

# 第3章Spark编程实训

---

## 第3章Spark编程实训

### CH03实训1：使用spark-shell进行词频统计

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

环境准备参考

数据准备

在spark-shell上执行

spark-shell预期输出：

在hdfs上观察

学习视频

### CH03实训2：通过spark编程统计某月份的客户总消费金额

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

数据准备

spark-shell脚本

预期输出

### CH03实训3：通过Spark编程计算各城市的平均气温

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

数据准备

在spark-shell上执行

spark-shell预期输出：

---

## CH03实训1：使用spark-shell进行词频统计

---

### 训练要点

1. 掌握RDD创建方法，读数数据文件和存储文件的方法。
2. 掌握flatMap操作、map操作、reduceByKey操作、filter操作方法。

### 需求说明

数据文件words.txt如下文所示，文件中包含了多行句子，现在要求对文档中的单词计数，并把单词计数超过3的结果存储到HDFS上

---

words.txt

What is going on there?

I talked to John on email. We talked about some computer stuff that's it.

I went bike riding in the rain, it was not that cold.

We went to the museum in SF yesterday it was \$3 to get in and they had free food. At the same time was a SF Giants game, when we got done we had to take the train with all the Giants fans, they are 1/2 drunk.

---

可以在linux下使用wget获取该文件:

wget <http://bigdata.hddly.cn/b46488/file/chap3/words.txt>

## 实现思路及步骤

1. 通过textFile的方法读取数据。
2. 通过flatMap将字符串切分成单词。
3. 通过map将单词转化为(单词, 1)的形式。
4. 通过reduceByKey将同一个单词的所有值相加。
5. 通过filter将单词计数大于3的结果过滤出来。
6. 通过saveAsTextFile将结果写入到HDFS

## 作业要求

1. 将测试文件words.txt上传到hdfs上的/user/myname/目录下,并截图(myname要改为自己名字全拼)
2. 通过crt进入spark-shell命令窗口, 执行词频分析脚本, 其中输出结果文件存于/user/myname/output\_wordcount, 并截图运行过程;
3. 提交hdfs上生成的结果文件: /user/myname/output\_wordcount 的内容截图, 图中需显示完整的文件路径。

## 实现参考

### 环境准备参考

```
1 //启动hadoop集群
2 start-all.sh
3 //启动spark集群
4 cd /usr/local/spark/sbin
5 ./start-all.sh
```

使用jps检查master主机进程

```
1 [root@master sbin]# jps
2 2421 SecondaryNameNode
3 2134 NameNode
4 2710 ResourceManager
5 3150 Master
6 3902 Jps
7 [root@master sbin]
```

## 数据准备

在slave1下执行

```
1 cd /root/spark
2 wget https://biglab.site//b59510spark/file/spark_t_data.tar.gz
3 tar -xzvf ./spark_t_data.tar.gz
4 hdfs dfs -rm /user/myname/words.txt
5 hdfs dfs -rm -r /user/myname/output_wordcount
6 hdfs dfs -put ./spark_t_data/words.txt /user/myname/
```

## 在spark-shell上执行

```
1 spark-shell
```

```
1 val worddata=sc.textFile("/user/myname/words.txt").flatMap(x=>x.split("
")).map(x=>(x,1)).reduceByKey((x,y)=>x+y)
2
3 val worddata_3=worddata.filter(x=>x._2>=2).map(x=>x._1)
4
5 worddata_3.repartition(1).saveAsTextFile("/user/myname/output_wordcount")
```

spark-shell预期输出：

```
1 scala> val
  worddata=sc.textFile("/user/myname/words.txt").flatMap(x=>x.split("
  ")).map(x=>(x,1)).reduceByKey((x,y)=>x+y)
2 worddata: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[16] at
  reduceByKey at <console>:24
3
4 scala>
5
6 scala> val worddata_3=worddata.filter(x=>x._2>3).map(x=>x._1)
7 worddata_3: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[18] at map at
  <console>:25
8
9 scala>
10
11 scala>
  worddata_3.repartition(1).saveAsTextFile("/user/myname/output_wordcount")
12
13 scala>
```

## 在hdfs上观察

```
1 | /user/myname/output_wordcount
```

## 学习视频

[学习视频](#)

# CH03实训2：通过spark编程统计某月份的客户总消费金额

## 训练要点

1. 掌握数据读取和存储的方法。
2. 掌握RDD基本操作

## 需求说明

大数据和人工智能技术的应用，使得我国广大科技工作者能够把握大势、抢占先机、直面问题、迎难而上，瞄准世界科技前沿，引领科技发展方向。

本任务中，某互联网企业，创建了线上购物平台，开拓了新的商品销售渠道。现有一份某电商2020年12月份的订单数据文件online retail.csv，记录了每位顾客每笔订单的购物情况，包含3个数据字段，字段说明如表3.7所示。因为该电商准备给重要的客户发放购物津贴作为福利回馈，提高顾客满意度，所以需要统计每位客户的总消费金额，并筛选出消费金额排在前50名的客户。

某电商的订单数据字段说明：

字段名称	说明
Invoice	订单编号
Price	订单价格 (单位: 元)
Customer ID	客户编号

## 实现思路及步骤

1. 读取数据并创建RDD.
2. 通过map()方法分割数据, 选择客户编号和订单价格字段组成键值对数据.
3. 使用reduceByKey()方法计算每位客户的总消费金额.
4. 使用sortBy()方法对每位客户的总消费金额进行降序排序, 取出前50条数据.

