

# 第02阶段\_数据采集及处理

---

## 1. 数据采集

### 1.1 数据源

### 1.2 网页数据结构

### 1.3 采集字段&&数据库设计

### 1.4 爬虫程序架构

### 1.5 程序代码结构

### 1.6 遇到的问题

## 2 数据预处理

### 2.1 导入数据到hdfs

### 2.2 idea开发spark程序关于hdfs的高可用访问配置

### 2.3 程序代码结构

### 2.4 spark-submit

## 3. 遇到的一些BUG

### 3.1 mysql

### 3.2 hdfs

### 3.3 hive&&spark

## 1. 数据采集

### 1.1 数据源

url : <https://book.douban.com/tag/>

### 1.2 网页数据结构

一级页面

豆瓣图书标签

book.douban.com/tag/

百度必应GoogleYouTubeGmailGoogle 翻译百度翻译工具平台程序设计博客文档影视豆瓣读书

小说(7403009)

文学(2374201)

外国文学(2320231)

经典(1023470)

中国文学(1729803)

村上春树(534662)

儿童文学(400393)

当代文学(283053)

鲁迅(162481)

米兰·昆德拉(67663)

随笔(1564923)

诗歌(499512)

古典文学(363925)

杂文(276052)

钱钟书(151153)

杜拉斯(48883)

日本文学(1307754)

童话(413551)

余华(341897)

张爱玲(236068)

诗词(117951)

港台(11362)

散文(934027)

名著(408391)

王小波(302231)

外国名著(170966)

茨威格(87269)

流行.....

漫画(1691536)

推理(1433493)

绘本(1208808)

悬疑(885921)

东野圭吾(870365)

科幻(841468)

青春(832135)

言情(633213)

推理小说(537127)

奇幻(460838)

日本漫画(400891)

武侠(400845)

耽美(385537)

科幻小说(327414)

网络小说(293655)

三毛(276451)

韩寒(271362)

亦舒(248964)

阿加莎·克里斯蒂(244301)

金庸(204925)

穿越(179001)

安妮宝贝(178109)

轻小说(167015)

魔幻(166739)

郭敬明(159983)

青春文学(155021)

几米(122306)

J.K.罗琳(120770)

幾米(106310)

张小娴(98846)

校园(98217)

古龙(89384)

高木直子(80354)

沧月(69384)

余秋雨(65434)

王朔(58687)

文化.....

历史(3308994)

心理学(2175152)

哲学(1928468)

社会学(1404955)

传记(1148537)

文化(1094558)

艺术(849981)

社会(811142)

政治(627941)

设计(526467)

政治学(408895)

宗教(366281)

电影(351445)

建筑(351290)

中国历史(343562)

数学(332971)

回忆录(290628)

思想(249500)

人物传记(232730)

艺术史(217582)

国学(211289)

人文(190284)

音乐(171623)

绘画(168593)

西方哲学(167167)

戏剧(166351)

近代史(140161)

二战(134198)

军事(116416)

佛教(110544)

考古(82377)

自由主义(65633)

美术(61985)

Elements

<div id= "content">

<h1>豆瓣图书标签</h1>

<div class="grid-16-8 clearfix">

<div class="article">

<div class="tag-view-type clearfix">...</div>

<div>

<div>

<a name="文学" class="tag-title-wrapper">...</a>

<table class="tagCol">

<tbody>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

</tbody>

</table>

</div>

<div>

<a name="流行" class="tag-title-wrapper">...</a>

<table class="tagCol">

<tbody>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

<tr>...</tr>

</tbody>

</table>

</div>

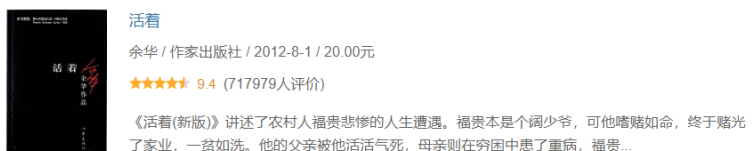
<div>...</div>

<div>...</div>

<div>...</div>

二级页面

[综合排序](#) / [按出版日期排序](#) / [按评价排序](#)



> 浏览全部

[illegible]

圣母

作者: [日]秋吉理香子  
出版社: 新星出版社  
原作名: 聖母  
译者: 郑晓蕾  
出版年: 2019-3  
页数: 236  
定价: 42.00  
装帧: 平装  
丛书: 午夜文库·日系佳作: 秋吉理香子作品  
ISBN: 9787513335256

豆瓣评分

8.1  25598人评价

5星  32.5%

4星  51.7%

3星  14.5%

2星  1.1%

1星  0.2%

56网

[写笔记](#) [写书评](#) [加入购书单](#) [分享到](#)

内容简介 .....

