

电工导 C 大作业报告

陈一鸣 黄浩初 罗宇晗 章凌恺

January 2024

目录

1	任务概览	2
2	实现细节和核心代码解释	2
2.1	爬虫实现	2
2.1.1	淘宝爬虫	2
2.1.2	京东爬虫	6
2.2	索引建立及情感分析	9
2.2.1	搜索引擎	9
2.2.2	评论内容情感分析	13
2.3	前端搭建	15
2.3.1	主页设计	15
2.3.2	品牌分类页	15
2.3.3	衣服种类页	16
2.3.4	搜索结果页	17
2.3.5	搜索结果页	17
2.3.6	以图搜图结果展示页面	18
2.3.7	chatbot 聊天返回结果界面	19
2.3.8	渲染界面	20
2.3.9	两个 style 函数的编写	22
2.4	图像识别	22
2.4.1	根据商品图片进行检索	22
2.4.2	基于 logo 的图像检索	23
2.5	其他创新功能实现	25
2.5.1	AI 智能服装推荐助手	25
3	可以改进之处	26
4	项目分工	27

1 任务概览

爬虫部分任务概览：本次爬虫使用了 selenium 进行爬取，并采用 BeautifulSoup 解析页面，根据网页的结构，得到我们需要的信息。但需要注意，由于存在登录，验证码等网页拦截的问题，我们需要采取一些反爬措施。

索引部分的任务概览：使用 Lucene 创建索引 + 全文检索即可对爬虫部分的内容建立对应的索引。

前端部分的任务概览：利用 flask 框架，css 等渲染，得到最终的搜索，结果展示等界面。其中包含关键词搜索，品牌名称限制以及以图搜图等相关的服装搜索，结果展示中包含根据评论内容打分进行的排序，根据价格排序，根据综合相关程度进行排序的结果。

图像识别模块任务概览：根据我们已有的知识，完成以图搜图有以下方式选择：提取颜色直方图，SIFT 图像描述子匹配，pyTorch 深度学习提取图像特征后 CNN 算法匹配，LSH 算法预处理检索。在我们小组自制的搜索引擎中，以图搜索的部分有根据商品图片检索，基于 logo 的图像检索，因此我们需要权衡各算法利弊，找到最合适的算法完成以上任务。

评论情感分析及 AI 助手部分任务概览：情感分析模块中使用 cnssenti 库进行中文情感分析。将所有的评论分析情感倾向，以实现根据评论内容对商品进行筛选。AI 助手部分概览：使用 gpt4all 库，利用本地部署的 gpt4all falcon LLM 模型进行服装推荐。为用户提供更多元化，更方便的购物体验。

2 实现细节和核心代码解释

2.1 爬虫实现

我们爬取的电商网站包括京东，淘宝，天猫，爬取到的基本商品信息包括商品价格，名称，图片链接，商品链接，品牌的信息。由于发现各大电商网站的评论一般和商品价格等属性所在的页面有所不同，因而我们采用的策略是，先爬取对应的基本商品信息，再根据我们爬取到的商品信息的 url，爬取对应的商品评论信息。

爬取商品信息我们采用的是 selenium 用来模拟用户的行为，并采集到对应的商品的信息。由于淘宝和京东等电商网站的结构不相同，并且反爬机制也有不同，因此我们对不同电商网站采取不同的爬取策略。

2.1.1 淘宝爬虫

淘宝网站的反爬虫策略主要在于滑块验证码和登录操作，因此我们写了一段代码，模拟登录

```
1 def login_to_taobao(driver, username, password):
2     try:
3         # 使用 Selenium 打开登录页面
4         login_url = 'https://login.taobao.com/member/login.jhtml'
5         driver.get(login_url)
6
7         # 输入用户名和密码
8         username_input = driver.find_element('id', 'fm-login-id')
9         username_input.send_keys(username)
10
11         password_input = driver.find_element('id', 'fm-login-password')
12         password_input.send_keys(password)
13
14         # 提交登录表单
15         password_input.send_keys(Keys.RETURN)
16
17         # 等待一段时间，确保登录成功
18         time.sleep(5)
```

```

19
20     # 处理滑块验证
21     handle_slider_verification(driver)
22
23 except Exception as e:
24     print(f'Error during login: {e}')

```

后面发现登录还需要处理滑块验证，我们写了如下滑块验证代码，以保证成功登录：

```

1 def handle_slider_verification(driver):
2     try:
3         switch_to_slider_iframe(driver)
4
5         # 使用WebDriverWait等待滑块元素可见
6         slider = WebDriverWait
7         (driver, 10).until(EC.visibility_of_element_located
8         ((By.XPATH, '//*[@id="nc_1_n1z"]'))))
9
10        # 创建ActionChains对象
11        actions = ActionChains(driver)
12
13        # 将鼠标移动到滑块
14        actions.move_to_element(slider).perform()
15
16        # 在滑块上按下左键
17        actions.click_and_hold().perform()
18
19        # 拖动鼠标
20        actions.move_by_offset(300, 0).perform()
21
22        # 释放鼠标
23        actions.release().perform()
24
25        # 切回主文档
26        driver.switch_to.default_content()
27
28    except Exception as e:
29        print(f'Error during slider verification: {e}')

```

后来发现这两个方法淘宝反爬虫的机制十分严格，这样子解决验证码的操作会验证失败。因此最终我们采用的是手动登录，发现这样子也可以实现，只要显式登录的过程中我们登录自己的账号，就不会出现这样的问题，因此我们关闭无头模式，进行登录，能够成功爬到对应的商品信息。核心的获取相应商品信息的代码如下：

```

1 def scrape_page(driver):
2     # 模拟滚动，加载更多的商品
3     for i in range(5):
4         driver.execute_script('window.scrollTo(0,
5         document.body.scrollHeight);')

```

```

6         time.sleep(2) # 等待加载完成，以应对反爬虫
7
8     # 获取完整的页面内容
9     page_source = driver.page_source
10
11    # 使用 BeautifulSoup 解析页面内容
12    soup = BeautifulSoup(page_source, 'html.parser')
13
14    # 提取标题、图片和超链接信息
15    items = soup.find_all('div', class_='Card—mainPicAndDesc—wvcDXaK')
16
17    # 使用列表存储标题、图片和超链接的对应关系
18    title_image_url_list = []
19
20    for item in items:
21        title_text = item.find('div', class_='Title—title—jCOPvpf')
22        .find('span').text.strip() if item
23        .find('div', class_='Title—title—jCOPvpf') else ''
24        image_url = item.find('div',
25                               class_='MainPic—mainPicWrapper—iv9Yv90')
26        .find('img')['src'].strip() if item
27        .find('div', class_='MainPic—mainPicWrapper—iv9Yv90') and item
28        .find('div', class_='MainPic—mainPicWrapper—iv9Yv90')
29        .find('img') else ''
30        href = item
31        .find_previous('a', class_='Card—doubleCardWrapper—L2XFE73')['href']
32        .strip()
33        if
34        item.find_previous('a',
35                           class_='Card—doubleCardWrapper—L2XFE73') else ''
36        price = item.find('div', class_='Price—priceWrapper—Q0Dn7pN')
37        .find('span', class_='Price—unit—VNGKLAP').text.strip()
38        +item.find('div', class_='Price—priceWrapper—Q0Dn7pN')
39        .find('span', class_='Price—priceInt—ZlsSi_M').text.strip()
40        +item.find('div', class_='Price—priceWrapper—Q0Dn7pN')
41        .find('span', class_='Price—priceFloat—h2RR0RK').text.strip()
42        href = urllib.parse.urljoin(url, href)
43    # 存储到列表
44    title_image_url_list.append((title_text, image_url, price, href))
45
46    return title_image_url_list

