



南開大學  
Nankai University

南开大学

---

## 信息检索系统原理——Web 搜索引擎

---

孟启轩

学号 : 2212452

专业 : 计算机科学与技术

指导教师 : 温延龙

2024 年 12 月 16 日

## 摘要

本实验基于信息检索系统原理课程所学，自主设计并实现了一个南开大学 Web 搜索引擎。结合 Elastic Search 实现了网页抓取与数据采集，采用先进的文本索引技术和链接分析算法（PageRank）对网页内容和结构进行深入解析与处理。搜索引擎支持多种查询服务，包括短语查询、通配符查询、查询日志等，显著提升了搜索的灵活性和精准度。此外，为了确保良好的用户交互，还设计了直观友好的 Web 页面界面以及个性化推荐功能，也尝试设计了个性化查询的功能。

**关键字：**信息检索，南开大学，搜索引擎，Python

## 目录

<b>一、 项目综述</b>	1
<b>二、 网页抓取</b>	2
(一) 设计思路 . . . . .	2
1. 爬虫目标 . . . . .	2
2. 爬虫功能需求 . . . . .	2
(二) 技术路线 . . . . .	2
1. 爬虫架构设计与实现细节 . . . . .	2
2. 数据清洗 . . . . .	3
(三) 完成结果 . . . . .	3
<b>三、 文本索引</b>	4
(一) 设计思路 . . . . .	4
(二) 技术路线 . . . . .	4
1. Elastic Search 连接 . . . . .	4
2. 创建索引并定义映射 . . . . .	4
3. 批量上传数据 . . . . .	5
4. 错误处理 . . . . .	5
(三) 完成结果 . . . . .	6
<b>四、 链接分析</b>	6
(一) 设计思路 . . . . .	7
(二) 技术路线 . . . . .	7
(三) 完成结果 . . . . .	8
<b>五、 查询服务</b>	8
(一) 实现功能 . . . . .	8
(二) 技术路线 . . . . .	9
1. 结果评分和排序 . . . . .	9
2. 短语查询 . . . . .	9
3. 通配符查询 . . . . .	10
4. 查询日志 . . . . .	10
(三) 完成结果 . . . . .	11

<b>六、 个性化查询</b>	<b>12</b>
(一) 设计思路 . . . . .	12
(二) 技术路线 . . . . .	13
<b>七、 Web 界面</b>	<b>13</b>
<b>八、 个性化推荐</b>	<b>14</b>
(一) 设计思路 . . . . .	14
(二) 技术路线 . . . . .	14
(三) 完成结果 . . . . .	15
<b>九、 总结与感悟</b>	<b>15</b>

## 一、项目综述

设计项目最初想要通过自己一步一步将整个 Web 搜索引擎搭建起来，但是在构建索引以及向量空间模型等步骤难度比较大，正确性也难以验证。所以经过几天的尝试后，决定结合 Elastic Search 这个强大的工具。系统架构如图1所示。



图 1: 系统架构

## 二、网页抓取

### (一) 设计思路

#### 1. 爬虫目标

- **目标域:** nankai.edu.cn
- **起始网站:** 南开大学官方网站 <http://www.nankai.edu.cn/>
- **抓取内容:** 网页标题 (title)、URL (url)、正文内容 (text)、页面内链接 (linksurl)。
- **数据存储:** 将抓取的数据存储至 CSV 文件，以便后续导入至 Elastic Search 进行索引构建。
- **反馈优化模块:** 通过用户反馈优化检索和生成策略。

#### 2. 爬虫功能需求

- **全面抓取:** 遍历目标网站的可访问页面，确保覆盖主要内容。
- **数据清洗:** 提取有效文本，去除无关内容（如脚本、样式）。
- **链接管理:** 处理相对路径链接，确保抓取范围限定在指定域名内。
- **效率优化:** 通过并发请求和自动限速，提高抓取速度，同时避免对目标服务器造成过大负载。
- **数据存储:** 结构化存储抓取数据，便于后续处理和分析。

### (二) 技术路线

#### 1. 爬虫架构设计与实现细节

- **爬虫框架:** Scrapy，高效、灵活，支持异步抓取，适合大规模数据抓取。
- **解析方法 (parse):** 主要负责从网页中提取并处理信息。具体来说，它使用 XPath 技术来选取网页的 <title> 标签内容作为标题，并挑选 <body> 标签内的文本节点（排除 <script> 和 <style> 标签）以获取正文内容，之后将这些文本合并并调用 clean\_text 方法去除多余的空格和换行符。同时，该方法会通过 XPath 提取所有 <a> 标签中的 href 属性来收集链接，处理其中的相对路径为绝对 URL，并过滤只保留属于 nankai.edu.cn 域名下的链接，最终以分号分隔的形式保存在 linksurl 字段中。所提取的数据将通过调用 save\_to\_csv 方法被写入 CSV 文件中进行存储。此外，对于每一个有效的链接，parse 方法还会递归地发起新请求继续抓取更多页面内容，确保了数据采集的深度与广度。
- **爬虫配置:** 请求间隔 (DOWNLOAD\_DELAY = 0.15)，爬取深度 (DEPTH\_LIMIT = 100)，并发配置 (最大并发请求数量 32, 单域并发请求数限制 16)，自动限速 (AUTOTHROTTLE)，请求头设置 (User-Agent 伪装成浏览器)
- **数据存储:** 将抓取的数据存储至 CSV 文件，以便后续导入至 Elasticsearch 进行索引构建。
- **允许爬取的域:** allowed\_domains = ["nankai.edu.cn"]

```

1 def parse(self, response):
2     # 提取网页标题
3     title = response.xpath("//title/text()").get(default="Untitled")
4     # 提取正文内容，仅保留有效文本
5     raw_text = response.xpath("//body//*[not(self::script or self::style)]/text()")
6     .getall()
7     clean_content = self.clean_text(" ".join(raw_text))
8     # 提取链接
9     links = response.xpath("//a[@href]/@href").extract()
10    full_links = []
11    for link in links:
12        if link.startswith("/"): # 处理相对路径
13            link = response.urljoin(link)
14        if self.allowed_domains[0] in link: # 确保只抓取允许的域
15            full_links.append(link)
16    linksurl = "; ".join(full_links) # 用分号分隔所有链接
17    # 保存数据到 CSV
18    self.save_to_csv(title, response.url, clean_content, linksurl)
19    # 继续爬取链接
20    for link in full_links:
21        yield scrapy.Request(link, callback=self.parse)

