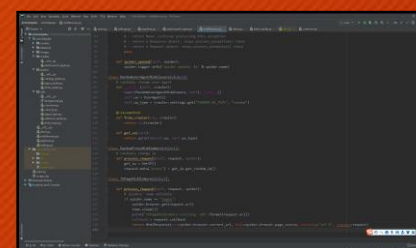
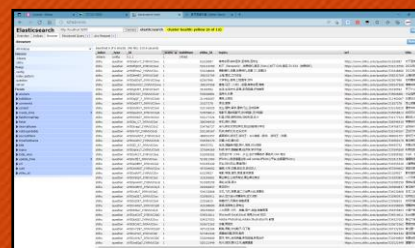
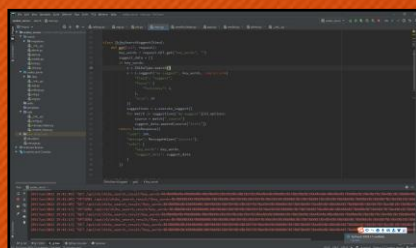
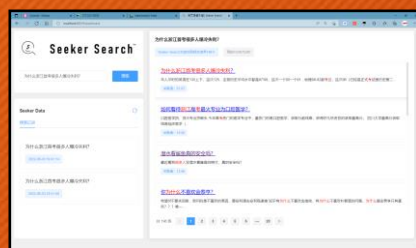




搜索引擎系统最佳实践

S
e
e
k
e
r
r
c
h

A Best Practice on Search Engine System



大数据应用开发课程设计作品

- 团队：The Seekers（付东源、谭斯谦）
- 汇报人：付东源



项目概述

项目定义

- 一个实现**搜索引擎**基本功能的软件系统

实现目标

- **数据爬取**
自动化登录、自动化爬虫脚本，目标网站**一键增量爬取**
- **索引入库**
将爬取数据按**索引**分类，并进一步**处理**，存入**数据仓库**
- **用户检索**
发送**HTTP请求**至项目服务器，请求**搜索关键词**相关数据
- **结果呈现**
开发用户易于接受的**网页应用**，简洁明了**展示搜索结果**



