

- \* To decide which one should we work for between Sensitivity & Specificity
- \* If identifying +ve is more important to us, then we will select algo that has high sensitivity
- \* If correctly identifying -ve is more important then we will select algo that has high specificity

## Precision & Recall $\Rightarrow$

$\rightarrow$  Used for Information Retrieval.

$\rightarrow$  Google Search Engine  $\Rightarrow$

$\rightarrow$  query fired

$\rightarrow$  Have millions of related records

$\rightarrow$  From these top 10-100 records are returned.

$$\text{Precision} = \frac{\text{Correctly predicted +ve}}{\text{Total +ve Predicted}} = \frac{TP}{TP + FP}$$

(Range  $(0 - 1)$ )

$$\text{Recall} = \frac{\text{Correctly Identified +ve}}{\text{Total Actual +ve}} = \frac{TP}{TP + FN}$$

(TPR) (Range  $0 - 1$ )

\* Ideally we want precision to be high ( $\approx 1$ ) for a good classifier

(range 0 - 1)

\* Ideally we want Precision to be high (*i.e.* 1) for a good classifier

$$\text{Precision} = 1 = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 \Rightarrow \text{when } \boxed{\text{FP} = 0}$$

\* Ideally we want Recall to be very high (*i.e.* 1) for a good classifier

$$\therefore \text{Recall} = 1 = \frac{\text{TP}}{\text{TP} + \text{FN}} \Rightarrow \text{when } \boxed{\text{FN} = 0}$$

So Ideally a good classifier has High Precision & Recall

But in reality there is trade-off

\* When we tweak our model to increase one, then the other decreases.  
(update)

Q) Explain

### F1-Score

- \* In reality we need a metric that takes into account both precision and recall.
- \* F1 score is a metric that takes into account both precision & recall.
- \* F1 score is harmonic mean of Precision & Recall.

$$\text{F1-Score} = \frac{\frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}}{\text{Precision} + \text{Recall}} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Harmonic mean of two variables  $\frac{2}{\frac{1}{a} + \frac{1}{b}} = \frac{2ab}{a+b}$

for n variables

$$H = \frac{n}{\frac{1}{n_1} + \frac{1}{n_2} + \dots + \frac{1}{n_n}} = \frac{n}{\sum_{i=1}^n \frac{1}{n_i}}$$

If F1-Score = 1  $\Rightarrow$  when Precision = 1  
Recall = 1

- \* When Precision and Recall both are high then F1-Score is high

When to Use F1-Score  $\rightarrow$

$\rightarrow$  Accuracy is not a good metric to use when we have class imbalance.

Ex  $\rightarrow$  let say 99% of people visiting site are onlookers and not purchasing anything.

$\rightarrow$  Suppose we have a model that predicts that 1% people visiting site are onlookers.

" 1% error is acceptable

people visiting site are onlookers.

