

Performance measurement of models

-1-

⇒ Accuracy ⇒ measure of a metric how good our model is.

$$\text{Accuracy} = \frac{\# \text{ correctly classified points}}{\# \text{ total number of points}}$$

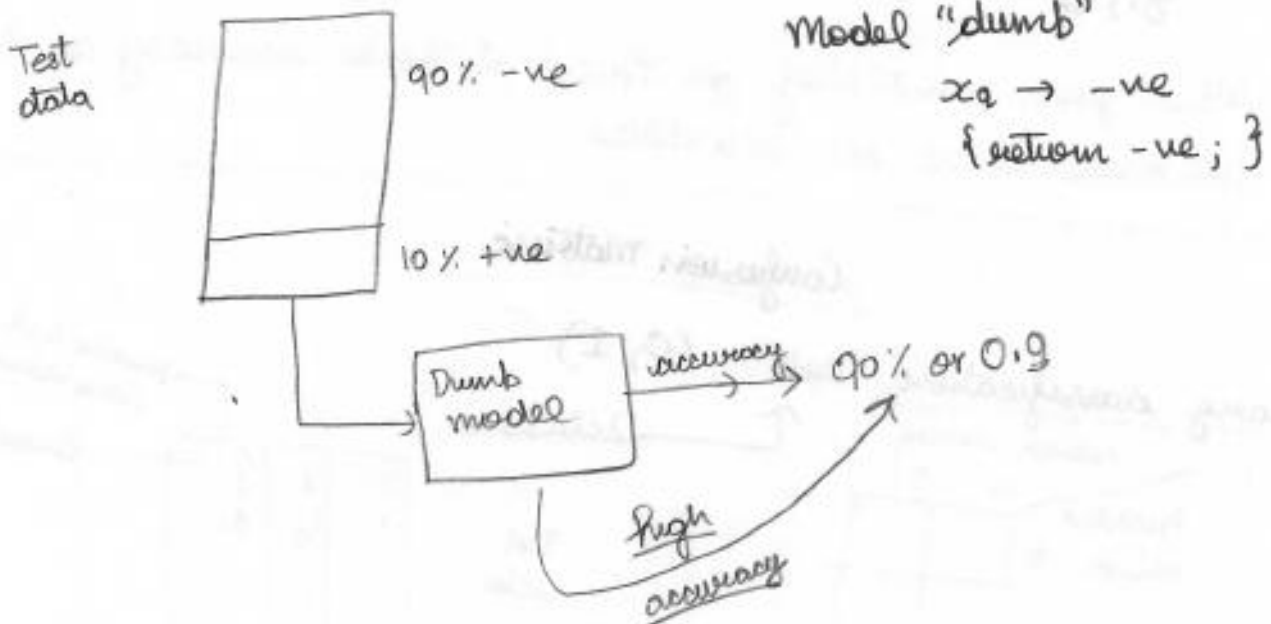
[0 to 1]
↙ ↘
Bad Best

100 pts → 60 +ve 40 -ve
(M) → 53 +ve, 7 -ve
→ 35 -ve, 5 +ve
Correctly = 88, Total errors = 12

$$\text{accuracy} = \frac{88}{100} = .88$$

⇒ Drawbacks of accuracy / Limitations of accuracy

Case 1: Imbalanced data :-



When you have imbalanced dataset never use accuracy.

Case 2:

predicted value of model 1 and model 2

	x	y	M_1	M_2	\hat{y}_1	\hat{y}_2
+ve	x_1	1	0.9	0.6	1	1
	x_2	1	0.8	0.65	1	1
-ve	x_3	0	0.1	0.45	0	0
	x_4	0	0.15	0.48	0	0
	x_5					

(Test-set)

Probability values

$M_1 \& M_2$

return a prob. score

KNN

$x_2 \rightarrow \text{prob}(\hat{y}_2 = 1)$

Probability score

Not/Can:

\hat{y} : predicted value
(\bar{x}, \bar{y})

predicted class label same for both models $M_1 \& M_2$.

But looking at probability score M_1 is better than M_2

But accuracy cannot use probability score.

$\hat{y}_1, \hat{y}_2 \rightarrow M_1 \& M_2$ have same accuracy.

