

★ Accuracy, Precision, Recall, F-1 score etc.
How many such metrics are there?

		Predicted condition		Sources: [12][13][14][15][16][17][18][19] view · talk · edit	
		Predicted positive (PP)	Predicted negative (PN)	Informedness, bookmaker informedness (BM) = $TPR + TNR - 1$	Prevalence threshold (PT) = $\frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
Actual condition	Positive (P) [a]	True positive (TP), hit ^[b]	False negative (FN), miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power = $\frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate type II error ^[c] = $\frac{FN}{P} = 1 - TPR$
	Negative (N) ^[d]	False positive (FP), false alarm, overestimation	True negative (TN), correct rejection ^[e]	False positive rate (FPR), probability of false alarm, fail-out type I error ^[f] = $\frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity = $\frac{TN}{N} = 1 - FPR$
Prevalence = $\frac{P}{P+N}$	Positive predictive value (PPV), precision = $\frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) = $\frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$	Negative likelihood ratio (LR-) = $\frac{FNR}{TNR}$	
Accuracy (ACC) = $\frac{TP + TN}{P + N}$	False discovery rate (FDR) = $\frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP (Ap) = $PPV + NPV - 1$	Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$	
Balanced accuracy (BA) = $\frac{TPR + TNR}{2}$	F ₁ score = $\frac{2 \cdot PPV \times TPR}{PPV + TPR}$ = $\frac{2 \cdot TP}{2 \cdot TP + FP + FN}$	Fowkes-Mallows index (FM) = $\sqrt{PPV \times TPR}$	Matthews correlation coefficient (MCC) = $\frac{\sqrt{TPR \times TNR \times PPV \times NPV}}{\sqrt{FNR \times FPR \times FOR \times FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index = $\frac{TP}{TP + FN + FP}$	

★ No need to remember because all of them are generated from TN, FP, FN, TP

★ Most imp. are: accuracy, precision, recall, F1 score

★ Other imp metrics are: TPR, FPR, TNR, FNR
Predicted

		N	P
Actual =	N _a N	TN	FP
	P _a P	FN	TP
		N _p	P _p

TPR = True Positive Rate

FPR = False " "

TNR = True Negative " "

FNR = False " "

★ Rule: Always divide by Actual

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

Recall

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

$$TP + FN$$

Recall

$$FP + TN$$

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

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

★ Recap Process of Logistic Regression

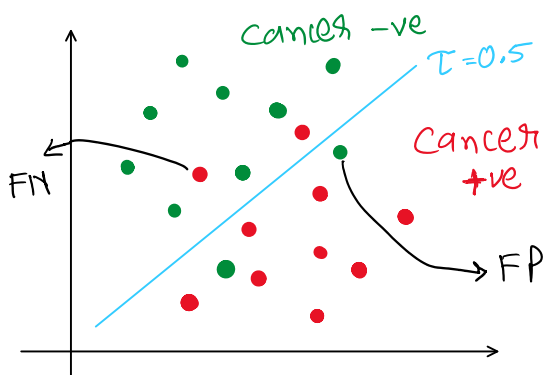
$$x_i \rightarrow w^T x + w_0 \rightarrow \sigma(z_i) \begin{cases} < 0.5 \text{ class-0} \\ > 0.5 \text{ class-1} \end{cases}$$

z_i
 $(-\infty, \infty)$
 $[0, 1]$

Threshold $\Rightarrow \tau = 0.5$

Should we always take $\tau = 0.5$? If not, why & when should we take a different τ ?

Example - Cancer detection

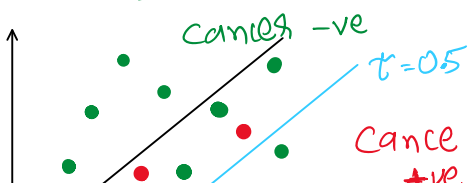


Let's consider FN:

Our model has declared these patients as 'cancer -ve' so, they were not recommended any further test/treatment. This may lead their condition worsen. **This is bad.**

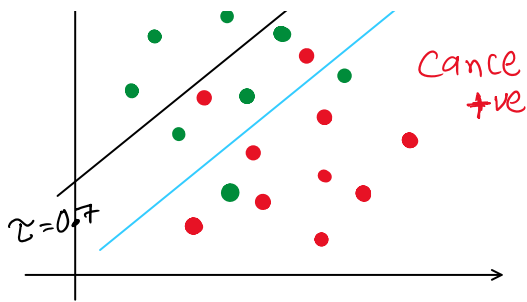
Lets talk about FP now: Our model has declared these patients as 'cancer +ve'. Hence we will recommend them the next round of test. **This is not as bad.**

In such cases when we want our model to be more sensitive to either FN or FP more than the other one we change τ .



