

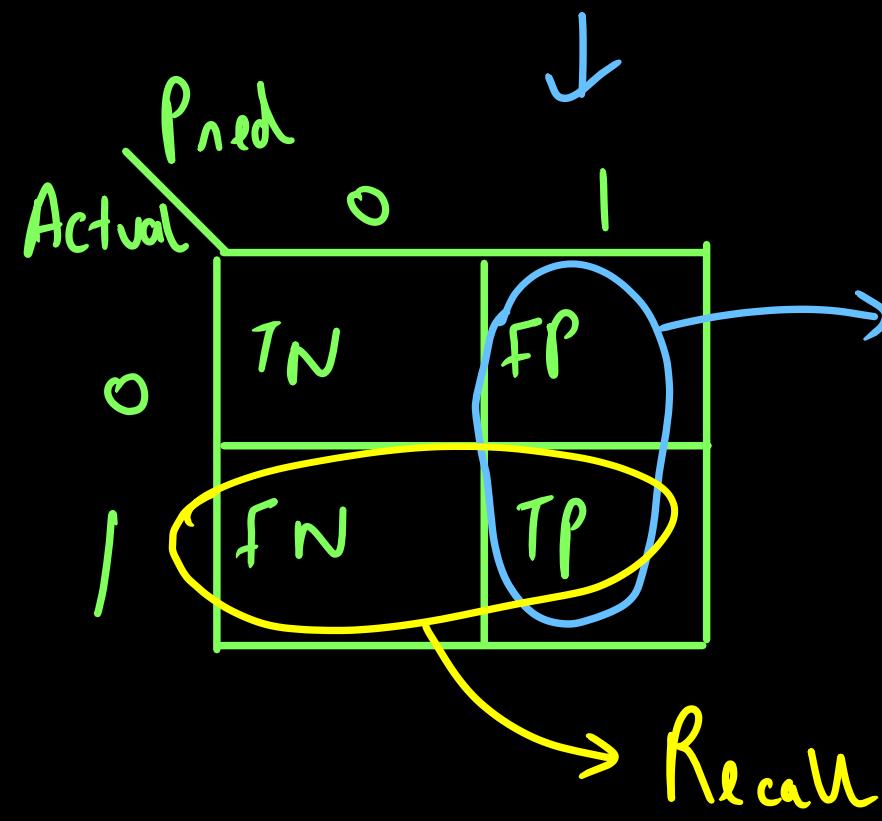
## Agenda

→ TPR, TNR, FPR, FNR

→ ROC AUC curve

→ PR curve

→ Imbalance data

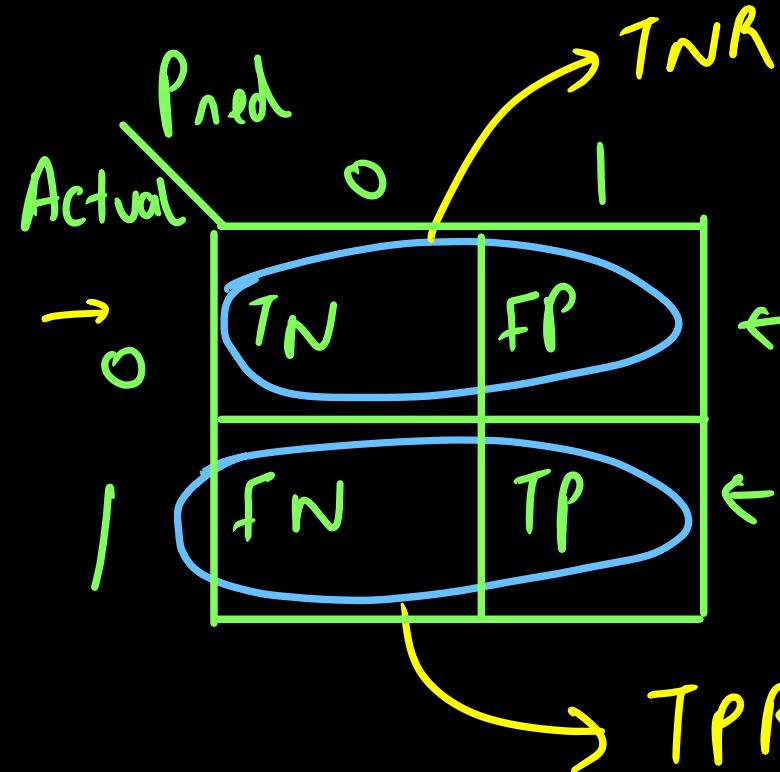


$$\uparrow P_{\text{Rec}} = \frac{TP}{TP + FP}$$

$$\uparrow R_{\text{ec}} = \frac{TP}{TP + FN}$$

$$F1 = \frac{2PR}{P+R}$$

True Positive Rate



$$\uparrow \checkmark TPR = \frac{TP}{\text{No. of Actual} = 1}$$

$$\uparrow \checkmark TNR = \frac{TN}{\text{No. of Actual} = 0}$$

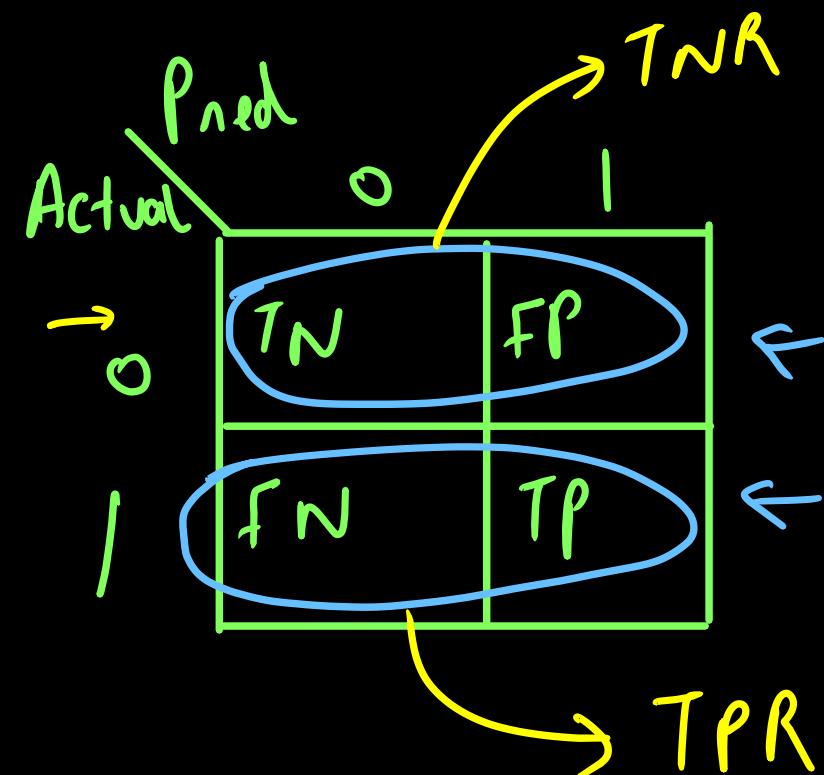
$$= \frac{TP}{TP + FN} = \text{Recall}$$

$$= \frac{TN}{TN + FP} = \text{Sensitivity}$$

$$= \frac{TN}{TN + FP} = \text{Specificity}$$

$\rightarrow FPR \quad FNR \rightarrow$  false Negative Rate

$\hookrightarrow$  false Positive rate



$$\downarrow FPR = \frac{FP}{FP + TN} = 1 - TNR$$

$$\downarrow FNR = \frac{FN}{FN + TP} = 1 - TPR$$

Min log-loss  $P, R, F1$

$$z^{(i)} = \omega^T x^{(i)} + b$$

$$\sigma(z^{(i)})$$

$$\hat{y} = P(y=1/x)$$

$$\hat{y} \geq \tau$$

$\tau = \text{threshold}$   
 $= 0.5$

$$\begin{cases} \hat{y} \geq \tau \rightarrow 1 \\ < \tau \rightarrow 0 \end{cases}$$

y	$\hat{y}$
1	0.9
1	0.8
0	0.7
0	0.6
0	0.3

$$\hat{y} \geq 0.5$$

1 ✓  
1 ✓  
0 ✗  
0 ✗  
0 ✓

$$\hat{y} \geq 0.8$$

1 ✓  
1 ✓  
0  
0 ✓  
0 ✓  
0 ✓

Pred

		Actual	
		0	1
0	0	TN 1	FP 2
	1	FN 0	TP 2

Pred

		Actual	
		0	1
0	0	TN 3	FP 0
	1	FN 0	TP 2

→ Problem of P, R, F1

↳ Dependent on threshold value

ROC (Receiver Operating Characteristics)

