

倒排索引和布尔检索

——莎士比亚全集搜索

1511274 朱钰颖 计算机科学与技术

一、实验环境

Java 1.8

IDE IntelliJ IDEA Ultimate 版本

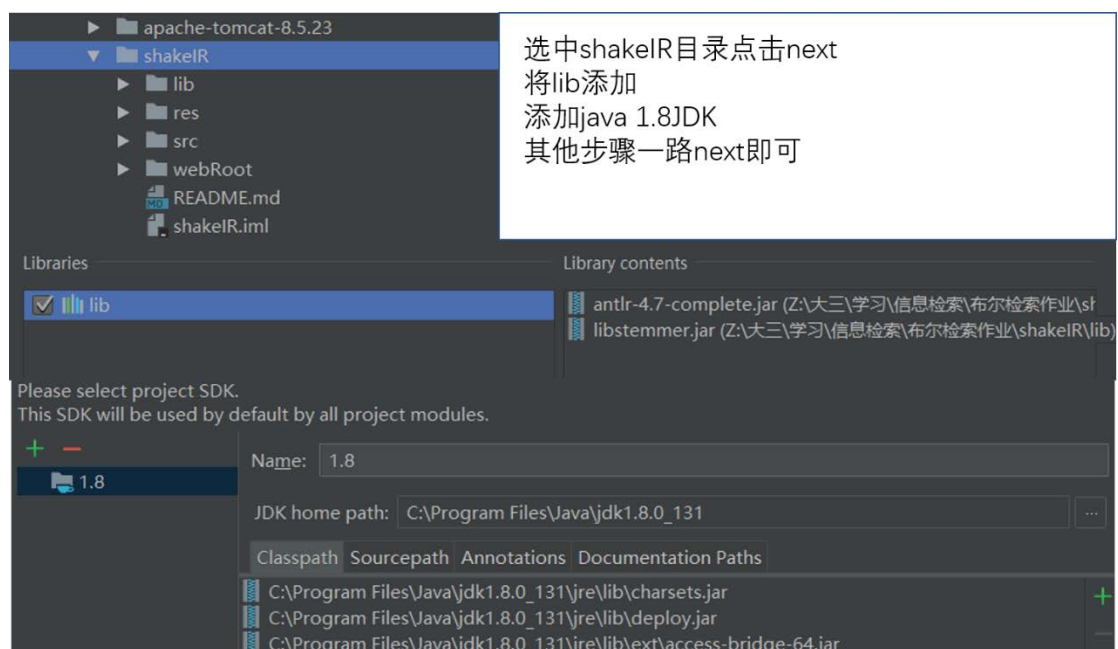
apache-tomcat-8.5.23

antlr-4.7-complete

snowball.stemmer

