实验四: 合并两表实验报告

小组成员: MF20330007 陈明远 MF20330003 陈迪 MG20330095 朱志威

● 实验要求:

本次实验要求将两个文件通过 mapreduce 先进行合并,将结果输入到 HDFS 上,然后通过 Hive 进行建表管理输出结果,最后通过 Hive 语句能查询到所建的表。

● 设计思路:

1. 数据类型

本次实验进行合并的两张表为 order. txt 和 product. txt, 所以本次实验选择定制自己的数据类型, 定义一个自己的 Tablebean 类, 然后实现 Writable 接口 (如果 key 值需要进行自定义排序,则实现 WritableComparable 接口)

2. Map:

本次实验有两种实现方式,Map 端进行合并或者 reduce 端进行合并。

1).在 Map 端进行合并

使用全局文件复制方法,利用 Hadoop 的Distirbuted Cache 机制,在 Mapper 的 setup 中,将需要进行复制的文件装入 Bufferreader 中,然后建立一个HashMap,将产品 id 作为 key,其余数据作为 value,然

后在 map 方法中直接进行合并,不需要通过 reduce 即可完成合并操作,此时 Map 输出的 key 为 Tablebean 类,value 为 Nullwritable 类型。这种方式适合一个表大一个表小的情况。

2). reduce 端进行合并

如果在 reduce 端进行合并,则只需要在 Map 中输出不同表的数据,然后送入 reduce 中,其中 map 输出的 k 为产品 id, value 为我们自定义的数据类型 Tablebean 类。这种合并方式会造成数据倾斜。

3. Reduce:

如果在 reduce 端进行合并,则需要将 map 输出的(k, v)进行整合,因为 reduce 中将同一 key 的不同 value 进行了 Merge,所以我们只需要将其中 order 中的数据存入一个 ArrayList 中,然后遍历这个数组,把所有的 order 数据与 product 中的数据进行合并,输出的 reduce 的 key 为 Tablebean 类,value 为 NullWriatable 类。

● 实现代码:

本次实验我们采用 reduce 端合并,代码如下:

- 1. 用户自制数据类型 Tablebean 类:
- 1. public static class Tablebean implements WritableComparable<Tablebean> {
- 2. //自定义自己的用户类,用来当作 key 或者 value,需要实现 writable 接口,如果需要自定义 map 的排序。则实现 writablecomparable
- 3. private String order_id; //订单id
- 4. private String order_date;//订单日期

```
5.
           private int product_id;//产品id
6.
           private String name;//产品名字
7.
           private int amount;//购买数量
           private int price;//单价
8.
9.
           private int flag;//标记,用来记录是第几个文件的数据
10.
           public Tablebean() {
11.
                super();
12.
           public void set(String order_id, String order_date, int product_id,
13.
   String name, int amount, int price, int flag) {
               this.order id = order id;
14.
               this.order date = order date;
15.
16.
               this.product_id = product_id;
17.
               this.name = name;
18.
               this.amount = amount;
               this.price = price;
19.
20.
               this.flag = flag;
21.
           }
22.
23.
           public String getOrder_id() {
                return this.order_id;
24.
25.
           }
26.
           public void setOrder id(String order id) {
27.
               this.order_id = order_id;
28.
29.
30.
           public String getOrder_date() {
31.
               return this.order_date;
32.
33.
           public void setOrder_date(String order_date) {
34.
               this.order_date = order_date;
35.
           }
36.
37.
           public int getProduct_id() {
               return this.product_id;
38.
39.
           public void setProduct id(int product id) {
40.
               this.product_id = product_id;
41.
42.
43.
44.
           public String getName() {
45.
               return this.name;
46.
           public void setName(String name) {
47.
```

```
48.
               this.name = name;
49.
           }
50.
           public int getAmount() {
               return this.amount;
51.
52.
53.
           public void setAmount(int amount) {
54.
               this.amount = amount;
55.
           }
56.
57.
           public int getPrice() {
58.
                return this.price;
59.
60.
           public void setPrice(int price) {
               this.price = price;
61.
62.
63.
           public int getFlag() {
64.
                return this.flag;
65.
           }
           public void setFlag(int flag) {
66.
67.
               this.flag = flag;
68.
69.
           @Override
70.
           public void write(DataOutput output) throws IOException {
               //任何需要定制的数据类型都至少要实现 write 和 readfields 两个接口
71.
72.
               output.writeUTF(this.order_id);
73.
               output.writeUTF(this.order_date);
74.
               output.writeInt(this.product_id);
75.
               output.writeUTF(this.name);
76.
               output.writeInt(this.amount);
77.
               output.writeInt(this.price);
78.
               output.writeInt(this.flag);
79.
           }
           @Override
80.
81.
           public void readFields(DataInput input) throws IOException {
82.
               this.order_id = input.readUTF();
83.
               this.order_date = input.readUTF();
               this.product_id = input.readInt();
84.
               this.name = input.readUTF();
85.
               this.amount = input.readInt();
86.
87.
               this.price = input.readInt();
88.
               this.flag = input.readInt();
89.
           }
90.
           @Override //toString 方法用来指定输入到 context 的 value 格式
           public String toString() {
91.
```

```
return this.order_id + " " + this.order_date + " " + this.produc
92.
   t_id + " " + this.name + " " + this.price + " " + this.amount;
93.
           }
94.
95.
           public int compareTo(Tablebean o) {
               return this.product_id - o.product_id;
96.
97.
98.
       //定义自己的 map,将产品 id 一样的(k,v)送到同一个 reduce
99.
         public static class DistributedMapper extends Mapper<LongWritable, Text</pre>
100.
   , Text, Tablebean> {
101.
             private Tablebean v = new Tablebean();
102.
             private Text k = new Text();
             @Override
103.
             protected void map(LongWritable key, Text value, Context context)
104.
                     throws IOException, InterruptedException {
105.
106.
                 String line = value.toString();
107.
                 FileSplit fs = (FileSplit) context.getInputSplit();
                 String filename = fs.getPath().getName();//获取文件名
108.
                 if(filename.startsWith("order")) {
109.
                     String[] fileds = line.split(" ");
110.
111.
                     v.setOrder_id(fileds[0]);
112.
                     v.setOrder date(fileds[1]);
                     v.setProduct_id(Integer.valueOf(fileds[2]));
113.
114.
                     v.setAmount(Integer.valueOf(fileds[3]));
115.
                     v.setName("");
116.
                     v.setPrice(∅);
117.
                     v.setFlag(₀);
118.
                     k.set(fileds[2]);//将 key 设置成产品 id
119.
                 }
120.
                 else {
121.
                     String[] fileds = line.split(" ");
122.
                     v.setProduct_id(Integer.valueOf(fileds[0]));
123.
124.
                     v.setName(fileds[1]);
                     v.setPrice(Integer.valueOf(fileds[2]));
125.
                     v.setOrder id("");
126.
127.
                     v.setOrder_date("");
128.
                     v.setFlag(1);
129.
                     v.setAmount(∅);
130.
131.
                     k.set(fileds[0]);
132.
                 }
133.
                 context.write(k, v);
```

```
134. }
135. }
```

2. Map 端:

```
1. public static class DistributedMapper extends Mapper<LongWritable, Text, Tex
    t, Tablebean> {
2.
            private Tablebean v = new Tablebean();
3.
            private Text k = new Text();
            @Override
4.
            protected void map(LongWritable key, Text value, Context context)
5.
                    throws IOException, InterruptedException {
6.
7.
                String line = value.toString();
                FileSplit fs = (FileSplit) context.getInputSplit();
8.
9.
                String filename = fs.getPath().getName();//获取文件名
10.
                if(filename.startsWith("order")) {
                    String[] fileds = line.split(" ");
11.
12.
                    v.setOrder_id(fileds[0]);
13.
                    v.setOrder_date(fileds[1]);
                    v.setProduct_id(Integer.valueOf(fileds[2]));
14.
15.
                    v.setAmount(Integer.valueOf(fileds[3]));
16.
                    v.setName("");
17.
                    v.setPrice(∅);
                    v.setFlag(∅);
18.
19.
                    k.set(fileds[2]);//将 key 设置成产品 id
20.
                }
21.
22.
                else {
                    String[] fileds = line.split(" ");
23.
24.
                    v.setProduct_id(Integer.valueOf(fileds[0]));
25.
                    v.setName(fileds[1]);
                    v.setPrice(Integer.valueOf(fileds[2]));
26.
27.
                    v.setOrder_id("");
28.
                    v.setOrder_date("");
29.
                    v.setFlag(1);
30.
                    v.setAmount(∅);
31.
32.
                    k.set(fileds[0]);
33.
                }
34.
                context.write(k, v);
35.
            }
36.
```

