

Last class (Aug 19)

- 1) Spam vs Non-Spam: Business Case
- 2) Issue with Accuracy
- 3) Confusion Matrix & Code
- 4) Precision & Code
- 5) Recall & Code
- 6) F1 Score & Code

Today's class

- 1) Recap — Quizzes
- 2) Sensitivity and Specificity
- 3) ROC curve
- 4) AUC under ROC curve
- 5) Precision Recall Curve
- 6) Handling Imbalance Data
 - class weights
 - Oversampling of minority
 - Undersampling of majority
 - SMOTE

Movie rec

→ FP → A bad movie recommended

→ /FN → A good movie not recommended
✓

high precision

Malware detection

FP: good software predicted as Malware

→ FN: malware undetected (detected as good software)

TP, TN, FP, FN

equally
imp

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

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

$$\begin{aligned} \left. \begin{array}{l} \text{Sensitivity} \\ \text{Specificity} \end{array} \right\} &= \frac{TP}{TP + FN} = \text{Recall} = \text{True Positive Rate} \\ &= \frac{TN}{TN + FP} = 1 - \frac{FP}{FP + TN} \\ &= 1 - \text{False Positive Rate} \end{aligned}$$

		Pred	
		0	1
Ground	0	TN	FP
	1	FN	TP

Balanced Data

500 → class 0

500 → class 1

TP, FP, FN not sufficient

Bank Fraud Detection → class 0 class 1
1000 1000

Cancer Detection → class 0 class 1
600 800 ✓ ✓

Precision → $\frac{TP}{TP + FP}$

Recall → $\frac{TP}{TP + FN}$

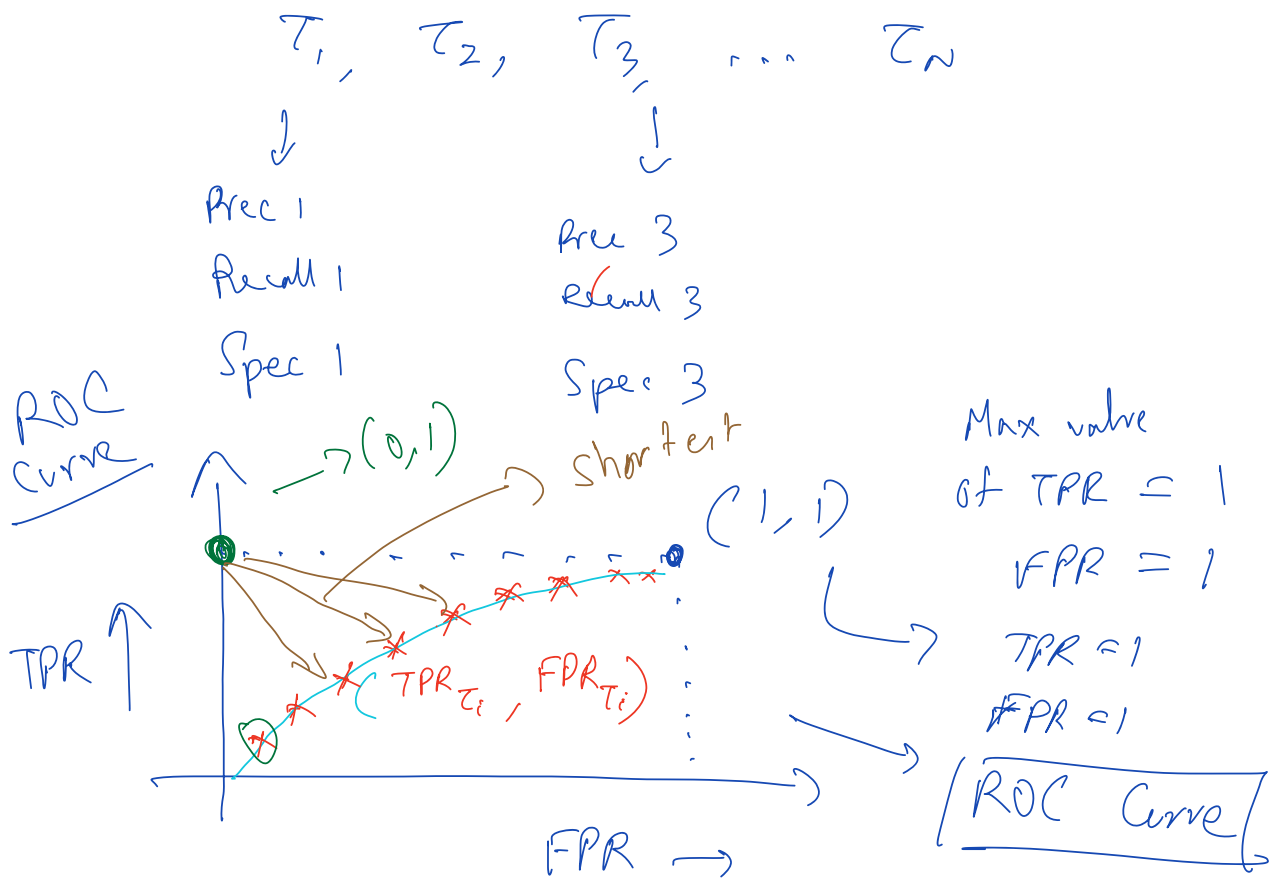
Acc = $\frac{TP + TN}{FP + FN + TP + TN}$

