

Date Pre-processing

- Handling missing value - Measure of central tendency
- Handling categorical variable - Encoding concept
- Handling Outliers - Transformation approach
- Scaling the data - Normalization / Standardization
- Handling Imbalanced dataset - Imbalance
 - ↳ applicable only with classification Problem

RNN Implementation & Interactive-MICE

Univariate Analysis

* Imbalanced vs Balanced dataset

2 class

(+)ve = Yes
(-)ve = No

$$n = 1000$$

$$\begin{aligned} n_1 &= (+)ve \\ n_2 &= (-)ve \end{aligned}$$

$$n = n_1 + n_2$$

Case I :- if $n_1 \approx n_2$ = Balance dataset

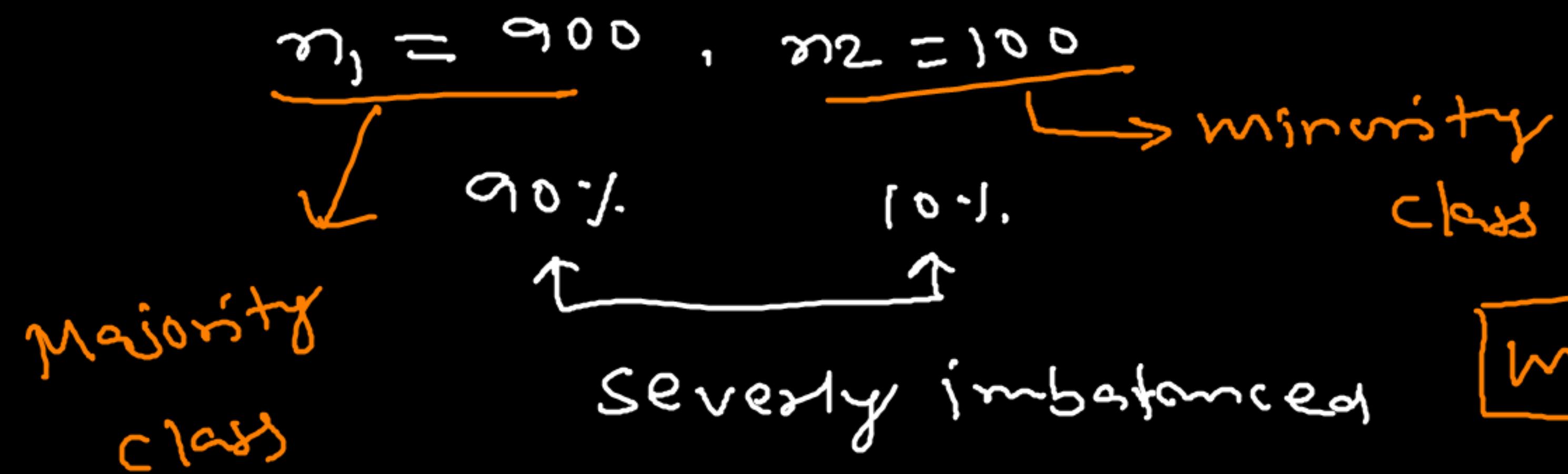
$$n_1 = 500, n_2 = 500$$

$$n_1 = n_2$$

$$\begin{aligned} n_1 &= 580, n_2 = 420 \\ &\quad \boxed{n_1 \approx n_2} \quad \underline{n_1 \neq n_2} \end{aligned}$$

58% - 42%.

Case II :- If $n_1 \ll n_2$ or $n_1 \gg n_2 \rightarrow$ [Imbalanced dataset]