```

加入了页面的滚动，以适应淘宝中动态加载的机制。之后的一些代码可以参考 crawler_scroll.py 中的代码，主要是实现了不断下一页的爬取，以获得更多的商品信息。最终爬取到的结果如下：

Title: Theory西装连衣裙明星同款2023秋冬新品羊毛混纺双排扣长袖收腰领
Image URL: https://img.alicdn.com/imgextra/i4/129970823/O1CN01Ru7ypD1HwyqO4LYdk !I0-saturn solar.jpg 460x460q90.jpg .webp
Brand: MANGO
Price: ¥538.00
Product URL: https://click.simba.taobao.com/cc_im?p=MANGO&s=504358029&k=10698&e=m8%2B7wtKXcj%2FZrgL90a1DGXpIS8XKwvlyXfkCyGbiMtBRdUTbPjgS%2FwuvpOVNikKjw%2B3A43hl5uuFHg5EttlDwA3bvAQgkomg2czn6qSt%2Fb6JPQMqY5IK%2FIUDjd9Hgc757N0ojhMKmoGBLB2AhkFDJFz29TnealeQjB%2Furky1ZzPvVC8CUGK%2FbO4MG6hzM0ktCoGPNJYakytSgRwgTdYnB910Kl15H8xGAhOtTjoO1iHn0%2F6Lam7%2BYUuUewODCDkuF6VagRXnLuQSRuAYVQJPFjvMkx0tRFKco1zSpicmT3IYf1WZMhNaniyAjb3dh%2BbOcotn%2F96tEQ9QknikVtexclOwo%2BwZ4BmU9n7DF%2FedBOXRen4JeeKTOuuBgpoRb5m5c3pIAX6yk35XurCRabTmHoHoQcrrtpvbC%2FnWEdvAPRFQSMqaa4JjVmMq%2FnqAKKO73XLTu5cP8vTdiEZg2D%2F96Ta159YBJ7F%2BCC6DlVRiVweJxCP3bTwewPukgl5Pg1cs8dq5GaEjnf6AlwV4A3pA%2B9zs%2BcLSOhtXB6osufqBxGZ9MkEFeXmaYSLya0V8tNyXlal209sx8lVADTUpPRBkQyy1l8wkrW4OITQNFrbw8ZmODv1Xhu1Aj%2Bw6UObIgnWdykHUExQPBvNBlaKe0oQzsE5TfAjNstPZ2z%2BdhZbplPaaqfdBWmiXx%2FfrVvgusTm4Z05kiHEfN8sY4panpRdKWEJ%2BLCMMt6lx2bstu8dFeroYgbS3%2BtcSKDZcV5BjYgmdn48pl1ZOzkcLL9uSLWRfW8bbxwPHETDTst4KKZ3gQpWQNRiUfj5rOXWKKNFuM0OseULpCulpe79lJDTmwcVeijUW3qZMSXRfC3Mvrw0%2BzD6tlhVMI3x89LqWnukW2v%2B6RnBFvdWIK2YilznCCUEMVzLVAmW4#B79lIC2gSDH%2FfpZ8AxM8JErGwjAWfvZMHWcYf8RcnsqQAFBfr%2BTb8PoWQ0GsFDQpckOc28ilc1L1LGZGfe%2FR%2FwyxG8BuAYsU2v2dvIccyTbF40Y%2FUGg8q8lfrLTWBX0yJlCQA%3D#detail

淘宝爬取评论内容是一个比较麻烦的事情，我们采取的方式是根据页面检查中的 network 模块，得到我们登录信息的 cookies 以及反爬虫所需要的 headers 和 parameters，利用以上三者信息，以及 requests.get 方法访问淘宝中 comment 所在的文件，即可得到对应的 json 文件内容，再根据网页中 comment 所在的位置，利用 json 文件的索引功能，找到对应的 comment 内容。具体代码结构如下：

```
1 cookies = {...} #从网页检查窗口中 network 位置得到
2 headers = {...}
3 param = {...}
4 response = requests.
5 get('https://rate.taobao.com/feedRateList.htm', params=params,
6 cookies=cookies, headers=headers)
7 # print(response.text)
8 json_match = re.search(r'\((.*)\)', response.text)
9 if json_match:
10     json_data = json_match.group(1)
11     # 将提取到的 JSON 数据转换为 Python 对象
12     json_data = json.loads(json_data)
13 # json_data = response.json()
14 print(json_data)
15 for com in json_data['comments']:
16     content = com['content']
17     href = 'https://item.taobao.com/item.htm?id=706407719891&ns=
18     1&abbucket=9'
19     validname = valid_filename(href)
20     with open(f'comment/{validname}.txt', 'a', encoding='utf-8') as file:
21         file.write(f'{content}\n')
22 else:
23     print("No JSON data found in the response.")
```

代码中前面一部分用于处理 requests.get 到的内容中非 json 文件的部分，我们发现得到的内容会包含如同 jsonp[.] 的格式，不能利用这个格式的内容得到评论内容，因此我们利用正则表达式处理掉这部分内容，使我们能够成功爬取到对应评论。这么做是可行的，并且通过改变 id 的方式，我们可以爬取到更多的商品评论。爬取的结果如下：

我个人很喜欢网购，所以每次都要写很长的评论，一条一条的写太浪费时间，因此想到了这个偷懒的方法，当您看到我的评价的时候，证明我对这个商品很满意，喜欢的亲可以购买，这个商品符合性价比，味道好，质量好的。
真的自发热啊，柔软舒服，贴身却没有束缚感。弹性大，所以S和M穿上没觉不同。
软嘴，舒适。袖子比较长适合打底
快递把这件衣服弄丢了，又找不到，跟客服联系了好久才补发。
高领很保暖但不扎人，很贴合身体，也很百搭，第二次购买了。
衣服没有以前厚，领子也比以前的高一点。
优衣库的物流我超爱
虽然薄但是确实非常温暖，很好
超好穿，贴身很温暖，弹力大
静电太大，物超所值
166 100 拍的s 超舒服 好穿
好看舒服
买了两件。
超级舒服。
等着降价，花都谢了。穿着很舒服的感觉
物美价廉
颜色超好看
此用户没有填写评价。
情绪管理才是最好的化妆品。不快乐比任何劳累都要消耗人。情绪一旦崩塌，即使你貌若天仙，外貌和气质被拖垮也就是瞬间的事情。谦卑和感恩是解决一切问题的万灵丹。了解自己、看清小我，进而在生活中磨练，让自己更加自在、解脱、快乐。而感恩，正是你用谦卑的心去体会一切之后，自然而然发生的。感恩会带来更多的谦卑、更多的福分、更多的快乐，这样就形成了一个非常好的良性循环谦卑和感恩是解决一切问题的万灵丹。了解自己、看清小我，进而在生活中磨练，让自己更加自在、解脱、快乐。而感恩，正是你用谦卑的心去体会一切之后，自然而然发生的。感恩会带来更多的谦卑、更多的福分、更多的快乐，这样就形成了一个非常好的良性循环。
太长了。。。。。。也不保准

