

实验四 MapReduce 2.x 编程

4.1 实验目的

- 学习编写简单的基于 Java API 的 MapReduce 程序
- 学习在单机集中式、单机伪分布、分布式部署方式下运行 MapReduce 程序

4.2 实验任务

- 完成基于 Java API 编写的“WordCount”的 MapReduce 程序
- 在单机集中式、单机伪分布式、分布式部署方式下调试、提交并运行该应用程序

4.3 实验环境

- 操作系统: Ubuntu 18.04
- JDK 版本: 1.8
- Hadoop 版本: 2.10.1
- IDEA 版本: 2020.2.3 (Ultimate 版)

4.4 实验步骤

4.4.1 编写 MapReduce 应用程序

(1) 新建 Maven 项目并添加依赖

- 新建名为“MapReduceDemo”的 Maven 项目
- 编辑项目根目录下的 pom.xml 文件, 在 <dependencies> 标签中添加 Hadoop 相关的依赖包 hadoop-common, hadoop-client, hadoop-hdfs

```
1 <dependencies>
2   <!-- https://mvnrepository.com/artifact/org.apache.hadoop/hadoop-common
3       -->
4   <dependency>
5     <groupId>org.apache.hadoop</groupId>
6     <artifactId>hadoop-common</artifactId>
7     <version>2.10.1</version>
8   </dependency>
   <!-- https://mvnrepository.com/artifact/org.apache.hadoop/hadoop-client
```

```

-->
9  <dependency>
10     <groupId>org.apache.hadoop</groupId>
11     <artifactId>hadoop-client</artifactId>
12     <version>2.10.1</version>
13 </dependency>
14 <!-- https://mvnrepository.com/artifact/org.apache.hadoop/hadoop-hdfs -->
15 <dependency>
16     <groupId>org.apache.hadoop</groupId>
17     <artifactId>hadoop-hdfs</artifactId>
18     <version>2.10.1</version>
19 </dependency>
20 </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 Java 程序

编写一个 Java 应用程序，实现对给定文本中的单词进行计数。

- 在项目的 `src/main/java` 目录下，选择 New->Package，输入名称 `cn.edu.ecnu.mapreduce.example.java.wordcount`
- 右键单击建好的包，选择 New->Java Class，输入名称 `WordCountMapper`

```

1 package cn.edu.ecnu.mapreduce.example.java.wordcount;
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 /* 步骤1: 确定输入键值对 [K1,V1] 的数据类型为 [LongWritable,Text]，输出键值对
11    [K2,V2] 的数据类型为 [Text,IntWritable] */
12 public class WordCountMapper extends Mapper<LongWritable, Text, Text,
13     IntWritable> {
14
15     @Override
16     protected void map(LongWritable key, Text value, Context context)
17         throws IOException, InterruptedException {
18         /* 步骤2: 编写处理逻辑将 [K1,V1] 转换为 [K2,V2] 并输出 */
19         // 以空格作为分隔符拆分成单词

```



```

18     String[] datas = value.toString().split(" ");
19     for (String data : datas) {
20         // 输出分词结果
21         context.write(new Text(data), new IntWritable(1));
22     }
23 }
24 }

```

- 右键单击建好的包，选择 New->Java Class，输入名称 *WordCountReducer*

```

1 package cn.edu.ecnu.mapreduce.example.java.wordcount;
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 /* 步骤1: 确定输入键值对[K2,List(V2)]的数据类型为[Text,
10    IntWritable], 输出键值对[K3,V3]的数据类型为[Text,IntWritable] */
11 public class WordCountReducer extends Reducer<Text, IntWritable, Text,
12    IntWritable> {
13     @Override
14     protected void reduce(Text key, Iterable<IntWritable> values, Context
15        context)
16        throws IOException, InterruptedException {
17         /* 步骤2: 编写处理逻辑将[K2,List(V2)]转换为[K3,V3]并输出 */
18         int sum = 0;
19         // 遍历累加求和
20         for (IntWritable value : values) {
21             sum += value.get();
22         }
23         // 输出计数结果
24         context.write(key, new IntWritable(sum));
25     }
26 }

```

- 右键单击建好的包，选择 New->Java Class，输入名称 *WordCount*

