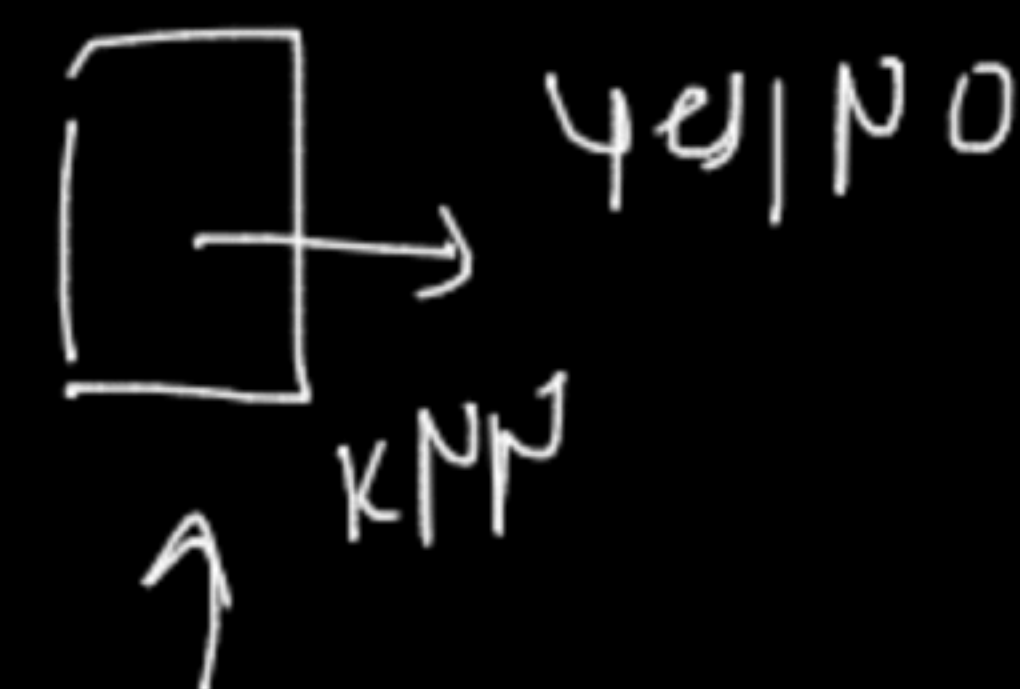
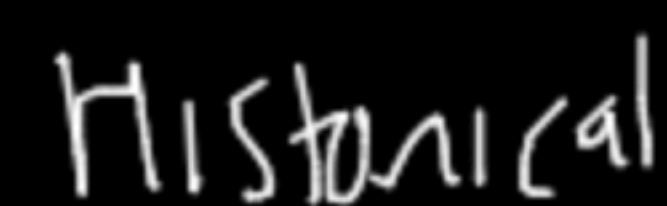
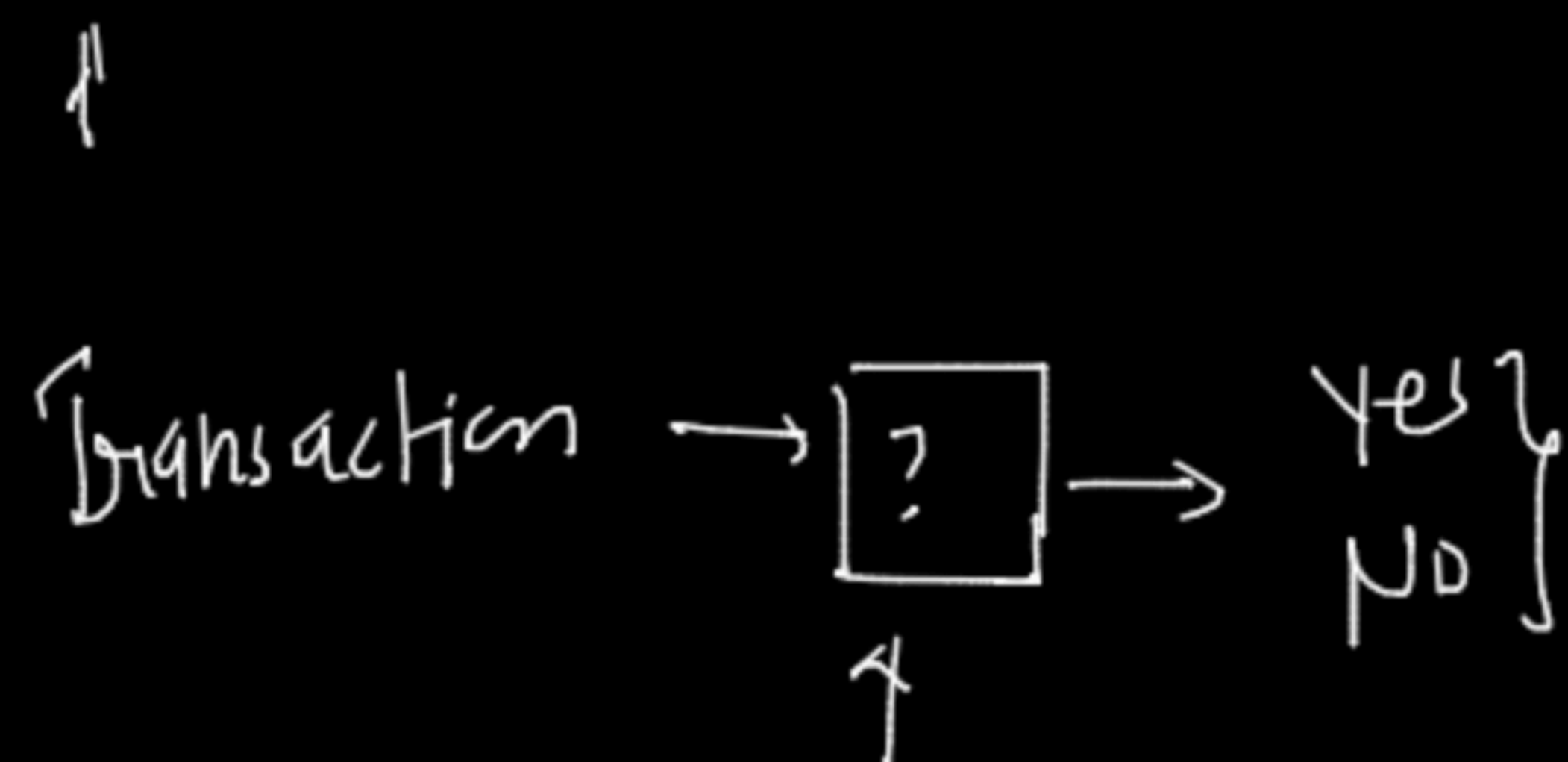


