

# Agenda

- ⇒ Recap
  - ⇒ Sensitivity and Specificity
  - ⇒ FNR and FPR
  - ⇒ ROC and AU-ROC
  - ⇒ PR Curve
  - ⇒ Data Imbalance
  - ⇒ Techniques to Handle
    - ⇒ Adjusting Weights of Loss Function
    - ⇒ UnderSampling and OverSampling
    - ⇒ SMOTE

# Recap

$$\text{precision} = \frac{TP}{TP + FP}$$

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

- ⇒ Precision ⇒ How precise is model in pred  $y=1$
  - ⇒ Recall ⇒ Out of all points belonging to  $y=1$ , How many model was able to predict.

## F-Beta Score

### F1-Score

When Both Precision and Recall are equally important, we can take mean of P and R

$$\textcircled{1} \text{ Simple-mean} = \frac{P+R}{2}$$

X

$$\textcircled{2} \text{ Harmonic Mean} = \frac{2 \times P \times R}{P+R}$$

\* F1-Score for bad Model  $\Rightarrow$  ?

$$P \Rightarrow 0 \\ R \Rightarrow 0$$

$$\text{F1-Score} = \frac{0 \times 0 \times 2}{0+0+0.0000} \Rightarrow 0$$

\* F1-Score for ideal Model  $\Rightarrow$  ?

$$P=1 \\ R=1$$

$$\text{F1-Score} = \frac{1 \times 1 \times 2}{1+1+0.0000} \Rightarrow \frac{2}{2+e}$$

Why Harmonic Mean?

	Precision	Recall	Mean	H.M
m1	0.30	0.80	0.55	0.436
m2	0.20	0.90	0.55	0.327
m3	0.70	0.40	0.55	0.51

$$\text{mean} = \frac{P+R}{2}$$

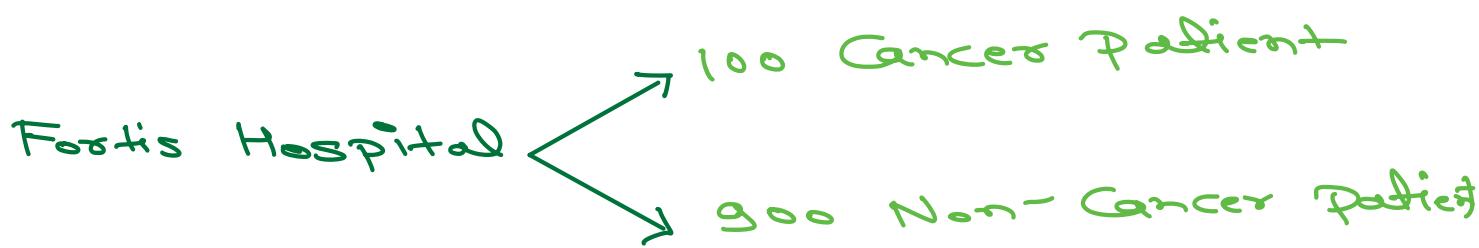
$$\text{H.M} \Rightarrow \frac{2 \times (P \times R)}{P+R}$$

$$\frac{2 \times 0.3 \times 0.8}{0.3 + 0.8}$$

Conclusion :

- F1-Score gives equal importance to both precision and Recall
- It adds a lot of penalty when either of P or R is very low

## Sensitivity



Aim: Correctly identify all Cancer patients

- TP ↑ : Catch as many cancer patients as possible
- FN ↓ : Reduce mis-classification of cancer patients

TP ↑ and FN ↓ : Sensitivity

$$\frac{TP}{TP + FN}$$

→ TPR  
→ Recall  
→ Sensitivity

Why sensitivity is important?

High Sensitivity is crucial as failing to detect disease at early stage can be Fatal

## Specificity

Role of TN and FP in Above scenario?

- $TN \uparrow$  : Model should be able to identify Non-Cancer Patient
- $FP \downarrow$  : Model should keep mis-classification of Non-Cancer LOW

Specificity  $\Rightarrow$

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

$TNR$   
(True Negative)  
Rate

Specificity is Sensitivity for Class 0

Why specificity is important?

