

obs: ~~forgetten~~

{ ML-2 = }

Unsupervised Learning

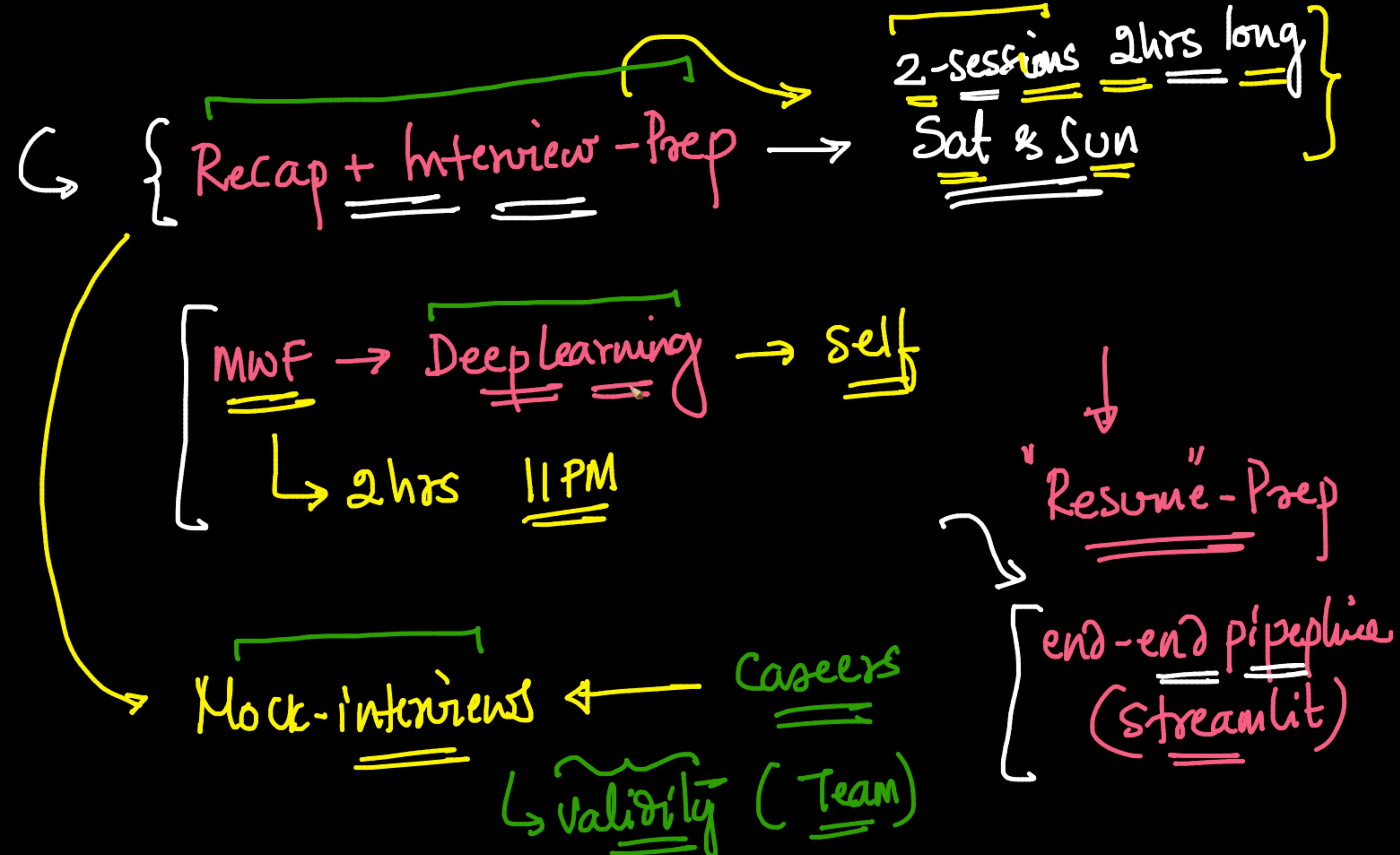
High-dim Data viz

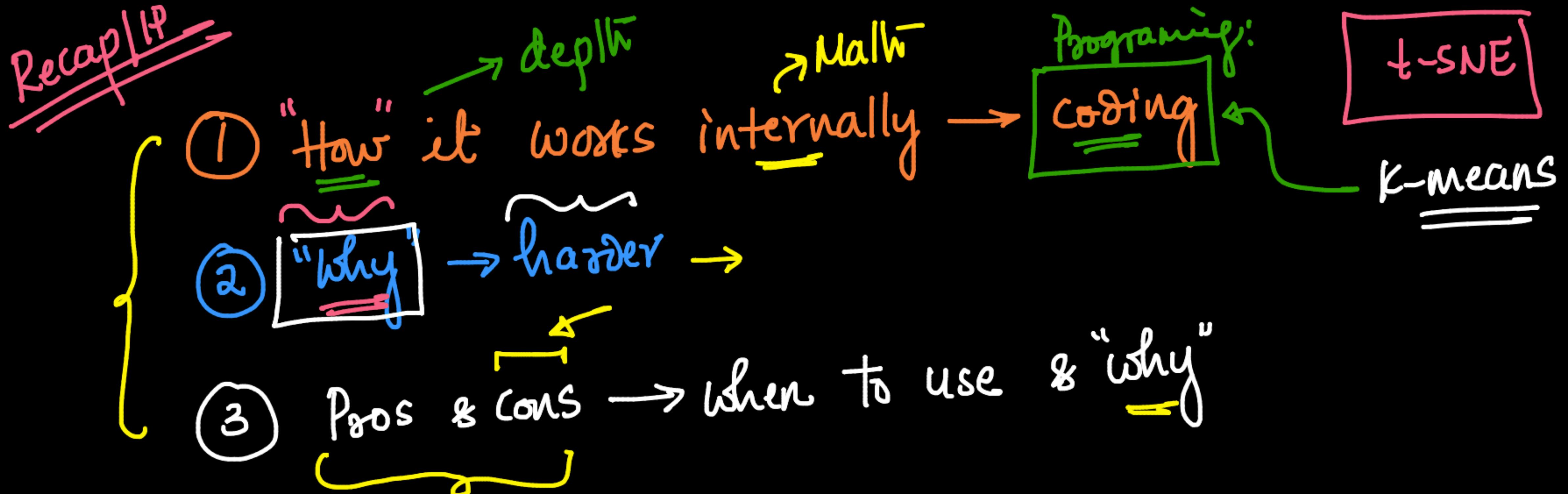
RecSys

Time Series Analysis & Forecasting

Agenda: =

Scenarios; recap; few interview Ques (why?)
✓ ✓ =





↳ Stackoverflow / Q&A / Wikipedia

Clustering:

- K-means, Kmeans++
Lloyd's
- hierarchical → linkages
- DBSCAN
- GMM
- High-dim
data viz
- UMAP, tSNE
- PCA,

RecSys:

- Content based: $U \cdot U^T$; $I \cdot I^T$; $U \cdot I^T$
- Coll. Filt → MF
 - ↳ $MF + \mu + b_u^T b_i$
- Hybrid-model & cold-start
- MF → clustering
 - ↳ PCA ; NMF
 - ↳ SVD
- Matrix Completion

TS Analysis & Forecasting $\xrightarrow{(+, *)}$

- ACF, PACF; decomposition
- imputation; TB-splitting; stationarity
- AR, MA, ARIMA; SARIMA,
SARIMA
- FB Prophet
- SMA, EWMA; DES; TES
- C-l m TS
- Change detection

- Anomaly detection
outlier,
- Isolation forest
 - DBSCAN; LOF
 - elliptical envelope
 - one-class SVM
 - RANSAC
 - { - Isotonic-reg (Misc..)
Platt scaling; elastic Net

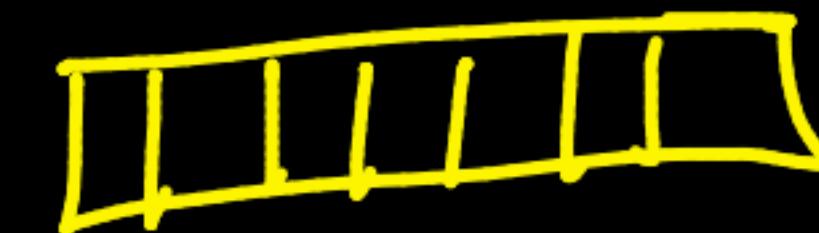
Q1

$n = \underline{10M}$ e-commerce
products

cluster

$d = \underline{40}$ (let)

