

# 作业二

## 一、实验目的

通过使用 Hadoop MapReduce 和 Spark 框架，实现数据处理的常见任务，包括数据去重、数据排序、计算平均值和单表关联，从而加深对分布式计算框架的理解和应用能力。

## 二、实验环境

- 操作系统：Ubuntu 22.04 LTS
- 虚拟机：UTM
- 集群节点：1 个 master 节点，2 个 slave 节点（slave1, slave2）
- 网络配置：静态 IP 地址
- 软件版本：
  - Hadoop: 2.7.1
  - Spark: 2.2.0
  - Scala: 2.11.8

## 三、hadoop MapReduce任务

### 3.1 数据去重

#### 3.3.1 代码实现

```
1 import org.apache.hadoop.fs.Path
2 import org.apache.hadoop.io.{Text, NullWritable}
3 import org.apache.hadoop.mapreduce.{Job, Mapper, Reducer}
4 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat
5 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat
6 import org.apache.hadoop.conf.Configuration
7
8 object DedupJob {
9
10    // Mapper: 直接输出每一行文本作为 key, value 是 NullWritable
11    class DedupMapper extends Mapper[Object, Text, Text, NullWritable] {
12        private val line = new Text()
13
14        override def map(key: Object, value: Text, context: Mapper[Object, Text, Text,
15          NullWritable]#Context): Unit = {
16            if (value != null && value.toString.trim.nonEmpty) {
17                line.set(value.toString.trim)
18                context.write(line, NullWritable.get())
19            }
20        }
21    }
```

```

22 // Reducer: 对每个唯一的 key (即原始行文本) 只输出一次
23 class DedupReducer extends Reducer[Text, NullWritable, Text, NullWritable] {
24     override def reduce(key: Text, values: java.lang.Iterable[NullWritable],
25                         context: Reducer[Text, NullWritable, Text,
26                         NullWritable]#Context): Unit = {
27         // 只输出一次, 达到去重目的
28         context.write(key, NullWritable.get())
29     }
30 }
31
32 // 主程序入口
33 def main(args: Array[String]): Unit = {
34     if (args.length < 2) {
35         System.err.println("Usage: DedupJob <input1> [<input2> ...] <output>")
36         System.exit(1)
37     }
38
39     val conf = new Configuration()
40     val job = Job.getInstance(conf, "Scala Hadoop Dedup Job")
41
42     job.setJarByClass(DedupJob.getClass)
43     job.setMapperClass(classOf[DedupMapper])
44     job.setReducerClass(classOf[DedupReducer])
45
46     job.setOutputKeyClass(classOf[Text])
47     job.setOutputValueClass(classOf[NullWritable])
48
49     // 添加所有输入路径
50     for (i <- 0 until args.length - 1) {
51         FileInputFormat.addInputPath(job, new Path(args(i)))
52     }
53
54     // 设置输出路径
55     FileOutputFormat.setOutputPath(job, new Path(args.last))
56
57     // 提交作业并退出
58     System.exit(if (job.waitForCompletion(true)) 0 else 1)
59 }

```

### 3.3.2 实验步骤

- 在本地创建并编辑文件file1、file2：

```

1 file1
2 2012-3-1 a
3 2012-3-2 b
4 2012-3-3 c
5 2012-3-4 d
6 2012-3-5 a
7 2012-3-6 b

```

```

8  2012-3-7 c
9  2012-3-3 c
10
11 file2
12 2012-3-1 b
13 2012-3-2 a
14 2012-3-3 b
15 2012-3-4 d
16 2012-3-5 a
17 2012-3-6 c
18 2012-3-7 d
19 2012-3-3 c
20
21 # 冗余数据
22 2012-3-3 c
23 2012-3-4 d
24 2012-3-5 a

```

使用命令 `hdfs dfs -put file1 file2 \user\lfl\input\dedup` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

## Browse Directory

/user/lfl/input/dedup

Go!

| Permission | Owner | Group      | Size | Last Modified      | Replication | Block Size | Name  |
|------------|-------|------------|------|--------------------|-------------|------------|-------|
| -rw-r--r-- | lfl   | supergroup | 88 B | 2025/6/28 17:39:11 | 2           | 128 MB     | file1 |
| -rw-r--r-- | lfl   | supergroup | 88 B | 2025/6/28 17:39:11 | 2           | 128 MB     | file2 |

### 2. 编译打包代码成 JAR 文件

实验过程已经把所有代码实现，因此四个源文件均一起编译并打包为 `MapReduce-assembly-0.1.jar`，提交作业时只需选定相应的类，如下：

```

lfl@master:~/MapReduce$ sbt 'show discoveredMainClasses'
[info] welcome to sbt 1.5.5 (Oracle Corporation Java 1.8.0_151)
[info] loading settings for project mapreduce-build from plugins.sbt ...
[info] loading project definition from /home/lfl/MapReduce/project
[info] loading settings for project mapreduce from build.sbt ...
[info] set current project to MapReduce (in build file:/home/lfl/MapReduce/)
[info] * DedupJob
[info] * SortedIndexJob
[info] * StudentAvgJob
[info] * TableReferJob
[success] Total time: 10 s, completed Jul 4, 2025 4:31:28 PM

```

后续任务直接在集群上运行作业

### 3. 在Hadoop集群上运行作业

在工程根目录下使用命令

```
1 hadoop jar target/scala-2.11/MapReduce-assembly-0.1.jar DedupJob \
2   /user/lfl/input/dedup/file1 \
3   /user/lfl/input/dedup/file2 \
4   /user/lfl/output/mapreddedup_result
```

### 3.1.3 实验结果

成功运行后可在Web页面查看输出文件是否存在:

/user/lfl/output/mapreddedup\_result Go!

