

Java API 文档搜索引擎

认识搜索引擎

先观察, 百度搜索的搜索结果页中, 包含了若干条结果.

每一个结果中, 又包含了图标, 标题, 描述, 展示url, 时间, 子链, 图片等.

搜索引擎的本质

输入一个查询词, 得到若干个搜索结果. 每个搜索结果包含了标题, 描述, 展示 url 和点击 url.

搜索的核心思路

当前我们有很多的网页(假设上亿个), 每个网页我们称为是一个**文档**

如何高效进行检索? 查找出有哪些网页是和查询词具有一定的相关性呢?

我们可以认为, 网页中包含了查询词(或者查询词的一部分), 就认为具有相关性.

那么我们就有了一个直观的解决思路

方案一 -- 暴力搜索

每次处理搜索请求的时候, 拿着查询词去所有的网页中搜索一遍, 检查每个网页是否包含查询词字符串.

这个方法是否可行?

显然, 这个方案的开销非常大. 并且随着文档数量的增多, 这样的开销会线性增长. 而搜索引擎往往对于效率的要求非常高.

方案二 -- 倒排索引

这是一种专门针对搜索引擎场景而设计的数据结构.

文档(doc): 被检索的html页面(经过预处理)

正排索引: "一个文档包含了哪些词". 描述一个文档的基本信息, 包括文档标题, 文档正文, 文档标题和正文的分词

/断句结果

倒排索引: "一个词被哪些文档引用了". 描述了一个词的基本信息, 包括这个词都被哪些文档引用, 这个词在该文档

中的重要程度, 以及这个词的出现位置等.

项目目标

实现一个 Java API 文档的简单的搜索引擎.

Java API 文档线上版本参见 <https://docs.oracle.com/javase/8/docs/api/index.html>

下载版本参见 <https://www.oracle.com/technetwork/java/javase/downloads/index.html>

下载之后看到一个 docs 目录, 里面存在一个 api 目录. 这里的 html 就和线上版本的文档是一一对应的.

例如, 我们熟悉的 java.util.Collection 类,

在线文档的链接是 <https://docs.oracle.com/javase/8/docs/api/java/util/Collection.html>

下载的文档的目录是 `\api\java\util\Collection.html`

这两者存在一定的对应关系.

核心流程

1. 索引模块: 扫描下载到的文档, 分析数据内容构建正排+倒排索引.
2. 搜索模块: 根据输入的查询词, 基于正排+倒排索引进行检索, 得到检索结果.
3. 前端模块: 编写一个简单的页面, 展示搜索结果. 点击其中的搜索结果能跳转到对应的 Java API 文档页面.

关于分词

分词是搜索中的一个核心操作. 尤其是中文分词, 比较复杂(当然, 咱们此处暂不涉及中文分词)

我们可以使用现成的分词库 ansj.

官网网站: https://github.com/NLPchina/ansj_seg

使用的简单示例: https://blog.csdn.net/weixin_44112790/article/details/86756741

```
1 public class test.TestAnsj {
2     public static void main(String[] args) {
3         String str = "小明毕业于清华大学计算机专业," +
4                     "后来去蓝翔技校和新东方深造," +
5                     "擅长使用计算机控制挖掘机炒菜";
6         List<Term> terms = ToAnalysis.parse(str).getTerms();
7         for (Term term : terms) {
8             System.out.print(term.getName() + "/")
9         }
10    }
11 }
```

输出结果

```
1 小明/毕业/于/清华大学/计算机/专业/,/后来/去/蓝/翔/技校/和/新/东方/深造/,/擅长/使用/计算
   机/控制/挖掘机/炒菜/
```

注意: 当 anjs 对英文分词时, 会自动把单词转为小写.

公共模块: 创建 DocInfo 类

创建 common 包, 创建 DocInfo 类. 这个就作为一个文档对应的结构.

```
1 public class DocInfo {
2     private int docId;
3     private String title;
4     private String url;
5     private String content;
6 }
```

预处理模块

遍历目录下所有的文件, 并读取每个文件的内容, 把所有文件整理成一个行文本文件.

这个行文本文件, 每行对应一个 html .

每一行中有三列, 用 \3 分割. 分别是标题, url, 正文.

