

HBase

第一章 简介

Hadoop的局限

HBase 是一个构建在 Hadoop 文件系统之上的面向列的数据库管理系统。

要想明白为什么产生 HBase，就需要先了解一下 Hadoop 存在的限制？Hadoop 可以通过 HDFS 来存储结构化、半结构甚至非结构化的数据，它是传统数据库的补充，是海量数据存储的最佳方法，它针对大文件的存储，批量访问和流式访问都做了优化，同时也通过多副本解决了容灾问题。

但是 Hadoop 的缺陷在于它只能执行批处理，并且只能以顺序方式访问数据，这意味着即使是最简单的工作，也必须搜索整个数据集，无法实现对数据的随机访问。实现数据的随机访问是传统的关系型数据库所擅长的，但它们却不能用于海量数据的存储。在这种情况下，必须有一种新的方案来解决海量数据存储和随机访问的问题，HBase 就是其中之一 (HBase, Cassandra, couchDB, Dynamo 和 MongoDB 都能存储海量数据并支持随机访问)。

注：数据结构分类：

- 结构化数据：即以关系型数据库表形式管理的数据；
- 半结构化数据：非关系模型的，有基本固定结构模式的数据，例如日志文件、XML 文档、JSON 文档、Email 等；
- 非结构化数据：没有固定模式的数据，如 WORD、PDF、PPT、EXL，各种格式的图片、视频等。

HBase简介

HBase 是一个构建在 Hadoop 文件系统之上的面向列的数据库管理系统。

HBase 是一种类似于 **Google's Big Table** 的数据模型，它是 Hadoop 生态系统的一部分，它将数据存储在 HDFS 上，客户端可以通过 HBase 实现对 HDFS 上数据的随机访问。它具有以下特性：

- 不支持复杂的事务，只支持行级事务，即单行数据的读写都是原子性的；
- 由于是采用 HDFS 作为底层存储，所以和 HDFS 一样，支持结构化、半结构化和非结构化的存储；
- 支持通过增加机器进行横向扩展；
- 支持数据分片；
- 支持 RegionServers 之间的自动故障转移；
- 易于使用的 Java 客户端 API；
- 支持 BlockCache 和布隆过滤器；
- 过滤器支持谓词下推。

HBase Table

HBase 是一个面向 **列** 的数据库管理系统，这里更为确切的而说，HBase 是一个面向 **列族** 的数据库管理系统。表 schema 仅定义列族，表具有多个列族，每个列族可以包含任意数量的列，列由多个单元格（cell）组成，单元格可以存储多个版本的数据，多个版本数据以时间戳进行区分。

下图为 HBase 中一张表的：

- RowKey 为行的唯一标识，所有行按照 RowKey 的字典序进行排序；
- 该表具有两个列族，分别是 personal 和 office；
- 其中列族 personal 拥有 name、city、phone 三个列，列族 office 拥有 tel、address 两个列。

Hbase 的表具有以下特点：

- 容量大：一个表可以有数十亿行，上百万列；
- 面向列：数据是按照列存储，每一列都单独存放，数据即索引，在查询时可以只访问指定列的数据，有效地降低了系统的 I/O 负担；
- 稀疏性：空 (null) 列并不占用存储空间，表可以设计的非常稀疏；
- 数据多版本：每个单元中的数据可以有多个版本，按照时间戳排序，新的数据在最上面；
- 存储类型：所有数据的底层存储格式都是字节数组 (byte[])。

Phoenix

Phoenix 是 HBase 的开源 SQL 中间层，它允许你使用标准 JDBC 的方式来操作 HBase 上的数据。在 **Phoenix** 之前，如果你要访问 HBase，只能调用它的 Java API，但相比于使用一行 SQL 就能实现数据查询，HBase 的 API 还是过于复杂。**Phoenix** 的理念是 **we put sql SQL back in NOSQL**，即你可以使用标准的 SQL 就能完成对 HBase 上数据的操作。同时这也意味着你可以通过集成 **Spring Data JPA** 或 **Mybatis** 等常用的持久层框架来操作 HBase。

