

CLASSIFICATION

METRICS-2

CLASS IMBALANCE

Confusion Matrix

predicted (\hat{y})

-ve Not SPAM +ve SPAM

Actual (y)

NOT SPAM -ve SPAM +ve

0	TN	FP
1	FN	TP

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FN + FP}$$

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

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

		predicted (\hat{y})	
		-ve NOT SPAM	+ve SPAM
Actual (y)	0	TN	FP
	1	FN	TP
		N_p	P_p

NA
 NOT SPAM
 -ve
 PA
 SPAM
 +ve

✓ Recall

True Positive Rate = $\frac{TP}{P_A} = \frac{TP}{TP+FN}$

Sensitivity.

True Negative Rate = $\frac{TN}{N_A} = \frac{TN}{TN+FP}$

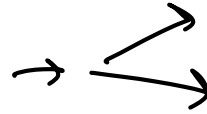
specificity

✓ False Positive Rate = $\frac{FP}{N_A} = \frac{FP}{FP+TN}$

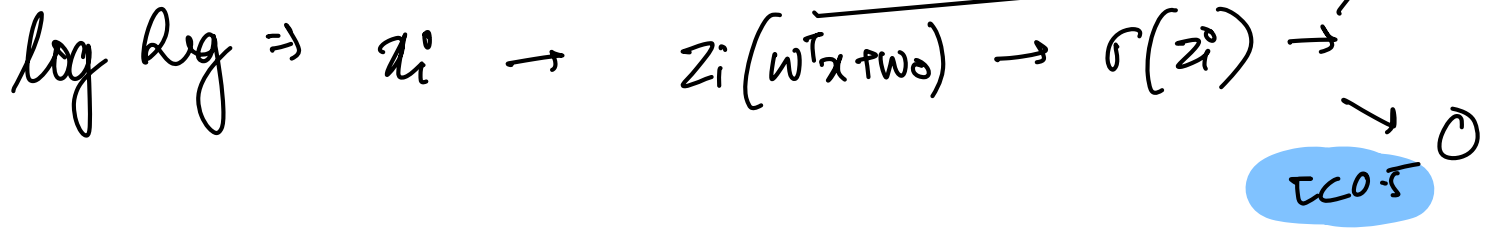
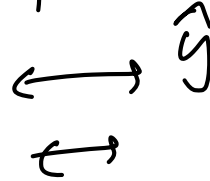
False Negative Rate = $\frac{FN}{P_A} = \frac{FN}{FN+TP}$

