# Hadoop数据分析平台 第7周

DATAGURU专业数据分析网站

# Hbase安装：单机模式

- 下载及解压hbase安装包

- 修改conf/hbase-env.sh脚本，设置环境变量

- 编辑hbase-site.xml进行配置

- 启动Hbase

- 验证Hmaster已经启动

- 进入shell

# 下载及解压Hbase安装包

# 修改hbase-env.sh

- 设置JAVA_HOME环境变量

```
# * See the License for the specific language govern
# * limitations under the License.
# */


# Set environment variables here.


# The java implementation to use.  Java 1.6 required
export JAVA_HOME=/usr/java/jdk1.6.0_26/


# Extra Java CLASSPATH elements.  Optional.
# export HBASE_CLASSPATH=


# The maximum amount of heap to use, in MB. Default
# export HBASE_HEAPSIZE=1000
```

**2012.10.23**

# 配置hbase-site.xml

■ 先创建用于存放数据的目录/home/grid/hbase-0.90.5/data

# 启动Hbase及验证

```
[grid@h1 hbase-0.90.5]$ bin/start-hbase.sh
starting master, logging to /home/grid/hbase-0.90.5/bin/../logs/hbase-grid-master-h1.out
[grid@h1 hbase-0.90.5]$ /usr/java/jdk1.6.0_26/bin/jsp
-bash: /usr/java/jdk1.6.0_26/bin/jsp: No such file or directory
[grid@h1 hbase-0.90.5]$ /usr/java/jdk1.6.0_26/bin/jps
5334 Jps
4150 SecondaryNameNode
4025 NameNode
5184 HMaster
4219 JobTracker
[grid@h1 hbase-0.90.5]$ bin/hbase shell
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 0.90.5, r1212209, Fri Dec  9 05:40:36 UTC 2011

hbase(main):001:0> quit
[grid@h1 hbase-0.90.5]$
```

**2012.10.23**

# Hbase安装：伪分布模式

- 在单点模式的基础上继续

- 编辑 hbase-env.sh增加HBASE_CLASSPATH环境变量

- 编辑hbase-site.xml打开分布模式

- 覆盖hadoop核心jar包

- 启动hbase

- 验证启动

# 编辑 hbase-env.sh增加 HBASE_CLASSPATH环境变量

■ 用于帮助hbase找到hadoop

```
# * See the License for the specific language governing permissi
# * limitations under the License.
# */


# Set environment variables here.


# The java implementation to use.   Java 1.6 required.
export JAVA_HOME=/usr/java/jdk1.6.0_26/


# Extra Java CLASSPATH elements.   Optional.
export HBASE_CLASSPATH=/home/grid/hadoop-0.20.2/conf


# The maximum amount of heap to use, in MB. Default is 1000.
# export HBASE_HEAPSIZE=1000
"hbase-env.sh" 76L, 3378C written
```

# 编辑hbase-site.xml打开分布模式

```
 * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, eithe
 * See the License for the specific language governing
 * limitations under the License.
 */
-->
<configuration>
<property>
<name>hbase.rootdir</name>
<value>file:///home/grid/hbase-0.90.5/data</value>
</property>
<property>
<name>hbase.cluster.distributed</name>
<value>true</value>
</property>
</configuration>
"hbase-site.xml" 33L, 1166C written
```