其次 **Phoenix** 的性能表现也非常优异，**Phoenix** 查询引擎会将 SQL 查询转换为一个或多个 HBase Scan，通过并行执行来生成标准的 JDBC 结果集。它通过直接使用 HBase API 以及协处理器和自定义过滤器，可以为小型数据查询提供毫秒级的性能，为千万行数据的查询提供秒级的性能。同时 Phoenix 还拥有二级索引等 HBase 不具备的特性，因为以上的优点，所以 **Phoenix** 成为了 HBase 最优秀的 SQL 中间层。

第二章 基本概念

Row Key

Row Key 是用来检索记录的主键。想要访问 HBase Table 中的数据，只有以下三种方式：

- 通过指定的 **Row Key** 进行访问；
- 通过 Row Key 的 range 进行访问，即访问指定范围内的行；
- 进行全表扫描。

Row Key 可以是任意字符串，存储时数据按照 **Row Key** 的字典序进行排序。这里需要注意以下两点：

- 因为字典序对 Int 排序的结果是
1,10,100,11,12,13,14,15,16,17,18,19,2,20,21,...,9,91,92,93,94,95,96,97,98,99。如果你使用整型的字符串作为行键，那么为了保持整型的自然序，行键必须用 0 作左填充。

- 行的一次读写操作时原子性的 (不论一次读写多少列)。

Column Family (列族)

HBase 表中的每个列，都归属于某个列族。列族是表的 Schema 的一部分，所以列族需要在创建表时进行定义。列族的所有列都以列族名作为前缀，例如 `courses:history`，`courses:math` 都属于 `courses` 这个列族。

Column Qualifier (列限定符)

列限定符，你可以理解为是具体的列名，例如 `courses:history`，`courses:math` 都属于 `courses` 这个列族，它们的列限定符分别是 `history` 和 `math`。需要注意的是列限定符不是表 Schema 的一部分，你可以在插入数据的过程中动态创建列。

Column(列)

HBase 中的列由列族和列限定符组成，它们由 `:` (冒号) 进行分隔，即一个完整的列名应该表述为 `列族名:列限定符`。

Cell

`Cell` 是行，列族和列限定符的组合，并包含值和时间戳。你可以等价理解为关系型数据库中由指定行和指定列确定的一个单元格，但不同的是 HBase 中的一个单元格是由多个版本的数据组成的，每个版本的数据用时间戳进行区分。

Timestamp(时间戳)

HBase 中通过 `row key` 和 `column` 确定的为一个存储单元称为 `Cell`。每个 `Cell` 都保存着同一份数据的多个版本。版本通过时间戳来索引，时间戳的类型是 64 位整型，时间戳可以由 HBase 在数据写入时自动赋值，也可以由客户显式指定。每个 `Cell` 中，不同版本的数据按照时间戳倒序排列，即最新的数据排在最前面。

第三章 常用Shell命令

基本命令

打开 Hbase Shell:

```
1 # hbase shell
```

获取帮助

```
1  # 获取帮助
2  help
3  # 获取命令的详细信息
4  help 'status'
```

查看服务器状态

```
1  status
```

查看版本信息

```
1  version
```

关于表的操作

查看所有表

```
1  list
```

创建表

命令格式： create '表名称', '列族名称 1','列族名称 2','列名称 N'

```
1  # 创建一张名为Student的表, 包含基本信息 (baseInfo)、学校信息 (schoolInfo) 两个列族
2  create 'Student', 'baseInfo', 'schoolInfo'
```

查看表的基本信息

命令格式： desc '表名'

```
1  describe 'Student'
```

表的启用/禁用

enable 和 disable 可以启用/禁用这个表,is_enabled 和 is_disabled 来检查表是否被禁用

```
1  # 禁用表
2  disable 'Student'
3  # 检查表是否被禁用
4  is_disabled 'Student'
5  # 启用表
6  enable 'Student'
7  # 检查表是否被启用
8  is_enabled 'Student'
```

检查表是否存在

```
1  exists 'Student'
```

删除表

```
1  # 删除表前需要先禁用表
2  disable 'Student'
3  # 删除表
4  drop 'Student'
```

增删改

添加列族

命令格式： alter '表名', '列族名'

```
1  alter 'Student', 'teacherInfo'
```

删除列族

命令格式： alter '表名', {NAME => '列族名', METHOD => 'delete'}

```
1  alter 'Student', {NAME => 'teacherInfo', METHOD => 'delete'}
```

