

Pyspark SQL,Dataframe,RDD 使用对比

2020 年 11 月 25 日

首先，启动 Hadoop 和 Spark，打开 Jupyter notebook。

1 创建 RDD

读取 HDFS 上的 buyer_favorite 文件创建 RDD 并查看文件内数据行数

```
[1]: RawUserRDD = sc.textFile("/input/wordcount/buyer_favorite")
```

```
[2]: RawUserRDD.count()
```

```
[2]: 30
```

查看前 5 行数据

```
[3]: RawUserRDD.take(5)
```

```
[3]: ['10181\t1000481\t2010-04-04 16:54:31',  
      '20001\t1001597\t2010-04-07 15:07:52',  
      '20001\t1001560\t2010-04-07 15:08:27',  
      '20042\t1001368\t2010-04-08 08:20:30',  
      '20067\t1002061\t2010-04-08 16:45:33']
```

以 Tab 为分隔符，获取前 5 行的每个字段

```
[4]: userRDD=RawUserRDD.map(lambda line:line.split("\t"))  
userRDD.take(5)
```

```
[4]: [['10181', '1000481', '2010-04-04 16:54:31'],  
      ['20001', '1001597', '2010-04-07 15:07:52'],  
      ['20001', '1001560', '2010-04-07 15:08:27'],
```

```
['20042', '1001368', '2010-04-08 08:20:30'],  
['20067', '1002061', '2010-04-08 16:45:33']]
```

2 创建 DataFrame

通过 userRDD 创建 DataFrame, 导入 Row 模块, 定义 DataFrames 的每一个字段名与数据类型

```
[5]: from pyspark.sql import Row  
user_Rows = userRDD.map(lambda p: Row(  
    ↪Row(buyer_id=int(p[0]),good_id=int(p[1]),dt=p[2]))  
user_Rows.take(5)
```

```
[5]: [Row(buyer_id=10181, dt='2010-04-04 16:54:31', good_id=1000481),  
      Row(buyer_id=20001, dt='2010-04-07 15:07:52', good_id=1001597),  
      Row(buyer_id=20001, dt='2010-04-07 15:08:27', good_id=1001560),  
      Row(buyer_id=20042, dt='2010-04-08 08:20:30', good_id=1001368),  
      Row(buyer_id=20067, dt='2010-04-08 16:45:33', good_id=1002061)]
```

创建了 user_Rows 之后, 使用 sqlContext.createDataFrame() 方法传入 user_Rows 数据, 创建 DataFrame, 然后使用 printSchema() 方法查看 DataFrames 的 Schema

```
[6]: user_df = sqlContext.createDataFrame(user_Rows)  
user_df.printSchema()
```

```
root  
 |-- buyer_id: long (nullable = true)  
 |-- dt: string (nullable = true)  
 |-- good_id: long (nullable = true)
```

```
[7]: user_df.show(5)
```

```
+-----+-----+-----+  
|buyer_id|          dt|good_id|  
+-----+-----+-----+  
|  10181|2010-04-04 16:54:31|1000481|  
|   20001|2010-04-07 15:07:52|1001597|  
|   20001|2010-04-07 15:08:27|1001560|
```

```
|    20042|2010-04-08 08:20:30|1001368|
|    20067|2010-04-08 16:45:33|1002061|
+-----+-----+-----+
only showing top 5 rows
```

也可以使用 `alias()` 方法来为 `DataFrame` 创建别名，例如 `user_df.alias("df")`，后续我们就可以使用 `df` 这个别名执行命令

```
[8]: df = user_df.alias("df")
      df.show(5)
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
|    10181|2010-04-04 16:54:31|1000481|
|    20001|2010-04-07 15:07:52|1001597|
|    20001|2010-04-07 15:08:27|1001560|
|    20042|2010-04-08 08:20:30|1001368|
|    20067|2010-04-08 16:45:33|1002061|
+-----+-----+-----+
only showing top 5 rows
```

3 创建 Spark SQL

使用 `registerDataFrameAsTable` 方法将 `DataFrame` 转成 `table`

```
[9]: sqlContext.registerDataFrameAsTable(df, "buyer_table")
```

使用 `sqlContext.sql()` 输入 `sql` 语句，使用 `select` 关键字查询文件内容行数，并使用 `from` 关键字指定要查询的表，最后使用 `show()` 方法显示查询结果

```
[10]: sqlContext.sql("select count(*) counts from buyer_table").show()
```

```
+-----+
|counts|
+-----+
|     30|
```

```
+-----+
```