但我们发现这个方式爬取到一定量的评论之后，并不能继续爬取了，似乎淘宝识别到了我的爬虫行为，并且禁止我的账号访问评论内容，这给我们带来了不小的冲击。因此在爬取到了一定数量的评论之后，我们尝试去京东网站爬取对应的信息。

2.1.2 京东爬虫

京东爬虫与淘宝爬虫类似，我们采用 selenium 进行京东网站的爬取，需要注意的是，我们需要将自动检测关闭，以消除反爬虫的影响：

```
1 chrome_driver_path =
2 'C:\\Users\\25458\\Downloads\\chromedriver-win64
3 \\chromedriver-win64\\chromedriver.exe'#这个路径是我chrome_drive
4 r 所在的路径
5 chrome_options = ChromeOptions()
6 # 这一部分将自动检测关闭，以保证爬虫的正常进行
7 chrome_options.add_argument('--disable-blink-features=Automation
8 \\Controlled')
9 chrome_options.add_argument('--disable-extensions')
10 chrome_options.add_argument('--disable-dev-shm-usage')
11 chrome_options.add_argument('--disable-gpu')
12 chrome_options.add_argument('--no-sandbox')
13 chrome_options.add_experimental_option('excludeSwitches',
14
15 ['enable-automation'])
16 chrome_options.add_experimental_option('useAutomationExtension', False)
17
18 driver = webdriver.Chrome(service=Service(chrome_driver_path),
19
20 options=chrome_options)
```

其他的包括翻页，网页结构分析得到对应的内容都与淘宝类似，这里不再赘述。

京东爬虫的评论获取：我们惊喜地发现，我们可以直接利用 selenium 访问商品链接，进而得到对应的商品评论信息，由于评论信息较多，爬取速度较慢，因此我们采取多线程爬取的形式，同时注意我们的爬虫文件中包含有 get_url.py 和 delete_txt.py 文件，作用分别是从小txt文件中提取出商品url到新的txt文件中，便于进行评论爬取，同时 delete_txt.py 的作用就是将已经爬取过的网页在文件夹中进行删除。防止重复爬取浪费时间。

```
1 with ThreadPoolExecutor(max_workers=1) as executor:
2     # 启动线程
```

```

3     threads = []
4     for _ in range(5):
5         thread = threading.Thread(target=worker)
6         thread.start()
7         threads.append(thread)
8
9     # 等待队列中的所有任务完成
10    url_queue.join()
11
12    # 停止线程
13    for _ in range(5):
14        url_queue.put(None)
15    for thread in threads:
16        thread.join()

```

每一个线程对应工作的 worker 函数如下：

```

1 def worker():
2     driver = webdriver.Chrome(service=Service(chrome_driver_path),
3
4     options=chrome_options)
5     while True:
6         try:
7             url = url_queue.get()
8             if url is None:
9                 break
10            driver.get(url)
11            get_fav_rate(driver, url)
12            scrape_page(driver, url)
13            for page in range(2, 11): # 限制爬取的页数
14                try:
15                    next_page_button = driver
16                        .find_element(By.CLASS_NAME, 'ui-pager-next')
17                    ActionChains(driver).move_to_element(next_page_button)
18                    .perform()
19                    time.sleep(2)
20                    next_page_button.click()
21                    scrape_page(driver, url)
22                except Exception as err:
23                    continue
24            url_queue.task_done()
25        except:
26            continue

```

其中加入了 time.sleep() 等函数防止反爬虫机制的检测。其中 get_fav_rate 函数和 scrape_page 函数分别实现了获得好评率和得到评论内容的功能。具体的实现如下：

```

1 # 得到好评率信息

```

```

2 def get_fav_rate(driver, url):
3     for i in range(2):
4         driver.execute_script('window.scrollTo(0,
5         document.body.scrollHeight);')
6         time.sleep(5)
7         page_source = driver.page_source
8         soup = BeautifulSoup(page_source, 'html.parser')
9         element = WebDriverWait(driver, 15).until(
10            EC.visibility_of_element_located((By.CLASS_NAME, "percent-con")))
11        )
12        fav_rate = soup.find('div', class_='percent-con')
13        validname = valid_filename(url)
14        with open(f'comment/{validname}.txt', 'w', encoding='utf-8') as file:
15            file.write(f'好评率: {fav_rate.text}\n')
16
17 # 爬取一页中所有的评论信息
18 def scrape_page(driver, url):
19     for i in range(1):
20         driver.execute_script('window.scrollTo(0,
21         document.body.scrollHeight);')
22         time.sleep(3)
23
24         element = WebDriverWait(driver, 15).until(
25            EC.visibility_of_element_located((By.CLASS_NAME, "comment-con")))
26        )
27
28        page_source = driver.page_source
29        soup = BeautifulSoup(page_source, 'html.parser')
30        comments = soup.find_all('p', class_='comment-con')
31
32        for comment in comments:
33            validname = valid_filename(url)
34            with open(f'comment/{validname}.txt', 'a', encoding='utf-8') as file:
35                file.write(f'{comment.text.strip()}\n')

```

分别采用了也采用了 `time.sleep()` 函数防止反爬虫，同时根据对应的标签的信息得到对应的内容。并且写入以对应的 url 为名的文件中。最终爬取的评论结果如下：

```

好评率: 92%
衣服很合身, 挺满意的.
衣服很好, 码偏小一点, 各位亲注意喽!
布料质量好, 穿着也舒服, 值得推荐
站在消费者的角度说: 卖家粗心大意, 有缺陷的衣服也出货。(有图有真相, 像素不好莫见怪)
尺码大小: 有点大了, 还好能穿
衣服偏短.....
不怎么样
一般吧
质量可以, 穿起好帅