```

Listing 1: 核心爬虫代码

## 2. 数据清洗

在爬取到数据后，发现有许多无法使用的数据，比如很多需要登陆的网页，会显示“跳转提示：访问的链接无效”等，所以对原始数据进行简单的清洗。

```

1 # 删除标题为 "提示信息" 或 "跳转提示" 或 "Untitled" 的行
2 df = df[df["title"] != "提示信息"]
3 df = df[df["title"] != "跳转提示"]
4 df = df[df["title"] != "Untitled"]
5 # 删除正文为空的行
6 df = df[df["text"].notnull()]
7 # 删除包含特定无效关键词的行
8 invalid_keywords = ["您使用的浏览器", "系统提示", "附件下载", "错误的模板参数", "暂不支
持移动端", " ", "跳转提示", "无效的文章地址", "访问地址无效"]
9 df = df[~df["text"].apply(lambda x: any(keyword in x for keyword in invalid_keywords))]
10 # 去重操作：删除完全相同的数据行
11 df = df.drop_duplicates()

```

Listing 2: 核心数据清洗代码

## (三) 完成结果

网页抓取得到的数据如图2所示，经过清洗后，共有 134879 条数据。

### 三、文本索引

data > cleanedfkuoutput.csv
1 title,url,text,linksurl 2 南开大学,https://www.nankai.edu.cn/,信息公开 图书馆 服务指南 登录邮箱 办公室 校友入校   English 学校概况 学校简介 现任领导 历届领导 历史回顾 南开新闻 3 南开大学研究生招生网,https://yzb.nankai.edu.cn/,收藏本站 导航 首页 硕士招生 博士招生 港澳台研究生招生 留学研究生招生 院系信息 师生风采 联系我们 报考指南 4 国际经贸关系专业,南开大学经济与社会发展研究院,https://esd.nankai.edu.cn/jxwm/gjhzbx/xmjk/gjgngzxy.htm,网站首页 学院概况 院长寄语 学院简介 组织机构 5 区域政策研究专业,https://chinareal.nankai.edu.cn/,“百闻”机构介绍 发展目标 主要功能 突出特色 新闻资讯 研究前沿 国际前沿 前沿文献导读 经典著作与文献 6 专题,https://www.nankai.edu.cn/zt/list.htm,信息公告 图书馆 服务指南 登录邮箱 办公室 校友入校   English 首页 学院概况 学校简介 现任领导 历届领导   7 南开大学招投标管理中心 - 首页,https://nkbb.nankai.edu.cn/,首页 部门概况 领导简介 办公室 职责 公告通知 荣誉公告 货物类 工程类 服务类 其他类 结果公告 8 南开大学中国区域经济应用实验室,https://nku-chinareal.nankai.edu.cn/,首页 部门概况 领导简介 办公室 职责 公告通知 荣誉公告 货物类 工程类 服务类 其他类 结果公告 9 南开校友网,https://nkuau.nankai.edu.cn/,> 南开校友网 育人 新闻动态 公告 南开新闻 总会快讯 分会动态 校友会 会长致辞 总会简介 总会章程 组织机构 校友 10 南开大学活动,https://www.nankai.edu.cn/xshd/list.htm,信息公开 图书馆 服务指南 登录邮箱 办公室 校友入校   English 首页 学院概况 学校简介 现任领导 历届 11 南开大学附属医院,https://fy.nankai.edu.cn/,“百闻”机构介绍 基层党组织 主要职能 部门领导 基层党组织 主要职能 部门领导 机构设置 后勤服务处   12 南开大学后勤保障服务,https://hq.nankai.edu.cn/,首页 机构概况 后勤党委 主要职能 部门领导 基层党组织 主要职能 部门领导 机构设置 后勤服务处   13 南开大学经济与社会发展研究院,https://esd.nankai.edu.cn/,“百闻”机构介绍 院长寄语 联系我们 教学项目 全日制硕士 全日制硕 14 2024年暑假后各单位服务指南,https://www.nankai.edu.cn/2024/9626/c17471a546267/page.htm,信息公告 图书馆 服务指南 登录邮箱 办公网 校友入校   Eng 15 南开大学爱国主义教育基地,https://aiguo.nankai.edu.cn/,导航 首页 南开简介 历史沿革 南开参观点 南开故事 南开人物 南开精神 参观点爱国主义路线 参观点介绍 16 南开大学滨海学院,https://bihai.nankai.edu.cn/,首页 学院概况 学院简介 现任领导简介 学院新闻 网站新闻 滨海学院报 就业工作 教育教学 教务部 精品课网 17 南开大学滨海学院,https://bihai.nankai.edu.cn/,首页 学院概况 学院简介 现任领导简介 学院新闻 网站新闻 滨海学院报 就业工作 教育教学 教务部 精品课网 134857 网内子,http://www.cim.nankai.edu.cn/2021/0419/c6729a353075/page.htm,自贡 English 版本 天下找我 数字所简介 顾问委员会和学术委员会名单 历届领导   134858 连增,http://www.cim.nankai.edu.cn/2021/0423/c6729a354123/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134859 黄文,http://www.cim.nankai.edu.cn/2021/0423/c6729a354124/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134860 杨志英,http://www.cim.nankai.edu.cn/2021/0603/c6729a370028/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134861 冯克勤,http://www.cim.nankai.edu.cn/2021/0610/c6729a371925/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134862 李海中,http://www.cim.nankai.edu.cn/2021/0624/c6729a371914/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134863 高速,http://www.cim.nankai.edu.cn/2021/0923/c6729a397678/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134864 董景英,http://www.cim.nankai.edu.cn/2021/1014/c6729a403503/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134865 丁剑,http://www.cim.nankai.edu.cn/2021/1021/c6729a405951/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134866 蒋萍澜,http://www.cim.nankai.edu.cn/2022/0505/c6729a448225/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134867 尚名久,http://www.cim.nankai.edu.cn/2018/1115/c6729a13657/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134868 张晓,http://www.cim.nankai.edu.cn/2018/1129/c6729a15586/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134869 林秉甫,http://www.cim.nankai.edu.cn/2018/1129/c6729a15587/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134870 张功斌,http://www.cim.nankai.edu.cn/2018/1213/c6729a16706/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134871 Akito Futaki,http://www.cim.nankai.edu.cn/2019/0325/c6729a12132/page.htm,首页 English 版本 关于我们 数学所简介 顾问委员会和学术委员会名单 历届领导   134872 实践活动 南开大学社会实践平台,http://shsj.nankai.edu.cn/index/huodong/sxtz/p/2035.aspx,实践活动 师生同行 基地 书屋 新闻 首页 参加活动 师生同行   134873 实践活动 南开大学社会实践活动平台,http://shsj.nankai.edu.cn/index/huodong/p/1.aspx,实践活动 师生同行 基地 书屋 新闻 首页 参加活动 师生同行   134874 一起做年夜饭 实践活动 南开大学社会实践活动平台,http://shsj.nankai.edu.cn/index/hdview/1d/22423.aspx,实践活动 师生同行 基地 书屋 新闻 首页 参加活动 师生同行   134875 敬老孝亲,传承良好年夜饭 实践活动 南开大学社会实践活动平台,http://shsj.nankai.edu.cn/index/hdview/1d/22660.aspx,实践活动 师生同行 基地 书屋 新闻 首页 参加活动 师生同行   134876 在除夕这天为家人做一顿年夜饭 实践活动 南开大学社会实践活动平台,http://shsj.nankai.edu.cn/index/hdview/1d/22493.aspx,实践活动 师生同行 基地 书屋 新闻   134877 余姚特色年夜饭 实践活动 南开大学社会实践活动平台,http://shsj.nankai.edu.cn/index/hdview/1d/22525.aspx,实践活动 师生同行 基地 书屋 新闻 首页 参加活动   134878 我的家,我们的年 实践活动 南开大学社会实践活动平台,http://shsj.nankai.edu.cn/index/hdview/1d/22610.aspx,实践活动 师生同行 基地 书屋 新闻 首页 参加活动   134879 134880

