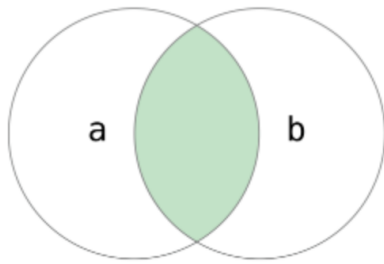# 第二节——HiveSQL基础技能2

## 课堂目标

1、掌握inner join

2、掌握left join

3、掌握full join和union all

4、掌握各种连接的组合

## 表连接

# 【inner join】（25min）



内连接：返回两个表的交集

举例说明：

表1：user_list_1

| user_id | user_name |
| --- | --- |
| 10001 | Abby |
| 10002 | Ailsa |
| 10003 | Alice |
| 10004 | Alina |
| 10005 | Allison |
| 10006 | Angelia |

表2：user_list_2

| user_id | user_name |
| --- | --- |
| 10001 | Abby |
| 10003 | Alice |
| 10004 | Alina |
| 10007 | Amanda |
| 10008 | Anne |
| 10009 | Ann |

找出既在**user_list_1**也在**user_list_2**的用户：

```
1   --既在user_list_1也在user_list_2的用户--
2   SELECT *
3   FROM user_list_1 a JOIN user_list_2 b ON a.user_id=b.user_id;
```

```
10001   Abby    10001   Abby
10003   Alice   10003   Alice
10004   Alina   10004   Alina
Time taken: 26.062 seconds, Fetched: 3 row(s)
```

**注意：**

- 表连接时，必须进行重命名
- on后面使用的连接条件必须起到唯一键值的作用
- inner可省略不写，效果一样

练习：

| user_trade列名 | 举例 |
|---|---|
| user_name | Amy, Dennis |
| piece | 购买数量 |
| price | 价格 |
| pay_amount | 支付金额 |
| goods_category | food, clothes, book, computer, electronics, shoes |
| pay_time | 1323308943，时间戳 |
| dt | partition，'yyyy-mm-dd' |

| user_refund列名 | 举例 |
|---|---|
| user_name | Amy, Dennis |
| refund_piece | 退款件数 |
| refund_amount | 退款金额 |
| refund_time | 1323308943，时间戳 |
| dt | partition，'yyyy-mm-dd' |

**1、在2019年购买后又退款的用户：**

```
1    --在2019年购买后又退款的用户--
2    SELECT a.user_name
3    FROM
4         (SELECT distinct user_name
5          FROM user_trade
6          WHERE year(dt)=2019)a
7        JOIN
8          (SELECT distinct user_name
9          FROM user_refund
10         WHERE year(dt)=2019)b on a.user_name=b.user_name;
```

**常见错误：没有写别名**

```
1    SELECT user_name
2    FROM
3         (SELECT distinct user_name
4          FROM user_trade
5          WHERE year(dt)=2019)a
6        JOIN
7          (SELECT distinct user_name
8          FROM user_refund
9          WHERE year(dt)=2019)b on a.user_name=b.user_name;
```

```
hive> SELECT user_name
    > FROM
    > (SELECT distinct user_name
    > FROM user_trade
    > WHERE year(dt)=2019)a
    > JOIN
    > (SELECT distinct user_name
    > FROM user_refund
    > WHERE year(dt)=2019)b on a.user_name=b.user_name;
FAILED: SemanticException Column user_name Found in more than One Tables/Subqueries
```

如果不进行去重的话，会怎么样?

```
1    SELECT count(a.user_name),
2         count(distinct a.user_name)
3    FROM
4         (SELECT user_name
5          FROM user_trade
6          WHERE year(dt)=2019)a
7        JOIN
8          (SELECT user_name
9          FROM user_refund
10         WHERE year(dt)=2019)b on a.user_name=b.user_name;
```

```
61        31
Time taken: 38.773 seconds, Fetched: 1 row(s)
```

