

实验四: Spark SQL 基础编程方法

胡嘉鑫

102102145

2023 年 11 月 22 日

目录

1 实验目的	1
2 实验平台	2
3 实验步骤	2
3.1 Spark SQL 基本操作	2
3.1.1 Problem Description	2
3.1.2 Code	3
3.1.3 Result	5
3.2 咖啡供应链统计	5
3.2.1 Problem Description	5
3.2.2 Code	6
3.2.3 Result	6
3.3 分析校运动会数据	6
3.3.1 Problem Description	6
3.3.2 Code	6
3.3.3 Result	6
4 出现的问题及其解决方案	6

1 实验目的

- 理解 SPARK 工作流程;
- 掌握 SPARK SQL 基础编程方法.

2 实验平台

- OS: Linux
- Hadoop v3.1.3
- Spark v3.4.0
- JDK v1.8

3 实验步骤

3.1 Spark SQL 基本操作

3.1.1 Problem Description

将下列 JSON 格式数据复制到 Linux 系统中，并保存命名为 employee.json。

```
{ "id":1 , "name":" Ella" , "age":36 }  
{ "id":2, "name":"Bob", "age":29 }  
{ "id":3 , "name":"Jack", "age":29 }  
{ "id":4 , "name":"Jim", "age":28 }  
{ "id":4 , "name":"Jim", "age":28 }  
{ "id":5 , "name":"Damon" }  
{ "id":5 , "name":"Damon" }
```

为 employee.json 创建 DataFrame，并写出 Scala 语句完成下列操作：

1. 查询所有数据;
2. 查询所有数据，并去除重复的数据;
3. 查询所有数据，打印时去除 id 字段;
4. 筛选出 age>30 的记录;
5. 将数据按 age 分组;
6. 将数据按 name 升序排列;
7. 取出前 3 行数据;
8. 查询所有记录的 name 列，并为其取别名为 username;
9. 查询年龄 age 的平均值;
10. 查询年龄 age 的最小值。

3.1.2 Code

```
package net.homework

import org.apache.spark.SparkContext._
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.functions.{collect_list, expr}
import org.apache.spark.sql.types.{
    StructField, StructType,
    StringType, LongType, DoubleType
}
import org.apache.spark.sql.Row

import java.nio.file.Paths

object App {
    def main(args : Array[String]) : Unit = {
        val spark = SparkSession
            .builder()
            .appName("pro1")
            .getOrCreate()

        val cwd = Paths.get("").toAbsolutePath.toString
        val inputPath = s"file://${cwd}/input"

        val df = spark.read.format("json").load(s"${inputPath}/employee.json")
        df.cache()

        // 查询所有数据;
        println(" 查询所有数据")
        df.selectExpr("*").show()
        println("=====")

        // 查询所有数据，并去除重复的数据;
        println(" 查询所有数据，并去除重复的数据")
        df.selectExpr("*").distinct().show()
    }
}
```

```

println("=====")

// 查询所有数据，打印时去除 id 字段；
println(" 查询所有数据，打印时去除 id 字段")
df.select("name", "age").distinct().show()
println("=====")

// 筛选出 age>30 的记录；
println(" 筛选出 age>30 的记录")
df.where("age > 30").show()
println("=====")

// 将数据按 age 分组；
println(" 将数据按 age 分组")
df.groupBy("age").agg(collect_list("name")).show()
println("=====")

// 将数据按 name 升序排列；
println(" 将数据按 name 升序排列")
df.orderBy("name").show()
println("=====")

// 取出前 3 行数据；
println(" 取出前 3 行数据")
df.show(3)
println("=====")

// 查询所有记录的 name 列，并为其取别名为 username；
println(" 查询所有记录的 name 列，并为其取别名为 username")
df.selectExpr("name AS username").show()
println("=====")

// 查询年龄 age 的平均值；
println(" 查询年龄 age 的平均值")
df.selectExpr("AVG(age)").show()
println("=====")

```

```

    // 查询年龄 age 的最小值。
    println(" 查询年龄 age 的最小值")
    df.selectExpr("MIN(age)").show()
    println("=====")

    spark.stop()
  }
}

```

3.1.3 Result

```

App.scala
26 object App {
27   def main(args: Array[String]) : Unit = {
28     val spark = SparkSession
29       .builder()
30       .appName("pro1")
31       .getOrCreate()
32
33     val cwd = Paths.get("").toAbsolutePath.toString
34     val inputPath = s"file://${cwd}/input"
35
36     val df = spark.read.format("json").load(s"${inputPath}/employe.json")
37     df.cache()
38
39     // 查询所有数据；
40     println("查询所有数据")
41     df.selectExpr("*").show()
42     println("=====")
43
44     // 查询所有数据，并去除重复的数据；
45     println("查询所有数据，并去除重复的数据")
46     df.selectExpr("*").distinct().show()
47     println("=====")
48
49     // 查询所有数据，打印时去除 id 字段；
50     println("查询所有数据，打印时去除 id 字段")
51     df.select("name", "age").distinct().show()
52     println("=====")
53
54     // 筛选出 age>30 的记录；
55     println("筛选出 age>30 的记录")
56
57   }
58 }

