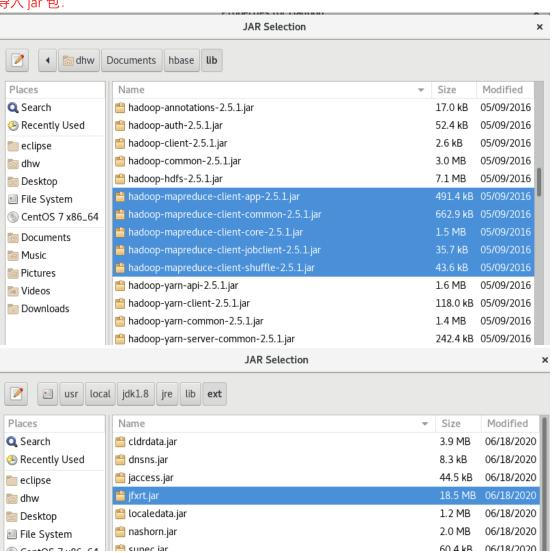
启动 Hadoop:

·	×
File Edit View Search Terminal Help	
File Edit View Search Terminal Help [dhw@localhost ~]\$ cd /usr/local/hadoop/ [dhw@localhost hadoop]\$./sbin/start-dfs.sh Starting namenodes on [localhost] dhw@localhost's password: localhost: starting namenode, logging to /usr/local/hadoop/logs/hadoop-dhw-nam ode-localhost.localdomain.out dhw@localhost's password: localhost: starting datanode, logging to /usr/local/hadoop/logs/hadoop-dhw-dat ode-localhost.localdomain.out Starting secondary namenodes [0.0.0.0] dhw@0.0.0.0's password: 0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop/logs/hadoop- w-secondarynamenode-localhost.localdomain.out [dhw@localhost hadoop]\$ jps 5956 Jps	an
5239 NameNode	
5735 SecondaryNameNode 5468 DataNode [dhw@localhost hadoop]\$	

导入 jar 包:



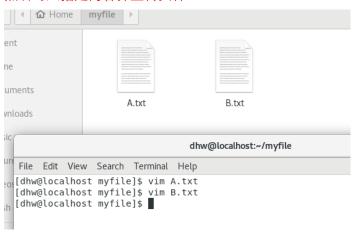
对于其他 jar 包,鼠标放上去按 ctrl 键点击左键,即可提示安装 jar 名称

实验

(一) 编程实现文件合并和去重操作

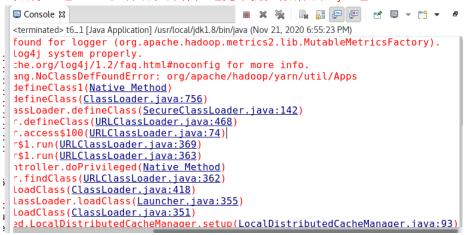
首先创建 input 目录

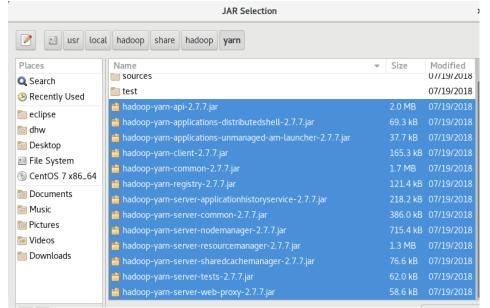
然后写入指定内容并上传文件



```
dhw@localhost:/usr/local/hadoop
                                                                              File Edit View Search Terminal Help
[dhw@localhost hadoop]$ ./bin/hdfs dfs -put /home/dhw/myfile/A.txt input
[dhw@localhost hadoop]$ ./bin/hdfs dfs -put /home/dhw/myfile/B.txt input
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls input
Found 3 items
-rw-r--r--
            1 dhw supergroup
                                       66 2020-11-21 17:44 input/A.txt
-rw-r--r--
                                       55 2020-11-21 17:45 input/B.txt
             1 dhw supergroup
- rw - r - - r - -
            3 dhw supergroup
                                       49 2020-10-31 14:27 input/t1.txt
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat input/A.txt
20180101 x
20180102 y
20180103 x
20180104 y
20180105 z
20180106 x
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat input/B.txt
20180101 v
20180102 v
20180103 x
20180104 z
20180105 y
[dhw@localhost hadoop]$
```

