

## 实验六 Spark 编程

### 6.1 实验目的

- 学习编写简单的基于 RDD API 的 Spark 程序
- 掌握在 IDEA 中调试 Spark 相关程序，以及在单机伪分布式、分布式部署方式下运行 Spark 相关程序的方法

### 6.2 实验任务

- 完成 WordCount 示例程序的编写
- 在单机伪分布式和分布式部署方式下运行 WordCount 示例程序

### 6.3 实验环境

- 操作系统: Ubuntu 18.04
- JDK 版本: jdk 1.8
- Hadoop 版本: 2.10.1
- Spark 版本: 2.4.7
- Scala 版本: 2.11.12
- IDEA 版本: 2020.2.3 (Ultimate 版)

### 6.4 实验步骤

#### 6.4.1 编写 Spark 应用程序

(1) 新建 Maven 项目并添加依赖

- 新建名为 “SparkDemo” 的 Maven 项目
- 编辑项目根目录下的 pom.xml 文件，分别在 <dependencies> 和 <build> 标签中添加 Spark 相关的插件 maven-compiler-plugin 和依赖包 spark-core

```
1 <dependencies>
2   <dependency>
3     <groupId>org.apache.spark</groupId>
4     <artifactId>spark-core_2.11</artifactId>
5     <version>2.4.7</version>
6   </dependency>
```

```
7 </dependencies>
```

修改完成后，在菜单界面选择 View->Tool Windows->Maven，在弹出的界面中点击 Reload All Maven Projects 加载依赖文件，第一次加载此过程可能耗时较长。

## (2) 编写代码

- Java 版: 新建 Java 类 *WordCount*

在项目的 *src/main/java* 目录下，选择 New -> Package，输入名称 *cn.edu.ecnu.spark.example.java.wordcount*，之后右键单击新建好的包，选择 New -> Java Class，输入名称 *WordCount*。

```
1 package cn.edu.ecnu.spark.example.java.wordcount;
2
3 import org.apache.spark.SparkConf;
4 import org.apache.spark.api.java.JavaPairRDD;
5 import org.apache.spark.api.java.JavaRDD;
6 import org.apache.spark.api.java.JavaSparkContext;
7 import org.apache.spark.api.java.function.*;
8 import scala.Tuple2;
9
10 import java.util.Arrays;
11 import java.util.Iterator;
12
13 public class WordCount {
14
15     public static void run(String[] args) {
16         /* 步骤1: 通过SparkConf设置配置信息，并创建SparkContext */
17         SparkConf conf = new SparkConf();
18         conf.setAppName("WordCountJava");
19         conf.setMaster("local"); // 仅用于本地进行调试，如在集群中运行则删除本行
20         JavaSparkContext sc = new JavaSparkContext(conf);
21
22         /* 步骤2: 按应用逻辑使用操作算子编写DAG，其中包括RDD的创建、转换和行动等 */
23         // 读入文本数据，创建名为lines的RDD
24         JavaRDD<String> lines = sc.textFile(args[0]);
25
26         // 将lines中的每一个文本行按空格分割成单个单词
27         JavaRDD<String> words =
28             lines.flatMap(
29                 new FlatMapFunction<String, String>() {
30                     @Override
31                     public Iterator<String> call(String line) throws Exception {
```



```
32         return Arrays.asList(line.split(" ")).iterator();
33     }
34 });
35 // 将每个单词的频数设置为1, 即将每个单词映射为[单词, 1]
36 JavaPairRDD<String, Integer> pairs =
37     words.mapToPair(
38         new PairFunction<String, String, Integer>() {
39             @Override
40             public Tuple2<String, Integer> call(String word) throws Exception
41             {
42                 return new Tuple2<String, Integer>(word, 1);
43             }
44         });
45 // 按单词聚合, 并对相同单词的频数使用sum进行累计
46 JavaPairRDD<String, Integer> wordCounts =
47     pairs
48     .groupByKey()
49     .mapToPair(
50         new PairFunction<Tuple2<String, Iterable<Integer>>, String,
51             Integer>() {
52             @Override
53             public Tuple2<String, Integer> call(Tuple2<String,
54                 Iterable<Integer>> t)
55                 throws Exception {
56                 Integer sum = Integer.valueOf(0);
57                 for (Integer i : t._2) {
58                     sum += i;
59                 }
60                 return new Tuple2<String, Integer>(t._1, sum);
61             }
62         });
63 // 合并机制
64 /*JavaPairRDD<String, Integer> wordCounts =
65     pairs.reduceByKey(
66         new Function2<Integer, Integer, Integer>() {
67             @Override
68             public Integer call(Integer t1, Integer t2) throws Exception {
69                 return t1 + t2;
70             }
71         });*/
72 // 输出词频统计结果到文件
```

```

71     wordCounts.saveAsTextFile(args[1]);
72     /* 步骤3: 关闭SparkContext */
73     sc.stop();
74 }
75
76 public static void main(String[] args) {
77     run(args);
78 }
79 }

