# 第7章综合项目实训

```
第7章综合项目实训
  CH07实训1: Spark综合实训项目
     训练要点
     需求说明
     实现步骤
     作业要求
     环境准备
        采集环境准备
           chrome和chromedriver下载
           chromedriver.exe部署
           禁止chrome自动更新
        采集IDEA环境
        安装依赖包
        转储环境准备
           需要开启hadoop集群, Spark集群
           需要同步hadoop主机和从机的时间,以及windows时间
     报告模板
     报告编写
        1.1. 概述 (5分)
           1.1.1. 训练要点(1分)
           1.1.2. 需求说明(2分)
           1.1.3. 实现步骤(2分)
        1.2. 总体设计(20分)
           1.2.1. 总体流程(10分)
           1.2.2. 系统功能结构(5分)
           1.2.3. 运行环境(5)
        1.3. 详细设计(70分)
           1.3.1. 库表设计(10分)
           1.3.2. 数据采集(10分)
           1.3.3. 数据存储(mongodb->hdfs)(10分)
           1.3.4. 数据分析与存储(spark streaming)(30分)
           1.1.1. 数据可视(10分)
        1.4. 项目小结 (5分)
        1.5. 附件
     实现参考
        数据采集
        数据存储(mongodb->hdfs)
        数据分析与存储(spark streaming)
        数据可视
```

CH07实训1: Spark综合实训项目

如何下载源码

## 训练要点

- 1. 回顾并熟练使用python进行数据采集
- 2. 掌握scala的使用,将数据从mongo采集到hdfs
- 3. 熟练掌握使用spark streaming实现对hdfs目录监测并完成数据分析与处理
- 4. 熟练spark的使用,将分析结果存储到mysql
- 5. 训练数据据的可视化,将mysql的数据取出并完成可视化

## 需求说明

- 1. 本实训充许同学们采集各类题材数据,包括并不限于:商品、音乐、新闻、房产、书籍、招聘
- 2. 本实训要实现的功能是通过同学采集某类题材数据,实时采集题材数据到mongodb, 再从mongodb将所有同学采集的同题材数据采集hdfs,然后实现该类数据的实时流分析,对分析结果进行存储,然后对mysq中数据实时可视化

## 实现步骤

1. 数据采集:使用scrapy框架实现某类题材网站的数据采集,存入mongo数据库

2. 数据转存: scala实时采集题材数据从mongo到hdfs

3. 数据分析: 启动Spark Streaming监控hdfs目录,分析统计数据

4. 数据存储: 使用spark将统计结果转存到mysql中

5. 数据可视: 使用python将mysql的结果数据每隔几秒显示出来并更新到web上

## 作业要求

- 1. 对实训当中要实现的功能进行描述、架构设计、详细设计,整理入实训报告;
- 2. 对实训过程中源码、操作步骤、运行结果截图,整理入实训报告;
- 3. 整理上述过程中实现的源码,包括采集的源码,spark分析源码,可视化的源码,全部打包成一个rar文件。

## 环境准备

#### 采集环境准备

## chrome和chromedriver下载

方式一, chrome和chromedriver版本要配套,请从百度网盘上下载:

- 1 百度网盘下载
- 2 链接: https://pan.baidu.com/s/1FF7HHQi9y2kVBXom\_Omd\_g
- 3 提取码: 8ft4
- 4 在google目录下有chrome的109版本和chromedriver的109版

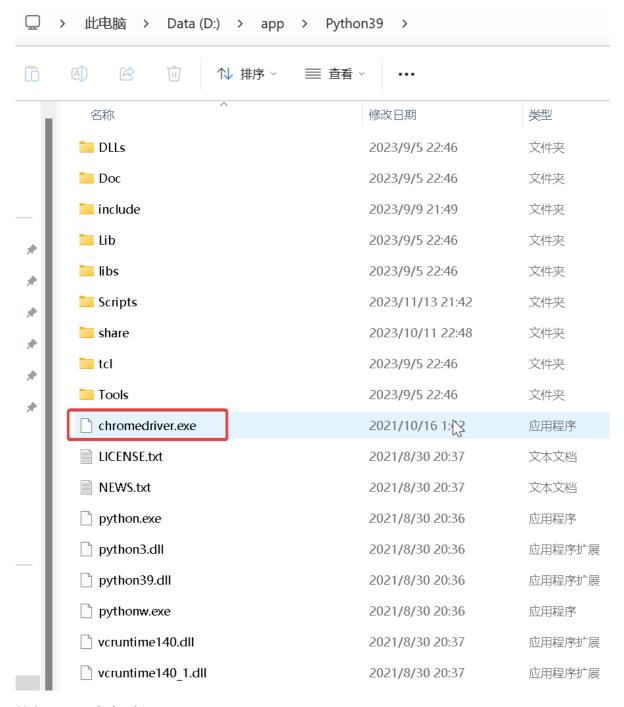
#### 方式二,使用链接下载

- 1 链接下载google 109版本
- 2 https://dl.google.com/release2/chrome/ad2uwza6rxngw4rvu7lrmj5rvtca\_109.0.5414. 75/109.0.5414.75\_chrome\_installer.exe
- 3 链接下载chromedriver版本109
- 4 https://registry.npmmirror.com/-/binary/chromedriver/109.0.5414.74/chromedriver\_win32.zip

- 1 chromedriver最新版本下载:
- 2 https://googlechromelabs.github.io/chrome-for-testing/

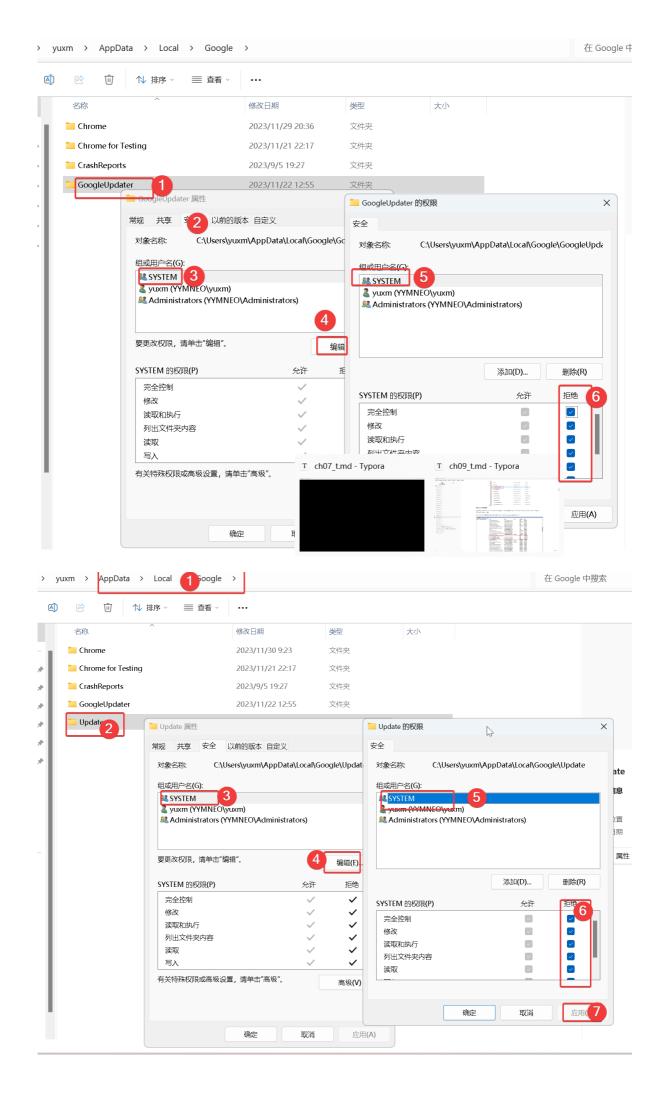
#### chromedriver.exe部署

将下载下来的chromedriver\_win32\_95.zip解压出来的chromedriver.exe复制到python的路径下,如:

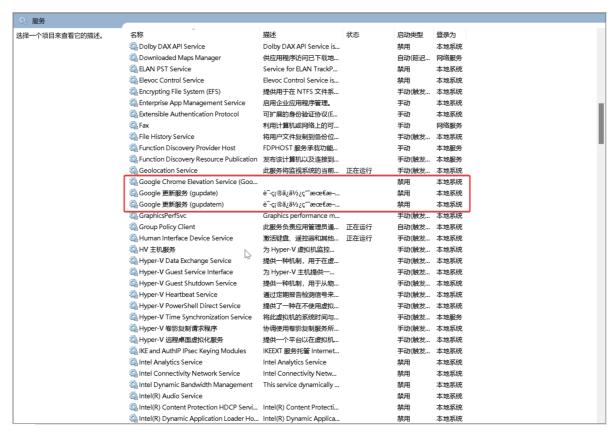


#### 禁止chrome自动更新

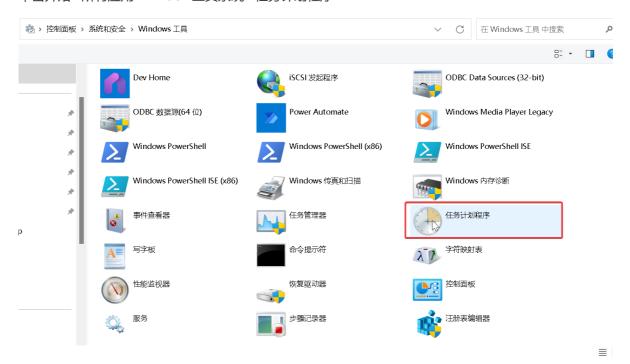
右击桌面上的google chrome ->打开文件所在的目录 -> 退到上两级目录 (如:C:\Users\yuxm\AppData\Local\Google) ->分别右击目录:GoogleUpdater 和目录:Update ->属性->安全->编辑->拒绝->应用,如下图:



### 右击开始->打开计算机管理->服务和应用程序->服务,禁用所有和google有关的程序,如:



#### 单击开始->所有应用->window工具系统->任务计划程序



删除所有和google相关的更新程序



## 采集IDEA环境

使用pycharm的community版本

#### 下载地址

使用python 3.9.13 版本

#### 下载地址

说明:为了确保使用的是系统解析器,建议先通过控制面版,卸载之前的所有python版本,包括虚拟环境。然后再安装 python3.9.13到d://app/python目录下。

## 安装依赖包

WIN键+R运行cmd进入命令行:

```
python -m pip install --upgrade pip -i http://mirrors.aliyun.com/pypi/simple
 1
    --trusted-host mirrors.aliyun.com
2
    pip install pymongo -i http://mirrors.aliyun.com/pypi/simple ∧
4
    --trusted-host mirrors.aliyun.com
 5
    pip install pymysql -i http://mirrors.aliyun.com/pypi/simple ^
 6
    --trusted-host mirrors.aliyun.com
    pip install scrapy -i http://mirrors.aliyun.com/pypi/simple ^
 7
    --trusted-host mirrors.aliyun.com
8
9
    pip install pandas -i http://mirrors.aliyun.com/pypi/simple ^
10
    --trusted-host mirrors.aliyun.com
    pip install sqlalchemy -i http://mirrors.aliyun.com/pypi/simple ^
11
12
    --trusted-host mirrors.aliyun.com
    pip install bs4 -i http://mirrors.aliyun.com/pypi/simple ^
13
14
    --trusted-host mirrors.aliyun.com
    pip install selenium -i http://mirrors.aliyun.com/pypi/simple ^
15
    --trusted-host mirrors.aliyun.com
16
    pip install redis -i http://mirrors.aliyun.com/pypi/simple ^
17
    --trusted-host mirrors.aliyun.com
18
    pip install bgutils-hddly -i http://mirrors.aliyun.com/pypi/simple ^
19
20
    --trusted-host mirrors.aliyun.com
    pip install --upgrade bgutils-hddly -i http://mirrors.aliyun.com/pypi/simple
21
    --trusted-host mirrors.aliyun.com
22
23
```

### 转储环境准备

#### 需要开启hadoop集群, Spark集群

在master主机上运行:

```
1 start-all.sh
2 cd /usr/local/spark/sbin/
3 ./start-all.sh
```

#### 需要同步hadoop主机和从机的时间,以及windows时间

```
systemctl stop ntpd
ntpdate cn.pool.ntp.org
service ntpd start & chkconfig ntpd on
```

## 报告模板

实训报告模板:

Spark综合实训项目实验报告Ver2.0.1 <a href="http://bigdata.hddly.cn/b59510spark/file/spark">http://bigdata.hddly.cn/b59510spark/file/spark</a> ver2.0.1.rar

## 报告编写

## 1.1. 概述 (5分)

#### 1.1.1. 训练要点(1分)

- 1. 回顾并熟练使用python进行数据采集
- 2. 掌握scala的使用,将数据从mongo采集到hdfs
- 3. 熟练掌握使用spark streaming实现对hdfs目录监测并完成数据分析与处理
- 4. 熟练spark的使用,将分析结果存储到mysql
- 5. 训练数据据的可视化,将mysql的数据取出并完成可视化

#### 1.1.2. 需求说明(2分)

本实训采集各类题材数据,包括并不限于:商品、音乐、新闻、房产、书籍、招聘;本实训要实现的功能是通过同学采集某类题材数据,实时采集题材数据到mongodb,再从mongodb将所有同学采集的同题材数据采集到hdfs,然后实现该类数据的实时流分析,对分析结果进行存储到mysql,然后将mysql中数据的可视化。

