

hadoop和spark使用

学生：张帅豪 学号：18030100101

hadoop和spark使用

1.已有配置：java python hadoop spark

2.实验目的：

3.实验题目

题目一

1.计算每个学生必修课的平均成绩。

2.按科目统计每个班的平均成绩。

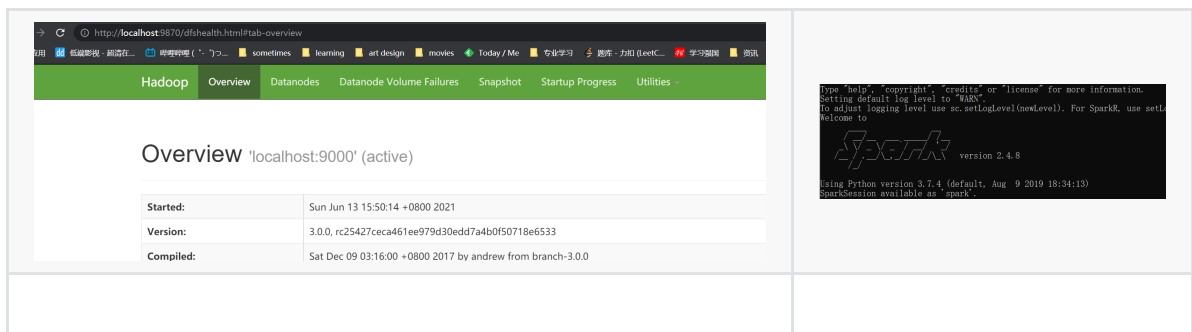
题目二

题目三

1.已有配置：java python hadoop spark

根据老师的安装手册，一步一步安装即可

得到以下界面



2.实验目的：

1. 学习基于MapReduce框架的分布式计算程序设计方法。
2. 学习基于Spark框架的分布式计算程序设计方法。

3.实验题目

- 准备事件：我们先将需要处理的文件上传到浏览器
input （child-parent.txt grades.txt input_file.txt）
input1 （grades.txt）
input2 （child-parent.txt）

题目一

输入文件为学生成绩信息，包含了必修课与选修课成绩，格式如下：

班级1, 姓名1, 科目1, 必修, 成绩1

班级2, 姓名2, 科目1, 必修, 成绩2

班级1, 姓名1, 科目2, 选修, 成绩3

.....,,,

grades.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

170315班,史伦泰,计算机图形学,选修,82

160301班,邱丰琪,计算机图形学,选修,90

160315班,崔和,工程制图,必修,71

170315班,沈明勤,离散数学,必修,61

160315班,田颖,数据结构,必修,97

160312班,孟思,高等数学,必修,53

160315班,江蔼博,工程优化,选修,93

170313班,熊佳,体育,必修,71

170311班,唐洋,软件工程,选修,96

170301班,龙睿,线性代数,必修,65

170211班,江小,数据结构,必修,77

编写两个Hadoop平台上的MapReduce程序，分别实现如下功能：

1. 计算每个学生必修课的平均成绩。
2. 按科目统计每个班的平均成绩。

1.计算每个学生必修课的平均成绩。

- 编写map代码

从数据中选择出**必修**的信息条：key:name value:score

eg: 170315班,史伦泰,计算机图形学,选修,82 史伦泰 82

160301班,邱丰琪,计算机图形学,选修,90 邱丰琪 90

.....

```
1  @Override
2      protected void map(LongWritable key, Text value,
3      Mapper<LongWritable, Text, Text, LongWritable>.Context context)
4          throws IOException, InterruptedException {
5      String line = value.toString();
6      String[] splited = line.split(",");
7      String name = splited[1]; //名字
8      String kind = splited[3]; //必修or选修
9      Integer score = Integer.parseInt(splited[4]); //分数
10     if(kind.equals("必修")){
11         context.write(new Text(name), new LongWritable(score));
12     }
```