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

Classification metrics



spam/NS



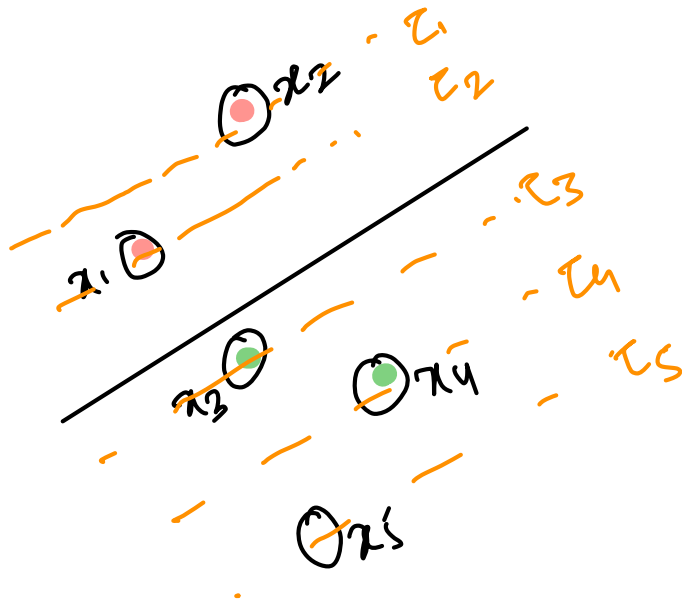
ROC Curve

$$\underline{T=0.5}$$

$$\underline{T=0.6}$$

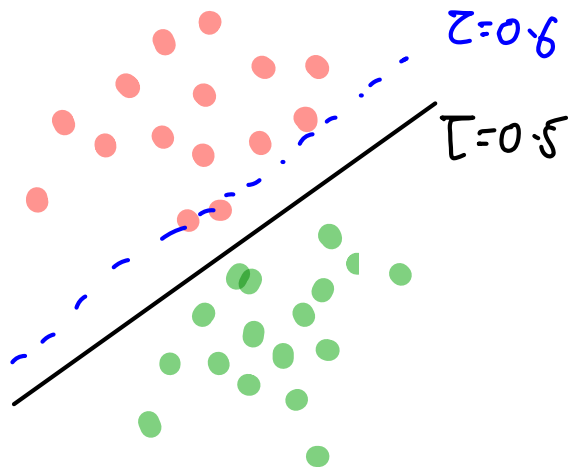
Receiver's operating characteristics.

Context : Binary Classification.



$$\sigma(z_i) \geq 0.6 \Rightarrow 1$$

$$\underline{\sigma(z_i) < 0.6 \Rightarrow 0}$$

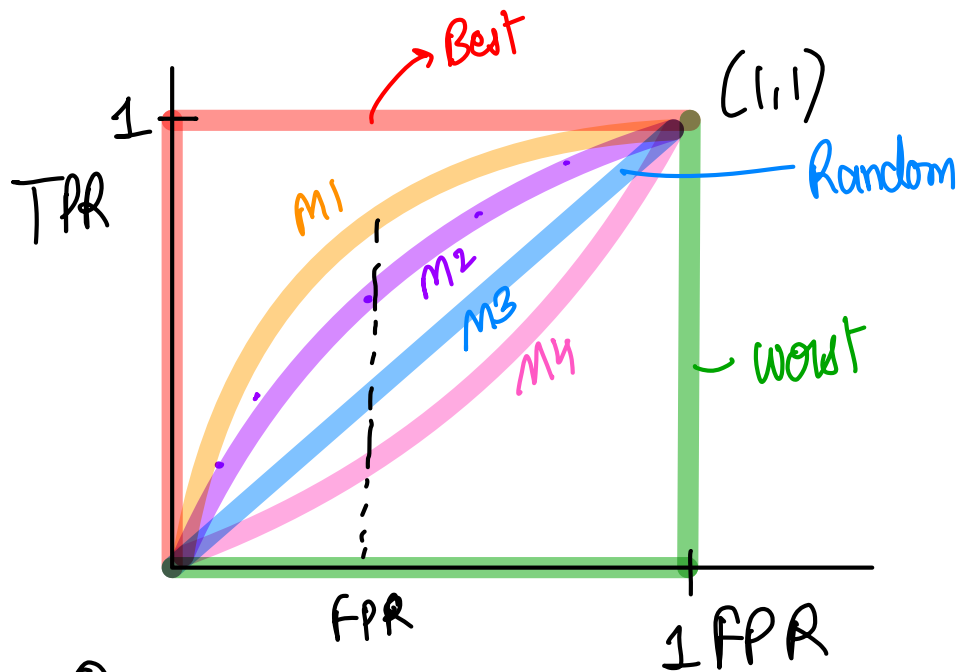


X	Y	$\sigma(z_i)$	$\tau=0.7$	$\tau=0.6$	$\tau=0.4$	$\tau=0.3$	$\tau=0.2$	
x_2	1	0.7	1	1	1	1	1	
x_1	1	0.6	0	1	1	1	1	
x_3	1	0.4	0	0	1	1	1	
x_4	0	0.3	0	0	0	1	1	
x_5	0	0.2	0	0	0	0	1	

$$\tau = 0.5$$

- ① Sort Records based on $\sigma(z_i)$ in descending order
- ② Find TPR and FPR for every threshold.

