

# ch5: Performance Measurement of Models

## Classification Metrics

### VI. Accuracy

Some classification Metrics:  
Accuracy is performed on Test.

Accuracy =  $\frac{\text{Total } \# \text{ of pts correctly classified}}{\text{Total } \# \text{ of pts. in } D_{\text{Test}}}$

Sensitivity. Pos +  $\frac{\text{Pos} + \text{neg}}{\text{Postneg}}$

Specificity. neg  $\frac{\text{Pos} + \text{neg}}{\text{Postneg}}$

(Check soln @ end)  $\frac{TP + FN}{D_{\text{Test}}}$

Accuracy  $\in [0, 1]$

worst — best

It is measured on test set usually.

let test pts = 100

Correct errors

60 +ve  $\longrightarrow$  53 +ve 7 -ve

40 -ve  $\longrightarrow$  35 -ve 5 +ve

errors = 12

correctly classified  
= 88

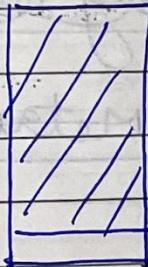
Accuracy =  $\frac{88}{100} = 0.88$

## Problems with accuracy :-

- Doesn't work well with ~~im~~ imbalanced data.

Suppose,

Test data has



90% -ve

10% +ve

~~$$\begin{matrix} \text{TP} & \times \text{TN} \\ + \text{FP} & \text{FP} + \text{TN} \end{matrix}$$~~

→ dumb model  $x_q \xrightarrow{\text{all pk.}} -\text{ve.}$

$$\text{Accuracy} = 0.9.$$

So, never use accuracy as a performance measure for imbalanced data.

- Assume  $M_1$  &  $M_2$  return probability score (Refer KNN to see how we calculate prob. score)

Test set has:-

$x$	$y$	Model 1	Model 2	$\hat{y}_{M_1}$	$\hat{y}_{M_2}$
$x_1$		0.9	0.6	1	0
$x_2$	1	0.8	0.65	1	1
$x_3$	0	0.1	0.45	0	0
$x_4$	0	0.16	0.48	0	0

- ① we can clearly see that  $M_1$  is better than  $M_2$  by looking at the prob. score.
- ② But since predicted class labels are same, we get same accuracy  
~~So, the two can have prob. score~~  
 its evaluating Accuracy understood done ✓
- V2) Confusion matrix, TPR, FPR, FNR, TNR

It doesn't process probability scores.

Consider a binary classification task. 2 classes, either 0 or 1.

		Actual		Predicted		Desired	
		0	1	0	1	$x_i$	$y_i$
Actual	0	a	b	0	$\hat{y}_i$	$x_i$	$\hat{y}_i$
	1	c	d	1	$\hat{y}_i$	$x_i$	$\hat{y}_i$
						$x_1$	$y_1$
						$x_2$	$y_2$
						$x_n$	$y_n$

a: # of pts. s.t.  $y_i = 0$  &  $\hat{y}_i = 0$   
 b: # of pts. s.t.  $y_i = 1$  &  $\hat{y}_i = 0$   
 c: # of pts. s.t.  $y_i = 0$  &  $\hat{y}_i = 1$   
 d: # of pts. s.t.  $y_i = 1$  &  $\hat{y}_i = 1$

↓  
 data point  
 ↓  
 actual class label  
 ↓  
 predicted class label

Actual on top  
 Predicted on left

This can be extended to multiclass classification.

For a sensible model, principal diagonal elements are higher and off diagonal elements are small.

$y_i$	0	1
0	$TN$	$FN$
1	$FP$	$TP$

$N \Rightarrow P = (\text{column total})$

eg  $\frac{TP}{P}$

Actual | Predicted  
 $(y_i)$  |  $(y_i)$

Negative

$N = \text{Total Negatives}$  then Positive row wise  
 $P = " \text{Positives}$   
 $n = N + P$

True on diagonal  
False on off diagonal

✓  $TPR = \text{True Positive Rate} = \frac{TP}{P}$   
(Sensitivity)

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

✓  $TNR = \text{True Negative Rate} = \frac{TN}{N}$   
(Specificity)

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

✓  $FPR = \text{False Positive Rate} = \frac{FP}{N}$

$$\frac{FP}{FN + FP}$$

✓  $FNR = \text{False Negative Rate} = \frac{FN}{P}$

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

eg -

	$\hat{y}_i$	
$\hat{y}_i$	0	1
0	850	6
1	50	94
	900	100
	N	P

$$TPR = 94\% = \frac{94}{100} \quad FPR = \frac{50}{900}$$

$$TNR = \frac{850}{900} \quad FNR = 6\% = \frac{6}{100}$$

So we have high TPR & TNR,  
which is desired.