(删除上述 t1.txt 文件后在操作, 避免干扰) 报错添加 jar 包:





```
_ _
1 t6_1.java ≅
   import org.apache.hadoop.io.Text;
   import org.apache.hadoop.mapred.Merger;
   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 t6 1 {
       public static void main(String[] args) throws Exception {
            Configuration conf = new Configuration();
conf.set("fs.defaultFS", "hdfs://localhost:9000");
conf.set("fs.hdfs.impl", "org.apache.hadoop.hdfs" +
                      '.DistributedFileSystem");
            String[] path = new String[] {"input", "output"};
            Job job = Job.getInstance(conf, "Merge and duplicate removal");
            job.setJarByClass(Merger.class); // name
job.setMapperClass(MyMap.class); // Map class
            job.setReducerClass(MyReduce.class); // add Reduce class
job.setOutputKeyClass(Text.class); // set output type
            job.setOutputValueClass(Text.class); // set input type
            FileInputFormat.addInputPath(job, new Path(path[0])); // input file
            FileOutputFormat.setOutputPath(job, new Path(path[1])); // output file
            System.exit(job.waitForCompletion(true) ? 0 : 1);
  public static class MyMap extends Mapper<Object, Text, Text> {
      private static Text text = new Text();
      public void map(Object key, Text value, Context context)
    throws IOException, InterruptedException {
           this.text = value;
           context.write(text, new Text(""));
  }
  public static class MyReduce extends Reducer<Text, Text, Text, Text>{
      public void reduce(Text key, Iterable<Text> values, Context context)
           throws IOException, InterruptedException {
context.write(key, new Text(""));
      }
  }
 ■ Console \( \mathbb{Z} \)
                                                                <terminated> t6_1 [Java Application] /usr/local/jdk1.8/bin/java (Nov 21, 2020 7:11:12 PM)
 log4j:WARN No appenders could be found for logger (org.apache.hadoop.metric
 log4j:WARN Please initialize the log4j system properly.
 log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for mc
```

```
dhw@localhost:/usr/local/hadoop
                                                                           ×
File Edit View Search Terminal Help
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls output
Found 2 items
-rw-r--r--
            3 dhw supergroup
                                       0 2020-11-21 19:11 output/ SUCCESS
                                     108 2020-11-21 19:11 output/part-r-00000
-rw-r--r--
            3 dhw supergroup
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat output/part-r-00000
20180101 x
20180101 y
20180102 y
20180103 x
20180104 y
20180104 z
20180105 y
20180105 z
20180106 x
[dhw@localhost hadoop]$
```

(二) 实现对输入文件的排序

创建目录 input2 用来存3个文件,同上创建文件

```
[dhw@localhost hadoop]$ ./bin/hdfs dfs -mkdir -p /usr/dhw
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls .
Found 2 items
drwxr-xr-x - dhw supergroup
                                      0 2020-11-21 19:10 input
drwxr-xr-x - dhw supergroup
                                      0 2020-11-21 19:11 output
[dhw@localhost hadoop]$ ./bin/hdfs dfs -mkdir -p /usr/dhw input2
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls .
Found 3 items
drwxr-xr-x - dhw supergroup
                                     0 2020-11-21 19:10 input
           - dhw supergroup
drwxr-xr-x
                                    0 2020-11-21 19:23 input2
drwxr-xr-x

    dhw supergroup

                                     0 2020-11-21 19:11 output
[dhw@localhost hadoop]$
```