搜索引擎

根据用户需求与一定算法，运用特定策略从互联网检索出指定信息反馈给用户的一门检索技术。

核心模块



爬虫



索引



检索



排序

依托技术

网络爬虫

网页处理

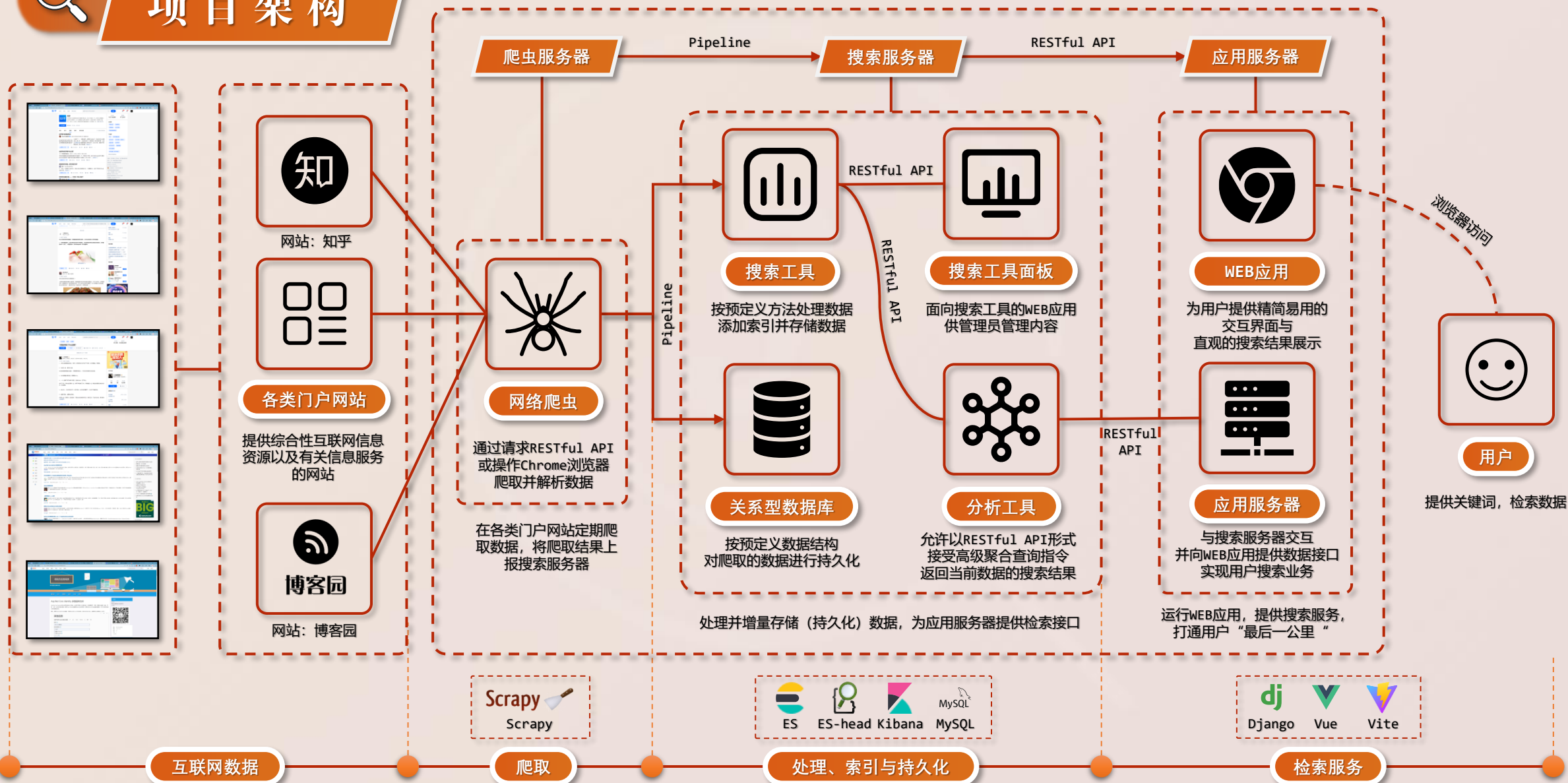
检索排序

NLP

...



项目架构



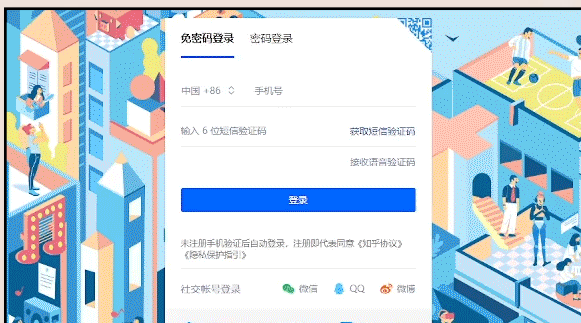


数据爬取

Scrapy



模拟登录模块



项目基于cv2库, 实现知乎网站爬虫
拼图验证码识别与自动拖动验证功能,
另集成百度智能云在线识别验证码功能。

基于HTML解析的爬取方式

使用ChromeDriver模拟浏览器操作, 获取HTML文档, 再通过**CSS选择器**定位并获取所需内容。所见即所得, 实现简单, 但效率较低。

```
item_loader_question.add_css(
    "topics", ".QuestionHeader-
    topics .Popover div::text")
```

知乎: 基于HTML解析进行爬取
zhihu_spider.py 139:9

基于AJAX请求的爬取方式

分析网站代码, 提取**AJAX**请求, 获取网站Cookie与信息id列表后, 直接**Request**请求获取数据。效率较高, 但有时涉及逆向, 难度较大。

```
Request(url=parse.urljoin(res
    ponse.url,
    '/NewsAjax/GetAjaxNewsInfo?co
    ntentId={}'.format(post_id)))
```