When given probability you should not choose accuracy as it gives both models same importance

Confusion Matrix

Binary classification

Predicted values	Actual values	
	0	1
0		
1		

task (0, 1)

2 classes

Test data

x_1	y_1	\hat{y}_1
x_2	y_2	\hat{y}_2
\vdots	\vdots	\vdots
x_n	y_n	\hat{y}_n

predicted class label

Binary

actual class label

data

NOTE Confusion matrix cannot take probability scores.

		actual	
		0	1
Predicted	0	a	b
	1	c	d

a: # pts s.t $y_i = 0$ (actual class label)
 $\hat{y}_i = 0$ (predicted class label)

b: # pts s.t $\hat{y}_i = 0$ (predicted)
 $y_i = 1$ (actual)

c: # pts s.t $y_i = 1$ (actual)
 $\hat{y}_i = 1$ (predicted)

d: # pts s.t $y_i = 0$ (actual)
 $\hat{y}_i = 1$ (predicted)

Multi Classification (C-Class)

		Actual →			
		0	1	...	C-1
pred ↓	0				
	1				
	...				
	C-1				

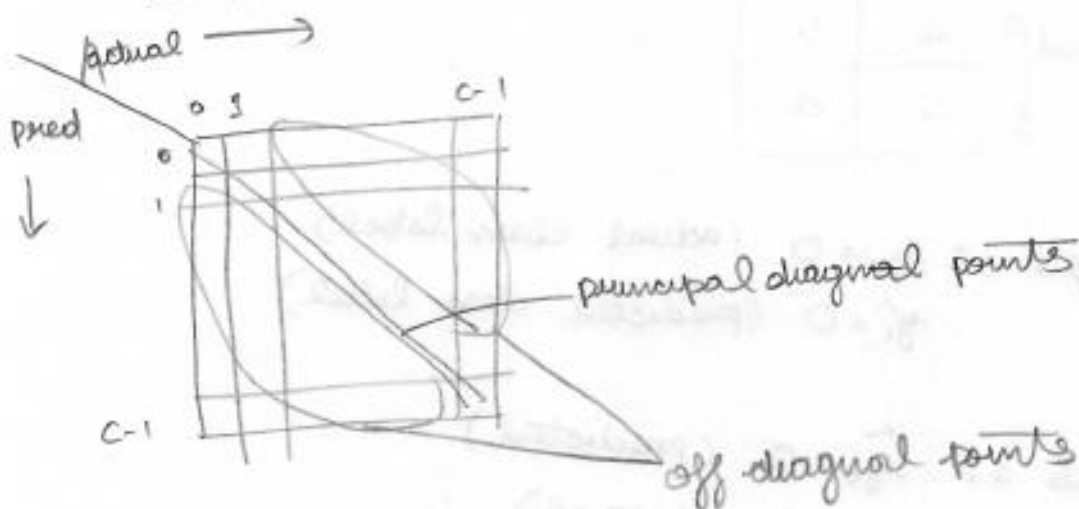
CXC

Binary setting

		Actual	
		0	1
predicted	0		
	1		

Model is sensible

		Actual	
		0	1
Pred	0	↑	small
	1	small	↑



Aim is that principal diagonal points must be larger than off diagonal points.

		actual	
		0	1
predicted	0	TN	FN
	1	FP	TP

(N) Total number of -ve N

(P) Total number of +ve

		Actual	
		0	1
Pred	0	-N	-ve
	1	+P	+P

True Positive

what is the predicted label?

Are you correct?

True Positive Rate

		Actual	
		0	1
Pred	0	TN	FN
	1	FP	TP

N

P

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

or Sensitivity

$$N + P = n$$

True Negative Rate = $\frac{\text{True Negative}}{\text{Negative}}$
or
Specificity

	Actual	
	0	1
Predicted	0	TN FN
	1	FP TP

False Positive Rate = $\frac{\text{False Positive}}{\text{Negatives}}$

False negative Rate = $\frac{\text{False Negative}}{\text{Positives}}$

Example

actual →

	0	1
predicted ↓	0 (TN) 850	(FN) 6
	(FP) 50	(TP) 94

N = 900 P = 100

Test :- 900 -ve
100 +ve } 9mbalanced data

Model :

TPR = 94%

TNR = $\frac{850}{900}$

FPR = $\frac{50}{900}$

FNR = $\frac{6}{100} = 0.06 = 6\%$

Model said to be good if TPR and TNR were high.