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |
|------------|-------|------------|------|-------------------|-------------|------------|------------|
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/4 22:42:41 | 2           | 128 MB     | _SUCCESS   |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/4 22:42:43 | 2           | 128 MB     | part-00000 |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/4 22:42:44 | 2           | 128 MB     | part-00001 |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/4 22:42:44 | 2           | 128 MB     | part-00002 |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/4 22:42:44 | 2           | 128 MB     | part-00003 |

查看并验证输出结果

```
lfl@master:~/MapReduce$ hdfs dfs -cat /user/lfl/output/sparkdedup_result/part*
2012-3-1 a
2012-3-1 b
2012-3-2 a
2012-3-2 b
2012-3-3 b
2012-3-3 c
2012-3-4 d
2012-3-5 a
2012-3-6 b
2012-3-6 c
2012-3-7 c
2012-3-7 d
```

## 3.2 数据排序

### 3.2.1 代码实现

```
1 import org.apache.hadoop.conf.Configuration
2 import org.apache.hadoop.fs.Path
3 import org.apache.hadoop.io.{IntWritable, LongWritable, Text}
4 import org.apache.hadoop.mapreduce.{Job, Mapper, Reducer}
5 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat
6 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat
7
8 import java.lang.{Long => JLong}
9 import scala.collection.JavaConverters._
10
```

```
11 object SortedIndexJob {
12
13     class SortMapper extends Mapper[LongWritable, Text, IntWritable, IntWritable] {
14         private val num = new IntWritable()
15
16         override def map(key: LongWritable, value: Text, context: Mapper[LongWritable,
17             Text, IntWritable, IntWritable]#Context): Unit = {
18             val line = value.toString.trim
19             if (line.nonEmpty) {
20                 try {
21                     num.set(line.toInt)
22                     context.write(num, new IntWritable(1)) // 只使用 key 参与排序
23                 } catch {
24                     case _: NumberFormatException => // 忽略非法行
25                 }
26             }
27         }
28     }
29
30     class SortReducer extends Reducer[IntWritable, IntWritable, Text, IntWritable] {
31         private var index = 1
32         private val outputKey = new Text()
33
34         override def reduce(key: IntWritable, values: java.lang.Iterable[IntWritable],
35                             context: Reducer[IntWritable, IntWritable, Text,
36                             IntWritable]#Context): Unit = {
37             for (_ <- values.iterator().asScala) {
38                 outputKey.set(index.toString)
39                 context.write(outputKey, key)
40                 index += 1
41             }
42         }
43     }
44
45     def main(args: Array[String]): Unit = {
46         if (args.length < 2) {
47             System.err.println("Usage: SortedIndexJob <input1> [<input2> ...] <output>")
48             System.exit(1)
49         }
50
51         val conf = new Configuration()
52         val job = Job.getInstance(conf, "SortedIndexJob")
53
54         job.setJarByClass(SortedIndexJob.getClass)
55         job.setMapperClass(classOf[SortMapper])
56         job.setReducerClass(classOf[SortReducer])
57
58         job.setMapOutputKeyClass(classOf[IntWritable])
59         job.setMapOutputValueClass(classOf[IntWritable])
60         job.setOutputKeyClass(classOf[Text])
61         job.setOutputValueClass(classOf[IntWritable])
```

```

61     for (i <- 0 until args.length - 1) {
62         FileInputFormat.addInputPath(job, new Path(args(i)))
63     }
64     FileOutputFormat.setOutputPath(job, new Path(args.last))
65
66     System.exit(if (job.waitForCompletion(true)) 0 else 1)
67 }
68 }
```

### 3.2.2 实验步骤

#### 1. 准备包含整数数据的输入文件

在本地创建并编辑文件file3、file4、file5：

```

lfl@master:~/data$ cat file3 file4 file5
2
32
654
32
15
756
65223
5956
22
650
92
26
54
6
```

使用命令 `hdfs dfs -put file3 file4 file5 \user\lfl\input\sort` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

| /user/lfl/input/sort |       |            |      |                    |             |            |       | Go! |
|----------------------|-------|------------|------|--------------------|-------------|------------|-------|-----|
| Permission           | Owner | Group      | Size | Last Modified      | Replication | Block Size | Name  |     |
| -rw-r--r--           | lfl   | supergroup | 25 B | 2025/6/28 17:39:12 | 2           | 128 MB     | file3 |     |
| -rw-r--r--           | lfl   | supergroup | 15 B | 2025/6/28 17:39:12 | 2           | 128 MB     | file4 |     |
| -rw-r--r--           | lfl   | supergroup | 8 B  | 2025/6/28 17:39:12 | 2           | 128 MB     | file5 |     |

#### 2. 在Hadoop集群上运行作业

在工程根目录下使用命令

```
1 hadoop jar target/scala-2.11/MapReduce-assembly-0.1.jar DedupJob \
2   /user/lfl/input/sort/file3 \
3   /user/lfl/input/sort/file4 \
4   /user/lfl/input/sort/file5 \
5   /user/lfl/output/mapredsorted_result
```

### 3.2.3 实验结果

成功运行后可在Web页面查看输出文件是否存在

|                                      |     |
|--------------------------------------|-----|
| /user/lfl/output/mapredsorted_result | Go! |
|--------------------------------------|-----|