# 覆盖hadoop核心jar包

■ 这是关键一步，主要目的是防止因为hbase和hadoop版本不同出现兼容问题，造成hmaster启动异常

```
[grid@h1 lib]$ mv hadoop-core-0.20-append-r1056497.jar hadoop-core-0.20-append-r1056497.sav
[grid@h1 lib]$ ls ../../hadoop-0.20.2/
bin              conf                    hadoop-0.20.2-core.jar       ivy              LICENSE.txt
build.xml        contrib                 hadoop-0.20.2-examples.jar   ivy.xml          logs
c++              docs                    hadoop-0.20.2-test.jar       lib              NOTICE.txt
CHANGES.txt      hadoop-0.20.2-ant.jar   hadoop-0.20.2-tools.jar      librecordio      pig_1340564601586.log
```
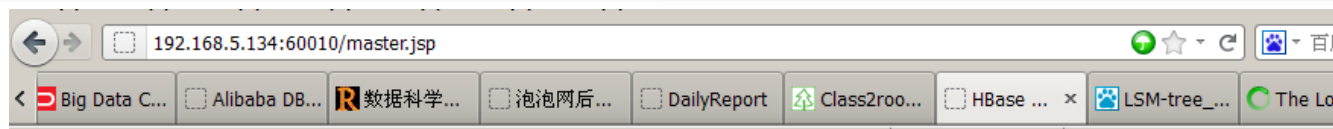
```
[grid@h1 lib]$ cp ../../hadoop-0.20.2/hadoop-0.20.2-core.jar .
[grid@h1 lib]$ ls
activation-1.1.jar          guava-r06.jar                    jersey-core-1.4.jar
asm-3.1.jar                 hadoop-0.20.2-core.jar           jersey-json-1.4.jar
avro-1.3.3.jar              hadoop-core-0.20-append-r1056497.sav  jersey-server-1.4.jar
commons-cli-1.2.jar         jackson-core-asl-1.5.5.jar       jettison-1.1.jar
commons-codec-1.4.jar       jackson-jaxrs-1.5.5.jar          jetty-6.1.26.jar
commons-el-1.0.jar          jackson-mapper-asl-1.4.2.jar     jetty-util-6.1.26.jar
```

# 启动hbase并验证

```
[grid@h1 lib]$ cd ..
[grid@h1 hbase-0.90.5]$ bin/start-hbase.sh
localhost: starting zookeeper, logging to /home/grid/hbase-0.90.5/bin/../logs/hbase-grid-zookeeper-h1.out
starting master, logging to /home/grid/hbase-0.90.5/bin/../logs/hbase-grid-master-h1.out
localhost: starting regionserver, logging to /home/grid/hbase-0.90.5/bin/../logs/hbase-grid-regionserver-h1.out
[grid@h1 hbase-0.90.5]$ /usr/java/jdk1.6.0_26/bin/jps
6022 Jps
4150 SecondaryNameNode
5895 HRegionServer
5789 HMaster
5747 HQuorumPeer
4025 NameNode
4219 JobTracker
[grid@h1 hbase-0.90.5]$
```

# Hbase安装：完全分布模式

- 配置hosts，确保涉及的主机名均可以解析为ip

- 编辑hbase-env.xml

- 编辑hbase-site.xml

- 编辑regionservers文件

- 把Hbase复制到其它节点

- 启动Hbase

- 验证启动

# Web管理界面

# Shell



```
hadoop@vincent-laptop:/usr/hbase-0.20.6$ cd bin
hadoop@vincent-laptop:/usr/hbase-0.20.6/bin$ ls
add_table.rb     hbase            hbase-daemons.sh  loadtable.rb       set_m
copy_table.rb    hbase-config.sh  HBase.rb          regionservers.sh   start-
Formatter.rb     hbase-daemon.sh  hirb.rb           rename_table.rb    stop-l
hadoop@vincent-laptop:/usr/hbase-0.20.6/bin$ ./hbase shell
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Version: 0.20.6, r965666, Mon Jul 19 16:54:48 PDT 2010
hbase(main):001:0>
```

# Shell命令帮助

```
hbase(main):001:0> help
HBase Shell, version 0.91.0-SNAPSHOT, r1130916, Sat Jul 23 12:44:34 CEST 2011
Type 'help "COMMAND"', (e.g. 'help "get"' -- the quotes are necessary) for
help on a specific command. Commands are grouped. Type 'help "COMMAND_GROUP"',
(e.g. 'help "general"') for help on a command group.

COMMAND GROUPS:
  Group name: general
  Commands: status, version

  Group name: ddl
  Commands: alter, create, describe, disable, drop, enable, exists,
            is_disabled, is_enabled, list
```

hbase(main):024:0>status

3 servers, 0 dead,1.0000 average load

hbase(main):025:0>version

0.90.4, r1150278,Sun Jul 24 15:53:29 PDT 2011

hbase(main):011:0>create 'member','member_id','address','info'

0 row(s) in 1.2210seconds

```
hbase(main):012:0>list

TABLE

member

1 row(s) in 0.0160seconds

hbase(main):006:0>describe 'member'

DESCRIPTION                                          ENABLED

{NAME => 'member', FAMILIES => [{NAME=> 'address', BLOOMFILTER => 'NONE', REPLICATION_SCOPE =>
    '0', true

VERSIONS => '3', COMPRESSION => 'NONE',TTL => '2147483647', BLOCKSIZE => '65536', IN_MEMORY =>
    'false', BLOCKCACHE => 'true'}, {NAME =>'info', BLOOMFILTER => 'NONE', REPLICATION_SCOPE => '0',
    VERSI

ONS => '3', COMPRESSION => 'NONE', TTL=> '2147483647', BLOCKSIZE => '65536', IN_MEMORY =>
    'false',

 BLOCKCACHE => 'true'}]}

1 row(s) in 0.0230seconds
```