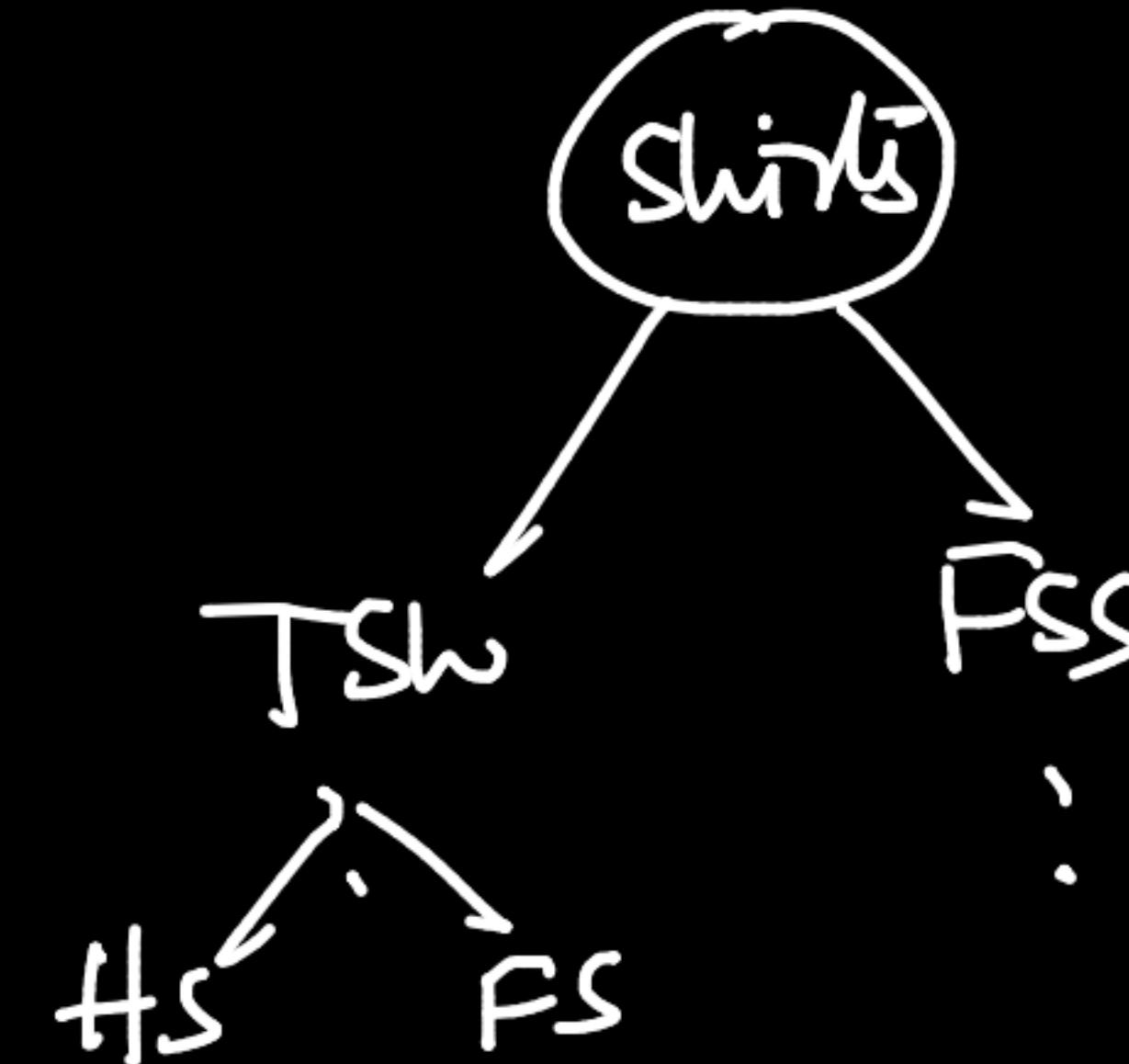
(open-ended)



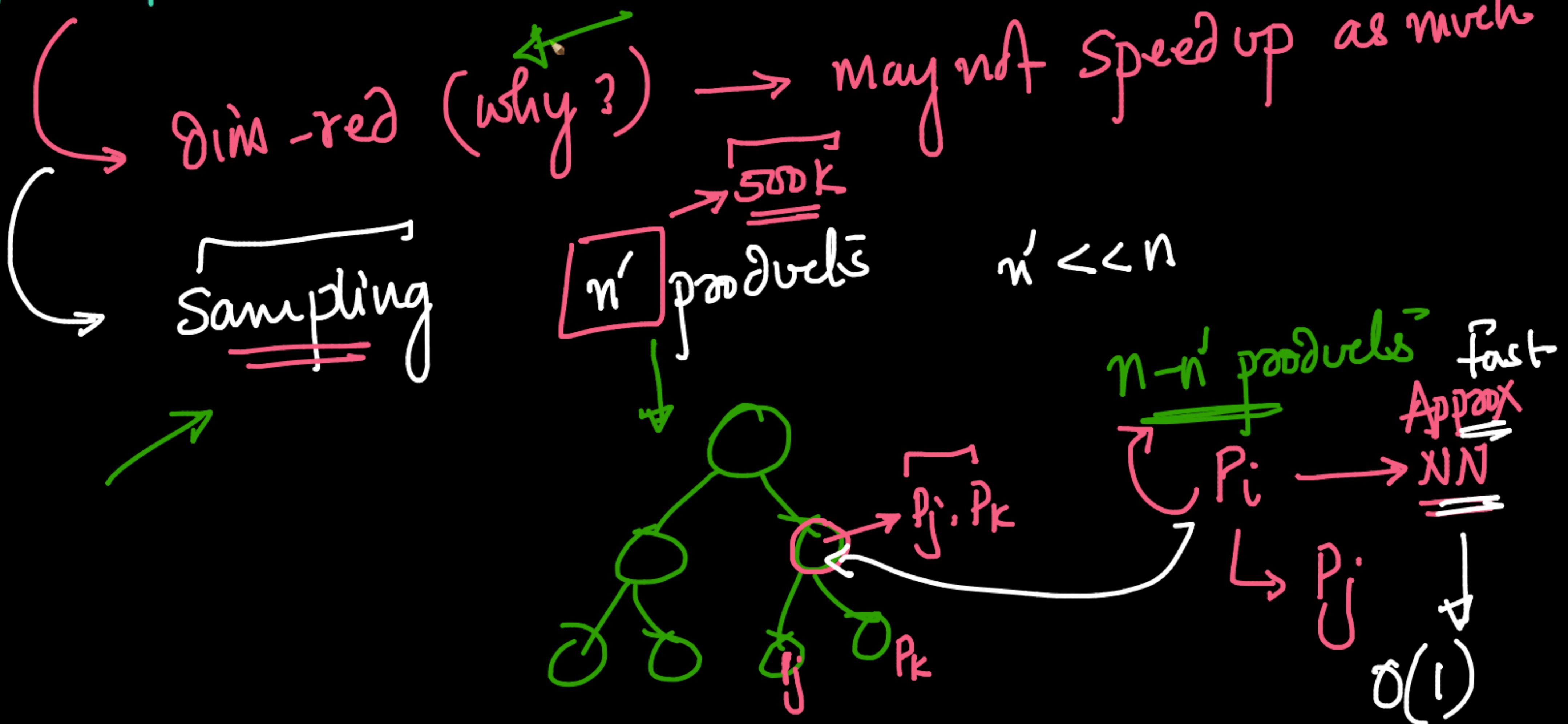
: item-vec

Hierarchical → comp-exp → 12 hrs
natural

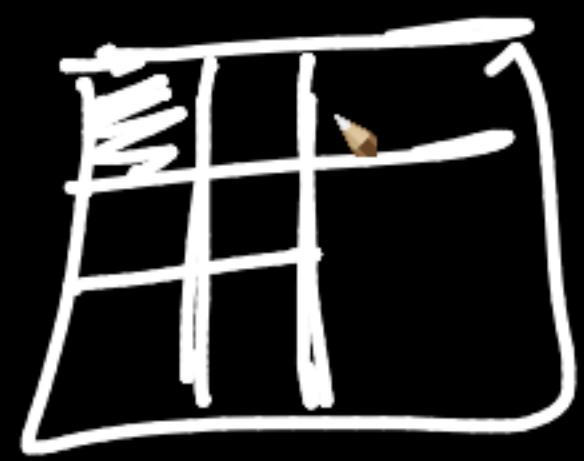
→ K-Means



Speed-up (ex SDE)