```

• Scala 版: 新建 Scala 类 *WordCountScala*

在项目的 *src/main* 目录下, 选择 New -> Directory, 输入 *scala*, 接着右键单击目录 *scala*, 选择 Mark Directory as -> Sources Root。在 *scala* 目录下选择 New->Package, 输入名称 *cn.edu.ecnu.spark.example.scala.wordcount*, 之后右键单击新建好的包, 选择 New -> Scala Class, 输入名称 *WordCount*。

```

1 package cn.edu.ecnu.spark.example.scala.wordcount
2
3 import org.apache.spark.{SparkConf, SparkContext}
4
5 object WordCount {
6
7     def run(args: Array[String]): Unit = {
8         /* 步骤1: 通过SparkConf设置配置信息, 并创建SparkContext */
9         val conf = new SparkConf
10         conf.setAppName("WordCount")
11         conf.setMaster("local") // 仅用于本地进行调试, 如在集群中运行则删除本行
12         val sc = new SparkContext(conf)
13
14         /* 步骤2: 按应用逻辑使用操作算子编写DAG, 其中包括RDD的创建、转换和行动等 */
15         // 读入文本数据, 创建名为lines的RDD
16         val lines = sc.textFile(args(0))
17         // 将lines中的每一个文本行按空格分割成单个单词
18         val words = lines.flatMap { line => line.split(" ") }
19         // 将每个单词的频数设置为1, 即将每个单词映射为[单词, 1]
20         val pairs = words.map { word => (word, 1) }
21
22         // 按单词聚合, 并对相同单词的频数使用sum进行累计
23         val wordCounts = pairs.groupByKey().map(t => (t._1, t._2.sum))
24         // 如需使用合并机制则将上一行替换为下行
25         // val wordCounts = pairs.reduceByKey(_+_ )
26

```



```

27 // 输出词频统计结果到文件
28 wordCounts.saveAsTextFile(args(1))
29
30 /* 步骤3: 关闭SparkContext */
31 sc.stop()
32 }
33
34 def main(args: Array[String]): Unit = {
35     run(args)
36 }
37 }

```

## 6.4.2 调试 Spark 应用程序

使用 IDEA 调试 “WordCount” 应用程序。

### (1) 准备数据

```

1 cd ~
2 mkdir spark_input
3 cp ~/Chinese-Cloze-RC-master/people_daily/pd/pd.test ~/spark_input/pd.test
4 # 之前实验的数据集，如果删除了可以通过wget
   https://github.com/ymcui/Chinese-Cloze-RC/archive/master.zip下载

```

### (2) 配置运行环境

点击菜单栏 Run -> Edit Configuration, 在弹出的界面中点击 + 号选择 Application, 新建 Application 配置, Name 为 WordCountJava (Scala 版则为 WordCountScala), 配置界面如图6.1所示。

- 配置 Main Class 为 *cn.edu.ecnu.spark.example.java.wordcount.WordCount* (Scala 版则配置为 *cn.edu.ecnu.spark.example.scala.wordcount.WordCount*)
- 配置 Program arguments 为 输入路径 输出路径  
如 */home/dase-local/spark\_input /home/dase-local/IdeaProjects/SparkDemo/spark\_output/*, 其中前者表示输入路径, 后者表示输出路径。