**2012.10.23**

# 删除列族：alter、disable、enable命令

hbase(main):003:0>alter 'member',{NAME=>'member_id',METHOD=>'delete'}

ERROR: Table memberis enabled. Disable it first before altering.

hbase(main):004:0>disable 'member'

0 row(s) in 2.0390seconds

hbase(main):005:0>alter'member',{NAME=>'member_id',METHOD=>'delete'}

0 row(s) in 0.0560seconds

hbase(main):008:0> enable 'member'

0 row(s) in 2.0420seconds

# 列出所有的表

hbase(main):028:0>list

TABLE

member

temp_table

2 row(s) in 0.0150seconds

# 删除表

hbase(main):029:0>disable 'temp_table'

0 row(s) in 2.0590seconds


hbase(main):030:0>drop 'temp_table'

0 row(s) in 1.1070seconds

# 查询一个表是否存在

hbase(main):021:0>exists 'member'

Table member

doesexist


0 row(s) in 0.1610seconds

# 判断表是否enable或disable

hbase(main):034:0>is_enabled 'member'

true


0 row(s) in 0.0110seconds


hbase(main):032:0>is_disabled 'member'

false


0 row(s) in 0.0110seconds

put'member','scutshuxue','info:age','24'

put'member','scutshuxue','info:birthday','1987-06-17'

put'member','scutshuxue','info:company','alibaba'

put'member','scutshuxue','address:contry','china'

put'member','scutshuxue','address:province','zhejiang'

put'member','scutshuxue','address:city','hangzhou'

 put'member','xiaofeng','info:birthday','1987-4-17'

put'member','xiaofeng','info:favorite','movie'

put'member','xiaofeng','info:company','alibaba'

put'member','xiaofeng','address:contry','china'

put'member','xiaofeng','address:province','guangdong'

put'member','xiaofeng','address:city','jieyang'

put'member','xiaofeng','address:town','xianqiao'

# 获取一个行健的所有数据

```
hbase(main):001:0>get 'member','scutshuxue'

COLUMN                    CELL

address:city              timestamp=1321586240244,value=hangzhou


 address:contry           timestamp=1321586239126,value=china


address:province          timestamp=1321586239197,value=zhejiang


 info:age                 timestamp=1321586238965,value=24


 info:birthday            timestamp=1321586239015, value=1987-06-
    17
 info:company             timestamp=1321586239071,value=alibaba


6 row(s) in 0.4720seconds
```

hbase(main):002:0>get 'member','scutshuxue','info'

COLUMN                              CELL

info:age                            timestamp=1321586238965,value=24

 info:birthday                       timestamp=1321586239015, value=1987-06-

    17

info:company                         timestamp=1321586239071,value=alibaba

3 row(s) in 0.0210seconds

hbase(main):002:0>get 'member','scutshuxue','info:age'

COLUMN                          CELL

 info:age                       timestamp=1321586238965,value=24


1 row(s) in 0.0320seconds

# 更新一条记录

hbase(main):004:0>put 'member','scutshuxue','info:age' ,'99'

0 row(s) in 0.0210seconds


hbase(main):005:0>get 'member','scutshuxue','info:age'

COLUMN                          CELL


 info:age                       timestamp=1321586571843,value=99


1 row(s) in 0.0180seconds

# 通过timestamp来获取数据

hbase(main):010:0>get 'member','scutshuxue',{COLUMN=>'info:age',TIMESTAMP=>1321586238965}

COLUMN                          CELL

info:age                        timestamp=1321586238965,value=24

1 row(s) in 0.0140seconds

hbase(main):011:0>get 'member','scutshuxue',{COLUMN=>'info:age',TIMESTAMP=>1321586571843}

COLUMN                          CELL

info:age                        timestamp=1321586571843,value=99

1 row(s) in 0.0180seconds

hbase(main):013:0>scan 'member'


结果略

hbase(main):016:0>delete 'member','temp','info:age'

