

Class Starts at 9:05 PM

## Agenda :

- Aggregates
  - Group By
  - Having
    - count()
    - min()
    - max()
    - avg()
    - sum()
- 

→ Till now whatever SQL Queries we have written, were to get few rows of the table, (Filtering out rows using ON and WHERE)

Sometimes, instead of getting few rows of the table, we have to get analysis of the entries of table.

Q1) Get the average PSP of every batch of scalar.

Batches	
id	name

Students			
id	name	batch-id	PSP

Approaching this as a DFA problem:

Step 1: Iterate over the Join of Batches and Students.

Step 2: Use maps to store data of Students batch-wise.

Step 3: Get the sum of PSP / sum of Students

To answer such type of Queries, we need the Aggregate Function.

Count(), SUM(), MIN(), MAX(),  
Avg().

Behaviour of every Agg. Fn



Count

Students			
id	name	age	batch_id
1	A	20	1
2	B	21	1
3	C	22	Null
4	D	23	2

Q) Give the count of students that have a batch?

Count → takes a lot of values and combines them into a single value equal to the number of elements in the list.

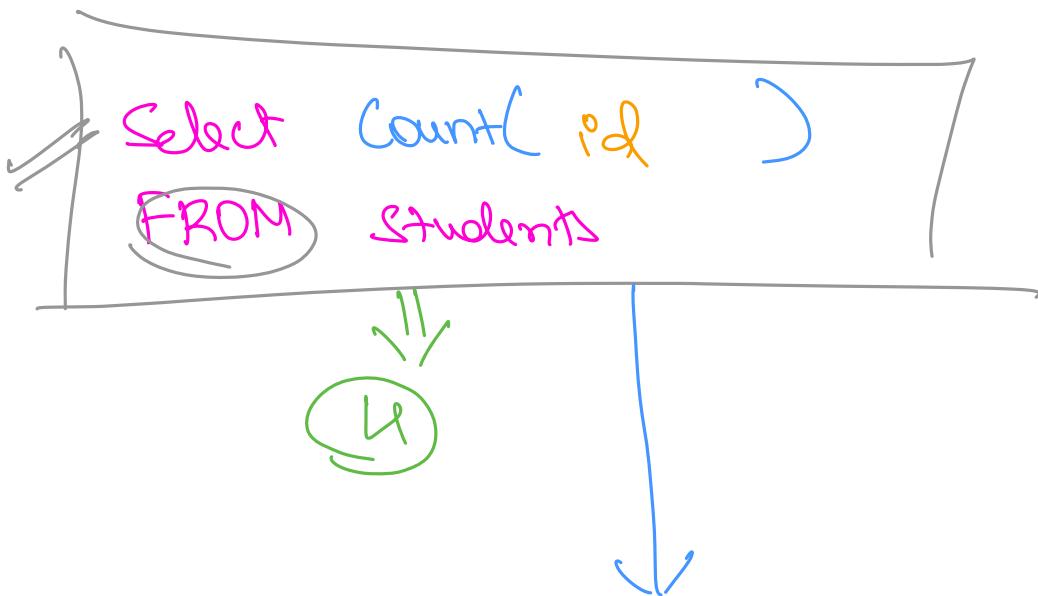
Eg  $\text{Count}(1, 2, 3, 5, 6) \rightarrow 5$

Select Count(batch\_id)  
FROM Students

Count (1, 1, Null, 2)  $\Rightarrow$  3

All aggregate F's ignore Null when computing their result.

what if I use id in Count.



Select Count(id)  
FROM Students  
WHERE batch\_id IS NOT NULL

Students			
id	name	age	batch_id
1	A	20	1
2	B	21	1
3	C	22	NULL
4	D	23	2

Q) Give me those students that have a batch and their age is less than 23.

