# Hadoop数据分析平台 第9周
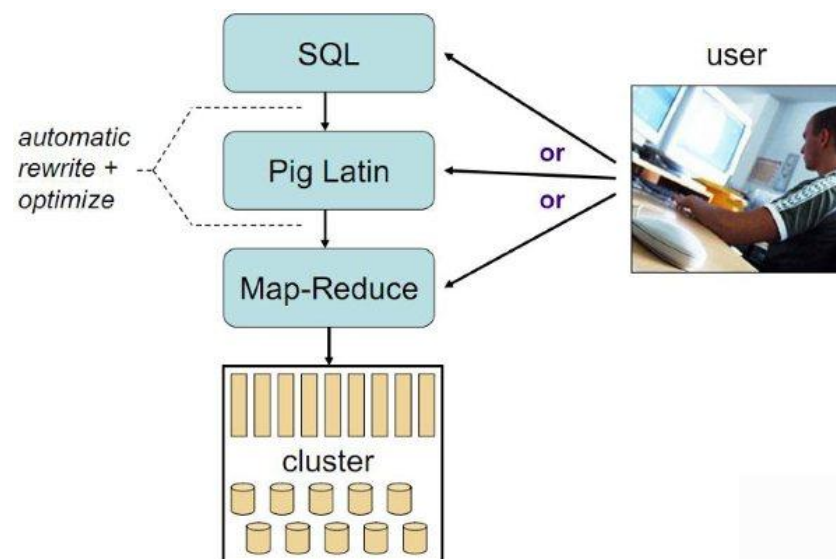
# Pig

- Hadoop客户端

- 使用类似于SQL的面向数据流的语言Pig Latin

- Pig Latin可以完成排序，过滤，求和，聚组，关联等操作，可以支持自定义函数

- Pig自动把Pig Latin映射为Map-Reduce作业上传到集群运行，减少用户编写Java程序的苦恼

- 三种运行方式：Grunt shell，脚本方式，嵌入式

# Hadoop流：最简便的M-R

```
% cat input/ncdc/sample.txt | ch02/src/main/ruby/max_temperature_map.rb | \
  sort | ch02/src/main/ruby/max_temperature_reduce.rb
1949    111
1950    22
```

```
% hadoop jar $HADOOP_INSTALL/contrib/streaming/hadoop-*-streaming.jar \
  -input input/ncdc/sample.txt \
  -output output \
  -mapper ch02/src/main/ruby/max_temperature_map.rb \
  -reducer ch02/src/main/ruby/max_temperature_reduce.rb
```

# Wordcount的例子

bin/hadoop  jar  contrib/streaming/hadoop-0.20.2-streaming.jar  -input
  input  -output  output  -mapper  /bin/cat  -reducer  /usr/bin/wc

**注意，命令一定要写完整的路径**

# 一个案例：生物数据库

# BLAST

# BLAST

# BLAST的Map-Reduce化

- BLAST比对算法，只涉及独立的一条基因信息，没有交叉计算，非常适合M-R

- BLAST算法用c实现，代码庞大，修改困难

- 权宜之计可以使用hadoop stream快速实现

# Hive

- 数据仓库工具。可以把Hadoop下的原始结构化数据变成Hive中的表

- 支持一种与SQL几乎完全相同的语言HiveQL。除了不支持更新、索引和事务，几乎SQL的其它特征都能支持

- 可以看成是从SQL到Map-Reduce的映射器

- 提供shell、JDBC/ODBC、Thrift、Web等接口

# Hive简介

- 起源自facebook由Jeff Hammerbacher领导的团队

- 构建在Hadoop上的数据仓库框架

- 设计目的是让SQL技能良好，但Java技能较弱的分析师可以查询海量数据

- 2008年facebook把hive项目贡献给Apache

# Hive的组件与体系架构

- 用户接口：shell, thrift, web等

- Thrift服务器

- 元数据库 "Derby, Mysql等

- 解析器

- Hadoop

# Hive安装

- 内嵌模式：元数据保持在内嵌的Derby模式，只允许一个会话连接

- 本地独立模式：在本地安装Mysql，把元数据放到Mysql内

- 远程模式：元数据放置在远程的Mysql数据库

# Hive安装：内嵌模式

1.下载

http://apache.dataguru.cn/hive/hive-0.8.1/hive-0.8.1.tar.gz

2.安装

(1)上传hive安装包到机器上,使用root用户登陆:

 tar -xvf hive-0.8.1.tar.gz

(2)将解压的hive分别移动并改名为/usr/local/hive

rm -rf /usr/local/hive mv hive-0.8.1  /usr/local/hive

3.配置hive

(1)修改/usr/local/hive/bin/hive-config.sh

在文件末尾加入

export JAVA_HOME=/usr/local/jdk export HIVE_HOME=/usr/local/hive export HADOOP_HOME=/usr/local/hadoop

(2) 根据hive-default.xml复制hive-site.xml

cp /usr/local/hive/conf/hive-default.xml /usr/local/hive/conf/hive-site.xml

(3)配置hive-site.xml,主要配置项如下:

hive.metastore.warehouse.dir：（HDFS上的）数据目录

hive.exec.scratchdir：（HDFS上的）临时文件目录

hive.metastore.warehouse.dir默认值是/user/hive/warehouse

hive.exec.scratchdir默认值是/tmp/hive-${user.name}

**2012.11.11**

以上是默认值，暂时不改。

(4)改变 /usr/local/hive的目录所有者为hadoop

chown -R hadoop:hadoop /usr/local/hive

(5)配置hive的log4j:

cp /usr/loca/hive/conf/hive-log4j.properties.template

/usr/loca/hive/conf/hive-log4j.properties

修改/usr/loca/hive/conf/hive-log4j.properties将

org.apache.hadoop.metrics.jvm.EventCounter改为

org.apache.hadoop.log.metrics.EventCounter

(6)启动hive

使用hadoop用户登陆,执行/usr/local/hive/bin/hive

# Hive安装：独立模式

- 安装Mysql并启动服务

- 在Mysql中为hive建立账号，并授予足够的权限，例如hive账号，授予all privileges

- 用上述账号登陆mysql，然后创建数据库，比如名叫hive，用于存放hive的元数据

- 在本地安装mysql客户端

- 配置hive-site.xml文件，指出使用本地Mysql数据库，已经连接协议，账号、口令等

- 把mysql-connector-java-x.x.x.jar复制到hive的lib目录下

- 启动hive能进入shell表示安装成功

# Hive安装：远程模式

- 在本地模式的基础上修改hive-site.xml文件，设置hive.metastore.local为false，并指向远程mysql数据库即可

# hive-site.xml文件内容

```xml
<property>
  <name>hive.metastore.local</name>
  <value>false</value>
  <description>controls whether to connect to remove metastore server or open a new metastore server in Hive Client JVM</description>
</property>


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

    <value>jdbc:mysql://mysql_server_host:3306/hivedb?createDatabaseIfNotExist=true&useUnicode=true&characterEncoding=latin1</value>
  <description>JDBC connect string for a JDBC metastore</description>
</property>
```

```
<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>mysql_username</value>
  <description>username to use against metastore database</description>
</property>


<property>
  <name>javax.jdo.option.ConnectionPassword</name>
  <value>mysql_password</value>
  <description>password to use against metastore database</description>
</property>
```

**2012.11.11**

# hive-site.xml文件内容

```xml
<property>
  <name>hive.stats.dbconnectionstring</name>

    <value>jdbc:mysql://mysql_server_host:3306/hive_stats?useUnicode=true&characterEncoding=latin1&user=mysql
    _username&password=mysql_password&createDatabaseIfNotExist=true</value>
  <description>The default connection string for the database that stores temporary hive statistics.</description>
</property>


<property>
  <name>hive.stats.dbconnectionstring</name>

    <value>jdbc:mysql://mysql_server_host:3306/hive_stats?useUnicode=true&characterEncoding=utf8&user=mysql_
    username&password=mysql_password&createDatabaseIfNotExist=true</value>
  <description>The default connection string for the database that stores temporary hive statistics.</description>
</property>
```

