

一 内容回顾（列举前一天重点难点内容）

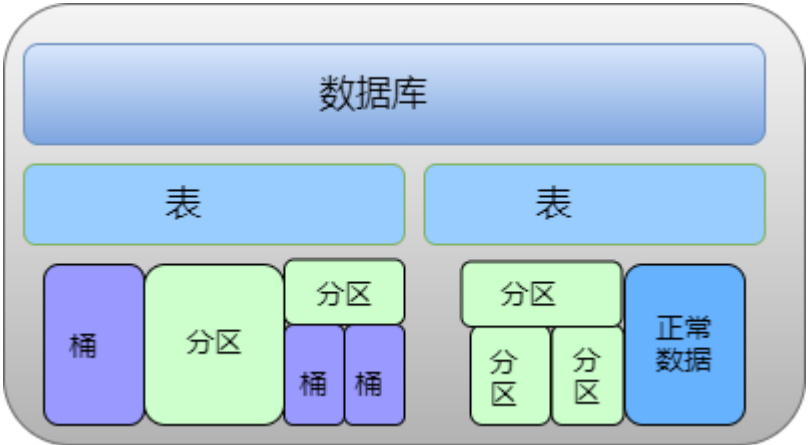
教学重点

1	Hive的简介
2	Hive的三种安装模式
3	表属性的修改
4	Hive的远程模式
5	Hive的架构
6	Hive的工作原理
7	Hive的基本操作
8	Hive加载数据的方式
9	表类型详解
10	分区表的相关概念
11	分区分类
12	分区操作

二 教学目标

1	分桶表来历
2	分桶表设置
3	抽样查询
4	分区分桶联合使用
5	Join连接
6	查询子句
7	排序
8	合并结果集
9	数据类型
10	系统内置函数

三 教学导读



3.1 为什么要分桶

- 1 当单个的分区或者表的数据量过大，分区不能更细粒度的划分数据，就需要使用分桶技术将数据划分成更细的粒度。

四 教学内容

4.1 分桶

4.1.1 分桶技术

```
1 [CLUSTERED BY (COLUMNNAME COLUMNTYPE [COMMENT 'COLUMN COMMENT'],...)  
2 [SORTED BY (COLUMNNAME [ASC|DESC])...] INTO NUM_BUCKETS BUCKETS]
```

4.1.2 关键字及其原理

bucket

分桶的原理:跟MR中的HashPartitioner原理一样:都是key的hash值取模reduce的数量

- MR中: 按照key的hash值除以reductTask取余
- Hive中: 按照分桶字段的hash值去模除以分桶的个数

4.1.3 分桶的意义

- 1 为了保存分桶查询的分桶结构 (数据按照分桶字段进行保存hash散列)
- 2 分桶表的数据进行抽样和JOIN时可以提高查询效率,一般是用来抽样查询

4.2 分桶表的创建示例

分桶表创建流程:

1. 首先要创建带分桶定义的表 (分桶表)
2. 然后创建一个临时表 (普通表)
3. 从临时表中使用分桶查询将查询到的数据插入到分桶表中

4.2.1 创建分桶表

```
1 create table if not exists buc1(  
2 id int,  
3 name string,  
4 age int  
5 )  
6 clustered by (id) into 4 buckets  
7 row format delimited fields terminated by ',';  
8  
9 数据: 从(data/buc1.txt加载如下格式数据)  
10 id,name,age
```

4.2.2 创建临时表

```
1 create table if not exists temp_buc1(  
2   id int,  
3   name string,  
4   age int  
5 )  
6 row format delimited fields terminated by ',';
```

4.2.3 分桶使用load方式加载数据不能体现分桶

```
1 load data local inpath '/root/hivedata/buc1.csv' into table buc1;
```

4.2.4 加载数据到临时表

```
1 load data local inpath '/root/hivedata/buc1.csv' into table temp_buc1;
```

4.2.5 使用分桶查询将数据导入到分桶表

```
1 insert overwrite table buc1  
2 select id,name,age from temp_buc1  
3 cluster by (id);
```

这里还是无法分桶,需要进行下面的强制设置

4.2.6 设置强制分桶的属性

```
1 <!-- 如果要分桶,就要打开分桶的强制模式 -->  
2 set hive.enforce.bucketing=false/true  
3 <name>hive.enforce.bucketing</name>  
4 <value>false</value>  
5 <description>  
6   Whether bucketing is enforced. If true, while inserting into the table, bucketing is  
   enforced.</description>  
7  
8 最好把强制排序也打开  
9 set hive.enforce.sorting=true;
```

4.2.7 如果设置了reduces的个数和总桶数不一样, 请手动设置

```
1 set mapreduce.job.reduces=-1 #-1表示可以根据实际需要来定制reduces的数量
```

4.2.8 创建指定排序字段的分桶表

```

1 create table if not exists buc3(
2   id int,
3   name string,
4   age int
5 )
6 clustered by (id)
7 sorted by (id desc) into 4 buckets
8 row format delimited fields terminated by ',';

```

4.2.9 导入数据

```

1 insert overwrite table buc3
2 select id,name,age from temp_buc1
3 distribute by (id) sort by (id asc);
4 和下面的语句效果一样
5 insert overwrite table buc3
6 select id,name,age from temp_buc1
7 cluster by (id);
8
9 注意:定义和导入数据时同时设置了分桶和排序,以导入数据时优先

```

4.3 分桶表查询案例

```

1 select * from buc3;
2 select * from buc3 tablesample(bucket 1 out of 1 on id);

```

4.3.1 查询第1桶的数据

```

1 语法解释:
2
3 tablesample(bucket x out of y on id)
4 x:代表从第几桶开始查询
5 y:可以是总桶数,总桶数的倍数或者是因子。x不能大于y。
6 y如果是4 ,已有的总桶数是4 那么要取的数据 4/4= 1个桶
7 y如果是2 ,已有的总桶数是4,那么要取的数据4/y 4/2=2个桶
8 on后面跟的是分桶的字段
9
10 hive根据y的大小, 决定抽样的比例。例如, table总共分了64份, 当y=32时, 抽取(64/32=)2个bucket的数据, 当y=128时, 抽取(64/128=)1/2个bucket的数据。
11
12 x表示从哪个bucket开始抽取。例如, table总bucket数为32, tablesample(bucket 3 out of 16), 表示总共抽取 (32/16=) 2个bucket的数据, 分别为第3个bucket和第 (3+16=) 第19个bucket的数据。
13
14 案例分析:
15 # 一共4桶,取得第一桶的数据
16 select * from buc3 tablesample(bucket 1 out of 4 on id);
17 # 取得第二桶的数据
18 #y如果是2 ,已有的总桶数是4,那么要取的数据4/y 4/2=2个桶
19 select * from buc3 tablesample(bucket 2 out of 2 on id);