```


2.2 索引建立及情感分析

2.2.1 搜索引擎

使用 Lucene 创建索引 + 全文检索。通过将搜索得到的每个商品的各种信息打包成 dict 进行排序。排序方法包括相关度, 好评率, 综合评价, 价格升序和价格降序。其中综合评价是在相关度的基础上, 通过好评、差评数量占比进行修正, 从而得到更加全面的评价分数。将检索到的信息排序好再打包发送到页面。对于 Lucene 中文分词性能欠佳的问题, 我们使用了 jieba 中文分词库, 利用 Lucene 总是会依照空格分词, 在两个中文词汇之间插入空格, 以提高 Lucene 的中文分词性能。而这会造成在进行 site-search 时, Lucene 识别进行 site-search 的命令会被分词库切割的问题。比如: 输入: img_url: https://... 分词库会把 img_url: 这个本该被 Lucene 整体识别的指令切割, 导致 Lucene 报错。所以我们采用对搜索结果进行过滤的方式, 实现针对品牌等属性的搜索。具体建立索引代码实现如下:

```
1 def indexDocs(self, root, writer):
2     t1 = FieldType()
3     t1.setStored(True)
4     t1.setTokenized(False)
5     t1.setIndexOptions(IndexOptions.NONE) # Not Indexed
6
7     t2 = FieldType()
8     t2.setStored(False)
9     t2.setTokenized(True)
10    # Indexes documents, frequencies and positions.
11    t2.setIndexOptions(IndexOptions.DOCS_AND_FREQS_AND_POSITIONS)
12
13
14    for root, dirnames, filenames in os.walk(root):
15        for filename in filenames:
16            if not filename.endswith('.txt'):
17                continue
18            print("adding", filename)
19            path = os.path.join(root, filename)
20            file = open(path, encoding='utf-8', errors='ignore')
21            contents = file.read()
22            file.close()
23            doc = Document()
24            doc.add(StringField("name", filename, Field.Store.YES))
25            doc.add(StringField("path", path, Field.Store.YES))
26            if len(contents) > 0:
27                title = self.getTxtAttribute(contents, 'Title')
28                title = jieba.cut(title)
29                title = " ".join(title) # put space between
30                segmented words
31                doc.add(TextField('title', title, Field.Store.YES))
32                # title of the product
33                brand = self.getTxtAttribute(contents, 'Brand')
34                brand = jieba.cut(brand)
35                brand = " ".join(brand) # put space between
36                segmented words
```

```

37         doc.add(TextField('brand', brand, Field.Store.YES))
38         # brand of the product
39         img_url = self.getTxtAttribute(contents, 'Image_URL')
40         doc.add(TextField('img_url', img_url, Field.Store.YES))
41         # image url of the product
42         price = self.getTxtAttribute(contents, 'Price')
43         doc.add(TextField('price', price, Field.Store.YES))
44         # price of the product
45         prd_url = self.getTxtAttribute(contents, 'Product_URL')
46         doc.add(TextField('prd_url', prd_url, Field.Store.YES))
47         # product url of the product
48         pos = self.getTxtAttribute(contents, 'pos')
49         doc.add(TextField('pos', pos, Field.Store.YES))
50         # number of positive comments
51         neg = self.getTxtAttribute(contents, 'neg')
52         doc.add(TextField('neg', neg, Field.Store.YES))
53         # number of negative comments
54         neu = self.getTxtAttribute(contents, 'neu')
55         doc.add(TextField('neu', neu, Field.Store.YES))
56         # number of neutral comments
57         pos_rate = self.getTxtAttribute(contents, 'pos_rate')
58         doc.add(TextField('pos_rate', pos_rate, Field.Store.YES))
59         # positive rate of the product
60         neg_rate = self.getTxtAttribute(contents, 'neg_rate')
61         doc.add(TextField('neg_rate', neg_rate, Field.Store.YES))
62         # negative rate of the product
63     else:
64         print("warning: no content in %s" % filename)
65     writer.addDocument(doc)

```

有关于搜索的代码如下:

```

1  def filter_by_brand(products, brand_range):
2      # products: a list of products
3      # brand_range: a list of brands that the user wants to search
4      # filter the products by the brands the user wants to search
5
6      if brand_range == []:
7          return products
8      filtered_products = []
9      for product in products:
10         for brand in brand_range:
11             if brand == product['brand']:
12                 filtered_products.append(product)
13                 break
14     return filtered_products
15

```

```

16 def filter_by_price(products, price_range):
17     # products: a list of products
18     # price_range: a list of two numbers, the price range
19     # filter the products by the price range
20
21     if price_range == []:
22         return products
23     filtered_products = []
24     for product in products:
25         if float(product['price']) >= price_range[0]
26         and float(product['price']) <= price_range[1]:
27             filtered_products.append(product)
28     return filtered_products
29
30 def get_comprehensive_score(score, pos, neg, neu, penalty = 3):
31     # score: search score
32     # pos: positive comment count
33     # neg: negative comment count
34     # neu: neutral comment count
35     # penalty: penalty per negative comment
36     # calculate the comprehensive score of a product
37
38     if pos + neg + neu == 0:
39         return score
40     comment_score = float(pos - neg * penalty) / (pos + neg + neu)
41     comprehensive_score = score * (1 + comment_score)
42     return comprehensive_score
43 # this is only a half-guessing method to evaluate the product
44 # you can change the penalty to get a better result
45
46 def sort(products, sort_method = 1):
47     # sort_method: 1. search score
48     #                2. pos_rate
49     #                3. comprehensive_score
50     #                4. price up
51     #                5. price down
52     # sort the result by the sort method
53
54     if sort_method == 1:
55         products.sort(key = lambda x: x['search_score'], reverse = True)
56     elif sort_method == 2:
57         products.sort(key = lambda x: x['pos_rate'], reverse = True)
58     elif sort_method == 3:
59         products.sort(key = lambda x: x['comprehensive_score']
60                       , reverse = True)
61     elif sort_method == 4:

```

```

62         products.sort(key = lambda x: x['price'], reverse = False)
63     elif sort_method == 5:
64         products.sort(key = lambda x: x['price'], reverse = True)
65     else:
66         print('Error: sort method not found!')
67         return
68     return products
69
70 def run(searcher, analyzer, command):
71     # while True:
72     if deletespace(command) == '':
73         print('Error: empty command!')
74         return
75
76     products = []
77     score_threshold = 0 # threshold for score of each token
78
79     print()
80     print("Searching for:", command)
81
82     command = jieba.cut(command)
83     command = list(command)
84     cmd_tokens = len(command) # number of tokens in command
85     command = " ".join(command) # put space between segmented words
86     query = QueryParser("title", analyzer).parse(command)
87     scoreDocs = searcher.search(query, 50000).scoreDocs
88     print("%s total matching documents." % len(scoreDocs))
89
90     for i, scoreDoc in enumerate(scoreDocs):
91         product = {}
92         if scoreDoc.score < score_threshold * cmd_tokens:
93             continue
94         doc = searcher.doc(scoreDoc.doc)
95         product['title'] = deletespace(doc.get("title"))
96         product['prd_url'] = doc.get("prd_url")
97         product['price'] = doc.get("price").rstrip('\n')
98         product['img_url'] = doc.get("img_url").rstrip('\n')
99         product['brand'] = doc.get("brand").rstrip('\n')
100        product['file_path'] = doc.get("path").rstrip('\n')
101        product['search_score'] = scoreDoc.score
102        product['pos_rate'] = doc.get("pos_rate").rstrip('\n')
103        product['comprehensive_score'] =
104        get_comprehensive_score(scoreDoc.score,
105        int(doc.get("pos").rstrip('\n')),
106
107        int(doc.get("neg").rstrip('\n')),