↳ Defined as plot b/w TPR (Recall)

and FPR as decision threshold



TPR



✓

$y$	$\hat{y}$
1	0.9
0	0.6
1	0.7
0	0.3
1	0.8

①

Sort in decreasing prob

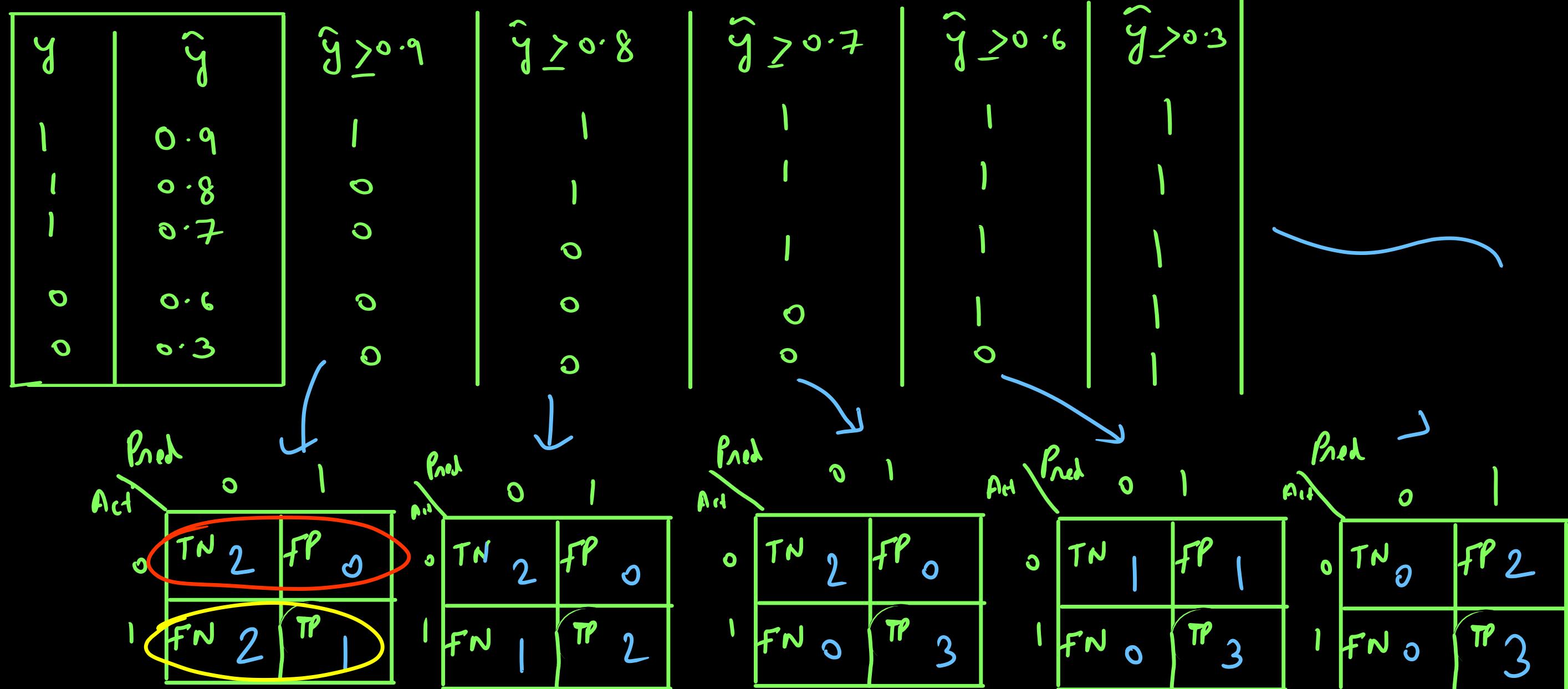


$y$	$\hat{y}$
1	0.9
1	0.8
1	0.7
0	0.6
0	0.3

2. Choose diff. threshold and calculate prediction

$y$	$\hat{y}$	$\hat{y} \geq 1$	$\hat{y} \geq 0.9$	$\hat{y} \geq 0.8$	$\hat{y} \geq 0.7$	$\hat{y} \geq 0.6$	$\hat{y} \geq 0.3$	$\hat{y} \geq 0$
1	0.9	0	1	1	1	1	1	1
1	0.8	0	0	1	1	1	1	1
1	0.7	0	0	0	1	1	1	1
0	0.6	0	0	0	0	1	1	1
0	0.3	0	0	0	0	0	1	1

