

编译原理 实验二：语法分析

09019106 牟倪

目录

一、实验目标.....	2
二、实验内容.....	2
三、基本思路.....	2
四、假设与名词定义.....	2
五、构造 Parsing Table.....	3
5.1 消除左递归与左公共因子.....	3
5.2 计算 First 与 Follow 集合	4
5.3 构造 Parsing Table.....	5
六、核心数据结构.....	6
6.1 用栈实现自顶而下分析.....	6
6.2 用树记录 parsing tree 信息.....	6
七、核心算法.....	7
7.1 用 map 记录产生式.....	7
7.2 硬编码构造 parsing table	7
7.3 解析 token 串并构造 parsing tree	8
7.4 parsing tree 的前序和后序遍历.....	10
八、实验.....	10
九、出现问题与解决方案.....	11
十、实验体会.....	11
附录：源代码.....	12

一、实验目标

- 理解语法分析过程。
- 对 LL parsing, 熟练掌握消除左递归→消除左公共因子→得到 First 和 Follow 集合→构造 parsing table 的过程。
- 对 LR parsing, 熟练掌握构造 NFA→化为 DFA→构造 parsing table 的过程。
- 自主实现 LL 或 LR 语法分析程序。

二、实验内容

- 自定义 token 和 context-free grammar。
- 完成词法分析程序：
 - 输入 token 串, 输出解析结果的 parsing tree 的前序和后序遍历。
 - 对不合法的输入 token 串, 报错。

三、基本思路

我选择实现 LL 分析程序。首先对原 grammar 消除左递归和左公共因子, 然后得到 First 和 Follow 集合, 最后构造 parsing table。将 token 序列成功解析后, 我们用树来记录 parsing tree 的信息, 并输出 parsing tree 的遍历序列。

四、假设与名词定义

终结符定义:

- id: 表示变量的 identifier;
- =: 赋值的等于号。
- num: 表示数字。
- if: 表示“如果”, 应用格式为 if then else。
- then: 表示“然后”, 应用格式为 if then else。
- else: 表示“其他”, 应用格式为 if then else。
- cmp: 表示比较的运算符, 如 <、>、≤、≥ 等。
- op: 运算符号, 如 +、-、*、/ 等。
- rop: 逻辑运算符号, 如 &&、|| 等。

非终结符定义:

- S: 文法开始符号, 表示 statement。

- E: 表示 expression, 返回值为数值。
- B: 表示 boolean expression, 返回 true/false。

context-free grammar 定义:

- $S \rightarrow \text{if } B \text{ then } S \text{ else } S$
- $S \rightarrow \text{id} = E$
- $B \rightarrow E \text{ cmp } E$
- $B \rightarrow E \text{ cmp } E \text{ rop } B$
- $E \rightarrow \text{num}$
- $E \rightarrow \text{num op } E$

五、构造 Parsing Table

5.1 消除左递归与左公共因子

grammar 中没有左递归, 但 B 和 E 的产生式有左公共因子。

把所有的 E 替换为 num Enext,

- $\text{Enext} \rightarrow \text{op num Enext}$
- $\text{Enext} \rightarrow \varepsilon$

消除 E 产生式中的左公共因子后, 文法如下所示:

- $S \rightarrow \text{if } B \text{ then } S \text{ else } S$
- $S \rightarrow \text{id} = \text{num Enext}$
- $B \rightarrow \text{num Enext cmp num Enext}$
- $B \rightarrow \text{num Enext cmp num Enext rop } B$
- $\text{Enext} \rightarrow \text{op num Enext}$
- $\text{Enext} \rightarrow \varepsilon$

