

12-OCT-2025

## Agenda:

- ROC
- AUC
- cross-validation
- k-nearest Neighbour (ML Algo)

## ROC Curve:

Classifier → Score →  $\text{prob}(0, 1) \rightarrow$   
1 (positive)  
0 (negative)

→ To make a decision we need a threshold

$\text{Score} \geq \text{threshold} \rightarrow 1 \text{ else } 0$

→ TPR - True Positive Rate

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

spam classifier:

→ how many actual spams we caught correctly

→ FPR - False Positive Rate

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

→ how many not spams were wrongly marked as spam.

ROC curve is TPR vs FPR for all threshold.  
 (Receiver Operating Characteristic)

$y$	$\hat{y}$	$\text{Threshold}(0.2) \geq$	$\text{Threshold}(0.3) \geq$
1	0.8	1	1
0	0.1	0	0
0	0.9	1	1
1	0.56	1	1
1	0.7	1	1
0	0.2	1	0
		↓	
		$\text{TPR} \& \text{FPR}$	$\text{TPR} \& \text{FPR}$

ROC Curve: —

for each threshold ( $0.1 \rightarrow 0.9$ ), we calculate two things:

→ TPR

→ FPR

→ we plot TPR & FPR for all threshold.

email	True label ( $y$ )	Model-score ( $\hat{y}$ )
A	1	0.95
B	0	0.70
C	1	0.60
D	0	0.40
E	1	0.20

Actual spam : A, C, E

Actual negative : B, D

	True label ( $y$ )	Model-score ( $\hat{y}$ )	Threshold ( $0.9$ )
A	1	0.95	1
B	0	0.70	0
C	1	0.60	0
D	0	0.40	0
E	1	0.20	0

Threshold: 0.8

TP: A (1)      FP: 0

TN: B,D (2)      FN: C,E (2)

TP: model predict positive, actual positive  
 TN: model predict negative, actual negative  
 FP: model predict positive, actual negative  
 FN: model predict negative, actual positive