→ The model is 1% wrong, Generally 1% error is acceptable

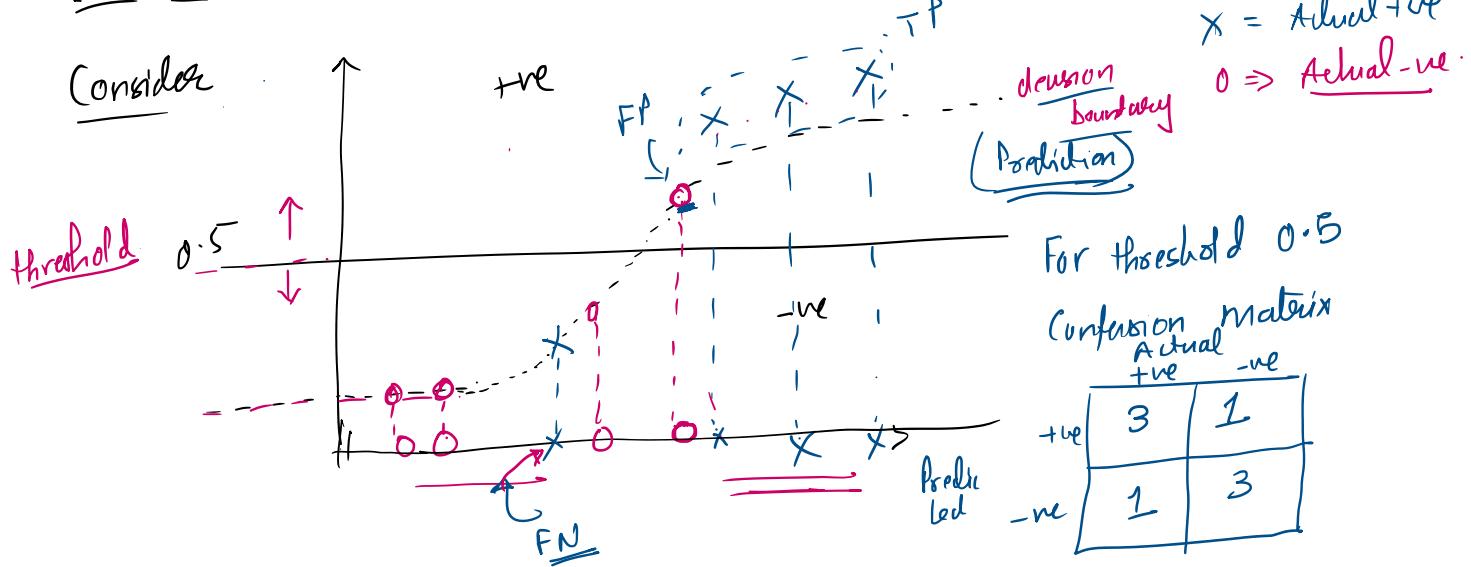
→ But such model in this case is useless

→ In such case instead of accuracy, we will prefer

F1-Score.

## ROC [ Receiver Operator characteristic ]

Consider

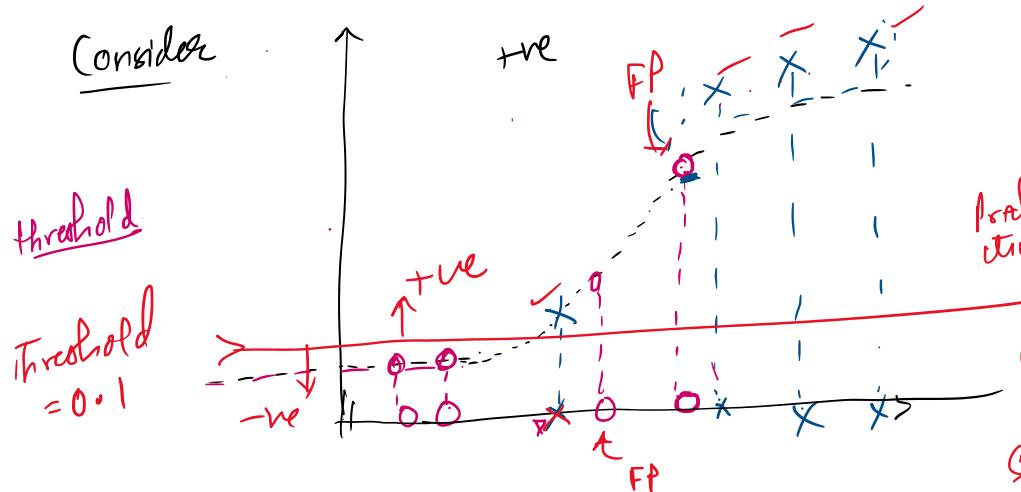


$$\text{Sensitivity} = \frac{3}{4} = \text{TPR}$$

$$\text{Specificity} = \frac{3}{4} = \text{TNR}$$

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

Consider



Now threshold = 0.1

Confusion matrix

		Actual	
		+ve	-ve
True	+ve	4	2
	-ve	0	2

$$\text{Sensitivity} = \frac{4}{4} = 1$$

$$\text{Specificity} = \frac{2}{4} = \frac{1}{2}$$

- \* Consider for logistic regression, where we identify a threshold point and prepare confusion matrix and calculate Sensitivity & Specificity.

