

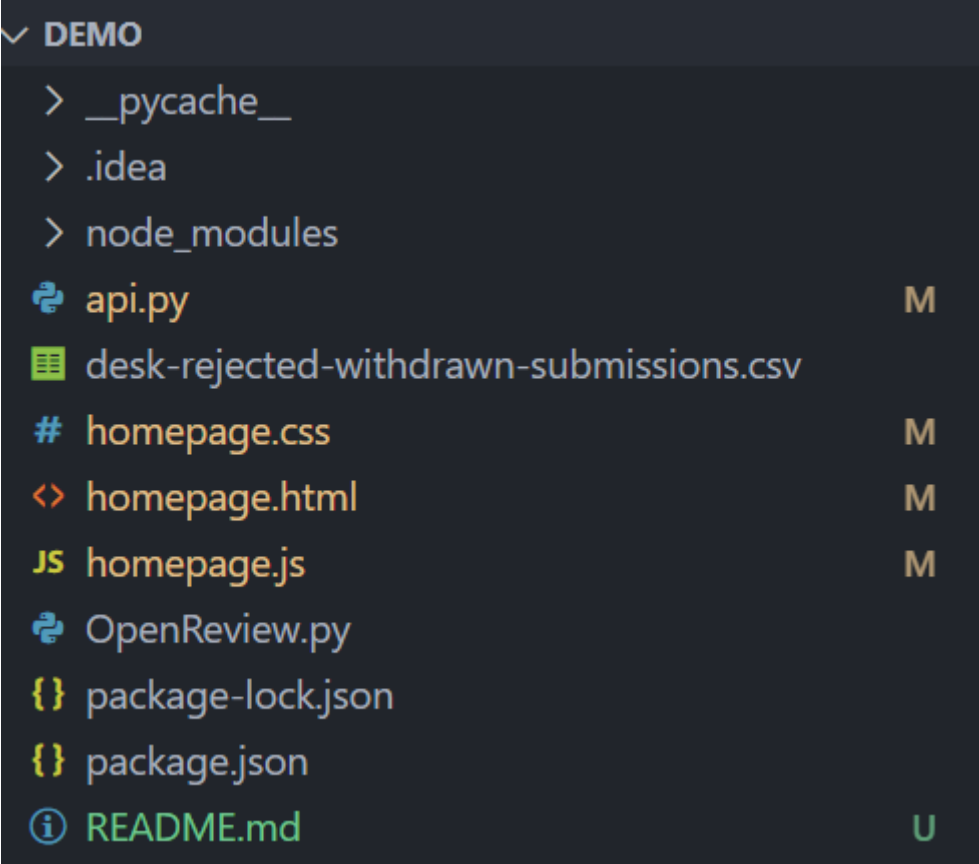
期末大作业实验报告

Python第四小组

1. 作业要求

```
1 # 题目三：爬虫
2
3 ## 题A：
4
5 具体描述：
6
7 * 从以下网站爬取`ICLR2023`的论文数据，输出所有论文的pdf链接；
8 * 爬取每篇论文的标题、作者和论文摘要，用这些信息制作一个新的网页（离线网页即可）；
9 * 针对论文的关键词绘制词云。
10
11 \[目标链接\](https://openreview.net/group?id=ICLR.cc/2023/Conference)
```

2. 项目结构



项目主要包括三个部分：

- 前端页面部分：
 - homepage.html
 - homepage.js
 - homepage.css
- api接口部分：
 - api.py
- 爬虫部分：
 - OpenReview.py

其中，针对关键词绘制词云，以及扩展的搜索功能都在api接口部分完成，爬取论文数据，并清洗数据，获得每篇论文的标题、作者和论文摘要以及下载链接，这一部分在爬虫部分完成，而调用api接口，进行页面渲染，这部分在前端页面部分完成。在开发期间，我们使用git进行版本管理，所有源码可以在[这里](#)找到。

3. 项目细节

3.1. 爬虫实现

我们共导入了以下包，其中第三行导入了线程池用来实现多线程爬虫，第四行用于处理url：

```
1 import requests
2 import pandas as pd
3 from concurrent.futures import ThreadPoolExecutor
4 from urllib.parse import urlencode, urlunparse, urlparse, parse_qs
```