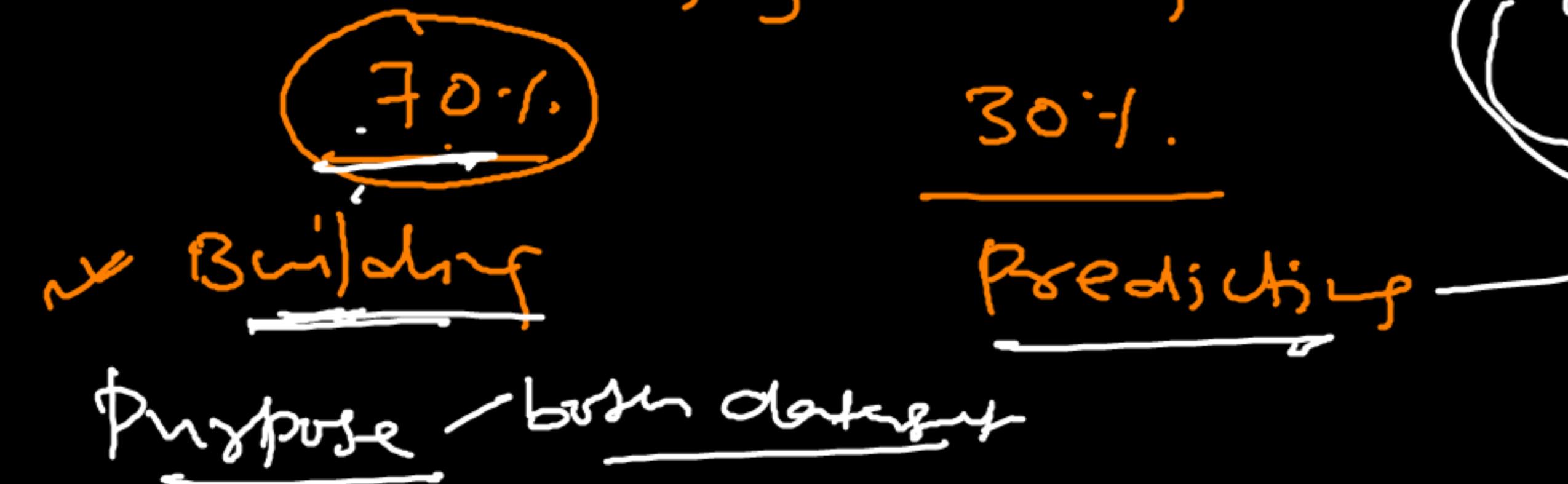
Minority * 2 \leq Majority

$$\frac{200}{2} < 900$$

$n_1 = 200, n_2 = 800$

0 class - Real form = $\sqrt{284315} = n_1 = 99.7\% \times$
1 class - found = 498 $= n_2 = 0.01\%$

Split the data into training & testing



$$\begin{aligned}n_1 &= \underline{50\%} \\n_2 &= \underline{50\%}\end{aligned}$$

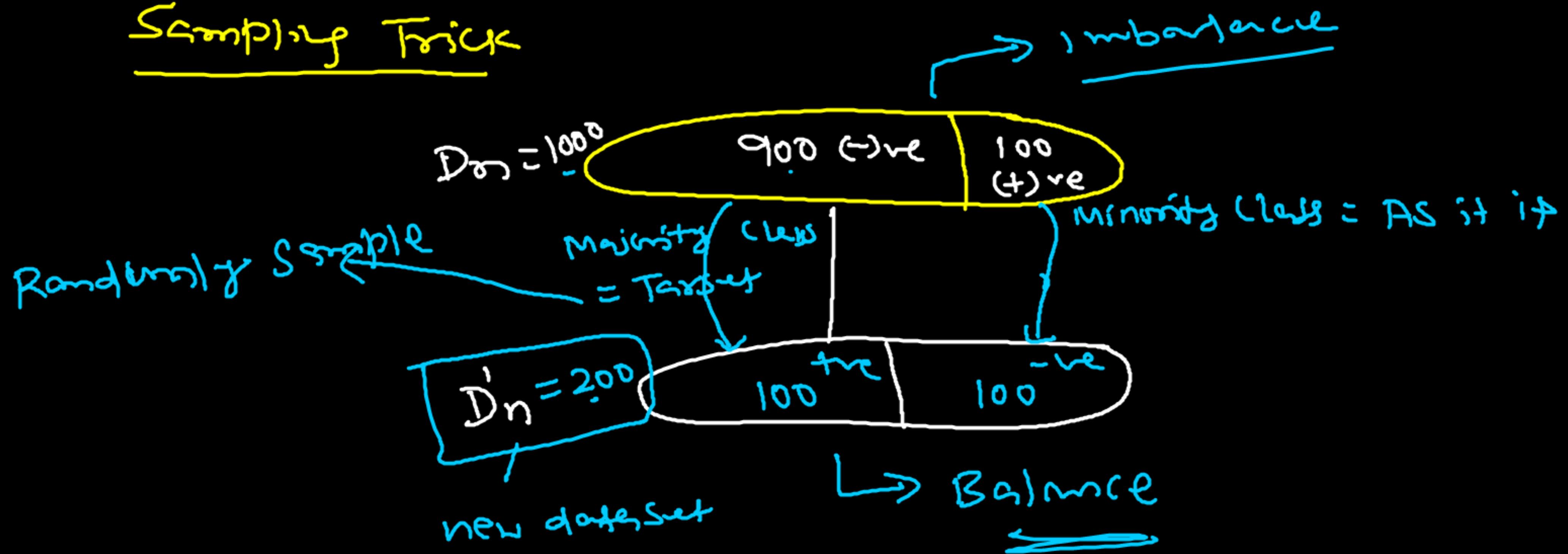
* How to work around imbalance dataset issue

- ① undersampling method → Nearest
- ② oversampling method.
- ③ Artificial/synthetic technique (SMOTE)
- ④ Class weight - Techniques

* undersampling method - Target (Majority class)

D_n $n_1 = 100 (+)ve$ — Minority Class
 $n_2 = 900 (-)ve$ — Majority Class

Sampling Trick



Problem with undersampling

$D_n = 1000 \rightarrow$ undersampling $\rightarrow 200 \rightarrow$ what about 800 datapoint

