# **SQL**

基础:选择select; 表连接join/left join/union

常用:去重distinct;聚合max/min/sum/count/group by;筛选having/where;排序order

by/sort by;条件case when...end;字符串substr/concat/split

进阶:日期函数to\_date/dateoff;分组排除row\_number();取百分比percentile

高级:窗口函数(分组排序row number/rank/dense rank;分组取最大最小

first\_value/last\_value;累积百分比cume\_dist/sum()over;错位lead/lag);日期函数(时间和日期戳转化unix timestamp/from unixtime;日期操作

to\_date/datesub/date\_add/dateoff); 字符串函数 (substr,concat,concat\_ws,length,upper这类基础函数; regexp\_replace正则/get\_json\_object); 其他常用函数 (lateral view explode行列转换; rand等随机抽样)

SQL是一种操作数据库的语言,主要包括CRUD,对于数据分析师来说,一般只要查找数据的,不能对数据库进行增删改

- 1 SELECT comlumns name --查找一列或多了之间用逗号隔开
- 2 FROM Table
- --目标表
- 3 WHERE condition --过滤条件
- 4 GROUP BY columns\_name --按列值分组,可以1个或多个列
- 5 HAVING condition --分组后的筛选条件, HAVING和WHERE的区别在于前者表达式中可包含函数 6 ORDER BY columns name --按列排序
- 7 LIMIT start, row count --对结果进行限定, start表示从哪行开始, row count表示结果行数

# 一、基础用法

#### 1.简单查询

#### 2.条件查询

有一张学生表student,包括学生id,姓名,年龄,班级,分数字段,下面用SQL来查找数据吧。

stu_id	stu_name	age	class	grade
1	小明	17	1	60
2	小花	18	2	90
3	小王	16	2	70
4	小李	20	1	20
5	小陈	17	2	60
6	小赵	15	和平@	80元效居

student表

- 1 -- 查找id. 姓名, 成绩列
- 2 | SELECT stu\_id, stu\_name, grade FROM student;

stu_id stu_name grade  1 小明 60 2 小花 90 3 小王 70 4 小李 20 5 小陈 60 6 小赵 80	信息 结	果1 概况	状态
2 小花     90       3 小王     70       4 小李     20       5 小陈     60	stu_id	stu_name	grade
3 小王     70       4 小李     20       5 小陈     60	1	小明	60
4 小李 20 5 小陈 60	2	小花	90
5 小陈 60	3	小王	70
	4	小李	20
6 小赵 80	5	小陈	60
	6	小赵	80

- 1 -- 查找成绩大于70分的学生id, 姓名, 成绩列
- 2 | SELECT s.stu id, s.stu name, s.grade FROM student s WHERE s.grade >=70;

#### 执行结果



# 3.聚合函数

- 1 --计算班级人数,平均成绩,最高成绩,最低成绩
- 2 SELECT s.class, COUNT (s.stu id) as

count1,AVG(s.grade),MAX(s.grade),MIN(s.grade) FROM student s GROUP BY s.class;

class         count1         AVG(s.grade)         MAX(s.grade)         MIN(s.grade)           1         3         53.33333333333333333333333333333333333	1	信息	结	果1	概况	状态			
		class		coun	ıt1	AVG(s.gr	ade)	MAX(s.grade)	MIN(s.grade)
2 3 73 232323232323 00 60	Þ	1			3	53.33333	3333333336	80	20
2 3 73.33333333333		2			3	73.33333	33333333	90	60

注释:计算班级平均成绩用到了group by语句,同样的class被分为一组。当使用分组语句时,select语句后面只能出现分组字段和聚合函数。

## 4.分组后筛选

```
1 | --Having语句筛选班级平均分大于60分的班级
2 | SELECT s.class,COUNT(s.stu_id) as count1,AVG(s.grade),MAX(s.grade),MIN(s.grade) FROM student s GROUP BY s.class HAVING AVG(s.grade)>60;
```

#### 执行结果

	息	结果1	概况	状态			
	class	cour	nt1	AVG(s.g	rade)	MAX(s.grade)	MIN(s.grade)
١	2		3	73.3333	333333333	90	60

#### 5.模糊查询

```
1 | --like进行模糊匹配
2 | where name like '陈%' | --找出姓陈的人
3 | where name like '%铭%' | --找出名字中有铭字的人
4 | where name like '陈_' | --找出姓陈且名字为两个字的人
```