更改列族存储版本的限制

默认情况下，列族只存储一个版本的数据，如果需要存储多个版本的数据，则需要修改列族的属性。修改后可通过 `desc` 命令查看。

```
1  alter 'Student', {NAME=>'baseInfo', VERSIONS=>3}
```

插入数据

命令格式： put '表名', '行键','列族:列','值'

注意：如果新增数据的行键值、列族名、列名与原有数据完全相同，则相当于更新操作

```
1  put 'Student', 'rowkey1', 'baseInfo:name', 'tom'
2  put 'Student', 'rowkey1', 'baseInfo:birthday', '1990-01-09'
3  put 'Student', 'rowkey1', 'baseInfo:age', '29'
4  put 'Student', 'rowkey1', 'schoolInfo:name', 'Harvard'
5  put 'Student', 'rowkey1', 'schoolInfo:location', 'Boston'
6
7  put 'Student', 'rowkey2', 'baseInfo:name', 'jack'
8  put 'Student', 'rowkey2', 'baseInfo:birthday', '1998-08-22'
9  put 'Student', 'rowkey2', 'baseInfo:age', '21'
10 put 'Student', 'rowkey2', 'schoolInfo:name', 'Yale'
11 put 'Student', 'rowkey2', 'schoolInfo:location', 'New Haven'
12
13 put 'Student', 'rowkey3', 'baseInfo:name', 'maike'
14 put 'Student', 'rowkey3', 'baseInfo:birthday', '1995-01-22'
15 put 'Student', 'rowkey3', 'baseInfo:age', '24'
16 put 'Student', 'rowkey3', 'schoolInfo:name', 'Yale'
17 put 'Student', 'rowkey3', 'schoolInfo:location', 'New Haven'
18
19 put 'Student', 'rowkey4', 'baseInfo:name', 'maike-jack'
```

获取指定行、指定行中的列族、列的信息

```
1  # 获取指定行中所有列的数据信息
2  get 'Student','rowkey3'
3  # 获取指定行中指定列族下所有列的数据信息
4  get 'Student','rowkey3','baseInfo'
5  # 获取指定行中指定列的数据信息
6  get 'Student','rowkey3','baseInfo:name'
```

删除指定行、指定行中的列

```
1  # 删除指定行
2  deleteall 'Student','rowkey3'
3  # 删除指定行中指定列的数据
4  delete 'Student','rowkey3','baseInfo:name'
```

查询

hbase 中访问数据有两种基本的方式：

- 按指定 rowkey 获取数据：get 方法；
- 按指定条件获取数据：scan 方法。

scan 可以设置 begin 和 end 参数来访问一个范围内所有的数据。get 本质上就是 begin 和 end 相等的一种特殊的 scan。

Get查询

```
1  # 获取指定行中所有列的数据信息
2  get 'Student','rowkey3'
3  # 获取指定行中指定列族下所有列的数据信息
4  get 'Student','rowkey3','baseInfo'
5  # 获取指定行中指定列的数据信息
6  get 'Student','rowkey3','baseInfo:name'
```

查询整表数据

```
1  scan 'Student'
```

查询指定列簇的数据

```
1  scan 'Student', {COLUMN=>'baseInfo'}
```

条件查询

```
1  # 查询指定列的数据
2  scan 'Student', {COLUMNS=> 'baseInfo:birthday'}
```

除了列（COLUMNS）修饰词外，HBase 还支持 Limit（限制查询结果行数），STARTROW（ROWKEY 起始行，会先根据这个 key 定位到 region，再向后扫描）、STOPROW（结束行）、TIMERANGE（限定时间戳范围）、VERSIONS（版本号）、和 FILTER（按条件过滤行）等。

如下代表从 rowkey2 这个 rowkey 开始，查找下两个行的最新 3 个版本的名字列的数据：

```
1 scan 'Student', {COLUMNS=> 'baseInfo:name', STARTROW => 'rowkey2', STOPROW => 'wrowkey4', LIMIT=>2,
  VERSIONS=>3}
```

条件过滤

Filter 可以设定一系列条件来进行过滤。如我们要查询值等于 24 的所有数据：

```
1 scan 'Student', FILTER=>"ValueFilter(=, 'binary:24')"
```