核心流程:

1. 枚举出文档目录下所有的 html 文件
2. 遍历每个文件, 把文件格式进行转换
3. 把最终结果写入到一个结果文件中.

在 searcher 包中创建 Parser 类. 这个类是一个单独的可执行程序.

```
1 public class Parser {
2     private static final String INPUT_PATH = "D:\\jdk1.8\\docs\\api";
3     private static final String OUTPUT_PATH = "D:/raw_data.txt";
4
5     public static void main(String[] args) {
6         try {
7             File resultFile = new File(OUTPUT_PATH);
8             FileWriter resultFileWriter = new FileWriter(resultFile);
9
10            // 1. 枚举出这个目录下的所有文件
11            ArrayList<File> fileList = new ArrayList<>();
12            enumFile(INPUT_PATH, fileList);
13            for (File f : fileList) {
14                System.out.println("converting " + f.getAbsolutePath());
15                // 2. 针对每个文件, 打开, 并读取内容, 进行转换
16                String line = convertLine(f);
17                // 3. 把转换结果写入到最终输出文档中
18                resultFileWriter.write(line);
19            }
20            resultFileWriter.close();
21            System.out.println("converting done!");
22        } catch (IOException e) {
23            e.printStackTrace();
24        }
25    }
26 }
```

```

27 // 递归完成目录枚举过程
28 private static void enumFile(String rootPath, ArrayList<File> fileList)
29 {
30     File rootFile = new File(rootPath);
31     File[] files = rootFile.listFiles();
32     for (File f : files) {
33         if (f.isDirectory()) {
34             enumFile(f.getAbsolutePath(), fileList);
35         } else if (f.getAbsolutePath().endsWith(".html")) {
36             fileList.add(f);
37         }
38     }
39 }
40 private static String convertLine(File f) {
41     // 1. 转换出标题
42     String title = convertTitle(f);
43     // 2. 转换出 url
44     String url = convertUrl(f);
45     // 3. 转换出正文(正文需要去除 html 标签)
46     String content = convertContent(f);
47     return title + "\3" + url + "\3" + content + "\n";
48 }
49
50 private static String convertTitle(File f) {
51     // 直接使用文件名作为标题
52     String name = f.getName();
53     return name.substring(0, name.length() - ".html".length());
54 }
55
56 private static String convertUrl(File f) {
57     // 这个 url 是指在线文档对应的链接。
58     // url 由两个部分构成。
59     // 第一部分是 https://docs.oracle.com/javase/8/docs/api
60     // 第二部分是 文件路径中 api 之后的部分。
61     String part1 = "https://docs.oracle.com/javase/8/docs/api";
62     String part2 = f.getAbsolutePath().substring(INPUT_PATH.length());
63     return part1 + part2;
64 }
65
66 private static String convertContent(File f) {
67     // 读取文件内容，并去除其中的 html 标签和换行
68     try {
69         FileReader fileReader = new FileReader(f);
70         // 是否当前读的字符是正文
71         boolean isContent = true;
72         StringBuilder output = new StringBuilder();
73         while (true) {
74             int ret = fileReader.read();
75             if (ret == -1) {
76                 break;
77             }
78             char c = (char)ret;
79             if (isContent) {
80                 if (c == '<') {
81                     isContent = false;
82                     continue;
83                 }

```

```

84         if (c == '\n' || c == '\r') {
85             c = ' ';
86         }
87         output.append(c);
88     } else {
89         if (c == '>') {
90             isContent = true;
91         }
92     }
93 }
94 fileReader.close();
95 return output.toString();
96 } catch (FileNotFoundException e) {
97     e.printStackTrace();
98 } catch (IOException e) {
99     e.printStackTrace();
100 }
101 return "";
102 }
103 }

```

索引模块

在 searcher 包中创建 Index 类. 这个类的目的是把刚才的 raw_data 文件做成索引.