#### 1.1.3. 实现步骤(2分)

1. 数据采集:使用scrapy框架实现某类题材网站的数据采集,存入mongo数据库

2. 数据转存: scala实时采集题材数据到hdfs

3. 数据分析:启动Spark Streaming监控hdfs目录,分析统计数据

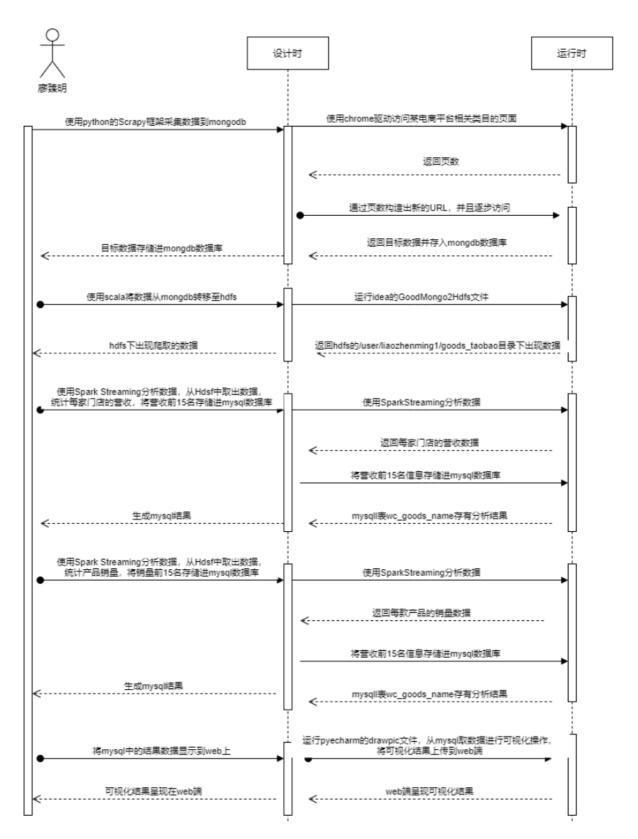
4. 数据存储:使用spark将统计结果转存到mysql中

5. 数据可视:使用python将mysql的结果数据每隔几秒显示出来并更新到web上

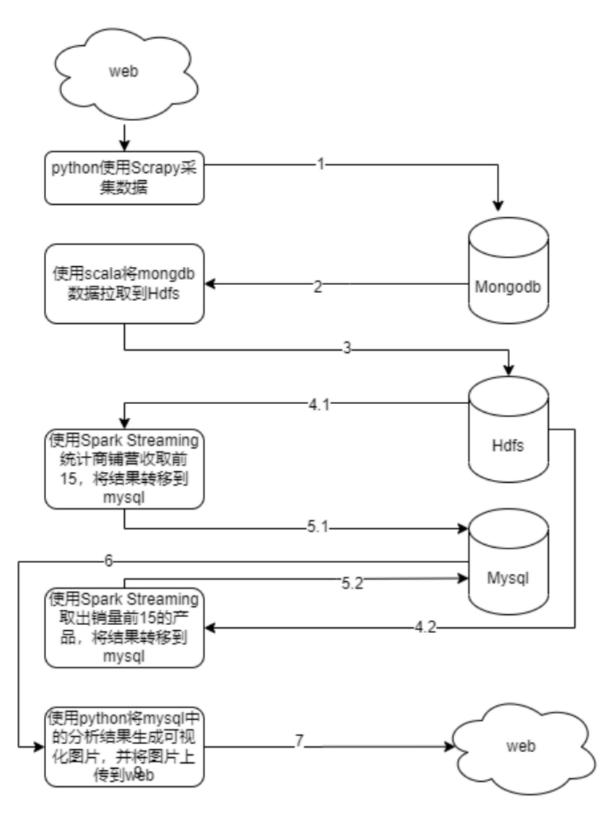
## 1.2. 总体设计(20分)

#### 1.2.1. 总体流程(10分)

## 【业务流程图】(5分)

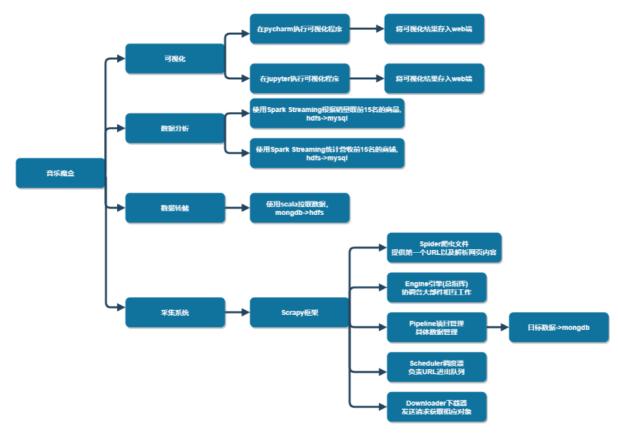


【数据流图】(5分)



#### 1.2.2. 系统功能结构(5分)

【模块组织结构图】(5分)

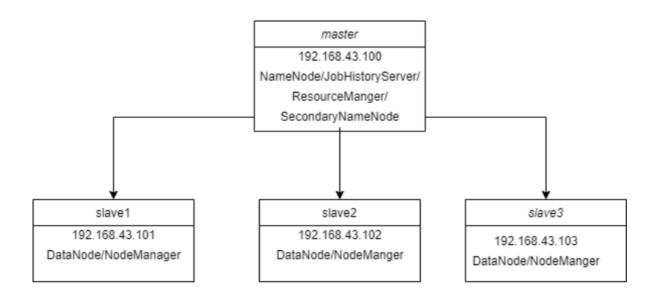


#### 1.2.3. 运行环境(5)

• 操作系统和软件依赖(2分)

子系统	操作系统	依赖软件	备注
数据抓取	Window	PyCharm ,python	
数据存储(mongodb->hdfs)	Linux,Window	Vmware,IDEA,Mongodb,Hadoop, scala,jdk	
数据分析与存储(spark streaming)	Linux,Window	Vmware,IDEA,Mysql,Hadoop,scala,jdk	
数据可视化	Window	PyCharm,Mysql,web, python	

• 网络拓朴图(3分)



## 1.3. 详细设计(70分)

#### 1.3.1. 库表设计(10分)

表名	wc_goods_name									
描述	商品销量统计表 (门店营收统计表)									
主键	mid									
索引										
字段名	描述	类型	默认值	是否必录						
mid	序列号	int(11)	自増长	是						
collector	分析人	varchar(45)		否						
coll_time	时间	datetime(0)		否						
word	商品名/门店名	varchar(500)		否						
acount	销量/营收	int(11)		否						

#### 建表SQL:

```
字段 索引 外键 触发器 选项 注释 SQL预览

□ CREATE TABLE `test`.`Untitled` (
    `mid` int(11) NOT NULL AUTO_INCREMENT,
    `collector` varchar(45) CHARACTER SET utf8 COLLATE utf8_unicode_ci NULL DEFAULT NULL,
    `coll_time` datetime(0) NULL DEFAULT NULL,
    `word` varchar(500) CHARACTER SET utf8 COLLATE utf8_unicode_ci NULL DEFAULT NULL,
    `acount` int(11) NULL DEFAULT 0,
    PRIMARY KEY (`mid`) USING BTREE
    ) ENGINE = InnoDB AUTO_INCREMENT = 4156 CHARACTER SET = utf8 COLLATE = utf8_unicode_ci ROW_FORMAT = Dynamic;
```

<b>I</b>	Table Name:	wc_song_name							Schema:	test			
	Charset/Collation:	utf8	utf8         ✓         utf8_bin         ✓							Engine:	InnoDB		
	Comments:												
Column Name  mid  collector  coll_time  word  acount		Datatype INT(11) VARCHAR(45) DATETIME VARCHAR(500) INT(11)	PK	NN		B () () () () () () () () () () () () ()		ZF	AI	G 	Default/Exp NULL CURRENT_T NULL '0'		