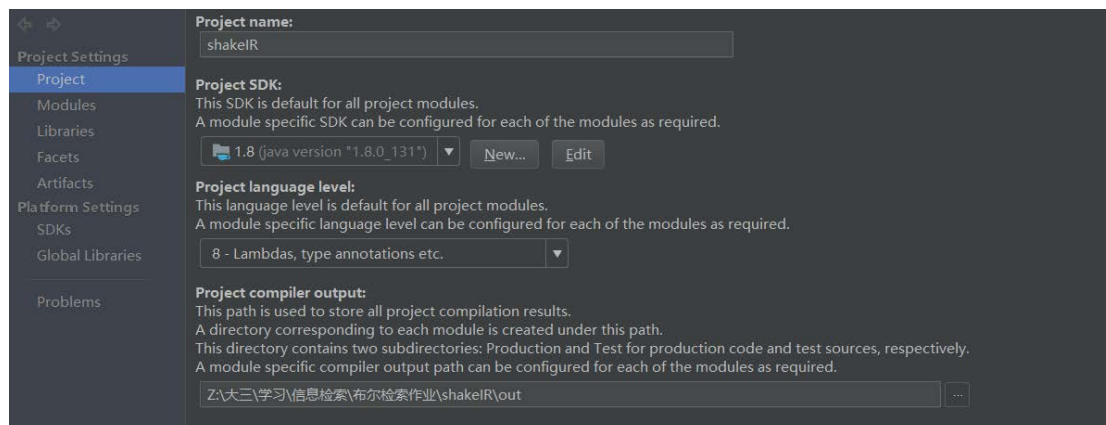
二、环境配置步骤

1、打开 intellij 后，import project，初始化步骤如下图所示：

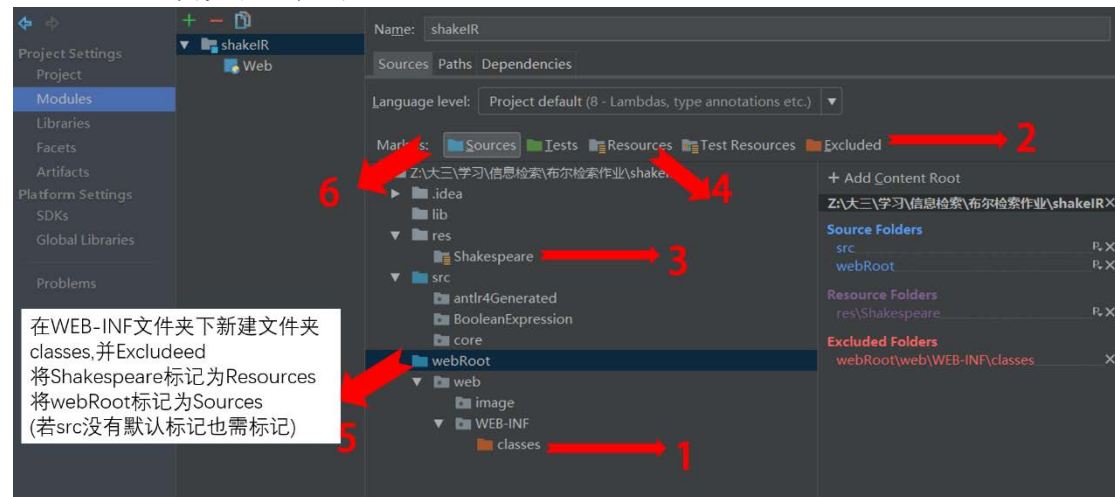


2、成功进入 intellij 后，点击 File->Project Structure

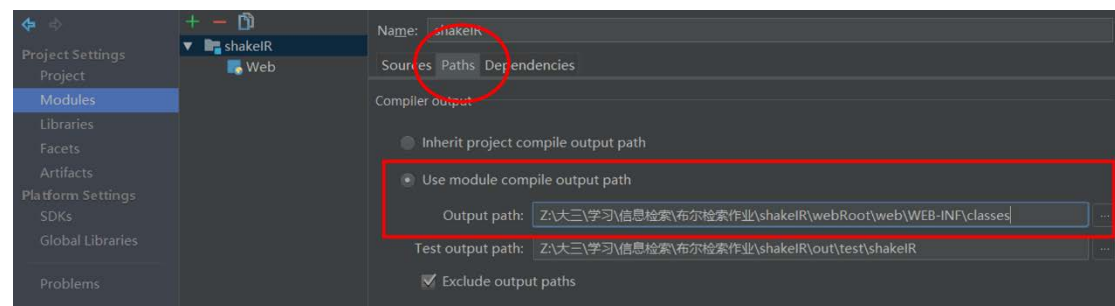
