

• 创建Maven项目

- 设置相关依赖 pom.xml

```
1  <dependencies>
2      <dependency>
3          <groupId>org.apache.hadoop</groupId>
4          <artifactId>hadoop-common</artifactId>
5          <version>2.7.5</version>
6      </dependency>
7      <dependency>
8          <groupId>org.apache.hadoop</groupId>
9          <artifactId>hadoop-client</artifactId>
10         <version>2.7.5</version>
11     </dependency>
12     <dependency>
13         <groupId>org.apache.hadoop</groupId>
14         <artifactId>hadoop-hdfs</artifactId>
15         <version>2.7.5</version>
16     </dependency>
17     <dependency>
18         <groupId>org.apache.hadoop</groupId>
19         <artifactId>hadoop-mapreduce-client-core</artifactId>
20         <version>2.7.5</version>
21     </dependency>
22     <dependency>
23         <groupId>junit</groupId>
24         <artifactId>junit</artifactId>
25         <version>4.13</version>
26     </dependency>
27 </dependencies>
28 <build>
29     <plugins>
30         <plugin>
31             <groupId>org.apache.maven.plugins</groupId>
32             <artifactId>maven-compiler-plugin</artifactId>
33             <version>3.0</version>
34             <configuration>
35                 <source>1.8</source>
36                 <target>1.8</target>
```

```

37         <encoding>UTF-8</encoding>
38         <!--      <verbal>true</verbal>-->
39     </configuration>
40 </plugin>
41 <plugin>
42     <groupId>org.apache.maven.plugins</groupId>
43     <artifactId>maven-shade-plugin</artifactId>
44     <version>2.4.3</version>
45     <executions>
46         <execution>
47             <phase>package</phase>
48             <goals>
49                 <goal>shade</goal>
50             </goals>
51             <configuration>
52                 <minimizeJar>true</minimizeJar>
53             </configuration>
54         </execution>
55     </executions>
56 </plugin>
57 </plugins>
58 </build>
59 </project>

```

• Map阶段编写

```

1  package sz.base.java;
2
3  import org.apache.hadoop.io.IntWritable;
4  import org.apache.hadoop.io.LongWritable;
5  import org.apache.hadoop.io.Text;
6  import org.apache.hadoop.mapreduce.Mapper;
7
8  import java.io.IOException;
9
10 /**
11  * map方法的生命周期： 框架每传一行数据就被调用一次
12  * MapReduce: 继承Mapper(分的类型) 形参为<k1,v1><k2,v2>
13  * LongWritable:需要序列化, offset 偏移量, 每行的偏移量    k1
14  * Text: 每行的数据, 整行数据    v1

```

```

15  * Text: 输出的每个单词      k2
16  * IntWritable: 1          v2
17  */
18  public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
19      //实现map
20      @Override
21      protected void map(LongWritable key, Text value, Mapper<LongWritable, Text, Text,
22      IntWritable>.Context context) throws IOException, InterruptedException {
23          //1.读取每行数据
24          String line = value.toString();
25          //2.逗号切分为每行数据
26          if (line != null || line.length() != 0) {
27              String[] words = line.split(",");
28              //3.转换成<单词,1>
29              for (String word : words) {
30                  //context写出到下一步
31                  context.write(new Text(word), new IntWritable(1));
32              }
33          }
34      }
35  }

```

• reduce阶段编写

```

1  package sz.base.java;
2
3  import org.apache.hadoop.io.IntWritable;
4  import org.apache.hadoop.io.Text;
5  import org.apache.hadoop.mapreduce.Reducer;
6
7  import java.io.IOException;
8
9  /**
10   * //生命周期: 框架每传递进来一个kv 组, reduce方法被调用一次
11   * reduce:
12   * Text: k2 单词
13   * IntWritable: v2 : 1
14   * Text: k3 : 单词
15   * IntWritable : v3 : 单词个数

```

```

16  */
17  public class WordCountReduce extends Reducer<Text, IntWritable, Text, IntWritable> {
18      @Override
19      protected void reduce(Text key, Iterable<IntWritable> values, Reducer<Text,
IntWritable, Text, IntWritable>.Context context) throws IOException,
InterruptedException {
20          //定义累加的和
21          int count = 0;
22          for (IntWritable value : values) {
23              //累加v2中的每个集合中的值
24              count += value.get();
25          }
26          //context写出去
27          context.write(key, new IntWritable(count));
28      }
29  }
30

```

• main方法编写

```

1  package sz.base.java;
2
3  import org.apache.hadoop.conf.Configuration;
4  import org.apache.hadoop.fs.Path;
5  import org.apache.hadoop.io.IntWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Job;
8  import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
9  import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
10
11  import java.io.IOException;
12
13  public class WordCountMain {
14      public static void main(String[] args) {
15          try {
16              //1.实例化job ， 每个MapReduce都是一个job 创建一个job任务对象,并指定job的名字
17              Job job = Job.getInstance(new Configuration(), "wordcount");
18              //定义job的jar包参数需要设置一下
19              job.setJarByClass(WordCountMain.class);
20              //2.设置一下当前输入的格式（数据格式）和输入的文件路径
21              job.setInputFormatClass(TextInputFormat.class);

```

```

22      TextInputFormat.addInputPath(job, new Path("D:\\logs\\input\\words.txt"));
23      //3.设置mapper类和map输出的key和value的类型
24      job.setMapperClass(WordCountMapper.class);
25      job.setOutputKeyClass(Text.class);
26      job.setOutputValueClass(IntWritable.class);
27      //4.设置自定义的分区、排序、规约、分组规则，如果不写按照默认
28
29      //5.设置reducer类 和reduce输出的key 和value 的类型
30      job.setReducerClass(WordCountReduce.class);
31      job.setOutputKeyClass(Text.class);
32      job.setOutputValueClass(IntWritable.class);
33      //6.设置输出的格式，输出的路径
34      job.setOutputFormatClass(TextOutputFormat.class);
35      TextOutputFormat.setOutputPath(job, new Path("D:\\logs\\output"));
36      //7.等待这个任务的执行
37      boolean status = job.waitForCompletion(true);
38      //8.如果成功了退出当前job程序
39      System.exit(status ? 0 : 1);
40  } catch (IOException | InterruptedException | ClassNotFoundException e) {
41      e.printStackTrace();
42  }
43
44  }
45  }
46

```

- 集群运行则把路径修改为hdfs路径

```

1  TextInputFormat.addInputPath(job, new Path("hdfs://node1:8020/input/wordcount"));
2  TextOutputFormat.setOutputPath(job, new Path("hdfs://node1:8020/output/wordcount"));

```

- 将程序打成JAR包，然后在集群的任意一个节点上用hadoop命令启动

```

1  hadoop jar wordcount.jar cn.itcast.WordCountDriver

```

- 如果出现(null) entry in command string: null

```

1  解决方法：
2  下载hadoop.dll文件，拷贝到c:\windows\system32目录中即可hadoop.dll
3  可以在github上下载：https://github.com/4ttty/winutils
4  各个版本的hadoop.dll好像是通用的。

```

