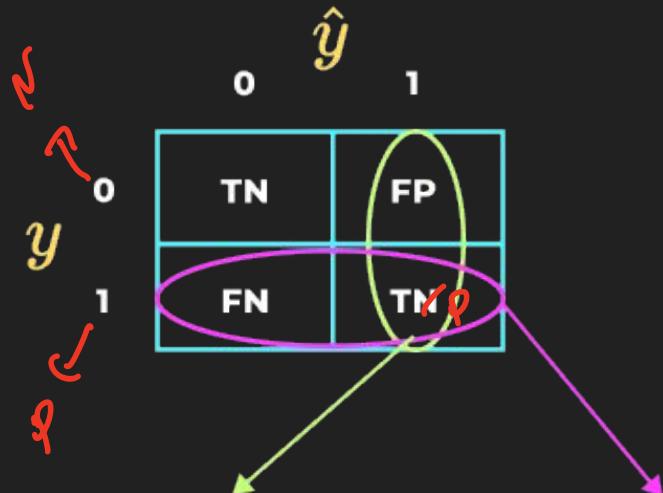


RECAP

Problem with Accuracy - Doesn't work for imbalance data

Confusion Matrix :



$$\text{Precision} = \frac{\text{Correct +ve Predictions}}{\text{All positive prediction}}$$

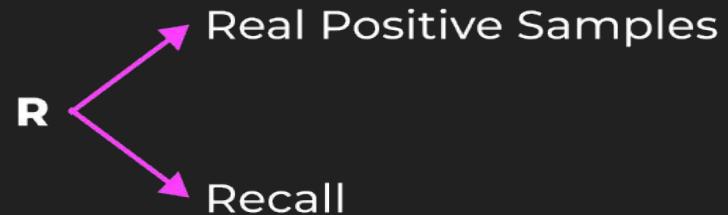
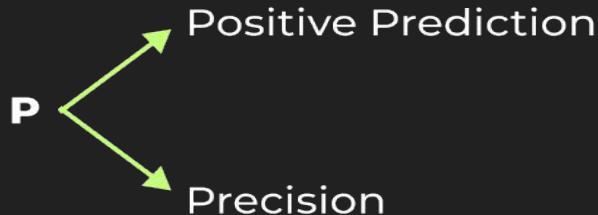
$$\frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{Correct +ve Predictions}}{\text{Real prediction}}$$

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



Quick way for leaving Precision & Recall :



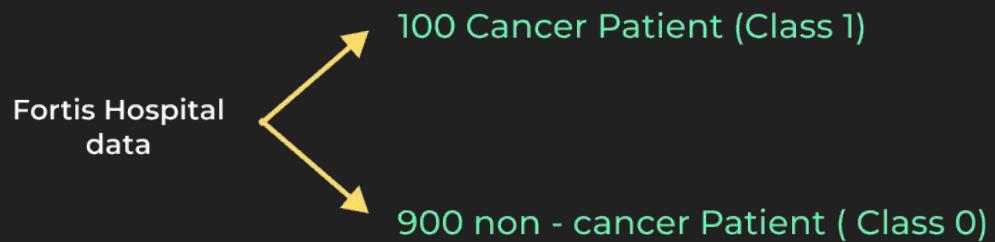
When FP & FN both
important.

F1_Score : Harmonic Mean between precision & recall

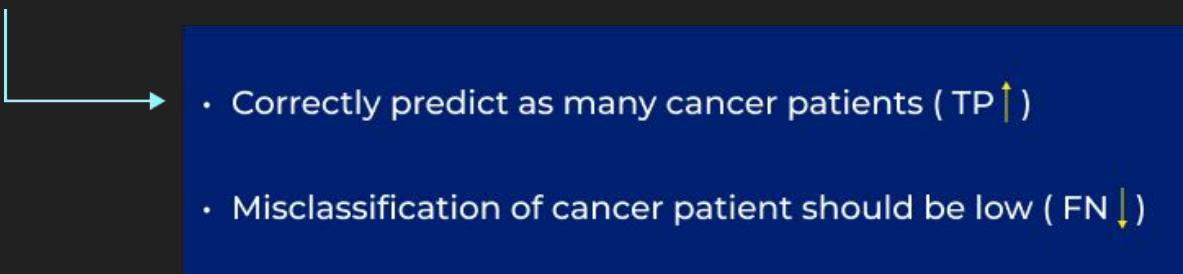
$$F1 = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2 \text{ (Precision)}(\text{recall})}{\text{Precision} + \text{Recall}}$$

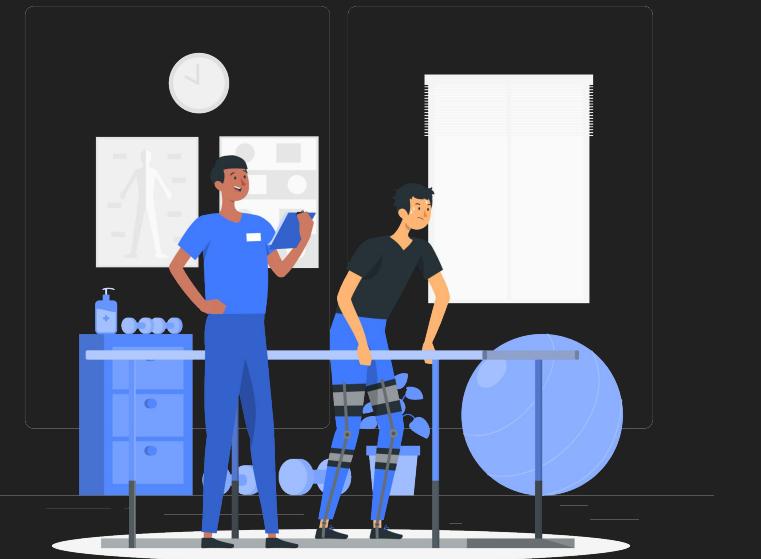
SENSITIVITY

Suppose you are a Data Scientist at Fortis :



AIM ? - Correctly identify all cancer patients

- 
- Correctly predict as many cancer patients ($TP \uparrow$)
 - Misclassification of cancer patient should be low ($FN \downarrow$)



This need of TP ↑ & FN ↓ is called Sensitivity. ≈ Recall

Sensitivity Important ? - Yes, since here failing of model to identify cancer patient can cause death

Formula ? -
$$\left[\frac{TP}{TP + FN} \right]$$

TP → $\frac{TP}{TP + FN}$

TPR → $\frac{TP}{TP + FN}$

Same as Recall

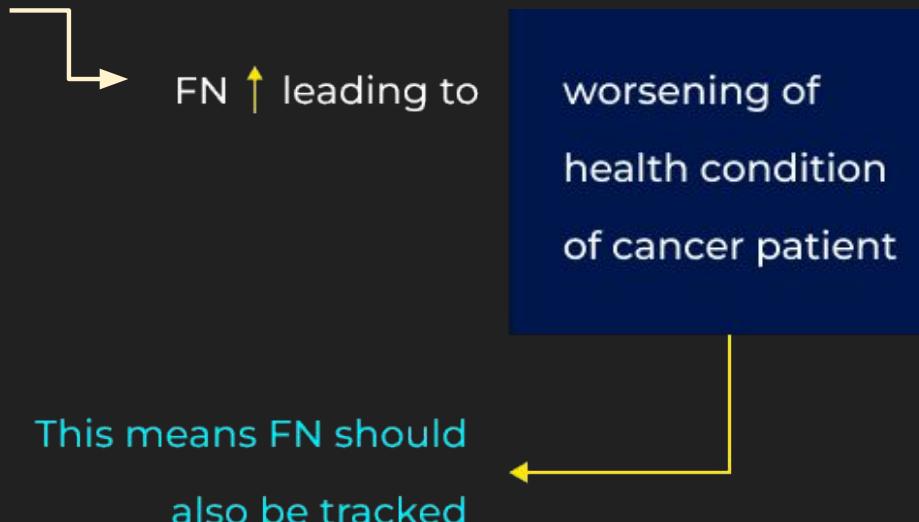
Tracks TP, hence also called TPR



$$TPR = \frac{\text{sensitivity}}{1 - \text{sensitivity}}$$

False Negative Rate (FNR) / Miss rate

What happens if model is insensitive ?