Imagine Dumb model :- all -ve

	act	
	0	1
pre	0	900 100
	1	0 0

N=900 P=100

900 -ve
100 +ve

↑ TPR = 0, ↑ TNR = $\frac{900}{900} = 1$
↓ FPR = 0, ↓ FNR = $\frac{100}{100} = 1$

Continued

↑ $\text{TPR} = 0\%$

↑ $\text{TNR} = 100\%$

↓ $\text{FPR} = 0\%$

↓ $\text{FNR} = 100\%$

From this we can predict that the model is dumb.

Which number is more important among 4?

This is domain specific.

eg.

diagnose cancer (not)

actual →
pred ↓

	0	1
0	TN	FN
1	FP	TP
	N	P

→ very low False negative rate (close to zero)
(missing someone who has cancer to nothing)

→ For cancer → we want high True positive rate

Precision and recall, F1-score

actual →

	0	1
pred 0	TN	FN
1	FP	TP
	N	P

P_r, R_e - Information retrieval

$$\text{Precision} = \frac{\text{True positive}}{\text{True positive} + \text{False positive}}$$

Information retrieval → million of documents ~~use~~ objective to find some document which are relevant.

$$\text{Recall} = \frac{\text{True positive}}{\# \text{ positive}}$$

NOTE: Precision and Recall values when you care about the true class.

Precision and Recall into one measure

Precision \uparrow
(0-1)

Recall \uparrow
(0-1)

$$F1\text{-Score} = \left(2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \right)$$

NOTE:

Harmonic mean \rightarrow

$$F1 = \frac{2}{\left(\frac{1}{\text{recall}} + \frac{1}{\text{pr}} \right)} \times \left[\text{avg} \left[\text{inv-recall}, \text{inv-precision} \right] \right]^{-1}$$

NOTE:

Harmonic mean:

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

Used in kaggle competitions

Lies between 0 to 1

Medical domain, we want high recall. We don't want any patient who actually has a disease to be misdiagnosed and the model say that he/she doesn't have it.

But, we are allowed to have a low precision. Meaning, maybe, the patient actually doesn't have a disease but we are okay on saying that he/she has. That is not life threatening.

$$\text{Error rate} = \frac{FP + FN}{P + N}$$

Why F-measures uses Harmonic mean?

It measures Test accuracy

Harmonic mean is applicable, if the numerators of the averaged values are the same. eg resistance in parallel

Arithmetic mean is applicable, if the denominators of the averaged values are the same. eg resistor in series

Harmonic mean \rightarrow ?

precision \times recall = $0.1 \times 0.1 = 0.01$

Case I
$$F1\text{ score} = \frac{2}{\frac{1}{0.1} + \frac{1}{0.1}} = \frac{2}{20} = 0.1$$

accurate value

Case II $P=0.8$
 $r=0.1$

$P \times r = 0.08$

$$F1\text{ score} = \frac{2}{\frac{1}{0.8} + \frac{1}{0.1}} = 0.177$$

F1-Score gather information more accurately than PXR

Evaluating Classifier Performance Receiver Operating Characteristics Curve (ROC)

\rightarrow electronics & radio engineers discovered by

x	y	\hat{y} (probability score)
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

decreasing of \hat{y}

binary
classification
Model +1
0
 \hookrightarrow score
prob score

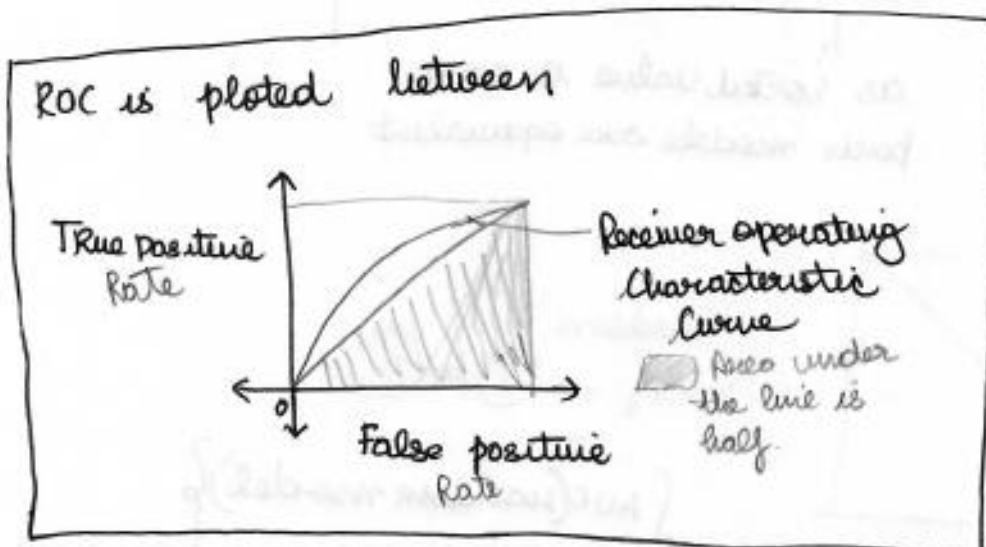
\uparrow score
 $\Rightarrow 1$

ROC curve

Step-1 Sort in decreasing order of \hat{y}

Step-2 Thresholds (Z)