Req → [80% Precision, 90% Recall]

Logistic Regression

$$\sigma(\vec{z}) = P(y^{(i)} = 1 / x^{(i)}) = p$$

class 0 or class 1 → $p \geq 0.5 \rightarrow \text{class 1}$
 $p < 0.5 \rightarrow \text{class 0}$



Ideal model: $TPR = 1 = \frac{TP}{TP + FN} \rightarrow 0$

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

Model M1, Model M2

↓
Logistic
Regression

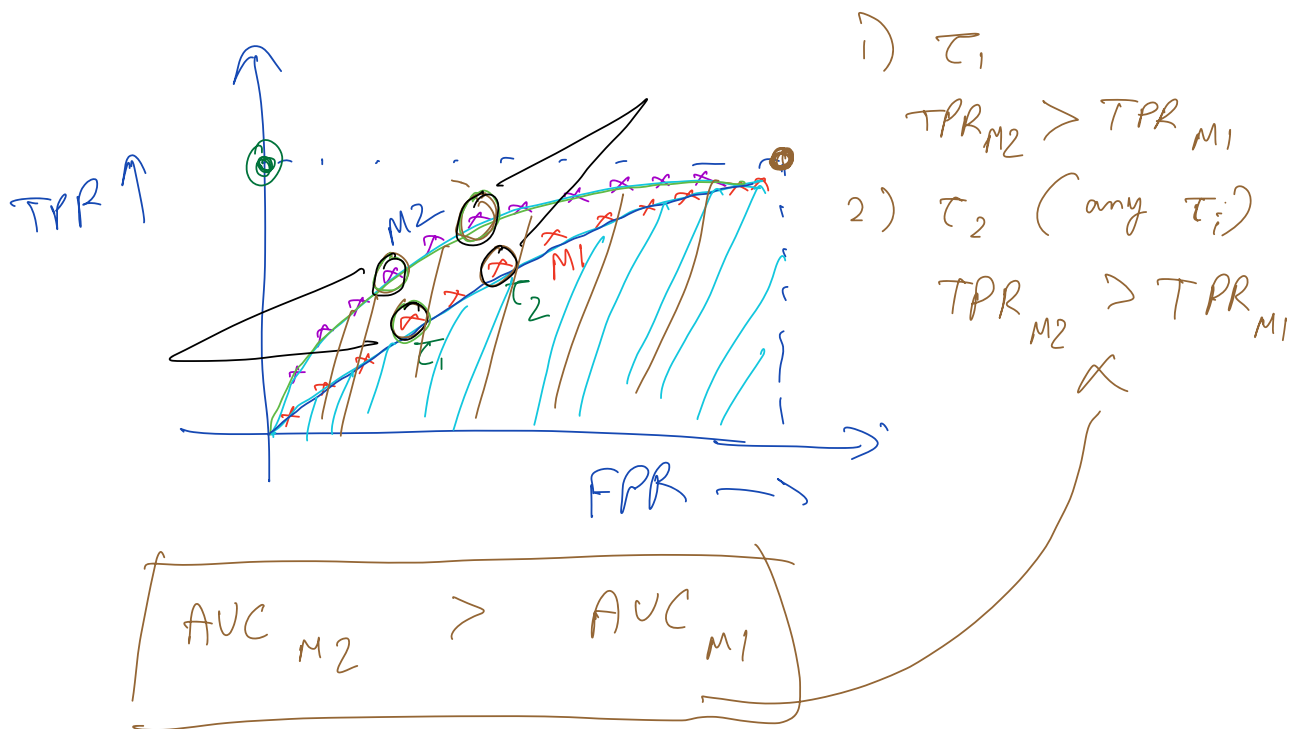
↓
Neural
network

$\tau = 0.5$

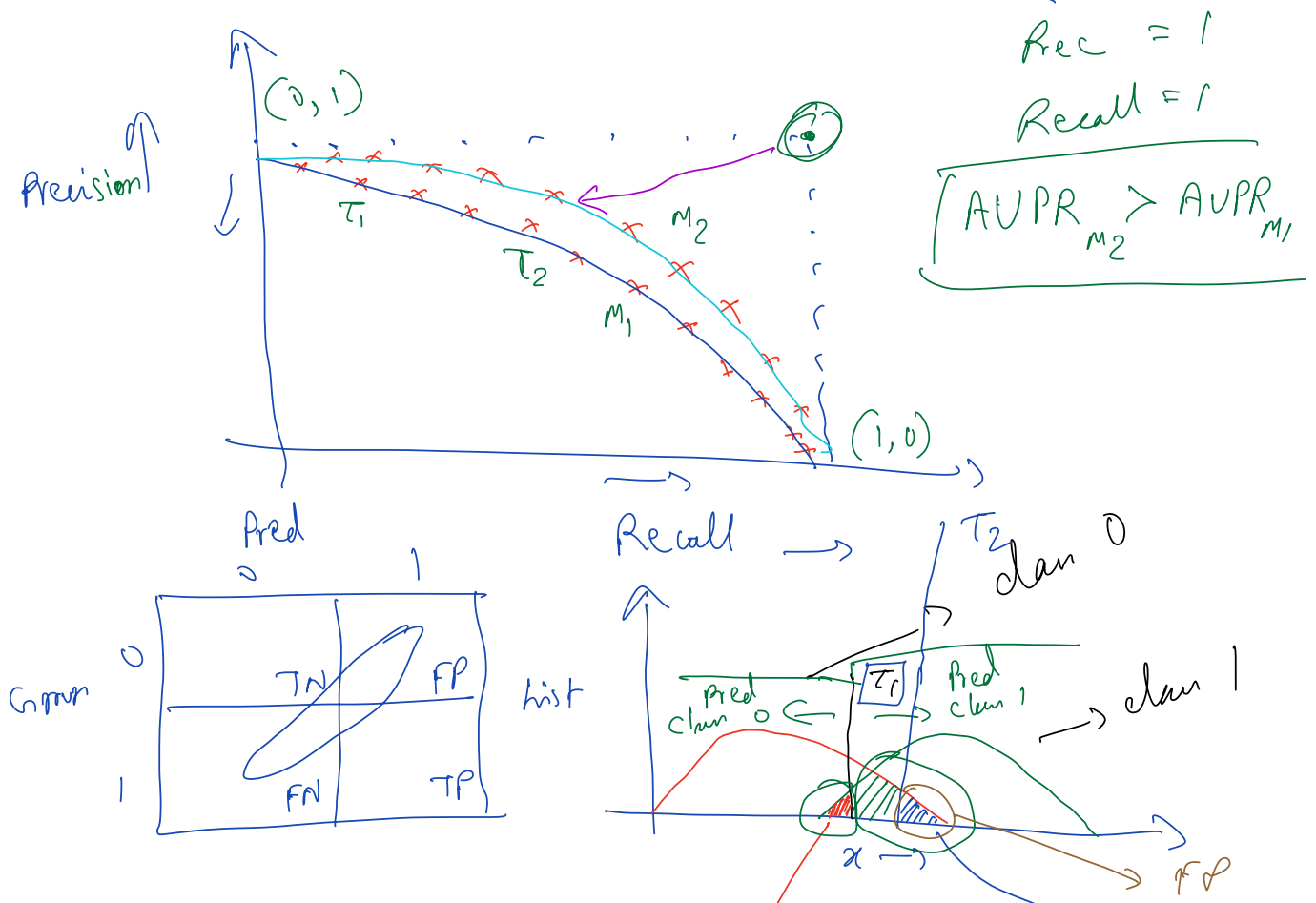
$Acc_1 = 0.8$

$Acc_1 = 0.8$

Balanced
Dataset



Imbalanced Dataset



