

# 实验报告

报告标题: MR-Hbase-Hive

学号: 21190630

姓名: 黄艺杰

日期: 2022 年 12 月 11 日

## 一、实验环境

1. 操作系统: Windows 10、Linux
2. 相关软件 (含版本号): VMWare15pro、FinalShell3.9、IntelliJ IDEA2019Professional
3. 其它工具: JDK1.8

## 二、实验内容及其完成情况

(针对上述实验内容逐一详述实验过程)

### 1. MapReduce 编程题

要求: 基于天气数据文件 (weather.txt), 计算月平均气温。要求有上传气象数据的代码, MR 程序代码, 运行过程, 下载并查看结果的代码验证方法: 在虚拟机中查看运行过程与结果。

先在 IDEA 上使用 put 函数将 weather.txt 文件上传到自定义文件夹 homework 中

### Browse Directory

/homework								Go!
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
-rw-r--r--	hyj	supergroup	166.96 KB	Thu Dec 08 21:34:43 +0800 2022	3	128 MB	weather.txt	

图 1.1HDFS 中 homework 文件夹

上传代码如下:

```
1. put("/D:/qiyesixun/wenjian/weather.txt", "/homework/weather.txt");
```

```
1. private static void put(String src,String dest) throws IOException {
2.     fileSystem.copyFromLocalFile(new Path(src),new Path(dest));
3.     System.out.println("put ok");
4. }
```

MR 代码如下:

### WeatherMapper:

```
1.  public class WeatherMapper extends Mapper<LongWritable, Text, Text, LongWritable> {
2.      @Override
3.      protected void map(
4.          LongWritable key,
5.          Text value,
6.          Mapper<LongWritable, Text, Text, LongWritable>.Context context
7.      ) throws IOException, InterruptedException {
8.          String[] ss = value.toString().split(",");
9.          if (ss[0].equals("DATE")) {
10.             return;
11.          }
12.
13.          String yymm = ss[0].substring(0, 6);
14.          try {
15.              context.write(
16.                  new Text(yymm),
17.                  new LongWritable(Long.parseLong(ss[5]))
18.              );
19.          } catch (Exception ex) {
20.              }
21.      }
22.  }
```

### WeatherReducer:

```
1.  public class WeatherReducer extends Reducer<Text, LongWritable, Text, DoubleWritable> {
2.      @Override
3.      protected void reduce(
4.          Text key,
5.          Iterable<LongWritable> values,
6.          Reducer<Text, LongWritable, Text, DoubleWritable>.Context context
7.      ) throws IOException, InterruptedException {
8.          long sum = 0L;
9.          long cnt = 0L;
10.         for (LongWritable lw : values) {
11.             cnt++;
12.             sum += lw.get();
13.         }
14.         context.write(
15.             key,
16.             new DoubleWritable(sum * 1.0 / cnt)
17.         );
18.     }
19. }
```

WeatherApp:

```
1.  public class WeatherApp {
2.      public static void main(String[] args) throws IOException {
3.          Configuration conf = new Configuration();
4.          Job job = Job.getInstance(conf);
5.          job.setJarByClass(WeatherApp.class);
6.          job.setMapperClass(WeatherMapper.class);
7.          job.setReducerClass(WeatherReducer.class);
8.          job.setMapOutputKeyClass(Text.class);
9.          job.setMapOutputValueClass(LongWritable.class);
10.         job.setOutputKeyClass(Text.class);
11.         job.setOutputValueClass(DoubleWritable.class);
12.         FileInputFormat.setInputPaths(job, new Path(args[0]));
13.         FileOutputFormat.setOutputPath(job, new Path(args[1]));
14.         try {
15.             job.waitForCompletion(true);
16.         } catch (Exception ex) {
17.             System.out.println(ex);
18.         }
19.     }
20. }
```

将以上文件打包后，上传到虚拟机中，并在虚拟机中运行结果，将结果保存到 HDFS 的 out1 文件夹中

The screenshot displays a virtual machine interface. On the left, a sidebar shows system information for 'FinalShell 3.9.3.4' with IP 192.168.110.101. The main area is divided into two panes. The top pane shows system metrics: running time 1:26, load 0.24, 0.13, 0.14, CPU 3%, memory 74% (1.3G/1.8G), and swap 0% (1M/2G). Below this is a table of running processes:

内存	CPU	命令
2.9M	2	sshd
210.6K	1.3	java
2.9M	0.7	top
0	0.3	migration/0

The bottom pane shows a terminal window with the command `[hyj@master soft]$ hadoop jar hadoop-demo-1.0.0.jar edu.nnu.weather_mr.WeatherApp /homework/weather.txt /out1`. The terminal output shows various Hadoop logs, including warnings about native code loading and deprecated session IDs, and informational messages about job submission and execution. The job ID is `job_local874523496_0001`.