Formula ? - FNR = 1 - Sensitivity

$$\begin{aligned}
 &= 1 - \frac{TP}{TP + FN} \\
 &= \frac{FN}{TP + FN}
 \end{aligned}$$



As patients are missed by the MODEL,

How ? - Measure change in FN which is called as FNR.

- FNR is called Miss Rate.



SPECIFICITY

TP

TN

Do TN & FP not play any role in medical firms ?

TN & FP plays important role

Model should correctly predict non - cancer patient (TN ↑)

Model should keep misclassification of non - cancer patient low (FP ↓)

Need of TN ↑ & FP ↓ is called **Specificity**

Analogy : Specificity is sensitivity for negative class (CLASS 0)

Specificity Important ?

- High specificity avoids
 - 1. Fruitless experience treatments
 - 2. Reduces social stigma anxiety for a non - cancer patient



Formula ? - Specificity =

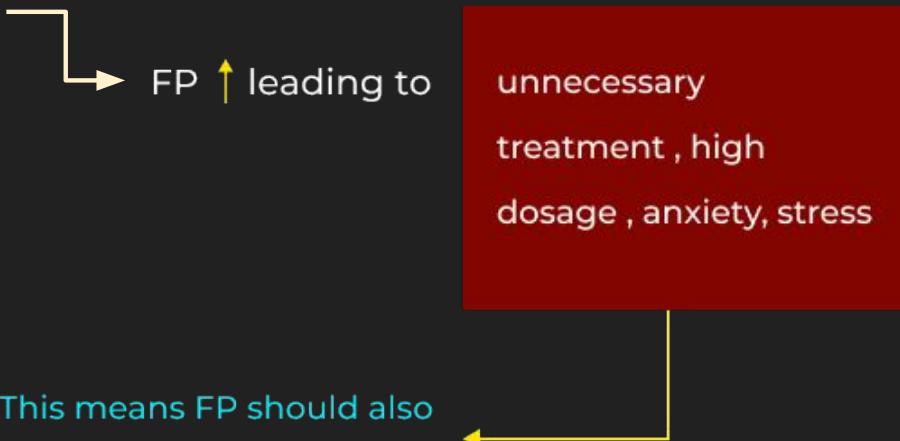
$$\frac{TN}{TN + FP}$$

→ Tracks TN , Hence called TNR

→ **TNR**

False Positive Rate (FPR)

What happens if model is not specific ?



Formula ? - FPR = **1 - Specificity**

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

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

How to track ? - Using FPR to measure change in FP

Points to Remember

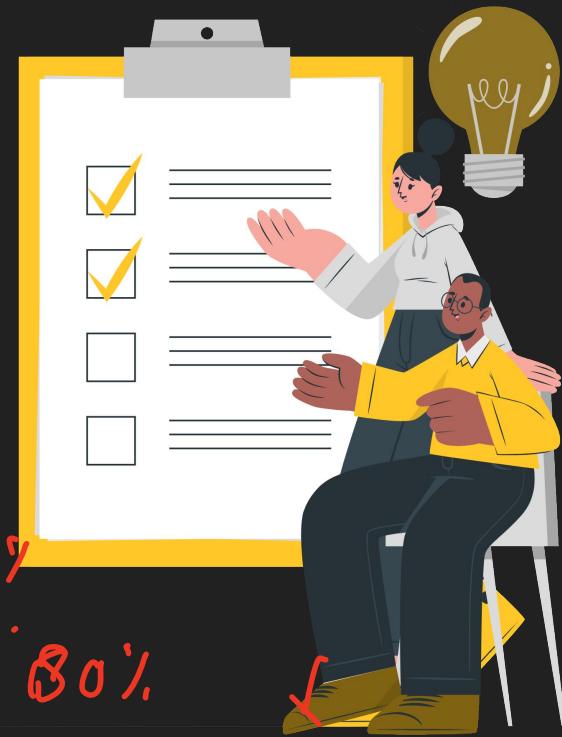
Sensitivity : TP ↑ & FN ↓ = $\frac{TP}{TP + FN}$ $\approx .92\%$

FNR : tracks FN = 1 - Sensitivity

Specificity : TN ↑ & FP ↓ = $\frac{TN}{TN + FP}$
(Sensitivity for Negative class)

FPR : tracks FP = 1 - Specificity

Accuracy : .80%
99%
70%
780%



RECEIVER OPERATING CURVE

How to improve performance of our model ?

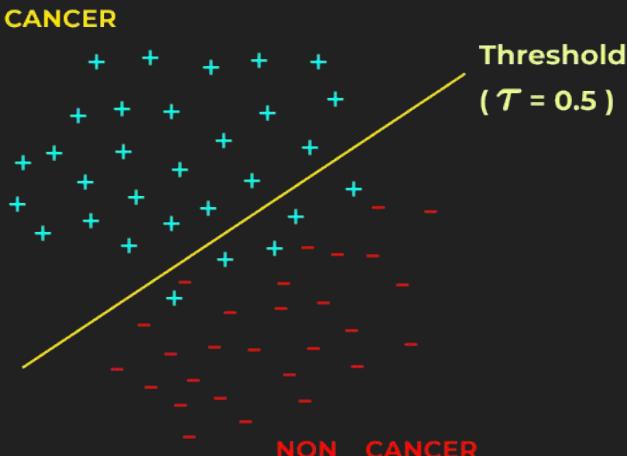


HYPERPARAMETER TUNING

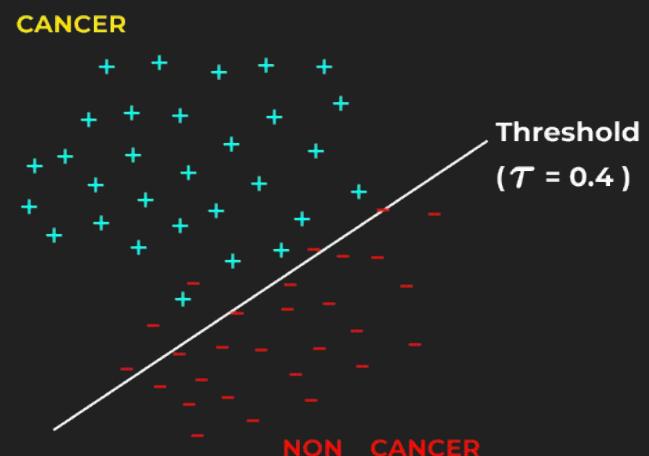


Find the correct threshold (τ)

- (the model misclassifies cancer patients
- need of a correct threshold)



Find the correct
threshold (τ)



How to find the correct threshold τ ?