Index 类中要包含正排索引和倒排索引, 结构如下

```

1  class Weight {
2      public String word;
3      public int docId;
4      public int weight;
5  }
6
7  public class Index {
8      // 正排索引, 下标对应 docId
9      private ArrayList<DocInfo> forwardIndex = new ArrayList<>();
10     // 倒排索引, key 是分词结果, value 是这个分词 term 对应的倒排拉链(包含一堆 docid)
11     private HashMap<String, ArrayList<Weight>> invertedIndex = new HashMap<>();
12 }

```

Index 类中要支持的接口如下:

```

1  // 根据 docId 查正排
2  public DocInfo getDocInfo(int docId) {
3      return forwardIndex.get(docId);
4  }
5
6  // 根据 分词结果 查倒排
7  public ArrayList<Weight> getInverted(String term) {
8      return invertedIndex.get(term);

```

```

9     }
10
11    // 构建索引, 根据 raw_data 文件在内存中构造索引结构
12    public void build(String inputPath) {
13        // TODO
14    }

```

接下来重点实现 Index.build 方法, 这也是实现索引模块最复杂的方法.

```

1    // 构建索引, 根据 raw_data 文件在内存中构造索引结构
2    public void build(String inputPath) {
3        System.out.println("build start! ");
4        long startTime = new Date().getTime();
5        try {
6            // 1. 按行读取文件内容(每行是一个 html)
7            FileReader fileReader = new FileReader(inputPath);
8            BufferedReader bufferedReader = new BufferedReader(fileReader);
9            while (true) {
10                String line = bufferedReader.readLine();
11                if (line == null) {
12                    break;
13                }
14                // 2. 构造 DocInfo 并更新正排索引
15                DocInfo docInfo = buildForward(line);
16                if (docInfo == null) {
17                    continue;
18                }
19                // 3. 构造 weight 并更新倒排索引
20                buildInverted(docInfo);
21                System.out.println("Build " + docInfo.getTitle() + " done!");
22            }
23            System.out.println("Build Done! " + (new Date().getTime() -
startTime));
24        } catch (IOException e) {
25            e.printStackTrace();
26        }
27    }

```

实现 buildForward (比较简单)

```

1    private DocInfo buildForward(String line) {
2        String[] tokens = line.split("\\3");
3        if (tokens.length != 3) {
4            System.err.println("tokens 长度不为 3: " + line);
5            return null;
6        }
7        DocInfo docInfo = new DocInfo();
8        docInfo.setDocId(forwardIndex.size());
9        docInfo.setTitle(tokens[0]);
10       docInfo.setUrl(tokens[1]);
11       docInfo.setContent(tokens[2]);
12       forwardIndex.add(docInfo);
13       return docInfo;

```

实现 buildInverted (比较复杂)

```

1  private void buildInverted(DocInfo docInfo) {
2      // 构造 weight 对象，并更新倒排索引。
3      // 此处 "权重" 简单粗暴的认为词出现的次数。
4      // weight = 10 * 这个词标题中出现的次数 + 1 * 这个词正文中出现的次数
5      // 核心流程：
6      // 1. 对标题进行分词
7      // 2. 遍历分词结果，统计标题中每个词出现的次数
8      // 3. 对正文进行分词
9      // 4. 遍历分词结果，统计正文中每个词出现的次数
10     // 5. 把以上内容都整理到一个 HashMap 中
11     // 6. 遍历 HashMap，可以得到 词 -> 权重 这样的映射关系，更新到倒排索引中
12     // 这个类用于辅助统计词出现的次数
13     class WordCnt {
14         public int titleCount;
15         public int contentCount;
16
17         public WordCnt(int titleCount, int contentCount) {
18             this.titleCount = titleCount;
19             this.contentCount = contentCount;
20         }
21     }
22     HashMap<String, WordCnt> wordCntMap = new HashMap<>();
23
24     // 1. 对标题进行分词
25     List<Term> terms = ToAnalysis.parse(docInfo.getTitle()).getTerms();
26     // 2. 遍历分词结果，统计标题中每个词出现的次数
27     for (Term term : terms) {
28         // 注意，此时的 word 已经被 anjs 转成小写了。
29         String word = term.getName();
30         WordCnt wordCnt = wordCntMap.get(word);
31         if (wordCnt == null) {
32             wordCntMap.put(word, new WordCnt(1, 0));
33         } else {
34             wordCnt.titleCount++;
35         }
36     }
37     // 3. 对正文进行分词
38     terms = ToAnalysis.parse(docInfo.getContent()).getTerms();
39     // 4. 统计正文中出现的词的个数
40     for (Term term : terms) {
41         // 注意，此时的 word 已经被 anjs 转成小写了。
42         String word = term.getName();
43         WordCnt wordCnt = wordCntMap.get(word);
44         if (wordCnt == null) {
45             wordCntMap.put(word, new WordCnt(0, 1));
46         } else {
47             wordCnt.contentCount++;
48         }
49     }
50     // 5. 把以上内容都整理到一个 HashMap 中
51     for (HashMap.Entry<String, WordCnt> entry : wordCntMap.entrySet()) {

