# 181840356-周娴静-实验4报告

## 一、题目

1. 分别编写MapReduce程序和Spark程序统计双十一最热门的商品和最受年轻人(age<30)关注的商家（"添加购物车+购买+添加收藏夹"前100名）；

2. 编写Spark程序统计双十一购买了商品的男女比例，以及购买了商品的买家年龄段的比例；

3. 基于Hive或者Spark SQL查询双十一购买了商品的男女比例，以及购买了商品的买家年龄段的比例；

4. 预测给定的商家中，哪些新消费者在未来会成为忠实客户，即需要预测这些新消费者在6个月内再次购买的概率。基于Spark MLlib编写程序预测回头客，评估实验结果的准确率。

| 用户行为日志 | user_log_format1 |
| --- | --- |
| 用户画像 | user_info_format1 |
| 训练数据 | train_format1 |
| 测试数据 | test_format1 |

## 二、设计思路及程序关键代码

### 1. Task1

#### （1）Task11（Mapreduce）

- 类比作业5 wordcount shakespeare
- 【限定时间为双十一】首先筛选出time_stamp=1111的商品
- 【对商品计数】筛选action_type !=0的数据，以商品item_id作为它的key

```
1    public void map(Object key, Text value, Context context
2    ) throws IOException, InterruptedException {
3        StringTokenizer itr = new StringTokenizer(value.toString());
4        while (itr.hasMoreTokens()) {
5            word.set(itr.nextToken());
6            //把每一行拆解成一个一个属性,下标从0开始
7            String[] temp=word.toString().split(",");
8            //首先筛选出time_stamp=1111的商品    [5]
9            //对action_type，如果是1，2，3（不等于0）则计数   [6]
10           if(temp[5].equals("1111"))
11               if(!temp[6].equals("0"))
```

```
12              {
13                  //以商品id   item_id作为它的key    [1]
14                  context.write(new Text(temp[1]), one);
15              }
16          }
17      }
18  }
19 }
```

## (2) Task12（Mapreduce）

- 与Task11类似
- 仿照作业5（wordcount shakespeare）设置停词文件user_info_skip.csv
- 存储patternsToSkip的HashSet，供后续使用

```
1 public void setup(Context context) throws IOException,
  InterruptedException {
2     conf = context.getConfiguration( );
3     //caseSensitive = conf.getBoolean ( "wordcount.case.sensitive",true);
4     Path patternsPath=new
  Path("hdfs://localhost:9000/Task12/user_info_skip.csv");
5     String patternsFileName = patternsPath.getName( ).toString();
6     parseSkipFile(patternsFileName);
7     //     }
8     //}
9 }
```
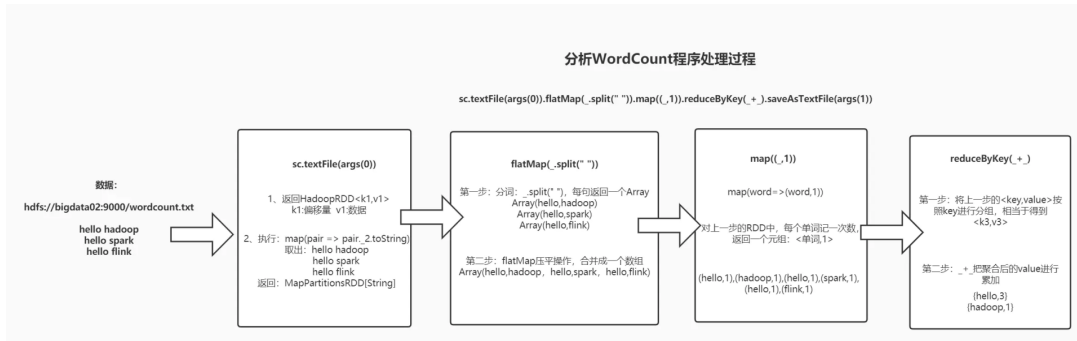
```
1 while(itr.hasMoreTokens()){
2     word.set(itr.nextToken());
3     //把每一行拆解成一个一个属性,下标从0开始
4     String[] temp=word.toString().split(",");
5     //首先筛选出time_stamp=1111的商品    [5]
6     //对action_type，如果是1，2，3（不等于0）则计数   [6]
7     //购物者是年轻人     [0]
8     if(temp[5].equals("1111"))
9         if(!temp[6].equals("0"))
10            if(! patternsToSkip.contains(temp[0])){
11                //以商家id   merchant_id作为它的key    [3]
12                context.write(new Text(temp[3]),one);        //用
  context.write收集<key,value>对
13                Counter counter =
  context.getCounter(CountersEnum.class.getName(),CountersEnum.INPUT_WORDS.t
  oString());  //2.0
14                counter.increment(1);   //2.0
15    }
```