- ⇒ Costly and fruitless Tests / procedures
- ⇒ Stress and anxieties to patients

## FNR and FPR

FNR  $\Rightarrow$  False Negative Rate

⇒ Miss Rate

⇒  $1 - \text{Sensitivity}$

⇒  $1 - \frac{TP}{TP + FN}$

$$\frac{TP + FN - TP}{TP + FN}$$

⇒  $\frac{FN}{TP + FN}$

$FN \uparrow$

\* In-sensitive Model  $\Rightarrow$  Model with high FNR is considered in-sensitive Model

$$FPR = 1 - \text{Specificity}$$

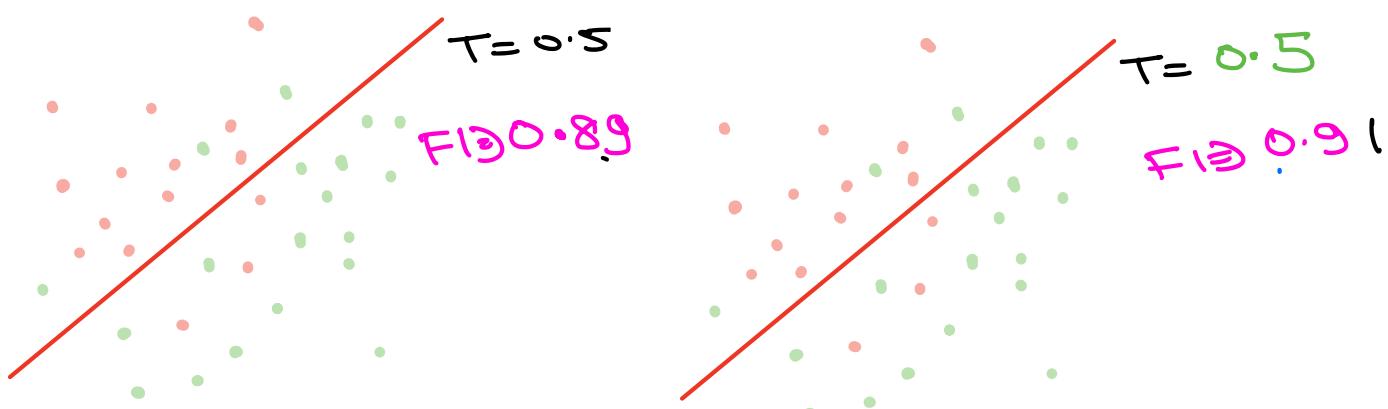
$$\Rightarrow 1 - \frac{TN}{TN+FP} \Rightarrow \frac{TN+FP-TN}{TN+FP}$$

False Positive Rate  $\Rightarrow \frac{FP}{TN+FP}$

\* Not Specific Model : When model has high FPR

### Scenario

$\text{Q: Suppose you have a model}$



$m_1 \Rightarrow \text{Threshold } (0.5)$

$m_2 \Rightarrow \text{Threshold } 0.5$

To improve Model Performance

→ Hyper parameter tuning (L1, L2R)

\* → Find the best value of Threshold

Find the best value of Threshold

Compare different Model independent of Threshold.

Q-1 How do we find the Best Threshold

Q-2 How do we compare two different type of Models?

⇒ Receiver operating Curve

TPR vs FPR

⇒ PR Curve

Precision vs Recall

# ROC

Step 1 : Get all probabilities and Sort

$x$	$y \in (0,1)$	$P = P(y_i=1/x_i)$
$x_1$	$y_1$	$P_1$
$x_2$	$y_2$	$P_2$
$x_3$	$y_3$	$P_3$
$\vdots$	$\vdots$	$P_{\dots}$
$x_n$	$y_n$	$P_n$

n Logn  
 Sort  
 the table  
 based  
 on  
 $P_{y=1/x}$   
 (descending)

