# 华中科技大学

# 大数据处理实验报告

## 实验三：MapReduce 的基本操作

姓　　名：刘日星

学　　号：X2020I1007

院　　系：计算机科学与技术

专　　业：计科与金融

年　　级：CS1804-交换生

指导教师：石宣化

2021 年 12 月 15 日

# 一：实验目的

1、了解 MapReduce 的用途

2、掌握 MapReduce 的基本命令

# 二：实验要求

1、第四节中的实验内容要附上完整的实验过程截图以及必要的文字说明，每个人的 IP 地址等不同，不能直接套用样例的截图。

2、请同学们在完成报告后，将报告的 pdf 版本命名为：大数据实验三+姓名+学号.pdf，并在这周五前,发到邮箱:
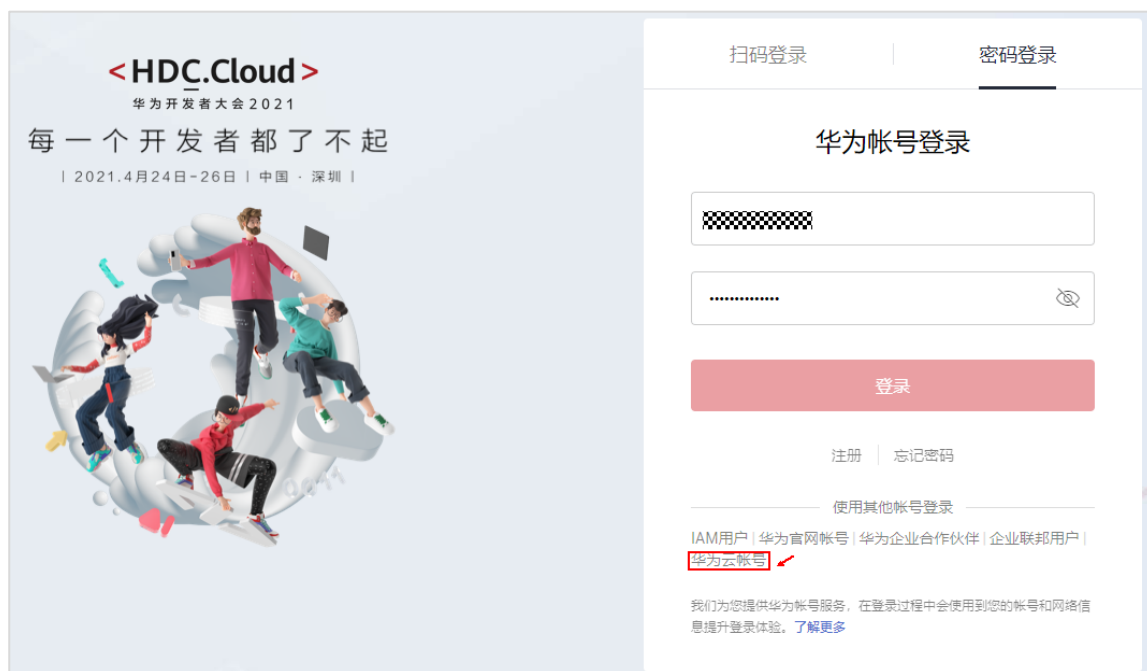
aaaaltaaaa@126.com(石老师班)

798792873@qq.com(郑老师班)

# 三：实验环境配置

步骤 1 登录华为云网站

https://www.huaweicloud.com



点击右上角登录，输入账号和密码

注意：华为云已统一登录入口，若仍不能登录则点击下方华为云账号进行登录。

步骤 2 点击"EI 企业智能"选择"MapReduce 服务"



步骤 3 点击"立即购买"

选择"自定义购买"(这里还要选上 Spark)



点击下一步，进入硬件配置

选择"按需计费","可用区 2",点击"弹性公网 IP",如下图:



点击"购买弹性公网 IP",选择"按需计费","按流量计费","5M",点击

"立即购买",如下图:



点击"提交",如下图:

| 产品类型 | 产品规格 | | 计费模式 | 数量 | 价格 |
|---|---|---|---|---|---|
| 弹性公网IP | 区域 | 北京四 | 按需计费 | 1 | ¥0.02/小时 |
| | 类型 | 全动态BGP | | | |
| | IPv6转换 | 停用 | | | |
| | 标签 | -- | | | |
| 带宽 | 带宽名称 | bandwidth-818e | 按需计费 | 1 | ¥0.80/GB |
| | 带宽类型 | 独享带宽 | | | |
| | 计费方式 | 按流量计费 | | | |
| | 带宽大小 | 5 Mbit/s | | | |

弹性公网IP费用 ¥0.02/小时 ＋ 公网流量费用 ¥0.80/GB ⑦

上一步　提交

购买成功，如下图：

返回 MapReduce 服务自定义购买界面绑定 EIP

选择"鲲鹏计算"，关闭高可用，调整 core 节点数为 1,如下图：

点击"下一步"

高级配置项参考如下：



点击"确认授权"

点击"立即购买"

步骤 4  点击"返回集群列表",如下图:



创建过程需要等待几分钟,待状态变为"运行中"集群创建完成



步骤 5  点击集群名称

步骤 6  配置安全组

点击集群名称



选择"节点管理"，点击含有"master1"的节点



在弹出页面中选择"安全组"，点击"更改安全组规则"，如下图所示：



选择"入方向规则"，点击"一键放通"，确认即可。

步骤 7  远程登录 master 节点

在安全组配置项，点击右上方"远程登录"，选择 cloudshell 登录。



输入密码，点击连接即可。

## 四：实验内容及步骤、实验的详细记录、实验结果分析

**请附上实验过程截图（截图需包含指令）以及必要的文字分析**

## 4.1 MapReduce

## 4.1.1 进入 hadoop(5')

cd /opt/client/HDFS/Hadoop



## 4.1.2 添加环境变量(10')

export HADOOP="/opt/client/HDFS/hadoop/share/hadoop"