图 2: 数据

### 三、文本索引

最初尝试使用 Whoosh 和 jieba 等工具对爬取到的数据构建倒排索引，不仅构建时间长，而且结果异常庞大，压缩后的结果也不易查看，正确性也难以验证。最终使用 Elastic Search 构建索引，不仅快速，正确性也得到保证，并且还有不错的可视化。

#### (一) 设计思路

- 连接:** 与 Elastic Search 集群交互。
- 高效上传数据:** 由于数据量非常大，而 Elastic Search 直接导入数据有大小限制，所以要确保大量数据能够快速、稳定地上传。
- 准确构建索引:** 设计合理的索引映射，确保数据的准确性和查询效率。
- 错误处理与日志记录:** 有效处理上传过程中可能出现的错误，并记录失败的文档以便后续分析。

#### (二) 技术路线

##### 1. Elastic Search 连接

使用 Elastic Search 类连接到指定的主机 (ES\_HOST = 'http://localhost:9200')。

##### 2. 创建索引并定义映射

- title:** text 类型，用于全文搜索。
- url:** keyword 类型，用于精确匹配和过滤。

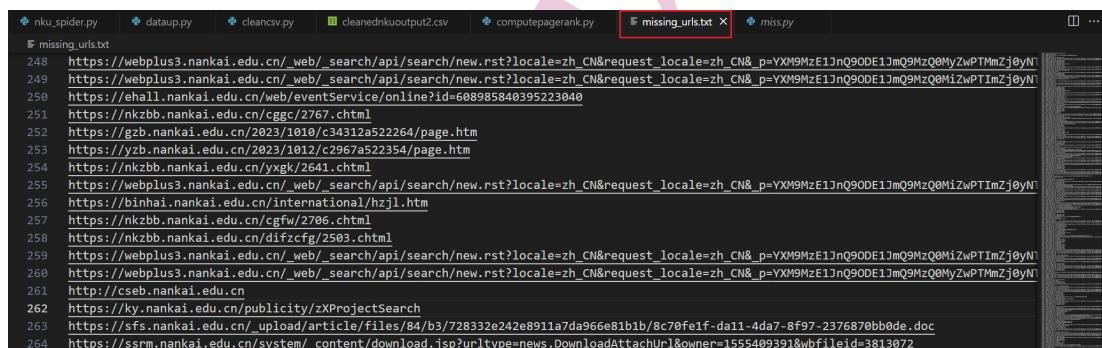
- **text:** text 类型，用于全文搜索。
- **linksurl:** keyword 类型，用于存储页面内的链接。
- **pagerank:** float 类型，用于存储页面的权重评分。(在链接分析部分详细解释)
- **suggest:** completion 类型，用于实现联想搜索功能。(在个性化推荐部分详细解释)

### 3. 批量上传数据

使用 Pandas 的 `read_csv` 方法以分块方式读取 CSV 文件，设置 `chunksize` 为 1000，以分批处理数据，提升内存利用率和上传效率。并且使用 `tqdm` 库显示数据上传的进度条，提供实时反馈。如图4。

### 4. 错误处理

- **捕获上传错误:** 在批量上传过程中，可能会遇到部分文档上传失败的情况，导致 CSV 中数据数量和索引中的文档数量不匹配。对于每个失败的文档，提取错误原因和文档 url，并将其记录到 `failed_documents.json` 文件中。
- **异常处理:** 通过捕捉到的错误的 url (如图3)，再次清洗数据，然后通过 Elastic 控制台 `DELETE /xxjs` 删除之前的索引，并重新构建索引上传数据。

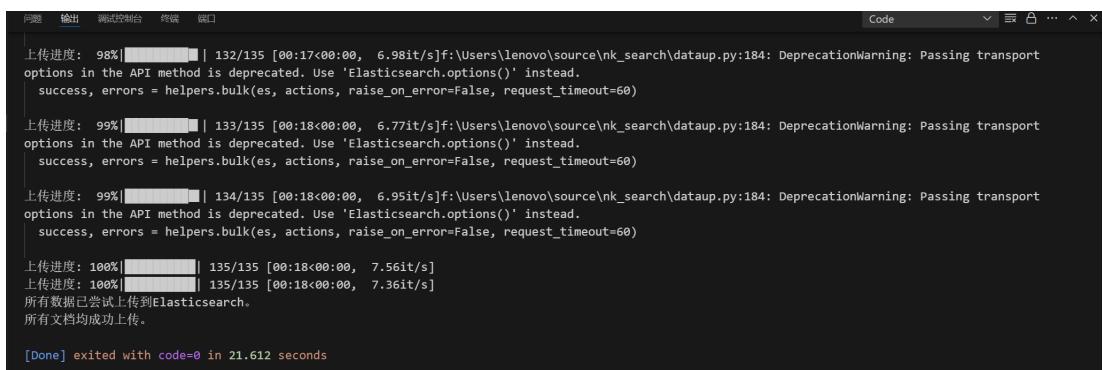


```

nku_spider.py    dataup.py    cleancsv.py    cleanednkuoutput2.csv    computepagerank.py    missing_urls.txt    miss.py
missing_urls.txt
248 https://webplus3.nankai.edu.cn/_web/_search/api/search/new.rst?locale=zh_CN&request_locale=zh_CN&p=YXM9MzE1JnQ90DE1JmQ9MzQ0MyZwPTImZj0yN
249 https://webplus3.nankai.edu.cn/_web/_search/api/search/new.rst?locale=zh_CN&request_locale=zh_CN&p=YXM9MzE1JnQ90DE1JmQ9MzQ0MyZwPTImZj0yN
250 https://ehall.nankai.edu.cn/web/eventsService/online?id=608985840395223040
251 https://nkzbb.nankai.edu.cn/cggc/2767.shtml
252 https://gzb.nankai.edu.cn/2023/1010/c34312a522264/page.htm
253 https://yzb.nankai.edu.cn/2023/1012/c2967a522354/page.htm
254 https://nkzbb.nankai.edu.cn/yxgk/2641.shtml
255 https://webplus3.nankai.edu.cn/_web/_search/api/search/new.rst?locale=zh_CN&request_locale=zh_CN&p=YXM9MzE1JnQ90DE1JmQ9MzQ0MyZwPTImZj0yN
256 https://binhai.nankai.edu.cn/international/hzjl.htm
257 https://nkzbb.nankai.edu.cn/cgfw/2706.shtml
258 https://nkzbb.nankai.edu.cn/difzcfg/2503.shtml
259 https://webplus3.nankai.edu.cn/_web/_search/api/search/new.rst?locale=zh_CN&request_locale=zh_CN&p=YXM9MzE1JnQ90DE1JmQ9MzQ0MyZwPTImZj0yN
260 https://webplus3.nankai.edu.cn/_web/_search/api/search/new.rst?locale=zh_CN&request_locale=zh_CN&p=YXM9MzE1JnQ90DE1JmQ9MzQ0MyZwPTImZj0yN
261 http://cseb.nankai.edu.cn
262 https://ky.nankai.edu.cn/publicity/zxprojectSearch
263 https://sfs.nankai.edu.cn/upload/article/files/84/b3/72833e242e8911a7da966e81b1b/8c70fe1f-da11-4da7-8f97-2376878bbde.doc
264 https://ssrm.nankai.edu.cn/system/_content/download.jsp?urltype=news.DownloadAttachUrl&owner=1555409391&wbfileid=3813072