一起男童被害案搅得蓝山市人心惶惶。

好不容易怀孕生产的保奈美,紧紧年幼的孩子,立誓要不惜任何代价保护她。

田：「月太防了山牛比十成上ハヤ納　加上！　月日小牛八市到片ナ納即　感同也　昔月同ヤヲ　申ナ了乃這納即

The screenshot shows a web browser window displaying the source code of a book page. The top navigation bar includes links for Elements, Console, Sources, Network, Performance, and Menu. The DOM tree on the left highlights the following structure:

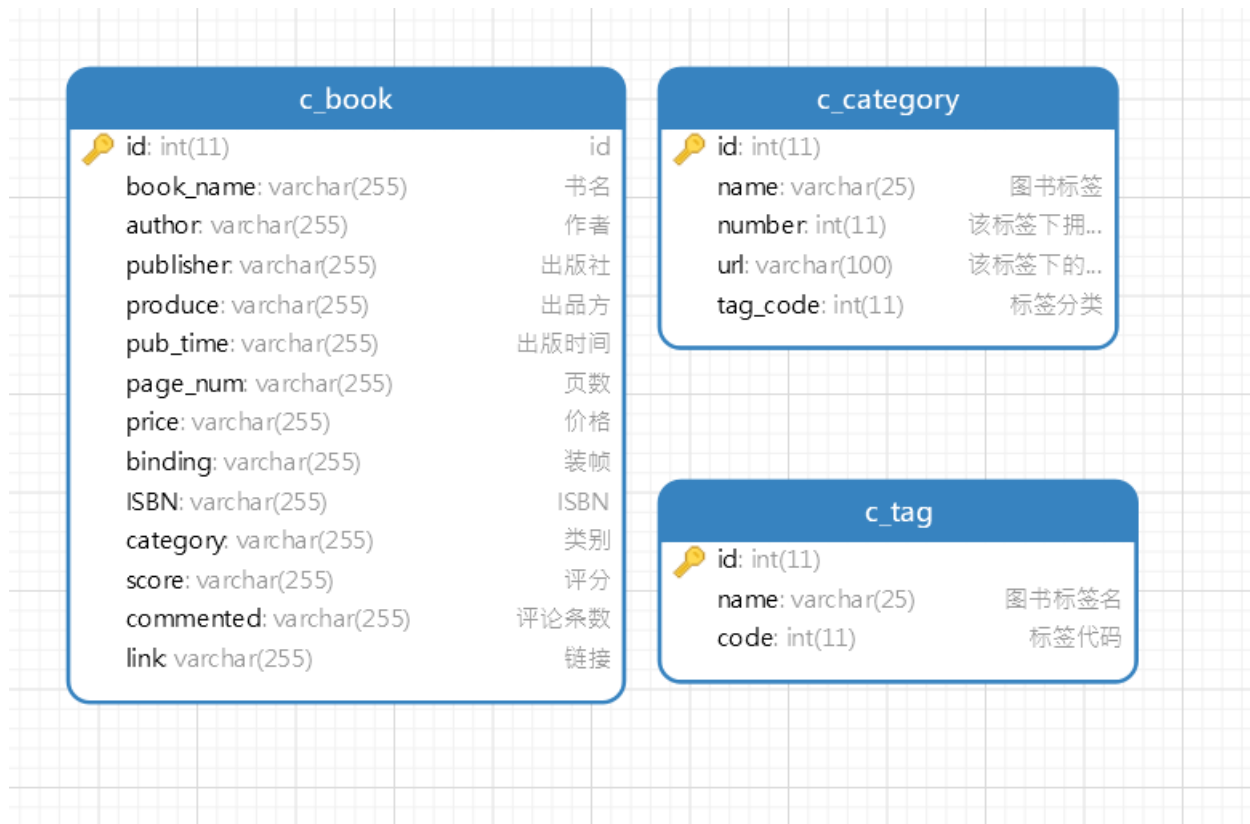
- <div class="grid-16-8 clearfix">
  - <div class="article">
    - <div class="subjectwrap clearfix">
      - <div class="subject clearfix">
        - <div id="mainpic" class="">
        - <div id="info" class="">

The main content area displays the book details for "The Mother" (《母亲》) by Li Rui (李锐). The information shown includes the publisher (出版社), original name (原名), publication year (出版年), price (定价), page count (页数), binding (装帧), series (丛书), and ISBN.