```

```

108
109         int(doc.get("neu").rstrip('\n')))
110         print('title:', deletespace(doc.get("title")), 'score:',
111             scoreDoc.score)
112         products.append(product)
113         # print 'explain:', searcher.explain(query, scoreDoc.doc)
114     print("Search finished!")
115
116     # prd_url = prodcut url
117     # img_url = image url
118     # there's no sales volume cuz we didn't crawl it
119
120     # there're 4 ways to sort the result
121     # 1. sort by search score
122     # 2. sort by pos_rate
123     # 3. sort by comprehensive_score
124     # 4. sort by price
125     products = sort(products, 3)
126
127     return products # all return in a list, each element is a dict

```

这其中关于核心代码的解释已经放在了注释当中，这里不再赘述。

2.2.2 评论内容情感分析

这一部分我们使用 `cnssenti` 库进行中文情感分析。将所有的评论分析情感倾向，如果 `positive > negative` 则将其识别为好评。类似地，将所有评论识别为 `positive`, `negative`, `neutral` 中的一种之后，按照 `positive` 和 `negative` 占比计算好评率和差评率。将结果写入原始数据以便于建立索引。由于从网页上爬取的评论包含大量的 `unicode` 表情符号，所以在进行情感分析之前，对评论文件进行了清洗，使用 `emoji` 库，遍历所有文件，去除了所有的表情符号。由于原理较为简单，在这里不做赘述。具体代码实现如下：

```

1  def analysis(comment_path, result_path):
2      # args: comment_path: the path of the comment files
3      #         result_path: the path of the analysis result files
4      # process all the files in the comment_path, and write the emo-analysis
5
6      # result to the result_path
7
8      for root, dirs, files in os.walk(comment_path):
9          for file in files:
10             print("processing file: ", file)
11             com_path = os.path.join(root, file)
12             res_path = os.path.join(result_path, file)
13             com_file_lines = open(com_path, 'r', encoding='utf-8').readlines()
14             res_file = open(res_path, 'w', encoding='utf-8')
15             # open the comment file and the result file
16
17             for c in com_file_lines:

```

```

18         print(">>>_" + c)
19         emot = senti.sentiment_calculate(c)
20         if emot['pos'] > emot['neg']: # positive
21             emot = 1
22         elif emot['pos'] < emot['neg']: # negative
23             emot = -1
24         else: # neutral
25             emot = 0
26         res_file.write(str(emot) + '\n')
27         print(emot, "_written_to_file")
28         # write the analysis result to the result file
29         in the form of 1, -1, 0
30
31 def conclude(analysis_path, result_path):
32     # args: analysis_path: the path of the analysis result files
33     #       result_path: the path of the conclusion files
34     # process all the files in the analysis_path, and write the
35     # conclusion to the result_path
36
37     for root, dirs, files in os.walk(analysis_path):
38         for file in files:
39             print("processing_file:", file)
40             com_path = os.path.join(root, file)
41             com_file_lines = open(com_path, 'r', encoding='utf-8').readlines()
42             # open the analysis result file
43
44             pos = 0
45             neg = 0
46             neu = 0
47             for c in com_file_lines:
48                 if c == '1\n':
49                     pos += 1
50                 elif c == '-1\n':
51                     neg += 1
52                 else:
53                     neu += 1
54             print("pos:", pos, "neg:", neg, "neu:", neu)
55             # count the number of positive, negative and neutral comments
56
57             res_path = os.path.join(result_path, file)
58             res_file = open(res_path, 'w', encoding='utf-8')
59             res_file.write("pos:" + str(pos) + '\n')
60             res_file.write("neg:" + str(neg) + '\n')
61             res_file.write("neu:" + str(neu) + '\n')
62             total = pos + neg + neu
63             # write the number of positive, negative and neutral

```

```

64         # comments to the conclusion file
65
66         if total != 0:
67             pos_rate = float(pos) / (pos + neg + neu)
68             neg_rate = float(neg) / (pos + neg + neu)
69         else:
70             pos_rate = 0
71             neg_rate = 0
72         res_file.write("pos_rate:␣" + str(pos_rate) + '\n')
73         res_file.write("neg_rate:␣" + str(neg_rate) + '\n')
74         # write the rate of positive and negative comments
75         # to the conclusion file
76
77         res_file.close()

```

2.3 前端搭建

这一部分主要进行前端部分实现的展示以及核心代码的呈现和解释。

2.3.1 主页设计

Taobao goods Search.html 实现了主页的设计，使用了 style.css 作为样式，核心代码的实现如下：

```

1  <div id="header">
2      <div id="menu">
3          <ul>
4              <li class="current_page_item"><a href="/home">Home</a></li>
5              <li><a href="/about">Brand</a></li>
6              <li><a href="/contact">Kinds</a></li>
7          </ul>
8      </div>
9      <div id="search">
10         <form id="searchform" method="get" action="/result">
11             <fieldset>
12                 <input id="s" type="text" name="keyword" value="" class="text">
13                 <input id="x" type="submit" value="Search" class="button">
14             </fieldset>
15         </form>
16     </div>
17 </div>
18 <div id="logo">
19 </div>

```

2.3.2 品牌分类页

about.html 为品牌分类页，提供了 logo 搜索，以图搜图，部分品牌介绍等功能，使用了 style2.css 作为样式，具体前端展示代码如下：

```

<div id="pageBrand">
  <div id="contentBrand">
    <div class="post">
      <h1 class="title">Brand</h1>
      <div class="entryBrand">
        <h1>logo搜索</h1>
        <form method="POST" action="/logo" enctype="multipart/form-data">
          <input type="file" name="image" accept="image/*">
          <input type="submit" value="上传图片">
        </form>
        <h1>以图搜图</h1>
        <form method="POST" action="/pic" enctype="multipart/form-data">
          <input type="file" name="image" accept="image/*">
          <input type="submit" value="上传图片">
        </form>
        <h1 class="name">常见品牌推荐</h1>
        <ul>
          <li>
            <h2><a href="{{ url_for('result', keyword='李宁') }}">李宁</a></h2>
            <p>李宁是“体操王子”李宁在1990年创立的专业体育品牌。
              李宁公司拥有完善的品牌营销、研发、设计、制造、经销及销售能力，以经营李宁品牌专业及休闲运动鞋、服装、器材
            </p>
          </li>
          <li>
            <h2><a href="{{ url_for('result', keyword='ANTA') }}">ANTA</a></h2>
            <p>安踏始创于一九九一年，海安踏体育用品有限公司（股份代号：2020（港币柜台）及82020（人民币柜台））在二零零
          </p>
          </li>
        </ul>
      </div>
    </div>
  </div>
</div>