为避免输入的 `sql` 语句过长，我们可以使用三个引号 `"""`，来将 `sql` 拆分成多行

```
[11]: sqlContext.sql("""
      select
      count(*) counts
      from
      buyer_table
      """).show()
```

```
+-----+
```

```
| counts |
```

```
+-----+
```

```
|      30 |
```

```
+-----+
```

4 查询部分字段

4.1 使用 RDD 查询部分数据

当我们使用 RDD 查询部分字段时，因为没有 `Schema`，未定义字段名，所以只能指定位置，这里我们查询 `buyer_id`、`good_ids` 和 `dt` 字段

```
[12]: userRDDnew = userRDD.map(lambda x:(x[0],x[1],x[2]))
      userRDDnew.take(5)
```

```
[12]: [('10181', '1000481', '2010-04-04 16:54:31'),
      ('20001', '1001597', '2010-04-07 15:07:52'),
      ('20001', '1001560', '2010-04-07 15:08:27'),
      ('20042', '1001368', '2010-04-08 08:20:30'),
      ('20067', '1002061', '2010-04-08 16:45:33')]
```

4.2 使用 DataFrame 查询部分数据

当使用 `DataFrame` 时，因为已经定义了 `Schema`，所以可以使用 `select` 方法输入字段名，使用 `DataFrame` 查询部分数据时，有 4 种语句，执行结果一样

4.2.1 select 字段名查询

```
[13]: user_df.select("buyer_id", "good_id", "dt").show(5)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
|  10181|1000481|2010-04-04 16:54:31|
|  20001|1001597|2010-04-07 15:07:52|
|  20001|1001560|2010-04-07 15:08:27|
|  20042|1001368|2010-04-08 08:20:30|
|  20067|1002061|2010-04-08 16:45:33|
+-----+-----+-----+
only showing top 5 rows
```

4.2.2 select (dataframe. 字段名) 查询

```
[14]: user_df.select(user_df.buyer_id,user_df.good_id,user_df.dt).show(5)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
|  10181|1000481|2010-04-04 16:54:31|
|  20001|1001597|2010-04-07 15:07:52|
|  20001|1001560|2010-04-07 15:08:27|
|  20042|1001368|2010-04-08 08:20:30|
|  20067|1002061|2010-04-08 16:45:33|
+-----+-----+-----+
only showing top 5 rows
```

4.2.3 select (df 别名. 字段名) 查询

```
[15]: user_df.select(df.buyer_id,df.good_id,df.dt).show(5)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
```

```
| 10181|1000481|2010-04-04 16:54:31|
| 20001|1001597|2010-04-07 15:07:52|
| 20001|1001560|2010-04-07 15:08:27|
| 20042|1001368|2010-04-08 08:20:30|
| 20067|1002061|2010-04-08 16:45:33|
+-----+-----+-----+
only showing top 5 rows
```

4.2.4 使用 [] 查询部分数据

```
[16]: df[df['buyer_id'],df['good_id'],df['dt']].show(5)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
| 10181|1000481|2010-04-04 16:54:31|
| 20001|1001597|2010-04-07 15:07:52|
| 20001|1001560|2010-04-07 15:08:27|
| 20042|1001368|2010-04-08 08:20:30|
| 20067|1002061|2010-04-08 16:45:33|
+-----+-----+-----+
only showing top 5 rows
```

4.3 使用 Spark SQL 查询部分数据

在 sql 语句中，可以使用 select 关键词来查询指定数据

```
[17]: sqlContext.sql("select buyer_id,good_id,dt from buyer_table").show(5)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
| 10181|1000481|2010-04-04 16:54:31|
| 20001|1001597|2010-04-07 15:07:52|
| 20001|1001560|2010-04-07 15:08:27|
| 20042|1001368|2010-04-08 08:20:30|
```

```
| 20067|1002061|2010-04-08 16:45:33|
+-----+-----+-----+
only showing top 5 rows
```

5 增加计算问题

当我们查询数据时，有些字段需要经过计算，现在我们使用 RDD、DataFrame、Spark SQL 三种方式，对 `buyer_id` 字段进行增加 100000 操作

5.1 使用 RDD 增加计算字段

```
[18]: userRDDnew = userRDD.map(lambda x: (x[0],x[1],x[2],int(x[0])+100000))
      userRDDnew.take(5)
```

```
[18]: [('10181', '1000481', '2010-04-04 16:54:31', 110181),
      ('20001', '1001597', '2010-04-07 15:07:52', 120001),
      ('20001', '1001560', '2010-04-07 15:08:27', 120001),
      ('20042', '1001368', '2010-04-08 08:20:30', 120042),
      ('20067', '1002061', '2010-04-08 16:45:33', 120067)]
```