```

```

20 # 从x开始取,取的数量是 (4/2)=2,当x==1时代表取的是1,3桶
21 select * from buc3 tablesample(bucket 1 out of 2 on id);
22 # 4/y=1/2,代表每一桶被分成2份,当x=1和5时表示第一桶,x=1表示取的是第一桶的前半部分.
23 # 小技巧:区分id是那一部分的方法:id/y值,例如当前的案例:id=8,y=8,id/y=0,因为商+1=1=x,桶中的部分就是
    前半部分,比如:id=12,y=8,id/y=4,商+1=5=x,所以是桶的后半部分
24 select * from buc3 tablesample(bucket 1 out of 8 on id);

```

4.3.2 查询

(注意: tablesample一定是放在from后面)

查询id为基数:

```

1  select *
2  from buc1 tablesample(bucket 2 out of 2 on id)
3  where name = "aa1";
4
5  select *
6  from buc2
7  where uname = "aa1";
8
9  查询:
10 select * from temp_buc1 limit 3;
11 select * from temp_buc1 tablesample(3 rows);#行数
12 select * from temp_buc1 tablesample(30 percent);#数量的百分比
13 select * from temp_buc1 tablesample(6B);B k M G T P #具体的数量
14
15 SELECT * from temp_buc1 order by rand() limit 3; #rand()是随机排序
16 https://www.aboutyun.com/thread-27093-1-1.html

```

4.4 分区分桶联合案例

4.4.1 案例

(先熟悉Select各种查询)

```

1  需求: 按照性别分区 (1男2女) , 在分区中按照id的奇偶进行分桶:
2
3  id,name,sex
4  数据请参考data/user.txt
5
6  建表:
7  create table if not exists stu(
8  id int,
9  name string
10 )
11 partitioned by (sex string)
12 clustered by (id) into 2 buckets
13 row format delimited fields terminated by ',';
14
15 建临时表:

```

```

16 create table if not exists stu_temp(
17 id int,
18 name string,
19 sex string
20 )
21 row format delimited fields terminated by ',';
22
23
24 -- 主要要把动态分区设为非严格模式,如下:
25 set hive.exec.dynamic.partition.mode=nonstrict;
26
27 -- 加载临时表的数据
28 load data local inpath '/root/hivedata/user.txt' into table stu_temp;
29
30
31 将数据导入到分区分桶表
32
33 insert overwrite table stu partition(sex)
34 select id,name,sex from stu_temp cluster by id ;
35
36 查询性别为女, 学号为奇数的学生
37
38 select * from stu tablesample(bucket 2 out of 2 on id)
39 where sex = '2';

```

4.4.2 分桶表总结

```

1 1. 定义
2     clustered by (id)      ---指定分桶的字段
3     sorted by (id asc|desc) ---指定数据的排序规则, 表示咱们预期的数据是以这种规则进行的排序
4
5 2. 导入数据
6     cluster by (id)  ---指定getPartition以哪个字段来进行hash, 并且排序字段也是指定的字段, 排序是以asc排列
7     distribute by (id)  ---- 指定getPartition以哪个字段来进行hash
8     sort by (name asc | desc) ---指定排序字段
9     order by --- 全域排序(只能有一个reduce)
10
11 区别: cluster by 这种方式可以分别指定getPartition和sort的字段
12
13 导数据时:
14     insert overwrite table buc3
15     select id,name,age from temp_buc1
16     distribute by (id) sort by (id asc);
17     和下面的语句效果一样
18     insert overwrite table buc4
19     select id,name,age from temp_buc1
20     cluster by (id) ;

```

4.4.3 注意事项

- 1 分区使用的是表外字段，分桶使用的是表内字段
- 2 分桶更加用于细粒度的管理数据，更多的是使用来做抽样、join

4.5 查询语句基本语法

4.5.1 select查询结构基本语法

下面是一个sql查询语句的基本结构

```
1  select selection_list # 查询的列
2  from table           # 要查询的表
3  join on              # 连接的表
4  where                # 查询条件
5  group by            # 分组查询
6  having               # 分组条件过滤
7  order by            # 字段排序
8  sort by             # 结果排序
9  limit               # 限制结果数
10 union/union all     # 合并表
```

4.5.2 sql语句的执行顺序

```
1  FROM
2  <left_table>
3  ON
4  <join_condition>
5  <join_type>
6  JOIN
7  <right_table>
8  WHERE
9  <where_condition>
10 GROUP BY
11 <group_by_list>
12 HAVING
13 <having_condition>
14 SELECT
15 DISTINCT
16 <select_list>
17 ORDER BY
18 <order_by_condition>
19 LIMIT
20 <limit_number>
```