Credit card fraud



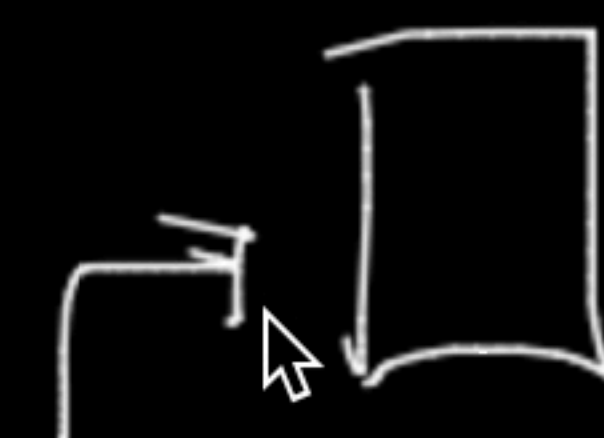
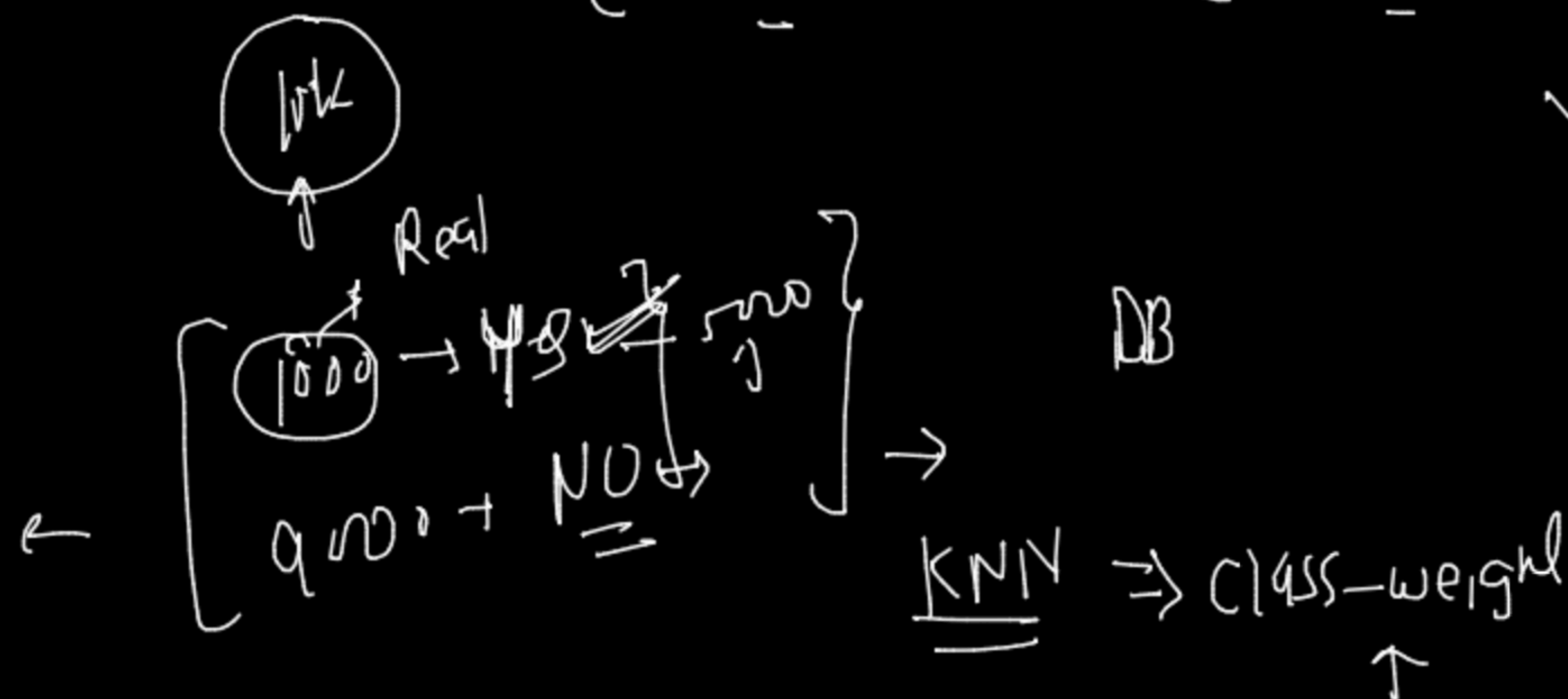
	transa-history	F2	F3	statue
c1	-	1	1	NO
	-	1	1	NO
	-	1	1	NO
	-	1	1	NO
	-	1	1	NO
	-	1	1	44
	-	1	1	44
	-	1	1	44

Real

→ 5000

* * [17 oversampling
2) unders

↑
balance



D_x (X, y) $\rightarrow \in [yes, no] \rightarrow$ binary classification
 input \rightarrow output

(+w)

KNN \rightarrow

$\left. \begin{matrix} c_1 \\ c_2 \\ c_3 \end{matrix} \right\} +ve$
 $\begin{matrix} c_1 \\ c_2 \\ c_3 \end{matrix}$

ABC

Reviews \rightarrow [?] \rightarrow Neutral $\left. \begin{matrix} +ve \\ -ve \end{matrix} \right\} = 3$
 text

$\begin{matrix} +ve & -ve & -ve \\ 3 & 2 & 2 \end{matrix}$

multiclass-classifications

$\frac{7}{2} \cdot \frac{7}{2}$

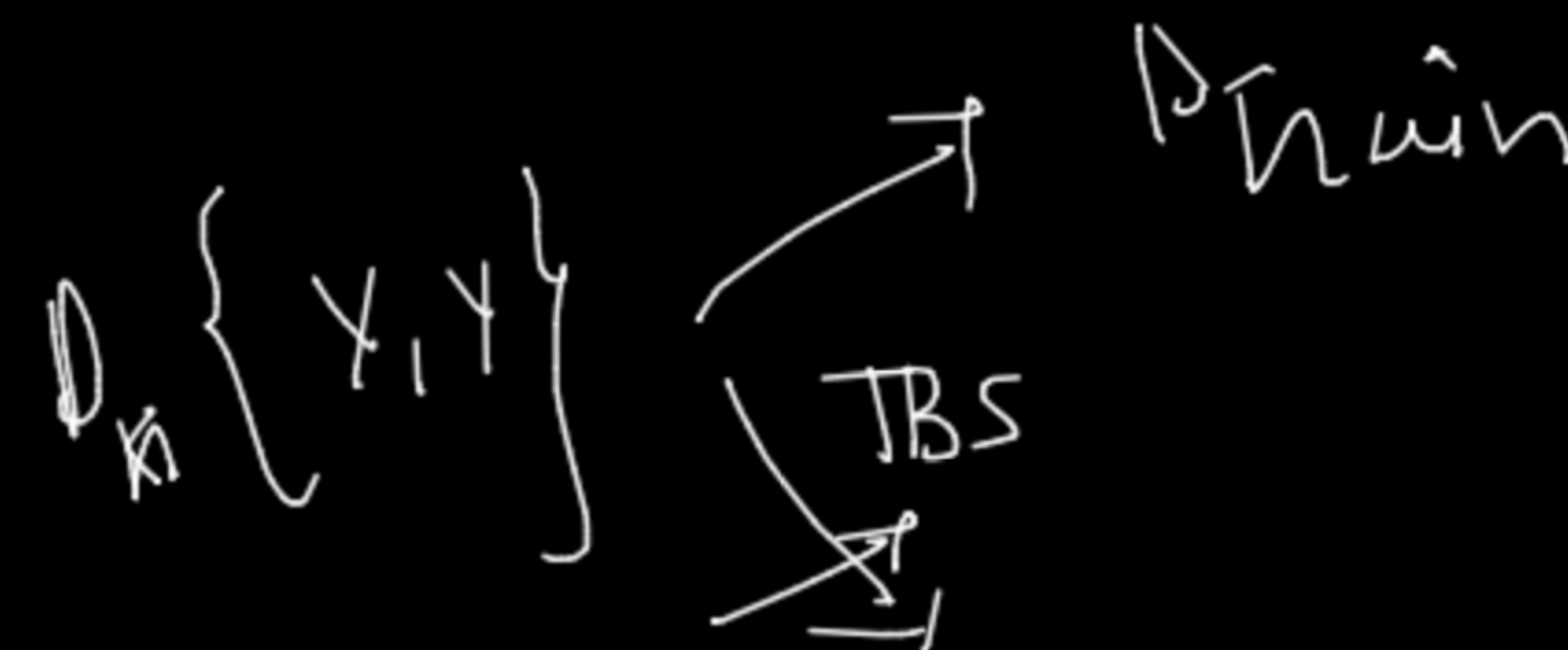
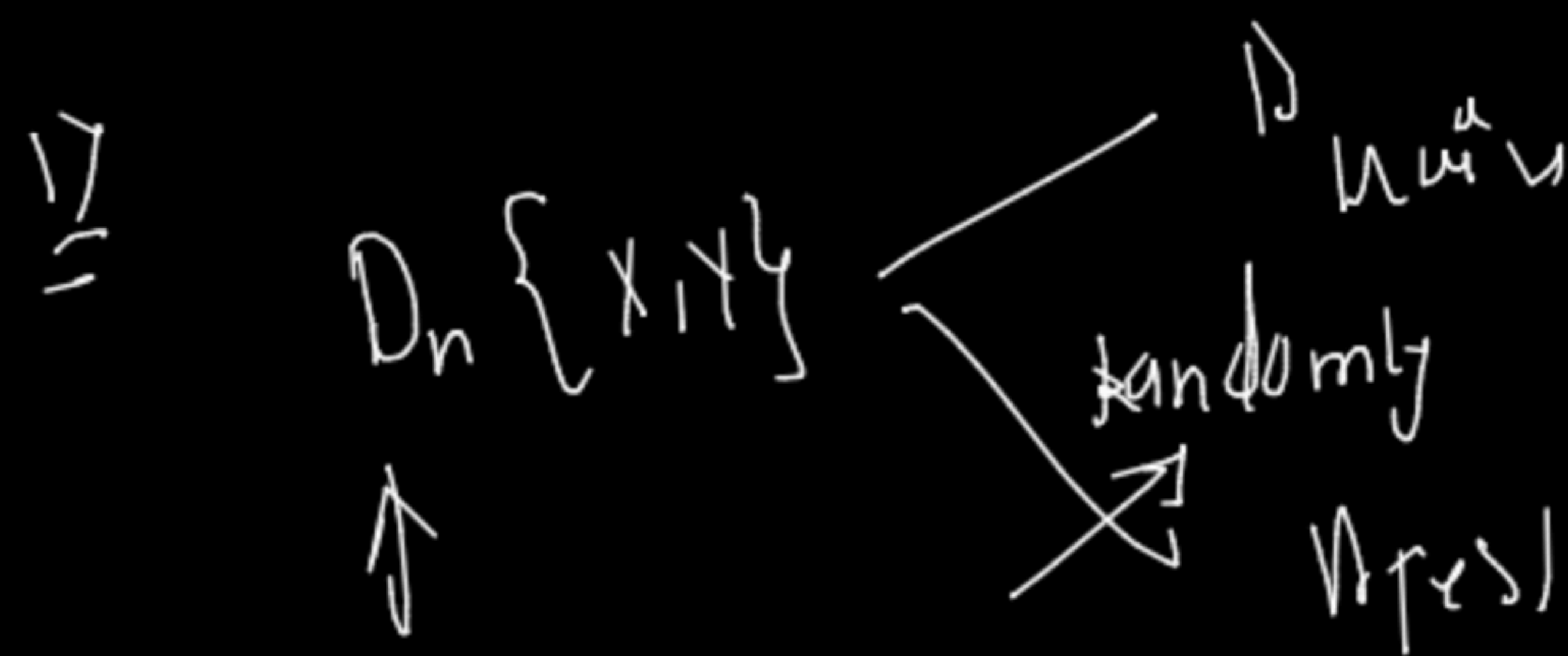
$$D_n = \{X, y \mid x \in R^d, y_i \in [0, 1, 2, \dots, L]\}$$