Step 2 : Calculate  $\hat{y}^i$  for all prob-scores

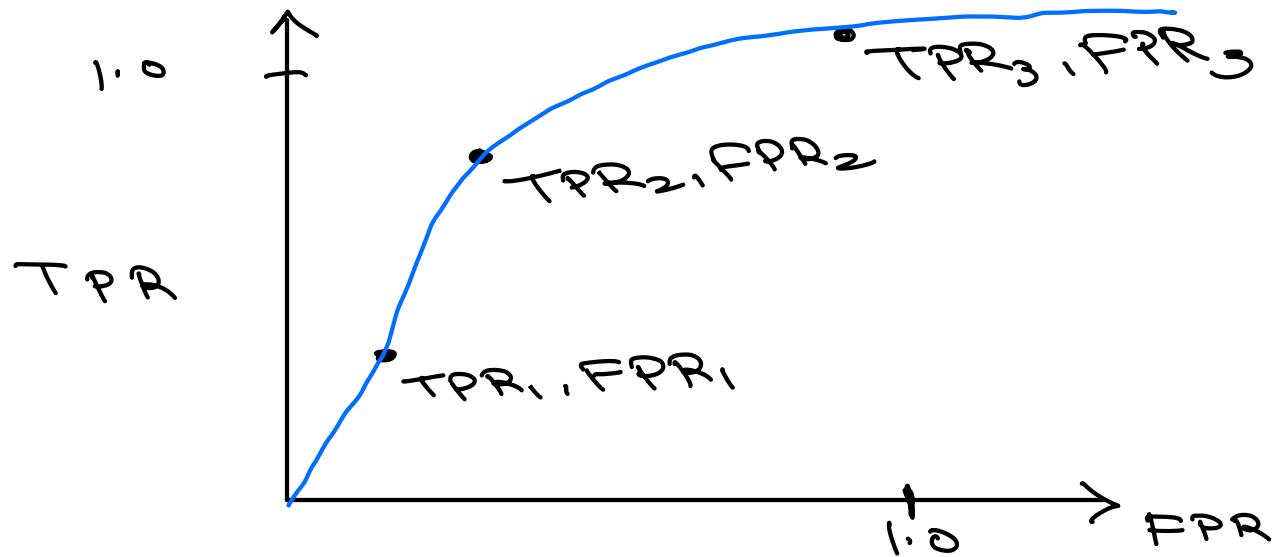
$x$	$y$	$P$	$\hat{y} \rightarrow P_1$	$\hat{y} \rightarrow P_2$	$\hat{y} \rightarrow P_3$	$\hat{y} \rightarrow P_n$
$x_1$	$y_1$	$P_1$	-	-	-	-
$x_2$	$y_2$	$P_2$	0	-	-	-
$\vdots$	$\vdots$	$\vdots$	0	0	-	-
$x_n$	$y_n$	$P_n$	0	0	0	-

$$TPR \Rightarrow \frac{TP}{TP + FN}$$

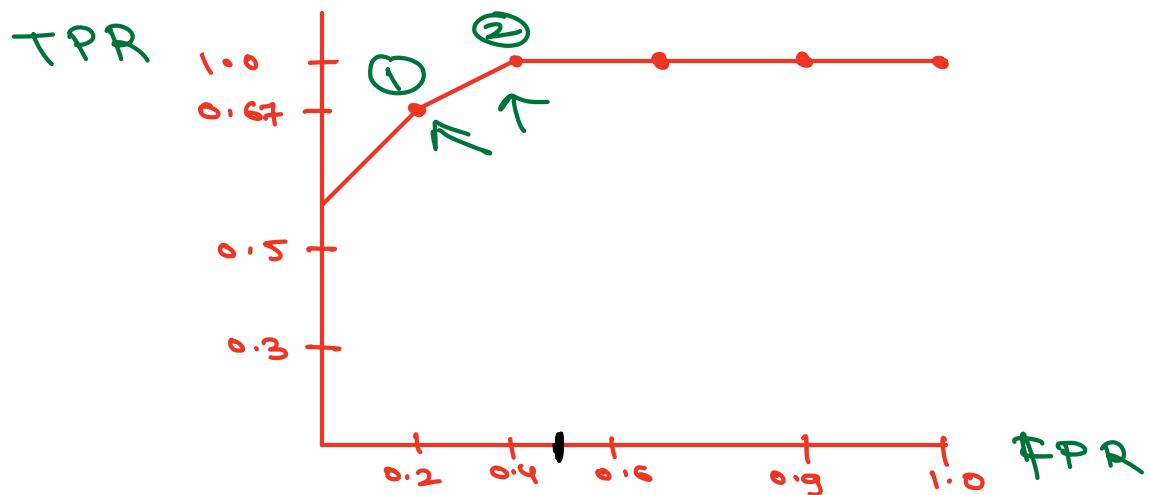
$$FPR \Rightarrow \frac{FP}{FP + TN}$$

$P$	$TPR$	$FPR$
$P_1$	$TPR_1$	$FPR_1$
$P_2$	$TPR_2$	$FPR_2$
$P_3$	$TPR_3$	$FPR_3$
$\vdots$	$\vdots$	$\vdots$
$P_n$	$TPR_n$	$FPR_n$

### Step 3: Plot TPR vs FPR



$\Leftrightarrow$  Ideal Threshold:



HomeWork: Plot ROC and Find opt T

X	Y	P
$x_1$	-1	0.65
$x_2$	-1	0.94
$x_3$	0	0.3
$x_4$	-1	0.92
$x_5$	0	0.7
$x_6$	0	0.2

P	T = 0.94
0.94	1
0.92	0
0.7	0
0.65	0
0.3	0
0.2	0
0	1

P	T = 0.92
0.94	1
0.92	0
0.7	0
0.65	0
0.3	0
0.2	0
0	1

$$TPR = \frac{TP}{TP + FN}$$

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

$$TPR_{0.94} \Rightarrow$$

$$\frac{1}{1+2} \Rightarrow$$

$$0.33$$

$$3 \Rightarrow 0.5$$

$$1.0$$

$$FPR_{0.94} \Rightarrow$$

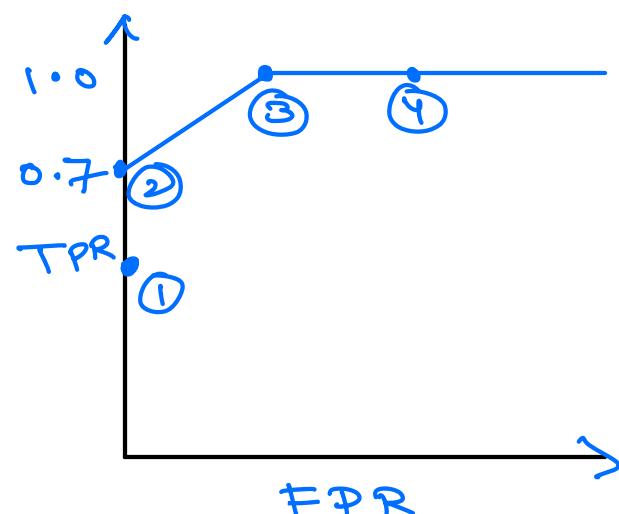
$$\Rightarrow$$

$$0$$

$$2 \Rightarrow 0 FPR$$

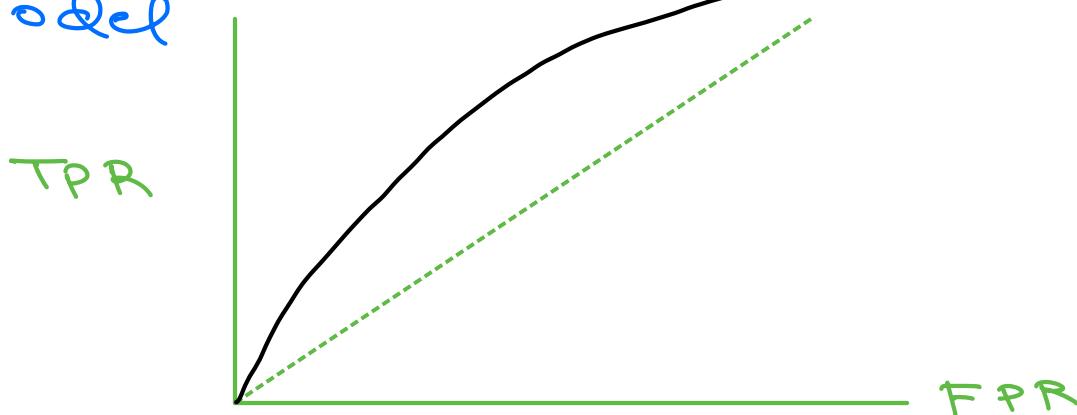
$$0.9 TPR$$

Thresh	TPR	FPR
0.94	0.33	0



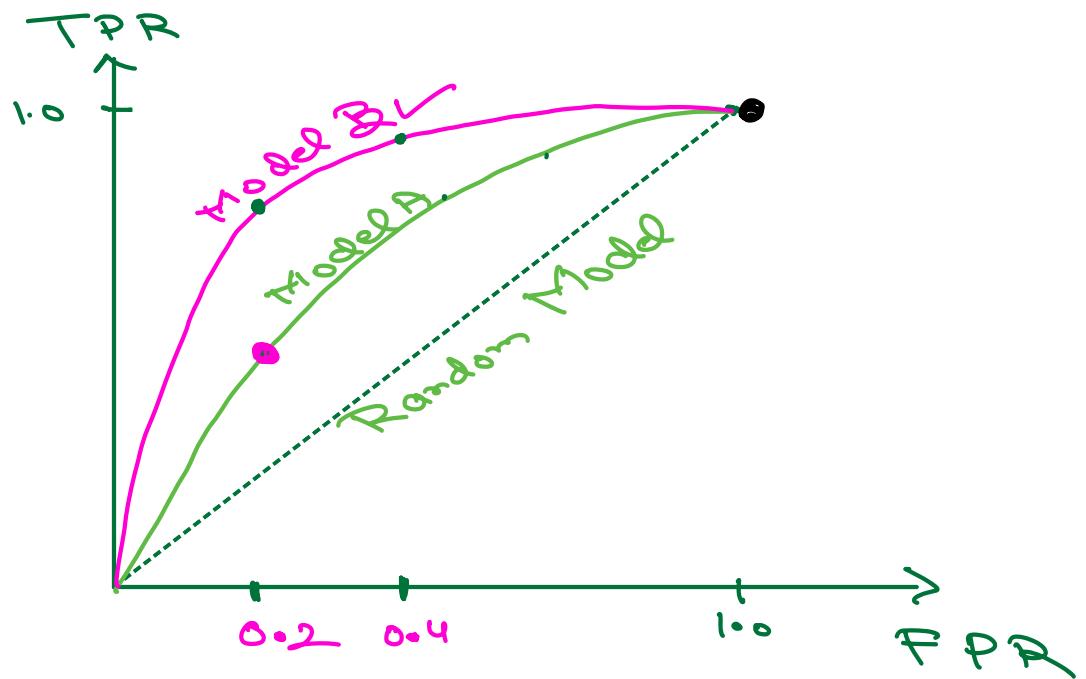
## ② Comparing different Models

ROC curve for a Random Model



