
Session 5

Mohamed Emary

December 15, 2024

All the examples are on the 'ITI' Database we have worked on before.

1 GROUP BY and HAVING

1.1 GROUP BY

In the Instructor table, to find the minimum salary between all the instructors, we can use:

```
1 SELECT
2     MIN(Salary)
3 FROM
4     Instructor;
```

But what if we want to find the minimum salary for each department? We can use the **GROUP BY** clause to group the rows based on the **Dept_name** and then apply the **MIN** function to each group.

```
1 SELECT
2     MIN(Salary) "Min Salary Per Dept",
3     Dept_Id
4 FROM
5     Instructor
6 GROUP BY
7     Dept_Id;
```

You may see a **NULL** value in the min salary column. This appears if there is one or more instructors in the department all having **NULL** salary. If there is at least one instructor with a salary, the **NULL** value will not appear (**NULL** is not a value and it only appears when there is no other value).

To avoid having **NULL** values in the result, we can use the **WHERE** clause to filter out the **NULL** values.

```
1 SELECT
2     MIN(Salary) "Min Salary Per Dept",
3     Dept_Id
```

```
4 FROM
5   Instructor
6 WHERE
7   Salary IS NOT NULL
8 GROUP BY
9   Dept_Id;
```

The GROUP BY clause make the aggregation function (MIN in this case) apply to each group instead of the whole table.

Another way to apply the operation above is to use PARTITION BY in the OVER clause. This way, we can apply the aggregation function to each group without using the GROUP BY clause.

```
1 SELECT
2   Dept_Id,
3   MIN(Salary) OVER (
4     PARTITION BY
5       Dept_Id
6   ) "Min Salary Per Dept"
7 FROM
8   Instructor
9 WHERE
10  Salary IS NOT NULL;
```

When using group by use it on a column that it's values are the same in multiple rows. Using it on a unique column like the primary key will be useless:

```
1 SELECT
2   St_Id,
3   COUNT(*)
4 FROM
5   Student
6 GROUP BY
7   St_Id;
```

You also can't group by *. For this to work you need to have at least two rows with the same values in all columns which shouldn't happen from the beginning since the primary key is unique. It's also not possible in code, you will get an error:

```
1 SELECT
2   St_Id,
3   COUNT(*)
4 FROM
5   Student
6 GROUP BY
7   * -- Error
```

To count the number of students in each department, the COUNT aggregate function would work on each group from the GROUP BY clause.

```
1 SELECT
2   Dept_Id,
3   COUNT(*) AS 'Number of Dep Students'
4 FROM
5   Student
```

```
6  GROUP BY
7    Dept_Id;
8
9  -- To ignore the NULL values in the Dept_Id column
10 SELECT
11     Dept_Id,
12     COUNT(*) AS 'Number of Dep Students'
13 FROM
14     Student
15 WHERE
16     Dept_Id IS NOT NULL
```

To count the number of students in the whole table:

```
1  SELECT
2    -- Here the COUNT function will work on the whole table
3    -- since there is no GROUP BY clause
4    COUNT(*) AS 'Total Number of Students'
5 FROM
6     Student;
```

Anything being selected next to the aggregate function and it's not an aggregate function should be in the GROUP BY:

```
1  SELECT
2    St_Lname,
3    COUNT(*)
4 FROM
5     Student
6 GROUP BY
7     St_Lname;
8
9  -- Another example
10 SELECT
11     St_Fname,
12     MAX(St_Age)
13 FROM
14     Student
15 GROUP BY
16     St_Fname;
```

Grouping by multiple columns have a similar idea to cross join:

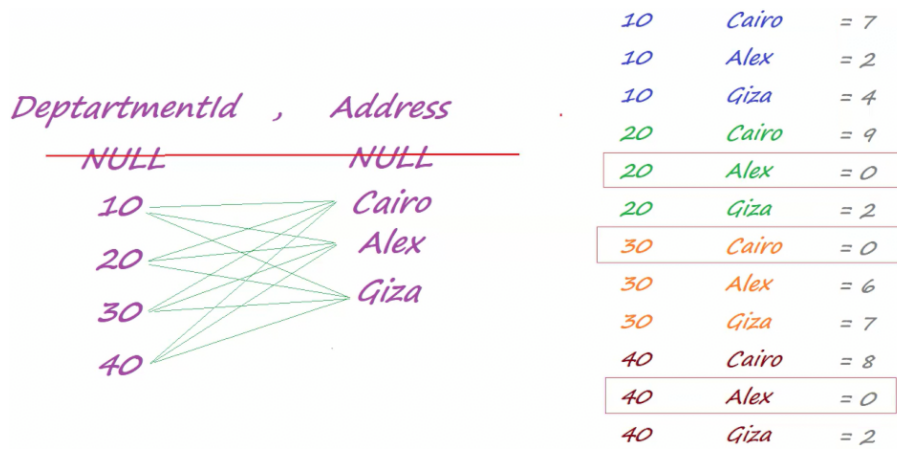


Figure 1: Grouping by multiple columns

```

1 SELECT
2   Dept_Id,
3   St_Address,
4   COUNT(*) AS 'Number of Students'
5 FROM
6   Student
7 WHERE
8   Dept_Id IS NOT NULL
9   AND St_Address IS NOT NULL
10  -- Grouping here is done by address first (the second column)
11  -- then by dept_id (the first column)
12 GROUP BY
13   Dept_Id,
14   St_Address;

```

1.2 HAVING

SELECT and WHERE work on table *record by record* while COUNT aggregate function works on the whole table, that is why you can't use AND COUNT(*) > 2 in the WHERE clause in the statement below

If you want to get the number of students in each department that have more than 2 students:

```

1 -- This will not work
2 SELECT
3   Dept_Id,
4   COUNT(*) AS 'Number of Students'
5 FROM
6   Student
7 WHERE
8   Dept_Id IS NOT NULL
9   AND COUNT(*) > 2
10 GROUP BY
11   Dept_Id

```

To fix that you will need to use HAVING keyword, HAVING works on the groups created by the GROUP BY clause HAVING is mostly used with aggregate functions.

```

1  SELECT
2      Dept_Id,
3      COUNT(*) AS 'Number of Dep Students'
4  FROM
5      Student
6  WHERE
7      Dept_Id IS NOT NULL
8  GROUP BY
9      Dept_Id
10 HAVING
11     COUNT(*) > 2;

```

In general, we use **HAVING** without **GROUP BY** if we want to apply a condition on an aggregate function and we are not selecting that aggregate function in the **SELECT** clause **because we can't use aggregate functions in the WHERE clause**.

With **HAVING** we always have a condition with aggregation that condition works on the groups created by the **GROUP BY** clause or on the whole table if there is no **GROUP BY** clause and we are using an aggregate function in the **SELECT** clause.

Some general rules:

1. Aggregate functions like **COUNT** work on the whole table.
2. **WHERE** works on the table record by record.
3. **HAVING** works on the groups created by the **GROUP BY** clause.

Example of using **HAVING** without **GROUP BY**, we want to get the sum of all instructors salaries if there is more than 10 instructors in the table:

```

1  SELECT
2      SUM(Salary) AS 'Total Salaries'
3  FROM
4      Instructor
5  HAVING
6      COUNT(*) > 10;

```

If we want to get the sum of salaries for each department, we can use one of the two statements below.

The two statement below is similar to each other but the second one uses join. The performance in the second one is worse than the first since we are applying operations on two tables

Generally this is not a good case for using join. Join is used when we want to get data from two tables

```

1  SELECT
2      Dept_Id,
3      SUM(Salary) AS 'SumOfSalaries'
4  FROM
5      Instructor
6  WHERE
7      Dept_Id IS NOT NULL
8  GROUP BY
9      Dept_Id;
10

```

```

11  -- Using join
12  SELECT
13      I.Dept_Id,
14      SUM(I.Salary) AS 'SumOfSalaries'
15  FROM
16      Instructor I,
17      Department D
18  WHERE
19      I.Dept_Id = D.Dept_Id
20  GROUP BY
21      I.Dept_Id;

```

Here each department have a different name, and to show the department name here we still need to group by Dept_Name to make the query work:

```

1  SELECT
2      I.Dept_Id,
3      D.Dept_Name,
4      SUM(I.Salary) AS 'SumOfSalaries'
5  FROM
6      Instructor I,
7      Department D
8  WHERE
9      I.Dept_Id = D.Dept_Id
10 GROUP BY
11     I.Dept_Id,
12     D.Dept_Name;

```

To select the students who act as supervisors and the number of students they supervise, we had to group by both Supr.St_Fname, Supr.St_Id since we are selecting both of them in the SELECT clause. If we are selecting only one of them we can group by only that column.

NOTE: If you group by only the St_Fname column, and we have two supervisors with the same St_Fname, the query will group them together and show the total number of students they supervise together.

```

1  SELECT
2      Supr.St_Fname 'Supervisor',
3      Supr.St_Id 'Supervisor ID',
4      COUNT(*) 'No of Students'
5  FROM
6      Student Stud,
7      Student Supr
8  WHERE
9      Supr.St_Id = Stud.St_super
10 GROUP BY
11     Supr.St_Fname,
12     Supr.St_Id;

```