把所有的 B 替换为 num Enext cmp num Enext Bnext,

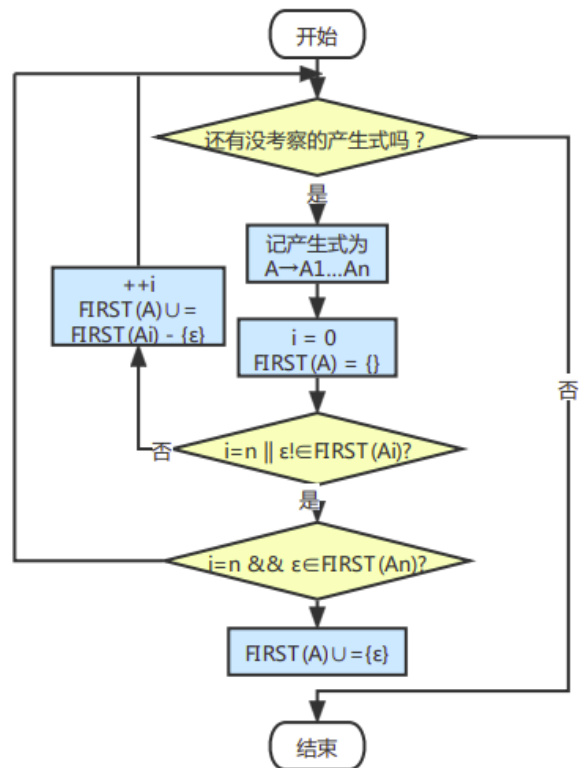
- $\text{Bnext} \rightarrow \text{rop num Enext cmp num Enext Bnext}$
- $\text{Bnext} \rightarrow \varepsilon$

消除 B 产生式中的左公共因子后, 文法如下所示:

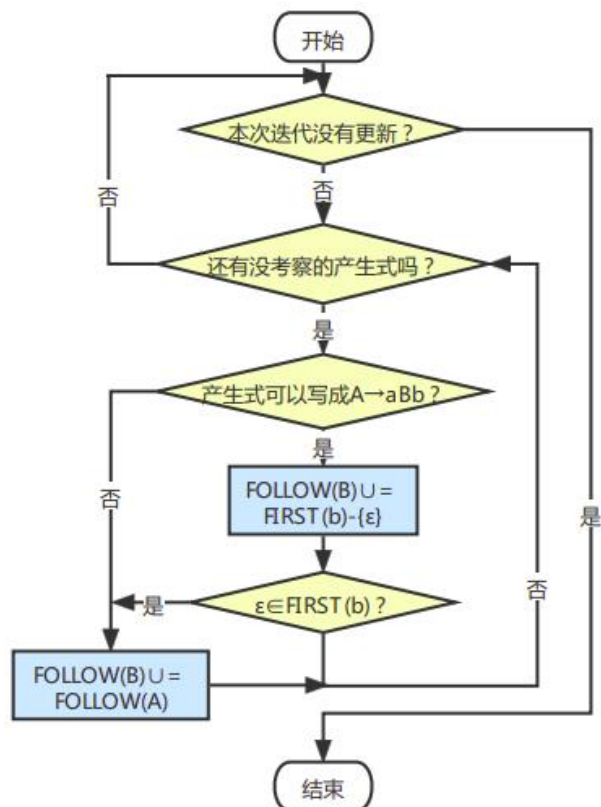
- $S \rightarrow \text{if num Enext cmp num Enext Bnext then } S \text{ else } S$
- $S \rightarrow \text{id} = \text{num Enext}$
- $\text{Bnext} \rightarrow \text{rop num Enext cmp num Enext Bnext}$
- $\text{Bnext} \rightarrow \varepsilon$
- $\text{Enext} \rightarrow \text{op num Enext}$
- $\text{Enext} \rightarrow \varepsilon$