A) Select Count (batch-id)

FROM	Students
WHERE	age < 23

Students		
id	name	batch-id
1	x	1
2	y	1
3	z	2

Batch	
id	name
1	A
2	B
3	C

Q) Let count of students with batch-name  
'A'

A) Select Count (s.id)

FROM Students s

JOIN batch b

ON s.batch-id = b.id

AND b.name = 'A'

```
Select      Count (s.id)
FROM       Students  s
JOIN       batches   b
ON         s.batchid = b.id
WHERE      b.name = 'A'
```

## CODE

```
Select
FROM A
JOIN B
ON condn 1
WHERE condn 2
```

$A = \{ \{ \} \dots \{ \} \dots \}$   
 $B = \{ \{ \} \dots \{ \} \dots \}$   
 $ans1 = [ ]$

for row1 in A :

    for row2 in B :

        if row1 and row2 match ON cond " :  
 $ans1.add(row1 + row2 \text{ (stitch)})$

JOIN



$ans2 = [ ]$

for row3 in  $ans1$  :

    if row3 matches the WHERE cond " :  
 $ans2.add[row3]$  :

virtual table.

Count-S.id = 0

for row in  $ans2$  :

    if row[S.id] is Not Null :

        Count-S.id += 1;

filtered table

Aggregate  
W  
F

print (Count-S.id) ;

Other aggregate f<sup>n</sup>

→ you can print multiple aggregate f<sup>n</sup> in the same query.

Select count(s.id), sum(s.pop), max(s.pop)  
min(s.age)  
FROM students

Max ( - - - - - )  
Min ( - - - - - )

Both these functions can take a col<sup>n</sup> which is comparable.

{  
    → int  
    → double  
    → Strings (lexicographic)

Average

~~Avg( 1, 2, 3, Null )~~ →

$$\rightarrow \frac{6}{3} = 2 \quad \checkmark$$

$$\rightarrow \frac{6}{4} = 1.5 \quad \times$$

Student	
id	psp
1	1
2	2
3	3
4	Null

Select (  $\frac{\text{Sum}(\text{psp})}{\text{Count}(\text{psp})}$  ) =  $\overline{\text{Avg}(\text{psp})}$

True?  $\checkmark$

False?

Print

Select

$$\frac{\underline{6}}{\underline{4}} = \overline{\text{Avg}(\text{psp})}$$

False ✓  
 $\frac{6}{7} \neq 2$

Q) Tell me how many students are there?

Students		
id	name	batch-id
1	x	1
2	y	1
3	z	2

Select count(id)  
FROM Students;

→ what if id could also be Nulls.

A) Select count(\*)  
FROM Students;

Count(\*)

$A = \{ \{ \} \dots \} \}$   
 $B = \{ \{ \} \dots \} \}$   
 $ans1 = [ ]$

for row1 in A:

    for row2 in B:

        if row1 and row2 match ON word<sup>n</sup>:  
            ans1.add(row1 + row2 (stitch))

JOIN

ans2 = [ ]

for row3 in ans1:

    if row3 matches the WHERE cond<sup>n</sup>:

        ans2.add[row3];

virtual table.

Count-Sid = 0

for row in 1  
    if

filtered table

ans2 is

Not Null's

Count-Sid += 1;

Aggregate  
W  
F

print (Count-Sid);

\* means any value

writing Query:

```
Select Count(*)  
FROM Students
```

exactly same as.

```
Select Count(1)  
FROM Students
```

⇒ Output?

whatever we give inside the Count(x) fn  
It will check for x is not Null.

Count ( sid ) → checks if sid is not Null  
Count ( batch\_id ) → checks if batch\_id is not Null.

Count (\*)  
Count (1)  
Count ('Hello') } all are same as \*

Checks if 1 is Not Null for every iteration.  
→ Always true.

Select Count(1) }  
Select Count(2) } same as  
Select Count('Hello') } Select Count(\*)

Which one is faster?

Select Count(1) OR Select Count(\*)

Select Count(1) is faster  $\Rightarrow$  No memory access.

