

天津大学



编译原理 词法分析器/语法分析器 开发报告

学 院 智能与计算/求是学部
专 业 计算机科学与技术
组 长 杨亦凡 3019234258
组 员 李自安 3019207257
组 员 石昊 3019208051
组 员 华溢 3019244091

目录

1	项目简介	4
1.1	开发环境	4
1.2	实现语言	4
1.3	项目结构	4
2	词法分析器	5
2.1	需求分析	5
2.2	数据结构描述	6
2.2.1	符号表项	6
2.2.2	符号表	7
2.2.3	状态	7
2.2.4	结点	7
2.2.5	边	8
2.2.6	图	8
2.3	算法描述	9
2.3.1	正则表达式转 NFA 算法	9
2.3.2	NFA 确定化算法	9
2.3.3	DFA 最小化算法	10
2.3.4	去除注释算法	10
2.3.5	词法分析器解析	13
2.4	输出格式说明	15
2.4.1	正则表达式转 NFA 算法	15
2.4.2	NFA 确定化算法	15
2.4.3	DFA 最小化算法	16
2.4.4	Token 序列	16
2.5	源程序编译步骤	16
3	语法分析器	17
3.1	数据结构描述	17
3.1.1	grammar 相关数据结构	17
3.1.2	FIRST & FOLLOW 集合	18
3.1.3	项目集规范族相关数据类型	19
3.1.4	LR1 分析表相关数据类型	21
3.1.5	LR1 语法分析器	22
3.2	算法描述	23
3.2.1	计算语法四元组	23
3.2.2	求 FIRST 集	24

3.2.3	求 FOLLOW 集	24
3.2.4	求闭包	25
3.2.5	求项目集规范族与 GO	26
3.2.6	求 ACTION 表与 goto 表	26
3.2.7	LR1 parse 解析	27
3.3	LR1 分析表描述	28
3.4	输出格式说明	29
3.4.1	FIRST 集与 FOLLOW 集输出	29
3.4.2	项目集规范族输出	30
3.4.3	LR1 分析表输出	31
3.4.4	规约序列输出	32
3.5	源程序编译步骤	32
4	主程序分析流程	33

1 项目简介

本项目针对 SQL-- 语言开发了一个编译器前端，包括词法分析器和语法分析器：

- (1) 使用自动机理论编写词法分析器
- (2) 使用自下而上的语法分析方法编写语法分析器

1.1 开发环境

macOS Monterey 12.3.1

1.2 实现语言

Python 3.9

1.3 项目结构



./src/lexer/lexical_analyzer.py 是词法分析器主程序
 ./src/LR1_parser/LR1.py 是语法分析器主程序
 ./tests 是测试用例文件夹
 ./output/token_table 是 token 序列输出
 ./output/parse_result 是规约序列输出

2 词法分析器

词法分析器首先用 SQL - - 语言所有单词符号对应的正则表达式通过正则表达式转 NFA 算法构建相应的 NFA，再通过 NFA 确定化算法构建相应的 DFA，然后通过 DFA 最小化算法构建有限自动机，最后根据有限自动机编程实现 SQL - - 语言的词法分析器。

2.1 需求分析

词法分析器的输入为 SQL - - 语言源代码，词法分析器对其预处理去除注释，识别单词的二元属性并生成符号表，将待分析代码转化为语法分析器可接受的序列。单词符号的类型包括关键字，标识符，界符，运算符，整数，浮点数，字符串。每种单词符号的具体要求如下：

TABLE 1: 关键字 (KW, 规定为大写)

类别	语法关键字
查询表达式	(1) SELECT, (2) FROM, (3) WHERE, (4) AS, (5) *
插入表达式	(6) INSERT, (7) INTO, (8) VALUES, (9) VALUE, (10) DEFAULT
更新表达式	(10) UPDATE, (11) SET
删除表达式	(13) DELETE
连接操作	(14) JOIN, (15) LEFT, (16) RIGHT, (17) ON
聚合操作	(18) MIN, (19) MAX, (20) AVG, (21) SUM
集合操作	(22) UNION, (23) ALL
组操作	(24) GROUP BY, (25) HAVING, (26) DISTINCT, (27) ORFER BY
条件语句	(28) TRUE, (29) FALSE, (30) UNKNOWN, (31) IS, (32) NULL

TABLE 2: 运算符 (OP, 规定为大写)

类别	语法关键字
比较运算符	(1) =, (2) >, (3) <, (4) >=, (5) <=, (6) !=, (7) <=>
逻辑运算符	(8) AND, (9) &&, (10) OR, (11) , (12) XOR, (13) NOT, (14) !
算术运算符	(13) -
属性运算符	(13) .

TABLE 3: 界符 (SE)

类别	语法关键字
界符	(1) (, (2)), (3) ,

标识符 (IDN) 为字母、数字和下划线 () 组成的不以数字开头的串;

整数 (INT)、浮点数 (FLOAT) 的定义与 C 语言正数相同, 负数通过在正数前面加算术运算符 (-) 实现;

字符串 (STRING) 定义与 C 语言相同, 使用双引号包含的任意字符串。

2.2 数据结构描述

2.2.1 符号表项

ADT 1: TokenLine 抽象数据类型

ADT TokenLine {

数据对象: D = {

待测代码中的单词符号 word_content,

单词符号大类 word_type,

单词符号种别 token_type,

单词符号内容 token_content

}

数据关系: R = {

<word_content, word_type>,

<word_type, token_type>,

<token_type, token_content>

}

基本操作:

input_raw_line(word_content, word_type)

初始条件: word_content 存在, word_type 合法

操作结果: 产生 token_type, 产生 token_content

output_token_line()

初始条件: word_content 存在, word_type 合法

操作结果: 产生 token_type, 产生 token_content

_process_raw_line()

初始条件: word_content 存在, word_type 合法

操作结果: 产生 token_type, 产生 token_content

_get_token_id(token_type)

初始条件: token_type 存在

操作结果：返回 word_content 在 token_type 表中对应的序号

```
} ADT TokenLine
```

2.2.2 符号表

ADT 2: TokenTable 抽象数据类型

```
ADT TokenTable {
    数据对象：D = {TokenLine i | TokenLine i 是符号表项, i = 0,1,...,n, n >= 0}
    数据关系：R = {<Ti, Ti+1> | Ti, Ti+1 属于D, i = 0,1,...,n-1}
    基本操作：
        add_token_line(toke_line)
            初始条件：toke_line 存在
            操作结果：将 toke_line 加入到 TokenTable 末尾
        print()
            初始条件：TokenTable 存在
            操作结果：格式化打印 TokenTable
        save(test_name)
            初始条件：TokenTable 存在
            操作结果：格式化存储 TokenTable, 文件名为 test_name
        reset()
            操作结果：重置 TokenLine
    } ADT TokenTable
```