值包含 yale 的所有数据：

```
1 scan 'Student', FILTER=>"ValueFilter(=, 'substring:yale')"
```

列名中的前缀为 birth 的：

```
1 scan 'Student', FILTER=>"ColumnPrefixFilter('birth')"
```

FILTER 中支持多个过滤条件通过括号、AND 和 OR 进行组合：

```
1 # 列名中的前缀为birth且列值中包含1998的数据
2 scan 'Student', FILTER=>"ColumnPrefixFilter('birth') AND ValueFilter
  ValueFilter(=, 'substring:1998')"
```

PrefixFilter 用于对 Rowkey 的前缀进行判断：

```
1 scan 'Student', FILTER=>"PrefixFilter('wr')"
```

第四章 HBase的JavaAPI

1 单例连接

```
1 import org.apache.flink.table.filesystem.FileSystemOutputFormat;
2 import org.apache.hadoop.conf.Configuration;
3 import org.apache.hadoop.hbase.client.AsyncConnection;
4 import org.apache.hadoop.hbase.client.Connection;
5 import org.apache.hadoop.hbase.client.ConnectionFactory;
6
7 import java.io.IOException;
8 import java.util.concurrent.CompletableFuture;
9
10 /**
11  * Project:  BigDataCode
12  * Create date:  2023/6/20
13  * Created by fujiahao
14  */
15
16 /**
17  * 单线程使用连接
18  */
19
20 public class HBaseConnection {
21
22     // 声明静态属性
23     public static Connection connection = null;
```

```

24     static {
25         // 创建连接
26         // 默认使用同步连接
27         try {
28             // 使用读取本地文件的形式添加参数 hbase-site.xml
29             connection = ConnectionFactory.createConnection();
30         } catch (IOException e) {
31             e.printStackTrace();
32         }
33     }
34
35     public static void closeConnection() throws IOException {
36         // 判断连接是否为空值
37         if (connection != null) {
38             connection.close();
39         }
40     }
41
42     public static void main(String[] args) throws IOException {
43
44         /*// 创建连接配置对象
45         Configuration conf = new Configuration();
46         conf.set("hbase.zookeeper.quorum", "master,slave1,slave2");
47
48         // 创建连接
49         // 默认使用同步连接
50         Connection connection = ConnectionFactory.createConnection(conf);
51
52         // 可以使用异步连接
53         // 不推荐
54         // CompletableFuture<AsyncConnection> asyncConnection =
55         ConnectionFactory.createAsyncConnection(conf);
56
57         // 使用连接
58         System.out.println(connection);
59
60         // 关闭连接
61         connection.close();*/
62
63         // 直接使用创建好的链接
64         // 不要在main线程里面单独创建
65         System.out.println(HBaseConnection.connection);
66
67         // 在main线程的最后记得关闭连接
68         HBaseConnection.closeConnection();
69     }
70

```

2 插入数据

```

1     /**
2         * 插入数据
3         * @param namespace 命名空间
4         * @param tableName 表名
5         * @param rowKey 行号
6         * @param columnFamily 列族

```



```

7      * @param columnName 列明
8      * @param value 值
9      */
10     public static void putCell(String namespace, String tableName, String rowKey, String
columnFamily, String columnName, String value) throws IOException {
11         // 1. 获取table
12         Table table = connection.getTable(TableName.valueOf(namespace, tableName));
13
14         // 2. 创建put对象
15         Put put = new Put(Bytes.toBytes(rowKey));
16
17         // 3. 添加列
18         put.addColumn(Bytes.toBytes(columnFamily), Bytes.toBytes(columnName),
Bytes.toBytes(value));
19         try {
20             table.put(put);
21         } catch (IOException e) {
22             e.printStackTrace();
23         }
24
25         // 关闭table
26         table.close();
27     }