2.1Project 设置如下



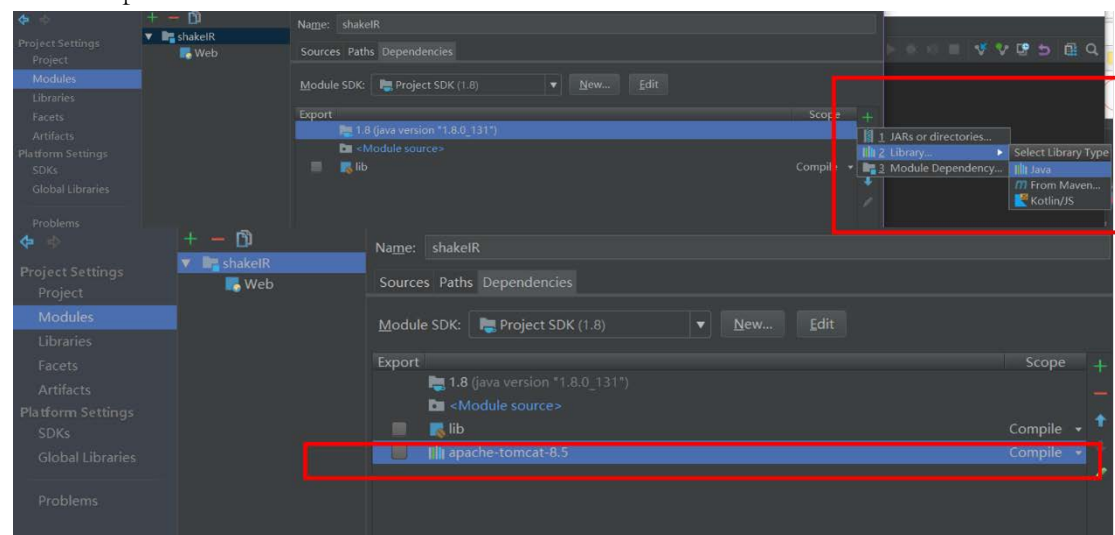
2.2 Modules 模块设置如下



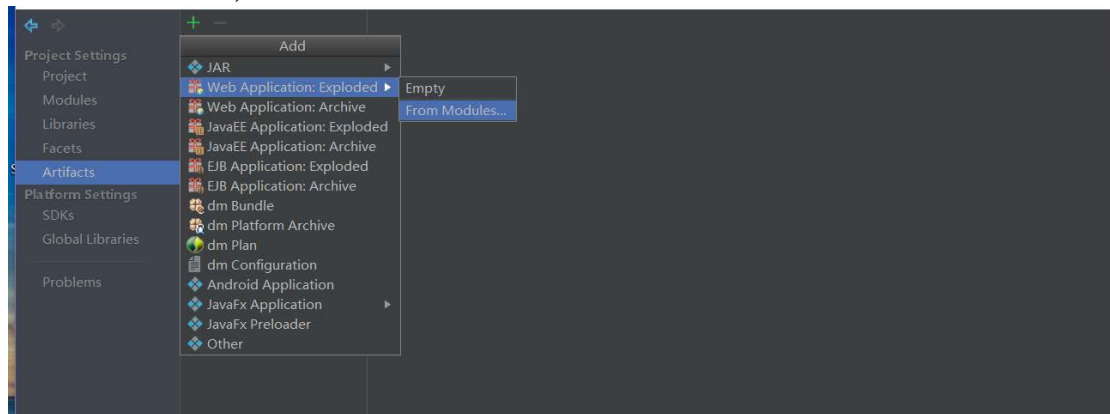
※更正：将 Shakespeare 的上层目录 res 标记为 Resources，而不是 Shakespeare 文件夹设置 Path