Test set has 900 -ve & 100 +ve

Ideally,  $TPR \uparrow$   $FPR \downarrow$   
 $TNR \uparrow$   $FNR \downarrow$

Let us now assume, a dumb model  
which gives opp label always -ve.  
Then

	$\hat{y}_i$	0	1
$\hat{y}_i$	0	900	0
0	900	0	FN
1	0	0	TP

$$TPR = 0\% = \frac{0}{100}$$

$$FPR = 0\% = \frac{0}{900}$$

$$TNR = 100\% = \frac{900}{900}$$

$$FNR = 100\% = \frac{100}{100}$$

So, we can detect dumb models by using these 4 metrics together.

Now, which of the 4 metrics is more important?

It is domain specific.

Eg- Diagnosing Cancer.

We must have very near or zero FNR as we can't afford to say an actual patient that he doesn't have cancer if he has.

Why, Little high FPR is acceptable as even if we say a non-cancer patient that he has cancer, later through powerful tests that can be classified.

Don't miss a cancer patient!

$TN$	$FN$	$TP$	$FP$	$TN$	$FP$	$TP$	$FN$
<del>under</del>	<del>good</del>	<del>done</del>	<del>likely</del>	<del>NEG</del>	<del>POS</del>	<del>POS</del>	<del>NEG</del>
<del>Yi</del>	<del>0</del>	<del>1</del>	<del>0</del>	<del>0</del>	<del>1</del>	<del>1</del>	<del>0</del>

### V3. Precision, Recall and F1-score

This is based on Information Retrieval.

e.g. - collecting only 10% out of Millions of datapoints

<del>Yi</del>	<del>0</del>	<del>1</del>	<del>TP</del>	<del>FP</del>	<del>TN</del>	<del>FN</del>	<del>Precision</del>
<del>Yi</del>	<del>0</del>	<del>1</del>	<del>TP</del>	<del>FP</del>	<del>TN</del>	<del>FN</del>	$\frac{TP}{TP+FP}$
<del>Yi</del>	<del>0</del>	<del>1</del>	<del>TP</del>	<del>FP</del>	<del>TN</del>	<del>FN</del>	$\frac{TP}{TP+FP}$
<del>Yi</del>	<del>0</del>	<del>1</del>	<del>TP</del>	<del>FP</del>	<del>TN</del>	<del>FN</del>	$\frac{TP}{TP+FP}$

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

(only in context of positive class)

of all the points the model declared to be +ve, what %age of them are actually +ve.

(Sensitivity)

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

(only in context of positive class)

of all the points which actually belong to class +ve, how many are predicted to be +ve.

<del>Yi</del>	<del>0</del>	<del>1</del>	<del>TP</del>	<del>FP</del>	<del>TN</del>	<del>FN</del>
<del>Yi</del>	<del>0</del>	<del>1</del>	<del>TP</del>	<del>FP</del>	<del>TN</del>	<del>FN</del>

We want both Pr & Re to be high.

[0,1], [0,1]

F1-Score combines them both:-

$$\text{F1-Score} = 2 * \left[ \frac{\text{Pre} * \text{Rec}}{\text{Pre} + \text{Rec}} \right]$$

\* Num  
+ Deno

Formula is derived from P.T.O

Specificity =  $\frac{\text{True Negative}}{\text{True Neg} + \text{False Positive}}$   
 $(\text{Negative})$

Harmonic mean of  $P_r$  &  $R_e \Rightarrow F_1$ -score

$$\frac{1}{F_1} = \frac{1}{2} \left( \frac{1}{P_r} + \frac{1}{R_e} \right) \text{ lies b/w } 0 \& 1$$

V4: Receiver Operating characteristic Curve

& AUC (Area Under <sup>this</sup> ~~Curve~~)

Used for binary classification only.  
 Assume model outputs a score which  
 can be interpreted similar to probability  
 of class 1 score.

