

# 作业五

## 一、设计思路：

### 1.开发环境与工具：

基于vscode，安装Java开发插件和Maven插件，借助maven在vscode上搭建hadoop开发环境

### 2.编写MapReduce程序：

- 读取输入文件A和B，将它们作为输入数据。
- 编写Mapper阶段，Mapper的任务是将输入数据按照指数编号和成分股代码进行映射，然后输出（指数编号，成分股代码）作为键值对。
- 编写Reducer阶段，Reducer的任务是接收Mapper输出的键值对，并在Reducer内部进行数据合并和去重操作。
- 最后，Reducer将合并后的数据写入输出文件。

**3.运行MapReduce程序：** 使用Hadoop集群来运行MapReduce程序，指定输入文件A和B以及输出文件的路径。

## 二、程序运行结果说明：

### 1. 在vscode上新建Maven项目，并在pom.xml文件中新增有关hadoop的相关依赖配置

```
<dependencies>
  <dependency>
    <groupId>junit</groupId>
    <artifactId>junit</artifactId>
    <version>4.11</version>
    <scope>test</scope>
  </dependency>

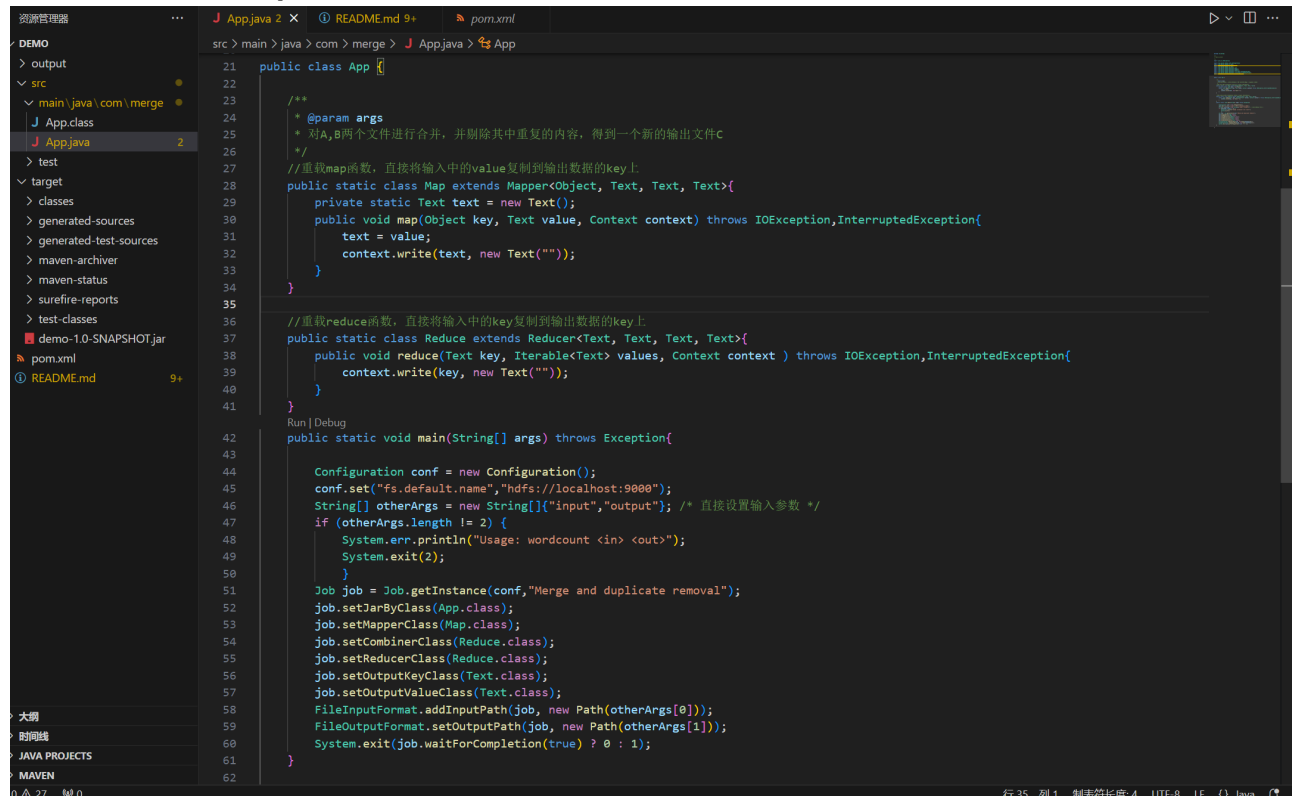
  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-common</artifactId>
    <version>3.3.6</version>
  </dependency>

  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-hdfs</artifactId>
    <version>3.3.6</version>
  </dependency>

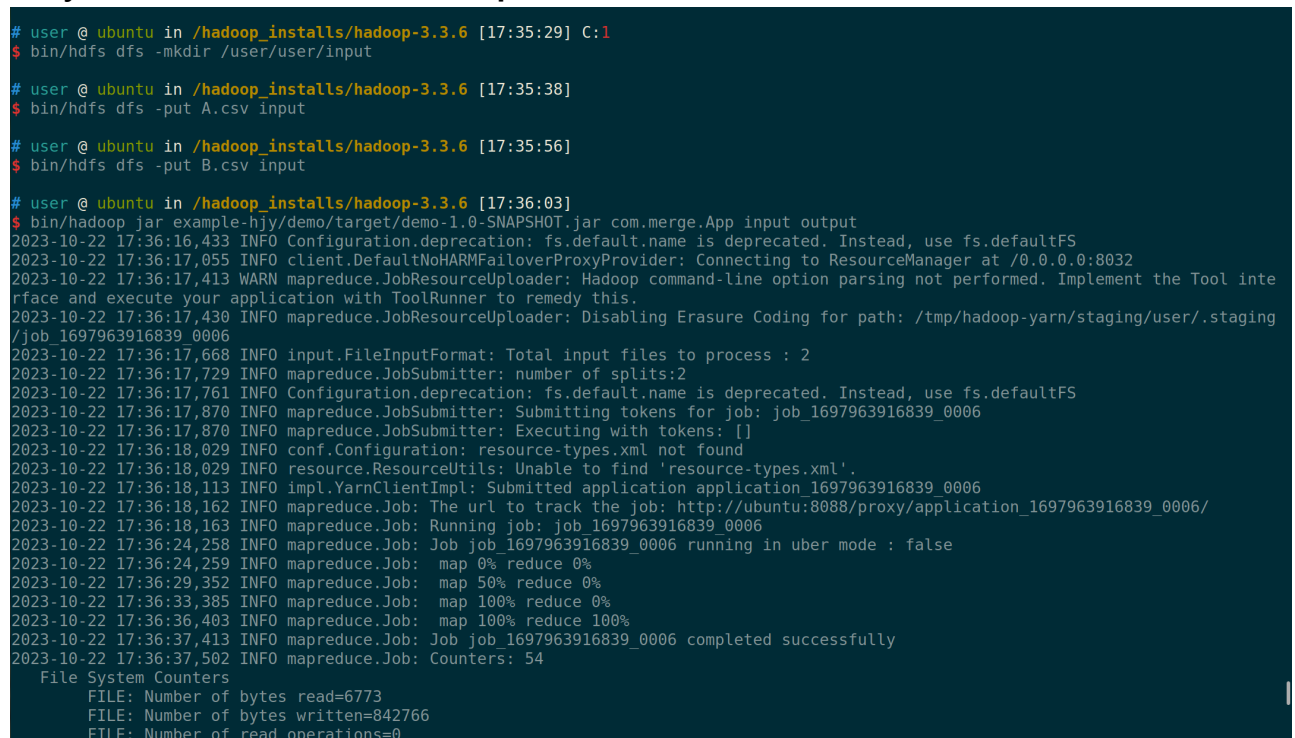