## (3) Task11（Spark）



分析WordCount程序处理过程

sc.textFile(args(0)).flatMap(_.split(" ")).map((_,1)).reduceByKey(_+_).saveAsTextFile(args(1))

- 读取数据文件，去除第一行
- 首先筛选出双十一（"1111"）的数据
- 进一步筛选action_type !=0的数据
- 实现二次排序

```
1  val
   data=sc.textFile("hdfs://localhost:9000/Task12/input2/user_log_format1.csv
   ").flatMap(_.split("\n"))
2  val arr = data.take(1)
3  val data1 = data.filter(!arr.contains(_)).filter(line=>line.split(",")
   (5).equals("1111")).map{
6          line => (line.split(",")(1), line.split(",")(6))
7        }.mapValues(_.toInt)
8  val data2 = data1.filter(value => value._2>0)
10 val data3 = data2.countByKey().toSeq.sortWith(_._2>_._2).take(100)
12 val data4 = sc.parallelize(data3)
13 data4.saveAsTextFile("hdfs://localhost:9000/Task12/output1-1")
16 :quit
19 hdfs dfs -cat /Task12/output1-1/*
```

## (4) Task12（Spark）

```
1  val log =
   spark.read.format("csv").option("header","true").load("hdfs://localhost:90
   00/Task12/input2/user_log_format1.csv")
3  val info =
   spark.read.format("csv").option("header","true").load("hdfs://localhost:90
   00/Task12/user_info_format1.csv")
5  val info1 = info.filter("age_range<4 and
   age_range>0").select("user_id","age_range")
6  val log1 = log.filter("time_stamp=1111 and
   action_type!=0").select("user_id","seller_id","action_type")
8  val dfjoin = info1.join(log1,"user_id")
10 val dfss = dfjoin.groupBy("seller_id").count()
12   //orderby count in desc and take first 100
```

```
13  val rddss = dfss.orderBy(dfss("count").desc).rdd.map(x=>
    (x(0),x(1))).take(100)
14    //save as testFile
15  sc.parallelize(rddss).saveAsTextFile("hdfs://localhost:9000/Task12/output1
    -2")
16  :quit
18  hdfs dfs -cat /Task12/outp
```

## 2. Task2

### (1) Task21

- 基本数据筛选，筛选出双十一购买商品的数据条以及gender为（0，1）的数据
- 进行dataframe的join
- 对gender作进一步分组，并统计

```
1   val log
    =spark.read.format("csv").option("header","true").load("hdfs://localhost:9
    000/Task12/input2/user_log_format1.csv")
3   val info
    =spark.read.format("csv").option("header","true").load("hdfs://localhost:9
    000/Task12/user_info_format1.csv")
5   val log1 =
    log.select("user_id","time_stamp","action_type").filter("time_stamp=1111
    and action_type=2")
6   val info1 = info.select("user_id","gender").filter("gender=1 or gender=0")
8   val dfjoin1 = info1.join(log1,"user_id")
10  val dfcount1 = dfjoin1.groupBy("gender").count()
12  dfcount1.withColumn("ratio",dfcount1("count")/dfjoin1.count).show
15  val rddr1 =
    dfcount1.withColumn("ratio",dfcount1("count")/dfjoin1.count).select("gende
    r","ratio").rdd.map(x=>(x(0),x(1))).collect()
16  sc.parallelize(rddr1).saveAsTextFile("hdfs://localhost:9000/Task12/output2
    -1")
```