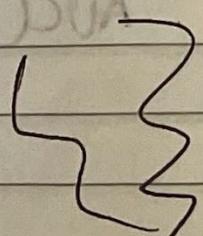
$x_i$	$y_i$	$\hat{y}_i$
$x_1$	1	0.95
$x_2$	1	0.92
$x_3$	0	0.80
$x_4$	1	0.76
$x_5$	1	0.71

Steps for ROC :

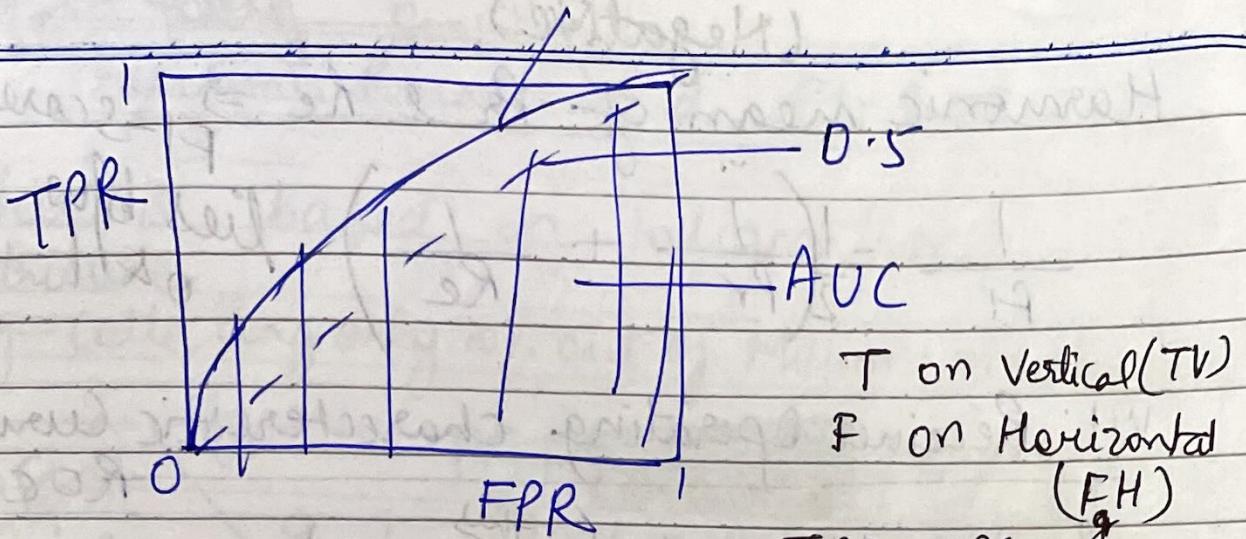
- 1) Sort entries in decreasing order of  $\hat{y}_i$
- 2) Thresholding ( $T$ )
- 3) Find TPR & FPR for each  $T$   
 $\text{So, } T_{1,0.95} \text{ if } \hat{y}_i \geq T \text{ label = 1}$   
 $\text{else label = 0.}$

$x_i$	$y_i$	$\hat{y}_i$	$\hat{y}_{T=0}$ as $y_{T=0}$	$y_{T=1}$
-------	-------	-------------	------------------------------	-----------

$x_1$	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



## ROC curve

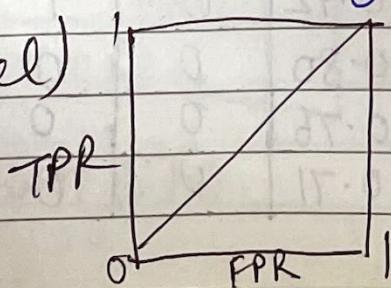


For a sensible model this will have ~~have~~ the curve like above.

### AUC properties:-

- 1) For imbalanced data, AUC can be high, ~~and~~ also for dumb/simple model also.
- 2) AUC is not dependent on the  $\hat{y}$  score. It ~~is~~ only depends on ordering of data.
- 3)  $\text{AUC}(\text{random model}) = 0.5$  which might be  $0.1$
- 4) If  $\text{AUC}(M) \in [0, 0.5]$ , then, simply compliment the o/p of the model. (swapping). change the class label or reverse the order of output's swap

$\text{AUC}(\text{Random Model})$



VS: Log loss (uses probability score)

On Test set of n-pts :-

$$\text{log loss} = -\frac{1}{n} \sum_{i=1}^n \left[ (\log(p_i) * y_i) + (\log(1-p_i) * (1-y_i)) \right]$$

