

云凡教育大数据学院

Hive 的数据存储模式

Hive 的数据分为表数据和元数据，表数据是 Hive 中表格（table）具有的数据；而元数据是用来存储表的名字，表的列和分区及其属性，表的属性（是否为外部表等），表的数据所在目录等。下面分别来介绍。

一、Hive 的数据存储

在《[Hive 到底是什么](#)》博文中我们提到 Hive 是基于 Hadoop 分布式文件系统的，它的数据存储在 Hadoop 分布式文件系统中。Hive 本身是没有专门的数据存储格式，也没有为数据建立索引，只需要在创建表的时候告诉 Hive 数据中的列分隔符和行分隔符，Hive 就可以解析数据。所以往 Hive 表里面导入数据只是简单的将数据移动到表所在的目录中（如果数据是在 HDFS 上；但如果数据是在本地文件系统中，那么是将数据复制到表所在的目录中）。

Hive 中主要包含以下几种数据模型：Table（表），External Table（外部表），Partition（分区），Bucket（桶）（本博客会专门写几篇博文来介绍分区和桶）。

1、表：Hive 中的表和关系型数据库中的表在概念上很类似，每个表在 HDFS 中都有相应的目录用来存储表的数据，这个目录可以通过 `$(HIVE_HOME)/conf/hive-site.xml` 配置文件中的 `hive.metastore.warehouse.dir` 属性来配置，这个属性默认的值是 `/user/hive/warehouse`（这个目录在 HDFS 上），我们可以根据实际的情况来修改这个配置。如果我有一个表 `wyp`，那么在 HDFS 中会创建 `/user/hive/warehouse/wyp` 目录（这里假定 `hive.metastore.warehouse.dir` 配置为 `/user/hive/warehouse`）；`wyp` 表所有的数据都存放在这个目录中。这个例外是外部表。

2、外部表：Hive 中的外部表和表很类似，但是其数据不是放在自己表所属的目录中，而是存放到别处，这样的好处是如果你要删除这个外部表，该外部表所指向的数据是不会被删除的，它只会删除外部表对应的元数据；而如果你要删除表，该表对应的所有数据包括元数据都会被删除。

3、分区：在 Hive 中，表的每一个分区对应表下的相应目录，所有分区的数据都是存储在对应的目录中。比如 `wyp` 表有 `dt` 和 `city` 两个分区，则对应 `dt=20131218,city=BJ` 对应表的目录为 `/user/hive/warehouse/dt=20131218/city=BJ`，所有属于这个分区的数据都存放在这个目录中。

4、桶：对指定的列计算其 hash，根据 hash 值切分数据，目的是为了并行，每一个桶对应一个文件（注意和分区的区别）。比如将 `wyp` 表 `id` 列分散至 16 个桶中，首先对 `id` 列的值计算 hash，对应 hash 值为 0 和 16 的数据存储的 HDFS 目录为：

`/user/hive/warehouse/wyp/part-00000`；而 hash 值为 2 的数据存储的 HDFS 目录为：

`/user/hive/warehouse/wyp/part-00002`。

来看下 Hive 数据抽象结构图

Hive的那些事：<http://www.iteblog.com/archives/tag/hive的那些事>

数据库(Database)



Hive 数据抽象

从上图可以看出，表是在数据库下面，而表里面又要分区、桶、倾斜的数据和正常的数据等；分区下面也是可以建立桶的。

二、Hive 的元数据

Hive 中的元数据包括表的名字，表的列和分区及其属性，表的属性（是否为外部表等），表的数据所在目录等。由于 Hive 的元数据需要不断的更新、修改，而 HDFS 系统中的文件是多读少改的，这显然不能将 Hive 的元数据存储存储在 HDFS 中。目前 Hive 将元数据存储存储在数据库中，如 MySQL、Derby 中。我们可以通过以下的配置来修改 Hive 元数据的存储方式

```
<property>

<name>javax.jdo.option.ConnectionURL</name>

<value>jdbc:mysql://localhost:3306/hive_hdp?characterEncoding=UTF-8

&createDatabaseIfNotExist=true</value>

<description>JDBC connect string for a JDBC metastore</description>

</property>
```

云凡教育大数据学院

```
<property>

  <name>javax.jdo.option.ConnectionDriverName</name>

  <value>com.mysql.jdbc.Driver</value>

  <description>Driver class name for a JDBC metastore</description>

</property>

<property>

  <name>javax.jdo.option.ConnectionUserName</name>

  <value>root</value>

  <description>username to use against metastore database</description>

</property>

<property>

  <name>javax.jdo.option.ConnectionPassword</name>

  <value>123456</value>

  <description>password to use against metastore database</description>

</property>
```

当然，你还需要将相应数据库的启动复制到`$(HIVE_HOME)/lib` 目录中，这样才能将元数据存储在对应的数据库中。

云凡教育大数据学院

云凡教育大数据学院 www.cloudyhadoop.com

直击 30 万年薪 业界首播：大数据企业面试+企业大数据实战项目公开课：

- 1) HDFS、YARN 相关面试题
- 2) MapReduce 高级编程相关面试题
- 3) Hbase、Hive、Storm、Spark、Solr 相关面试题
- 4) 项目经验

12 月 16 日，晚 21: 00 在 YY: 20483828 直线开讲
点击即可进入课堂: <http://www.yy.com/20483828>
详情请加入 QQ 群: 374152400 ， 咨询课程顾问！



关注云凡教育微信公众号 **yfteach**，第一时间获取公开课信息。

实时在线授课，一线研发技术
www.yfteach.com