	P	TPR	FPR	
CM1	0.7	TPR ₁	FPR ₁	(— , —)
CM2	0.6	:		(— , —)
CM3	0.4	:		(— , —)
CM4	0.3	:		(— , —)
CM5	0.2	TPR ₅	FPR ₅	(— , —)



- ① \mathbb{I}
- ② Models Comparison
- ③ $[0, 1]$

$$1 \rightarrow \text{TPR} = \frac{TP}{TP + FN}$$

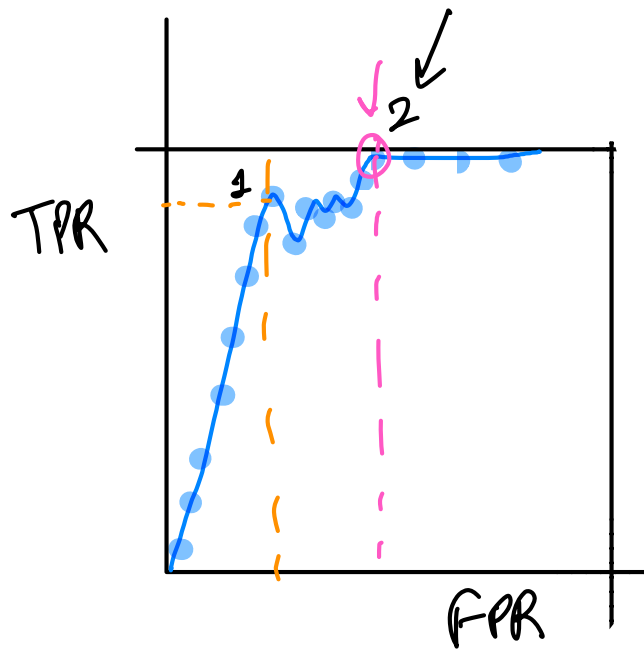
$$\text{FPR} = \frac{FP}{N} = \frac{FP}{FP + TN}$$

↓
0

Best \Rightarrow $\text{TPR} = 1$
 $\text{FPR} = 0$

Worst \Rightarrow $\text{TPR} = 0$
 $\text{FPR} = 1$

Random model:



How to choose z for model?