```

```

<45_exp4_src_spark_sql/pro1/src/main/scala/net/homework//502838:/bin/bash
> 0 WuJiaxin@102102145:~/workspace/spark_demo/102102145_exp4_src_spark
1 sql/pro1$ ./run.sh
2 查询所有数据
3 +-----+-----+
4 | age | id | name |
5 +-----+-----+
6 | 36 | 1 | Ella |
7 | 29 | 2 | Bob |
8 | 29 | 3 | Jack |
9 | 28 | 4 | Jim |
10 | 28 | 4 | Jim |
11 | null | 5 | Damon |
12 | null | 5 | Damon |
13 +-----+-----+
14
15
16 查询所有数据，并去除重复的数据
17 +-----+-----+
18 | age | id | name |
19 +-----+-----+
20 | 36 | 1 | Ella |
21 | 29 | 3 | Jack |
22 | null | 5 | Damon |
23 | 29 | 2 | Bob |
24 | 28 | 4 | Jim |
25 +-----+-----+
26
27
28 查询所有数据，打印时去除 id 字段
29 +-----+
30 | name | age |

```

图 1: 运行结果

```

App.scala
3      println("查询所有数据，打印时去除 id 字段")
4      df.select("name", "age").distinct().show()
5      println("=====")
6
7      // 筛选出 age>30 的记录；
8      println("筛选出 age>30 的记录")
9      df.where("age > 30").show()
10     println("=====")
11
12     // 将数据按 age 分组；
13     println("将数据按 age 分组")
14     df.groupBy("age").agg(collect_list("name")).show()
15     println("=====")
16
17     // 将数据按 name 升序排列；
18     println("将数据按 name 升序排列")
19     df.orderBy("name").show()
20     println("=====")
21
22     // 取出前 3 行数据；
23     println("取出前 3 行数据")
24     df.show(3)
25     println("=====")
26
27     // 查询所有记录的 name 列，并为其取别名为 username；
28     println("查询所有记录的 name 列，并为其取别名为 username")
29     df.selectExpr("name AS username").show()
30     println("=====")
31
32     // 查询年龄 age 的平均值；
33     println("查询年龄 age 的平均值")
34
<45_exp4_src_spark_sql/proj/src/main/scala/net/homework//502838:/bin/b
3      =====
4      查询所有数据，打印时去除 id 字段
5      1 +-----+
6      0 | name | age |
7      1 +-----+
8      2 |Damon| null |
9      3 | Bob | 29 |
10     4 | Jim | 28 |
11     5 | Ella | 36 |
12     6 | Jack | 29 |
13     7 +-----+
14     8
15     9
16     10 筛选出 age>30 的记录
17     11 +-----+
18     12 | age | id | name |
19     13 +-----+
20     14 | 36 | 1 | Ella |
21     15 +-----+
22
23     16
24     17 将数据按 age 分组
25     18 +-----+
26     19 | age | collect_list(name) |
27     20 +-----+
28     21 | 29 | [Bob, Jack] |
29     22 | null | [Damon, Damon] |
30     23 | 28 | [Jim, Jim] |
31     24 | 36 | [Ella] |
32     25 +-----+
33     26
34     27
App.scala
L:41/79 0 51% unix utf-8[scala]
main/scala/net/homework//502838:/bin/bash[-] L:31/116 11 26% unix utf-