This leads to another question:

Which value of τ is



Which value of τ is the 'best choice'?

Answer to this question is:

R.O.C. curve

Receiver's Operating Characteristic Curve

→ How to create an ROC curve?

X_i	Y_i	$\text{Sig}(Z_i)$
X1	1	0.3
X2	1	0.2
X3	1	0.7
X4	0	0.6
X5	0	0.5

step-1: Sort in descending order of $\sigma(Z_i)$

step-2: For each unique value of $\sigma(Z_i)$, take τ & compute \hat{y} (y_{pred})

X_i	Y_i	$\text{Sig}(Z_i)$	$T=0.7$	$T=0.6$	$T=0.5$	$T=0.3$	$T=0.2$
			y_{pred}	y_{pred}	y_{pred}	y_{pred}	y_{pred}
X3	1	0.7	1	1	1	1	1
X4	0	0.6	0	1	1	1	1
X5	0	0.5	0	0	1	1	1
X1	1	0.3	0	0	0	1	1
X2	1	0.2	0	0	0	0	1

step-3: Calculate TPR & FPR for each τ

$\tau = 0.7$

TN	FP
2	0
FN	TP
2	1

$\tau = 0.6$

TN	FP
1	1
FN	TP
2	1

$\tau = 0.5$

TN	FP
0	2
FN	TP
2	1

$\tau = 0.3$

TN	FP
0	2
FN	TP
1	2

$\tau = 0.2$

TN	FP
0	2
FN	TP
0	3

$$\text{TPR} = \frac{1}{3} = 0.3$$

$$\text{TPR} = \frac{1}{3} = 0.3$$

$$\text{TPR} = \frac{1}{3} = 0.3$$

$$\text{TPR} = \frac{2}{3} = 0.6$$

$$\text{TPR} = 1$$

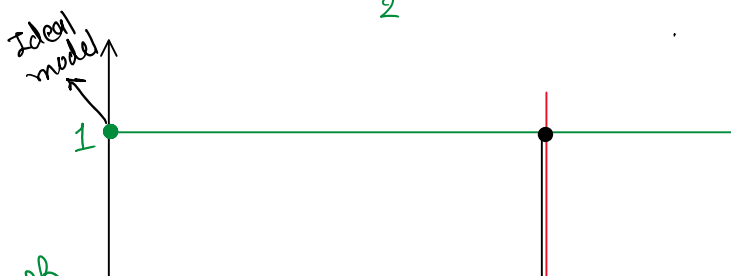
$$\text{FPR} = 0$$

$$\text{FPR} = \frac{1}{2} = 0.5$$

$$\text{FPR} = 1$$

$$\text{FPR} = 1$$

$$\text{FPR} = 1$$

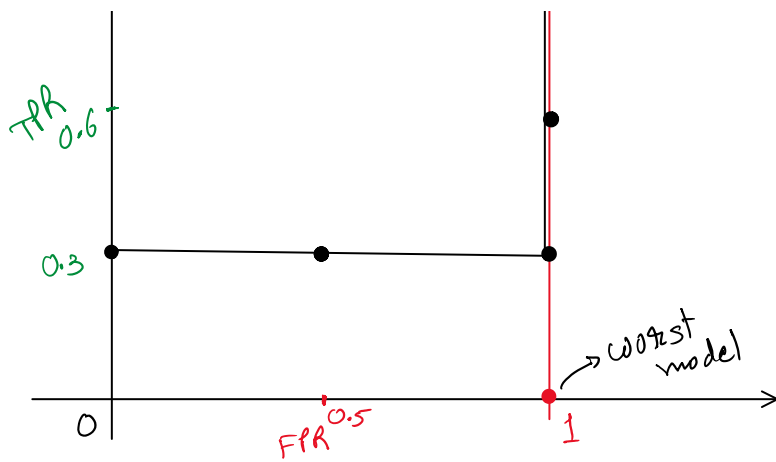


For a worst model:

$$TP = 0, TN = 0, FP = P, FN = Q$$

where $P + Q = N$

$$\text{TPR} = \frac{TP}{(TP + FN)} = 0$$



$$TPR = TP / (TP + FN) = 0$$

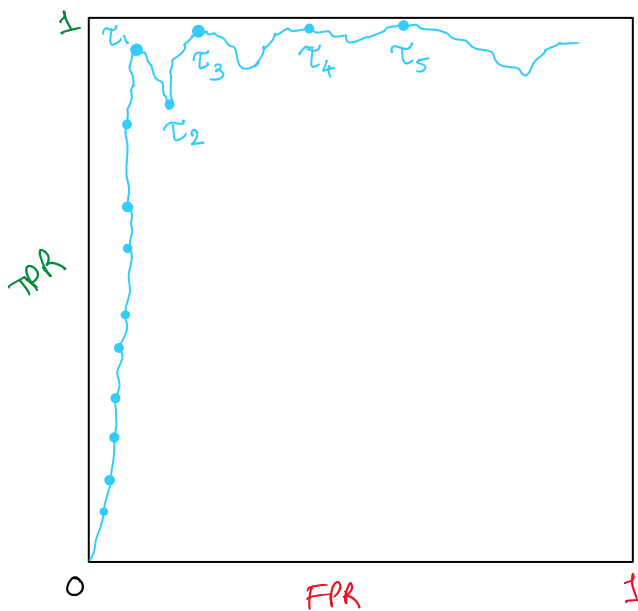
$$FPR = FP / (FP + TN) = 1$$

$$TPR = 0 \text{ \& \; } FPR = 1$$

For a best model:

$$TPR = 1, FPR = 0$$

A real world ROC curve looks like:



Clearly, τ_1 is better than τ_2

But, which one will you select from τ_1

\& \; \tau_3? Usually, classifying points

correctly (TPR) is more important than

caring about mis-classified points (FPR)

(usually, not always) therefore we will select τ_3 .

If the difference in TPR is not significant

then we choose τ with low FPR. eg,

we will choose τ_3 from τ_3, τ_4 \& \; τ_5 .

Imagine working with 15-20 confusion matrices to figure out the best choice of τ vs. pointing it out from an ROC curve!

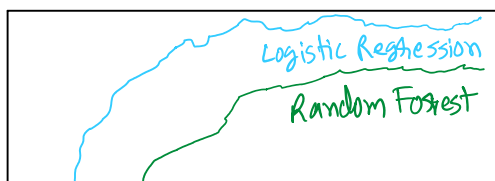
In future, we are going to use multiple different techniques to solve one problem. eg, to classify obese

\& \; non-obese people, we are trying out two ML

techniques (1) Logistic Regression \& \; (2) Random Forest

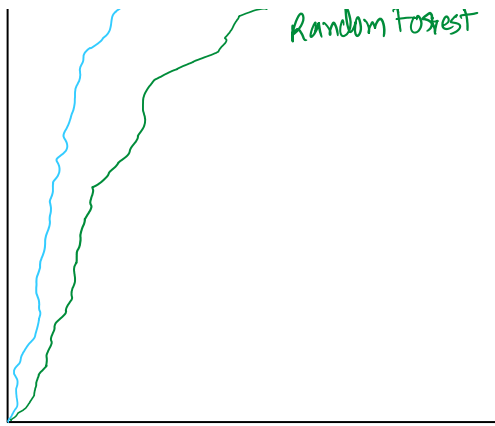
\therefore we will get two ROC curves one for each. Let's say

they are like below:



Looking at the ROC curves, which technique seems better?