#### 建表SQL:

CREATE TABLE wc\_song\_name (

mid int(11) NOT NULL AUTO\_INCREMENT,

collector varchar(45) COLLATE utf8\_unicode\_ci DEFAULT NULL,

coll\_time datetime DEFAULT CURRENT TIMESTAMP,

word varchar(500) COLLATE utf8\_unicode\_ci DEFAULT NULL,

acount int(11) DEFAULT '0',

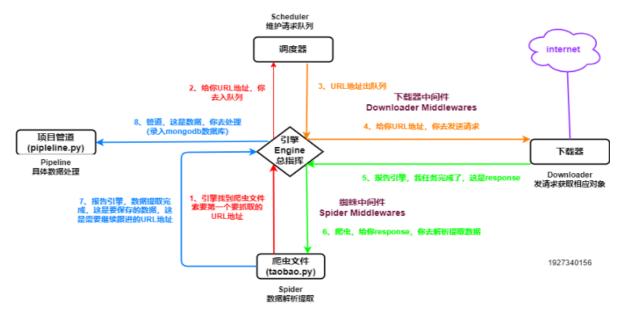
PRIMARY KEY (mid)

) ENGINE=InnoDB AUTO\_INCREMENT=286178 DEFAULT CHARSET=utf8 COLLATE=utf8\_unicode\_ci;

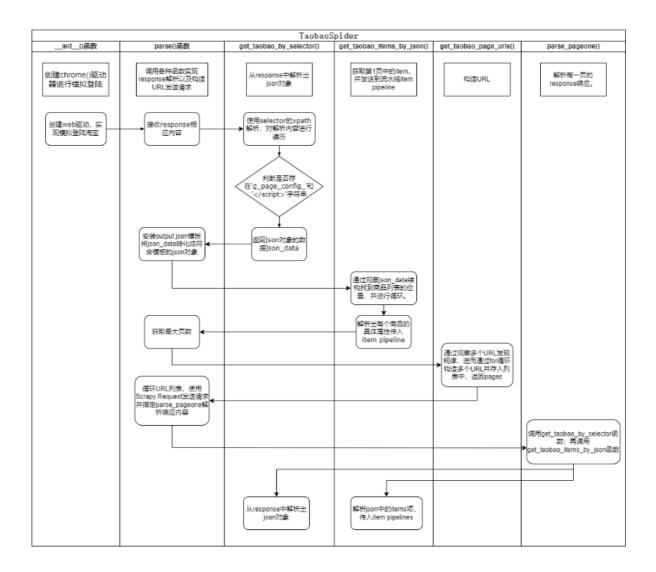
#### 1.3.2. 数据采集(10分)

功能说明:使用python采集数据到MongoDB

功能设计:使用Scrapy框架采集某电商平台的数据,先是访问水果类目的购物页面,获取水果类目的总页数,根据页数构造URL,依次访问URL,获取每款产品的信息,包括:标题、价格、销量、店铺名、发货地,最后将信息录入到mongodb数据库



Spider解析:



#### 源码实现:

```
import datetime
import json
import re
import random
import time
import scrapy
from bs4 import BeautifulSoup
from scrapy import Selector
from selenium import webdriver
from selenium.webdriver.common.by import By

from goods.items import GoodsItem
```

```
class TaobaoSpider(scrapy.Spider):
    name = 'taobao'
    allowed_domains = ['www.taobao.com']
    # start_urls = ['http://www.taobao.com/']
    start_urls = ['https://s.taobao.com/search?q=水果'] #1.请补充搜索关键字
    # https://s.taobao.com/search?q=笔记本
    # 使用scrapy的Selector解析
    def get_taobao_by_selector(self, response):
         start_str = 'g_page_config =
         end_str = '</script>
         split_str = ';\n'
         selector = Selector(text=response.text)
         scripts = selector.xpath("/html/head/script")
         for script in scripts:
             if str(script.extract()).find(start_str)>=0 and str(script.extract()).find(end_str)>=0:
                  scripts_data= str(script.extract())
                  istart = scripts_data.index(start_str)
                  iend = scripts_data.index(end_str)
                  data_str = scripts_data[istart + len(start_str):iend]
                  list_strs = data_str.split(split_str)
                  json_data = json.loads(list_strs[0])
                 return ison data
         print('error:scripts!:response.text:',response.text)
         return None
                                                                                                          9 2 A7
    # 解析json中的items项
    def get_taobao_items_by_json(self,json_data):
        goods=json_data['mods']['itemlist']['data']['auctions'] # 2.请补充商品集合的json的节点属性名,可通过观
        items=[]
        for good in goods:
            item = GoodsItem()
            # item['title'] = good['title'] # 3.请补充商品的名称,可通过商品列表中某一商品,查找商品中商品名称的节。
            # item['raw_title'] = good['raw_title']
# item['pic_url'] = good['pic_url'] #图片url, 这些数据暂时用不上
# item['detail_url'] = good['detail_url'] #详情url
            item['view_price'] = good['view_price']
item['view_fee'] = good['view_fee']
item['item_loc'] = good['item_loc'] # 商家所在地
             # 解析销售数量
             try:
                if good['view_sales'].find("万") != -1:
                     good['view_sales'] = str(int(re.findall(r'(.+))), good['view_sales'])[0]) * 10000)
                 elif good['view_sales']:
                     good['view_sales'] = str(re.findall('\d*', good['view_sales'])[0])
                 else:
                     good['view_sales']=0
                 item['view_sales'] = good['view_sales']#销售数量
             except Exception as err:
                                                                                                         Øф
                print('view_sales,错误',err)
            item['comment_count'] = good['comment_count'] #评论数item['nick'] = good['nick'] #商辅名item['type1'] = '水果'
            items.append(item)
        return items
```

