

# 第4章MapReduce编程入门实训

---

(源自:<https://biglab.site>)

(版本:Ver1.4-20231024)

## 第4章MapReduce编程入门实训

### 实训1 获取成绩表的最高分记录

训练要点

需求说明

实现思路及步骤

作业要求

实现参考

准备测试数据

进入服务器/root/hadoop目录:

获取实验数据

上传数据到hdfs

验证数据是否已上传

编写脚本

FindMaxScore类

FindMaxMapper类

FindMaxReducer类

编译与导出jar

编译生成jar包

上传jar包到master

运行MR程序

使用crt或xshell进入master

运行结果如:

查看HDFS结果输出

## 实训1 获取成绩表的最高分记录

---

### 训练要点

1. 掌握MapReduce的编程的基本操作
2. 掌握MapReduce编程的一些常用的数据处理方法, 如最大值, 去重等
3. 掌握对MapReduce的执行流程
4. 掌握MapReduce程序的输入输出格式

### 需求说明

样例文件Subject\_score成绩表,每行包含两个字段 科目和分数, 要求获取成绩列表中每个科目成绩的最高记录, 输出到最高成绩表中

## 实现思路及步骤

1. 准备测试数据，从服务器上下载测试文件subject\_score.txt到linux系统; 然后从linux系统上传测试文件到hdfs
2. 编写脚本，实现成绩表最高分统计的MR脚，包括FindMaxScore驱动类， FindMaxMapper类和 FindMaxReducer类
3. 编译和导出hadoop.jar并上传到集群的master服务器上
4. 使用hadoop jar hadoop.jar命令运行MR程序，观察运行过程
5. 观察hdfs上输出文件内容是否达到预期，内容中应该包含每科的最高分。

## 作业要求

1. 提供集群网络拓扑图，图中需说明master,slave的名称、IP、JPS进程名。
2. 在<http://master:9870>上截取本小组集群中本成员目录下/user/myname中上传的文件,需包含 subject\_score.txt文件的截图
3. 在eclipse中，分别截图 map类， reduce类， main方法的源码图
4. 使用hadoop jar hadoop.jar命令运行MR程序，运行过程截图
5. 在<http://master:9870>的文件系统中，打开运行输出结果:/user/myname/output\_maxscore/下的文件内容，截图

## 实现参考

### 准备测试数据

进入服务器/root/hadoop目录:

```
1 | cd /root/hadoop
```

获取实验数据

```
1 | wget https://biglab.site/b37066/file/subject_score.txt
```

上传数据到hdfs

```
1 | hdfs dfs -put ./subject_score.txt /user/myname/
```

验证数据是否已上传

```
1 | hdfs dfs -ls /user/myname/sub*
```

### 编写脚本

FindMaxScore类

```
1 | package chap4_score;  
2 |  
3 | import java.io.IOException;  
4 |  
5 | import org.apache.hadoop.conf.Configuration;  
6 | import org.apache.hadoop.fs.FileSystem;
```

```

7  import org.apache.hadoop.fs.Path;
8  import org.apache.hadoop.io.IntWritable;
9  import org.apache.hadoop.io.Text;
10 import org.apache.hadoop.mapreduce.Job;
11 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
12 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
13 import org.apache.hadoop.util.GenericOptionsParser;
14
15 import utils.ConfUtil;
16
17 public class FindMaxScore {
18
19     public static void main(String[] args) throws IOException,
ClassNotFoundException, InterruptedException {
20         Configuration conf = ConfUtil.GetConf(FindMaxScore.class);
21         String[] otherArgs = new GenericOptionsParser(conf,
args).getRemainingArgs();
22         if (otherArgs.length < 2) {
23             otherArgs = new String[] { "/user/limm/subject_score.txt",
"/user/limm/output_subject_score" };
24         }
25
26         if (otherArgs.length < 2) {
27             System.err.println("Usage: FindMaxScore <in> [<in>...] <out>");
28             System.exit(2);
29         }
30         Job job = Job.getInstance(conf, "FindMaxScore");
31         job.setJarByClass(FindMaxScore.class);
32         job.setMapperClass(FindMaxMapper.class);
33         // job.setCombinerClass(FindMaxReducer.class);
34         job.setReducerClass(FindMaxReducer.class);
35
36         job.setMapOutputKeyClass(Text.class);
37         job.setMapOutputValueClass(IntWritable.class);
38
39         job.setOutputKeyClass(Text.class);
40         job.setOutputValueClass(IntWritable.class);
41         job.setNumReduceTasks(1);
42         for (int i = 0; i < otherArgs.length - 1; ++i) {
43             FileInputFormat.addInputPath(job, new Path(otherArgs[i]));
44         }
45         FileSystem.get(conf).delete(new Path(otherArgs[otherArgs.length -
1]), true);
46         FileOutputFormat.setOutputPath(job, new
Path(otherArgs[otherArgs.length - 1]));
47         System.exit(job.waitForCompletion(true) ? 0 : 1);
48
49     }
50
51 }
52