设置 Dependencies

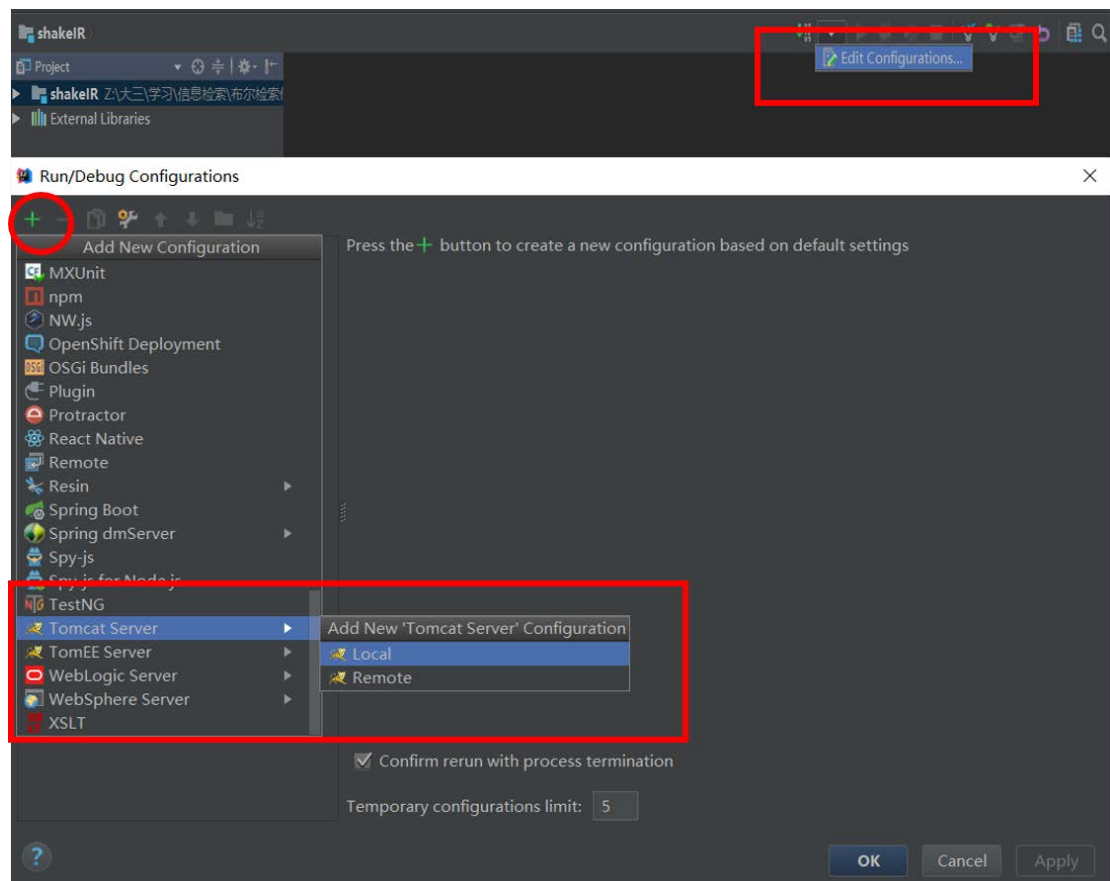


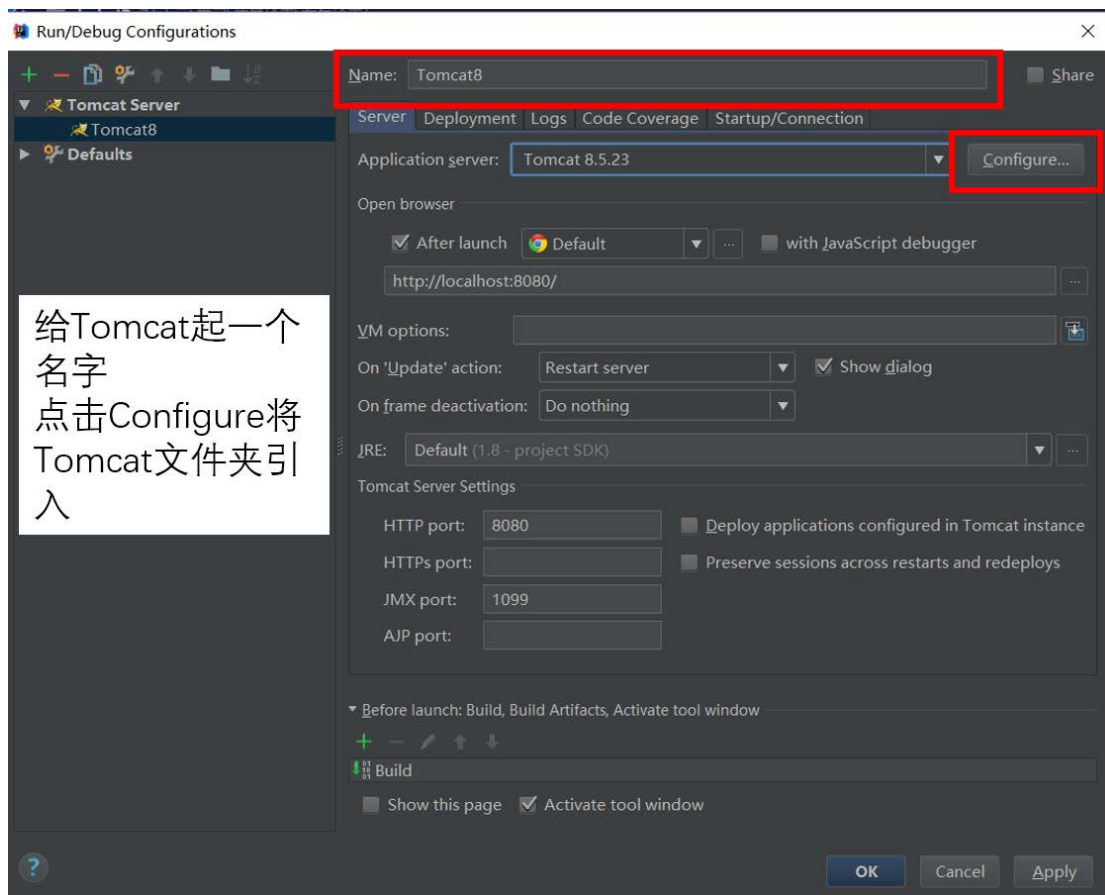
2.3 设置 Artifacts，如下图所示



3、配置 Tomcat

IDEA 右上角点击 Edit Configuration，按照图中所示添加一个本地 Tomcat

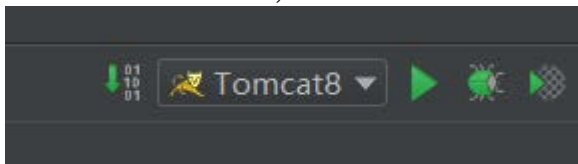




至此，环境配置完毕

三、程序运行

将选项调至 Tomcat，点击运行。

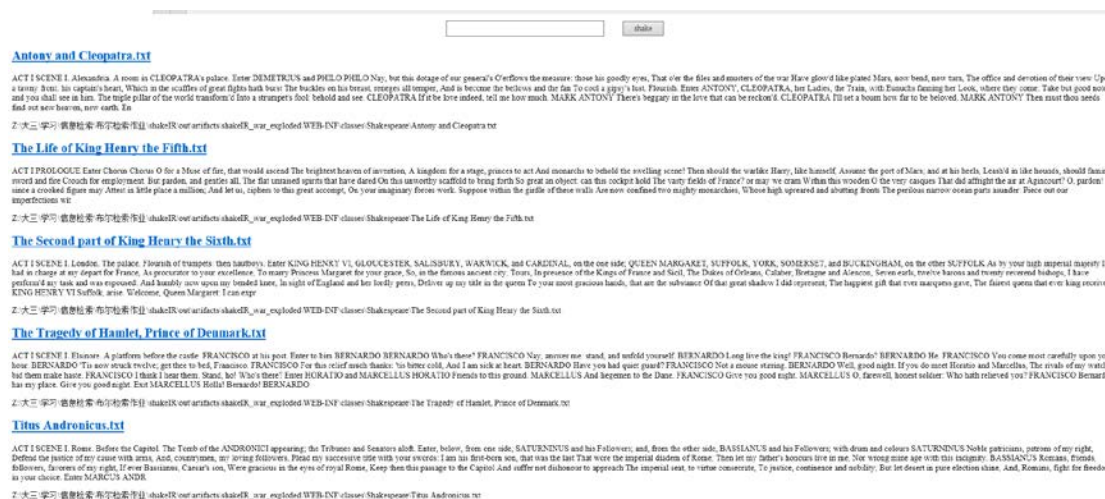


不出意外的话网页会自动跳出，若没有自动跳出，输入 `http://localhost:8080` 即可打开网页。

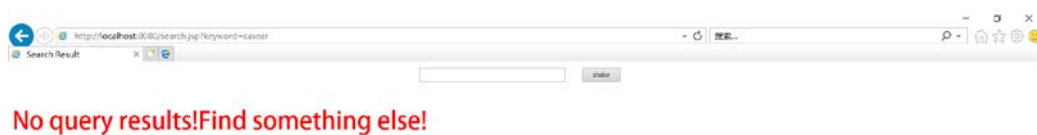
页面如下图所示，在底部文本框输入查询语句，并点击 shake 即可查询。



结果界面：输入 **Brutus AND Caesar AND NOT Calpurnia** 得到查询结果如下



在顶部文本框可继续进行查询，若查询不到则显示以下界面：



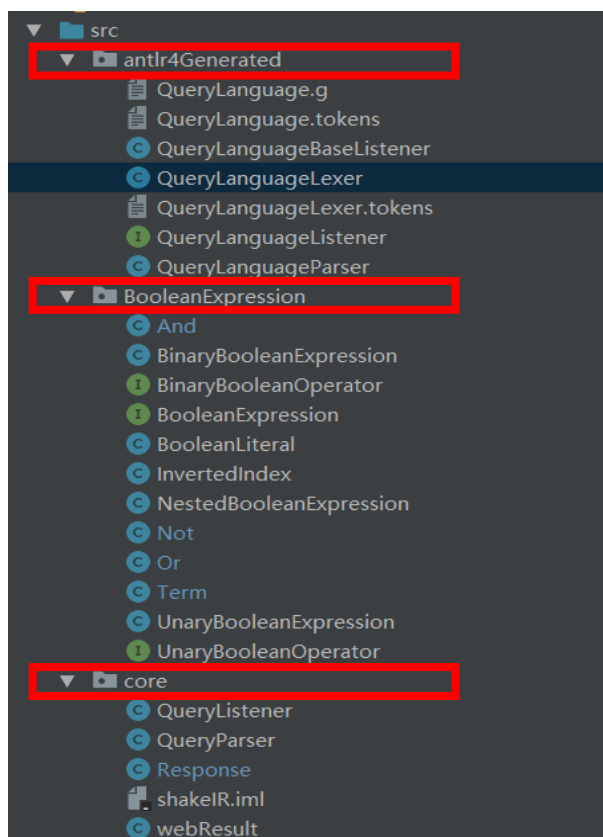
四、程序实现

4.1 程序结构

antlr4Generated 文件夹中是编写的语法文件和利用 antlr 语法工具生成的代码。