0 row(s) in 0.0150seconds

hbase(main):018:0>get 'member','temp'

COLUMN                          CELL


0 row(s) in 0.0150seconds

# 删除整行

hbase(main):001:0>deleteall 'member','xiaofeng'

0 row(s) in 0.3990seconds

# 查询表中有多少行

hbase(main):019:0>count 'member'

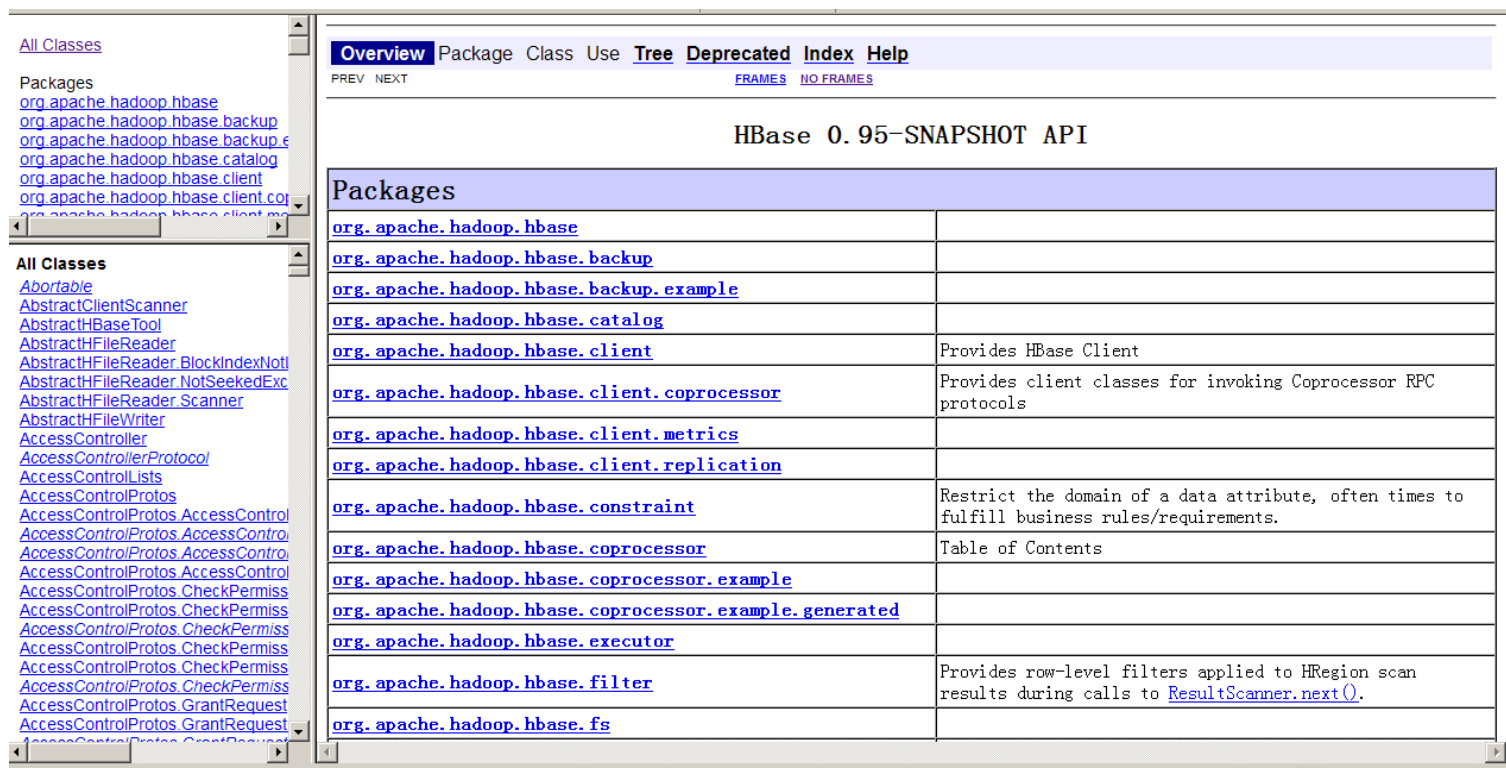2 row(s) in 0.0160seconds

hbase(main):035:0>truncate 'member'

Truncating 'member'table (it may take a while):