### (2) Task22

- 与Task21类似

```
1   val log
    =spark.read.format("csv").option("header","true").load("hdfs://localhost:9
    000/Task12/input2/user_log_format1.csv")
3   val info
    =spark.read.format("csv").option("header","true").load("hdfs://localhost:9
    000/Task12/user_info_format1.csv")
```

```
5  val info1 = info.select("user_id","age_range").filter("age_range>0 and
   age_range<9")
6  val log1 =
   log.select("user_id","time_stamp","action_type").filter("time_stamp=1111
   and action_type=2")
8  val dfjoin = info1.join(log1,"user_id")
10 val dfcount = dfjoin.groupBy("age_range").count()
13 dfcount.withColumn("ratio",dfcount("count")/dfjoin.count).show
15 /*将答案保存到本地*/
16 val dfratio =
   dfcount.withColumn("ratio",dfcount("count")/dfjoin.count).select("age_rang
   e","ratio").orderBy("age_range")
18 val rddr = dfratio.rdd.map(x=>(x(0),x(1)))
20 val rddr = dfratio.rdd.map(x=>(x(0),x(1))).collect()
22 sc.parallelize(rddr).saveAsTextFile("hdfs://localhost:9000/Task12/output2-
   2")
```

## 3. Task3

### (1) Task31

- 将数据存储为dataframe以后进行createOrReplaceTempView试图转化，方便后续进行表的连接
- 采用嵌套查询方法，内部select初步筛选出购买了商品且gender为（0，1）的数据，设置distinct保证最多指读取该用户id一次
- 外部select进一步根据gender分组，并进行统计

```
1  //terminal
2  cd $SPARK_HOME
3  bin/spark-shell
4
5
6  //scala
7  val user_log =
   spark.read.format("csv").option("header","true").load("hdfs://localhost:90
   00/Task12/input2/user_log_format1.csv")
8
9  val user_info = spark.read.format("csv").option("header","true").load("hdf
   s://localhost:9000/Task12/user_info_format1.csv")
10
11 user_log.createOrReplaceTempView("table_log")
12
13 user_info.createOrReplaceTempView("table_info")
15 //3-1
17 val people1=spark.sql("select gender,count(*) as num from (select distinct
   a.user_id,gender from table_log a,table_info b where a.user_id=b.user_id
   and a.action_type='2' and gender in('0','1')) group by gender")
```

```
18  people1.show
19  //下面显示ratio
20  people1.agg("num"-> "sum").show //结果为407308
21  people1.withColumn("ratio",people1("num")/407308).show
```

## (2) Task32

- 与Task31类似

```
1  //3-2
2  val people2=spark.sql("select age_range,count(*) as num from (select
   distinct a.user_id,age_range from table_log a,table_info b
   where a.user_id=b.user_id and a.action_type='2'
   and age_range in('1','2','3','4','5','6','7','8')) group
   by age_range order by age_range")
3  people2.show
4  // //下面显示ratio(但不是sql语句
5  people2.agg("num"-> "sum").show //结果为329039
6  people2.withColumn("ratio",people2("num")/329039).show
```