```

图 3: 丢失的数据



```

上传进度: 98% [██████████] | 132/135 [00:17<00:00,  6.98bit/s] f:\Users\lenovo\source\nk_search\dataup.py:184: DeprecationWarning: Passing transport options in the API method is deprecated. Use 'Elasticsearch.options()' instead.
    success, errors = helpers.bulk(es, actions, raise_on_error=False, request_timeout=60)

上传进度: 99% [██████████] | 133/135 [00:18<00:00,  6.77it/s] f:\Users\lenovo\source\nk_search\dataup.py:184: DeprecationWarning: Passing transport options in the API method is deprecated. Use 'Elasticsearch.options()' instead.
    success, errors = helpers.bulk(es, actions, raise_on_error=False, request_timeout=60)

上传进度: 99% [██████████] | 134/135 [00:18<00:00,  6.95it/s] f:\Users\lenovo\source\nk_search\dataup.py:184: DeprecationWarning: Passing transport options in the API method is deprecated. Use 'Elasticsearch.options()' instead.
    success, errors = helpers.bulk(es, actions, raise_on_error=False, request_timeout=60)

上传进度: 100% [██████████] | 135/135 [00:18<00:00,  7.56it/s]
上传进度: 100% [██████████] | 135/135 [00:18<00:00,  7.36it/s]
所有数据已尝试上传到Elasticsearch。
所有文档均成功上传。

[Done] exited with code=0 in 21.612 seconds

```

图 4: 上传数据

```

1 INDEX_NAME = 'xxjs'                                # Elasticsearch索引名称
2 CHUNK_SIZE = 1000                                 # 每批处理的行数
3
4 # 读取CSV并分批上传
5 chunks = pd.read_csv(CSV_FILE_PATH, chunksize=CHUNK_SIZE, encoding='utf-8')
6 for chunk in tqdm(chunks, total=total_lines//CHUNK_SIZE + 1, desc="上传进度"):
7     # 构建批量操作
8     actions = []
9     for _, row in chunk.iterrows():
10         action = {
11             "_index": INDEX_NAME,
12             "_source": {
13                 "title": row.get('title', ''),
14                 "url": row.get('url', ''),
15                 "text": row.get('text', ''),
16                 "linksurl": row.get('linksurl', ''),
17                 "pagerank": row.get('pagerank', ''),
18                 "suggest": row.get('suggest', '')
19             }
20         }
21         actions.append(action)

```

Listing 3: 索引构建与数据上传

### (三) 完成结果

The screenshot shows the Elasticsearch UI for the 'xxjs' index. At the top, it displays '文档 (128,310)' and '文档数 134737'. Below this, the 'Discover' tab is active, showing a table of search results. The table columns include '\_id', '\_index', '\_score', '\_type', '\_version', '标题', '摘要', '文本', 'url', '链接', and '建议'. Each row contains a snippet of text from the document, such as '习近平在山西考察工作结束时强调，要牢固树立绿水青山就是金山银山的理念，坚定不移走绿色发展的道路。' and '习近平指出，山西要牢固树立绿水青山就是金山银山的理念，坚定不移走绿色发展的道路。'.

图 5: 索引

## 四、链接分析

由于我的 CSV 文件数量和索引中的文档数量不一致（上传 Elastic 时有些文档出错），并且通过 CSV 计算 pagerank 的文档 id 是行号，而 Elastic 上会为每个文档生成唯一的随机 id，使得文档到 id 的映射不对应，导致所有文档的 pagerank 都是 0。

那么先计算 pagerank，然后将 pagerank 属性添加到每个文档后面，再构建索引上传数据，问题不就解决了吗。

## (一) 设计思路

- **构建图结构：**基于 CSV 中的链接数据，构建一个有向图（Directed Graph），每个网页视为一个节点，链接关系视为边。
- **计算 PageRank：**使用 NetworkX 的 pagerank 算法计算每个节点（网页）的 pagerank 值。
- **添加 pagerank 属性到索引：**为每个文档添加 pagerank 属性，然后重新构建索引上传数据。

## (二) 技术路线

- **处理链接数据：**在每一行数据中，提取网页的 url 和它链接的 linksurl 字段。对每个链接进行检查，若非空，则将该链接作为目标网页的指向网页，建立有向边。
- **图的构建：**使用 NetworkX 库中的 DiGraph（有向图）数据结构。对于每一对（网页 A，链接 B），我们在图中添加一条从网页 A 到网页 B 的有向边。
- **图的优化与限制：**对于每条链接，进行有效性检查，确保其不为空并且不重复。若图中某些节点（网页）没有任何出链（即没有任何网页链接到它），也需在图中保留这些节点，防止图的计算出错。
- **PageRank 算法：**使用 NetworkX 中的 pagerank 函数来计算每个节点的 pagerank 值。PageRank 算法基于以下公式：

$$PR(A) = \frac{1-d}{N} + d \sum_{i \in In(A)} \frac{PR(i)}{L(i)}$$

其中：

- $PR(A)$  是网页 A 的 PageRank 值。
- $d$  是阻尼因子，通常设置为 0.85。
- $In(A)$  是指向网页 A 的网页集合。
- $L(i)$  是网页  $i$  的出链数。
- $N$  是图中网页的总数。

该算法基于迭代计算，直到所有网页的 PageRank 值收敛。

```

1 def compute_pagerank(csv_file_path):
2     G = nx.DiGraph()      # 创建有向图
3
4     # 读取 CSV 数据
5     with open(csv_file_path, 'r', encoding='utf-8') as f:
6         reader = csv.DictReader(f)
7
8         for row in reader:
9             url = row['url']