**2、在2017年和2018年都购买的用户：**

```
1   --在2017年和2018年都购买的用户--
2   SELECT a.user_name
3   FROM
4        (SELECT distinct user_name
5        FROM user_trade
6        WHERE year(dt)=2017)a
7     JOIN
8        (SELECT distinct user_name
9        FROM user_trade
10       WHERE year(dt)=2018)b on a.user_name=b.user_name;
```

| trade_2017&2018&2019列名 | 举例 |
|---|---|
| **user_name** | Amy, Dennis |
| **amount** | 金额 |
| **trade_time** | 交易时间，2017-02-05 06:31:50 |

**3、在2017年、2018年、2019都有交易的用户：**

```
1    --在2017年、2018年、2019都有交易的用户--
2    ##第一种写法
3    SELECT distinct a.user_name
4    FROM trade_2017 a
5    JOIN trade_2018 b on a.user_name=b.user_name
6    JOIN trade_2019 c on b.user_name=c.user_name;
7
8    ##第二种写法
9    SELECT a.user_name
10   FROM
11       (SELECT distinct user_name
12       FROM trade_2017)a
13     JOIN
14       (SELECT distinct user_name
15       FROM trade_2018)b on a.user_name=b.user_name
16     JOIN
17       (SELECT distinct user_name
18       FROM trade_2019)c on b.user_name=c.user_name;
```

- 在表的数据量级很大时，推荐第二种写法
- 为什么写a.user_name呢？b.user_name或者c.user_name可以吗？——当然可以。

因为JOIN是取表的交集，所以不管取哪个表的这个字段都是一样的结果。

**初学者常犯错误写法：**

```
1    ----错误写法----
2    SELECT distinct a.user_name
3    FROM trade_2017 a
4    JOIN trade_2018 b
5    JOIN trade_2019 c on a.user_name=b.user_name=c.user_name;
```

# 【left join】（25min）

首先，我们先来看一下，对表1和表2进行左连接后，发生了什么。

```
1    SELECT *
2    FROM user_list_1 a LEFT JOIN user_list_2 b ON a.user_id=b.user_id;
```

```
10001   Abby    10001   Abby
10002   Ailsa   NULL    NULL
10003   Alice   10003   Alice
10004   Alina   10004   Alina
10005   Allison NULL    NULL
10006   Angelia NULL    NULL
Time taken: 25.652 seconds, Fetched: 6 row(s)
```
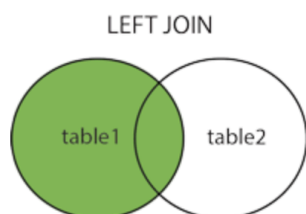
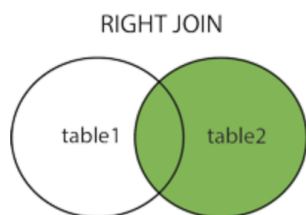进行左连接后，以左边的表**1**为全集，返回能够匹配上的右边表**2**的匹配结果，没有匹配上的则显示**NULL**。

拓展：

**right join**：以右边的表为全集，返回能够匹配上的左表的匹配结果，没有匹配上的则显示NULL

但其完全可以由left join改写出同样的结果，所以较少使用

**LEFT JOIN** 关键字从左表（**table1**）返回所有的行，即使右表（**table2**）中没有匹配。如果右表中没有匹配，则结果为 **NULL**。
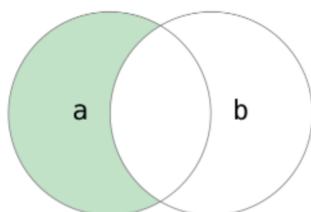


LEFT JOIN

**RIGHT JOIN** 关键字从右表（**table2**）返回所有的行，即使左表（**table1**）中没有匹配。如果左表中没有匹配，则结果为 **NULL**。



RIGHT JOIN

举例说明：

如何取出，在**user_list_1**表中但是不在**user_list_2**的用户？



```
1  --在user_list_1表中但是不在user_list_2的用户--
2  SELECT a.user_id,
3         a.user_name
4  FROM user_list_1 a LEFT JOIN user_list_2 b ON a.user_id=b.user_id
5  WHERE b.user_id is null;
```