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |
|------------|-------|------------|------|-------------------|-------------|------------|------------|
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/4 22:43:13 | 2           | 128 MB     | _SUCCESS   |
| -rw-r--r-- | lfl   | supergroup | 23 B | 2025/7/4 22:43:15 | 2           | 128 MB     | part-00000 |
| -rw-r--r-- | lfl   | supergroup | 27 B | 2025/7/4 22:43:15 | 2           | 128 MB     | part-00001 |
| -rw-r--r-- | lfl   | supergroup | 31 B | 2025/7/4 22:43:15 | 2           | 128 MB     | part-00002 |

查看并验证输出结果

```
lfl@master:~/MapReduce$ hdfs dfs -cat /user/lfl/output/sparksorted_result/part*
1      2
2      6
3     15
4     22
5     26
6     32
7     32
8     54
9     92
10    650
11    654
12    756
13    5956
14   65223
```

可观察到输出文件中的数据按升序排列

## 3.3 平均

### 3.3.1 代码实现

```
1 import org.apache.hadoop.conf.Configuration
2 import org.apache.hadoop.fs.Path
3 import org.apache.hadoop.io.{IntWritable, Text}
4 import org.apache.hadoop.mapreduce.{Job, Mapper, Reducer}
5 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat
6 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat
```

```
7
8 import scala.jdk.CollectionConverters._
9
10 object StudentAvgJob {
11
12     // Mapper 输出: key 为学生名, value 为该科成绩
13     class AvgMapper extends Mapper[Object, Text, Text, IntWritable] {
14         private val name = new Text()
15         private val scoreWritable = new IntWritable()
16
17         override def map(key: Object, value: Text, context: Mapper[Object, Text, Text,
18             IntWritable]#Context): Unit = {
19             val line = value.toString.trim
20             if (line.nonEmpty) {
21                 val parts = line.split("\\s+")
22                 if (parts.length == 2) {
23                     try {
24                         name.set(parts(0))
25                         scoreWritable.set(parts(1).toInt)
26                         context.write(name, scoreWritable)
27                     } catch {
28                         case _: NumberFormatException => // 忽略非法数据
29                     }
30                 }
31             }
32         }
33
34     // Reducer 输出: key 为学生名, value 为平均成绩
35     class AvgReducer extends Reducer[Text, IntWritable, Text, IntWritable] {
36         override def reduce(key: Text, values: java.lang.Iterable[IntWritable],
37             context: Reducer[Text, IntWritable, Text,
38             IntWritable]#Context): Unit = {
39             val scores = values.asScala.map(_.get)
40             val sum = scores.sum
41             val count = scores.size
42             val avg = if (count > 0) sum / count else 0
43             context.write(key, new IntWritable(avg))
44         }
45     }
46
47     def main(args: Array[String]): Unit = {
48         if (args.length < 2) {
49             System.err.println("Usage: StudentAvgJob <input1> [<input2> ...] <output>")
50             System.exit(1)
51         }
52
53         val conf = new Configuration()
54         val job = Job.getInstance(conf, "Scala Hadoop StudentAvgJob")
55
56         job.setJarByClass(this.getClass)
57         job.setMapperClass(classOf[AvgMapper])
```

```

57     job.setReducerClass(classOf[AvgReducer])
58
59     job.setMapOutputKeyClass(classOf[Text])
60     job.setMapOutputValueClass(classOf[IntWritable])
61     job.setOutputKeyClass(classOf[Text])
62     job.setOutputValueClass(classOf[IntWritable])
63
64     // 多个输入路径
65     for (i <- 0 until args.length - 1) {
66       FileInputFormat.addInputPath(job, new Path(args(i)))
67     }
68
69     // 设置输出路径
70     FileOutputFormat.setOutputPath(job, new Path(args.last))
71
72     System.exit(if (job.waitForCompletion(true)) 0 else 1)
73   }
74 }
```

### 3.3.2 实验步骤

- 准备包含学生姓名和成绩的数据文件，每行格式为“姓名 成绩”。

在本地创建并编辑文件math、chinese、english:

```

1 # math
2 张三 88
3 李四 99
4 王五 66
5 赵六 77
6
7 #chinese
8 张三 78
9 李四 89
10 王五 96
11 赵六 67
12
13 #english
14 张三 80
15 李四 82
16 王五 84
17 赵六 86
```

使用命令 `hdfs dfs -put math chinese english \user\lfl\input\average` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

/user/lfl/input/average

Go!

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name    |
|------------|-------|------------|------|-------------------|-------------|------------|---------|
| -rw-r--r-- | lfl   | supergroup | 40 B | 2025/7/3 11:23:50 | 2           | 128 MB     | chinese |
| -rw-r--r-- | lfl   | supergroup | 40 B | 2025/7/3 11:23:50 | 2           | 128 MB     | english |
| -rw-r--r-- | lfl   | supergroup | 42 B | 2025/7/3 11:23:50 | 2           | 128 MB     | math    |

## 2. 在Hadoop集群上运行作业

在项目根目录下运行命令

```
1 hadoop jar target/scala-2.11/MapReduce-assembly-0.1.jar SortedIndexJob \
2   /user/lfl/input/sort/math \
3   /user/lfl/input/sort/chinese \
4   /user/lfl/input/sort/english \
5   /user/lfl/output/mapredaverage_result
```

### 3.3.3 实验结果

运行成功后可在Web页面查看输出文件是否存在：

/user/lfl/output/mapredaverage\_result

Go!