## 4. Task4

**参考:**

https://www.jianshu.com/p/aa6cb1ef6f69

https://tianchi.aliyun.com/notebook-ai/detail?postId=143593

- 参考Baseline做法，选取age_range、gender、total_logs、unique_item_ids、categories、browse_days、one_clicks、shopping_carts、purchase_times、favourite_times 建立特征工程
- 对建立好的特征作缺失值处理
- 利用pyspark-ml随机森林进行预测

```
1   #jupyter notebook
2   # In[1]:
5   from pyspark.sql import SparkSession
6   spark = SparkSession.builder.appName('random_forest').getOrCreate()
7   df=spark.read.csv('hdfs://localhost:9000/Task12/feature.csv',inferSchema=True,header=True)
8   # In[2]:
10  from pyspark.ml.feature import VectorAssembler
11  df_assembler = VectorAssembler(inputCols=
    ['age_range','gender','total_logs','unique_item_ids','categories','browse_days','one_clicks','shopping_carts','purchase_times','favourite_times'], outputCol="features")
12  df = df_assembler.transform(df)
```

```
15   -测试集和训练集划分比例：0.75：0.25
16   -随机森林基学习器数量：120
17   # In[3]:
18   model_df = df.select(['features','label'])
19   train_df,test_df=model_df.randomSplit([0.75,0.25])
20   # In[4]:
22   from pyspark.ml.classification import RandomForestClassifier
23   rf_classifier=RandomForestClassifier(labelCol='label',numTrees=120).fit(tr
     ain_df)
24   rf_predictions=rf_classifier.transform(test_df)
26   # In[5]:
27   from pyspark.ml.evaluation import BinaryClassificationEvaluator
28   from pyspark.ml.evaluation import MulticlassClassificationEvaluator
30   rf_accuracy=MulticlassClassificationEvaluator(labelCol='label',metricName=
     'accuracy').evaluate(rf_predictions)
31   print('MulticlassClassificationEvaluator 随机森林测试的准确
     性：{0:.0%}'.format(rf_accuracy))
32   # In[6]:
34   rf_auc=BinaryClassificationEvaluator(labelCol='label').evaluate(rf_predict
     ions)
35   print('BinaryClassificationEvaluator 随机森林测试的准确
     性：{0:.0%}'.format(rf_auc))
36
```

## 三、Ubuntu命令行关键代码

### 1. Mapreduce指令

```
1    cd ~/hadoop_installs/hadoop-3.3.0
2    sbin/start-all.sh
3    jps
5    cd /mnt/hgfs/ubuntu_share/Experiment4/No1.2
6    bin/hdfs dfs -mkdir /Task12
8    #input1放的是mytest.CSV测试样例
9    bin/hdfs dfs -mkdir /Task12/input1
10   bin/hdfs dfs -put
     /mnt/hgfs/ubuntu_share/Experiment4/No1.2/user_info_skip.csv /Task12/
13   bin/hdfs dfs -put /mnt/hgfs/ubuntu_share/Experiment4/No1.2/mytest11.CSV
     /Task12/input1
15   bin/hadoop jar /mnt/hgfs/ubuntu_share/Experiment4/No1.2/popular_merchant-
     7.0-SNAPSHOT.jar project.code /Task12/input2 /Task12/output7
```

### 2. 安装Spark

```
1    tar zxvf spark-3.0.1-bin-hadoop2.7.tgz  -C /usr/local/
2    cd /usr/local/
```

```
3   mv spark-3.0.1-bin-hadoop2.7 spark
5   cd /mnt/hgfs/ubuntu_share/
6   tar zxvf scala-2.12.11.tgz -C /usr/local/
7   cd /usr/local/
8   mv scala-2.12.11 scala
9   vim ~/.bashrc
10  # ~/.bashrc内容
12  #SCALA
13  export SCALA_HOME=/usr/local/scala
14  export PATH=$PATH:$HOME/bin:$JAVA_HOME/bin:$SCALA_HOME/bin
15  #SPARK
16  export SPARK_HOME=/usr/local/spark
17  export PATH=$PATH:$SPARK_HOME/bin
18  #localhost:8080
20  cd /usr/local/spark/conf
22  cp spark-env.sh.template spark-env.sh
23  vim spark-env.sh
25  #内容
26  export JAVA_HOME=/usr/java/jdk1.8.0_261
27  export HADOOP_HOME=/usr/local/hadoop
28  export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
29  export SCALA_HOME=/usr/local/scala
30  export SPARK_HOME=/usr/local/spark
31  export SPARK_MASTER_IP=127.0.0.1
32  export SPARK_MASTER_PORT=7077
33  export SPARK_MASTER_WEBUI_PORT=8099
34  export SPARK_WORKER_CORES=3
35  export SPARK_WORKER_INSTANCES=1
36  export SPARK_WORKER_MEMORY=5G
37  export SPARK_WORKER_WEBUI_PORT=8081
38  export SPARK_EXECUTOR_CORES=1
39  export SPARK_EXECUTOR_MEMORY=1G
40  export LD_LIBRARY_PATH=${LD_LIBRARY_PATH}:$HADOOP_HOME/lib/native
```