```

3 创建表格

```

1     /**
2      * 创建表格
3      * @param namespace 命名空间名称
4      * @param tableName 表格名称
5      * @param columnFamilies 列族名称(可以有多个)
6      */
7     public static void createTable(String namespace, String tableName, String... columnFamilies)
throws IOException {
8         // 判断是否有至少一个列族
9         if (columnFamilies.length == 0) {
10             System.out.println("创建表格至少有一个列族");
11             return;
12         }
13
14         // 判断表格是否存在
15         if (isTableExists(namespace, tableName)) {
16             System.out.println("表格已经存在");
17             return;
18         }
19
20         // 1. 获取admin
21         Admin admin = connection.getAdmin();
22
23         // 2. 创建表格
24         // 2.1 创建表格描述 => 建造者模式
25         TableDescriptorBuilder builder =
TableDescriptorBuilder.newBuilder(TableName.valueOf(namespace, tableName));
26
27         // 2.2 添加参数
28         for (String columnFamily : columnFamilies) {
29

```

```

30         // 2.3 创建列族描述的建造者
31         ColumnFamilyDescriptorBuilder columnFamilyDescriptorBuilder =
ColumnFamilyDescriptorBuilder.newBuilder(Bytes.toBytes(columnFamily));
32
33         // 2.4 对应当前的列族添加参数
34         // 添加版本参数
35         columnFamilyDescriptorBuilder.setMaxVersions(5);
36
37         // 2.5 创建添加完参数的列族描述
38         builder.setColumnFamily(columnFamilyDescriptorBuilder.build());
39     }
40
41     // 2.3 创建对应的表格描述
42     try {
43         admin.createTable(builder.build());
44     } catch (IOException e) {
45         e.printStackTrace();
46     }
47
48     // 3.关闭admin
49     admin.close();
50 }

```

4 创建命名空间

```

1     /**
2      * 创建命名空间
3      * @param namespace 命名空间名称
4      */
5     public static void createNamespace(String namespace) throws IOException {
6         // 1.获取admin
7         // 此处异常先不抛出，最后统一处理
8         // admin连接：轻量级，不是线程安全的，不推荐池化或者缓存该连接
9         Admin admin = connection.getAdmin();
10
11         // 2.调用方法创建命名空间
12         // 代码相对shell更底层，所以shell能够实现的功能，代码一定能实现
13         // 所以需要填写完整的命名空间描述
14
15         // 2.1 创建命名空间描述的建造者 => 设计师
16         NamespaceDescriptor.Builder builder = NamespaceDescriptor.create(namespace);
17
18         // 2.2 给命名空间添加需求
19         builder.addConfiguration("user", "atguigu");
20
21         // 2.3 使用builder构造出对应的添加完参数的对象
22         // 创建命名空间出现的问题，都属于本方法自身问题
23         try {
24             admin.createNamespace(builder.build());
25         } catch (IOException e) {
26             System.out.println("命名空间已经存在");
27             e.printStackTrace();
28         }
29
30         // 3.关闭admin
31         admin.close();
32     }

```

5 带过滤的扫描

```

1  /**
2      * 带过滤的扫描
3      * @param namespace 命名空间
4      * @param tableName 表格名
5      * @param startRow 开始
6      * @param stopRow 结束
7      * @param columnFamily 列族
8      * @param columnName 列明
9      * @param value 值
10     * @throws IOException
11     */
12     public static void filterScan(String namespace, String tableName, String startRow,
13                                   String stopRow, String columnFamily, String
columnName, String value) throws IOException {
14         // 1. 获取table
15         Table table = connection.getTable(TableName.valueOf(namespace, tableName));
16
17         // 2. 创建scan对象
18         Scan scan = new Scan();
19         // 如果直接调用，就扫描整张表
20         // 添加参数来控制扫描数据
21         scan.withStartRow(Bytes.toBytes(startRow));
22         scan.withStopRow(Bytes.toBytes(stopRow));
23
24         // 可以添加多个过滤
25         FilterList filterList = new FilterList();
26
27         // 创建过滤器
28         // (1) 结果只保留当前列的数据
29         ColumnValueFilter columnValueFilter = new ColumnValueFilter(
30             // 列族名称
31             Bytes.toBytes(columnFamily),
32             // 列名
33             Bytes.toBytes(columnName),
34             // 比较关系
35             CompareOperator.EQUAL,
36             // 值
37             Bytes.toBytes(value)
38         );
39
40         // (2) 结果保留整行数据
41         // 结果会同时保留没有当前列的数据
42         // 比如：“bigdata:student”中“rowKey”为“1001”的“info:age”值是18,但是“info:name”为空值
43         // 也会保留该行数据
44         SingleColumnValueFilter singleColumnValueFilter = new SingleColumnValueFilter(
45             // 列族名称
46             Bytes.toBytes(columnFamily),
47             // 列名
48             Bytes.toBytes(columnName),
49             // 比较关系
50             CompareOperator.EQUAL,
51             // 值
52             Bytes.toBytes(value)