若出现如下错误则:

```
at org.apache.hadoop.fs.shell.Command.run(Command.java:165)
at org.apache.hadoop.fs.FsShell.run(FsShell.java:287)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:70)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:84)
at org.apache.hadoop.fs.FsShell.main(FsShell.java:340)
put: Checksum error: file:/home/dhw/myfile/t1.txt at 0 exp: -1197748375 got: -44
4126375
[dhw@localhost hadoop]$
```

```
dhw@localhost:~/myfile _ _ _ x

File Edit View Search Terminal Help

[dhw@localhost myfile]$ ls

A.txt B.txt t1.txt t2.txt t3.txt
[dhw@localhost myfile]$ ls -a

. . . A.txt B.txt .t1_temp.txt.crc t1.txt .t1.txt.crc t2.txt t3.txt
[dhw@localhost myfile]$ rm .t1_temp.txt.crc
[dhw@localhost myfile]$ rm .t1.txt.crc
[dhw@localhost myfile]$ ls -a

. . . A.txt B.txt t1.txt t2.txt t3.txt
[dhw@localhost myfile]$
```

然后重新上传

```
[dhw@localhost hadoop]$ ./bin/hdfs dfs -put /home/dhw/myfile/t1.txt input2
[dhw@localhost hadoop]$ ./bin/hdfs dfs -put /home/dhw/myfile/t2.txt input2
[dhw@localhost hadoop]$ ./bin/hdfs dfs -put /home/dhw/myfile/t3.txt input2
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls input2/t1.txt
-rw-r--r-- 1 dhw supergroup
                                      12 2020-11-21 19:30 input2/t1.txt
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat input2/t1.txt
37
12
40
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat input2/t2.txt
16
39
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat input2/t3.txt
1
45
25
[dhw@localhost hadoop]$
```

```
🕖 t6_2.java 🛭
                                                                                    _
   public class t6_2 {
      public static void main(String[] args) throws Exception {
          String[] path = new String[] {"input2", "output2"};
          String[] other = new GenericOptionsParser(conf, path).getRemainingArgs();
          Job job = Job.getInstance(conf, "mergesort");
job.setJarByClass(MergeSort.class); // name
job.setMapperClass(MyMap.class); // Map class
          job.setReducerClass(MyReduce.class); // add Reduce class
          job.setOutputKeyClass(IntWritable.class); // set output type
          job.setOutputValueClass(IntWritable.class); // set input type
          FileInputFormat.addInputPath(job, new Path(other[0]));
          FileOutputFormat.setOutputPath(job, new Path(other[1])); // output file
          System.exit(job.waitForCompletion(true) ? 0 : 1);
    public static class MyMap extends Mapper<Object, Text, IntWritable, IntWritable> {
        private static IntWritable data = new IntWritable();
        String line = value.toString();
            this.data.set(Integer.parseInt(line));
            context.write(data, new IntWritable(1));
        }
    }
    private static IntWritable linenum = new IntWritable(1);
        public void reduce(IntWritable key, Iterable<IntWritable> values, Context context)
                throws IOException, InterruptedException {
            for (IntWritable num : values) {
                context.write(linenum, key);
                linenum = new IntWritable(linenum.get() + 1);
        }
    }
```

```
🖳 Console 🖾 🛚
 <terminated> t6_2 [Java Application] /usr/local/jdk1.8/bin/java (Nov 21, 2020 7:39:44 PM)
 log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.S
 log4j:WARN Please initialize the log4j system properly.
 log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for mc
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls .
Found 4 items
drwxr-xr-x - dhw supergroup
                                     0 2020-11-21 19:26 input
           - dhw supergroup
drwxr-xr-x
                                     0 2020-11-21 19:31 input2
drwxr-xr-x

    dhw supergroup

                                    0 2020-11-21 19:11 output
drwxr-xr-x - dhw supergroup
                                    0 2020-11-21 19:40 output2
[dhw@localhost hadoop]$ ./bin/hdfs dfs -ls output2
Found 2 items
-rw-r--r-- 3 dhw supergroup
                                     0 2020-11-21 19:40 output2/_SUCCESS
           3 dhw supergroup
-rw-r--r--
                                   54 2020-11-21 19:40 output2/part-r-00000
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat output2/part-r-00000
        5
3
4
       12
5
       16
6
        25
        33
8
        37
9
        39
10
       40
11
       45
[dhw@localhost hadoop]$
```