# hive-site.xml文件内容

```
<property>
  <name>hive.stats.dbclass</name>
  <value>jdbc:mysql</value>
  <description>The default database that stores temporary hive statistics.</description>
</property>


<property>
  <name>hive.stats.jdbcdriver</name>
  <value>com.mysql.jdbc.Driver</value>
  <description>The JDBC driver for the database that stores temporary hive statistics.</description>
</property>


<property>
  <name>hive.metastore.uris</name>
  <value>thrift://127.0.0.1:9083</value>
</property>
```

**2012.11.11**

# Hive shell

- 执行HiveQL（大约相当于SQL 92标准）

- 查看或临时设置Hive参数，只对当前会话有效

- 创建函数

- 导入jar包

**2012.11.11**

# 创建表

# 创建表

# 插入数据

```
hive> INSERT OVERWRITE TABLE RESULT
    > SELECT
    >   IMSI, IMEI, SUBSTR ( CGI, 8 ), STARTTIME, NULL, UPDATETYPE, NULL, NULL,
3
    > FROM LOC
    > WHERE IMSI IS NOT NULL;
Total MapReduce jobs = 2
Launching Job 1 out of 2
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_201206262230_0011, Tracking URL = http://localhost:50030/jobd
etails.jsp?jobid=job_201206262230_0011
Kill Command = /home/james/hadoop/bin/../bin/hadoop job  -Dmapred.job.tracker=lo
calhost:9001 -kill job_201206262230_0011
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 0
2012-06-27 00:45:25,408 Stage-1 map = 0%,  reduce = 0%
2012-06-27 00:45:28,421 Stage-1 map = 100%,  reduce = 0%
2012-06-27 00:45:31,432 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_201206262230_0011
Ended Job = 1152434878, job is filtered out (removed at runtime).
Moving data to: hdfs://localhost:9000/tmp/hive-james/hive_2012-06-27_00-45-20_50
1_6276951200625837126/-ext-10000
Loading data to table default.result
Deleted hdfs://localhost:9000/user/hive/warehouse/result
Table default.result stats: [num_partitions: 0, num_files: 1, num_rows: 0, total
_size: 538, raw_data_size: 0]
7 Rows loaded to result
MapReduce Jobs Launched:
Job 0: Map: 1   HDFS Read: 1002 HDFS Write: 538 SUCESS
Total MapReduce CPU Time Spent: 0 msec
```

# 表连接

```
hive> SELECT RESULT.IMSI, LOC.INSTIME
    > FROM RESULT JOIN LOC ON ( RESULT.IMSI = LOC.IMSI );
Total MapReduce jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapred.reduce.tasks=<number>
Starting Job = job_201206262230_0013, Tracking URL = http://localhost:50030/jobd
etails.jsp?jobid=job_201206262230_0013
Kill Command = /home/james/hadoop/bin/../bin/hadoop job  -Dmapred.job.tracker=lo
calhost:9001 -kill job_201206262230_0013
Hadoop job information for Stage-1: number of mappers: 2; number of reducers: 1
2012-06-27 00:48:28,570 Stage-1 map = 0%,   reduce = 0%
2012-06-27 00:48:31,592 Stage-1 map = 100%,   reduce = 0%
2012-06-27 00:48:40,624 Stage-1 map = 100%,   reduce = 100%
Ended Job = job_201206262230_0013
MapReduce Jobs Launched:
Job 0: Map: 2  Reduce: 1   HDFS Read: 1540 HDFS Write: 252 SUCESS
Total MapReduce CPU Time Spent: 0 msec
OK
460000722940589 2012-03-16 00:00:00
460000940196027 2012-03-16 00:00:00
460020202346902 2012-03-16 00:00:00
460022676514472 2012-03-16 00:00:00
460023173370082 2012-03-16 00:00:00
460027157683337 2012-03-16 00:00:00
460029146542227 2012-03-16 00:00:00
Time taken: 20.828 seconds
```

