

词频统计

idea

src

main

java

cn.itcast.hadoop.mr

hadoop.mr

WordCountCombiner

WordCountDriver

WordCountMapper

WordCountReducer

mr

dedup

InvertedIndex

topn

com.itcast.hdfsdemo

42

43

44

45

46

47

48

49

50

51

// 指定要处理的数据所在的位置

FileInputFormat.setInputPaths(wcjob, commaSeparatedPaths: "E:/Idea/mr/{input/*}");

// 指定处理完成之后的结果所保存的位置

FileOutputFormat.setOutputPath(wcjob, new Path(pathString: "E:/Idea/mr/output"));

// 提交程序并且监控打印程序执行情况

boolean res = wcjob.waitForCompletion(verbose: true);

System.exit(res ? 0 : 1);

WordCountDriver : main()

WordCountDriver (1)

Reduce shuffle bytes=84

Reduce input records=6

Reduce output records=6

Spilled Records=12

Shuffled Maps =1

Failed Shuffles=0

Merged Map outputs=1

GC time elapsed (ms)=6

Total committed heap usage (bytes)=385875968

Shuffle Errors

BAD_ID=0

电脑 > Workspaces (E:) > Idea > mr > output

名称

修改日期

类型

大小

._SUCCESS.crc

2019/9/26 15:56

CRC 文件

1 KB

.part-r-00000.crc

2019/9/26 15:56

CRC 文件

1 KB

._SUCCESS

2019/9/26 15:56

文件

0 KB

part-r-00000

2019/9/26 15:56

文件

1 KB

E:\Idea\mr\output\part-r-00000 - Notepad++ [Administrator]

文件(F) 编辑(E) 搜索(S) 视图(V) 编码(N) 语言(L) 设置(T) 工具(O) 宏(M) 运行(R) 插件(P)

part-r-00000

1 hadoop 2

2 hello 3

3 itcast 1

4 mapreduce 2

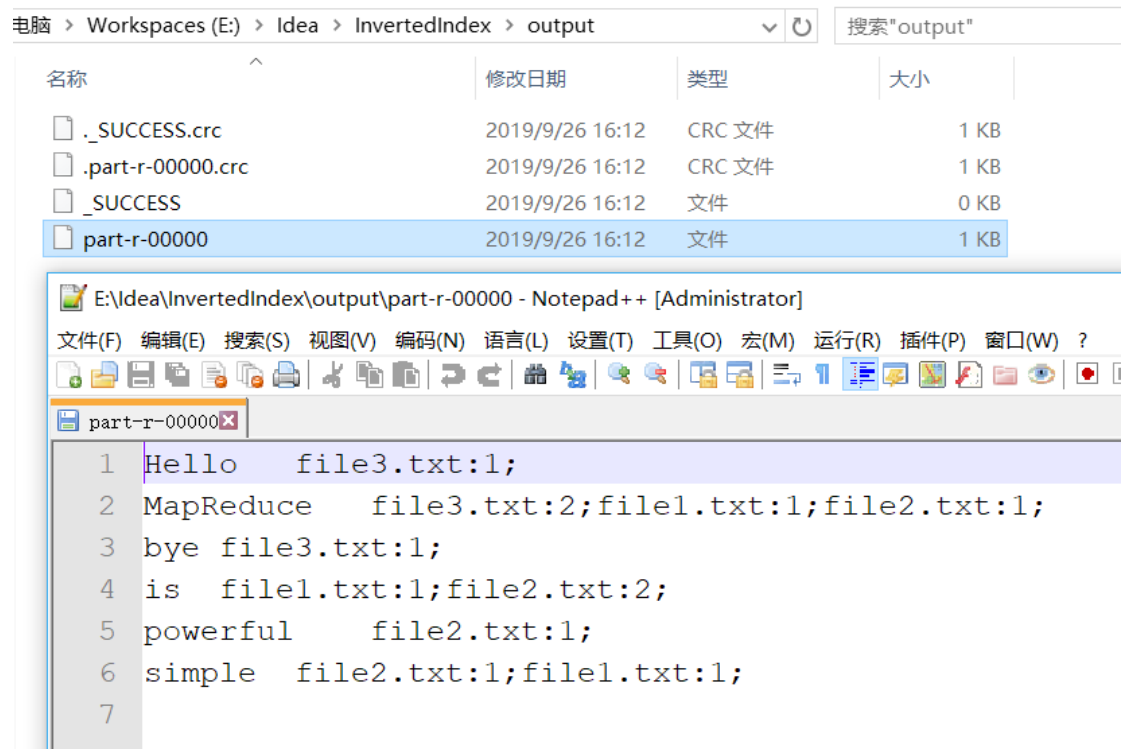
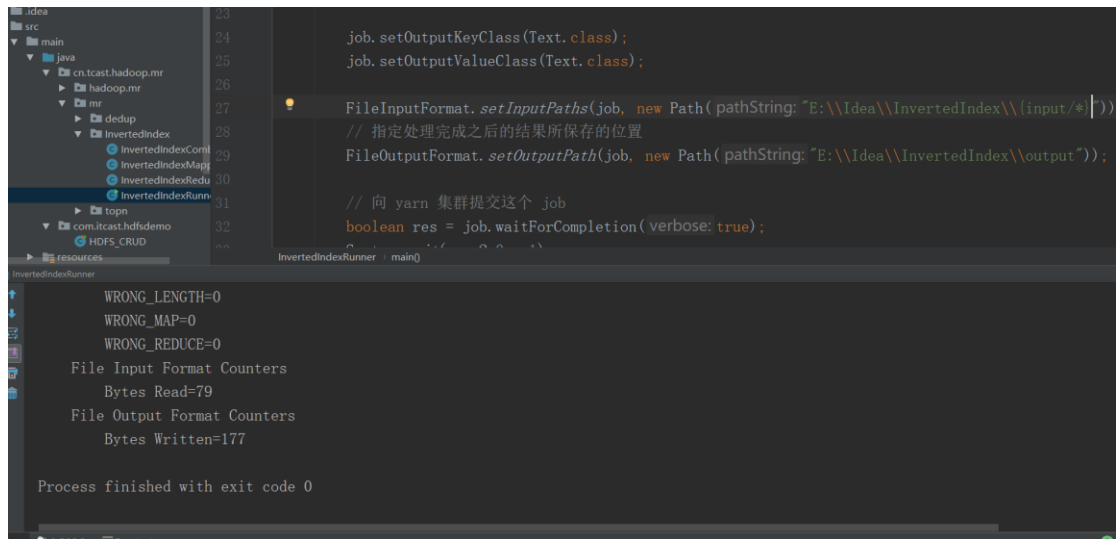
5 spark 1