3. Reduce 端:

```
4. public static class MyReducer extends Reducer<Text, Tablebean, Tablebean, Nu
   llWritable> {
5.
           @Override
           public void reduce(Text key, Iterable<Tablebean> values, Context con
6.
   text)
7.
                   throws IOException, InterruptedException {
8.
               Tablebean pt = new Tablebean();
               ArrayList<Tablebean> arrayList = new ArrayList<Tablebean>();
9.
10.
               for(Tablebean value: values) {//一个 values 对应一个 key,所以
11.
   values 中所有数据的产品 id 是一样的
                   if(value.getFlag() == 1) {//如果是 product.txt 中的数据,则只进
12.
   行一次复制,用来后面和 order 中的数据合并
13.
                       try {
14.
                           BeanUtils.copyProperties(pt, value);
15.
                       } catch (Exception e){
16.
                           e.printStackTrace();
17.
                   }else {
18.
                      Tablebean other = new Tablebean();
19.
20.
21.
                          BeanUtils.copyProperties(other, value);
22.
                          arrayList.add(other);//将 order 中的数据存储
                      }catch (Exception e) {
23.
24.
                          e.printStackTrace();
25.
26.
                   }
27.
28.
               for(Tablebean tablebean: arrayList) {//将 order 中的数据与
   product 中的数据合并
29.
                   tablebean.setName(pt.getName());
30.
                   tablebean.setPrice(pt.getPrice());
                   context.write(tablebean, NullWritable.get());
31.
32.
               }
33.
34.
           }
```

4. Main 函数:

```
1. public static void main(String[] args) throws Exception {
```

```
Configuration conf = new Configuration();
3.
        String[] otherargs = (new GenericOptionsParser(conf, args)).getRemaining
    Args();
        if (otherargs.length != 2) {
4.
5.
             System.err.println("Usage:hadoop jar MyJoin2.jar MyJoin2 main input
     output");
6.
             System.exit(2);
7.
        }
8.
        Job job = Job.getInstance(conf, "MyJion2");
9.
10.
        job.setJarByClass(MyJoin2.class);
        job.setMapperClass(DistributedMapper.class);
11.
12.
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(Tablebean.class);
13.
        job.setReducerClass(MyReducer.class);
14.
        job.setOutputKeyClass(Tablebean.class);
15.
16.
        job.setOutputValueClass(NullWritable.class);
17.
        job.setNumReduceTasks(1);
        \label{linear} File Input Format. add Input Path (job, \begin{subarray}{c} new \\ Path (other args[\begin{subarray}{c} 0 \\ \end{bmatrix})); \\
18.
19.
        FileOutputFormat.setOutputPath(job, new Path(otherargs[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
20.
21. }
```

● 实验结果

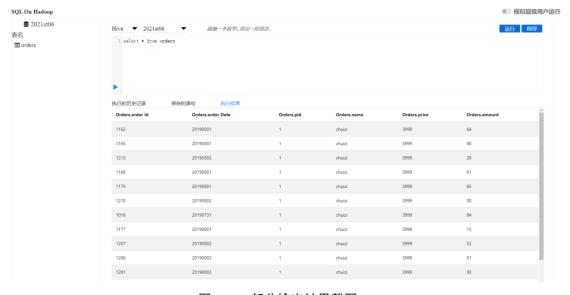


图 1.hive 部分输出结果截图(1)

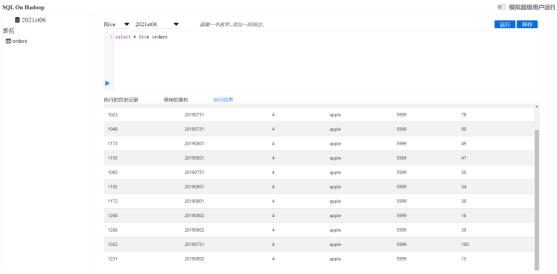


图 2.hive 部分输出结果(2)



图 3.集群监控页面



图 4.Job 执行情况页面

● 实验总结

本次实验通过不同方式合并两张表,对 Hadoop 编程更加熟悉,同时掌握了 hive 查表操作,更加体会到了 hadoop 处理数据的强大之处。