BooleanExpression 文件夹中是倒排索引和布尔检索所需计算程序

core 是核心代码



4.2 倒排索引

4.2.1 数据结构

InvertedIndex 类中有两个变量 docIndex 和 InvertedIndex

docIndex 记录键值对<文件编号,文件名字>,使用 HashMap

InvertedIndex 记录<词语, 词语所带的信息 Info>, 使用 TreeMap

Info 是一个内部类, 包含词语的一些信息。

在本程序中主要用到的变量是倒排表 posting, 其他变量虽有记录但在布尔查询时并未使用 (起初是为了进行优化功能所设)。

```
public class InvertedIndex {
    class Info {
        private int docfreq;//how many doc has this word
        private int totalfreq;//word total frequency in all files

        ArrayList<Integer> posting;
        TreeMap<Integer, Integer> posIndex;//docID+word position(last appearance)
        TreeMap<Integer, Integer> wordfreq;//word in each files,<docID,freq>

        public Info() {}

        public Info(int docfreq, TreeMap<Integer, Integer> posIndex, TreeMap<Integer, Integer> wordfreq) {}
    }

    public HashMap<Integer, String> docIndex;
    public TreeMap<String, Info> InvertedIndex;
}
```

4.2.2 方法

getIndex 为文件编号，并储存到数据结构中，且写到了 txt 文件中。

getWordFrequency 获得倒排索引表和词语频率

preprocess 利用 snowball 分词工具对文本进行预处理

```
public void getFileIndex() {...}

public void getWordFrequency() throws IOException {...}

private void getWordFrequencyHelp(String wordOfDoc, Integer docID, Integer ByteRecord) {...}

public static ArrayList<String> preprocess(String[] words) {...}
```

