

第5章Spark SQL 实训

第5章Spark SQL 实训

CH05实训1：统计分析某公司每年的产品销售量及销售额

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

上传测试数据

加载数据

使用临时表

实现查询和统计

CH05实训2：基于DataFrame实现老师教学质量统计分析

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

测试数据文件准备

Hive数据库准备

Hive运行结果

Spark脚本参考

Spark运行结果

CH05实训3：基于DataFrame实现学生成绩统计分析

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

测试数据准备

脚本参考

运行过程结果

CH05实训1：统计分析某公司每年的产品销售量及销售额

训练要点

1. RDD转化为DataFrame的方法
2. DataFrame的查询操作

需求说明

以某公司的交易数据，该数据包含3个.txt文档数据，分别为日期数据tbData.txt，订单表头数据tbStock.txt和订单明细tbStockDetail.txt。

将tbData.txt,tbStock.txt,tbStockDetail.txt上传到HDFS，通过sparkContext读取该数据得到一个RDD，将该RDD数据转换成DataFrame进行查询分析。分析的目标如：

1. 计算所有订单中每年的销售单数，销售总额
2. 计算所有订单每年最大金额订单的销售额
3. 计算所有订单中每年最畅销货品

实现思路及步骤

1. 将数据上传到HDFS
2. 使用SparkContext读取HDFS上的数据
3. 将2中得到的RDD数据转为DataFrame，并注册成为临时表
4. 针对分析目标1，使用join连接3个临时表进行查询
5. 针对分析目标2，首先求出每份订单的销售额以其发生时间，然后以第1步的查询结果和tbData连接，求出每年最大金额订单的销售订单的销售额
6. 针对分析目标3，首先求出每年每个货品的销售利息额，然后求出每年单品销售的最大金额，最后求出每年与销售额最大相符的货品就是最畅销货品

作业要求

1. 提供hdfs上/user/myname/sparkSql/train5data目录下文件列表的截图，应包含文件tbStock.txt, tbStockDetail.txt, tbDate.txt
2. 提供分析目标1,2,3的运行过程及运行结果的截图

实现参考

上传测试数据

```
1 cd /root/spark
2 wget -c http://bigdata.hddly.cn/b46488/file/train5data.tar.gz
3 tar -xzvf ./train5data.tar.gz
4 hdfs dfs -mkdir -p /user/myname/sparkSql
5 hdfs dfs -put -d ./train5data /user/myname/sparkSql/
6 hdfs dfs -ls /user/myname/sparkSql/train5data
```

加载数据

打开spark-shell，在命令行中执行如下：

```
1 case class tbStock(ordernumber:String, locationid:String, dateid:String)
  extends Serializable
2 val tbStockRdd =
  spark.sparkContext.textFile("/user/myname/sparkSql/train5data/tbStock.txt")
3 val tbStockDS = tbStockRdd.map(_.split(",")).map(attr=>
  tbStock(attr(0), attr(1), attr(2))).toDS
4 tbStockDS.show
5
6 case class tbStockDetail(ordernumber:String, rownum:Int, itemid:String,
  number:Int, price:Double, amount:Double) extends Serializable
7 val tbStockDetailRdd =
  spark.sparkContext.textFile("/user/myname/sparkSql/train5data/tbStockDetail.txt")
8 val tbStockDetailDS = tbStockDetailRdd.map(_.split(",")).map(attr=>
  tbStockDetail(attr(0), attr(1).trim().toInt, attr(2), attr(3).trim().toInt, attr(
  4).trim().toDouble, attr(5).trim().toDouble)).toDS
9 tbStockDetailDS.show()
```

```

10
11 case class tbDate(dateid:String, years:Int, theyear:Int, month:Int, day:Int,
    weekday:Int, week:Int, quarter:Int, period:Int, halfmonth:Int) extends
    Serializable
12 val tbDateRdd =
    spark.sparkContext.textFile("/user/myname/sparkSql/train5data/tbDate.txt")
13 val tbDateDS = tbDateRdd.map(_.split(",")).map(attr=>
    tbDate(attr(0),attr(1).trim().toInt,
    attr(2).trim().toInt,attr(3).trim().toInt, attr(4).trim().toInt,
    attr(5).trim().toInt, attr(6).trim().toInt, attr(7).trim().toInt,
    attr(8).trim().toInt, attr(9).trim().toInt)).toDS
14 tbDateDS.show()