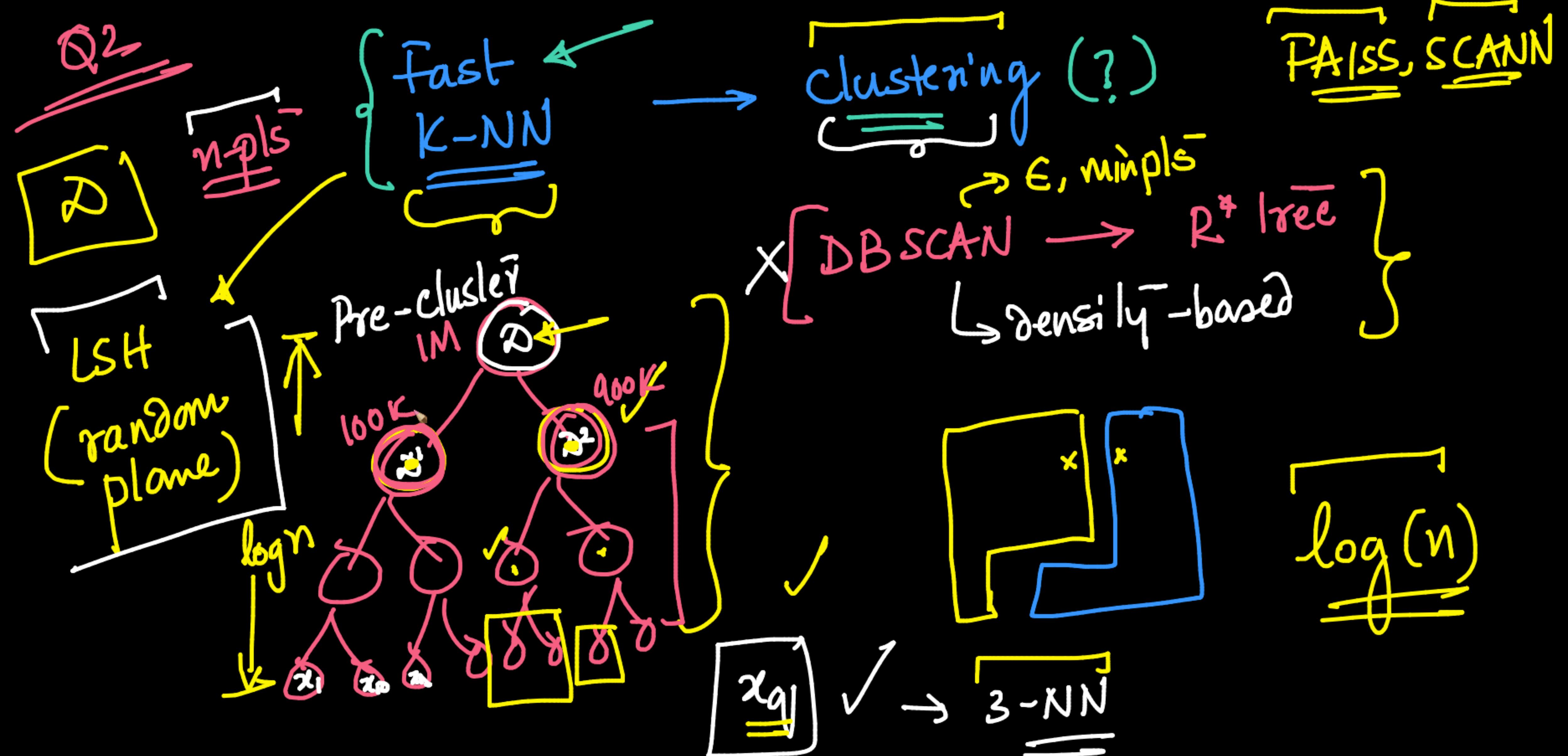
C

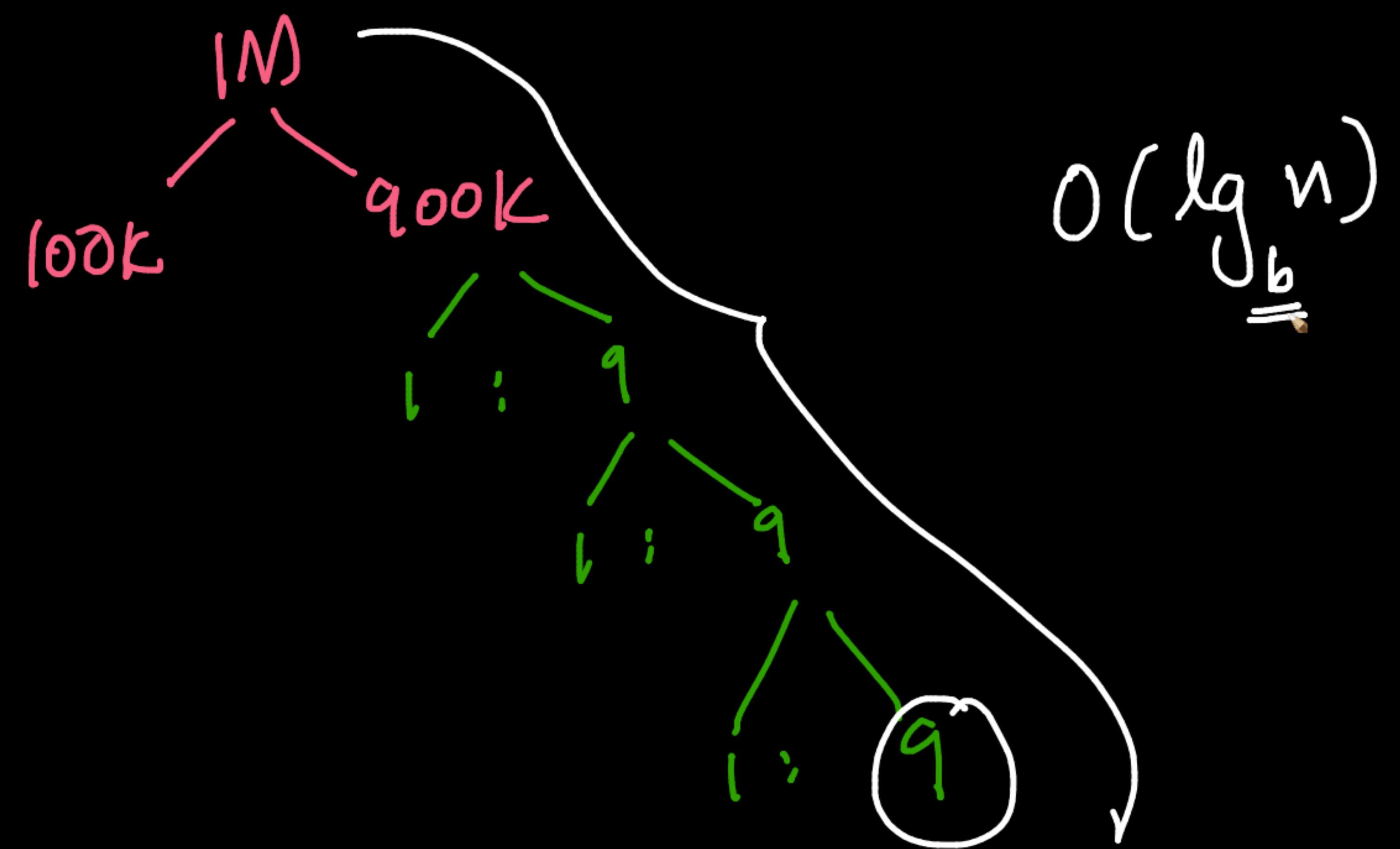


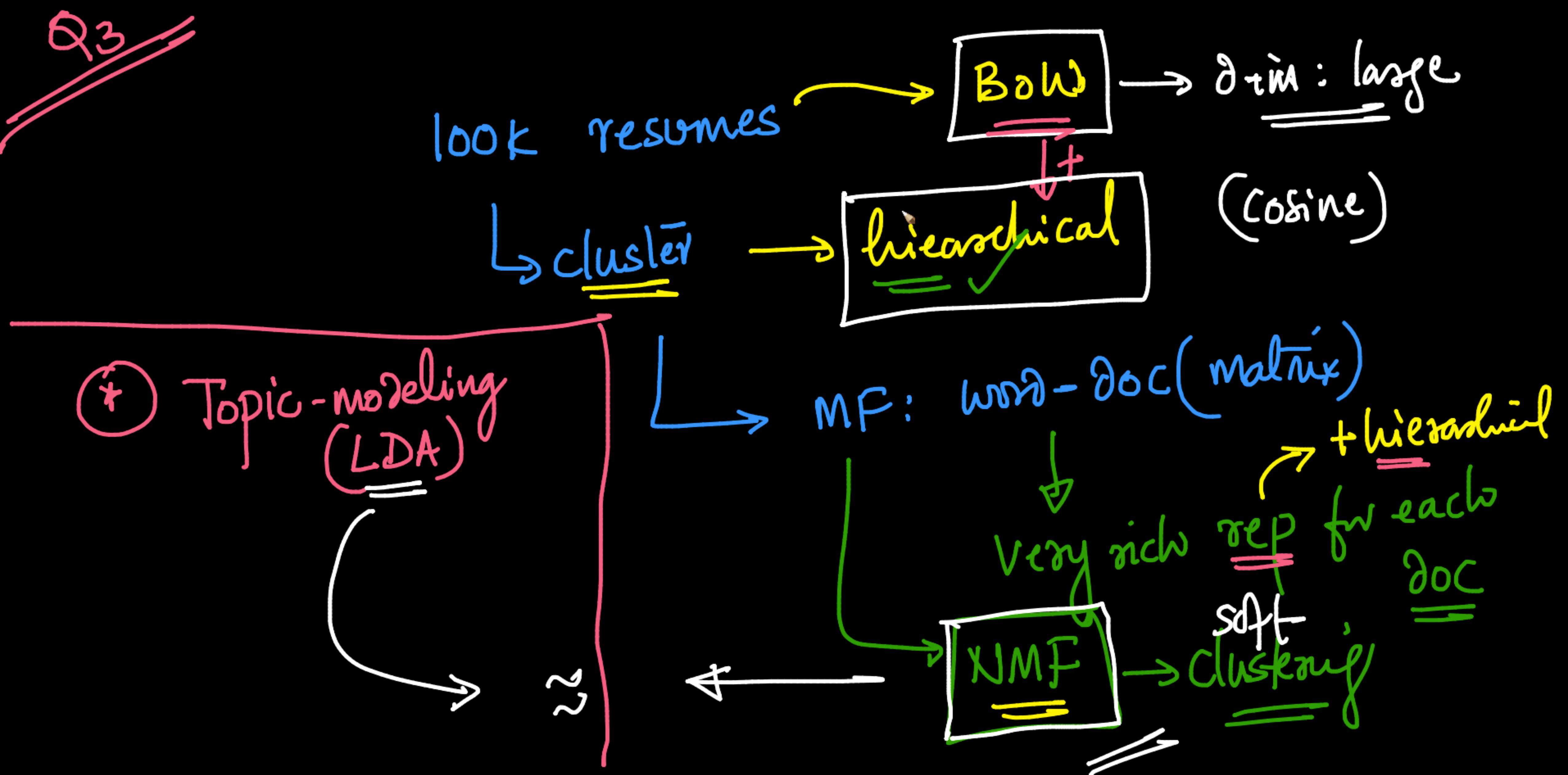
c_1

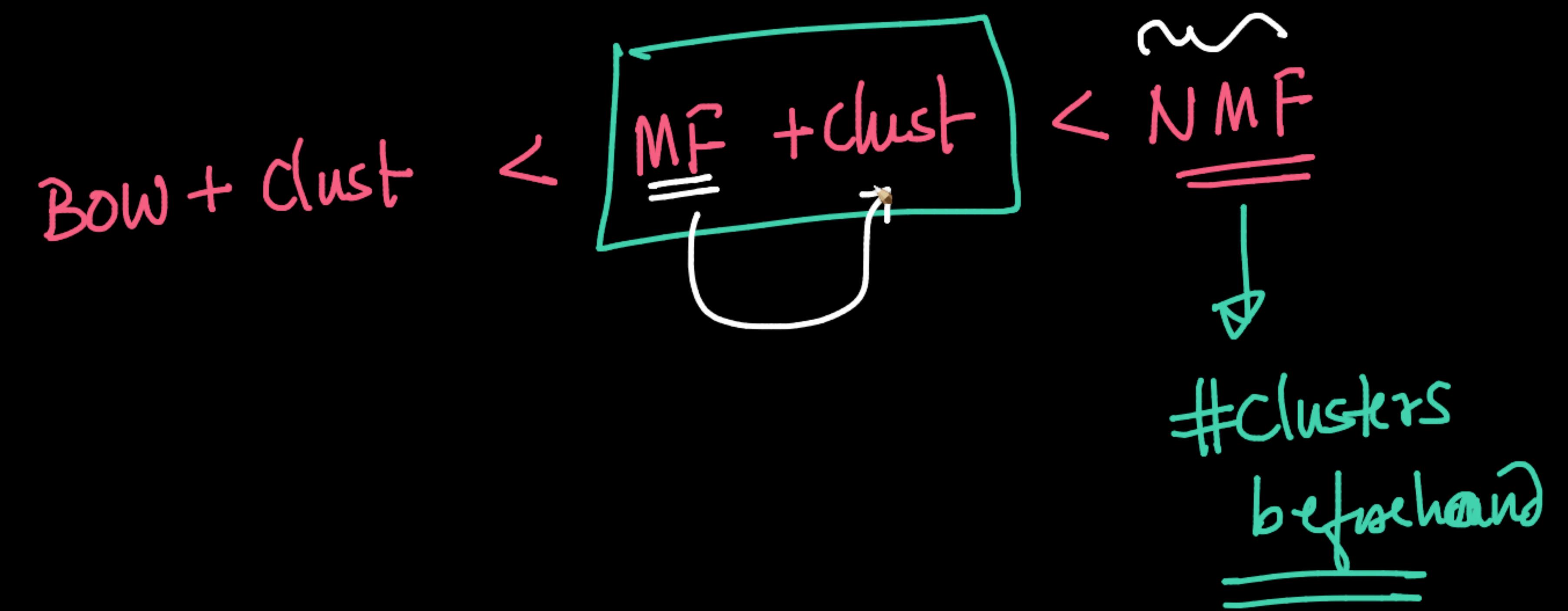
o o o o o o o
 { {

c_2



 $O(\lg n)$
 \underline{b}





Q4
n = 1MM
 $d = 2$



lat & long
Src, dest