2012.11.11

# JDBC/ODBC接口

- 用户可以像连接传统关系数据库一样使用JDBC或ODBC连接Hive

- 目前还不成熟

1.使用jdbc的方式连接Hive，首先做的事情就是需要启动hive的Thrift Server,否则连接
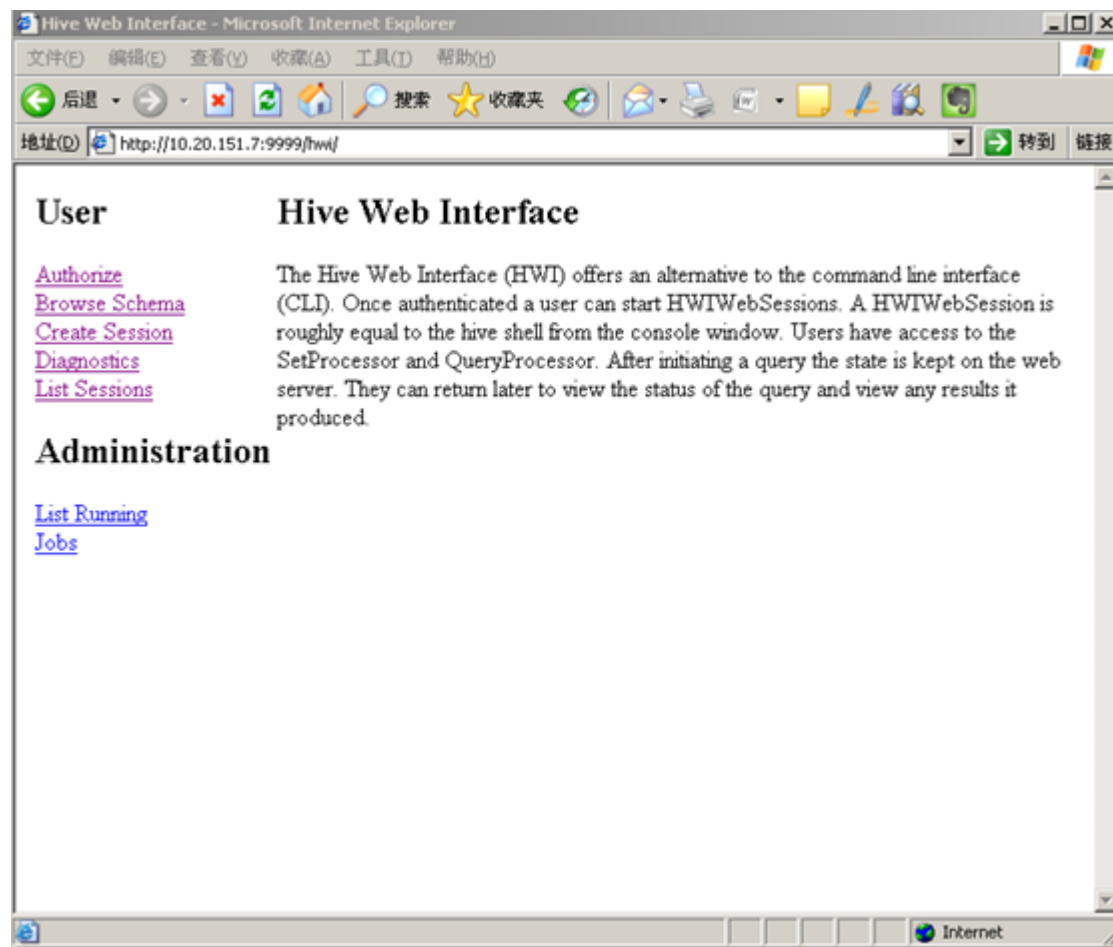   hive的时候会报connection refused的错误。

启动命令如下：

hive --service hiveserver


2.新建java项目，然后将hive/lib下的所有jar包和hadoop的核心jar包hadoop-0.20.2-
   core.jar添加到项目的类路径上。