6 world 1

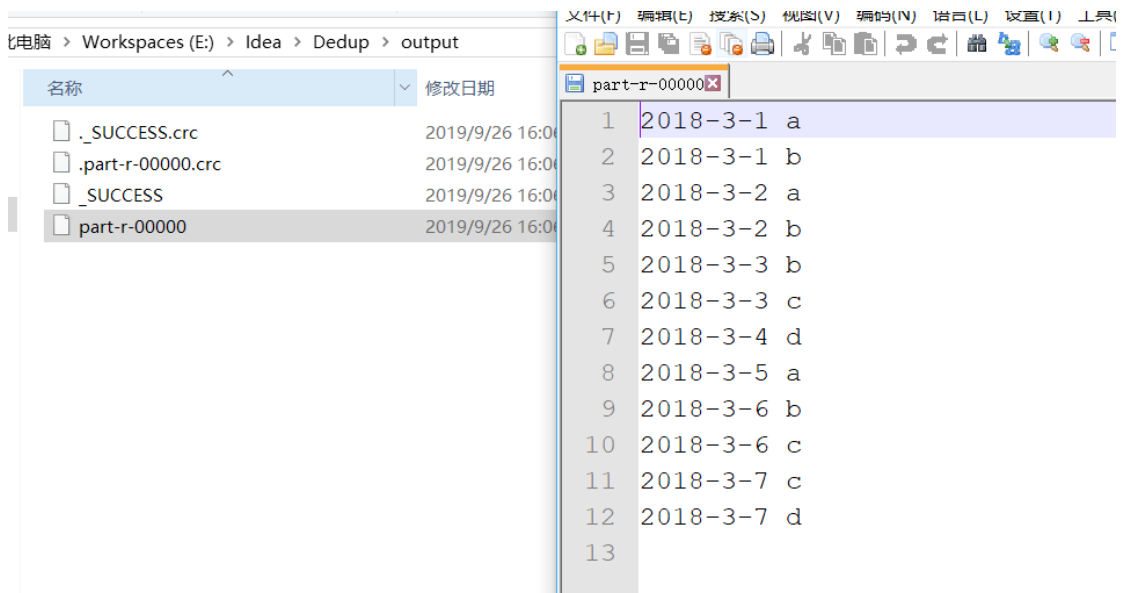
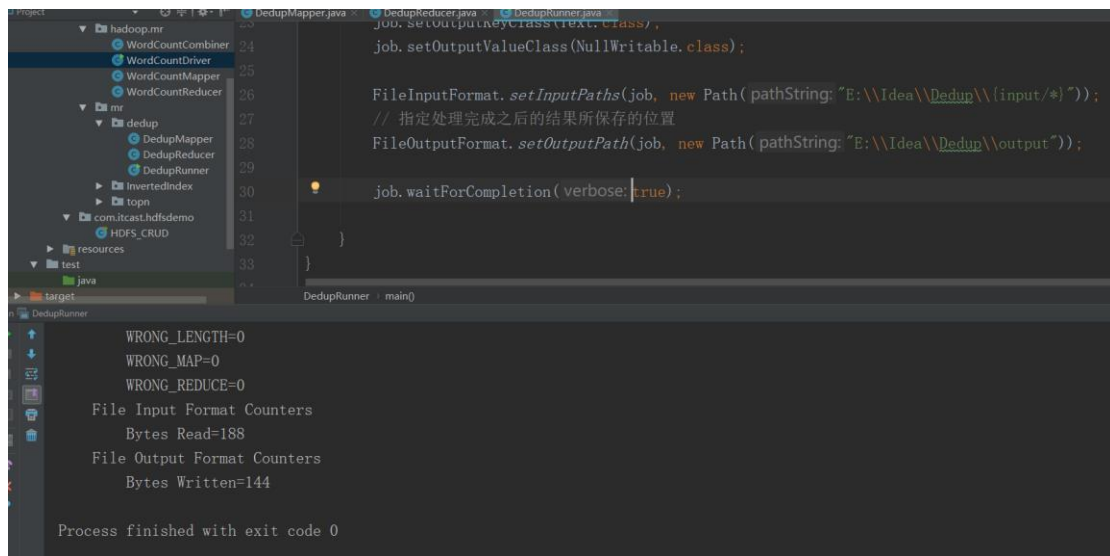
7