```

```

52 // 6. 遍历 HashMap, 可以得到 词 -> 权重 这样的映射关系, 更新到倒排索引中
53 weight weight = new weight();
54 weight.setDocId(docInfo.getDocId());
55 weight.setWeight(entry.getValue().titleCount * 10
56                 + entry.getValue().contentCount);
57 weight.setWord(entry.getKey());
58 // 这个逻辑也可以使用 Map.putIfAbsent. 此处为了直观, 还是直接用 get 吧
59 ArrayList<weight> invertedList = invertedIndex.get(entry.getKey());
60 if (invertedList == null) {
61     invertedList = new ArrayList<>();
62     invertedIndex.put(entry.getKey(), invertedList);
63 }
64 invertedList.add(weight);
65 }
66 }

```

测试下索引的构建, 在 Index 中写一个 main 方法, 运行一下.

```

1 public static void main(String[] args) {
2     Index index = new Index();
3     index.build("d:/raw_data.txt");
4     ArrayList<weight> weights = index.getInverted("arraylist");
5     for (weight weight : weights) {
6         System.out.println("docId: " + weight.getDocId());
7         System.out.println("weight: " + weight.getWeight());
8         System.out.println("title: " +
9         index.getDocInfo(weight.getDocId()).getTitle());
10    }
11 }

```

挑两个 case 来看看 weight 计算的是否正确即可.

搜索模块

在 searcher 包中, 创建 DocSearcher 类. 实现搜索的核心流程; 创建 Result 类, 表示结果

```

1 // 不同于之前的 DocInfo, 这个 Result 是用来表示最终展示结果的.
2 public class Result {
3     private String title;
4     private String url;
5     // desc 是正文的一段摘要
6     private String desc;
7
8     // getter setter 略
9 }

```

DocSearcher 需要实现的方法:

1. 构造方法, 要构建索引
2. search 方法, 搜索的核心流程


```

1 public class DocSearcher {
2     private Index index = new Index();
3
4     public DocSearcher() {
5         index.build("d:/raw_data.txt");
6     }
7
8     // 根据查询词，进行搜索，得到搜索结果集合。
9     // 结果集中包含若干条记录，每个记录中包含搜索结果的标题，描述，url
10    public List<Result> search(String query) {
11
12    }
13 }

```

实现 search 方法

```

1 // 根据查询词，进行搜索，得到搜索结果集合。
2 // 结果集中包含若干条记录，每个记录中包含搜索结果的标题，描述，url
3 // 1. [分词] 对查询词进行分词
4 // 2. [触发] 对每个分词结果查找倒排索引，得到一个倒排拉链
5 // 3. [排序] 针对结果集合进行排序，按权重降序排序即可
6 // 4. [返回] 构造返回结果
7 public List<Result> search(String query) {
8     // 1. [分词] 对查询词进行分词
9     List<Term> terms = ToAnalysis.parse(query).getTerms();
10    // 2. [触发] 对每个分词结果查找倒排索引，得到一个倒排拉链
11    ArrayList<Weight> allTokenResult = new ArrayList<>();
12    for (Term term : terms) {
13        String word = term.getName();
14        ArrayList<Weight> invertedList = index.getInverted(word);
15        // 倒排拉链可能找不到，比如这个词根本就在索引中不存在。
16        if (invertedList == null) {
17            continue;
18        }
19        allTokenResult.addAll(invertedList);
20    }
21    // 3. [排序] 针对结果集合进行排序，按权重降序排序即可
22    allTokenResult.sort(new Comparator<Weight>() {
23        @Override
24        public int compare(Weight o1, Weight o2) {
25            // 如果升序就是 o1 - o2，降序是 o2 - o1
26            return o2.getWeight() - o1.getWeight();
27        }
28    });
29    // 4. [返回] 构造返回结果
30    ArrayList<Result> results = new ArrayList<>();
31    for (Weight weight : allTokenResult) {
32        DocInfo docInfo = index.getDocInfo(weight.getDocId());
33        Result result = new Result();
34        result.setTitle(docInfo.getTitle());
35        result.setUrl(docInfo.getUrl());
36        result.setDesc(GenDesc(docInfo.getContent(), weight.getword()));
37        results.add(result);
38    }
39    return results;
40 }