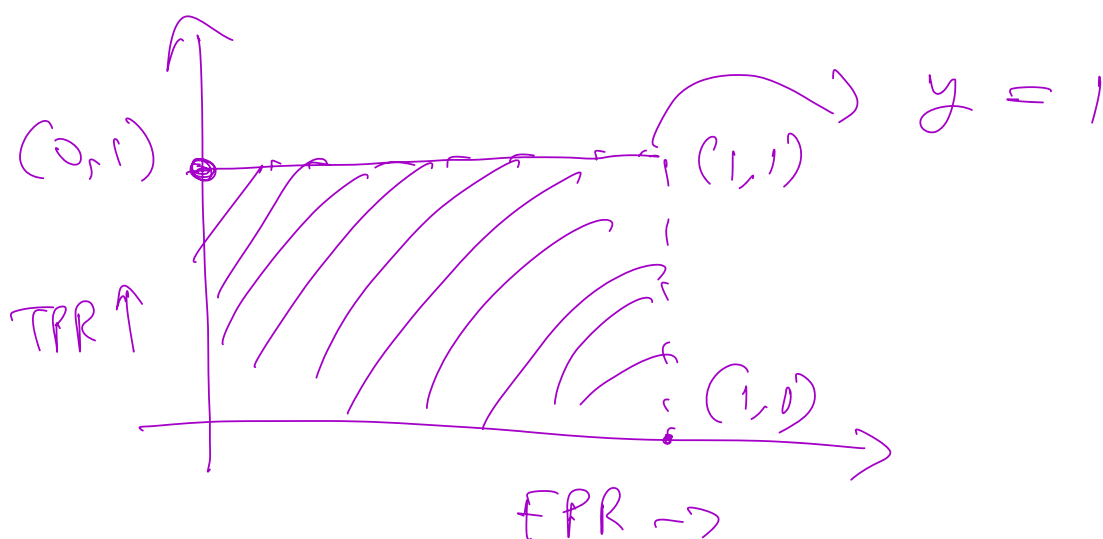
$\tau_1 : FN \ll FP$
 $\tau_2 : FP \ll FN$

$G.D : \text{class 1 } \begin{matrix} +ve \\ -ve \end{matrix}$
 $Pred : \text{class 0}$

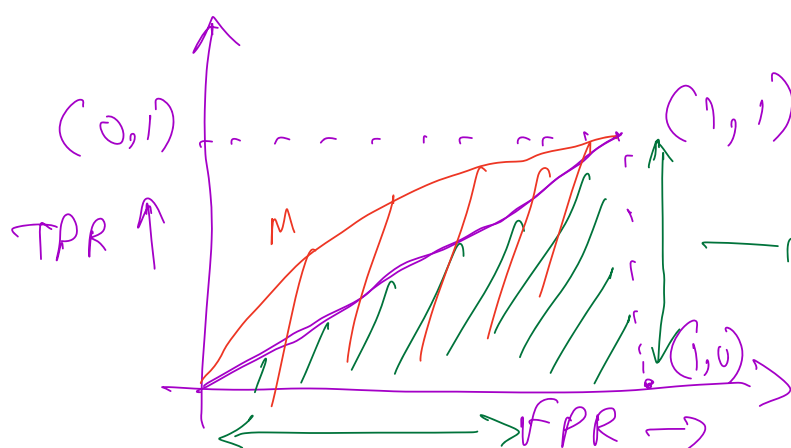
$G.D : \text{class 0}$
 $Pred : \text{class 1}$