5.2 计算 First 与 Follow 集合

构造 First 集的算法如下所示：



构造 Follow 集的算法如下所示：

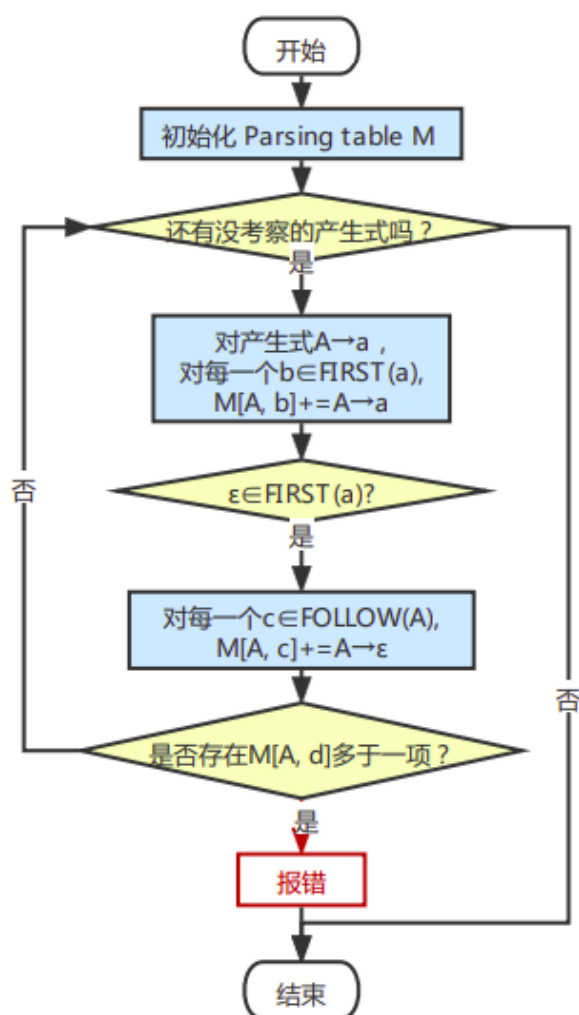


目前，非终结符有：S, Bnext, Enext；终结符有：if, then, else, id, =, num, cmp, rop, op。我们构造 First 集和 Follow 集如下：

		FISRT	FOLLOW
S	$S \rightarrow \text{if num Enext cmp num Enext Bnext then S else S}$	{if}	{else, #}
	$S \rightarrow \text{id = num Enext}$	{id}	
Bnext	$\text{Bnext} \rightarrow \text{rop num Enext cmp num Enext Bnext}$	{rop}	{then}
	$\text{Bnext} \rightarrow \epsilon$	{ ϵ }	
Enext	$\text{Enext} \rightarrow \text{op num Enext}$	{op}	{cmp, then, rop, else, #}
	$\text{Enext} \rightarrow \epsilon$	{ ϵ }	

5.3 构造 Parsing Table

构造 parsing table 的算法如下所示：



构造 parsing table 如下：

	if	then	else	id	=	num	cmp	rop	op	#
S	S → if num Enext cmp num Enext Bnext then S else S			S → id = num Enext						
Bnext		Bnext → ε						Bnext → rop num Enext cmp num Enext Bnext		
Enext		Enext → ε	Enext → ε				Enext → ε	Enext → ε	Enext → op num Enext	Enext → ε

六、核心数据结构

6.1 用栈实现自顶而下分析

在最初，栈底为语法开始符号 **S**。当读入 token 后，如果栈顶符号为终结符（即 token），比较两个 token 是否相等，若相等则弹栈、继续读入 token，若不相等则报错。如果栈顶符号为非终结符，则根据 parsing table 进行推导，即将产生式左部从栈中弹出，然后将产生式右部压入栈中，读入 token 的进度不变化。

6.2 用树记录 parsing tree 信息

我们构造节点类 **Node**，包含：

- **string value**，代表节点存储的值。如果是叶节点，**value** 为终结符；如果是非叶节点，**value** 为非终结符。
- **vector<Node*> children**，代表指向子节点的指针。

每完成一次推导，我们就把弹出栈的产生式左部（即非终结符）拿到，并相应去构造产生式右部的节点，令产生式左部的节点的 **children** 指向产生式右部的节点，并将这些产生式右部节点压栈。

七、核心算法

7.1 用 map 记录产生式

我们将 token 以 string 的形式记录，用 `map<int, vector<string>>` 的形式记录产生式，其中 map 的 key 为整型，是约定好的产生式编号。代码如下：

```
1. map<int, vector<string>> productions;
2.
3. void init(){
4.     // productions
5.     string temp1[11] = {"if", "num", "Enext", "cmp", "num", "Enext", "Bnext", "then", "S", "else", "S"};
6.     productions[1] = vector<string>(temp1, temp1 + 11);
7.
8.     string temp2[4] = {"id", "=", "num", "Enext"};
9.     productions[2] = vector<string>(temp2, temp2 + 4);
10.
11.    string temp3[7] = {"rop", "num", "Enext", "cmp", "num", "Enext", "Bnext"};
12.    productions[3] = vector<string>(temp3, temp3 + 7);
13.
14.    productions[4] = vector<string>();
15.
16.    string temp5[3] = {"op", "num", "Enext"};
17.    productions[5] = vector<string>(temp5, temp5 + 3);
18.
19.    productions[6] = vector<string>();
20.}
```

