

1. 作业5说明

211275009 陈铭浩

1.1. maven创建项目后，更改pom.xml相关配置

```
1  #使用maven-archetype-quickstart(version=1.4)
2  1. 修改编译使用的jdk版本（本机为JDK1.8）
3      <properties>
4          <project.build.sourceEncoding>UTF-
5              8</project.build.sourceEncoding>
6          <maven.compiler.source>1.8</maven.compiler.source>
7          <maven.compiler.target>1.8</maven.compiler.target>
8      </properties>
9  2. 添加项目用到的Java依赖包
10     <dependencies>
11         <dependency>
12             <groupId>junit</groupId>
13             <artifactId>junit</artifactId>
14             <version>4.11</version>
15             <scope>test</scope>
16         </dependency>
17         <dependency>
18             <groupId>org.apache.hadoop</groupId>
19             <artifactId>hadoop-client</artifactId>
20             <version>3.3.6</version>
21         </dependency>
22         <dependency>
23             <groupId>org.apache.hadoop</groupId>
24             <artifactId>hadoop-common</artifactId>
25             <version>3.3.6</version>
26         </dependency>
27         <dependency>
28             <groupId>org.apache.hadoop</groupId>
```

```
28     <artifactId>hadoop-hdfs</artifactId>
29     <version>3.3.6</version>
30 </dependency>
31 <dependency>
32     <groupId>org.apache.hadoop</groupId>
33     <artifactId>hadoop-mapreduce-client-core</artifactId>
34     <version>3.3.6</version>
35 </dependency>
36 <dependency>
37     <groupId>org.apache.hadoop</groupId>
38     <artifactId>hadoop-mapreduce-client-common</artifactId>
39     <version>3.3.6</version>
40 </dependency>
41 </dependencies>
```

1.2. 程序设计思路

1. 实验要求

对于给定的两个输入文件A和B（如附件所示），文件内包含了纳斯达克100指数、标普500指

数、道琼斯工业指数的部分成分股数据，第一列为三大指数编号（简化为101、102、103），第

二列为成分股代码。请编写Mapreduce程序，对两个文件内容进行合并，并对重复内容进行剔

除，最后合并为一个包含三大指数成分股的输出文件。输出文件格式与文件A和B格式相同。

2. 设计思路

首先将xlsx文件转为csv文件便于后续处理，而后在UnionCalc主类中包括UnionCalcMapper 类，UnionCalcReducer 类和main方法

- UnionCalcMapper 类：

这是一个Mapper类，用于处理输入数据。

在 map 方法中，将每一行文本数据（CSV记录）转换为字符串，然后使用逗号分隔符拆分记录，提取索引和股票代码。

为了去除不必要的空格，使用 .trim() 方法来修整索引和股票代码。

最后，将索引作为键，股票代码作为值，发送到Reducer以分组相同索引的记录。

- UnionCalcReducer 类:

这是一个Reducer类，用于处理来自Mapper的键值对数据，去除重复记录。

在 reduce 方法中，创建了一个 Set 集合（stockCodeSet）来存储唯一的股票代码。

然后，遍历相同索引下的所有股票代码，将它们添加到 stockCodeSet 中，从而去除重复。

最后，将去除重复后的记录输出。

- main 方法:

main 方法是程序的入口点，用于配置Hadoop作业。

通过 Job 类，指定了Map和Reduce的类，以及输入和输出数据的格式。

设置了输入文件路径和输出文件路径。

最后，启动Hadoop作业，等待作业完成，然后退出。

```
public class UnionCalc{
    public static class UnionCalcMapper extends Mapper<Object, Text, Text, Text>{
        protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
            String line = value.toString();

            String[] fields = line.split(",");

            if (fields.length >= 2) {
                String index = fields[0].trim();
                String stockCode = fields[1].trim();
                context.write(new Text(index), new Text(stockCode));
            }
        }
    }

    public static class UnionCalcReducer extends Reducer<Text, Text, Text, Text>{
        protected void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException {
            Set<String> stockCodeSet = new HashSet<>();

            // 遍历某个index下所有相同成分股代码的记录，将其添加到集合中
            for (Text value : values) {
                stockCodeSet.add(value.toString());
            }

            // 输出去除重复后的记录
            for (String stockCode : stockCodeSet) {
                context.write(key, new Text(stockCode));
            }
        }
    }
}
```

```
public static void main(String[] args) throws Exception{
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "UnionCalc");

    job.setJarByClass(UnionCalc.class);
    job.setMapperClass(UnionCalcMapper.class);
    job.setReducerClass(UnionCalcReducer.class);

    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(Text.class);

    job.setInputFormatClass(org.apache.hadoop.mapreduce.lib.input.TextInputFormat.class);
    job.setOutputFormatClass(org.apache.hadoop.mapreduce.lib.output.TextOutputFormat.class);

    FileInputFormat.addInputPath(job, new Path(args[0])); // 输入文件A和B的路径input
    FileOutputFormat.setOutputPath(job, new Path(args[1])); // 输出文件路径output

    System.exit(job.waitForCompletion(verbose:true) ? 0 : 1);
}
```