## 四、实验结果及截图

### 1. Task1截图

Task11截图

```
root@zxj13-virtual-machine:/usr/local/spark#hdfs dfs -cat /Task12/output1-1/*
(191499,2494)
(353560,2250)
(1059899,1917)
(713695,1754)
(655904,1674)
(67897,1572)
(221663,1547)
(1039919,1511)
(454937,1387)
(81360,1361)
(514725,1356)
(783997,1351)
(823766,1343)
(107407,1319)
(889095,1272)
(936203,1270)
(770668,1257)
(698879,1235)
(349999,1218)
(671759,1167)
(186456,1162)
(315345,1067)
(729259,1021)
(946001,1015)
(181387,1002)
(926069,1002)
(28895,983)
(89953,975)
(413046,965)
(944554,948)
(617878,927)
(676215,873)
(213297,864)
(15207,859)
(513855,842)
(49881,842)
(48664,831)
(179830,827)
(981145,815)
(441588,814)
(939915,805)
```

```
(764906,787)
(1025501,778)
(3001,765)
(147751,763)
(343432,758)
(141675,757)
(100215,745)
(722301,742)
(678194,732)
(526229,730)
(1100222,730)
(952198,725)
(846908,711)
(1102377,708)
(605027,707)
(201405,704)
(15173,698)
(758374,698)
(335720,695)
(110347,693)
(950987,687)
(834362,687)
(399879,685)
(81901,676)
(487805,668)
(203050,659)
(1075577,659)
(173776,659)
(796566,653)
(276750,650)
(209821,647)
(735931,647)
(779070,645)
(235204,640)
(318890,635)
(986262,634)
(886674,622)
(386646,616)
(717309,616)
(28186,615)
```
```
(376482,610)
(554408,603)
(772645,601)
(992011,598)
(784134,597)
(472166,595)
(825218,590)
(566407,585)
(918348,580)
(982357,580)
(293244,577)
(419724,571)
(256896,559)
(893999,554)
(870470,549)
(1042707,549)
(82431,548)
(1093758,545)
(1112049,543)
```

Task12截图

```
root@zxj13-virtual-machine:/usr/local/spark#hdfs dfs -cat /Task12/output1-2/*
(4044,7278)
(3491,3661)
(1102,3588)
(3828,3434)
(4173,3348)
(3734,3303)
(2385,3214)
(4976,3064)
(798,2997)
(422,2893)
(1892,2792)
(1393,2774)
(4282,2737)
(1535,2720)
(4760,2669)
(4644,2607)
(3760,2600)
(184,2406)
(598,2400)
(3698,2111)
(375,2078)
(4043,2077)
(2537,1949)
(1760,1935)
(2482,1920)
(4659,1843)
(2138,1841)
(606,1818)
(1257,1775)
(4218,1760)
(4538,1639)
(141,1633)
(4918,1596)
(3826,1588)
(2031,1584)
(420,1579)
(2813,1553)
(173,1496)
(962,1495)
(2223,1471)
(4079,1467)
(1816,1462)
(2217,1442)
(2403,1436)
(66,1432)
(742,1423)
(2468,1360)
(1056,1356)
(4845,1344)
(2418,1339)
(2669,1335)
(4048,1333)
(4648,1263)
(1861,1253)
(4818,1220)
(4766,1203)
(2273,1198)
(3022,1194)
(2336,1186)
(361,1176)
(474,1174)
(4798,1167)
(2545,1159)
(1087,1133)
(4847,1126)
(2954,1119)
(1346,1106)
(4871,1098)
(3578,1088)
(2387,1076)
(3971,1074)
(4605,1056)
(2676,1051)
(1480,1049)
(2823,1047)
(4127,1040)
```

```
(1,1030)
(2206,1027)
(4257,1001)
(4129,969)
(4160,968)
(2193,963)
(2664,954)
(1727,951)
(310,948)
(1310,930)
(1487,918)
(1200,914)
(3859,914)
(2928,897)
(3163,885)
(3623,884)
(4287,875)
(1867,868)
(3173,867)
(2318,865)
(2677,865)
(4427,864)
(786,839)
(643,827)
```