KNN \rightarrow 'vs Rest'

let T-KNN \rightarrow

$\begin{matrix} 1 & 2 & 4 & \dots & C \\ \hline 1 & 6 & & & \end{matrix}$

Train test set difference



eg. sales

Day	discount	BMI	Volume	
1	-	-	-	2 months 60 days
2	-	-	-	
3	-	-	-	
1	-	-	-	
<hr/>				
1	-	-	-	<u>60 days</u> → test
	-	-	-	
	-	-	-	
	-	-	-	

120

Training

$D_{\text{train}}(x_{\text{train}}, y_{\text{train}}) \rightarrow \text{plot it}$

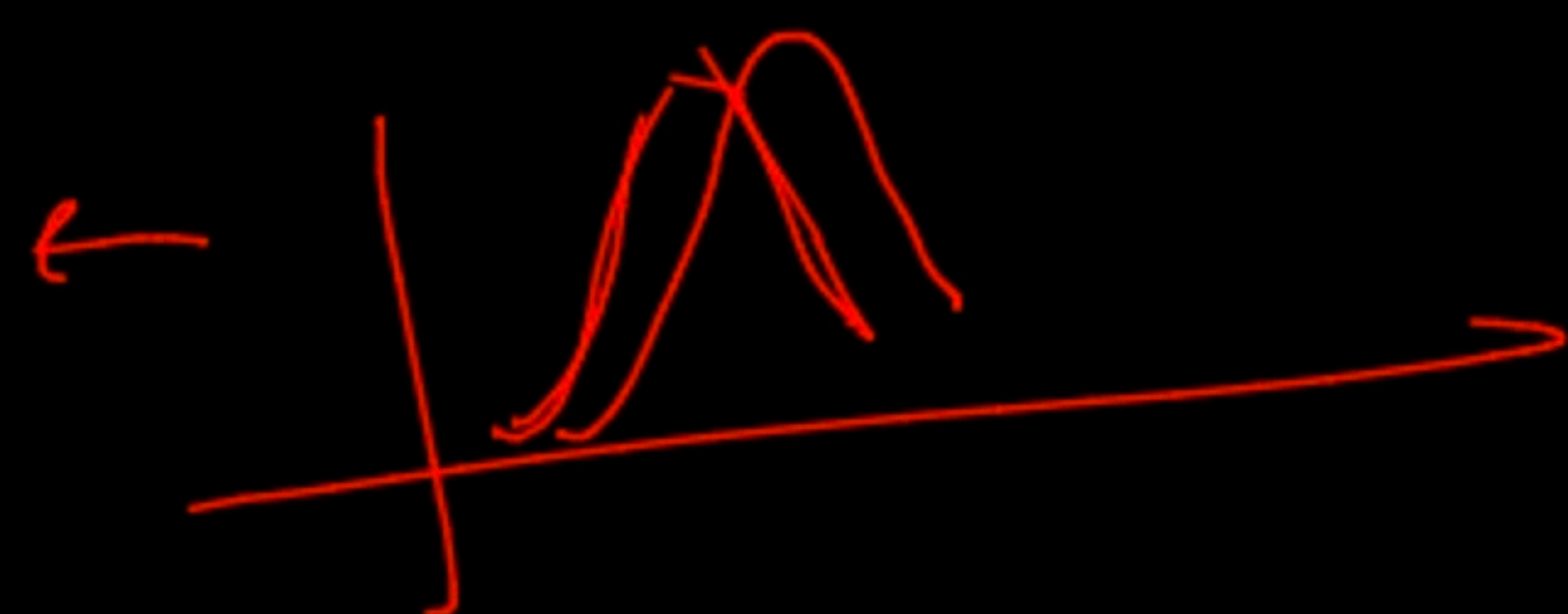
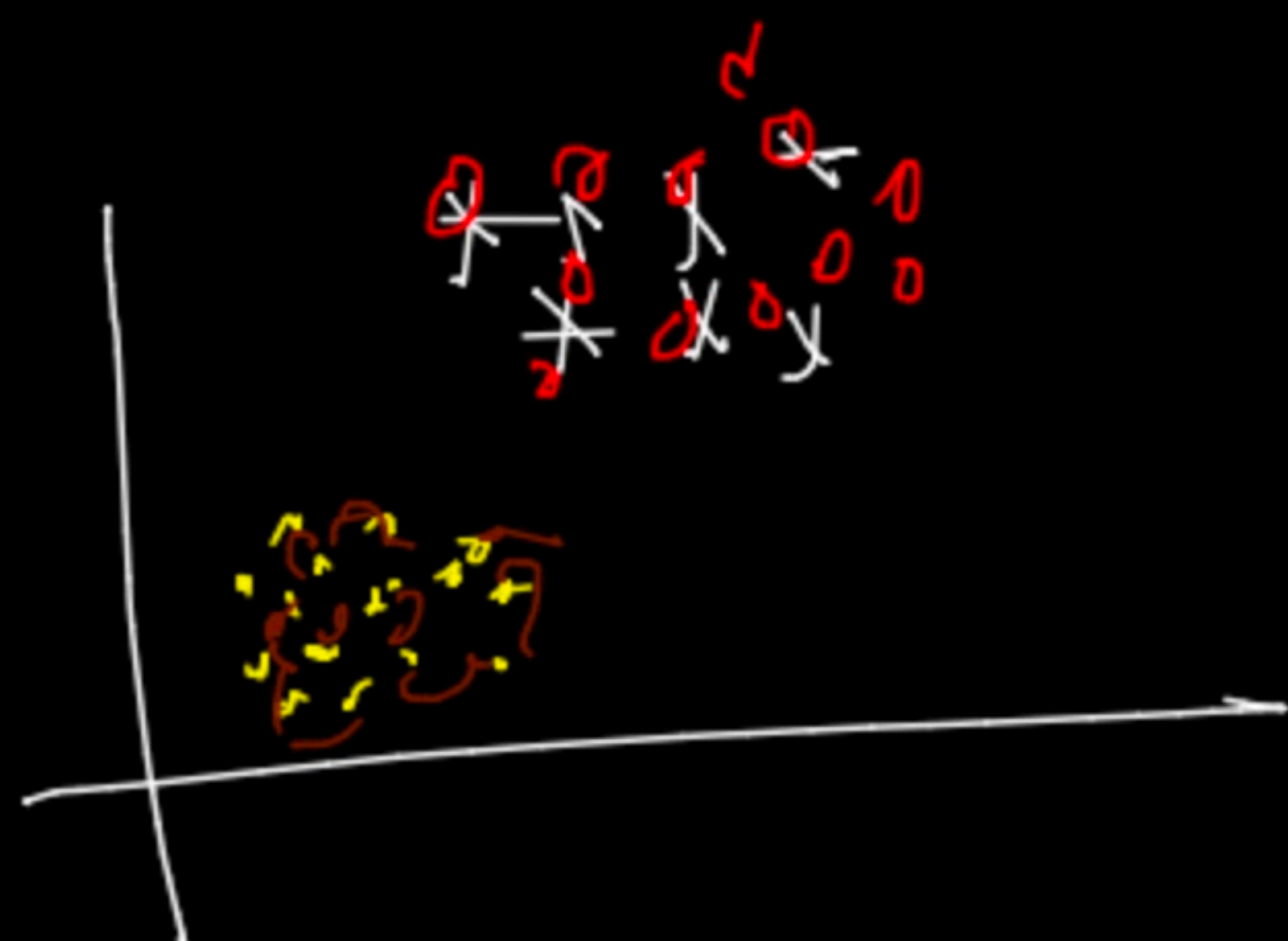
$\begin{cases} +ve \\ -ve \end{cases} \rightarrow \begin{matrix} 40\% \\ 60\% \end{matrix}$