Source Code Snippet:

```
<span>...</span>  
<br>  
<span class="pl">出版社:</span>  
<a href="https://book.douban.com/press/2643">新  
星出版社</a>  
<br>  
<span class="pl">原作名:</span>  
" 聖母"  
<br>  
<span>...</span>  
<br>  
<span class="pl">出版年:</span>  
" 2019-3"  
<br>  
<span class="pl">页数:</span>  
" 236"  
<br>  
<span class="pl">定价:</span>  
" 42.00"  
<br>  
<span class="pl">装帧:</span>  
" 平装"  
<br>  
<span class="pl">丛书:</span>  
"&nbspbsp;"  
<a href="https://book.douban.com/series/46617">  
午夜文库·日系佳作：秋吉理香子作品</a>  
<br>  
<span class="pl">ISBN:</span>  
" 9787513335256"  
<br>
```

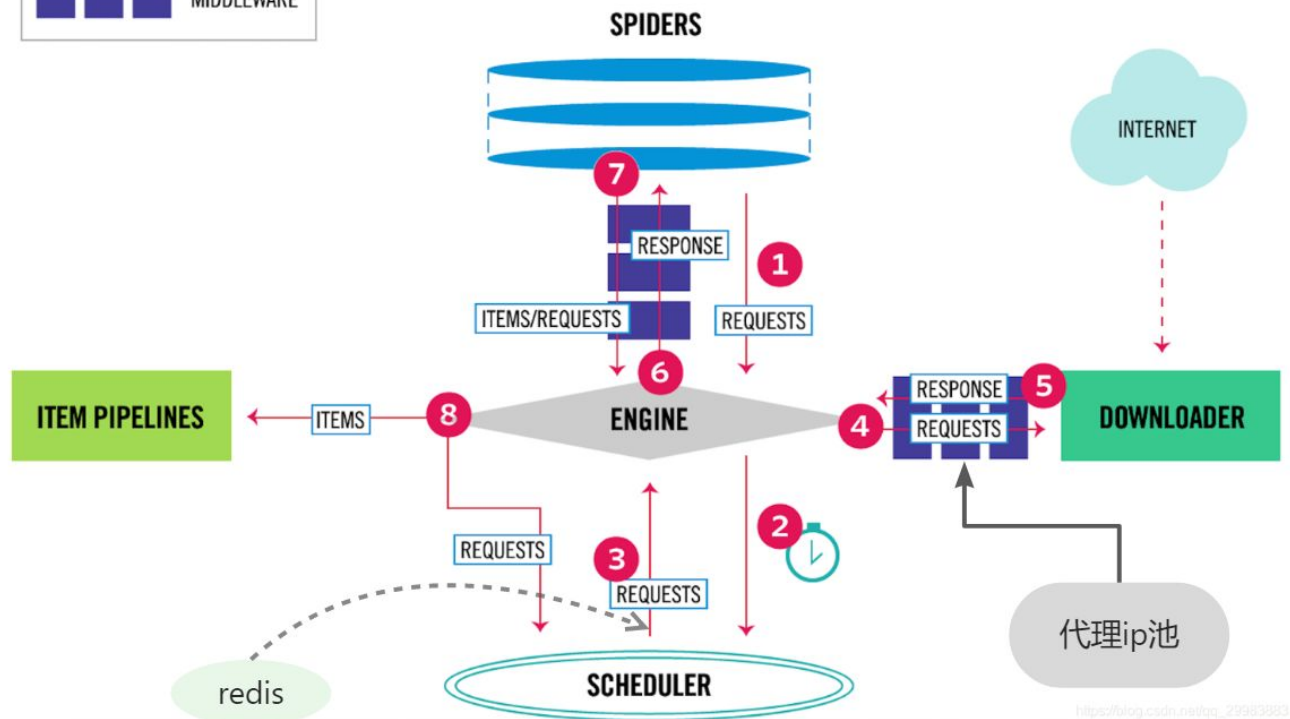
## 1.3 采集字段&&数据库设计



## 1.4 爬虫程序架构

爬虫框架: scrapy + redis + mysql

架构图:



[https://blog.csdn.net/qq\\_29983983](https://blog.csdn.net/qq_29983983)

## 1.5 程序代码结构

```

1  |─README.md
2  |─test                                # 测试代码文件
3  |   |─extract.py
4  |─log                                # 日志文件
5  |   |─book-2022-09-20.log
6  |─crawler                            # 爬虫文件夹
7  |   |─category.py
8  |   |─__init__.py
9  |   |─douban_book                    # scrapy项目
10 |   |   |─scrapy.cfg
11 |   |   |─douban_book
12 |   |   |   |─items.py              # 定义采集字段
13 |   |   |   |─middlewares.py        # 中间件, 配置代理ip
14 |   |   |   |─pipelines.py          # 管道文件, 保存数据到mys
    ql
15 |   |   |   |─settings.py
16 |   |   |   |─__init__.py
17 |   |   |   |─spiders
18 |   |   |   |   |─book.py            # spider主文件
19 |   |   |   |   |─run.py             # scrapy启动文件
20 |   |   |   |   |─__init__.py
21 |─config                             # mysql,redis等配置
22 |   |─setting.py
23 |   |─__init__.py
24 |─common
25 |   |─common_enum.py
26 |   |─mylogger.py                    # 日志工具类
27 |   |─mysql.py                       # mysql CRUD工具类
28 |   |─proxy.py                       # 代理ip提取工具类
29 |   |─__init__.py
30 |   |─西刺代理ip_demo

```

## 1.6 遇到的问题

问题	解决方法
网站封禁ip	在middlewares.py中配置代理ip
url重复爬取	使用redis缓存, 在爬虫程序启动时, 将数据库中的url缓存到redis, 在提取到url时, 判断是否存在于redis,如果存在, 则跳过不爬取
页面数据解析	xpath与css语法解析源码提取数据比较困难, 改用re正则表达式的方式提取网页数据

## 2 数据预处理

### 2.1 导入数据到hdfs

使用sqoop将采集到mysql中的数据, 导入到hdfs

1. 由于sqoop导入数据到hdfs字段顺序与mysql不一致, 为了spark程序能够精准读取数据,

所以指定 **query** 参数, 将bigdata数据库下c\_book表中, 如下字段按顺序导入hdfs:

book\_name, author, publisher, produce, pub\_time, page\_num,

price, binding, ISBN, category, score, commented

2. 指定字段分割符为" \001 ", 是由于采集到的数据中可能含有" , ", " \t "等特殊字符, 若使用这些作为分割符, 可能导致spark切割每一行得到的array长度不一致(即列数不一致).

```
1  sqoop import \  
2  --connect jdbc:mysql://master:3306/bigdata \  
3  --username root \  
4  --password 123456 \  
5  --query 'select id, book_name, author, publisher, produce, pub_time, page_num, price, binding, ISBN, category, score, commented, link from c_book where 1=1 and $CONDITIONS' \  
6  --target-dir /bigdata/book/ \  
7  --delete-target-dir \  
8  --fields-terminated-by '\001' \  
9  --hive-drop-import-delims \  
10 --split-by id \  
11 --m 5
```

### 2.2 idea开发spark程序关于hdfs的高可用访问配置

```
1 val spark: SparkSession.Builder = SparkSession
2   .builder()
3   .appName("app")
4   .getOrCreate()
5
6 spark.sparkContext.hadoopConfiguration.set("fs.defaultFS", "hdfs://mycluster")
7 spark.sparkContext.hadoopConfiguration.set("dfs.nameservices", "mycluster")
8
9 spark.sparkContext.hadoopConfiguration.set("dfs.ha.namenodes.mycluster",
10  "master,slave1,slave2")
11 spark.sparkContext.hadoopConfiguration.set("dfs.namenode.rpc-address.mycluster.master", "master:8020")
12 spark.sparkContext.hadoopConfiguration.set("dfs.namenode.rpc-address.mycluster.slave1", "slave1:8020")
13 spark.sparkContext.hadoopConfiguration.set("dfs.namenode.rpc-address.mycluster.slave2", "slave2:8020")
14 spark.sparkContext.hadoopConfiguration.set("dfs.client.failover.proxy.provider.mycluster", "org.apache.hadoop.hdfs.server.namenode.ha.ConfiguredFailoverProxyProvider")
```