FN
 FP

Ideal Model



50% +ve, 50% -ve | Random Model

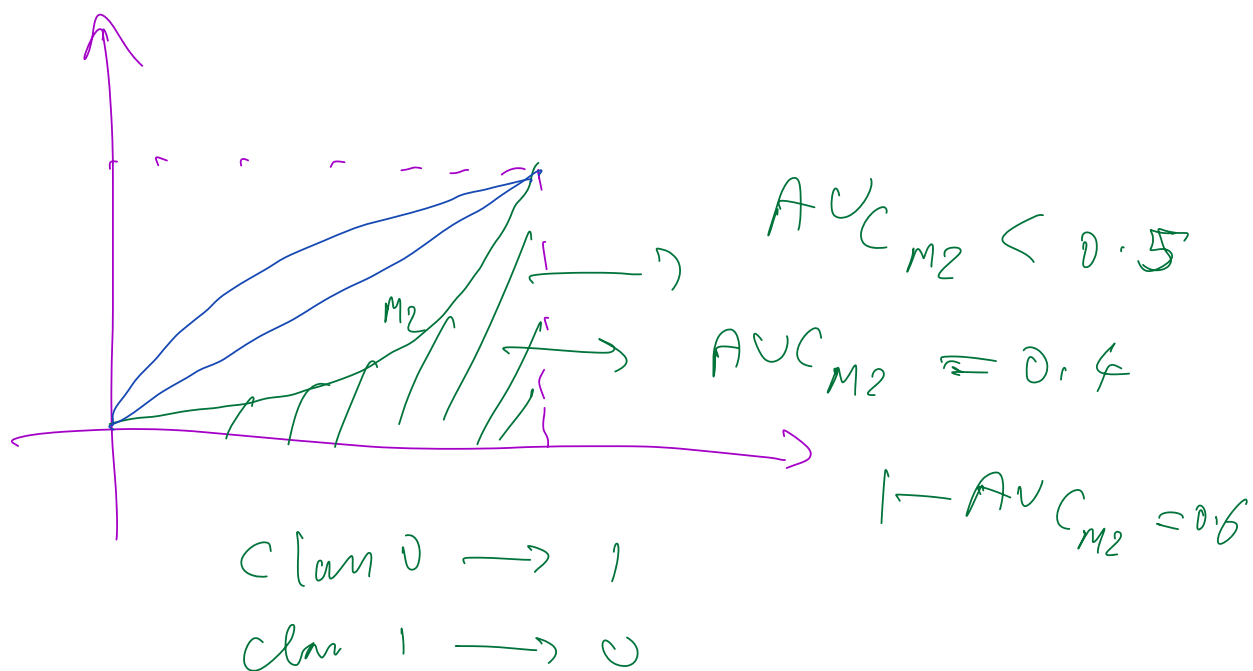


\downarrow
 class 0 : 50% prob
 class 1 : 50% prob

$$AUC = \frac{1}{2} \times 1 \times 1$$

$$= 0.5$$

$$AUC_M > 0.5$$



$$X = [x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}, x^{(5)}, x^{(6)}]$$

$$Y = [\underline{1}, \underline{1}, \underline{0}, \underline{1}, \underline{0}, \underline{0}] \rightarrow \begin{matrix} 3 \text{ +ve} \\ 3 \text{ -ve} \end{matrix}$$

$$P = [0.65, 0.94, 0.3, 0.92, 0.7, 0.2]$$

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

Sort them
descending order

X	Y	\hat{y}_{T_1}	\hat{y}_{T_5}	P
x ⁽²⁾	→ 1	1	1	0.94 = τ_1
x ⁽⁴⁾	→ 1	0	1	0.92 = τ_2
x ⁽⁵⁾	→ 0	0	1	0.7 = τ_3
x ⁽¹⁾	→ 1	0	1	0.65 = τ_4
→ x ⁽³⁾	→ 0	0	1	0.3 = τ_5
x ⁽⁶⁾	→ 0	0	0	0.2 = τ_6