$\{x_{\text{test}}, y_{\text{test}}\}$

D_{test}

val/weights

good

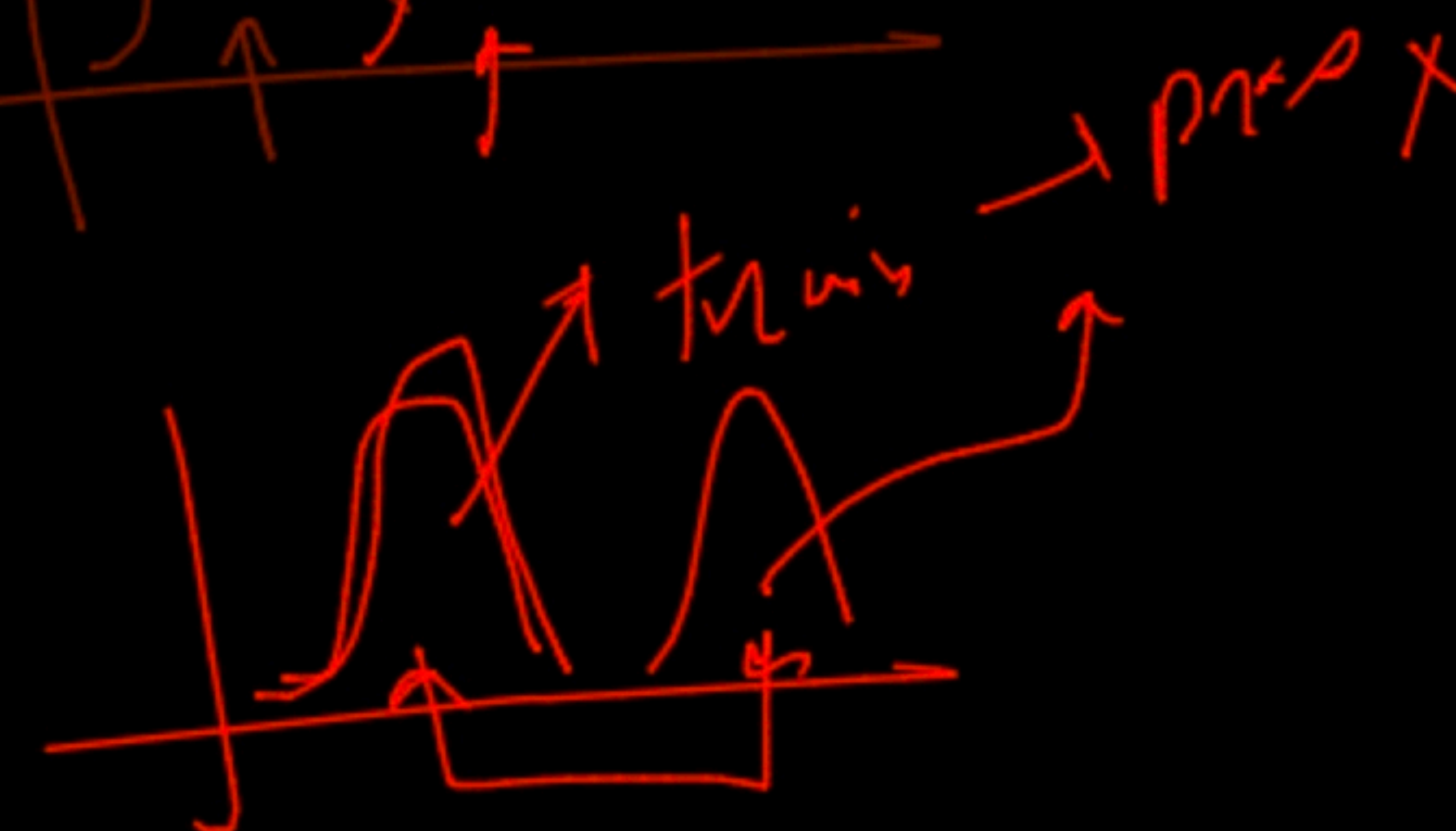
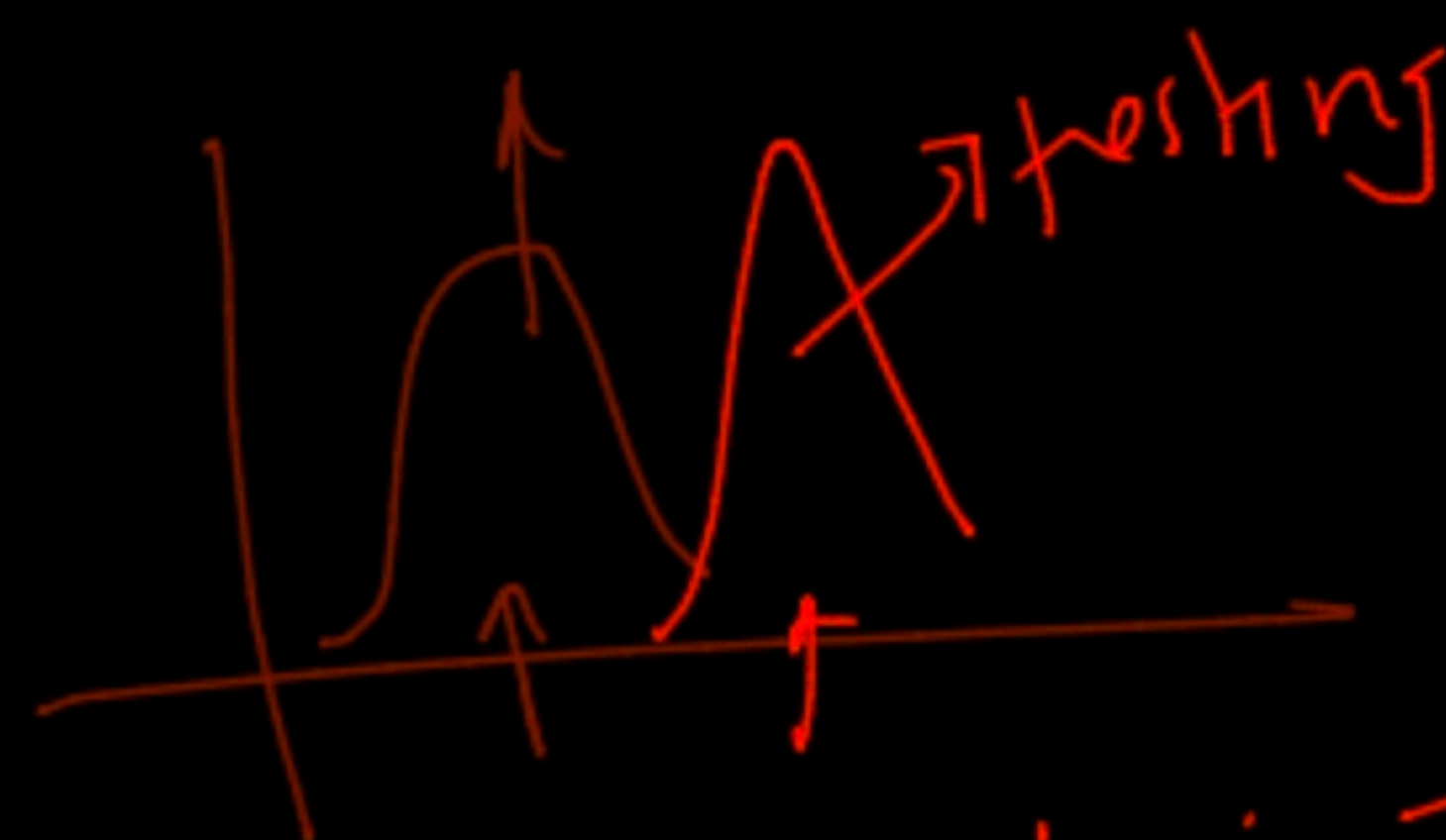
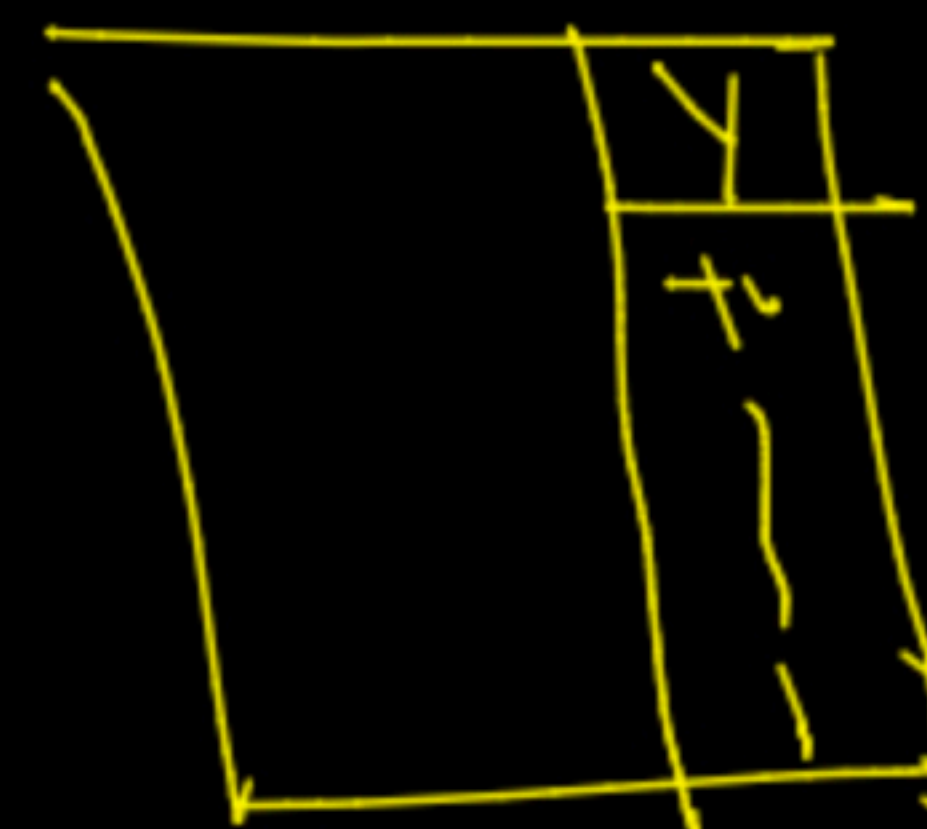


