



大数据技术基础课程实验报告

实验三：基于华为云安装 HBase、ZooKeeper 及 HBase 应用实践

付容天

学号 2020211616

班级 2020211310

计算机学院（国家示范性软件学院）

2023 年 4 月 3 日

1. 实验描述与实验目的

在本次实验中，我们需要在实验一、二搭建完毕的集群环境上，继续安装 HBase、ZooKeeper，并实践 HBase 的基本使用。

本次实验的目的是：掌握 HBase、ZooKeeper 等组件的安装与使用，以及使用 MapReduce 批量将 HBase 表上的数据导入到 HDFS 中，学习本实验能快速掌握 HBase 数据库在分布式计算中的应用，理解 Java API 读取 HBase 数据机制等相关内容。

2. 实验过程与实验分析

这部分的实验中，我按照实验指导书，先后完成了：

- (1) 实验环境的检查，保证实验一、二搭建的集群环境运行正常；
- (2) 下载、安装、配置 ZooKeeper 组件，包括：建立软链接、配置文件修改、修改数据目录、ZooKeeper 的拷贝等内容；
- (3) 在四个节点上启动 ZooKeeper，检查是否启动成功，并通过-status 选项查看各节点是 follower 还是 leader；
- (4) 下载、安装、配置 HBase 组件，包括：建立软链接、环境变量修改、配置文件修改、HBase 的拷贝等内容；
- (5) 在 node1 上启动 HBase，进入 HBase Shell，使用 create 和 put 命令创建了如下图所示的数据库表格：

```
root@frt-2020211616-0001:/
Took 0.0086 seconds
hbase(main):001:0> create '2020211616_frt','cf1'
Created table 2020211616_frt
Took 2.0826 seconds
=> Hbase::Table - 2020211616_frt
hbase(main):002:0> put '2020211616_frt','2020211616_frt-001','cf1:keyword','applicate'
Took 0.2319 seconds
hbase(main):003:0> put '2020211616_frt','2020211616_frt-002','cf1:keyword','OnePlus 5'
Took 0.0092 seconds
hbase(main):004:0> put '2020211616_frt','2020211616_frt-003','cf1:keyword','iphone 6s'
Took 0.0215 seconds
hbase(main):005:0> scan '2020211616_frt'
ROW                                COLUMN+CELL
2020211616_frt-001                 column=cf1:keyword, timestamp=1680425252713, value=applicate
2020211616_frt-002                 column=cf1:keyword, timestamp=1680425280347, value=OnePlus 5
2020211616_frt-003                 column=cf1:keyword, timestamp=1680425299969, value=iphone 6s
3 row(s)
Took 0.0593 seconds
hbase(main):006:0>
```

图 1：数据库表格（学号 2020211616，名字缩写 frt）

上图 1 中显示了数据库表格的结构，共有三行，且三行分别命名（即 row key）为 2020211616_frt-001、2020211616_frt-002、2020211616_frt-003，并且定义了列族 cf1（具有 keyword）。直观来看，这个表格可以表示为：

行键	cf1 (列族)
	keyword
2020211616_frt-001	Applicate
2020211616_frt-002	OnePlus 5
2020211616_frt-003	iphone 6s

表 1: 创建的数据库表格的形式与内容

随后，在 IDEA 中创建名为 MyHBase 的工程，进行下面的操作：

- (6) 配置 pom.xml 文件，添加依赖，注意处理三个缺少依赖 jar 的包；
- (7) 新建名为 org.frt2020211616.hbase.inputSource 的 package；
- (8) 在新建 package 中创建类 MemberMapper，编写相应代码；
- (9) 在新建 package 中创建类 Main，编写相应代码。

下面展示了两个类的源代码：

```

1 package org.frt2020211616.hbase.inputSource;
2
3 import org.apache.hadoop.hbase.Cell;
4 import org.apache.hadoop.hbase.client.Result;
5 import org.apache.hadoop.hbase.io.ImmutableBytesWritable;
6 import org.apache.hadoop.hbase.mapreduce.TableMapper;
7 import org.apache.hadoop.hbase.util.Bytes;
8 import org.apache.hadoop.io.Writable;
9 import org.apache.hadoop.io.Text;
10 import java.io.IOException;
11
12 usage

