

目录

| | | |
|----------|---------------------------|-----------|
| 1 | 实验介绍 | 2 |
| 1.1 | 实验内容 | 2 |
| 1.1.1 | 实验六 | 2 |
| 1.1.2 | 实验七 | 2 |
| 1.2 | 实验环境 | 2 |
| 1.3 | 所需技能 | 2 |
| 2 | 实验过程 | 3 |
| 2.1 | 准备实验环境 | 3 |
| 2.2 | 后端构建 | 3 |
| 2.2.1 | 设置 JVM 与 Lucene | 3 |
| 2.2.2 | 设置路由 | 4 |
| 2.2.3 | 网页搜索 | 5 |
| 2.2.4 | 图片搜索 | 5 |
| 2.3 | 前端构建 | 6 |
| 2.3.1 | Index 页面 | 6 |
| 2.3.2 | 网页搜索结果页面 | 8 |
| 2.3.3 | 图片搜索结果页面 | 10 |
| 3 | 总结与分析 | 12 |

1 实验介绍

1.1 实验内容

1.1.1 实验六

使用 `web.py`, 结合前面学习的 `HTML`, `Lucene`, 中文分词等知识点, 根据上次实验爬取的网页, 建立一个简单的搜索引擎. 搜索结果中要求包含: 标题, 超链接, 关键词上下文以及网址.

1.1.2 实验七

制作一个图片加文字的搜索引擎, 作为中期整合. 即, 在上次的基础上, 加入图片搜索, 使用 `css` 制定样式.

自主调节 `css` 中的各项属性参数与 `div` 布局, 不存在标准答案, 只需根据个人喜好决定页面表示形式和显示哪些细节.

1.2 实验环境

操作系统: Ubuntu 18.04 LTS

Python 版本: 3.7

PyLucene 版本: 7.4.0

`web.py` 版本: 0.40.dev1

1.3 所需技能

实验六主要涉及使用 `Python` 框架搭建简单的网站, 用到了 `web.py`, 需要学习其基本的工作方式与 `API`. 同时, 需要将之前实验中实现的多个功能模块, 如中文分词与 `Lucene` 索引, 进行整合使用.

实验七主要在于整合前期工作, 同时合理使用 `HTML` 和 `CSS` 布局来使得网页整齐美观. 这一过程需要综合考虑多个模块, 多次尝试调整, 以达到满意的效果.

2 实验过程

2.1 准备实验环境

相较于之前的实验, 本次实验仅需额外安装 web.py 即可准备好实验环境. 但需要说明的是, 由于 web.py 目前处于欠维护状态, 其对 Python 3 的支持并非十分完善. 如果通过 pip 安装, 需手动指定版本, 否则所安装的版本将仅对 Python 2 可用:

```
1 pip install web.py==0.40-dev1
```

更进一步的, 由于我所使用的 Python 版本为最新的 3.7, 这样安装后的 web.py 依旧存在问题: 框架报错而无法运行, 静态文件访问出错等. 这些问题主要是由于 Python 3.7 的一些变动导致的, 为解决这些问题, 需要手动从 Github 上下载最新的 web.py 源码安装, 如下:

```
1 git clone git://github.com/webpy/webpy.git
2 ln -s `pwd`/webpy/web .
```

2.2 后端构建

2.2.1 设置 JVM 与 Lucene

由于 PyLucene 仅仅是 Lucene 的 Python 接口, Lucene 本身是运行于 Java 虚拟机环境 (JVM) 中, 需要在 Python 程序中启动 JVM, 并且在每次执行搜索时调用. 对于 web.py, 在主程序中启动 Web 应用时, 会将当前所有的全局变量, 即 `globals()`, 打包传递给 web.py 程序, 并且在之后的每次响应时从中获取所需变量.