博客园: 基于AJAX请求进行爬取
cnblogs_spider.py 94:19



索引入库



MySQL

MySQL

Crawled Data Item Object

网站数据模型实例

yield



ElasticSearchPipeline

ElasticSearch数据流处理器

MySQLTwistedPipeline

MySQL数据流处理器

JSONExporterPipeline

数据流JSON文件生成器

ArticleImagePipeline

图片批量下载器

成员变量

Item初始化时，需要提供该数据对象所需的必要字段，生成数据Model，以进行后续操作。

```
zhihu_id = scrapy.Field()
topics = scrapy.Field()
url = scrapy.Field()
title = scrapy.Field()
content = scrapy.Field()
create_time = scrapy.Field()
update_time = scrapy.Field()
answers = scrapy.Field()
comments = scrapy.Field()
total_view = scrapy.Field()
clicks = scrapy.Field()
```

Item成员变量 items.py 98:5

基类方法

所有Item都应覆写的方法，包括存储至ES、生成DML语句等。

get_insert_sql

方法：生成SQL语句

save_to_es

方法：存储至ElasticSearch



数据入库

```
insert_sql = "insert into zhihu_question(zhihu_id, ...)
VALUES (%s, ...) ON DUPLICATE KEY UPDATE content =
VALUES(content)..."
```

DML-INSERT模板语句 items.py 111:9

```
params = (zhihu_id, topics, url, title, content,
answers, comments, total_view)
```

参数表 items.py 121:9

提供DML-INSERT模板语句与参数表，通过pymysql数据库驱动，完成MySQL数据增量入库。



数据索引

```
class ZhihuType(DocType):
    suggest = Completion(analyzer=ik_analyzer)
    topics = Text(analyzer="ik_max_word")
```

ElasticSearch索引类，使用最大分词analyzer
elasticsearch_types.py 17:1

```
if __name__ ==
 "__main__":
    ZhihuType.init()
```

创建ES索引
elasticsearch_types.py 35:1

```
def save_to_es(self):
    zhihu = ZhihuType()
    zhihu.save()
```

Item基类方法
items.py 125:5

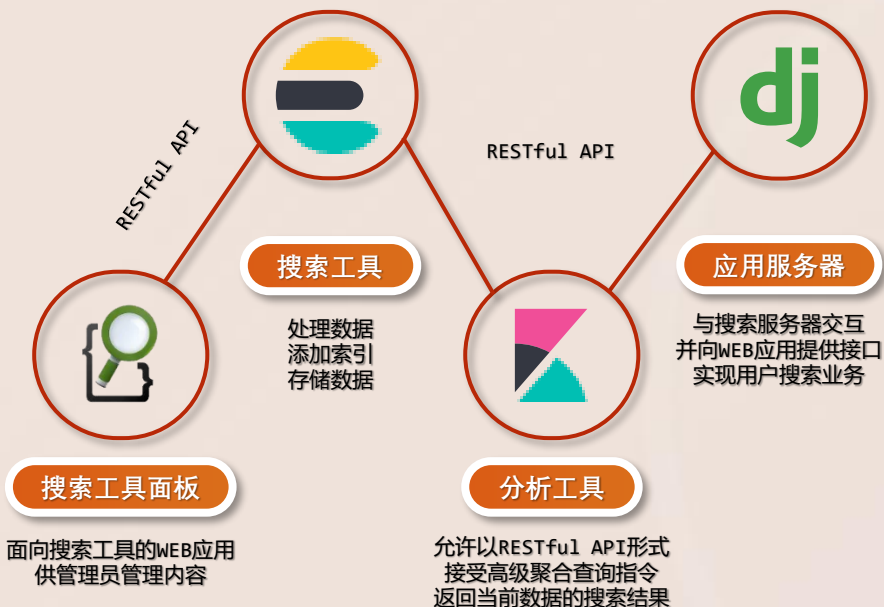
```
class ElasticSearchPipeline(object):
    def process_item(self, item, spider):
        item.save_to_es()
        return item
```