4.3 查询语句解析

参考资料:

<http://yijun1171.github.io/2015/03/30/ANTLR4> 学习笔记-语法字典-Grammar-Lexicon/

<https://theendian.com/blog/antlr-4-lexer-parser-and-listener-with-example-grammar/>

下载 ANTLR4 包后，编写文法文档.g 文件，使用 **java -jar path/ antlr-4.7-complete.jar /path/xx.g** 进行编译，即可自动生成该文法的解析器。

文法文档如下图所示:


```

1      grammar QueryLanguage;
2
3      @header {
4      package antlr4Generated;
5      }
6
7      TRUE:'true';
8      FALSE:'false';
9      AND:'and';
10     OR:'or';
11     NOT:'not';
12     LPAR:'(';
13     RPAR:')';
14
15     parse:expr EOF;
16
17     expr:nestedExpr|unaryExpr|expr binaryOp expr|literalExpr|termExpr;
18
19     nestedExpr:LPAR expr RPAR;
20
21     unaryExpr:unaryOp expr;
22
23     literalExpr:boolean_literal;
24
25     termExpr:TERM;
26
27     unaryOp:NOT;
28
29     binaryOp:AND|OR;
30
31     boolean_literal:TRUE|FALSE;
32     🧠
33     TERM:[a-zA-Z_][a-zA-Z_0-9]*;
34
35     WS:(' ' | '\r' | '\t' | '\u000C' | '\n')* | EOF->skip;

```

在自动生成的文件中我们只关心 QueryLanguageBaseListener, 继承这个文件, 重写里面的 enterxx 和 exitxx 方法, 编写 QueryListener 类。
编写 QueryParser 类, 将 query 传入此类中, 由 parse () 函数进行解析, match () 函数进行查询。

```

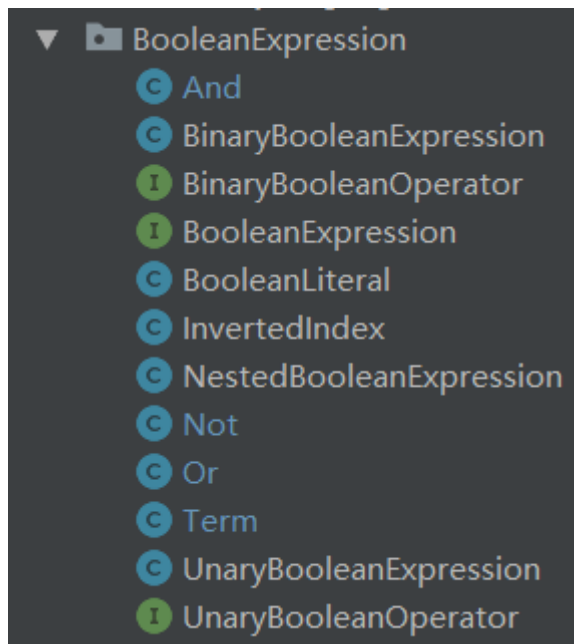
🧠
public class QueryParser {
    private BooleanExpression booleanExpression;

    public BooleanExpression parse(String query) {
        QueryLanguageLexer lexer=new QueryLanguageLexer(new ANTLRInputStream(query.toLowerCase()));
        QueryLanguageParser parser=new QueryLanguageParser(new CommonTokenStream(lexer));
        QueryListener queryListener=new QueryListener();
        parser.addParseListener(queryListener);
        parser.parse();
        booleanExpression=queryListener.getOperand();
        System.out.println("PARSED QUERY: " + booleanExpression + "\n");
        return booleanExpression;
    }

    public ArrayList<Integer> match(InvertedIndex invertedIndex){
        if(invertedIndex !=null&&invertedIndex.invertedIndex.size()>0){
            return booleanExpression.evaluate(invertedIndex);
        }
        return null;
    }
}