1.3. 打包maven项目

```
1 1. 运行 Maven 编译和打包命令：
2 #在终端中，使用以下 Maven 命令来编译和打包你的项目：
3 mvn clean package
4 #这个命令会执行 clean 生命周期（清理旧构建），然后执行 package 生命周期，将
  项目打包为 JAR 文件。构建好的 JAR 文件将位于项目目录的 target 子目录下。
5 2. 查找 JAR 文件：
6 #构建成功后，你可以在项目目录的 target 文件夹中找到生成的 JAR 文件，通常以项
  目名称和版本号命名，如 my-java-project-1.0-SNAPSHOT.jar。
```

1.4. 在hadoop中运行程序

1. 准备伪分布式hadoop

```
1 start-all.sh
```

2. 将jar文件和输入文件上传到Hadoop中

3. 使用 hadoop jar 命令来运行 JAR 文件

```
1 hadoop jar your-project.jar main-class input-path output-path
2 #your-project.jar 是你上传到 Hadoop 的 JAR 文件
3 #main-class 是包含 main 方法的 Java 类的全名，用于启动你的作业
4 #input-path 是输入数据的 HDFS 路径
5 #output-path 是输出结果的 HDFS 路径
```

1.5. 运行结果

```
2023-10-24 21:06:17,622 INFO mapreduce.Job: map 0% reduce 0%
2023-10-24 21:06:21,664 INFO mapreduce.Job: map 100% reduce 0%
2023-10-24 21:06:25,683 INFO mapreduce.Job: map 100% reduce 100%
2023-10-24 21:06:25,690 INFO mapreduce.Job: Job job_1698151280234_0002 completed successfully
2023-10-24 21:06:25,749 INFO mapreduce.Job: Counters: 54
```



Cluster

About

Nodes

Node Labels

Applications

NEW

NEW SAVING

SUBMITTED

ACCEPTED

RUNNING

FINISHED

FAILED

KILLED

Scheduler

Tools

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Used Resources
2	0	0	2	0	<memory:0 B, vCores:0>

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes
1	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation
Capacity Scheduler	[memory-mb (unit=M), vcores]	<memory:1024, vCores:1>	<memory:8192, vCores:1>