5.2 使用 DataFrame 增加计算字段

```
[19]: df.select("buyer_id","good_id","dt",df.buyer_id+100000).show(5)
```

```
+-----+-----+-----+-----+
|buyer_id|good_id|dt|(buyer_id + 100000)|
+-----+-----+-----+-----+
| 10181|1000481|2010-04-04 16:54:31|110181|
| 20001|1001597|2010-04-07 15:07:52|120001|
| 20001|1001560|2010-04-07 15:08:27|120001|
| 20042|1001368|2010-04-08 08:20:30|120042|
| 20067|1002061|2010-04-08 16:45:33|120067|
+-----+-----+-----+-----+
only showing top 5 rows
```

还可以使用 `alias()` 方法为计算字段取一个别名，这里我们取名为 `new_buyer_id`

```
[20]: df.select("buyer_id", "good_id", "dt", (df.buyer_id+100000).alias("new_buyer_id")).
      →show(5)
```

```
+-----+-----+-----+-----+
|buyer_id|good_id|          dt|new_buyer_id|
+-----+-----+-----+-----+
|   10181|1000481|2010-04-04 16:54:31|    110181|
|   20001|1001597|2010-04-07 15:07:52|    120001|
|   20001|1001560|2010-04-07 15:08:27|    120001|
|   20042|1001368|2010-04-08 08:20:30|    120042|
|   20067|1002061|2010-04-08 16:45:33|    120067|
+-----+-----+-----+-----+
```

only showing top 5 rows

5.3 使用 Spark SQL 增加计算字段

```
[21]: sqlContext.sql("""
      select
      buyer_id,
      good_id,
      dt,
      buyer_id+100000
      from
      buyer_table
      """).show(5)
```

```
+-----+-----+-----+-----+
|buyer_id|good_id|          dt|(buyer_id + CAST(100000 AS BIGINT))|
+-----+-----+-----+-----+
|   10181|1000481|2010-04-04 16:54:31|    110181|
|   20001|1001597|2010-04-07 15:07:52|    120001|
|   20001|1001560|2010-04-07 15:08:27|    120001|
|   20042|1001368|2010-04-08 08:20:30|    120042|
|   20067|1002061|2010-04-08 16:45:33|    120067|
+-----+-----+-----+-----+
```

only showing top 5 rows

6 筛选数据

6.1 使用 RDD 筛选数据

```
[22]: userRDD.filter(lambda r:r[0]=='20056').take(5)
```

```
[22]: [['20056', '1003289', '2010-04-12 10:50:55'],
       ['20056', '1003290', '2010-04-12 11:57:35'],
       ['20056', '1003292', '2010-04-12 12:05:29'],
       ['20056', '1002420', '2010-04-15 11:24:49'],
       ['20056', '1003066', '2010-04-15 11:43:01']]
```

6.2 使用 DataFrame 筛选数据

使用 `filter()` 方法筛选数据

```
[23]: user_df.filter("buyer_id=20056").show()
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
|  20056|2010-04-12 10:50:55|1003289|
|  20056|2010-04-12 11:57:35|1003290|
|  20056|2010-04-12 12:05:29|1003292|
|  20056|2010-04-15 11:24:49|1002420|
|  20056|2010-04-15 11:43:01|1003066|
|  20056|2010-04-15 11:43:06|1003055|
|  20056|2010-04-15 11:45:24|1010183|
|  20056|2010-04-15 11:45:49|1002422|
|  20056|2010-04-15 11:45:54|1003100|
|  20056|2010-04-15 11:45:57|1003094|
|  20056|2010-04-15 11:46:04|1003064|
|  20056|2010-04-15 16:15:20|1010178|
+-----+-----+-----+
```

使用（**dataframe. 字段名**）筛选

```
[24]: df.filter(df.buyer_id=="20056").show()
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
|  20056|2010-04-12 10:50:55|1003289|
|  20056|2010-04-12 11:57:35|1003290|
|  20056|2010-04-12 12:05:29|1003292|
|  20056|2010-04-15 11:24:49|1002420|
|  20056|2010-04-15 11:43:01|1003066|
|  20056|2010-04-15 11:43:06|1003055|
|  20056|2010-04-15 11:45:24|1010183|
|  20056|2010-04-15 11:45:49|1002422|
|  20056|2010-04-15 11:45:54|1003100|
|  20056|2010-04-15 11:45:57|1003094|
|  20056|2010-04-15 11:46:04|1003064|
|  20056|2010-04-15 16:15:20|1010178|
+-----+-----+-----+
```

