

## → "Precision, Recall & F1-score" :-

There are some isolated measures called "Precision & Recall" & they are often used in information retrieval problems.

They are often used when we have a large corpus of text. They are used in search engines.

Actual ↓	0	1
0	TN	FN
1	FP	TP

'N'   'P'

The formula for precision is

$$\left\{ Pr = \frac{TP}{TP + FP} \right\}$$

Let's first understand what does precision mean intuitively :-

Precision (Pr) :- Of all the points, the model declared / predicted to be true, what %age of them are actually true.

↓  
In precision, we are not worried too much about the negative class, same is the case with recall also.

Note :- "Precision & recall" take care about the "positive-class" & not about the "negative-class".

Recall :- It is nothing but the TPR (True positive rate).

$$TPR = \frac{TP}{P}$$

In both the cases, Recall & Precision basically revolves around "TP" only.

So, "Precision & Recall" are the metrics, which are very useful when we care about mostly the "positive class".

Let's understand recall intuitively, Recall basically means of all the points which actually belong to "positive class", how many, the model detected to be of "positive class".

"Precision"  $\Rightarrow$  of all the points that the model predicted to be positive, how many are actually positive.

"Recall"  $\Rightarrow$  of all the actually positive points, how of them are predicted to be positive by the model.

Now given these two measures, is there any way for us to combine ["Precision"] & ["Recall"] into one measure (single measure).

Note  $\Rightarrow$  we always want "Precision" & "Recall" to be high.

$$\left. \begin{array}{l} \text{Ps} \uparrow \\ (0 \text{ to } 1) \end{array} \right\} \quad \left. \begin{array}{l} \text{Re} \uparrow \\ (0 \text{ to } 1) \end{array} \right\}$$

There is a measure called "F1-score" which combines both these metrics ("Precision" & "Recall").

$$\text{'F1-score'} = \left( 2 * \frac{\text{Ps} * \text{Re}}{\text{Ps} + \text{Re}} \right)$$



Interpreting "F1-score" is much more harder than interpreting "Precision" & "Recall".

Given Precision, & Recall, it is much easier to understand, but there will be some instances, where everything has to be converted into just a single metric, that's where we use "F1-score".

✓ The formula is based on "Harmonic Mean"

$$F1 = \frac{2}{\left( \frac{1}{\text{recall}} + \frac{1}{\text{precision}} \right)} = \left[ \text{Avg} \left( \text{inv-recall} ; \text{inv-prec} \right) \right]^{-1}$$

$\downarrow$

$$\frac{1}{2} \left( \frac{1}{\text{re}} + \frac{1}{\text{ps}} \right)$$

F1-score is often used in lots of "kaggle competitions" (3)  
 & F1 score is high, if "Precision" & "Recall" both are high.

Precision & Recall in F1-score are represented as:

$$F1\text{-score} = \left( 2 \times \frac{Pr \times Re}{Pr + Re} \right)$$

These numbers lie b/w "0 & 1"  
 F1 score also lies b/w "0 & 1"  
 [ "1" being very good & "0" being very bad ]

→ Receiver operating characteristic curve & AUC :-  
 (ROC) ↓  
"Area under the curve"

Other very interesting metric is called "ROC" curve & "AUC" {Area under the curve}

This curve has very nice interesting history, This curve has very nice interesting history. Actually, it was designed by "Electronics & Radio Engineers" during "second world war" to predict how well their missiles are working.

Let's understand, what "ROC" is :-

Imagine if we have a dataset like (comprising of five datapoints) for

"x"	"y"	"G"
x <sub>1</sub>	1	
x <sub>2</sub>	1	
x <sub>3</sub>	0	
x <sub>4</sub>	1	
x <sub>5</sub>	1	

Each datapoint, we have corresponding class labels

The class labels (Actual class labels) associated with these data points are (1, 1, 0, 1, 1)

Now let's assume that our model, which does "binary classification" not only gives us

the class labels associated with the datapoints, It also gives us a score like a probability score. It could be any score here, a score that represents as:- More the score, more is the chance that it belongs to class 1,



or vice versa.

Let's assume that our model gives the probability score associated with each point as:

$x$	$y$	$\hat{y}$
$x_1$	1	0.95
$x_2$	1	0.92
$x_3$	0	0.80
$x_4$	1	0.76
$x_5$	1	0.71

Now the probability score that  $x_1$  belongs to class '1' is let's say 0.95

Similarly for  $x_2$  let's say it is 0.92

$x_3$	"	"	"	0.80
$x_4$	"	"	"	0.76
$x_5$	"	"	"	0.71

Decreasing order of  $\hat{y}$

First thing that you would notice here is, we have sorted our data in decreasing order of  $\hat{y}$ .

1) So, first step in 'ROC' plotting is,

Take all your data & sort it in decreasing order of " $\hat{y}$ ".

2) Once we have sorted our data in decreasing order of " $\hat{y}$ ", we will do "Thresholding". ( $\tau$ )

So, "Thresholding" works as follows:  $\rightarrow$

We can take any value of  $\hat{y}$  for thresholding.