```

```

10     links = row['linksurl'].split(';') # 链接是以分号分隔的
11     # 将每个网页与其链接添加到图中
12     for link in links:
13         link = link.strip()
14         if link: # 确保链接不为空
15             G.add_edge(url, link)
16
17     pagerank = nx.pagerank(G, alpha=0.85) # 计算 PageRank 值
18

```

Listing 4: 链接分析

### (三) 完成结果

```

[running] python -u "f:\Users\lenovo\source\nk_search\computepagerank.py"
读取CSV文件...
CSV文件读取完成，共有 135657 条记录。
创建URL到文档ID的映射...
总共有 139337 个文档被映射。
构建图结构...
构建图结构: 0% | 0/135657 [00:00<?, ?it/s]
构建图结构: 1% | 1958/135657 [00:00<00:12, 10519.37it/s]
构建图结构: 2% | 2265/135657 [00:00<00:11, 11413.19it/s]
构建图结构: 3% | 3605/135657 [00:00<00:10, 12297.31it/s]
构建图结构: 4% | 4989/135657 [00:00<00:10, 12891.47it/s]
构建图结构: 5% | 6279/135657 [00:00<00:10, 12270.72it/s]
构建图结构: 6% | 7511/135657 [00:00<00:10, 11948.75it/s]
构建图结构: 6% | 8710/135657 [00:00<00:11, 11336.68it/s]
构建图结构: 7% | 9851/135657 [00:00<00:14, 8839.14it/s]
构建图结构: 8% | 10835/135657 [00:01<00:13, 9078.27it/s]
构建图结构: 9% | 11900/135657 [00:01<00:13, 9467.88it/s]
构建图结构: 10% | 12953/135657 [00:01<00:12, 9741.45it/s]
构建图结构: 10% | 13964/135657 [00:01<00:12, 9700.94it/s]
构建图结构: 11% | 14960/135657 [00:01<00:12, 9321.93it/s]
构建图结构: 12% | 16180/135657 [00:01<00:11, 10094.87it/s]
构建图结构: 13% | 17950/135657 [00:01<00:09, 12243.59it/s]
构建图结构: 15% | 19917/135657 [00:01<00:08, 14373.77it/s]

```

图 6: 链接分析

```

dataup.py    pagerank2.py    pagerankedoutput.csv ×
pagerankedoutput.csv
149 te.nankai.edu.cn; https://graduate.nankai.edu.cn/_s11/byjj/1
150 职工生育津贴政策 最新招聘及引进 • 南开大学2025年学生工作辅导员岗位
151 .ist.htm; https://bs.nankai.edu.cn/gltd/list.htm; https://bs
152 .nankai.edu.cn/7264/list.htm; https://sfs.nankai.edu.cn/7265
153 e.nankai.edu.cn/2024/0909/c34572a550319/page.htm; https://fi
154 7D3BB5C4603CD7BE8A4042F2C2FE3 [1.893046976412916e-06]
155 学与文化》连续入选“复印... [ 2024-09-14 ] « 南开大学团队发布《
156 ankai.edu.cn/gltd/list.htm; https://bs.nankai.edu.cn/zpxx/li
157 队与花开岭公益基地举行了签约仪式，建立南开大学学生社会实践基地。实
158 j.nankai.edu.cn/index/jidi.aspx; http://shsj.nankai.edu.cn/:

```

图 7: pagerank

## 五、 查询服务

### (一) 实现功能

- 站内查询:** 基本的查询操作。
- 短语查询:** 基本的查询操作，支持对多个 Term 的查询。

- **通配查询：**用户可能并不知道自己要查什么，支持通配符（正则）查询操作。
- **查询日志：**查询日志（历史），能够知道用户之前查了什么。
- **结果评分和排序：**根据 final\_score 将查询到的结果排序，返回得分高的结果。

## (二) 技术路线

### 1. 结果评分和排序

- **查询执行与数据提取：**执行查询，获取匹配的文档，并提取所需的数据（score 和 pagerank 值）。每个文档的评分（score）是 Elastic Search 根据查询条件自动计算的，表示该文档与查询的相关度。
- **标准化处理：**因为 score 和 pagerank 并不在同一量级，所以无法直接加权组合，需要对 score 和 pagerank 值进行标准化处理，以便将它们放在相同的量纲下进行加权计算。
- **加权组合得到 final\_score：**设置了 TFIDF\_WEIGHT = 0.7 和 PAGERANK\_WEIGHT = 0.3，使用加权平均的方式，将标准化后的 score 和 pagerank 结合起来得到 final\_score，根据最终评分返回得分靠前的结果。

```

1 scores = [hit['_score'] for hit in hits]
2 pageranks = [hit['_source'].get('pagerank', 0) for hit in hits]
3 normalized_scores = self.normalize_values(scores)
4 normalized_pageranks = self.normalize_values(pageranks)
5 TFIDF_WEIGHT = 0.7
6 PAGERANK_WEIGHT = 0.3
7
8 results = []
9 for idx, hit in enumerate(hits):
10     source = hit['_source']
11     final_score = TFIDF_WEIGHT * normalized_scores[idx] + PAGERANK_WEIGHT *
12         normalized_pageranks[idx]
13     results.append({
14         'title': source.get('title', ''),
15         'url': source.get('url', ''),
16         'text': source.get('text', '')[:200], # snippet
17         'pagerank': source.get('pagerank', 0),
18         'final_score': final_score
19     })
20 results.sort(key=lambda x: x['final_score'], reverse=True)

```

Listing 5: 结果评分和排序

### 2. 短语查询

- **查询类型：**使用 Elastic Search 的 match\_phrase 查询来查找包含精确短语的文档。
- **查询过程：**用户输入查询短语，函数将其传递给 match\_phrase 查询模块。然后通过 Elasticsearch 查询文档中是否包含该短语。最后返回匹配文档的结果。
- **注意：**“南开大学”和“南开 大学”的区别。

```

1 def search_phrase(self, phrase, user=None):
2     query = {
3         "query": {
4             "match_phrase": {
5                 "text": {
6                     "query": phrase
7                 }
8             }
9         }
10    }
11    return self.execute_query(query, user)

```

Listing 6: 短语查询

### 3. 通配查询

- 查询类型:** 使用 Elastic Search 的 wildcard 查询来实现通配符查询。
- 查询过程:** 用户输入查询字符串，允许模糊匹配。然后将查询传递给 wildcard 查询模块，搜索文本字段中符合条件的文档。最后返回匹配结果。

```

1 def search_wildcard(self, wildcard_query, user=None):
2     query = {
3         "query": {
4             "wildcard": {
5                 "text": {
6                     "value": wildcard_query.lower(),
7                     "case_insensitive": True
8                 }
9             }
10            }
11        }
12        return self.execute_query(query, user)