图 1.2 在虚拟机上运行 WeatherApp

可在虚拟机上远程查看结果内容

```
[hyj@master soft]$ hadoop fs -cat /out1/part-r-00000
22/12/08 22:39:17 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your plat
form... using builtin-java classes where applicable
201601 -3.808199121522694
201602 1.680161943319838
201603 9.741197183098592
201604 16.772661870503597
201605 21.534170153417016
201606 25.697014925373136
[hyj@master soft]$
```

图 1.3 虚拟机上远程查看结果

将实验结果下载到本地 result 文件中的 weather.txt，并在本地查看，下载代码如下：

```
1. get("/out1/part-r-00000", "/D:/qiyeshixun/result/weather.txt");

1. private static void get(String src, String dest) throws IOException {
2.     fileSystem.copyToLocalFile(false, new Path(src), new Path(dest), true);
3.     System.out.println("get ok!");
4. }
```

在本地打开下载的 weather.txt 文件查看结果：

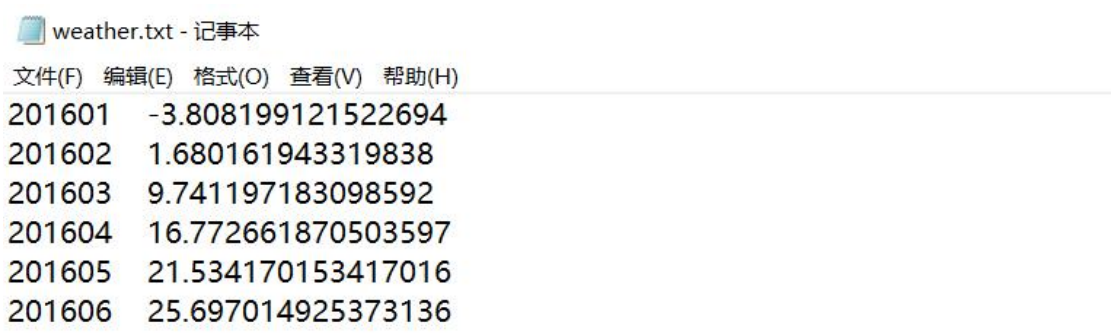


图 1.4 在本地查看下载的实验结果

2. MapReduce 编程题

要求：统计在四大名著中出现总次数最多的 20 个词（2 个字或以上的词）。要求有文件上传的代码、MR 程序、运行过程，下载结果的代码，排序数据并打印最终结果（输出：词和出现次数）的代码。验证方法：在虚拟机中查看运行过程，在 IDEA 中查看最终结果。

将四大名著上传到 HDFS 的 demo 文件夹中并合并成一个文件 Four\_masterpieces.txt

Browse Directory

/demo								Go!
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
-rw-r--r--	hyj	supergroup	657 B	Sat Dec 03 22:14:16 +0800 2022	3	128 MB	2.txt	
-rw-r--r--	hyj	supergroup	8.53 MB	Sat Dec 03 20:56:23 +0800 2022	3	128 MB	Four_masterpieces.txt	

图 2.1 上传四大名著并合并成一个文件

上传代码如下:

```
1. writeAllFilesOfDir("/D:/qiyeshixun/book", "/demo/Four_masterpieces.txt");
```

```
1. private static void writeAllFilesOfDir(String dir,String dest) {
2.     File path = new File(dir);
3.     if(path.exists() && path.isDirectory()){
4.         File[] files=path.listFiles();
5.         byte[] bts=new byte[4096];
6.         int size = 0;
7.         FSDataOutputStream outputStream=null;
8.         try {
9.             outputStream = fileSystem.create(new Path(dest));
10.            for (int i = 0; i < files.length; i++) {
11.                FileInputStream fis = null;
12.                if (files[i].isFile()) {
13.                    try {
14.                        fis = new FileInputStream(files[i]);
15.                        while ((size = fis.read(bts)) > 0) {
16.                            outputStream.write(bts,0,size);
17.                        }
18.                    } catch (Exception ex) {
19.                        System.out.println(ex);
20.                    } finally {
21.                        try {
22.                            fis.close();
23.                        } catch (Exception e) {
24.                        }
25.                    }
26.                }
27.            }
28.        }
29.        catch(Exception ex){
30.            System.out.println(ex);
31.        }
32.        finally{
33.        }
34.    }
35.    else{
36.        System.out.println("not found");
37.    }
38. }
```

MR 代码如下:

WordCountMapper 代码:

```

1.  public class WordCountMapper extends Mapper<LongWritable, Text,Text,LongWritable> {
2.      @Override
3.      protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
4.          JiebaSegmenter seg=new JiebaSegmenter();
5.          List<SegToken> tokenList=seg.process(value.toString(),JiebaSegmenter.SegMode.SEARCH);
6.          LongWritable one=new LongWritable(1);
7.          for(SegToken segToken : tokenList){
8.              if(segToken.word.matches("[\u4e00-\u9fa5]{2,}")){
9.                  context.write(new Text(segToken.word),one);
10.             }
11.         }
12.     }
13. }

```

WordCountReducer 代码:

```

1.  public class WordCountReducer extends Reducer<Text, LongWritable,Text,LongWritable> {
2.      @Override
3.      protected void reduce(Text key, Iterable<LongWritable> values, Context context) throws IOException, InterruptedException {
4.          //super.reduce(key, values, context);
5.          long sum=0;
6.          for (LongWritable value : values) {
7.              sum+=value.get();
8.          }
9.          context.write(key,new LongWritable(sum));
10.     }
11. }

```

WordCountApp 代码:

```

1.  public class WordCountApp {
2.      public static void main(String[] args) throws Exception {
3.          Configuration config=new Configuration();
4.          Job job=Job.getInstance(config);
5.
6.          job.setMapperClass(WordCountMapper.class);
7.          job.setReducerClass(WordCountReducer.class);
8.          job.setJarByClass(WordCountApp.class);
9.
10.         job.setMapOutputKeyClass(Text.class);
11.         job.setMapOutputValueClass(LongWritable.class);
12.         job.setOutputKeyClass(Text.class);
13.         job.setOutputValueClass(LongWritable.class);
14.