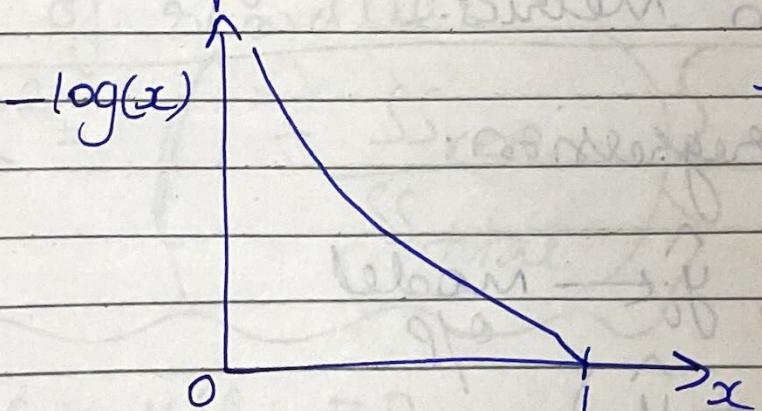
either of the term  
is valid, the other becomes  
zero

negative  
of average  
of log  
loss

Average of negative of log of probabilities  
of correct class label.

All  $\log$  loss, smaller the better.  
" prob., higher " "

Log loss is very powerful because  
unlike all previous metrics, it makes use of actual  
class probabilities.



So, more closer  
the correct class  
probability as to 1,  
lesser the loss.

For multiclass classification,  
total  $\sum_{j=1}^n p_{ij}$  c-classes

$$\text{log loss} = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c y_i * \log(p_{ij})$$

if  $x_i \in \text{class } j$   
0 otherwise

Prob. of  $x_i$   
 $\in \text{class } j$

- 1) One key disadvantage of log loss is it's hard to interpret, as it can go as high as  $\infty$ .
- 2) No upperbound, can't determine how bad.

### Regression Metrics

V6:  $R^2$  or Coefficient of determination

Some regression metrics.

$y_i \in \mathbb{R}$  is regression.

Test  
Set  
 $i=1 \text{ to } n$

$x_i, y_i, \hat{y}_i \leftarrow$  model  
op

$$e_i = (y_i - \hat{y}_i)$$

error  
for  
 $x_i$

Just like in classification, a simple model would be one that returns the label of the majority class for all inputs.

Similarly, for regression, a simple model would be one that returns the mean of  $y_i$ 's.

Sum of squared error for the simple mean model

$$\text{Sum of squared error} = \sum_{i=1}^n (y_i - \bar{y}_i)^2$$

Note:  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

Also called  $SS_{\text{Total}}$

$$\text{Residue} = \hat{e}_i = (y_i - \hat{y}_i)$$

$$\text{Sum of squares of residue} = SS_{\text{res.}} = \sum_{i=1}^n (\hat{e}_i)^2$$

$$R^2 = \left( 1 - \frac{SS_{\text{res.}}}{SS_{\text{total}}} \right)$$

Actual metric

$$\text{Case 1: } \hat{e}_i = 0$$

$$\Rightarrow e_i = 0$$

$$\Rightarrow R^2 = 1 \quad (\text{Best value})$$

$SS_{\text{res.}} \rightarrow$  closer to 0 better it is

Case 2:  $SS_{\text{res.}} < SS_{\text{total}}$

$$\Rightarrow R^2 \in [0, 1]$$

Case 3:  $SS_{\text{res.}} = SS_{\text{total}}$

$$\Rightarrow R^2 = 0$$

$\Rightarrow$  Same as simple mean model.

Case 4:  $SS_{\text{res.}} > SS_{\text{total}}$

$$\Rightarrow R^2 = 1 - (\underbrace{\text{grs} > 1}_{\text{greater than some value}})$$

$\Rightarrow$  Negative.

$\Rightarrow$  Means it is worse than a simple mean model.

Read Wikipedia  $\rightarrow$  Coefficient of determination

V7: Median Absolute Deviation of Errors [MAD]

$R^2 \Rightarrow$  Relies on mean, hence not very robust to outliers (if one is very large it becomes an outlier)

For each  $x_i \xrightarrow{\text{we have}} y_i > \hat{y}_i \rightarrow e_i$

If  $e_i$  = random variable.

median( $e_i$ ) = central value of errors  
(is a robust measure of mean)

MAD( $e_i$ ) = Median of  $(|e_i - \text{median}(e_i)|)$   
(is a robust measure of std dev.)