## AU-ROC

⇒ Area Under ROC curve

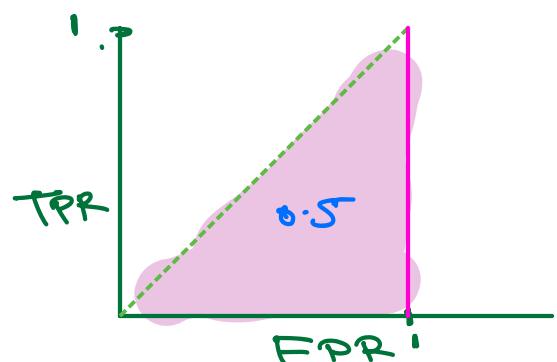


AU of Model A < AU of Model B

\* For Random Model

$$\text{TPR} \times b \times t$$

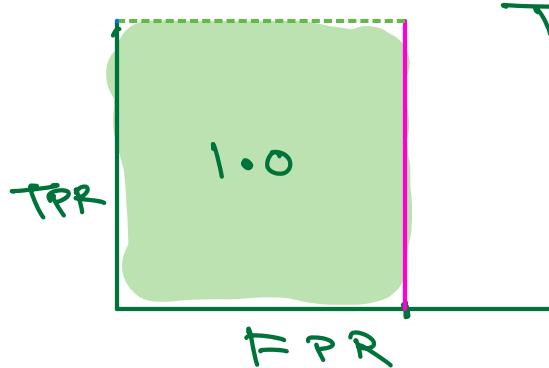
$$\text{TPR} \times 1 \times 1 \Rightarrow 0.5$$



\* For ideal Model

$$b \times t$$

$$1 \times 1 \Rightarrow 1.0$$



TPR always

1.0

## Issue with AU-ROC

③ Dependent on Order of probabilities

$$y \Rightarrow [1, 1, 0]$$

$$m_1 \Rightarrow [0.95, 0.92, 0.80]$$

$$m_2 \Rightarrow [0.20, 0.10, 0.08]$$

$$y_{pred} = 0.8$$

$y$	$P$	$\hat{y} \geq 0.95$	$\hat{y} \geq 0.92$	$\hat{y} \geq 0.80$
1	0.95	1	1	1
1	0.92	0	1	1
0	0.80	0	0	1

$m_1$

$M_1$

$y$	$P$	$\hat{y} \geq 0.2$	$\hat{y} \geq 0.10$	$\hat{y} = 0.08$
1	0.2	1	1	1
1	0.1	0	1	1
0	0.08	0	0	1

