

★ Performance measurements:-

① In this chapter we will look at how to measure the performance of models.

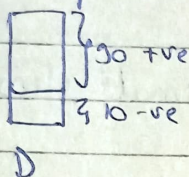
① Accuracy :- $\frac{\text{No. of correctly classified points}}{\text{Total No. of points in Test set.}}$

Accuracy lies b/w 0 to 1
bad \rightarrow 0, good \rightarrow 1

There are some problems with Accuracy

① Imbalanced dataset :- Accuracy may be high even for wrong prediction (bad model)

suppose Ex



And model predicts all points as +ve.

Accuracy = 90% even the model is not doing anything.

② Accuracy cannot use prob score

② Confusion matrix :-

Let's start with binary classification task : (0,1)

Predict \ Actual	0	1
0	a	b
1	c	d

Dataset	x_1	y_1	g_1
	x_2	y_2	g_2
	x_3	y_3	g_3
	\vdots	\vdots	\vdots
	x_n	y_n	g_n

• Confusion matrix cannot process prob scores

a = No. of points which are actually zero and predicted as zero.

b = No. of points which are actually 1 and predicted as zero

c = " " " " " " " 0 " " " "

d = " " " " " " " 1 " " " "

In a multi-class classifier, we can draw a $C \times C$ matrix

Actual y_i

Predict

	0	1	...	C
0				
1				
...				
C-1				

multiclass

If a model is good. then diagonal value of confusion matrix should be high.

	0	1
0	↑	
1		↑

Actual

Predict

	0	1
0	TN	FN
1	FP	TP

FN: Predicted = Negative → incorrect ∵ actual = 0

TP: Predict = 1 = Positive is it correct? Yes! T ∴ TP

what is predicted label

are you correct

FP: Predicted positive = 1, "incorrect predicts b/w actual = 0"

A

P

	0	1
0	TN	FN
1	FP	TP

$$N + P = n$$

Total
Negati
(N)

Total positive (P)

$$\text{True positive rate} = \frac{TP}{P}$$

$$\text{True negative rate} = \frac{TN}{N}$$

$$\text{False positive rate} = \frac{FP}{N}$$

$$\text{False negative rate} = \frac{FN}{P}$$

⑩ Precision, Recall & F1-score

Precision, Recall are very useful in information Retrieval.

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

(of all the points the model predicted to be +ve what percentage of them are actually +ve.)

$$\text{Recall} = \text{True positive rate} = \frac{TP}{P}$$

* Precision & Recall are only P rating about +ve class.

~~Can we~~ we want Precision \uparrow , Recall \uparrow
 can we combine the idea of precision & recall.

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

(ii) Receiver operating characteristic Curve (ROC curve) & AUC
 (Area under the curve)

	y	\hat{y} predicted
x_1	1	0.95
x_2	1	0.92
x_3	0	0.88
x_4	1	0.76
x_5	1	0.71

Let binary ~~sc~~ classifier not only give class label, it also gives a score.
 more score mean more chance of point to belong to that class.

step 1: ~~take~~ sort the data in decreasing order of \hat{y} .

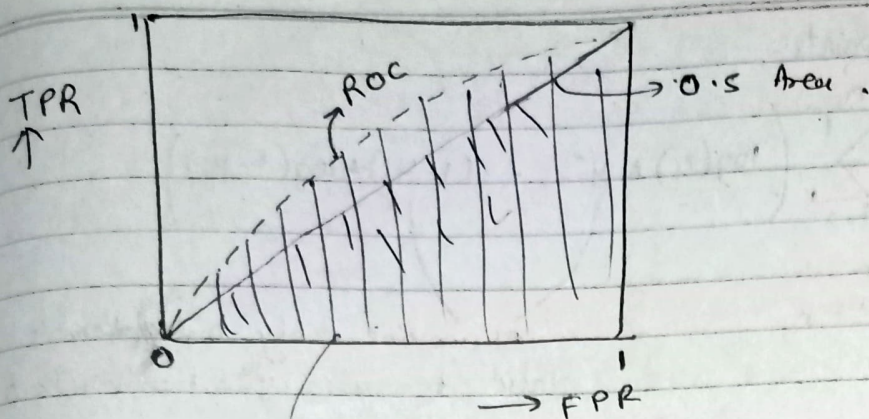
step 2: Thresholding:- Initially the first value (max value) is threshold, T .