Show 20 entries

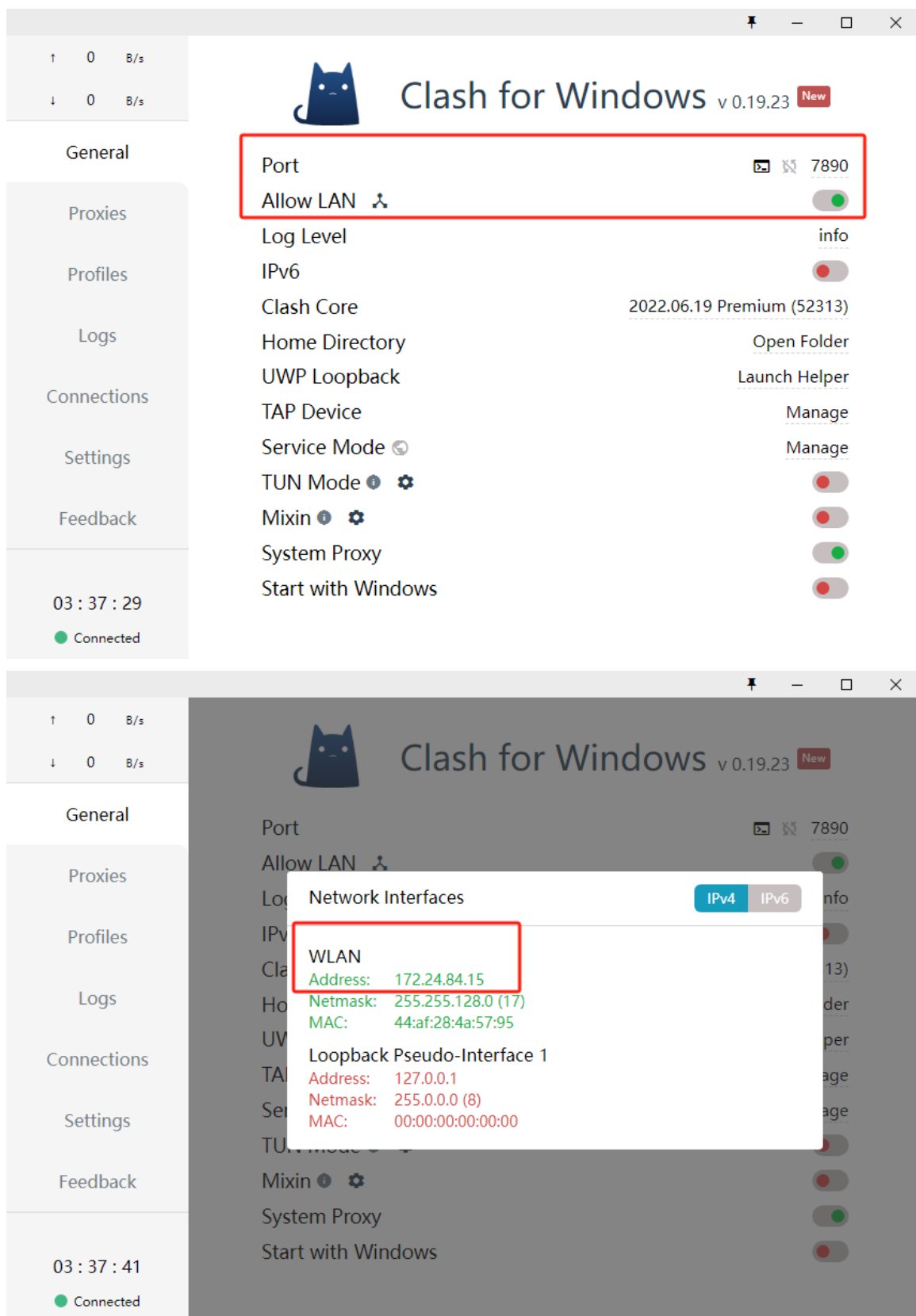
ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Resource
application_1698151280234_0002	cmh	UnionCalc	MAPREDUCE		default	0	Tue Oct 24 21:06:12 +0800 2023	Tue Oct 24 21:06:13 +0800 2023	Tue Oct 24 21:06:24 +0800 2023	FINISHED	SUCCEEDED	N/A