\sqrt{x} \uparrow $(+ve/-ve)$

Bad

x

train U
data



Impact of scale and column standardization

	f_1	f_2
x_1	—	
x_i		
x_j	—	
x_k		

3-5
0

(0-100) (0-1)
uniform → same units

KNN

test

data

f_1 f_2
 x_1 [23, 0.2]
 x_2 [28, 0.2]
 x_3 [23, 1.0]

0.8
0.1
=



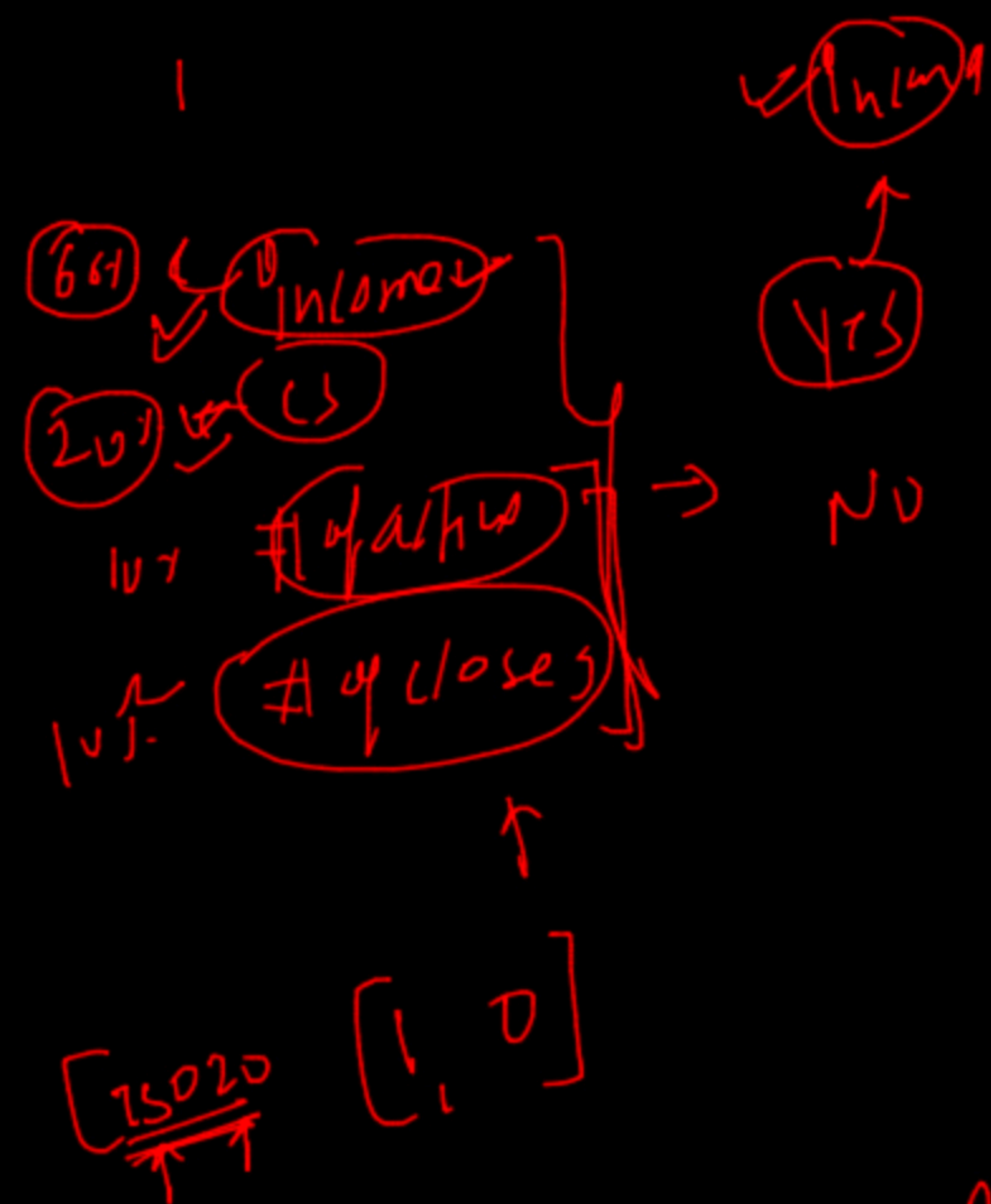
distance (x_1, x_2) = 5

(0-1)
↑

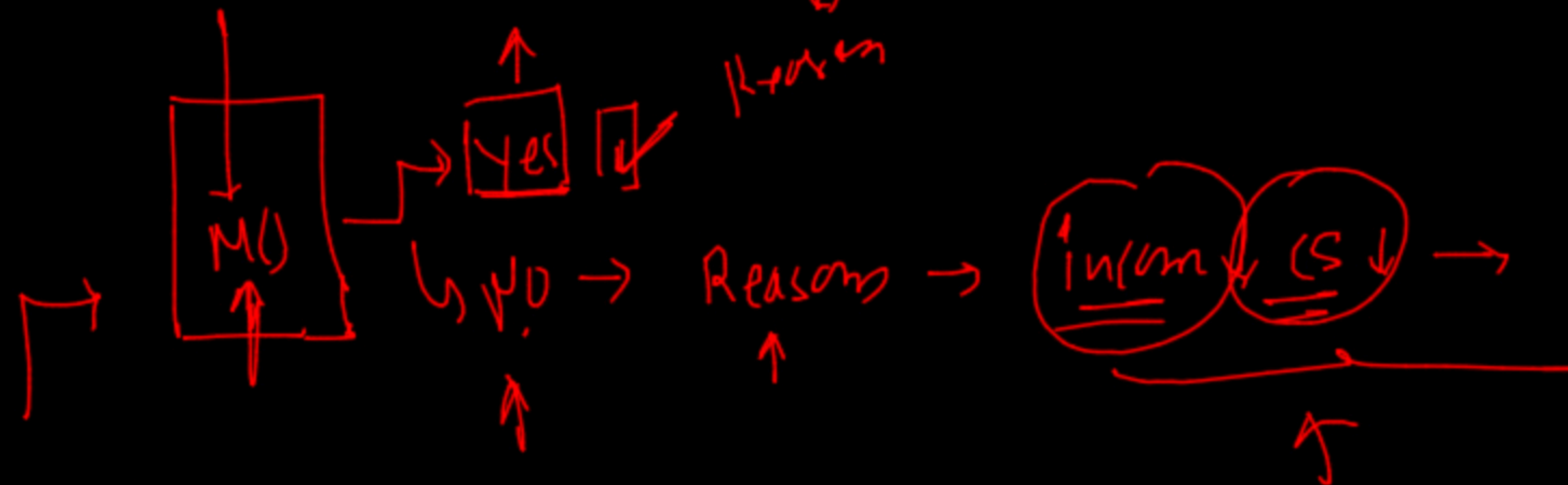
distance (x_1, x_3) = 0.8

logically →

→ model interpretability and Black Box