```

2.3.3 衣服种类页

contact.html 为衣服种类页，提供衣服穿搭 gpt 推荐，按衣服适用人群场合，具体分类搜索等功能，使用了 style2.css 作为样式，具体实现代码如下：

```

<div id="pageBrand">
  <div id="contentBrand">
    <div class="post">
      <h1 class="title">Kinds</h1>
      <div id="search-g">
        <form id="searchform" method="get" action="/gpt">
          <fieldset>
            <input id="s-g" type="text" name="question" value="" class="text">
            <input id="x-g" type="submit" value="Ask GPT For Advice" class="button">
          </fieldset>
        </form>
      </div>
      <h1>Kinds of Clothes</h1>
      <div class="kinds-container">
        <!-- Contact info start -->
        <div class="kind-item">
          <a href="{{ url_for('result', keyword='男') }}">
            
          </a>
        </div>
        <div class="kind-item">
          <a href="{{ url_for('result', keyword='女') }}">
            
          </a>
        </div>
        <div class="kind-item">
          <a href="{{ url_for('result', keyword='童') }}">
            
          </a>
        </div>
      </div>
    </div>
  </div>
</div>

```


2.3.4 搜索结果页

search.html 为搜索结果页, 采用白色背景以避免信息过多影响阅读, 提供排序方式选择, 价格区间, 品牌筛选, 分页展示等多种搜索结果筛选形式, 实现了以图片为主要呈现形式的商品展示方式, 使用了 style2.css 作为样式, 核心代码如下:

```
1 <div class="product-container">
2     {% for i in range(0,length)%}
3     <div class="product-item">
4         <a href="{{u[i]}}" target="_blank">
5             
6         </a>
7         <p>商品名称<p>
8         <p><strong>{{t[i]}}</strong></p>
9         <p>品牌&nbsp;;&nbsp;{{brand[i]}}</p>
10        <p>好评率&nbsp;;&nbsp;{{pos[i]}}</p>
11        <p class="price">综合评价指数&nbsp;;&nbsp;{{com[i]}}</p>
12        <p class="price">价格&nbsp;;&nbsp;{{p[i]}}</p>
13    </div>
14    {%endfor%}
15</div>
```

2.3.5 搜索结果页

logoSearch.html 为 logo 搜索结果页, 提供 logo 搜索结果与相关的商品, 以相对商品指数进行排序, 使用了 style2.css 作为样式, 具体代码如下:

```
1 <div id="content">
2     <div class="post">
3         <h1 class="title">Search Results</h1>
4         <h1 class="name">upload logo might be {{brand[0]}}</h1>
5         <div class="entry">
6             <h2 style="color:purple;">Reupload your logo picture
7             </h2>
8             <form method="POST" action="/logo"
9             enctype="multipart/form-data">
10                <input type="file" name="image"
11                accept="image/*">
12                <input type="submit" value="上传图片">
13            </form>
14            <h2 style="color:purple;">
15            Displaying {{length}} results for your search:</h2>
16            <!-- List of search results start -->
17            <div class="product-container">
18                {% for i in range(0,length)%}
19                <div class="product-item">
20                    <a href="{{u[i]}}" target="_blank">
21                        
```


[illegible]

2.3.7 chatbot 聊天返回结果界面

gpt.html 为 Chatbot 返回结果页，提供 bot 推荐的五个搜索关键词，并实现一键搜索功能，使用 style2.css 作为样式，核心代码实现如下：

```

1 <div id="contentBrand">
2     <div class="post">
3         <h1 class="title">Kinds</h1>
4         <div class="kinds-contrainer">
5             <!-- Contact info start -->
6             <h1>GPT's recommend for your question</h1>
7             <ul>
8                 <li><a href="{ {url_for('result',
9 keyword=command[0]) }}" style="font-size:
10 2em;" target="_blank">{{command[0]}}</a></li>
11                 <li><a href="{ {url_for('result',
12 keyword=command[1]) }}" style="font-size:
13 2em;" target="_blank">{{command[1]}}</a></li>
14                 <li><a href="{ {url_for('result',
15 keyword=command[2]) }}" style="font-size:
16 2em;" target="_blank">{{command[2]}}</a></li>
17                 <li><a href="{ {url_for('result',
18 keyword=command[3]) }}" style="font-size:
19 2em;" target="_blank">{{command[3]}}</a></li>
20                 <li><a href="{ {url_for('result',
21 keyword=command[4]) }}" style="font-size:
22 2em;" target="_blank">{{command[4]}}</a></li>
23             </ul>
24             <!-- Contact info end -->
25         </div>
26     </div>
27 </div>

```

2.3.8 渲染界面

app.py 主要是使用 flask 进行渲染, 将搜索函数返回的结果进行进一步处理之后通过变量传递进入 html 之中进行进一步处理, 具体代码如下:

主页代码:

```
1 @app.route('/home', methods=['POST', 'GET'])
2 def bio_data_form():
3     if request.method == "POST":
4         keyword = request.form['keyword']
5         return redirect(url_for('result', keyword=keyword, method = 3))
6 return render_template("Taobao_goods_Search.html")
```

搜索页:

```
1 @app.route('/result', methods=['POST', 'GET'])
2 def result():
3     t = []
4     u = []
5     p = []
6     img = []
7     b = []
8     com = []
9     pos = []
10    name = ['搜索指数', '好评率', '综合评价指数', '价格']
11
12    STORE_DIR = "result/index"
13    jieba.load_userdict("userdict.txt")
14    print('lucene', lucene.VERSION)
15    #base_dir = os.path.dirname(os.path.abspath(sys.argv[0]))
16    vm_env.attachCurrentThread()
17    directory = SimpleFSDirectory(File(STORE_DIR).toPath())
18    searcher = IndexSearcher(DirectoryReader.open(directory))
19    # set a new similarity computing method
20    searcher.setSimilarity(SimpleSimilarity()) # simplesimilarity func
21    analyzer = StandardAnalyzer()#Version.LUCENE_CURRENT)
22    if 'keyword' in request.args:
23        keyword = request.args.get('keyword')
24    elif 'keyword' in request.form:
25        keyword = request.form.get('keyword')
26    else:
27        return redirect(url_for('bio_data_form'))
28
29    method = int(request.args.get('method', 3))
30    selected_brands = []
31    if 'brands' in request.form:
32        selected_brands = request.form.getlist('brands')
33    print(selected_brands)
```

```

34 price = []
35 if 'price1' in request.args:
36     price1 = int(request.args.get('price1'))
37     price2 = int(request.args.get('price2'))
38     price.append(price1)
39     price.append(price2)
40 if keyword == '':
41     return redirect(url_for('bio_data_form'))
42 keyword = ' '.join(jieba.cut(keyword))
43
44 products = run(searcher, analyzer, keyword, method)
45 products = filter_by_price(products, price)
46 products = filter_by_brand(products, selected_brands)
47 for i in range(len(products)):
48     t.append(products[i]['title'])
49     u.append(products[i]['prd_url'])
50     p.append(products[i]['price'])
51     img.append(products[i]['img_url'])
52     b.append(products[i]['brand'])
53     com.append(round(products[i]['comprehensive_score'], 4))
54     pos.append(round(float(products[i]['pos_rate']), 4))
55 page = int(request.args.get('page', 1))
56 # 当前页，默认为第一页面
57 per_page = 50
58 # 每页显示的结果数量，可以自己设定
59 start = (page - 1) * per_page
60 end = start + per_page
61 paged_results = t[start:end] # 取当前页面的结果
62 paged_urls = u[start:end]
63 paged_prices = p[start:end]
64
65 del searcher
66 length = min(len(t), len(u))
67 length = min(length, len(p))
68 total_results = length
69 total_pages = (total_results + per_page - 1) // per_page
70 current_t = len(paged_results)
71
72 return render_template("search.html", keyword=keyword,
73 length=current_t, img = img, brand = b, price = price,
74 u=paged_urls, t=paged_results,
75 p=paged_prices, name = name, com = com, pos = pos,
76 page=page, total_pages=total_pages,
77 total_results = total_results, method = method)