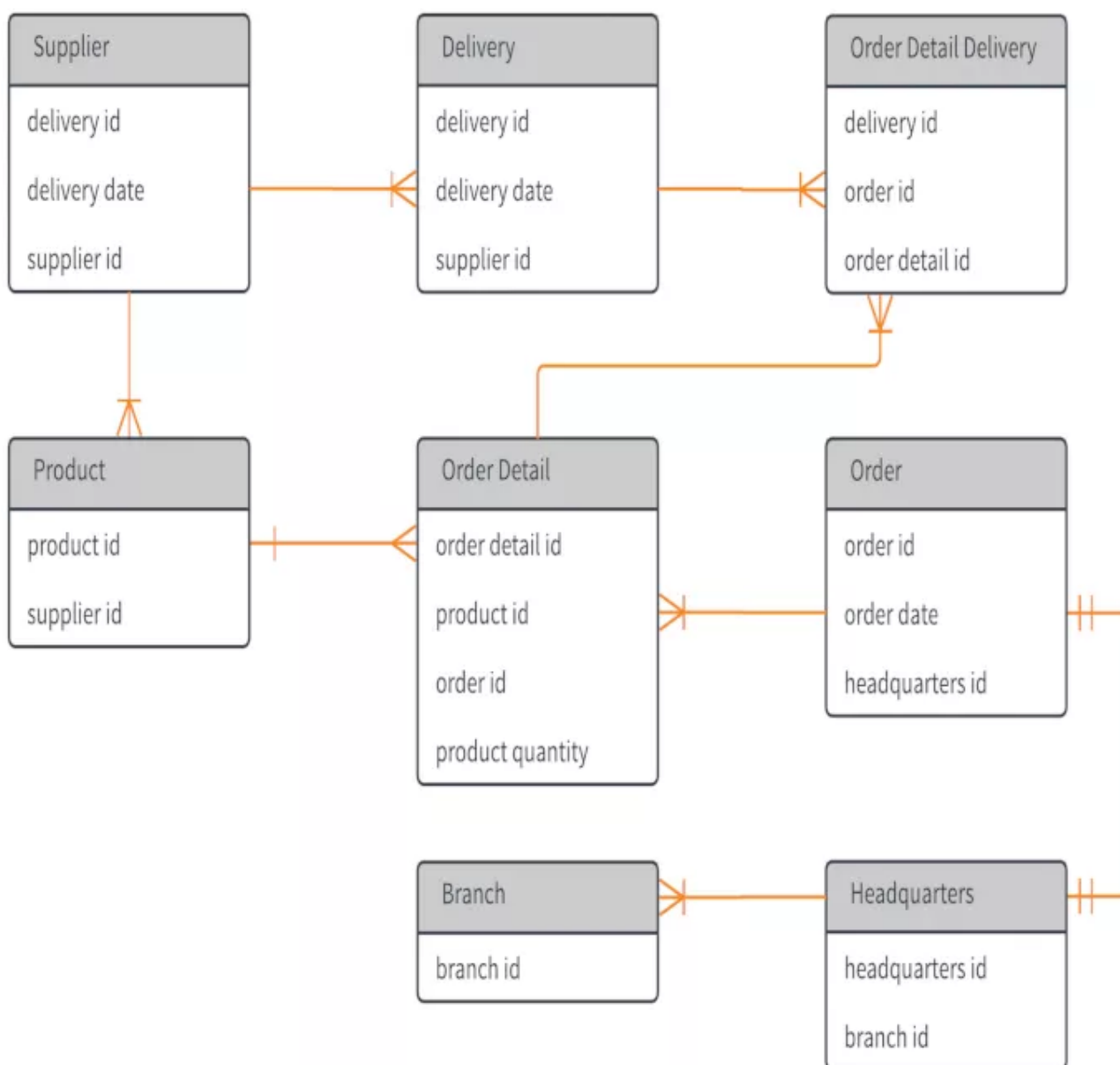
4.5.3 查询注意事项

```
1 尽量不要使用子查询、尽量不要使用 in not in
2  select * from aa1
3  where id in (select id from bb);
4  查询尽量避免join连接查询，但是这种操作咱们是永远避免不了的。
5  查询永远是小表驱动大表（永远是小结果集驱动大结果集）
```

4.5.4 数据库建模

关系型数据库最难的地方，就是建模（model）。

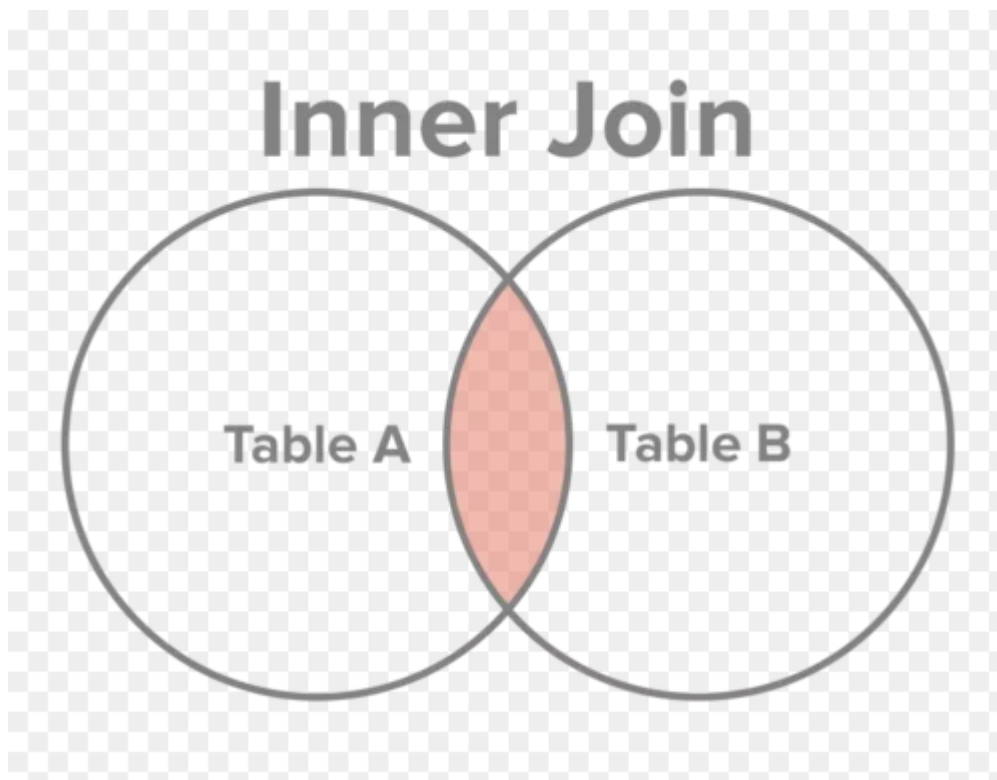
错综复杂的数据，需要建立模型，才能储存在数据库。所谓"模型"就是两样东西：实体（entity）+ 关系（relationship）ER图。实体指的是那些实际的对象，带有自己的属性，可以理解成一组相关属性的容器。关系就是实体之间的联系，通常可以分成"一对一"、"一对多"和"多对多"等类型。