Cluster

Median to deliver

$\leq K$

deliveries (Zepto)

dark stores

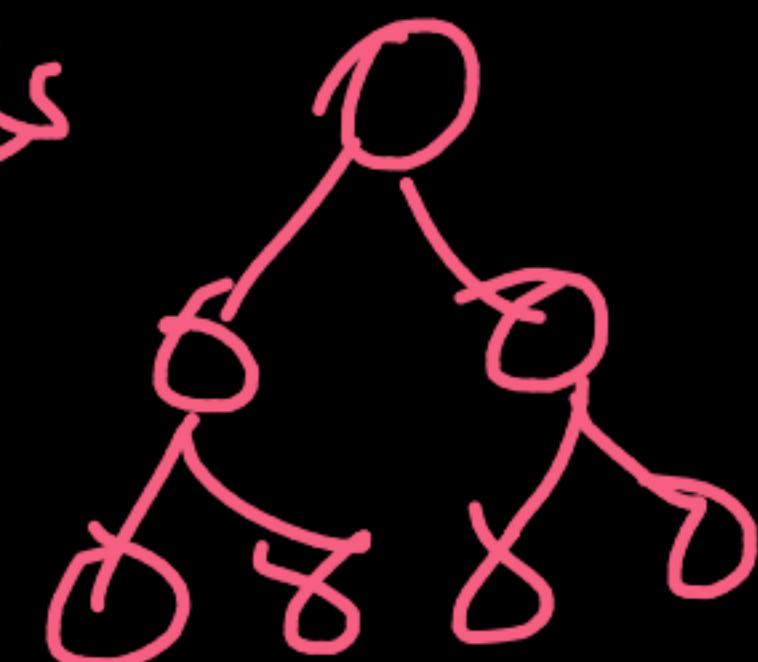
DBSCAN (ϵ, minpts)
density difference
K-Means
 \approx equal size
clusters





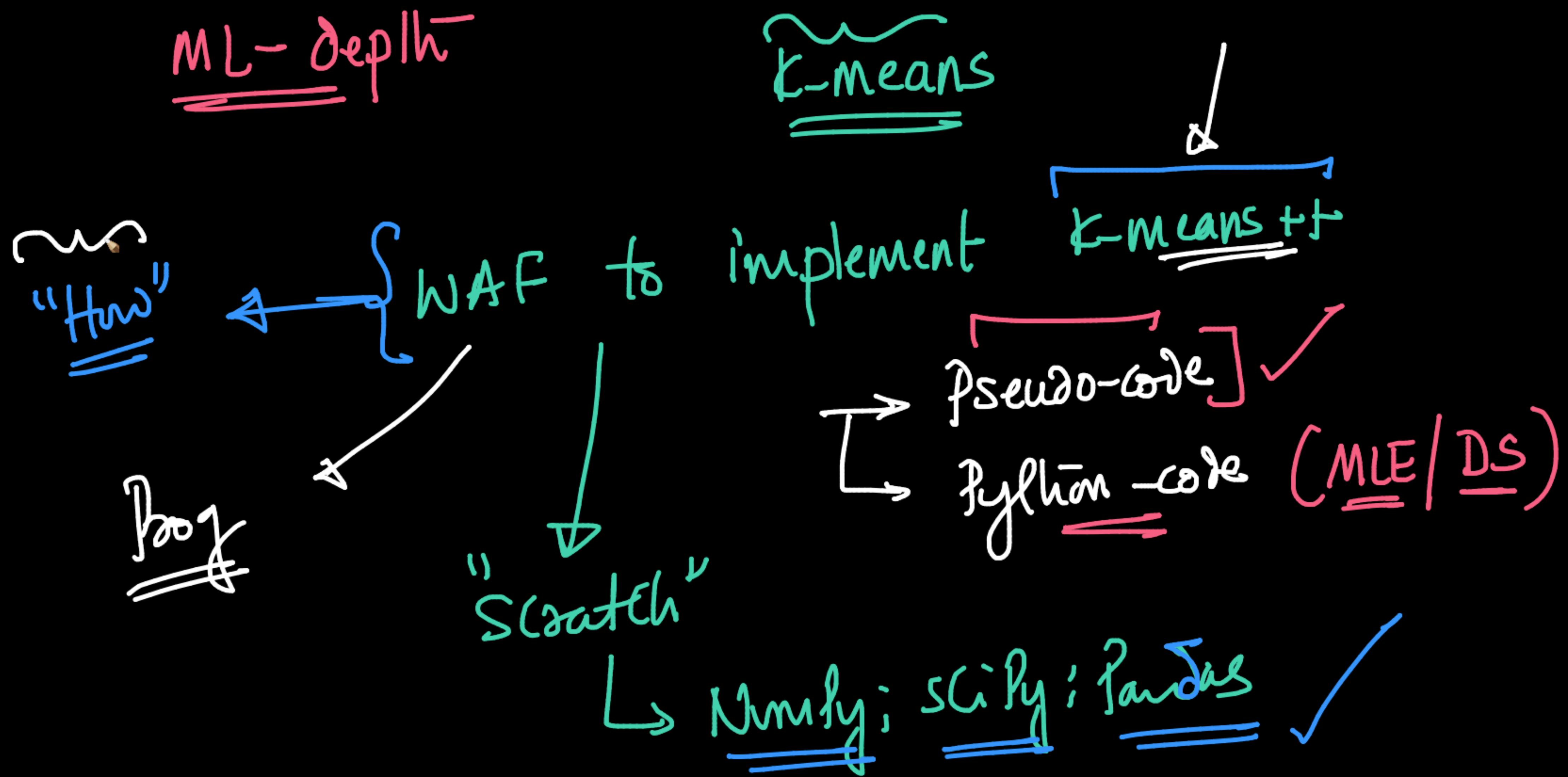
adv & disadv of each technique + "Hui"