3. Plot confusion matrix 4 get TPR and FPR



$$TNR = \frac{TN}{TN+FP} = \frac{2}{3} \rightarrow 0.67$$

$$FPR = \frac{FP}{TN+FP} = 0$$

$$\frac{2}{3} \rightarrow 0.67$$

$$\frac{3}{3+0} = 1$$

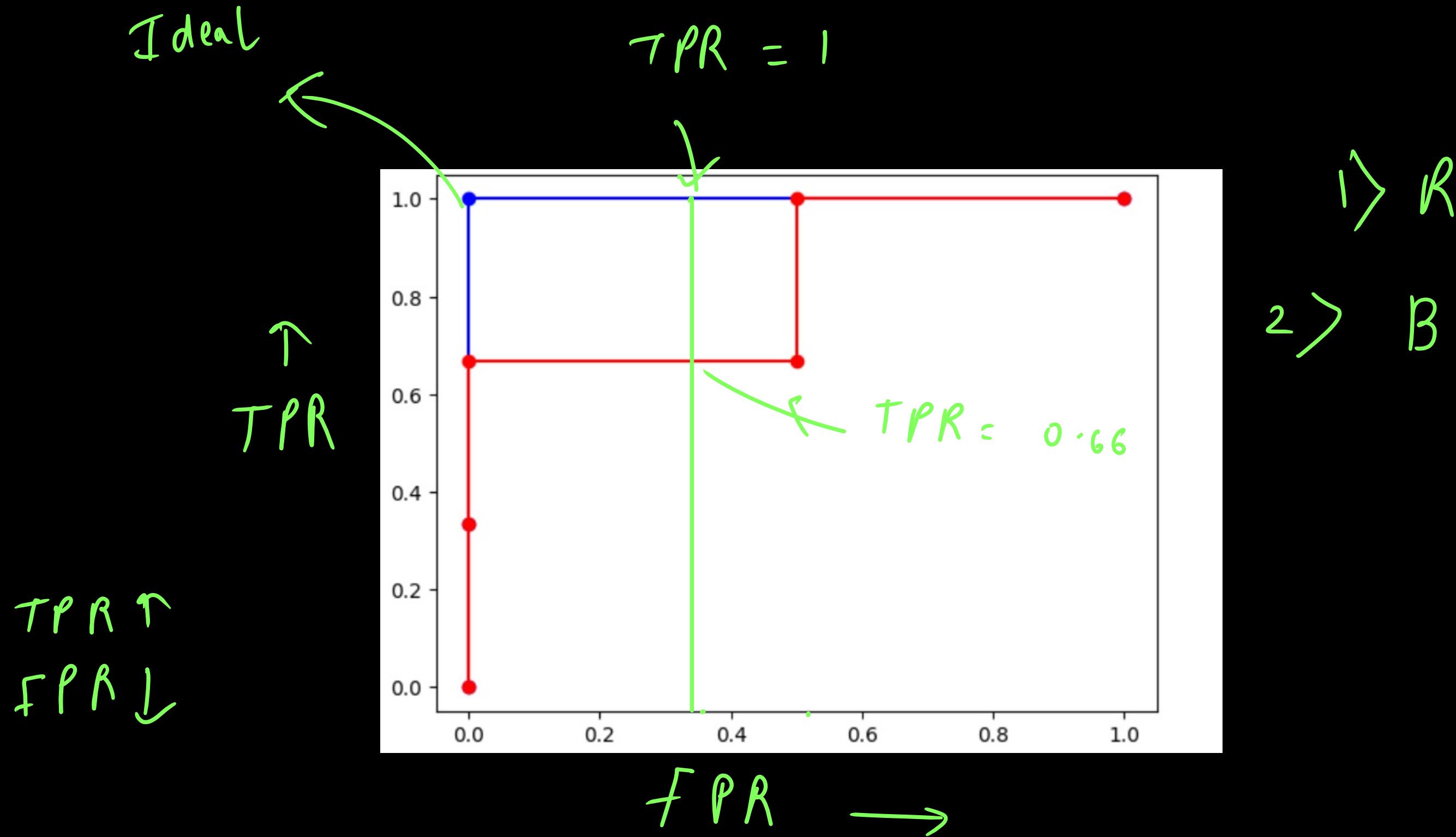
$$\frac{3}{3+0} = 1$$

$$\frac{0}{2+0} = 0$$

$$\frac{3}{3+0} = 1$$

$$\frac{2}{2+0} = 1$$

$$\frac{1}{1+1} = 0.5$$



$$TRR = \frac{1}{1+2} = \frac{1}{3} \rightarrow 0.33$$