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |
|------------|-------|------------|------|-------------------|-------------|------------|------------|
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/4 22:43:31 | 2           | 128 MB     | _SUCCESS   |
| -rw-r--r-- | lfl   | supergroup | 20 B | 2025/7/4 22:43:33 | 2           | 128 MB     | part-00000 |
| -rw-r--r-- | lfl   | supergroup | 10 B | 2025/7/4 22:43:34 | 2           | 128 MB     | part-00001 |
| -rw-r--r-- | lfl   | supergroup | 10 B | 2025/7/4 22:43:34 | 2           | 128 MB     | part-00002 |

查看并验证运行结果：

```
lfl@master:~$ hdfs dfs -cat /user/lfl/output/sparkaverage_result/part*
张三 82
李四 90
王五 82
赵六 76
```

可观察到输出文件中包含每个学生的平均成绩，验证了计算平均功能。

## 3.4 单表关联

### 3.4.1 代码实现

```
1 import org.apache.hadoop.conf.Configuration
2 import org.apache.hadoop.fs.Path
3 import org.apache.hadoop.io.{Text}
4 import org.apache.hadoop.mapreduce.{Job, Mapper, Reducer}
```

```

5 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat
6 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat
7
8 import scala.jdk.CollectionConverters._
9
10 object TableReferJob {
11
12     class RelationMapper extends Mapper[Object, Text, Text, Text] {
13         override def map(key: Object, value: Text, context: Mapper[Object, Text, Text,
14 Text]#Context): Unit = {
15             val line = value.toString.trim
16             if (line.nonEmpty && !line.startsWith("child")) { // 跳过表头
17                 val parts = line.split("\\s+")
18                 if (parts.length == 2) {
19                     val child = parts(0)
20                     val parent = parts(1)
21                     // 发出 parent -> child
22                     context.write(new Text(parent), new Text("C:" + child))
23                     // 发出 child -> parent
24                     context.write(new Text(child), new Text("P:" + parent))
25                 }
26             }
27         }
28
29     class RelationReducer extends Reducer[Text, Text, Text, Text] {
30         override def reduce(key: Text, values: java.lang.Iterable[Text],
31                             context: Reducer[Text, Text, Text, Text]#Context): Unit = {
32
33             val grandchildren = scala.collection.mutable.ListBuffer[String]()
34             val grandparents = scala.collection.mutable.ListBuffer[String]()
35
36             for (v <- values.asScala) {
37                 if (v.toString.startsWith("C:")) {
38                     grandchildren += v.toString.substring(2)
39                 } else if (v.toString.startsWith("P:")) {
40                     grandparents += v.toString.substring(2)
41                 }
42             }
43
44             for (gc <- grandchildren; gp <- grandparents) {
45                 context.write(new Text(gc), new Text(gp))
46             }
47         }
48     }
49
50     def main(args: Array[String]): Unit = {
51         if (args.length < 2) {
52             System.err.println("Usage: TableReferJob <input1> [<input2> ...] <output>")
53             System.exit(1)
54         }
55     }

```

```

56     val conf = new Configuration()
57     val job = Job.getInstance(conf, "Scala Hadoop TableRefer Job")
58
59     job.setJarByClass(this.getClass)
60     job.setMapperClass(classOf[RelationMapper])
61     job.setReducerClass(classOf[RelationReducer])
62
63     job.setMapOutputKeyClass(classOf[Text])
64     job.setMapOutputValueClass(classOf[Text])
65     job.setOutputKeyClass(classOf[Text])
66     job.setOutputValueClass(classOf[Text])
67
68     // 添加输入路径
69     for (i <- 0 until args.length - 1) {
70       FileInputFormat.addInputPath(job, new Path(args(i)))
71     }
72
73     // 设置输出路径
74     FileOutputFormat.setOutputPath(job, new Path(args.last))
75
76     System.exit(if (job.waitForCompletion(true)) 0 else 1)
77   }
78 }
```

### 3.4.2 实验步骤

- 准备包含亲子关系的数据文件，第一行为表头“child parent”

在本地创建并编辑文件child-parent:

```

1 child parent
2 Tom Lucy
3 Tom Jack
4 Jone Lucy
5 Jone Jack
6 Lucy Mary
7 Lucy Ben
8 Jack Alice
9 Jack Jesse
10 Terry Alice
11 Terry Jesse
12 Philip Terry
13 Philip Alma
14 Mark Terry
15 Mark Alma
```

使用命令 `hdfs dfs -put child-parent \user\lfl\input\tablerefer` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

- 在spark集群上运行作业

在项目根目录下运行命令

```
1 | hadoop jar target/scala-2.11/MapReduce-assembly-0.1.jar tablereferJob \
2 |   /user/lfl/input/tablerefer/child-parent \
3 |   /user/lfl/output/mapredtablerefer_result
```

### 3.4.3 实验结果

运行成功后可在Web页面查看输出文件是否存在：

| /user/lfl/output/mapredtablerefer_result |       |            |      |                   |             |            |            | Go! |
|--|-------|------------|------|-------------------|-------------|------------|------------|-----|
| Permission                               | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |     |
| -rw-r--r--                               | lfl   | supergroup | 0 B  | 2025/7/4 22:43:51 | 2           | 128 MB     | _SUCCESS   |     |
| -rw-r--r--                               | lfl   | supergroup | 0 B  | 2025/7/4 22:43:51 | 2           | 128 MB     | part-00000 |     |
| -rw-r--r--                               | lfl   | supergroup | 23 B | 2025/7/4 22:43:53 | 2           | 128 MB     | part-00001 |     |
| -rw-r--r--                               | lfl   | supergroup | 48 B | 2025/7/4 22:43:53 | 2           | 128 MB     | part-00002 |     |
| -rw-r--r--                               | lfl   | supergroup | 0 B  | 2025/7/4 22:43:54 | 2           | 128 MB     | part-00003 |     |
| -rw-r--r--                               | lfl   | supergroup | 0 B  | 2025/7/4 22:43:54 | 2           | 128 MB     | part-00004 |     |
| -rw-r--r--                               | lfl   | supergroup | 78 B | 2025/7/4 22:43:54 | 2           | 128 MB     | part-00005 |     |