```

```

15.         FileInputFormat.setInputPaths(job, new Path(args[0]));
16.         FileOutputFormat.setOutputPath(job,new Path(args[1]));
17.
18.         job.waitForCompletion(true);
19.     }
20. }

```

将以上代码打包在虚拟机上运行，将结果保存在 HDFS 的 out2 文件夹，运行过程如下：

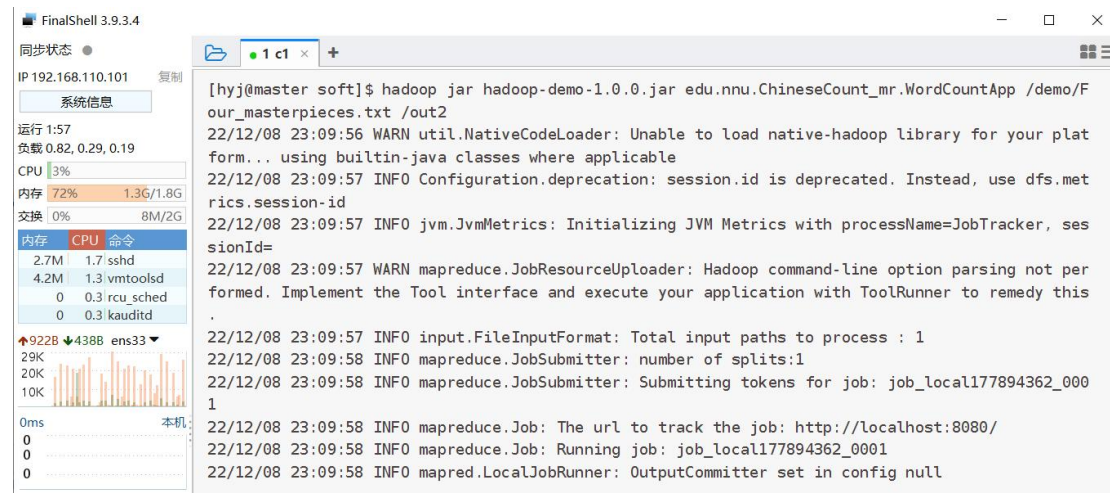


图 2.2 WordCountApp 在虚拟机上运行过程

下载结果的代码如下：

```

1.     get("/out2/part-r-00000", "/D:/qiyeshixun/result/Four_masterpieces.txt");

1.     1. private static void get(String src, String dest) throws IOException {
2.         2.         fileSystem.copyToLocalFile(false, new Path(src), new Path(dest), true);
3.         3.         System.out.println("get ok!");
4.         4.     }

```

最后在 IDEA 排序数据并打印最终的结果，其代码如下：

```

1.     public class WordCountResultApp {
2.         static class Word{
3.             public String key;
4.             public Integer count;
5.         }
6.         public static void main(String[] args){
7.             File file =new File("D:/qiyeshixun/result/Four_masterpieces.txt");
8.             if(!file.exists()){
9.                 System.out.println("file not found");
10.                return;
11.            }
12.            InputStream is=null;

```

```

13.         Reader reader=null;
14.         BufferedReader br=null;
15.         List<Word> words=new ArrayList<Word>();
16.         try{
17.             is=new FileInputStream(file);
18.             reader =new InputStreamReader(is);
19.             br=new BufferedReader(reader);
20.             String line=null;
21.             while((line=br.readLine())!=null){
22.                 String []tmp=line.split("\\t");
23.                 Word word=new Word();
24.                 word.key=tmp[0];
25.                 word.count=Integer.parseInt(tmp[1]);
26.                 words.add(word);
27.             }
28.         }catch(Exception ef){
29.             System.out.println(ef);
30.         }finally{
31.             try{
32.                 br.close();
33.             }catch (Exception e){}
34.             try{
35.                 reader.close();
36.             }catch(Exception e){}
37.             try{
38.                 is.close();
39.             }catch(Exception e){}
40.         }
41.         words.sort((a,b)->b.count-a.count);
42.         for(int i=0;i<20;i++){
43.             System.out.println(words.get(i).key+" "+words.get(i).count);
44.         }
45.     }
46. }