(三) 对给定的表格进行信息挖掘

查看准备文件

```
dhw@localhost:/usr/local/hadoop
                                                                        File Edit View Search Terminal Help
[dhw@localhost hadoop]$ ./bin/hdfs dfs -mkdir -p /usr/dhw input3
[dhw@localhost hadoop]$ ./bin/hdfs dfs -put /home/dhw/myfile/child parent.txt in
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat input3/child parent.txt
child parent
Steven Lucy
Steven Jack
Jone Lucy
Jone Jack
Lucy Mary
Lucy Frank
Jack Alice
Jack Jesse
David Alice
David Jesse
Philip David
Philip Alma
Mark David
Mark Alma
[dhw@localhost hadoop]$
```

```
↓
1 t6_3.java 

□
1 t6_3.
                                                                                                                                                                                                                                                                   _ _
         public class t6_3 {
                     public static void main(String[] args) throws Exception {
                               Configuration conf = new Configuration();
conf.set("fs.defaultFS", "hdfs://localhost:9000");
conf.set("fs.hdfs.impl", "org.apache.hadoop.hdfs" -
".DistributedFileSystem");
                                String[] path = new String[] {"input3", "output3"};
                                Job job = Job.getInstance(conf, "Single table join");
job.setJarByClass(t6_3.class); // name
job.setMapperClass(MyMap.class); // Map class
job.setReducerClass(MyReduce.class); // add Reduce class
                                job.setOutputKeyClass(Text.class); // set output type
                                job.setOutputValueClass(Text.class); // set input type
                                \label{linear_putPath} File InputFormat. \textit{addInputPath}(job, \ \textit{new} \ Path(path[0])); \ \textit{// input file}
                                FileOutputFormat.setOutputPath(job, new Path(path[1])); // output file
                                System.exit(job.waitForCompletion(true) ? 0 : 1);
                public static int time = 0;
                public static class MyMap extends Mapper<Object, Text, Text, Text> {
                            public void map(Object key, Text value, Context context)
                                                   throws IOException, InterruptedException {
                                       String line = value.toString();
                                       String[] childAndParent = line.split(" ");
                                       List<String> list = new ArrayList<String>(2);
                                       for (String childOrParent : childAndParent) {
   if (! "".equals(childOrParent)) {
                                                              list.add(childOrParent);
                                       }
                                       if (!"child".equals(list.get(0))) {
                                                   String childName = list.get(0);
                                                   String parentName = list.get(1);
                                                   String relationType = "1";
                                                   context.write(new Text(parentName), new Text(relationType + "+"
                                                                         + childName + "+" + parentName));
                                                   relationType = "2";
                                                   context.write(new Text(childName), new Text(relationType + "+"
                                                                         + childName + "+" + parentName));
                                       }
```

```
public static class MyReduce extends Reducer<Text, Text, Text, Text> {
           public void reduce(Text key, Iterable<Text> values, Context context)
                   throws IOException, InterruptedException {
               if (time == 0) {
                   context.write(new Text("grand_child"), new Text("grand_parent"));
                   time ++;
               }
               List<String> grandChild = new ArrayList<String>();
               List<String> grandParent = new ArrayList<String>();
               for (Text text : values) {
                   String s = text.toString();
                   String[] relation = s.split("\\+");
String relationType = relation[0];
                   String childName = relation[1];
                   String parentName = relation[2];
                   if ("1".equals(relationType)) {
                       grandChild.add(childName);
                     else {
                       grandParent.add(parentName);
               int grandParentNum = grandParent.size();
             int grandParentNum = grandParent.size();
             int grandChildNum = grandChild.size();
             if (grandParentNum != 0 && grandChildNum != 0) {
                 for (int m = 0; m < grandChildNum; m++) {</pre>
                     for (int n = 0; n < grandParentNum; n++) {</pre>
                          context.write(new Text(grandChild.get(m)),
                                  new Text(grandParent.get(n)));
                }
           }
       }
    }
}
```

实验结果:

```
[dhw@localhost hadoop]$ ./bin/hdfs dfs -cat output3/part-r-00000
grand child
                grand_parent
Mark
       Jesse
Mark
        Alice
Philip Jesse
Philip Alice
Jone
        Jesse
Jone
       Alice
Steven Jesse
Steven Alice
Steven Frank
Steven Mary
        Frank
Jone
        Mary
Jone
[dhw@localhost hadoop]$
```