查看并验证运行结果：

```
lfl@master:~/data$ hdfs dfs -cat /user/lfl/output/sparktablerefer_result/part*
grandchild      grandparent
Philip  Alice
Philip  Jesse
Mark   Alice
Mark   Jesse
Tom    Mary
Tom    Ben
Jone   Mary
Jone   Ben
Tom    Alice
Tom    Jesse
Jone   Alice
Jone   Jesse
```

输出文件中包含每个孙子和祖父母的关系。

## 四、Spark任务

### 4.1 数据去重

#### 4.1.1 代码实现

```
1 | import org.apache.spark.{SparkConf, SparkContext}
2 |
3 | object SparkDedupJob {
```

```

4  def main(args: Array[String]): Unit = {
5    if (args.length < 2) {
6      System.err.println("Usage: DedupApp <input1> [<input2> ... <inputN>] <output>")
7      System.exit(1)
8    }
9
10   val output = args.last
11   val inputPaths = args.slice(0, args.length - 1)
12
13   val conf = new SparkConf().setAppName("DedupApp")
14   val sc = new SparkContext(conf)
15
16   val inputRDDs = inputPaths.map(path => sc.textFile(path))
17   val mergedRDD = sc.union(inputRDDs)
18
19   val result = mergedRDD.distinct().sortBy(x => x)
20
21   result.saveAsTextFile(output)
22
23   sc.stop()
24 }
25 }
```

## 4.1.2 实验步骤

### 1. 准备包含重复数据的输入文件

在本地创建并编辑文件file1、file2：

```

1  file1
2  2012-3-1 a
3  2012-3-2 b
4  2012-3-3 c
5  2012-3-4 d
6  2012-3-5 a
7  2012-3-6 b
8  2012-3-7 c
9  2012-3-3 c
10
11 file2
12 2012-3-1 b
13 2012-3-2 a
14 2012-3-3 b
15 2012-3-4 d
16 2012-3-5 a
17 2012-3-6 c
18 2012-3-7 d
19 2012-3-3 c
20
21 # 冗余数据
22 2012-3-3 c
23 2012-3-4 d
```

使用命令 `hdfs dfs -put file1 file2 \user\lfl\input\dedup` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

## Browse Directory

/user/lfl/input/dedup

Go!

| Permission | Owner | Group      | Size | Last Modified      | Replication | Block Size | Name  |
|------------|-------|------------|------|--------------------|-------------|------------|-------|
| -rw-r--r-- | lfl   | supergroup | 88 B | 2025/6/28 17:39:11 | 2           | 128 MB     | file1 |
| -rw-r--r-- | lfl   | supergroup | 88 B | 2025/6/28 17:39:11 | 2           | 128 MB     | file2 |

### 2. 编译打包代码成 JAR 文件

实验过程已经把所有代码实现，因此四个源文件均一起编译并打包为 `sparkjob_2.11-0.1.jar`，提交作业时只需选定相应的类，如下：

```
lfl@master:~/SparkJob$ sbt 'show discoveredMainClasses'
[info] welcome to sbt 1.5.5 (Oracle Corporation Java 1.8.0_151)
[info] loading settings for project sparkjob-build from plugins.sbt ...
[info] loading project definition from /home/lfl/SparkJob/project
[info] loading settings for project sparkjob from build.sbt ...
[info] set current project to SparkJob (in build file:/home/lfl/SparkJob/)
[info] * SparkAvgJob
[info] * SparkDedupJob
[info] * SparkSortedIndexJob
[info] * SparkTableReferJob
[success] Total time: 13 s, completed Jul 3, 2025 9:54:56 PM
```

后续任务直接在集群上运行作业

### 3. 在spark集群上运行作业

在工程根目录下使用命令

```
1 spark-submit \
2   --class SparkDedupJob \
3   --master spark://master:7077 \
4   target/scala-2.11/sparkjob_2.11-0.1.jar \
5   hdfs:/user/lfl/input/dedup/file1 \
6   hdfs:/user/lfl/input/dedup/file2 \
7   hdfs:/user/lfl/output/sparkdedup_result
```

## 4.1.3 实验结果

成功运行后可在Web页面查看输出文件是否存在

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |
|------------|-------|------------|------|-------------------|-------------|------------|------------|
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/3 10:38:22 | 2           | 128 MB     | _SUCCESS   |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/3 10:38:21 | 2           | 128 MB     | part-00000 |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/3 10:38:21 | 2           | 128 MB     | part-00001 |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/3 10:38:21 | 2           | 128 MB     | part-00002 |
| -rw-r--r-- | lfl   | supergroup | 33 B | 2025/7/3 10:38:21 | 2           | 128 MB     | part-00003 |

查看并验证输出结果

```
lfl@master:~/MapReduce$ hdfs dfs -cat /user/lfl/output/sparkdedup_result/part*
2012-3-1 a
2012-3-1 b
2012-3-2 a
2012-3-2 b
2012-3-3 b
2012-3-3 c
2012-3-4 d
2012-3-5 a
2012-3-6 b
2012-3-6 c
2012-3-7 c
2012-3-7 d
```