```

运行后输出四大名著出现总次数最多的 20 个词：



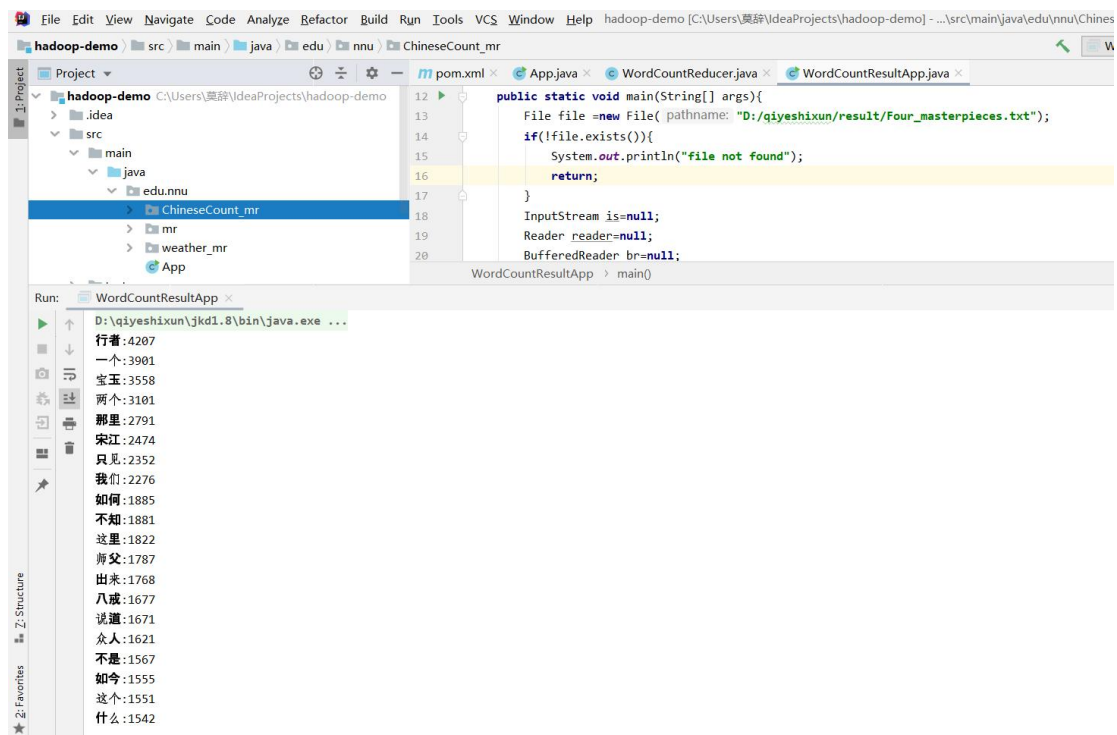


图 2.3 统计四大名著热词

### 3. Hbase 集群的安装配置（包括 Zookeeper 的分布式配置）：

验证方法：在宿主机中通过浏览器查看 CentOS 7 的 60010 端口

先在虚拟机上将 zookeeper 的压缩包解压到根目录并改名为 zookeeper

```
[hyj@master ~]$ cd soft
[hyj@master soft]$ ll
总用量 27408
-rw-rw-r--. 1 hyj hyj 28065455 12月  9 12:43 zookeeper-3.4.5-cdh5.7.0.tar.gz
[hyj@master soft]$ tar xvf zookeeper-3.4.5-cdh5.7.0.tar.gz -C ~/
```