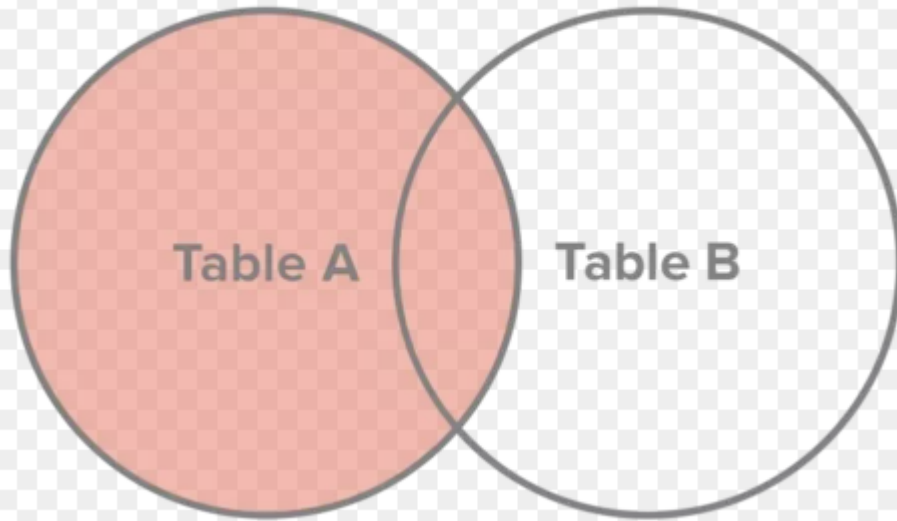
4.6 Join的语法与特点

4.6.1 表之间关系

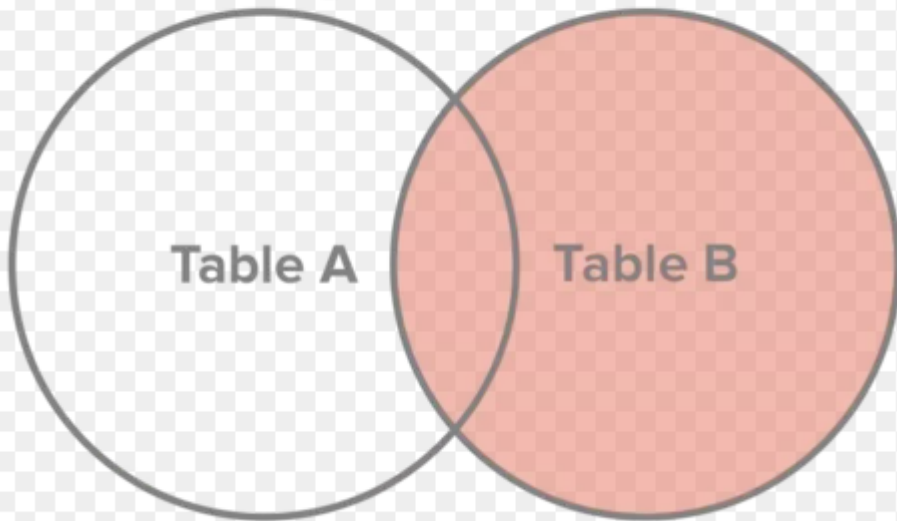
- 1 在关系型数据库里面，每个实体有自己的一张表 (table)，所有属性都是这张表的字段 (field)，表与表之间根据关联字段“**连接**” (join) 在一起。所以，表的连接是关系型数据库的核心问题。
- 2
- 3 表的连接分成好几种类型。
- 4 内连接 (inner join)
- 5 外连接 (outer join)
- 6 左连接 (left join)
- 7 右连接 (right join)
- 8 全连接 (full join)
- 9
- 10 以前，很多文章采用维恩图 (两个圆的集合运算)，解释不同连接的差异。



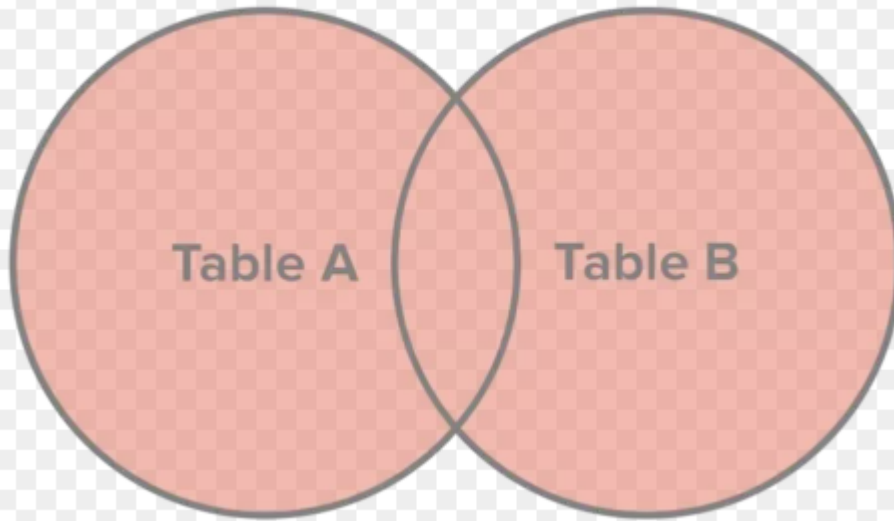
Left Join



Right Join



Full Join



1 所谓"连接", 就是两张表根据关联字段, 组合成一个数据集。问题是, 两张表的关联字段的值往往是不一致的, 如
2 果关联字段不匹配, 怎么处理? 比如, 表 A 包含张三和李四, 表 B 包含李四和王五, 匹配的只有李四这一条记
3 录。

4

5 很容易看出, 一共有四种处理方法。

6 只返回两张表匹配的记录, 这叫内连接 (inner join) 。

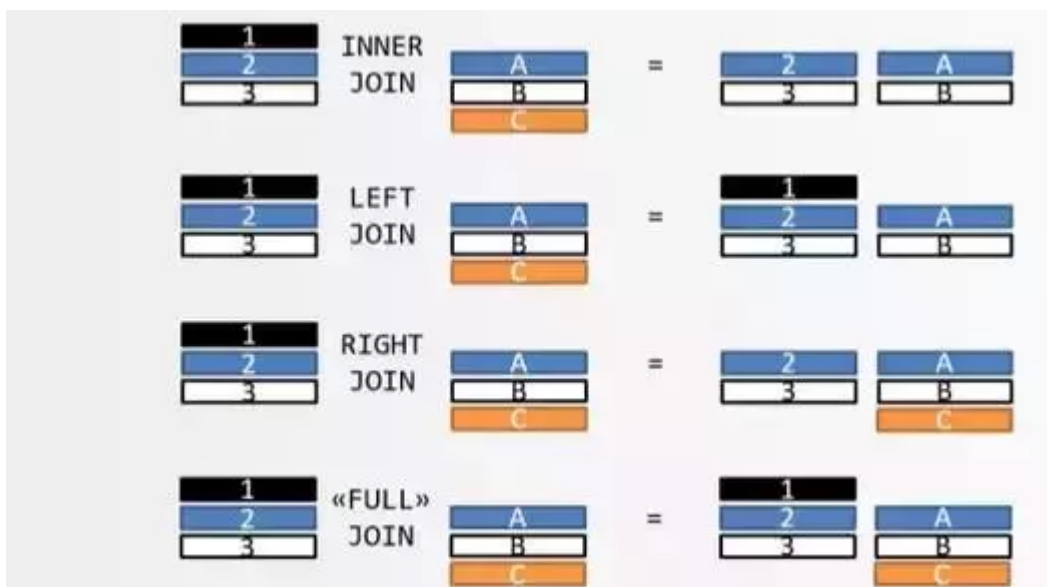
7 返回匹配的记录, 以及表 A 多余的记录, 这叫左连接 (left join) 。

