

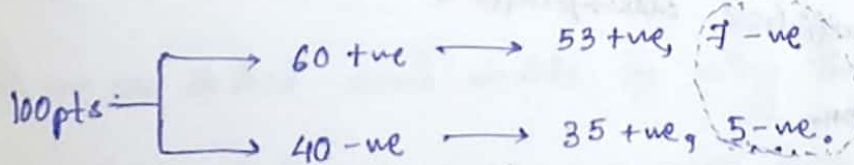
Performance Measurement of Models

some classification metrics

$$\text{Accuracy} = \frac{\text{\#pts correctly classified}}{\text{Total \#pts in } D_{\text{test}}}$$

easy to understand measure.

measured on test set usually.



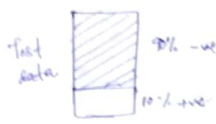
$$\text{error} = 12, \quad \text{correctly} = 88.$$

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

Problems with accuracy

- ① Doesn't work well with unbalanced data.

Accuracy $\in [0, 1]$
↓ (bad) ↓ (good)



dumb any $x_i \rightarrow -ve$
 Model
 Accuracy = 0.9
 So, never use accuracy as performance measure for imbalanced data

x	y	Model 1	Model 2	y_1	y_2
x_1	1	0.9	0.6	1	1
x_2	1	0.8	0.65	1	1
x_3	0	0.1	0.15	0	0
x_4	0	0.16	0.15	0	0

* We can clearly say that M_1 is better than M_2 by looking at the probability values.
 * But since predicted class labels are same, we get same accuracy.

* Accuracy \rightarrow can't use probability scores.

Confusion Matrix (doesn't process probability scores).

consider a binary classification task (0, 1).
 2 classes

Actual \rightarrow

	0	1
Pred 0	a	b
Pred 1	c	d

x_i	y_i	\hat{y}_i
x_1	y_1	\hat{y}_1
x_2	y_2	\hat{y}_2
x_n	y_n	\hat{y}_n

predicted class label

as #pts x_i s.t. $y_i = 0$ & $\hat{y}_i = 0$.
 as #pts x_i s.t. $y_i = 1$ & $\hat{y}_i = 1$.

This can be also extended to multiclass classifications

consider c-class classification

Actual \rightarrow

	0	1	...	c-1
Pred 0				
Pred 1				
Pred c-1				

for a sensible model,
 principal diagonal elements \uparrow .
 off diagonal elements - small.

act \rightarrow

	0	1
pred 0	TN	FN
pred 1	FP	TP

(N) (P)

$\frac{T}{T+P}$ what is the predicted label.
 is the prediction correct.

N : # negatives.
 P : # positives.
 $n = N + P$.

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

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

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

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

act \rightarrow

	0	1
pred 0	850	6
pred 1	50	91

800N 100P

$TPR = 91\%$
 $TNR = \frac{850}{900}$
 $FAR = \frac{50}{900}$
 $FNR = 6\%$

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

Test Set: 900 -ve, 100 +ve } imbalanced

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

X $TPR = 0\%$, $FPR = 0\%$
 $TNR = 100\%$, $FNR = 100\%$

let us now assume a dumb model which gives off label always -ve.

act \rightarrow

	0	1
pred 0	900	100
pred 1	0	0

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

* Next, which of the 4 metrics is more important?

\rightarrow domain specific.

ex \rightarrow diagnosing cancer \rightarrow We must have very low FNR as we can't afford to say an actual patient that he doesn't have cancer.
 * little high FPR is acceptable as even if we say a non-patient that he may have cancer, later through powerful tests that can be identified.

③ Precision, Recall and F1-Score (Information Retrieval).

Pred	0	1
0	TN	FN
1	FP	TP
	N	P

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

of all pts the model declared to be 1's, what %age of them are actually 1's.

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

of all the ~~models~~ pts which actually belong to class 1, how many are predicted to be 1's.

We want both Pr & Re to be high.
[0, 1].

F1 score combines them both

$$\text{F1-Score} = 2 * \left(\frac{\text{Pr} * \text{Re}}{\text{Pr} + \text{Re}} \right) \quad \left\{ \begin{array}{l} \text{Harmonic mean of Pr \& Re.} \\ \frac{1}{F1} = \frac{1}{Pr} + \frac{1}{Re} \end{array} \right.$$

used as metric in

many Kaggle competitions. Pr & Re — more interpretable — can understand in simple english

④ Receiver Operating Characteristic Curve (ROC) & AUC (Area Under Curve).

— for binary classification only.

Assume model outputs a score which can be interpreted similar to probability of class '1' score.

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	0.71

Steps for ROC

① Sort entries in decreasing order of \hat{y}

② Thresholding (z).