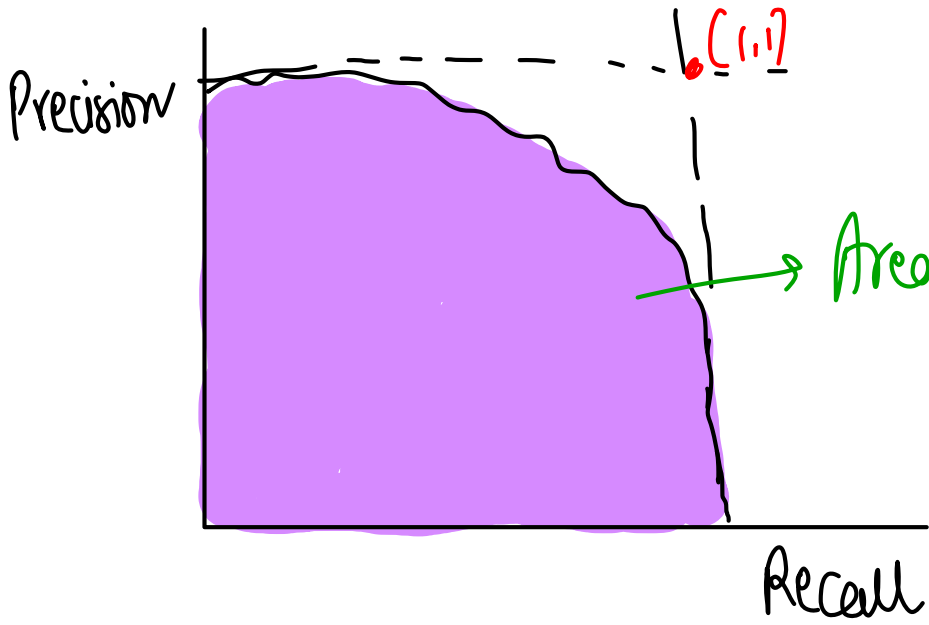
TPR \uparrow

FPR \downarrow

② ✓

problems with ROC-AUC

① doesn't work with imbalanced dataset



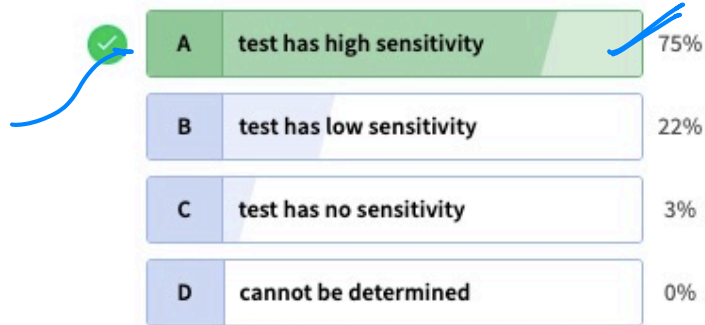
Area under PRC → Curve.
↓
Precision Recall

Quiz time!

🕒 Quiz Ended!

What to say when screening test identifies 92 Cancer patients out of 100?

32 users have participated



100
92
✓

Quiz time!

🕒 Quiz Ended!

How many points are typically used to plot an ROC curve?

28 users have participated

- | | | | |
|---|--|-----|-----|
| A | 2 points (0,0) and (1,1) | 18% | ✓ X |
| B | 3 points representing the thresholds 0.25, 0.5, and 0.75 | 11% | X |
| C | 10 points equally spaced between 0 and 1 | 4% | X |
| D | Depends on the number of unique threshold values | 68% | ✓ |
- Handwritten blue annotations: A checkmark and an 'X' next to option A; an 'X' next to option B; an 'X' next to option C; a green checkmark and a blue circle around option D, with a blue wavy line underneath it.

Quiz time!

🕒 Quiz Ended!

Which of the following metrics can be directly derived from the ROC curve?

31 users have participated

A Accuracy 6%

B Precision 3%

C Recall 3%

✓ D Area Under the Curve (AUC) 87%

Class Imbalance

One class \rightarrow dominated.

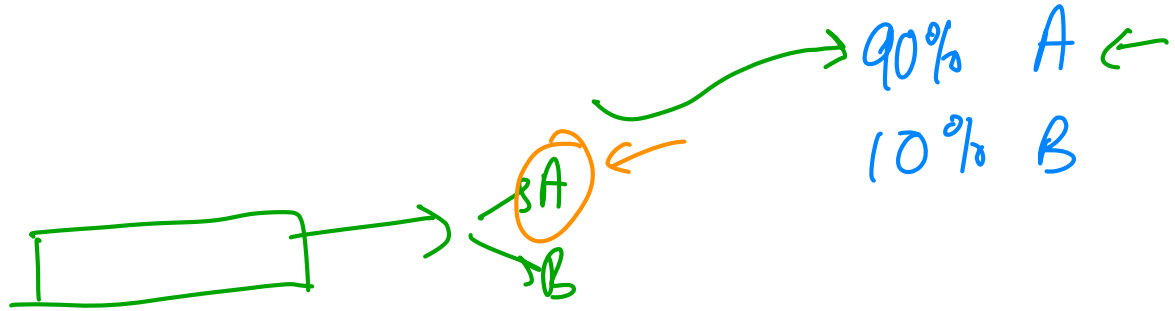
A	B	
50%	50%	\rightarrow Balanced
60%	40%	\rightarrow Slightly Balanced
70%	30%	\rightarrow slightly imbalanced
80%	20%	\rightarrow imbalanced
90%	10%	\rightarrow } highly imbalanced.
95%	5%	
99%	1%	