- and calculate Sensitivity & Specificity
- \* If threshold changes then the confusion matrix and accordingly the Sensitivity and Specificity changes.
  - \* We can have many such thresholds bet<sup>n</sup> 0 → 1
  - \* We want to analyze the performance at diff threshold and want to identify the best of it.
  - \* For this we will plot TPR and FPR for diff threshold.
- 

→ From the above ROC curve for Random Forest model, it will help us to find the best threshold.

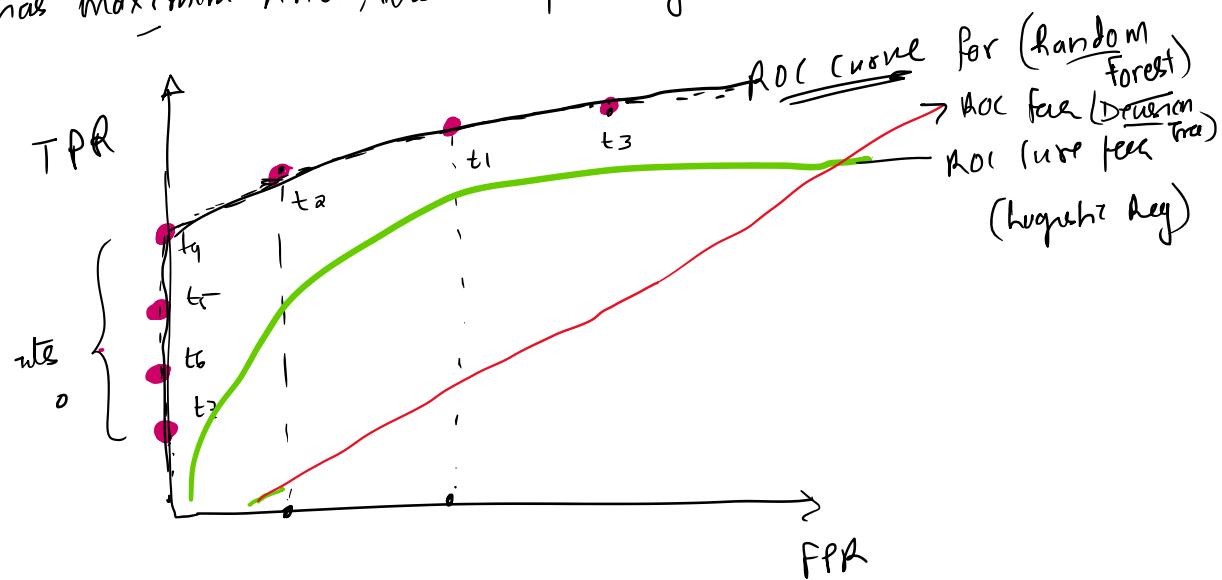
## Q) Explain

AUC [Area Under Curve] → It is a method to compare ROC for more than one method and will help to judge which one is better.

→ Here will find Area Under ROC for each method.  
→ ROC that has maximum Area, the corresponding method will be best.

ROC curve for (Random Forest)

→ ROC that has maximum Area, the better it is.



From above ROC Curves, we find that there is maximum area under ROC Curve for Random forest, hence the method Random forest will be the best method.

Q) Explain Kappa Statistic?

- \* Kappa Statistic or Cohen's Kappa is statistical measure of inter-rater reliability for categorical variable.
- \* It is used when two/more raters apply a criteria based on a tool to assess whether or not some condition occurs.
- \* Ex: let say Two doctors rates whether or not each of 20 patients has diabetes based on symptoms.
- \* If two raters uses same criteria on same target to evaluate and then their agreement is very high then we will have evidence of reliable rating.
- \* If their agreement is not very high then →
  - either criterion tool is not useful
  - or raters are not trained enough.

- \* Kappa Statistics correct for chance agreement and not percent agreement

Evaluator A

	Yes	No
Yes	35	20
No	15	40

Evaluator B

- ⇒ 35 times both agreed - said Yes
- ⇒ 40 times both agreed ⇒ said No
- ⇒ 20 time A said NO but B said Yes
- ⇒ 15 times A said Yes but B said No.