20% - Model builders

$$|D'_n| \ll |D_n|$$

80% \rightarrow Throwing away

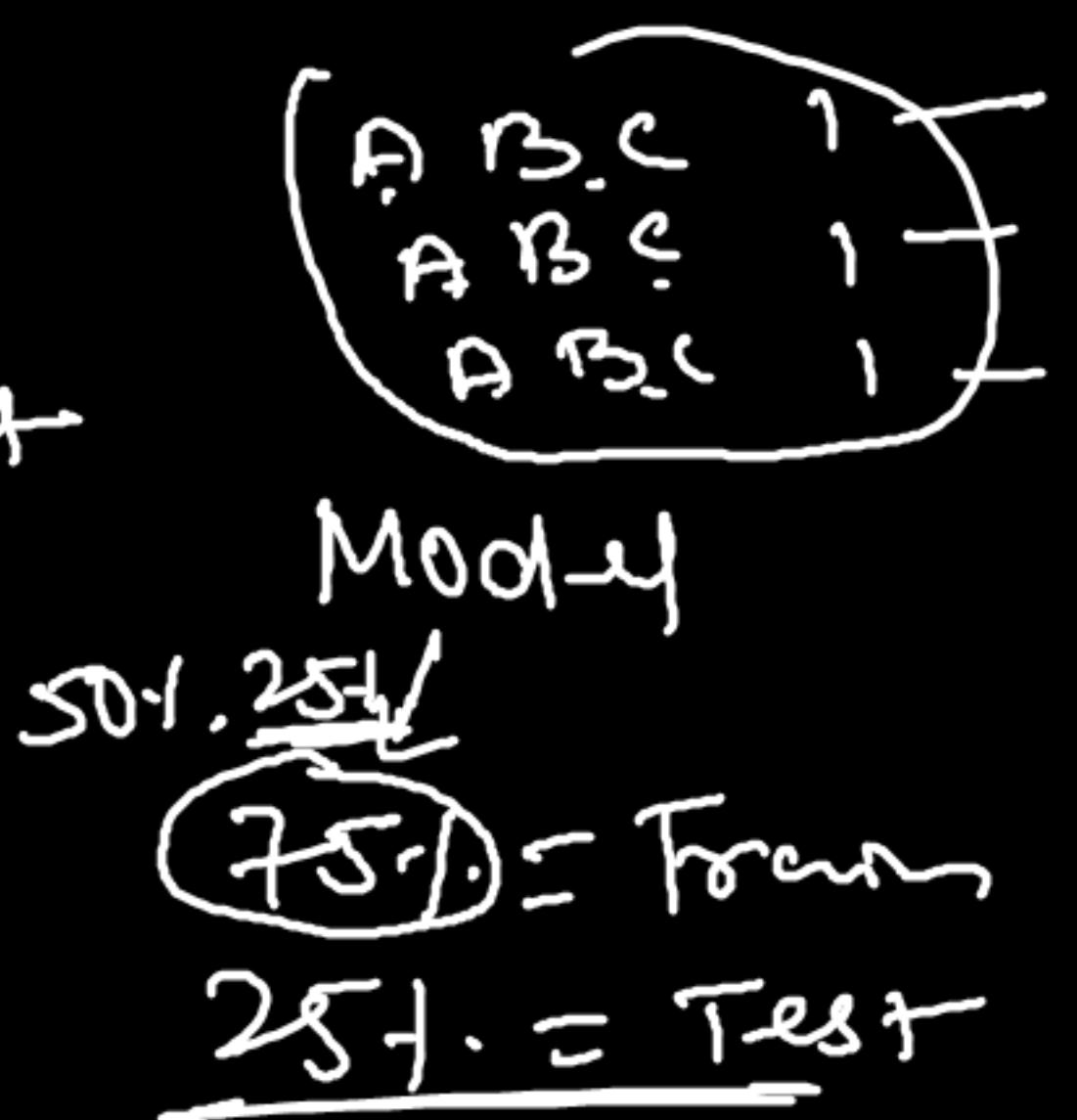
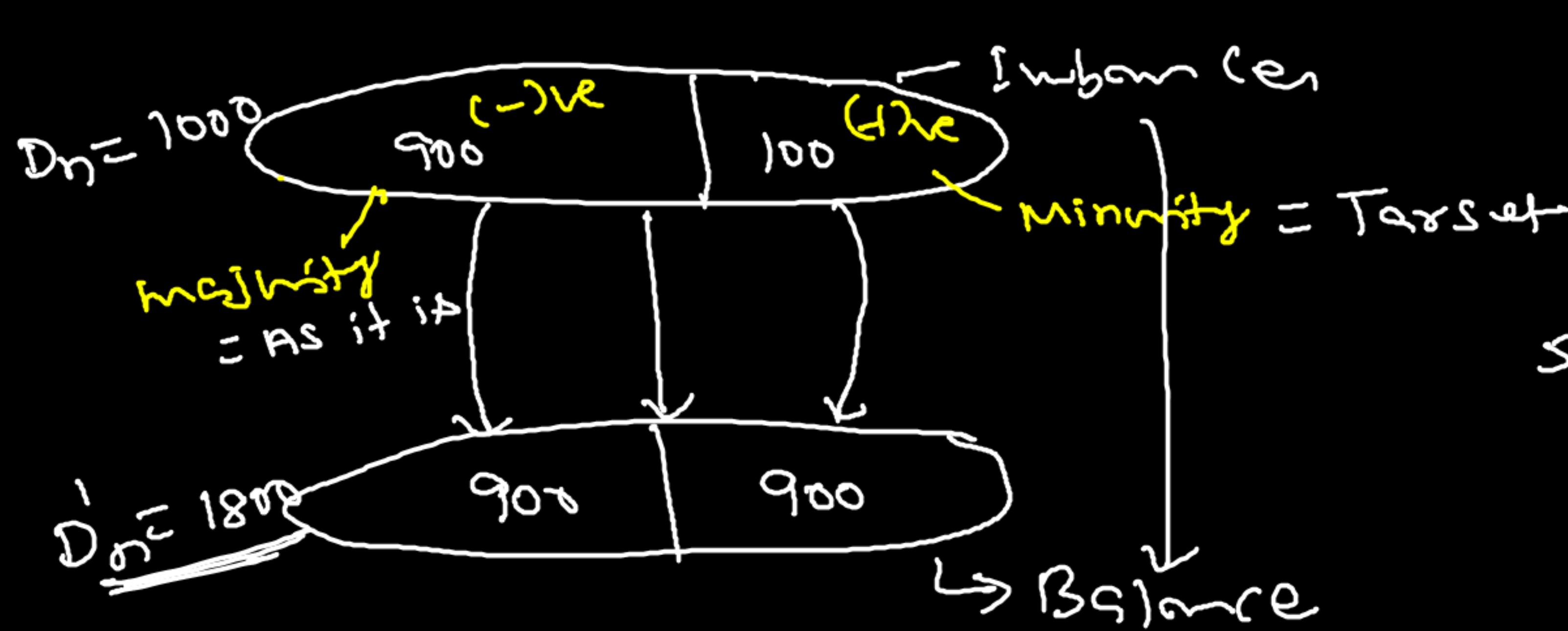
\hookrightarrow Throwing away data