```

实现下 GenDesc 方法, 生成一下描述信息

```
1 private String GenDesc(String content, String word) {
2     // 先找一下 word 在 content 中第一次出现的位置
3     int firstPos = content.indexOf(word);
4     if (firstPos == -1) {
5         return "";
6     }
7     // 直接截取 firstPos 周围的文本
8     // 从 firstPos 往前找 60 个字符作为描述开始,
9     // 然后从描述开始位置往后找 160 个字符作为整个描述
10    // 注意: 此处的 60, 160 都是拍脑门出来的
11    int descBeg = firstPos < 60 ? 0 : firstPos - 60;
12    if (descBeg + 160 > content.length()) {
13        return content.substring(descBeg);
14    }
15    // 正文长度充足, 在最后加上 ...
16    return content.substring(descBeg, descBeg + 160) + "...";
17 }
```

写个 main 方法测试一下 DocSearcher 类

```
1 public static void main(String[] args) {
2     DocSearcher docSearcher = new DocSearcher();
3     List<Result> results = docSearcher.search("ArrayList");
4     for (Result result : results) {
5         System.out.println(result);
6     }
7 }
```

实现 Servlet

接口设计

```
1 请求:
2  /search?query=ArrayList
3
4 响应:
5  {
6      ok: true,
7      data: [
8          {
9              title: "title1",
10             url: "url1",
11             desc: "desc1",
12         },
13         {
14             title: "title2",
15             url: "url2",
16             desc: "desc2",
17         },
18     ]
19 }
```

关于 Json

Json 是一种常见是数据格式组织方式. 源于 JavaScript, 是一种键值对风格的数据格式.

Java 中可以使用 Gson 库来完成 Json 的解析和构造.

在 Maven 中新增 Gson 的依赖

```
1 <dependency>
2   <groupId>com.google.code.gson</groupId>
3   <artifactId>gson</artifactId>
4   <version>2.8.2</version>
5 </dependency>
```

简单示例 (创建一个 TestGson 类)

```
1 public class TestGson {
2     public static void main(String[] args) {
3         HashMap<String, Object> data = new HashMap<>();
4         data.put("name", "曹操");
5         data.put("skill1", "剑气");
6         data.put("skill2", "三段跳");
7         data.put("skill3", "吸血加攻击");
8         data.put("skill4", "释放技能加攻速");
9         Gson gson = new GsonBuilder().create();
10        String jsonData = gson.toJson(data);
11        System.out.println(jsonData);
12    }
13 }
```

实现 DocSearcherServlet

创建 servlet 包, 创建 DocSearcherServlet 类, 继承自 HttpServlet . 此处需要使用 Gson 作为 Json 解析库.

```
1 public class DocSearcherServlet extends HttpServlet {
2     private DocSearcher docSearcher = new DocSearcher();
3
4     @Override
5     protected void doGet(HttpServletRequest req, HttpServletResponse resp)
6     throws ServletException, IOException {
7         resp.setContentType("application/json; charset=utf-8");
8         HashMap<String, Object> data = new HashMap<>();
9         Gson gson = new GsonBuilder().create();
10        String query = req.getParameter("query");
11        if (query == null || query.equals("")) {
12            data.put("ok", false);
13            data.put("reason", "query is empty");
14            String respData = gson.toJson(data);
15            resp.getWriter().write(respData);
16            return;
17        }
18    }
19 }
```

```

17         List<Result> results = docSearcher.search(query);
18         data.put("ok", true);
19         data.put("data", results);
20         String respData = gson.toJson(data);
21         resp.getWriter().write(respData);
22     }
23 }

```

修改 web.xml

```

1 <servlet>
2     <servlet-name>DocSearcherServlet</servlet-name>
3     <servlet-class>servlet.DocSearcherServlet</servlet-class>
4 </servlet>
5 <servlet-mapping>
6     <servlet-name>DocSearcherServlet</servlet-name>
7     <url-pattern>/search</url-pattern>
8 </servlet-mapping>

```

此时发现一个问题. 索引制作的过程是在第一次收到请求的时候才开始的. 这显然不科学.