练习：

**1、在2019年购买，但是没有退款的用户：**

```
1   --在2019年购买，但是没有退款的用户--
2   SELECT a.user_name
3   FROM
4       (SELECT distinct user_name
5       FROM user_trade
6       WHERE year(dt)=2019)a
7     LEFT JOIN
8       (SELECT distinct user_name
9       FROM user_refund
10      WHERE year(dt)=2019)b on a.user_name=b.user_name
11  WHERE b.user_name is null;
```

**2、在2019年购买用户的学历分布：**

| user_info 列名 | 举例 |
| --- | --- |
| user_id | 10001,10002 |
| user_name | Amy, Dennis |
| sex | [male, female] |
| age | [13,70] |
| city | beijing, shanghai |
| firstactivetime | 2019-04-19 15:40:00 |
| level | [1,10] |
| extra1 | string类型：{"systemtype":"ios","education":"master","marriage_status":"1","phonebrand":"iphone X"} |
| extra2 | map<string,string>类型：{"systemtype":"ios","education":"master","marriage_status":"1","phonebrand":"iphone X"} |

```
1   --在2019年购买用户的学历分布--
2   SELECT b.education,
3         count(a.user_name)
4   FROM
5       (SELECT distinct user_name
6       FROM user_trade
7       WHERE year(dt)=2019)a
8     LEFT JOIN
9       (SELECT user_name,
10              get_json_object(extra1, '$.education') as education
11      FROM user_info)b on a.user_name=b.user_name
12  GROUP BY b.education;
```

注意：get_json_object(extra1, '$.education')可以换成extra2['education']

### 3、在2017和2018年都购买，但是没有在2019年购买的用户：

| trade_2017&2018&2019列名 | 举例 |
| --- | --- |
| **user_name** | Amy, Dennis |
| **amount** | 金额 |
| **trade_time** | 交易时间，2017-02-05 06:31:50 |

```
1   --在2017和2018年都购买，但是没有在2019年购买的用户--
2   SELECT a.user_name
3   FROM
4        (SELECT distinct user_name
5        FROM trade_2017)a
6      JOIN
7        (SELECT distinct user_name
8        FROM trade_2018)b on a.user_name=b.user_name
9      LEFT JOIN
10       (SELECT distinct user_name
11       FROM trade_2019)c on b.user_name=c.user_name
12  WHERE c.user_name is null;
```

- a.user_name换成b.user_name也可以

这种写法行不行？——可以的！

```
1   SELECT c.user_name
2   FROM
3        (SELECT distinct a.user_name
4        FROM trade_2017 a JOIN trade_2018 b on a.user_name=b.user_name)c
5      LEFT JOIN
6        (SELECT distinct user_name
7        FROM trade_2019)d on c.user_name=d.user_name
8   WHERE d.user_name is null;
```

这种写法行不行？——可以，但不推荐！

```
1  SELECT distinct a.user_name
2  FROM trade_2017 a
3  JOIN trade_2018 b on a.user_name=b.user_name
4  LEFT JOIN trade_2019 c on b.user_name=c.user_name
5  WHERE c.user_name is null;
```

注意：如果表比较小的时候，这样写影响不大。但是有分区的大表，这样写执行速度很慢。

# 【full join】（10min）

## 举例说明：

首先，我们先来看一下，对表1和表2进行全连接后，发生了什么。

```
1  SELECT *
2  FROM user_list_1 a FULL JOIN user_list_2 b ON a.user_id=b.user_id;
```
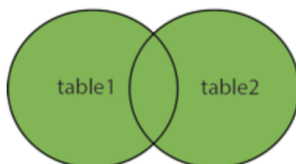


**FULL OUTER JOIN** 关键字只要左表（**table1**）和右表（**table2**）其中一个表中存在匹配，则返回行。**FULL OUTER JOIN** 关键字结合了 **LEFT JOIN** 和 **RIGHT JOIN** 的结果。



## 练习：

**1、user_list_1和user_list_2的所有用户：**

```
1   --user_list_1和user_list_2的所有用户--
2   SELECT coalesce(a.user_name,b.user_name)
3   FROM user_list_1 a FULL JOIN user_list_2 b on a.user_id=b.user_id;
```

```
Abby
Ailsa
Alice
Alina
Allison
Angelia
Amanda
Anne
Ann
Time taken: 26.578 seconds, Fetched: 9 row(s)
```

