# SQL for Data Scientists/MLE/SDE

Applied AI Course.com

InterviewPrep.AppliedCourse.com

Why ?

Data scientist $\longrightarrow$ mostly one round

MLE $\longrightarrow$ Programming + SQL

SDE $\longrightarrow$ Programming + SQL

$$\text{Data} \xrightarrow{\text{SQL}} \text{Code}$$

Services / startups / product - based

Level 1 & 2

Level 2 & 3 } examples ahead!

rare

Applied AI Course $\longrightarrow$ SQL on IMDB dataset

[ 10 queries ]

$\downarrow$

Level 2 &3

nested -queries

①

Schema: ⟨Id, Marks⟩ ←— Students

Find the second highest marks

↳ very common in many companies

Sub-problems   $\longrightarrow$   sort   descending   order

$\longrightarrow$   pick the   2nd from top

edge-cases   $\longrightarrow$

1, 89

2, 89     [duplicates]

3, 86

Select distinct Marks from Students

order by Marks desc

LIMIT 1    OFFSET 1

**Alt:**

(a) Select MAX(Marks) from Students

(b) Select max(Marks) AS marks from

Students

WHERE marks < ( select max(marks)
from students)

② Points : $\langle x, y \rangle$

find the shortest distance between the points

[Amazon, Microsoft]

Python/Java/C++:

for each $P_i$ in points
   for each $P_j$ in points
      if $P_i \neq P_j$
         $d(P_i, P_j)$

SELECT MIN (SQRT ( POW( P1.x - P2.x , 2) + POW ( P1.y - P2.y, 2)
))
as minDist

FROM Points P1 JOIN Points P2
ON P1.x != P2.x OR P1.y != P2.y

③

orders: cust, item, date

Select customers who purchased atleast two items and on two different dates

→ e-commerce companies

group-by COUNT $\longrightarrow$ item Cnt, date Cnt
cust $\hookrightarrow$ distinct

where $\longrightarrow$ filtering $(\geqslant 2)$

Select Cust, COUNT ( DISTINCT item) as itemCnt

COUNT ( DISTINCT date) as dateCnt

from Orders

GROUP BY Cust

HAVING itemCnt >= 2 AND dateCnt >= 2

④

scores: player, country, goals

Find players who scored more than all Spanish players and more than atleast one german player

Find players who scored more than all Spanish players and more than atleast one german player

goals > ALL ( Spanish player goals)

goals > ANY (german player goals)

Select t.player from scores AS t

where t.goals >ALL ( SELECT t1.goals
from scores AS t1
where t1.country = 'spain')

AND

t.goals >ANY ( SELECT t2.goals from
scores as t2 where
t2.country = 'germany' )

(5)

| water_schemes | | |
|---|---|---|
| scheme_no | district_name | capacity |
| 1 | Ajmer | 20 |
| 1 | Bikaner | 10 |
| 2 | Bikaner | 10 |
| 3 | Bikaner | 20 |
| 1 | Churu | 10 |
| 2 | Churu | 20 |
| 1 | Dungargarh | 10 |

Print names of districts whose total capacity $\geq$ average capacity of all districts

total :  distrct_name, capacity $\xleftarrow[\text{sum}]{\text{groupby}}$  water-schemes

total_avg :  average-Capacity $\xleftarrow[]{\text{avg}}$  total

result :  distrct-name $\xleftarrow[X]{\text{where}}$  total, total-avg

*name of the sub-query block*

```sql
with total(name, capacity) as

    select district_name, sum(capacity)
    from water_schemes
    group by district_name
```

} → *sub-query block*

```sql
with total_avg(capacity) as
    select avg(capacity)
    from total


select name
    from total, total_avg
    where total.capacity ≥ total_avg.capacity
```

[ sub-query - refactoring ]

Key focus

→ Basic syntax

→ nested Queries

→ Self - joins   [ Select ... from Emp E₁,
                               Emp E₂ ... ]

→ Break a problem into Sub-problems

For more practice:

① GATE CS    https://www.geeksforgeeks.org/dbms-gq/sql-gq/

② Online resources: " SQL interview Questions Amazon"

↳ SOLVE actual Query-related Questions

↳ Theory is easy!