```

图 2: 运行结果

```

App.scala
9      df.show(3)
10     println("=====")
11
12     // 查询所有记录的 name 列，并为其取别名为 username；
13     println("查询所有记录的 name 列，并为其取别名为 username")
14     df.selectExpr("name AS username").show()
15     println("=====")
16
17     // 查询年龄 age 的平均值；
18     println("查询年龄 age 的平均值")
19     df.selectExpr("AVG(age)").show()
20     println("=====")
21
22     // 查询年龄 age 的最小值；
23     println("查询年龄 age 的最小值")
24     df.selectExpr("MIN(age)").show()
25     println("=====")
26
27     spark.stop()
28 }
29 }
App.scala
L:68/79 42 86% unix utf-8[scala]
<45_exp4_src_spark_sql/proj/src/main/scala/net/homework//502838:/bin/bash
27 将数据按 name 升序排列
26 +-----+
25 | age | id | name |
24 +-----+
23 | 36 | 1 | Ella |
22 | 29 | 2 | Bob |
21 | null | 5 | Damon |
20 | null | 5 | Damon |
19 | 29 | 3 | Jack |
18 | 28 | 4 | Jim |
17 | 28 | 4 | Jim |
16 +-----+
15
14
13 取出前 3 行数据
12 +-----+
11 | age | id | name |
10 +-----+
9 | 36 | 1 | Ella |
8 | 29 | 2 | Bob |
7 | 29 | 3 | Jack |
6 +-----+
5 only showing top 3 rows
4
3
2 查询所有记录的 name 列，并为其取别名为 username
1 +-----+
0 | username |
1 +-----+
2 | Ella |
3 | Bob |
App.scala
L:68/79 42 86% unix utf-8[scala]
main/scala/net/homework//502838:/bin/bash[-] L:87/116 10 75% unix utf-8

```

图 3: 运行结果

```
App.scala
9 df.show(3)
10 println("=====")
11
12 // 查询所有记录的 name 列，并为其取别名为 username;
13 println("查询所有记录的 name 列，并为其取别名为 username")
14 df.selectExpr("name AS username").show()
15 println("=====")
16
17 // 查询年龄 age 的平均值;
18 println("查询年龄 age 的平均值")
19 df.selectExpr("AVG(age)").show()
20 println("=====")
21
22 // 查询年龄 age 的最小值。
23 println("查询年龄 age 的最小值")
24 df.selectExpr("MIN(age)").show()
25 println("=====")
26
27 spark.stop()
28 }
29 }
```

```
<45_exp4_src_spark_sql/pro1/src/main/scala/net/homework//502838:/
21 =====
22 查询所有记录的 name 列，并为其取别名为 username
23 +-----+
24 |username|
25 +-----+
26 | Ella |
27 |  Bob |
28 |  Jack |
29 |   Jim |
30 |   Jim |
31 | Damon |
32 | Damon |
33 +-----+
34
35 =====
36 查询年龄 age 的平均值
37 +-----+
38 |avg(age)|
39 +-----+
40 |   30.0 |
41 +-----+
42
43 =====
44 查询年龄 age 的最小值
45 +-----+
46 |min(age)|
47 +-----+
48 |    28 |
49 +-----+
50
51 =====
```

图 4: 运行结果

3.2 咖啡供应链统计

3.2.1 Problem Description

Coffee Chain.csv 数据字段如下:

Area Code 区号

Ddate 统计日期

Market 市场位置

Market Size 市场规模

Product 产品

Product Type 产品类别

State 所在州

Type 产品属性

Budget Cogs 预算成本

Budget Margin 预算盈余

Budget Profit 利润预算

Budget Sales 销售预算

Coffee Sales 实际销售

Cogs 实际成本

Inventory 库存

Margin 实际盈余

Marketing 销售量

Number of Records 记录数

Profit 实际利润

Total Expenses 其他成本

1. 查看咖啡连锁店的销售量排名，按照销售量降序排列。
2. 查看咖啡销售量和所在州的关系，按降序排列。
3. 查询咖啡的平均利润和售价，按平均利润降序排列。
4. 查询市场规模、市场地域与销售量的关系。按总销量降序排列。
5. 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系。

3.2.2 Code

```
package net.homework

import org.apache.spark.SparkContext._
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.Row
import org.apache.spark.sql.functions.{desc, expr}

import java.nio.file.Paths

object App {
  def main(args : Array[String]) : Unit = {
    val spark = SparkSession.builder().appName("pro2").getOrCreate()

    val cwd = Paths.get("").toAbsolutePath.toString
    val inputPath = s"file://${cwd}/input"

    val df = spark.read
      .option("header", "true")
      .option("inferSchema", "true")
      .format("csv")
      .load(s"${inputPath}/Coffee_Chain.csv")
    df.cache()
```



```

//df.printSchema()

// 查看咖啡连锁店的销售量排名，按照销售量降序排列。
println(" 查看咖啡连锁店的销售量排名，按照销售量降序排列")
df.select("product", "marketing")
  .groupBy("product")
  .agg(expr("SUM(marketing) AS number"))
  .orderBy(desc("number"))
  .show()
println("=====")

// 查看咖啡销售量和所在州的关系，按降序排列。
println(" 查看咖啡销售量和所在州的关系，按降序排列")
df.select("state", "marketing")
  .groupBy("state")
  .agg(expr("SUM(marketing) AS number"))
  .orderBy(desc("number"))
  .show()
println("=====")

// 查询咖啡的平均利润和售价，按平均利润降序排列。
println(" 查询咖啡的平均利润和售价，按平均利润降序排列")
df.select("product", "coffee sales", "profit")
  .groupBy("product")
  .agg(expr("avg(`coffee sales`)"), expr("avg(profit) as avg_profit"))
  .orderBy(desc("avg_profit"))
  .show()
println("=====")

// 查询市场规模、市场地域与销售量的关系。按总销量降序排列。
println(" 查询市场规模、市场地域与销售量的关系。按总销量降序排列")
df.select("market", "market size", "coffee sales")
  .groupBy("market", "market size")
  .agg(expr("SUM(`coffee sales`) AS total_sales"))
  .orderBy(desc("total_sales"))