```

使用临时表

```

1 tbStockDS.createOrReplaceTempView("tbStock")
2 tbDateDS.createOrReplaceTempView("tbDate")
3 tbStockDetailDS.createOrReplaceTempView("tbStockDetail")

```

实现查询和统计

```

1 //计算所有订单中每年的销售单数、销售总额
2 spark.sql("SELECT c.theyear, COUNT(DISTINCT a.ordernumber), SUM(b.amount)
    FROM tbStock a JOIN tbStockDetail b ON a.ordernumber = b.ordernumber JOIN
    tbDate c ON a.dateid = c.dateid GROUP BY c.theyear ORDER BY c.theyear").show
3 //计算所有订单每年最大金额订单的销售额
4 spark.sql("SELECT a.dateid, a.ordernumber, SUM(b.amount) AS SumOfAmount FROM
    tbStock a JOIN tbStockDetail b ON a.ordernumber = b.ordernumber GROUP BY
    a.dateid, a.ordernumber").show
5 spark.sql("SELECT theyear, MAX(c.SumOfAmount) AS SumOfAmount FROM (SELECT
    a.dateid, a.ordernumber, SUM(b.amount) AS SumOfAmount FROM tbStock a JOIN
    tbStockDetail b ON a.ordernumber = b.ordernumber GROUP BY a.dateid,
    a.ordernumber ) c JOIN tbDate d ON c.dateid = d.dateid GROUP BY theyear ORDER
    BY theyear DESC").show
6 //计算所有订单中每年最畅销货品
7 spark.sql("SELECT c.theyear, b.itemid, SUM(b.amount) AS SumOfAmount FROM
    tbStock a JOIN tbStockDetail b ON a.ordernumber = b.ordernumber JOIN tbDate c
    ON a.dateid = c.dateid GROUP BY c.theyear, b.itemid").show
8
9 spark.sql("SELECT d.theyear, MAX(d.SumOfAmount) AS MaxOfAmount FROM (SELECT
    c.theyear, b.itemid, SUM(b.amount) AS SumOfAmount FROM tbStock a JOIN
    tbStockDetail b ON a.ordernumber = b.ordernumber JOIN tbDate c ON a.dateid =
    c.dateid GROUP BY c.theyear, b.itemid ) d GROUP BY d.theyear").show
10
11 spark.sql("SELECT DISTINCT e.theyear, e.itemid, f.maxofamount FROM (SELECT
    c.theyear, b.itemid, SUM(b.amount) AS sumofamount FROM tbStock a JOIN
    tbStockDetail b ON a.ordernumber = b.ordernumber JOIN tbDate c ON a.dateid =
    c.dateid GROUP BY c.theyear, b.itemid ) e JOIN (SELECT d.theyear,
    MAX(d.sumofamount) AS maxofamount FROM (SELECT c.theyear, b.itemid,
    SUM(b.amount) AS sumofamount FROM tbStock a JOIN tbStockDetail b ON
    a.ordernumber = b.ordernumber JOIN tbDate c ON a.dateid = c.dateid GROUP BY
    c.theyear, b.itemid ) d GROUP BY d.theyear ) f ON e.theyear = f.theyear AND
    e.sumofamount = f.maxofamount ORDER BY e.theyear").show

```

CH05实训2：基于DataFrame实现老师教学质量统计分析

训练要点

1. 熟悉spark SQL的配置方法
2. 掌握在spark SQL中操作Hive表
3. 掌握创建DataFrame的方法。
4. 掌握DataFrame的查询操作
5. 掌握DataFrame的输出操作

需求说明

在某学校的一个微型班级中，有4个表记录着老师与学生的情况，利用大数据技术进行一定处理，方便老师分析自身的教学质量及学生的学习情况，分别是student表、score表、course表和teacher表。student表有4个数据字段，分别为s_id、s_name、s_birth和s_sex, 字段说明如表5—9所示