```

图 2: MemberMapper 类的头文件

```

13 usage
14 public class MemberMapper extends TableMapper<Writable, Writable> {
15     2 usages
16     private Text k = new Text();
17     2 usages
18     private Text v = new Text();
19     3 usages
20     public static final String FIELD_COMMON_separator = "\u0001";
21     @Override
22     protected void setup(Context context) throws IOException, InterruptedException {}
23     @Override
24     protected void map(ImmutableBytesWritable row, Result columns,
25         Context context) throws IOException, InterruptedException {
26         String value = null;
27         String rowkey = new String(row.get());
28         byte[] columnFamily = null;
29         byte[] columnQualifier = null;
30         long ts = 0L;
31         try{
32             for(Cell cell : columns.listCells()){
33                 value = Bytes.toStringBinary(cell.getFamilyArray());
34                 columnFamily = cell.getQualifierArray();
35                 ts = cell.getTimestamp();
36                 k.set(rowkey);
37                 v.set(Bytes.toString(columnFamily)+FIELD_COMMON_separator+Bytes.toString(columnQualifier)
38                     +FIELD_COMMON_separator+value+FIELD_COMMON_separator+ts);
39                 context.write(k, v);
40             }
41         } catch (Exception e) {
42             e.printStackTrace();
43             System.err.println("Error:"+e.getMessage()+" ,Row:"+Bytes.toString(row.get())+"Value:"+value);
44         }
45     }
46 }
47 }

```

图 3: MemberMapper 类的源代码

```

1 package org.frt2020211616.hbase.inputSource;
2
3 import org.apache.commons.logging.Log;
4 import org.apache.commons.logging.LogFactory;
5 import org.apache.hadoop.conf.Configuration;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.fs.Path;
8 import org.apache.hadoop.fs.FileSystem;
9 import org.apache.hadoop.hbase.HBaseConfiguration;
10 import org.apache.hadoop.hbase.client.Scan;
11 import org.apache.hadoop.hbase.util.Bytes;
12 import org.apache.hadoop.hbase.mapreduce.TableMapReduceUtil;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14 import org.apache.hadoop.mapreduce.Job;
15 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
16
17 3 usages

```

图 4: Main 类的头文件

```

17 3 usages
18 public class Main {
19     static final Log LOG = LogFactory.getLog(Main.class);
20     public static final String NAME = "Member Test1";
21     public static final String TEMP_INDEX_PATH = "hdfs://node1:8020/tmp/2020211616_frt";
22     public static String inputTable = "2020211616_frt";
23
24     public static void main(String[] args) throws Exception {
25         Configuration conf = HBaseConfiguration.create();
26         Scan scan = new Scan();
27         scan.setBatch(0);
28         scan.setCaching(10000);
29         scan.setMaxVersions();
30         scan.setTimeRange(System.currentTimeMillis() - 3*24*3600*1000L, System.currentTimeMillis());
31         scan.addColumn(Bytes.toBytes("cf1"), Bytes.toBytes("keyword"));
32
33         conf.setBoolean("mapred.map.tasks.speculative.execution", false);
34         conf.setBoolean("mapred.reduce.tasks.speculative.execution", false);
35         Path tmpIndexPath = new Path(TEMP_INDEX_PATH);
36         FileSystem fs = FileSystem.get(conf);

```

图 5: Main 类的源代码

```

37         if(fs.exists(tmpIndexPath)) {
38             fs.delete(tmpIndexPath, true);
39         }
40
41         Job job = new Job(conf, NAME);
42         job.setJarByClass(Main.class);
43
44         TableMapReduceUtil.initTableMapperJob(inputTable, scan, MemberMapper.class, Text.class, Text.class, job);
45         job.setNumReduceTasks(0);
46         job.setOutputFormatClass(TextOutputFormat.class);
47         FileOutputFormat.setOutputPath(job, tmpIndexPath);
48
49         boolean success = job.waitForCompletion(true);
50         System.exit(success ? 0 : 1);
51     }
52 }

```

图 6: Main 类的源代码（续）

下面对代码进行简单的分析。图 2 展示了 MemberMapper 类依赖的头文件，主要来自 org.apache.hadoop；图 3 展示了 MemberMapper 类的源代码，该类扩展自 Mapper 类，以 HBase 中的表作为输入源进行相应的处理：图中 21 到 25 行定义了用到的变量，26 到 35 行为主要的处理逻辑（通过 for 循环读入 columns 中的元素，处理后写入到 context 中）。

图 4 展示了 Main 类依赖的头文件，主要来自 org.apache.commons.logging 和 org.apache.hadoop 两处。图 5 和图 6 展示了 Main 类的源代码，其中 18-21