7.2 硬编码构造 parsing table

我们用函数 `int parsing_table(string stacktop, string nextval)` 来完成查 parsing table 的动作，参数 1 为栈顶 token，参数 2 为下一个 token，返回要使用的产生式编号。代码如下：

```
1. int parsing_table(string stacktop, string nextval){
2.     if(stacktop == "S"){
3.         if(nextval == "if") return 1;
4.         if(nextval == "id") return 2;
5.         fail();
6.     }
7. }
```

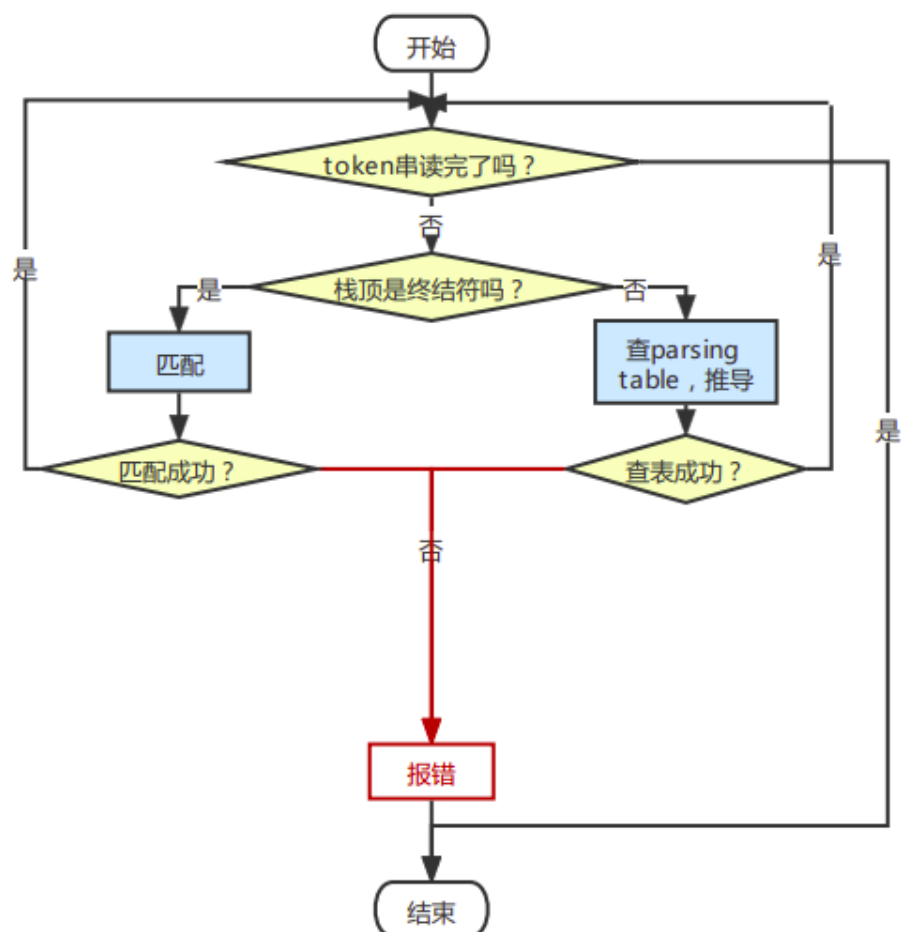
```

6.     }
7.     else if(stacktop == "Bnext"){
8.         if(nextval == "rop") return 3;
9.         if(nextval == "then") return 4;
10.        fail();
11.    }
12.    else if(stacktop == "Enext"){
13.        if(nextval == "op") return 5;
14.        if(nextval == "then" || nextval == "else" ||
15.           nextval == "cmp" || nextval == "rop" || nextval == "#
    ") return 6;
16.        fail();
17.    }
18.    return -1;
19.}

```

7.3 解析 token 串并构造 parsing tree

我们用函数 `void parse(vector<string> tokens)` 来完成 token 串的解析。该函数会修改作为全局变量的 parsing tree 数据结构。流程图如下：