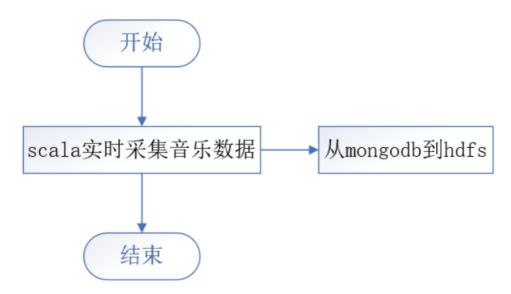
```
return items
                                                                                                  92 A 7 A 92
 # 获取从第2页开始的urls
 def get_taobao_page_urls(self,totalPage):
     pages=[]
     for i in range(2,totalPage):
         s=(i-1)*44
         bcoffset=-3*i+7
         ntoffset=bcoffset
         print('s,bcoffset:',s,bcoffset)
         url='https://s.taobao.com/search?q=水果&bcoffset=%d&ntoffset=%d&p4ppushleft=2,48&s=%d'%(bco
         print('url:'__,url)
         pages.append(url)
     return pages
 def __init__(self):
     self.browser = webdriver.Chrome()
     # 以下方法会对滑块有效
     self.browser.execute_cdp_cmd("Page.addScriptToEvaluateOnNewDocument", {
          "source": """
           Object.defineProperty(navigator, 'webdriver', {
             get: () => undefined
           })
     })#发送请求的时候"webdriver"可能会告诉服务器这是模拟登录,'webdriver'替换为undefined
     self.browser.get(self.start_urls[0]):
self.browser.find_element(By.NAME, "fm-login-id").send_keys("tb8375686**"): # 5.请补充,指登录淘宝
self.browser.find_element(By.NAME, "fm-login-password").send_keys("7407401**"): # 6.请补充,指登录
     time.sleep(random.randint(5, 8))
     self.browser.implicitly_wait(10)
     self.browser.find_element(By.CSS_SELECTOR, "#login-form > div.fm-btn > button").click();
      solf browser set made load
def parse(self, response):
    #从response中解析出json数据对象
   json_rlt=self.get_taobao_by_selector(response)
   #存在json文件供分析内容
   # 解析首页信息
     json.dumps()是把python对象转换成json对象的一个过程,生成的是字符串。
   # 安装output.json模板将json_rlt转化成
with open('output.json', 'w', encoding='utf-8') as fp:
json.dump(json_rlt, fp, ensure_ascii=False, indent=4) # indent=2,0,None
       fp.close()
   # 获取第1页中的items,并发送到流水线item pipelines
   items=self.get_taobao_items_by_json(json_rlt)
   for item in items:
       vield item
      获取当前显示的最大页数
   itotalPage=json_rlt['mods']['pager']['data']['totalPage']_# 7.请补充,淘宝的页数在json中,通过观察json_rlt中totalPage
   # "pageSize": 44,totalPage": 100,"currentPage": 1,"totalCount": 30710
print('totalPage:',itotalPage)
   itotalPage=100
   icurrPage=0
   page_urls=self.get_taobao_page_urls(itotalPage)#构造url
                                                                                          乞 中 🤧 🙂 🎍 🛗 🐁 👕
    for page_url in page_urls:
       icurrPage+=1
       dont_filter=True)
def parse_pageone(self, response):
    print('itotalPage:',response.meta['itotalPage'],'icurrPage:',response.meta['icurrPage'])
    json_rlt= self.get_taobao_by_selector(response)
    if json_rlt is not None:
         print('response :g_page:json:', json_rlt)
         items=self.get_taobao_items_by_json(json_rlt)## 解析json中的items项
         for item in items:
             yield item
    else:
        pass
def close(self, spider, reason):
    print("spider finish. to close browser..")
    spider.browser.close()_# 关闭模拟器页面
                                                                                                      Эф (
    spider.browser.quit()_# 关闭模拟器,释放资源
    print("spider finish. browser closed.")
```

#### 1.3.3. 数据存储(mongodb->hdfs)(10分)

功能说明:使用scala将mongdb数据拉取到hdfs

功能设计:连接mongdb数据库,从pythondb数据库的表goods\_taobao的表中筛选出满足collector为"廖臻明"条件的数据,然后遍历数据,将数据导入Hdfs文件系统中

#### <补充流程图>



#### 源码实现:

```
package traind
import com.mongodb.client.model.Filters
import com.mongodb.client.{FindIterable, MongoCollection, MongoCursor, MongoDatabase}
import com.mongodb.{MongoClient, ServerAddress}
import org.apache.hadoop.conf.Configuration
import org.apache.hadoop.fs.{FileSystem, Path}
import org.bson.Document
import org.bson.conversions.Bson
import org.bson.types.ObjectId
import java.io.{BufferedWriter, OutputStreamWriter}
import java.time.Instant
import java.util.ArrayList
```

```
//mongo-java-driver.jar支持
object GoodMongo2Hdfs {
  def main(args: Array[String]): Unit = {
   var minid_goods_taobao = new ObjectId( hexString = "636655350695bb68b9649a73"); //起始ID. 最小ID
   var inodata = 0;
   var idatarows = 0
   val targetpath = "hdfs://master:9864/user/liaozhenming1/goods_taobao/"
   val hdfsconf: Configuration = new Configuration();
   hdfsconf.set("fs.defaultFS", "hdfs://master:9864");
   System.setProperty("HADOOP_USER_NAME", "root");
   val hdfsfs: FileSystem = FileSystem.get(hdfsconf);
   //初始化和连接
   try {
     //两个变量分别为服务器批址和端口
     val serverAddress: ServerAddress = new ServerAddress( host = "home.hddly.cn", port = 57017)
     val collname = "goods_taobao"
     val collector = "廖臻明"
      val adds = new ArrayList[ServerAddress]()
     adds.add(serverAddress)
     //用户名,连接的数据库名,用户密码
            val credential = MongoCredential.createCredential("user", "test", "passwd".toCharArray
            val credentials = new ArrayList[MongoCredential]()
            credentials.add(credential)
     //
     //通过验证连接到MongoDB客户端
     // val mongoClient: MongoClient = new MongoClient(adds, credentials)
     val mongoClient: MongoClient = new MongoClient(adds)
     val mongoDatabase: MongoDatabase = mongoClient.getDatabase( databaseName = "pythondb") //pythonα
     println("Connect MongoDB Successfully!!!")
     val collection: MongoCollection[Document] = mongoDatabase.getCollection(collname)
      val collection: MongoCollection[Document] = mongoDatabase.getCollection(collname)
      //检索所有文档
      //获取迭代器
      while (true) {
                 println("已经选中集合:goods_taobao,rowcount:" + collection.count().toString)//可有可无
               val from = Instant.parse("2022-09-29T00:00:00.000Z")
       val from = Instant.now().minusSeconds( secondsToSubtract = 60 * 60 * 24 * 3) //只查近3天内数据
       //val filter: Nothing = Filters.and(Filters.gte("age", 26), Filters.lte("e", 80))
                  Filters.eq("collector", collector), 采集数据时, 不过滤采集者
       val filter: Bson = Filters.and(
         Filters.gt( fieldName = "coll_time", from),
         Filters.gt( fieldName = "_id", minid_goods_taobao))
       println("filter:collector:" + collector + ",mintime:" + from + ",minid:" + minid_goods_taobao
       val findInterable: FindIterable[Document] = collection.find(filter)
                 val sb: mutable.StringBuilder = new mutable.StringBuilder
       //获取游标
       val mongoCursor: MongoCursor[Document] = findInterable.iterator()
       if (mongoCursor.hasNext) {
         val newPath = new Path(targetpath + minid_goods_taobao.toString)
         hdfsfs.delete(newPath, true)
         val os = hdfsfs.create(newPath)
         val bw = new BufferedWriter(new OutputStreamWriter(os, charsetName = "utf-8"))
          while (mongoCursor.hasNext) {
           val document = mongoCursor.next()
           bw.write(document.toJson())
           bw.write( str = "\n")
           if (minid_goods_taobao.toString.compareTo(document.get("_id").toString) < 0) {</pre>
              minid_goods_taobao = new ObjectId(document.get("_id").toString)
```

```
minid_goods_taobao = new ObjectId(document.get("_id").toString)
            idatarows += 1
          }
          bw.close
          println("finish:" + minid_goods_taobao.toString + ",datarows:" + idatarows)
        else {
         inodata += 1
         println("no data." + inodata + ",datarows:" + idatarows)
        Thread sleep 5000
    } catch {
      case e: Exception =>
       println(e.getClass.getName + ": " + e.getMessage)
    }
    hdfsfs.close
  }
}
```