```

1 package cn.edu.ecnu.mapreduce.example.java.wordcount;
2
3 import org.apache.hadoop.conf.Configured;
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.FileInputFormat;
9 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
10 import org.apache.hadoop.util.Tool;
11 import org.apache.hadoop.util.ToolRunner;
12
13 public class WordCount extends Configured implements Tool {
14
15     @Override
16     public int run(String[] args) throws Exception {
17         /* 步骤1: 设置作业的信息 */
18         Job job = Job.getInstance(getConf(), getClass().getSimpleName());
19         // 设置程序的类名
20         job.setJarByClass(getClass());
21
22         // 设置数据的输入输出路径
23         FileInputFormat.addInputPath(job, new Path(args[0]));
24         FileOutputFormat.setOutputPath(job, new Path(args[1]));
25
26         // 设置map和reduce方法
27         job.setMapperClass(WordCountMapper.class);
28         job.setReducerClass(WordCountReducer.class);
29
30         // 设置map方法的输出键值对数据类型
31         job.setMapOutputKeyClass(Text.class);
32         job.setMapOutputValueClass(IntWritable.class);
33         // 设置reduce方法的输出键值对数据类型
34         job.setOutputKeyClass(Text.class);
35         job.setOutputValueClass(IntWritable.class);
36
37         return job.waitForCompletion(true) ? 0 : 1;
38     }
39
40     public static void main(String[] args) throws Exception {