$$FPR = \frac{0}{2+0} = 0$$

$$\frac{2}{3} \rightarrow 0.66$$

$$\frac{0}{2} = 0$$

$$\frac{3}{3+0} = 1$$

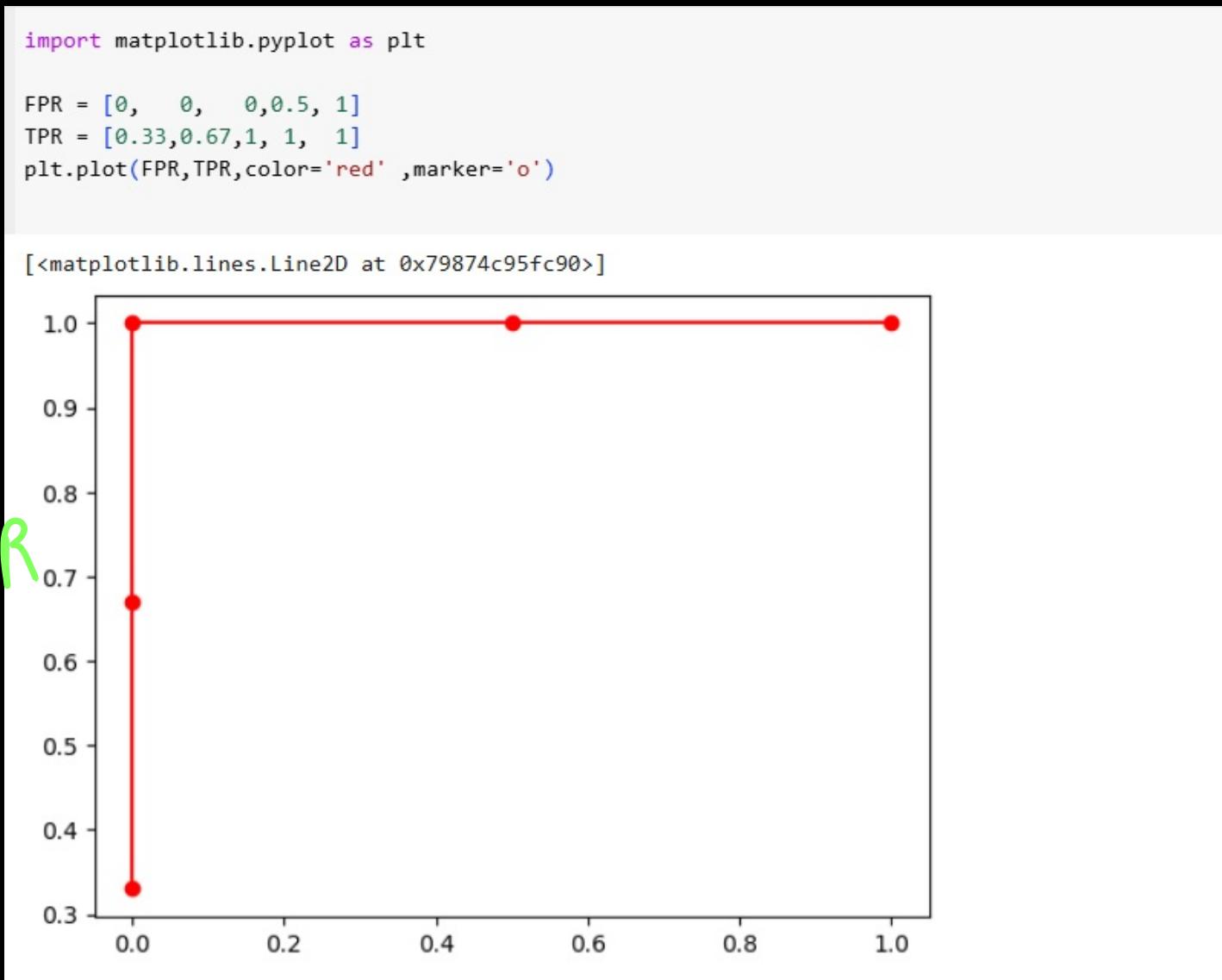
$$\frac{0}{2+0} = 0$$

$$\frac{3}{3+0} = 1$$

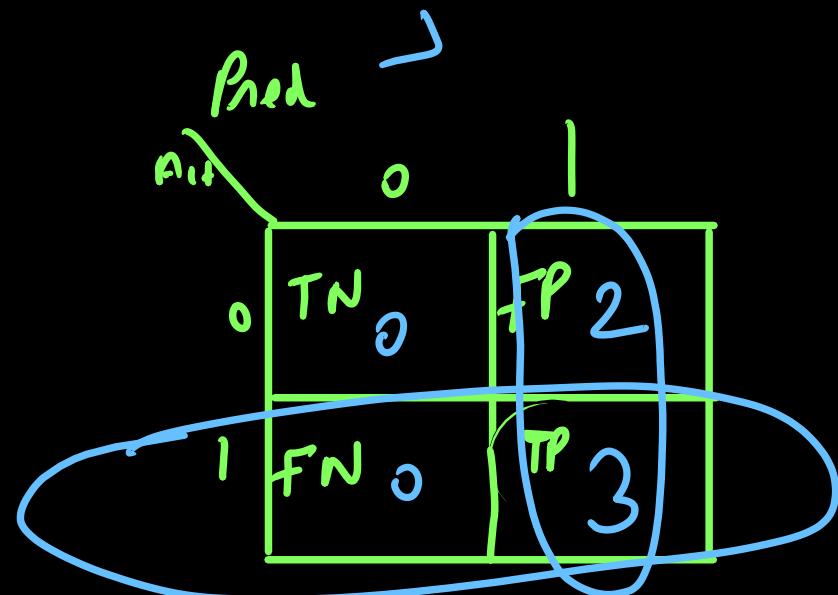
$$\frac{1}{1+1} = 0.5$$

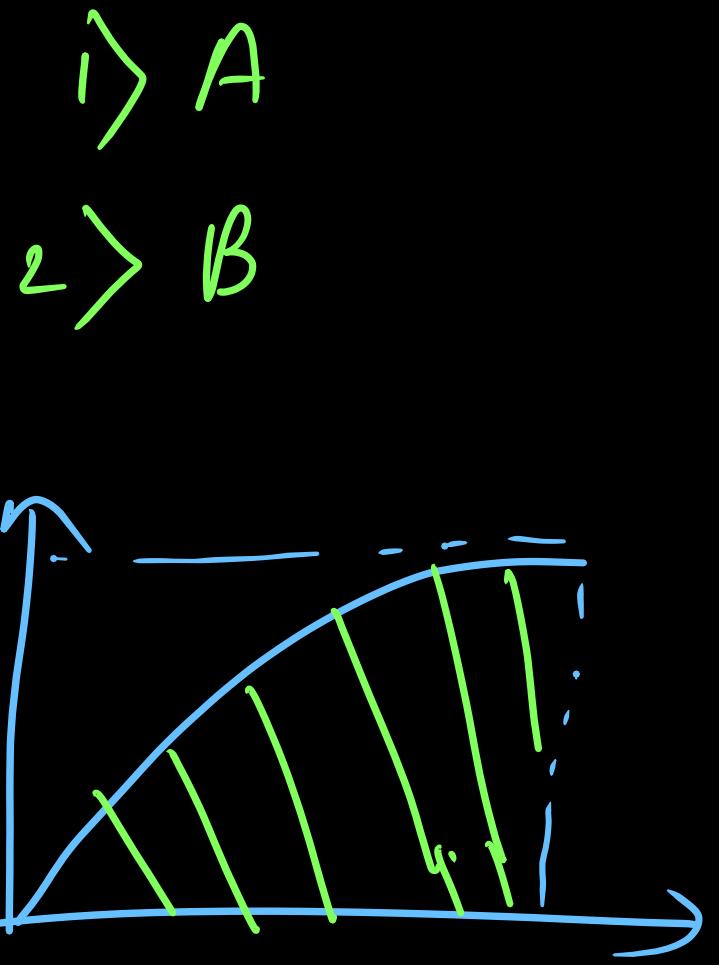
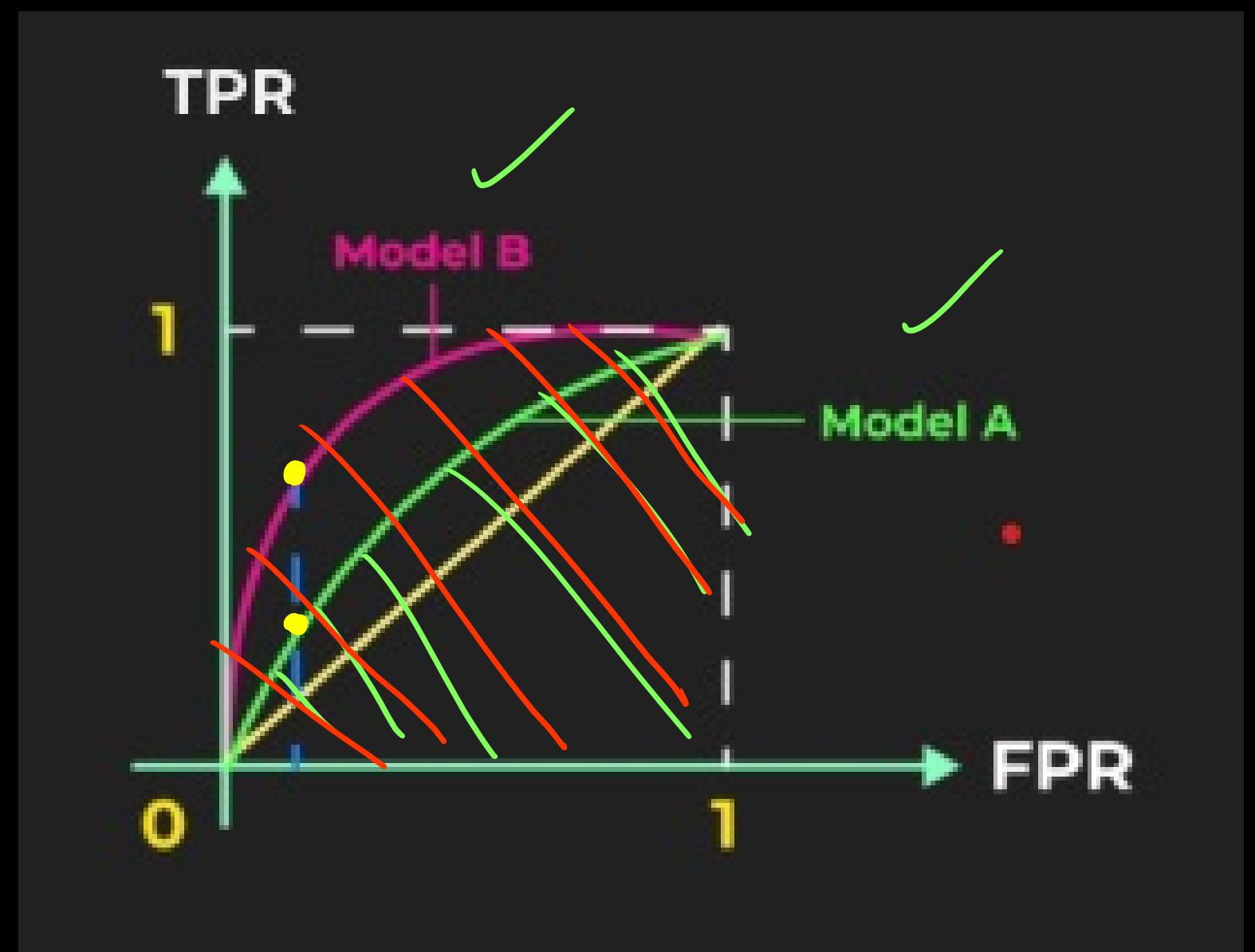
$$\frac{3}{3+0} = 1$$