$m_2$

$M_2$

⇒ Will not work very well for Data Imbalance

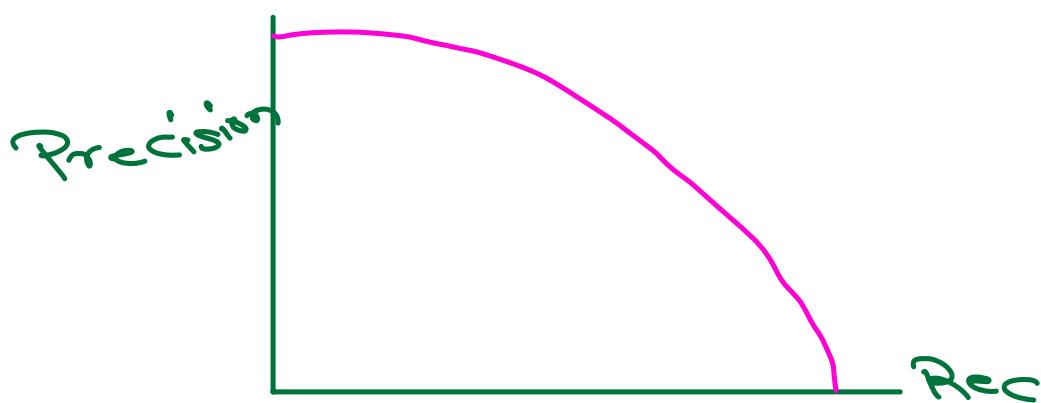
$$TPR = \frac{TP}{TP+FN}$$

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

⇒ PR-Curve and AU-PR are better choice

### Precision Recall Curve

⇒ Same steps as AU-ROC

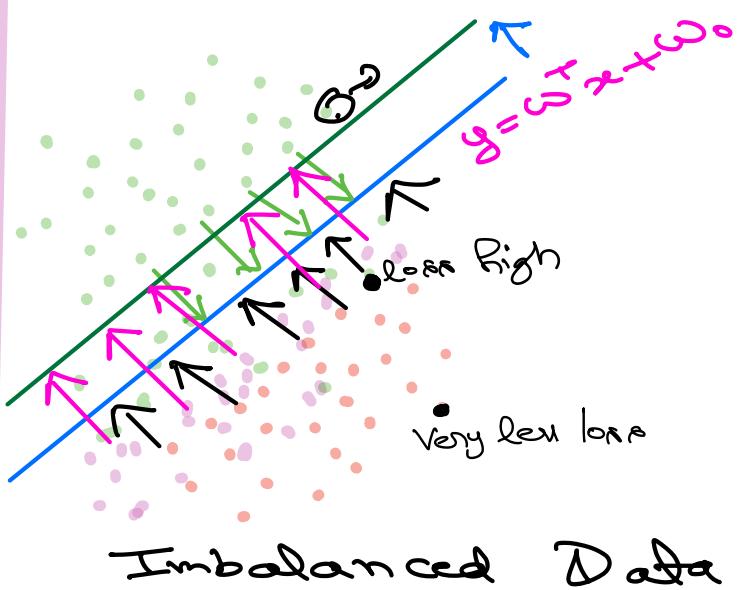
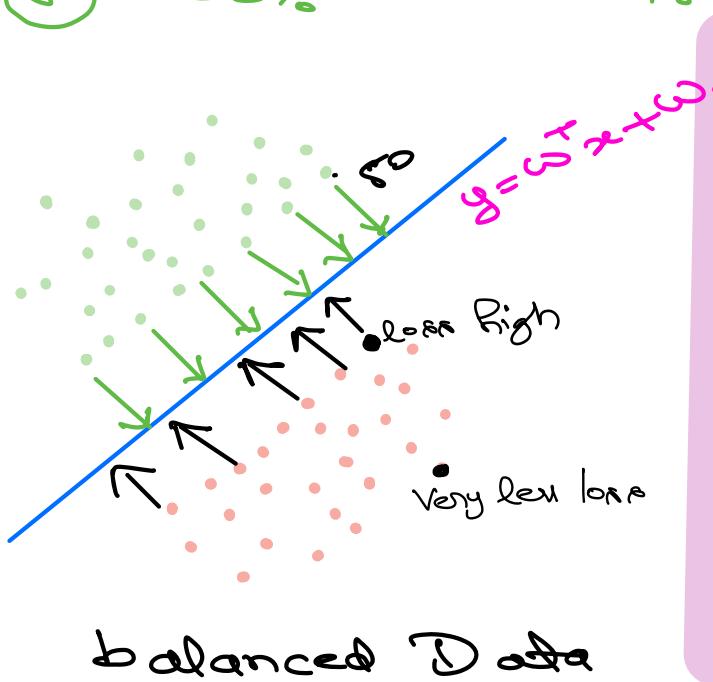


