \alpha \gamma \beta$: 经过0步或多步派生

$\alpha A \beta \xRightarrow{i}_G \alpha \gamma \beta$: 恰好经过 i 步派生

$\Rightarrow_{lm}, \xRightarrow{*}_{lm}$: 最左派生, 即派生过程中仅替换符号串最左边的变元

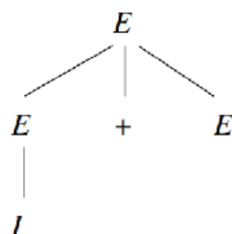
$\Rightarrow_{rm}, \xRightarrow{*}_{rm}$: 最右派生

$A \xRightarrow{*} w$ 当且仅当 $A \xRightarrow{*}_{lm} w$, 当且仅当 $A \xRightarrow{*}_{rm} w$

语法分析树

用来表示派生, 可以从树中看出整个派生过程和最终产生的符号串

例:



• 定义

1. 每个内节点是变元符号 $\in V$
2. 每个叶子节点 $\in V \cup T \cup \{\epsilon\}$
3. 若某个内节点标记是 A ，其子节点从左到右为

$$X_1, X_2, X_3, \dots, X_n$$

则 $A \rightarrow X_1, X_2, X_3, \dots, X_n$ 是 P 的一个产生式

• 相关定理和结论

语法分析树的全部**叶节点**从左到右连接起来，称为该树的**产物**或**结果**，如果根节点是**初始符号** S ，叶节点是**终结符**，那么该树的产物属于 $L(G)$

若某个 X_i 为 ϵ ，则 X_i 一定是 A 唯一的子节点，且 $A \rightarrow \epsilon$ 是一个产生式 (*)

• **等价定义**，对 $CFG\ G = (V, T, P, S)$, $A \in V$ ，以下命题等价：

- 文法 G 中存在以 A 为根节点，产物为 α 的语法分析树
- $A \xRightarrow{*} w$
- $A \xRightarrow{lm} w$
- $A \xRightarrow{rm} w$

文法和语言的歧义性

- 若 $CFG\ G$ 使某些符号串有两棵**不同**的语法分析树，则该文法是歧义的，有些文法的歧义性可以通过**重新设计文法**来消除
- 如果上下文无关语言 L 的所有文法都是有歧义的，则称 L 是固有歧义的
 - 例： $L = \{a^i b^j c^k \mid i = j \text{ or } j = k\}$ 是固有歧义的
- 「判定任何给定 $CFG\ G$ 是否是有歧义的」是一个不可判定问题

文法的化简

- **定理**：每个不带 ϵ 的 CFL 都可以由一个不带无用符号， ϵ -产生式和单元产生式的文法来定义
- 文法化简的可靠顺序
 - 消除 ϵ -产生式
 - 消除单元产生式
 - 消除非产生的无用符号
 - 消除非可达的无用符号

消除 ε -产生式

注意, 对于CFG L , 消除了 ε -产生式之后, 该文法对应的语言为 $L - \{\varepsilon\}$, 若空串本身不在该语言中, 则 $L = L - \{\varepsilon\}$

- 确定可空变元
 - 若 $A \rightarrow \varepsilon$, 则 A 是可空的
 - 若 $B \rightarrow \alpha$, 且 α 中每个符号都是可空的, 则 B 是可空的
- 替换产生式, 对含有可空变元的一条产生式 $A \rightarrow X_1, X_2, \dots, X_n$, 用一组产生式 $A \rightarrow Y_1, Y_2, \dots, Y_n$ 代替
 - 若 X_i 不是可空的, 则 Y_i 为 X_i
 - 若 X_i 是可空的, 则 Y_i 为 X_i 或 ε
 - Y_i 不全为 ε

注意最后要先删除所有的 ε -产生式, 如果题目明确要求了化简后还要保留语言中的空串, 则再加个产生式 $S \rightarrow \varepsilon$ 即可, 根据题意决定操作

消除单元产生式

- 确定单元对: 如果有 $A \xRightarrow{*} B$, 则称 $[A, B]$ 为单元对 (注意这里的 A 和 B 都是单个变元)
 - 若 $A \rightarrow B \in P$ 则 $[A, B]$ 是单元对
 - 若 $[A, B]$ 和 $[B, C]$ 都是单元对, 则 $[A, C]$ 是单元对
- 消除单元产生式
 - 删除全部形为 $A \rightarrow B$ 的单元产生式
 - 对每个单元对 $[A, B]$, 将 B 的产生式复制给 A