使用中括号 [] 筛选

```
[25]: df.filter(df["buyer_id"]=="20056").show()
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
|  20056|2010-04-12 10:50:55|1003289|
|  20056|2010-04-12 11:57:35|1003290|
|  20056|2010-04-12 12:05:29|1003292|
|  20056|2010-04-15 11:24:49|1002420|
|  20056|2010-04-15 11:43:01|1003066|
|  20056|2010-04-15 11:43:06|1003055|
|  20056|2010-04-15 11:45:24|1010183|
|  20056|2010-04-15 11:45:49|1002422|
|  20056|2010-04-15 11:45:54|1003100|
|  20056|2010-04-15 11:45:57|1003094|
```

```
| 20056|2010-04-15 11:46:04|1003064|
| 20056|2010-04-15 16:15:20|1010178|
+-----+-----+-----+
```

6.3 使用 Spark SQL 筛选数据

```
[26]: sqlContext.sql("""
select
*
from
buyer_table
where buyer_id = "20056"
""").show(5)
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
| 20056|2010-04-12 10:50:55|1003289|
| 20056|2010-04-12 11:57:35|1003290|
| 20056|2010-04-12 12:05:29|1003292|
| 20056|2010-04-15 11:24:49|1002420|
| 20056|2010-04-15 11:43:01|1003066|
+-----+-----+-----+
only showing top 5 rows
```

7 数据排序

7.1 RDD 按单个字段给数据排序

- 在 RDD 中可以使用 `takeOrdered(num,key=None)` 方法对数据进行排序:
- 参数解释:
 - `num`: 要显示的项数
 - `key`: 使用 `lambda` 语句设置要排序的字段

使用 RDD 按 `buyer_id` 的升序给数据排序

```
[27]: userRDD.takeOrdered(10,key=lambda x:int(x[0]))
```

```
[27]: [['10181', '1000481', '2010-04-04 16:54:31'],
      ['20001', '1001597', '2010-04-07 15:07:52'],
      ['20001', '1001560', '2010-04-07 15:08:27'],
      ['20042', '1001368', '2010-04-08 08:20:30'],
      ['20054', '1002420', '2010-04-14 15:24:12'],
      ['20054', '1010675', '2010-04-14 15:23:53'],
      ['20054', '1002429', '2010-04-14 17:52:45'],
      ['20054', '1003326', '2010-04-20 12:54:44'],
      ['20054', '1003103', '2010-04-15 16:40:14'],
      ['20054', '1003100', '2010-04-15 16:40:16']]
```

使用 RDD 按 buyer_id 的降序给数据排序

```
[28]: userRDD.takeOrdered(10,key=lambda x:-1*int(x[0]))
```

```
[28]: [['20076', '1002427', '2010-04-14 19:35:39'],
      ['20076', '1003101', '2010-04-15 16:37:27'],
      ['20076', '1003103', '2010-04-15 16:37:05'],
      ['20076', '1003100', '2010-04-15 16:37:18'],
      ['20076', '1003066', '2010-04-15 16:37:31'],
      ['20067', '1002061', '2010-04-08 16:45:33'],
      ['20064', '1002422', '2010-04-15 11:35:54'],
      ['20056', '1003289', '2010-04-12 10:50:55'],
      ['20056', '1003290', '2010-04-12 11:57:35'],
      ['20056', '1003292', '2010-04-12 12:05:29']]
```

7.2 DataFrame 按单个字段给数据排序

使用.orderBy() 方法来进行排序，因为默认为升序，所以我们不需要注明 ascending

```
[29]: user_df.select('buyer_id','good_id','dt').orderBy('buyer_id').show(10)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
|  10181|1000481|2010-04-04 16:54:31|
```

```
| 20001|1001560|2010-04-07 15:08:27|
| 20001|1001597|2010-04-07 15:07:52|
| 20042|1001368|2010-04-08 08:20:30|
| 20054|1003103|2010-04-15 16:40:14|
| 20054|1010675|2010-04-14 15:23:53|
| 20054|1003326|2010-04-20 12:54:44|
| 20054|1002420|2010-04-14 15:24:12|
| 20054|1002429|2010-04-14 17:52:45|
| 20054|1003100|2010-04-15 16:40:16|
+-----+-----+-----+
only showing top 10 rows
```