#### 运行截图:

```
GoodsHdfs2Mysql × GoodHdfs2Mysql (1)
🔁 GoodMongo2Hdfs (2) 🗵
 no data.23,datarows:1632
 filter:collector:廖臻明,mintime:2022-11-26T03:10:46.382Z,minid:6382d52c9fda80d1174de088
 no data.24,datarows:1632
 filter:collector:廖臻明,mintime:2022-11-26T03:10:51.430Z,minid:6382d52c9fda80d1174de088
 no data.25.datarows:1632
 filter:collector:廖臻明,mintime:2022-11-26T03:10:56.523Z,minid:6382d52c9fda80d1174de088
 no data.26,datarows:1632
 filter:collector:廖臻明,mintime:2022-11-26T03:11:01.530Z,minid:6382d52c9fda80d1174de088
 finish:6382d5439fda80d1174de089,datarows:1633
 filter:collector:廖臻明,mintime:2022-11-26T03:11:06.586Z,minid:6382d5439fda80d1174de089
 finish:6382d54a9fda80d1174de08c,datarows:1636
 filter:collector:廖臻明,mintime:2022-11-26T03:11:11.705Z,minid:6382d54a9fda80d1174de08c
 finish:6382d54e9fda80d1174de18e,datarows:1894
 filter:collector:廖臻明,mintime:2022-11-26T03:11:16.876Z,minid:6382d54e9fda80d1174de18e
 no data.27,datarows:1894
 filter:collector:廖臻明,mintime:2022-11-26T03:11:21.885Z,minid:6382d54e9fda80d1174de18e
 no data.28, datarows:1894
 filter:collector:廖臻明,mintime:2022-11-26T03:11:26.892Z,minid:6382d54e9fda80d1174de18e
 no data.29,datarows:1894
```

在hadoop集群上运行截图:

#### 1.3.4. 数据分析与存储(spark streaming)(30分)

功能说明:使用spark streaming分析数据,根据销量取前15名的商品,统计营收前15名的商铺。

#### 功能设计:

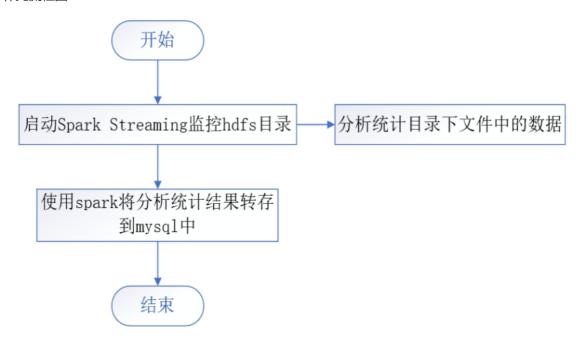
#### 1、根据销量,取前15名的商品:

第一步,连接Hdfs文件系统,遍历每条数据,取出商品标题、销量。第二步,根据店铺对数据进行分组求和。第三步,根据商品的销量进行排序,取前15名。第四步,连接mysql数据库,将分析出的15条结果插入mysql数据库。

#### 2、统计营收前15名的商铺:

第一步,连接Hdfs文件系统,遍历每条数据,取出店铺字段,计算店铺内的某款产品的收入。第二步,根据店铺对数据进行分组求和,求得每家店铺的营收。第三步,根据店铺的营收进行排序,取前15名。第四步,连接mysql数据库,将分析出的15条结果插入mysql数据库。

#### <补充流程图>



1、根据销量取前15名的商品:

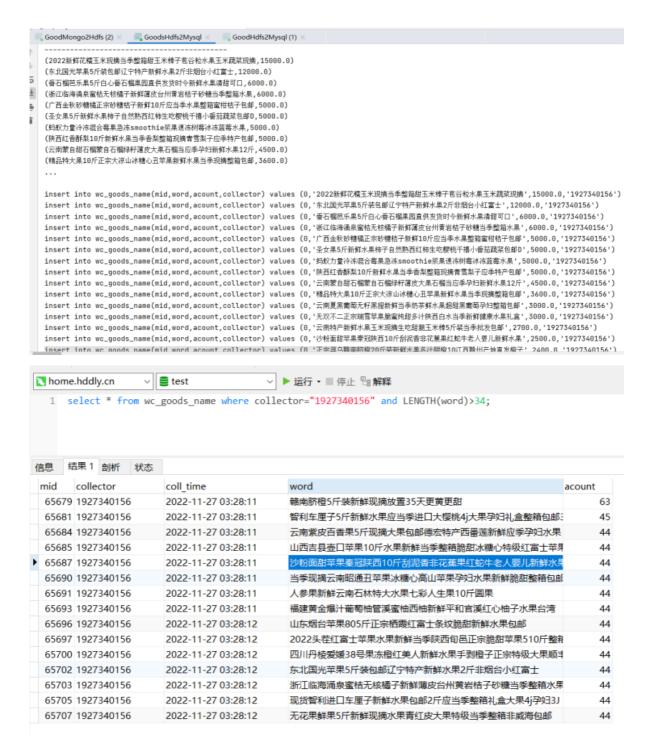
```
package traind
import ...
 * 统计产品销量前15名
Jobject GoodsHdfs2Mysql {
def main(args: Array[String]): Unit = {
    //hdfs dfs -rm -r /user/liaozhenming1/goods_taobao1/
     val streamingpath = "hdfs://master:9864/user/liaozhenming1/goods_taobao1/"
    System.setProperty("HADOOP_USER_NAME", "root");
    val sparkconf = new SparkConf().setMaster("local[*]").setAppName("goodscount") //local[*],spark://master:7077
    val ssc = new StreamingContext(sparkconf, Seconds(10))
     println("streamingpath:" + streamingpath)
    val lines = ssc.textFileStream(streamingpath)
    lines.print()
     val html = lines.map { line => val words = line.split( regex = ","); words(θ)}
     println(html)
     val ItemPairs = lines.map(line => {
      val document: Document = Document.parse(line)
      println("line:" + line);
      var strGoodsName = document.get("raw_title").toString()//"nick":是店铺名
      var floatSales = document.get("view_sales").toString().toDouble
      strGoodsName = StrUtil.replaceSumbo(strGoodsName)
      println("goods:" + strGoodsName);
      (strGoodsName, floatSales)
    1)
    ItemPairs.print()
    val htmlCount = ItemPairs.reduceByKeyAndWindow((v1: Double, v2: Double) => v1 + v2, Seconds(60 * 60 * 24), Seconds(20))
         窗口长度为1天,滑动时间为20秒,近一天内采集排名
    htmlCount.print()
    //根据门店营收进行排名,取前10名
     val hottestHtml = htmlCount.transform(itemRDD => {
      val top10 = itemRDD.map(pair => (pair._2, pair._1)).sortByKey( ascending = false).map(pair => (pair._2, pair._1)).take( num = 15)
      ssc.sparkContext.makeRDD(top10).repartition( numPartitions = 1)
             ssc.sparkContext.makeRDD(top10)
    1-)
     hottestHtml.print()
     val url = "jdbc:mysql://home.hddly.cn:53386/test?useSSL=false&useUnicode=true&characterEncoding=UTF-8&autoReconnect=true&failOve
     val user = "test"
     val password = "test"
    val mycode = "1927340156" //请将1927100100改为自己的学号
     hottestHtml.foreachRDD(rdd => {
       rdd.foreachPartition(partitionOfRecords => {
         val conn = ConnectionPool.getConn(url, user, password)
         conn.setAutoCommit(false)
         val stmt = conn.createStatement()
         var i = 1
         //将数据插入数据库
         partitionOfRecords.foreach(record => {
           val strsql = "insert into wc_goods_name(mid,word,acount,collector) values (0,'" +
            record._1 + "'," + record._2 + ",'" + mycode + "')"
           println(strsql)
           stmt.addBatch(strsql)
          i += 1
         })
         stmt.executeBatch()
         if (i > 1) {
          ConnectionPool.delete(conn, table = "wc_goods_name", mycode)
         conn.commit()
                  conn.close() //这里不能关闭,返回池中使用
      })
     1)
     ssc.start()
     ssc.awaitTermination()
     ssc.stop()
```