```

Listing 7: 通配查询

### 4. 查询日志

- 查询日志功能:** 在每次用户执行查询时，记录查询时间、用户 ID、查询的短语或通配符、相关文档（如标题、URL、摘要）等信息，保存为日志文件（history.txt）。
- 日志查看:** 返回用户的搜索历史。

```

1 def log_query(self, user, query, results, history_file='history.txt'):
2     timestamp = datetime.now().strftime('%Y-%m-%d %H:%M:%S')
3     with open(history_file, 'a', encoding='utf-8') as f:
4         f.write(f"{timestamp} | {user} | {query}\n")
5         for res in results:
6             snippet = res['text'].replace('\n', ' ').replace('\r', ' ')[:100]
7             f.write(f"{timestamp} | {user} | {res['title']} | {res['url']} | {snippet}\n")

```

Listing 8: 查询日志

### (三) 完成结果

短语查询，通配查询，查询日志结果如下：

```
===== 搜索引擎 =====
选择查询类型：
1. 短语查询
2. 通配查询
3. 联想关联（搜索建议）
4. 查看查询历史
5. 退出登录
输入数字选择操作：1
请输入短语查询：南开大学
F:\Users\lenovo\source\nk_search\search.py:147: DeprecationWarning: Received 'size' via a specific parameter in the presence of a 'body' parameter, which is deprecated and will be removed in a future version. Instead, use only 'body' or only specific parameters.
response = self.es.search(index=self.index, body=function_score_query, size=4)

===== 搜索结果 =====
Rank 1:
Title: 南开大学
URL: https://www.nankai.edu.cn/
Final Score: 0.189863
Snippet: 信息公开 图书馆 服务指南 登录邮箱 办公网 校友入校 | English 学校概况 学校简介 现任领导 历届领导 历史回眸 南开新闻 南开新闻网 南开大学报 院系机构 专业学院 职能部门 直属和后勤单位 研究机构 人才师资 人事人才 人才招聘 博士后 教育教学 本科教育 研究生教育 本科教学质量报告 远程教育 继续教育 科学研究 自然科学 社会科学 自然科学研究院 社会科学研究院 仪器设备 ...

Rank 2:
Title: 规章制度
URL: https://xcb.nankai.edu.cn/gzzd/list.htm
Final Score: 0.189863
Snippet: 导航 首 页 机构设置 部门简介 现任领导 规章制度 常用下载 联系我们 规章制度 当前位置： 首页 规章制度 规章制度 南开大学 2023 年宣传思想工作要点 2023-04-13 中共南开大学委员会关于学习宣传贯彻党的二十大精神的安排意见 2023-04-13 南开大学 2022 年宣传思想工作要点 2023-04-13 关于深入学习宣传贯彻第十次党代会精神的安排意见 2023-04-13 南 ...

Rank 3:
Title: 常用下载
URL: https://xcb.nankai.edu.cn/11347/list.htm
Final Score: 0.122589
Snippet: 导航 首 页 机构设置 部门简介 现任领导 规章制度 常用下载 联系我们 常用下载 当前位置： 首页 常用下载 常用下载 南开大学教师到校外举办讲座、报告会申请 2020-05-19 南开大学报告会、研讨会、讲座、论坛和读书会、学术沙龙申请 2020-05-19 南开大学网上直播活动备案表 2020-05-19 每页 14 记录 总共 3 记录 第一页 <<上一页 下一页>> 尾页 页码 1 / ...

Rank 4:
Title: 宣传部
URL: https://xcb.nankai.edu.cn/main.htm
Final Score: 0.100000
Snippet: 导航 首 页 机构设置 部门简介 现任领导 规章制度 常用下载 联系我们 重要论述 更多>> 习近平寄语南开师生：只有把小我融入... 17日上午，习近平在天津南开大学参观了百年校史主题展览，察看了化学学院和元素有机化... 习近平总书记视察南开大学 2019-01-17 习近平：开辟马克思主义中国化时代化新境界 2023-10-17 习近平对宣传思想文化工作作出重要指示 2023-10-10 习 ...
=====

```

图 8: 短语查询 1

```
===== 搜索引擎 =====
选择查询类型：
1. 短语查询
2. 通配查询
3. 联想关联（搜索建议）
4. 查看查询历史
5. 退出登录
输入数字选择操作：1
请输入短语查询：大学 南开
F:\Users\lenovo\source\nk_search\search.py:147: DeprecationWarning: Received 'size' via a specific parameter in the presence of a 'body' parameter, which is deprecated and will be removed in a future version. Instead, use only 'body' or only specific parameters.
response = self.es.search(index=self.index, body=function_score_query, size=4)

===== 搜索结果 =====
Rank 1:
Title: 南开大学校史网
URL: https://xs.nankai.edu.cn/
Final Score: 1.00000
Snippet: 首页 走进南开 活动要闻 纪念专题 校史研究 南开人物 校史文萃 更多 » 活动要闻 / News 教室搬上257米高空！天塔迎来首堂南开校史课 南开校史文化沙龙分享老校友口述史 “百年南开校史文化图书专架”揭幕 张元龙作专场讲座 叶嘉莹先生追思会举行 “叶先生，一路走好！”南开大学师生追念叶嘉莹先生 国际著名教育家、诗人、中国古典文学研究泰斗、南开大学讲席教授叶嘉莹先生逝世 “西南联大时期民族...

Rank 2:
Title: 欢迎光临陈省身数学研究所数学图书馆！
URL: http://www.mathlib.nankai.edu.cn/
Final Score: 0.220993
Snippet: English Version 本馆概况 入馆指南 读者守则 赔偿规则 证籍管理 读者服务 ProQuest公司赠送American Mathematical Society出版社电子书53册 德古意特出版社和普林斯顿出版社数学专著电子图书开通通知 全球科学出版社(Global Science Press)数学电子期刊库开通通知 关于规范电子资源数据库使用的通知 馆藏目录检索 到馆新书目录 中文订...

Rank 3:
Title: 欢迎关注公众号
URL: https://law.nankai.edu.cn/2023/0910/c33629a519482/page.htm
Final Score: 0.100000
Snippet: 返回首页 会议室预定 服务指南 下载专区 English 旧版入口 学院概况 学院简介 学院党委 行政领导 历任领导 专业委员会 教代会、工会 管理团队 党团建设 风采展示 工作流程 常用下载 师资队伍 在职教师 荣休教师 永远的怀念 人才培养 教学培养 毕业学位 奖励资助 常用下载 科研动态 科研机构 科研成果 科研项目 学术活动 合作交流 国际交流 招生工作 本科生 学术硕士 专业硕...