- 编写reduce代码

将新k2, v2----> k3,v3

eg 张哥<100,90,80> ----> 求平均数 张哥 90= (100+90+80) /3

```

1  @Override
2      protected void reduce(Text k2, Iterable<LongWritable> v2s,
3                          Reducer<Text, LongWritable, Text,
4                          DoubleWritable>.Context context) throws IOException, InterruptedException {
5          long count = 0L;
6          long sum = 0L;
7          double aver;
8          for (LongWritable v2 : v2s) { //遍历求和
9              count += v2.get();
10             sum = sum + 1;
11         }
12         aver = count*1.0/sum; //防止结果取整数
13         DoubleWritable v3 = new DoubleWritable(aver);
14         context.write(k2, v3);
15     }

```

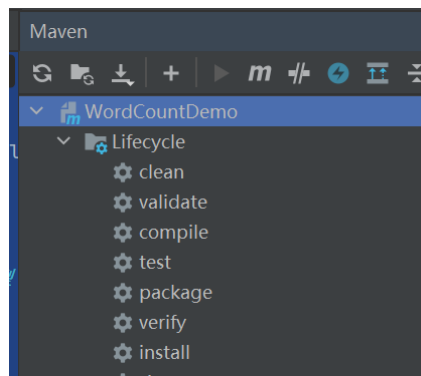
- 主函数（一般不怎么变化）

```

1  public static void main(String[] args) throws Exception {
2      Configuration conf = new Configuration();
3      //1: 创建了一个job任务对象
4      Job job = Job.getInstance(conf,
5      MapReduceWordCountDemo.class.getSimpleName());
6      //打成jar执行
7      job.setJarByClass(MapReduceWordCountDemo.class);
8
9      //2配置job任务对象（八大步骤）
10     //第一步：指定文件的读取方式和读取路径
11     FileInputFormat.setInputPaths(job, args[0]);
12     //FileInputFormat.setInputPaths(job, new
13     Path("http://localhost:9870/explorer.html#/input"));
14
15     //第二步：指定map阶段的处理方式和数据类型
16     job.setMapperClass(MyMapper.class);
17     //设置map阶段k2的类型
18     job.setMapOutputKeyClass(Text.class);
19     //设置map阶段v2的类型
20     job.setMapOutputValueClass(LongWritable.class);
21
22     //第三四五六采用默认方式
23
24     //第7步：指定reduce阶段的处理方式和数据类型
25     job.setReducerClass(MyReducer.class);
26     //设置map阶段k3的类型
27     job.setOutputKeyClass(Text.class);
28     //设置map阶段v3的类型
29     job.setOutputValueClass(DoubleWritable.class);
30
31     //第八步：设置输出类型
32     FileOutputFormat.setOutputPath(job, new Path(args[1]));
33
34     //交给yarn去执行，直到执行结束才退出本程序
35     job.waitForCompletion(true);
36 }

```

- 打包jar包--->11.jar（先clean再package）



- 运行：如图

第一次实验

1. 计算每个学生必修课的平均成绩。

已有11.jar 输入数据在input1 reference--》com.org.xidian.MapReduceWordCountDemo

命令行：hadoop jar F:\11.jar com.org.xidian.MapReduceWordCountDemo /input1 /output11

- 运行结果

Size: 223962

Availability:

- 192.168.91.1

File contents

```

丁一 75.57142857142857
丁一石 71.28571428571429
丁丰 77.42857142857143
丁丰欣 71.85714285714286
丁丽 80.85714285714286
丁丽焕 82.14285714285714
丁义 77.14285714285714
丁伟子 67.0
          
```

2.按科目统计每个班的平均成绩。

- 编写map代码

从数据中选择出信息条：key:subject+"\t"+banji+"\t" value:score

eg: 170315班,史伦泰,计算机图形学,选修,82 计算机图形学 170315班 82
 160301班,邱丰琪,计算机图形学,选修,90 计算机图形学 160301班 90


```

1  @Override
2      protected void map(LongWritable key, Text value,
3      Mapper<LongWritable, Text, Text, LongWritable>.Context context)
4          throws IOException, InterruptedException {
5
6      String line = value.toString();
7      String[] splited = line.split(",");//分离
8      String banji = splited[0];
9      String subject = splited[2];
10     String feature = subject+"\t"+banji+"\t";
11     Integer score = Integer.parseInt(splited[4]);
12     context.write(new Text(feature), new LongWritable(score));
13 }