2. 统计营收前15名的商铺:

```
package traind
import ...
//统计门店营收前15名
object GoodHdfs2Mysql {
  def main(args: Array[String]): Unit = {
    //请将下方链接中myname 替换为本人姓名全拼
    //hdfs dfs -rm -r /user/liaozhenming1/goods_taobao1/
    val streamingpath = "hdfs://master:9864/user/liaozhenming1/goods_taobao1/"
    System.setProperty("HADOOP_USER_NAME", "root");
    val sparkconf = new SparkConf().setMaster("local[*]").setAppName("goodscount") //local[*],spark://master:7077
    val ssc = new StreamingContext(sparkconf, Seconds(10))
    println("streamingpath:" + streamingpath)
    val lines = ssc.textFileStream(streamingpath)
    lines.print()
    val html = lines.map { line => val words = line.split( regex = ","); words(θ)}
    println(html)
    val ItemPairs = lines.map(line => {
     val document: Document = Document.parse(line)
     println("line:" + line);
     var strGoodName = document.get("nick").toString()//"nick":是店铆名
     strGoodName = StrUtil.replaceSymbo(strGoodName)
     var intSales:Integer =Integer.valueOf(document.get("view_sales").toString())//该店铺某款产品的销售量
     var floatPrice = document.get("view_price").toString().toDouble//该店铺某款产品的价格
     val add4=(x:Int,y:Double)=>x * y
     var totalSale = add4(intSales,floatPrice)//该店某款产品的营收
     println("goods:" + strGoodName);
     (strGoodName, totalSale)
    })
  ItemPairs.print()
  val htmlCount = ItemPairs.reduceByKeyAndWindow((v1: Double, v2: Double) => v1 + v2, Seconds(60 * 60 * 24), Seconds(20))
  // 窗口长度为1天,滑动时间为20秒,近一天内采集排名
  htmlCount.print()
  val hottestHtml = htmlCount.transform(itemRDD => {
   val top10 = itemRDD.map(pair => (pair._2, pair._1)).sortByKey( ascending = false).map(pair => (pair._2, pair._1)).take( num = 15)
    ssc.sparkContext.makeRDD(top10).repartition( numPartitions = 1)
           ssc.sparkContext.makeRDD(top10)
  3-)
  hottestHtml.print()
  val url = "jdbc:mysql://home.hddly.cn:53306/test?useSSL=false&useUnicode=true&characterEncoding=UTF-8&autoReconnect=true&failOver
  val user = "test"
  val password = "test"
  val mycode = "1927340156" //请将1927100100改为自己的学号
    hottestHtml.foreachRDD(rdd => {
     rdd.foreachPartition(partitionOfRecords => {
       val conn = ConnectionPool.getConn(url, user, password)
        conn.setAutoCommit(false)
       val stmt = conn.createStatement()
       var i = 1
        //将数据插入数据库
        partitionOfRecords.foreach(record => {
         val strsql = "insert into wc_goods_name(mid,word,acount,collector) values (θ,'" +
          record._1 + "'," + record._2 + ",'" + mycode + "')"
         println(strsql)
         stmt.addBatch(strsql)
        1)
        stmt.executeBatch()
         if (i > 1) {
           ConnectionPool.delete(conn, "wc_goods_name", mycode)
         }-
        conn.commit()
                 conn.close() //这里不能关闭, 返回池中使用
      1)
    ssc.start()
    ssc.awaitTermination()
    ssc.stop()
```

#### 运行截图:

1、根据销量,取前15名的商品:

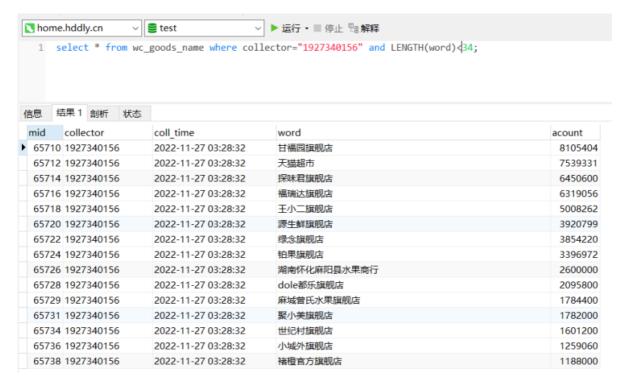


#### 在hadoop集群上运行:

spark-submit --master yarn --class traind.GoodsHdfs2Mysql /home/hadoop/train.jar

#### 2、统计营收前15名的商铺:

```
GoodMongo2Hdfs (2) ×
                                                 🏿 GoodsHdfs2Mysql 🗴 🛮 🗐 GoodHdfs2Mysql (1) 🗵
   (OOLe都水旗舰店,1/96488.8)
   (铂果旗舰店,768000.0)
   (jojo1987312,750000.0)
   (王小二旗舰店,605952.0)
   ( 小蝴蝶, 527340.0)
   (曼溪小铺,523200.0)
   (果味颂旗舰店,498240.0)
   (鞍山南果梨旗舰店,479600.0)
   (植都旗舰店,426240.0)
   (爱果精神旗舰店,363000.0)
   insert into wc_goods_name(mid,word,acount,collector) values (0,'dole都乐旗舰店',1796400.0,'1927340156')
   insert into wc_goods_name(mid,word,acount,collector) values (0,'铂果旗舰店',768000.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (0,'jojo1987312',750000.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (0,'王小二旗舰店',605952.0,'1927340156')
   insert into wc_goods_name(mid,word,acount,collector) values (0,'\ 小蝴蝶',527340.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (θ,'曼溪小铺',523200.0,'1927340156')
   insert into wc_goods_name(mid,word,acount,collector) values (0,'果味颂旗躵店',498240.0,'1927340156')
   insert into wc_goods_name(mid,word,acount,collector) values (0,'鞍山南果梨旗舰店',479600.0,'1927340156')
   insert into wc_goods_name(mid,word,acount,collector) values (0,'植都旗舰店',426240.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (0,'愛果精神旗舰店',363000.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (0,'鲜蜂队旗舰店',354239.9999999994,'1927340156')
   insert into wc_goods_name(mid,word,acount,collector) values (0,'顾一旗舰店',309600.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (0,'果鲜你水果旗舰店',302400.0,'1927340156')
  insert into wc_goods_name(mid,word,acount,collector) values (0,'a402811915',301680.0,'1927340156') insert into wc_goods_name(mid,word,acount,collector) values (0,'a402811915',acount,collector) insert into wc_goods_name(mid,word,acount,collector) insert into wc_goods_name(mid,word,acount,collector) values (0,'a402811915',acount,collector) insert into wc_goods_name(mid,word,acount,collector) insert insert into wc_goods_name(mid,word,acount,collector) insert into wc_goods_name(mid,word,acount,collector) insert into wc_goods_name(mid,word,acount,collector) insert into wc_goods_name(mid,word,acount,co
  insert into wc_goods_name(mid,word,acount,collector) values (0,'生果乐旗舰店',287040.0,'1927340156')
   Time: 1669518780000 ms
```