```

```

53         );
54
55         // 本身可以添加多个过滤器
56         filterList.addFilter(singleColumnValueFilter);
57
58         // 添加过滤
59         scan.setFilter(filterList);
60
61         try {
62             // 读取多行数据
63             ResultScanner scanner = null;
64             scanner = table.getScanner(scan);
65             // result来记录一行数据
66             // ResultScanner来记录多行数据
67             for (Result result : scanner) {
68                 Cell[] cells = result.rawCells();
69                 for (Cell cell : cells) {
70                     System.out.print(new String(CellUtil.cloneRow(cell)) + "-"
71                                     + new String(CellUtil.cloneFamily(cell)) + "-"
72                                     + new String(CellUtil.cloneQualifier(cell)) + "-"
73                                     + new String(CellUtil.cloneValue(cell)) + "\t");
74                 }
75                 System.out.println();
76             }
77         } catch (IOException e) {
78             e.printStackTrace();
79         }
80
81         table.close();
82     }

```

6 读取数据

```

1  /**
2      * 读取数据
3      * @param namespace 命名空间
4      * @param tableName 表名
5      * @param rowKey 行号
6      * @param columnFamily 列族
7      * @param columnName 列明
8      */
9      public static void getCells(String namespace, String tableName, String rowKey, String
columnFamily, String columnName) throws IOException {
10         // 1. 获取table对象
11         Table table = connection.getTable(TableName.valueOf(namespace, tableName));
12
13         // 2. 创建get对象
14         Get get = new Get(Bytes.toBytes(rowKey));
15
16         // 如果直接调用get方法读取数据，此时读一整行数据
17         // 如果想要读取某一列数据，需要添加对应的参数
18         get.addColumn(Bytes.toBytes(columnFamily), Bytes.toBytes(columnName));
19         // 设置读取数据的版本
20         get.readAllVersions();
21
22         // 读取数据，得到result对象
23         try {

```

```

24         Result result = null;
25         result = table.get(get);
26         // 处理数据
27         Cell[] cells = result.rawCells();
28         for (Cell cell : cells) {
29             // cell存储数据比较底层
30             // 会乱码
31             String value = new String(CellUtil.cloneValue(cell));
32             System.out.println(value);
33         }
34     } catch (IOException e) {
35         e.printStackTrace();
36     }
37
38     table.close();
39 }

```

7 判断表格是否存在

```

1  /**
2      * 判断表格是否存在
3      * @param namespace 命名空间名称
4      * @param tableName 表格名称
5      * @return true表示存在
6      */
7  public static boolean isTableExists(String namespace, String tableName) throws IOException {
8      // 1. 获取admin
9      Admin admin = connection.getAdmin();
10
11      // 2. 判断表格是否存在
12      boolean b = false;
13      try {
14          b = admin.tableExists(TableNames.valueOf(namespace, tableName));
15      } catch (IOException e) {
16          e.printStackTrace();
17      }
18
19      // 3. 关闭连接
20      admin.close();
21
22      // 4. 返回结果
23      return b;
24  }

```

8 扫描数据

```

1  /**
2      * 扫描数据
3      * @param namespace 命名空间
4      * @param tableName 表名
5      * @param startRow 开始行(包含)
6      * @param stopRow 结束行(不包含)
7      */
8  public static void scanRows(String namespace, String tableName, String startRow, String
stopRow) throws IOException {
9      // 1. 获取table