消除无用符号

若 $S \xRightarrow{*} \alpha X \beta$, 则 X 是可达的, 若 $\alpha X \beta \xRightarrow{*} w (w \in T^*)$, 则 X 是产生的

如果 X 既是产生的, 也是可达的, 则 X 是有用的, 注意这里的 $X \in (V \cup T)$

- 计算"产生的"符号集
 - T 中符号都是产生的
 - 对产生式 $A \rightarrow \alpha \in P$ 且 α 中符号都是产生的, 则 A 是产生的
- 计算"可达的"符号集
 - 符号 S 是可达的
 - 对产生式 $A \rightarrow \alpha \in P$ 且 A 是可达的, 则 α 中符号都是可达的
- 删除全部含有"非产生的"和"非可达的"符号的产生式

文法的范式

乔姆斯基范式CNF

- 定义: 每个不带 ε 的CFL都可以由这样的CFG G 定义, G 中每个产生式的形式都为:

$$A \rightarrow BC \text{ 或 } A \rightarrow a$$

其中 A, B, C 是变元, a 是终结符

- 转换方法

1. 设文法不带无用符号, ε -产生式和单元产生式 (任何文法都可以化简成这种形式)

- 考虑文法每个形式为 $A \rightarrow X_1, X_2, \dots, X_m (m \geq 2)$ 的产生式, 若 X_i 为终结符 a , 则引入新变元 C_a 替换 X_i 并增加新产生式 $C_a \rightarrow a$
- 所有产生式的形式变成 $A \rightarrow B_1, B_2, \dots, B_m (m \geq 2)$ 和 B_i (即 C_a) $\rightarrow a$
- 对于 $A \rightarrow B_1 B_2, \dots, B_m (m \geq 2)$, 引入新变元 D_1, D_2, \dots, D_{m-2} , 即替换为一组级联产生式

$$A \rightarrow B_1 D_1, D_1 \rightarrow B_2 D_2, \dots, D_{m-2} \rightarrow B_{m-1} B_m$$

格雷巴赫范式GNF

- 定义: 每个不带 ϵ 的 CFL 都可以由这样的 CFG G 定义, G 中每个产生式的形式都为

$$A \rightarrow a\alpha$$

其中 A 是变元, a 是终结符, α 是零或多个变元的串

关于 GNF , 要记住的是, 其每个产生式都会引入一个终结符, 长度为 n 的串的派生恰好是 n 步

下推自动机

形式定义

- 定义下推自动机 (PDA) P 为七元组

$$P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$$

- Q : 有穷状态集
- Σ : 有穷输入符号集
- Γ : 有穷栈符号集
- $\delta: Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \rightarrow 2^{Q \times \Gamma^*}$: 状态转移函数
- $q_0 \in Q$: 初始状态
- Z_0 : 初始栈底符号
- $F \subseteq Q$: 接受状态集

- PDA 的动作

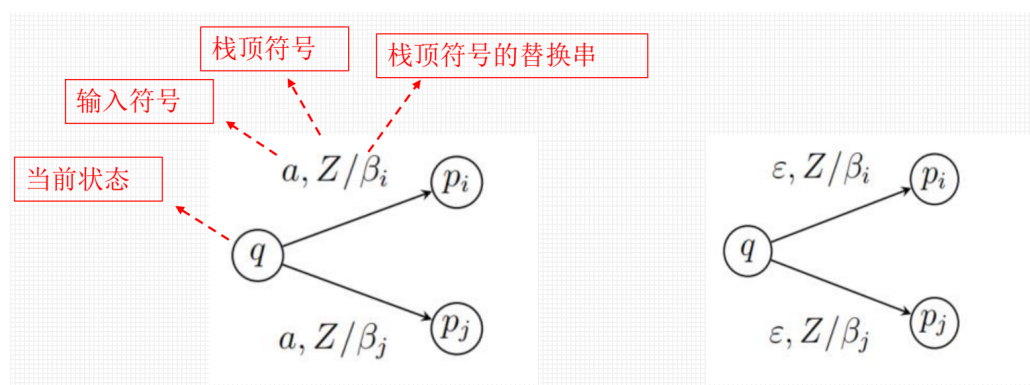
若 $q, p_i \in Q (1 \leq i \leq m), a \in \Sigma \cup \{\epsilon\}, Z \in \Gamma, \beta_i \in \Gamma^*$, 则 $\delta(q, a, Z)$ 表示如下的动作:

$$\delta(q, a, Z) = \{(p_1, \beta_1), (p_2, \beta_2), \dots, (p_m, \beta_m)\}$$

或

$$\delta(q, \epsilon, Z) = \{(p_1, \beta_1), (p_2, \beta_2), \dots, (p_m, \beta_m)\}$$

- 状态转移图



• 瞬时描述

$(q, aw, Z\alpha) \vdash_p (p, w, \beta\alpha)$: 在状态 q 从当前符号串 aw 中消耗 a (可为 ε), 转移到状态 p , 同时用 β (可为 ε)替换栈顶的 Z , 余下的符号串为 w

$(q, aw, Z\alpha) \vdash (p, w, \beta\alpha)$: 所指PDA明确时, 可以忽略角标 p

$(q, aw, Z\alpha) \vdash^* (p, w, \beta\alpha)$: 经过零个或多个动作进行转移

终态方式和空栈方式的转换

首先明确两种方式的别:

- **终态方式:** $L(P) = \{w | (q_0, w, Z_0 \vdash^* (p, \varepsilon, \gamma), p \in F)\}$

- 检查符号串是否消耗完
- 检查是否进入了接受状态

- **空栈方式:** $N(P) = \{w | (q_0, w, Z_0 \vdash^* (p, \varepsilon, \varepsilon))\}$

- 检查符号串是否消耗完
- 检查栈是否已清空

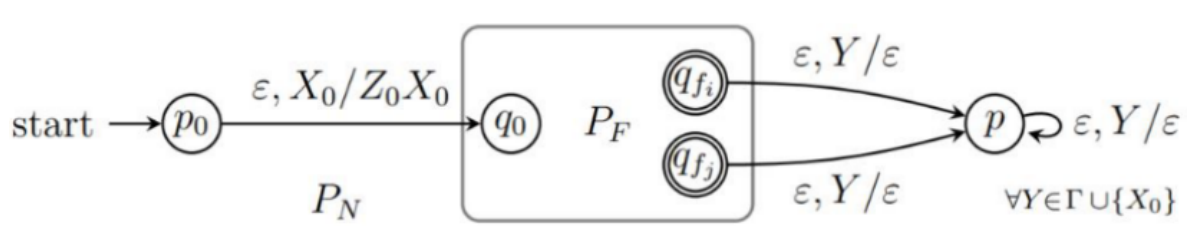
• 定理

如果PDA P_F 以终态方式接受语言 L , 那么一定存在PDA P_N 以空栈方式接受 L

如果PDA P_N 以空栈方式接受语言 L , 那么一定存在PDA P_F 以终态方式接受 L

证明略

从终态方式到空栈方式



- 设 $P_F = (Q, \Sigma, \Gamma, \delta_F, q_0, Z_0, F)$

构造 $P_N = (Q \cup \{p_0, p\}, \Sigma, \Gamma \cup \{X_0\}, \delta_N, p_0, X_0, \emptyset)$ 即添加了新的开始状态 q_0 、消除状态 q^* 和一个栈符号 X_0