* value counts

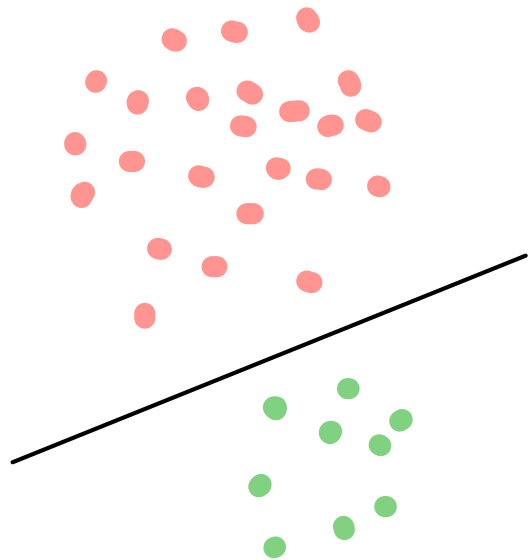
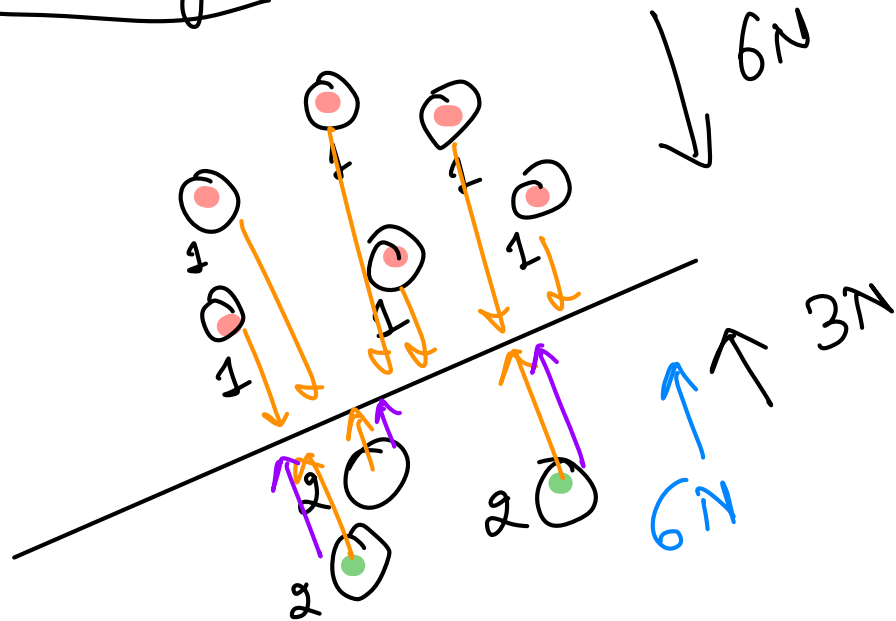
* Count plot

Problems with imbalanced dataset

- ① Accuracy / Metrics become unreliable
- ② Model becomes biased towards majority class



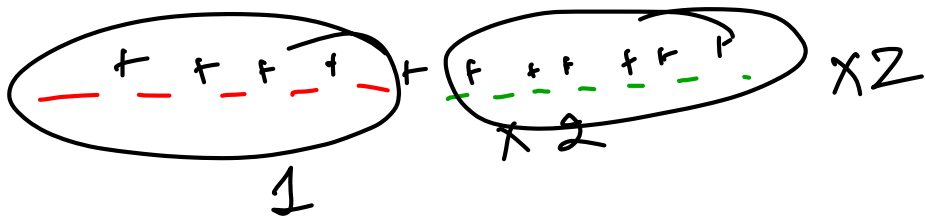
Class Weights



6 Red 3 Green

1 " " " " " " 2 " " " " 1 " "