拓展：coalesce是一个函数，(expression_1, expression_2, …,expression_n)依次参考各参数表达式，遇到非null值即停止并返回该值。如果所有的表达式都是空值，最终将返回一个空值。


## 【union all】 （20min）

表1：user_list_1

| user_id | user_name |
| --- | --- |
| 10001 | Abby |
| 10002 | Ailsa |
| 10003 | Alice |
| 10004 | Alina |
| 10005 | Allison |
| 10006 | Angelia |

表3：user_list_3

| user_id | user_name |
| --- | --- |
| 10290 | Michael |
| 10291 | Avery |
| 10292 | Reilly |
| 10293 | Dillon |
| 10294 | Walton |

**union all**：联合所有

举例说明：

将**user_list_1**和**user_list_3**合并在一起：

```
1  SELECT user_id,
2         user_name
3  FROM user_list_1
4  UNION ALL
5  SELECT user_id,
6         user_name
7  FROM user_list_3;
```

```
10001   Abby
10002   Ailsa
10003   Alice
10004   Alina
10005   Allison
10006   Angelia
10290   Michael
10291   Avery
10292   Reilly
10293   Dillon
10294   Walton
Time taken: 22.26 seconds, Fetched: 11 row(s)
```

注意：

- 字段名称必须一致

- 字段顺序必须一致
- 没有连接条件

```
1   ----错误写法----
2   SELECT user_name,
3          user_id
4   FROM user_list_1
5   UNION ALL
6   SELECT user_id,
7          user_name
8   FROM user_list_3;
```

```
Abby     10001
Ailsa    10002
Alice    10003
Alina    10004
Allison  10005
Angelia  10006
10290    Michael
10291    Avery
10292    Reilly
10293    Dillon
10294    Walton
Time taken: 17.473 seconds, Fetched: 11 row(s)
```

练习：

**1、2017-2019年有交易的所有用户数**：

| trade_2017&2018&2019列名 | 举例 |
| --- | --- |
| **user_name** | Amy, Dennis |
| **amount** | 金额 |
| **trade_time** | 交易时间，2017-02-05 06:31:50 |

```
1   --2017-2019年有交易的所有用户数--
2   ##写法一
3   SELECT count(distinct a.user_name),
4           count(a.user_name)
5   FROM
6       (
7           SELECT user_name
8           FROM trade_2017
9       UNION ALL
10          SELECT user_name
11          FROM trade_2018
12      UNION ALL
13          SELECT user_name
14          FROM trade_2019)a;
15
16
17  ##写法二
18  SELECT count(distinct a.user_name),
19          count(a.user_name)
20  FROM
21      (
22          SELECT user_name
23          FROM trade_2017
24      UNION
25          SELECT user_name
26          FROM trade_2018
27      UNION
28          SELECT user_name
29          FROM trade_2019)a;
```

- **UNION ALL和UNION的区别：**

| 对比 | UNION ALL | UNION |
|------|-----------|-------|
| 对重复结果的处理 | 不会去除重复记录 | 在进行表连接后会筛选掉重复的记录 |
| 对排序的处理 | 只是简单的将两个结果合并后就返回 | 将会按照字段的顺序进行排序 |
| 效率 | 更快 | 更慢 |
| 总述 | 不去重不排序 | 去重且排序 |

PS：如果表很大时推荐先去重，再进行union all。

- 常见错误——没有对union all后的表进行重命名:

```
1   ----错误写法----
2   SELECT count(distinct user_name)
3   FROM
4     (
5        SELECT user_name
6        FROM trade_2017
7      UNION ALL
8        SELECT user_name
9        FROM trade_2018
10     UNION ALL
11       SELECT user_name
12       FROM trade_2019);
```

- 常见错误——直接对表进行UNION ALL:

```
1   ----错误写法----
2   SELECT count(distinct user_name)
3   FROM trade_2017
4   UNION ALL trade_2018
5   UNION ALL trade_2019;
```