54 字节

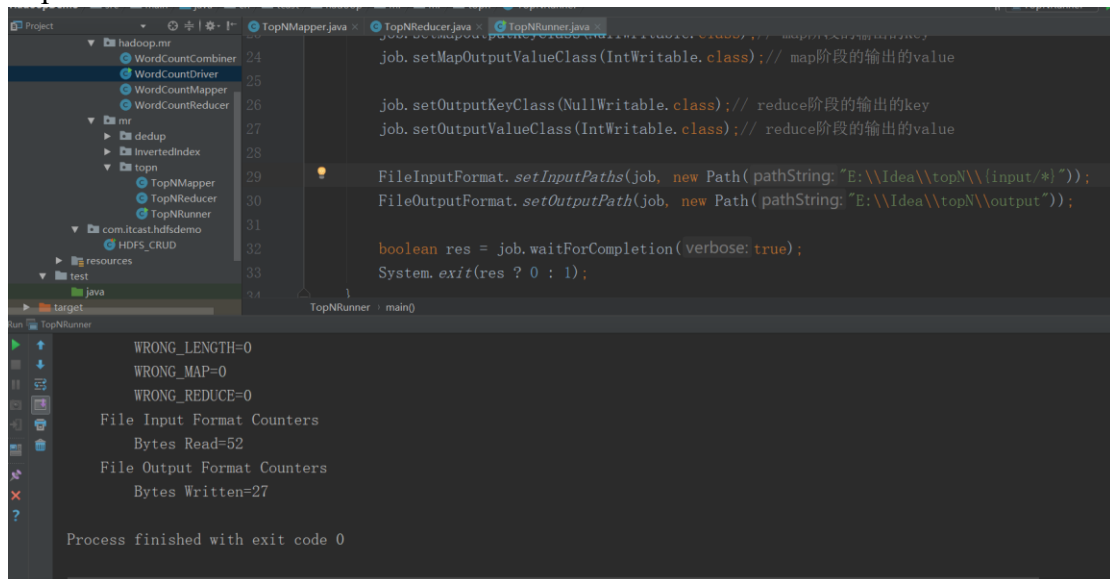
倒排索引



数据去重



TopN



电脑 > Workspaces (E:) > Idea > topN > output

名称	修改日期	类型
._SUCCESS.crc	2019/9/26 16:03	CRC 文件
.part-r-00000.crc	2019/9/26 16:03	CRC 文件
._SUCCESS	2019/9/26 16:03	文件
part-r-00000	2019/9/26 16:03	文件

E:\Idea\topN\output\part-r-00000 - Notepad++ [Administrator]

文件(E) 编辑(E) 搜索(S) 视图(V) 编码(N) 语言(L) 设置(T) 工具(O)

part-r-00000

```
1 20
2 19
3 18
4 17
5 16
6
```