```

```

        .show()
println("=====")

// 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系。
println(" 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系")
df.select("type", "profit", "coffee sales", "total expenses",
    "marketing", "budget sales")
    .groupBy("type")
    .agg(
        expr("AVG(`coffee sales`) AS avg_sales"),
        expr("AVG(profit) AS avg_profit"),
        expr("SUM(`coffee sales`) AS total_sales"),
        expr("AVG(`total expenses`) AS other_cost"))
    .show()
println("=====")

spark.stop()
}
}

```

3.2.3 Result

```

27 // 查看咖啡连锁店销售量排名，按照销售量降序排列。
26 println("查看咖啡连锁店销售量排名，按照销售量降序排列")
25 df.select("product", "marketing")
24   .groupBy("product")
23   .agg(expr("SUM(marketing) AS number"))
22   .orderBy(desc("number"))
21   .show()
20 println("=====")
19
18 // 查看咖啡销售量和所在州的关系，按降序排列。
17 println("查看咖啡销售量和所在州的关系，按降序排列")
16 df.select("state", "marketing")
15   .groupBy("state")
14   .agg(expr("SUM(marketing) AS number"))
13   .orderBy(desc("number"))
12   .show()
11 println("=====")
10
9 // 查询咖啡的平均利润和售价，按平均利润降序排列。
8 println("查询咖啡的平均利润和售价，按平均利润降序排列")
7 df.select("product", "coffee sales", "profit")
6   .groupBy("product")
5   .agg(expr("avg(`coffee sales`)"), expr("avg(profit) as avg_profit"))
4   .orderBy(desc("avg_profit"))
3   .show()
2 println("=====")
1
0 // 查询市场规模，市场地域与销售量的关系，按总销量降序排列。
1 println("查询市场规模，市场地域与销售量的关系，按总销量降序排列")
2 df.select("market", "market size", "coffee sales")
3   .groupBy("market", "market size")

```

```

0 查看咖啡连锁店销售量排名，按照销售量降序排列
1 +-----+
2 | product|number|
3 +-----+
4 | Caffe Mocha| 19686|
5 | Colombian| 17346|
6 | Lemon| 15880|
7 | Chamomile| 12166|
8 | Decaf Irish Cream| 11362|
9 | Decaf Espresso| 10112|
10 | Earl Grey| 9846|
11 | Darjeeling| 9766|
12 | Green Tea| 7126|
13 | Mint| 6108|
14 | Caffe Latte| 5472|
15 | Amaretto| 4658|
16 | Regular Espresso| 2946|
17 +-----+
18
19 =====
20 查看咖啡销售量和所在州的关系，按降序排列
21 +-----+
22 | state|number|
23 +-----+
24 | California| 16062|
25 | New York| 13044|
26 | Nevada| 12056|
27 | Illinois| 8962|
28 | Iowa| 7594|
29 | Washington| 7244|
30 | Colorado| 6810|

```

~/src/main/scala/net/homework/App.scala L:53/77 1 68% unix utf-8|[scala] result.txt L:1/91 1 1% unix utf-8|[text]

图 5: 运行结果

```

18 // 查看咖啡销售量和所在州的关系，按降序排列。
17 println("查看咖啡销售量和所在州的关系，按降序排列")
16 df.select("state", "marketing")
15   .groupBy("state")
14   .agg(expr("SUM(marketing) AS number"))
13   .orderBy(desc("number"))
12   .show()
11 println("=====")
10
9 // 查询咖啡的平均利润和售价，按平均利润降序排列。
8 println("查询咖啡的平均利润和售价，按平均利润降序排列")
7 df.select("product", "coffee sales", "profit")
6   .groupBy("product")
5   .agg(expr("avg(`coffee sales`)"), expr("avg(profit) as avg_profit"))
4   .orderBy(desc("avg_profit"))
3   .show()
2 println("=====")
1
0 // 查询市场规模、市场地域与销售量的关系，按总销量降序排列。
1 println("查询市场规模、市场地域与销售量的关系，按总销量降序")
2 df.select("market", "market size", "coffee sales")
3   .groupBy("market", "market size")
4   .agg(expr("SUM(`coffee sales`) AS total_sales"))
5   .orderBy(desc("total_sales"))
6   .show()
7 println("=====")
8
9 // 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系。
10 println("查询咖啡属性与平均售价、平均利润、销售量与其他成本")
11 df.select("type", "profit", "coffee sales", "total expenses")
12   .groupBy("type")
13   .agg(
14     expr("AVG(`coffee sales`) AS avg_sales"),
15     expr("AVG(profit) AS avg_profit"),
16     expr("SUM(`coffee sales`) AS total_sales"),
17     expr("AVG(`total expenses`) AS other_cost")
18   )
19   .show()
20 println("=====")
21