字段名称	说明
s_id	学生学号
s_name	学生姓名
s_birth	出生日期
s_sex	学生性别

score表有3个数据字段，分别为s_id、c_id和s_score,字段说明如表

字段名称	说明
s_id	学生学号
c_id	课程编号
s_score	分数

course表有3个数据字段，分别表示c_id、c_name和t_id,字段说明如

字段名称	说明
c_id	课程编号
c_name	课程名称
t_id	教师工号

teacher表有2个数据字段，分别为t_id和t_name，字段说明如

字段名称	说明
t_id	教师工号
t_name	教师姓名

为了分析老师的教学质量，请根据score表、course表和teacher表的数据，查询不同老师所教不同课程的平均分。

实现思路及步骤

1. 检查是否已配置spark SQL，若没有，则先配置spark SQL
2. 启动spark-shell、启动Hive元数据服务、访问Hive。
3. 在Hive中创建student数据库，并在student数据库下创建score表、course表和teacher表，并导入数据
4. 使用spark.sql()方法分别读取Hive中的score表、course表和teacher表的数据。
5. 使用groupBy()、withColumn()、count()、sum()、join()、drop()和cast()方法，按教师工号和课程进行分组，聚合查询不同老师所教不同课程的平均分!
6. 将查询结果输出到Hive的student数据库中的teacher_courseAvg表中

作业要求

1. 提供hdfs上/user/myname/目录下文件列表的截图，应包含文件student.txt，score.txt，course.txt，teacher.txt
2. 提供在Hive中创建student数据库，并在student数据库下创建score表、course表和teacher表，并导入数据的过程截图
3. 提供分析运行过程及运行结果的截图

实现参考

测试数据文件准备

在master上操作数据文件

```
1 | cd /root/spark
2 | rm -f spark_t_data.tar.gz*
3 | wget https://biglab.site//b59510spark/file/spark_t_data.tar.gz
4 | tar -xzvf ./spark_t_data.tar.gz
5 | hdfs dfs -rm /user/myname/score.txt
6 | hdfs dfs -rm /user/myname/course.txt
7 | hdfs dfs -rm /user/myname/student.txt
8 | hdfs dfs -rm /user/myname/teacher.txt
9 | hdfs dfs -put ./spark_t_data/score.txt /user/myname/
10 | hdfs dfs -put ./spark_t_data/course.txt /user/myname/
11 | hdfs dfs -put ./spark_t_data/student.txt /user/myname/
12 | hdfs dfs -put ./spark_t_data/teacher.txt /user/myname/
```

Hive数据库准备

在master主机上的hive上执行

```
1 | --创建student数据库
2 | create database student;
```

```

3  --切换student数据库
4  use student;
5  --创建score表
6  create table score(
7  s_id String,
8  c_id String,
9  s_score Int
10 )
11 row format delimited fields terminated by ',';
12 load data local inpath '/root/spark/spark_t_data/score.txt' into table score;
13 --创建teacher表
14 create table teacher(
15 t_id String,
16 t_name String
17 )
18 row format delimited fields terminated by ',';
19 load data local inpath '/root/spark/spark_t_data/teacher.txt' into table
teacher;
20
21 --创建course表
22 create table course(
23 c_id String,
24 c_name String,
25 t_id String
26 )
27 row format delimited fields terminated by ',';
28 load data local inpath '/root/spark/spark_t_data/course.txt' into table
course;
29
30 --创建student表
31 create table student(
32 s_id String,
33 s_name String,
34 s_birth String,
35 s_sex String
36 )
37 row format delimited fields terminated by ',';
38 load data local inpath '/root/spark/spark_t_data/student.txt' into table
student;

```

Hive运行结果

```

1  hive> --创建student数据库
2  hive> create database student;
3  116548 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.q1.session.SessionState - METASTORE_FILTER_HOOK will
be ignored, since hive.security.authorization.manager is set to instance of
HiveAuthorizerFactory.
4  116557 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored

```

```
5 116589 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
6 116608 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore - Failed to get database
hive.student, returning NoSuchObjectException
7 OK
8 Time taken: 0.429 seconds
9 hive> --切换student数据库
10 hive> use student;
11 OK
12 Time taken: 0.048 seconds
13 hive> --创建score表
14 hive> create table score(
15     > s_id String,
16     > c_id String,
17     > s_score Int
18     > )
19     > row format delimited fields terminated by ',';
20 OK
21 Time taken: 1.064 seconds
22 hive> load data local inpath '/root/spark/spark_t_data/score.txt' into table
score;
23 Loading data to table student.score
24 542020 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
25 542761 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
26 OK
27 Time taken: 1.178 seconds
28 hive> --创建teacher表
29 hive> create table teacher(
30     > t_id String,
31     > t_name String
32     > )
33     > row format delimited fields terminated by ',';
34 OK
35 Time taken: 0.129 seconds
36 hive> load data local inpath '/root/spark/spark_t_data/teacher.txt' into
table teacher;
37 Loading data to table student.teacher
38 700252 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
39 700499 [f960459d-ecbf-4bdc-a460-b76bcca31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
40 OK
```

```
41 Time taken: 0.417 seconds
42 hive>
43   > --创建course表
44   > create table course(
45   > c_id String,
46   > c_name String,
47   > t_id String
48   > )
49   > row format delimited fields terminated by ',';
50 OK
51 Time taken: 0.099 seconds
52 hive> load data local inpath '/root/spark/spark_t_data/course.txt' into table
course;
53 Loading data to table student.course
54 700797 [f960459d-ecbf-4bdc-a460-b76bcc31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
55 701052 [f960459d-ecbf-4bdc-a460-b76bcc31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
56 OK
57 Time taken: 0.462 seconds
58 hive>
59   > --创建student表
60   > create table student(
61   > s_id String,
62   > s_name String,
63   > s_birth String,
64   > s_sex String
65   > )
66   > row format delimited fields terminated by ',';
67 OK
68 Time taken: 0.101 seconds
69 hive> load data local inpath '/root/spark/spark_t_data/student.txt' into
table student;
70 Loading data to table student.student
71 702611 [f960459d-ecbf-4bdc-a460-b76bcc31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
72 702844 [f960459d-ecbf-4bdc-a460-b76bcc31823 main] WARN
org.apache.hadoop.hive.metastore.ObjectStore -
datanucleus.autoStartMechanismMode is set to unsupported value null . Setting
it to value: ignored
73 OK
74 Time taken: 0.414 seconds
75 hive>
```

使用exit;退出hive环境。

Spark脚本参考

在spark-shell环境下运行

```
1 val teacher = spark.sql("select * from student.teacher")
2 val score = spark.sql("select * from student.score")
3 val course = spark.sql("select * from student.course")
4 // 查询不同老师所教不同课程平均分
5 val df1 =
  score.groupBy("c_id").sum().withColumnRenamed("sum(s_score)","sums")
6 val df2 = score.groupBy("c_id").count()
7 val df3 = df1.join(df2,"c_id")
8 val df4=df3.withColumn("avg", df3.col("sums")/df3.col("count").cast("int"))
9 val df5 =
  df4.join(course,"c_id").join(teacher,"t_id").drop("sums","count","t_id")
10 df5.show
11 df5.write.mode(org.apache.spark.sql.SaveMode.Overwrite).saveAsTable("teacher_
  courseAvg")
12 val r1t = spark.sql("select * from teacher_courseAvg")
13 r1t.show
```

Spark运行结果

```
1 scala> val teacher = spark.sql("select * from student.teacher")
2 teacher: org.apache.spark.sql.DataFrame = [t_id: string, t_name: string]
3
4 scala> val score = spark.sql("select * from student.score")
5 score: org.apache.spark.sql.DataFrame = [s_id: string, c_id: string ... 1
  more field]
6
7 scala> val course = spark.sql("select * from student.course")
8 course: org.apache.spark.sql.DataFrame = [c_id: string, c_name: string ... 1
  more field]
9
10 scala> // 查询不同老师所教不同课程平均分
11
12 scala> val df1 =
  score.groupBy("c_id").sum().withColumnRenamed("sum(s_score)","sums")
13 df1: org.apache.spark.sql.DataFrame = [c_id: string, sums: bigint]
14
15 scala> val df2 = score.groupBy("c_id").count()
16 df2: org.apache.spark.sql.DataFrame = [c_id: string, count: bigint]
17
18 scala> val df3 = df1.join(df2,"c_id")
19 df3: org.apache.spark.sql.DataFrame = [c_id: string, sums: bigint ... 1 more
  field]
20
21 scala> val df4=df3.withColumn("avg",
  df3.col("sums")/df3.col("count").cast("int"))
22 df4: org.apache.spark.sql.DataFrame = [c_id: string, sums: bigint ... 2 more
  fields]
23
24 scala> val df5 =
  df4.join(course,"c_id").join(teacher,"t_id").drop("sums","count","t_id")
```

```

25 df5: org.apache.spark.sql.DataFrame = [c_id: string, avg: double ... 2 more
fields]
26
27 scala> df5.show
28 +-----+-----+-----+-----+
29 |c_id|          avg|c_name|t_name|
30 +-----+-----+-----+-----+
31 | 01|          64.5| 语文| 李四|
32 | 03|          68.5| 英语| 王五|
33 | 02|72.66666666666667| 数学| 张三|
34 +-----+-----+-----+-----+
35
36
37 scala>
df5.write.mode(org.apache.spark.sql.SaveMode.Overwrite).saveAsTable("teacher_
courseAvg")
38
39 scala> val r1t = spark.sql("select * from teacher_courseAvg")
40 r1t: org.apache.spark.sql.DataFrame = [c_id: string, avg: double ... 2 more
fields]
41
42 scala> r1t.show
43 +-----+-----+-----+-----+
44 |c_id|          avg|c_name|t_name|
45 +-----+-----+-----+-----+
46 | 01|          64.5| 语文| 李四|
47 | 03|          68.5| 英语| 王五|
48 | 02|72.66666666666667| 数学| 张三|
49 +-----+-----+-----+-----+
50
51
52 scala>
53