2.2.3 状态

ADT 3: State 抽象数据类型

```
ADT State {
    数据对象：D = {
        状态编号 id,
        结点集合 node_set,
        起始状态标识位 is_begin,
        结束状态标识位 is_end
    }
    } ADT State
```

2.2.4 结点

ADT 4: Node 抽象数据类型

```

ADT Node {
    数据对象: 结点编号 num
} ADT Node

```

2.2.5 边

ADT 5: Edge 抽象数据类型

```

ADT Edge {
    数据对象: D = {
        弧尾结点编号 node_tail,
        弧头结点编号 node_head,
        边的类型 node_type
    }
    数据关系: R = {<Ni, Nj> | Ni, Nj 属于 Node, i, j = 0,1,...,n}
} ADT Edge

```

2.2.6 图

ADT 6: Graph 抽象数据类型

```

ADT Graph {
    数据对象: D = {
        源点编号 node_tail,
        汇点编号 node_head,
        源点集合 node_heads = {Ni | Ni 属于 Node, i = 0,1,...,n}
        边集合 edges = {Ei | Ei 属于 Edge, i = 0,1,...,n}
    }
    基本操作:
        add_star(obj)
            初始条件: obj 存在
            操作结果: 构建一个带有 obj 类型的自环的结点,
                    头尾各通过空弧接一个结点, 并将整体加入到图中。
        add_union(obj1, obj2)
            初始条件: obj1, obj2 存在
            操作结果: 构建一对通过两条反向弧连接的结点,
                    弧的类型分别为 obj1, obj2,
                    并将整体加入到图中。
        add_concat(obj1, obj2)
            初始条件: obj1, obj2 存在

```


操作结果：构建一对通过两条同向弧连接的结点，
弧的类型分别为 `obj1`, `obj2`,
并将整体加入到图中。

} ADT Graph

2.3 算法描述

2.3.1 正则表达式转 NFA 算法

Algorithm 1 REGEX TO NFA Algorithm

Inputs:

正则表达式 $regex$

Outputs:

NFA $M = \{S, \Sigma, \delta, S_0, S_t\}$

- 1: 对输入的正则表达式预处理，将正则表达式中隐含的 `&` 补充
 - 2: 初始化操作符栈 $OPTR$ 和操作数栈 $OPND$
 - 3: **while** $OPTR[-1] == \text{'\#'} \text{ and } regex[i] == \text{'\#'} \text{ do}$
 - 4: 扫描处理后的正则式，读入一个字符 ch
 - 5: **while** ch 为操作数 **do**
 - 6: 将 ch 压入 $OPND$ ，读入一个字符 ch
 - 7: **end while**
 - 8: 比较当前输入操作符和操作符栈顶的操作符的优先级
 - 9: **if** 当前输入操作符优先级更高 **then**
 - 10: 将 ch 压入 $OPTR$ ，读入一个字符 ch
 - 11: **end if**
 - 12: **if** 两者优先级相同 **then**
 - 13: 弹出 $OPTR$ 栈顶元素，读入一个字符 ch
 - 14: **end if**
 - 15: **if** 当前输入操作符优先级更低 **then**
 - 16: 弹出 $OPTR$ 栈顶元素，根据其具体类别弹出 $OPND$ 1 个或 2 个操作数，生成相应的状态图并压入 $OPND$
 - 17: **end if**
 - 18: **end while**
-

2.3.2 NFA 确定化算法

Algorithm 2 NFA TO DFA Algorithm

Inputs:

NFA $M = \{S, \Sigma, \delta, S_0, S_t\}$

Outputs:

DFA $M' = \{S, \Sigma, \delta, S_0, S_t\}$

- 1: 从 M 的 S_0 出发, 将仅经过任意条 ε 弧能到达的状态组成的集合 I 作为 M' 的初态 q_0
 - 2: **while** 不再产生新的状态 **do**
 - 3: 从 I 中元素出发, 将经过任意 $a \in \Sigma$ 的 a 弧转换 I_a 所组成的集合作为 M' 的状态
 - 4: **end while**
-

2.3.3 DFA 最小化算法

Algorithm 3 DFA Minimization Algorithm

Inputs:

DFA $M = \{S, \Sigma, \delta, S_0, S_t\}$

Outputs:

DFA $M' = \{S, \Sigma, \delta, S_0, S_t\}$

- 1: 构造状态集的初始划分 Π , 包含终态 S_t 和非终态 $S - S_t$ 两组
 - 2: **while not done do**
 - 3: 对 P_i 使用传播性原则构造新划分 Π_{new}
 - 4: **if** $\Pi_{new} == \Pi$ **then**
 - 5: $\Pi_{final} == \Pi$ **done**
 - 6: **else**
 - 7: $\Pi == \Pi_{new}$
 - 8: **end if**
 - 9: **end while**
 - 10: **for all** Π_{final} **do**
 - 11: 选一代表元素 s , 加入到 M' 的 S
 - 12: **if** 代表 s 满足 $\delta(s, a) == t$ **then**
 - 13: 令 r 作为 t 组的代表, 将转换 $\delta(s, a) = r$ 加入到 M' 的 δ
 - 14: **end if**
 - 15: 将含有 S_0 的组的代表作为 M' 的开始状态
 - 16: 将含有 S_t 的组的代表作为 M' 的终态
 - 17: **end for**
 - 18: 去除 M' 的 S 中死状态
-

2.3.4 去除注释算法

Algorithm 4 Input Preprocessing Algorithm

Inputs:

带注释的 SQL -- 语言源代码 $Text_{raw}$

Outputs:

去除注释的 SQL -- 语言源代码 $Text_{preprocessed}$

- 1: 当前字符串 $current_string \leftarrow ""$

```

2: 当前状态  $state \leftarrow$  代码
3: for all  $ch \in Text_{raw}$  do
4:   if  $state ==$  代码 then
5:     if  $ch == \text{'\text{'}}$  then
6:        $state \leftarrow$  斜杠
7:     else if  $ch == \text{'\"}}$  then
8:        $state \leftarrow$  短横线
9:     else
10:       $current\_string += ch$ 
11:      if  $ch == \text{'\textbackslash'}}$  then
12:         $state \leftarrow$  字符
13:      else if  $ch == \text{'\"}}$  then
14:         $state \leftarrow$  字符串
15:      end if
16:    end if
17:  else if  $state ==$  短横线 then
18:    if  $ch == \text{'\text{'}}$  then
19:       $state \leftarrow$  单行注释
20:    else
21:       $current\_string += \text{'\"}}$  +  $ch$ 
22:       $state \leftarrow$  代码
23:    end if
24:  else if  $state ==$  斜杠 then
25:    if  $ch == \text{'\text{'}}$  then
26:       $state \leftarrow$  多行注释
27:    else if  $ch == \text{'\text{'}}$  then
28:       $state \leftarrow$  单行注释
29:    else
30:       $current\_string += \text{'\text{'}}$ 
31:       $current\_string += ch$ 
32:       $state \leftarrow$  代码
33:    end if
34:  else if  $state ==$  多行注释 then
35:    if  $ch == \text{'\text{'}}$  then
36:       $state \leftarrow$  多行注释遇到 *
37:    else
38:      if  $ch == \text{'\textbackslashn'}}$  then
39:         $current\_string += \text{'\textbackslashr\textbackslashn'}}$ 
40:       $state \leftarrow$  多行注释