Loss: $w_j \times \log \text{loss} + \nearrow \text{Reg}$

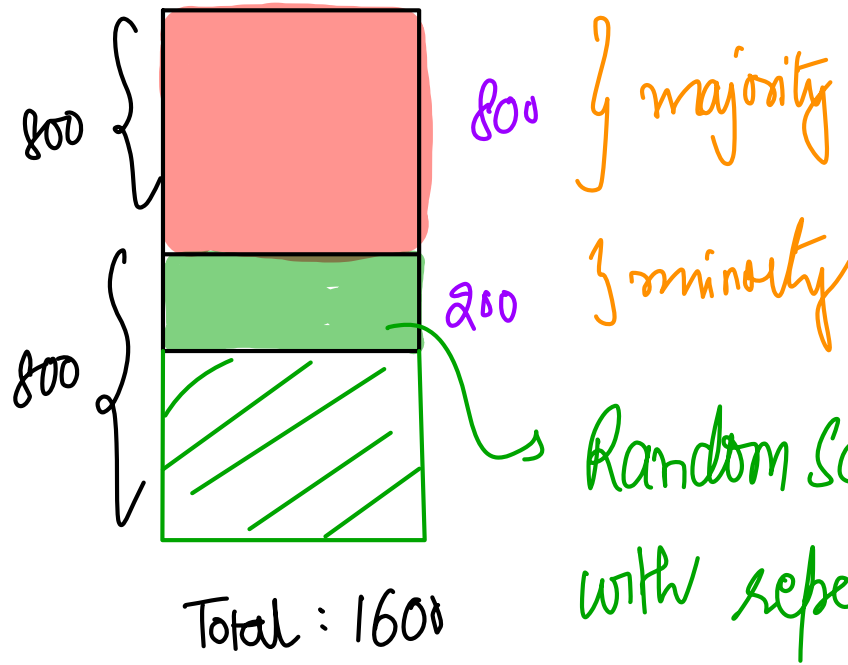


$$-\frac{1}{n} \sum_{i=1}^n [y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)]$$

Over Sampling

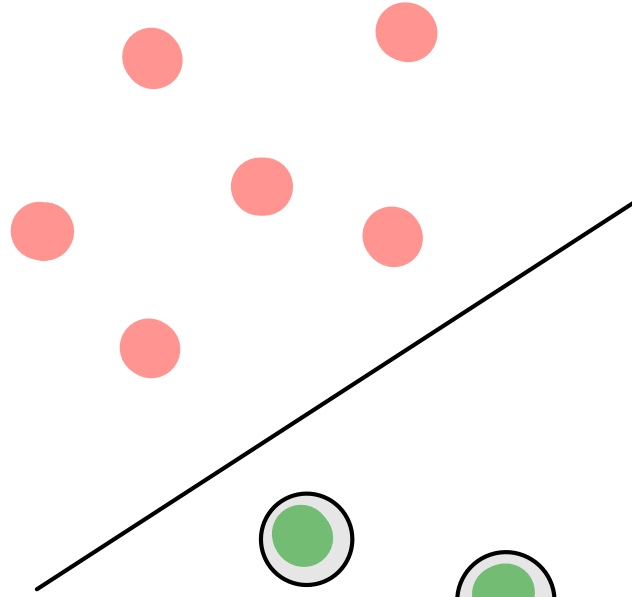
disadvantage

- ① possibility of overfitting
- ② Duplication of data



Advantage:

- ① No data loss
- ② Better than Under sampling



repeating the
same data points

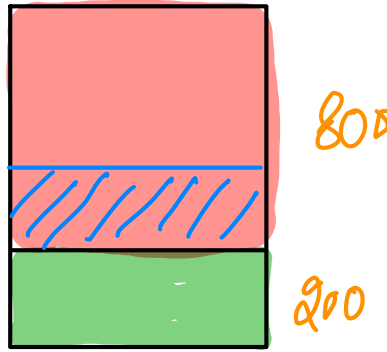
Under Sampling

disadvantage

- ① Data loss
Possibility of losing
some very imp data point
- ② Sample may be biased.