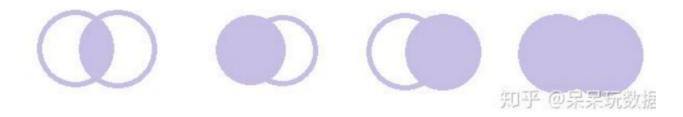
注释: like进行模糊匹配时, %可以匹配任意多个任意字符, \_匹配任意一个字符。

# 二、表连接

表连接是指两张表可以通过想用的关键字段进行连接,包括内连接和外连接。外连接包括左外连接、右 外连接和全外连接。

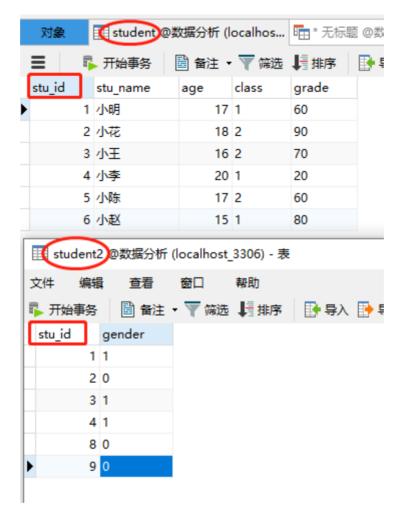
- 1.内连接指结果表只包含同时存在两张表中的连接字段
- 2.左外连接是指既包含内连接的连接字段还包含左表未连接的字段
- 3.右外连接是指既包含内连接的连接字段还包含右表未连接的字段
- 4.全外连接是指包含两张表的所有连接字段.

#### 韦恩图



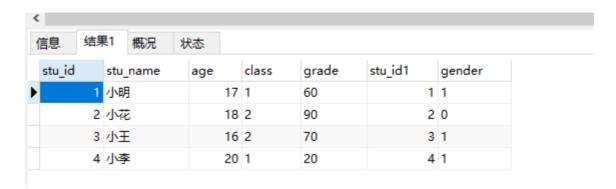
从左到右,依次是内连接,左外连接,右外连接,全外连接

表连接,两张表student, student2, stu\_id为连接字段



## 1.内连接

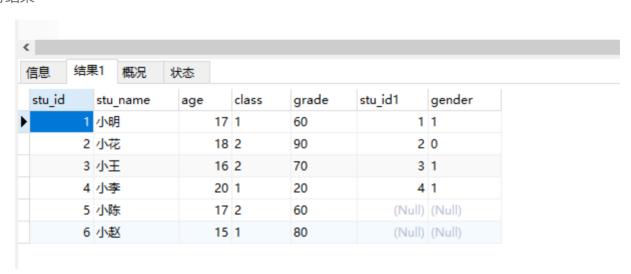
```
1 | -- 内连接, *是通配符, 取出表中所有变量
2 | SELECT s1.*,s2.* FROM student s1 INNER JOIN student2 s2 ON s1.stu_id=s2.stu_id;
```



# 2.左外连接

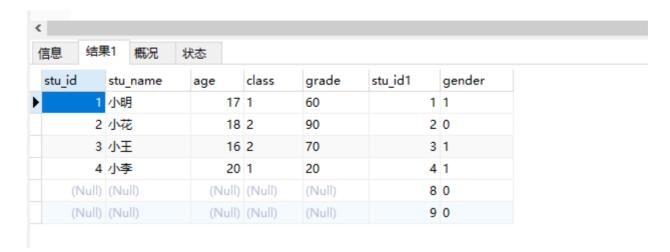
```
1 | -- 左外连接
2 | SELECT s1.*,s2.* FROM student s1 LEFT JOIN student2 s2 ON s1.stu_id=s2.stu_id;
```