- Considering Validation data

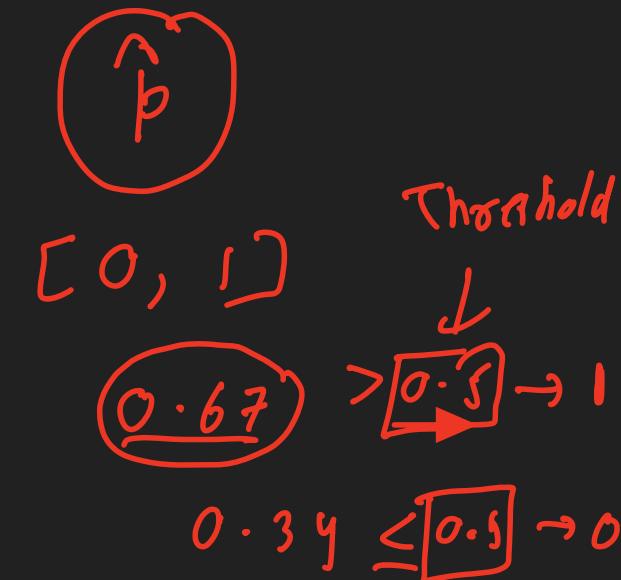
Step : 1

X	Y	$P = P(Y^{(i)} = 1 X^{(i)})$
$x^{(1)}$	$y^{(1)}$	$p^{(1)}$
$x^{(2)}$	$y^{(2)}$	$p^{(2)}$
$x^{(n)}$	$y^{(n)}$	$p^{(n)}$

↑
n
↓

Sort samples in
descending order
based on value $P^{(i)}$

Complexity : $O(n \log n)$



Step : 2 Use every $p^{(i)}$ as threshold ($\tau^{(i)}$) for predicting $\hat{y}^{(i)}$

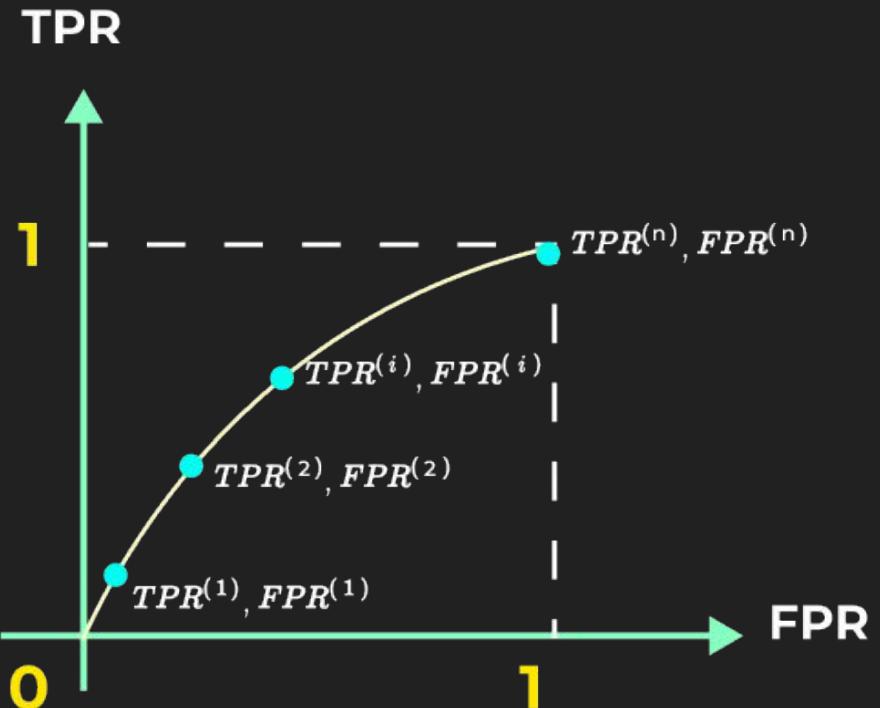
X	Y	P	$\hat{y}_\tau = p^{(1)}$	$\hat{y}_\tau = p^{(2)}$	$\hat{y}_\tau = p^{(n)}$
$x^{(1)}$	$y^{(1)}$	$p^{(1)}$	1	1	1
$x^{(2)}$	$y^{(2)}$	$p^{(2)}$	0	1	1
.
.
$x^{(n)}$	$y^{(n)}$	$p^{(n)}$	0	0	1

Based on $\hat{y}_\tau = p^{(i)}$, we find TPR & FPR

P	TPR	FPR
$p^{(1)}$	$TPR^{(1)}$	$FPR^{(1)}$
$p^{(2)}$	$TPR^{(2)}$	$FPR^{(2)}$
.	.	.
.	.	.
$p^{(n)}$	$TPR^{(n)}$	$FPR^{(n)}$

$$TPR = \frac{TP}{TP + FN}, \quad FPR = \frac{FP}{FP + TN}$$

Step : 3 - Plot TPR vs FPR



In spam classifier, we want TP to increase and FP to decrease

Therefore , pick hat threshold where TPR increases & FPR decreases

Understanding ROC with example

Step :1

X	Y	P
$x^{(1)}$	1	0.65
$x^{(2)}$	1	0.94
$x^{(3)}$	0	0.30
$x^{(4)}$	1	0.92
$x^{(5)}$	0	0.70
$x^{(6)}$	0	0.20

Sort according
to P

X	Y	P
$x^{(2)}$	1	0.94
$x^{(4)}$	1	0.92
$x^{(5)}$	0	0.70
$x^{(1)}$	1	0.65
$x^{(3)}$	0	0.30
$x^{(6)}$	0	0.20

Step : 2 Create $\hat{y}_\tau = p^{(i)}$

X	Y	P	$\hat{y}_\tau = 0.94$	$\hat{y}_\tau = 0.92$	$\hat{y}_\tau = 0.20$
$x^{(2)}$	1	0.94	1	1	1
$x^{(4)}$	1	0.92	0	1	1
$x^{(5)}$	0	0.70	0	0	1
$x^{(1)}$	1	0.65	0	0	1
$x^{(3)}$	0	0.30	0	0	1
$x^{(6)}$	0	0.20	0	0	1

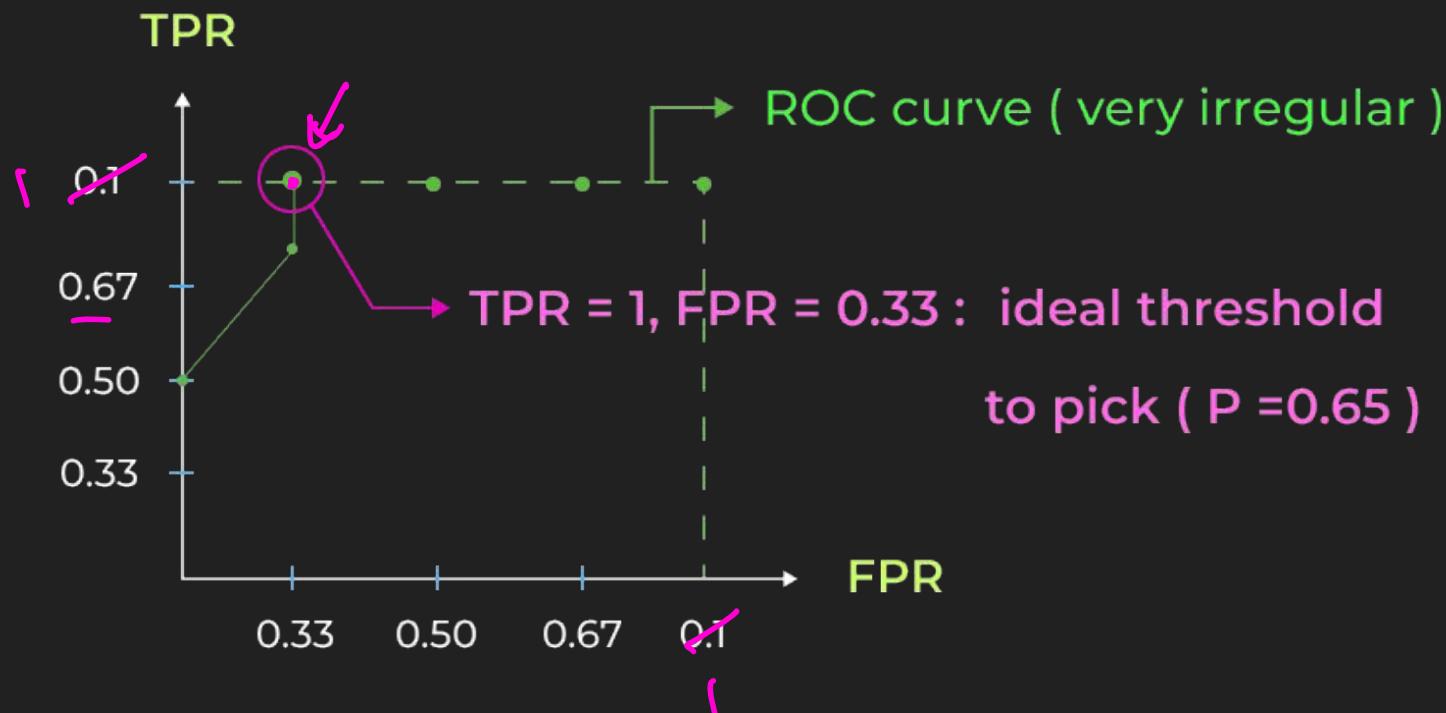
Find TPR , FPR - TP = 1 , FP = 1, FN = 2 , TN = 3

$$\therefore TPR = \frac{1}{1+2} = \frac{1}{3} = 0.33 \quad FPR = \frac{0}{(0+3)} = 0$$