## 2、2019年每个用户的支付和退款金额汇总:

| user_trade列名 | 举例 |
| --- | --- |
| user_name | Amy, Dennis |
| piece | 购买数量 |
| price | 价格 |
| pay_amount | 支付金额 |
| goods_category | food, clothes, book, computer, electronics, shoes |
| pay_time | 1323308943,时间戳 |
| dt | partition,'yyyy-mm-dd' |

| user_refund列名 | 举例 |
|---|---|
| user_name | Amy, Dennis |
| refund_piece | 退款件数 |
| refund_amount | 退款金额 |
| refund_time | 1323308943，时间戳 |
| dt | partition，'yyyy-mm-dd' |

```
1   --2019年每个用户的支付和退款金额汇总--
2   SELECT a.user_name,
3          sum(a.pay_amount),
4          sum(a.refund_amount)
5   FROM
6     (
7       SELECT user_name,
8              sum(pay_amount) as pay_amount,
9              0 as refund_amount
10      FROM user_trade
11      WHERE year(dt)=2019
12      GROUP BY user_name
13    UNION ALL
14      SELECT user_name,
15             0 as pay_amount,
16             sum(refund_amount) as refund_amount
17      FROM user_refund
18      WHERE year(dt)=2019
19      GROUP BY user_name
20    )a
21  GROUP BY a.user_name;
```

- 如何用**full join**来实现该题目？

```sql
SELECT coalesce(a.user_name,b.user_name),
       a.pay_amount,
       b.refund_amount
FROM
    (SELECT user_name,
            sum(pay_amount) as pay_amount
     FROM user_trade
     WHERE year(dt)=2019
     GROUP BY user_name)a
  FULL JOIN
    (SELECT user_name,
            sum(refund_amount) as refund_amount
     FROM user_refund
     WHERE year(dt)=2019
     GROUP BY user_name)b on a.user_name=b.user_name;
```

与第一种写法有什么结果差异？——没有退款的人，退款金额显示NULL

**如果用第二种写法，如何把NULL都变成0呢？**

```sql
SELECT coalesce(a.user_name,b.user_name),
       if(a.pay_amount is null,0,a.pay_amount),
       if(b.refund_amount is null,0,b.refund_amount)
FROM
    (SELECT user_name,
            sum(pay_amount) as pay_amount
     FROM user_trade
     WHERE year(dt)=2019
     GROUP BY user_name)a
  FULL JOIN
    (SELECT user_name,
            sum(refund_amount) as refund_amount
     FROM user_refund
     WHERE year(dt)=2019
     GROUP BY user_name)b on a.user_name=b.user_name;
```

问题变形：**2019年每个支付用户的支付金额和退款金额**

```
1   SELECT a.user_name,
2          a.pay_amount,
3          b.refund_amount
4   FROM
5       (SELECT user_name,
6              sum(pay_amount) as pay_amount
7        FROM user_trade
8        WHERE year(dt)=2019
9        GROUP BY user_name)a
10     LEFT JOIN
11       (SELECT user_name,
12              sum(refund_amount) as refund_amount
13        FROM user_refund
14        WHERE year(dt)=2019
15        GROUP BY user_name)b on a.user_name=b.user_name;
```

# 重点练习（**20min**）

**1**、首次激活时间在**2017年**，但是一直没有支付的用户年龄段分布(使用**user_trade**和**user_info**两个表)

| user_info 列名 | 举例 |
|---|---|
| **user_id** | 10001,10002 |
| **user_name** | Amy, Dennis |
| **sex** | [male, female] |
| **age** | [13,70] |
| **city** | beijing, shanghai |
| **firstactivetime** | 2019-04-19 15:40:00 |
| **level** | [1,10] |
| **extra1** | string类型：{"systemtype":"ios","education":"master","marriage_status":"1","phonebrand":"iphone X"} |
| **extra2** | map<string,string>类型：{"systemtype":"ios","education":"master","marriage_status":"1","phonebrand":"iphone X"} |