\rightarrow Should be avoided

\rightarrow Not good idea to remove 80% data

* Oversampling method = Target (Minority) $\frac{100}{100}$ ① - very very less
② - never $\underline{\underline{=}}$

D_n $n_1 = 100 = 10\%$.
 $n_2 = 900 = 90\%$. \rightarrow Enhance - how to handle it ?



Oversampling

→ Placing more points from the minority class
of the dataset



Advanced Techniques → Not used very frequently

→ extrapolation → synthesis | synthetic approach
SMOTE

You are screen sharing Stop Share

Class weight

100 +ve 900 -ve

$$Y = mx + c$$

Slope/co-effici
weight

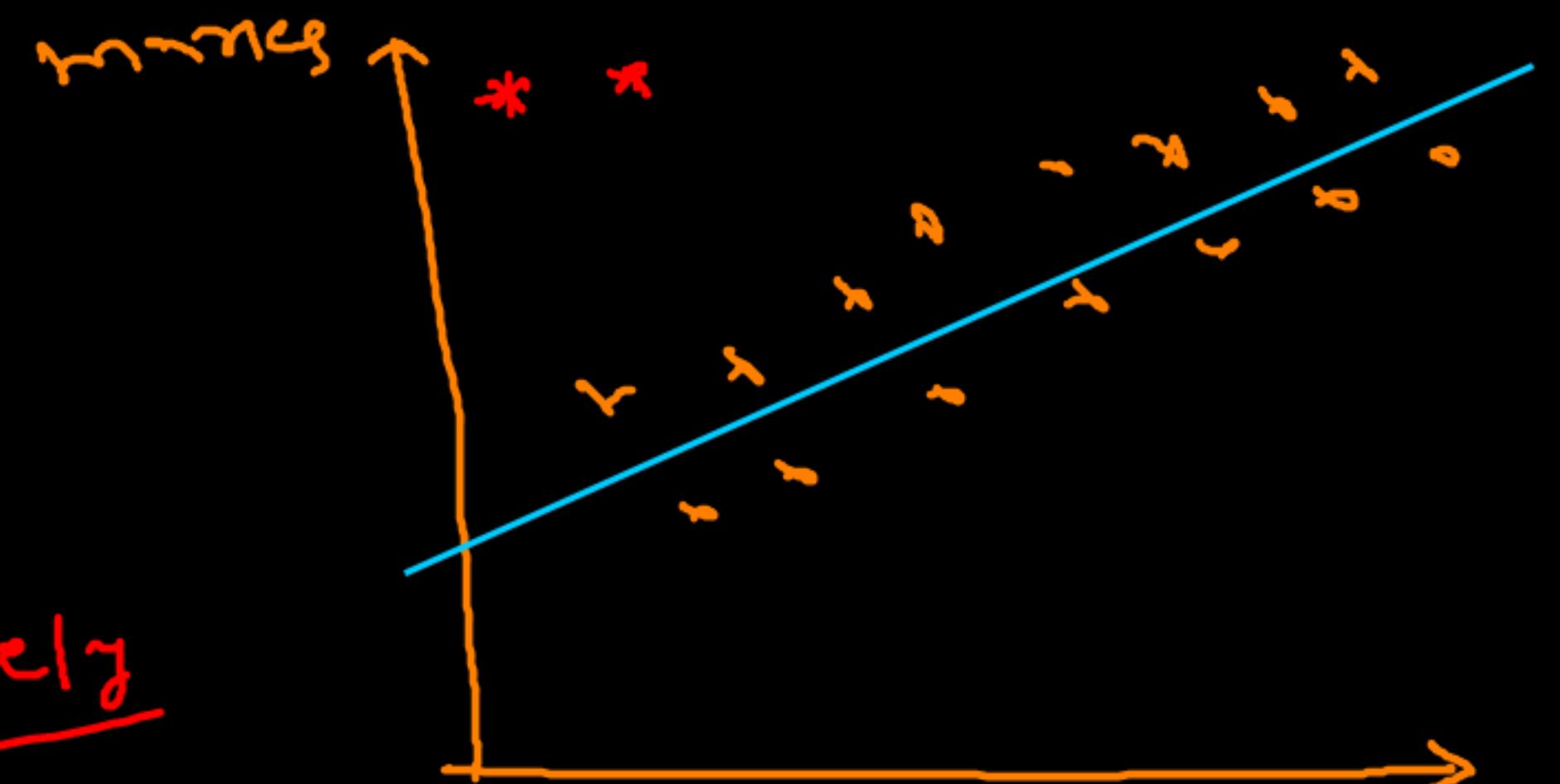
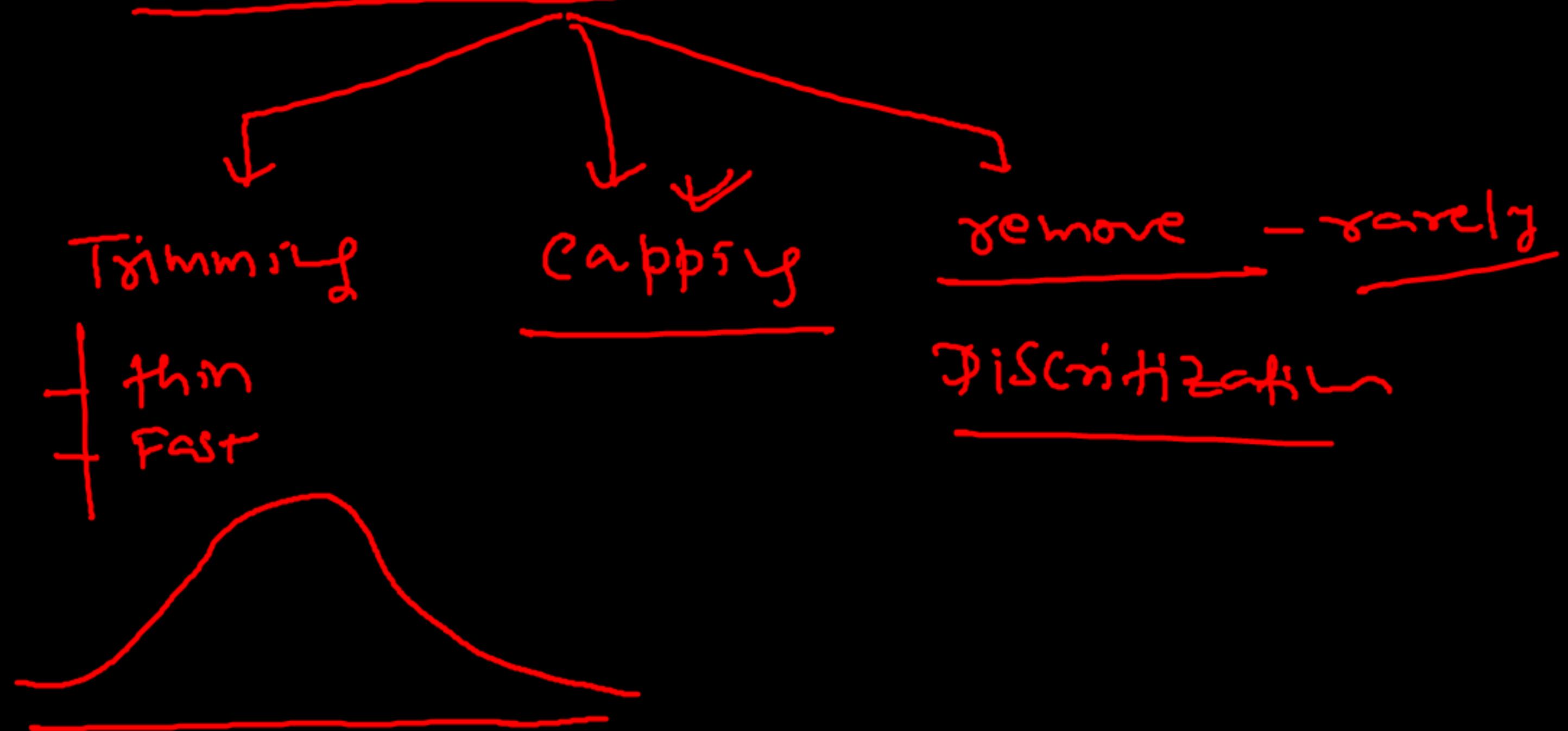
$w_+ = 9 \rightarrow$ more weight in the minority class