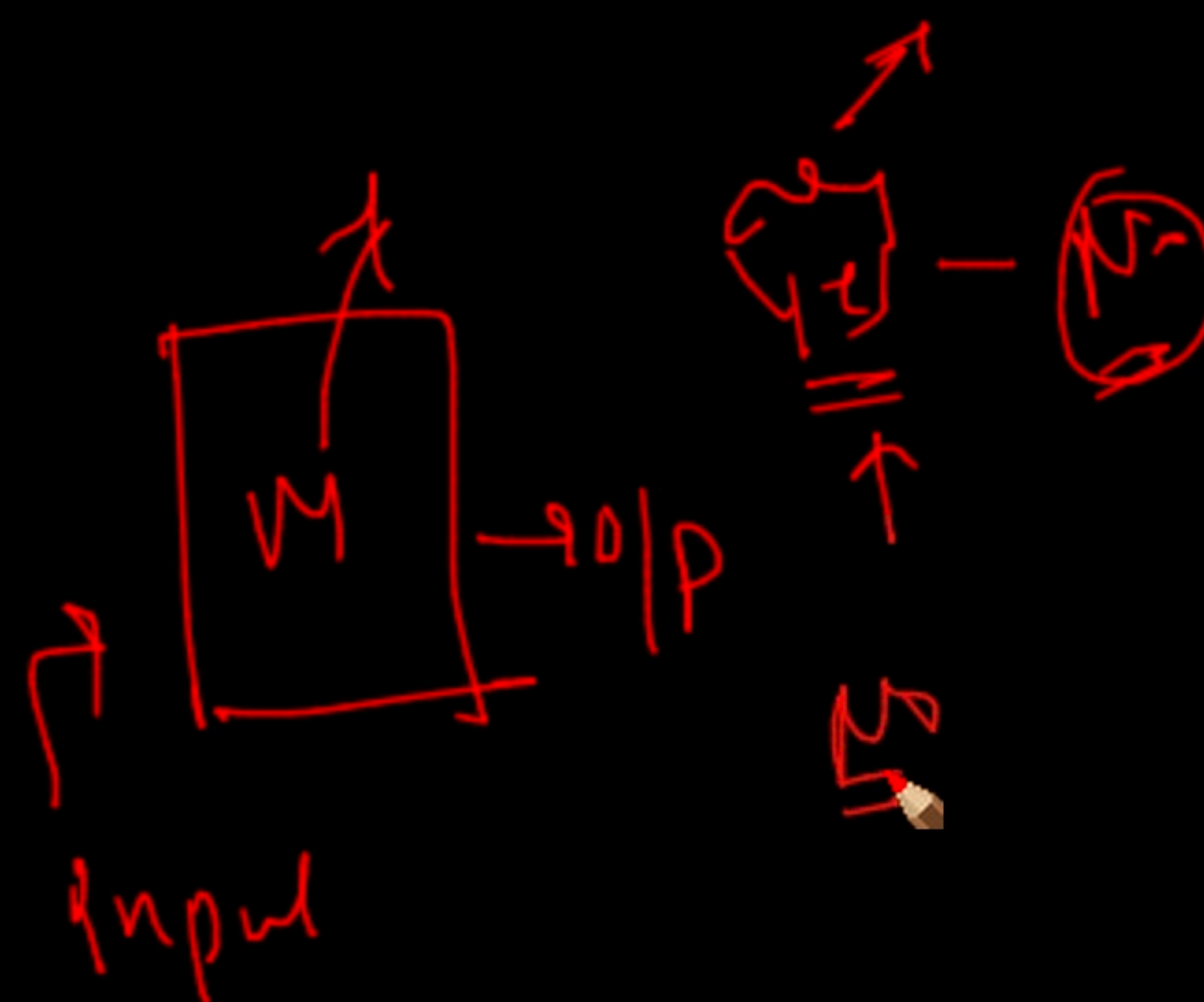


Black Box
 ↑
 'KNN' → $\begin{pmatrix} C1 \\ C2 \\ C3 \end{pmatrix} \rightarrow \text{Yes}$ → Income → Banking / CS



Input
 ↑
 $(\text{Income}, \text{CS}, \# \text{ of active-len}, \# \text{ of closes})$

KNN → $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$



Interview
 ↑
KNN → Hard

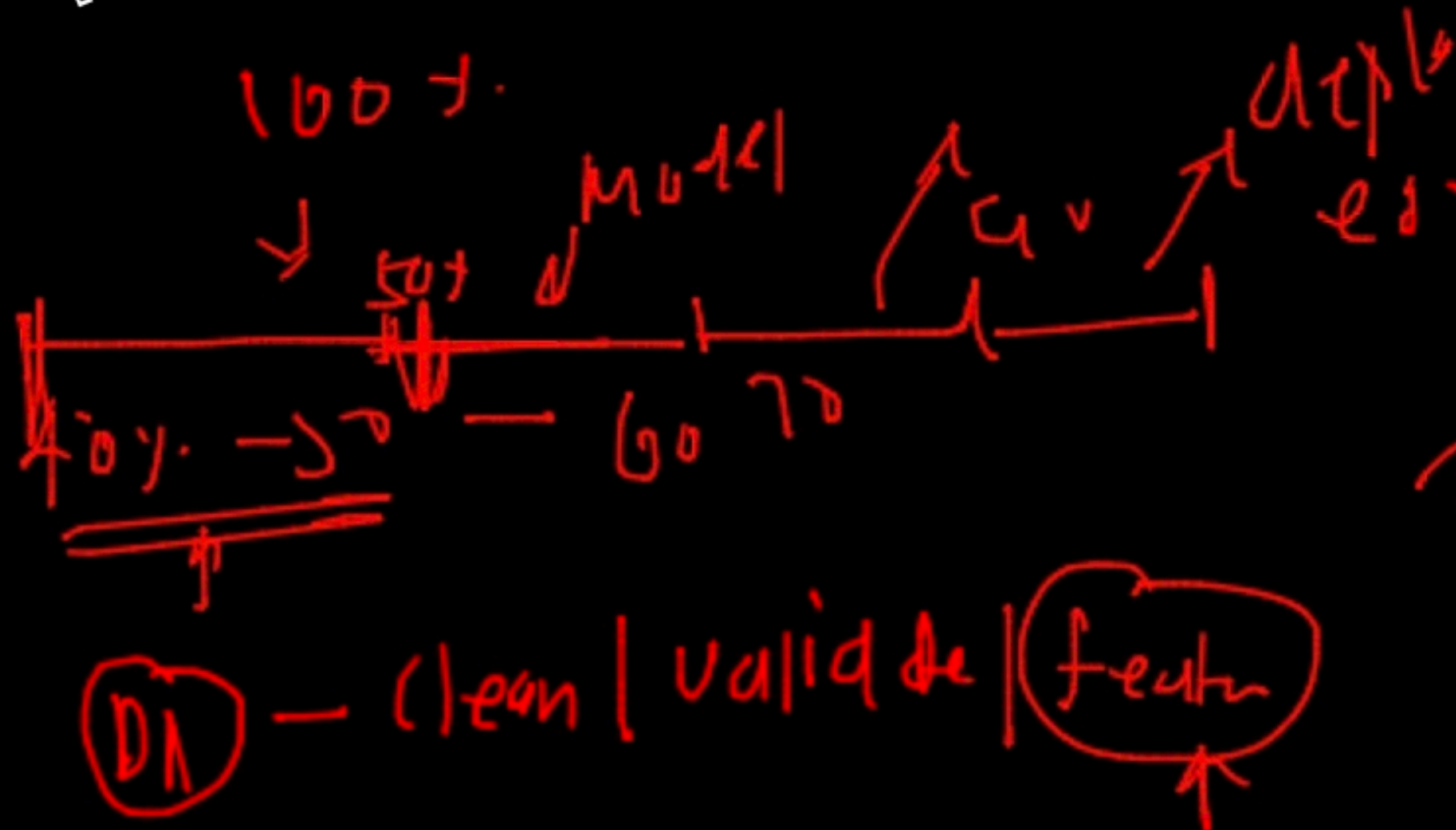
Feature Importance and Forward selection