| user_trade列名 | 举例 |
|---|---|
| user_name | Amy, Dennis |
| piece | 购买数量 |
| price | 价格 |
| pay_amount | 支付金额 |
| goods_category | food, clothes, book, computer, electronics, shoes |
| pay_time | 1323308943，时间戳 |
| dt | partition，'yyyy-mm-dd' |

```
1   --首次激活时间在2017年，但是一直没有支付的用户年龄段分布--
2   SELECT a.age_level,
3          count(a.user_name)
4   FROM
5        (SELECT user_name,
6                case when age<20 then '20岁以下'
7                     when age>=20 and age<30 then '20-30岁'
8                     when age>=30 and age<40 then '30-40岁'
9                     else '40岁以上' end as age_level
10        FROM user_info
11        WHERE year(firstactivetime)=2019)a
12      LEFT JOIN
13        (SELECT distinct user_name
14        FROM user_trade
15        WHERE dt>0)b on a.user_name=b.user_name
16   WHERE b.user_name is null
17   GROUP BY a.age_level;
```

```
20-30岁  6
20岁以下           8
30-40岁  6
40岁以上          22
Time taken: 86.026 seconds, Fetched: 4 row(s)
```

**常见错误——没有对子查询中的字段进行重命名：**

```
1    ----错误写法----
2    SELECT a.age,
3           count(a.user_name)
4    FROM
5        (SELECT user_name,
6                case when age<20 then '20岁以下'
7                     when age>=20 and age<30 then '20-30岁'
8                     when age>=30 and age<40 then '30-40岁'
9                     else '40岁以上' end
10       FROM user_info
11       WHERE year(firstactivetime)=2019)a
12     LEFT JOIN
13       (SELECT distinct user_name
14       FROM user_trade
15       WHERE dt>0)b on a.user_name=b.user_name
16   WHERE b.user_name is null
17   GROUP BY a.age;
```

**2、2018、2019年交易的用户，其激活时间段分布(使用trade_2018、trade_2019和user_info三个表)**

| trade_2017&2018&2019列名 | 举例 |
| --- | --- |
| **user_name** | Amy, Dennis |
| **amount** | 金额 |
| **trade_time** | 交易时间，2017-02-05 06:31:50 |

| user_info 列名 | 举例 |
| --- | --- |
| user_id | 10001,10002 |
| user_name | Amy, Dennis |
| **sex** | [male, female] |
| **age** | [13,70] |
| **city** | beijing, shanghai |
| **firstactivetime** | 2019-04-19 15:40:00 |
| **level** | [1,10] |
| **extra1** | string类型：{"systemtype":"ios","education":"master","marriage_status":"1","phonebrand":"iphone X"} |
| **extra2** | map<string,string>类型：{"systemtype":"ios","education":"master","marriage_status":"1","phonebrand":"iphone X"} |

```sql
--2018、2019年交易的用户，其激活时间段分布--
SELECT hour(firstactivetime),
       count(a.user_name)
FROM
    (
        SELECT user_name
        FROM trade_2018
    UNION
        SELECT user_name
        FROM trade_2019)a
    LEFT JOIN user_info b on a.user_name=b.user_name
GROUP BY hour(firstactivetime);
```

```
0     10
1     4
2     5
3     7
4     7
5     5
6     6
7     7
8     5
9     15
10    3
11    5
12    10
13    6
14    6
15    4
16    10
17    10
18    6
19    13
20    14
21    12
22    9
23    6
Time taken: 53.367 seconds, Fetched: 24 row(s)
```

## 总结

1. 在实际业务场景中，熟练选择JOIN、LEFT JOIN来解决具体问题
2. 区分好FULL JOIN和UNION ALL的使用场景
3. 在多表连接时，注意各种细节和业务逻辑
4. 避免常见错误

# 作业

作业1：在2019年购买后又退款的用户性别分布(使用user_trade、user_refund)

作业2：在2018年购买，但是没在2019年购买的用户城市城市分布(使用user_trade、user_refund)

作业3：2017-2019年，有交易但是没退款的用户的手机品牌分布(使用trade_2017、trade_2018、trade_2019、user_refund、user_info)