$w_- = 1 \rightarrow$ less " " " " majority "

AI

Outlier Treatment

How to handle outlier



$$(\text{+ve outlier}) = Q_3 + 1.5 \times IQR - \text{isn}$$

$$(\text{-ve outlier}) = Q_1 - 1.5 \times IQR - \text{isn}$$

Normal Distribution

Empirical Rule

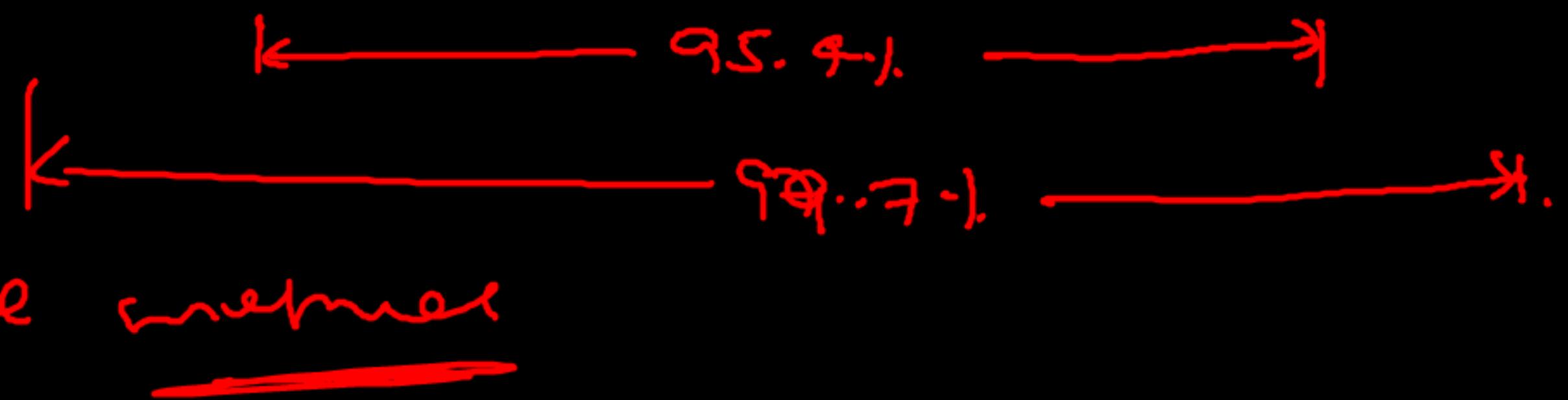
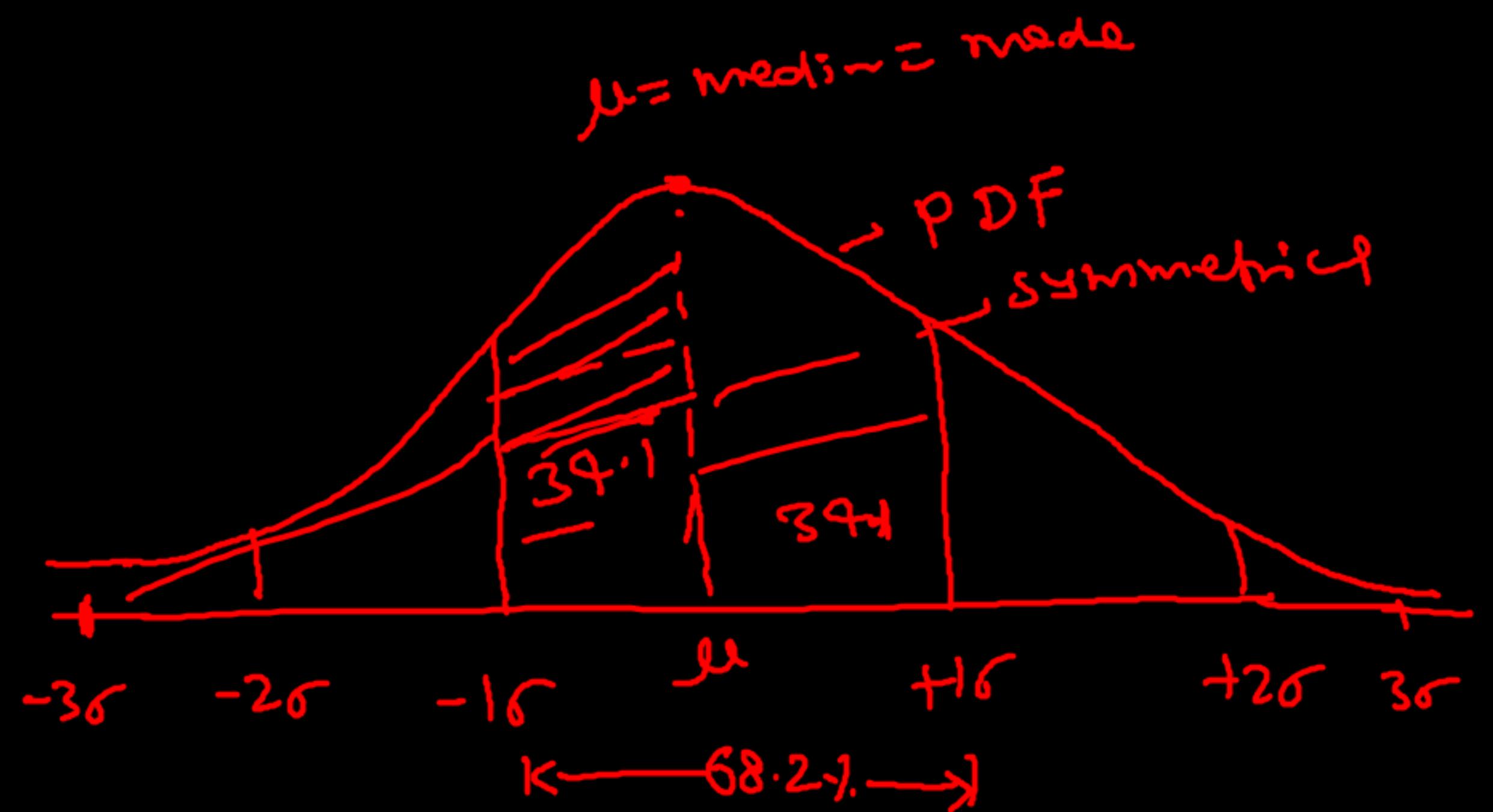
$$\text{Prob}[-\mu \leq x \leq \mu] \approx 68.2\%$$