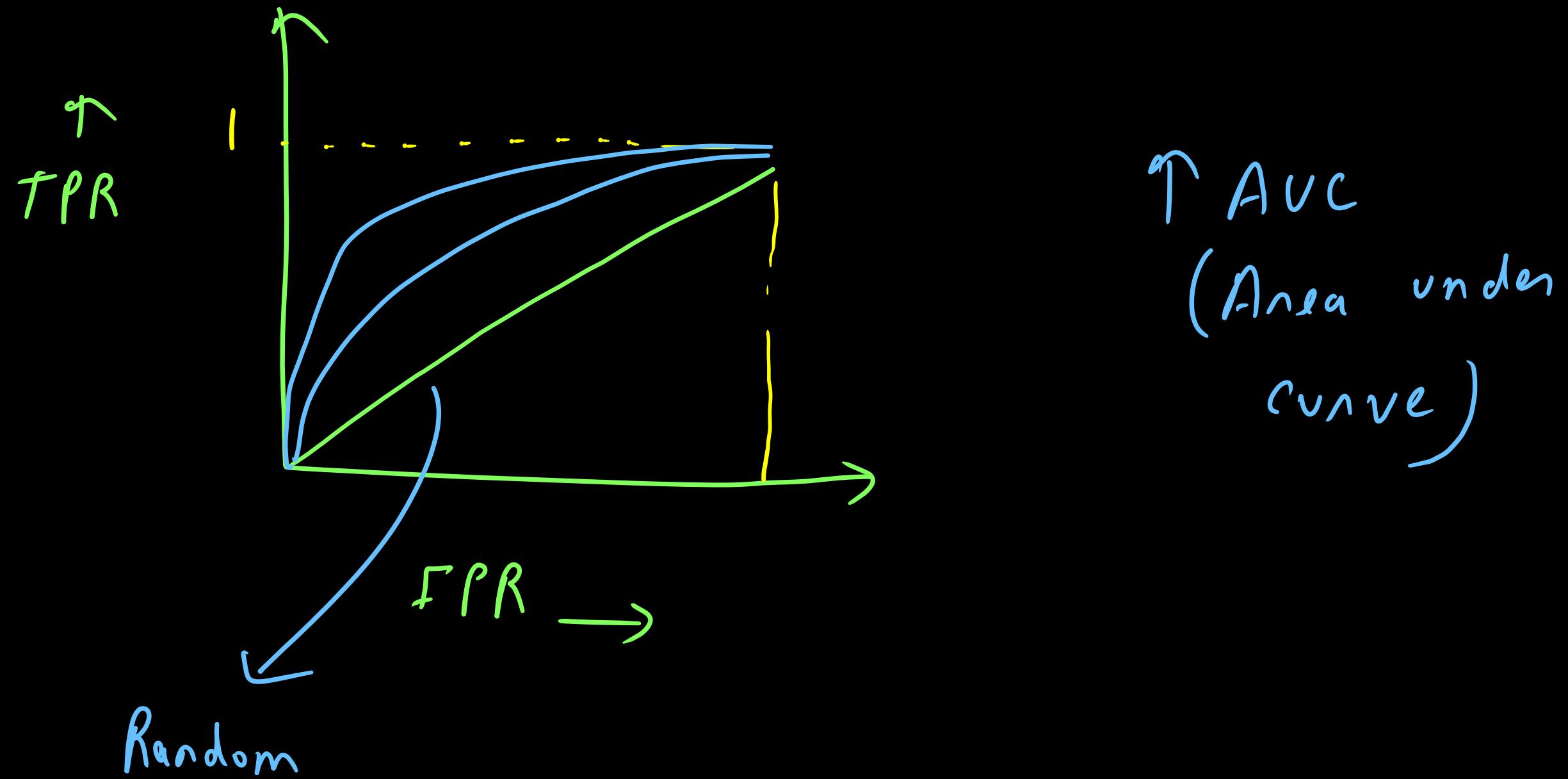
$$\frac{2}{2+0} = 1$$



$FPR \longrightarrow$







## Problem for AUC - ROC

---

1. Does not work for imbalanced dataset

✓ 2. It takes "order of Probab" and not <sup>actual probab</sup>  


Model 1



$y$	$P_{m1}$	$\hat{y}_\tau \geq 0.95$	$\hat{y}_\tau = 0.92$
1	0.95	1	1
1	0.92	1	0
0	0.80	0	0
1	0.76	0	0
1	0.77	0	0

MI

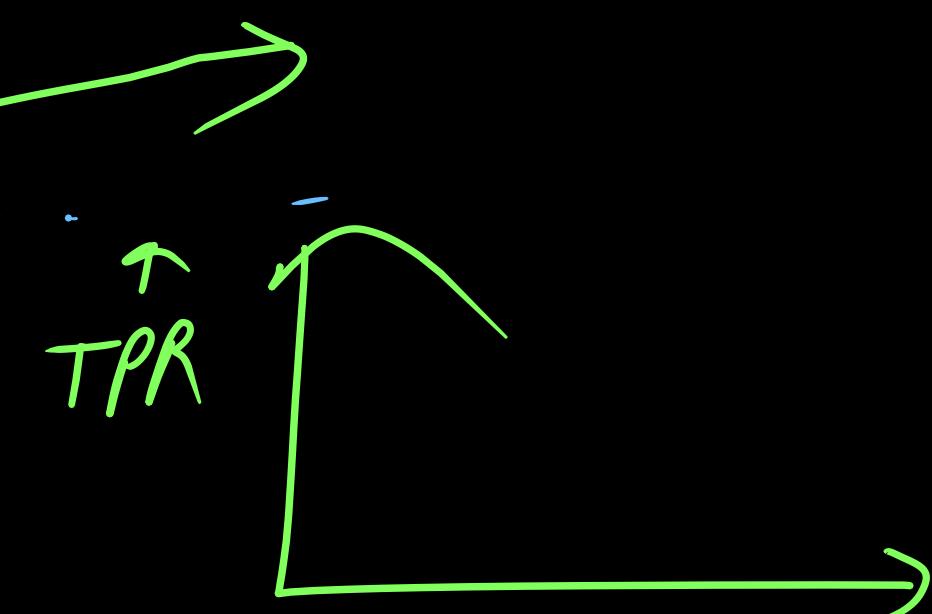
Model 2



$y$	$P_{m2}$	$\hat{y}_\tau \geq 0.2$	$\hat{y}_\tau = 0.1$
1	0.2	1	1
1	0.1	1	0
0	0.08	0	0
1	0.06	0	0
1	0.02	0	0