41         /* 步骤2: 运行作业 */
42         int exitCode = ToolRunner.run(new WordCount(), args);
43         System.exit(exitCode);
44     }
45 }
```


4.4.2 调试 MapReduce 应用程序

使用 IDEA 调试 WordCount 应用程序¹。

(1) 准备数据

使用之前准备的数据集 pd.train，数据集的位置为 `/home/dase-local/input/`。

(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 WordCount, 配置界面如图 4.1 所示。

a) 配置 Main Class 为 `cn.edu.ecnu.mapreduce.example.java.wordcount.WordCount`

b) 配置 Program arguments 为 `/home/dase-local/input/ /home/dase-local/IdeaProjects/MapReduceDemo/output`

该程序需要传入两个参数，第一个参数为输入文件路径，即要统计单词数的文本文件或文件夹的路径，若输入的是文件夹的路径，输出的是文件夹中所有文本文件的单词数统计结果；第二个参数为输出文件路径，即存放统计结果的文件路径。输出文件路径无需在程序运行前创建好，程序运行过程中会生成该输出文件路径。

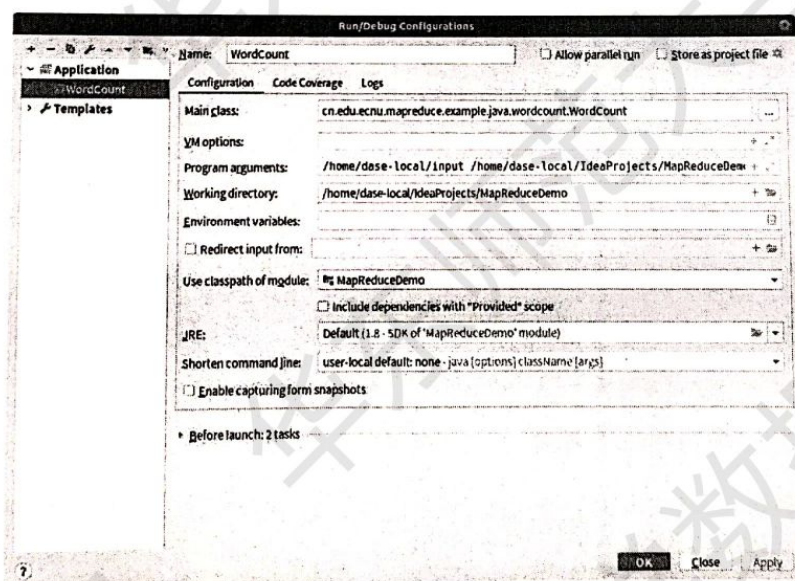


图 4.1 运行参数配置界面

(3) 运行应用程序

在菜单栏点击 Run->Run 'WordCount', 在 IDEA 中直接运行程序，待运行完毕后，项目根目录下会出现“output”文件夹，其中“part-r-00000”文件就是运行结果，如图 4.2 所示。

¹若在 Windows 下调试 MapReduce 应用程序，则需按照附录 A 先进行相关配置

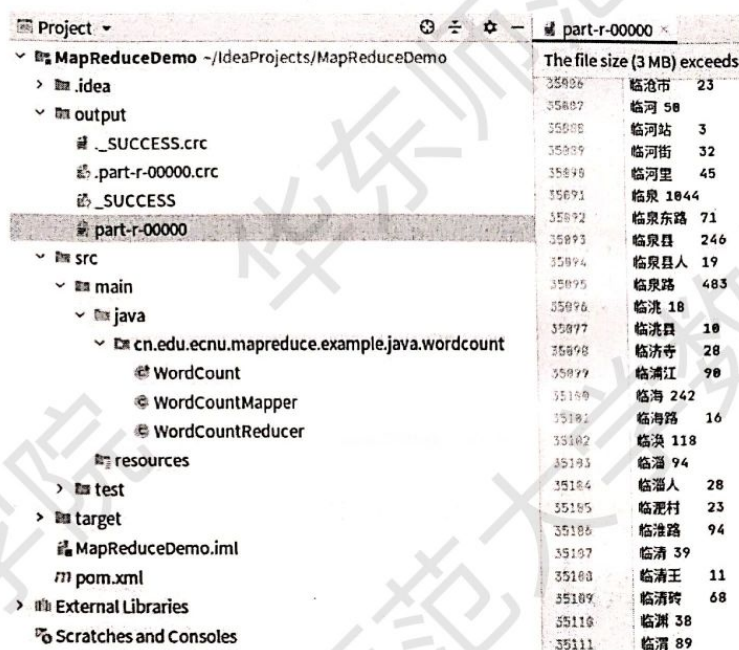


图 4.2 IDEA 中应用程序调试结果

4.4.3 运行 MapReduce 应用程序

在部署的 Hadoop 上运行 WordCount 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 将项目打包成可执行 jar 包

jar 包名称为 WordCount.jar, 打包路径为 `/home/dase-local/IdeaProjects/HDFSFileOP/out/artifacts/WordCount/`, 配置界面如图 4.3 所示。

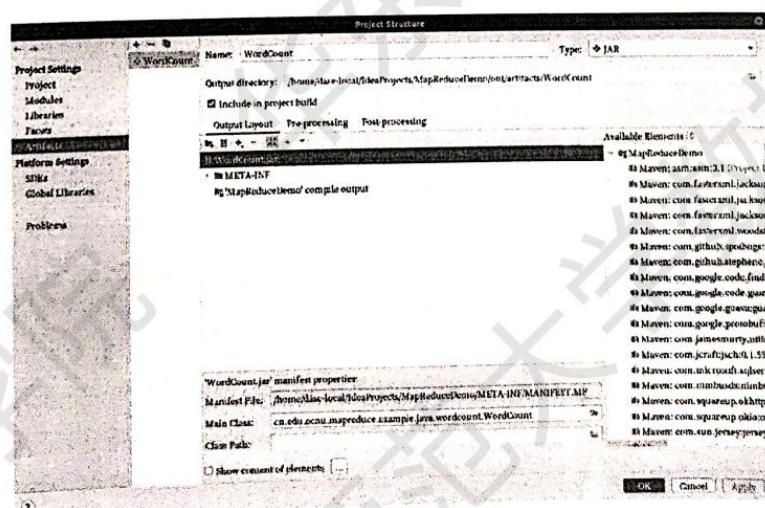


图 4.3 将项目打包为可执行 jar 包

- 复制 jar 包到指定目录

在客户端节点执行以下命令，将之前打好的包拷贝至 `/home/dase-local/hadoop-2.10.1/myApp` 路径下。