$$\text{Prob}[-2\mu \leq x \leq 2\mu] \approx 95.4\%$$

$$\text{Prob}[-3\mu \leq x \leq 3\mu] \approx 99.7\%$$

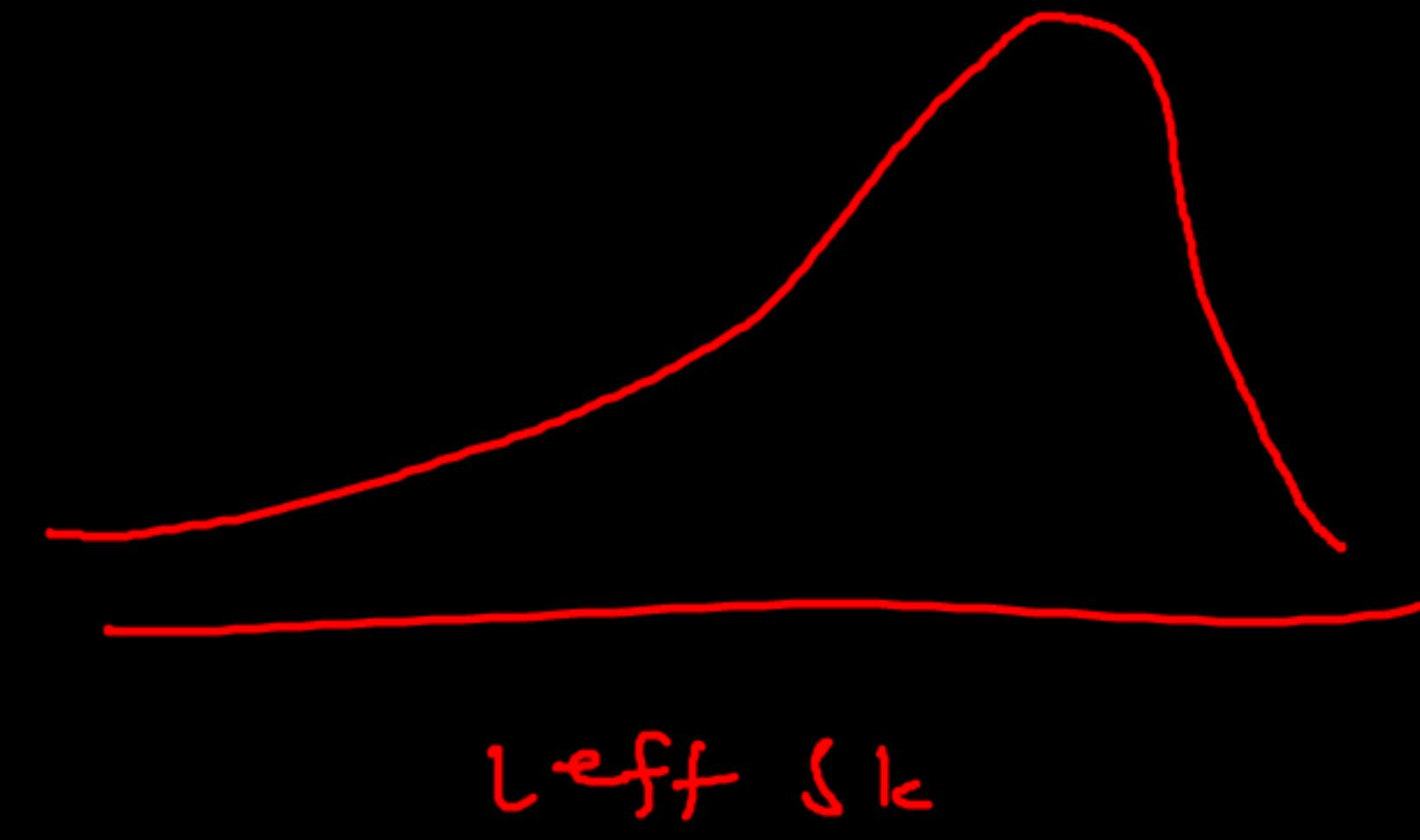