Select Count(\*) I have to check for non actual values for all columns

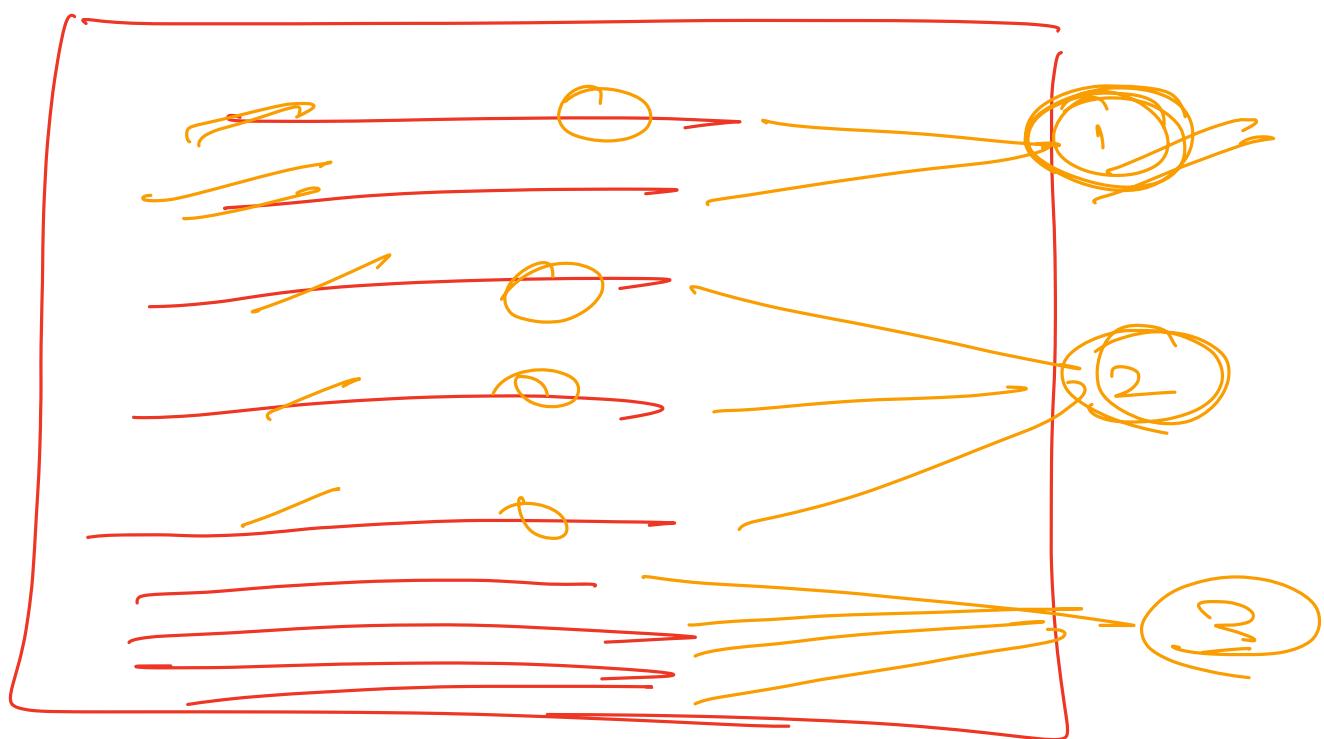
Break till 10:28 PM

Students			
id	name	psp	batch_id

Q) Get the count of students for every batch.

batch\_id      count of Students.

1	50
2	30
3	140
4	60



All Students of  
batch 1 → Count(\*) → answer  
for b1.

All Students of  
batch 2 → Count(\*) → answer for  
batch 2.

I can make use of aggregates, Not on the entire table, but on the sub-parts (groups) of the table.

Group By: Allows you to break your table into smaller groups, so as to be used by the aggregate functions

Eg → i) Group by batch\_id.  
// It will bring all the rows of same batch together.

Select batchid , Count( id )  
 FROM Students  
 Group By batchid

1	50
2	30
3	40
4	60

- Q) 2) Group by on 2 diff columns.  
 Q) How many different groups will be there?  
 Group by b-id, ins-id.

b-id	ins-id	st-id
1	100	20
1	100	21
1	150	30
1	150	31
1	150	32
1	100	32

1	100	24
2	200	40
2	200	45
1	100	50

Groups

1, 100 → 5 Students

1, 100 → 3

2, 200 → 2

Students

id	name	batch_id	psp
1	A	1	10
2	B	1	20
3	C	2	30
4	D	1	40

Q) For each batch, get me the no. of Students.