```
1 mkdir ~/hadoop-2.10.1/myApp/ #在hadoop目录下新建myApp/目录
2 cp ~/IdeaProjects/MapReduceDemo/out/artifacts/WordCount/WordCount.jar
   ~/hadoop-2.10.1/myApp/
   #将打包好的可执行jar包拷贝到hadoop安装路径下的myApp/目录下
```

(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 Yarn 和 HDFS 服务

```
1 su dase-local
2 ~/hadoop-2.10.1/sbin/start-yarn.sh
3 ~/hadoop-2.10.1/sbin/mr-jobhistory-daemon.sh start historyserver
4 ~/hadoop-2.10.1/sbin/start-dfs.sh
```

- 上传输入文件

输入文件 `pd.train` 在上一章中已上传至 `hdfs://localhost:9000/user/dase-local/input/`。

- 通过提交 jar 包运行应用程序

```
1 cd hadoop-2.10.1/
2 ./bin/hadoop jar ./myApp/WordCount.jar input/pd.train output
   #使用hadoop命令在Yarn模式下运行jar包，并在运行时输入参数
```

该程序实现了统计 HDFS “`/user/dase-local/input/`” 目录下的 “`README.txt`” 文件单词数量的操作。运行结果存储在 HDFS 的 “`/user/dase-local/output/`” 目录下的 “`part-r-00000`” 文件中，如图4.4所示。

- 查看运行结果

```
1 cd hadoop-2.10.1/
2 ./bin/hdfs dfs -tail output/part-r-00000 #查看输出结果
```

通过命令行查看的输出结果如图 4.5所示。

- 停止相关服务

```
1 ~/hadoop-2.10.1/sbin/stop-yarn.sh
2 ~/hadoop-2.10.1/sbin/mr-jobhistory-daemon.sh stop historyserver
3 ~/hadoop-2.10.1/sbin/stop-dfs.sh
```

(3) 分布式部署方式下运行应用程序

Browse Directory

/user/dase-local/output										Go			
Show	25	entries	Search										
<input type="checkbox"/>	dl	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name				
<input type="checkbox"/>		-rw-r--r--	dase-local	supergroup	0 B	Apr 26 14:38	1	128 MB	_SUCCESS				
<input type="checkbox"/>		-rw-r--r--	dase-local	supergroup	2.86 MB	Apr 26 14:38	1	128 MB	part-r-00009				
Showing 1 to 2 of 2 entries										Previous			

图 4.4 单机伪分布式部署下应用程序运行结果存储位置

```

M B A      2
M H        26
M K 2      68
M P 3      39
M P V      12
N          390
N B A      30
P 订制     21
P D A      38
P M        640
P M F      10
P P        20
P S A      23
P o w e r  10

```

图 4.5 单机伪分布式部署下应用程序运行结果

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

在客户端节点执行以下命令，将之前打好的包拷贝至 dase-dis 用户/home/dase-dis/hadoop-2.10.1/myApp 目录下。

```

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

```

- 启动 Yarn 和 HDFS 服务

在主节点执行以下命令：

```

1 su dase-dis
2 ~/hadoop-2.10.1/sbin/start-yarn.sh
3 ~/hadoop-2.10.1/sbin/mr-jobhistory-daemon.sh start historyserver

```



```
4 ~/hadoop-2.10.1/sbin/start-dfs.sh
```

- 上传输入文件

输入文件 pd.train 在上一章中已上传至 `hdfs://ecnu01:9000/user/dase-dis/input/`。

- 通过提交 jar 包运行应用程序

在客户端节点执行以下命令：

```
1 cd ~/hadoop-2.10.1/  
2 ./bin/hadoop jar ./myApp/WordCount.jar input/pd.train output  
3 ./bin/hdfs dfs -tail output/part-r-00000 #查看输出结果
```

该程序实现了统计 HDFS “/user/dase-dis/input/” 目录下的 “README.txt” 文件单词数量的操作。运行结果存储在 HDFS 的 “/user/dase-dis/output/” 目录下的 “part-r-00000” 文件中。

- 停止服务

在主节点执行以下命令：