可以看到两个文件的重复数据均已删除，输出文件中只包含去重后的数据

## 4.2 数据排序

### 4.1.1 代码实现

```
1 import org.apache.spark.sql.SparkSession
2
3 object SparkSortedIndexJob {
4     def main(args: Array[String]): Unit = {
5         if (args.length < 2) {
6             System.err.println("Usage: SparkSortedIndexJob <input1> [<input2> ...]
<output>")
7             System.exit(1)
8         }
9
10        val inputPaths = args.slice(0, args.length - 1)
11        val outputPath = args.last
12
13        val spark = SparkSession.builder()
14            .appName("Spark Sorted Index Job")
15            .getOrCreate()
16
17        val sc = spark.sparkContext
```

```

18
19 // 读取多个输入路径并合并
20 val inputRDD = sc.textFile(inputPaths.mkString(","))
21
22 // 清理空行并尝试转换为 Int
23 val numbersRDD = inputRDD
24   .map(_.trim)
25   .filter(_.nonEmpty)
26   .flatMap(line =>
27     try {
28       Some(line.toInt)
29     } catch {
30       case _: NumberFormatException => None
31     }
32   )
33
34 // 排序并 zip 上从 1 开始的索引
35 val sortedIndexedRDD = numbersRDD
36   .sortBy(num => num)
37   .zipWithIndex()
38   .map { case (value, idx) => (idx + 1, value) } // index 从 1 开始
39
40 // 格式化为文本: index \t value
41 val outputRDD = sortedIndexedRDD.map { case (index, value) => s"$index\t$value" }
42
43 // 写入输出
44 outputRDD.saveAsTextFile(outputPath)
45
46 spark.stop()
47 }
48 }
```

## 4.2.2 实验步骤

1. 准备包含整数数据的输入文件

在本地创建并编辑文件file3、file4、file5：

```
lfl@master:~/data$ cat file3 file4 file5
2
32
654
32
15
756
65223
5956
22
650
92
26
54
6
```

使用命令 `hdfs dfs -put file3 file4 file5 \user\lfl\input\sort` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

| /user/lfl/input/sort |       |            |      |                    |             |            | Go!                   |
|----------------------|-------|------------|------|--------------------|-------------|------------|-----------------------|
| Permission           | Owner | Group      | Size | Last Modified      | Replication | Block Size | Name                  |
| -rw-r--r--           | lfl   | supergroup | 25 B | 2025/6/28 17:39:12 | 2           | 128 MB     | <a href="#">file3</a> |
| -rw-r--r--           | lfl   | supergroup | 15 B | 2025/6/28 17:39:12 | 2           | 128 MB     | <a href="#">file4</a> |
| -rw-r--r--           | lfl   | supergroup | 8 B  | 2025/6/28 17:39:12 | 2           | 128 MB     | <a href="#">file5</a> |

## 2. 在spark集群上运行作业

在工程根目录下使用命令

```
1 spark-submit \
2   --class SparkSortedIndexJob \
3   --master spark://master:7077 \
4   target/scala-2.11/sparkjob_2.11-0.1.jar \
5   hdfs:/user/lfl/input/sort/file3 \
6   hdfs:/user/lfl/input/sort/file4 \
7   hdfs:/user/lfl/input/sort/file5 \
8   hdfs:/user/lfl/output/sparksorted_result
```

### 4.2.3 实验结果

运行成功后可在Web页面查看输出文件是否存在：

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |
|------------|-------|------------|------|-------------------|-------------|------------|------------|
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/3 10:44:23 | 2           | 128 MB     | _SUCCESS   |
| -rw-r--r-- | lfl   | supergroup | 23 B | 2025/7/3 10:44:22 | 2           | 128 MB     | part-00000 |
| -rw-r--r-- | lfl   | supergroup | 27 B | 2025/7/3 10:44:22 | 2           | 128 MB     | part-00001 |
| -rw-r--r-- | lfl   | supergroup | 31 B | 2025/7/3 10:44:22 | 2           | 128 MB     | part-00002 |

查看并验证运行结果：

```
lfl@master:~/MapReduce$ hdfs dfs -cat /user/lfl/output/sparksorted_result/part*
1      2
2      6
3     15
4     22
5     26
6     32
7     32
8     54
9     92
10    650
11    654
12    756
13   5956
14  65223
```

可观察到输出文件中的数据按升序排列

## 4.3 平均

### 4.3.1 代码实现

```
1 import org.apache.spark.sql.SparkSession
2
3 object SparkAvgJob {
4   def main(args: Array[String]): Unit = {
5     if (args.length < 2) {
6       System.err.println("Usage: SparkAvgJob <input1> [<input2> ...] <output>")
7       System.exit(1)
8     }
9
10    val inputPaths = args.slice(0, args.length - 1)
11    val outputPath = args.last
12
13    val spark = SparkSession.builder()
14      .appName("SparkAvgJob")
15      .getOrCreate()
16
17    val sc = spark.sparkContext
```

```

18
19     val inputRDD = sc.textFile(inputPaths.mkString(","))
20
21     val cleanedRDD = inputRDD
22         .map(_.replaceAll("[\u00A0\u2000-\u200B\u3000\uFEFF]", " "))
23         // 替换全角空格
24         .map(_.replaceAll("\s+", " ").trim)           // 替换多个空格或 tab 为一个空格
25         .filter(_.nonEmpty)
26
27     val parsed = cleanedRDD.flatMap { line =>
28         val parts = line.split(" ")
29         if (parts.length == 2) {
30             try Some((parts(0), parts(1).toInt)) catch {
31                 case _: NumberFormatException => None
32             }
33         } else None
34     }
35
36     val averageByStudent = parsed
37         .mapValues(score => (score, 1)) // (name, (score, 1))
38         .reduceByKey { case ((s1, c1), (s2, c2)) => (s1 + s2, c1 + c2) }
39         .mapValues { case (sum, count) => sum / count } // 取整平均
40
41     val outputRDD = averageByStudent
42         .sortByKey() // 可选: 按姓名排序输出
43         .map { case (name, avg) => s"$name $avg" }
44
45     outputRDD.saveAsTextFile(outputPath)
46
47     spark.stop()
48 }
```