代码如下：

```
1. void parse(vector<string> tokens){
2.     root = new Node("S");
3.     st.push(root);
4.     int cur = 0;
5.     while(!st.empty() && tokens[cur] != "#"){
6.         string stacktop = st.top()->value;
7.         if(!is_non_terminal(stacktop)){ // terminal
8.             if(stacktop == tokens[cur]){
9.                 st.pop();++cur;
10.            }
11.            else fail();
12.        }
13.        else{ // non terminal
14.            string nextval = tokens[cur];
15.            int idx = parsing_table(stacktop, nextval);
16.            Node *temp = st.top();st.pop();
17.            int len = productions[idx].size();
18.            for(int i=0; i<len; ++i){
19.                Node* newNode = new Node(productions[idx][i]);
20.                temp->children.push_back(newNode);
21.            }
22.            for(int i=len-1; i>=0; --i){
23.                st.push(temp->children[i]);
24.            }
25.        }
26.    }
27.}
```

7.4 parsing tree 的前序和后序遍历

我们用递归的方式实现 parsing tree 的前序和后序遍历。代码如下：

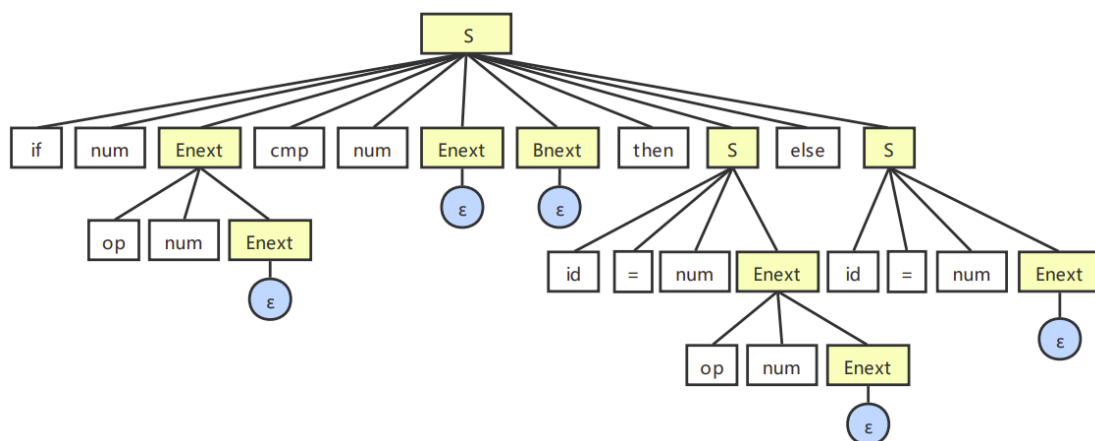
```
1. void preorder_traversal(Node* node){
2.     cout<<node->value<<" ";
3.     for(Node* i:node->children) preorder_traversal(i);
4. }
5.
6. void postorder_traversal(Node* node){
7.     for(Node* i:node->children) postorder_traversal(i);
8.     cout<<node->value<<" ";
9. }
```

八、实验

测试用例为 "if num op num cmp num then id = num op num else id = num #". 这样看可能比较抽象，我们将其实例化一下：

```
1. if a + b > c then x = a + b else x = c
```

这就是一个以上 token 序列的实例化。parsing tree 如下所示，其中黄色节点为非终结符，白色节点为终结符。



其前序遍历序列为：S if num Enext op num Enext cmp num Enext Bnext then S id = num Enext op num Enext else S id = num Enext

后序遍历序列为：if num op num Enext Enext cmp num Enext Bnext then id = num op num Enext Enext S else id = num Enext S S

程序执行结果：

```
preorder:
S if num Enext op num Enext cmp num Enext Bnext then S id = num Enext op num Enext else S id = num Enext
postorder:
if num op num Enext Enext cmp num Enext Bnext then id = num op num Enext Enext S else id = num Enext S S
```

前序遍历与后序遍历都输出正确。

我们再来测试一个失败的用例 "if then else #", 程序报错：

```
parsing failed!
```

可以看出，实验结果正确，程序编写合理。

九、出现问题与解决方案

没有出现问题。

十、实验体会

经过本次实验，我复习了自顶而下 LL 语法分析的内容，消除左递归、左公共因子、计算 First 和 Follow 集合、构造 parsing table 也更熟练了。