### (3) 运行应用程序

在菜单栏点击 Run -> Run ‘WordCountJava’ (Scala 版为 Run ‘WordCountScala’), 运行结果如图6.2所示, 程序运行结束在 IDEA 中产生的文件如图6.3所示。

注意: 运行前删除目录 */home/dase-local/IdeaProjects/SparkDemo/spark\_output/*。

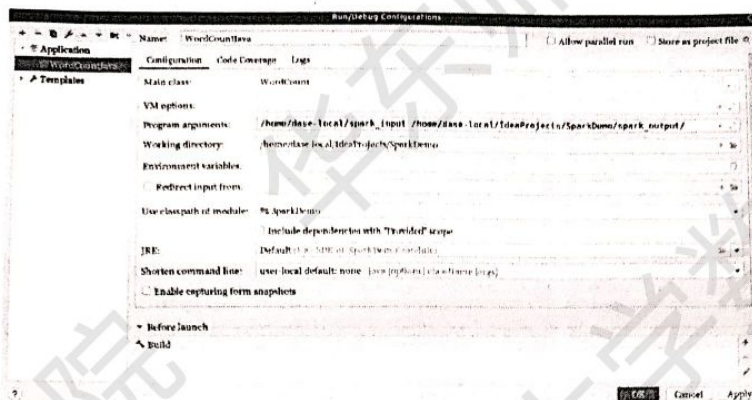


图 6.1 运行应用程序配置界面

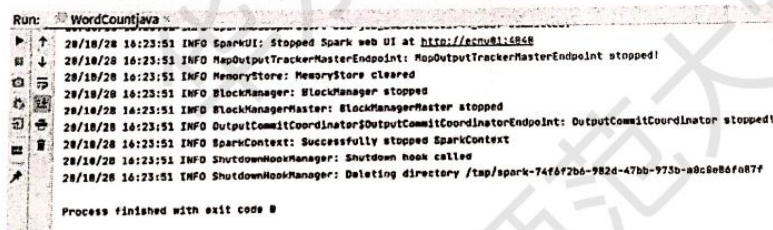


图 6.2 在 IDEA 中的运行结果

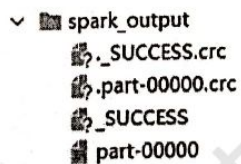


图 6.3 程序运行产生的文件



### 6.4.3 运行 Spark 应用程序

以 Java 版本为例，展示如何运行 Spark 应用程序。

#### (1) 准备工作

以下操作在各节点均以 dase-local 用户身份进行。

- 使用 IDEA 将程序打成 jar 包

注意：打包前注释掉代码中的 `conf.setMaster("local");`。

jar 包名称为 `RddWordCountJava.jar`，打包路径为 `/home/dase-local/IdeaProjects/SparkDemo/out/artifacts/RddWordCountJava/`，配置界面如图 6.4 所示。