### 执行结果



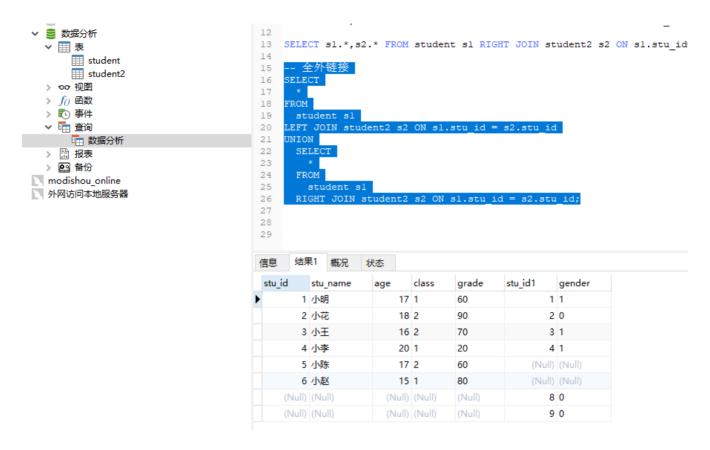
# 3.右外连接

```
1 -- 右外连接
2 SELECT s1.*,s2.* FROM student s1 RIGHT JOIN student2 s2 ON s1.stu_id=s2.stu_id;
```



## 4.全外连接

```
-- 全外连接
2
   SELECT
3
4
   FROM
5
    student s1
   LEFT JOIN student2 s2 ON s1.stu id = s2.stu id
6
7
   UNION
8
      SELECT
9
      FROM
11
          student s1
       RIGHT JOIN student2 s2 ON s1.stu_id = s2.stu_id;
12
```



# 三、子查询

子查询是指把一个查询的结果放到另一个查询里面使用,功能十分丰富,既可以放在from语句中作为临时表供另一个查询使用,也可以放在where语句子句后面进行过滤。

放在where子句后面的子查询返回结果可以是单行,也可以是多行。

单行子查询一般与>,<,=,<=,>=等比较符一起使用。

而多行子查询与in(等于子查询返回结果的任意一个), any(比较符满足子查询中的任意一个情况), all(比较符满足子查询中所有情况)等配套使用

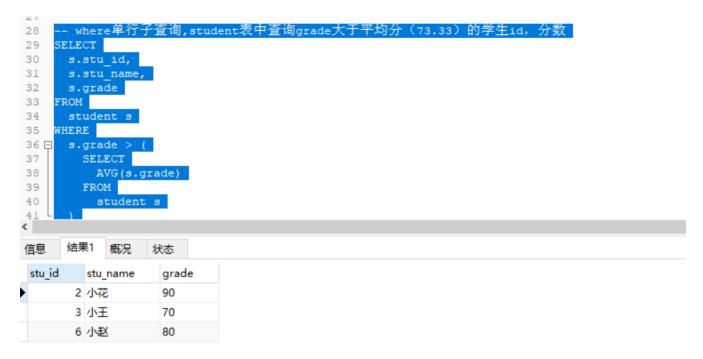
# 1.where单行子查询

student表中查询grade大于平局分的学生id, 姓名, 分数

```
-- where单行子查询
 2
    SELECT
 3
        s.stu_id,
4
        s.stu name,
 5
        s.grade
 6
    FROM
 7
        student s
8
    WHERE
9
        s.grade > (
10
             SELECT
11
                 AVG(s.grade)
             FROM
13
                 student s
```

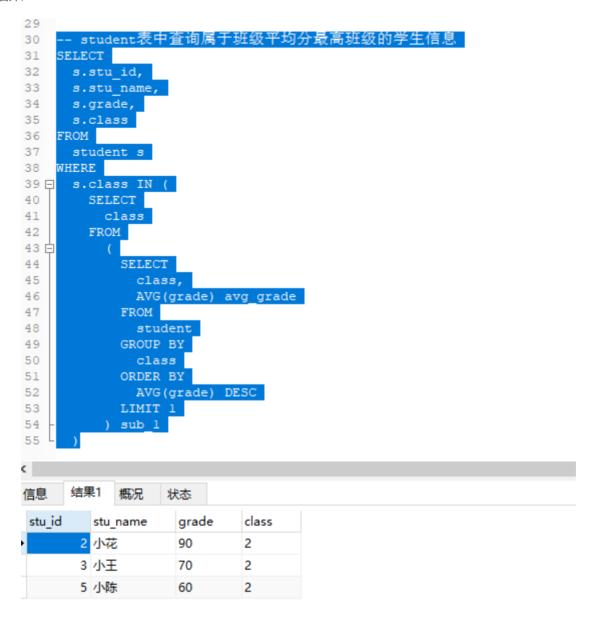
14 )