11

$\Rightarrow$  do time A said NO but B said YES  
 $\Rightarrow$  IC time A said YES but B said NO.

Cohen suggested following Statistics.  $\Rightarrow$

value $\leq 0$	$\Rightarrow$ No agreement
0.01 $\rightarrow$ 0.20	$\Rightarrow$ as <u>none</u> to <u>slight</u>
0.21 $\rightarrow$ 0.40	$\Rightarrow$ as <u>fair</u>
0.41 $\rightarrow$ 0.60	$\Rightarrow$ as <u>moderate</u>
0.60 $\rightarrow$ 0.80	$\Rightarrow$ as <u>substantial</u>
0.81 $\rightarrow$ 1.00	$\Rightarrow$ <u>perfect agreement</u>

\* Rather than calculating the percentage of items, the raters agreed on Cohen's Kappa attempts to account the fact that rater may happen to agree on some items purely by chance

Ex Two curators asked to rate 70 paintings.

		Curator C2	
		Yes	No
Curator C1	Yes	25	10
	No	15	20

Step 1) Calculate Relative Agreement bet<sup>n</sup> Curators.

$$P_o = \frac{\text{Both Said Yes} + \text{Both Said No}}{\text{Total}} = \frac{25+10}{70} = 0.6429$$

Step 2) Calculate hypothetical probabilities of chance Agreement bet<sup>n</sup> Curators.

$$P(\text{Yes}) = \underline{C_1(\text{Yes})} \times \underline{C_2(\text{Yes})} = \frac{(25+10)}{70} \times \frac{(25+15)}{70} = 0.2857$$

$$\checkmark P(Y_{10}) = \frac{C_1(Y_{10})}{\text{Total Ans}} * \frac{C_2(Y_{10})}{\text{Total Ans}} = \frac{(25+10)}{70} * \frac{(25+15)}{70} = \underline{\underline{0.2857}}$$

$$\checkmark P(N_0) = \frac{C_1(N_0)}{\text{Total Ans}} * \frac{C_2(N_0)}{\text{Total Ans}} = \frac{(15+20)}{70} * \frac{(10+20)}{70} = \underline{\underline{0.214285}}$$

$$\underline{\underline{P_e}} = P(Y_{10}) + P(N_0) = 0.2857 + 0.214285 = \underline{\underline{0.5}}$$

Calculate Cohen's Kappa =  $K = \frac{P_o - P_e}{1 - P_e} = \frac{(0.6429 - 0.5)}{1 - 0.5} = \underline{\underline{0.2857}}$

gt is in range  $0.21 \rightarrow 0.40$  so the agreement bet<sup>n</sup> two  
 Cuthors is fair

## Module - 3

### Ensemble learning →

#### 3.1 Understand Ensembles

K-fold Cross Validation

⇒ Boosting

Stumping (Adaboost)