8 返回匹配的记录, 以及表 B 多余的记录, 这叫右连接 (right join) 。

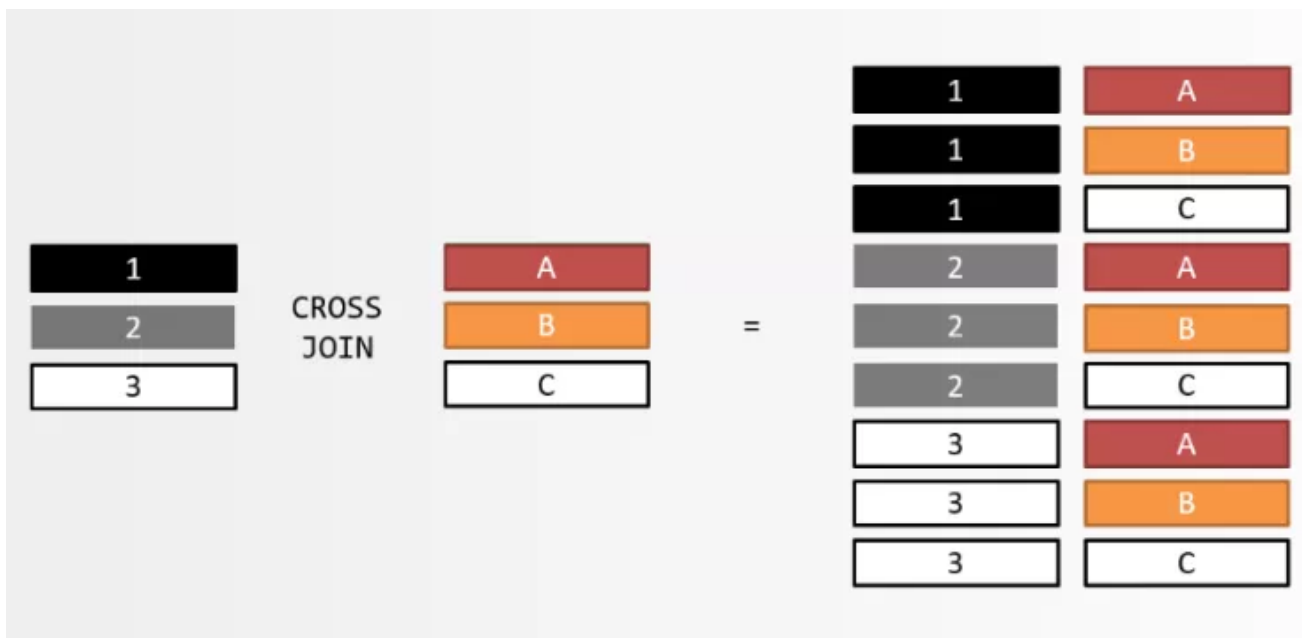
9 返回匹配的记录, 以及表 A 和表 B 各自的多余记录, 这叫全连接 (full join) 。

10

11 下图就是四种连接的图示。



- 1 上图中，表 A 的记录是 123，表 B 的记录是 ABC，颜色表示匹配关系。返回结果中，如果另一张表没有匹配的记
录，则用 null 填充。
- 2
- 3 这四种连接，又可以分成两大类：内连接（inner join）表示只包含匹配的记录，外连接（outer join）表示还包
含不匹配的记录。所以，左连接、右连接、全连接都属于外连接。
- 4
- 5 此外，还存在一种特殊的连接，叫做“交叉连接”（cross join），指的是表 A 和表 B 不存在关联字段，这时表
A（共有 n 条记录）与表 B（共有 m 条记录）连接后，会产生一张包含 $n \times m$ 条记录的新表（见下图）。



4.6.2 Join案例-语句特点演示

4.6.2.1 连接分类

- 1 类型: (left join\left outer join)
- 2 right join
- 3 right outer join
- 4 inner join
- 5 full outer join
- 6
- 7 特殊类型:
- 8 left semi join

准备数据

- 1 请参照/data/u1.txt ,/data/u2.txt

4.6.2.2 建两张表

```
1 create table if not exists u1(  
2 id int,  
3 name string  
4 )  
5 row format delimited fields terminated by ',';  
6  
7 create table if not exists u2(  
8 id int,  
9 name string  
10 )  
11 row format delimited fields terminated by ',';
```

4.6.2.3 加载数据

```
1 load data local inpath '/root/hivedata/u1.txt' into table u1;  
2 load data local inpath '/root/hivedata/u2.txt' into table u2;
```

4.6.2.4 内连接

使用关键字join、inner join、多表用逗号分开

不加任何的on 或者是where过滤条件时，称之为笛卡尔积。

```
1 select *  
2 from u1  
3 join u2 on u1.id =u2.id;  
4  
5 select *  
6 from u1  
7 inner join u2 on u1.id =u2.id;  
8  
9 select *  
10 from u1,u2  
11 where u1.id=u2.id;
```

4.6.2.5 左连接

数据以左表的数据为准，左表存在的数据都查询出，右表的数据关联上就出来，关联不上以NULL代替

```

1 left join/left outer join/left semi join :
2
3 left join 从hive0.8版本开始有
4 left join 和 left outer join 几乎差不多
5 select *
6 from u1
7 left join u2 on u1.id =u2.id;
8
9 select *
10 from u1
11 left outer join u2 on u1.id =u2.id;
12
13 上面的 SQL 语句还可以加上where条件从句，对记录进行筛选，比如只返回表 A 里面不匹配表 B 的记录。

```

4.6.2.6 右连接

以右表为准，来匹配左表信息，如果匹配不上，使用NULL来代替。

```

1 right join /right outer join
2
3 select *
4 from u1
5 right join u2 on u1.id =u2.id;
6
7
8 select *
9 from u1
10 right outer join u2 on u1.id =u2.id;
11 hive不支持right semi join:
12 这两个得到的结果也是一样的。

```

4.6.2.7 全连接

```

1 full outer join
2 (相互进行连接，如果有一张表的数据连接不上来使用NULL来代替)
3 select *
4 from u1
5 full outer join u2 on u1.id =u2.id;
6
7 另一个例子，返回表 A 或表 B 所有不匹配的记录。

```

4.6.3 Hive专有Join特点

4.6.3.1 left semi join

在Hive中，有一种专有的join操作,left semi join,我们称之为半开连接。它是left join的一种优化形式，只能查询左表的信息，主要用于解决Hive中左表的数据是否存在的问题。相当于exists关键字的用法。

```

1  select *
2  from u1
3  left semi join u2 on u1.id =u2.id;
4
5  select u1.* ,u2.*
6  from u1
7  left semi join u2 on u1.id =u2.id;
8
9  select u1.*
10 from u1
11 where exists (select 1 from u2 where u2.id =u1.id);
12
13 # left semi join 叫做半开连接，通常是left join的一种优化，只能查询左表的信息，然后主要解决Hive中存
    在不存在的问题。

```

4.6.3.2 子查询

```

1  # 演示员工表
2  create table emp
3  (
4      EMPNO    string,
5      ENAME    string,
6      JOB      string,
7      MGR      string,
8      HIREDATE string,
9      SAL      string,
10     COMM      string,
11     DEPTNO    string
12 ) row format delimited fields terminated by ','
13
14 # 演示部门表
15 create table dept
16 (
17     DEPTNO    string,
18     DNAME     string,
19     LOC       string
20 ) row format delimited fields terminated by ',';