- 先将 P_F 的栈底符号 Z_0 压栈

$$\delta_N(p_0, \varepsilon, X_0) = \{(q_0, Z_0 X_0)\}$$

- 模拟 P_F 的动作

$$\delta_N(q, a, Y) \text{ 包含 } \delta_F(q, a, Y) \text{ 的全部元素}$$

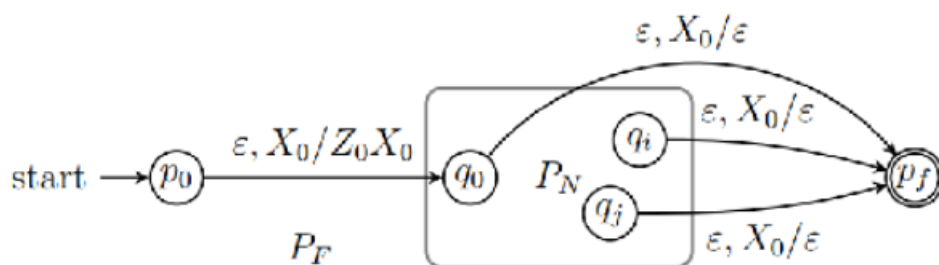
- 从 $q_f \in F$ 开始弹出栈中符号

$$\delta_N(q_f, \varepsilon, Y) \text{ 包含 } (p, \varepsilon)$$

- 状态 p 时, 弹出全部栈中符号, 即 $\forall Y \in \Gamma \cup \{X_0\}$

$$\delta_N(p, \varepsilon, Y) = \{(p, \varepsilon)\}$$

从空栈方式到终态方式



- 设 $P_N = (Q, \Sigma, \Gamma, \delta_N, q_0, Z_0, \emptyset)$
构造 $P_F = (Q \cup \{p_0, p_f\}, \Sigma, \Gamma \cup \{X_0\}, \delta_F, p_0, X_0, \{p_f\})$
- 开始时, 将 P_N 栈底符号 Z_0 压栈
 $\delta_F(p_0, \epsilon, X_0) = \{(q_0, Z_0 X_0)\}$
- 模拟 P_N 的动作
 $\delta_F(q, a, Y)$ 包含 $\delta_N(q, a, Y)$ 的全部元素
- $\forall q \in Q$, 看到 P_F 的栈底 X_0 , 则转移到新的终态 p_f
 $\delta_F(q, \epsilon, X_0) = \{(p_f, \epsilon)\}$

下推自动机和文法的等价性

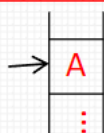
任意CFG到PDA

- 定理: 对任何CFL L , 一定存在PDA P , 使得 $L = N(P)$
- 构造方法
给定CFG $G = (V, T, P, S)$
构造PDA $P_N = (\{q\}, T, V \cup T, \delta, q, S, \emptyset)$

其中 δ :

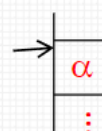
1. 对每个变元 A , 输入空串, 压入产生式: $\delta(q, \epsilon, A) = \{(q, \alpha) | A \rightarrow \alpha \in P\}$
2. 对每个终结符 a , 输入字符, 弹出字符 a : $\delta(q, a, a) = \{(q, \epsilon)\}$

变化前:

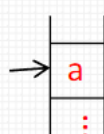


$$\delta(q, \epsilon, A) = \{(q, \alpha) | A \rightarrow \alpha \in P\}$$

变化后:

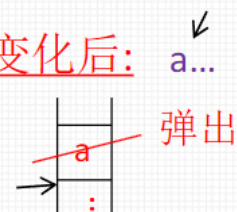


变化前:



$$\delta(q, a, a) = \{(q, \epsilon)\}$$

变化后:



这种方法构造出来的PDA在实际识别的时候会产生大量的dead paths, 识别时只针对正确的路径分析即可

GNF到PDA

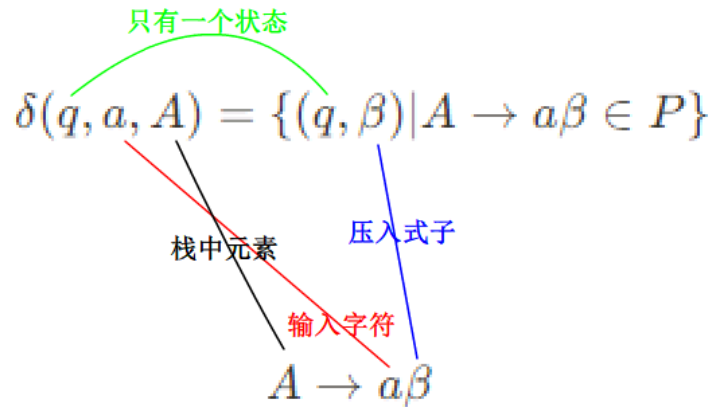
- 构造方法

如果有GNF格式的CFG $G = (V, T, P, S)$, 那么构造PDA $P_N = (\{q\}, T, V, \delta, q, S, \emptyset)$