```

```

10         Table table = connection.getTable(TableName.valueOf(namespace, tableName));
11
12         // 2. 创建scan对象
13         Scan scan = new Scan();
14         // 如果直接调用，就扫描整张表
15         // 添加参数来控制扫描数据
16         scan.withStartRow(Bytes.toBytes(startRow));
17         scan.withStopRow(Bytes.toBytes(stopRow));
18
19         try {
20             // 读取多行数据
21             ResultScanner scanner = null;
22             scanner = table.getScanner(scan);
23             // result来记录一行数据
24             // ResultScanner来记录多行数据
25             for (Result result : scanner) {
26                 Cell[] cells = result.rawCells();
27                 for (Cell cell : cells) {
28                     System.out.print(new String(CellUtil.cloneRow(cell)) + "-"
29                                     + new String(CellUtil.cloneFamily(cell)) + "-"
30                                     + new String(CellUtil.cloneQualifier(cell)) + "-"
31                                     + new String(CellUtil.cloneValue(cell)) + "\t");
32                 }
33                 System.out.println();
34             }
35         } catch (IOException e) {
36             e.printStackTrace();
37         }
38
39         table.close();
40     }

```

9 删除表格

```

1     /**
2         * 删除表格
3         * @param namespace 命名空间
4         * @param tableName 表名
5         * @return 是否删除
6         */
7     public static boolean deleteTable(String namespace, String tableName) throws IOException {
8         // 判断表格是否存在
9         if (isTableExists(namespace, tableName)) {
10             System.out.println("表格不存在，无法删除");
11             return false;
12         }
13
14         // 1. 获取admin
15         Admin admin = connection.getAdmin();
16
17         // 2. 删除表格
18         try {
19             // HBase删除表格之前一定要先标记表格为不可用
20             admin.disableTable(TableName.valueOf(namespace, tableName));
21             admin.deleteTable(TableName.valueOf(namespace, tableName));
22         } catch (IOException e) {
23             e.printStackTrace();

```

```

24         }
25
26         // 3. 关闭连接
27         admin.close();
28
29         return true;
30     }

```

10 删除一行中的一列数据

```

1  /**
2      * 删除一行中的一列数据
3      * @param namespace 命名空间
4      * @param tableName 表名
5      * @param rowKey 行号
6      * @param columnFamily 列族
7      * @param columnName 列名
8      */
9      public static void deleteColumn(String namespace, String tableName, String rowKey, String
columnFamily, String columnName) throws IOException {
10         Table table = connection.getTable(TableName.valueOf(namespace, tableName));
11
12         Delete delete = new Delete(Bytes.toBytes(rowKey));
13         // 删除一个版本
14         delete.addColumn(Bytes.toBytes(columnFamily), Bytes.toBytes(columnName));
15         // 删除所有版本
16         delete.addColumns(Bytes.toBytes(columnFamily), Bytes.toBytes(columnName));
17
18         try {
19             table.delete(delete);
20         } catch (IOException e) {
21             e.printStackTrace();
22         }
23
24         table.close();
25     }

```

11 修改表格

```

1  /**
2      * 修改表格
3      * @param namespace 命名空间
4      * @param tableName 表名
5      * @param columnFamily 列族
6      * @param version 版本
7      */
8      public static void modifyTable(String namespace, String tableName, String columnFamily, int
version) throws IOException {
9
10         // 判断表格是否存在
11         if (!isTableExists(namespace, tableName)){
12             System.out.println("表格不存在");
13             return;
14         }
15
16         // 1. 获取admin

```

```

17         Admin admin = connection.getAdmin();
18
19         try {
20             // 2. 修改表
21             // 2.0 获取之前的表格描述!!!!!!!(这个很重要!!!!!!)
22             TableDescriptor descriptor = admin.getDescriptor(TableName.valueOf(namespace,
tableName));
23
24             // 2.1 创建一个表格描述
25             // 如果直接填写tableName，相当于创建了一个新的表格描述建造者，没有之前的信息。
26             // 如果想要修改之前的信息，必须调用方法填写一个旧的表格描述
27             TableDescriptorBuilder tableDescriptorBuilder =
TableDescriptorBuilder.newBuilder(descriptor);
28
29             // 2.2 对应建造者进行表格数据的修改
30             // 需要填写旧的列族描述
31             ColumnFamilyDescriptor columnFamily1 =
descriptor.getColumnFamily(Bytes.toBytes(columnFamily));
32             ColumnFamilyDescriptorBuilder columnFamilyDescriptorBuilder =
ColumnFamilyDescriptorBuilder.newBuilder(columnFamily1);
33             // 修改对应版本
34             columnFamilyDescriptorBuilder.setMaxVersions(version);
35             // 此处修改的时候，如果填写新创建的，那么别的参数会初始化
36
37             tableDescriptorBuilder.modifyColumnFamily(columnFamilyDescriptorBuilder.build());
38
39             admin.modifyTable(tableDescriptorBuilder.build());
40         } catch (IOException e) {
41             e.printStackTrace();
42         }
43
44         // 关闭admin
45         admin.close();
46     }