ElasticSearch数据流处理器 pipelines.py 90:1

- 通过设置settings.py中的ITEM_PIPELINES属性，接入指定的pipeline，灵活切换数据处理方式，结合yield，实现并行处理（MySQL+ElasticSearch）
- 设计索引类型，规定索引字段与analyzer，运行init()方法，在ES中创建索引，此后ElasticSearchPipeline调用Item.save_to_es()，完成该条数据的处理、索引



用户检索

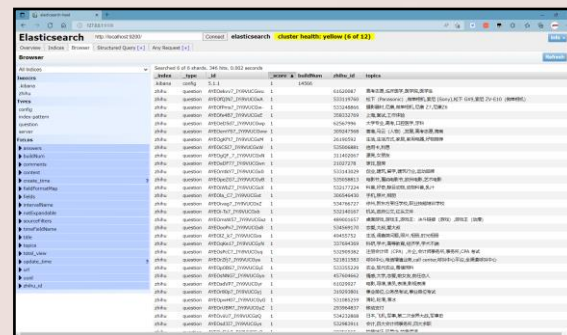


搜索工具: Elasticsearch

```
cd <BASE_DIR>/elastic_search/bin elasticsearch.bat
```

通过cmd命令启动
ElasticSearch

- 以服务器形式运行在9200端口，为项目搜索核心功能提供支撑
- 集成 ik_max_word、ik_smart 等 analyzer（分析器），可对文本进行分词等分析处理
- 命中搜索结果局部高亮（Highlight）



ElasticSearch-Head: ES管理控制台 (Port 9300)



分析工具: Kibana

```
cd <BASE_DIR>/kibana/bin kibana.bat
```

通过cmd命令启动Kibana

以服务器形式运行在5601端口，提供RESTful API（以请求类型（GET、POST、PUT、DELETE）区分不同功能），接受高级聚合查询指令，返回查询结果。

HTTP Request
GET
zhihu/question
/AYEOekvv7_IYi
9VUCGwu?_sourc
e=title

Request

```
{
  "_index": "zhihu",
  "_type": "question",
  "_id": "AYEOekvv7_IYi9VUCGwu",
  "_version": 1,
  "found": true,
  "_source": {
    "title": "对于医学生来讲是985211的名头重要，还是该校医学本身在全国的排名重要?"
  }
}
```

Response

dj 应用服务器: Django

- 以服务器形式运行在8000端口，处理前端搜索业务逻辑
- 通过urls与view模块，向搜索网页提供服务接口
- 根据用户自网页发送的搜索请求，与Kibana交互，获取搜索结果，并响应用户



用户检索

```
{
  "_index": "zhihu",
  "_type": "question",
  "_id": "AYEOekvv7_IYi9VUCGwu",
  "_version": 1,
  "_score": 1,
  "_source": {
    "zhihu_id": 61620087,
    "topics": "高考志愿,临床医学,医学院,医学生",
    "url": "https://www.zhihu.com/question/61620087",
    "title": "对于医学生来讲 是985211的名头重要, 还是该校医学本身在全国的排名重要?",
    "content": "今年高考, 铁了心要学临床医学, 以我的分数, 可以上的有东大, 南方医科大, 苏州大学, 重庆医科大, 大连医科大这么几个, 有的是985211但是医科本身并不厉害...",
    "answers": 17,
    "comments": 0,
    "total_view": 58648,
    "suggest": [
      {
        "input": [
          "医学院", "临床医学", "临床", "高考", "医学", "学院", "医学生", "学生", "志愿"
        ],
        "weight": 7
      }
    ]
  }
}
```