```

2.3.9 两个 style 函数的编写

对页面各部分元素的字体大小，背景颜色，子元素间距等等进行了规范，对交互精细度进行了美化，比如搜索按钮在 hover 时会有放大效果，图片链接在鼠标悬停时会缩放，并显示提示框提示点击后下一步的商品名或者商品种类名，在开源网站设计了背景图和项目主页图，设计了 slogan。

2.4 图像识别

2.4.1 根据商品图片进行检索

基本思路：考虑到实际情况中用户输入商品图片的目的是检索同类商品，我们排除了提取颜色直方图算法以及对颜色匹配敏感的 LSH 算法，而 CNN 算法需要建立我们已有一万余张图片的图像特征库，这不符合搜索引擎的速度要求。因此我们选择 SIFT 算法，并利用 threading 多线程方法加快搜索速度。

代码实现：search_by_picture.py 主要功能分为三个函数，url_to_img 函数用于打开爬虫数据库中图片的 URL，读入图片信息并转化为三维或二维的数组形式。由于存在 URL 网址无法打开的情况，我们加入异常处理机制，防止程序报错中止。

```
1 def url_to_img(img_url):
2     try:
3         resp = urllib.request.urlopen(img_url)
4         image = np.asarray(bytearray(resp.read()), dtype="uint8")
5         image = cv2.imdecode(image, cv2.IMREAD_COLOR)
6         return image
7     except:
8         return np.array([])
```

process_folder 函数用于对已打开的文件夹遍历爬到的 txt 文件数据，并提取爬虫数据库中图片的关键点和描述符，将其与用户输入的商品图片进行比较匹配，对于满足比例要求的好的匹配点计数，当数量达到一定阈值时，可以认为该图片与用户输入图片相似度高，并用好的匹配点个数模拟相似度评分返回结果。这一函数单独存在的意义是方便各个线程单独独立工作。核心代码如下

```
1 kp2, des2 = sift.detectAndCompute(goods, None)
2 matches = flann.knnMatch(des1, des2, k=2)
3 good_matches = [m for m, n in matches if m.distance < 0.75*n.distance]
4 if (len(good_matches) >= 20):
5     with lock:
6         product = getProductInfo(data, filename, len(good_matches)/40)
7         products.append(product)
8         results.append(products)
9         if len(results) >= 50:
10             stop_event.set()
```

picture_detection 函数用于遍历文件夹中的各个子文件夹，创建子线程对它们完成上述 process_folder 函数，其中也完成了创建 FLANN 匹配器帮助匹配的功能。核心代码如下

```
1 # 创建FLANN匹配器
2 FLANN_INDEX_KDTREE = 1
3 index_params = dict(algorithm = FLANN_INDEX_KDTREE, trees = 5)
4 search_params = dict(checks=50)
5 flann = cv2.FlannBasedMatcher(index_params, search_params)
```

```

6     # 遍历文件夹中各商品图片
7     results = []
8     lock = threading.Lock()
9     stop_event = threading.Event()
10    threads = []
11    for brand in os.listdir(folder_path):
12        thread = threading.Thread(target=process_folder,
13                                  args=(brand, folder_path, img, sift, flann, des1, results,
14                                        lock, stop_event))
15        threads.append(thread)
16        thread.start()
17    for thread in threads:
18        thread.join()

```

2.4.2 基于 logo 的图像检索

基本思路：同任务 2.1，考虑到实际情况中用户输入 logo 图片的目的是限定商品类别，我们排除了提取颜色直方图算法以及对颜色匹配敏感的 LSH 算法。并且由于商品图片中 logo 部分所占比例小且模糊，与用户输入的 logo 图片匹配精度低，因此我们选择从网上爬取我们已有品牌数据的 12 张 logo 图片作比较。比较次数缩小到 12 次，我们采用的是 CNN 算法和 SIFT 算法双重评价标准，得到结果的精确度得到显著提升。

代码实现：search_by_logo.py 主要部分分为一个类和三个函数，完成准备工作的 getProductInfo 函数实现提取匹配好的商品信息的功能。部分代码如下

```

1  def getProductInfo(data, filename, score):
2      product = dict()
3      pos, neg, neu = 0, 0, 0
4      for content in data:
5          if(re.match("Title:", content)):
6              content = content.lstrip("Title:").rstrip('\n')
7              product["title"] = content
8              continue
9          if(re.match("Product_URL:", content)):
10             content = content.lstrip("Product_URL").rstrip('\n')
11             product["prd_url"] = content
12             continue

```

MyCBIR 类通过已有的训练模型提取图像特征并比较两张图片的相似度。在相似度计算方面，我们采用的方法是计算两特征向量夹角的余弦值。核心代码如下

```

1  # 图像预处理
2  def preprocess(self, img):
3      input_image = default_loader(img)
4      input_image = trans(input_image)
5      input_image = torch.unsqueeze(input_image, 0)
6      #print(input_image)
7      return input_image
8

```

```

9      # 抽取图像特征
10     def features(self, x):
11         x = model.conv1(x)
12         x = model.bn1(x)
13         x = model.relu(x)
14         x = model.maxpool(x)
15         x = model.layer1(x)
16         x = model.layer2(x)
17         x = model.layer3(x)
18         x = model.layer4(x)
19         x = model.avgpool(x)
20         return x
21
22     # 计算相似度
23     def compare(self, other):
24         v1 = self.features
25         v2 = other.features
26         normalized_v1 = v1 / np.linalg.norm(v1)
27         normalized_v2 = v2 / np.linalg.norm(v2)
28         return 1 - np.sqrt(np.sum((normalized_v1-normalized_v2)**2))