使用 **DataFrame** 按 **buyer_id** 的降序给数据排序，我们需要使用 **.desc()** 方法或者指定 **ascending=0**

```
[30]: user_df.select("buyer_id", "good_id", "dt").orderBy(user_df.buyer_id.desc()).
      ↪ show(10)
```

```
+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
| 20076|1003066|2010-04-15 16:37:31|
| 20076|1002427|2010-04-14 19:35:39|
| 20076|1003100|2010-04-15 16:37:18|
| 20076|1003101|2010-04-15 16:37:27|
| 20076|1003103|2010-04-15 16:37:05|
| 20067|1002061|2010-04-08 16:45:33|
| 20064|1002422|2010-04-15 11:35:54|
| 20056|1003290|2010-04-12 11:57:35|
| 20056|1003066|2010-04-15 11:43:01|
| 20056|1003292|2010-04-12 12:05:29|
+-----+-----+-----+
only showing top 10 rows
```

```
[31]: user_df.select('buyer_id', 'good_id', 'dt').orderBy('buyer_id', ascending=0).
      ↪ show(10)
```

```

+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
|  20076|1003066|2010-04-15 16:37:31|
|  20076|1002427|2010-04-14 19:35:39|
|  20076|1003100|2010-04-15 16:37:18|
|  20076|1003101|2010-04-15 16:37:27|
|  20076|1003103|2010-04-15 16:37:05|
|  20067|1002061|2010-04-08 16:45:33|
|  20064|1002422|2010-04-15 11:35:54|
|  20056|1003290|2010-04-12 11:57:35|
|  20056|1003066|2010-04-15 11:43:01|
|  20056|1003292|2010-04-12 12:05:29|
+-----+-----+-----+
only showing top 10 rows

```

7.3 Spark SQL 按单个字段给数据排序

使用 Spark SQL 按 `buyer_id` 的升序给数据排序

```

[32]: sqlContext.sql("""
select
buyer_id,
good_id,
dt
from
buyer_table
order by buyer_id
""").show(10)

```

```

+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
|  10181|1000481|2010-04-04 16:54:31|
|  20001|1001560|2010-04-07 15:08:27|
|  20001|1001597|2010-04-07 15:07:52|
|  20042|1001368|2010-04-08 08:20:30|

```

```

| 20054|1003103|2010-04-15 16:40:14|
| 20054|1010675|2010-04-14 15:23:53|
| 20054|1003326|2010-04-20 12:54:44|
| 20054|1002420|2010-04-14 15:24:12|
| 20054|1002429|2010-04-14 17:52:45|
| 20054|1003100|2010-04-15 16:40:16|
+-----+-----+-----+
only showing top 10 rows

```

使用 Spark SQL 按 buyer_id 的降序给数据排序

```

[33]: sqlContext.sql("""
select
buyer_id,
good_id,
dt
from
buyer_table
order by buyer_id desc
""").show(10)

```

```

+-----+-----+-----+
|buyer_id|good_id|          dt|
+-----+-----+-----+
| 20076|1003066|2010-04-15 16:37:31|
| 20076|1002427|2010-04-14 19:35:39|
| 20076|1003100|2010-04-15 16:37:18|
| 20076|1003101|2010-04-15 16:37:27|
| 20076|1003103|2010-04-15 16:37:05|
| 20067|1002061|2010-04-08 16:45:33|
| 20064|1002422|2010-04-15 11:35:54|
| 20056|1003290|2010-04-12 11:57:35|
| 20056|1003066|2010-04-15 11:43:01|
| 20056|1003292|2010-04-12 12:05:29|
+-----+-----+-----+
only showing top 10 rows

```

8 按多个字段给数据排序

8.1 RDD 按多个字段给数据排序

```
[34]: userRDD.takeOrdered(10,key=lambda x:(x[0],-int(x[1])))
```

```
[34]: [['10181', '1000481', '2010-04-04 16:54:31'],
      ['20001', '1001597', '2010-04-07 15:07:52'],
      ['20001', '1001560', '2010-04-07 15:08:27'],
      ['20042', '1001368', '2010-04-08 08:20:30'],
      ['20054', '1010675', '2010-04-14 15:23:53'],
      ['20054', '1003326', '2010-04-20 12:54:44'],
      ['20054', '1003103', '2010-04-15 16:40:14'],
      ['20054', '1003100', '2010-04-15 16:40:16'],
      ['20054', '1002429', '2010-04-14 17:52:45'],
      ['20054', '1002420', '2010-04-14 15:24:12']]
```