图 3.1 解压 zookeeper 的压缩包

接下来在 .bash\_profile 中配置环境，配置 zoo.cfg 文件

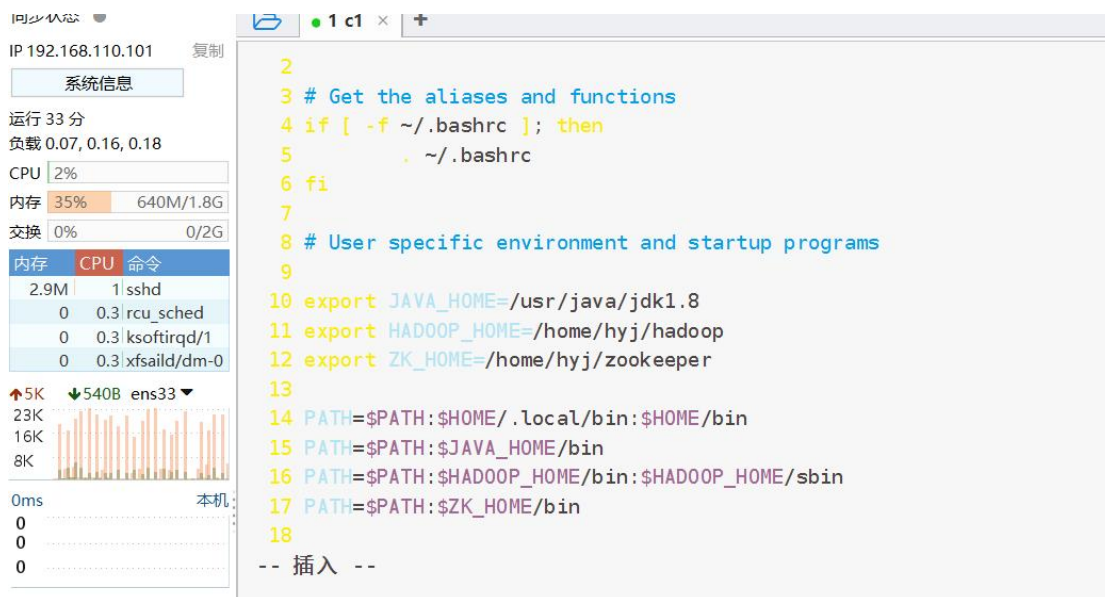


图 3.2 编辑.bash\_profile 环境配置文件

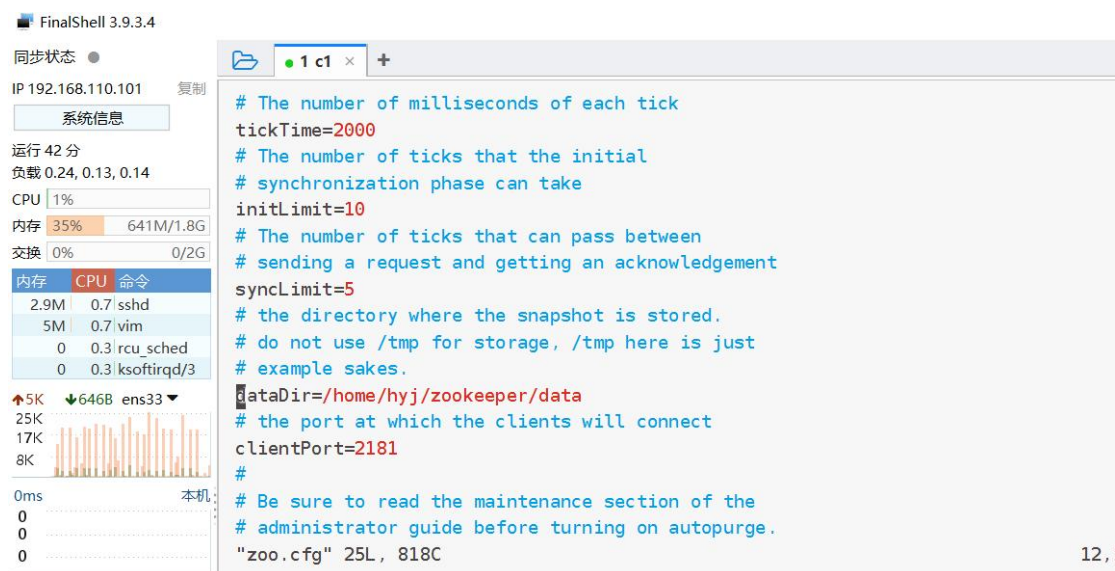


图 3.3 编辑 zoo.cfg 文件

接下来上传 Hbase 的文件，与上述过程基本类似，先上传文件并解压缩，改名，添加环境变量，修改 hbase-env.sh 文件。

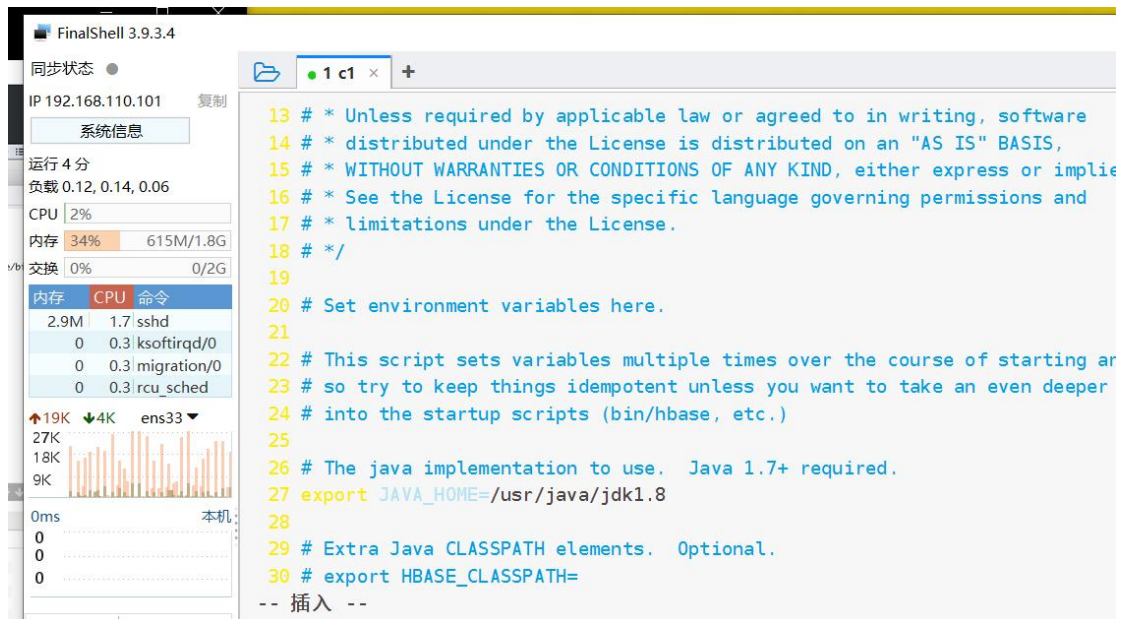


图 3.4 配置 hbase-env.sh 文件

### 配置 hbase-site.xml

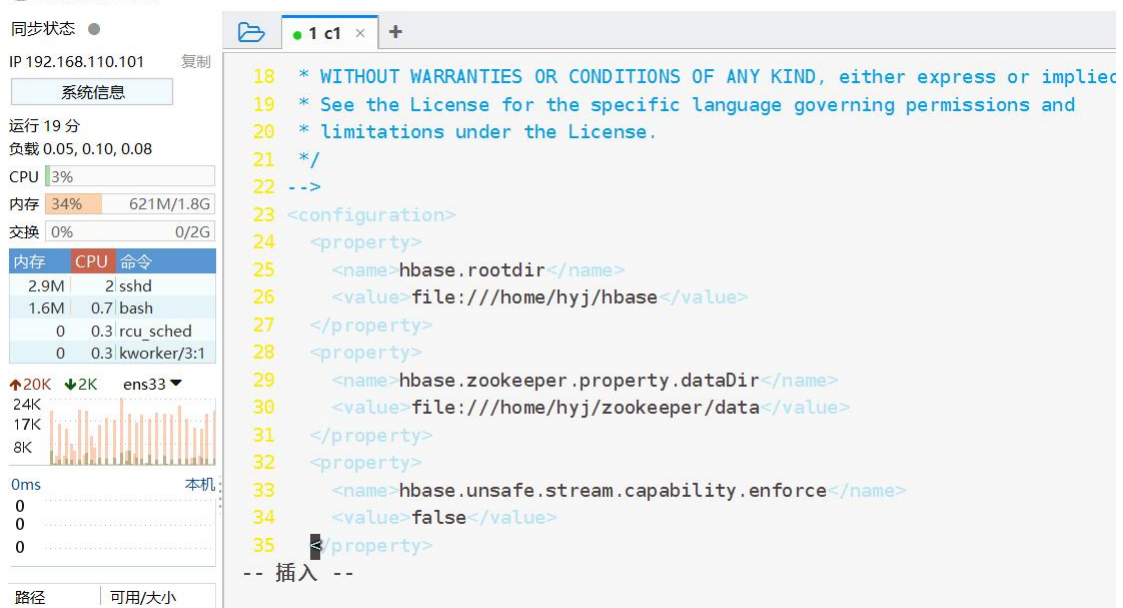


图 3.5 配置 hbase-site.xml 文件

打开 hbase.sh, 在 jps 中看到 HMaster

```

[hyj@master conf]$ start-hbase.sh
starting master, logging to /home/hyj/hbase/logs/hbase-hyj-master-master.out
[hyj@master conf]$ jps
37633 HMaster
38409 Jps