首先，规定了实体格式和目标网页每页文章数量：

```
1 pd_dic = {"titles": [], "authors": [], "keywords": [], "pdfs": []}
2 limit = 25
```

接下来定义了须爬取的五大类型对应的基础url：

```
1 # 所有ICLR前%5的论文
2 url_0 = "https://api.openreview.net/notes?
content.venue=ICLR+2023+notable+top+5%25&details=replyCount&offset=25&limit=25&invitation=ICLR.cc%2F2023%2FConference%2F-
%2FBlind_Submission"
3 temp_resp_0 = requests.get(url_0)
4 count_0 = temp_resp_0.json()["count"]
5 # 所有ICLR前%25的论文
6 url_1 = "https://api.openreview.net/notes?
content.venue=ICLR+2023+notable+top+25%25&details=replyCount&offset=0&limit=25&invitation=ICLR.cc%2F2023%2FConference%2F-
%2FBlind_Submission"
7 temp_resp_1 = requests.get(url_1)
8 count_1 = temp_resp_1.json()["count"]
9 # 所有发表的论文
10 url_2 = "https://api.openreview.net/notes?
content.venue=ICLR+2023+poster&details=replyCount&offset=0&limit=25&invitation=ICLR.cc%2F2023%2FConference%2F-%2FBlind_Submission"
11 temp_resp_2 = requests.get(url_2)
12 count_2 = temp_resp_2.json()["count"]
13 # 所有提交的论文
14 url_3 = "https://api.openreview.net/notes?
content.venue=Submitted+to+ICLR+2023&details=replyCount&offset=0&limit=25&invitation=ICLR.cc%2F2023%2FConference%2F-
%2FBlind_Submission"
15 temp_resp_3 = requests.get(url_3)
16 count_3 = temp_resp_3.json()["count"]
17 # desk-rejected-withdrawn-submissions
18 url_4 = "https://api.openreview.net/notes?
details=replyCount%2Cinvitation%2Coriginal&offset=0&limit=25&invitation=ICLR.cc%2F2023%2FConference%2F-%2FWithdrawn_Submission"
19 temp_resp_4 = requests.get(url_4)
20 count_4 = temp_resp_4.json()["count"]
```

然后实现了爬取每一页，并将数据存入dict字典的函数：

```
1 def get_one_page(_url, _dict):
2     try:
3         resp = requests.get(_url)
4         resp.raise_for_status() # 如果响应状态码不是200，就主动抛出异常
5     except requests.RequestException as e:
6         print(f"请求{_url}时发生错误: {e}")
7         return
8
9     notes = resp.json()["notes"]
10    for note in notes:
11        id = note["id"]
12        content = note["content"]
13        pdf = f"https://openreview.net/pdf?id={id}"
14        _dict["titles"].append(content["title"])
15        _dict["authors"].append(content["authors"])
16        _dict["keywords"].append(content["keywords"])
17        _dict["pdfs"].append(pdf)
18        print("ok")
```

接下来是运用多线程，对于需要爬取的每一页，提交一个任务，用最多50个线程并行爬取。

```
1 def Thread_Method(__url, _count, _dict = pd_dic):
2     with ThreadPoolExecutor(50) as t:
3         for i in range(0, _count, limit):
4             # 解析URL
5             url_parts = list(urlparse(__url))
6             # 解析查询参数
7             query = dict(parse_qs(url_parts[4]))
8             # 更新offset参数
9             query.update({"offset": str(i)})
10            # 重新生成查询参数字符串
11
12            url_parts[4] = urlencode(query)
```

```
12         # 重新生成URL
13         _url = urlunparse(url_parts)
14         t.submit(get_one_page, _url, _dict)
```

最后，封装为核心功能，并测试文件功能是否正常。

```
1 def func(mode):
2     dic = {"titles": [], "authors": [], "keywords": [], "pdfs": []}
3     url = globals()['url_' + str(mode)]
4     count = globals()['count_' + str(mode)]
5     Thread_Method(url, count, dic)
6     return dic
7
8 if __name__ == "__main__":
9     pd_dic = func(4)
10    db = pd.DataFrame(pd_dic)
11    db.to_csv("desk-rejected-withdrawn-submissions.csv")
```