$$\begin{array}{l} \text{a.) } Z_1 = 0.95 \\ \text{if } \hat{y} \geq Z_1 \\ \quad 1 \\ \text{else} \\ \quad 0 \end{array} \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \begin{array}{l} \text{TPR} \\ \text{FPR} \end{array}$$



Thresholds are used for plotting different ROC curves.
ROC plots is plotted using different thresholds.

By setting different thresholds you get new True positive Rate and False positive Rate that is called ROC curve.

Area under the ROC curve can be 0 to 1

\uparrow \uparrow
 terrible V. good

★ AUC brilliant properties

① Imbalanced data AUC can be high for dumb model

② AUC did not care about actual values. It depends on the ordering. (Sorting)

		M_1	M_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

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

↓
As sorted value is same.
both models are equivalent

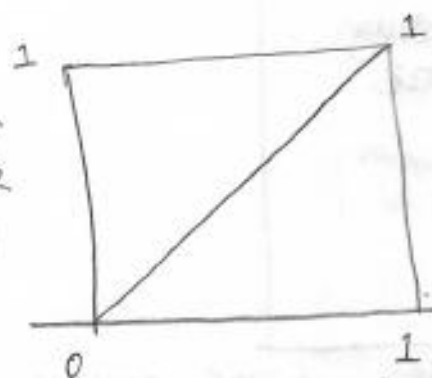
5) Classifier A is better than Classifier B if A's AUC is larger than B's AUC

6) A perfect classifier has AUC of 1

7) A classifier making random guesses has AUC of 0.5

③ ↑
TPR

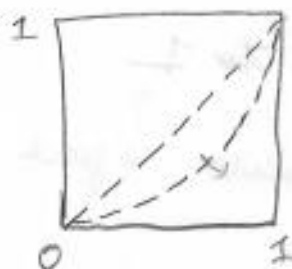
ROC



Random
→ $x_a \rightarrow \textcircled{1} \text{ or } \textcircled{0}$

$$\left\{ \begin{array}{l} AUC(\text{random model}) \\ = 0.5 \end{array} \right\}$$

④ Model M:



$$AUC(M) : = 0.2$$

↑
worse than random

AUC value 0.5 to 1 → normal

0.5 → random

0.0 to 0.5 → guess class labels

Change the
class label to
opposite to it

$$\left\{ \begin{array}{l} \hat{y}_i = 0 \rightarrow 1 \\ \hat{y}_i = 1 \rightarrow 0 \end{array} \right.$$

⚡ Log Loss : probability Scores

Binary classification :

x	y	$\hat{y} = p$	
x_1	1	0.9	$-\log(0.9) \rightarrow 0.0457$
x_2	1	0.6	$-\log(0.6) \rightarrow 0.22$
x_3	0	0.1	$-\log(0.9) \rightarrow 0.0457$
x_4	0	0.4	$-\log(0.6) \rightarrow 0.22$

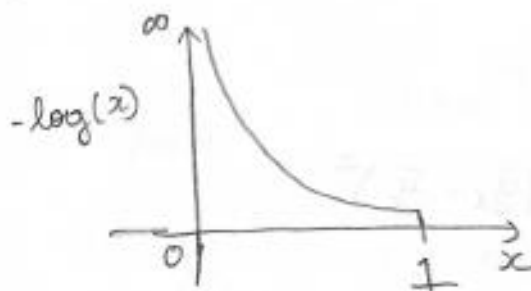
Test-set of n points :

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

If model is good p_i values should be equal to 1 that means model is good.

log loss should be as small as possible.

log loss \div avg. neg log (probability of correct class label)



$$0 \leq p \leq 1$$

smaller the better
all loss function min/objective
is to be small

Multi Class log-loss :-

$$-\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^C y_{ij} \log(P_{ij})$$

↓

= 1 if $x_i \in \text{class } j$
0 o/w.