if $\hat{y} \geq T$
 label = 1
 else
 label = 0

calculate the TPR, & FPR corresponding to the predicted class label

	y	\hat{y}	$\hat{y}, T=0.95$	$\hat{y}, T=0.92$	$\hat{y}, T=0.88$	$\hat{y}, T = \text{least value}$
x_1	1	0.95	1	1	1	...
x_2	1	0.92	0	1	1	
x_3	0	0.88	0	0	1	
x_4	1	0.76	0	0	0	
x_5	1	0.71	0	0	0	
			TPR ₁ FPR ₁	TPR ₂ FPR ₂	TPR ₃ FPR ₃	TPR _N FPR _N

plot these TPR_i, FPR_i



→ AUC = Area under ROC curve

↳ 0 to 1
bad ← → v. good

This is only useful to binary classification.

AUC have following property

- (i) AUC can be impacted by imbalanced data: AUC have large value for even dumb model.
- (ii) AUC didnot dependent on the \hat{y} scores themselves, it is only dependent on ordering.

Ex:-

	x_i	y_i	\hat{y}_1	\hat{y}_2
	x_1	1	0.95	0.2
	x_2	1	0.92	0.1
	x_3	0	0.80	0.08
	x_4	1	0.76	0.07
	x_5	1	0.71	0.06

↳ Order of the value are same, i.e. they are already in sorted order, hence their ~~for~~ AUC will be similar.

(iii) AUC of a random model will be exactly 0.5.

(iv) If AUC < 0.5 just reverse the model outputs.

log-loss :-

It uses the probability scores.

x	y	\hat{y}	
x_1	1	0.9	→ p_1
x_2	1	0.6	→ p_2
x_3	0	0.1	→ p_3
x_4	0	0.4	→ p_4

Given a Test-set of n -points: $\text{log-loss} = -\frac{1}{n} \sum \left(\log(p_i) y_i + (1-y_i) (1-p_i) \right)$

Test-set of n points:-

$$\text{log-loss} = -\frac{1}{n} \sum_{i=0}^n \left(\log(p_i) * y_i + (1-y_i) * \log(1-p_i) \right)$$

-ve of average

At one time only one of them is valid bcz when $y_i = 1 \Rightarrow (1-y_i) = 0$

$$y_i = 0 \Rightarrow (1-y_i) = 1$$

We want log-loss to as small as possible. it can lie b/w 0 to ∞ .

We can easily extend log loss to multi-class.

Let we have c classes.

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

\uparrow
 $= 1$ if x_{ij} b/w class j
 $= 0$ otherwise.

Only draw of log loss is, it's hard to interpret.

(41) R^2 or coefficient of determination

We have looked at various classification model measurement. Now look at regression performance measurement.

$$\text{error}, e_i = y_i - \hat{y}_i$$

$$\text{Total sum of square} = \sum_{i=1}^n (y_i - \bar{y})^2$$

average value of all y_i in data.

Simple - mean model:- it ~~is~~ for any query x_i it return the mean of \bar{y} .

$$SS_{\text{total}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SS_{\text{residuals}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

actual predicted value

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

Case 1: $SS_{\text{res}} = 0$, if there is no error
then $R^2 = 1$ (but value)

Case 2: $SS_{\text{res}} < SS_{\text{total}}$; $R^2 \equiv 0$ to 1

Case 3: $SS_{\text{res}} = SS_{\text{total}}$; $R^2 \equiv 0$ { It is a simple mean model }

Case 4: $SS_{\text{res}} > SS_{\text{total}}$; $R^2 < 0$ { model is worse than simple mean model }

to calculate R^2 we use SS_{res} which uses mean
hence R^2 is not very robust to outliers.

If we can image error as R.V we can calculate median of e_i 's

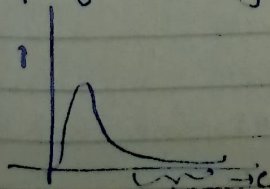
$$MAD(e_i) = \text{media}(|e_i - \text{median}(e_i)|)$$

So, instead of using R^2 we can use Median of e_i
or $MAD(e_i)$ as a measure.

If e_i median(e_i) is less model is good.
or $MAD(e_i)$ is less model is good.

→ Distribution of errors:-

we can also use pdf, cdf of e_i 's to calculate the performance of regression model.



↳ very few e_i are large hence model is good?