```

准备数据:

请参照

```

1  data/emp.txt
2  data/dept.txt

```

导入数据:

```

1 load data local inpath  "/root/hivedata/emp.txt" into  table emp;
2 select * from emp;
3
4 load data local inpath  "/root/hivedata/dept.txt" into  table dept;
5 select * from dept;

```

```

1 # Hive对子查询支持不是很友好，特别是 "="问题较多
2 select
3 e.*
4 from emp e
5 where e.deptno = (
6 select
7 d.deptno
8 from dept d
9 limit 1
10 );
11
12 select
13 e.*
14 from emp e
15 where e.deptno in (
16 select
17 d.deptno
18 from dept d
19 );

```

```

1 # inner join 和outer join的区别:
2 # 分区字段对outer join 中的on条件是无效，对inner join 中的on条件有效
3
4 有inner join 但是没有full inner join
5 有full outer join但是没有outer join
6 所有join连接，只支持等值连接(= 和 and )。不支持 != 、 < 、 > 、 <> 、 >=、 <= 、 or

```

4.6.3.3 map-side join

如果所有的表中有小表，将会把小表缓存内存中，然后在map端进行连接关系查找。Hive在map端 查找时将减小查询量，从内存中读取缓存小表数据，效率较快，还省去大量数据传输和shuffle耗时

```

1 #注意使用下面属性打开Map-Join:
2 set hive.auto.convert.join=true
3 select
4 e.*
5 from u1 d
6 join u2 e
7 on d.id = e.id;

```



```

1  # 以前的老版本，需要添加(/+MAPJOIN(小表名)/)来标识该join为map端的join。hive 0.7以后hive已经废弃，
   但是仍然管用：
2
3  select
4  /+MAPJOIN(d)/
5  e.*
6  from u1 d
7  join u2 e
8  on d.id = e.id
9  ;
10 到底小表多大才会被转换为map-side join：
11  set hive.mapjoin.smalltable.filesize=25000000  约23.8MB

```

```

1  on : 所有on只支持等值连接。
2
3  where : where后面通常是表达式、还可以是非聚合函数表达式(但是不能是聚合函数表达式)
4
5  select
6  d.*
7  from dept d
8  where length(d.dname) > 5 ;

```

4.6.3.4 表达式别名

对有些表达式来说,查询的结果字段名比较难理解,这时候可以给表达式起一个别名,方便查看结果的表头.使用as把复杂表达式命名为一个容易懂的别名

```

1  select
2  d.dname,
3  length(d.dname) as nameLength
4  from dept d ;

```

4.7 查询子句

4.7.1 Where语句特点

where后不能跟聚合函数

```

1  #下面代码不能正确执行
2  select
3  e.deptno,
4  count(e.deptno) ct
5  from emp e
6  where count(e.deptno) > 3
7  group by e.deptno;

```

4.7.2 Group By语句特点

group by : 分组, 通常和聚合函数搭配使用

查询的字段要么出现在group by 后面，要么出现在聚合函数里面

```
1 select
2 e.deptno,
3 count(e.ename) ct
4 from emp e
5 group by e.deptno;
```

4.7.3 Having语句

Having是对分组以后的结果集进行过滤。

```
1 select
2 e.deptno,
3 count(e.deptno) ct
4 from emp e
5 group by e.deptno
6 having ct > 3;
```

4.7.4 Limit语句

限制从结果集中取数据的条数,一般用于分页

```
1 # 将set hive.limit.optimize.enable=true 时, limit限制数据时就不会全盘扫描, 而是根据限制的数量进行
   抽样。
2
3 同时还有两个配置项需要注意:
4 hive.limit.row.max.size          这个是控制最大的抽样数量
5 hive.limit.optimize.limit.file  这个是抽样的最大文件数量
6
7 select
8     e.deptno,
9     count(e.deptno) ct
10 from emp e
11 group by e.deptno
12 having ct > 3
13 limit 1;
```

4.7.5 排序

4.7.5.1 order by

```
1 全局排序:引发全表扫描,reduce数量一个,不建议使用
```

4.7.5.2 sort by

- 1 默认分区内排序,当reduceTask数量是1时候,那么效果和order by 一样
- 2 一般和distribute by 搭配使用

4.7.5.3 distribute by

- 1 用来确定用哪个列(字段)来分区,一般要写在sort by的前面

4.7.5.4 cluster by

- 1 分区的列和分区内排序的列相同时,那么可以用cluster by deptno来代替 distribute by deptno sort
- 2 cluster by : 兼有distribute by以及sort by的升序功能。
- 3 排序只能是升序排序(默认排序规则),不能指定排序规则为asc 或者desc。

4.7.5.5 排序详细区别

- 1 distribute by : 根据by后的字段和reducer个数,决定map的输出去往那个reducer。
- 2 默认使用查询的第一列的hash值来决定map的输出去往那个reducer。如果reducer的个数为1时没有任何体现。
- 3
- 4 sort by:局部排序,只保证单个reducer有顺序。
- 5 order by:全局排序,保证所有reducer中的数据都是有顺序。
- 6 如果reducer个数只有一个,两者都差不多。
- 7 两者都通常和 desc 、 asc 搭配。默认使用升序asc。
- 8
- 9 order by的缺点:
- 10 由于是全局排序,所以所有的数据会通过一个Reducer 进行处理,当数据结果较大的时候,一个Reducer 进行处理十分影响性能。
- 11 注意事项:
- 12 当开启MR 严格模式的时候ORDER BY 必须要设置 LIMIT 子句 , 否则会报错
- 13
- 14 -- 手动设置reducer个数:
- 15 set mapreduce.job.reduces=3;
- 16 select
- 17 e.empno
- 18 from emp e
- 19 order by e.empno desc ;
- 20
- 21 只要使用order by , reducer的个数将是1个。
- 22

4.8 表合并 Union

如果把两张表 结果联合在一起,可以用union,有下面两种用法

- union : 将多个结果集合并,去重,排序
- union all : 将多个结果集合并,不去重,不排序。

下面是具体代码演示

- 1 select

```

2  d.deptno as deptno,
3  d.dname as dname
4  from dept d
5  union
6  select
7  e.deptno as deptno,
8  e.ename as dname
9  from emp e;
10
11 select
12 d.deptno as deptno,
13 d.dname as dname
14 from dept d
15 union all
16 select
17 d.dname as dname,
18 d.deptno as deptno
19 from dept d;
20
21 -- 单个union 语句不支持: orderBy、clusterBy、distributeBy、sortBy、limit
22 -- 单个union语句字段的个数要求相同，字段的顺序要求相同。
23
24 如果想去重查询中的某列数据,可以使用 distinct关键字