## 2. Task2截图

### Task21截图

```
scala> val log =spark.read.format("csv").option("header","true").load("hdfs://local
host:9000/Task12/input2/user_log_format1.csv")
log: org.apache.spark.sql.DataFrame = [user_id: string, item_id: string ... 5 more
fields]

scala> val info =spark.read.format("csv").option("header","true").load("hdfs://loca
lhost:9000/Task12/user_info_format1.csv")
info: org.apache.spark.sql.DataFrame = [user_id: string, age_range: string ... 1 mo
re field]

scala> val log1 = log.select("user_id","time_stamp","action_type").filter("time_sta
mp=1111 and action_type=2")
log1: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [user_id: string, ti
me_stamp: string ... 1 more field]

scala> val info1 = info.select("user_id","gender").filter("gender=1 or gender=0")
info1: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [user_id: string, g
ender: string]

scala> val dfjoin1 = info1.join(log1,"user_id")
dfjoin1: org.apache.spark.sql.DataFrame = [user_id: string, gender: string ... 2 mo
re fields]

scala> val dfcount1 = dfjoin1.groupBy("gender").count()
dfcount1: org.apache.spark.sql.DataFrame = [gender: string, count: bigint]

scala> dfcount1.withColumn("ratio",dfcount1("count")/dfjoin1.count).show
+------+------+------------------+
|gender| count|             ratio|
+------+------+------------------+
|     0|846054|0.7232596926427983|
|     1|323725|0.2767403073572017|
+------+------+------------------+


scala> val rddr1 = dfcount1.withColumn("ratio",dfcount1("count")/dfjoin1.count).sel
ect("gender","ratio").rdd.map(x=>(x(0),x(1))).collect()
rddr1: Array[(Any, Any)] = Array((0,0.7232596926427983), (1,0.2767403073572017))

root@zxj13-virtual-machine:/usr/local/spark#hdfs dfs -cat /Task12/output2-1/*
(0,0.7232596926427983)
(1,0.2767403073572017)
```