```

MySIFT 函数实现思路与任务 2.1 相同，不加以赘述。

logo_detection 为两种算法的整合函数，用于遍历文件夹中的各个子文件夹，对每个品牌子文件夹中的 logo 图片分别用上述两种算法同用户输入的 logo 图片进行匹配，找到最佳匹配品牌。在两种算法结合的方式上，我们考虑的是 CNN 相似度评分 +SIFT 好的匹配点个数/1000 作为最终评分的机制。这样的匹配结果经多次实验证明效果较好。核心代码如下

```

1  # 遍历文件夹中各品牌的 logo 图片
2  for brand in os.listdir(folder_path):
3      match_score[brand] = 1
4      for filename in os.listdir(os.path.join(folder_path, brand)):
5          if (filename.endswith(".jpg") or filename.endswith(".png")):
6              path = os.path.join(folder_path, brand, filename)
7              print("Processing {}".format(filename))
8
9              # 用CBIR算法计算相似度
10             logo_tool = MyCBIR(path)
11             match_score[brand] *= target_tool.compare(logo_tool)
12             # print(match_score)
13
14             # 用SIFT算法计算相似度
15             match_score[brand] += MySIFT(target, path)/10
16             # print(match_score)
17
18     # 对比匹配结果，找到最合适匹配
19     match_score = sorted(match_score.items(), key=lambda x : x[1],
20                           reverse=True)

```


2.5 其他创新功能实现

2.5.1 AI 智能服装推荐助手

使用 gpt4all 库，利用本地部署的 gpt4all falcon LLM 模型进行服装推荐。由于在建立页面内聊天窗功能时遇到技术困难，我们最后选择针对用户输入的需求，进行单次推荐。通过适当的 prompt，让模型输出符合需求的 5 种服装类别，然后再分别进行搜索，达到推荐的效果。GPT4ALL Falcon 下载地址：https://gpt4all.io/models/gguf/gpt4all-falcon-newbpe-q4_0.gguf 由于本地部署的 gpt 性能有限，中文语言能力较差，因此使用了 translate 库进行输入和输出的翻译，以实现更好的效果。i/o 样例：

```
>>> 我想买一件耐克的T恤，最好是速干的
感谢您的申请。 根据您的要求，以下五个关键词可能符合您的需求：

1. 快干T恤
2. Dri-FIT T恤
3. 性能T恤
4. 透气T恤
5. 轻便T恤
['快干T恤', 'Dri-FIT T恤', '性能T恤', '透气T恤', '轻便T恤']
```

返回一个包含 5 个建议的列表，然后搜索就可以得到推荐的商品。由于容器没有安装 GLIBCXX_3.4.26 库，导致无法在容器内运行 gpt，且经过多次尝试无法安装，所以我们采用文件 i/o 的方式，通过在容器外提前运行 gpt，通过读取和写入 txt 文件来进行输入输出，既解决了容器内无法运行 gpt 的问题，又有效规避了 gpt 预热带来的延迟。人工智能 bot (gpt) 提供购买建议的实现代码如下：

```
1 gpt 建议生成代码
2 def chat(model, hints, max_tokens=200, temp=0):
3     # args: model: the gpt4all model
4     #       hints: the prompts
5     #       max_tokens: the max tokens of the output, higher
6     # max_tokens means longer output
7     #       temp: the temperature of the output,
8     # higher temp means more randomness
9     # return: the attributions of the request
10
11     CACHE_FILE_PATH = "container/samples4/Group_Task/cache/chat_input.txt"
12     OUTPUT_FILE_PATH = "container/samples4/Group_Task/cache/chat_output.txt"
13     ATTRS_FILE_PATH = "container/samples4/Group_Task/cache/chat_attrs.txt"
14     # cache file. read from it pre 1 second. if has content, use it
15     as input, and clear it.
16
17     with model.chat_session():
18         # prompts
19         for hint in hints:
20             outp = model.generate(hint, max_tokens, n_batch=256, temp=temp)
21             file = open(CACHE_FILE_PATH, 'w', encoding='utf-8')
22             file.write('')
23             file.close()
24             # pre-clear the cache file
```

```

25
26     print("这里是人工智能服装推荐系统，您可以输入您的需求，
27     我将为您推荐最适合您的服装。")
28
29     while(True):
30         file = open(CACHE_FILE_PATH, 'r', encoding='utf-8',
31                     errors='ignore')
32         output_file = open(OUTPUT_FILE_PATH, 'w', encoding='utf-8')
33         attrs_file = open(ATTRS_FILE_PATH, 'w', encoding='utf-8')
34         command = file.read()
35         while(command == ''):
36             sleep(1)
37             print("waiting for input...")
38             file = open(CACHE_FILE_PATH, 'r', encoding='utf-8',
39                         errors='ignore')
40             command = file.read()
41         # detect contents in the input cache file
42
43         print("command: " + command)
44         hint = command
45         file = open(CACHE_FILE_PATH, 'w', encoding='utf-8')
46         file.write('')
47         file.close()
48         if hint == 'exit':
49             break
50
51         outp = model.generate('request: ' + translate(hint,
52                                                         src='zh', tgt='en'), max_tokens, n_batch=256, temp=temp)
53         # llm generate suggestions
54         outp = translate(outp)
55         output_file.write(outp + '\n')
56         attributions = search_attributions(outp + '\n')
57         attrs_file.writelines(attributions)
58         # write the output and attributions to the cache files
59
60         print(outp)
61         print(attributions)
62     print("chat session ended!")

```

具体的代码解析在注释中已经比较清楚，总体上实现了人机的交互。

3 可以改进之处

我们开发此次项目用了不少的时间，但仍有许多可以改进的地方：

1. 本次爬虫爬取时忘记爬取销量这一重要信息，并且由于电商网站网页不断地变更，我们的代码可能不一定一直适用，因而如果想要爬取更多有效的信息，应该爬取更多的服装网站信息，并且增强代码的鲁棒性，稳定性，以适应

网页更多的变化。

2. 本次项目没有引入数据库工作，如果可以的话，我们应该进一步实现根据用户的搜索记录来得到更加优化的搜索结果展示。

3. 前端部分前期经验不足，没有充分利用 javascript，后面开始学习想要使用的时候已经写好的代码量已经救不回去了，框架已定，导致许多功能实现比较复杂。因此后续可以通过 javascript 对页面进行优化, 对代码进行简化等操作。

4.style 设计较为杂乱，查询函数较为困难，容易出现素材复用难度大的问题。因此后续应该进行素材的高效利用以及整合包装。

4 项目分工

陈一鸣 (522030910055): 爬虫实现及网页异常处理

黄浩初 (522030910076): 前端搭建及对接工作

罗宇晗 (522030910023): 以图搜图及 logo 匹配工作

章凌恺 (522030910056): 索引建立以及创新功能实现