附录：源代码

```
1. #include<iostream>
2. #include<stack>
3. #include<vector>
4. #include<map>
5. #include<string>
6. #include<cstring>
7.
8. using namespace std;
9.
10. class Node{
11. public:
12.     string value;
13.     vector<Node*> children;
14.     Node(string v):value(v){}
15.     Node(){}
16. };
17.
18. map<int, vector<string>> productions;
19. Node* root = NULL;
20. stack<Node*> st;
21.
22. void fail(){
23.     cout<<"parsing failed!\n";
24.     exit(1);
25. }
26.
27. void init(){
28.     // productions
29.     string temp1[11] = {"if", "num", "Enext", "cmp", "num", "Enext", "Bnext", "then", "S", "else", "S"};
30.     productions[1] = vector<string>(temp1, temp1 + 11);
31.
32.     string temp2[4] = {"id", "=", "num", "Enext"};
33.     productions[2] = vector<string>(temp2, temp2 + 4);
34. }
```

```

35.     string temp3[7] = {"rop", "num", "Enext", "cmp", "num", "Enex
    t", "Bnext"};
36.     productions[3] = vector<string>(temp3, temp3 + 7);
37.
38.     productions[4] = vector<string>();
39.
40.     string temp5[3] = {"op", "num", "Enext"};
41.     productions[5] = vector<string>(temp5, temp5 + 3);
42.
43.     productions[6] = vector<string>();
44. }
45.
46. bool is_non_terminal(string token){
47.     return (token == "S" || token == "Bnext" || token == "Enext")
    ;
48. }
49.
50. int parsing_table(string stacktop, string nextval){
51.     if(stacktop == "S"){
52.         if(nextval == "if") return 1;
53.         if(nextval == "id") return 2;
54.         fail();
55.     }
56.     else if(stacktop == "Bnext"){
57.         if(nextval == "rop") return 3;
58.         if(nextval == "then") return 4;
59.         fail();
60.     }
61.     else if(stacktop == "Enext"){
62.         if(nextval == "op") return 5;
63.         if(nextval == "then" || nextval == "else" ||
64.             nextval == "cmp" || nextval == "rop" || nextval == "#
            ") return 6;
65.         fail();
66.     }
67.     return -1;
68. }
69.
70. void parse(vector<string> tokens){
71.     root = new Node("S");
72.     st.push(root);
73.     int cur = 0;
74.     while(!st.empty() && tokens[cur] != "#"){
75.         string stacktop = st.top()->value;

```

```

76.         if(!is_non_terminal(stacktop)){ // terminal
77.             if(stacktop == tokens[cur]){
78.                 st.pop();++cur;
79.             }
80.             else fail();
81.         }
82.         else{ // non terminal
83.             string nextval = tokens[cur];
84.             int idx = parsing_table(stacktop, nextval);
85.             Node *temp = st.top();st.pop();
86.             int len = productions[idx].size();
87.             for(int i=0; i<len; ++i){
88.                 Node* newNode = new Node(productions[idx][i]);
89.                 temp->children.push_back(newNode);
90.             }
91.             for(int i=len-1; i>=0; --i){
92.                 st.push(temp->children[i]);
93.             }
94.         }
95.     }
96.
97. }
98.
99. void preorder_traversal(Node* node){
100.     cout<<node->value<<" ";
101.     for(Node* i:node->children) preorder_traversal(i);
102. }
103.
104. void postorder_traversal(Node* node){
105.     for(Node* i:node->children) postorder_traversal(i);
106.     cout<<node->value<<" ";
107. }
108.
109. int main(){
110.     init();
111.     string temptoken[17] = {"if", "num", "op", "num", "cmp", "nu
        m",
112.         "then", "id", "=", "num", "op", "num",
113.         "else", "id", "=", "num", "#"};
114. };
115.     vector<string> token_to_parse = vector<string>(temptoken, te
        mptoken + 17);
116.     //string badtoken[4] = {"if", "then", "else", "#"}; // 不合法
        的 token 序列, 报错

```

```
117.    //vector<string> token_to_parse = vector<string>(badtoken, b
    adtoken + 4);
118.
119.    parse(token_to_parse);
120.    cout<<"preorder:\n";
121.    preorder_traversal(root); cout<<"\n\npostorder:\n";
122.    postorder_traversal(root); cout<<"\n\n";
123. }
```