```

4.9 数据类型

前面已经讲过基本数据类型,以下是复杂基本类型:

分类	类型	描述	字面量示例
复杂类型	ARRAY	有序的的同类型的集合	array(1,2)
	MAP	key-value,key必须为原始类型, value 可以任意类型	map('a',1,'b',2)
	STRUCT	字段集合,类型可以不同	struct('1',1,1.0), named_struct('col1','1','col2',1,'col3',1.0)
	UNION	在有限取值范围内的一个值	create_union(1,'a',63)

4.9.1 特殊基本类型

```

1  Java中有的而Hive中没有的:
2  long
3  char
4  short
5  byte

```

简单数据类型创表案例:

```

-- 简单数据类型创表案例:

```

```

1  -- 创建一个函数多种基本类型的表
2  create table if not exists bs1(
3  id1 tinyint,
4  id2 smallint,
5  id3 int,
6  id4 bigint,
7  sla float,
8  sla1 double,
9  isok boolean,
10 content binary,
11 dt timestamp
12 )
13 row format delimited fields terminated by ',';
14
15 -- 准备要插入的数据:
16 23, 12,342523,455345345,30000,600005,nihao, helloworld2 ,2017-06-02 11:41:30
17 12, 13,342526,455345346,80000,100000, true, helloworld1, 2017-06-02 11:41:30
18
19 load data local inpath '/root/hivedata/bs1.csv' into table bs1;

```

4.9.2 复杂的数据类型

复杂类型分为三种,分别是数组array,键值对map,和结构体struct

```

1 array : col array<基本类型> ,下标从0开始, 越界不报错, 以NULL代替
2 map   : column map<string,string>
3 struct: col struct

```

4.9.2.1 array示例

```

1  -- 数据如下: 注意下面列之间是通过TAB来分隔的
2  zhangsan    78,89,92,96
3  lisi        67,75,83,94
4
5  # 注意terminated顺序,新建数组类型
6  create table if not exists arr1
7  (
8      name string,
9      score array<string>
10 )
11 row format delimited fields terminated by '\t'
12 collection items terminated by ',';
13
14 # 导入数据
15 load data local inpath '/root/hivedata/arr1.csv' into table arr1;
16
17 # 查询:
18 select * from arr1;
19 select name,score[1] from arr1 where size(score) > 3;

```

4.9.2.1.1 列转行

就是把一列数据转化成多行,如下:

```
1 #原始数据:
2 zhangsan 90,87,63,76
3
4 #转化后数据
5 zhangsan 90
6 zhangsan 87
7 zhangsan 63
8 zhangsan 76
```

内嵌查询:

- **explode** : 展开
列表中的每个元素生成一行

```
1 select explode(score) score from arr1;
```

- **lateral view** : 虚拟表

侧视图的意义是配合explode, 一个语句生成把单行数据拆解成多行后的数据结果集。

```
1 select name,cj from arr1 lateral view explode(score) score as cj;
```

统计每个学生的总成绩:

```
1 select name,sum(cj) as totalscore from arr1 lateral view explode(score) score as cj group by
   name;
```

4.9.2.1.2 行转列

就是把多行数据转化成一行数据:

```
1 #原始数据:
2 zhangsan 90
3 zhangsan 87
4 zhangsan 63
5 zhangsan 76
6
7 #转化后数据
8 zhangsan 90,87,63,76
```

- 准备数据:

```

1 create table arr_temp
2 as
3 select name,cj
4 from arr1 lateral view explode(score) score as cj;

```

- collect_set函数:

它们都是将分组中的某列转为一个数组返回

```

1 create table if not exists arr3(
2 name string,
3 score array<string>
4 )
5 row format delimited fields terminated by ' '
6 collection items terminated by ',';

```

- 将数据写成array格式:

```

1 insert into arr3 select name,collect_set(cj) from arr_temp group by name;

```

4.9.2.2 map示例

有数据如下:

```

1 zhangsan chinese:90,math:87,english:63,nature:76
2 lisi chinese:60,math:30,english:78,nature:0
3 wangwu chinese:89,math:25,english:81,nature:9

```

- 创建map类型的表

```

1 create table if not exists map1(
2 name string,
3 score map<string,int>
4 )
5 row format delimited fields terminated by ' '
6 collection items terminated by ','
7 map keys terminated by ':';

```

- 加载数据

```

1 load data local inpath '/root/hivedata/map1.txt' into table map1;

```

- map格式数据查询

```

1  #查询数学大于35分的学生的英语和自然成绩:
2  select
3  m.name,
4  m.score['english'] ,
5  m.score['nature']
6  from map1 m
7  where m.score['math'] > 35;

```

4.9.2.2.1 map列转行

使用上面的数据

- explode 展开

```

1  select explode(score) as (m_class,m_score) from map1;

```

- Lateral view

Lateral View和split, explode等一起使用, 它能够将一行数据拆成多行数据, 并在此基础上对拆分后的数据进行聚合。

```

1  select name,m_class,m_score from map1 lateral view explode(score) score as m_class,m_score;

```

4.9.2.2.2 map行转列

准备数据:

```

1  name7,38,75,66
2  name6,37,74,65
3  name5,36,73,64
4  name4,35,72,63
5  name3,34,71,62
6  name2,33,70,61
7  name1,32,69,60

```

- 创建临时表,并加载数据

```

1  create table map_temp(
2  name string,
3  score1 int,
4  score2 int,
5  score3 int
6  )
7  row format delimited fields terminated by ',';
8  # map_temp.csv
9
10 load data local inpath '/root/hivedata/map_temp.csv' into table map_temp;

```

- 创建要导入数据map表


```

1 create table if not exists map2(
2   name string,
3   score map<string,int>
4 )
5 row format delimited fields terminated by ' '
6 collection items terminated by ','
7 map keys terminated by ':';

```

- 导入数据:

```

1 insert into map2
2 select name,map('chinese',score1,'math',score2,'english',score3) from map_temp;

```

4.9.2.3 struct

```

1 create table if not exists str1
2 (
3   name string,
4   score struct<chinese:int,math:int,english:int>
5 )
6   row format delimited fields terminated by '\t'
7   collection items terminated by ',';
8 #导入数据:
9 load data local inpath '/opt/data/arr1.csv' into table str1;
10
11 查询数据:
12 # 查询数学大于35分的学生的英语和语文,数学成绩:
13 select name,
14         score.english,
15         score.chinese,
16         score.math
17 from str1
18 where score.math > 35;

```

4.9.2.4 复杂数据类型案例

```

1 uid uname belong tax addr
2 1 xdd ll,lw,lg,lm wuxian:300,gongjijin:1200,shebao:300 北京,西城区,中南海
3 2 lkq lg,lm,lw,ll,mm wuxian:200,gongjijin:1000,shebao:200 河北,石家庄,中山路