测试后，成功获取了csv文件 `desk-rejected-withdrawn-submissions.csv`，文件内容正确。

3.2. api接口实现

我们共导入了以下包，包括刚刚爬虫部分封装的func()函数。

```
1 from flask import Flask, jsonify, request, send_file
2 from flask_cors import cross_origin
3 from OpenReview import func
4 from wordcloud import WordCloud
5 from io import BytesIO
```

我们采用Flask框架作为服务端架构，首先创建实例app，并初始化data为None。

```
1 app = Flask(__name__)
2 data = None
```

在路由 `/api/data` 处封装api，并允许跨源访问。

```
1 @app.route('/api/data', methods=['GET'])
2 @cross_origin()
3 def get_data():
4     mode = int(request.args.get('mode'))
5
6     global data
7     data = func(mode)
8
9     return jsonify(data) # 返回指定页的数据
```

封装词云生成api，先把所有文章的关键词展平，再调用WordCloud库绘制词云，并使用flask库的send_file()方法将图片发送到前端。

```
1 @app.route('/api/keywords-wordcloud', methods=['GET'])
2 @cross_origin()
3 def generate_keywords_wordcloud():
4     global data
5     flat_list = []
6     for item in data['keywords']:
7         for words in item:
8             flat_list.append(words)
9     keywords_text = ",".join(flat_list)
10
11    # 生成词云
12    wordcloud = WordCloud(width=800, height=400, background_color='white').generate(keywords_text)
13
14    # 保存词云为图片
15    img_buffer = BytesIO()
16    wordcloud.to_image().save(img_buffer, format='PNG')
17    img_buffer.seek(0)
18
19    return send_file(img_buffer, mimetype='image/png')
```

拓展了搜索功能，时间有限没有做基于正则表达式的部分匹配，仅支持完全匹配，且仅可搜索当前页面上显示的文章。

```
1 @app.route('/api/search', methods=['GET'])
2
3 @cross_origin()
```

```
3 def search_data():
4     keyword = request.args.get('keyword')
5     # 在全局变量 data 中搜索关键词并返回结果
6     result = search_in_data(keyword)
7     return jsonify(result)
8
9 def search_in_data(keyword):
10     global data
11
12     # 在 data 中搜索关键词，找到匹配的索引
13     result = []
14     for i, sublist in enumerate(data['keywords']):
15         if keyword in sublist:
16             result.append(i)
17     return result
```

启动Flask实例，监听5000端口。

```
1 if __name__ == '__main__':
2     app.run(port=5000)
```

3.3. 前端代码实现

前端代码共包含三个部分，其中页面的标签栏和搜索栏是静态网页，而词云与论文部分是动态加载的。

3.3.1. 静态页面结构

```
1 <!DOCTYPE html>
2 <html lang="en">
3 <head>
4     <meta charset="UTF-8">
5     <title>论文聚合平台</title>
6     <link rel="stylesheet" type="text/css" href="homepage.css">
7 </head>
8 <body onload="initPage()">
9     <header class="header">
10         
11         <a href="homepage.html" class="home">论文聚合平台</a>
12         <div class="search">
13             <input id="searchInput" type="text" placeholder="搜索">
14             <button onclick="search()">
15                 <svg width="18" height="18" viewBox="0 0 24 24" fill="currentColor">
16                     <g fill-rule="evenodd" clip-rule="evenodd">
17                         <path d="M11.5 18.389c3.875 0 7-3.118 7-6.945 0-3.826-3.125-6.944-7-6.944s-7 3.118-7 6.944 3.125 6.945 7
18 6.945Zm0 1.5c4.694 0 8.5-3.78 8.5-8.445C20 6.781 16.194 3 11.5 3S3 6.78 3 11.444c0 4.664 3.806 8.445 8.5 8.445Z"></path>
19                         <path d="M16.47 16.97a.75.75 0 0 1 1.06 0l3.5 3.5a.75.75 0 0 1 1-1.06 1.061-3.5-3.5a.75.75 0 0 1 0-1.06Z">
20                     </path>
21                 </g>
22             </svg>
23         </button>
24     </div>
25 </header>
26 <h1>The Eleventh International Conference on Learning Representations</h1>
27 <div style="box-shadow: 0 1px 3px hsla(0,0%,7%,.1);display: flex;flex-direction: column;align-items: center">
28     <div class="tabCardContainer">
29         <button id="tab_0" class="tabCard" onclick="clickTab(0)" style="color: #121212">Notable-top-5%</button>
30         <button id="tab_1" class="tabCard" onclick="clickTab(1)" style="color: #121212">Notable-top-25%</button>
31         <button id="tab_2" class="tabCard" onclick="clickTab(2)" style="color: #121212">Poster</button>
32         <button id="tab_3" class="tabCard" onclick="clickTab(3)" style="color: #121212">Submitted</button>
33         <button id="tab_4" class="tabCard" onclick="clickTab(4)" style="color: #121212">Desk Rejected/Withdrawn
34 Submissions</button>
35     </div>
36     <div id="loading"><p>加载中...</p></div>
37     <div class="img-container"></div>
38     <div class="article-container"></div>
39 </div>
40 <script src="homepage.js"></script>
41 </body>
42 </html>
```