```

```

9 =====
8 查看咖啡销售量和所在州的关系，按降序排列
7
6 +-----+-----+
5 | state|number|
4 +-----+-----+
3 | California| 16062|
2 | New York| 13044|
1 | Nevada| 12056|
0 | Illinois| 8962|
1 | Iowa| 7594|
2 | Washington| 7244|
3 | Colorado| 6810|
4 | Wisconsin| 6552|
5 | Florida| 6284|
6 | Oregon| 6026|
7 | Utah| 5670|
8 | Ohio| 5166|
9 | Oklahoma| 4984|
10 | Texas| 4746|
11 | Missouri| 4328|
12 | Connecticut| 4180|
13 | Louisiana| 3790|
14 | Massachusetts| 3428|
15 | New Mexico| 3042|
16 | New Hampshire| 2506|
17 +-----+-----+
18
19 =====
20 查询咖啡的平均利润和售价，按平均利润降序排列
21 +-----+-----+

```

/src/main/scala/net/homework/App.scala L:53/77 1 68% unix utf-8|[scala] result.txt L:28/91 1 30% unix utf-8|[text]

图 6: 运行结果

```

9 // 查询咖啡的平均利润和售价，按平均利润降序排列。
8 println("查看咖啡的平均利润和售价，按平均利润降序排列")
7 df.select("product", "coffee sales", "profit")
6   .groupBy("product")
5   .agg(expr("avg(`coffee sales`)"), expr("avg(profit) as avg_profit"))
4   .orderBy(desc("avg_profit"))
3   .show()
2 println("=====")
1
0 // 查询市场规模、市场地域与销售量的关系，按总销量降序排列。
1 println("查询市场规模、市场地域与销售量的关系，按总销量降序")
2 df.select("market", "market size", "coffee sales")
3   .groupBy("market", "market size")
4   .agg(expr("SUM(`coffee sales`) AS total_sales"))
5   .orderBy(desc("total_sales"))
6   .show()
7 println("=====")
8
9 // 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系。
10 println("查询咖啡属性与平均售价、平均利润、销售量与其他成本")
11 df.select("type", "profit", "coffee sales", "total expense")
12   .groupBy("type")
13   .agg(
14     expr("AVG(`coffee sales`) AS avg_sales"),
15     expr("AVG(profit) AS avg_profit"),
16     expr("SUM(`coffee sales`) AS total_sales"),
17     expr("AVG(`total expenses`) AS other_cost")
18   )
19   .show()
20 println("=====")
21

```

```

3 =====
2 查询咖啡的平均利润和售价，按平均利润降序排列
1
0 +-----+-----+
1 | product| avg(coffee sales)| avg_profit|
2 +-----+-----+
3 | Regular Espresso| 333.7638888888889| 139.79166666666666|
4 | Colombian| 267.31458333333336| 116.25833333333334|
5 | Earl Grey| 231.84722222222223| 83.90277777777777|
6 | Darjeeling| 190.49739583333334| 75.65885416666667|
7 | Decaf Espresso| 191.5735294117647| 72.30882352941177|
8 | Chamomile| 196.81770833333334| 70.0140625|
9 | Lemon| 199.84583333333333| 62.22708333333333|
10 | Caffe Latte| 166.19907407407408| 52.66203703703704|
11 | Caffe Mocha| 176.88333333333333| 36.829166666666666|
12 | Decaf Irish Cream| 162.10416666666666| 36.4296875|
13 | Mint| 185.98958333333334| 32.052083333333336|
14 | Amaretto| 136.81770833333334| 25.46875|
15 | Green Tea| 114.0625|-0.8020833333333334|
16 +-----+-----+
17
18 =====
19 查询市场规模、市场地域与销售量的关系，按总销量降序排列
20 +-----+-----+
21 | market| market size| total_sales|
22 +-----+-----+
23 | West| Small Market| 175372|
24 | Central| Major Market| 152579|
25 | East| Major Market| 138260|
26 | Central| Small Market| 112466|
27 | West| Major Market| 96892|