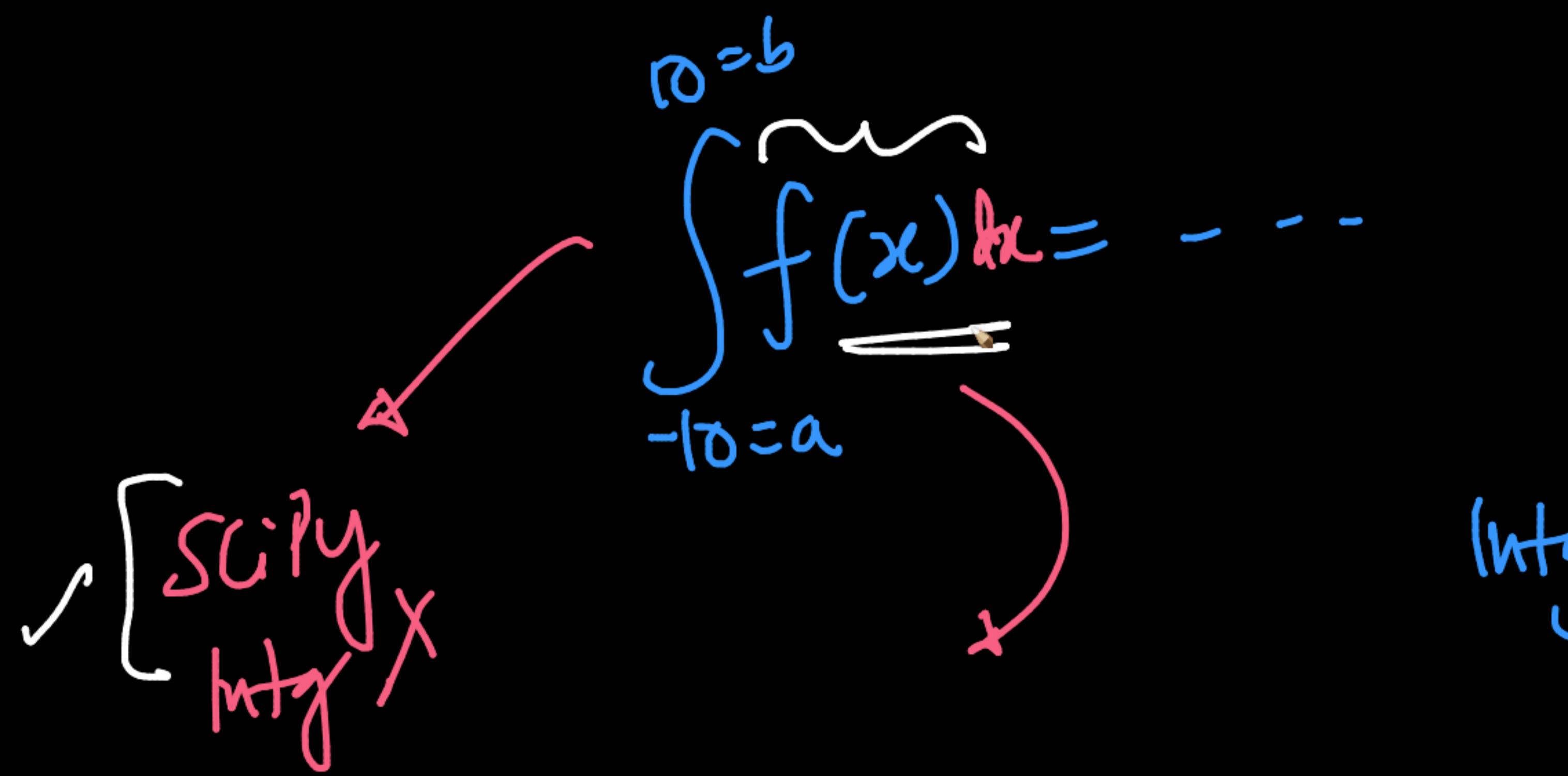
Scenarios → constraints
→ Properties



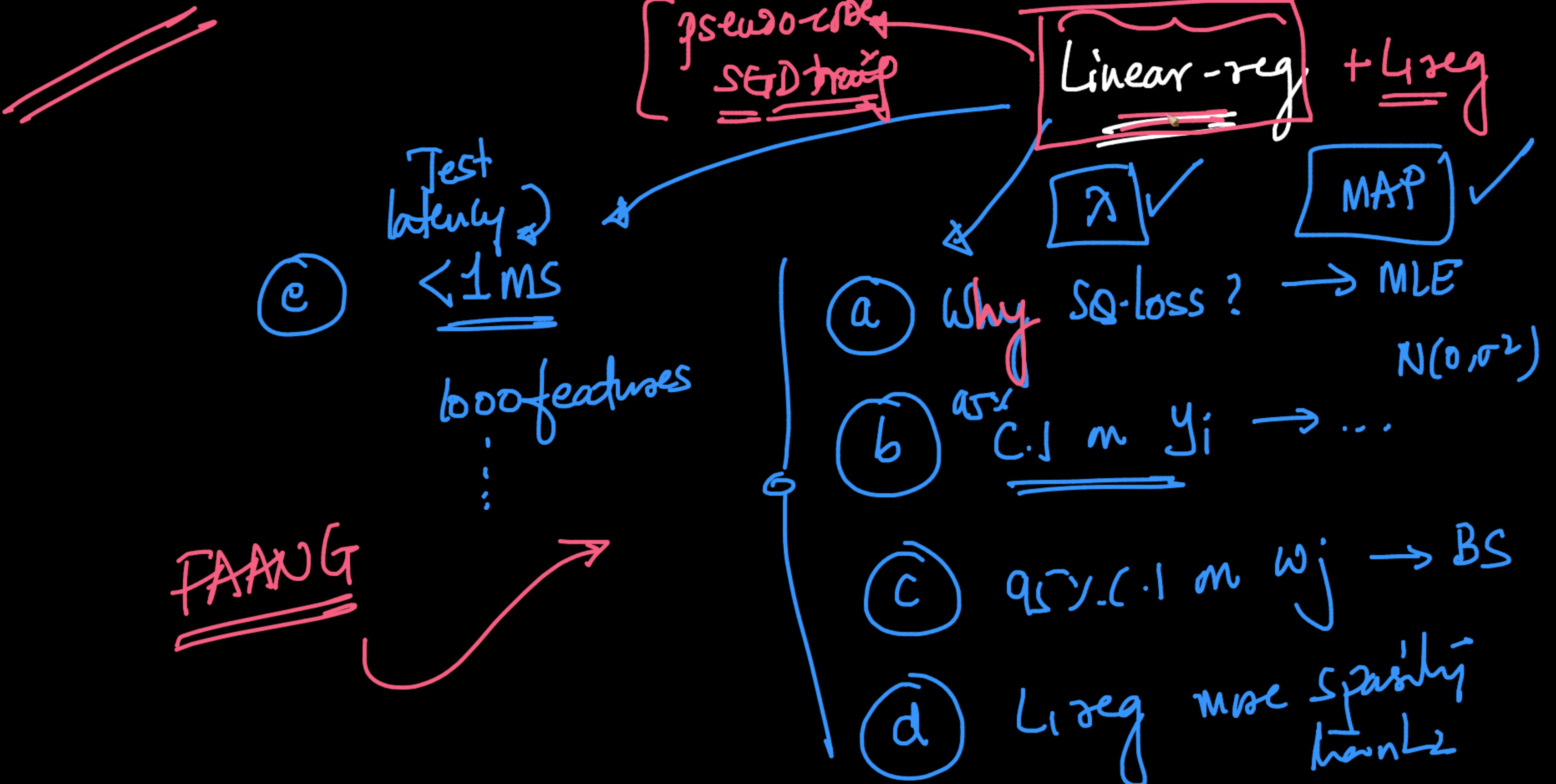
QS

ML-depth





code
 $\text{Intg}(f, a, b)$





Why tSNE

not preserve "global" struc.