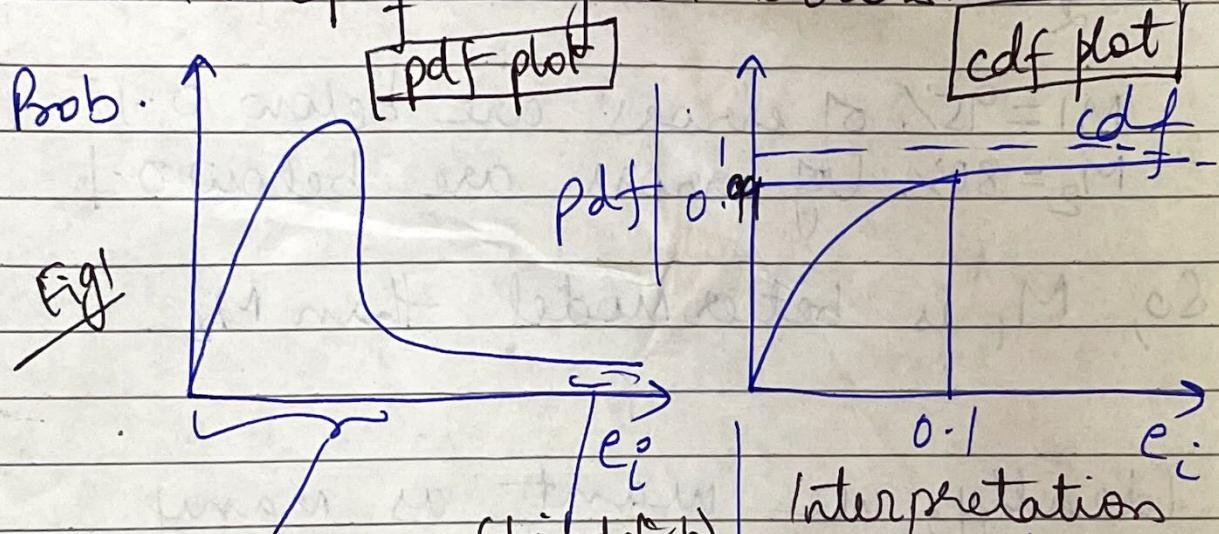
Defined for a set.

Use median & MAD if outliers are present.

Use mean & std. dev. if no outliers

### V 8 : Distribution of Errors

We can use pdf & cdf of errors.



Most  $e_i$   
are  
small

(tailed dist.)  
very  
few  
 $e_i$  are  
large

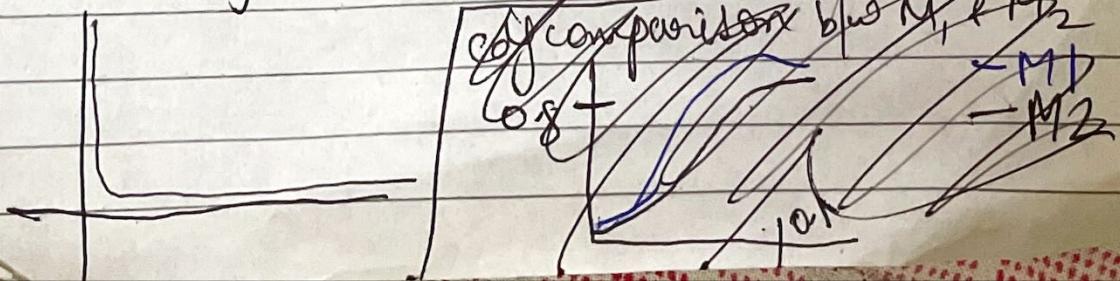
Interpretation  
99% of  $e_i$  are  
 $< 0.1$   
1% of  $e_i \geq 0.1$

This tells that large model is good.

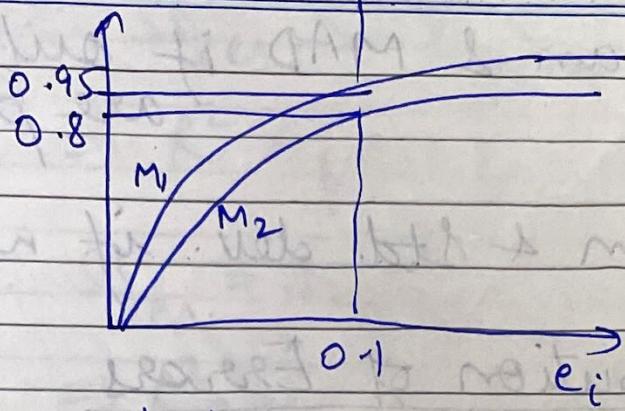
No cdf can be  $>$  than 1.

So, we can use these standard statistical measures to get an idea of distribution of errors.

Ideal pdf in Dist of errors should be :-



e.g -



$M_2$  is below  $M_1$ ,

$M_1 = 95\%$  of errors are below 0.1

$M_2 = 80\%$  of errors are below 0.1

So,  $M_1$  is better model than  $M_2$ .

Ideally, we want as many errors close to 0 as possible.