```
1 ~/hadoop-2.10.1/sbin/stop-yarn.sh  
2 ~/hadoop-2.10.1/sbin/mr-jobhistory-daemon.sh stop historyserver  
3 ~/hadoop-2.10.1/sbin/stop-dfs.sh
```

4.5 思考题

1 如何从 Hadoop JobHistory 的 Web UI 中查看某个 MapReduce 应用程序分别启动了多少个 Map 任务和 Reduce 任务？

2 如何从 Hadoop JobHistory 的 Web UI 中查看某个 MapReduce 应用程序启动的 Map 任务和 Reduce 任务分别是在哪些 NodeManager 上执行的？

3 对于一个 MapReduce 应用程序，是否存在某一时刻 Map 任务和 Reduce 任务同时运行？请结合 Hadoop JobHistory 的 Web UI 给出例证。

4 编写一个利用 Combine 机制实现的 WordCount 程序，使用实验二中的 pd.train 数据集，在分布式部署方式下分别运行本实验中的 WordCount 程序和使用 Combine 机制的 WordCount 程序，观察二者的时间差并分析原因。

附录

A Windows 下调试 MapReduce 应用程序

(1) 下载 winutils

在 Windows 下, Hadoop 需要一些 native 库来访问本地文件系统, 否则无法正确运行 MapReduce 程序。因此, 需要下载所需的 native 库, 即 winutils²。将下载完成的 winutils 进行解压, 并将解压文件夹下的 hadoop.dll、winutils.exe 等文件放置在 C:\Hadoop\bin 目录下。

