

Agenda:

1. Roles
2. Skills for each role
3. Interview rounds
4. Sample Interview Questions
5. How to prepare & ace them
6. Projects
7. Non-tech & Non-CS backgrounds

Freshers (approx..)

12-ISLP*1*

Data Analyst →

BA

PA;

Op Analyst

RA; FA

SCA

Roles

20-28L*2*

Data scientists

20-28L*3*

Machine learning Engineer

MLOps

SDE(ML); CV-Engineer

DevOps

25-35L*4*

Applied Scientist

Data Engineer

↳ SDE + Data pipelines

5

Research Scientist

Skills:

1. Python & libraries ; basic DS
2. SQL
3. Data Analysis
 - ↳ Probability
 - Applied stats
 - Plotting

→ Pandas; NumPy; scipy; Matplotlib; Seaborn
inbuilt

Data
Analyst
(must)

4: Dashboarding/Viz Tools
(not a must)

- DA (good to have)

5: Math: Optimization;
Linear Algebra,
High-dim geom

6: Classical ML
↳ Supervised: GB DT; RF; ...
↳ Unsup: Clustering
→ RecSys; Time series

DS

7: Deep learning

↳ NN

→ CNN (CV)

→ NLP (RNN, Transformer)

DS

8: ML ops

→ API; Docker / Kubernetes

→ Spark

→ (how to use)

→ AWS / Azure / GCP (not must)

→ ML System design

→ basics of NoSQL

MLE

↳ part SDE → DS & algo (some companies)
↳ part DS ✓

Applied Scientist

↳ DS + Cutting edge research

PhD ↳ { Research Scientist → GB/DeepMind/FAIR/OpenAI/MSR
⇒ ↳ designs new algo from scratch

Interview
guides



lectures
(Medi, some
hard)

Data Analyst

1. SQL → Subqueries; joins
2. Python → Libraries; loops; ^{nested} recursion
3. Data analysis → scenarios / data based

Data scientist

1. Programming → Python +SQL
2. Maths →
 - Pools, stats, optimization, LA;
 - Data Analysis; high dim geom

3. breadth → ML & DL
4. depth → internals of techniques
5. scenario based → MLOPS / CoeML / DL ...
=

MLE:

1. Programming → DSA }
2. SQL / Programming
3. Breadth } scenario-based
=
4. depth



Interview Questions

ML breath:



scenarios

Train Data: \mathcal{D}_{TY}

Test Data: \mathcal{D}_{Te}

(Q)