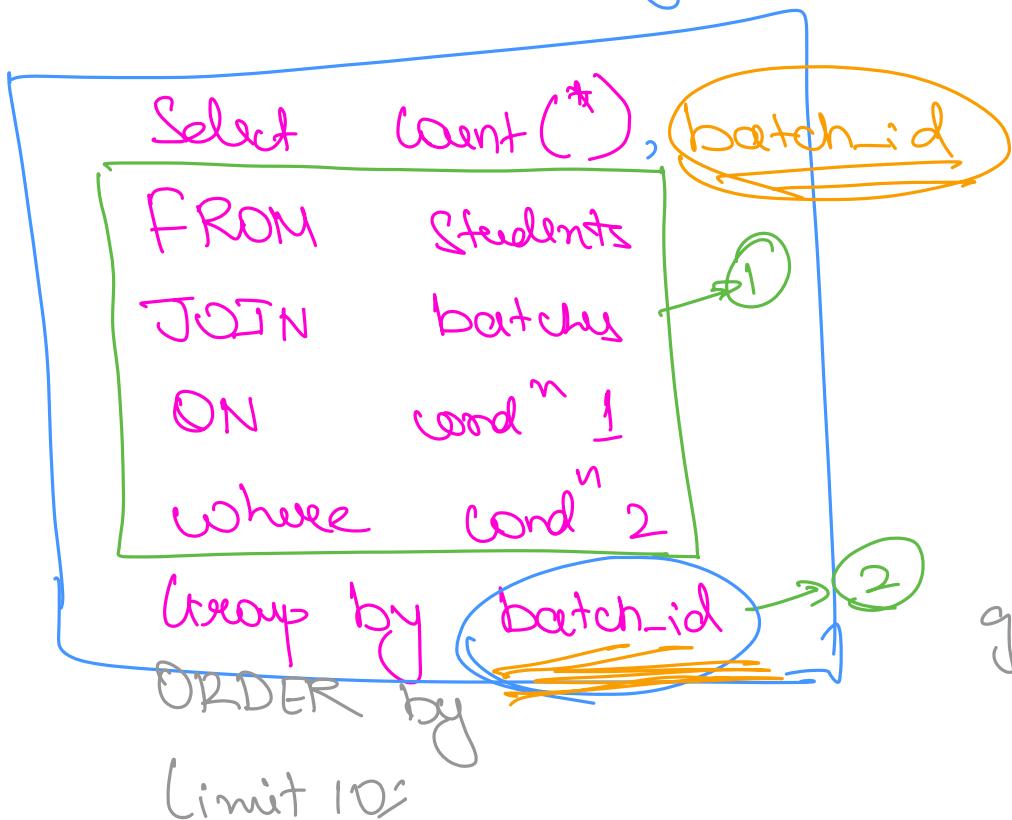
Select count(id)

from Students

Group By batch\_id



General SQL Query.



Select count, b-name

group by b-name

Output :

(Count*)	batch_id	
3	1	Nov22
1	2	Nov22

→ You can only use those columns that are present in Group By.

↳ Select will have those columns, that are mentioned in Group By

Q) Print the batch names with more than 100 students.

Students			
id	name	batch-id	psp

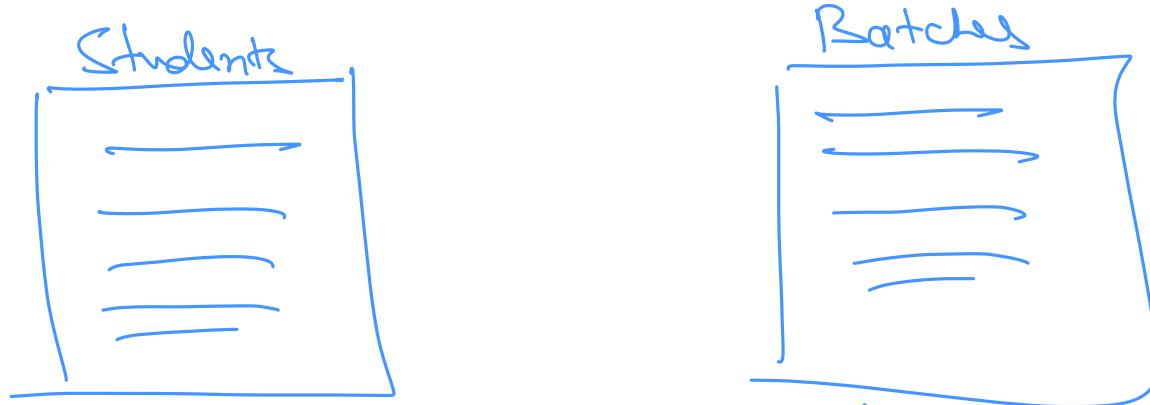
Batches	
id	name

1) Step 1 : ①) Print batch-name along with  
      ②) the count of students.