```

- 编写reduce代码

将新k2, v2----> k3,v3

eg 数学 11班<100,90,80> ----> 求平均数 数学 11班 90= (100+90+80) /3

```

1  @Override
2      protected void reduce(Text k2, Iterable<LongWritable> v2s,
3      Reducer<Text, LongWritable, Text,
4      DoubleWritable>.Context context) throws IOException, InterruptedException {
5
6      long count = 0L;
7      long sum = 0L;
8      double aver;
9      for (LongWritable v2 : v2s) { //遍历相加
10         count += v2.get();
11         sum = sum + 1;
12     }
13     aver = count*1.0 /sum; //取平均
14     DoubleWritable v3 = new DoubleWritable(aver);
15     context.write(k2, v3);
16 }

```

- 主函数 (一般不怎么变化)

```

1  public static void main(String[] args) throws Exception {
2      Configuration conf = new Configuration();
3      //1: 创建了一个job任务对象
4      Job job = Job.getInstance(conf,
5      MapReduceWordCountDemo.class.getSimpleName());
6      //打成jar执行
7      job.setJarByClass(MapReduceWordCountDemo.class);
8
9      //2配置job任务对象 (八大步骤)
10     //第一步: 指定文件的读取方式和读取路径
11     FileInputFormat.setInputPaths(job, args[0]);
12     //FileInputFormat.setInputPaths(job, new
13     Path("http://localhost:9870/explorer.html#/input"));
14
15     //第二步: 指定map阶段的处理方式和数据类型
16     job.setMapperClass(MyMapper.class);
17     //设置map阶段k2的类型
18     job.setMapOutputKeyClass(Text.class);
19 }

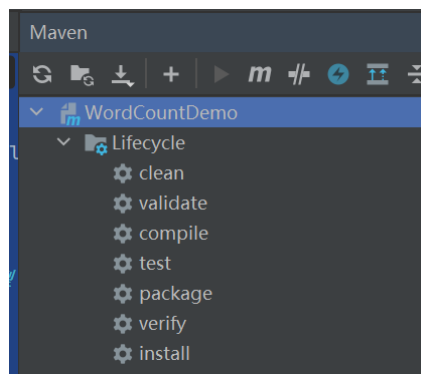
```

```

17 //设置map阶段v2的类型
18 job.setMapOutputValueClass(LongWritable.class);
19
20 //第三四五六采用默认方式
21
22 //第7步：指定reduce阶段的处理方式和数据类型
23 job.setReducerClass(MyReducer.class);
24 //设置map阶段k3的类型
25 job.setOutputKeyClass(Text.class);
26 //设置map阶段v3的类型
27 job.setOutputValueClass(DoubleWritable.class);
28
29 //第八步：设置输出类型
30 FileOutputFormat.setOutputPath(job, new Path(args[1]));
31
32 //交给yarn去执行，直到执行结束才退出本程序
33 job.waitForCompletion(true);
34 }

```

- 打包jar包--->12.jar (先clean再package)



- 运行：如图

第二次实验

2. 按科目统计每个班的平均成绩。

已有12.jar 输入数据在input1 reference--》com.org.xidian.MapReduceWordCountDemo

命令行：hadoop jar F:\12.jar com.org.xidian.MapReduceWordCountDemo /input1 /output12

- 运行结果

File contents

体育 160301班	74.9068493150685
体育 160311班	74.5564738292011
体育 160312班	74.78428927680798
体育 160313班	73.83185840707965
体育 160314班	74.60559796437659
体育 160315班	73.85034965034966
体育 170301班	75.28333333333333
体育 170311班	73.79004037685061

题目二

输入文件的每一行为具有父子/父女/母子/母女/关系的一对人名，例如：

Tim, Andy

Harry, Alice

Mark, Louis

Andy, Joseph

.....,

假定不会出现重名现象。

child-parent.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

Tim, Andy
Harry, Alice
Mark, Louis
Andy, Joseph
Lareina,Dennis
Louisa,Sally
Katherine,Wesley
Lynn,Glendon
Samantha,Hayden
Louise,Nancy
Austin,Sean
Ray,Candy
Claire,Kate
Ellen,Maureen
Iris,Dennis
Terence.Joanna