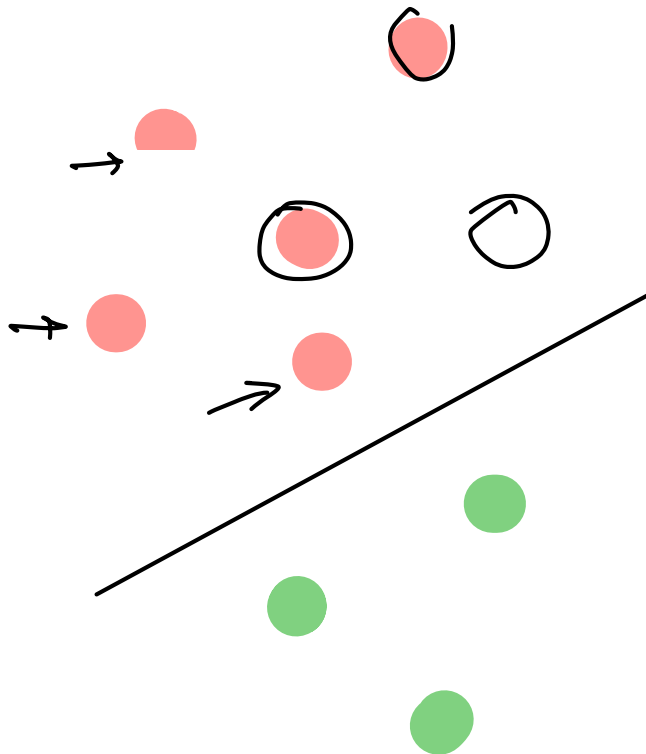
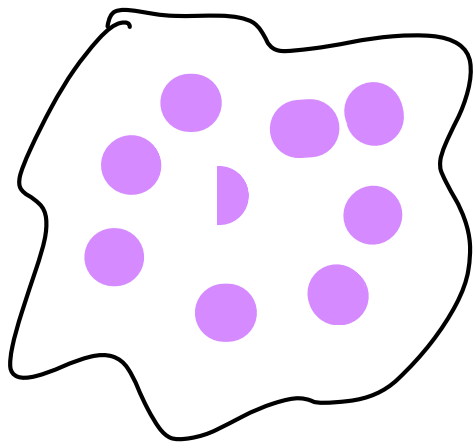
Advantages :

- ① time / compute / cost effective



Random Selection
Sampling

\neq # datapoints
in the minority
class.



SMOTE → Synthetic Minority Oversampling
Technique

Quiz time!

🕒 Quiz Ended!

Which among the following is a balanced data ?