#### 执行结果



### 2.where/from多行子查询

```
-- student表中查询属于班级平均分最高班级的学生信息
 2
    SELECT
 3
       s.stu_id,
 4
       s.stu_name,
 5
       s.grade,
 6
       s.class
 7
    FROM
 8
       student s
 9
    WHERE
10
       s.class IN (
11
            SELECT
                class
13
            FROM
14
15
                    SELECT
16
                        class,
17
                        AVG(grade) avg grade
18
                    FROM
19
                       student
20
                    GROUP BY
21
                        class
                    ORDER BY
23
                        AVG(grade) DESC
24
                   LIMIT 1
25
                ) sub 1
26
        )
```



上面这个嵌套子查询,一共用了两层子查询:

首先最里层的子查询计算了每个班级的平均分,然后用order by子句按班级平均分倒序排列(最高分排第一个),然后limit 1限制输出一行数据,得到了最高分的班级行,接着把这层子查询的结果放到外面一层子查询的from语句后面。

第二层子查询得到class,最后外查询的where子句判断class是不是等于第二层子查询中的某一个

#### 3.all

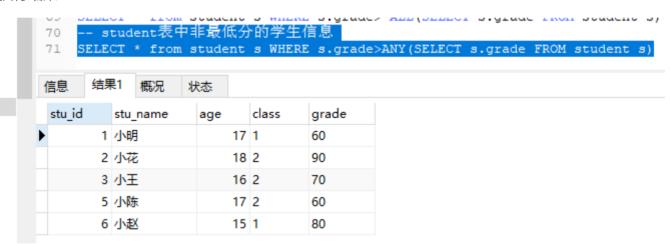
```
1 | -- student表中最高分的学生信息
2 | SELECT * from student s WHERE s.grade>=ALL(SELECT s.grade FROM student s)
```



#### 4.any

```
1 | -- student表中非最低分的学生信息
2 | SELECT * from student s WHERE s.grade>ANY(SELECT s.grade FROM student s)
```

#### 执行结果



注释: 只有最低分的学生成绩不大于任何一个人的分数, 被排除在外

### 四、case when

case when语句时sql中的一个非常重要的功能,可以完成很多复杂的计算,相当于一个表达式。可以放在任何可放表达式的地方。

#### 语法:

case when 条件 then 结果 when 条件 then 结果 else end

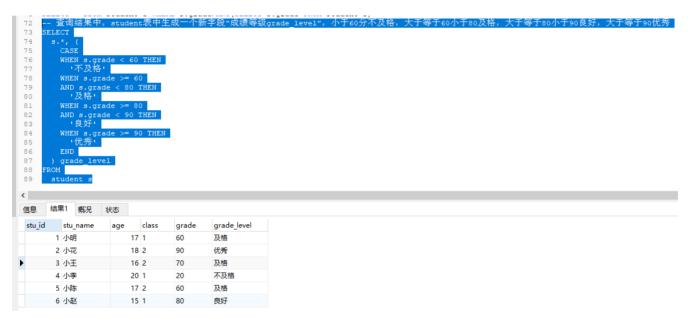
else可不加,是缺省条件下的值,如果不加,有缺省情况则为null。

case when还可以和group by语句搭配使用,用在sum,count,max等聚合函数内部

#### 1.例1

查询结果中, student表中生成一个新字段"成绩等级grade\_level", 小于60分不及格, 大于等于60小于80及格, 大于等于80小于90良好, 大于等于90优秀

```
WHEN s.grade < 60 THEN
 4
 5
             '不及格'
           WHEN s.grade >= 60
 6
 7
           AND s.grade < 80 THEN
               '及格'
 8
9
           WHEN s.grade >= 80
           AND s.grade < 90 THEN
10
               '良好'
11
12
           WHEN s.grade >= 90 THEN
               '优秀'
13
           END
14
15
      ) grade_level
16 FROM
17
      student s
```

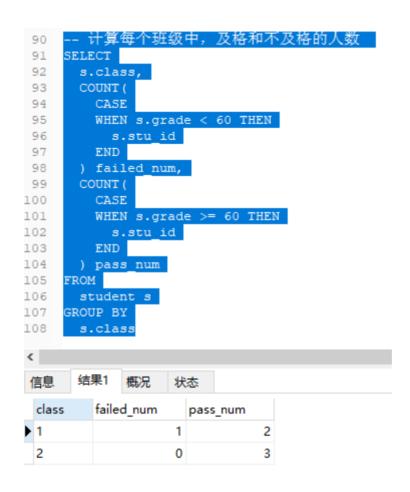


### 2.例2

计算每个班级中, 及格和不及格的人数

```
SELECT
 2
       s.class,
 3
       COUNT (
 4
           CASE
 5
           WHEN s.grade < 60 THEN
 6
            s.stu id
 7
           END
8
       ) failed_num,
9
       COUNT (
10
          CASE
           WHEN s.grade >= 60 THEN
11
            s.stu id
13
          END
```

```
14 ) pass_num
15 FROM
16 student s
17 GROUP BY
18 s.class
```