$D_{tr}(X, Y)$ \rightarrow [yes, no]
 \uparrow

CRD
F

100 feature / input $\rightarrow M()$

Data analysis \rightarrow 50 important

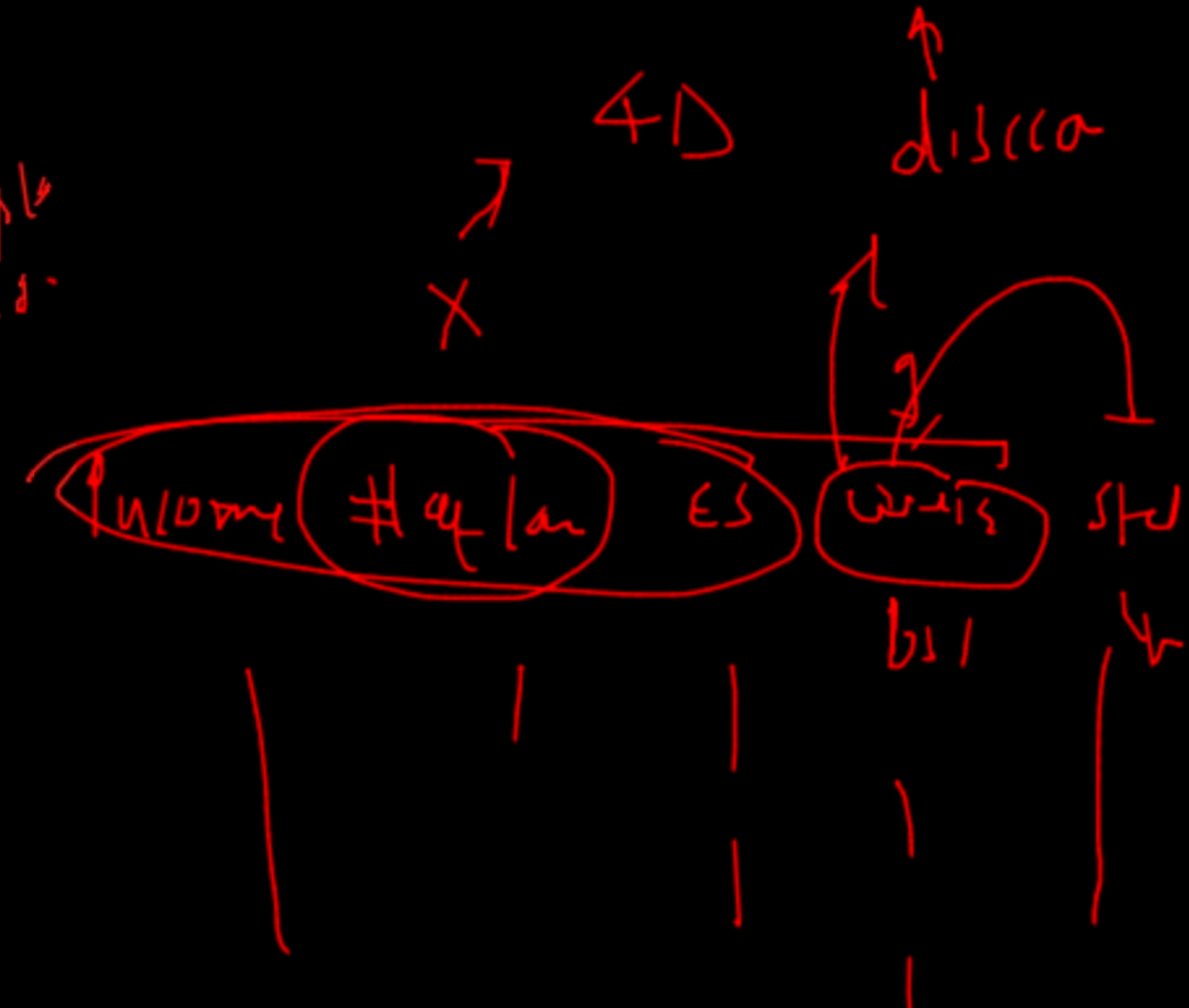


$KNN(X, Y)$
 \uparrow

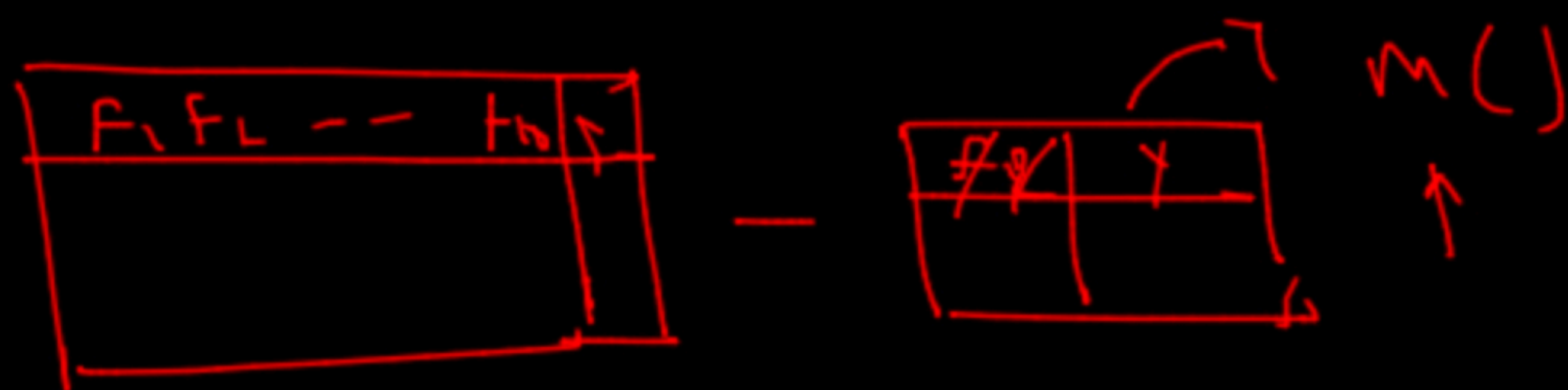
\rightarrow Data analysis
 \uparrow



$KNN(X, Y)$
 \uparrow



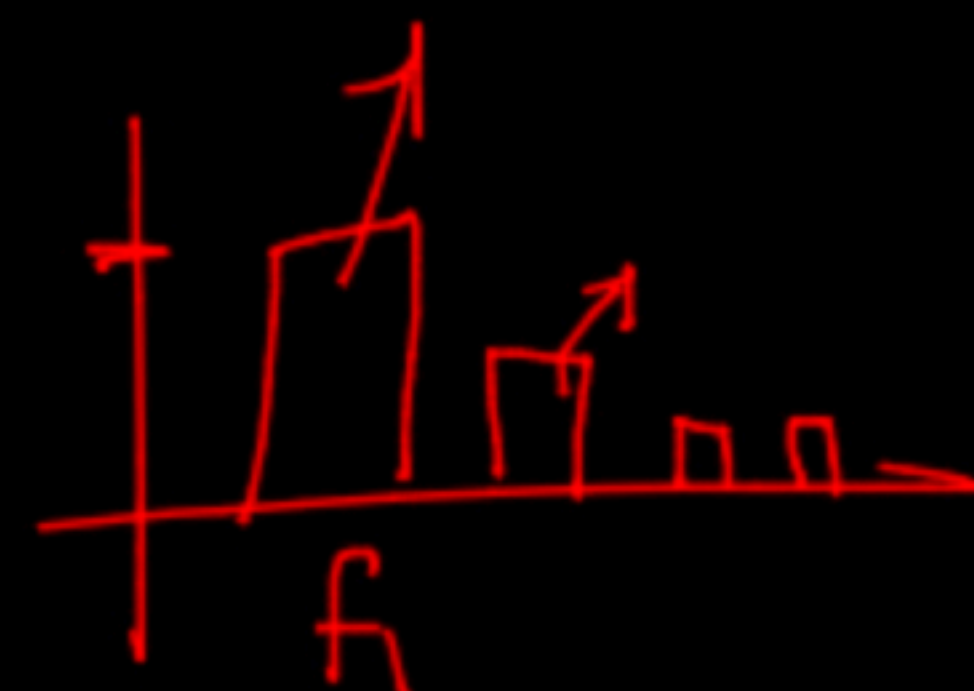
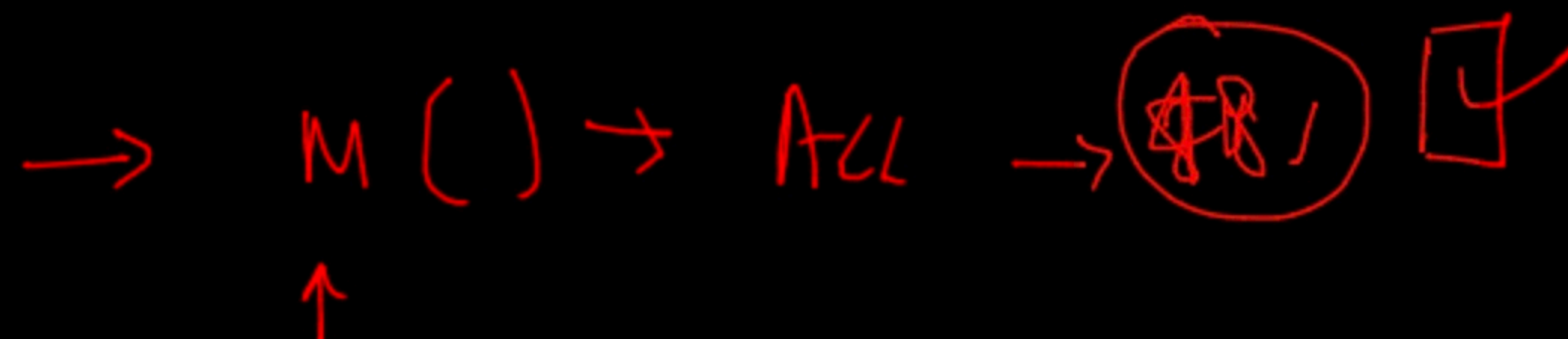
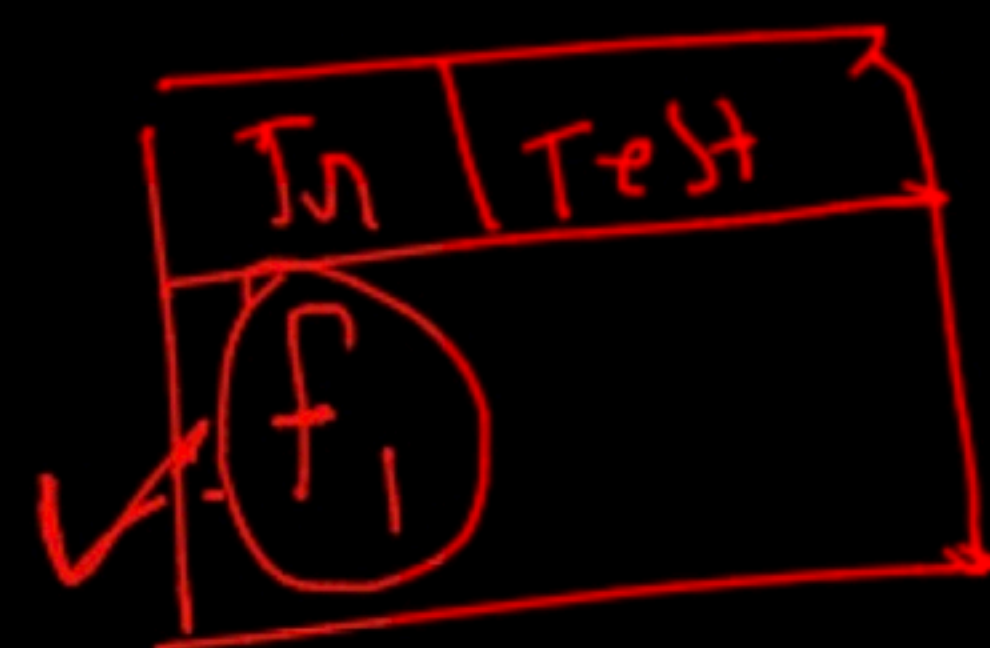
$f_1, f_2, f_3, \dots, f_{10}$