  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-client</artifactId>
    <version>3.3.6</version>
  </dependency>

  <dependency>
    <groupId>org.apache.hadoop</groupId>
    <artifactId>hadoop-yarn-api</artifactId>
    <version>3.3.6</version>
  </dependency>
</dependencies>
```

## 2. 编译源代码，实现MapReduce的合并去重思想



## 3. 导出jar文件，将程序复制到本地Hadoop系统的执行目录，在伪分布式环境下进行测试



## 三、运行成功的WEB页面截图：

1. part-r-00000和\_SUCCESS截图

part-r-00000 - 记事本  
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

101,AAPL  
101,AMGN  
101,AXP  
101,BA  
101,CAT  
101,CRM  
101,CSCO  
101,DIS  
101,DOW  
101,GS  
101,HD  
101,HON  
101,INTC  
101,JPM  
101,KO  
101,MMM  
101,MSFT  
101,NKE  
101,UNH  
101,V  
102,AAPL  
102,ACXP  
102,ADAF  
102,AHG  
102,ALXO  
102,AONC  
102,COYA  
102,CRVO  
102,CSCO  
102,EEIQ  
102,GRTS  
102,INTC  
102,JOAN  
102,MSFT  
102,NFTG  
102,OMGA  
102,ORGS  
102,PRZO  
102,SBFM  
102,TSBX

2. C.xlsx截图