# 五、组合查询

sql中可用union或union all将多个查询结果拼接起来,两者的区别是union会对结果重新排序,而union all不会,所以一般为了节省内存在不需要排序的情况下用union all更好

#### eg.

```
1 -- 组合查询,查询两个班都及格的学生信息
2 SELECT * FROM student s1 WHERE s1.grade >=60 ORDER BY s1.class ASC
3 
4 SELECT * FROM student s1 WHERE s1.class=1 AND s1.grade>=60
5 UNION ALL
6 SELECT * FROM student s1 WHERE s1.class=2 AND s1.grade>=60
```

```
-- 查询两个班都及格的学生信息
109
110
     SELECT * FROM student sl WHERE sl.grade >=60 ORDER BY sl.class ASC
111
     SELECT * FROM student sl WHERE sl.class=1 AND sl.grade>=60
112
     UNION ALL
113
     SELECT * FROM student sl WHERE sl.class=2 AND sl.grade>=60
114
            概况
                   状态
信息
                           class
 stu id
         stu name
                    age
                                  grade
        1 小明
                        17 1
                                  60
        6 小赵
                        15 1
                                  80
        2 小花
                        18 2
                                  90
        3 小王
                        16 2
                                  70
        5 小陈
                        17 2
                                  60
```

# 六、常用函数

```
-- 字符串函数
1
2
   -- 截取字符串string, 从start开始的length个字符, 类似excel的mid
   substring(string, start, length)
3
   -- 截取字符串string, 从最左边开始的length个字符, 类似excel的left
5
   left(string,length)
   -- 截取字符串string, 从最右边开始的length个字符, 类似excel的right
   right (string, length)
7
   -- mysql中查找stringl在string2中出现的位置
8
   instr(string1,string2[start location])
9
   -- mysql中删除指定个数字符,并在指定位置处插入新字符
10
   insert(string1, start, length, string2)
11
```

### eg.

```
SELECT
 2
       substring(
           '中国浙江省杭州市',
 4
           3,
 5
 6
       ) AS province,
 7
       LEFT (
8
           '中国浙江省杭州市',
9
           2
10
       ) AS country,
       RIGHT (
11
           '中国浙江省杭州市',
13
           3
14
       ) AS city,
15
       instr(
            '中国浙江省杭州市',
16
17
           '杭州市'
18
       ) AS index_city,
```

```
19 INSERT (
20 '中国浙江省杭州市',
21 6,
22 3,
23 '宁波市'
24 ) AS replace_city
```

```
115
116 SELECT
117 = substring(
         中国浙江省杭州市,
 118
 119
 120
       ) AS province,
 121
122 ঢ়
       LEFT (
         中国浙江省杭州市,
 123
 124
         2
       ) AS country,
 125 L
       RIGHT ( 中国浙江省杭州市',
 126 □
 127
       3
) AS city,
 128
 129
      instr(
中国浙江省杭州市,
 130 □
 131
         '杭州市'
 132
 133
       ) AS index city,
       INSERT (
中国浙江省杭州市',
 134 🗐
 135
 136
         3,
'宁波市'
 137
 138
       ) AS replace city
 139
 <
       结果1
             概况
                  状态
                                    replace city
  province
            country
                    city
                          index_city
 ▶ 浙江省
                                   6 中国浙江省宁波市
            中国
                    杭州市
```

# 七、日期函数

```
-- 日期函数,不同数据库会有一些区别,此处以mysql为例
1
   -- 当前日期
2
3
   current date()
   -- date加减expr的unit (年月日周)的date
4
   date add(date, interval expr unit)
   -- expr1减去expr2<mark>的天数</mark>
   datediff(expr1,expr2)
7
   -- expr1减去expr2的时间
8
   timediff(expr1,expr2)
9
   --格式化日期
10
11 date format(date, format)
12
   -- 将字符转换成日期
13 str to date(str, format)
```