8.2 DataFrame 按多个字段给数据排序

使用 `orderBy(["buyer_id", "good_id"],ascending=[0,1])` 按多个字段给数据排序

参数解释：- 第一个置要排序的字段：["buyer_id" , "good_id"] - 第二个数：设置排序字段的升序/降序：ascending=[1,0]，其中第一个 buyer_id 设置为 1，表示升序；第二个字段 good_id 设置为 0，表示降序

```
[35]: user_df.orderBy(["buyer_id","good_id"],ascending=[1,0]).show(10)
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
|  10181|2010-04-04 16:54:31|1000481|
|   20001|2010-04-07 15:07:52|1001597|
|   20001|2010-04-07 15:08:27|1001560|
|   20042|2010-04-08 08:20:30|1001368|
|   20054|2010-04-14 15:23:53|1010675|
|   20054|2010-04-20 12:54:44|1003326|
|   20054|2010-04-15 16:40:14|1003103|
|   20054|2010-04-15 16:40:16|1003100|
|   20054|2010-04-14 17:52:45|1002429|
```



```
|    20054|2010-04-14 15:24:12|1002420|
+-----+-----+-----+
only showing top 10 rows
```

使用.desc() 方法表示降序

```
[36]: user_df.orderBy(user_df.buyer_id,user_df.good_id.desc()).show(10)
```

```
+-----+-----+-----+
|buyer_id|          dt|good_id|
+-----+-----+-----+
|   10181|2010-04-04 16:54:31|1000481|
|   20001|2010-04-07 15:07:52|1001597|
|   20001|2010-04-07 15:08:27|1001560|
|   20042|2010-04-08 08:20:30|1001368|
|   20054|2010-04-14 15:23:53|1010675|
|   20054|2010-04-20 12:54:44|1003326|
|   20054|2010-04-15 16:40:14|1003103|
|   20054|2010-04-15 16:40:16|1003100|
|   20054|2010-04-14 17:52:45|1002429|
|   20054|2010-04-14 15:24:12|1002420|
+-----+-----+-----+
only showing top 10 rows
```

8.3 Spark SQL 按多个字段给数据排序

```
[37]: sqlContext.sql("""
select
buyer_id,
good_id,
dt
from buyer_table
order by good_id desc,buyer_id
""").show(10)
```

```
+-----+-----+-----+
```

buyer_id	good_id	dt
20054	1010675	2010-04-14 15:23:53
20056	1010183	2010-04-15 11:45:24
20056	1010178	2010-04-15 16:15:20
20054	1003326	2010-04-20 12:54:44
20056	1003292	2010-04-12 12:05:29
20056	1003290	2010-04-12 11:57:35
20056	1003289	2010-04-12 10:50:55
20054	1003103	2010-04-15 16:40:14
20076	1003103	2010-04-15 16:37:05
20076	1003101	2010-04-15 16:37:27

only showing top 10 rows

9 查询不重复的数据

9.1 RDD 查询不重复数据

使用`.distinct()` 方法查询 `buyer_id` 不重复数据

```
[38]: userRDD.map(lambda x:x[0]).distinct().collect()
```

```
[38]: ['20001',
      '20067',
      '20056',
      '20076',
      '20064',
      '10181',
      '20042',
      '20054',
      '20055']
```

查询 `buyer_id` 和 `good_id` 都不重复的数据，并限制查看 10 条

```
[39]: userRDD.map(lambda x:(x[0],x[1])).distinct().take(10)
```

```
[39]: [('20001', '1001597'),
      ('20056', '1003292'),
      ('20054', '1010675'),
      ('20056', '1002420'),
      ('20056', '1003066'),
      ('20056', '1010183'),
      ('20056', '1003100'),
      ('20056', '1003094'),
      ('20056', '1003064'),
      ('20056', '1010178')]
```

9.2 DataFrame 查询不重复数据

使用`.distinct()` 方法查询 `buyer_id` 不重复数据

```
[40]: user_df.select("buyer_id").distinct().show()
```

```
+-----+
|buyer_id|
+-----+
|  20064|
|  20056|
|  20042|
|  20001|
|   10181|
|  20067|
|  20055|
|  20076|
|  20054|
+-----+
```

查询 `buyer_id` 和 `good_id` 都不重复的数据，并限制查看 10 条

```
[41]: user_df.select("buyer_id", "good_id").distinct().show(10)
```

```
+-----+-----+
|buyer_id|good_id|
```

```

+-----+-----+
|  20054|1003100|
|  20055|1001679|
|  20056|1003055|
|  20076|1003066|
|  20056|1002422|
|  20056|1002420|
|  20056|1003100|
|  20064|1002422|
|  20056|1003094|
|  20054|1003103|
+-----+-----+

```

only showing top 10 rows

9.3 Spark SQL 查询不重复数据

使用 `distinct` 关键字查询 `buyer_id` 不重复数据

```

[42]: sqlContext.sql("""
      select
      distinct buyer_id
      from
      buyer_table
      """).show()