## 作业要求

提交源码和脚本运行过程截图

## 实现参考

### 数据准备

在slave1上运行:

```
1 | cd /root/spark
2 | wget https://biglab.site//b59510spark/file/spark_t_data.tar.gz
3 | tar -xzvf ./spark_t_data.tar.gz
4 | hdfs dfs -put ./spark_t_data/online_retail.txt /user/myname/
```

### spark-shell脚本

```
1 | //在Spark-shell中 读取上传到HDFS上的数据
2 | val input = sc.textFile("/user/myname/online_retail.txt")
3 |
4 | //去除csv数据文件的第一行
5 | val cutinput = input.mapPartitionsWithIndex((ix, it) => {
6 |     if (ix == 0) it.drop(1)
7 |     it
8 | })
9 | cutinput.collect
10 |
11 | val re = cutinput.map(line => {val data = line.split(",");(data(0),
12 |     data(1).toDouble)}).reduceByKey(_+_ )
13 | val re_sort = re.filter(x=> !(x._1=="")).sortBy(x=>x._2,false)
14 |
15 | re_sort.take(50)
16 |
```

## 预期输出

```
1 scala> //在Spark-shell中 读取上传到HDFS上的数据
2
3 scala> val input = sc.textFile("/user/myname/online_retail.txt")
4 input: org.apache.spark.rdd.RDD[String] = /user/myname/online_retail.txt
  MapPartitionsRDD[1] at textFile at <console>:24
5
6 scala>
7
8 scala> //去除csv数据文件的第一行
9
10 scala> val cutinput = input.mapPartitionsWithIndex((ix, it) => {
11     |   if (ix == 0) it.drop(1)
12     |   it
13     | })
14 cutinput: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at
  mapPartitionsWithIndex at <console>:25
15
16 scala> cutinput.collect
17 res0: Array[String] = Array(17850,2.55,536365, 17850,3.39,536365,
  17850,2.75,536365, 17850,3.39,536365, 17850,3.39,536365, 17850,7.65,536365,
  17850,4.25,536365, 17850,1.85,536366, 17850,1.85,536366, 13047,4.25,536368,
  13047,4.95,536368, 13047,4.95,536368, 13047,4.95,536368, 13047,1.69,536367,
  13047,2.1,536367, 13047,2.1,536367, 13047,3.75,536367, 13047,1.65,536367,
  13047,4.25,536367, 13047,4.95,536367, 13047,9.95,536367, 13047,5.95,536367,
  13047,5.95,536367, 13047,7.95,536367, 13047,7.95,536367, 13047,5.95,536369,
  12583,3.75,536370, 12583,3.75,536370, 12583,3.75,536370, 12583,0.85,536370,
  12583,0.65,536370, 12583,0.85,536370, 12583,1.25,536370, 12583,2.95,536370,
  12583,2.95,536370, 12583,1.95,536370, 12583,1.95,536370, 12583,1.95,536370,
  12583,0.85,536370, 1258...
18
19 scala>
20
21 scala> val re = cutinput.map(line => {val data = line.split(",");(data(0),
  data(1).toDouble)}).reduceByKey(_+_ )
22 re: org.apache.spark.rdd.RDD[(String, Double)] = ShuffledRDD[4] at
  reduceByKey at <console>:25
23
24 scala>
25
26 scala> val re_sort = re.filter(x=> !(x._1=="")).sortBy(x=>x._2,false)
27 re_sort: org.apache.spark.rdd.RDD[(String, Double)] = MapPartitionsRDD[10] at
  sortBy at <console>:25
28
29 scala>
30
31 scala> re_sort.take(50)
```

```
32 res1: Array[(String, Double)] = Array((12748,1618.1500000000015),
    (14911,1573.1600000000014), (17850,1176.2299999999982),
    (17841,1073.1299999999997), (14606,828.5199999999996),
    (16607,726.3199999999999), (14527,666.1399999999996),
    (17340,613.5099999999999), (15311,582.7899999999997),
    (15044,545.4599999999998), (13174,519.7500000000001),
    (14667,506.5499999999998), (15727,456.4499999999993),
    (17961,414.82000000000005), (14030,413.75), (18116,410.9799999999973),
    (15039,404.6699999999998), (16873,401.4299999999999),
    (18118,389.7099999999998), (15574,386.2499999999999),
    (14180,382.2499999999994), (16713,377.5199999999987),
    (18055,377.3799999999999), (14505,373.2999999999956),
    (15498,369.6699999999996), (15808,366.2999999999998),
    (15570,363.5099999999993), (12...
33
34 scala>
```