→ preserving
neigh

↳ crowding problem ✓
↳ stochasticity → may not
different solns for different runs

$x_i^s \rightarrow y_i^s$
high low

[HINT: Obj.-fn in OPTIMZN-problem]

$\min_{P,Q} KL(P, Q)$

$d \text{-dim } p_{ij} = \text{Prob that } x_i \text{ & } x_j \text{ are neighbors}$

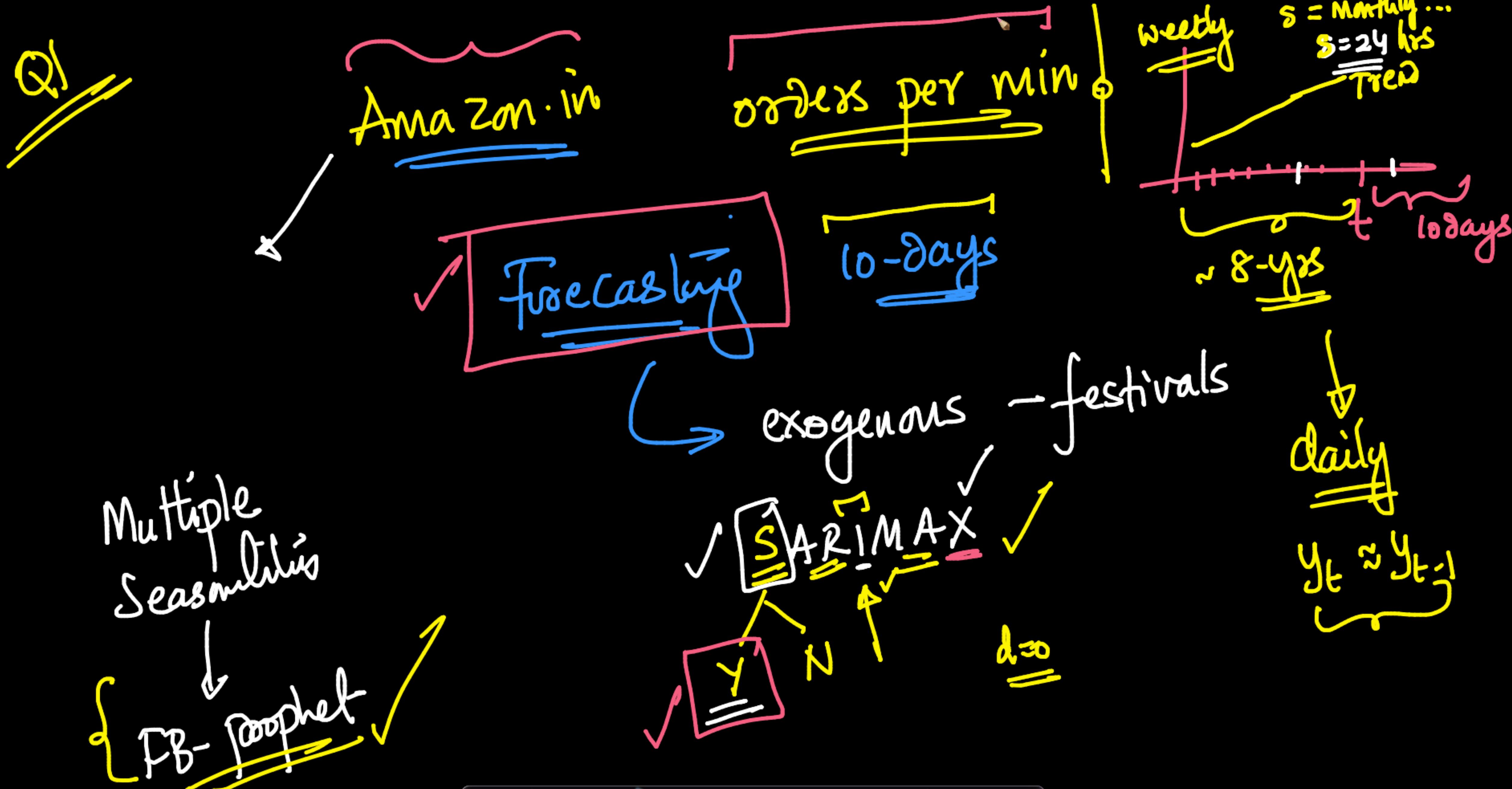
$d' \text{-dim } q_{ij} = y_i \text{ & } y_j \text{ are neighbors}$

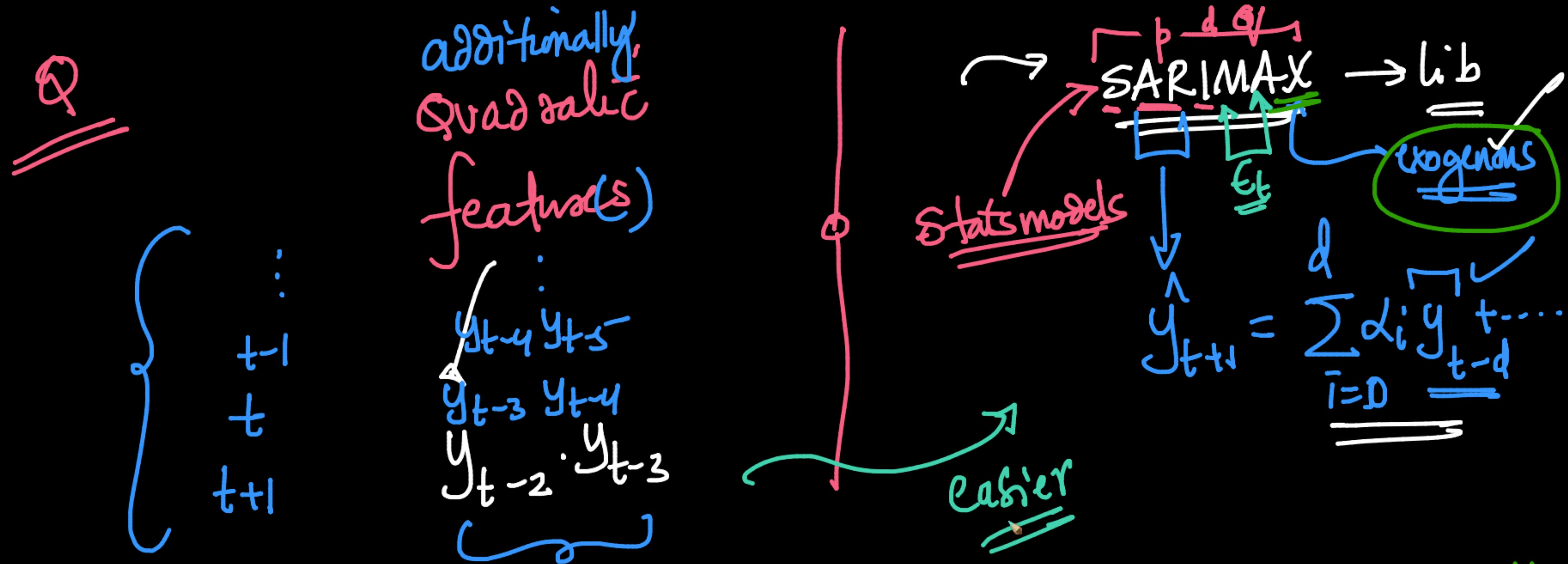
v-small

Gaussian $\rightarrow P_{ij} \approx Q_{ij} \leftarrow \text{disb}_1$

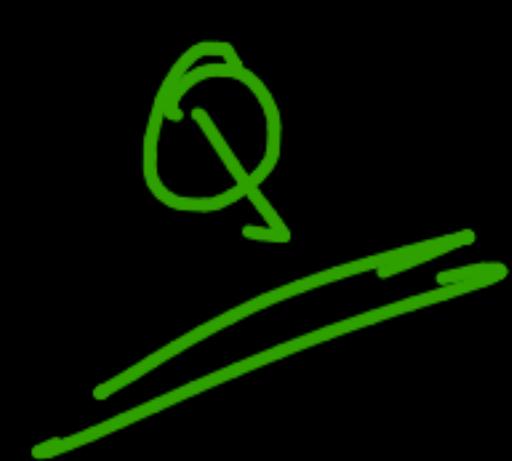
x_i $\xrightarrow[=4\sigma]{>3\sigma} x_j$