export

CLASSPATH="$HADOOP/common/hadoop-common-2.8.3-mrs-1.9.0.jar:$HADOOP/

mapreduce/hadoop-mapreduce-client-core-2.8.3-mrs-1.9.0.jar:$HADOOP/common

/lib/commons-cli-1.2.jar:$CLASSPATH"

```
[root@node-master1jgZo hadoop]# export HADOOP="/opt/client/HDFS/hadoop/share/hadoop"
[root@node-master1jgZo hadoop]# export CLASSPATH="$HADOOP/common/hadoop-common-2.8.3-mrs-1.9.
0.jar:$HADOOP/mapreduce/hadoop-mapreduce-client-core-2.8.3-mrs-1.9.0.jar:$HADOOP/common/lib/c
ommons-cli-1.2.jar:$CLASSPATH"
[root@node-master1jgZo hadoop]# _
```
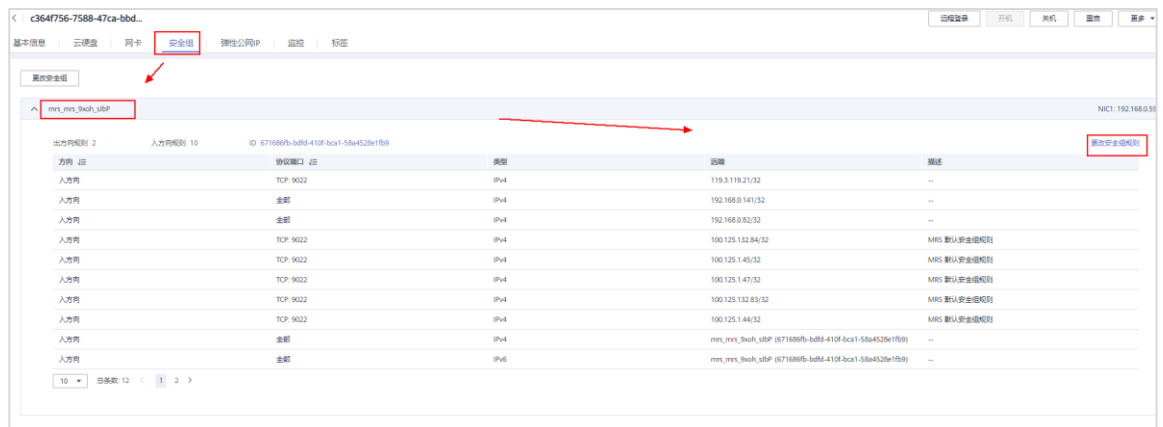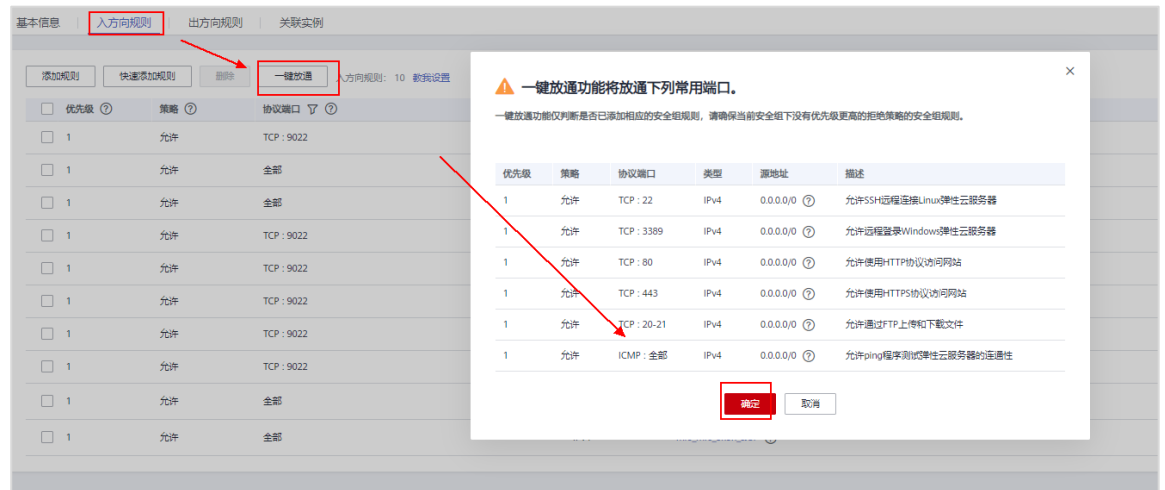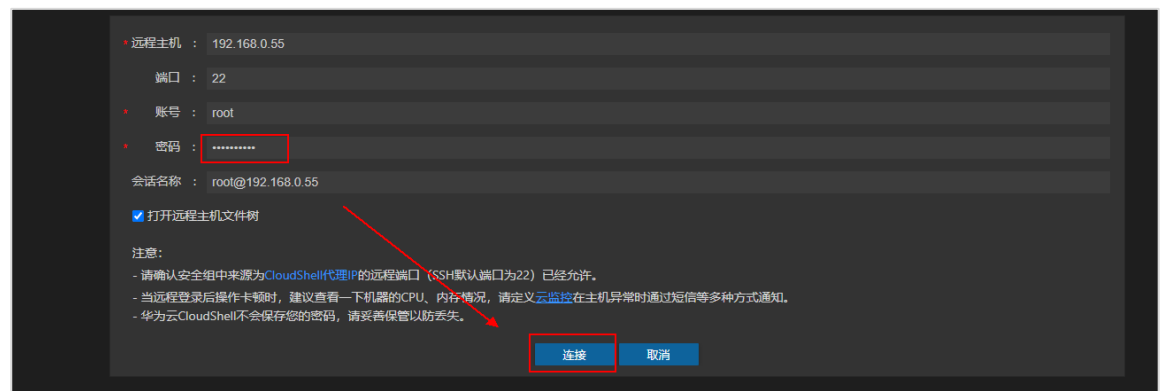
配置 java 环境

# 4.1.3 创建 java 程序 WordCount.java,在里面输入以下代码(5')

import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

  public static class TokenizerMapper
       extends Mapper<Object, Text, Text, IntWritable>{