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

$$= 1/3 = 0.33$$

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

$$= 0/(0+2) = 0$$

Threshold = 0.6

$$\text{TPR} = 0.67$$

$$\text{FPR} = 0.5$$

Threshold = 0.4

$$\text{TPR} = 0.67$$

$$\text{FPR} = 1.0$$

Threshold = 0.0

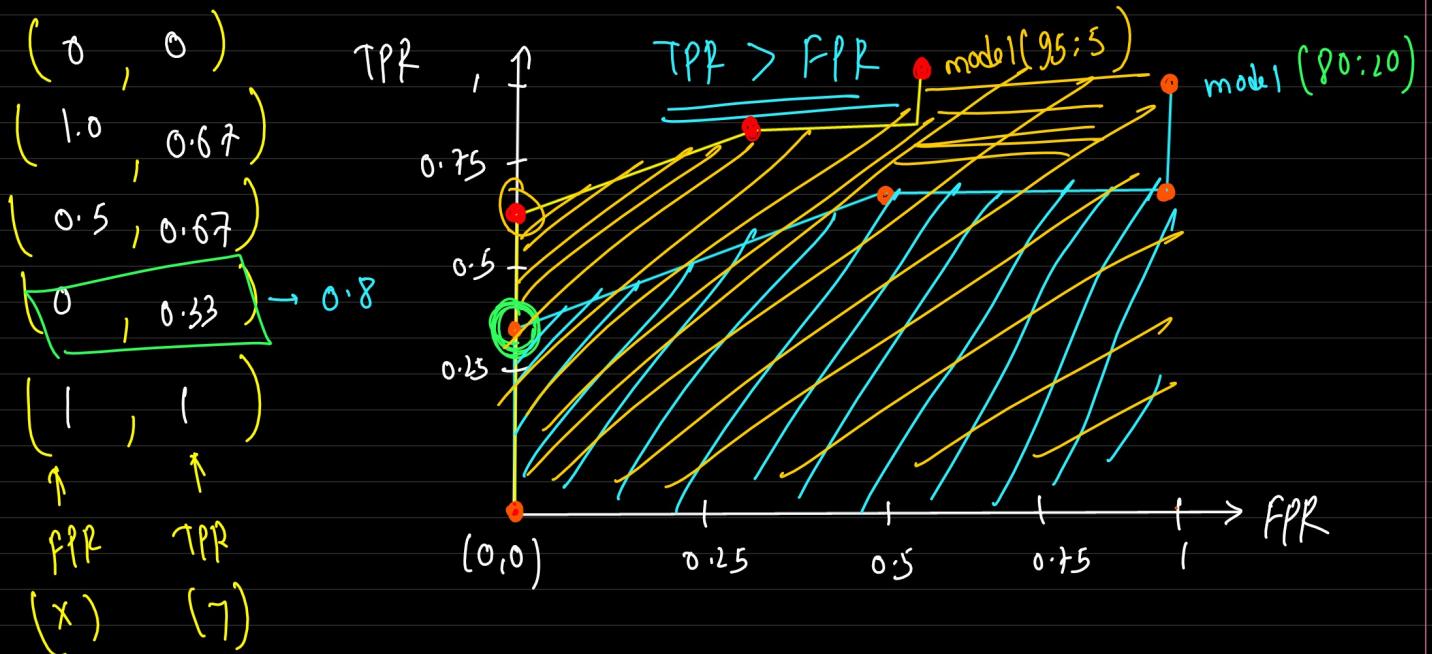
$$\text{TPR} = 1$$

$$\text{FPR} = 1$$

Threshold = 1.0

$$\text{TPR} = 0$$

$$\text{FPR} = 0$$



AUC - Area under ROC curve

- numerical value of the area under ROC curve

: prob that the model ranks a random positive higher than random negative.

$$\begin{aligned} \text{AUC} &\rightarrow 1 \\ &= 0.7 \\ &= 0.5 \end{aligned}$$

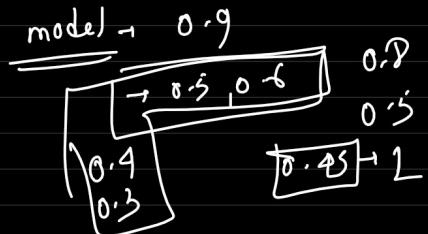
: model is better at separating classes

$x_1 \ x_2 \ x_3 \ x_4$  cancer

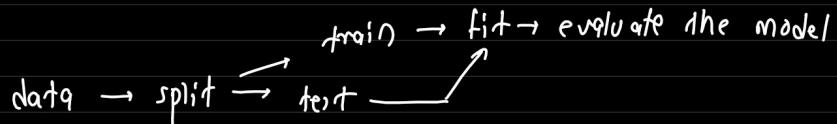
1  
0

1  
0

$\xrightarrow{5-10\% \rightarrow 1}$   
 $90-95 \rightarrow 0$



## Cross-validation :-



\* Maybe our test set is too hard.

\* Maybe our test set is too easy.

→ reduce this dependency of data split.

→ solution:

Instead of one fix split, we train multiple times, each time with a different split.

Index	Data
1	A
2	B
3	C
4	D
5	E
6	F

: 3 fold cross validation:

: Each data becomes test data once

fold	train	test
1	[C, D, E, F]	[A, B] → ML.fit → Metrics
2	[A, B, E, F]	[C, D] → ML.fit → Metrics
3	[A, B, C, D]	[E, F] → ML.fit → Metrics

1 → 96%  
 2 → 84%  
 3 → 90%

↓ average → **bias variance** ) not for model improvement.

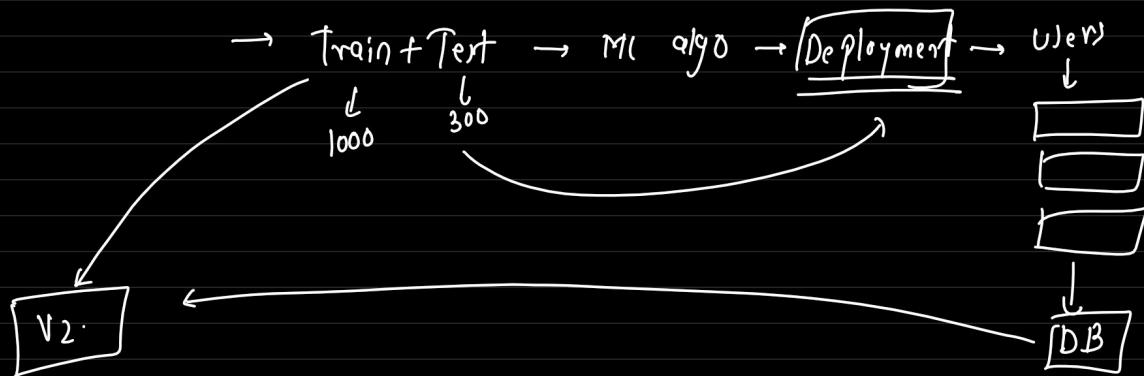
$k$ -fold  $\rightarrow$  split data into  $k$  parts, each used once for testing

stratified  $k$ -fold  $\rightarrow$  ensure class balance in each fold

Data  $\rightarrow$  1000 rows

EDA  $\rightarrow$  FE

$\rightarrow$  (row-validation) split  $\rightarrow$  general  $\xrightarrow{\text{Train}}$   
accuracy  $\xrightarrow{\text{Test}} \begin{matrix} 1 \\ 2 \\ 3 \end{matrix}$



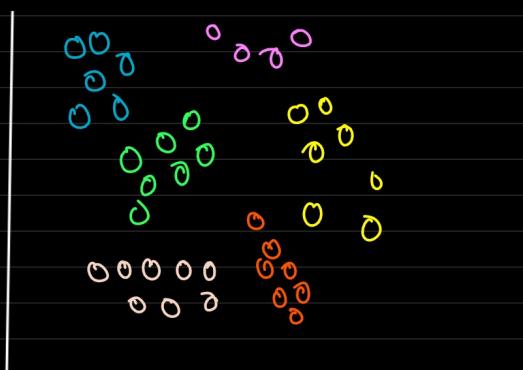
$k$ NN  $\rightarrow$   $k$ -Nearest Neighbour.