#### 在hadoop集群上运行:

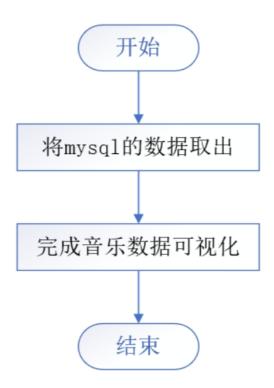
spark-submit --master yarn --class traind.GoodHdfs2Mysql /home/hadoop/train.jar

#### 1.1.1. 数据可视(10分)

功能说明:使用python语言绘制可视化图片,并上传到web端

功能设计:使用python语言的barplot()根据存储在mysql的分析结果绘制直方图图片,并将这些图片上传到指定的web网页。

<补充流程图>



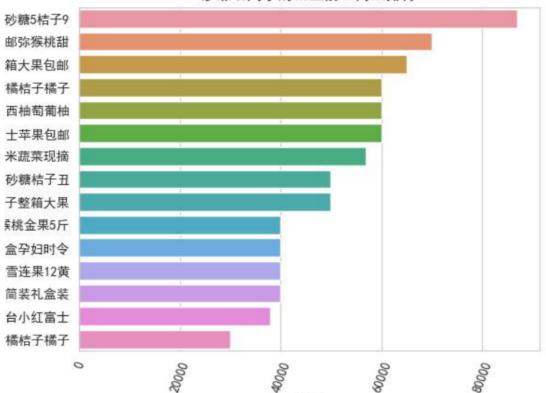
#### 源码实现:

```
# coding=utf-8
from bgutils.ftpUtil import ftpUtil
from bgutils.mysqlUtil import mysqlUtil
from matplotlib import pyplot as plt
≙import seaborn as sns
# 数据可视化
class drawpic:
    def draw(self,myname,mycode):
        mutil= mysqlUtil()
        ftputil=ftpUtil()
        # SELECT * FROM test.wc_goods_name where word!='' order by acount desc;
        SQL = '''select * from wc_goods_name where word!="" and collector="''+ \
mycode +''' and LENGTH(word)<34 order by acount desc_'''
        # 调用 mysqlconn 类的 query() 方法
        df_data = mutil.query(sql=SQL)
        # df_data.to_csv
        # 使用seaborn库绘图
        sns.set_style('whitegrid', {'font.sans-serif':['simhei', 'Arial']})
        # 设置中文字体
        #plt.rcParams['font.sans-serif'] = font
        # 设置正常显示负号
        plt.rcParams['axes.unicode_minus']=False
        count = df_data['acount']
        index = df_data['word']
        sns.barplot(x=count, y=index)
        plt.xticks(rotation=70)
        plt.xlabel('数量')
        plt.ylabel('商品名称')
        plt.title(myname+'同学的商铺营收前15名排行')
        # plt.show()
        idxproject= "22002" #22001是项目编号,长度5位
        idxpic="06" #04是图片在该项目中的顺序号,长度2位
        filepath=mycode+"_"+idxpic +".jpg" #生成的文件以学号+序号+后缀组成
        plt.savefig(filepath)
        ftputil.putfile_stud(filepath, idxproject,mycode,idxpic);
```

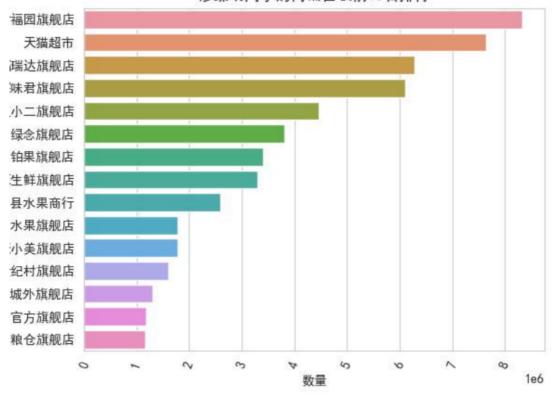
```
| if __name__ == "__main__":
| draw=drawpic();
| myname="廖臻明" #请将引号中的myname替换为本人姓名,如: 张三
| mycode="1927340156" #请将引号中的mycode替换为本人学号,如:2019001001
| draw.draw(myname,mycode);
```

#### 运行截图:

## 廖臻明同学的销量前15商品排行



## 廖臻明同学的商铺营收前15名排行



## 1.4. 项目小结 (5分)

本次实验主要是将Spark专业课程的学习内容进行巩固。在基础环境搭建完成的情况下,先是在pyecharm上搭建Scrapy框架,使用Scrapy框架采集电商平台水果类目的数据,并且将数据存储进mongodb数据库;为了方便后续的分析使用scala将mongodb数据库的数据转移至Hdfs文件系统当中;再使用spark streaming分析数据,根据销量取前15名的商品,统计营收前15名的商铺,并将分析结果存储进mysql数据库;最后使用pyecharm或jupyter工具运用python将存储在mysql的分析结果进行可视化,并将可视化结果上传至web端以便将来进行展示。

总体过程比较复杂,但难度简单,加深了对专业知识的巩固和认识。主要难度集中在各种环境的搭建上,可能配置不当会不停报错,在老师和同学的帮助下顺利完成。

通过这一次实训让我更加熟练的使用Spark Streaming进行数据处理和数据分析,相信这次实训对日后的工作有较大的帮助。

#### 1.5. 附件

整个过程中实现的源码,包括采集的源码,spark分析源码,可视化的源码,全部打包成rar文件另行提交。

## 实现参考

### 数据采集

源码参考:

https://jihulab.com/biglab-share/scrapy/-/tree/main/b59510/chap9/music\_scrapy?ref\_type=heads

## 数据存储(mongodb->hdfs)

源码参考:

https://jihulab.com/biglab-share/spark/-/blob/main/b59510/src/traina/MusicMongo2Hdfs.scala?ref\_type=heads\_

## 数据分析与存储(spark streaming)

源码参考

https://jihulab.com/biglab-share/spark/-/blob/main/b59510/src/traina/MusicHdfs2Mysql.scala?ref\_t\_wpe=heads

## 数据可视

源码参考一(pycharm):

https://jihulab.com/biglab-share/scrapy/-/tree/main/b59510/chap9/music view?ref type=heads

### 如何下载源码

以上源码包含在两个git项目:

- 1 https://jihulab.com/biglab-share/scrapy.git
- 2 https://jihulab.com/biglab-share/spark.git

由于 jihulab.com进入收费,我们即将切换到自建git服务器上,切换方法如:

使用资源管理器进入scrapy目录,在地址栏中输入cmd进入命令行,运行:

git remote set-url origin http://home.hddly.cn:8093/biglab-share/scrapy.git

同样,使用资源管理器进入spark目录,在地址栏中输入cmd进入命令行,运行:

1 git remote set-url origin http://home.hddly.cn:8093/biglab-share/spark.git

GIT安装和下载源码视频

学习视频