```

第五章 Phoenix

简介

Phoenix 是 HBase 的开源 SQL 中间层，它允许你使用标准 JDBC 的方式来操作 HBase 上的数据。在 Phoenix 之前，如果你要访问 HBase，只能调用它的 Java API，但相比于使用一行 SQL 就能实现数据查询，HBase 的 API 还是过于复杂。Phoenix 的理念是 `we put sql SQL back in NOSQL`，即你可以使用标准的 SQL 就能完成对 HBase 上数据的操作。同时这也意味着你可以通过集成 Spring Data JPA 或 Mybatis 等常用的持久层框架来操作 HBase。

其次 Phoenix 的性能表现也非常优异，Phoenix 查询引擎会将 SQL 查询转换为一个或多个 HBase Scan，通过并行执行来生成标准的 JDBC 结果集。它通过直接使用 HBase API 以及协处理器和自定义过滤器，可以为小型数据查询提供毫秒级的性能，为千万行数据的查询提供秒级的性能。同时 Phoenix 还拥有二级索引等 HBase 不具备的特性，因为以上的优点，所以 Phoenix 成为了 HBase 最优秀的 SQL 中间层。

Phoenix 简单使用

创建表

```
1 CREATE TABLE IF NOT EXISTS us_population (  
2     state CHAR(2) NOT NULL,  
3     city VARCHAR NOT NULL,  
4     population BIGINT  
5     CONSTRAINT my_pk PRIMARY KEY (state, city));
```

新建的表会按照特定的规则转换为 HBase 上的表，关于表的信息，可以通过 Hbase Web UI 进行查看

插入数据

Phoenix 中插入数据采用的是 **UPSERT** 而不是 **INSERT**，因为 Phoenix 并没有更新操作，插入相同主键的数据就视为更新，所以 **UPSERT** 就相当于 **UPDATE + INSERT**

```
1 UPSERT INTO us_population VALUES(' NY', 'New York', 8143197);  
2 UPSERT INTO us_population VALUES(' CA', 'Los Angeles', 3844829);  
3 UPSERT INTO us_population VALUES(' IL', 'Chicago', 2842518);  
4 UPSERT INTO us_population VALUES(' TX', 'Houston', 2016582);  
5 UPSERT INTO us_population VALUES(' PA', 'Philadelphia', 1463281);  
6 UPSERT INTO us_population VALUES(' AZ', 'Phoenix', 1461575);  
7 UPSERT INTO us_population VALUES(' TX', 'San Antonio', 1256509);  
8 UPSERT INTO us_population VALUES(' CA', 'San Diego', 1255540);  
9 UPSERT INTO us_population VALUES(' TX', 'Dallas', 1213825);  
10 UPSERT INTO us_population VALUES(' CA', 'San Jose', 912332);
```

改数据

```
1 -- 插入主键相同的数据就视为更新  
2 UPSERT INTO us_population VALUES(' NY', 'New York', 999999);
```

删除数据

```
1 DELETE FROM us_population WHERE city='Dallas';
```

查询数据

```
1 SELECT state as "州", count(city) as "市", sum(population) as "热度"  
2 FROM us_population  
3 GROUP BY state  
4 ORDER BY sum(population) DESC;
```

退出命令

```
1 !quit
```

扩展

从上面的操作中可以看出，Phoenix 支持大多数标准的 SQL 语法。关于 Phoenix 支持的语法、数据类型、函数、序列等详细信息，因为涉及内容很多，可以参考其官方文档，官方文档上有详细的说明：

- **语法 (Grammar)** : <https://phoenix.apache.org/language/index.html>
- **函数 (Functions)** : <http://phoenix.apache.org/language/functions.html>
- **数据类型 (Datatypes)** : <http://phoenix.apache.org/language/datatypes.html>
- **序列 (Sequences)** : <http://phoenix.apache.org/sequences.html>

- 联结查询 (Joins) : <http://phoenix.apache.org/joins.html>