```

## FindMaxMapper类

```
1 package chap4_score;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.LongWritable;
7 import org.apache.hadoop.io.Text;
8 import org.apache.hadoop.mapreduce.Mapper;
9
10 public class FindMaxMapper extends
Mapper<LongWritable,Text,Text,IntWritable>
11 {
12     Text course = new Text();
13     IntWritable score = new IntWritable();
14     @Override
15     protected void map(LongWritable key, Text value, Context context)
16         throws IOException, InterruptedException {
17         String line = new
String(value.getBytes(),0,value.getLength(),"GBK"); //为了处理乱码
18         String[] values= line.toString().trim().split(" ");
19         course.set(values[0]);
20         score.set(Integer.parseInt(values[1]));
21         context.write(course, score); //这里可能出现乱码,如果不进行 GBK转码的
22     }
23
24
25
26 }
27
```

## FindMaxReducer类

```
1 package chap4_score;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Reducer;
8
9 public class FindMaxReducer extends
Reducer<Text,IntWritable,Text,IntWritable>{
10
11     @Override
12     protected void reduce(Text key, Iterable<IntWritable> values,
13         Context context)
14         throws IOException, InterruptedException {
15         int maxScore=-1;
16         Text course = new Text();
17         for (IntWritable score:values) {
18             if (score.get()>maxScore) {
19                 maxScore= score.get();
20                 course=key;

```

```

21         }
22     }
23     context.write(course,new IntWritable(maxScore));
24 }
25
26
27 }
28

```

## 编译与导出jar

### 编译生成jar包

菜单-》 Build -》 Build Artifacts -》 hadoop-》 Build, 参考4.3节编译生成jar包图例

### 上传jar包到master

1. 在Hadoop项目, 左侧树图中-》 out-》 artifacts-》 hadoop -》 hadoop.jar, 右击hadoop.jar, 菜单中选择复制
2. 打开xftp, 进入master-》 root-》 hadoop目录, 右击选择粘贴

## 运行MR程序

(以下myname需换成本人姓名的拼音, 全拼或简拼)

### 使用crt或xshell进入master

```

1 cd /root/hadoop
2 hadoop jar hadoop.jar chap4_score.FindMaxScore \
3 /user/myname/subject_score.txt \
4 /user/myname/output_subject_score

```

### 运行结果如:

```

1 [root@master ~]# cd /root/hadoop
2 [root@master hadoop]# hadoop jar hadoop.jar chap4_score.FindMaxScore \
3 > /user/myname/subject_score.txt \
4 > /user/myname/output_subject_score
5 SLF4J: Class path contains multiple SLF4J bindings.
6 SLF4J: Found binding in [jar:file:/usr/local/hadoop-
7 3.1.4/share/hadoop/common/lib/slf4j-log4j12-
8 1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
9 SLF4J: Found binding in [jar:file:/usr/local/hadoop-
10 3.1.4/share/hadoop/common/slf4j-log4j12-
11 1.7.30.jar!/org/slf4j/impl/StaticLoggerBinder.class]
12 SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an
13 explanation.
14 SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
15 class name:chap4_score.FindMaxScore
16 2023-10-24 14:56:12,755 INFO client.RMProxy: Connecting to ResourceManager
17 at master/192.168.128.130:8032
18 2023-10-24 14:56:13,930 INFO mapreduce.JobResourceUploader: Disabling
19 Erasure Coding for path: /tmp/hadoop-
20 yarn/staging/root/.staging/job_1698130548927_0001
21 2023-10-24 14:56:15,365 INFO input.FileInputFormat: Total input files to
22 process : 1