```

图 3.6 开启 hbase.sh

最后进入浏览器中 60010 端口

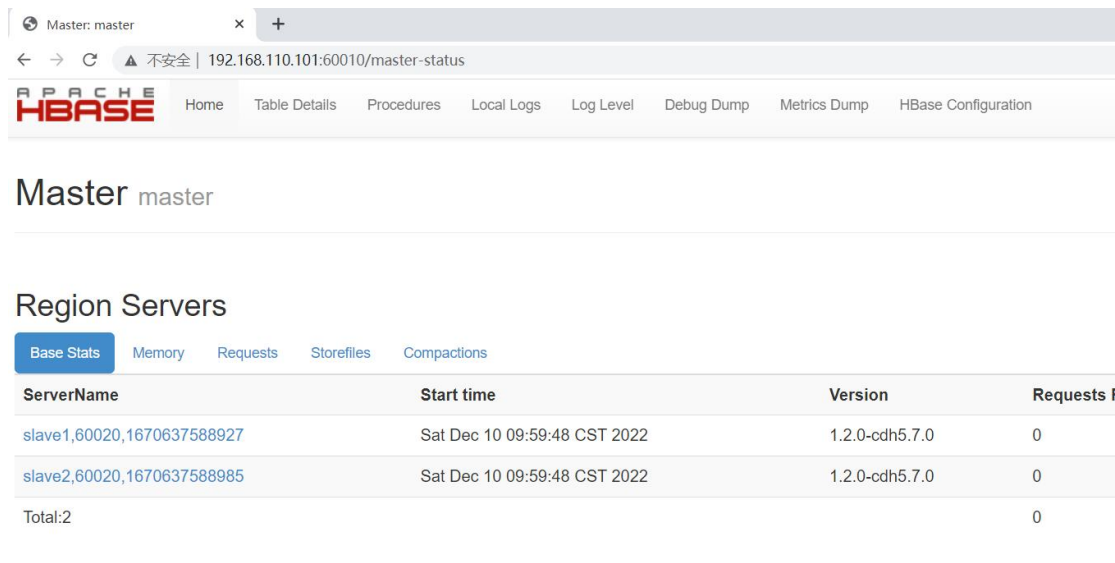


图 3.7 浏览器查看 60010 端口

#### 4. Hbase 编程题

基于课堂案例，定义一个方法，参数为课程名称，功能为完成根据该课程名称查询出该课程所有学生的成绩总分和均分，在 `main` 方法中测试该方法（课程名自行指定，如果是不存在的课程名，则输出提示信息后结束）。要求有完整的代码。

验证方法：在 IDEA 中输出运行结果

该方法代码如下：

```

1. private static void CourseScores(String courseName) throws Exception {
2.     String single="";
3.     switch(courseName){
4.         case "Java":single="C001";
5.         break;
6.         case "大数据":single="C002";
7.         break;
8.         case "MySQL":single="C003";
9.         break;
10.        case "数学":single="C004";
11.        break;
12.    }
13.    if(single.equals("")){System.out.println("无此课程");return;}
14.    Table table=connection.getTable(TableName.valueOf("student"));
15.    Scan scan=new Scan();
16.    Filter filter=new ColumnPrefixFilter(single.getBytes());
17.    scan.setFilter(filter);
18.    ResultScanner result= table.getScanner(scan);
19.    Result rs=null;

```

```

20.         int sum=0;
21.         int cnt=0;
22.         while((rs=result.next())!=null) {
23.             Cell[] cells = rs.rawCells();
24.             for (int i = 0; i < cells.length; i++) {
25.                 Cell cell = cells[i];
26.                 String val = new String(CellUtil.cloneValue(cell));
27.                 sum+=Integer.parseInt(val);
28.                 cnt++;
29.             }
30.         }
31.         System.out.println(coursename+"课程学生的总分为:"+sum+"\t 平均分为:"+sum*1.0/cnt));
32.         table.close();
33.     }

```

运行完整代码，查询 MySQL 课程，结果如下：

```

Run: StudentApp x
D:\qiyeshixun\jdk1.8\bin\java.exe ...
log4j:WARN No appenders could be found for logger (org.apache.hadoop.security.Groups).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
表【student】创建成功，共有【2】个列族，分别为：[info, scores]
表【course】创建成功，共有【1】个列族，分别为：[info]
数据添加成功！
数据添加成功！
数据添加成功！
行键【002】 scores:C001 69
MySQL 课程学生的总分为:5322 平均分为:73.91666666666667
Process finished with exit code 0

```

图 4.1MySQL 课程学生总分和平均分

## 5. Hive 的安装配置（包括 MySQL 的安装配置）

验证方法：在 CentOS 7 中使用 beeline 方式连接 Hive

首先是上传并解压 hive，之后改名为 hive

```

[hyj@master soft]$ ls
hive-1.1.0-cdh5.7.0.tar.gz
[hyj@master soft]$ tar xvf hive-1.1.0-cdh5.7.0.tar.gz -C ~/
hive-1.1.0-cdh5.7.0/