```

```

+-----+
|buyer_id|
+-----+
|  20064|
|  20056|
|  20042|
|  20001|
|  10181|
|  20067|
|  20055|
|  20076|
|  20054|

```

+-----+

使用 `distinct` 关键字查询 `buyer_id` 和 `good_id` 不重复数据，并限制查看 10 条

```
[43]: sqlContext.sql("""
      select
      distinct buyer_id,good_id
      from
      buyer_table
      """).show(10)
```

+-----+-----+

|buyer_id|good_id|

+-----+-----+

| 20054|1003100|

| 20055|1001679|

| 20056|1003055|

| 20076|1003066|

| 20056|1002422|

| 20056|1002420|

| 20056|1003100|

| 20064|1002422|

| 20056|1003094|

| 20054|1003103|

+-----+-----+

only showing top 10 rows

10 分组统计数据

10.1 RDD 分组查询

按照 `buyer_id` 分组统计数据，我们需要用到 `map/reduce`，此方法可用作 `wordcount`

```
[44]: userRDD.map(lambda x:(x[0],1)).reduceByKey(lambda x,y:x+y).collect()
```

```
[44]: [('20001', 2),
      ('20067', 1),
      ('20056', 12),
      ('20076', 5),
      ('20064', 1),
      ('10181', 1),
      ('20042', 1),
      ('20054', 6),
      ('20055', 1)]
```

10.2 DataFrame 分组查询

使用 `groupby()` 方法和 `count()` 方法对 `buyer_id` 进行分组查询

```
[45]: user_df.select("buyer_id").groupby("buyer_id").count().show()
```

```
+-----+-----+
|buyer_id|count|
+-----+-----+
|  20064|    1|
|  20056|   12|
|  20042|    1|
|  20001|    2|
|  10181|    1|
|  20067|    1|
|  20055|    1|
|  20076|    5|
|  20054|    6|
+-----+-----+
```

10.3 Spark SQL 分组查询

```
[46]: sqlContext.sql("""
      select
      buyer_id,count(*) counts
      from buyer_table
```

```
group by buyer_id
""").show()
```

```
+-----+-----+
|buyer_id|counts|
+-----+-----+
|  20064|    1|
|  20056|   12|
|  20042|    1|
|  20001|    2|
|  10181|    1|
|  20067|    1|
|  20055|    1|
|  20076|    5|
|  20054|    6|
+-----+-----+
```

11 Join 连接数据

```
[47]: buyer_log = sc.textFile("file:/data/buyer_log")
      buyer_log.take(5)
```

```
[47]: ['462\t10262\t2010-03-26 19:55:10\t123.127.164.252\t1',
      '463\t20001\t2010-03-29 14:28:02\t221.208.129.117\t2',
      '464\t20001\t2010-03-29 14:28:02\t221.208.129.117\t1',
      '465\t20002\t2010-03-30 10:56:35\t222.44.94.235\t2',
      '466\t20002\t2010-03-30 10:56:35\t222.44.94.235\t1']
```

以 **Tab** 为分隔符，分隔每个字段

```
[48]: userRDD2=buyer_log.map(lambda line:line.split("\t"))
      userRDD2.first()
```

```
[48]: ['462', '10262', '2010-03-26 19:55:10', '123.127.164.252', '1']
```

```
[49]: from pyspark.sql import Row
user_Rows2 = userRDD2.map(lambda p:
    ↳Row(id=int(p[0]),buyer_id=int(p[1]),dt=p[2],ip=p[3],opt_type=p[4]))
user_Rows2.take(5)
```

```
[49]: [Row(buyer_id=10262, dt='2010-03-26 19:55:10', id=462, ip='123.127.164.252',
opt_type='1'),
Row(buyer_id=20001, dt='2010-03-29 14:28:02', id=463, ip='221.208.129.117',
opt_type='2'),
Row(buyer_id=20001, dt='2010-03-29 14:28:02', id=464, ip='221.208.129.117',
opt_type='1'),
Row(buyer_id=20002, dt='2010-03-30 10:56:35', id=465, ip='222.44.94.235',
opt_type='2'),
Row(buyer_id=20002, dt='2010-03-30 10:56:35', id=466, ip='222.44.94.235',
opt_type='1')]
```

创建了 `user_Rows2` 之后，使用 `sqlContext.createDataFrame()` 方法传入 `user_Rows2` 数据，创建 `DataFrame`，然后使用 `printSchema()` 方法查看 `DataFrames` 的 `Schema`

```
[50]: buyerlog_df = sqlContext.createDataFrame(user_Rows2)
buyerlog_df.printSchema()
```

```
root
|-- buyer_id: long (nullable = true)
|-- dt: string (nullable = true)
|-- id: long (nullable = true)
|-- ip: string (nullable = true)
|-- opt_type: string (nullable = true)
```