Let's take the first largest value as threshold.

$$\text{as } \tau_1 = 0.95,$$

After taking  $\tau_1 = 0.95$ , we will say that, if  $\hat{y} \geq \tau_1$ , then declare it to be "class 1".

else declare it to be "class 0".

$$\tau_1 = 0.95.$$

$$\text{if } \hat{y} \geq \tau_1$$

else  
0

for every "datapoint" in our dataset, we will get "class label" & we can compute "TPR & FPR".

So, when we have  $\tau_1 = 0.95$ , all the points for which  $\hat{y} \geq 0.95$ , will have a final class label of 1 else it is '0'.

$x$	Actual $y$	$\hat{y}$	$\tilde{y}_{T_1=0.95}$
$x_1$	1	0.95	1
$x_2$	1	0.92	0
$x_3$	0	0.80	0
$x_4$	1	0.76	0
$x_5$	1	0.71	0

Here we are finding associated class label when  $T_1=0.95$  (Threshold=0.95) we represent it with  $\tilde{y}$

So, Corresponding  $\tilde{y}$  given  $T_1=0.95$  we can compute "TPR" & "FPR"

Next we will take  $T_2=0.92$  & we compute  $\tilde{y}$  given  $T_2=0.92$

$x$	$y$	$\hat{y}$	$\tilde{y}_{T_1=0.95}$	$\tilde{y}_{T_2=0.92}$
$x_1$	1	0.95	1	1
$x_2$	1	0.92	0	1
$x_3$	0	0.80	0	0
$x_4$	1	0.76	0	0
$x_5$	1	0.71	0	0

Corresponding to  $T_2=0.92$ , we can even compute "TPR" & "FPR".

Similarly, we can keep changing the threshold values, so if we have "n" points in our dataset we could have "n" thresholds.

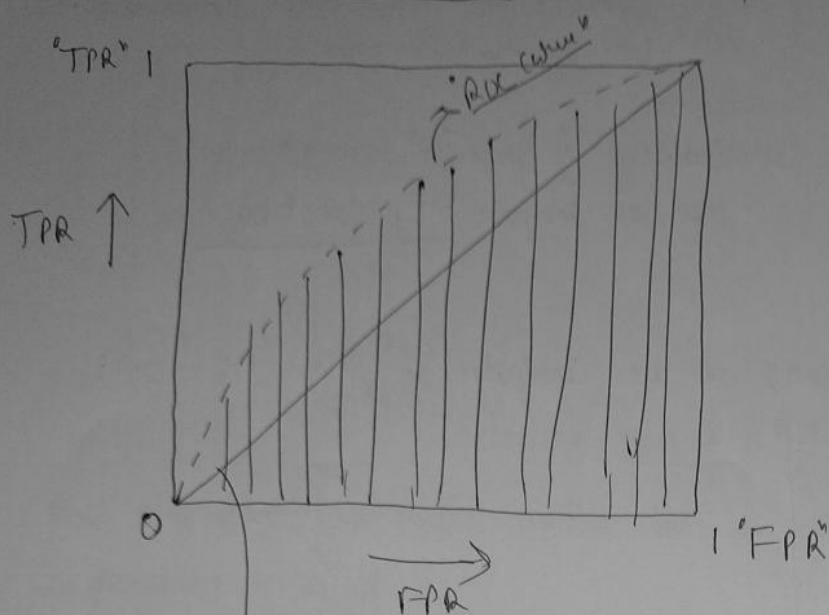
& if we have "n" thresholds ( $T_1, T_2, T_3, \dots, T_n$ )

Now for each " $T$ " we can get correspondingly "FPR" & "TPR".

$T_1$	$T_2$	$T_3$	$T_4$	-----	$T_n$
↓	↓	↓	↓	-----	↓
FPR <sub>1</sub>	FPR <sub>2</sub>	FPR <sub>3</sub>	FPR <sub>4</sub>	-----	FPR <sub>n</sub>
TPR <sub>1</sub>	TPR <sub>2</sub>	TPR <sub>3</sub>	TPR <sub>4</sub>	-----	TPR <sub>n</sub>

Now we will draw a plot, where x-axis is "FPR" & y-axis is "TPR"

Value of 'FPR + TPR' can lie b/w '0.41'



Now we will plot the  $(TPR, FPR)$  pairs.

Typically, if our model is sensible, we could get a line/curve as shown above. This line/curve is called 'Receiver operating characteristic curve'.

Let's draw another line as shown above. This straight line divides the entire region into two parts.  
The area under this straight line is  $= '0.5'$

Now 'AUC' = It basically the total area under this whole 'ROC curve'.  
This total area under the curve (AUC) can lie between

0.41 Higher the value of 'AUC' it is better.

[1 means very good  
0 means terrible]

AUC is only useful for 'binary classification' task.

There are extensions to 'ROC' for multiclass classification also, but in general, we will typically use it for binary classification only.