```

/src/main/scala/net/homework/App.scala L:53/77 1 68% unix utf-8|[scala] result.txt L:49/91 1 53% unix utf-8|[text]

图 7: 运行结果

```
9 // 查询咖啡的平均利润和售价，按平均利润降序排列。
8 println("查询咖啡的平均利润和售价，按平均利润降序排列")
7 df.select("product", "coffee_sales", "profit")
6 .groupBy("product")
5 .agg(expr("avg('coffee_sales')"), expr("avg(profit) as avg_profit"))
4 .orderBy(desc("avg_profit"))
3 .show()
2 println("=====")
1
0 // 查询市场规模、市场地域与销售量的关系，按总销量降序排列。
1 println("查询市场规模、市场地域与销售量的关系，按总销量降序排列")
2 df.select("market", "market_size", "coffee_sales")
3 .groupBy("market", "market_size")
4 .agg(expr("SUM('coffee_sales') AS total_sales"))
5 .orderBy(desc("total_sales"))
6 .show()
7 println("=====")
8
9 // 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系
10 println("查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系")
11 df.select("type", "profit", "coffee_sales", "total_expense")
12 .groupBy("type")
13 .agg(
14   expr("AVG('coffee_sales') AS avg_sales"),
15   expr("AVG(profit) AS avg_profit"),
16   expr("SUM('coffee_sales') AS total_sales"),
17   expr("AVG('total_expenses') AS other_cost")
18 )
19 .show()
20 println("=====")
21

3 =====
2 查询市场规模、市场地域与销售量的关系，按总销量降序排列
1 +-----+-----+-----+
0 || market| market size|total_sales|
1 +-----+-----+-----+
2 | West|Small Market| 175372|
3 | Central|Major Market| 152579|
4 | East|Major Market| 138260|
5 | Central|Small Market| 112466|
6 | West|Major Market| 96892|
7 | South|Small Market| 66516|
8 | East|Small Market| 40316|
9 | South|Major Market| 37410|
10 +-----+-----+-----+
11
12 =====
13 查询咖啡属性与平均售价、平均利润、销售量与其他成本的关系
14 +-----+-----+-----+-----+
15 | type| avg_sales| avg_profit|total_sales|
16 +-----+-----+-----+-----+
17 | Regular|196.74458333333334| 63.66583333333333| 28701840|55.48958|
18 | Decaf| 188.1082251082251|57.762445887445885| 21955584|52.21158|
19 +-----+-----+-----+-----+
20
21 =====
```

图 8: 运行结果

3.3 分析校运动会数据

3.3.1 Problem Description

有一份运动会竞赛结果的数据集，其中包括比赛项目、班级、运动员、成绩和名次，使用 Spark SQL 进行分析。

1. 计算平均成绩:
 - 计算所有比赛项目的平均成绩
2. 统计每个班级的名次总数:
 - 统计每个班级获得的第一名、第二名和第三名的次数。
 - 列出获得第一名次数超过 2 次的班级。
3. 筛选并统计特定项目的成绩:
 - 筛选出所有田径项目（如 100 米短跑、200 米短跑）的比赛结果。
 - 统计在这些田径项目中获得前三名的个人数量。

提示:通过 StructType 创建模式; 利用 spark.read.csv() 读取数据; 使用 createOrReplaceTempView 创建临时表; 使用 filter 进行过滤, 使用 groupby, agg 进行聚合.

3.3.2 Code

```
package net.homework

import org.apache.spark.{SparkConf, SparkContext}
import org.apache.spark.SparkContext._
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.types.{
    StructField, StructType,
    StringType, LongType, DoubleType
}
import org.apache.spark.sql.Row

import java.nio.file.Paths

object App {
    def main(args : Array[String]) : Unit = {
        val spark = SparkSession
            .builder()
            .appName("pro3")
            .getOrCreate()

        val cwd = Paths.get("").toAbsolutePath.toString
        val inputPath = s"file://${cwd}/input"

        val textData = spark
            .sparkContext
            .textFile(s"${inputPath}/运动会数据.txt")
            .filter(
                x => {
                    !x.startsWith(" 比赛") && x.length > 0
                }
            )
        val rddInfo = textData
            .map(
                x => {
                    val arr = x.split(",")