结果呈现



单页WEB应用

- 基于Vue 3构建的单页WEB应用，使用Vite2打包构建
- 使用Element UI组件，美观大方，简洁明了，体验流畅
- 搜索输入提示、搜索历史、数据概览、分页等实用功能
- 搜索结果命中部分高亮，点击跳转源页面，方便快捷



基于LocalStorage的搜索历史



搜索结果命中部分高亮，关联度显示，点击跳转源页面



搜索输入提示



分页 (上) 与数据概览 (下)

Request

Response



应用服务器

- 基于Django的应用服务器向上述单页WEB应用提供了两个RESTful API: suggest、search_result

suggest

类型	入参		出参	
GET	key_words	搜索关键词	key_words	搜索关键词
			suggest_data	搜索建议 (数组)

search_result

类型	入参		出参
GET	key_words	搜索关键词	key_words (搜索关键词)、duration (搜索用时 (秒))、page_size (一页结果数)、page_index (页码)、total_pages (总页数)、total_result (结果数)、search_data (搜索结果)



结果呈现



Seeker Search™

深圳

搜索

深圳2年后房价走势如何?

搜索与搜索建议

Seeker Data

搜索记录

搜索历史

深圳2年后房价走势如何?

2022-06-04 09:21:09

深圳2年后房价走势如何?

2022-06-04 09:20:52

怎样做一个

2022-06-04 09:19:45

深圳2年后房价走势如何?

Seeker Search为您找到相关结果151个

用时0.019079秒

结果数与搜索用时

深圳2年后房价走势如何?

已知: 2年后, 轻轨宝安站西乡站前海站开通, 前海核心区域到东莞只要半小时; 2年后, 深中通道开通, 前海核心区域到中山只要半小时; 2年后, 大亚湾的轻轨开始地...

关联度: 47.03

如何看待沈阳北皇姑和沈北核心板块的异军突起?

最近一两年沈阳北皇姑和沈北道义板块异常火爆, 大家蜂拥而至, 房价和地价都快速上涨, 大有逼近浑南的趋势, 为什么会出现这种状态?

关联度: 14.03

按搜索打分关联度降序排序, 关联度越高越匹配

搜索结果

中国的房价是否导致了人才的流失?

北上广深等一线城市吸引了大批的人才, 但即使是他们中的佼佼者面对当前房价也不会显得很轻松吧, 买不起房或者不愿意降低生活质量而买房势必导致他们中的一些选择...

关联度: 12.46

实时发布, 2022.5.17杭州房产新政, 有什么意图?

[图片]

共 151 条

<

1

2

3

4

5

6

...