Serially from top to bottom, take each value of \hat{y} as

Then, if $\hat{y} \geq z$,

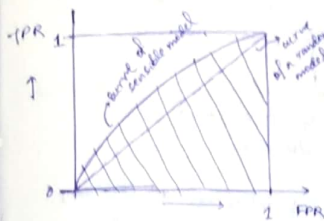
label = 1

else

label = 0.

Find TPR & FPR for each such threshold.

x	y	\hat{y}	$\hat{y} \geq z$	$\hat{y} < z$
x_1	1	0.95	1	
x_2	1	0.92	1	
x_3	0	0.80	0	
x_4	1	0.76	0	
x_5	1	0.71	0	
			TPR ₁	FPR ₁
			TPR ₂	FPR ₂



Area under curve (AUC) — 0 to 1.
(terrible) (v. good)

for a sensible model, this will look like the curve will look like.

AUC properties

① imbalanced data → AUC can be high for dumb-model / simple model.

② AUC is not dependent on the \hat{y} scores. Only depends on ordering.

Ex →

x_1	1	0.95	0.2
x_2	1	0.92	0.1
x_3	0	0.88	0.08
x_4	1	0.76	0.07
x_5	1	0.71	0.06

Both have completely different values. But the sorted order is same.

③ AUC (random model) = 0.5

If $\text{AUC}(M) \in [0, 0.5]$, → simply complement the op of the model (swapping).

⑤ Log-loss (uses probability scores).

consider binary classification case.

x y $\hat{y}=p$ loss

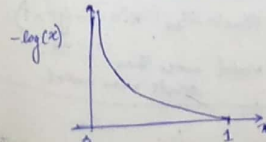
x_1	1	0.9 → 0.0457
x_2	1	0.6 → 0.22
x_3	0	0.1 → 0.0457
x_4	0	0.4 → 0.22

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

$$\text{Log-loss} = -\text{avg. avg. log (prob. correct class label)}.$$

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 is the loss.



for Multiclass classification,

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

prob. of x_i in class j .
1 if x_i in class j
0 otherwise

* One key disadvantage of log-loss is its hard to interpret, as it can go as high as ∞ , No upper-bound. Can't determine how bad.

Some regression metrics

$y \in \mathbb{R}$ in regression

Test Set x_i, y_i, \hat{y}_i model output
 $e_i = (y_i - \hat{y}_i)$

* Just like in classification, a simple model would be one that returns the label of the majority class for all y 's.
* Similarly, for regression, a simple model would be one that returns the mean of y 's.

Error for this simple model = $\sum_{i=1}^n (y_i - \bar{y})^2$ — also called SS_{tot} or sum of squares (using simple mean model).

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

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

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

aka search for coefficient of determination

Case 1: $SS_{\text{res}} = 0 \rightarrow e_i = 0 \rightarrow (R^2 = 1)$ best value.
Case 2: $SS_{\text{res}} < SS_{\text{tot}} \rightarrow R^2 \in [0, 1]$
Case 3: $SS_{\text{res}} = SS_{\text{tot}} \rightarrow R^2 = 0$ model same as simple mean.
Case 4: $SS_{\text{res}} > SS_{\text{tot}} \rightarrow R^2 = 1 - (y_i > 1)$ Model worse than simple mean model.

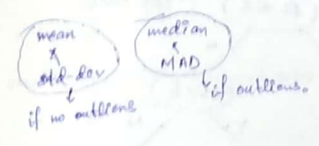
Median Absolute deviation of errors

R^2 - relies on mean, hence not very robust to outliers.

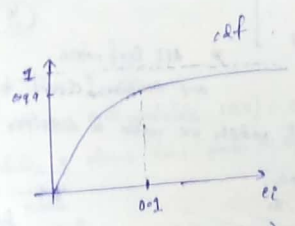
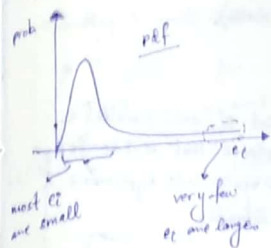
for each $x_i \xrightarrow{\text{model}} y_i, \hat{y}_i, e_i$

e_i is random variable.

median(e_i) = central value of errors, — defined for a set of pts.
 $MAD(e_i) = \text{median}(|e_i - \text{median}(e_i)|)$ — defined for a set

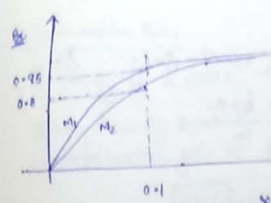


Distribution of errors



99% of e_i (errors) are < 0.1
1% of errors > 0.1

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



M_2 cdf below M_1 .

M_1 : 95% of errors are below 0.1

M_2 : 80% of errors are below 0.1

So, M_1 better model than M_2 .

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