# FBDP作业五

## maven的熟悉与操作

1. 首先在ide如vscode或者intellij上构建maven类，创建
2. 修改pom.xml，载入Hadoop，hdfs等插件；pom.xml内需要填写dependence以及build的插件信息
3. `mvn clean` 先对文件进行晴空
4. `mvn package` 对文档进行编译，得到jar文件
5. 发现报错 `Exception in thread "main" java.lang.ClassNotFoundException: cys.nju.edu.cn.WordCount` 原因是在pom里面要对主要的class进行定位，定位信息应该是 packageName.className，若定位不准确则会报错
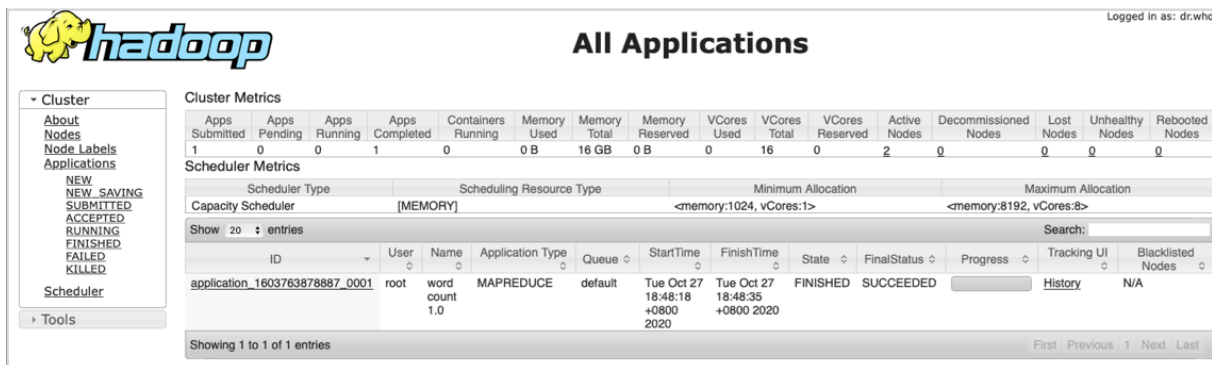
## 在bdkit上运行

Hdfs运行的几个代码（往后熟悉就会记住）

```
hdfs dfs -mkdir /wordcount

hdfs dfs -put /input /wordcount

hadoop jar target/jarFileName.jar /wordcount/input /wordcount/output
```

首先运行的是wordcount1.0版本，确定熟悉mapreduce的基本编写以及maven的用法、hadoop的运行等

## 2.0 版本

此为1.0版本结果截图，接着开始wordcount2.0版本，要求达到：忽略分词、停词，不区分大小写，忽略数字等。

由此，对1.0版本所作的修改为：

1. +setup() – 处理参数的输入，如caseSensitive的选择，skip文件的来源等
2. 在map()方法进行修改，利用正则表达式去除数字、字母以及排除停词，但是得到如下的结果，单词缺失字母。因为在stop-words-list中需要去除如'a'，'i'等短词，若使用replaceall方法，则会把长单词中此类字母也一起去掉。

```
3837    wth
2534    ths
10449   hve
10241   scene
9300    thou
9156    wll
5647    thy
 382    shll
文件夹中打开桌面的内容
 309    enter
4015    whch
3806    let
3388    lke
3350    love
3206    nry
2879    mke
2856    gve
2662    shkespe
2652    hopge
2458    pvo
2383    frst
2253    tke
2152    set
```

解决方法：利用正则表达式去掉字母和符号，利用tokenizer的特性可以分离不同的词，分离后判断单词长度且与停词比较是否要去掉或纳入。

```java
String num = "[0-9]+";

String p = "\\pP";

line = line.replaceAll(num,"");

line = line.replaceAll(p,"");

StringTokenizer itr = new StringTokenizer(line);

while (itr.hasMoreTokens()) {

    boolean founded = false;

    String s = itr.nextToken();

    for (String pattern : patternToSkip) {

        if (s.equals(pattern)) {

            founded = true;

            break;

        }

    }

if (s.length() >= 3 && !founded) {

    word.set(s);

    context.write(word, one);

}
```

3. 实现倒序排列输出。原本按照1.0的方法得到一个输出，存储在临时文件夹，重新利用map读取原结果，利用inversemapper class可以改变自动sort按照单词字典序排序，转为按照value即单词计数来排序。同时在比较时引入 `IntWritableDecreasingComparator` 的override版本，正向排序可变为逆向排序（比较结果改为负数即可）

4. 按格式要求输出，则需要修改reducer类别，使输出的的key为排名，value为单词与计数的string。并修改main的输出reducer class。



## 本次实验的收获

1. 熟悉了maven创建操作以及bdkit的使用
2. 熟悉了git仓库操作与使用
3. 熟悉了mapreduce程序编写的方法