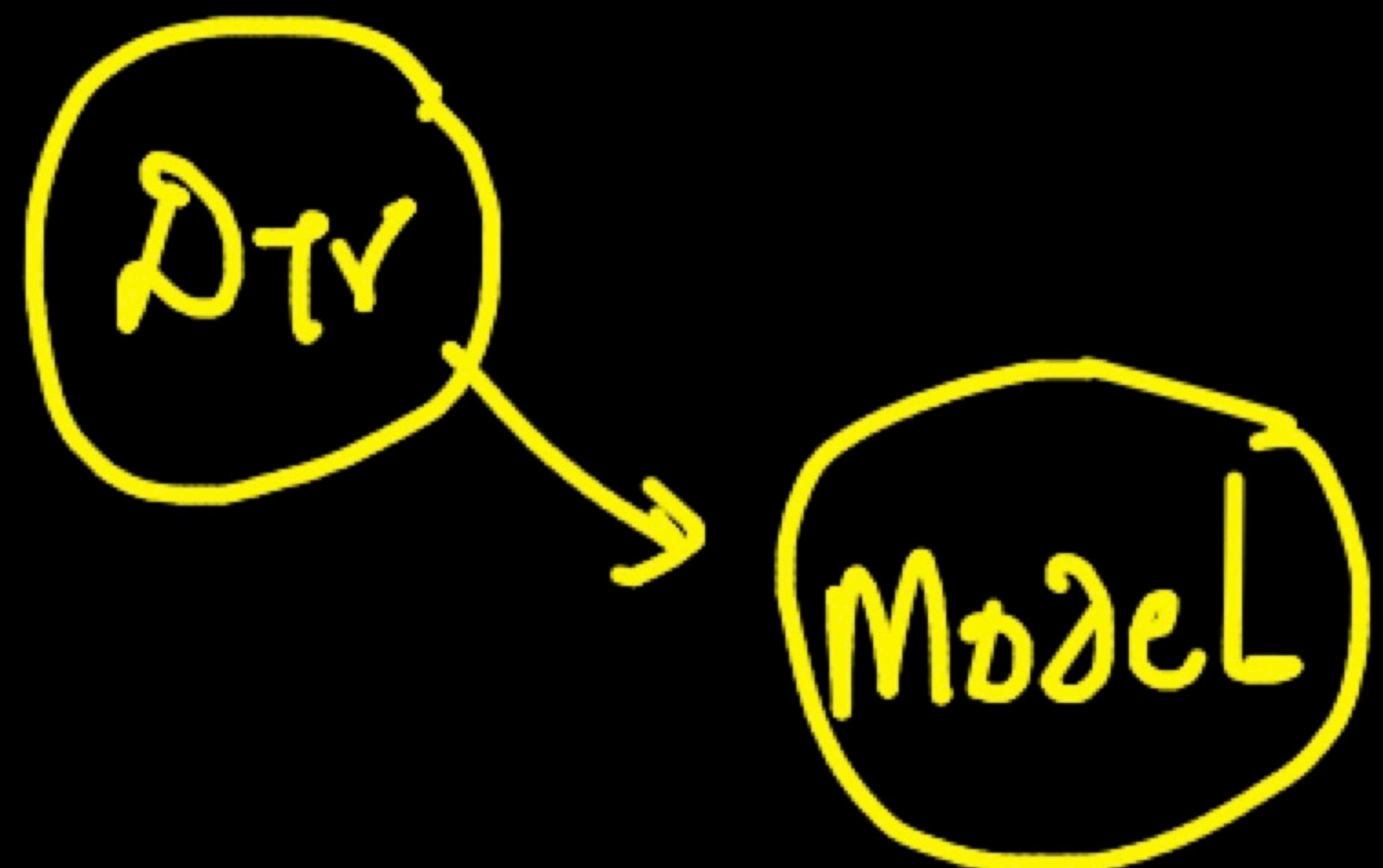
Determine if

\mathcal{D}_{TY} is similar \mathcal{D}_{Te}

=
Classification

EEDA:

range check y_i in D_{Tr} & D_{Te} $\rightarrow C_{1,2}$

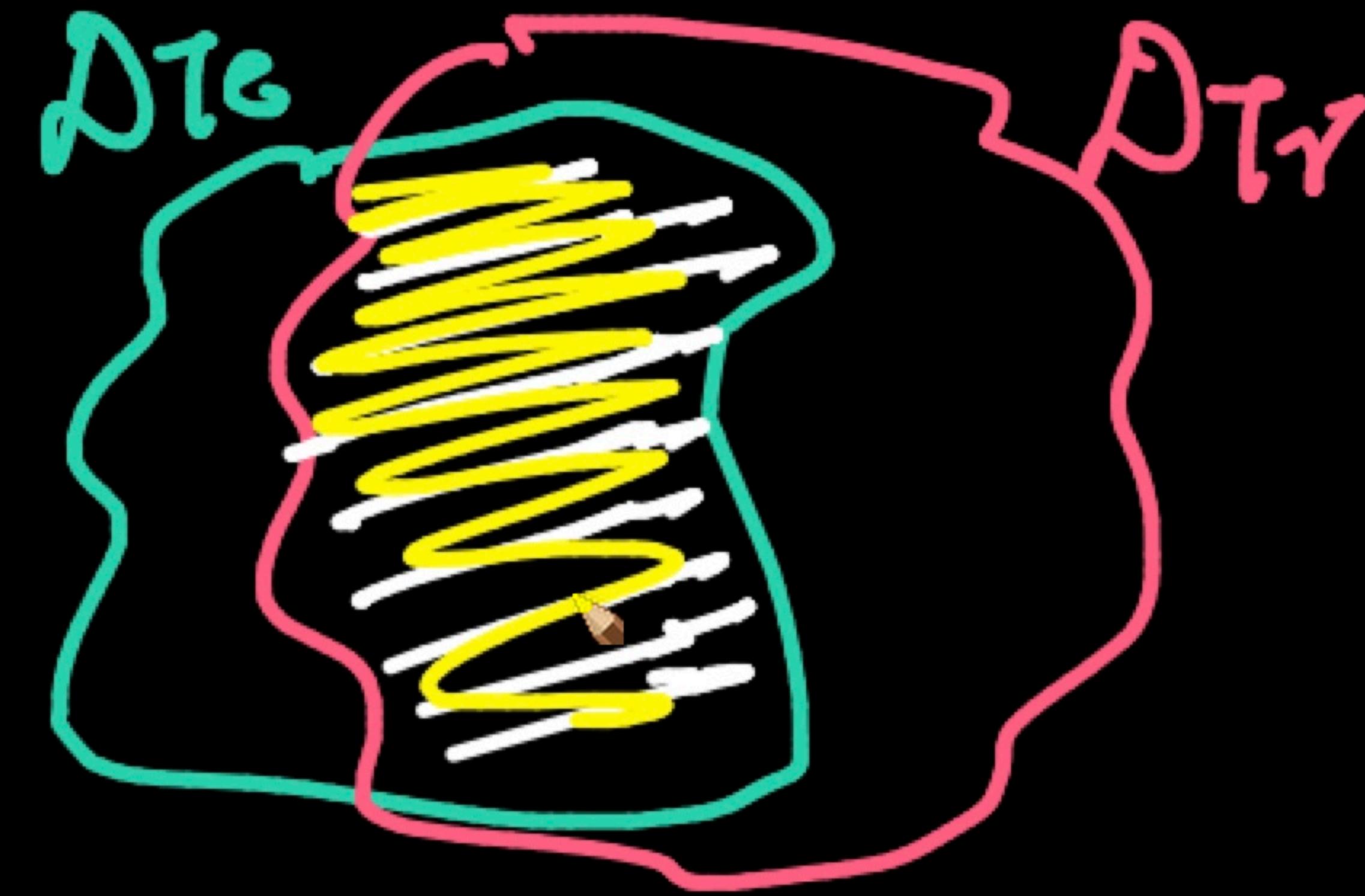


$D_{Te} \rightarrow \text{Model} \rightarrow F_{l_{Te}}$

Train: $F_l_{Score}_{Tr}$

$Vl \equiv \begin{cases} \text{if } F_{l_{Tr}} \approx F_{l_{Te}} \\ \quad - \text{then } D_{Tr} \& D_{Te} \text{ are similar} \\ \text{else} \\ \quad - \text{dissimilar} \end{cases}$

Follow-up:



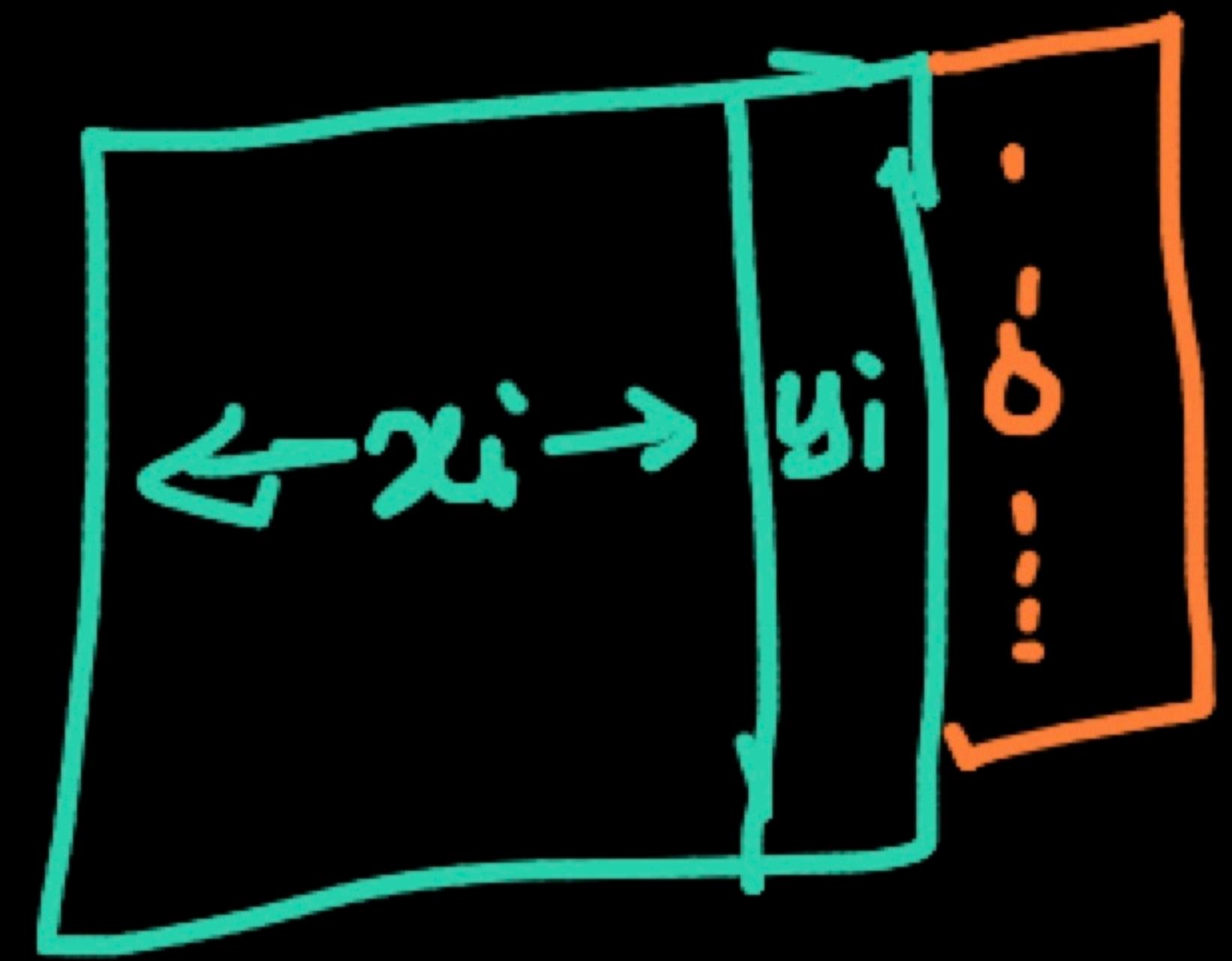
d-dim

$$Fl_{Tr} \approx Fl_{Te}$$

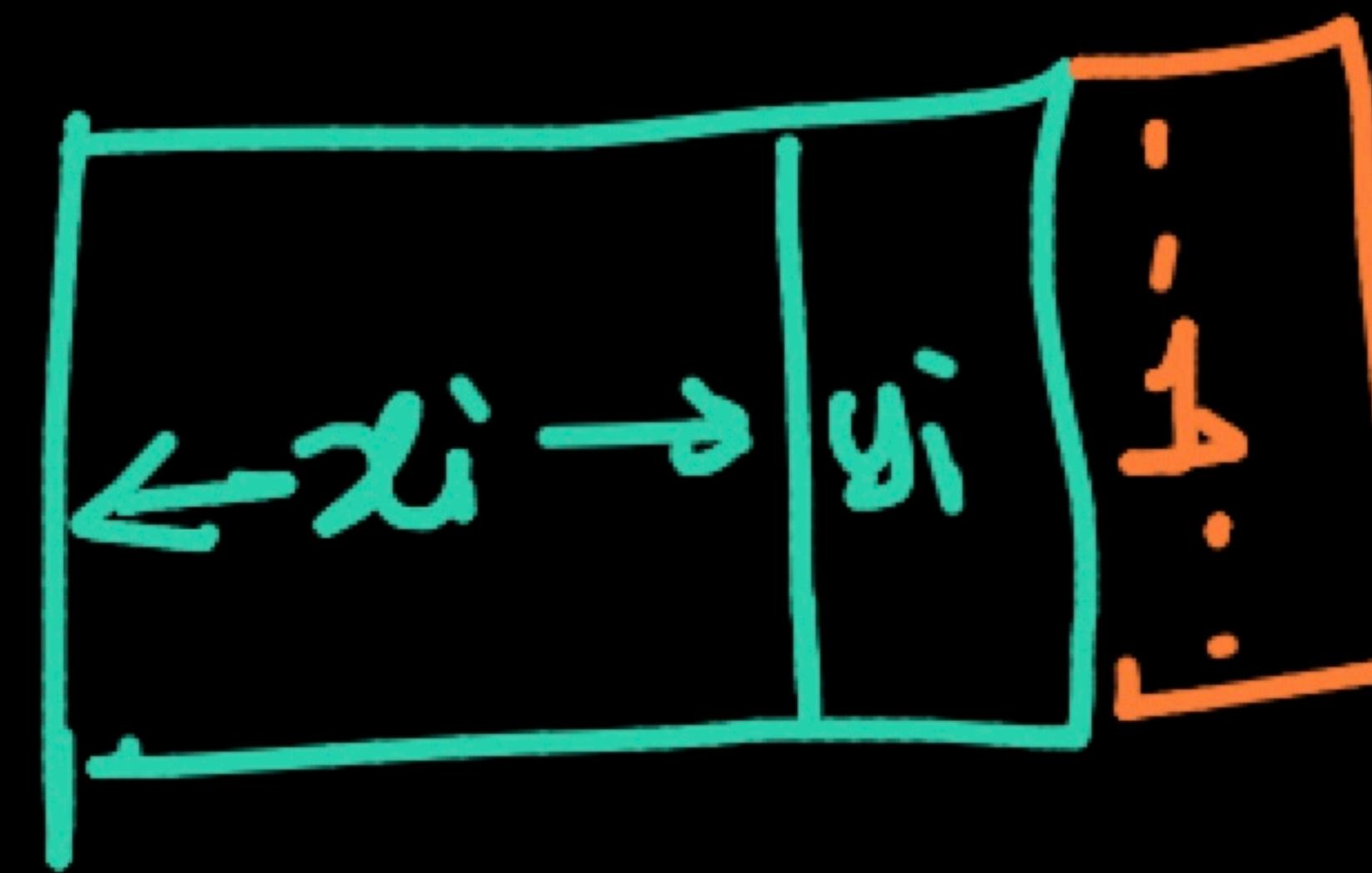
{ Is there a better technique
 + see if
 disb of DTv ↑ disb of DTe

$D \approx 1$

$D_{\text{new}} \rightarrow T_Y$



T_E



Complex but generic

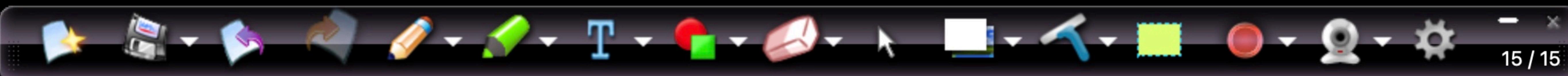
$\rightarrow \text{GBDT}$
(not overfit)

if GBDT can classify
then $D_T \neq D_E$
else similar

Alt: (Static) → ANDVA (variations)



Data Analyst



185: Solution with step by step

1st and Future - Player Contact



Code

The given problem is to find the high earners in each department. In other words, we need to find employees in each department with the top three unique salaries.

To solve this problem, we can use the following query:

```
SELECT d.Name as Department, e.Name as Employee, e.Salary as Salary  
FROM Department d, Employee e  
WHERE (  
    SELECT COUNT(DISTINCT Salary)  
    FROM Employee  
    WHERE Salary > e.Salary AND DepartmentId = d.Id  
) < 3 AND e.DepartmentId = d.Id  
ORDER BY d.Id, e.Salary DESC;
```

This query first performs a cross join between the Department and Employee tables to get all possible combinations of departments and employees.

Then, the WHERE clause filters out the employees who do not have a top three unique salary in their department. The subquery in the WHERE clause returns the count of distinct salaries greater than the employee's salary for the employee's department. If this count is less than 3, the employee's salary is in the top three unique salaries for the department.

Finally, the query selects the department name, employee name, and salary for the high earners and orders the result by department ID and employee salary in descending order.

Overall, this query efficiently finds the high earners in each department by utilizing a subquery and avoiding the use of expensive self-joins.

Related Solutions



SQL || DEPARTMENT TO...

MySQL



Beats 98% || EASY & SI...

MySQL



Using basic utilities of S...

MySQL

2+



185. Department Top T...

MySQL



Clean Code with Dense ...