## 4.3.2 实验步骤

- 准备包含学生姓名和成绩的数据文件，每行格式为“姓名 成绩”。

在本地创建并编辑文件math、chinese、english:

```

1 # math
2 张三 88
3 李四 99
4 王五 66
5 赵六 77
6
7 #chinese
8 张三 78
9 李四 89
10 王五 96
11 赵六 67
12
13 #english
```

```
14 | 张三 80  
15 | 李四 82  
16 | 王五 84  
17 | 赵六 86
```

使用命令 `hdfs dfs -put math chinese english \user\lfl\input\average` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

| /user/lfl/input/average |       |            |      |                   |             |            |         | Go! |
|-------------------------|-------|------------|------|-------------------|-------------|------------|---------|-----|
| Permission              | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name    |     |
| -rw-r--r--              | lfl   | supergroup | 40 B | 2025/7/3 11:23:50 | 2           | 128 MB     | chinese |     |
| -rw-r--r--              | lfl   | supergroup | 40 B | 2025/7/3 11:23:50 | 2           | 128 MB     | english |     |
| -rw-r--r--              | lfl   | supergroup | 42 B | 2025/7/3 11:23:50 | 2           | 128 MB     | math    |     |

## 2. 在spark集群上运行作业

在项目根目录下运行命令

```
1 spark-submit \  
2   --class SparkAvgJob \  
3   --master spark://master:7077 \  
4   target/scala-2.11/sparkjob_2.11-0.1.jar \  
5   hdfs://user/lfl/input/average/math \  
6   hdfs://user/lfl/input/average/chinese \  
7   hdfs://user/lfl/input/average/english \  
8   hdfs://user/lfl/output/sparkaverage_result
```

### 4.3.3 实验结果

运行成功后可在Web页面查看输出文件是否存在：

| /user/lfl/output/sparkaverage_result |       |            |      |                   |             |            |            | Go! |
|--------------------------------------|-------|------------|------|-------------------|-------------|------------|------------|-----|
| Permission                           | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |     |
| -rw-r--r--                           | lfl   | supergroup | 0 B  | 2025/7/3 11:55:13 | 2           | 128 MB     | _SUCCESS   |     |
| -rw-r--r--                           | lfl   | supergroup | 20 B | 2025/7/3 11:55:11 | 2           | 128 MB     | part-00000 |     |
| -rw-r--r--                           | lfl   | supergroup | 10 B | 2025/7/3 11:55:10 | 2           | 128 MB     | part-00001 |     |
| -rw-r--r--                           | lfl   | supergroup | 10 B | 2025/7/3 11:55:11 | 2           | 128 MB     | part-00002 |     |

查看并验证运行结果：

```
lfl@master:~$ hdfs dfs -cat /user/lfl/output/sparkaverage_result/part*  
张三 82  
李四 90  
王五 82  
赵六 76
```

可观察到输出文件中包含每个学生的平均成绩，验证了计算平均功能。

## 4.4 单表关联

### 4.4.1 代码实现

```
1 import org.apache.spark.{SparkConf, SparkContext}
2
3 object SparkTableReferJob {
4     def main(args: Array[String]): Unit = {
5         if (args.length != 2) {
6             System.err.println("Usage: SparkTableReferJob <input> <output>")
7             System.exit(1)
8         }
9
10        val inputPath = args(0)
11        val outputPath = args(1)
12
13        val conf = new SparkConf().setAppName("SparkTableReferJob")
14        val sc = new SparkContext(conf)
15
16        // 读取数据并过滤表头
17        val lines = sc.textFile(inputPath)
18            .filter(line => !line.trim.startsWith("child"))
19
20        // 拆分为 (child, parent)
21        val relations = lines.map(_.trim.split("\\s+")).map {
22            case Array(child, parent) => (child, parent)
23        }
24
25        // 模拟左表: key = parent, value = "L_child"
26        val left = relations.map { case (child, parent) => (parent, "L_" + child) }
27
28        // 模拟右表: key = child, value = "R_parent"
29        val right = relations.map { case (child, parent) => (child, "R_" + parent) }
30
31        // 合并并按 key 分组
32        val joined = left.union(right).groupByKey()
33
34        // Reduce 阶段做自连接: 分开左右表并做笛卡尔积
35        val grandchildGrandparent = joined.flatMap {
36            case (_, values) =>
37                val (leftValues, rightValues) = values.partition(_.startsWith("L_"))
38                val grandchildren = leftValues.map(_.substring(2))
39                val grandparents = rightValues.map(_.substring(2))
40                for {
41                    gc <- grandchildren
42                    gp <- grandparents
43                } yield (gc, gp)
44        }
```

```

46 // 添加表头并保存为文本
47 val header = sc.parallelize(Seq("grandchild\tgrandparent"))
48 val results = grandchildGrandparent.map { case (gc, gp) => s"$gc\t$gp" }
49
50 header.union(results).saveAsTextFile(outputPath)
51
52 sc.stop()
53 }
54 }
```