Task22截图

```scala
scala> val log =spark.read.format("csv").option("header","true").load("hdfs://local
host:9000/Task12/input2/user_log_format1.csv")
log: org.apache.spark.sql.DataFrame = [user_id: string, item_id: string ... 5 more
fields]

scala> val info =spark.read.format("csv").option("header","true").load("hdfs://loca
lhost:9000/Task12/user_info_format1.csv")
info: org.apache.spark.sql.DataFrame = [user_id: string, age_range: string ... 1 mo
re field]

scala> val info1 = info.select("user_id","age_range").filter("age_range>0 and age_r
range<9")
info1: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [user_id: string, a
ge_range: string]

scala> val log1 = log.select("user_id","time_stamp","action_type").filter("time_sta
mp=1111 and action_type=2")
log1: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [user_id: string, ti
me_stamp: string ... 1 more field]

scala> val dfjoin = info1.join(log1,"user_id")
dfjoin: org.apache.spark.sql.DataFrame = [user_id: string, age_range: string ... 2
more fields]

scala> val dfcount = dfjoin.groupBy("age_range").count()
dfcount: org.apache.spark.sql.DataFrame = [age_range: string, count: bigint]

scala> dfcount.withColumn("ratio",dfcount("count")/dfjoin.count).show
+---------+------+--------------------+
|age_range| count|               ratio|
+---------+------+--------------------+
|        7| 19363|0.019736191647869567|
|        3|327758| 0.33407502464093547|
|        8|  3476|0.003542994482672861|
|        5|133480| 0.13605261897214427|
|        6|103774| 0.10577408211878409|
|        1|    54|5.504076584129301E-5|
|        4|268549|  0.2737248634428407|
|        2|124637| 0.12703918392891178|
+---------+------+--------------------+

scala> val dfratio = dfcount.withColumn("ratio",dfcount("count")/dfjoin.count).sele
ct("age_range","ratio").orderBy("age_range")
dfratio: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [age_range: strin
g, ratio: double]

scala> val rddr = dfratio.rdd.map(x=>(x(0),x(1))).collect()
rddr: Array[(Any, Any)] = Array((1,5.504076584129301E-5), (2,0.12703918392891178),
(3,0.33407502464093547), (4,0.2737248634428407), (5,0.13605261897214427), (6,0.1057
7408211878409), (7,0.019736191647869567), (8,0.003542994482672861))

scala>

scala> sc.parallelize(rddr).saveAsTextFile("hdfs://localhost:9000/Task12/output2-2"
)
```

```
root@zxj13-virtual-machine:/usr/local/spark#hdfs dfs -cat /Task12/output2-2/*
(1,5.504076584129301E-5)
(2,0.12703918392891178)
(3,0.33407502464093547)
(4,0.2737248634428407)
(5,0.13605261897214427)
(6,0.10577408211878409)
(7,0.019736191647869567)
(8,0.003542994482672861)
```

3. Task3截图

Task31截图

```
scala> people1.show
+------+------+
|gender|   num|
+------+------+
|     0|285638|
|     1|121670|
+------+------+


scala> people1.agg("num"-> "sum").show
+--------+
|sum(num)|
+--------+
|  407308|
+--------+

scala> people1.withColumn("ratio",people1("num")/407308).show
+------+------+------------------+
|gender|   num|             ratio|
+------+------+------------------+
|     0|285638|0.7012825674919225|
|     1|121670|0.2987174325080774|
+------+------+------------------+
```