↳ prob that $x_i \in \text{class } j$

best case = 0

$$\text{log loss}(M) = 1.0$$

R² Squared / Coefficient of determination

$$y_i \in \mathbb{R}$$

Test :- x_i, y_i, \hat{y}_i
 $i \rightarrow 1 \text{ to } n.$ ↑
 model
 prediction/output.

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

↑
error for pt i

◇ Sum of Squares (SS)

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

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

{ (mean / average values of y_i 's)
in test data }

example predict height (EIR)

features - w, c, h, \dots

mean can be average value. that can be used to build simplest model.

$SS_{total} = \sum_{i=1}^n (y_i - \bar{y})^2$	\bar{y}_{mean} squared errors using simple mean model.
$SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$	\hat{y} actual value.

Residue = $e_i = y_i - \hat{y}_i$

\uparrow error \uparrow actual \nwarrow predicted

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

Case 1: $SS_{res} = 0 \leftarrow e_i = 0 \rightarrow R^2 = 1$ (best-value)

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

Case 3: $SS_{res} = SS_{total}$; $R^2 = 1 - 1 = 0 \rightarrow$ Model is same as simple mean model.

Case 4: $SS_{res} > SS_{total}$

$$R^2 = 1 - (gt > 1) = -ve$$

Model is worse than a simple mean model

Median absolute deviation (MAD)

$$SS_{\text{Residual}} = \sum_{i=1}^n e_i^2$$

one e_3 is very large
whole R^2 will be effected.

R^2 is not very robust to outliers.

MAD \rightarrow (Median absolute deviation)

e_i : random variable.

Median (e_i) = central - value of errors. = small

MAD(e_i) = Median ($|e_i - \text{Median}(e_i)|$) = small

no outlier

Central tendency mean or median of e_i s

Std
dev

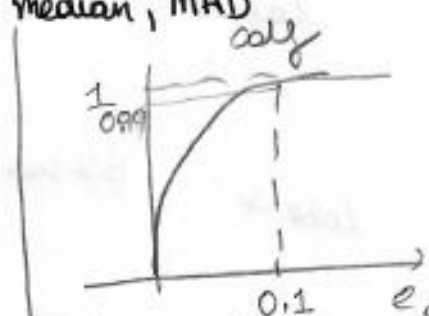
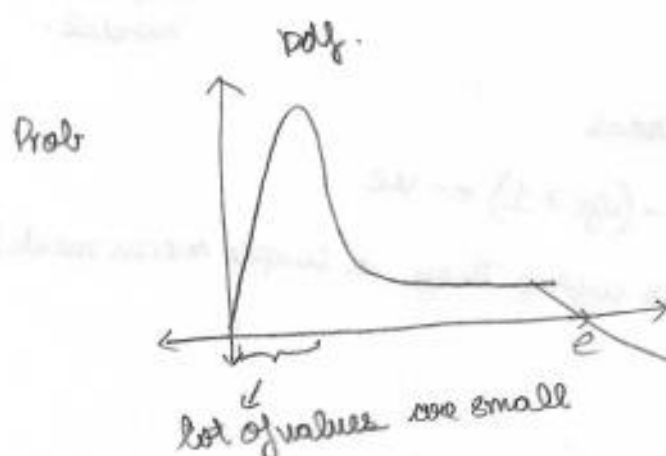
MAD

↓

Robust to outliers

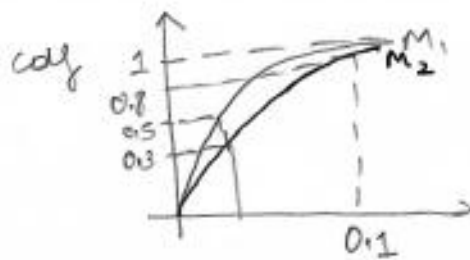
Distribution of errors:

$x_i \rightarrow y_i, \hat{y}_i, e_i \rightarrow$ mean, std - dev
median, MAD



99 % of errors < 0.1
1% error ≥ 0.1

Ideally what we want $e_i \rightarrow 0$

Model M_1 & M_2 

M_2 : cdf is below model 1

M_1 : # errors 95% errors are below 0.1

M_2 : 90% errors below 0.1

M_1 : 50% errors are below 0.01

M_2 : 30% errors are below 0.01

M_1 is better than M_2 .