```


4.4 布尔查询



三个接口： BooleanExpression, BinaryBooleanOperator, UnaryBooleanOperator 定义布尔表达式，二元运算和一元运算。其他类实现这三个接口。

三个基本布尔逻辑运算实现：

And

如果返回 list 含有-1（没有在文档集中找到此词），继续向上层返回-1

通过两个变量 i, j 对两个 list 进行合并

```

16 public ArrayList<Integer> evaluate(BooleanExpression left, BooleanExpression right, InvertedIndex invertedIndex) {
17     ArrayList<Integer> lefttrtn = left.evaluate(invertedIndex);
18     ArrayList<Integer> righttrtn = right.evaluate(invertedIndex);
19     ArrayList<Integer> rtn = new ArrayList<>();
20     /*have nonexistent word*/
21     if (lefttrtn.contains(-1) || righttrtn.contains(-1)) {
22         rtn.add(-1);
23     } else {
24
25         int i = 0;
26         int j = 0;
27
28         while (i < lefttrtn.size() && j < righttrtn.size()) {
29             if (lefttrtn.get(i) == righttrtn.get(j)) {
30                 rtn.add(lefttrtn.get(i));
31                 i++;
32                 j++;
33             } else if (lefttrtn.get(i) < righttrtn.get(j)) {
34                 i++;
35             } else {
36                 j++;
37             }
38         }
39     }
40     return rtn;
41 }

```

Or

同样考虑含有-1 的情况，通常返回两个 list 的并集。

```

14 public ArrayList<Integer> evaluate(BooleanExpression left, BooleanExpression right, InvertedIndex invertedIndex) {
15     ArrayList<Integer> lefttrtn = left.evaluate(invertedIndex);
16     ArrayList<Integer> righttrtn = right.evaluate(invertedIndex);
17     ArrayList<Integer> rtn = new ArrayList<>();
18     /*have nonexistent word*/
19     if (lefttrtn.contains(-1)) {
20         return righttrtn;
21     } else if (righttrtn.contains(-1)) {
22         return lefttrtn;
23     } else if (lefttrtn.contains(-1) && righttrtn.contains(-1)) {
24         rtn.add(-1);
25         return rtn;
26     } else {
27         int i = 0;
28         int j = 0;
29
30         while (i < lefttrtn.size() && j < righttrtn.size()) {
31             if (lefttrtn.get(i) == righttrtn.get(j)) {
32                 rtn.add(lefttrtn.get(i));
33                 i++;
34                 j++;
35             } else if (lefttrtn.get(i) < righttrtn.get(j)) {
36                 rtn.add(lefttrtn.get(i));
37                 i++;
38             } else {
39                 rtn.add(righttrtn.get(j));
40                 j++;
41             }
42         }
43         if (i < lefttrtn.size()) {
44             for (; i < lefttrtn.size(); i++) {
45                 rtn.add(lefttrtn.get(i));
46             }
47         }
48         if (j < righttrtn.size()) {
49             for (; j < righttrtn.size(); j++) {
50                 rtn.add(righttrtn.get(j));
51             }
52         }
53         return rtn;
54     }
}

```

Not

取全集的补集

```
15 public ArrayList<Integer> evaluate(BooleanExpression booleanExpression, InvertedIndex invertedIndex) {
16
17     ArrayList<Integer> posting = booleanExpression.evaluate(invertedIndex);
18     ArrayList<Integer> rtn = new ArrayList<>();
19
20     int flag = invertedIndex.docIndex.size();
21     if (posting.contains(-1)) {
22         for (int i = 0; i < flag; i++) {
23             rtn.add(i);
24         }
25     } else {
26         for (int i = 0; i < flag; i++) {
27             if (posting.contains(i)) {
28                 continue;
29             } else {
30                 rtn.add(i);
31             }
32         }
33     }
34     return rtn;
35 }
```

五、程序改进

- 1、没有对索引进行压缩
- 2、没有对搜索出的文档进行排序
- 3、网页搜索结果没有链接到文档本身，即无法打开链接结果
- 4、进行布尔逻辑运算时，没有首先判断大小，进行运算顺序的优化