```

```

41:      end if
42:  end if
43:  else if state == 多行注释遇到 then
44:      if ch == '/' then
45:          state ← 代码
46:      else if ch == '*' then
47:          state ← 多行注释遇到
48:      else
49:          state ← 多行注释
50:      end if
51:  else if state == 单行注释 then
52:      if ch == '\\' then
53:          state ← 拆行注释
54:      else if ch == '\n' then
55:          current_string += '\r\n'
56:          state ← 代码
57:      else
58:          state ← 单行注释
59:      end if
60:  else if state == 拆行注释 then
61:      if ch == '\\' OR ch == '\r' OR ch == '\n' then
62:          if ch == '\n' then
63:              current_string += '\r\n'
64:          end if
65:          state ← 拆行注释
66:      else
67:          state ← 单行注释
68:      end if
69:  else if state == 字符 then
70:      current_string += ch
71:      if ch == '\\' then
72:          state ← 字符中的转义字符
73:      else if ch == '\" then
74:          state ← 代码
75:      else
76:          state ← 字符
77:      end if
78:  else if state == 字符中的转义字符 then
79:      current_string += ch

```

```

80:   state ← 字符
81: else if state == 字符串 then
82:   current_string += ch
83:   if ch == '\\' then
84:     state ← 字符串中的转义字符
85:   else if ch == '\"' then
86:     state ← 代码
87:   else
88:     state ← 字符串
89:   end if
90: else if state == 字符串中的转义字符 then
91:   current_string += ch
92:   state ← 字符串
93: end if
94: end for
95: Text_preprocessed += current_string

```

2.3.5 词法分析器解析

词法分析器将去除注释的 SQL -- 语言源代码逐行输入到有限状态机逻辑 2.3.5，初步解析出单词和单词符号大类，再经过符号表项解析逻辑得到准确的单词符号种别和单词符号内容。

值得说明的是，对于复合单词“ORDER BY”与“GROUP BY”，当词法分析器识别到“ORDER”或“GROUP”后超前搜索 4 个字符，如果为“BY”则将两者合并输出，否则转到词法检查逻辑，根据具体情况决定是单独输出“ORDER”/“GROUP”或是抛出异常。

2.4 输出格式说明

2.4.1 正则表达式转 NFA 算法

```

Input: (12|345)*123
===== NFA =====
start_node:9
node_head:13
Edge 0 -> 1, node_type:1
Edge 1 -> 2, node_type:2
Edge 3 -> 4, node_type:3
Edge 4 -> 5, node_type:4
Edge 5 -> 6, node_type:5
Edge 7 -> 0, node_type:varepsilon
Edge 7 -> 3, node_type:varepsilon
Edge 2 -> 8, node_type:varepsilon
Edge 6 -> 8, node_type:varepsilon
Edge 9 -> 10, node_type:varepsilon
Edge 9 -> 7, node_type:varepsilon
Edge 8 -> 10, node_type:varepsilon
Edge 8 -> 7, node_type:varepsilon
Edge 10 -> 11, node_type:1
Edge 11 -> 12, node_type:2
Edge 12 -> 13, node_type:3

```

2.4.2 NFA 确定化算法

```

Input: (12|345)*123
===== DFA =====
start_node:0
node_heads:[5]
Edge 0 -> 1, node_type:1
Edge 0 -> 2, node_type:3
Edge 1 -> 3, node_type:2
Edge 2 -> 4, node_type:4
Edge 3 -> 1, node_type:1
Edge 3 -> 5, node_type:3
Edge 4 -> 6, node_type:5
Edge 5 -> 4, node_type:4
Edge 6 -> 1, node_type:1
Edge 6 -> 2, node_type:3

```

2.4.3 DFA 最小化算法

```

Input: (12|345)*123
===== Construct matrix =====
['1', '2', '3', '4', '5']
[1, -1, 2, -1, -1]
[-1, 3, -1, -1, -1]
[-1, -1, -1, 4, -1]
[1, -1, 5, -1, -1]
[-1, -1, -1, -1, 6]
[-1, -1, -1, 4, -1]
[1, -1, 2, -1, -1]
=====
===== Minimized DFA =====
start_node:0
node_heads:[5]
Edge 1 -> 3, node_type:2
Edge 2 -> 4, node_type:4
Edge 3 -> 1, node_type:1
Edge 3 -> 5, node_type:3
Edge 4 -> 0, node_type:5
Edge 5 -> 4, node_type:4
Edge 0 -> 1, node_type:1
Edge 0 -> 2, node_type:3

```

2.4.4 TOKEN 序列

Token 输出格式:

[待测代码中的单词符号] [TAB] <[单词符号种别],[单词符号内容]>

其中, 单词符号种别为 KW (关键字)、OP (运算符)、SE (界符)、IDN (标识符)、INT (整形数)、FLOAT (浮点数)、STRING (字符串); 单词符号内容 KW、OP、SE 为其编号, KW 见表 2.1, OP 见表 2.1, SE 见表 2.1; 其余类型为其值。

2.5 源程序编译步骤

本项目的词法分析器采用 python 编写, 无须编译, 可以直接运行。其中, 代码对环境的要求如下:

```
python3
```

在工程根目录终端下, 键入如下命令即可快速运行:

```
python main.py
```

其中, main.py 文件的输入为 ./tests 目录下的 "testcase-*.sql" 文件, 通过词法分析与语法分析分别输出 token 序列和规约序列。

若要单独查看与调试词法分析器本身的功能, 可以在工程根目录下键入以下命令:


```
python ./src/lexer/lexical_analyzer.py
```

若要查看与调试词法分析器子模块 X(可选项为 nfa2dfa, rex2nfa, dfa2min) 的功能, 可以在工程根目录下键入以下命令:

```
python ./src/lexer/X.py
```

3 语法分析器

3.1 数据结构描述

3.1.1 GRAMMAR 相关数据结构

语法分析前需要首先对文法进行预处理, 即将以文本文件形式存储的语法转换为 $[V_N, V_T, P, S]$ 的形式, 其中 V_N 为非终结符集, V_T 为终结符集, P 为产生式集, S 为开始符号。因此分别设计 **production** (产生式) 与 **grammar** (文法) 两个抽象数据类型, 来实现相应类型数据的存储与计算, 分别如下:

ADT 7: production 抽象数据类型

```
ADT production {
    数据对象: D = {
        左部符号 left,
        右部符号 right,
        产生式编号 index
    }
    基本操作:
        init(left, right, index)
            初始条件: left存在, right与index合法
            操作结果: 初始化production
} ADT production
```

ADT 8: grammar 抽象数据类型

```
ADT grammar {
    数据对象: D = {
        非终结符集 non_terminal,
        终结符集 terminal,
        产生式集 productions,
        初始符号 start
        grammar文本内容 grammar_context
    }
```

}

基本操作:

`init(grammar_file_path)`

初始条件: 路径存在

操作结果: 初始化`grammar`各变量`load_grammar(grammar_file_path)`

初始条件: 路径存在

操作结果: 根据`grammar`文本文件, 计算得到`non_terminal`, `terminal`, `productions`, `start`等成员对象`get_augmented_grammar()`初始条件: `grammar`加载完成操作结果: 将`grammar`转换为其对应的增广文法

} ADT grammar

3.1.2 FIRST & FOLLOW 集合

在计算项目集规范族之前, 需要首先根据文法计算出该文法对应的 FIRST 集与 FOLLOW 集 (事实上, 在计算项目集规范族时 FOLLOW 集并不是必要的, 但是为了完整, 我们仍旧设计该数据类型与算法)。为了能够规范地存储与计算 FIRST 集与 FOLLOW 集, 设计 FIRST 与 FOLLOW 两个抽象数据类型, 分别如下:

ADT 9: FIRST 抽象数据类型

ADT FIRST {

数据对象: $D = \{$ FIRST集合 `first_sets`,文法对象 `grammar_obj`

}

基本操作:

`init(grammar_obj)`初始条件: `grammar_obj`存在

操作结果: 初始化FIRST

`calculate_first_sets()`

初始条件: FIRST已初始化

操作结果: 根据文法计算该文法中各个符号的FIRST集合, 并存入`first_sets`中`calculate_follow_set(symbols)`初始条件: FIRST已初始化、`first_sets`已计算完成、`symbols`合法操作结果: 返回符号串`symbols`的FIRST集合`dump_first_sets_into_file(file_path)`初始条件: FIRST已初始化、`first_sets`已计算完成、`file_path`合法

操作结果：将FIRST按照规定格式写入文件file_path中

} ADT FIRST

ADT 10: FOLLOW 抽象数据类型

ADT FOLLOW {

数据对象：D = {

FOLLOW集合 follow_sets

}

基本操作：

init(grammar_obj, FIRST_obj)

初始条件：grammar_obj存在，FIRST_obj存在，并且FIRST_obj与grammar_obj相对应

操作结果：初始化FOLLOW

calculate_follow()

初始条件：FOLLOW已初始化

操作结果：根据文法grammar_obj与其对应的FIRST_obj计算该文法中各个符号（终结符）的FOLLOW集合，并存入follow_sets中

dump_follow_sets_into_file(file_path)

初始条件：FOLLOW已初始化、follow_sets已计算完成、file_path合法

操作结果：将FOLLOW按照规定格式写入文件file_path中

} ADT FOLLOW

3.1.3 项目集规范族相关数据类型

计算文法的项目集规范族时，涉及到了较多的算法和表示形式，因此分别设计如下三个抽象数据类型：

ADT 11: item 抽象数据类型

ADT item {

数据对象：D = {

项目产生式左端符号 left,

项目产生式右端符号 right,

项目dot位置 dot_pos,

项目对应的终结符号集合 terminals,

项目对应的产生式的编号 index

}

基本操作：

init(left, right, dot_pos, terminals, index)

初始条件: left, right, dot_pos, terminals, index合法

操作结果: 初始化item

go(symbol)

初始条件: item已初始化、symbol合法

操作结果: 计算move(item, symbol), 即根据输入的终结符得到下一个项目。

} ADT item

ADT: 12 itemSet 抽象数据类型

ADT itemSet {

数据对象: D = {

项目集合 item_set,

项目集合的编号 index

}

基本操作:

init(item_set, index)

初始条件: item_set合法、index合法

操作结果: 初始化itemSet

calculate_closure(grammar_obj, FIRST_obj)

初始条件: itemSet已初始化、grammar_obj存在、FIRST_obj存在,
并且FIRST_obj

与grammar_obj相对应

操作结果: 计算项目集合的闭包, 项目增加入item_set中

} ADT itemSet

ADT 13: itemSets 抽象数据类型

ADT itemSets {

数据对象: D = {

项目集规范族 item_sets

go关系 go

}

基本操作:

init()

初始条件: 无

操作结果: 初始化itemSets

calculate_itemSets(grammar_obj, FIRST_obj)

```

    初始条件: itemSets已初始化、grammar_obj存在、FIRST_obj存在,
    并且FIRST_obj与grammar_obj相对应
    操作结果: 根据语法计算其对应的项目集规范族, 并存入item_sets中
calculate_go(itemSet_obj, grammar_obj)
    初始条件: itemSet_obj合法、grammar_obj存在
    操作结果: 计算itemSet_obj在grammar_obj的文法下可以跳转到的项目集,
    并将相应的关系存入go中
convert_go()
    初始条件: itemSets已初始化
    操作结果: 将成员变量go从原来的列表类型转换为字典类型
merge_item_sets()
    初始条件: itemSets各成员变量已经计算完成
    操作结果: 将成员变量item_sets中各个项目集进行merge操作, 将项目集
    内具有相同产生式和dot位置的项目对应的终结符号集合合并起来, 作为一个项目
dump_into_file():
    初始条件: itemSets已计算完毕
    操作结果: 将itemSets中的数据按照规定格式写入文件
} ADT itemSets

```

3.1.4 LR1 分析表相关数据类型

构建 LR1 分析表时, 我们需要根据文法分别计算 `action_table` 以及 `goto_table`, 因此分别设计 `action` 与 `goto` 相关的两个抽象数据类型, 如下:

ADT 14: actionTable 抽象数据类型

```

ADT actionTable {
    数据对象: D = {
        动作表 action_table
    }
    基本操作:
        init()
            初始条件: 无
            操作结果: 初始化actionTable
        calculate_action_Table(grammar_obj, itemSets_obj)
            初始条件: actionTable已初始化、grammar_obj存在、itemSets_obj存在,
            并且itemSets_obj与grammar_obj相对应
            操作结果: 根据语法规则和项目集规范族计算其对应的动作表, 并存入
            action_table中
        dump_table_into_file(file_path)