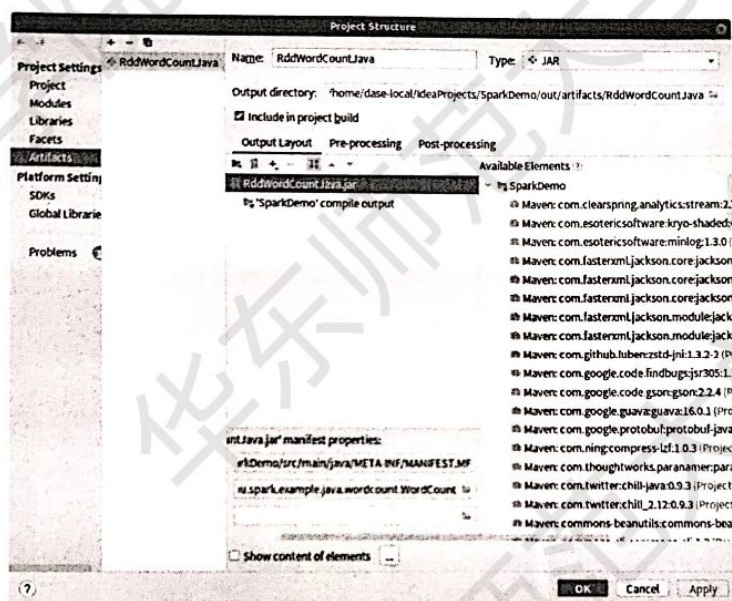


图 6.4 jar 包配置界面

- 复制 jar 包到指定目录

```
1 mkdir ~/spark-2.4.7/myApp
2 cp ~/IdeaProjects/SparkDemo/out/artifacts/RddWordCountJava/RddWordCountJava.jar
   ar ~/spark-2.4.7/myApp/
```

#### (2) 单机伪分布式部署方式下运行应用程序

以下操作在各节点均以 dase-local 用户身份进行。

- 启动服务

```
1 su dase-local
2 cd ~
3 ~/hadoop-2.10.1/sbin/start-dfs.sh
4 ~/spark-2.4.7/sbin/start-all.sh
5 ~/spark-2.4.7/sbin/start-history-server.sh
```

### • 准备数据

```
1 ~/hadoop-2.10.1/bin/hdfs dfs -cp ./input/pd.test ./spark_input
2 #该命令将hdfs://user/dase-local/input/pd.test拷贝到hdfs://user/dase-local/spar
3 k_input中
4 #其中, pd.test数据文件之前已经上传至hdfs://user/dase-local/input
5 #spark_input目录之前也已创建完成
```

### • 通过提交 jar 包运行应用程序

```
1 ~/hadoop-2.10.1/bin/hdfs dfs -rm -r spark_output #删除 HDFS 上的输出路径
2 ~/spark-2.4.7/bin/spark-submit \
3 --master spark://localhost:7077 \
4 --class cn.edu.ecnu.spark.example.java.wordcount.WordCount \
5 /home/dase-local/spark-2.4.7/myApp/RddWordCountJava.jar
   hdfs://localhost:9000/user/dase-local/spark_input
   hdfs://localhost:9000/user/dase-local/spark_output
   #将之前实验中上传到HDFS上的文件作为输入
```

```
21/02/28 15:00:17 INFO handler.ContextHandler: Stopped o.s.j.s.ServletContextHan
dler@62fe067{/jobs/job,null,UNAVAILABLE,@Spark}
21/02/28 15:00:17 INFO handler.ContextHandler: Stopped o.s.j.s.ServletContextHan
dler@3e34ace1{/jobs/json,null,UNAVAILABLE,@Spark}
21/02/28 15:00:17 INFO handler.ContextHandler: Stopped o.s.j.s.ServletContextHan
dler@6e1d8f9e{/jobs,null,UNAVAILABLE,@Spark}
21/02/28 15:00:17 INFO ui.SparkUI: Stopped Spark web UI at http://219.228.148.15
4:4040
21/02/28 15:00:17 INFO cluster.StandaloneSchedulerBackend: Shutting down all exe
cutors
21/02/28 15:00:18 INFO cluster.CoarseGrainedSchedulerBackend$DriverEndpoint: Ask
ing each executor to shut down
21/02/28 15:00:18 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMas
terEndpoint stopped!
21/02/28 15:00:18 INFO memory.MemoryStore: MemoryStore cleared
21/02/28 15:00:18 INFO storage.BlockManager: BlockManager stopped
21/02/28 15:00:18 INFO storage.BlockManagerMaster: BlockManagerMaster stopped
21/02/28 15:00:18 INFO scheduler.OutputCommitCoordinator$OutputCommitCoordinator
Endpoint: OutputCommitCoordinator stopped!
21/02/28 15:00:18 INFO spark.SparkContext: Successfully stopped SparkContext
21/02/28 15:00:18 INFO util.ShutdownHookManager: Shutdown hook called
21/02/28 15:00:18 INFO util.ShutdownHookManager: Deleting directory /tmp/spark-a
427bbda-9833-421c-8b34-52e2f6e4e1c2
```