```

图 9: 短语查询 2

```

===== 搜索引擎 =====
选择查询类型:
1. 短语查询
2. 通配查询
3. 联想关联（搜索建议）
4. 查看查询历史
5. 退出登录
输入数字选择操作: 2
请输入通配符查询: 商*
===== 搜索结果 =====

Rank 1:
Title: 南开大学商学院.Nankai Business School
URL: https://ibs.nankai.edu.cn/
Final Score: 1.00000
Snippet: 新网站入口 教工登录 南开大学 English 首页 | 学院简介 | 科学研究 | 师资队伍 | 学科专业 | 教学教务 | 实验教学 | 对外交流 | 学生工作 | 教育培训 | 校友工作 | 人才招聘 第三期领导力课程研修班火热招生中 DBA教育 EMBA教育 MBA教育 MEM教育 MPM教育 MPAcc教育 MV教育 MLIS教育 高层管理教育 使命与愿景 本着传承南开历史, 突出优势和...
Rank 2:
Title: 南开商学院-企业管理
URL: https://ibs.nankai.edu.cn/teacher/bm
Final Score: 0.100000
Snippet: 新网站入口 教工登录 南开大学 English 首页 | 学院简介 | 科学研究 | 师资队伍 | 学科专业 | 教学教务 | 实验教学 | 对外交流 | 学生工作 | 教育培训 | 校友工作 | 人才招聘 师资队伍 企业管理系 会计学系 市场营销系 财务管理系 人力资源管理系 管理科学与工程系 信息资源管理系 现代管理研究所 企业管理系 崔连广 教授 电话: 23498009 邮箱: cu...
Rank 3:
Title: 南开商学院 研究中心
URL: https://ibs.nankai.edu.cn/page/researchcenter
Final Score: 0.100000
Snippet: 新网站入口 教工登录 南开大学 English 首页 | 学院简介 | 科学研究 | 师资队伍 | 学科专业 | 教学教务 | 实验教学 | 对外交流 | 学生工作 | 教育培训 | 校友工作 | 人才招聘 首页 研究中心 获奖成果 研究项目 《南开管理评论》 资源链接 资料中心 研究中心 公司治理研究院 南开大学公司治理研究院是1997年11月在南开大学国际商学院现代管理研究所的公司治理研究...

```

图 10: 通配查询



```

===== 搜索引擎 =====
选择查询类型:
1. 短语查询
2. 通配查询
3. 联想关联（搜索建议）
4. 查看查询历史
5. 退出登录
输入数字选择操作: 4

===== 查询历史 =====
2024-12-16 21:20:32 | 南开大学
2024-12-16 21:20:32 | 南开大学
2024-12-16 21:20:32 | 规章制度
2024-12-16 21:20:32 | 常用下载
2024-12-16 21:20:32 | 宣传部
2024-12-16 21:21:48 | 商*
2024-12-16 21:21:48 | 南开大学商学院 .Nankai Business School
2024-12-16 21:21:48 | 南开商学院-企业管理
2024-12-16 21:21:48 | 南开商学院 研究中心
2024-12-16 21:21:48 | 南开商学院-学院介绍
2024-12-16 21:23:41 | 校友捐赠
2024-12-16 21:23:41 | 校友信息
2024-12-16 21:23:41 | 校友信息
2024-12-16 21:23:41 | 校友活动

```

图 11: 查询日志

## 六、个性化查询

### (一) 设计思路

- 用户系统:** 设计用户注册登录功能来存储不同用户的信息以及查询记录。
- 用户历史记录的利用:** 查询的历史记录包含了用户的兴趣偏好。通过从历史记录中提取关键词 (history\_terms)，在查询过程中，给予包含这些关键词的文档更高的权重。
- 动态加载用户历史:** 为了确保个性化查询的实时性，系统需要能够根据当前用户的身份，动态加载其历史查询记录。
- 搜索结果的重新排序:** 使用查询偏好得分和归一化加权操作，结合常规的搜索得分 (如 TF-IDF、PageRank 等) 对结果进行重新排序。对于包含用户历史关键词的文档，可以通过加大其得分来使其优先显示，从而增强个性化推荐的效果。

## (二) 技术路线

设计好了用户的注册登录系统来存储不同的用户信息和查询日志。

最初尝试从查询历史中的文档提取关键词和查询词，计算查询偏好得分，但是失败了。然后想要为每个文档定一个属性，比如科技、美食、文化等，根据用户查询的历史中哪一类别的文档最多，来推测用户的查询偏好，但是效果并不理想。所以个性化查询并没有得到完全的实现。

```

1 history_terms = self.load_user_history(user)      # 加载用户历史记录
2 if history_terms:
3     # 对包含历史关键词的文档加权
4     function_score_query["query"]["function_score"]["functions"].append({
5         "filter": {
6             "terms": {
7                 "text": history_terms
8             }
9         },
10        "weight": 1.5  # 提高包含用户历史关键词的文档的权重
11    })

```

Listing 9: 个性化查询

## 七、Web 界面

为了展示功能和方便查询以及个性化操作，设计了初始页面、登录页面、查询页面、搜索结果页面、退出页面等简单可视化页面。未展示过的页面如下所示：



```

○ PS F:\Users\lenovo\source\nk_search> python search.py
成功连接到Elasticsearch.

===== 欢迎使用南开大学搜索引擎 =====

请选择操作：
1. 注册
2. 登录
3. 退出
输入数字选择操作：1

```

图 12: 初始页面



```

===== 欢迎使用南开大学搜索引擎 =====

请选择操作：
1. 注册
2. 登录
3. 退出
输入数字选择操作：1

===== 注册 =====
请输入用户名: 1216
请输入密码:
请确认密码:
用户 '1216' 注册成功!

请选择操作：
1. 注册
2. 登录
3. 退出
输入数字选择操作：2

===== 登录 =====
请输入用户名: 1216
请输入密码:
用户 '1216' 登录成功!

===== 搜索引擎 =====

```

图 13: 注册登录页面

```

===== 搜索引擎 =====
选择查询类型:
1. 短语查询
2. 通配查询
3. 联想关联(搜索建议)
4. 查看查询历史
5. 退出登录
输入数字选择操作: 5
用户 '1216' 已登出。

请选择操作:
1. 注册
2. 登录
3. 退出
输入数字选择操作: 3
退出程序。

```

图 14: 退出页面

## 八、个性化推荐

### (一) 设计思路

- 联想建议的生成:** 通过标题中的前几个词来构造一个简单的联想短语。通过分割标题字符串并取前两个词作为建议。如果标题仅包含一个词，则该词会作为建议的一部分。同时为了避免建议过长，使用省略号 (...) 来截断建议短语，确保联想建议简洁明了。
- suggest 属性:** 添加一个新的属性 suggest，将得到的联想建议存入映射中。
- Elastic Search 联想功能:** 使用 Elastic Search 的 completion suggester 提供基于前缀的模糊匹配，返回与用户输入前缀相关的关键词或短语建议。

### (二) 技术路线

使用 Elastic Search 提供的 completion suggester API，根据前缀（如标题的部分文字）获取联想建议。然后根据已索引的文档和建议的字段进行模糊匹配，返回最相关的建议项。最后通过 wildcard\_suggest 方法查询并返回相关的建议。

```

1 def wildcard_suggest(self, prefix):
2     # 使用 Elasticsearch 的 completion suggester
3     suggest = {
4         "suggest": {
5             "simple_phrase": {
6                 "prefix": prefix,
7                 "completion": {
8                     "field": "suggest",
9                     "fuzzy": {
10                         "fuzziness": "AUTO"
11                     },
12                     "size": 5
13                 }
14             }
15         }
16     }
17     try:
18         response = self.es.search(index=self.index, body=suggest)
19         suggestions = response.get('suggest', {}).get('simple_phrase', [])
20         if suggestions:

```

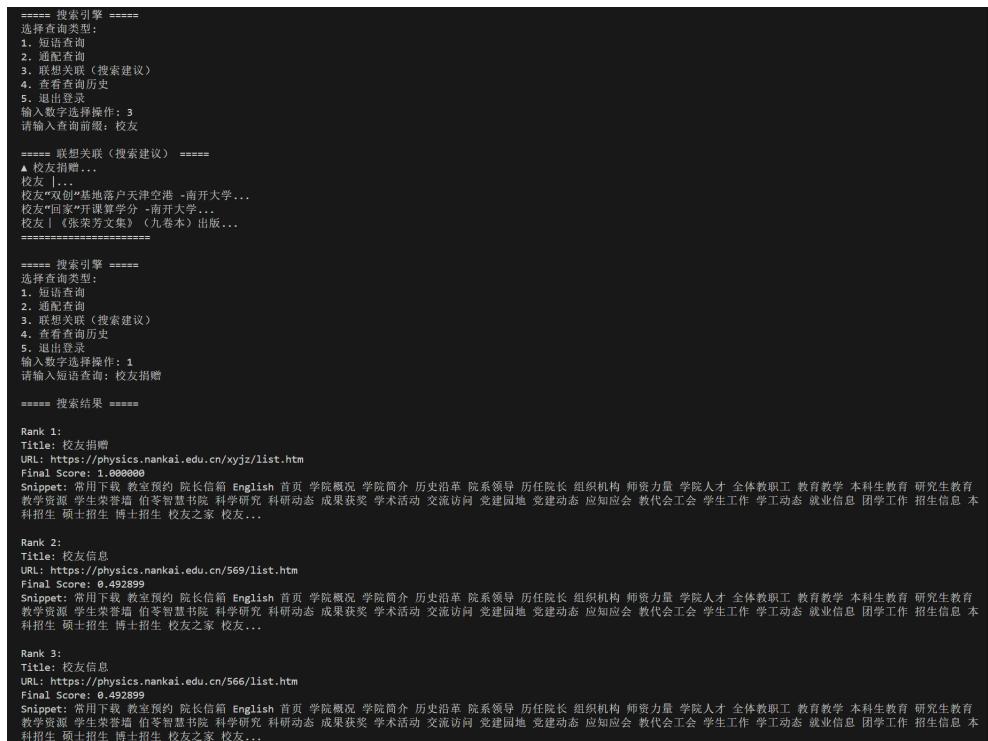
```

21     options = suggestions[0].get('options', [])
22     return [sugg['text'] for sugg in options]
23
24     return []
25 except Exception as e:
26     print(f"获取建议时发生错误: {e}")
27     return []

```

Listing 10: 联想关联（搜索建议）

### (三) 完成结果



```

===== 搜索引擎 =====
选择查询类型:
1. 短语查询
2. 通配查询
3. 联想关联 (搜索建议)
4. 查看查询历史
5. 退出登录
输入数字选择操作: 3
请输入查询前缀: 校友

===== 联想关联 (搜索建议) =====
▲ 校友捐赠...
校友 [...]
校友“双一流”基地落户天津空港 -南开大学...
校友“回家”开课助学分 -南开大学...
校友《张荣芳文集》(九卷本)出版...

===== 搜索引擎 =====
选择查询类型:
1. 短语查询
2. 通配查询
3. 联想关联 (搜索建议)
4. 查看查询历史
5. 退出登录
输入数字选择操作: 1
请输入短语查询: 校友捐赠

===== 搜索结果 =====

Rank 1:
Title: 校友捐赠
URL: https://physics.nankai.edu.cn/xyjz/list.htm
Final Score: 1.00000
Snippet: 常用下载 教室预约 院长信箱 English 首页 学院概况 学院简介 历史沿革 院系领导 历任院长 组织机构 师资力量 学院人才 全体教职工 教育教学 本科教育 研究生教育 教学资源 学生荣誉墙 伯苓智慧书院 科学研究 科研动态 成果获奖 学术活动 交流访问 党建园地 党建动态 应知应会 教代会工会 学生工作 学工动态 就业信息 团学工作 招生信息 本科招生 硕士招生 博士招生 校友之家 校友...

Rank 2:
Title: 校友信息
URL: https://physics.nankai.edu.cn/569/list.htm
Final Score: 0.49289
Snippet: 常用下载 教室预约 院长信箱 English 首页 学院概况 学院简介 历史沿革 院系领导 历任院长 组织机构 师资力量 学院人才 全体教职工 教育教学 本科教育 研究生教育 教学资源 学生荣誉墙 伯苓智慧书院 科学研究 科研动态 成果获奖 学术活动 交流访问 党建园地 党建动态 应知应会 教代会工会 学生工作 学工动态 就业信息 团学工作 招生信息 本科招生 硕士招生 博士招生 校友之家 校友...

Rank 3:
Title: 校友信息
URL: https://physics.nankai.edu.cn/566/list.htm
Final Score: 0.49289
Snippet: 常用下载 教室预约 院长信箱 English 首页 学院概况 学院简介 历史沿革 院系领导 历任院长 组织机构 师资力量 学院人才 全体教职工 教育教学 本科教育 研究生教育 教学资源 学生荣誉墙 伯苓智慧书院 科学研究 科研动态 成果获奖 学术活动 交流访问 党建园地 党建动态 应知应会 教代会工会 学生工作 学工动态 就业信息 团学工作 招生信息 本科招生 硕士招生 博士招生 校友之家 校友...

```

图 15: 联想关联（搜索建议）

## 九、总结与感悟

历时 8 天，终于完成了本次系统性的实验。每一步都充满挑战，尤其是爬虫爬取数据阶段，在刚开始时存在很多疑虑，10 万条数据怎么爬？要爬多久？爬虫怎么这么慢？爬下来的数据后面能用吗？前前后后爬了好几次，每次一爬就是 8、9 个小时，甚至第一版爬虫爬了 24 小时，还不敢干别的，生怕爬虫会断，真的是“万事开头难”。这次试验遇到了很多问题，也做过很多尝试，不断地试错、改错以及调整方向。

当然，通过应用信息检索系统原理课程所学知识，设计和实现这个 Web 搜索引擎，收获颇丰。我深刻理解了搜索引擎的工作原理及其背后的技术细节，尤其是索引构建和链接分析。这次实验还增强了我对海量数据的处理能力，也掌握了强大工具 Elastic Search 的使用方法，锻炼了我实际项目开发与问题解决的能力，可能也会让我以后在查询我想要的信息时能够更加精准。