```

初始条件: `action_table`已计算完毕、`file_path`合法
 操作结果: 将`action_table`中的数据按照规定格式写入文件

```
} ADT actionTable
```

ADT 15: gotoTable 抽象数据类型

```
ADT gotoTable {
  数据对象: D = {
    跳转表 goto_table
  }
  基本操作:
    init()
      初始条件: 无
      操作结果: 初始化gotoTable
    calculate_goto_Table(grammar_obj, itemSets_obj)
      初始条件: gotoTable已初始化、grammar_obj存在、itemSets_obj存在, 并且
      itemSets_obj与grammar_obj相对应
      操作结果: 根据语法规则和项目集规范族计算其对应的跳转表, 并存入
      goto_table中
    dump_table_into_file(file_path)
      初始条件: goto_table已计算完毕、file_path合法
} ADT gotoTable
```

3.1.5 LR1 语法分析器

计算出前面所有的数据类型后, 我们可以进行语法分析器的进一步设计。我们将语法分析器封装成一个类, 以方便调用与计算, 其对应的抽象数据类型如下:

ADT 16: LR1_parser 抽象数据类型

```
ADT LR1_parser {
  数据对象: D = {
    语法分析器对应文法 grammar_obj,
    语法分析器对应FIRST集 FIRST_obj,
    语法分析器对应FOLLOW集 FOLLOW_obj,
    语法分析器对应项目集规范族 itemSets_obj,
    语法分析器对应分析表 analysisTable_obj,
    语法分析器基本配置项 Configs
  }
  基本操作:
    init(Configs)
```

```

    初始条件: Configs合法
    操作结果: 根据Configs初始化LR1_parser
  LR1_init(grammar_path)
    初始条件: grammar_path合法
    操作结果: 根据grammar_path计算LR1_parser所有成员变量
  parse(input_stack, REDIRECT_STDOUT_TO_FILE, save_path)
    初始条件: input_stack合法、REDIRECT_STDOUT_TO_FILE合法、save_path合法
    操作结果: 对input_stack进行语法分析。若REDIRECT_STDOUT_TO_FILE
    为真, 则将语法分析过程输出到文件save_path中;
} ADT LR1_parser

```

3.2 算法描述

3.2.1 计算语法四元组

Algorithm 5 Load grammar

Inputs:

grammar.txt

```

1: 初始化:  $non\_terminals \leftarrow [], terminals \leftarrow [], productions \leftarrow [], start \leftarrow None$ 
2:  $all\_symbols \leftarrow [], cnt \leftarrow 1$ 
3: for grammar.txt 中的每一个产生式 do
4:   将产生式左端符号 append 到  $non\_terminals$  中
5:   将产生式左端符号 append 到  $all\_symbols$  中
6:   initialize production
7:    $production.left \leftarrow$  产生式左端符号
8:    $production.right \leftarrow$  产生式右端符号集
9:    $production.idx \leftarrow cnt$ 
10:   $cnt \leftarrow cnt + 1$ 
11:  for 产生式右端符号  $i$  do
12:    将产生式右端符号  $i$  append 到  $all\_symbols$  中
13:  end for
14: end for
15:  $Unique(non\_terminals, terminals)$ 
16:  $terminals = all\_symbols - non\_terminals$ 
17:  $start \leftarrow productions[0].left$ 
18: 初始化 grammar
19:  $grammar.non\_terminals \leftarrow non\_terminals$ 
20:  $grammar.terminals \leftarrow terminals$ 
21:  $grammar.productions \leftarrow productions$ 

```

```

22: grammar.start  $\leftarrow$  start
23: return grammar

```

3.2.2 求 FIRST 集

Algorithm 6 Calculate First Set

Inputs:

文法 *grammar*

```

1: while FIRST 仍在增大 do
2:   for N in grammar.non_terminals do
3:     First[N]  $\leftarrow$  []
4:   end for
5:   for T in grammar.terminals do
6:     First[T]  $\leftarrow$  [T]
7:   end for
8:   for P in grammar.productions do
9:     if P.right[0] in grammar.erminals then
10:      First[P.left].append(P.right[0])
11:    end if
12:    if P.right = [] then
13:      First[P.left].append( $\varepsilon$ )
14:    end if
15:    if P.right 的左侧含有连续 i 个符号位于 grammar.non_terminals 中, 且对任意 i,
      First[i] 均含  $\varepsilon$  then
16:      First[P.left].extend(FIRST[i + 1] - [ $\varepsilon$ ])
17:    end if
18:  end for
19: end while
20: return First

```

3.2.3 求 FOLLOW 集

Algorithm 7 Calculate Follow Set

Inputs:

文法 *grammar*, FIRST 集

```

1: while FOLLOW 仍在增大 do
2:   for N in grammar.non_terminals do
3:     Follow[N]  $\leftarrow$  []
4:   end for
5:   Follow[grammar.start]  $\leftarrow$  [#]

```



```

6:  for  $P$  in  $grammar.productions$  do
7:    if  $P$  的形式为  $A \rightarrow \alpha B \beta$ , 其中  $A, B$  位于  $grammar.non\_terminals$  中,  $\alpha, \beta$  为任意符号串 then
8:       $Follow[B].extend(FIRST[\beta] - [\epsilon])$ 
9:    end if
10:   if  $P$  的形式为  $A \rightarrow \alpha B$ , 其中  $A, B$  位于  $grammar.non\_terminals$  中,  $\alpha$  为任意符号串 then
11:      $Follow[B].extend(Follow[A])$ 
12:   end if
13:   if  $P$  的形式为  $A \rightarrow \alpha B \beta$ , 其中  $A$  位于  $grammar.non\_terminals$  中,  $\alpha, \beta$  为任意符号串, 且  $\epsilon$  含在  $First[\beta]$  中 then
14:      $Follow[B].extend(Follow[A])$ 
15:   end if
16: end for
17: end while
18: return  $Follow$ 

```

3.2.4 求闭包

Algorithm 8 Calculate Closure

Inputs:

文法 $grammar$, FIRST 集, FOLLOW 集, 项目集 I

```

1:  $closure[I] \leftarrow []$ 
2: for  $i$  in  $I$  do
3:    $closure[I].append(i)$ 
4: end for
5: while  $closure[I]$  仍在增大 do
6:   for  $item$  in  $closure[I]$  do
7:     if  $item$  的形式为  $A \rightarrow \alpha \bullet B \beta, a$ , 其中  $A, B$  位于  $grammar.non\_terminals$  中,  $\alpha, \beta$  为任意符号串,  $a$  位于  $grammar.terminals$  中 then
8:       for  $grammar.productions$  中所有满足  $P.left = B$  的  $P$  do
9:         for 所有  $b$  in  $grammar.terminals$ , 并且  $b$  位于  $First[\beta a]$  中 do
10:           $closure[I].append(B \rightarrow \alpha \gamma, b)$ 
11:        end for
12:      end for
13:    end if
14:  end for
15: end while
16: return  $closure[I]$ 