在搜索框和文章容器部分，参考了知乎的样式表，包括svg搜索图标等。

3.3.2. 前端代码

鉴于本次课程为Python，故不详细展开介绍js部分。

```
1  const loading = document.getElementById('loading');
2  let data;
3
4  function clickTab(x) {
5      let elem = document.getElementById("tab_" + x)
6      if (elem.style.color === "rgb(18, 18, 18)")
7          elem.style.color = "#056de8";
8      for (let i = 0; i < 5; i++) {
9          let temp = document.getElementById("tab_" + i)
10         if (i !== x && temp.style.color === "rgb(5, 109, 232)")
11             temp.style.color = "#121212";
12     }
13     loading.style.display = 'flex';
14     getData(x).then((d) => {
15         data = d;
16         console.log(data);
17         displayPage();
18     });
19 }
20
21 async function getData(mode) {
22     try {
23         const response = await fetch('http://localhost:5000/api/data?mode=' + mode);
24         const data = await response.json();
25         // 使用返回的数据渲染页面
26         loading.style.display = 'none';
27         return data;
28     } catch (error) {
29         console.error('Error:', error);
30     }
31 }
32
33 function displayPage() {
34     if (data.length === 0) console.log("获取文章失败");
35     const container = document.querySelector('.article-container');
36     if (container) {
37         container.innerHTML = '';
38         getKeywordsWordcloud()
39         // 遍历数据数组，渲染每个文章块
40         for (let i = 0; i < data.authors.length; i++) {
41             // 创建一个新的文章块元素
42             let article = document.createElement('div');
43             article.className = 'article';
44
45             // 填充文章块的内容
46             article.innerHTML = `
47                 <h2>${data.titles[i]}</h2>
48                 <p>Author: ${data.authors[i]}</p>
49                 <p>Keywords: ${data.keywords[i]}</p>
50                 <a href="${data.pdfs[i]}" target="_blank" class="download-button">Download PDF</a>
51             `;
52
53             // 将文章块添加到 article-container 中
54             container.appendChild(article);
55         }
56     } else {
57         console.error('Article container not found!');
58     }
59 }
60
61 function initPage() {
62     clickTab(0)
63 }
64
65 function getKeywordsWordcloud() {
66     const container = document.querySelector('.img-container');
67     container.innerHTML = '';
68     fetch(`http://localhost:5000/api/keywords-wordcloud`)
69         .then(response => response.blob())
70         .then(blob => {
71             // 创建一个表示图片的URL
72             let imageUrl = URL.createObjectURL(blob);
73
74             // 在页面上显示词云图片
75             let imgElement = document.createElement('img');
76             imgElement.src = imageUrl;
77             container.appendChild(imgElement);
78
79         })
80 }
```

```

79         .catch(error => console.error('Error fetching wordcloud:', error));
80     }
81
82     function search() {
83         // 获取输入的关键词
84         let keyword = document.getElementById('searchInput').value;
85         // 发送搜索请求到后端API
86         fetch(`http://localhost:5000/api/search?keyword=${encodeURIComponent(keyword)}`)
87             .then(response => response.json())
88             .then(data => {
89                 // 处理搜索结果
90                 console.log(data);
91                 showSearchResult(data)
92             })
93             .catch(error => {
94                 console.error('Error during search:', error);
95             });
96     }
97
98     function showSearchResult(indexs) {
99         const img_container = document.querySelector('.img-container');
100         img_container.innerHTML = '';
101         const art_container = document.querySelector('.article-container');
102         art_container.innerHTML = '';
103         for (let i = 0; i < indexs.length; i++) {
104             // 创建一个新的文章块元素
105             let article = document.createElement('div');
106             article.className = 'article';
107
108             // 填充文章块的内容
109             article.innerHTML = `
110                 <h2>${data.titles[indexs[i]]}</h2>
111                 <p>Author: ${data.authors[indexs[i]]}</p>
112                 <p>Keywords: ${data.keywords[indexs[i]]}</p>
113                 <a href="${data.pdfs[indexs[i]]}" target="_blank" class="download">Download PDF</a>
114             `;
115
116             // 将文章块添加到 article-container 中
117             art_container.appendChild(article);
118         }
119     }