```

CH05实训3：基于DataFrame实现学生成绩统计分析

训练要点

1. 掌握在spark SQL中操作Hive表。
2. 掌握创建DataFrame对象的方法
3. 掌握DataFrame的查询操作

需求说明

为了进行学生成绩统计分析，需要根据实训1中student表、course表和score表这3个表的数据，完成以下查询任务。

1. 查询所有学生的学生学号、学生姓名、选课总数和所有课程的总成绩；

2. 查询课程名称为“数学”，且分数低于60的学生姓名和分数
3. 查询学生的平均成绩及名次

实现思路及步骤

1. 在Hive的student数据库中创建student表并导入数据，因为score表和course表在实训1中已经创建过，所以无须再创建
2. 使用spark.sql()方法分别读取Hive中的score表、course表和student表的数据
3. 使用groupBy()、sum()、count()、join()和drop()方法，查询所有学生的学生学号、学生姓名、选课总数和所有课程的总成绩。
4. 使用join()和where()方法，查询课程名称为“数学”，且分数低于60的学生姓名和分数
5. 使用groupBy()、join()、sum()、count()、withColumn()、col()、selectExpr()、drop()和cast()方法，查询学生的平均成绩及名次

作业要求

1. 提供hdfs上/user/myname/sparkSql/train5data目录下文件列表的截图，应包含文件student.txt, score.txt, course.txt, teacher.txt
2. 提供在Hive中创建student数据库，并在student数据库下创建score表、course表和teacher表，并导入数据的过程截图
3. 提供分析步骤4, 5的运行过程及运行结果的截图
4. 提供hdfs上user/myname/sparkSql/student_courseAvg结果表的内容的截图

实现参考

测试数据准备

参考一，如果已完成测试数据准备则不用再执行。

脚本参考

启动spark-shell

```
[root@spark2-master-0 sbin]# spark-shell
```

```
1 // 读取数据
2 val student = spark.sql("select * from student.student")
3 val score = spark.sql("select * from student.score")
4 val course = spark.sql("select * from student.course")
5 // 查询所有同学的学生编号、学生姓名、选课总数、所有课程的总成绩
6 val df1 = score.groupBy("s_id").sum("s_score")
7 //val df2 = score.groupBy("s_id").count().join(df1,"s_id","left") //error
8 val df2 = score.groupBy("s_id").count().join(df1,Seq("s_id"),"left_outer")
9 //val df3 = student.drop("s_birth","s_sex").join(df2,"s_id","left") //error
10 val df3 = student.drop("s_birth","s_sex").join(df2,Seq("s_id"),"left_outer")
11 df3.show
12 //查询课程名称为"数学"，且分数低于60的学生姓名和分数
13 val df4 = score.join(course, "c_id").where("c_name='数学' and s_score<60")
14 df4.show
15 // 查询学生平均成绩及其名次
16 val df5 =
    score.groupBy("s_id").sum("s_score").join(score.groupBy("s_id").count(),"s_id")
```

```

17 val df6 = df5.withColumn("avg",
    (df5.col("sum(s_score)"/df5.col("count")).cast("int")).drop("sum(s_score)", "
    count")
18 val df7 = df6.selectExpr("s_id","avg","rank() over(order by avg desc) as
    ord")
19 df7.show
20 df7.write.mode(org.apache.spark.sql.SaveMode.Overwrite).saveAsTable("student_
    courseAvg")
21 val r1t = spark.sql("select * from student_courseAvg")
22 r1t.show