Some properties of AUC (some downfall of AUC)

(7)

① If we have 'imbalanced data'  $\rightarrow$  AUC can be high even for a dumb model / or improved model.  
Note:  $\Rightarrow$  AUC can be impacted by 'imbalanced data'.

② If we look at AUC we will note that it does not care about the actual 'y'. It is caring only about the sorting of 'y'.

$\therefore$  AUC is not dependent on the y score, it depends only on the ordering of 'y'.

Let's say we have five datapoints in our dataset, & we have two models operating on this dataset. The corresponding 'y's generated by these models 'M<sub>1</sub>' & 'M<sub>2</sub>' is as:-

'x'	'y'	M <sub>1</sub> y	M <sub>2</sub> y
x <sub>1</sub>	1	0.95	0.9
x <sub>2</sub>	1	0.98	0.1
x <sub>3</sub>	0	0.80	0.08
x <sub>4</sub>	1	0.76	0.07
x <sub>5</sub>	1	0.71	0.06

For both these models 'M<sub>1</sub>' & 'M<sub>2</sub>' the AUC curve will be same because sorted order of 'y' for both 'M<sub>1</sub>' & 'M<sub>2</sub>' is going to be same.

$$\underline{AUC(M_1) = AUC(M_2)}$$

because it doesn't care about actual scores;

Note:  $\Rightarrow$  'AUC' is often used in ML a lot especially for 'binary classification'.

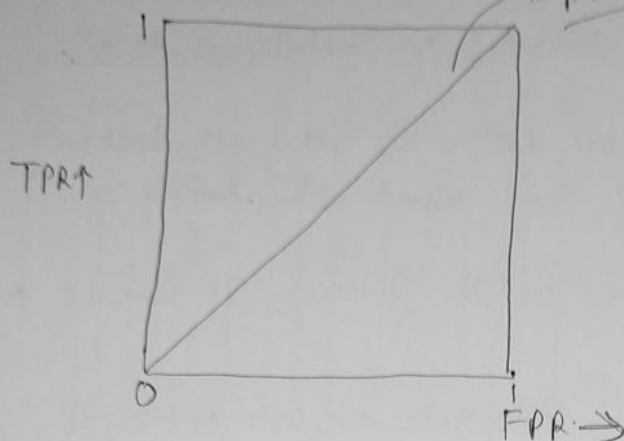
③ Suppose we have a model which randomly decides about the class, given a query point 'x<sub>q</sub>', this model randomly assigns a class label, such a model is called a random model.

Random model

$\hookrightarrow x_q \rightarrow \text{'1 or 0'}$

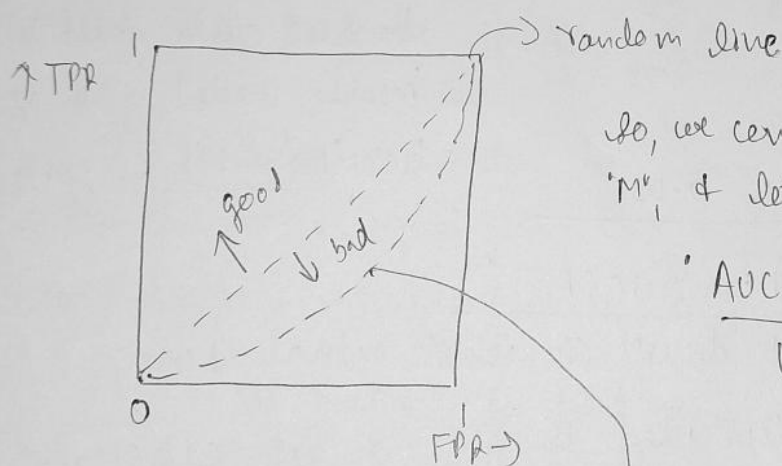


So, if a model is random, not dumb then ROC for that random model will look like a straight line



∴ AUC of a random model will be exactly "0.5".

(4) Suppose, we train a model 'M' & we plotted the 'ROC'



So, we computed the AUC of model 'M', & let's say it is 0.2

$$\underline{AUC(M) = 0.2}$$

Which is worst than random

0.2 basically means that ROC curve looks like as shown above. Everything above the straight line, means it is a good model & everything below that straight line means it is a bad model.

∴ If we see 'AUC' value b/w 0.5 to 1, it means a good model

& if we see 'AUC' value ~~of~~ <sup>at</sup> '0.5', we say that model is not doing anything sensible, it is just a random model.



if AUC lies b/w "0.0 to 0.5" then this is what we do. (9)  
Then if  $\hat{y}_i = 0$  <sup>class label</sup> simply change the class labels  
 $\hat{y}_i = 1$

[ If model o/p is '0' change it to '1'  
If it is '1' change it to '0' ]  
↳ simple swapping of class labels

Now if we do so,  
then the modified model will have an AUC of

$$\frac{1 - 0.2 = 0.8}{\text{which is good}}$$

So, if you ever see a model with an AUC less than 0.5  
simply switch the "class labels". & when we swap the class labels,  
we get an  $AUC = 1 - \text{original AUC}$

In this case " $1 - 0.2 = 0.8$ ".