如果想让索引能在服务启动的时候就进行制作, 也就需要让 Servlet 在 Tomcat 启动时就被创建和加载.

此时可以在 web.xml 中新增一个 load-on-startup 选项即可.

```

1 <servlet>
2     <servlet-name>DocSearcherServlet</servlet-name>
3     <servlet-class>servlet.DocSearcherServlet</servlet-class>
4     <load-on-startup>1</load-on-startup>
5 </servlet>

```

其中这个值 如果填成 ≤ 0 的, 就表示第一次收到请求再加载. 如果 > 0 的, 就表示服务器启动就加载. 其中数字越小, 就越先加载.

参考 <https://blog.csdn.net/lvxiangan/article/details/80582487>

实现前端页面

前端页面结构比较简单, 只要包含一个输入框和一个按钮即可. 另外要对返回的结果进行一个比较合理的组织.

创建一个 index.html. 这个就不讲了吧, 直接发给大家.

```

1 <html>
2 <head>
3     <!-- Bootstrap 文档: https://v3.bootcss.com/css/ -->
4     <!-- vue 文档: https://cn.vuejs.org/v2/guide/ -->
5     <!-- Required meta tags -->
6     <meta charset="utf-8">
7     <meta name="viewport" content="width=device-width, initial-scale=1,
8         shrink-to-fit=no">
9     <!-- Bootstrap CSS -->

```

```

10     <link rel="stylesheet"
href="https://cdn.jsdelivr.net/npm/bootstrap@3.3.7/dist/css/bootstrap.min.cs
s" integrity="sha384-
BVYiISIFeK1dGmJRAkycuHAHRg32OmUcww7on3RYdg4Va+PmSTsz/K68vbdEjh4u"
crossorigin="anonymous">
11
12     <title>Java API 搜索</title>
13     <style>
14     #app {
15         margin-left:50px;
16         margin-right:50px;
17     }
18     div button {
19         width:100%;
20     }
21     .row {
22         padding-top: 10px;
23     }
24     .col-md-5, .col-md-1 {
25         padding-left:2;
26         padding-right:2;
27     }
28     .title {
29         font-size: 22px;
30     }
31     .desc {
32         font-size: 18px;
33     }
34     .url {
35         font-size: 18px;
36         color: green;
37     }
38     </style>
39 </head>
40 <body>
41 <div id="app">
42     <div class="row">我是 logo</div>
43     <div class="row">
44         <div class="col-md-5">
45             <input type="text" class="form-control" placeholder="请输入关键字"
v-model="query">
46             </div>
47             <div class="col-md-1">
48                 <button class="btn btn-success" v-on:click="search()">搜索
</button>
49             </div>
50         </div>
51         <div class="row" v-for="result in results">
52             <!--用来存放结果-->
53             <div class="title"><a v-bind:href="result.url">{{result.title}}</a>
</div>
54             <div class="desc">{{result.desc}}</div>
55             <div class="url">{{result.url}}</div>
56         </div>
57     </div>
58 </body>
59 <script src="https://apps.bdimg.com/libs/jquery/2.1.4/jquery.min.js">
</script>

```

```
60 <script
    src="https://cdn.jsdelivr.net/npm/bootstrap@3.3.7/dist/js/bootstrap.min.js"
    integrity="sha384-
    Tc5IQib027qvyjSMfHjOMaLkfuwVxZxUPnCJA7l2mCWNIpG9mGCD8WGNIcPD7Txa"
    crossorigin="anonymous"></script>
61 <script src="https://cdn.jsdelivr.net/npm/vue"></script>
62 <script>
63     var vm = new Vue({
64         el: "#app",
65         data: {
66             query: "",
67             results: [ ]
68         },
69         methods: {
70             search() {
71                 $.ajax({
72                     url: "/java_doc_searcher/search?query=" + this.query,
73                     type: "get",
74                     context: this,
75                     success: function(respData, status) {
76                         this.results = respData.data;
77                     }
78                 })
79             },
80         }
81     })
82 </script>
83 </html>
```

后续改进

1. 数据量更大
2. 请求量更大
3. 业务更复杂(相关性策略)