Select b\_name, count(1)

FROM Students S  
JOIN batchel b  
ON S.batch\_id = b.id

GROUP BY b\_name



b_name	count(1)
A	110
B	90
C	120

I need to do some filtering on this table.

Do I need to filter on Rows? No.

I want to filter on the groups.

## HAVING

→ Allows you to filter on the groups.

Q) Print the batch names with more than 100 students.

```
Select b.name, count()
FROM Students S
JOIN batchel b
ON S.batch_id = b.id
GROUP BY b.name
HAVING count() > 100
```

b.name	count()
A	110
B	90
C	120

## Code

A = [ { } , ... , { } , ... ]

B = [ { } , ... , { } , ... ]

ans1 = [ ]

for row1 in A :

    for row2 in B :

        JOIN

        if row1 and row2 match ON cond<sup>n</sup> :  
            ans1.add( row1 + row2 (stitch) )

ans2 = [ ]

    virtual table.

    for row3 in ans1 :

        if row3 matches the WHERE cond<sup>n</sup> :

ans2.add[ row3 ] :

Map <( Group By (w<sup>m</sup>) , int > countMap

Map <( Group By (w<sup>m</sup>) , Double > avg Map

for row in ans2 :

```
// populate the maps  
CountMap.put(bname, emitting-value + 1)  
avgMap.update()
```

GROUP BY

generating the aggregates

for each group in CountMap:

if condition in Having clause is true:

print(CountMap.get(group), avgMap.get(group))

Count Map

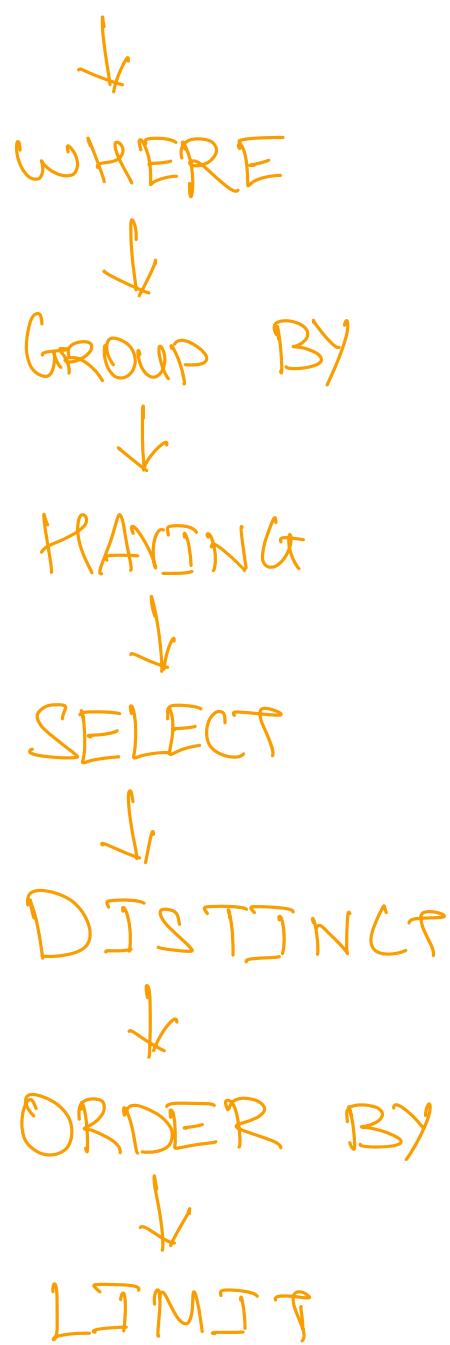
A	110
B	90
C	120

ORDER of Operations:

FROM



JOIN (ON)



Q) Can I apply WHERE after Group By?

No

HAVING will be applied.