```

```

1 #查询: 下属个数大于4个, 公积金小于1200, 省份在河北的数据
2
3 create table if not exists tax(
4   id int,
5   name string,
6   belong array<string>,
7   tax map<string,double>,

```

```

8  addr struct<province:string,city:string,road:string>
9  )
10 row format delimited fields terminated by ' '
11 collection items terminated by ','
12 map keys terminated by ':'
13 stored as textfile;
14
15 # 导入数据
16
17 load data local inpath '/root/hivedata/tax.csv' into table tax;
18
19 #查询：下属个数大于4个，公积金小于1200，省份在河北的数据
20 select id,
21 name,
22 belong[0],
23 belong[1],
24 tax['wuxian'],
25 tax['shebao'],
26 addr.road
27 from tax
28 where size(belong) > 4 and
29 tax['gongjijin'] < 1200 and
30 addr.province = '河北';

```

4.9.2.5 嵌套数据类型(了解)

在Hive中,有时候数据结构比较复杂,这时候可以用到嵌套类型,就是一个类型里面可以放置另外一个类型.这时候为了避免数据发送混乱,一般可以使用Hive自带的分隔符,

Hive共支持8个层级的分隔符, 依次是: \001,\002,\003,...\008

下面可以简单演示用一般分隔符来表示嵌套类型

```

1  -- map嵌套使用(数据为qt.txt)
2  uid uname belong tax addr
3  1   xdd wuxian:(300,300),gongjijin:1200,1500,shebao:300
4  2   lkq wuxian:(200,300),gongjijin:1000,1200,shebao:200
5
6  -- 定义嵌套类型数据表
7  create table qt(
8  id int,name string,
9  addr map<string,array<string>>)
10 row format delimited fields terminated by '\t'
11 collection items terminated by ','
12 map keys terminated by ':'
13
14 -- 加载数据类型,有些数据可能会丢失
15 load data local inpath '/root/hivedata/qt.txt' overwrite into table qt;

```

4.10 系统内置函数

4.10.1 函数查看

可以用下面两个命令查看Hive中的函数

```
1  -- 显示Hive中所有函数
2  show functions;
3
4  -- 查看某个函数的用法
5  desc function array;
```

4.10.2 日期函数

因为Hive的核心功能和海量数据统计分析,而在统计分析时日期时间是一个非常重要的维度,所以日期函数在Hive使用中尤为重要.

```
1  -- 时间戳转日期
2  select from_unixtime(1505456567);
3  select from_unixtime(1505456567, 'yyyyMMdd');
4  select from_unixtime(1505456567, 'yyyy-MM-dd HH:mm:ss');
5
6  -- 获取当前时间戳
7  select unix_timestamp();
8
9  -- 日期转时间戳
10 select unix_timestamp('2017-09-15 14:23:00');
11
12 -- 计算时间差
13 select datediff('2018-06-18', '2018-11-21');
14
15 -- 查询当月第几天
16 select dayofmonth(current_date);
17
18 -- 月末:
19 select last_day(current_date);
20
21 -- 当月第1天:
22 select date_sub(current_date, dayofmonth(current_date)-1);
23 -- 下个月第1天:
24 select add_months(date_sub(current_date, dayofmonth(current_date)-1), 1);
25
26 -- 当前日期
27 select current_date
28
29 -- 字符串转时间 (字符串必须为: yyyy-MM-dd格式)
30 select to_date('2017-01-01 12:12:12');
31
32 -- 日期、时间戳、字符串类型格式化输出标准时间格式:
33 select date_format(current_timestamp(), 'yyyy-MM-dd HH:mm:ss');
34 select date_format(current_date(), 'yyyyMMdd');
35 select date_format('2017-01-01', 'yyyy-MM-dd HH:mm:ss');
```

4.10.3 字符串函数

```
1 lower-- (转小写)
2 select lower('ABC');
3
4 upper-- (转大写)
5 select lower('abc');
6
7 length-- (字符串长度, 字符数)
8 select length('abc');
9
10 concat-- (字符串拼接)
11 select concat("A", 'B');
12
13 concat_ws -- (指定分隔符)
14 select concat_ws('-', 'a' , 'b', 'c');
15
16 substr-- (求子串)
17 select substr('abcde',3);
```

4.10.4 类型转换函数

```
1 cast(value as type) -- 类型转换
2 select cast('123' as int)+1;
```

4.10.5 数学函数

```
1 round --四舍五入((42.3 =>42))
2 select round(42.3);
3
4 ceil --向上取整(42.3 =>43)
5 select ceil(42.3);
6
7 floor --向下取整(42.3 =>42)
8 select floor(42.3);
```

4.10.6 判断为空函数

```
1 nvl (expr1, expr2)
2 #作用：将查询为Null值转换为指定值。
3 #若expr1为Null, 则返回expr2, 否则返回expr1。
4 select nvl(count,2);
```

五 实战应用

5.1 实战案例一

六 教学总结

6.1 教学重点

- 1 分桶表设置
- 2 抽样查询
- 3 分区分桶联合使用
- 4 查询语句基本语法
- 5 数据库建模
- 6 表之间关系
- 7 Join连接
- 8 查询子句
- 9 排序
- 10 合并结果集
- 11 数据类型:array,map,struct
- 12 系统内置函数:常见函数

七 课后作业

7.1

创建一个分桶表,并且练习分桶表的数据插入:

答案:

1

7.2

练习分桶表的抽样查询,总共分8个桶,要求分三种情况抽样,可以取一桶,取4桶,取8桶的情况下,写出SQL语句

答案:

1

7.3

注意: 下面题目的数据参考Hive第一天的数据

1 查询男生、女生人数:

1

2 查询出只选修一门课程的全部学生的学号和姓名:

1

3 查询1981年出生的学生名单

1

4 查询平均成绩大于80的所有学生的学号、姓名和平均成绩：

1

5 查询每门课程的平均成绩，结果按平均成绩升序排序，平均成绩相同时，按课程号降序排列：

1

6 查询课程名称为“数学”，且分数低于60的学生名字和分数：

1

7 查询所有学生的选课情况：

1

8 查询任何一门课程成绩在70分以上的姓名、课程名称和分数：

1

9 查询01课程比02课程成绩高的所有学生的学号

1

10 查询平均成绩大于60分的同学的学号和平均成绩

1

11 查询所有同学的学号、姓名、选课数、总成绩

1

12 查询所有课程成绩小于60的同学的学号、姓名：

1

13 查询没有学全所有课的同学的学号、姓名：

1

14 查询至少有一门课与学号为01同学所学相同的同学的学号和姓名：

1

15 查询至少学过学号为01同学所有一门课的其他同学学号和姓名；

1

八 解决方案

8.1 应用场景

8.2 核心面试题