TPR

FPR

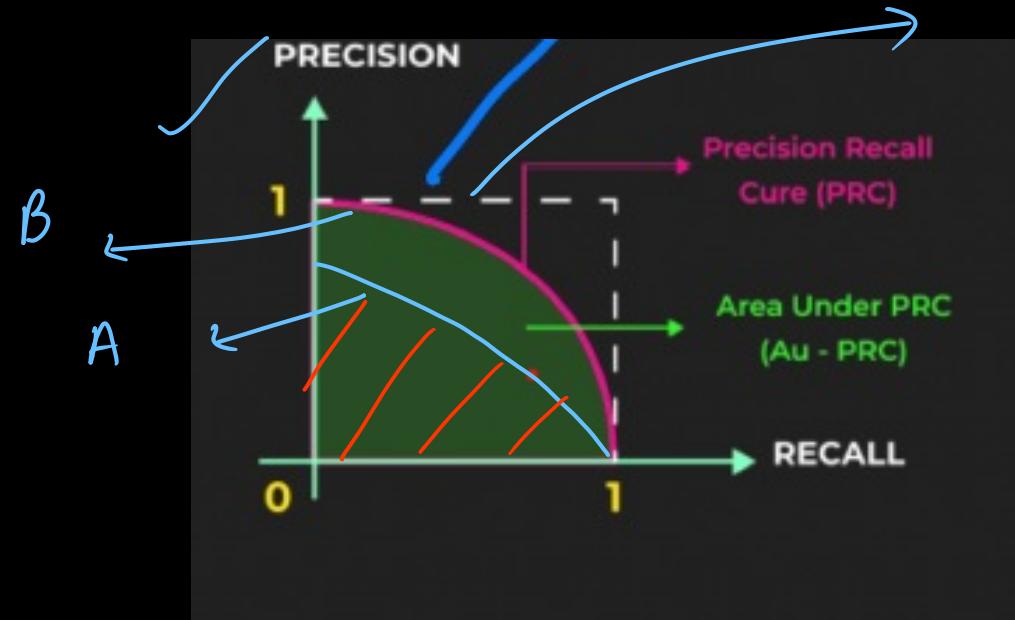


Sum

→ Break until 22:33 PM

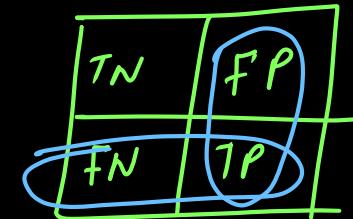
---

## PR - curve



Ideal case

AUC-PR curve



ROC → Balanced

PR → Imbalanced dataset & Balanced Dataset

-ve → 9990

+ve → 10

$$P = \frac{5}{5+5} = 0.5$$

-ve → 10000

1

2

3

4

5

0.94

0

1

2

3

4

5

R = TPR

1

ROC

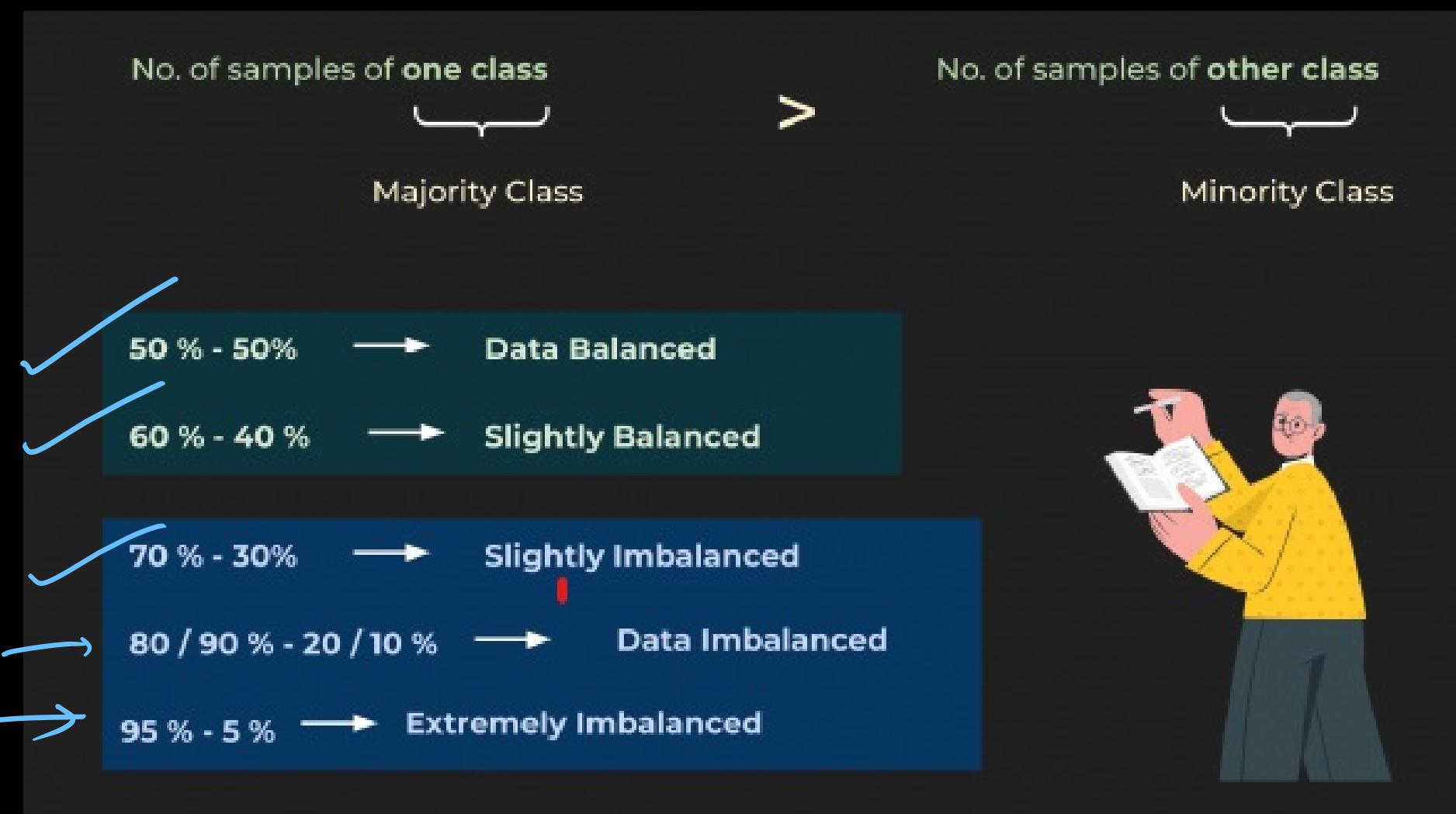
FPR

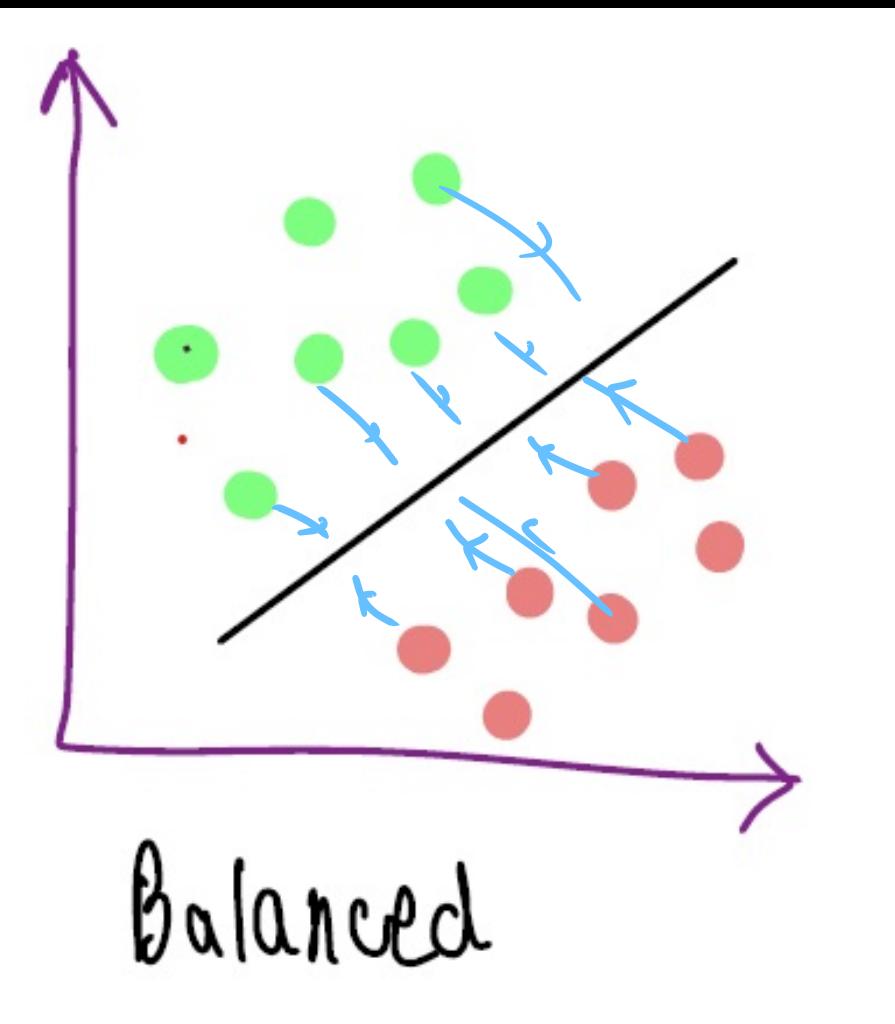
$\frac{FP}{FP+TN}$

$= \frac{5}{5+9990} \approx 0.001$

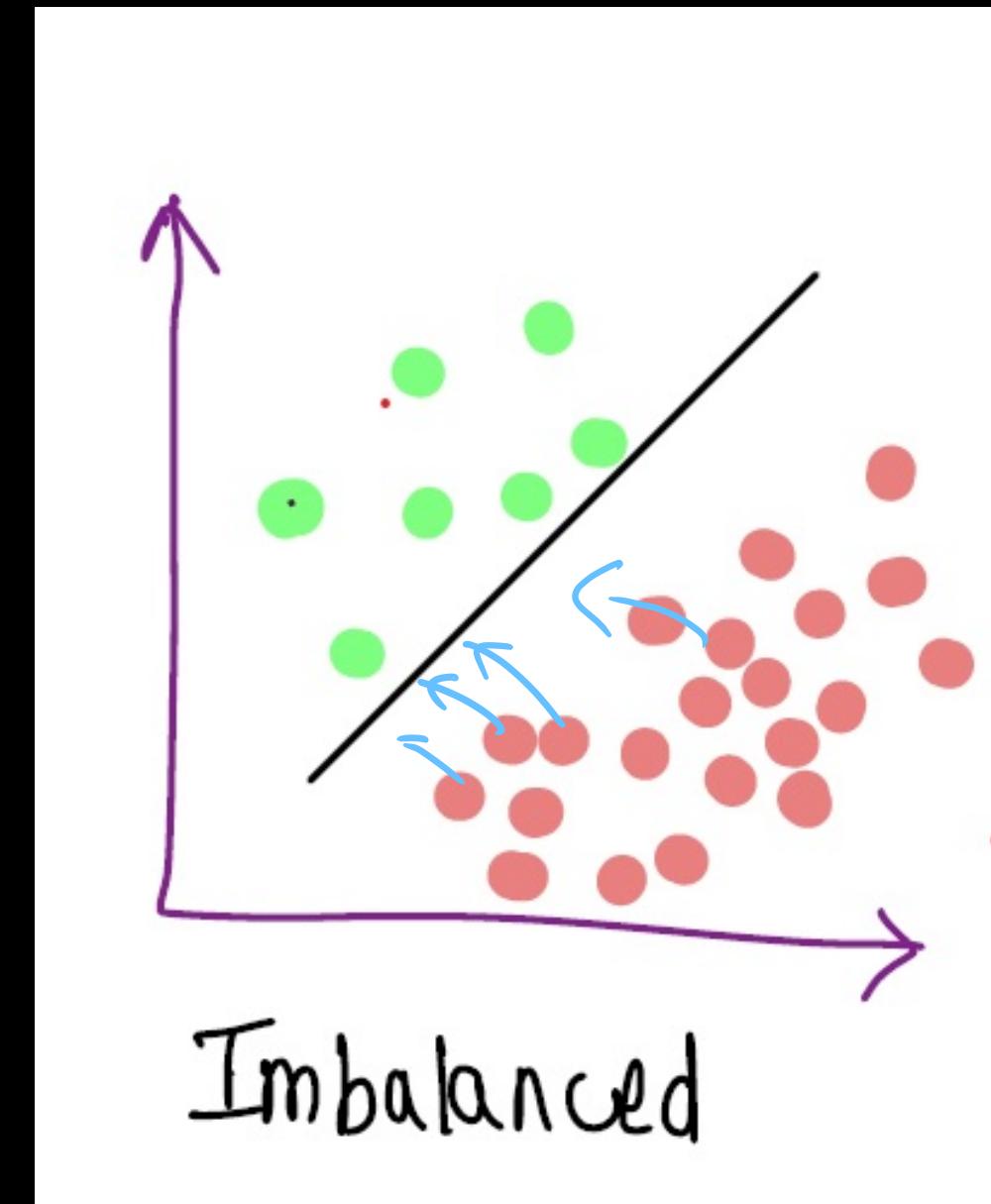
$$\text{TPR} (\text{Recall}) = \frac{5}{5+0} = 1$$

