

# Lecture 4 & Lab 4

---

## Distinct

### 1. **No duplicated** identifier.

adding `distinct` will help us to deduplicate the columns we select, i.e. eliminate the identical rows for the columns

**NULL** is treated as a distinct element

used when we are interested in different keys

### 2. syntax

```
1 | select distinct {col} from {tab}
```

then we will see the counts of some columns

```
1 | select distinct {col1}, {col2}, ... from {tab}
```

the tuple `(col1, col2, ...)` is distinct

### 3. Remark

We cannot use `distinct` to select both distinct and not distinct columns. Once `distinct` is used, the combination of the following columns should be distinct and there is no exception.

We can use `group by` to deal with this problem.

## Aggregate Functions

### 1. Aggregate functions will **aggregate all rows** that share a feature and **return a characteristic** of each group of aggregated rows.

### 2. **count**

#### 0. reference

[PostgreSQL COUNT Function: Counting Rows That Satisfy a Condition \(postgresqltutorial.com\)](https://www.postgresqltutorial.com/postgresql-count/)

#### 1. syntax

Summary: It **won't** consider the **NULL** if and only if you do the **single-column counting**

```

1  -- to count the number of rows including NULLs
2  select count(*) from {tab1}
3
4  -- to count the number of rows of col1 ignoring NULL
5  select count({col1}) from {tab1}
6
7  -- to count the number of different rows of col1 ignoring NULL
8  select count(distinct {col1}) from {tab1}
9
10 -- to count the number of different rows of (col1, col2, ...) including
    NULLs
11 select count(distinct({col1}, {col2}, ...)) from {tab1};

```

## 2. example

num1	num2	num3
1	2	null
1	null	null
null	2	1
null	null	null
null	null	1

```

1  select count(*) from table1; -- 5
2  select count(num3) from table1; -- 2
3  select count(distinct num3) from table1; -- 1
4  select count(distinct (num2, num3)) from table1; -- 4

```

## 3. group by

1. Aggregate all rows to several rows with the specified criteria.

Then some other aggregate functions can aggregate each group and return a specific value

2. syntax

**NULL is also a group in postgresQL**

```

1  select {col1} -- columns must also appear after group by
2  from {tab}
3  where {cond}
4  group by {col1}
5
6  select {col1}, {col2} -- columns must also appear after group by
7  from {tab}
8  where {cond}
9  group by {col1}, {col2}
10
11 -- or we can use aggregate functions

```

```

12 -- as can be omitted without changing the meaning
13 select {col1} aggregate_func({col2}) as {new_name}
14 from {tab}
15 where {cond}
16 group by {col1} -- not compulsory to appear col2 here

```

### 3. example

num1	num2
null	1
null	1
1	2
1	1
null	null

```

1 select num1, count(num2) from table1 group by num1;
2 /*
3 null 2
4 1    2
5 */
6 select num1, count(*) from table1 group by num1;
7 /*
8 null 3
9 1    2
10 */
11 select num1, num2 from table1 group by num1, num2;
12 /*
13 null 1
14 1    2
15 1    1
16 null null
17 */

```

### 4. other aggregate functions

they all ignore **NULL**

```

1 min(col)
2 max(col)
3 stddev(col) -- find the standard deviation of a column
4 avg(col)

```

### 5. composite relation

```

1  -- for the more complicated composite relation, use alias.col1 to specify
   the exact col
2  select {alias1.col1} from
3  {
4      select ...
5      from ...
6      where...
7  } {alias1} -- give the intermediate results an alias
8  where {cond}

```

examples:

num1	num2
null	1
null	1
1	2
1	1
null	null

```

1  select newta.std from(
2      select num1, stddev(num2) std
3      from testable1
4      group by num1
5  ) newta where newta.std > 0.58;

```

## 6. having

1. From actual meaning, rows firstly filtered by `where`, then grouped through `group by`, finally filtered by `having`
2. syntax

```

1 | select {} from {} where {cond1} group by {cols} having {cond1_}