那么, 我们需要在主程序中启动 JVM, 作为全局变量供之后使用. 同时, 为了避免每次搜索时创建 Lucene 搜索器的开销, 可以在程序启动时就创建好搜索器, 同样作为全局变量传递. 如下:

```
1 if __name__ == '__main__':
2     vm_env = lucene.initVM(vmargs=['-Djava.awt.headless=true'])
3     base_dir = os.path.dirname(os.path.abspath(sys.argv[0]))
4     directory = {
5         "html": SimpleFSDirectory(Paths.get(os.path.join(base_dir,
6             ↪ HTML_INDEX_DIR))),
7         "image": SimpleFSDirectory(Paths.get(os.path.join(base_dir,
8             ↪ IMAGE_INDEX_DIR))),
9     }
10    searcher = {
11        "html": IndexSearcher(DirectoryReader.open(directory["html"])),
12        "image": IndexSearcher(DirectoryReader.open(directory["image"])),
13    }
14    analyzer = StandardAnalyzer()
15    app = web.application(urls, globals(), autoreload=False) #attention here
16    app.run()
```

注意在上面的代码中, 需要禁用 web.py 的自动重载, 否则会导致上述代码失效. 可能的原因在于此种情况下主程序将不再是我们编写的该文件, 而是 web.py 自身, 则 `__name__ != '__main__'`, 而是实际的文件名.

在初始化 JVM 后, 只需在每次调用时:

```
1 vm_env.attachCurrentThread()
```

2.2.2 设置路由

由于需要区分网页搜索和图片搜索 (当然也可以不区分, 只需在搜索时额外传递一个类型参数), 则对二者的 Index 页面各自分配一个处理类 (原因在于这两个类都非常简洁), 但对于搜索结果页面, 我选择共用一个类, 通过 URL 来区分类别 (原因在于处理网页和图片的搜索请求流程基本相同, 可以通过共用类来减少冗余代码). 如下:

```
1 urls = (  
2     "/", "Index",  
3     "/image", "IndexImage",  
4     "/search/(.+)", "Search",  
5 )  
  
6 class Index:  
7     def GET(self):  
8         return render.index()  
  
9 class IndexImage:  
10     def GET(self):  
11         return render.index_image()  
  
12 class Search:  
13     def GET(self, search_type="html"):  
14         search_data = web.input()  
15         query_string = search_data.get("s", "")  
16         if search_type not in ("html", "image"):  
17             search_type = "html"  
18         if not query_string:  
19             return render.index_image() if search_type == "image" else  
20                 ↪ render.index()  
21         search_result = search_function[search_type](query_string)  
22         return result_template[search_type](query_string, search_result)
```

注意以上代码的最后两行, 通过表驱动法避免了 if-else 语句的泛滥, 对于目前仅有的两种搜索类型可能优势并不明显, 但可以方便地添加新的搜索类型而无需修改此处逻辑, 有利于日后的功能增加.

2.2.3 网页搜索

对网页的搜索,在往次实验的基础上,只需对原有程序作出一些修改,即可基本实现,故不再赘述.主要的问题在于,如何实现关键词上下文以及关键词高亮.

让我们稍做分析:Lucene 所创建的索引中,并未保存网页的完整内容 (`t.setStored(False)`),而是对其进行提取处理后保存了索引信息(如倒排索引).那么,自然是难以通过 Lucene 直接对关键词进行高亮处理,除非重新创建索引并且保存网页的内容信息.

那么,我采用了十分朴素的高亮方式:打开检索结果对应的本地 HTML 文件,对其内容进行去 HTML 标签及分词处理,得到一个词汇列表,然后将其与输入的搜索字符串分词后的词汇列表对比,找到第一个重合项,然后将此项的前后一定数量的词汇合并为所展示的内容摘要.如下:

```
1 with open(doc.get("path"), mode="r", encoding="utf8") as file:
2     content = file.read()
3     html2text = HTML2Text()
4     html2text.ignore_links = True
5     html2text.ignore_images = True
6     content = html2text.handle(content)
7     cutted_content = jieba.cut(content)
8     flag = False
9     if cutted_query:
10         word_num, cnt = 20, 0
11         for x in cutted_content:
12             if not x.strip():
13                 continue
14             if not flag and x in cutted_query:
15                 flag = True
16                 content = ""
17             if flag and cnt < word_num:
18                 cnt += 1
19                 content += (x if x not in cutted_query else
20                             "<span class='highlight'>{0}</span>".format(x)
21                             )
22             elif cnt >= word_num:
23                 break
24     single_result["content"] = content if flag else content[:100]
```

这样的做法可能有些粗糙,但效果尚可接受,在后续的实验,必要时需要对其做出修改.