输出文件（只展示了前37行）

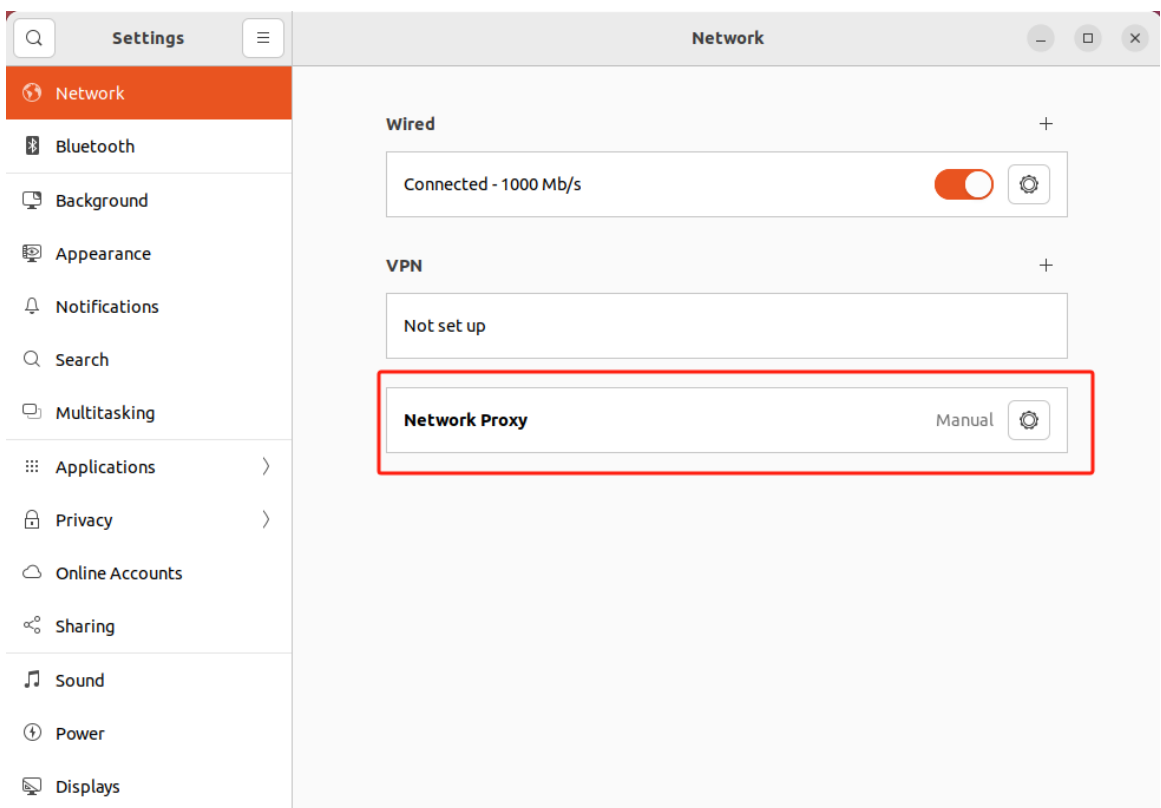
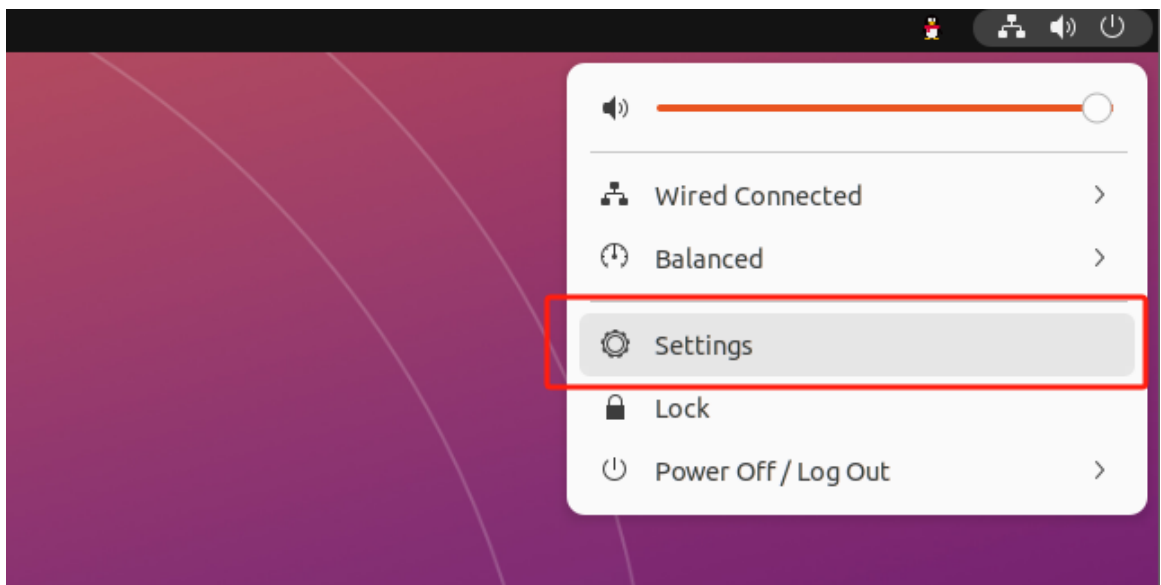
Open		part-r-00000		Save				
		~/hadoop/hadoop-3.3.6/homework5/output						
1	101	AAPL						
2	101	CSCO						
3	101	KO						
4	101	JPM						
5	101	HON						
6	101	NKE						
7	101	GS						
8	101	MMM						
9	101	AMGN						
10	101	DOW						
11	101	UNH						
12	101	DIS						
13	101	MSFT						
14	101	V						
15	101	INTC						
16	101	CAT						
17	101	AXP						
18	101	HD						
19	101	CRM						
20	101	BA						
21	102	AHG						
22	102	AAPL						
23	102	CSCO						
24	102	ACXP						
25	102	TT00						
26	102	GRTS						
27	102	NFTG						
28	102	MSFT						
29	102	JOAN						
30	102	PRZO						
31	102	ADAF						
32	102	WBA						
33	102	OMGA						
34	102	VCNX						
35	102	EELQ						
36	102	ORGS						
37	102	COYA						
Plain Text		Tab Width: 8		Ln 6, Col 12		INS		