sorted

$$T_1: \begin{array}{l} \text{pred +ve} = 1 \\ \text{pred -ve} = 5 \end{array}$$

TP	FP	TN	FN
1	0	3	2

$$TPR = \frac{1}{1+2} = 0.33$$

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

$$T_3: \begin{array}{l} \text{pred +ve} = 3 \\ \text{pred -ve} = 3 \end{array}$$

$$(0, 0.33)$$

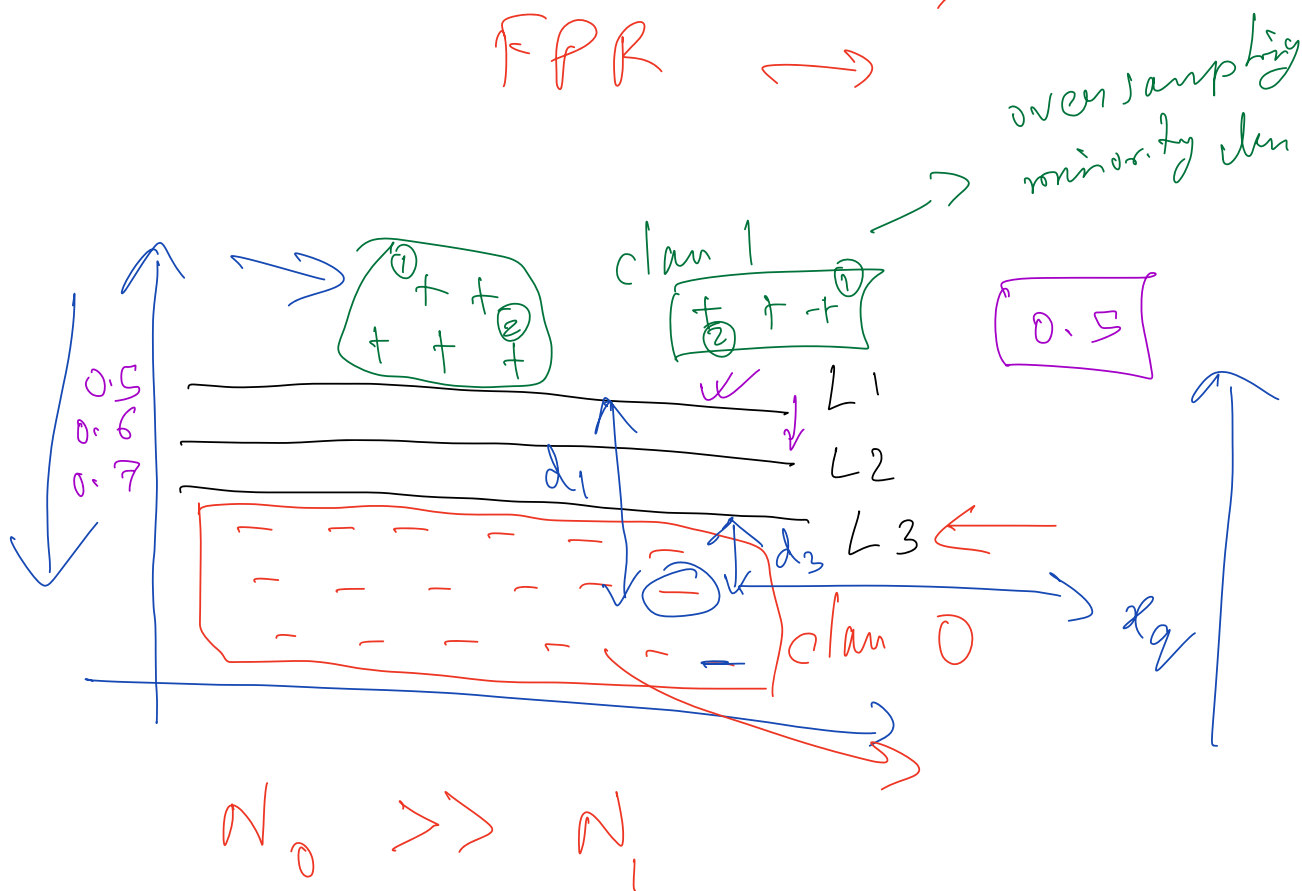
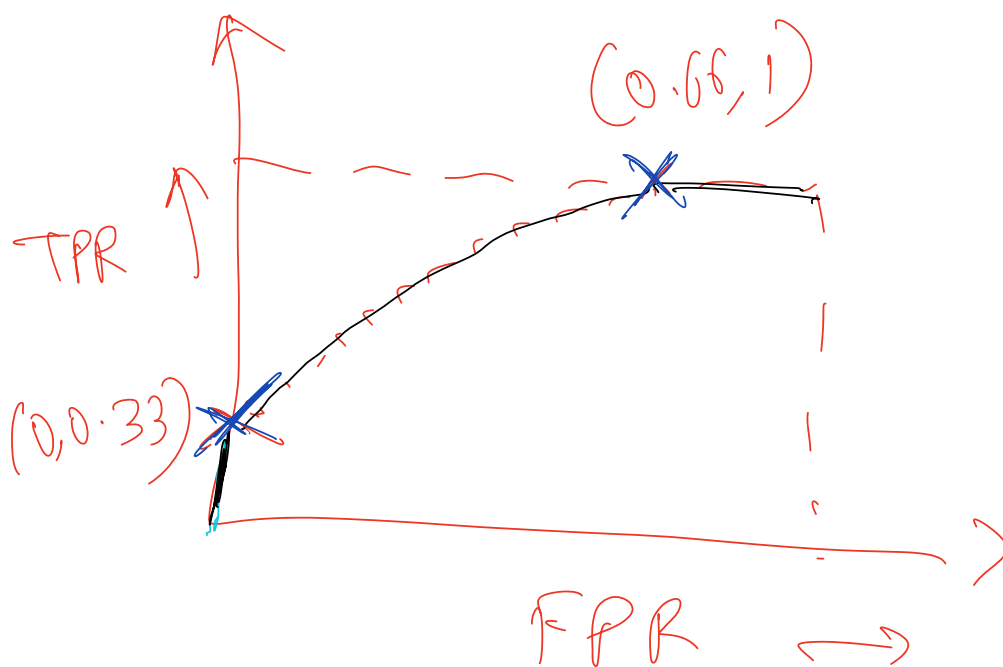
$$T_5: \begin{array}{l} \text{pred +ve} = 5 \\ \text{pred -ve} = 1 \end{array}$$