Ans- Logistic Regression

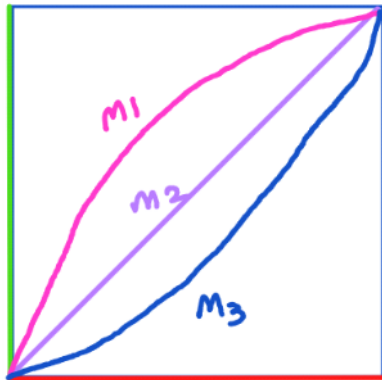


which is better than

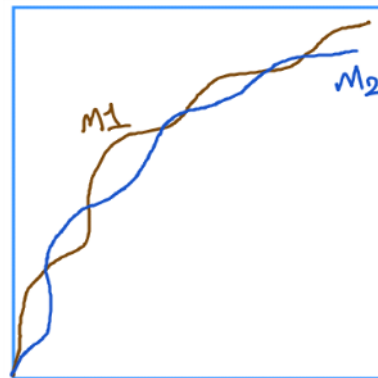
Ans - Logistic Regression

why - Because it has higher TPR compare to Random forest OR its ROC curve is closer to the ideal model.

Sometimes it is easy to find which method is having higher TPR / closer to the ideal ROC curve:



But sometimes it is not that easy to determine which one is better:



To solve this problem we need a **number** which we can use to compare the methods.

The solution is: Area Under the Curve - AUC

AUC - The method with higher AUC is better than the one with lower AUC.

This AUC is also called: **ROC-AUC**

The limitation of roc-auc is that it doesn't work well with imbalanced data.

For imbalanced data, we don't create ROC instead, we create **PRC** (Precision Recall Curve) & use prc-auc

Class Imbalanced

There is no clear distinction available about balanced/imbalanced data. But a widely accepted perception is:

Imbalance starts →

- 50-50 = balanced
- 60-40 = slightly balanced
- 70-30 = slightly imbalanced
- 80-20 = imbalanced
- 90-10
- 95-5
- 99-1
- ⋮

} Highly imbalanced

→ How do we find whether the data is balanced or not?

① `value_counts()`

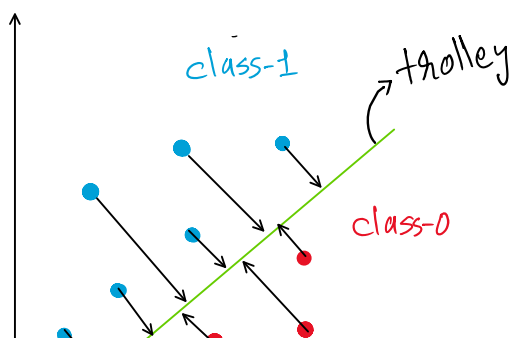
② `countplot`

→ What are the problems with imbalanced data?

① Accuracy and other a few metrics will no longer be reliable because

② Model starts becoming biased towards the majority class.

Understanding this problem:



solution-1: class-weights

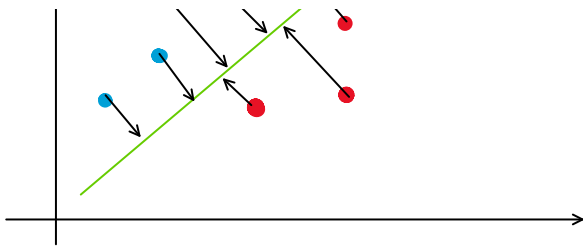
{ • : 1, • : 2 }

$L_1 L_2 L_3 L_4 L_5 L_6 L_7 L_8 L_9$

$\sum(L_i)$

→ 1*

→ 2*

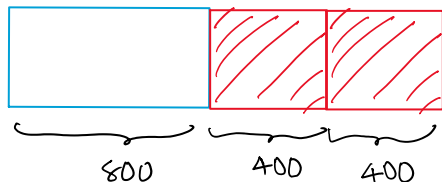


$$\sum(L_i)$$

$$-\frac{1}{n} \sum y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)$$

1* 2*
0 0

Solution-2: Over sampling - Duplicate the present minority



class data randomly until its size matches with the majority class.

Random sampling with repetition

Advantages: ① No data loss

Disadvantages:

① Duplication of data

② Chances of overfitting increases.

③ Under sampling - Randomly select as many data points from majority class as in the minority class & only keep them.

eg, if minority class has 100 datapoints & majority class has 900 then we will randomly select 100 data points from majority class & train our model on these 200 datapoints only.

Advantage:

① Computationally in-expensive

Disadvantage:

① Data Loss

② Our randomly selected sample may not capture all the edge cases as well as the complete pattern of the entire data.