# Data Imbalance

When frequency of one category is higher than other categories

Majority% - Minority% → Category

- ① 50% - 50% → Balanced
- ② 60% - 40% → Slightly Balance
- ③ 70% - 30% → Slightly Imbalance
- ④ 80% - 20% → Imbalance
- ⑤ 95% - 5% → Highly Imbalanced



# Handling Data Imbalance

Suppose → 1000 Emails

150 Spam  
850 Non-Spam

Balanced Data

$$-\frac{1}{n} \sum_i (y_i \times \log(y_{\hat{o}}) + (1-y_i) \times \log(1-y_{\hat{o}})) + \gamma w_i^2$$

balanced  $\Rightarrow 50\% 50\%$

Imbalanced 25% (1) 75% (0)

$$-\frac{1}{n} \sum_i (y_i \times \log(y_{\hat{o}}) + (1-y_i) \times \log(1-y_{\hat{o}})) + \gamma w_i^2$$

Strategy 1: Set Class Weights

$$\text{NLL} \times \text{Class-1} + \lambda w^2$$

$$\text{NLL} \times \text{Class-Weights}$$

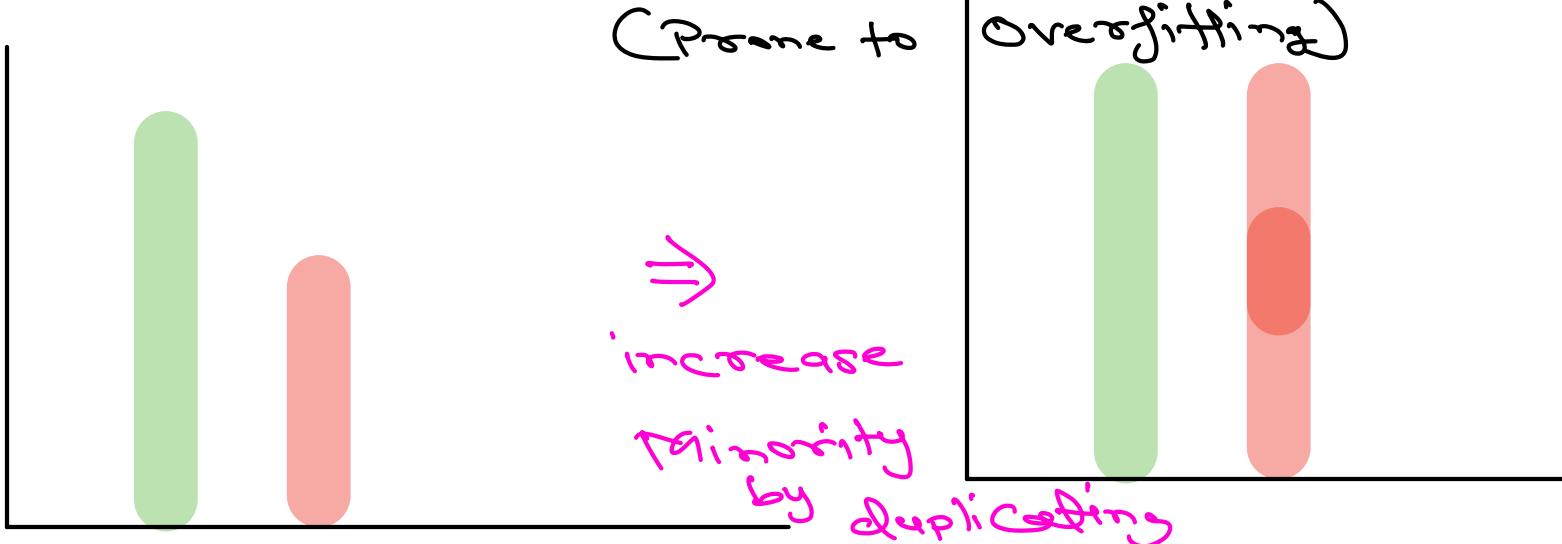
## Strategy 2: Modify the Data UnderSampling and OverSampling