    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(Object key, Text value, Context context
                    ) throws IOException, InterruptedException {

```
      StringTokenizer itr = new StringTokenizer(value.toString());
      while (itr.hasMoreTokens()) {
        word.set(itr.nextToken());
        context.write(word, one);
      }
    }
  }


  public static class IntSumReducer
        extends Reducer<Text,IntWritable,Text,IntWritable> {
    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context
                       ) throws IOException, InterruptedException {
      int sum = 0;
      for (IntWritable val : values) {
        sum += val.get();
      }
      result.set(sum);
      context.write(key, result);
    }
  }

  public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "word count");
    job.setJarByClass(WordCount.class);
    job.setMapperClass(TokenizerMapper.class);
    job.setCombinerClass(IntSumReducer.class);
    job.setReducerClass(IntSumReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
  }
}
```

```
[root@node-master1dlKK hadoop]# vim WordCount.java
[root@node-master1dlKK hadoop]#
```

```java
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {
        public static class TokenizerMapper
                extends Mapper<Object, Text, Text, IntWritable>{
                        private final static IntWritable one = new IntWritable(1);

                        private Text word = new Text();

                        public void map(Object key, Text value, Context context
                        ) throws IOException, InterruptedException {
                                StringTokenizer itr = new StringTokenizer(value.toString());
                                while (itr.hasMoreTokens()) {
                                        word.set(itr.nextToken());
                                        context.write(word, one);
                                }
                        }
                }
        public static class IntSumReducer
        extends Reducer<Text,IntWritable,Text,IntWritable> {
        private IntWritable result = new IntWritable();

        public void reduce(Text key, Iterable<IntWritable> values,
                        Context context
                                                                        1,1          Top
```

编写 java 文件

## 4.1.4 编译 WordCount.java(5')

```
[root@node-master1dlKK hadoop]# javac WordCount.java
[root@node-master1dlKK hadoop]# ls
 bin                LICENSE.txt          'WordCount$IntSumReducer.class'
 etc                NOTICE.txt           'WordCount$TokenizerMapper.class'
 hdfs-c-example     README.txt            WordCount.class
 include            sbin                  WordCount.java
 lib                share
 libexec            version.properties
[root@node-master1dlKK hadoop]#
```

```
[root@node-master1dlKK hadoop]# jar cf WordCount.jar WordCount*.class
[root@node-master1dlKK hadoop]# ls
 bin                NOTICE.txt           'WordCount$IntSumReducer.class'
 etc                README.txt           'WordCount$TokenizerMapper.class'
 hdfs-c-example     sbin                  WordCount.class
 include            share                 WordCount.jar
 lib                test1                 WordCount.java
 libexec            test2
 LICENSE.txt        version.properties
[root@node-master1dlKK hadoop]#
```

## 4.1.5 创建文件 test1,内容为 hello hust,文件 test2, 内容为 hello 学号,将他们放入 hdfs 的/input 文件夹 内。(方法见实验一)(5')

```
[root@node-master1dlKK hadoop]# vim test1
[root@node-master1dlKK hadoop]# vim test2
[root@node-master1dlKK hadoop]#
```

```
hello hust
~
~
~
~
```

```
hello X2020I1007
~
~
~
~
```

```
[root@node-master1dlKK hadoop]# hdfs dfs -mkdir /input
[root@node-master1dlKK hadoop]# hdfs dfs -put test1 /input
[root@node-master1dlKK hadoop]# hdfs dfs -put test2 /input
[root@node-master1dlKK hadoop]# hdfs dfs -ls /input
Found 2 items
-rw-r--r--   1 root ficommon          11 2021-12-14 20:27 /input/test1
-rw-r--r--   1 root ficommon          17 2021-12-14 20:28 /input/test2
[root@node-master1dlKK hadoop]#
```