TP	FP	TN	FN
3	2	1	0

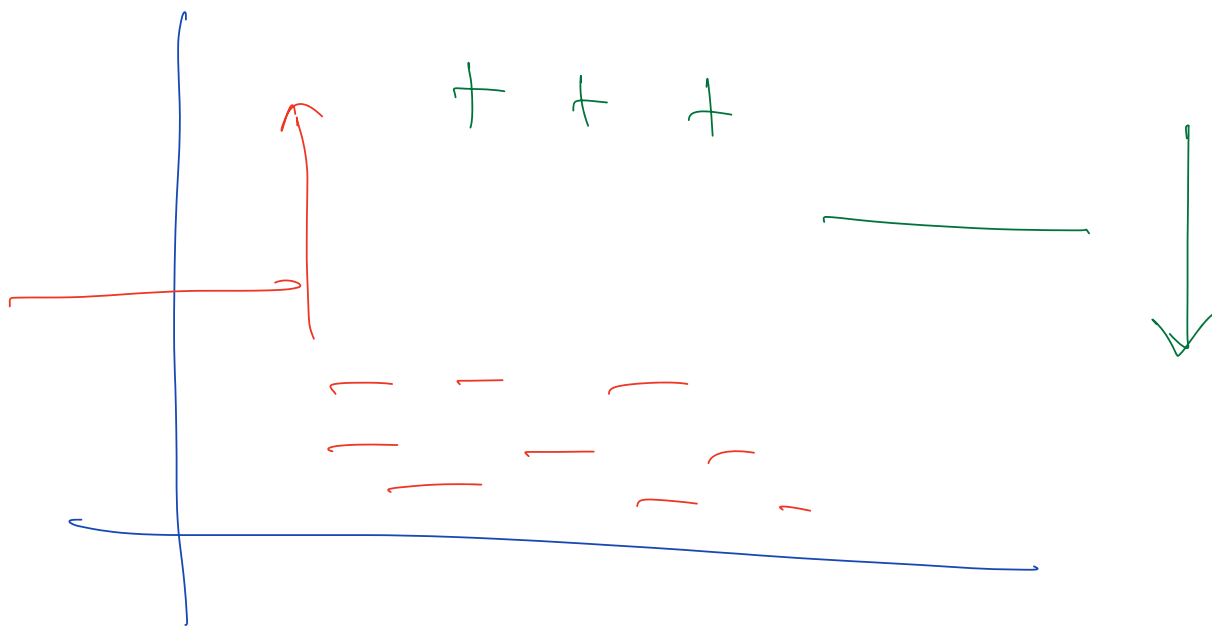
$$TPR = \frac{3}{3+0} = 1$$

$$FPR = \frac{2}{2+1} = 0.66$$

$$(0.66, 1)$$



$$\rightarrow d_1 > d_3 \quad C_1 > C_3$$



$$N_0 = 1000$$

$$N_1 = 100$$

$L(\) \rightarrow$ -ve sample $\rightarrow \frac{1}{10}$ or

$L(\) \rightarrow$ +ve sample \leftarrow

$$w = \frac{1}{10} \quad \left[\begin{array}{l} L_{+,1} \\ L_{+,2} \\ \vdots \\ L_{+,100} \end{array} \right] \rightarrow \left[\begin{array}{l} L_{-,1} \\ L_{-,2} \\ \vdots \\ L_{-,10} \end{array} \right] \quad w = 1$$