编写Hadoop平台上的MapReduce程序，找出所有具有grandchild-grandparent关系的人名组。

- 编写map代码

从数据中选择出信息条： key:child value:"<"+parent

 key:parent value:">"+child

eg: 儿a, 父a1 儿a <父a1 父a1 >儿a

.....

```
1  @Override
2      protected void map(LongWritable key, Text value,
3      Mapper<LongWritable, Text, Text, Text>.Context context)
4      throws IOException, InterruptedException {
5
6      String line = value.toString();
7      String[] splited = line.split(",");
8      String child = splited[0];
9      String parent = splited[1];
10     context.write(new Text(child), new Text("<"+parent));
11     context.write(new Text(parent), new Text(">"+child));
12 }
```

- 编写reduce代码

将新k2, v2-----> k3,v3

eg 儿a< <a1,<a2,>a3,<a4,>a5 > 转化1 a1,a2,a4 a3,a5 转化2 a3-a1 a3-a2 a3-a4 a5-a1
a5-a2 a5-a4

```
1  @Override
2      protected void reduce(Text k2, Iterable<Text> v2s,
3                          Reducer<Text, Text, Text, Text>.Context
context) throws IOException, InterruptedException {
4      ArrayList<Text> grandparent = new ArrayList<Text>();
5      ArrayList<Text> grandchild = new ArrayList<Text>();
6      for (Text v2 : v2s) { //对各个values中的值进行处理
7          String ss = v2.toString();
8          String s = ss.substring(0, 1);
9          if (s.equals("<")) {
10             grandchild.add(new Text(ss.substring(1)));
11         } else {
12             grandparent.add(new Text(ss.substring(1)));
13         }
14     }
15     //再将grandparent与grandchild中的东西，一一对应输出。
16     for (int i = 0; i < grandchild.size(); i++) {
17         for (int j = 0; j < grandparent.size(); j++) {
18             context.write(grandchild.get(i), grandparent.get(j));
19         }
20     }
21 }
```

- 主函数（一般不怎么变化）

```
1  public static void main(String[] args) throws Exception {
2      Configuration conf = new Configuration();
3      //1: 创建了一个job任务对象
4      Job job = Job.getInstance(conf,
MapReduceWordCountDemo.class.getSimpleName());
5      //打成jar执行
6      job.setJarByClass(MapReduceWordCountDemo.class);
7
8      //2配置job任务对象（八大步骤）
9      //第一步：指定文件的读取方式和读取路径
10     FileInputFormat.setInputPaths(job, args[0]);
11     //FileInputFormat.setInputPaths(job, new
Path("http://localhost:9870/explorer.html#/input"));
12
13     //第二步：指定map阶段的处理方式和数据类型
14     job.setMapperClass(MyMapper.class);
15     //设置map阶段k2的类型
16     job.setMapOutputKeyClass(Text.class);
17     //设置map阶段v2的类型
18     job.setMapOutputValueClass(LongWritable.class);
19
20     //第三四五六采用默认方式
21
22     //第7步：指定reduce阶段的处理方式和数据类型
23     job.setReducerClass(MyReducer.class);
24     //设置map阶段k3的类型
25     job.setOutputKeyClass(Text.class);
26     //设置map阶段v3的类型
```

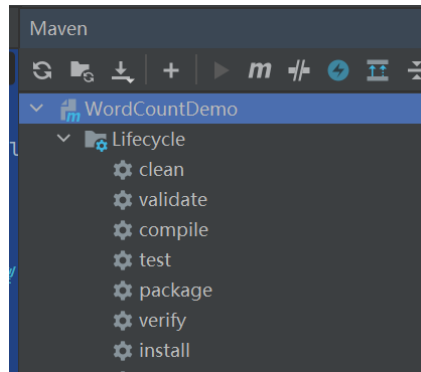


```

27     job.setOutputValueClass(DoubleWritable.class);
28
29     //第八步：设置输出类型
30     FileOutputFormat.setOutputPath(job, new Path(args[1]));
31
32     //交给yarn去执行，直到执行结束才退出本程序
33     job.waitForCompletion(true);
34 }