```

运行过程结果

```

1 scala> // 读取数据
2
3 scala> val student = spark.sql("select * from student.student")
4 student: org.apache.spark.sql.DataFrame = [s_id: string, s_name: string ... 2
    more fields]
5
6 scala> val score = spark.sql("select * from student.score")
7 score: org.apache.spark.sql.DataFrame = [s_id: string, c_id: string ... 1
    more field]
8
9 scala> val course = spark.sql("select * from student.course")
10 course: org.apache.spark.sql.DataFrame = [c_id: string, c_name: string ... 1
    more field]
11
12 scala> // 查询所有同学的学生编号、学生姓名、选课总数、所有课程的总成绩
13
14 scala> val df1 = score.groupBy("s_id").sum("s_score")
15 df1: org.apache.spark.sql.DataFrame = [s_id: string, sum(s_score): bigint]
16
17 scala> //val df2 = score.groupBy("s_id").count().join(df1,"s_id","left")
    //error
18
19 scala> val df2 =
    score.groupBy("s_id").count().join(df1,Seq("s_id"),"left_outer")
20 df2: org.apache.spark.sql.DataFrame = [s_id: string, count: bigint ... 1 more
    field]
21
22 scala> //val df3 = student.drop("s_birth","s_sex").join(df2,"s_id","left")
    //error
23
24 scala> val df3 =
    student.drop("s_birth","s_sex").join(df2,Seq("s_id"),"left_outer")
25 df3: org.apache.spark.sql.DataFrame = [s_id: string, s_name: string ... 2
    more fields]
26
27 scala> df3.show
28 +-----+-----+-----+-----+
29 |s_id|s_name|count|sum(s_score)|
30 +-----+-----+-----+-----+
31 | 01| 赵雷| 3| 269|
32 | 02| 钱电| 3| 210|

```

```

33 | 03| 孙风| 3| 240|
34 | 04| 李云| 3| 100|
35 | 05| 周梅| 2| 163|
36 | 06| 吴兰| 2| 65|
37 | 07| 郑竹| 2| 187|
38 | 08| 王菊| null| null|
39 +-----+
40
41
42 scala> //查询课程名称为"数学", 且分数低于60的学生姓名和分数
43
44 scala> val df4 = score.join(course, "c_id").where("c_name='数学' and
45 s_score<60")
46 df4: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [c_id: string,
47 s_id: string ... 3 more fields]
48
49 scala> df4.show
50 +-----+-----+-----+-----+-----+
51 |c_id|s_id|s_score|c_name|t_id|
52 +-----+-----+-----+-----+-----+
53
54
55 scala> // 查询学生平均成绩及其名次
56
57 scala> val df5 =
58 score.groupBy("s_id").sum("s_score").join(score.groupBy("s_id").count(),"s_id")
59 df5: org.apache.spark.sql.DataFrame = [s_id: string, sum(s_score): bigint ...
60 1 more field]
61
62 scala> val df6 = df5.withColumn("avg",
63 (df5.col("sum(s_score)"/df5.col("count")).cast("int")).drop("sum(s_score)","
64 count")
65 df6: org.apache.spark.sql.DataFrame = [s_id: string, avg: int]
66
67 scala> val df7 = df6.selectExpr("s_id","avg","rank() over(order by avg desc)
68 as ord")
69 df7: org.apache.spark.sql.DataFrame = [s_id: string, avg: int ... 1 more
70 field]
71
72 scala> df7.show
73 +-----+-----+-----+
74 |s_id|avg|ord|
75 +-----+-----+-----+
76 | 07| 93| 1|
77 | 01| 89| 2|
78 | 05| 81| 3|
79 | 03| 80| 4|
80 | 02| 70| 5|
81 | 04| 33| 6|
82 | 06| 32| 7|
83 +-----+-----+-----+

```

```
78
79
80 scala> df7.write.mode(org.apache.spark.sql.SaveMode.Overwrite).saveAsTable("student_
courseAvg")
81
82 scala> val r1t = spark.sql("select * from student_courseAvg")
83 r1t: org.apache.spark.sql.DataFrame = [s_id: string, avg: int ... 1 more
field]
84
85 scala> r1t.show
86 +-----+-----+-----+
87 |s_id|avg|ord|
88 +-----+-----+-----+
89 | 07| 93| 1|
90 | 01| 89| 2|
91 | 05| 81| 3|
92 | 03| 80| 4|
93 | 02| 70| 5|
94 | 04| 33| 6|
95 | 06| 32| 7|
96 +-----+-----+-----+
97
98
99 scala>
```