11.1 DataFrame 联接

`user_df` 通过 `buyer_id` 左外联接 `buyerlog_df`，联接结果将会创建另外一个 `DataFrame` (`joined_df`)

```
[51]: joined_df = user_df.join(buyerlog_df,user_df.buyer_id == buyerlog_df.
    ↳buyer_id,'left_outer')
joined_df.printSchema()
```



```

root
|-- buyer_id: long (nullable = true)
|-- dt: string (nullable = true)
|-- good_id: long (nullable = true)
|-- buyer_id: long (nullable = true)
|-- dt: string (nullable = true)
|-- id: long (nullable = true)
|-- ip: string (nullable = true)
|-- opt_type: string (nullable = true)

```

代码解释:

`joined_df = user_df.join(buyerlog_df : user_df 联接 buyerlog_df 创建 joined_df`

`user_df.buyer_id == buyerlog_df.buyer_id` : 设置联接条件

`'left_outer'` : 设置联接方式

`joined_df.printSchema()` : 打印 `joined_df` 的 Schema

[52]: `joined_df.show(20)`

```

+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+
|buyer_id|          dt|good_id|buyer_id|          dt|  id|
ip|opt_type|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+
|  20064|2010-04-15 11:35:54|1002422|    null|          null|null|
null|    null|
|  20056|2010-04-12 10:50:55|1003289|    null|          null|null|
null|    null|
|  20056|2010-04-12 11:57:35|1003290|    null|          null|null|
null|    null|
|  20056|2010-04-12 12:05:29|1003292|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:24:49|1002420|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:43:01|1003066|    null|          null|null|

```

```

null|    null|
|  20056|2010-04-15 11:43:06|1003055|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:45:24|1010183|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:45:49|1002422|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:45:54|1003100|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:45:57|1003094|    null|          null|null|
null|    null|
|  20056|2010-04-15 11:46:04|1003064|    null|          null|null|
null|    null|
|  20056|2010-04-15 16:15:20|1010178|    null|          null|null|
null|    null|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:14:10|
512|218.9.124.214|          2|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:14:11|
513|218.9.124.214|          1|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:25:38|
514|218.9.124.214|          5|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:26:30|
515|218.9.124.214|          4|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:31:28|
516|218.9.124.214|          1|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:31:43|
517|218.9.124.214|          5|
|  20042|2010-04-08 08:20:30|1001368|  20042|2010-04-08 08:46:09|
518|218.9.124.214|          1|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+
only showing top 20 rows

```

```
[53]: joined_df.filter("ip='123.127.164.252']").show()
```

```

+-----+-----+-----+-----+-----+-----+-----+-----+

```

```

-----+-----+
|buyer_id|          dt|good_id|buyer_id|          dt| id|
ip|opt_type|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-03-31
16:48:43|481|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-01
17:35:05|482|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-02
10:34:20|483|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-06
13:39:37|490|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-07
10:02:08|502|123.127.164.252|      1|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+

```

11.2 Spark SQL Join 操作

之前我们创建了 `buyer_table` 表，现在我们使用 `registerDataFrameAsTable()` 方法来将 `buyerlog_df` 转换为 `buyerlog_table`

```
[54]: sqlContext.registerDataFrameAsTable(buyerlog_df, "buyerlog_table")
```

```
[55]: sqlContext.sql("select count(*) counts from buyerlog_table").show()
```

```

+-----+
|counts|
+-----+
|    62|
+-----+

```

查看 `buyerlog_table` 有多少行记录

```
[56]: sqlContext.sql("select count(*) counts from buyerlog_table").show()
```

```

+-----+
|counts|
+-----+
|    62|
+-----+

```

使用 spark sql 将 buyer_table 和 buyerlog_table 进行左连接，并查询 ip 为 "123.127.164.252" 的数据

```

[57]: sqlContext.sql("""
select
b.*,l.*
from buyer_table b
left join buyerlog_table l on b.buyer_id = l.buyer_id
where l.ip='123.127.164.252'
""").show()

```

```

+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+
|buyer_id|          dt|good_id|buyer_id|          dt| id|
ip|opt_type|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-03-31
16:48:43|481|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-01
17:35:05|482|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-02
10:34:20|483|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-06
13:39:37|490|123.127.164.252|      1|
|  10181|2010-04-04 16:54:31|1000481|  10181|2010-04-07
10:02:08|502|123.127.164.252|      1|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+-----+

```