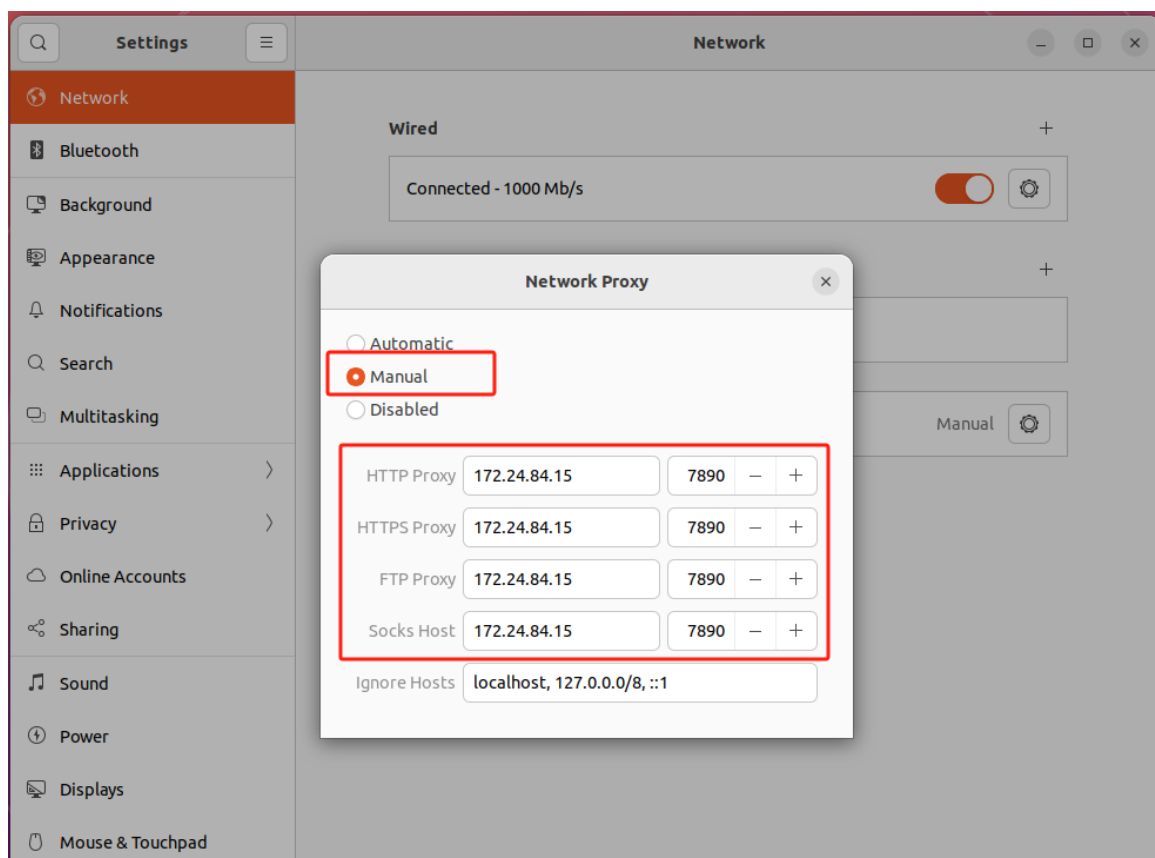
1.6. 附：由于在本作业中涉及到了虚拟机中github的使用，因此将在虚拟机中科学上网的方法记录如下

1. 打开clash中的Allow LAN按钮，点击三角符号，记录WLAN的IP地址及端口号



2. 打开虚拟机右上角的settings-Network，将Network Proxy修改为manual，并在四行中填入之前记录的IP及端口号





3. 设置完成 😊