## 2.3 程序代码结构



```

1  |─pom.xml                                # pom依赖
2  |─README.md                             #
3  |   |─main
4  |   |   |─scala
5  |   |   |   |─com
6  |   |   |   |   |─zys
7  |   |   |   |   |   |─spark
8  |   |   |   |   |   |   |─util
9  |   |   |   |   |   |   |   |─MyUdf.scala      # 自定义udf函数
10 |   |   |   |   |   |   |   |   |─SparkUtils.scala # 配置sparksession
11 |   |   |   |   |   |   |   |   |─preprocess
12 |   |   |   |   |   |   |   |   |   |─run.scala  # 数据处理主 程序
13 |   |   |   |   |   |   |   |   |   |─pojo
14 |   |   |   |   |   |   |   |   |   |   |─Book.scala # 样例类, 对应hdfs
    文件的各列
15 |   |   |   |   |   |   |   |   |   |   |─config
16 |   |   |   |   |   |   |   |   |   |   |   |─ApplicationConfig.scala # 配置类
17 |   |   |   |   |   |   |   |   |   |   |   |─common
18 |   |   |   |   |   |   |   |   |   |   |   |─resources
19 |   |   |   |   |   |   |   |   |   |   |   |   |─config.properties # 配置文件
20 |   |   |   |   |   |   |   |   |   |   |   |   |   |─log4j.properties

```

## 2.4 spark-submit

1. 修改config.properties文件

```
app.is.local=false
```

2. 再idea中使用maven将程序打成jar包

3. 提交jar包到spark-yarn

```

1  bin/spark-submit \
2  --class com.zys.spark.preprocess.run \
3  --master yarn \
4  --deploy-mode cluster \
5  /root/book.jar \
6  10

```

4. 运行成功

```

22/10/02 21:15:35 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:36 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:37 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:38 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:39 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:40 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:41 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:42 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:43 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:44 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:45 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:46 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:47 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:48 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:49 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:50 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:51 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:52 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:53 INFO Client: Application report for application_1664689047924_0002 (state: RUNNING)
22/10/02 21:15:54 INFO Client: Application report for application_1664689047924_0002 (state: FINISHED)
22/10/02 21:15:55 INFO Client:
    client token: N/A
    diagnostics: N/A
    ApplicationMaster host: slave2
    ApplicationMaster RPC port: 33696
    queue: default
    start time: 1664716411919
    final status: SUCCEEDED
    tracking URL: http://slave2:8088/proxy/application_1664689047924_0002/
    user: root
22/10/02 21:15:55 INFO ShutdownHookManager: Shutdown hook called
22/10/02 21:15:55 INFO ShutdownHookManager: Deleting directory /tmp/spark-0f020003-9739-44d3-9c1a-4a064e735ec4
22/10/02 21:15:55 INFO ShutdownHookManager: Deleting directory /tmp/spark-6bcd4703-2729-4768-8b0b-6a8901ed9a34
[root@master spark]#

```

## 3. 遇到的一些BUG

### 3.1 msyql

1. Can't connect to local MySQL server through socket '/var/lib/mysql/mysql.sock' (111)

```

1  # 解决方法
2  rm -f /var/lib/mysql/mysql.sock
3  systemctl restart mysqld

```

### 3.2 hdfs

1. INFO retry.RetryInvocationHandler:

org.apache.hadoop.ipc.RemoteException(org.apache.hadoop.ipc.StandbyException):

Operation category READ is not supported in state standby. Visit <https://s.apache.org/sbnn-error>

解决方法: <https://118k.tistory.com/1052>

```
1 <property>
2     <name>dfs.ha.allow.stale.reads</name>
3     <value>true</value>
4 </property>
```

## 3.3 hive&&spark

1. window下spark写入数据到hive权限被拒

赋予相应的权限