**Task32截图**

```
scala> people2.show
+---------+------+
|age_range|   num|
+---------+------+
|        1|    24|
|        2| 52871|
|        3|111654|
|        4| 79991|
|        5| 40777|
|        6| 35464|
|        7|  6992|
|        8|  1266|
+---------+------+


scala>

scala> people2.agg("num"-> "sum").show //结果为329039
+--------+
|sum(num)|
+--------+
|  329039|
+--------+


scala>

scala> people2.withColumn("ratio",people2("num")/329039).show
+---------+------+--------------------+
|age_range|   num|               ratio|
+---------+------+--------------------+
|        1|    24|7.293968192220375E-5|
|        2| 52871| 0.16068308012120144|
|        3|111654|  0.3393336352225724|
|        4| 79991| 0.24310492069329168|
|        5| 40777|  0.1239275587392376|
|        6| 35464| 0.10778053665370975|
|        7|  6992|0.021249760666668692|
|        8|  1266|0.003847568221396...|
+---------+------+--------------------+
```
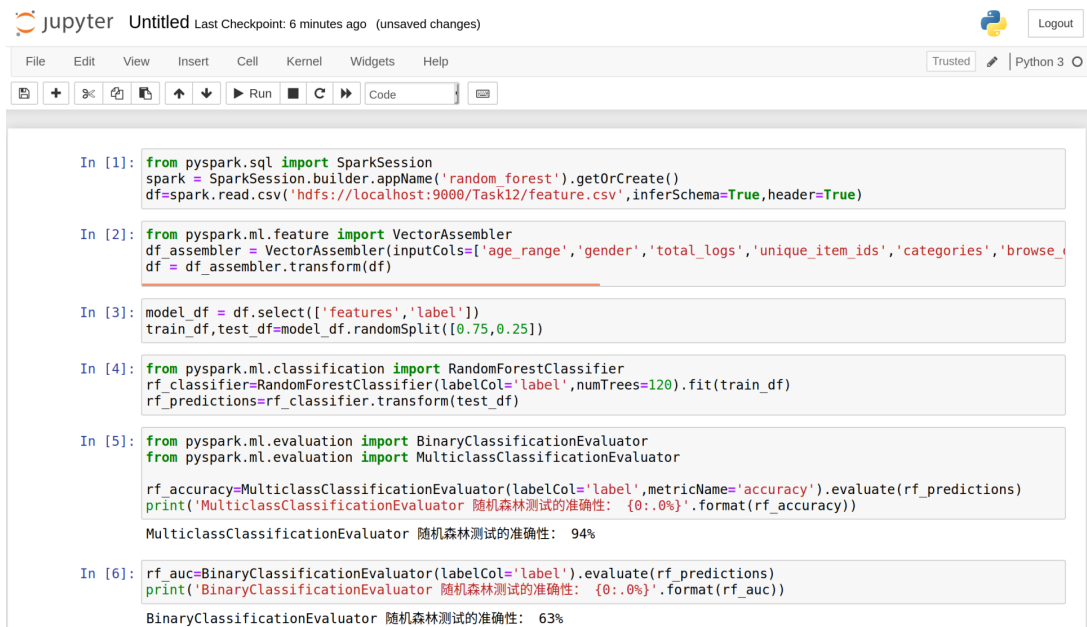
4. Task4截图

File    Edit    View    Insert    Cell    Kernel    Widgets    Help                                Trusted  ✏  │ Python 3  ○

```
In [1]: from pyspark.sql import SparkSession
        spark = SparkSession.builder.appName('random_forest').getOrCreate()
        df=spark.read.csv('hdfs://localhost:9000/Task12/feature.csv',inferSchema=True,header=True)

In [2]: from pyspark.ml.feature import VectorAssembler
        df_assembler = VectorAssembler(inputCols=['age_range','gender','total_logs','unique_item_ids','categories','browse_
        df = df_assembler.transform(df)

In [3]: model_df = df.select(['features','label'])
        train_df,test_df=model_df.randomSplit([0.75,0.25])

In [4]: from pyspark.ml.classification import RandomForestClassifier
        rf_classifier=RandomForestClassifier(labelCol='label',numTrees=120).fit(train_df)
        rf_predictions=rf_classifier.transform(test_df)

In [5]: from pyspark.ml.evaluation import BinaryClassificationEvaluator
        from pyspark.ml.evaluation import MulticlassClassificationEvaluator

        rf_accuracy=MulticlassClassificationEvaluator(labelCol='label',metricName='accuracy').evaluate(rf_predictions)
        print('MulticlassClassificationEvaluator 随机森林测试的准确性： {0:.0%}'.format(rf_accuracy))

        MulticlassClassificationEvaluator 随机森林测试的准确性： 94%

In [6]: rf_auc=BinaryClassificationEvaluator(labelCol='label').evaluate(rf_predictions)
        print('BinaryClassificationEvaluator 随机森林测试的准确性： {0:.0%}'.format(rf_auc))

        BinaryClassificationEvaluator 随机森林测试的准确性： 63%
```

## 五、问题总结

1. pyspark连接jupyter时出错

   解决方案：

   参考https://blog.csdn.net/hecongqing/article/details/85016154
   　　　https://blog.csdn.net/donaldsy/article/details/96194346

   修改配置文件 `vim ~/.jupyter/jupyter_notebook_config.py`
   找到 `#c.NotebookApp.allow_root = False` ，去掉#，并修改为True

## 六、可能的改进之处

1. 第四题预测回头客没来得及用其他方法进行预测对比分析，后续有时间可以尝试别的预测方法