## CH03实训3：通过Spark编程计算各城市的平均气温

### 训练要点

1. 掌握RDD创建方法
2. 掌握map()、groupBy()、mapvalues()、reduce()方法的使用

### 需求说明

天气预报每天都会显示各城市的温度，方便出行人士根据当天的温度穿上合适的衣物现有一份各城市的温度数据文件avgTemperature.txt,数据如表所示，记录了某段时间范围内各城市每天的温度，文件中每一行数据分别表示城市名和温度，现要求使用spark 编程计算出各城市的平均气温。

#### 各城市温度部分数据

```
1 | beijing 28.1
2 | shanghai 28.7
3 | guangzhou 32.0
4 | shenzhen 33.1
5 | beijing 27.3
6 | shanghai 30.1
7 | guangzhou 33.3
8 | shenzhen 28.6
9 | beijing 28.2
10 | shanghai 29.1
11 | guangzhou 32.0
12 | shenzhen 32.1
```

## 实现思路及步骤

1. 通过textFile()方法读取数据创建RDD
2. 使用map()方法将输入数据按制表符(TAB)进行分割，并转化成（城市，温度）的形式。
3. 使用groupBy()方法按城市分组，得到每个城市对应的所有温度。
4. 使用mapvalues()和reduce()方法计算各城市平均气温

## 作业要求

提交源码和脚本运行过程截图

## 实现参考

### 数据准备

```
1 cd /root/spark
2 wget https://biglab.site//b59510spark/file/spark_t_data.tar.gz
3 tar -xzvf ./spark_t_data.tar.gz
4 hdfs dfs -put ./spark_t_data/avgTemperature.txt /user/myname/
```

### 在spark-shell上执行

```
1 val input = sc.textFile("/user/myname/avgTemperature.txt")
2 // 利用map算子转换数据类型
3 val mapRDD = input.map(
4     data => {
5         val datas = data.split("\t")
6         (datas(0), datas(1).toDouble)
7     }
8 )
9 //按城市分组
10 val groupRDD = mapRDD.groupBy(_._1)
11 // 计算各城市平均气温
12 val avgTmp = groupRDD.mapValues(data => {
13     val a = data.reduce((x, y) => ("", x._2 + y._2))
14     a._2 / data.size
15 })
16 avgTmp.collect
```

### spark-shell预期输出：

```
1 scala> val
worddata=sc.textFile("/user/myname/words.txt").flatMap(x=>x.split("
")).map(x=>(x,1)).reduceByKey((x,y)=>x+y)
2 worddata: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[16] at
reduceByKey at <console>:24
3
4 scala>
5
6 scala> val worddata_3=worddata.filter(x=>x._2>3).map(x=>x._1)
```



```

7 worddata_3: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[18] at map at
  <console>:25
8
9 scala>
10
11 scala>
worddata_3.repartition(1).saveAsTextFile("/user/myname/output_wordcount")
12
13 scala> val input = sc.textFile("/user/myname/avgTemperature.txt")
14 input: org.apache.spark.rdd.RDD[String] = /user/myname/avgTemperature.txt
MapPartitionsRDD[25] at textFile at <console>:24
15
16 scala> // 利用map算子转换数据类型
17
18 scala> val mapRDD = input.map(
19     |     data => {
20     |     val datas = data.split("\t")
21     |     (datas(0), datas(1).toDouble)
22     |     }
23     |     )
24 mapRDD: org.apache.spark.rdd.RDD[(String, Double)] = MapPartitionsRDD[26] at
map at <console>:25
25
26 scala> //按城市分组
27
28 scala> val groupRDD = mapRDD.groupBy(_._1)
29 groupRDD: org.apache.spark.rdd.RDD[(String, Iterable[(String, Double)])] =
ShuffledRDD[28] at groupBy at <console>:25
30
31 scala> // 计算各城市平均气温
32
33 scala> val avgTmp = groupRDD.mapValues(data => {
34     |     val a = data.reduce((x, y) => ("", x._2 + y._2))
35     |     a._2 / data.size
36     |     })
37 avgTmp: org.apache.spark.rdd.RDD[(String, Double)] = MapPartitionsRDD[29] at
mapValues at <console>:25
38
39 scala> avgTmp.collect
40 res3: Array[(String, Double)] = Array((guangzhou,32.43333333333333),
(beijing,27.866666666666667), (shenzhen,32.6), (shanghai,29.3))
41
42 scala>

```