31 users have participated



A

50 -ve samples, 50 +ve samples

97%

B

100 -ve samples, 10 +ve samples

0%

C

10 -ve samples, 100 +ve samples

0%

D

2 -ve samples, 98 +ve samples

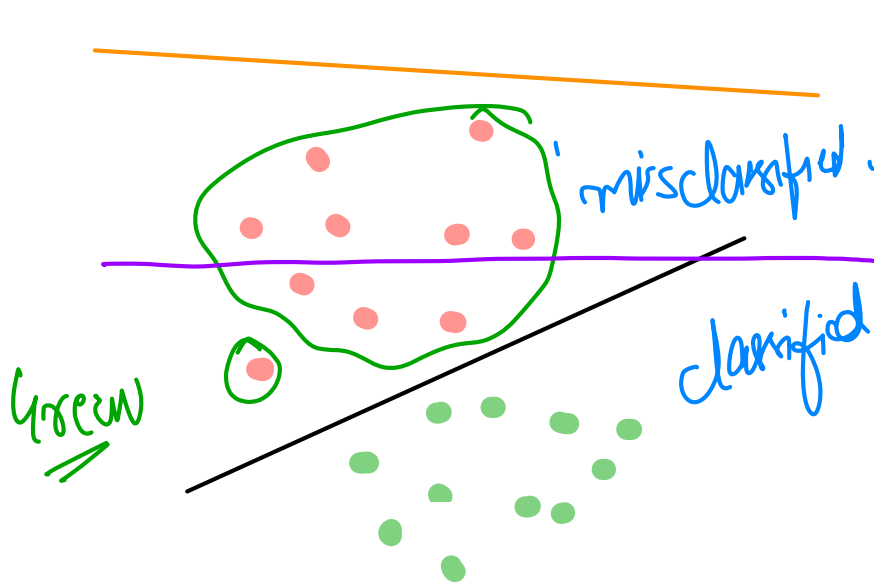
3%

① Assessment / Problem Solving Class

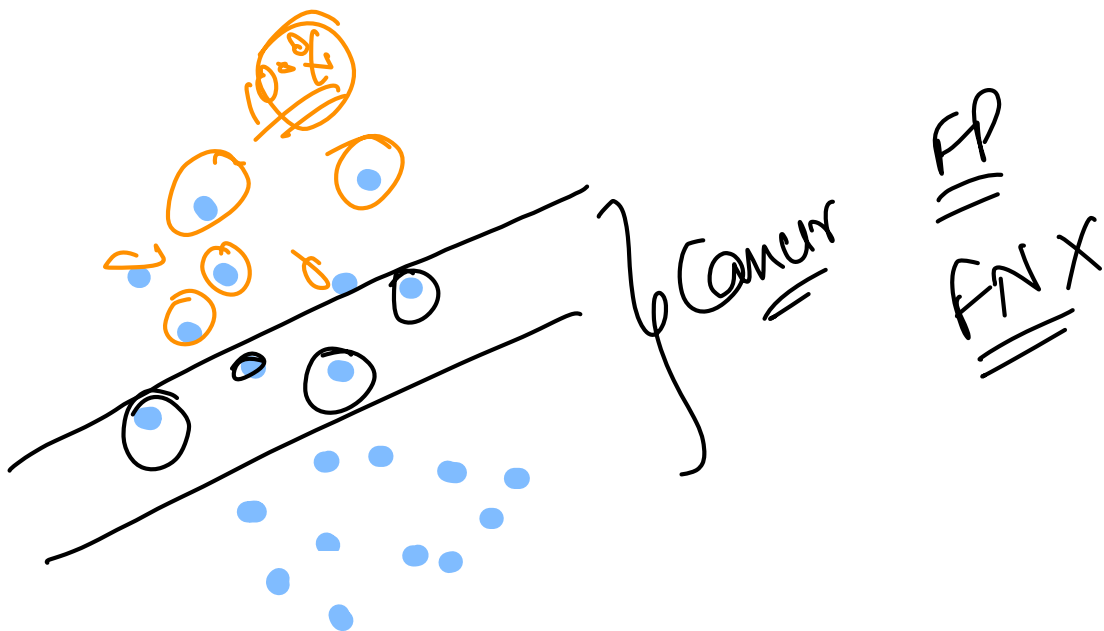
② Code →

$$\textcircled{x_i^0} \rightarrow W^T x_i^0 + w_0 \rightarrow s(z_i) \rightarrow \begin{cases} \text{if } \sigma(z_i) \geq 0.5 & 1 \\ \text{if } \sigma(z_i) < 0.5 & 0 \end{cases}$$

(0, 0.5)



Loss ↑↑↑



$$P[\text{non Cancer} \mid \text{patient}] = \underline{\underline{0.7}} \quad \checkmark$$

TPR $\uparrow\uparrow$
FPR $\downarrow\downarrow$

True Non C \rightarrow NonC

True C \rightarrow C
False C \rightarrow NonC