Repeating step - 2 for each P :

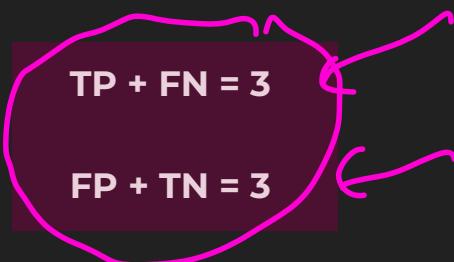
P	TPR	FPR
0.94	0.33	0
0.92	0.50	0
0.70	0.67	0.33
0.65	1	0.33
0.30	1	0.67
0.20	1	1

Pair of TPR & FPR for all 6 probabilities



What will be the ROC curve for Random Model ?

As our example data is **BALANCED**



3 +ve data
3 -ve data

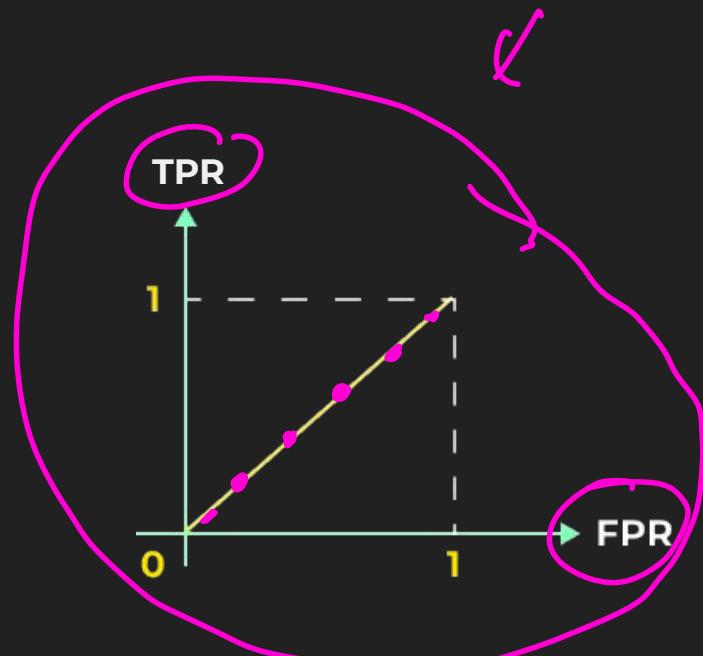
Also as Random Model generates (1, 0) randomly

Produces K TP & K FP points

$$TPR = \frac{K}{3}, \quad FPR = \frac{K}{3}$$

Hence, TPR = FPR

(y - axis) = (x - axis)

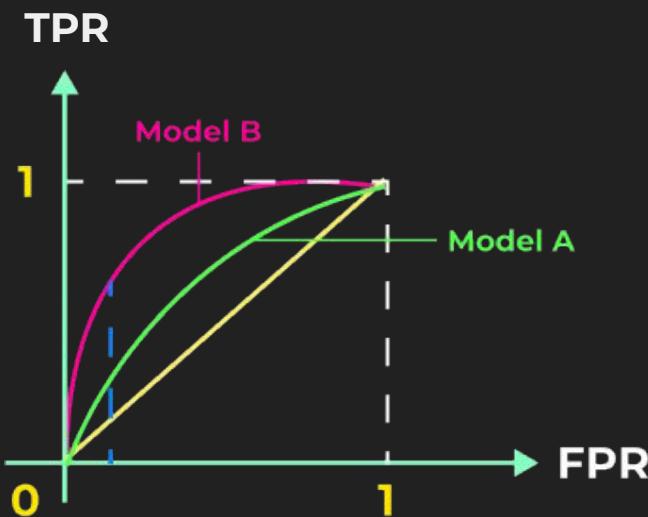


Suppose we have 2 models

Model A

Model B

Which is better ?

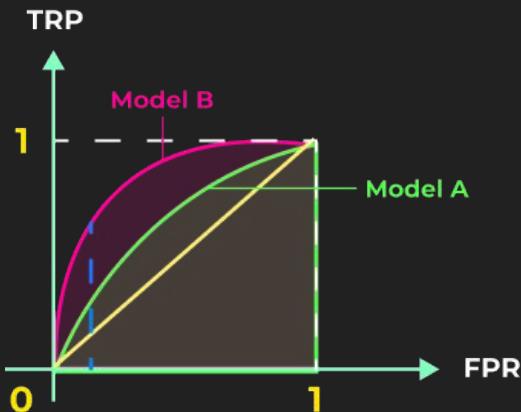


Intuitively,

$\text{TPR}(\text{B}) > \text{TPR}(\text{A})$ for same FPR.

Hence , MODEL B is better

How to mathematically show Model B better ?



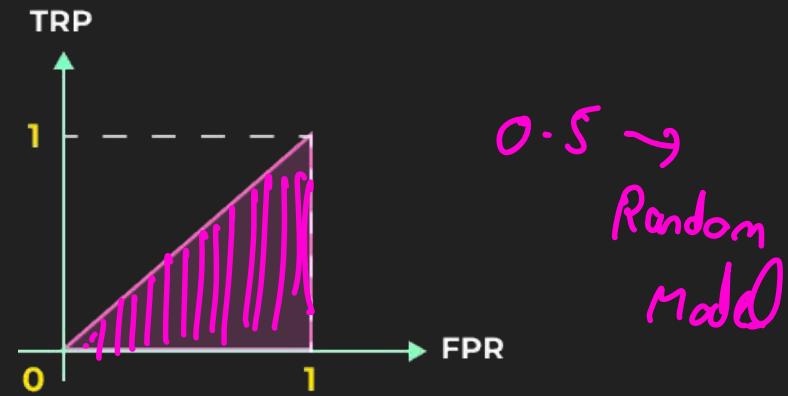
Let's take area under the ROC curve (All - ROC / Auc)

Clearly ,

Auc Model B > Auc Model A

Hence, Model B better

What will be Auc - Random Model ?



$$Auc = \frac{1}{2} * 1 * 1$$

$$\int = 0.5$$

Area under curve
(AUC)

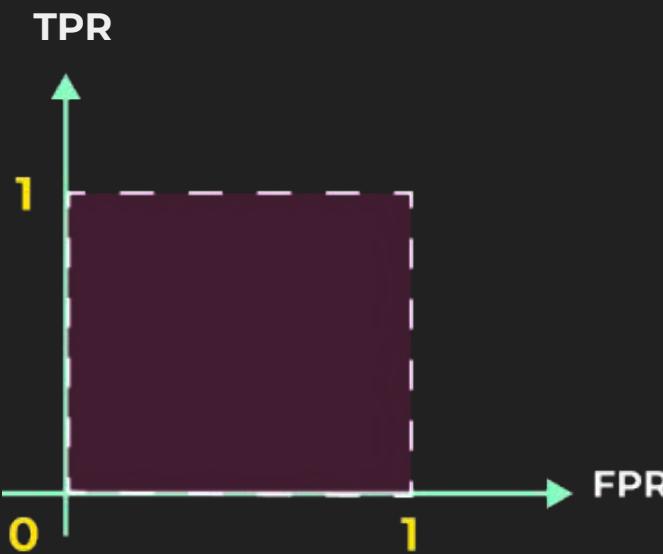
For ideal model

threshold ≈ 1 (*most* $\hat{y} = 0$)