```

- 打包jar包--->2.jar (先clean再package)



- 运行：如图

第三次实验

编写Hadoop平台上的MapReduce程序，找出所有具有grandchild-grandparent关系的人名组。
已有2.jar 输入数据在input2 reference--》com.org.xidian.MapReduceWordCountDemo

命令行：hadoop jar F:\2.jar com.org.xidian.MapReduceWordCountDemo /input2 /output2

- 运行结果

- 192.168.91.1

File contents

Gino	Wendy
Gino	Tasha
Katrina	Peggy
Katrina	Randy
Katrina	Olina
Katrina	Jean
Katrina	Charlotte
Katrina	Christina

题目三

输入文件为学生成绩信息，包含了必修课与选修课成绩，格式如下：

班级1, 姓名1, 科目1, 必修, 成绩1

班级2, 姓名2, 科目1, 必修, 成绩2

班级1, 姓名1, 科目2, 选修, 成绩3

```
grades.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
170315班,史伦泰,计算机图形学,选修,82
160301班,邱丰琪,计算机图形学,选修,90
160315班,崔和,工程制图,必修,71
170315班,沈明勤,离散数学,必修,61
160315班,田颖,数据结构,必修,97
160312班,孟思,高等数学,必修,53
160315班,江蔼博,工程优化,选修,93
170313班,熊佳,体育,必修,71
170311班,唐洋,软件工程,选修,96
170301班,龙睿,线性代数,必修,65
170311班,江小,数据结构,必修,77
```

编写一个Spark程序，同时实现如下功能：

1. 计算每个学生必修课的平均成绩。
2. 统计学生必修课平均成绩在：90~100,80~89,70~79,60~69和60分以下这5个分数段的人数。

- 先读入数据

```
1 from pyspark import SparkConf, SparkContext
2 conf = SparkConf().setMaster("local").setAppName("spark")
3 sc = SparkContext(conf=conf)
4 textData = sc.textFile("/input1/grades.txt")
```

- 筛选必修课信息条，并选择学生名字和分数列

```
1 def select(x):
2     m = x.split(",")
3     name = m[1]
4     score = int(m[4])
5     return (name, score)
6 lines = textData.filter(lambda line : "必修" in line).map(lambda x:select(x))
```

- 每个学生必修课的平均成绩

```
1 avgData = lines.mapValues(lambda x : (x,1)).reduceByKey(lambda x,y :
2     (x[0]+y[0],x[1] + y[1])).mapValues(lambda x : int(x[0] / x[1]))
3 avgData.saveAsTextFile("/result01")
```

- 统计分数段人数

```

1  def sort(x):
2      if(x>=90 and x<=100):
3          return ("A", 1)
4      if(x>=80 and x<=89):
5          return ("B", 1)
6      if(x>=70 and x<=79):
7          return ("C", 1)
8      if(x>=60 and x<=69):
9          return ("D", 1)
10     if(x<60):
11         return ("E", 1)
12  fData = avgData.map(lambda x:sort(x[1])).reduceByKey(lambda x,y : (x+y))
13  fData.saveAsTextFile("/result02")

```

- 运行：如图

第四次实验

1. 计算每个学生必修课的平均成绩。
2. 统计学生必修课平均成绩在：90~100,80~89,70~79,60~69和60分以下这5个分数段的人数。

已有spark 输入数据在input1

命令行：spark-submit spark.py

输出在result01 result02

- 运行结果

<div> <div>/result01</div> <div>Show 25 entries</div> <div> <div>Permission</div> <div>-rw-r--r--</div> </div> </div> <div> ('崔和', 71) ('沈明勤', 72) ('田颖', 76) ('孟思', 72) ('熊佳', 82) ('龙睿', 68) ('汪小', 76) ('孔秀妣', 66) </div>	<div> <div>/result02</div> <div>Show 25 entries</div> <div> <div>Permission</div> <div>-rw-r--r--</div> </div> </div> <div>File contents</div> <div> ('C', 5615) ('B', 1492) ('D', 1865) ('A', 19) ('E', 28) </div>
---	---