#### eg.

```
SELECT
 2
        CURRENT DATE () AS now,
        date add(
 3
 4
          CURRENT DATE (),
            INTERVAL - 1 DAY
 5
 6
       ) AS yesterday,
 7
        datediff(
            CURRENT DATE (),
 8
            DATE ADD(
 9
10
               CURRENT DATE (),
               INTERVAL - 1 DAY
11
12
            )
13
        ) AS date cha,
14
        date_format(CURRENT_DATE(), '%Y/%m') AS ym,
        str to date('2020-07-02', '%Y-%m-%d') AS strdate
15
```

```
141
142
      SELECT
143
        CURRENT DATE () AS now,
        date_add(
144 □
          CURRENT DATE (),
INTERVAL - 1 DAY
145
 146
147 L
        ) AS yesterday,
        datediff(
 148 👨
          CURRENT DATE (),
149
150 🖨
          DATE ADD(
             CURRENT DATE (),
 151
             INTERVAL - 1 DAY
 152
 153
 154 L
       ) AS date cha,
        date format(CURRENT DATE(), '%Y/%m') AS ym,
        str to date('2020-07-02', '%Y-%m-%d') AS strdate
 156
 <
        结果1
               概况
 信息
                     状态
  now
               yesterday
                             date cha
                                        ym
                                                 strdate
 ▶ 2020-07-02
                2020-07-01
                                      1 2020/07
                                                 2020-07-02
```

```
-- 其他一些常用函数
   -- 对某个数据列进行指定小数位四舍五入
   round(columns, decimals)
3
   -- 返回某个字段长短
4
   len(str)
5
   -- 将字符全部小写
6
7
   lower(str)
8
   -- 将字符全部大写
9
   upper(str)
   -- 返回第一行记录
   first()
11
   -- 返回最后一个记录的值
12
   last()
   -- 类型转化, 比如cast(str as bigint)
14
   cast(expr as stype)
15
   -- 将时间戳转换为时间
16
17
   from unixtime(timestamp)
```

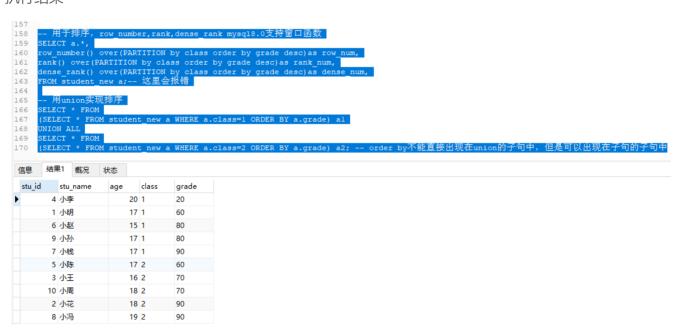
# 八、窗口函数

窗口函数也称为OLAP函数,可以对数据库数据进行实时分析处理。语法如下:

```
1 | <窗口函数> over ([partition by <分组columns>]
2 | order by <排序columns> [窗口子句])
```

窗口函数经常用来计算排序,也可以用在一些聚合函数上,下面举例说明,一张新的学生表 student\_new





### 【举个栗子 12】窗口函数-排序

```
--用于排序, row_number,rank,dense_rank

SELECT a.*,

row_number() over(partition by class order by grade desc) as row_num,

rank() over(partition by class order by grade desc) as rank_num,

dense_rank() over(partition by class order by grade desc) as dense_num,

FROM student_new a
```

stu_id	stu_name	age	class	grade	row_num	rank_number	dense rank
7	小钱	17	1	90	1	1	1
6	小赵	15	1	80	2	2	2
9	小孙	17	1	80	3	2	2
1	小明	17	1	60	4	4	3
4	小李	20	1	20	5	5	4
2	小花	18	2	90	1	1	1
8	小冯	19	2	90	2	1	1
	小王	16	2	70	3	3	2
10	小周	18	2	70	4	知乎 @素	呆玩数据 2
5	小陈	17	2	60	5	5	3

由上面这个例子可以清晰地看出三个窗口排序函数的区别,row\_number是在每个分组窗口中给定唯一序号,而rank碰到相等值序号一样,会跳过之后的位次,而dense\_rank碰到相等值序号也一致,但不跳过之后的位次。