$\therefore TPR \approx 1 \text{ & } FPR \approx 0$

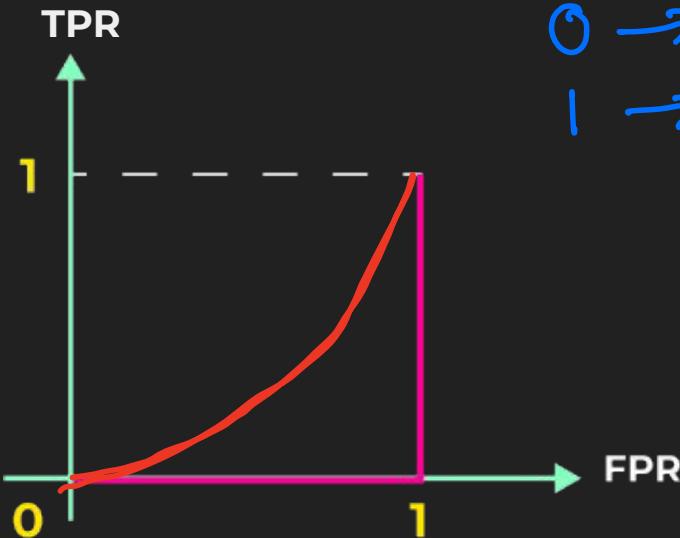
threshold ≈ 0 (*most* $\hat{y} = 1$)

$\therefore TPR \approx 1 \text{ & } FPR \approx 1$





What will be Auc - bad Model ?



$0 \rightarrow 1$
 $1 \rightarrow 0$
 $< 0.4 \rightarrow 0$
 $\geq 0.4 \rightarrow 1$

Would $(1 - p)$ fix bad model ?

$$AUC_{new} = 1 - AUC_{old}$$

$$= 1 - 0$$

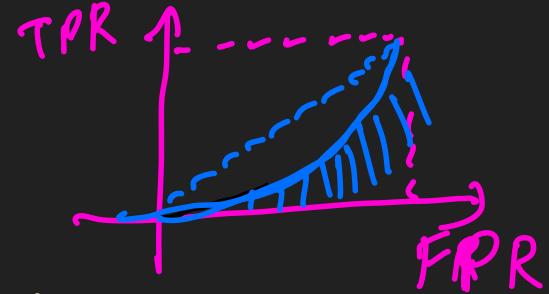
$$AUC_{new} = 1$$

Hence, for any model which has $ROC < 0.5$

Means bad model probabilities misclassified
every datapoint

Break: 8:11 AM

Can be fixed by
revising probability



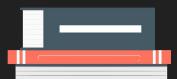
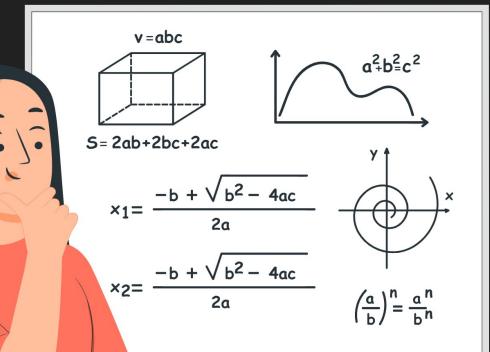
Au - ROC

We find AUC metric on every
possible threshold

Precision , Recall , F1

All find metric on a single threshold
(generally, $\mathcal{T} = 0.5$)

How Au - ROC different from Precision / Recall / F1 ?



Issue with Au - ROC curve :

Suppose $y = [1, 1, 0, 1, 1] \rightarrow$ Imbalance Data

$$P_{m1} = [0.95, 0.92, 0.80, 0.76, 0.71]$$

$$P_{m2} = [0.20, 0.10, 0.08, 0.06, 0.02]$$

then what will be \hat{y}_{m1} & \hat{y}_{m2} ?

M1

y	P_{m1}	$\hat{y}_\tau = 0.95$	$\hat{y}_\tau = 0.92$
1	0.95	1	1
1	0.92	0	1
0	0.80	0	0
1	0.76	0	0
1	0.71	0	0

M2

y	P_{m_2}	$\hat{y}_\tau = 0.2$	$\hat{y}_\tau = 0.1$
1	0.2	1	1
1	0.1	0	1
0	0.08	0	0
1	0.06	0	0
1	0.02	0	0

M1 : 0.95 > 0.92 > 0.80 > 0.76 > 0.71

M2 : 0.2 > 0.1 > 0.08 > 0.06 > 0.02

$$[TPR_{M1}^{(1)}, FPR_{M2}^{(1)}] \dots [TPR_{M1}^{(n)}, FPR_{M2}^{(n)}]$$

=

$$[TPR_{M2}^{(1)}, FPR_{M2}^{(1)}] \dots [TPR_{M2}^{(n)}, FPR_{M2}^{(n)}]$$

Hence, Auc(M1) = Auc(M2)

Observe

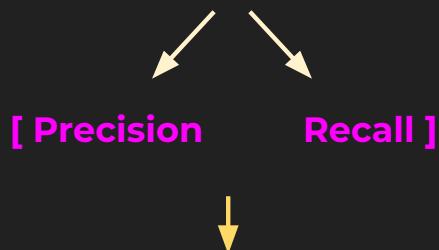
- AU-ROC considers how data is ordered rather than the value itself

Hence AUC(M1) and AUC(M2) same

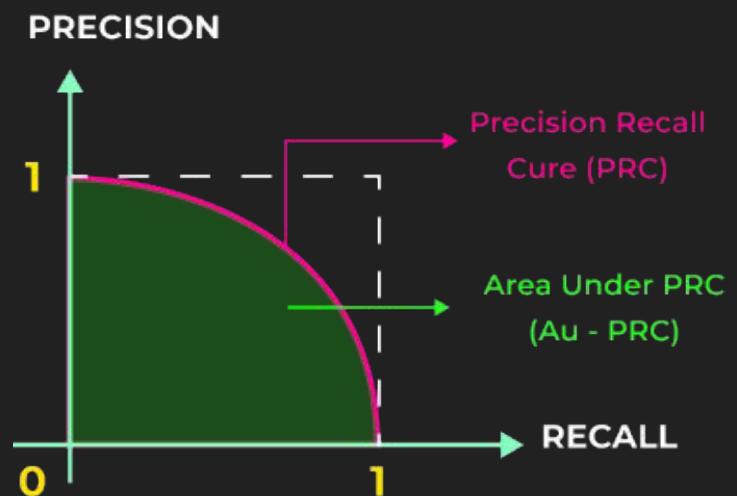
PR - ROC Curve

What can be used instead of Au - ROC when data imbalanced ?

A F1 score works well for imbalance data:



Let's take **Precision and Recall** values for every probability instead of **TPR and FPR**



Points to Remember

Sensitivity : TP \uparrow & FN \downarrow =
$$\frac{TP}{TP + FN}$$

FNR : tracks FN = 1 - Sensitivity

Specificity : TN \uparrow & FP \downarrow =
$$\frac{TN}{TN + FP}$$

(Sensitivity for Negative class)

FPR : tracks FP = 1 - Specificity

- Au - ROC for random model = **0.5**
- Au - ROC for ideal model = **1**
- Au - ROC does **not work in imbalance setting**
- Au - ROC depends on **order of probabilities.**

Understanding Imbalance Data



50 % - 50% → **Data Balanced**

60 % - 40 % → **Slightly Balanced**

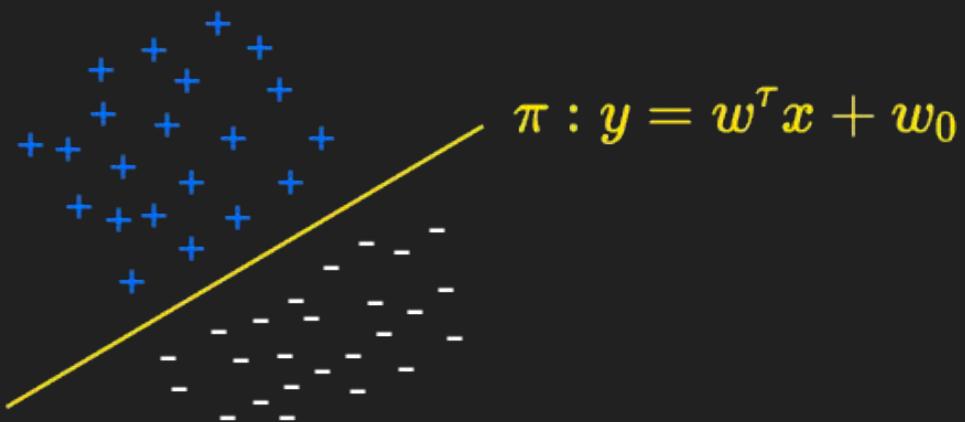
70 % - 30% → **Slightly Imbalanced**

80 / 90 % - 20 / 10 % → **Data Imbalanced**

95 % - 5 % → **Extremely Imbalanced**



How does Logistic - Regression look when data is balanced ?