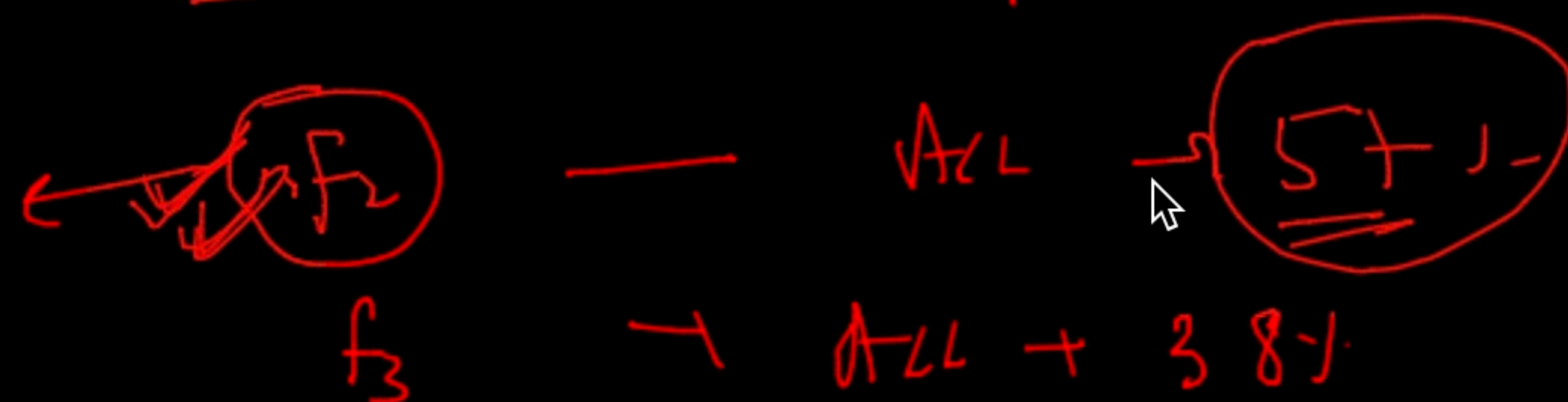
$$D_n = \{T_n, T_r\}$$



Step 1



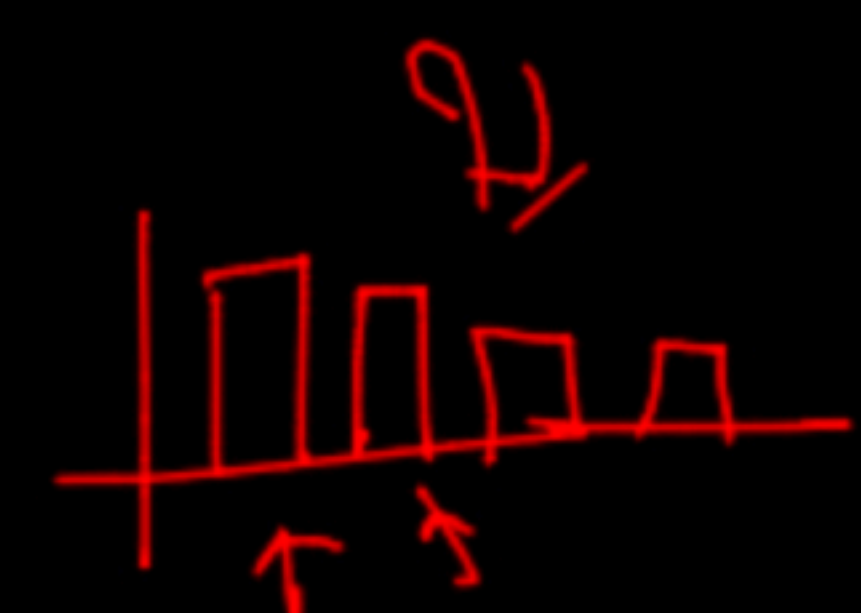
important



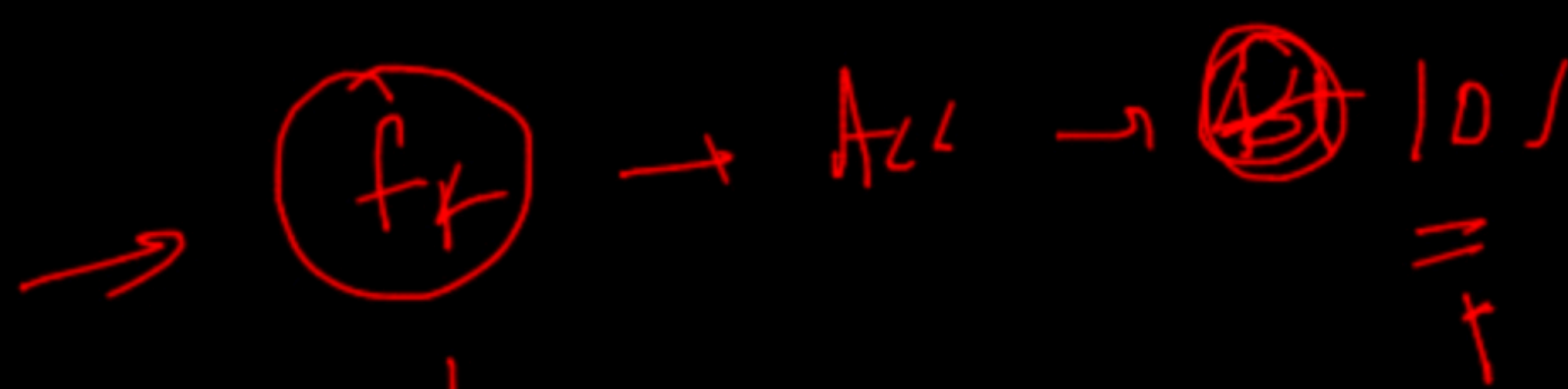
"Feature important"

shape

evidently



least



RFE \rightarrow to select Feature

\uparrow

sklearn

\uparrow

\rightarrow 100 col \rightarrow data analysis

\downarrow

50 col \rightarrow RFE (10 \uparrow) \rightarrow \rightarrow

$D \uparrow \rightarrow \downarrow$

11'

100 info
 \uparrow

$D \uparrow$ $D \downarrow$
 \uparrow
100 200

$D \uparrow$
 \uparrow
info

Tr	res
f_1	f_2

$\rightarrow M() \rightarrow 64\%$

f_1	f_2	f_3
-------	-------	-----------------------------

$\rightarrow M() \rightarrow 58\% \downarrow$

f_1	f_2	f_3
-------	-------	-------

$\rightarrow M() \rightarrow 74\%$