图 6.5 在伪分布式部署方式下应用程序的运行结果

### Browse Directory

hdfs://user/dase-local/spark\_output

Show 25 entries

ll	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	drwxr-xr-x	dase-local	supergroup	0 B	Apr 15 22:01	1	128 MB	SUCCESS
drwxr-xr-x	drwxr-xr-x	dase-local	supergroup	5120 KB	Apr 15 22:01	1	128 MB	part-00000
drwxr-xr-x	drwxr-xr-x	dase-local	supergroup	5123 KB	Apr 15 22:01	1	128 MB	part-00001

Showing 1 to 3 of 3 entries

Previous Next

图 6.6 程序运行产生的文件



执行结果如图6.5所示，程序运行产生的文件如图6.6所示。

- 停止服务

```
1 ~/hadoop-2.10.1/sbin/stop-dfs.sh
2 ~/spark-2.4.7/sbin/stop-all.sh
3 ~/spark-2.4.7/sbin/stop-history-server.sh
```

### (3) 分布式部署方式下运行应用程序

以下操作在各节点均以 dase-dis 用户身份进行。

- 复制 jar 包到指定目录

在客户端节点 (ecnu04) 执行以下命令，将之前打好的包拷贝至 dase-dis 用户下的 ~/spark-2.4.7/myApp 路径中。

```
1 su dase-dis
2 mkdir ~/spark-2.4.7/myApp/ #在hadoop目录下新建myApp/目录
3 scp dase-local@localhost:/home/dase-local/spark-2.4.7/myApp/
   /home/dase-dis/hadoop-2.10.1/myApp/
```

- 启动服务

在主节点执行以下命令：

```
1 su dase-dis
2 ~/hadoop-2.10.1/sbin/start-dfs.sh
3 ~/spark-2.4.7/sbin/start-all.sh
4 ~/spark-2.4.7/sbin/start-history-server.sh
```

- 准备数据

在客户端节点 (ecnu04) 执行以下命令：

```
1 ~/hadoop-2.10.1/bin/hdfs dfs -cp ./input/pd.test ./spark_input
2 #pd.test数据文件之前已经上传至hdfs:///user/dase-dis/input
3 #spark_input目录之前也已创建完成
```

- 通过提交 jar 包运行应用程序

在客户端节点 (ecnu04) 执行以下命令：

```
1 ~/hadoop-2.10.1/bin/hdfs dfs -rm -r spark_output
2 ~/spark-2.4.7/bin/spark-submit \
3   --master spark://ecnu01:7077 \
4   --class cn.edu.ecnu.spark.example.java.wordcount.WordCount \
5   /home/dase-dis/spark-2.4.7/myApp/RddWordCountJava.jar
   hdfs://ecnu01:9000/user/dase-dis/spark_input
   hdfs://ecnu01:9000/user/dase-dis/spark_output
```

分布式环境下程序运行产生的文件如图6.7所示。

Browse Directory

Ames/Ames-dfs/spark\_output

Show as entries

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	data-dis	supergroup	0 B	Apr 15 22:42	2	128 MB	SUCCESS
-rw-r--r--	data-dis	supergroup	54.02 KB	Apr 15 22:42	2	128 MB	part-00000
-rw-r--r--	data-dis	supergroup	51.23 KB	Apr 15 22:42	2	128 MB	part-00001

Showing 1 to 3 of 3 entries

Search

Previous Next

图 6.7 程序运行产生的文件

- 停止服务

在主节点执行以下命令：

```
1 ~/hadoop-2.10.1/sbin/stop-dfs.sh
2 ~/spark-2.4.7/sbin/stop-all.sh
3 ~/spark-2.4.7/sbin/stop-history-server.sh
```

## 6.5 思考题

- 1 第6.4.3节中编写的 WordCount 应用程序由多少个 Stage 组成？请结合 Spark 的 Web UI 进行说明。
- 2 对于第6.4.3节中编写的 WordCount 应用程序，根据实际运行情况，结合 Spark 的 Web UI 说明运行过程中 Shuffle 的数据量有多大。