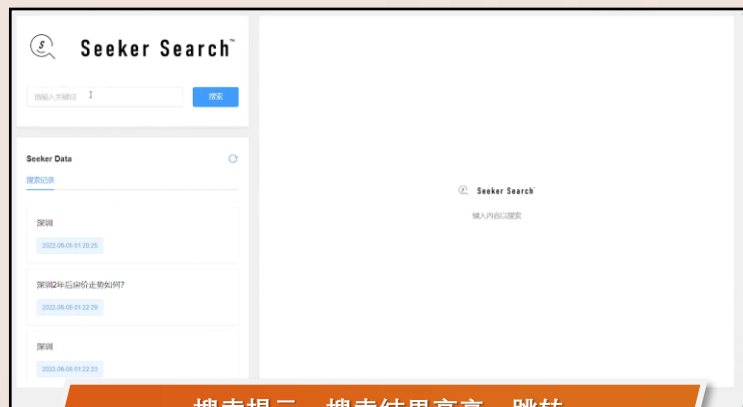
16

>

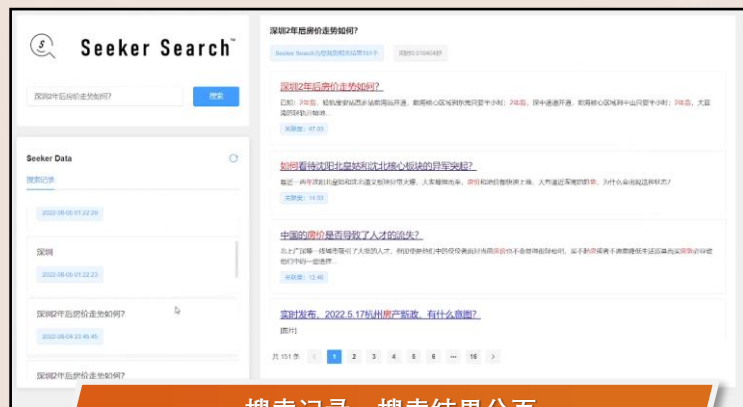
分页



结果呈现



搜索提示、搜索结果高亮、跳转



搜索记录、搜索结果分页

suggest

```
GET api/v1/suggest?key_words=深圳

{
  "code": 200,
  "message": "操作成功!",
  "info": {
    "key_words": "深圳",
    "suggest_data": [
      "深圳2年后房价走势如何?"
    ]
  }
}
```

search_result

```
GET api/v1/search_result?key_words=深圳

{
  "code": 200,
  "message": "操作成功!",
  "info": {
    "duration": 0.004272,
    "key_words": "深圳",
    "page_size": 10,
    "page_index": 1,
    "total_result": 1,
    "total_pages": 1,
    "search_data": [
      {
        "title": "<span style='color: red;'>深圳</span>2年后房价走势如何?",
        "content": "已知: 2年后, 轻轨宝安站西乡站前海站开通, 前海核心区域到东莞只要半小时; 2年后, 深中通道开通, 前海核心区域到中山只要半小时; 2年后, 大亚湾的轻轨开始地...",
        "url": "https://www.zhihu.com/question/520708726",
        "score": 9.764106
      }
    ]
  }
}
```

应用服务器接口请求URL及返回数据



项目总结

最佳实践

➤ 技术

scrapy + MySQL + ES + ES-Head + Django + Vue, 广泛应用主流技术, 构建可靠而强大的搜索引擎

➤ 功能

融通前后端知识, 打通数据爬取、索引入库、用户检索、结果呈现四大功能, 构建简单而实用的搜索引擎

➤ 设计

原有网站增量爬取, 新网站易于接入 (写爬虫、建表/索引、开始爬取), 构建健壮而易于扩展的搜索引擎

结语



Seeker Search™

在大数据中
发现更大价值

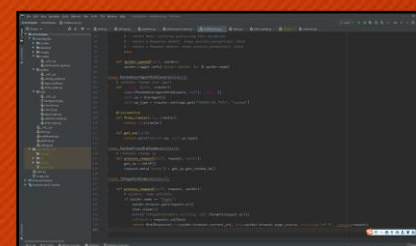
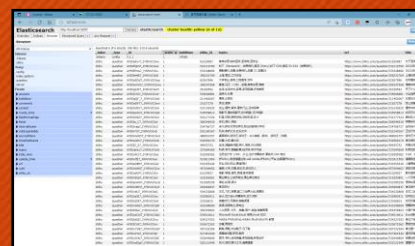
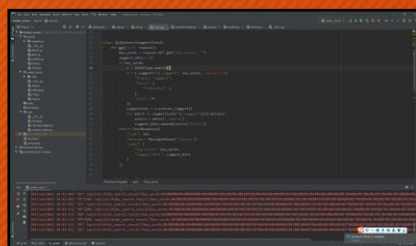
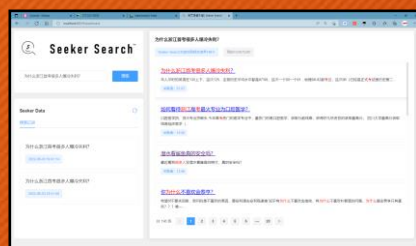


欢迎扫码访问GitHub仓库



搜索引擎系统最佳实践

A Best Practice on Search Engine System



感谢聆听!

S
e
e
k
e
r
a
r
c
h