```

```

        Row(arr(0).trim, arr(1).trim, arr(2).trim,
            arr(3).trim.toDouble, arr(4).trim.toLong)
    })

val schema = new StructType(Array(
    new StructField("category", StringType),
    new StructField("class", StringType),
    new StructField("player", StringType),
    new StructField("score", DoubleType),
    new StructField("rank", LongType)
))
//rddInfo.foreach(println)

val dfInfo = spark.createDataFrame(rddInfo, schema)
//dfInfo.show()

dfInfo.createOrReplaceTempView("info")

// 计算所有比赛项目的平均成绩
println(" 计算所有比赛项目的平均成绩:")
spark
    .sql("""SELECT ROUND(AVG(score), 2) AS Avg_Score
        FROM info
        GROUP BY category""")
    .show()
println("=====")

// 统计每个班级获得的第一名、第二名和第三名的次数。
println(" 统计每个班级获得的第一名、第二名和第三名的次数:")
spark
    .sql("""SELECT info.class, COUNT(*) AS First_Num
        FROM info
        WHERE rank=1
        GROUP BY class""")
    .show()
spark

```

```

.sql("""SELECT info.class, COUNT(*) AS Second_Num
      FROM info
      WHERE rank=2
      GROUP BY class""")
.show()
spark
.sql("""SELECT info.class, COUNT(*) AS Third_Num
      FROM info
      WHERE rank=3
      GROUP BY class""")
.show()
println("=====")

// 列出获得第一名次数超过 2 次的班级。
println(" 列出获得第一名次数超过 2 次的班级:")
spark
.sql("""SELECT info.class, COUNT(*) AS Num
      FROM info
      WHERE rank=1
      GROUP BY class having count(*)>2""")
.show()
println("=====")

// 筛选出所有田径项目（如 100 米短跑、200 米短跑）的比赛结果。
println(" 筛选出所有田径项目（如 100 米短跑、200 米短跑）的比赛结果:")
spark
.sql("""SELECT *
      FROM info
      ORDER BY category, rank""")
.show()

// 统计在这些田径项目中获得前三名的个人数量。
println(" 统计在这些田径项目中获得前三名的个人数量:")
spark
.sql("""SELECT DISTINCT COUNT(*)
      FROM info

```

```

WHERE rank<=3"")
.show()

spark.stop()
}
}

```

3.3.3 Result

```

3 val spark = SparkSession
4   .builder()
5   .appName("pro3")
6   .getOrCreate()
7
8 val cwd = Paths.get("").toAbsolutePath.toString
9 val inputPath = s"file://${cwd}/input"
10
11 val textData = spark
12   .sparkContext
13   .textFile(s"${inputPath}/运动会数据.txt")
14   .filter(
15     x => {
16       !x.startsWith("比赛") && x.length > 0
17     }
18   )
19 val rddInfo = textData
20   .map(
21     x => {
22       val arr = x.split(",")
23       Row(arr(0).trim, arr(1).trim, arr(2).trim,
24         arr(3).trim.toDouble, arr(4).trim.toLong)
25     }
26   )
27 val schema = new StructType(Array(
28   new StructField("category", StringType),
29   new StructField("class", StringType),
30   new StructField("player", StringType),
31   new StructField("score", DoubleType),
32   new StructField("rank", LongType)
33 ))
34 //rddInfo.foreach(println)

```

0 计算所有比赛项目的平均成绩：

Avg_Score
1.83
11.92
6.93
12.74
22.06

11 统计每个班级获得的第一名、第二名和第三名的次数：

class	First_Num
A班	3
B班	2

class	Second_Num
A班	3
B班	1
C班	1

class	Third_Num
A班	3
B班	1
C班	1

图 9: 运行结果


```

6 // rddInfo.foreach(println)
7
8 val dfInfo = spark.createDataFrame(rddInfo, schema)
9 // dfInfo.show()
10
11 dfInfo.createOrReplaceTempView("info")
12
13 // 计算所有比赛项目的平均成绩
14 println("计算所有比赛项目的平均成绩:")
15 spark
16   .sql("""SELECT ROUND(AVG(score), 2) AS Avg_Score
17         FROM info
18         GROUP BY category""")
19   .show()
20   println("=====")
21
22 // 统计每个班级获得的第一名、第二名和第三名的次数。
23 println("统计每个班级获得的第一名、第二名和第三名的次数:")
24 spark
25   .sql("""SELECT info.class, COUNT(*) AS First_Num
26         FROM info
27         WHERE rank=1
28         GROUP BY class""")
29   .show()
30   spark
31   .sql("""SELECT info.class, COUNT(*) AS Second_Num
32         FROM info
33         WHERE rank=2
34         GROUP BY class""")
35   .show()
36   spark
37   .sql("""SELECT info.class, COUNT(*) AS Third_Num
38         FROM info
39         WHERE rank=3
40         GROUP BY class""")
41   .show()
42   println("=====")
43
44 // 列出获得第一次数超过 2 次的班级:
45 spark
46   .sql("""SELECT class, COUNT(*) AS First_Count
47         FROM info
48         WHERE rank=1
49         GROUP BY class""")
50   .show()
51   println("=====")
52
53 // 筛选出所有田径项目 (如 100 米短跑、200 米短跑) 的比赛结果:
54 spark
55   .sql("""SELECT category, class, player, score, rank
56         FROM info
57         WHERE category='短跑'""")
58   .show()
59   println("=====")
60
61 // 统计在这些田径项目中获得前三名的个人数量。
62 println("统计在这些田径项目中获得前三名的个人数量:")
63 spark
64   .sql("""SELECT DISTINCT COUNT(*)
65         FROM info
66         WHERE rank <= 3""")
67   .show()
68   spark.stop()
69 }