```

```
14 2023-10-24 14:56:15,560 INFO mapreduce.JobSubmitter: number of splits:1
15 2023-10-24 14:56:15,906 INFO mapreduce.JobSubmitter: Submitting tokens for
job: job_1698130548927_0001
16 2023-10-24 14:56:15,909 INFO mapreduce.JobSubmitter: Executing with tokens:
[]
17 2023-10-24 14:56:16,274 INFO conf.Configuration: resource-types.xml not
found
18 2023-10-24 14:56:16,274 INFO resource.ResourceUtils: Unable to find
'resource-types.xml'.
19 2023-10-24 14:56:16,886 INFO impl.YarnClientImpl: Submitted application
application_1698130548927_0001
20 2023-10-24 14:56:16,956 INFO mapreduce.Job: The url to track the job:
http://master:8088/proxy/application_1698130548927_0001/
21 2023-10-24 14:56:16,957 INFO mapreduce.Job: Running job:
job_1698130548927_0001
22 2023-10-24 14:56:31,544 INFO mapreduce.Job: Job job_1698130548927_0001
running in uber mode : false
23 2023-10-24 14:56:31,545 INFO mapreduce.Job: map 0% reduce 0%
24 2023-10-24 14:56:41,704 INFO mapreduce.Job: map 100% reduce 0%
25 2023-10-24 14:56:48,790 INFO mapreduce.Job: map 100% reduce 100%
26 2023-10-24 14:56:48,809 INFO mapreduce.Job: Job job_1698130548927_0001
completed successfully
27 2023-10-24 14:56:48,982 INFO mapreduce.Job: Counters: 53
28     File System Counters
29         FILE: Number of bytes read=780006
30         FILE: Number of bytes written=2004973
31         FILE: Number of read operations=0
32         FILE: Number of large read operations=0
33         FILE: Number of write operations=0
34         HDFS: Number of bytes read=482714
35         HDFS: Number of bytes written=63
36         HDFS: Number of read operations=8
37         HDFS: Number of large read operations=0
38         HDFS: Number of write operations=2
39     Job Counters
40         Launched map tasks=1
41         Launched reduce tasks=1
42         Data-local map tasks=1
43         Total time spent by all maps in occupied slots (ms)=30800
44         Total time spent by all reduces in occupied slots (ms)=17756
45         Total time spent by all map tasks (ms)=7700
46         Total time spent by all reduce tasks (ms)=4439
47         Total vcore-milliseconds taken by all map tasks=7700
48         Total vcore-milliseconds taken by all reduce tasks=4439
49         Total megabyte-milliseconds taken by all map tasks=15769600
50         Total megabyte-milliseconds taken by all reduce
tasks=9091072
51     Map-Reduce Framework
52         Map input records=60000
53         Map output records=60000
54         Map output bytes=660000
55         Map output materialized bytes=780006
56         Input split bytes=113
57         Combine input records=0
58         Combine output records=0
```

```

59      Reduce input groups=6
60      Reduce shuffle bytes=780006
61      Reduce input records=60000
62      Reduce output records=6
63      Spilled Records=120000
64      Shuffled Maps =1
65      Failed Shuffles=0
66      Merged Map outputs=1
67      GC time elapsed (ms)=292
68      CPU time spent (ms)=4430
69      Physical memory (bytes) snapshot=326193152
70      Virtual memory (bytes) snapshot=7199862784
71      Total committed heap usage (bytes)=141520896
72      Peak Map Physical memory (bytes)=210395136
73      Peak Map Virtual memory (bytes)=3597443072
74      Peak Reduce Physical memory (bytes)=115798016
75      Peak Reduce Virtual memory (bytes)=3602419712
76      Shuffle Errors
77          BAD_ID=0
78          CONNECTION=0
79          IO_ERROR=0
80          WRONG_LENGTH=0
81          WRONG_MAP=0
82          WRONG_REDUCE=0
83      File Input Format Counters
84          Bytes Read=482601
85      File Output Format Counters
86          Bytes Written=63
87      [root@master hadoop]#

```

## 查看HDFS结果输出

查看hdfs:/user/myname/output\_maxscore目录下结果如：

The screenshot shows the Hadoop web interface with a modal window titled "File information - part-r-00000". The modal displays the following details:

- Block information: Block 0
- Block ID: 1073742393
- Block Pool ID: BP-1924713958-192.168.128.130-1695194852746
- Generation Stamp: 1569
- Size: 63
- Availability:
  - slave1
  - slave2

The "File contents" section shows the following text:

```

化学 99
数学 149
物理 99
生物 99
英语 144
语文 114

```

The background interface shows a "Browse Director" view for the path "/user/myname/output\_subject\_score" with a table of files. The file "part-r-00000" is highlighted, and the modal window is open over it.

####