Logistic Regression creates
a hyperplane



Which separates the data perfectly
by finding optional weights (w_0, w)
using Gradient Descent

What happens to logistic Regression when data is imbalance ?

$$\text{logloss} = -\frac{1}{n} \left[\sum_{i=1}^n y^{(i)} \log(p^{(i)}) + (1 - y^{(i)}) \log(1 - p^{(i)}) \right]$$

↓

$$p^{(i)} = P(y^{(i)} = 1 | \mathbf{x}^{(i)})$$

Suppose $p^{(i)} = 0.7$ Means model 70% confident that $\mathbf{x}^{(i)}$ has label = 1

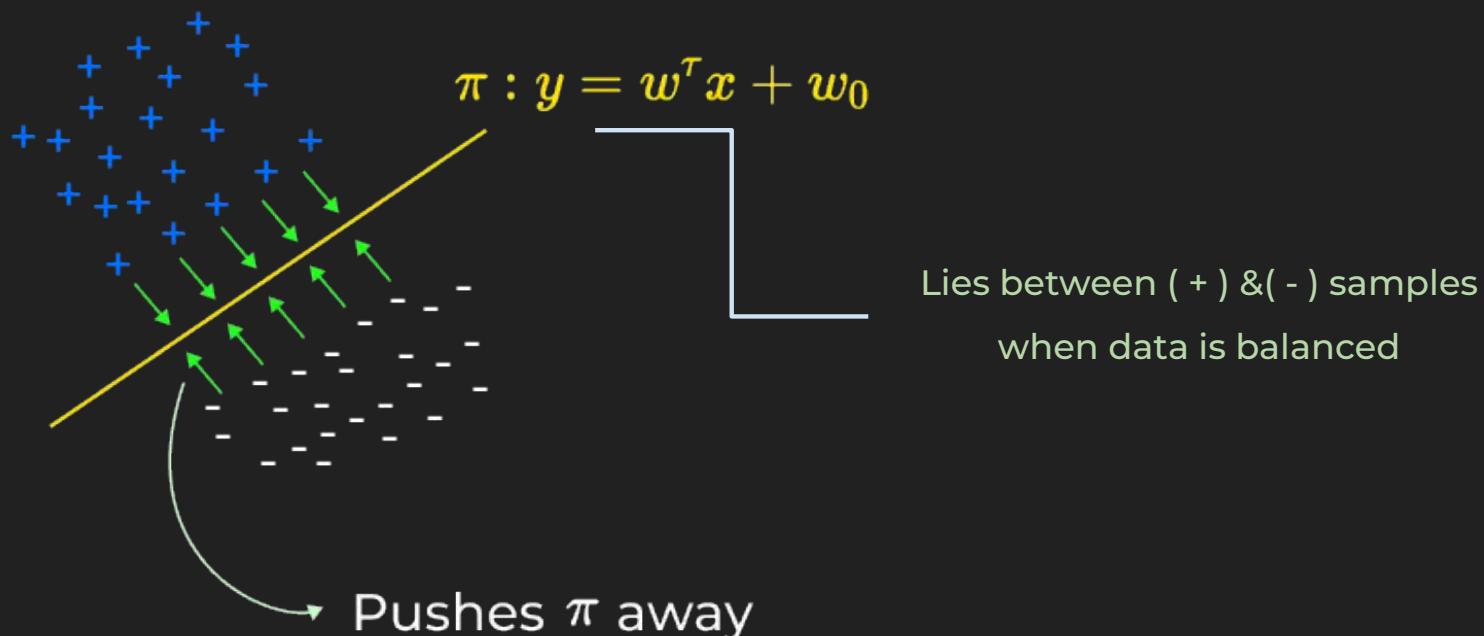
$$p^{(1)} = 0.5] \rightarrow \text{logloss} = -[(1 - 0) \log(1 - 0.5)] = 0.301$$

$$p^{(2)} = 0.1] \rightarrow \text{logloss} = -[(1 - 0) \log(1 - 0.1)] = 0.045$$

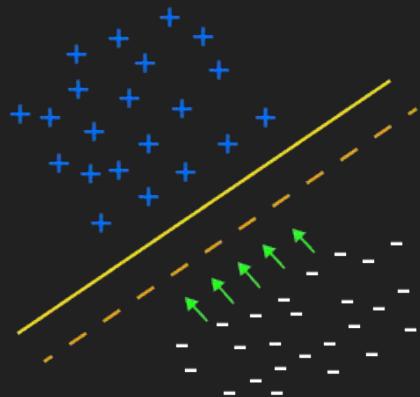
We can conclude, as logloss small for $p^{(2)}$, model is confident $p^{(2)} \neq \text{Class1}$

Hence , to conclude

Model is confident \rightarrow point is further away from it



Logistic Regression



Hyperplane shifts towards minority class

Uncertainty when predicting + ve class data

How to make imbalance data balanced ?

Suppose → **1000 data samples**

150 SPAM

850 NON - SPAM



Class Weight

∴ Non - Spam 5.67 times Spam

If 1 spam data has weightage of 5.67 non - spam

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

$W_i = 5.67$ when spam

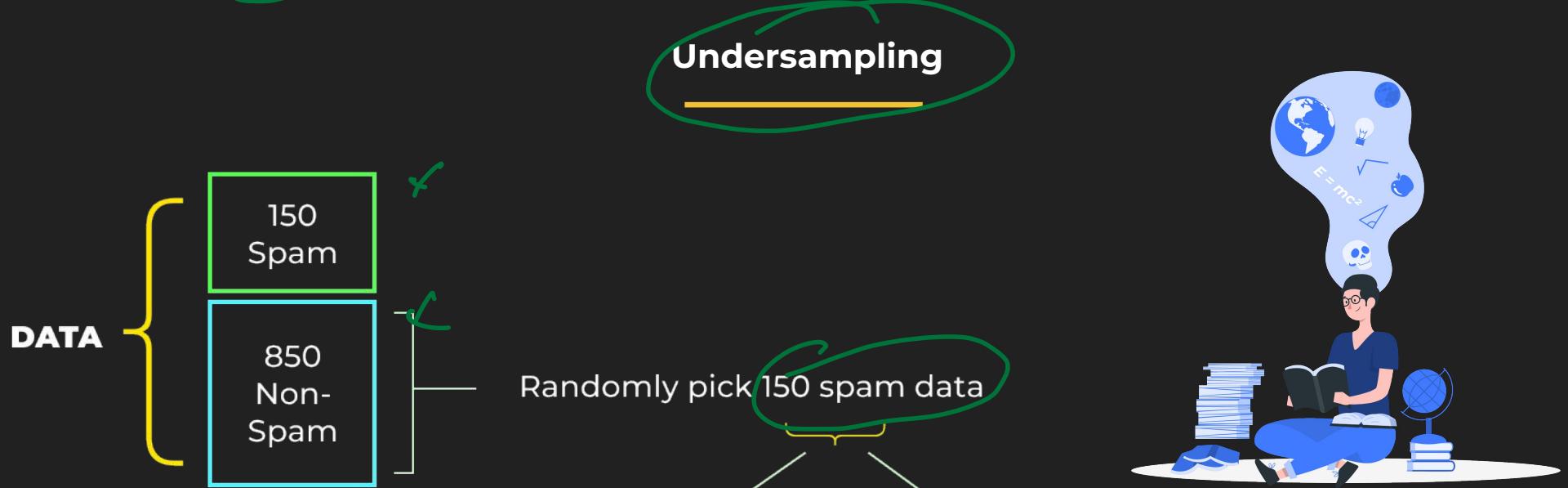
$W_i = 1$ when non - spam

$$w_i \cdot (-) + \text{---} \cdot w_j$$

$$\begin{array}{r} 1000 \\ 1000 \\ \hline 900 \\ \Rightarrow 9 \end{array}$$



$$\begin{array}{l} L_L \Rightarrow \\ 0 \rightarrow \text{---} \\ | \rightarrow \text{---} \end{array}$$