- Disabling table…

- Dropping table…

- Creating table…

0 row(s) in 4.3430seconds

# Hbase API

- 《Hbase权威指南》第3-5章

- http://hbase.apache.org/apidocs/index.html

# 什么情况下使用Hbase？

- 成熟的数据分析主题，查询模式已经确立并且不轻易改变

- 传统的关系型数据库已经无法承受负荷，高速插入，大量读取

- 适合海量的，但同时也是简单的操作（例如key-value）

# 场景一：浏览历史

# 关系型数据库的困难

- **简单的事情只要上了量就会变成无比复杂的事情**

- Order by耗费很多性能

- 大量发生，但又无法分布式处理

- 顾客需要实时看到自己的足迹，因此不能使用缓存技巧

# Hbase迎接挑战

- 天生就是面向时间戳查询

- 基于行键的查询异常快速，特别是最近的数据被放在内存的memstore里，完全没有IO开销

- 分布式化解负荷

# 模式设计

- 行键：userid

- 列族和列：book:bookid

- 为了充分利用分布式，可以用reverse key，hash等技巧改造行键

# 场景二：商品推荐

# 用关系型数据库实现

- [http://f.dataguru.cn/thread-84-1-1.html](http://f.dataguru.cn/thread-84-1-1.html)

- 拿ITPUB实验了一把。

  阅读推荐说白了，就是你打开一个帖子，看到有一个提示写着读了本帖的人有xx%读了xxxx贴，有xx%读了xxxx帖。。。等等，这项功能也可以推广到商品推荐，音乐推荐，下载推荐等等。

  在ITPUB中设置了一个log表，记录每次用户点击，有3个列，分别是时间戳，用户id，还有点击的主题id

  使用了一段时间的数据大约有1000万行，写了个sql搞定

```
01.    select A.threadid,count(distinct A.userid) from testtj A,testtj B where A.userid=B.userid and B.threadid=1479820 group by A.threadid
       order by 2 desc limit 10;
02.
03.    +----------+--------------------------+
04.    | threadid | count(distinct A.userid) |
05.    +----------+--------------------------+
06.    |  1479820 |                     1054 |
07.    |  1455924 |                      840 |
08.    |  1466253 |                      817 |
09.    |  1472481 |                      783 |
10.    |  1469262 |                      745 |
11.    |  1478790 |                      740 |
12.    |  1476679 |                      711 |
13.    |  1476821 |                      664 |
14.    |  1476860 |                      636 |
15.    |  1476068 |                      614 |
16.    +----------+--------------------------+
17.    10 rows in set (9.11 sec)
       复制代码
```

**2012.10.23**

- 两个表，一个是u-t，另一个是t-u

- U-t表的结构：行键为userid，列族和列为thread:threadid

- T-u表结构：行键为threadid，列族和列为user:userid

- 查询：先在u-t表从userid->threadid，再从t-u表从threadid->userid，在计算程序中 实现去重和统计功能

- 例子：学生表（学号，身份证号，姓名，性别，系，年龄），有时在学号上查询，有时在身份证号上查询

- 主表：行键为学号，列族为学生，下面的列是身份证号，姓名，性别，系，年龄

- 辅助（索引）表：行键为身份证号，列族和列为学号

**2012.10.23**

```
<userId> : <colfam> : <messageId> : <timestamp> : <email-message>

12345 : data : 5fc38314-e290-ae5da5fc375d : 1307097848 : "Hi Lars, ..."
12345 : data : 725aae5f-d72e-f90f3f070419 : 1307099848 : "Welcome, and ..."
12345 : data : cc6775b3-f249-c6dd2b1a7467 : 1307101848 : "To Whom It ..."
12345 : data : dcbee495-6d5e-6ed48124632c : 1307103848 : "Hi, how are ..."


<userId>-<messageId> : <colfam> : <qualifier> : <timestamp> : <email-message>

12345-5fc38314-e290-ae5da5fc375d : data : : 1307097848 : "Hi Lars, ..."
12345-725aae5f-d72e-f90f3f070419 : data : : 1307099848 : "Welcome, and ..."
12345-cc6775b3-f249-c6dd2b1a7467 : data : : 1307101848 : "To Whom It ..."
12345-dcbee495-6d5e-6ed48124632c : data : : 1307103848 : "Hi, how are ..."
```

**2012.10.23**

# 好处

- 便于分布

- 便于多条件伸缩查询

# Thanks

**FAQ时间**