MySQL

(Q2)

ML depth

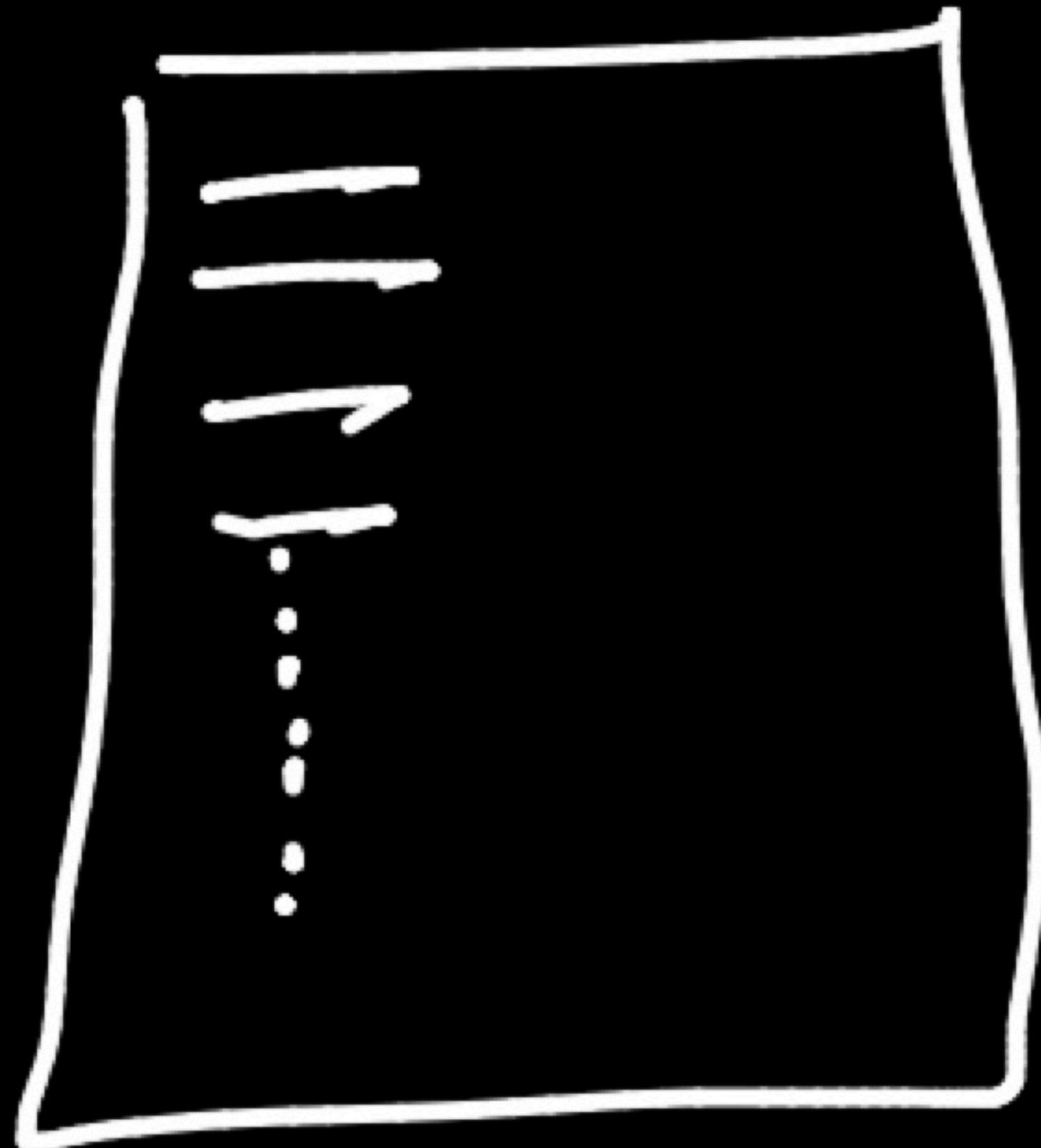
↳ SVM →

#SV; reduce #SV...
how do we make SVM
run fast @ evaluation
time

 $x_q \rightarrow \text{SVM} \rightarrow y_q$
(fast)

{ Random Forest → fewer & shallower trees
↳ GBDT preferred over RF? (Math..)
↳ loss fu.