2.2.4 图片搜索

类似的,图片搜索基本上可以利用之前实验的代码实现,仅需做出一些简单的修改.

```
1 def search_image(query_string, result_num=10):
2     vm_env.attachCurrentThread()
```

```

3      cutted = [x for x in jieba.cut_for_search(query_string) if x.strip()]
4      command = " ".join(cutted)
5      query = QueryParser("content", analyzer).parse(command)
6      scoreDocs = searcher["image"].search(query, result_num).scoreDocs
7      result = list()
8      for num, scoreDoc in enumerate(scoreDocs):
9          doc = searcher["image"].doc(scoreDoc.doc)
10         single_result = {
11             "url": doc.get("url"),
12             "description": doc.get("raw_content"),
13         }
14         result.append(single_result)
15     return result

```

2.3 前端构建

2.3.1 Index 页面

网页搜索的 Index 页面主体部分都是一个文本输入框加上表单提交的按钮. 为了使网页整齐, 使用 CSS+DIV 进行布局:

```

1  <!DOCTYPE html>
2  <html lang="zh">
3  <head>
4      <meta charset="UTF-8">
5      <link rel="stylesheet" href="/static/css/main.css">
6      <title> 网页搜索 </title>
7  </head>
8  <body>
9      <div id="magicWrapper">
10         <div id="superCenter" align="center">
11             <h1 style="font-size: 60px;" >
12                 <span>The</span><span class="beacon">Beacon</span>
13             </h1>
14             <div id="searchType">
15                 <span class="active"><a href="/"> 网页 </a></span>
16                 <span class="type"><a href="/image"> 图片 </a></span>
17             </div>
18             <div id="homeSearch" align="center">
19                 <form action="/search/html">
20                     <input type="text" id="homeInput" name="s" placeholder=" 搜点
                        ↳ 什么吧">

```

```
21         <input type="submit" id="homeSubmit" value=" 搜索">
22         <br>
23     </form>
24 </div>
25 </div>
26 </div>
27 </body>
28 </html>
```

具体的 CSS 属性不在此贴出, 最终实现的效果如下:



图 1: 网页搜索 Index 页面

而对于图片搜索, 其布局基本类似, 为了体现图片搜索的特点而有所区分, 对其设置了不同的背景, 效果如下:

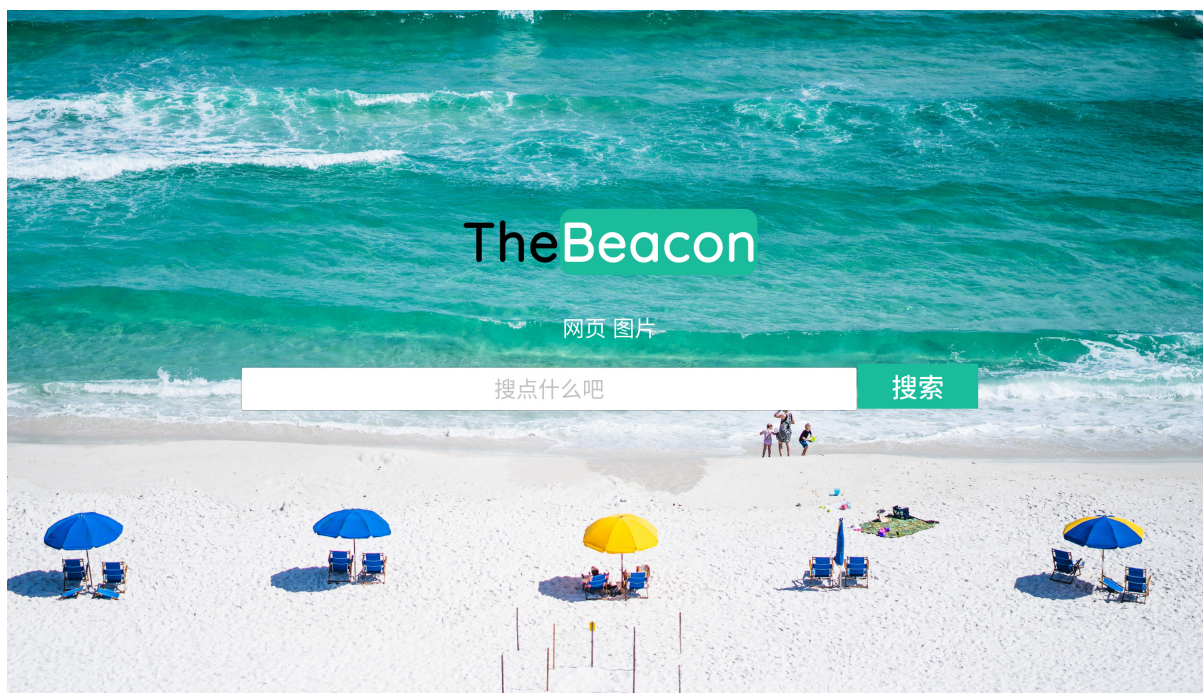


图 2: 图片搜索 Index 页面

2.3.2 网页搜索结果页面

搜索结果页面模板接受 Python 传来的参数: `query_string, search_result`, 类型分别为字符串以及字典组成的列表, 以此生成搜索页面. 则可布局为多个 DIV, 每个 DIV 为一条搜索结果. 同时, 为了方便使用, 可以设置一个始终悬浮于页面最上的输入框, 用于进行新的搜索.

```

1  $def with (query_string, search_result)
2  <!DOCTYPE html>
3  <html lang="zh">
4  <head>
5      <meta charset="UTF-8">
6      <link rel="stylesheet" href="/static/css/main.css">
7      <title>${query_string} 网页搜索结果 </title>
8  </head>
9  <body>
10     <div class = "floatingBar">
11         <form action="/search/html">
12             <input type="text" name="s" class="floatingInput" value=
13                 ↪ "${query_string}">
14             <input type="submit" class="floatingSubmit" value=" 搜索">
15         </form>
16     </div>
17     <div class = "resultContainer">
18         $for item in search_result:
19             <div class = "item">

```



```

19         <div class = "title">
20             <a href="$item['url']">$item["title"]</a>
21         </div>
22         <div class = "content">$.item["content"]</div>
23         <div class = "url">
24             <a href="$item['url']">$item["url"][:80]</a>
25         </div>
26     </div>
27 </div>
28 </body>
29 </html>

```

用到的主要 CSS 属性如下:

```

1  div.item{
2      margin-bottom: 10px;
3      padding: 20px 15px 15px;
4      border-radius: 5px;
5      background-color: #fff;
6      box-sizing: border-box;
7      box-shadow: 0 0 20px 2px rgba(0, 0, 0, .1);
8      -webkit-box-shadow: 0 0 20px 2px rgba(0, 0, 0, .1);
9      -moz-box-shadow: 0 0 20px 2px rgba(0, 0, 0, .1);
10 }

11 .floatingBar{
12     position: fixed;
13     top: 10px;
14     width:80%;
15     margin-left: 10%;
16     margin-right: 10%;
17 }

```

最终效果如图所示:



图 4: 图片搜索结果

3 总结与分析

本次实验将之前实验内容进行了整合,在这一过程中,我发现了以前所写代码的一些缺陷,诸如不能灵活配置参数,过度依赖上下文等,这些都阻碍着对其的整合.因此,在今后的编程中,更需要注意模块化的设计思想.

此外,在编写 HTML 和 CSS 代码时,由于没有正式地学习过前端知识,常常遇到一些难以处理的细节问题,需要反复尝试并查找各种资料,耗费了不少时间,但也加深了对其的理解与掌握.