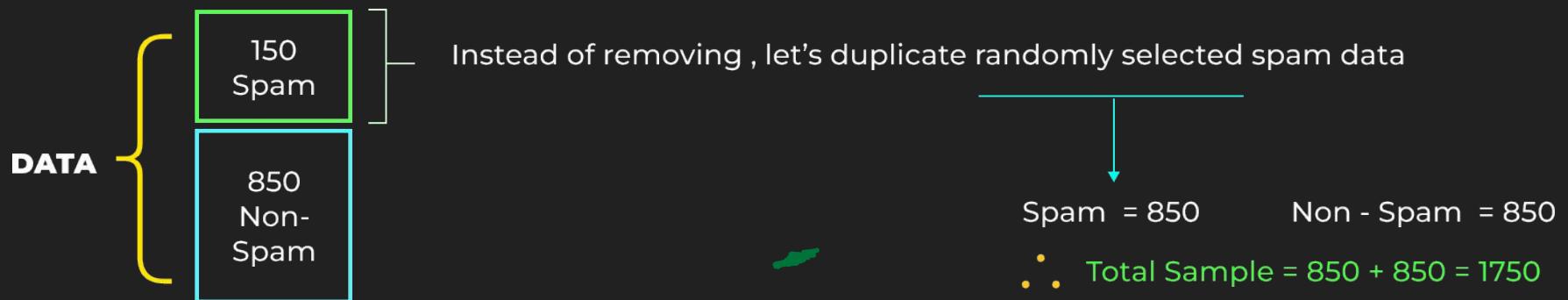
As model trained on smaller data

Model Reliability Decreases

When to use undersampling ?

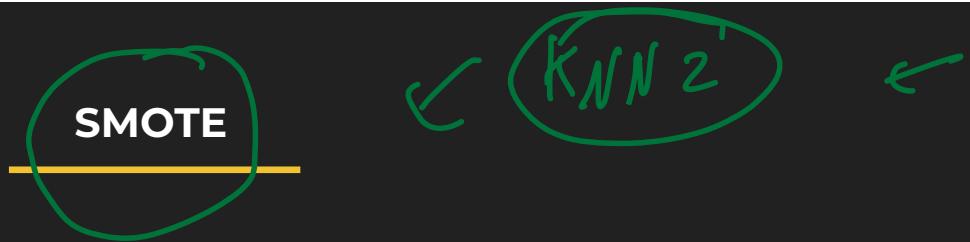
No. of samples \approx 1 Billion

Oversampling



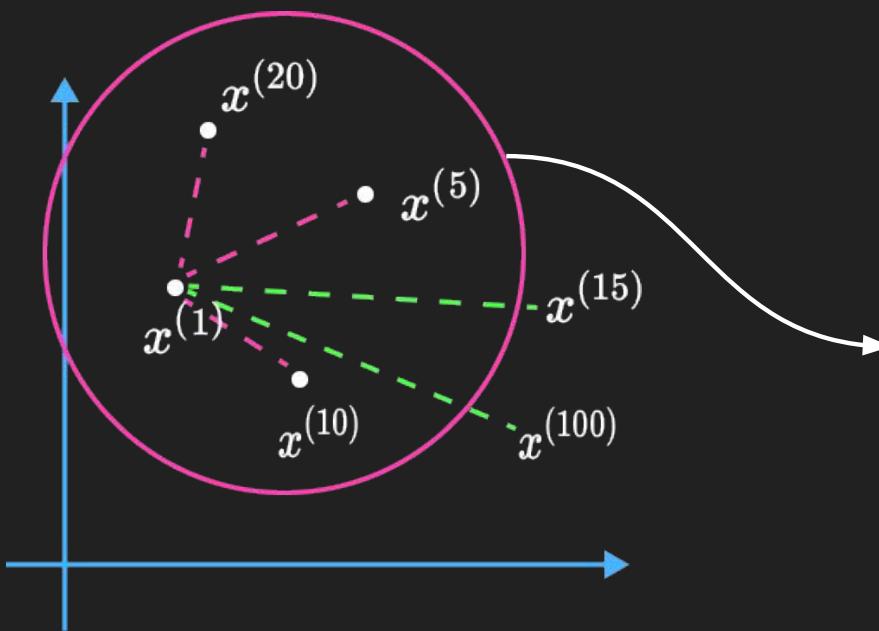
Why model overfit when using Oversampling ?





How to create synthetic data ? $x^{(1)} \rightarrow \text{random spam data}$

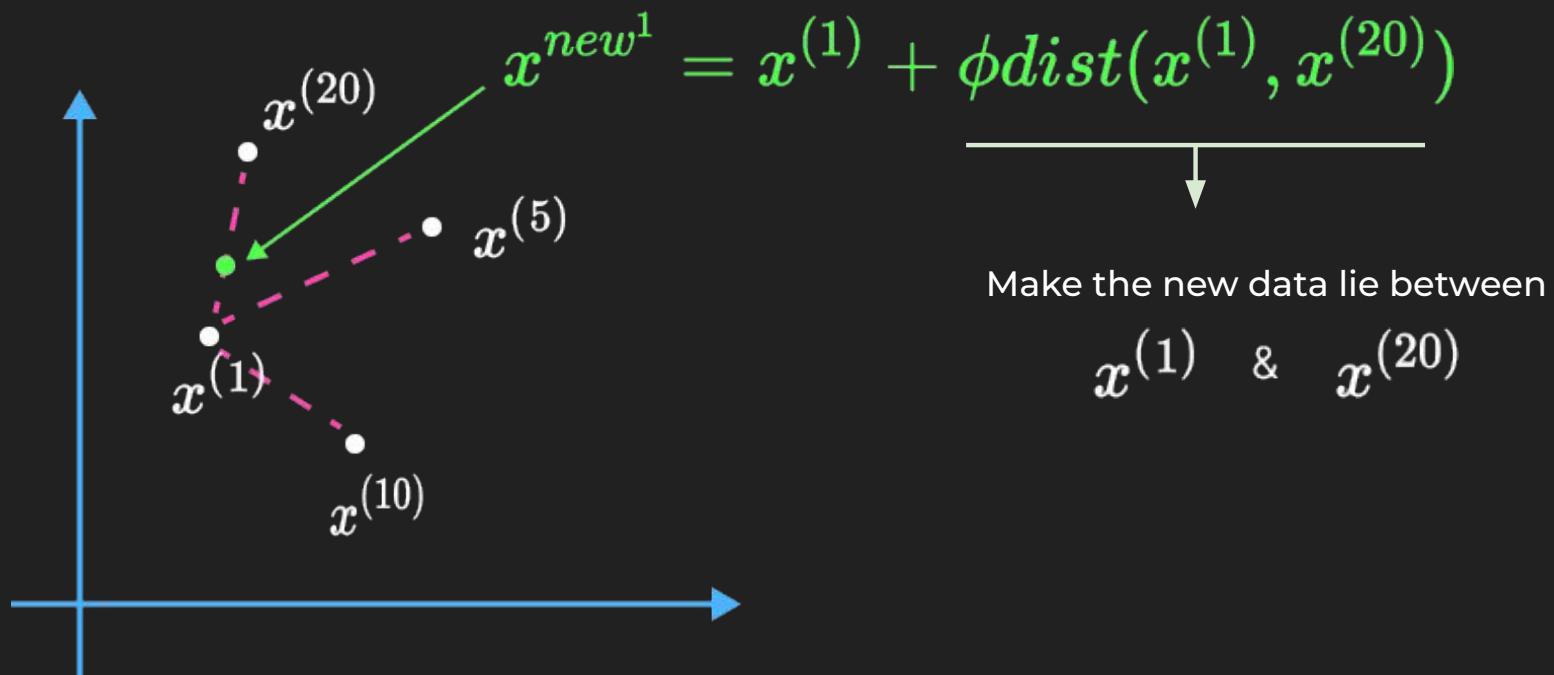
$$x^{(1)} = [f_1, f_2, f_3, \dots, f_d]$$