$$w_c = \frac{k}{N_c}$$

$$N_0 = 1000$$

$$N_1 = 100$$

$$\begin{cases} w_0 = \frac{k}{1000} \\ w_1 = \frac{k}{100} \end{cases}$$

$$w_0 + w_1 = 1$$

$$k \left(\frac{1}{1000} + \frac{1}{100} \right) = 1$$

$$\Rightarrow k \left(\frac{1100}{1000 \times 100} \right) = 1$$

$$N_0 = 100$$

$$N_1 = 1000$$

$$N_2 = 5000$$

$$w_0 = \frac{k}{N_0}$$

$$w_1 = \frac{k}{N_1}$$

$$w_2 = \frac{k}{N_2}$$

$$w_0 + w_1 + w_2 = 1$$

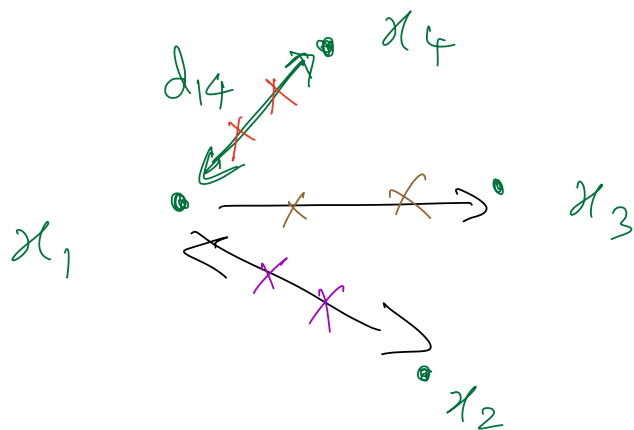
$$\text{log-loss} (L) = - \sum_{i=1}^m \left[y^{(i)} \log(\hat{y}^{(i)}) \right.$$

$$\left. + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$

Weighted log-loss

$$= - \sum_{i=1}^m \left[w_i y^{(i)} \log(\hat{y}^{(i)}) \right. \\ \left. + w_0 (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$

SMOTE



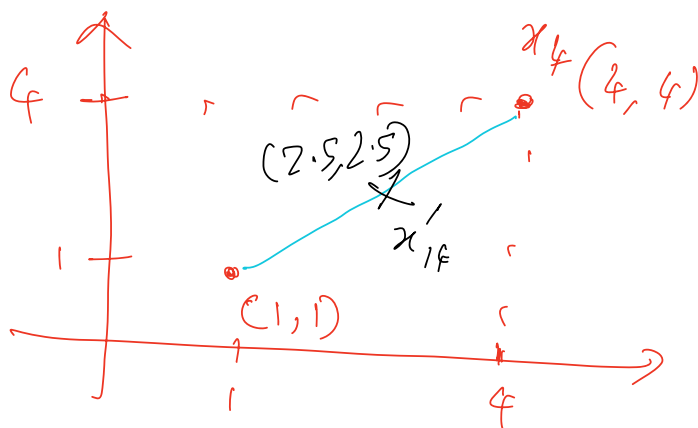
$$\epsilon = (0, 1)$$

$$0.4, 0.7$$

$$x'_{14} = x_1 + \epsilon \cdot d_{14}$$

$$x_1 = (1, 1)$$

$$x_4 = (4, 4)$$



$$\epsilon = 0.5$$

$$x'_{14} = x_1 + \epsilon \cdot d_{14}$$

$$= (1, 1) + 0.5 \times (4-1, 4-1)$$

$$= (1, 1) + 0.5 \times (3, 3)$$

$$= \left(1 + \frac{3}{2}, 1 + \frac{3}{2}\right)$$

$$= (2.5, 2.5)$$

360 +ve

40 -ve

$$\begin{aligned} P &= 1 \\ R &= 0.9 \end{aligned}$$

dumb model

TP TN FP FN

$$F1 = 0.947$$

$$Acc = 0.9$$

		Pred	
		0	1
GD	0	500 TN	10 FP
	1	10 FN	100 TP

$$Rec = \frac{TP}{TP + \textcircled{FN}} \rightarrow 0 = 1$$

$$Prec = \frac{TP}{TP + \textcircled{FP}} \rightarrow 1$$

$$\text{Prec} = \frac{100}{100+10} = \frac{10}{11} \approx 0.91$$

$$\text{Recall} = 0.91$$

$$\text{Acc} = \frac{600}{620} \approx 0.96$$

Balanced ^{data} is a sub-set of
Imbalanced data

If Prec & Recall is high,
 \Rightarrow Acc should be high.