```

3.3.3. 样式表

鉴于本次课程为Python，故不详细展开介绍css部分。

```

1  body {
2      font-family: -apple-system,BlinkMacSystemFont,Helvetica Neue,PingFang SC,Microsoft YaHei,Source Han Sans SC,Noto Sans CJK
      SC,WenQuanYi Micro Hei,sans-serif;
3      color: #121212;
4      font-size: 15px;
5      display: flex;
6      flex-direction: column;
7      padding: 80px;
8  }
9  h1 {
10     text-align: center;
11     padding: 20px;
12 }
13
14 .header {
15     display: flex;
16     align-items: center;
17     justify-content: space-between;
18     background-color: #fff;
19     height: 100px;
20     padding: 0 10px;
21     box-shadow: 0 1px 3px rgba(0, 0, 0, 0.1);
22 }
23
24 .logo {
25     width: 80px;
26     height: 80px;
27 }
28
29 .home {
30     margin-right: auto;
31     padding: 0 40px;

```

```
32     color: #121212;
33     text-align: center;
34     font-size: 25px;
35     font-weight: 600;
36     text-decoration:none;
37 }
38
39 .search {
40     display: flex;
41     align-items: center;
42     background-color: #f5f5f5;
43     border-radius: 20px;
44     padding: 0 20px;
45     height: 40px;
46     width: 300px;
47 }
48
49 .search input {
50     border: none;
51     outline: none;
52     background-color: transparent;
53     margin-left: 10px;
54     font-size: 16px;
55     width: 100%;
56 }
57
58 .search input::placeholder {
59     color: #999;
60 }
61
62 .search button {
63     border: none;
64     outline: none;
65     background-color: transparent;
66     cursor: pointer;
67 }
68
69 .tabCardContainer {
70     display: flex;
71     align-items: center;
72     justify-content: space-between;
73     height: 58px;
74     width: 100%;
75 }
76
77 .tabCard {
78     cursor: pointer;
79     margin: 0 22px;
80     border: none;
81     outline: none;
82     background-color: transparent;
83     font-size: 16px;
84 }
85
86 #loading {
87     display: flex;
88     align-items: center;
89     justify-content: center;
90     width:100%;
91     height:80px;
92     text-align: center;
93 }
94 #loading p {
95     font-size: 20px;
96 }
97
98 .article {
99     color: #121212;
100     font-family: -apple-system,BlinkMacSystemFont,Helvetica Neue,PingFang SC,Microsoft YaHei,Source Han Sans SC,Noto Sans CJK
    SC,WenQuanYi Micro Hei,sans-serif;
101     font-size: 15px;
102     -webkit-tap-highlight-color: rgba(18,18,18,0);
103     background: #fff;
104     box-sizing: border-box;
105     border-radius: 0;
106     outline: none;
107     overflow: initial;
108     position: relative;
109
110     padding: 20px;
```

```
110     border-bottom: 1px solid #f5f5f5;
111     box-shadow: none;
112     margin-bottom: 0;
113 }
114
115 .img-container {
116     display: flex;
117     justify-items: center;
118     align-items: center;
119 }
120
121 .download-button {
122     background-color: #999999;
123     color: white;
124     font-size: 20px;
125 }
```

4. 成员分工及贡献占比

本小组共5名成员，分工及贡献占比如下：

学号+姓名	分工	贡献占比
61522312万奕含	爬虫部分	100%
61522314吴清晏	api接口实现及文档撰写	100%
71122207董子翔	前端js代码完善	99%
58122327李家豪	前端页面搭建	99%
61522122李宗辉	前端css样式表调整及制作视频汇报	99%