$\sum_{ij} \tilde{P}_{ij} \log \frac{P_{ij}}{Q_{ij}}$





Introduce non-lr features into existing impletn
of SARIMAX



Pseudo code on

PACF

What

ACF

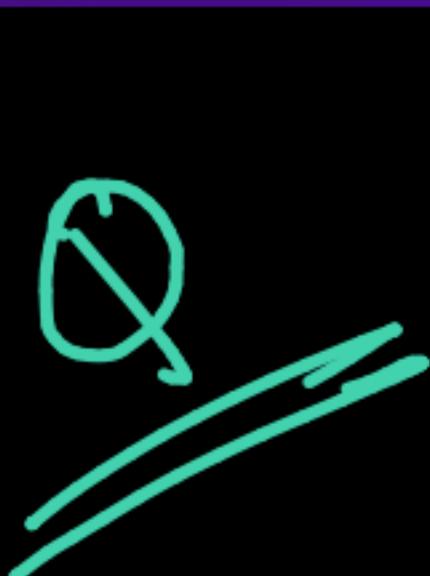
"How"

$$\delta_t = y_t - \underline{f_{AR}(y_{t-1})}$$

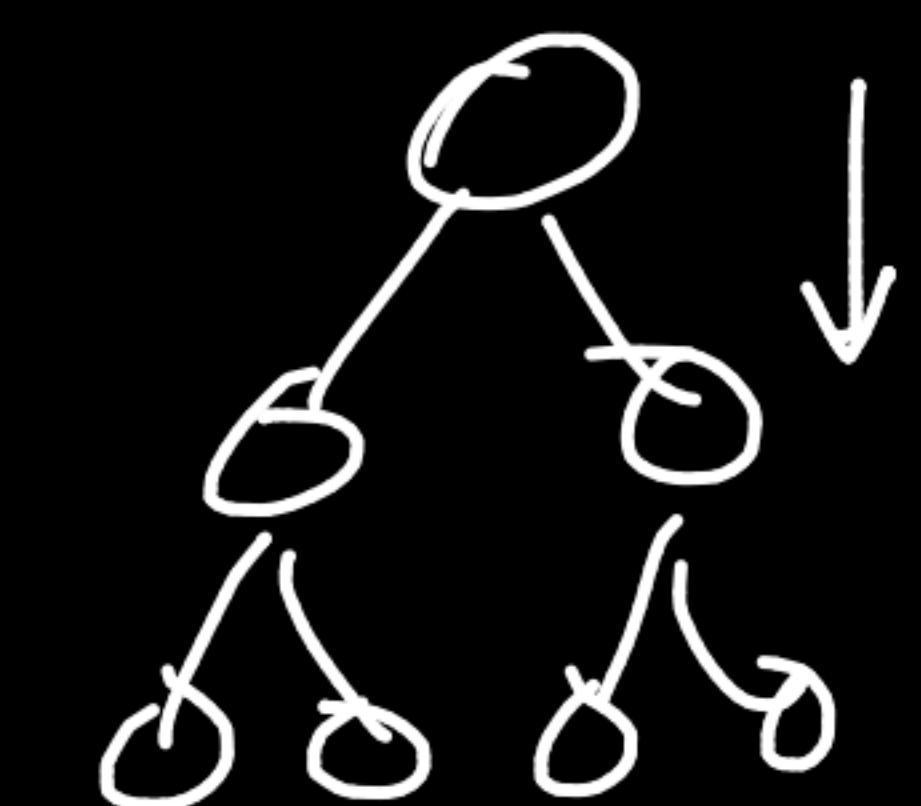
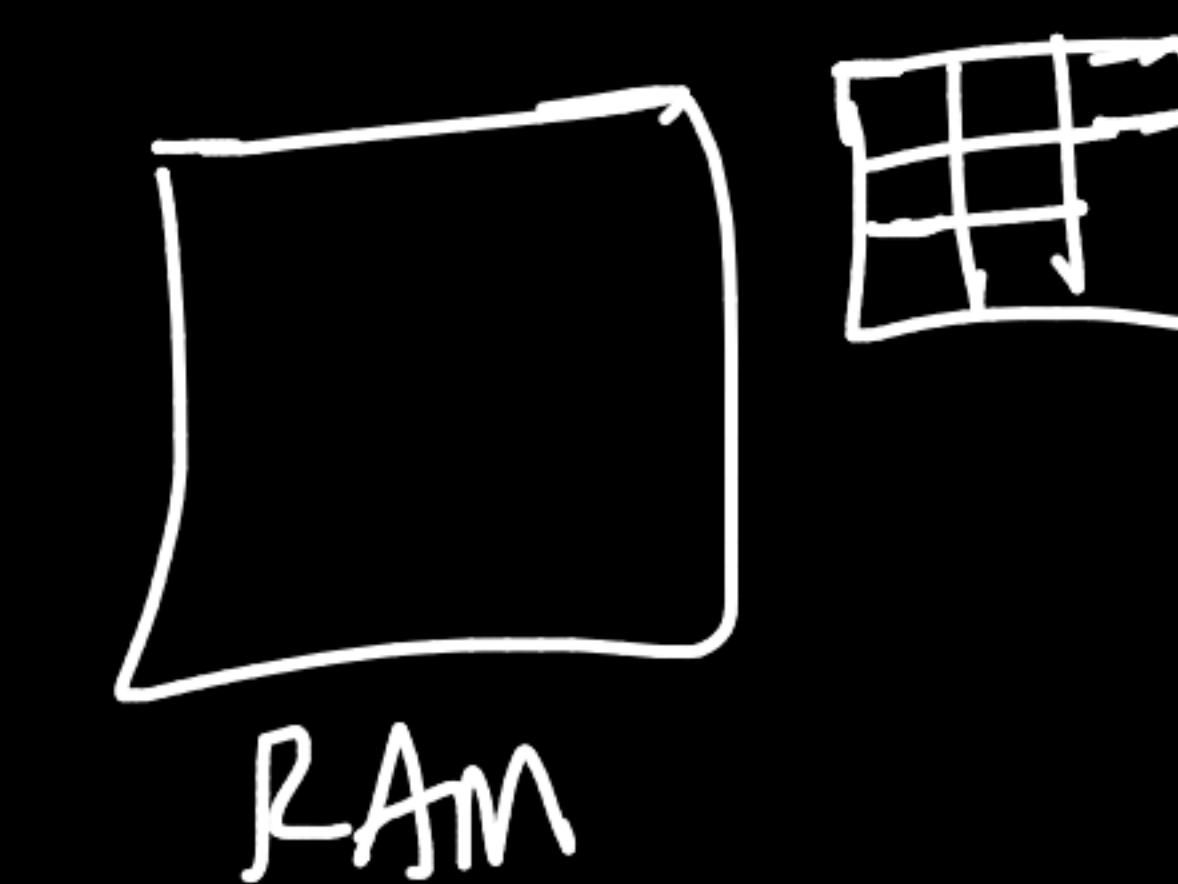