```

3.2.5 求项目集规范族与 GO

Algorithm 9 Calculate Item Set and GO

Inputs:

文法 *grammar*, FIRST 集, FOLLOW 集

```

1: 初始化:  $itemset \leftarrow []$ 
2:  $itemset.append(grammar.start \rightarrow \cdot grammar.right, \#)$ 
3:  $itemset[0] \leftarrow closure[itemset[0]]$ 
4: 计算  $itemset[0]$  可以 GO 到项目集  $I_s$ , 其中  $GO(itemset[0], A) = closure[J]$ , 其中  $J = \{\text{任何形如 } [A \rightarrow \alpha B \cdot \beta, a] \text{ 的项目 } |[A \rightarrow \alpha B \beta, a] \text{ 在 } itemset[0] \text{ 中}\}$ . 同时将相应的 GO 关系加入 GO 中。
5: for  $I$  in  $I_s$  do
6:    $itemset.append(I)$ 
7: end for
8: while  $itemset$  仍在增大 do
9:   for  $I$  in  $itemset$  do
10:    计算  $I$  可以 GO 到项目集  $I_s$ , 其中  $GO(itemset[0], A) = closure[J]$ , 其中  $J = \{\text{任何形如 } [A \rightarrow \alpha B \cdot \beta, a] \text{ 的项目 } |[A \rightarrow \alpha B \beta, a] \text{ 在 } itemset[0] \text{ 中}\}$ . 同时将相应的 GO 关系加入 GO 中。
11:    for  $I$  in  $I_s$  AND  $I$  不在  $itemset$  中 do
12:       $itemset.append(I)$ 
13:    end for
14:  end for
15: end while
16: return  $itemset, GO$ 

```

3.2.6 求 ACTION 表与 GOTO 表

Algorithm 10 Calculate Analysis Table

Inputs:

项目集规范族 $itemsets$, GO 关系, 文法 *grammar*

```

1: for  $itemset \in itemsets$  do
2:   for  $I \in itemset$  do
3:     if  $I$  的形式为  $[A \rightarrow \alpha a \beta, b]$ , 且  $Go[I, a] = I_j$ ,  $a$  在  $grammar.terminal$  中 then
4:        $action[itemset.idx, a] \leftarrow s_j$ 
5:     end if
6:     if  $I$  的形式为  $[A \rightarrow \alpha \cdot, a]$  then
7:        $action[itemset.idx, a] \leftarrow r_j$ , 其中  $r_j$  为  $I$  的编号
8:     end if
9:   if  $I$  的形式为  $[grammar.start_p \rightarrow grammar.start \cdot, \#]$  then

```

```

10:       $action[itemset.idx, \#] \leftarrow acc$ 
11:    end if
12:    if 不满足以上所有情况 then
13:      抛出错误
14:    end if
15:  end for
16: end for
17: for  $go$  in  $GO$  do
18:    $goto[go.idx, go.symbol] \leftarrow go.next$ 
19: end for
20: return  $action, goto$ 

```

3.2.7 LR1 PARSE 解析

Algorithm 11 Parse

Inputs:

待解析的 $input_stack(token \text{ 序列})$ 以及 $action, goto$, 文法 $grammar$

```

1: 初始化:  $state\_stack = Stack()$ ,  $symbol\_stack = Stack()$ 
2: 初始化:  $state\_stack.push(0)$ ,  $symbol\_stack.push(\#)$ 
3: while  $input\_stack$  不空 do
4:    $top\_state = state\_stack.top()$ ,  $top\_symbol = symbol\_stack.top()$ 
5:   if  $top\_symbol = \#$  or  $top\_symbol$  in  $grammar.terminal$  then
6:      $action\_now \leftarrow action[top\_state, top\_symbol]$ 
7:     if  $action\_now = acc$  then
8:       将  $top\_symbol$ ,  $top\_state$  与  $acc$  等信息打印
9:       解析成功
10:      return
11:    end if
12:    if  $action\_now = None$  then
13:      将  $top\_symbol$ ,  $top\_state$  与  $err$  等信息打印
14:      语法错误
15:      return
16:    end if
17:    if  $action\_now[0] = s$  then
18:      将  $top\_symbol$ ,  $top\_state$  与  $move$  等信息打印
19:      将  $action\_now.state$  压入  $state\_stack$ ,  $input\_stack.peak()$  压入  $symbol\_stack$ , 将
        $input\_stack$  栈顶元素弹出
20:    end if
21:    if  $action\_now[0] = r$  then

```

```

22:      将 top_symbol, top_state 与 reduction 等信息打印
23:      按照 action_now.state 对应的 grammar 将 symbol_stack 栈顶几个元素进行规约, 并
      将规约结果压入 symbol_stack
24:      根据 symbol_stack.top() 与当前状态 top_state 查找 goto 表, 将 goto 表中的结果压入
      state_stack
25:      end if
26:  end if
27:  if top_symbol in grammar.non_terminals then
28:    goto_now  $\leftarrow$  goto[top_state, top_symbol]
29:    if goto_now = None then
30:      将 top_symbol, top_state 与 err 等信息打印
31:      语法错误
32:      return
33:    end if
34:    if goto_now[0] = s then
35:      将 top_symbol, top_state 与 move 等信息打印
36:      将 goto_now.state 压入 state_stack, input_stack.peak() 压入 symbol_stack, 将
      input_stack 栈顶元素弹出
37:    end if
38:  end if
39: end while

```

3.3 LR1 分析表描述

LR1 分析表分为 *action* 表与 *goto* 表, 其对应的抽象数据类型详见 3.1.4, 其对应的构造算法详见 3.2.6, 输出示例详见 3.4.3。下面对该分析表结构及其功能进行简要描述:

LR1 分析表主要由两部分组成, 分别为 *action_table*(动作表) 与 *goto_table*(状态转移表)。

在实现上, *action_table* 是一个字典, 其中 *key* 为二元元组, 第一个元素为分析栈当前状态, 第二元素为当前面临输入栈的栈顶符号 (该符号的合法值为 *grammar* 的终结符或代表输入结束的“#”符号); *value* 为此时应该采取的动作。动作类型包括有“s*”, “r*”, “acc”, “err”, 其中“*”表示任意数字。分别对应的动作为:

- “s*”: 将当前输入栈栈顶符号放到分析栈 (符号栈) 栈顶, 弹出输入栈栈顶符号, 同时进行状态转移, 即将“*”对应的状态放到状态栈栈顶;
- “r*”: 利用序号为 * 的产生式对分析栈 (符号栈) 进行规约操作, 具体为: 将序号为 * 的产生式右端符号依次从分析栈 (符号栈) 栈顶弹出, 同时在状态栈中弹出其对应的状态; 最后将序号为 * 的产生式左端符号放到分析栈 (符号栈) 栈顶, 同时将 *goto* 表中对应的状态放置到状态栈顶。
- “acc”: 规约成功, 程序退出。

- “err”: 规约出错, 程序退出。

`goto_table` 也是以字典的方式实现, 字典的 `key` 与 `value` 和 `action_table` 的设置相同, 只不过 `key` 的第二个元素限定为非终结符。`goto` 表的动作类型包括有“`s*`”, “`err`”, 其中 “`*`” 表示任意数字。分别对应的动作为:

- “`s*`”: 将当前输入栈栈顶符号放到分析栈 (符号栈) 栈顶, 弹出输入栈栈顶符号, 同时进行状态转移, 即将“`*`”对应的状态放到状态栈栈顶;
- “`err`”: 规约出错, 程序退出。

进行语法分析时, 首先根据 3.2 中的算法计算语法对应的分析表, 然后分别对状态栈、分析栈与输入栈进行初始化操作 (细节详见 3.2.7), 根据当前栈的状态查询 LR1 分析表进行对应的动作即可。

3.4 输出格式说明

3.4.1 FIRST 集与 FOLLOW 集输出

FIRST 集每一行的输出格式:

[符号 (终结符或非终结符)] [=] [符号对应的 FIRST 集]

输出示例:

```
FIRST:
!=['!']
!=='['! !=']
#=['#']
&&=['&&']
(['(']
)=[')']
*=['*']
,=[',']
-=['-']
.=['.']
<=['<']
<==['< ==']
<=>=['< >']
==['=']
>=['>']
>==['> ==']
ALL=['ALL']
AND=['AND']
AS=['AS']
AVG=['AVG']
```

FOLLOW 集每一行的输出格式:

[符号 (终结符或非终结符)] [=] [符号对应的 FOLLOW 集]

输出示例:

FOLLOW:

```

expressionOrDefault=[''],'WHERE',' ','#']
expressionOrDefaultListRec=['')']
expressionRec=['#','HAVING','ORDER_BY','UNION',')']
expressionRight=['WHERE',' ','HAVING','ORDER_BY','UNION','GROUP_BY','LEFT','RIGHT',')','JOIN',' ','']
expressions=['#','HAVING','ORDER_BY','UNION',')']
expressionsWithDefaults=[')']
expressionsWithDefaultsListRec=['#']

```

3.4.2 项目集规范族输出

项目集规范族中每一个项目集的输出格式：

[itemset] [项目集在项目集规范族中的序号] [:] [\n]

[项目1] [\n]

[项目2] [\n]

.....

[项目n] [\n] [\n]

其中，项目的输出格式如下：

[左侧符号] [,] [右侧符号集] [,] [dot位置] [,] [终结符号集] [,] [产生式序号]

输出示例：

```

item_set0:
rootp, ['root'], 0, ['#'], 0,
root, ['dmlStatement'], 0, ['#'], 1,
dmlStatement, ['selectStatement'], 0, ['#'], 2,
dmlStatement, ['insertStatement'], 0, ['#'], 3,
dmlStatement, ['updateStatement'], 0, ['#'], 4,
dmlStatement, ['deleteStatement'], 0, ['#'], 5,
selectStatement, ['querySpecification', 'unionStatements'], 0, ['#'], 6,
insertStatement, ['insertKeyword', 'tableName', 'insertStatementRight'], 0, ['#'], 112,
updateStatement, ['UPDATE', 'tableName', 'elementNameAlias', 'SET', 'updatedElement', 'updatedElementL
deleteStatement, ['DELETE', 'FROM', 'tableName', 'deleteStatementRight'], 0, ['#'], 135,
querySpecification, ['SELECT', 'unionType', 'selectElements', 'selectClause'], 0, ['#', 'UNION'], 15,
querySpecification, ['(', 'querySpecification', ')'], 0, ['#', 'UNION'], 16,
insertKeyword, ['INSERT', 'into'], 0, ['IDN'], 115,

item_set1:
querySpecification, ['(', 'querySpecification', ')'], 1, ['#', 'UNION'], 16,
querySpecification, ['SELECT', 'unionType', 'selectElements', 'selectClause'], 0, [')'], 15,
querySpecification, ['(', 'querySpecification', ')'], 0, [')'], 16,

```

项目集规范族各项目集之间的转移关系 go 的输出格式：

[()][()][当前项目集序号][,][面临符号][()][,][转移后的项目集序号][)]

输出示例：

```

((0, '('), 1)
((0, 'DELETE'), 2)
((0, 'INSERT'), 3)
((0, 'SELECT'), 4)
((0, 'UPDATE'), 5)
((0, 'deleteStatement'), 6)
((0, 'dmlStatement'), 7)
((0, 'insertKeyword'), 8)
((0, 'insertStatement'), 9)
((0, 'querySpecification'), 10)
((0, 'root'), 11)
((0, 'selectStatement'), 12)
((0, 'updateStatement'), 13)
((1, '('), 14)
((1, 'SELECT'), 15)
((1, 'querySpecification'), 16)

```

3.4.3 LR1 分析表输出

LR1 分析表包括 action 表和 goto 表两个部分，二者分别输出。

其中：

action 表输出格式：

[(] [当前状态] [,] [面临输入符号] [)] [---->] [动作类型] [动作值]

在 action 表中，面临输入符号限定为非终结符，动作类型的可选值为“s”、“r”、“acc”、“None”（表示不存在该动作，与“err”等价）。

输出示例：

```

action table:
(0, '!') ----> None
(0, '!=') ----> None
(0, '#') ----> None
(0, '&&') ----> None
(0, '(') ----> s1
(0, ')') ----> None
(0, '*') ----> None
(0, ',') ----> None
(0, '-') ----> None
(0, '.') ----> None
(0, '<') ----> None
(0, '<=') ----> None
(0, '<=>') ----> None
(0, '=') ----> None
(0, '>') ----> None
(0, '>=') ----> None
(0, 'ALL') ----> None
(0, 'AND') ----> None
(0, 'AS') ----> None
(0, 'AVG') ----> None
(0, 'DEFAULT') ----> None
(0, 'DELETE') ----> s2
(0, 'DISTINCT') ----> None

```

goto 表输出格式：

[()] [当前状态] [,] [面临输入符号] [] [---->] [动作类型] [动作值]

在 `goto` 表中，面临输入符号限定为终结符，动作类型的可选值：“s”、“None”（表示不存在该动作，与“err”等价）。输出示例：

```
goto table:
(0, 'aggregateWindowedFunction') ----> None
(0, 'booleanLiteral') ----> None
(0, 'comparisonOperator') ----> None
(0, 'constant') ----> None
(0, 'decimalLiteral') ----> None
(0, 'deleteStatement') ----> s6
(0, 'deleteStatementRight') ----> None
(0, 'dmlStatement') ----> s7
```