# 4.1.6 运行 WordCount.jar 将 hdfs 的 /input 作为输入，/output 作为输出，并打印 /output 目录下的文件，显示出词频统计的结果(5')

export
HADOOP_CLASSPATH=$HADOOP_CLASSPATH:/opt/client/HDFS/hadoop/WordCount.jar

hadoop jar WordCount.jar WordCount hdfs:///input hdfs:///output

hdfs dfs -cat   /output/part-r-00000

```
[root@node-master1dlKK hadoop]# export HADOOP_CLASSPATH=$HADOOP_CLASSPATH:/o
pt/client/HDFS/hadoop/WordCount.jar
[root@node-master1dlKK hadoop]# hadoop jar WordCount.jar WordCount hdfs:///i
nput hdfs:///output
WARNING: Use "yarn jar" to launch YARN applications.
21/12/14 20:32:39 WARN mapreduce.JobResourceUploader: Hadoop command-line op
tion parsing not performed. Implement the Tool interface and execute your ap
plication with ToolRunner to remedy this.
21/12/14 20:32:40 INFO input.FileInputFormat: Total input files to process :
 2
21/12/14 20:32:42 INFO mapreduce.JobSubmitter: number of splits:2
21/12/14 20:32:43 INFO mapreduce.JobSubmitter: Submitting tokens for job: jo
b_1639483812097_0002
21/12/14 20:32:44 INFO impl.YarnClientImpl: Submitted application applicatio
n_1639483812097_0002
21/12/14 20:32:44 INFO mapreduce.Job: The url to track the job: http://node-
master1dlkk.mrs-fbva.com:8088/proxy/application_1639483812097_0002/
21/12/14 20:32:44 INFO mapreduce.Job: Running job: job_1639483812097_0002

[root@node-master1dlKK hadoop]# hdfs dfs -cat  /output/part-r-00000
X2020I1007      1
hello   2
hust    1
[root@node-master1dlKK hadoop]#
```

16

## 4.2 Spark

## 4.2.1 打开 spark(5')

Pyspark

```
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 2.2.2-mrs-1.9.0
      /_/

Using Python version 2.7.15 (default, Apr  1 2019 00:00:00)
SparkSession available as 'spark'.
>>>
```

## 4.2.2 读取 hdfs 文件内容(5')

lines = spark.read.text("hdfs:///input").rdd.map(lambda r: r[0])

```
>>> lines = spark.read.text("hdfs:///input").rdd.map(lambda r: r[0])
2021-12-14 20:36:58,279 | INFO | Thread-3 | resolve relation took 200 ms |
org.apache.spark.internal.Logging$class.logInfo(Logging.scala:54)
2021-12-14 20:36:58,659 | INFO | Thread-3 | skip CarbonOptimizer | org.apac
he.spark.sql.optimizer.CarbonLateDecodeRule.checkIfRuleNeedToBeApplied(Carbo
nLateDecodeRule.scala:95)
2021-12-14 20:36:58,659 | INFO | Thread-3 | skip CarbonOptimizer | org.apac
he.spark.sql.optimizer.CarbonLateDecodeRule.checkIfRuleNeedToBeApplied(Carbo
nLateDecodeRule.scala:95)
2021-12-14 20:36:58,659 | INFO | Thread-3 | Skip CarbonOptimizer | org.apac
he.spark.sql.optimizer.CarbonLateDecodeRule.apply(CarbonLateDecodeRule.scala
:72)
2021-12-14 20:36:58,691 | INFO | Thread-3 | Pruning directories with: | or
g.apache.spark.internal.Logging$class.logInfo(Logging.scala:54)
2021-12-14 20:36:58,695 | INFO | Thread-3 | Post-Scan Filters: | org.apach
```