→ Imbalanced Data





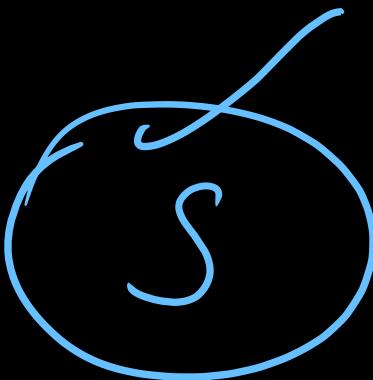
Balanced



Imbalanced

## Spam Classification

1. Spam  $\rightarrow$  150



2. NS  $\rightarrow$  850

$$\frac{NS}{S} = \frac{850}{150} = 5.67$$

$$NS = 5.67S$$

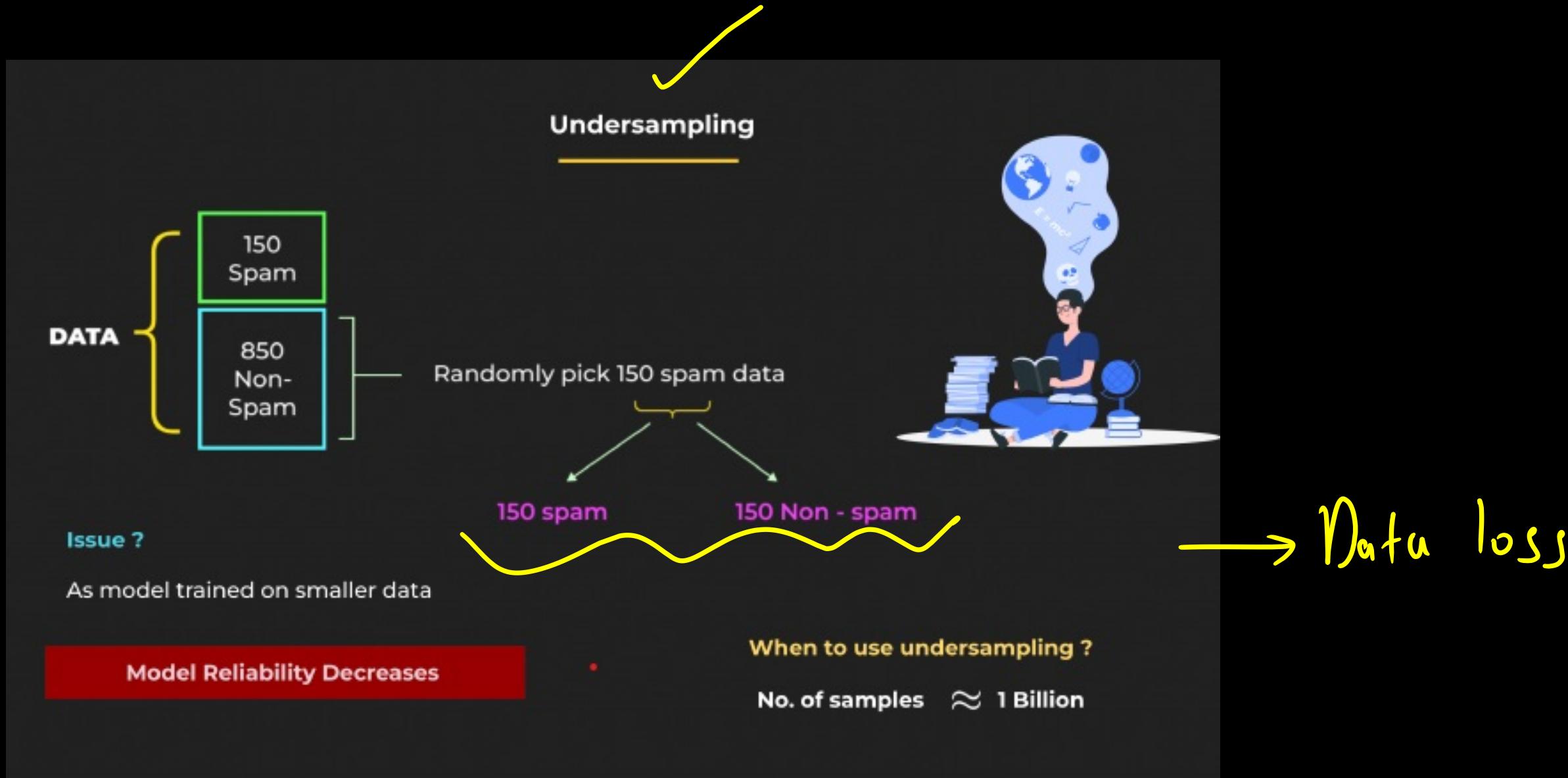
→ More weightage to minority class

$$\therefore \text{loss} = \sum_{i=1}^n \log loss_i W_i + \lambda \sum_{j=1}^d w_j^2$$

$W_i = 5.67 \text{ when spam}$      $\frac{N_S}{S} = \frac{0}{1} = 5.67$   
 $W_i = 1 \text{ when non-spam}$