XGBoost

3.2 ⇒ Bagging

Subagging

Random Forest

Comparison with Boosting

Different ways to

Combine Classifiers

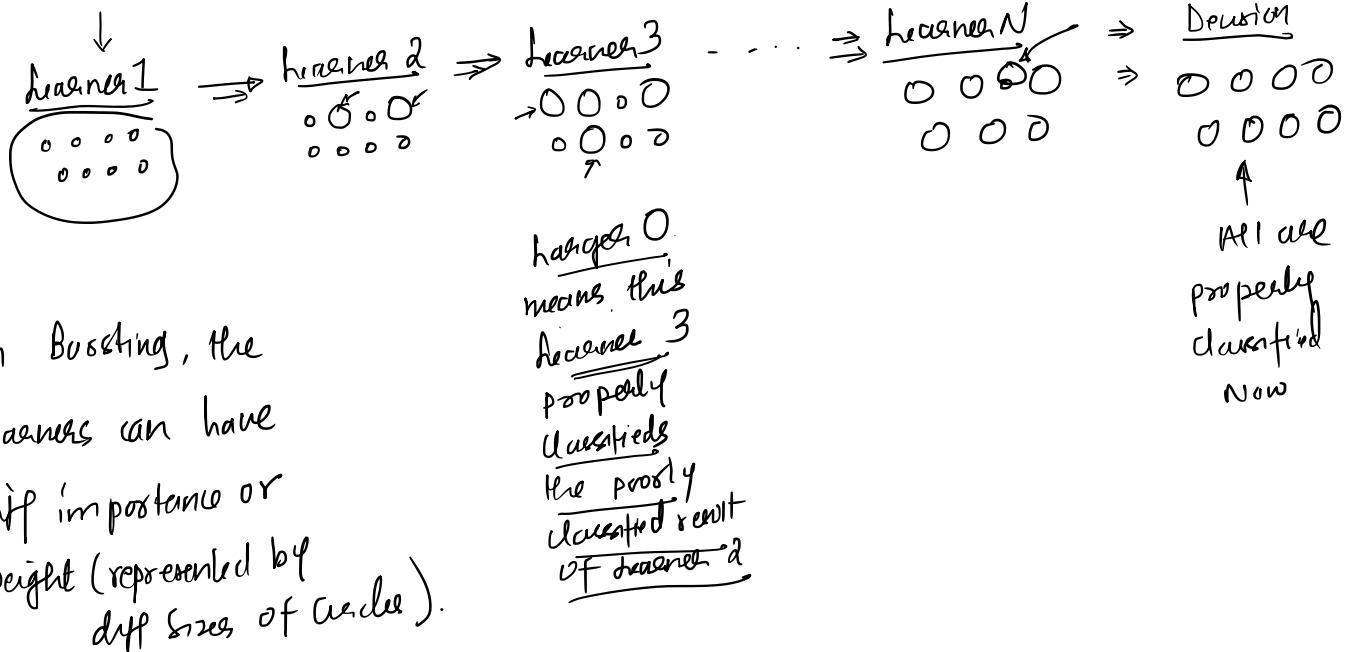
## \* Ensemble ->

- \* In ML, ensemble is a model that combines the prediction from two or more models.
- \* The models that contributes to Ensemble are known as ensemble members.
- \* The members may or may not be trained on same training data and they may be of same type or different type.
- \* It's very powerful method to improve the performance of the model.
- \* It's technique that uses group of weak learners in order to create a strong and aggregated learner.

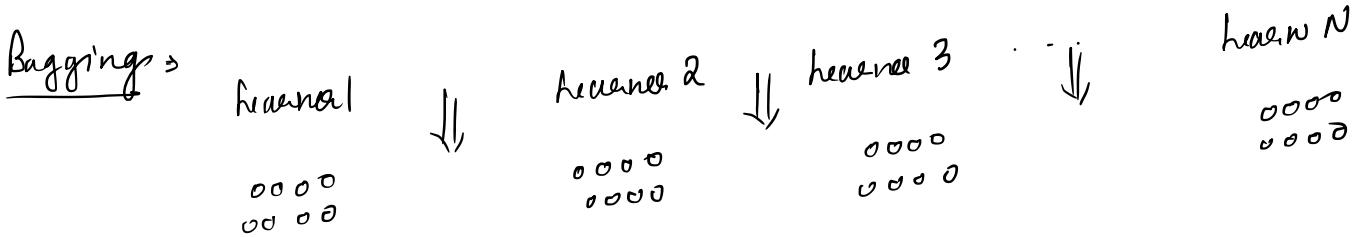
→ The Ensemble technique helps to reduce the Variance (By Bagging) and Bias (By Boosting) and thus helps in improving the predictions.

- Boosting model :-
- \* It falls inside family of Ensemble method.
  - \* It consists of filtering or weighting the data that is used to train team of Weak Learners, so that the new learner can give more weight on sample that is poorly classified by previous learner.
- In Boosting the learners are trained

## Sequentially



\* In Boosting, the learners can have diff importance or weight (represented by diff sizes of circles).



- ⇒ In Bagging the weak learners are trained in parallel using randomness
- ⇒ All learners have same weights.