## 4.2.3 词频统计(5')

counts = lines.flatMap(lambda x: x.split(' ')).map(lambda x: (x, 1)).reduceByKey(lambda x, y: x + y)
output = counts.collect()

```
>>> counts = lines.flatMap(lambda x: x.split(' ')).map(lambda x: (x, 1)).red
uceByKey(lambda x, y: x + y)
>>> output = counts.collect()
2021-12-14 20:41:04,911 | INFO | Thread-3 | Starting job: collect at <stdin
>:1 | org.apache.spark.internal.Logging$class.logInfo(Logging.scala:54)
2021-12-14 20:41:04,934 | INFO | dag-scheduler-event-loop | Registering RDD
 5 (reduceByKey at <stdin>:1) | org.apache.spark.internal.Logging$class.logI
nfo(Logging.scala:54)
2021-12-14 20:41:04,937 | INFO | dag-scheduler-event-loop | Got job 0 (coll
ect at <stdin>:1) with 2 output partitions | org.apache.spark.internal.Loggi
ng$class.logInfo(Logging.scala:54)
2021-12-14 20:41:04,938 | INFO | dag-scheduler-event-loop | Final stage: Re
sultStage 1 (collect at <stdin>:1) | org.apache.spark.internal.Logging$class
.logInfo(Logging.scala:54)
2021-12-14 20:41:04,938 | INFO | dag-scheduler-event-loop | Parents of fina
l stage: List(ShuffleMapStage 0) | org.apache.spark.internal.Logging$class.l
ogInfo(Logging.scala:54)
2021-12-14 20:41:04,941 | INFO | dag-scheduler-event-loop | Missing parents
: List(ShuffleMapStage 0) | org.apache.spark.internal.Logging$class.logInfo(
Logging.scala:54)
2021-12-14 20:41:04,947 | INFO | dag-scheduler-event-loop | Submitting Shuf
fleMapStage 0 (PairwiseRDD[5] at reduceByKey at <stdin>:1), which has no mis
sing parents | org.apache.spark.internal.Logging$class.logInfo(Logging.scala
:54)
2021-12-14 20:41:05,054 | INFO | dag-scheduler-event-loop | Block broadcast
_1 stored as values in memory (estimated size 13.0 KB, free 92.8 MB) | org.a
pache.spark.internal.Logging$class.logInfo(Logging.scala:54)
2021-12-14 20:41:05,080 | INFO | dag-scheduler-event-loop | Block broadcast
_1_piece0 stored as bytes in memory (estimated size 8.0 KB, free 92.8 MB) |
org.apache.spark.internal.Logging$class.logInfo(Logging.scala:54)
2021-12-14 20:41:05,087 | INFO | dispatcher-event-loop-3 | Added broadcast_
1_piece0 in memory on *.*.0.142:22873 (size: 8.0 KB, free: 93.3 MB) | org.ap
ache.spark.internal.Logging$class.logInfo(Logging.scala:54)
```

# 4.2.4 输出结果(10')

```
>>> print(output)
[(u'hust', 1), (u'X2020I1007', 1), (u'hello', 2)]
>>>
```

# 4.3 附加题

file1:

20210001 Math 90

20210002 Math 80

20210003 Math 70

file2:

20210001 English 80

20210002 English 70

20210003 English 60

# 1.将以上两个文件存入 hdfs(10')

```
[root@node-master1dlKK ~]# vim file1
[root@node-master1dlKK ~]# vim file2
[root@node-master1dlKK ~]# 
```

```
20210001 Math 90          20210001 English 80
20210002 Math 80          20210002 English 70
20210003 Math 70          20210003 English 60

~                         ~
~                         ~
~                         ~
~                         ~
~                         ~
~                         ~
```

```
[root@node-master1dlKK ~]# hdfs dfs -mkdir /hdfs
[root@node-master1dlKK ~]# hdfs dfs -put file1 /hdfs
[root@node-master1dlKK ~]# hdfs dfs -put file2 /hdfs
[root@node-master1dlKK ~]# hdfs dfs -ls /hdfs
Found 2 items
-rw-r--r--   1 root ficommon         52 2021-12-14 20:56 /hdfs/file1
-rw-r--r--   1 root ficommon         61 2021-12-14 20:56 /hdfs/file2
[root@node-master1dlKK ~]# []
```

将上述内容编写入文件并存入名为 hdfs 文件夹内

# 2.编写 mapreduce 的程序，输出每门课的平均成绩。(10')

# 3.编写 mapreduce 的程序，输出每位同学有多少门课成绩低于 75 分。(10')

# 五：实验总结(10′)

通过本次实验，让我学会了在 Hadoop 平台上配置 java 环境并且编译.java 文件对数据文件 test1 与 test2 进行处理以及在 Spark 平台上使用 python 代码对数据文件进行词频统计计算，从而实现 WordCount 功能（单词词频统计计算功能）。从中对 MapReduce 大数据并行处理的实现方法有了更多的认知。而本次最大的难点在于我还没学过 Java 这门语言，以至于在附加题最后两题上出现了卡顿，不知道如何实现那两个算法，同时还读不懂 Java 程序。因此我认识到了自己的短板在于缺少了 Java 编程语言方面的知识，在后面的实验开始前，我会好好地恶补一下这方面的知识，从而使实验完成得更加顺利。

## 附录:

1 MapReduce 官方教程

https://hadoop.apache.org/docs/r2.8.3/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html

2 Spark 官方教程

https://spark.apache.org/examples.html