

Performance Measurement of Models

* Accuracy: $\frac{\# \text{ correctly classified points}}{\text{Total \# of pts in Dist}}$

Ranges from 0 to 1.
 Bad \nearrow Good \nwarrow

It is easy to understand.

\Rightarrow It should not be used if Data is imbalanced.

\Rightarrow It can't be used to compare 2 models if prob-score is different.

X	y	Probability Scores		\hat{y}_1	\hat{y}_2
		M_1	M_2		
x_1	1	0.9	0.6	1	1
x_2	1	0.8	0.65	1	1
x_3	0	0.1	0.45	0	0
x_4	0	0.15	0.48	0	0

Accuracy $M_1 >$ Accuracy M_2

On the basis of probability scores, $M_1 > M_2$, M_1 model is better than M_2 .

* Confusion Matrix

In case of Binary Classification Task (0, 1)

		Actual \rightarrow	
		0	1
Pred \downarrow	0	a	b
	1	c	d

a \rightarrow no. of pts such that $y=0$ & $\hat{y}=0$
 b \rightarrow " " " " " " $y=1$ & $\hat{y}=0$

Multiclass Classification

Bred \ Actual →	0	1	2	3
0 ↓	↑			
1 ↓		↑		
2 ↓			↑	
3 ↓				↑

- Principal diagonal elements ↑
- Off diagonal element should be ↓

Bred \ Actual →	0	1
0 ↓	TN	FN
1 ↓	FP	TP
	N	P

T P
 ↗ ↘ What is predicted label
 ↘ Prediction is correct or not

N → No. of -ve pts

P → No. of +ve pts

True Negative Rate, $TNR = \frac{TN}{N}$

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

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

True Positive Rate, $TPR = \frac{TP}{P}$

TNR + TPR should be ↑

FNR + FPR should be ↓

Out of TNR and TPR, what should be Higher?

Ans. It is very domain specific.

* Precision, Recall and F1-Score (They are more focused on +ve class)

Bred \ Act →	0	1
0 ↓	TN	FN
1 ↓	FP	TP

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

(True predicted to be true out of labels that are actually true out of pts that are predicted to be true.)

$$\text{Recall} = \text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad \left(\frac{\text{Correctly predicted +ve pts}}{\text{total actual +ve pts}} \right)$$

F1-Score :- Harmonic mean of Precision and Recall.

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

Receiver Operating Characteristic Curve (ROC) and AUC

(Primarily used for Binary classification)

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	0.7

→, First y then \hat{y} are arranged in decreasing order

⇒ Second, thresholding (Z) is done with each value of \hat{y} .

i.e.,

if $\hat{y} \geq Z$,

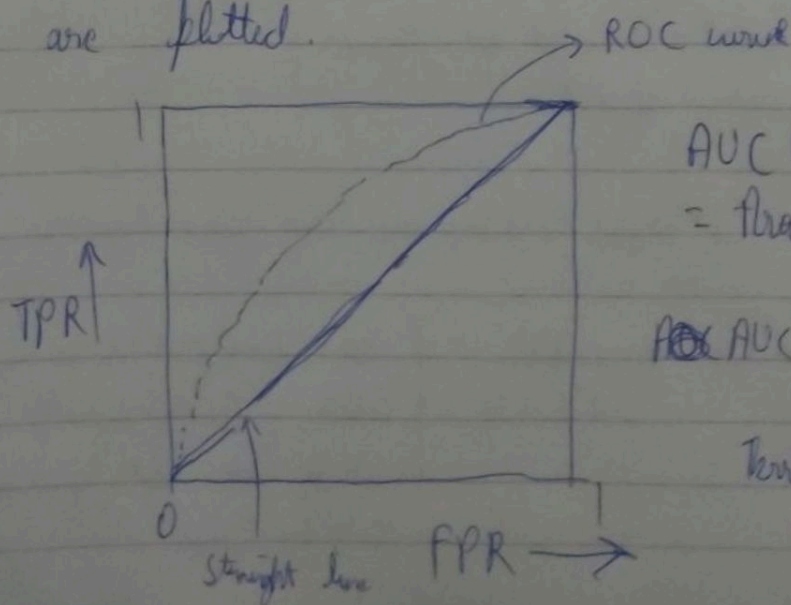
return 1

else

return 0

So, after this thresholding, ~~for~~ FPR & TPR are calculated for each value of Z or \hat{y} .

Then they are plotted.