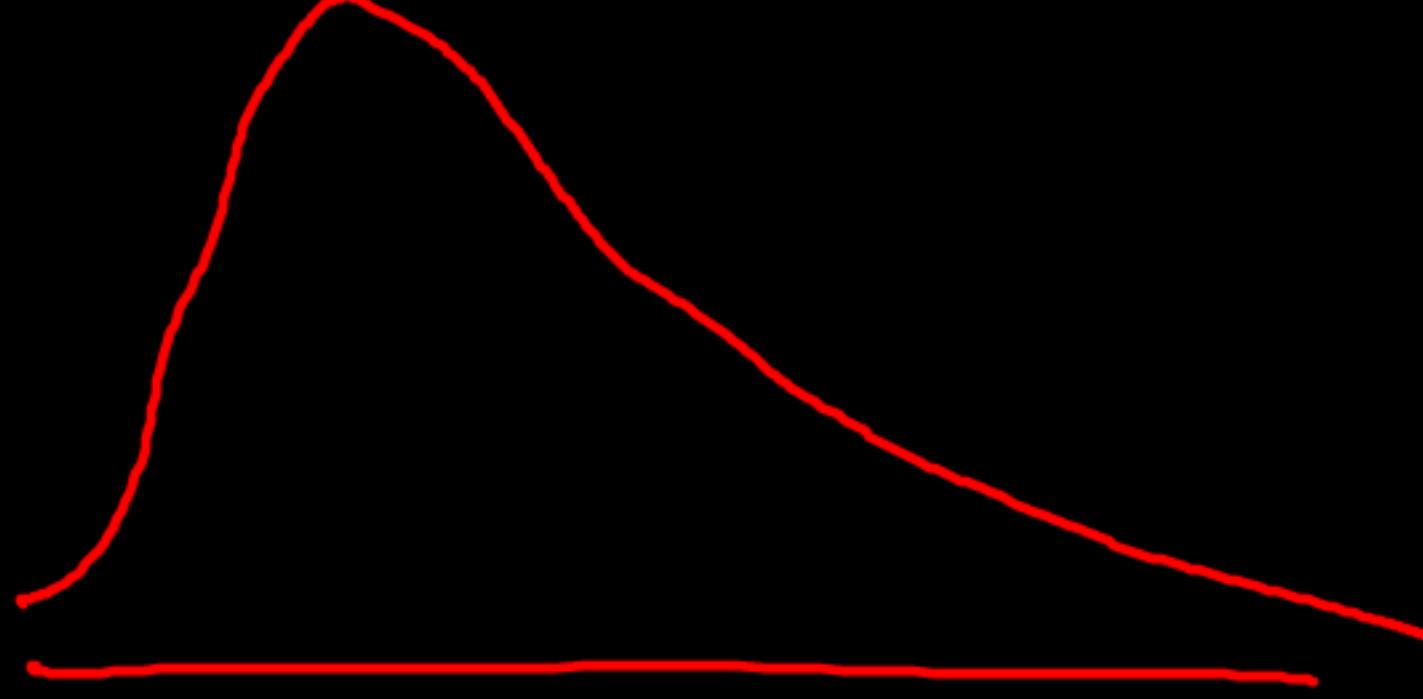
$(\mu + 3\sigma) >$ outliers

$(\mu - 3\sigma) <$ \hookrightarrow Z-Score method



②

Skewed dist - Box plot = IQR



right tail skewed :