3.4.4 规约序列输出

规约序列每一行的输出格式：

[序号] [TAB] [选用规则序号] [TAB] [栈顶符号]#[面临输入符号] [TAB] [执行动作]

其中，选用规则序号见附件文法规则；执行动作为“reduction”（归约），“move”（LR 分析的移进），“accept”（接受）或“error”（错误）输出示例：

```
1 1 / ##INSERT move
2 2 / INSERT#INTO move
3 3 116 INTO#IDN reduction
4 4 115 into#IDN reduction
5 5 / insertKeyword#IDN move
6 6 48 IDN#( reduction
7 7 44 uid#( reduction
8 8 / tableName#( move
9 9 / (#IDN move
10 10 48 IDN#, reduction
11 11 / uid#, move
12 12 / ,#IDN move
13 13 48 IDN#, reduction
14 14 / uid#, move
15 15 / ,#IDN move
16 16 48 IDN#) reduction
17 17 47 uid#) reduction
18 18 46 uidListRec#) reduction
19 19 46 uidListRec#) reduction
20 20 45 uidListRec#) reduction
21 21 / uidList#) move
22 22 / )#VALUES move
23 23 119 VALUES#( reduction
24 24 / insertFormat#( move
25 25 / (#STRING move
```

3.5 源程序编译步骤

本项目的 LR1 语法分析器采用 `python` 编写，无须编译，可以直接运行。其中，代码对环境的要求如下：


```
python3
matplotlib==3.4.3
pickle>=0.75
```

在工程根目录终端下，键入如下命令即可快速运行：

```
python main.py
```

其中，main.py 文件的输入为 ./tests 目录下的 "testcase-*.sql" 文件，通过词法分析与语法分析分别输出 token 序列和规约序列。

若要单独查看与调试 LR1 语法分析器本身的功能，可以在工程根目录下键入以下命令：

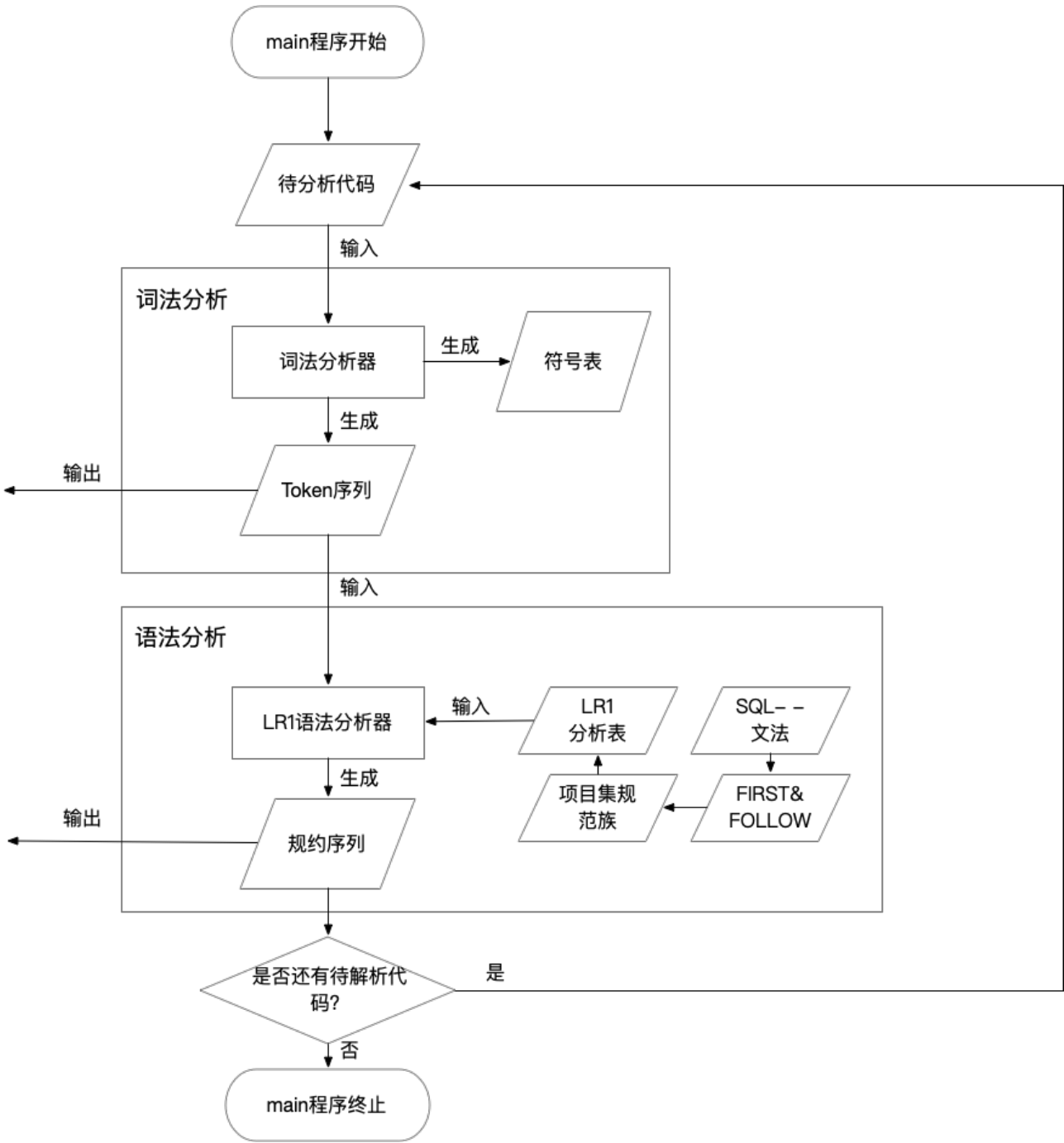
```
python ./src/LR1_parser/LR1.py
```

若要查看与调试 LR1 语法分析器子模块 X (可选项为 grammar, FF, item, analysisTable) 的功能，可以在工程根目录下键入以下命令：

```
python ./src/LR1_parser/ds/X.py
```

4 主程序分析流程

main 程序运行后，首先读取位于 ./tests 文件夹中的测试样例，然后依次调用词法分析器与语法分析器对测试样例进行分析。采用多遍编译的方法。支持连续测试多个测试样例并依次输出。相应流程图如下：



在使用前，通过修改 **Configs** 文件中的相应配置项来设置主程序的输入文法、输入测试样例、输出位置等参数，尔后在工程根目录终端下，键入如下命令即可快速运行：

```
python main.py
```