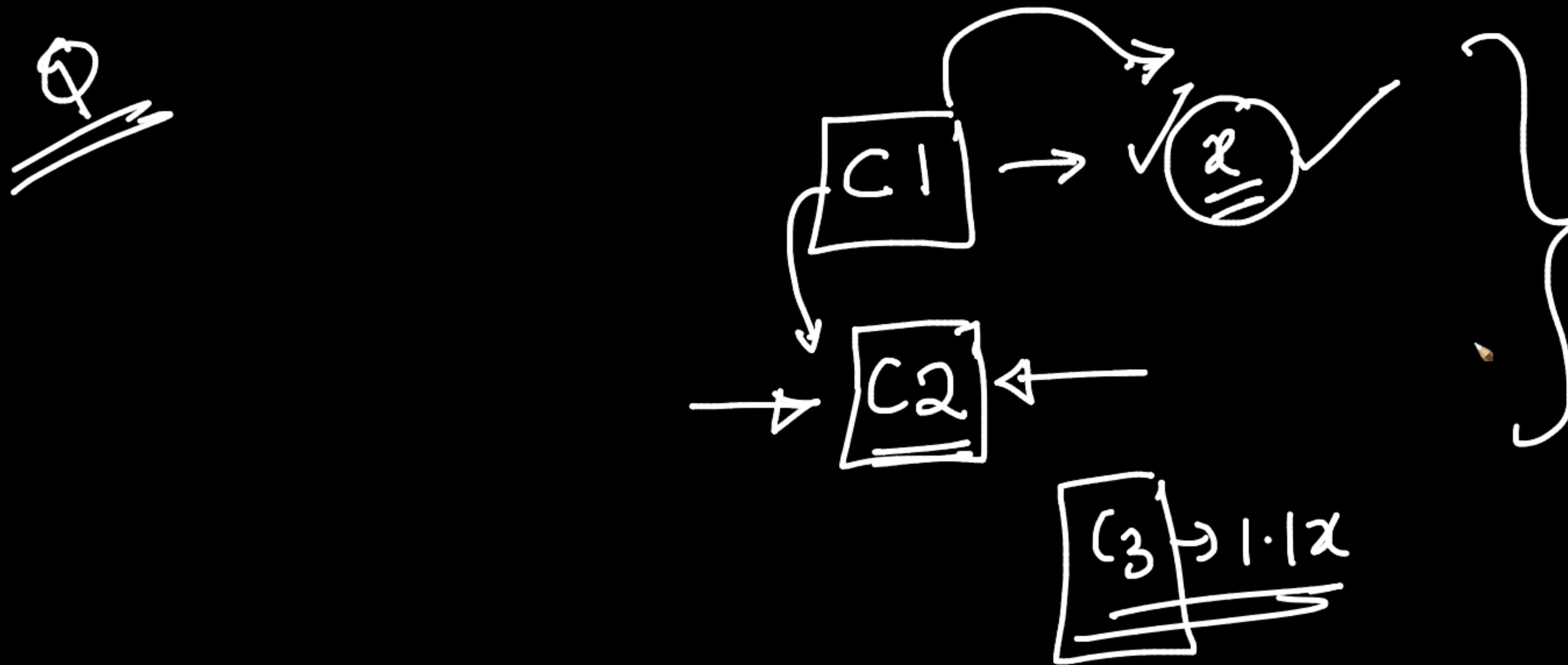
$$y_t \leftarrow f_{AR}(y_{t-1})$$

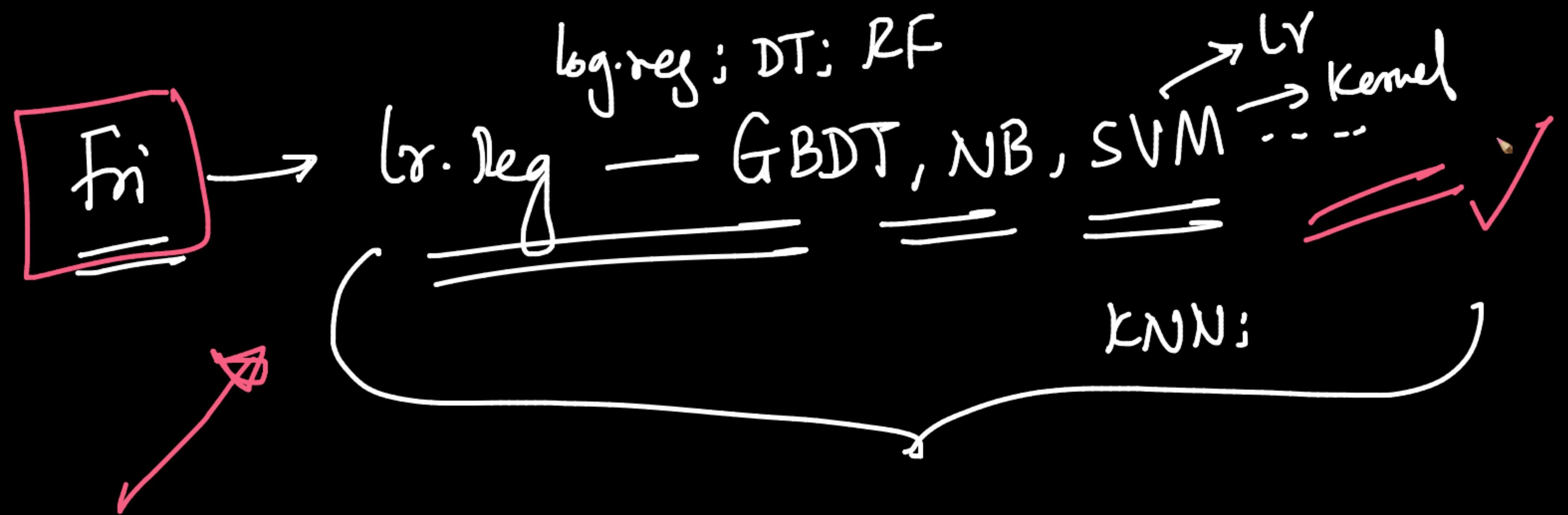
Corr. Coeff
(y_{t-2}, δ_t)

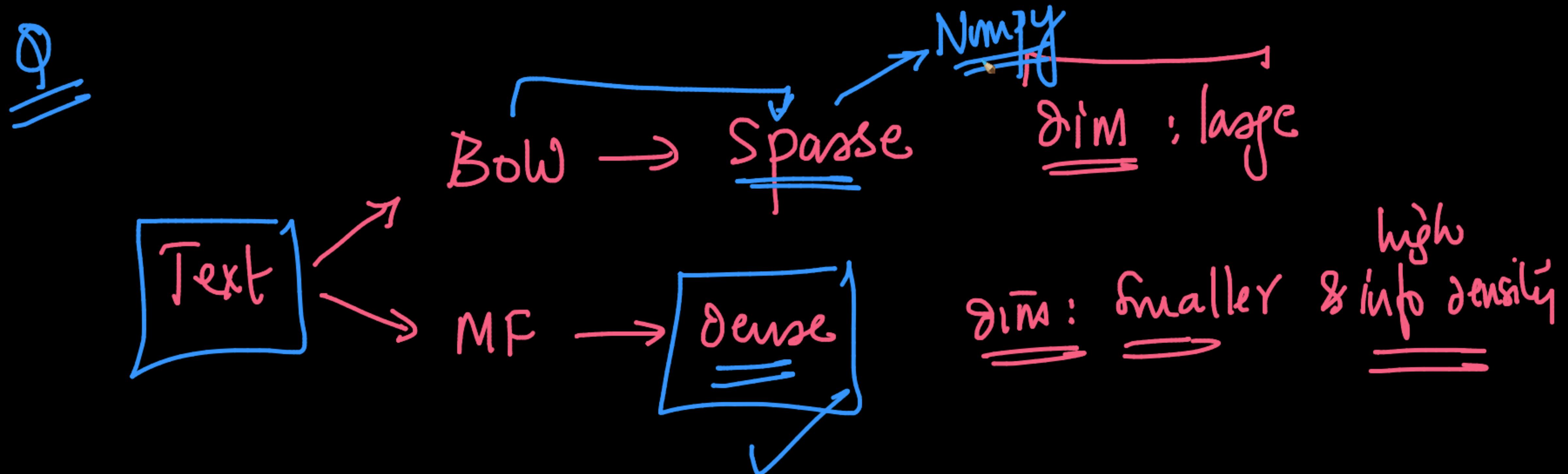


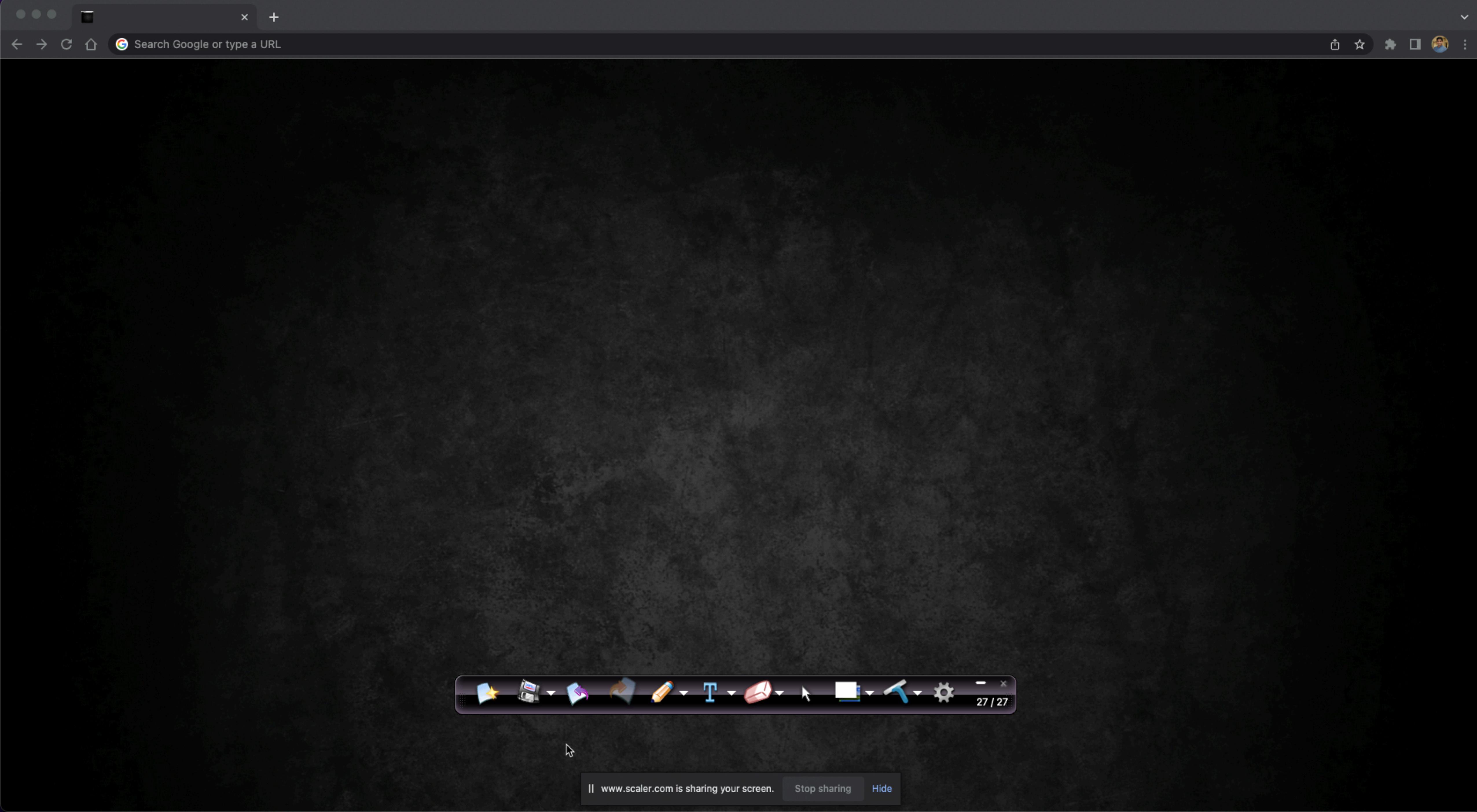
disb impletn of hier-clustering \rightarrow divisive











II www.scaler.com is sharing your screen.

Stop sharing

Hide