```

图 5.1 上传并解压 hive 的压缩文件

然后为 hive 配置环境路径

```

[hyj@master ~]$ echo $PATH
/usr/local/bin:/usr/bin:/usr/local/sbin:/usr/sbin:/home/hyj/.local/bin:/home/hyj/bin:/usr/java/jdk1.8/bin:/home/hyj/hbase/bin:/home/hyj/hadoop/bin:/home/hyj/hadoop/sbin:/home/hyj/zookeeper/bin:/home/hyj/hive/bin

```

图 5.2 为 hive 设置环境变量

之后配置 hive-env.sh 文件和 hive-site.xml 文件，上传 mysql 的 jar 包，在虚拟机上安装相关



## 依赖

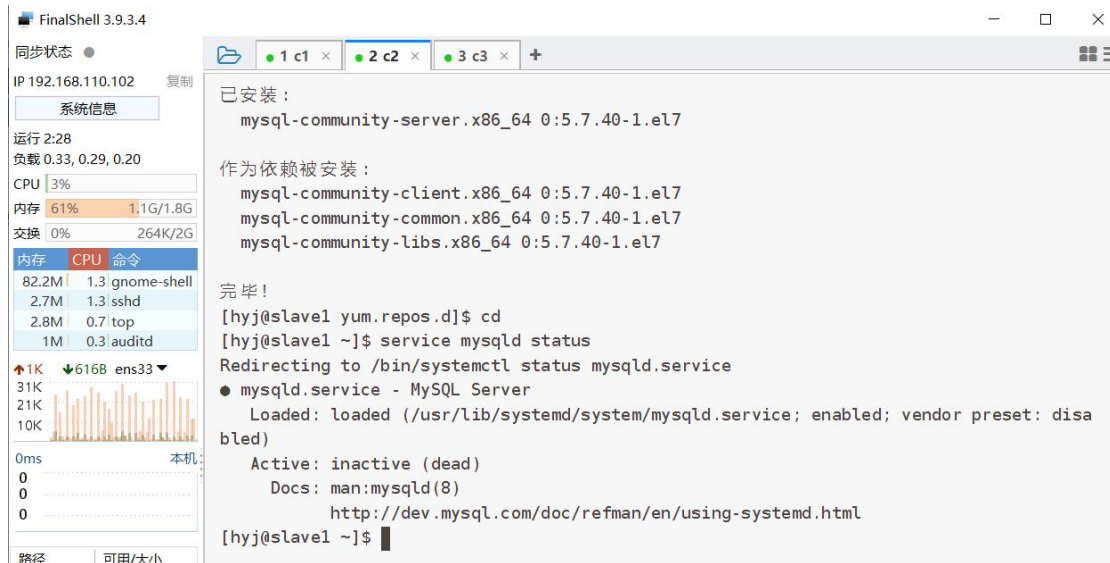


图 5.3 查看数据库状态

接着激活数据库并修改密码，修改为 Hyjnb@102

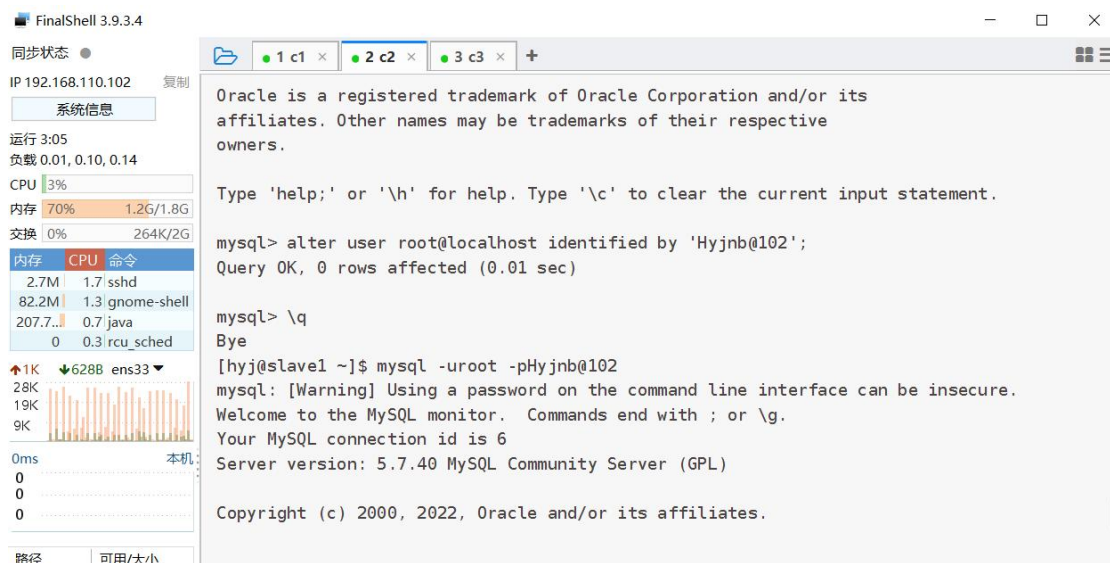


图 5.4 进入数据库并修改密码

最后用 beeline 命令进入 beeline 模式，利用 ! connect jdbc:hive2://master:10000 连接到 Hive

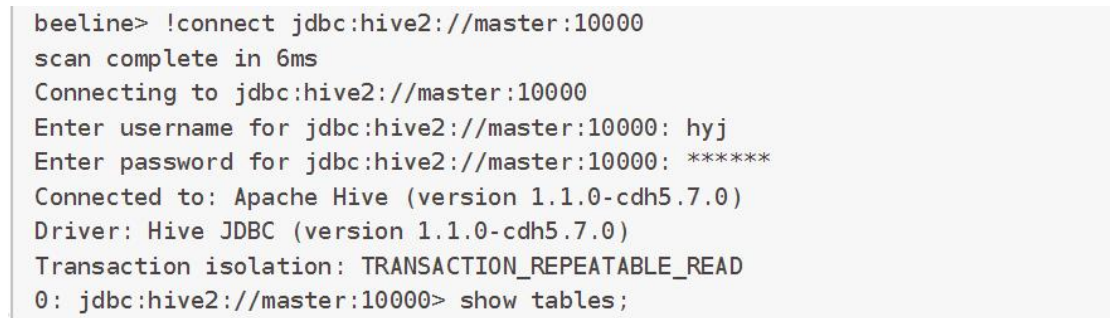


图 5.5 使用 beeline 方式连接 Hive

## 6. Hive 编程题

在 beeline 模式下创建数据库，命名为：hive\_+自己姓名拼音缩写，将该数据文件上传到 hdfs 后，导入 music 表；编写 SQL 语句，完成以下查询：学号尾数为单数的同学统计查询每首歌的播放次数学号尾数为双数的同学统计查询每种播放端的播放次数每个步骤均要求完整截图。

先通过 beeline 连接 Hive，进入 beeline 模式

```
[hyj@master ~]$ beeline
2022-12-11 15:35:53,403 WARN [main] mapreduce.TableMapReduceUtil: The hbase-prefix-tree module jar containing PrefixTreeCodec is not present. Continuing without it.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/hyj/hbase/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/hyj/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
2022-12-11 15:35:53,800 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Beeline version 1.1.0-cdh5.7.0 by Apache Hive
```

图 6.1 使用 beeline 连接

之后在 beeline 模式通过 nohup hiveserver2 > hs2.log 2>&1 & 语句连接数据库

```
beeline> !connect jdbc:hive2://master:10000
scan complete in 4ms
Connecting to jdbc:hive2://master:10000
Enter username for jdbc:hive2://master:10000: hyj
Enter password for jdbc:hive2://master:10000: *****
Connected to: Apache Hive (version 1.1.0-cdh5.7.0)
Driver: Hive JDBC (version 1.1.0-cdh5.7.0)
Transaction isolation: TRANSACTION_REPEATABLE_READ
```

图 6.2 在 beeline 模式中连接数据库

进入数据库后按要求创建并查看数据库

```
0: jdbc:hive2://master:10000> create database hive_hyj;
No rows affected (1.159 seconds)
0: jdbc:hive2://master:10000> show databases;
+-----+
| database_name |
+-----+