Finding euclidean distance from every spam
data & selecting 3 closet ones

$$[x^{(5)}, x^{(10)}, x^{(20)}]$$

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



Points to Remember

Sensitivity : TP \uparrow & FN \downarrow =
$$\frac{TP}{TP + FN}$$

FNR : tracks FN = 1 - Sensitivity

Specificity : TN \uparrow & FP \downarrow =
$$\frac{TN}{TN + FP}$$

(Sensitivity for Negative class)

FPR : tracks FP = 1 - Specificity

- Au - ROC for random model = **0.5**
- Au - ROC for ideal model = **1**
- Au - ROC does **not work in imbalance setting**
- Au - ROC depends on **order of probabilities.**

Points to Remember

Hyperplane of Logistic Regression is pushed away by majority class

Ways to balance , imbalance data :

- Class Weight : Add weight to minority class samples
- Undersampling : Repeat minority class samples
- SMOTE : Generate synthetic minority class samples



CME



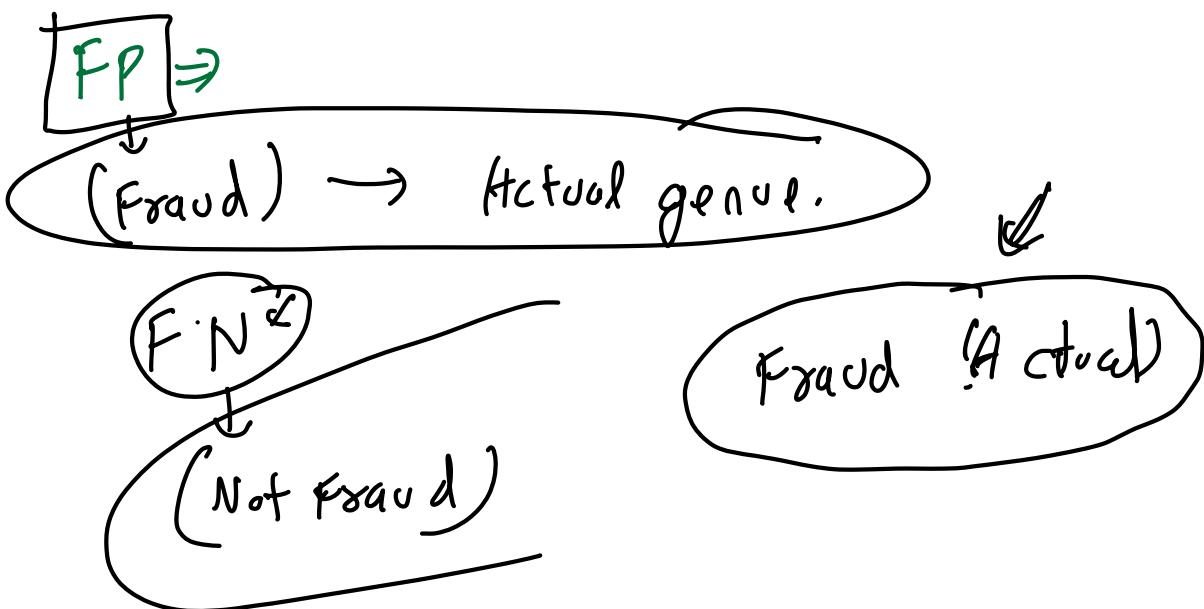
	0	1
0	TN	FP
1	FN	TP

No Fraud Fraud

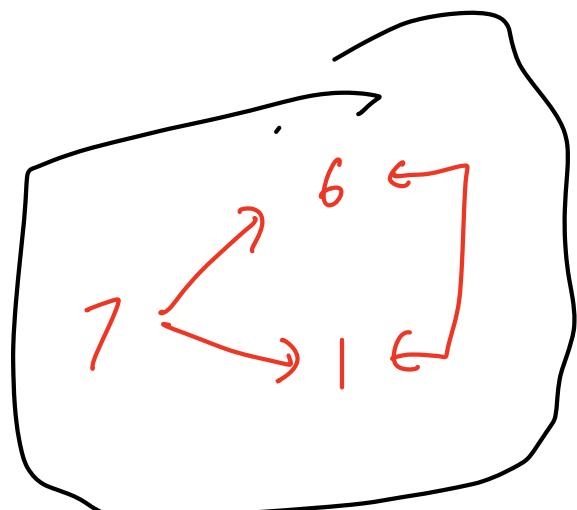
↓ FP ↓ FN

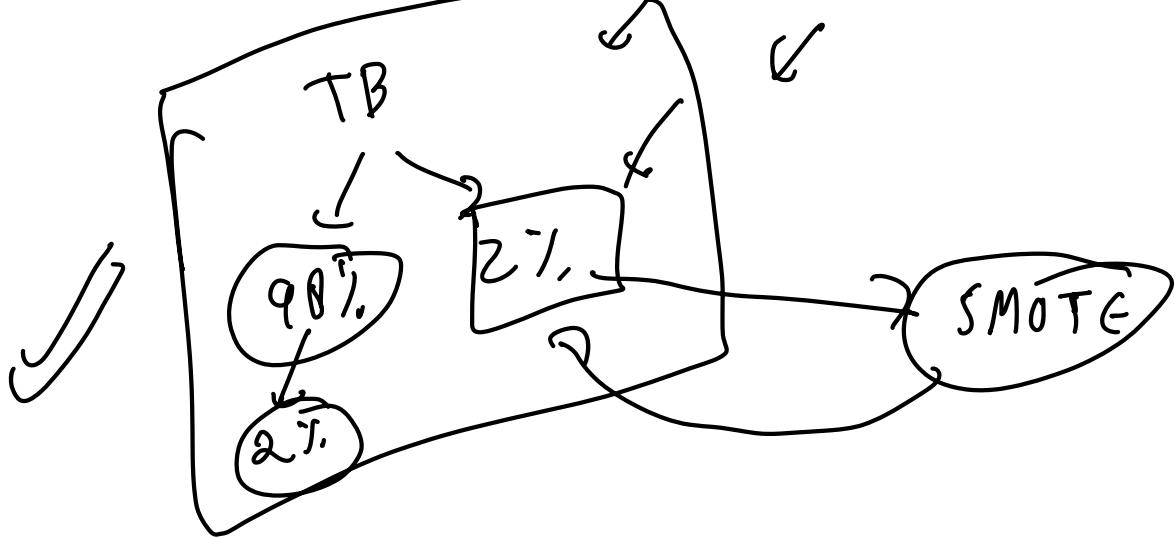
← Precision = 90% ← Recall = 90% ⇒

Accuracy ⇒



$P_1 \rightarrow 0$
 $P_2 \rightarrow 0$
 $P_3 \rightarrow 0$
 $P_4 \rightarrow 1$
 $P_5 \rightarrow 0$
 $P_6 \rightarrow 0$
 $P_7 \rightarrow 0$





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