GBDT



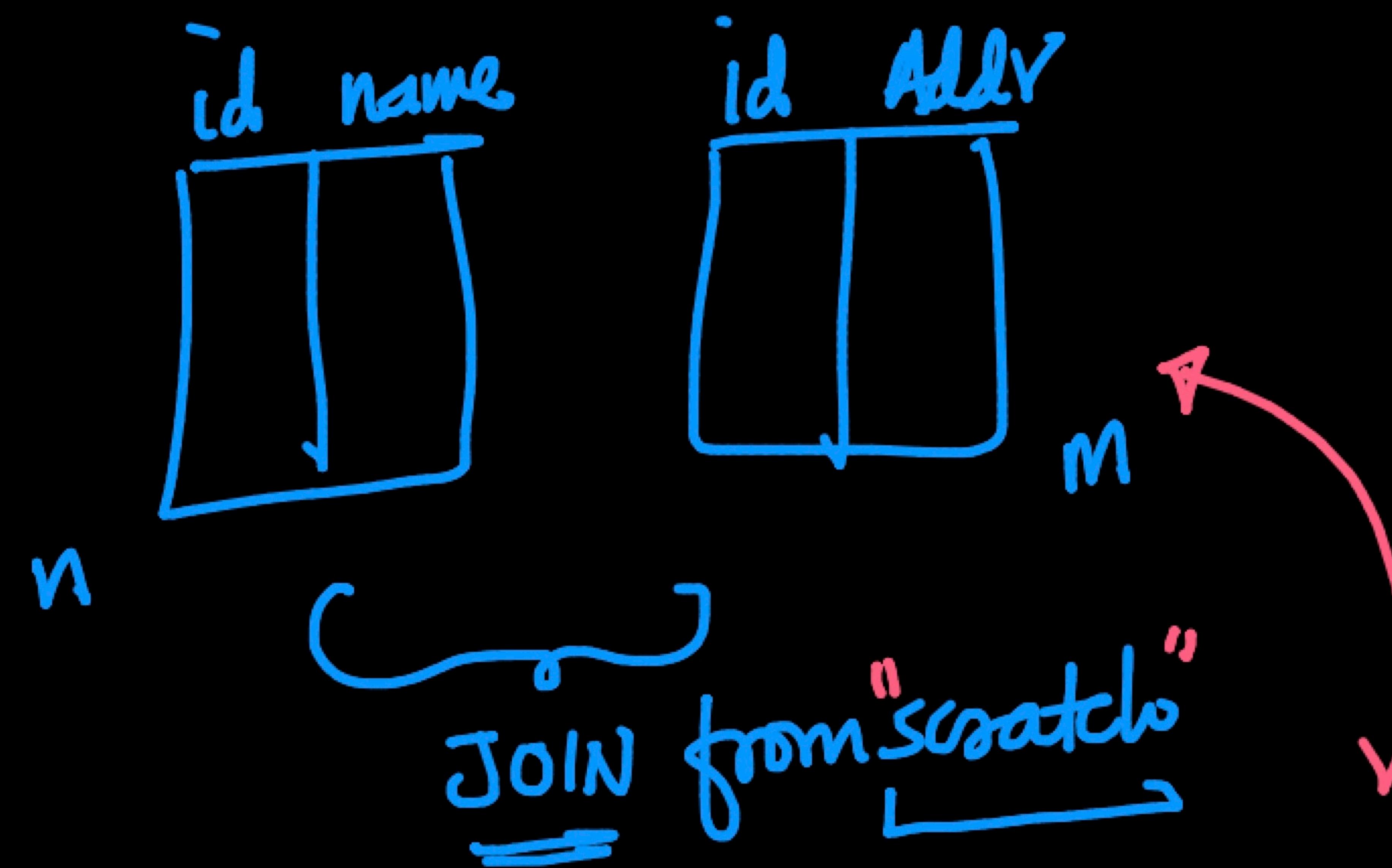
(1h)

Pseudo-code ..

(B3)

Python

DS MLE



implement library for [Pandas merge]

any library



pow(x,n) from scratch

$$x^n$$

\rightarrow n is integer

\rightarrow n is real

(8) [Transformers]

↳ most time taking aspect of training



how to speed up this?

GPT VS BERT

→ Speed up BERT →