```

## 3. examples

num1	num2
null	1
null	1
1	2
1	1

num1	num2
null	null

```

1 select num1, stddev(num2) std
2   from testable1
3  where num2 is not null
4  group by num1
5  having stddev(num) > 0.58
6  -- having std > 0.58 is wrong
7  -- having num2 > 1 is wrong
8  -- the expression must be aggregated or specify in the GROUP;

```

## Join

1. aim: we want to retrieve data from **multiple** table and display them in **one** table

2. mechanism:

make a **cartesian product** for the keys on table1 and table2,

then filter those `ta1_keyCol != ta2_keyCol`

rows with **key\_Col is null** is ignored

3. syntax:

```

1 select {ta1_col1}, {ta1_col2}...
2    {ta2_col1}, {ta2_col2}...
3   from {ta1} (alias1) join {ta2} (alias2)
4  on {alias1}.{ta1_keyCol} = {alias2}.{ta2_keyCol}
5  where {conditions}
6     ...;
7  -- though alias is not compulsory, it's a good habit
8  -- it's vital when ta1_keyCol have the same name as ta2_keyCol

```

4. natural join:(not recommended)

what if we don't specify the column?

```

1 select * from people natural join credits;
2
3  -- will automatically join by recognized key(same column name)
4  -- but sometimes we don't guarantee the same name means same semantic
5
6  -- The same as
7  select *
8  from people join credit
9  on people.peopleid = credits.peopleid
10 -- note that peopleid is a column in both people and credits

```

what if we want to simplify the query

```

1  -- better
2  select *
3  from people join credits using(peopleid)

```

## 5. self join:

join the same table together:

For example: How can we find how many **unordered pairs** of people with the same first name?  
(An inner pairing problem)

```

1  -- recall the mechanism: cartesian product
2
3  -- wrong solution 1:
4  select count(*)
5  from people p1 join people p2
6  on p1.first_name = p2.first_name
7  -- someone matching himself is counted
8
9  -- wrong solution 2:
10 select count(*)
11 from people p1 join people p2
12 on p1.first_name = p2.first_name
13 where p1.id != p2.id
14 -- A matches C and then C will matches A, that is, for one pair (A, C) will
   be counted twice
15
16 -- correct answer
17 select count(*) / 2
18 from people p1 join people p2
19 on p1.first_name = p2.first_name
20 where p1.id != p2.id

```

## 6. join in a subquery

syntax:

```

1  select {col1}, {col2}, {col3}...
2  from
3      {tab1} ({alias1})
4      join {tab2} ({alias2}) on {alias1}.{ta1_keyCol} = {alias2}.{ta2_keyCol}
5      join {tab3} ({alias3}) on {alias1}.{ta1_keyCol} = {alias3}.{ta3_keyCol}
6      ...
7  where {conditions}...

```

mechanism is similar, just filter more than one time.

# Quickly search some characteristics in a table

1. search the column names of a table

```
1 | SELECT column_name FROM information_schema.columns WHERE table_name =  
   | {tab1};
```

2. search the possible values for a column

```
1 | select distinct {col} from {tab1}
```

3. advanced: check the column names, types, not null, reference, primary key, foreign key.

**Command line** or **powershell** is needed in windows

Ensure that you add the **\bin** to the path

```
1 | # [optional]  
2 | chcp 437 # switch to English  
3 | chcp 936 # switch to Chinese  
4 |  
5 | # firstly log in to connect the database  
6 | psql -U {user_name} -d {database_name}  
7 |  
8 | # input your password  
9 | {your_password}  
10 |  
11 | # check characteristics of a table  
12 | \d {table_name}  
13 |  
14 | # search  
15 | select  
16 | distinct {col} from {tab}  
17 |  
18 | # terminate query  
19 | ctrl+c
```