对每个产生式 $A \rightarrow a\beta$, 定义 δ 为:

$$\delta(q, a, A) = \{(q, \beta) | A \rightarrow a\beta \in P\}$$

他们之间的关系是这样的：



PDA到CFG

如果PDA $P_N = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, \emptyset)$, 有 $L = N(P)$ (即 P_N 以空栈方式接受 L) 那么可以构造CFG $G = (V, \Sigma, P, S)$, 其中:

变元 $V = \{[qXp] \mid p, q \in Q, X \in \Gamma\} \cup \{S\}$

产生式集合 P 包括:

1. $\forall p \in Q$, 构造产生式 $S \rightarrow [q_0 Z_0 p]$
2. $\forall (p, Y_1, Y_2, \dots, Y_n) \in \delta(q, a, X)$, 构造 $|Q|^n$ 个产生式

$$[qXr_n] \rightarrow a[pY_1r_1][r_1Y_2r_2]\dots[r_{n-1}Y_nr_n]$$

其中 $a \in \Sigma \cup \{\varepsilon\}$; $X, Y_i \in \Gamma$; r_i 是 Q 中任意一个状态

3. $\forall (p, Y_1 Y_2 \dots Y_n) \in \delta(q, a, X)$, 若 $Y_1 Y_2 \dots Y_n = \varepsilon$, 则有产生式 $[qXp] \rightarrow a$ (这代表着从状态 q 完全弹出栈符号 X 到状态 p , 只需要接受一个符号 a)
- 通常我们在完成了上述步骤之后, 还需要**消除无用符号**, 这是因为, 在上面的过程中我们多次地**假定所有状态之间相互可达**, 假如某个状态 q_i 是不能到达 q_j 的, 但是我们在构造产生式时还是会用到状态 $[q_i X q_j]$, 因此是需要消除无用符号的
 - PPT上的例题是没有消除 ε -产生式的
 - 此外, 为了化简方便, 可以进一步使用 A, B, \dots 替换这些看起来复杂的变元

细节方面

在构造的过程中, 大量地使用了「穷举」的思想, 需要后续来消除无用的符号

确定型下推自动机

- 定义, 如果PDA $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$ 满足:
 - 对于任何的 (q, a, X) 组合, $q \in Q, a \in \Sigma \cup \{\varepsilon\}, X \in \Gamma, \delta(q, a, X)$ 至多有一个转移
 - 如果 $\delta(q, a, X)$ 对于某个 $a \in \Sigma$ 是有定义的, 则 $(q, \varepsilon, X) = \emptyset$
- 上面的定义可以从两个方面考虑:
 - 在确定某个状态 q 、某个栈顶符号 X 和输入字符 a 下, 状态转移的结果数量最多一个
 - 在确定某个状态 q 、某个栈顶符号 X 下, 不可能既能接受终结符 a 进行状态转移, 又可以接受空字符 ε 进行状态转移
- 表示语言的能力上 $DFA \Leftrightarrow NFA$, 但是 $DPDA$ 和 PDA 不等价

- 终态方式

$DPDA$ 以终态方式接受的语言也被称为 $DCFL$ (确定的上下文无关语言)

$RL \subset DCFL$ (正则语言真包含于 $DCFL$)

$DCFL$ 应用于:

1. 非固有歧义语言的真子集
2. 著名语法分析器YACC的基础 $LR(k)$ 文法是 $DCFL$ 的子集

- 空栈方式

定义: 语言 L 中不存在字符串 x 和 y , 使 x 是 y 的前缀, 则 L 具有前缀性质

定理: 如果 $L = N(P)$, 即被某个 $DPDA P$ 以空栈方式接受, 当且仅当 L 有前缀性质, 且存在以终态方式接受的 $DPDA P'$ 使 $L = L(P')$

$DPDA P$ 的 $N(P)$ 更有限, 不能接受 0^*

- $DPDA$ 与歧义文法

定理: $DPDA P$, 语言 $L = L(P)$, 那么 L 有无歧义的 CFG

定理: $DPDA P$, 语言 $L = N(P)$, 那么 L 有无歧义的 CFG