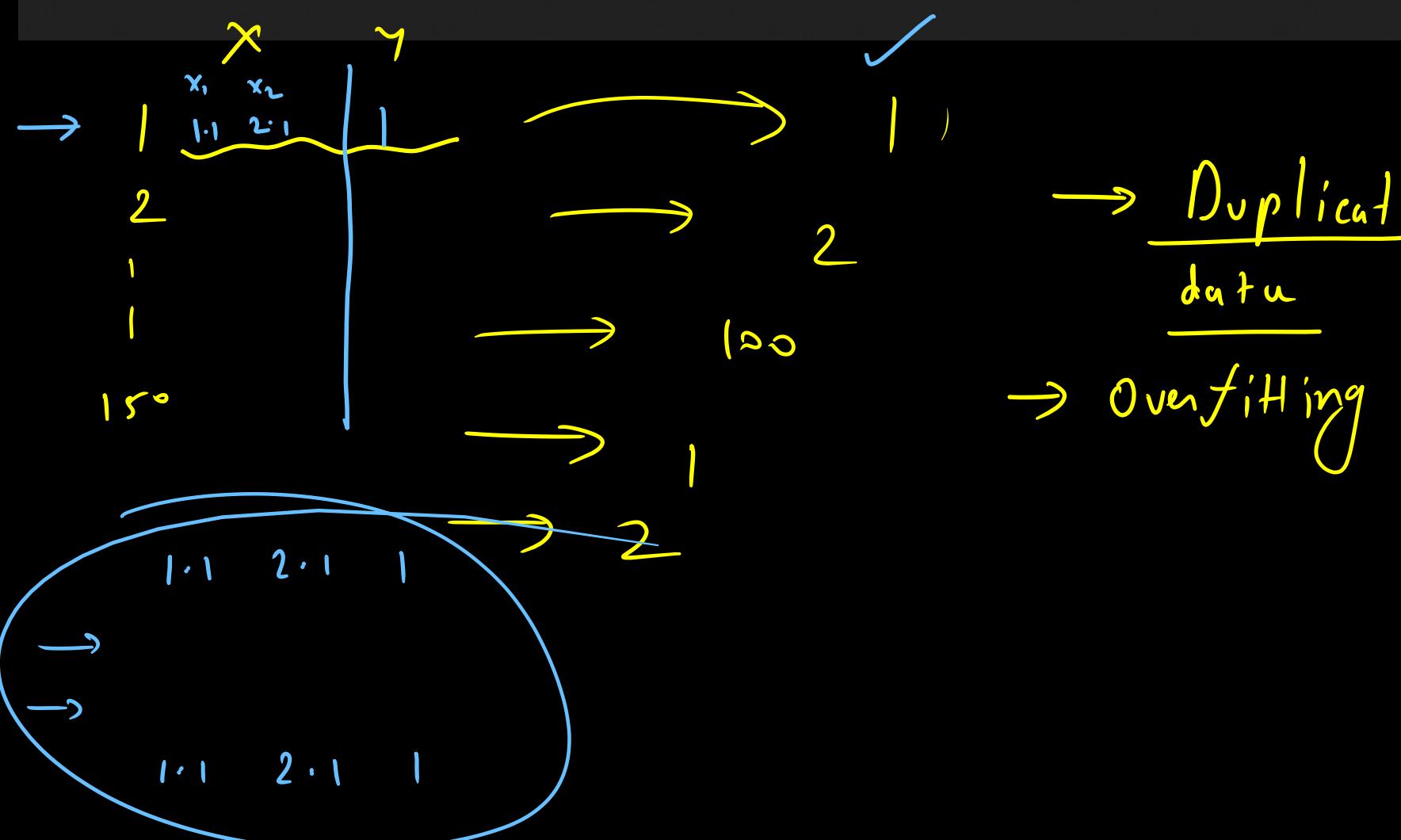
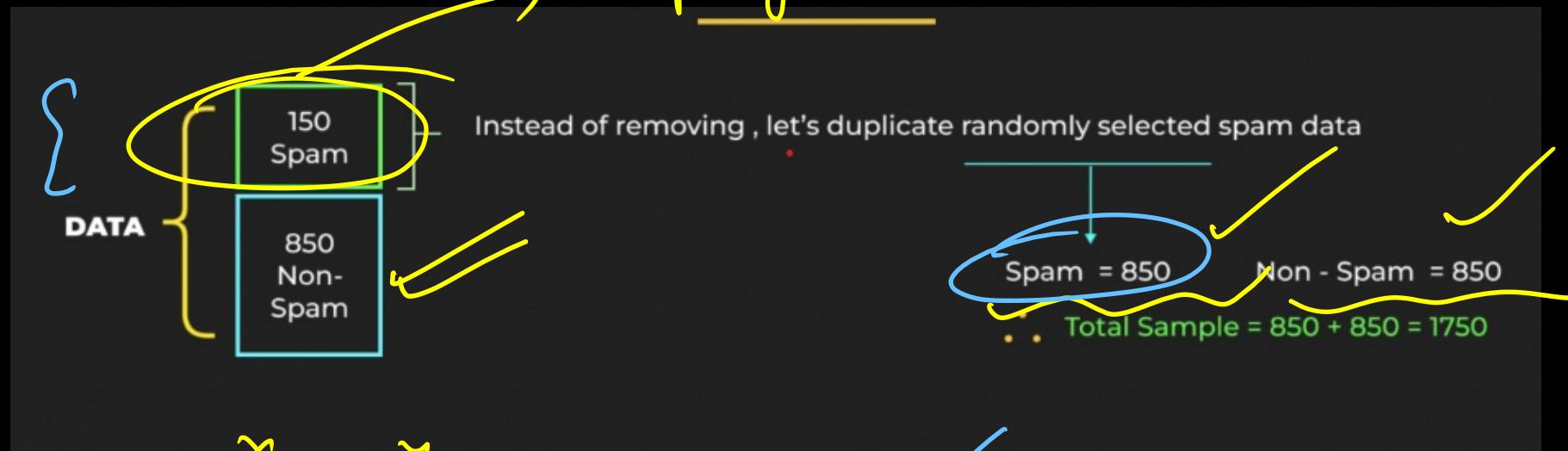
$$\{ 0 : 1 , 1 : 5.67 \}$$

## 2) Undersampling



### 3) Oversampling

Sampling with replacement



## SMOTE

↳ Synthetically Minority Oversample technique

↳ Generate Synthetic data

1. Select Random point from minority class
2. Find  $K = 3$  nearest point
3. Randomly select one of  $K$  - neighbor

$$x_{syn} = x_i + \lambda (x_j - x_i)$$

1. Select a random point from minority class

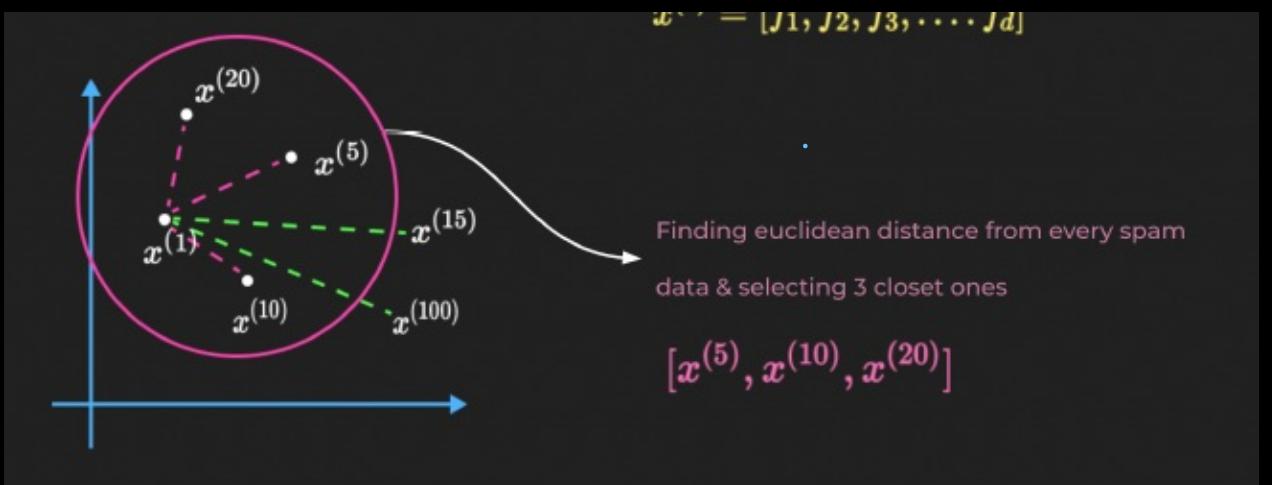
2. Find  $K = 3$  nearest point

3. Randomly select one of  $K$ -neighbors

$$x_{syn} = x_i + \lambda (x_j - x_i)$$

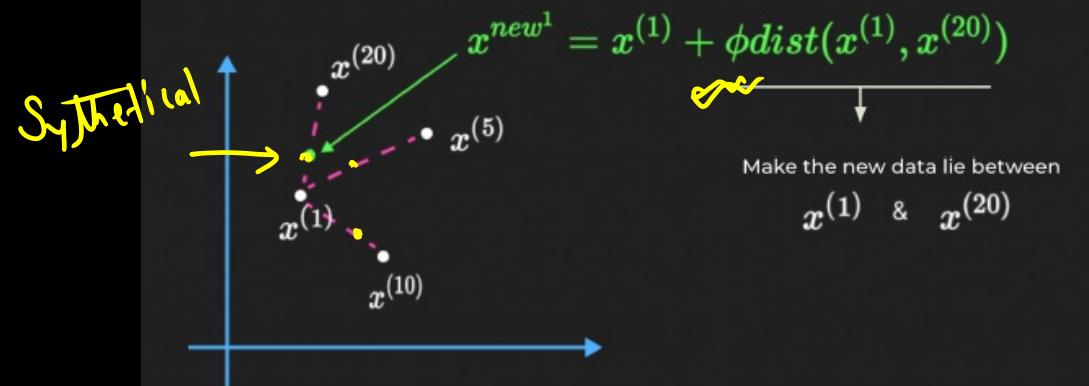
$\lambda \sim U[0, 1]$

$x^{(1)}$ ,  $x^{(5)}$ ,  $x^{(10)}$ ,  $x^{(20)}$



$X$	1
	1
	1
	0
	0

Taking a random number  $\phi \in [0, 1]$



0.1

0.2

0.8