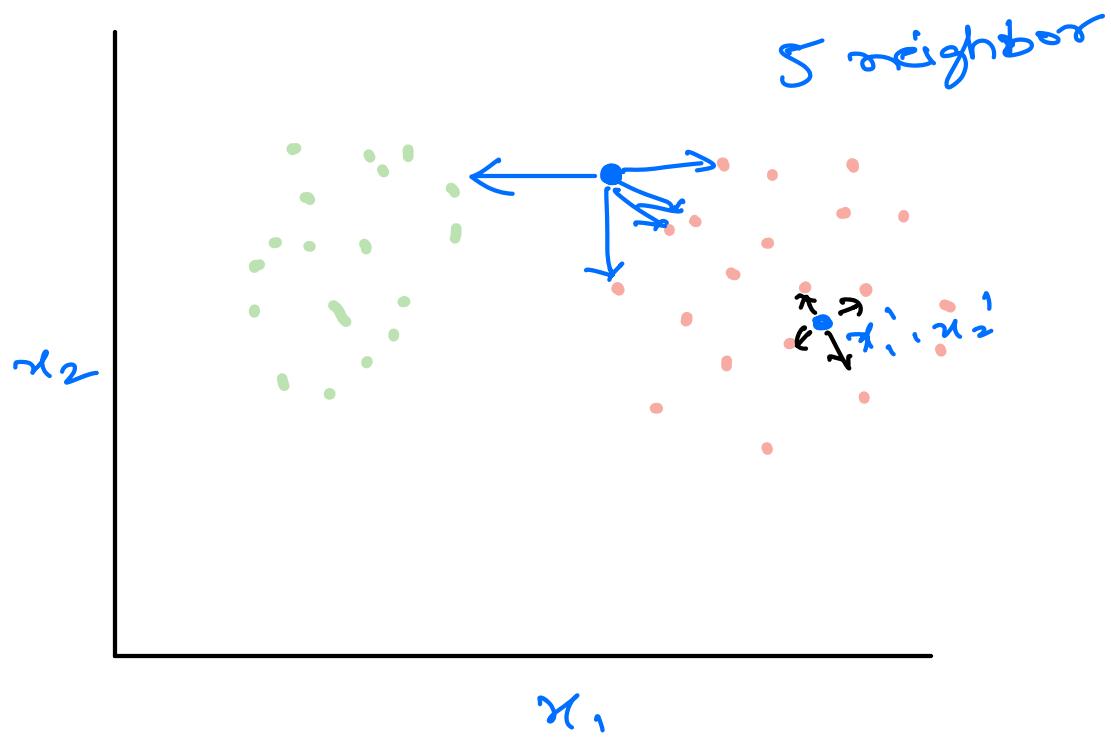
### ① UnderSampling (Prone to Underfitting)



### ② Oversampling by Duplication (Prone to Overfitting)



## Strategy 3: Synthetic Data Generation



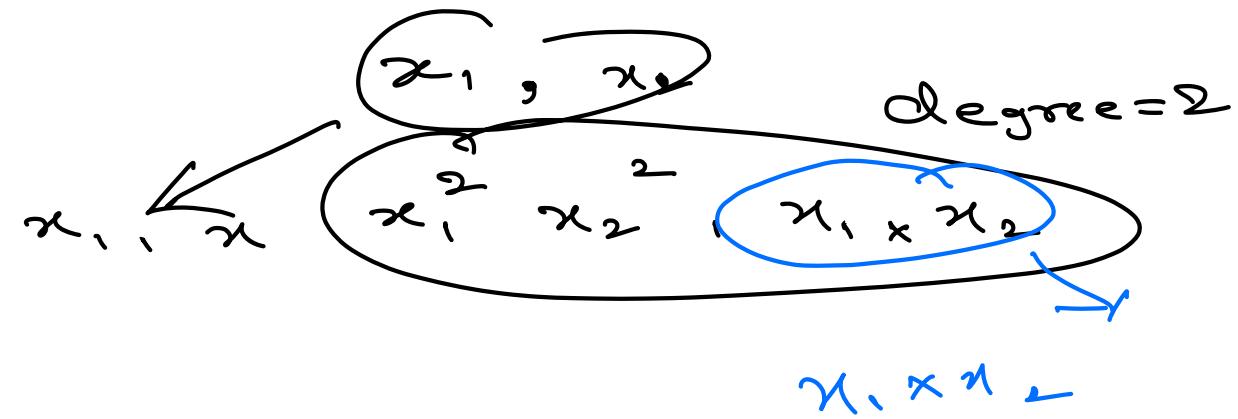
$$x_{\text{new}} \leftarrow x_i + \gamma \Delta(x_i, x_j)$$

Random Number  
between  
(0,1)

SMOTE  $\rightarrow$  Synthetic Minority Oversampling  
Technique



0.5



$x_1 \times x_2$