$\rightarrow$  classification

$\rightarrow$  regression

\* mathematical equation (line or curve)  $\times$

\* what if data doesn't form a nice line



Solution: KNN

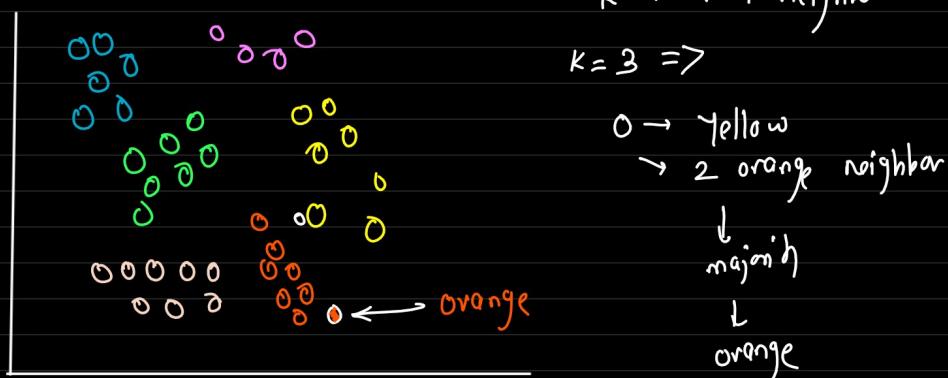
→ Instead of relying on equation, why not use the data itself to make decision.

→ we don't compute any parameters:  $\underbrace{m_1, m_2, \beta, M}$

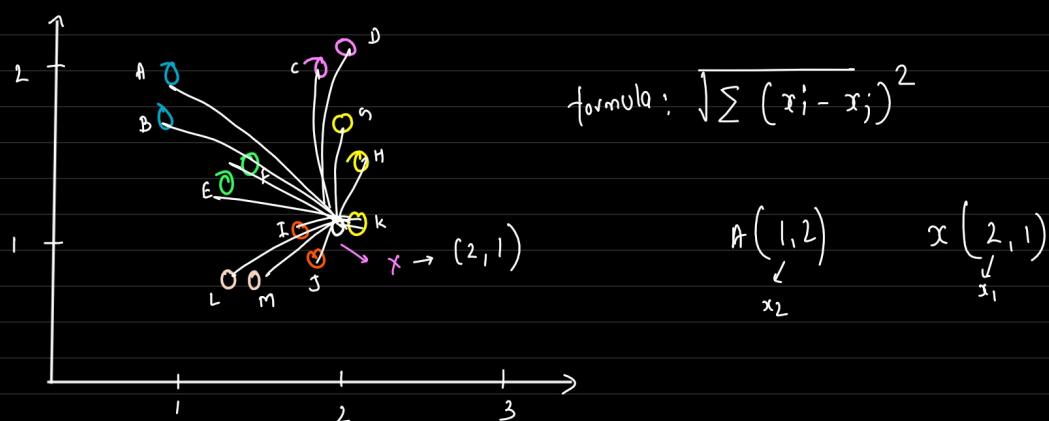
→ called as lazy learner, bcoz it doesn't learn any equation.

→ supervised algo

$k=1 \Rightarrow 1$  neighbor



$k \Rightarrow$  odd  $\rightarrow$  majority



point	coordinates	formula	Distance from $(2, 1)$	
A	$(1, 2)$	$\sqrt{(2-1)^2 + (1-2)^2}$	1.41	sort the distance in ascending order
B	$(0.9, 1.7)$		1	
C	$(x, y)$		2	
D			4	

→ Data point closer to  $x(2,1)$

$k=5 \rightarrow$  3 → orange → majority  
 $\rightarrow 2 \rightarrow$  yellow

Euclidean distance:

$$d(x_i, x_j) = \sqrt{\sum (x_i - x_j)^2}$$

$$\text{Manhattan distance} = \sum |x_i - x_j|$$

Optimal value of  $k$ :

Small ( $k=1, 2$ ) → High variance

Medium ( $k=5, 7, 9$ ) → Balanced

Large ( $k \rightarrow N$ ) → High bias