| default      |
| demo         |
| hive_hyj     |
```

图 6.3 查看新建的 hive\_hyj 数据库

接着进入创建的 hive\_hyj 数据库，按照给定数据类型和题目要求创建表格，将事先上传在 HDFS 上的 music.txt 导入到新创建的数据库中

```
0: jdbc:hive2://master:10000> use hive_hyj;
No rows affected (0.149 seconds)
0: jdbc:hive2://master:10000> create table music(
0: jdbc:hive2://master:10000> id int,
0: jdbc:hive2://master:10000> song string,
0: jdbc:hive2://master:10000> singer string,
0: jdbc:hive2://master:10000> sex string,
0: jdbc:hive2://master:10000> rhythm string,
0: jdbc:hive2://master:10000> player string)
0: jdbc:hive2://master:10000> row format delimited fields terminated by '\t';
No rows affected (1.582 seconds)
0: jdbc:hive2://master:10000> load data inpath 'hdfs://master:8020/homework/music.txt'
into table music;
No rows affected (1.395 seconds)
```

图 6.4 创建并导入 music 表

使用 SQL 语句查看表内容

```
0: jdbc:hive2://master:10000> select * from music;
+-----+-----+-----+-----+-----+-----+
+--+
| music.id | music.song | music.singer | music.sex | music.rhythm | music.player |
+-----+-----+-----+-----+-----+-----+
+--+
| 1        | song1      | singer1      | man       | slow         | pc           |
| 2        | song2      | singer2      | woman     | quick        | android      |
| 3        | song3      | singer3      | man       | quick        | ios          |
| 4        | song4      | singer4      | man       | slow         | pc           |
| 5        | song5      | singer5      | woman     | quick        | android      |
| 6        | song6      | singer6      | man       | quick        | ios          |
```

图 6.5 查看 music 表中内容

利用 SQL 语句查询每个播放器的使用次数

```
0: jdbc:hive2://master:10000> select player,count(1) from music group by player;
+-----+-----+
| player | _c1 |
+-----+-----+
| android | 11 |
| ios     | 10 |
| pc      | 9  |
+-----+-----+
3 rows selected (4.712 seconds)
```

图 6.6 查询不同播放器的使用次数

### 三、实验总结

（可以总结实验中出现问题以及解决的思路，也可以列出没有解决的问题）



1. 在第一次用 `mysql -uroot -p*****` 命令进入数据库时，密码中的符号前要加 `' \'` 转义，否则会报错
2. 在将 IDEA 上写的代码到虚拟机上运行时，每次都要先 `clean`，然后在 `package`，最后将打包好的文件上传