(2) 配置环境变量

- 在“此电脑”图标上右键选择属性, 然后选择“高级系统设置”, 如图 4.6 所示。

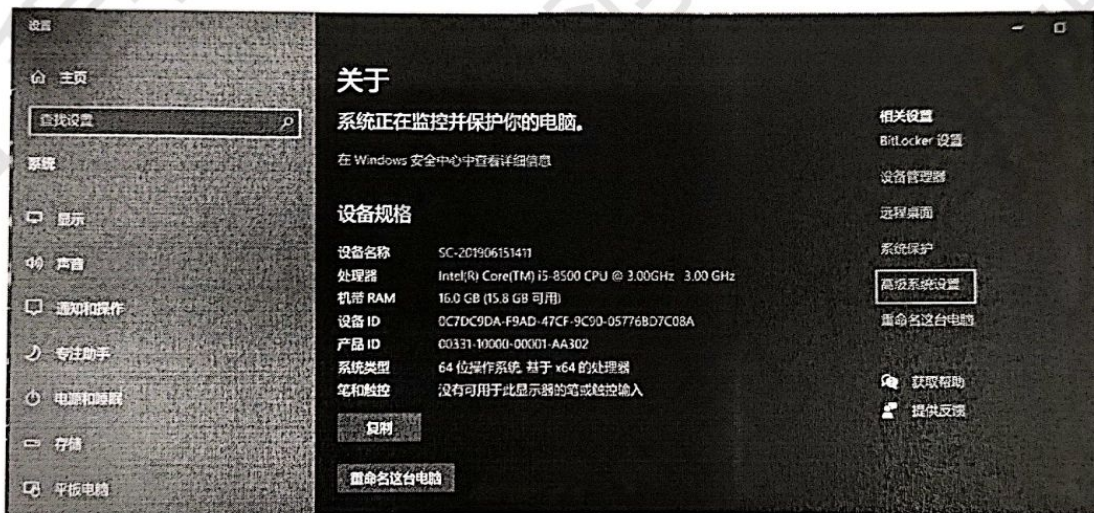


图 4.6 选择高级系统设置

- 在弹出的界面中选择“高级”并点击该界面的“环境变量”, 如图 4.7 所示。
- 点击系统变量下的“新建”按钮, 并添加名为“HADOOP_HOME”的环境变量, 如图 4.8 所示。
- 点击系统变量中名为“PATH”的环境变量, 并在其中添加 winutils 相关的路径, 如图 4.9 所示。

(3) 编程注意事项

- 环境变量配置完成之后需要重启 IDE。
- 请确保 winutils 版本与程序引入的 hadoop 依赖版本保持一致。

²https://github.com/stevcloughran/winutils/releases/download/tag_2017-08-29-hadoop-2.8.1-native/hadoop-2.8.1.zip

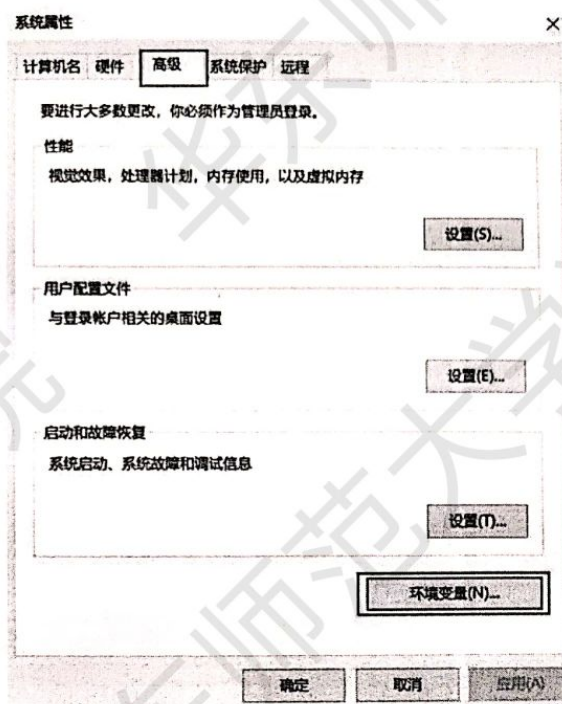


图 4.7 更改环境变量配置

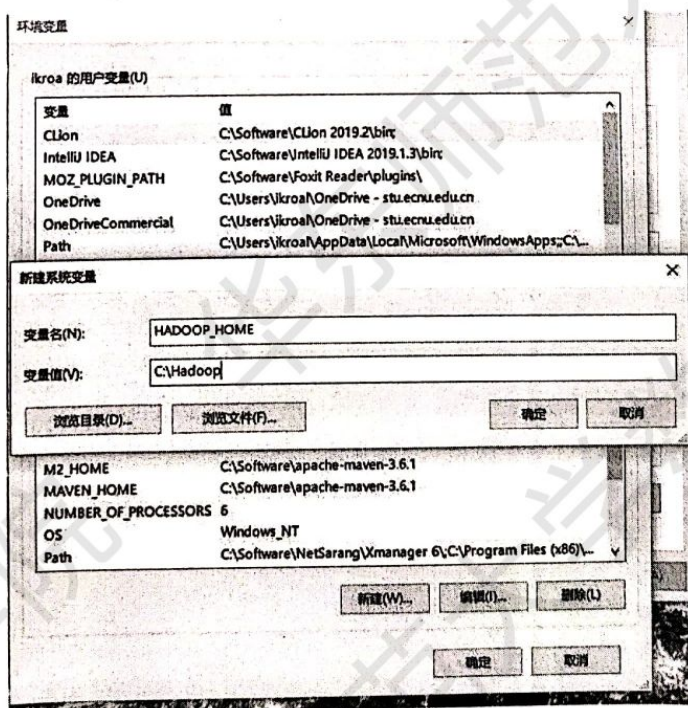


图 4.8 添加名为 HADOOP_HOME 的环境变量

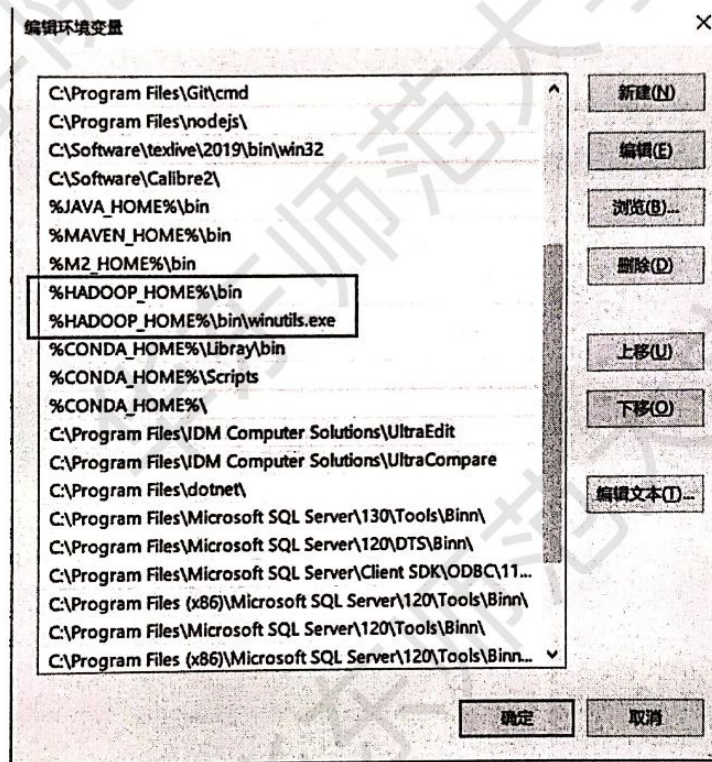


图 4.9 在 PATH 中添加 winutils 相关路径