

Search Engine Presentation

ZiHeng Wang

Directory

- ① Method/Optimization
- ② Some of the Methods I Explored
- ③ Actual Display
- ④ Summary

Title

Method/Optimization

Crawler

Crawlers are the foundation of search engines.

For this big assignment, my crawler is written in python2.7 and uses libraries like urlparse and bs4 to simplify the difficulties.

Crawler

Specific implementation:

From the very beginning of the school of Information website, each URL is traversed using BFS.

For each url traversed, use wget for crawling.

Crawler

Details to note:

- ① Need to judge all run out of the information institute web page.
- ② Remove duplicate sites(like index.php and the URL with '#').
- ③ Delete illegal links (such as download links, email and inaccessible urls).

Word Processing

Word processing is also a very important part of search engines.

It is also written in python2.7 and uses libraries such as jieba and HTMLParser to simplify the difficulties.

Word Processing

Specific implementation:

After you normalize the web page (which you can do using BeautifulSoup), get all the text through `handle_data(HTMLParser)`, and then cut the words with `jieba.cut`.

Word Processing

Details to note:

① Stop words need to be removed after word cutting.

② It is necessary to record whether each piece of text is necessary to obtain (for example, endnotes and other contents, we do not need to obtain).

TF - IDF Algorithm

Tf-idf algorithm is the core of search engine. It is basically the same as what is taught in class.

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \times \log_{10}(N / \text{df}_t)$$

TF - IDF Algorithm

I made some simple optimizations:

- ① The elements with large IDF are weighted.
- ② When there are too many query words, delete words which have small IDF.
- ③ Add weight when the current word is time.
- ④ It is assumed that the words in the title are more important.

Optimization can ensure the search speed and improve the accuracy

Web

Complete by using python's flask library.

Bootstrap was also used to make the HTML more beautiful.

It is a pity that I can't do better because of the time and my poor skill.

Web

这是一个清爽的搜索软件

Web

这是一个清爽的搜索软件

数据工程与知识工程教育部重点实验室

搜索

共展示100条结果

数据工程与知识工程教育部重点实验室 - 机构设置 - 中国人民大学信息学院

...数据工程与知识工程教育部重点实验室成立于2006年，是“985工程”二期“数据工程与知识工程”科技创新平台承担单位，中国人民大学建设的第一个省部级重点实验室。现任主任为教育部科技进步一等奖获得者杜小勇教授。实验室以中国人民大学的学科优势为依托，...

来源：http://info.ruc.edu.cn/overview_structure_dept.php?id=13

数据工程与知识工程研究所 - 机构设置 - 中国人民大学信息学院

...中国人民大学数据工程与知识工程研究所是一个集教学、科研与应用开发为一体的高科技学术机构。担负着计算机科学与技术专业博士点的建设任务。自1987年建所以来，始终站在学科前沿，跟踪国际先进技术，承担了大量国家重点研究项目，积极开展对外...

来源：http://info.ruc.edu.cn/overview_structure_dept.php?id=17

数据库与商务智能教育部工程研究中心 - 机构设置 - 中国人民大学信息学院

...商务智能教育部工程研究中心成立于2001年，致力于数据库与商务智能新技术研究、产品研发以及应用推广。目前，该中心成立数据库与商务智能工程研究中心研究实验室、研发中心、测试中心以及办公室，并以人大金仓公司市场运作实体成立了数据库...

来源：http://info.ruc.edu.cn/overview_structure_dept.php?id=15

Title

Some of the Methods I Explored

WAND

A Safe algorithm.

However, I have tried and found that the length of crawling web articles is usually short, and there are not many documents, so using the WAND algorithm is still likely to make all documents need to accurately calculate their vectors, or even slower than the plain one, so it is abandoned.

(This also took some time to study, so sad

Title

Actual Display

Title

Summary

Summary

In this is not a very short week, from the beginning of nothing, to now although the search engine has done, but the results are still so-so.

It is also a pity that I failed to finish many things I wanted to do because Of my weak basic skills and lack of knowledge in this field. Python was also the first time I came into contact with, so I spent a lot of time in searching materials and programming.

Summary

However, this course is at least a start. I learned a lot of knowledge and tried to put it into practice. I encountered many setbacks and had the help of teachers, teaching assistants and classmates.

In a word, I still have to work hard in the future.

THANK YOU FOR YOUR
LISTENING!