## 4.4.2 实验步骤

- 准备包含亲子关系的数据文件，第一行为表头“child parent”

在本地创建并编辑文件child-parent:

```

1 child parent
2 Tom Lucy
3 Tom Jack
4 Jone Lucy
5 Jone Jack
6 Lucy Mary
7 Lucy Ben
8 Jack Alice
9 Jack Jesse
10 Terry Alice
11 Terry Jesse
12 Philip Terry
13 Philip Alma
14 Mark Terry
15 Mark Alma
```

使用命令 `hdfs dfs -put child-parent \user\lfl\input\tablerefer` 将文件上传到Hadoop集群（路径内文件夹均已提前创建），上传成功后通过访问 `master:50070` Web页面查看：

- 在spark集群上运行作业

在项目根目录下运行命令

```

1 spark-submit \
2   --class SparkTableReferJob \
3   --master spark://master:7077 \
4   target/scala-2.11/sparkjob_2.11-0.1.jar \
5   hdfs:/user/lfl/input/tablerefer/child-parent \
6   hdfs:/user/lfl/output/sparktablerefer_result
```

## 4.4.3 实验结果

运行成功后可在Web页面查看输出文件是否存在：

| Permission | Owner | Group      | Size | Last Modified     | Replication | Block Size | Name       |
|------------|-------|------------|------|-------------------|-------------|------------|------------|
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/3 17:16:50 | 2           | 128 MB     | _SUCCESS   |
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/3 17:16:47 | 2           | 128 MB     | part-00000 |
| -rw-r--r-- | lfl   | supergroup | 23 B | 2025/7/3 17:16:48 | 2           | 128 MB     | part-00001 |
| -rw-r--r-- | lfl   | supergroup | 48 B | 2025/7/3 17:16:48 | 2           | 128 MB     | part-00002 |
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/3 17:16:45 | 2           | 128 MB     | part-00003 |
| -rw-r--r-- | lfl   | supergroup | 0 B  | 2025/7/3 17:16:45 | 2           | 128 MB     | part-00004 |
| -rw-r--r-- | lfl   | supergroup | 78 B | 2025/7/3 17:16:48 | 2           | 128 MB     | part-00005 |

查看并验证运行结果：

```
lfl@master:~/data$ hdfs dfs -cat /user/lfl/output/sparktablerefer_result/part*
grandchild      grandparent
Philip  Alice
Philip  Jesse
Mark   Alice
Mark   Jesse
Tom    Mary
Tom    Ben
Jone   Mary
Jone   Ben
Tom    Alice
Tom    Jesse
Jone   Alice
Jone   Jesse
```

输出文件中包含每个孙子和祖父母的关系。

## 五、遇到的问题及解决方案

### 5.1 任务运行时slave节点的NodeManager进程终止

- 问题描述：在Reduce任务的copy拉去阶段时slave节点的NodeManager进程终止，查看日志后发现：
 

```
2025-07-02 21:17:06,991 WARN org.apache.hadoop.yarn.server.nodemanager.containermanager.monitor.ContainersMonitorImpl: NodeManager configured with 8 G physical memory allocated to containers, which is more than 80% of the total physical memory available (3.8 G). Thrashing might happen.
```

```
2025-07-03 11:57:48,828 ERROR org.apache.hadoop.yarn.server.nodemanager.NodeManager: RECEIVED SIGNAL 15: SIGTERM
```
- 解决方案：目前认为的可能原因是内存溢出，导致ResourceManager终止了NodeManager进程，通过修改 `yarn-site.xml` 文件配置：

```
1 <property>
2   <name>yarn.nodemanager.resource.memory-mb</name>
3   <value>3072</value> <!-- 3GB 内存 -->
4 </property>
5 <property>
6   <name>yarn.nodemanager.resource.cpu-vcores</name>
7   <value>2</value> <!-- CPU核心数 2 -->
8 </property>
```

## 五、实验总结

通过本次实验，我成功使用 Hadoop MapReduce 和 Spark 框架实现了数据去重、数据排序、计算平均值和单表关联等常见数据处理任务。在实验过程中，深入理解了分布式计算框架的工作原理和编程模型，同时也掌握了在 Hadoop 集群和 Spark 系统上运行作业的方法。

在 Hadoop MapReduce 任务中，需要手动编写 Mapper 和 Reducer 类，对数据进行处理和聚合，代码相对复杂，但能够更好地控制数据处理的细节。而在 Spark 任务中，使用 RDD 和 DataFrame 等高级抽象，代码更加简洁，开发效率更高。

对比两种框架，Hadoop MapReduce 更适合处理大规模批处理任务，其编程模型虽然繁琐，但具有良好的容错性和稳定性；Spark 则在处理速度和灵活性上更具优势，尤其适合迭代计算和实时数据处理场景，其 DAG 执行引擎能显著提升复杂任务的执行效率。

通过本次实验，我的分布式系统开发能力得到了显著提升，不仅掌握了 Hadoop 和 Spark 的核心编程模型，还学会了如何在集群环境中部署和调试应用程序。同时，也对分布式计算的优势和挑战有了更深入的理解，为今后进一步学习和应用大数据技术奠定了坚实的基础。在未来的学习和工作中，我将继续探索更多大数据处理技术和框架，不断提升自己的专业技能。