AUC (Area under curve)

= Area under ROC curve

AUC → 0 to 1

Terrible

V. good

AUC of a model

Some Properties of AUC:-

- ⇒ Imbalanced data ⇒ AUC can be high for a dumb model too.
- ⇒ AUC is not dependent on \hat{y} scores, it only depends on the ordering of \hat{y} scores.
- ⇒ AUC of a random model is 0.5.
- ⇒ If AUC of a model is lesser than 0.5, just swap the predicted labels $0 \leftrightarrow 1$ to get $AUC_{\text{final}} = 1 - AUC_{\text{initial}}$

Log-Loss

- One of the best metric for classification problems.
Range :- $[0, \infty)$, lower the better

X	y_i	$\hat{y}_i = P_i$	log loss
x_1	1	0.9	$-\log(0.9) = 0.0457$
x_2	1	0.6	$-\log(0.6) = 0.22$
x_3	0	0.1	$-\log(0.9) = 0.0457$
x_4	0	0.4	$-\log(0.6) = 0.22$

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

↑ Actual values

It penalizes the small deviations in prob-score.

Multi-Class Logloss $\Rightarrow -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c y_{ij} \log(P_{ij})$

\downarrow
 1 if $x \in \text{class } j$
 0 otherwise

↗ Prob that $x \in \text{class } j$

Median Absolute Deviation of errors

R^2 - Coefficient of determination is not robust to outliers.

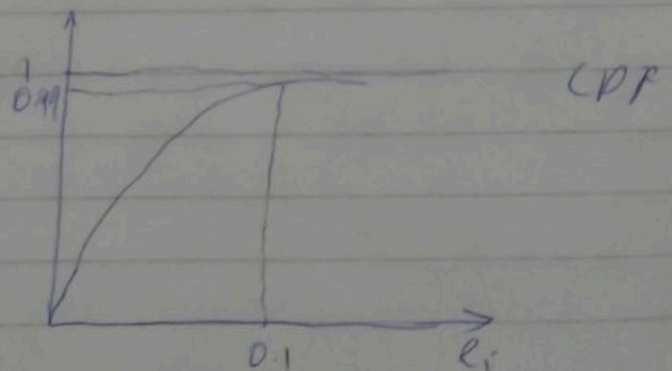
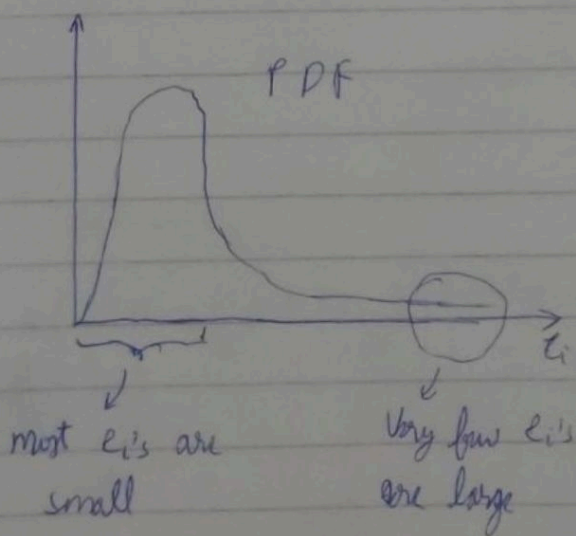
$$e_i \rightarrow y_i, \hat{y}_i$$

mean \leftarrow median(e_i) = Central value of errors
 $MAD(e_i) = \text{median}(|e_i - \text{median}(e_i)|)$

↓
Standard deviation

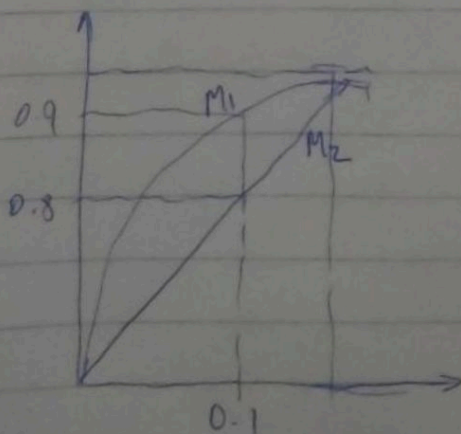
Distribution of errors.

y_i, \hat{y}_i, e_i (errors), we can also create PPF & CDF for e_i 's



99% of errors are < 0.1

Comparison of Models using CDF:-



For model M_1 :- 90% errors < 0.1

For model M_2 :- 80% errors < 0.1

M_1 is better than M_2 .