**2012.11.11**

```java
public static void main(String[] args) throws Exception {

    // TODO Auto-generated method stub

    Class.forName("org.apache.hadoop.hive.jdbc.HiveDriver");

    String dropSql="drop table pokes";

    String createSql="create table pokes (foo int,bar string)";

    String insertSql="load data local inpath '/home/zhangxin/hive/kv1.txt' overwrite into  table pokes";

    String querySql="select bar from pokes limit 5";

    Connection connection=DriverManager.getConnection("jdbc:hive://localhost:10000/default", "", "");

    Statement statement=connection.createStatement();

    statement.execute(dropSql);

    statement.execute(createSql);

    statement.execute(insertSql);

    ResultSet rs=statement.executeQuery(querySql);

     while(rs.next())

    {

        System.out.println(rs.getString("bar"));

    } }
```

**2012.11.11**

# Web接口

- 假设hive部署在
  10.20.151.7机器上，
  conf/hive-default.xml
  文件都是默认值，那么
  我们直接在浏览器中输
  入：
  http://10.20.151.7:999
  9/hwi/ 就可以访问了

NUCLEUS_TABLES

A

DBS

SEQUENCE_TABLE

SERDES

TBLS
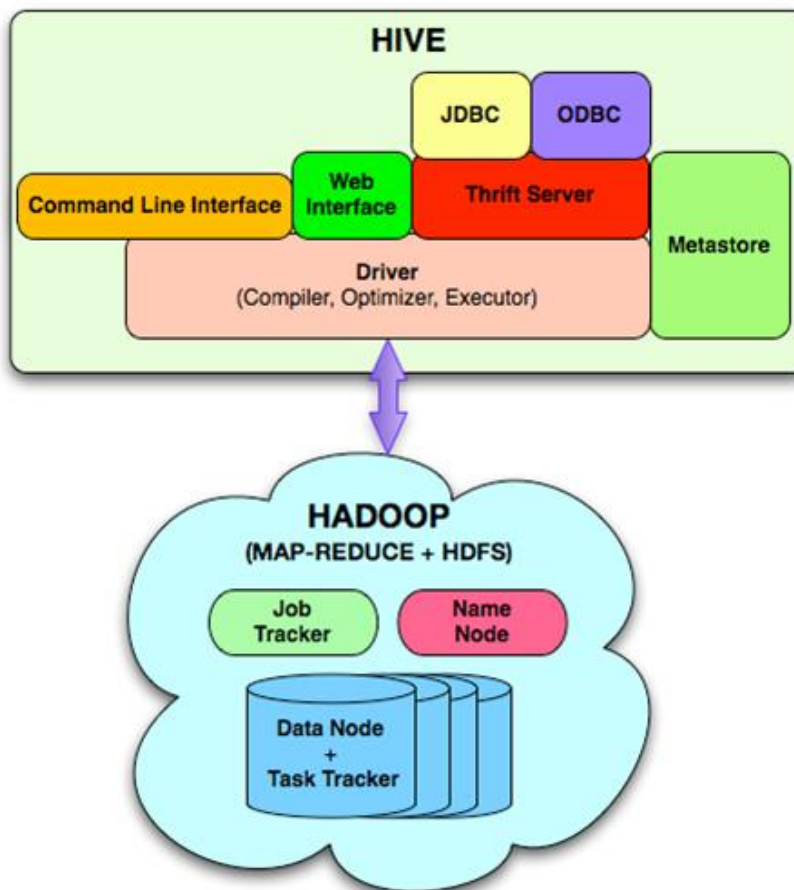
SDS

PARTITION_KEYS

COLUMNS

BUCKETING_COLS

SD_PARAMS

SORT_COLS

SERDE_PARAMS

TABLE_PARAMS



**2012.11.11**

# Hive的数据放在哪儿？

- 数据在HDFS的warehouse目录下，一个表对应一个子目录

- 桶与reduce

- 本地的/tmp目录存放日志和执行计划

# Hive的UDF

- 见刘鹏书P196

# Thanks

**FAQ时间**