```

```

1 class|Third_Num|
2 +-----+
3 | B班 | 1 |
4 | C班 | 4 |
5 +-----+
6
7 列出获得第一次数超过 2 次的班级:
8
9 class|Num|
10 +-----+
11 | A班 | 3 |
12 +-----+
13
14 筛选出所有田径项目 (如 100 米短跑、200 米短跑) 的比赛结果:
15
16 category|class|player|score|rank|
17 +-----+-----+-----+-----+-----+
18 |100米短跑| A班 | 张三 |12.45| 1 |
19 |100米短跑| B班 | 李四 |12.62| 2 |
20 |100米短跑| C班 | 王五 |12.75| 3 |
21 |100米短跑| A班 | 李华 |12.82| 4 |
22 |100米短跑| C班 | 李明 |13.05| 5 |
23 |200米短跑| A班 | 刘强 |21.1| 1 |
24 |200米短跑| A班 | 张三 |21.81| 2 |
25 |200米短跑| C班 | 王五 |22.35| 3 |
26 |200米短跑| C班 | 李明 |22.45| 4 |
27 |200米短跑| B班 | 李四 |22.6| 5 |

```

图 10: 运行结果

```

6 println("=====")
7
8 // 列出获得第一次数超过 2 次的班级。
9 println("列出获得第一次数超过 2 次的班级:")
10 spark
11   .sql("""SELECT info.class, COUNT(*) AS Num
12         FROM info
13         WHERE rank=1
14         GROUP BY class having count(*)>2""")
15   .show()
16   println("=====")
17
18 // 筛选出所有田径项目 (如 100 米短跑、200 米短跑) 的比赛结果。
19 println("筛选出所有田径项目 (如 100 米短跑、200 米短跑) 的比赛结果:")
20 spark
21   .sql("""SELECT *
22         FROM info
23         ORDER BY category, rank""")
24   .show()
25
26 // 统计在这些田径项目中获得前三名的个人数量。
27 println("统计在这些田径项目中获得前三名的个人数量:")
28 spark
29   .sql("""SELECT DISTINCT COUNT(*)
30         FROM info
31         WHERE rank <= 3""")
32   .show()
33   spark.stop()
34 }

```

```

1 筛选出所有田径项目 (如 100 米短跑、200 米短跑) 的比赛结果:
2
3 category|class|player|score|rank|
4 +-----+-----+-----+-----+-----+
5 |100米短跑| A班 | 张三 |12.45| 1 |
6 |100米短跑| B班 | 李四 |12.62| 2 |
7 |100米短跑| C班 | 王五 |12.75| 3 |
8 |100米短跑| A班 | 李华 |12.82| 4 |
9 |100米短跑| C班 | 李明 |13.05| 5 |
10 |200米短跑| A班 | 刘强 |21.1| 1 |
11 |200米短跑| A班 | 张三 |21.81| 2 |
12 |200米短跑| C班 | 王五 |22.35| 3 |
13 |200米短跑| C班 | 李明 |22.45| 4 |
14 |200米短跑| B班 | 李四 |22.6| 5 |
15 |跳远| B班 | 李四 |6.95| 1 |
16 |跳远| A班 | 张三 |7.05| 2 |
17 |跳远| B班 | 李华 |6.8| 3 |
18 |跳高| B班 | 李四 |1.9| 1 |
19 |跳高| A班 | 张三 |1.85| 2 |
20 |跳高| C班 | 王五 |1.75| 3 |
21 |铅球| A班 | 张三 |12.34| 1 |
22 |铅球| C班 | 李明 |11.92| 2 |
23 |铅球| C班 | 王五 |11.5| 3 |
24
25 统计在这些田径项目中获得前三名的个人数量:
26
27 count(1)|
28 +-----+
29 | 1 |
30 +-----+

```

图 11: 运行结果

4 出现的问题及其解决方案

没有问题.