NOTES

Classification accuracy measures:
Performance measure of models: **Accuracy**
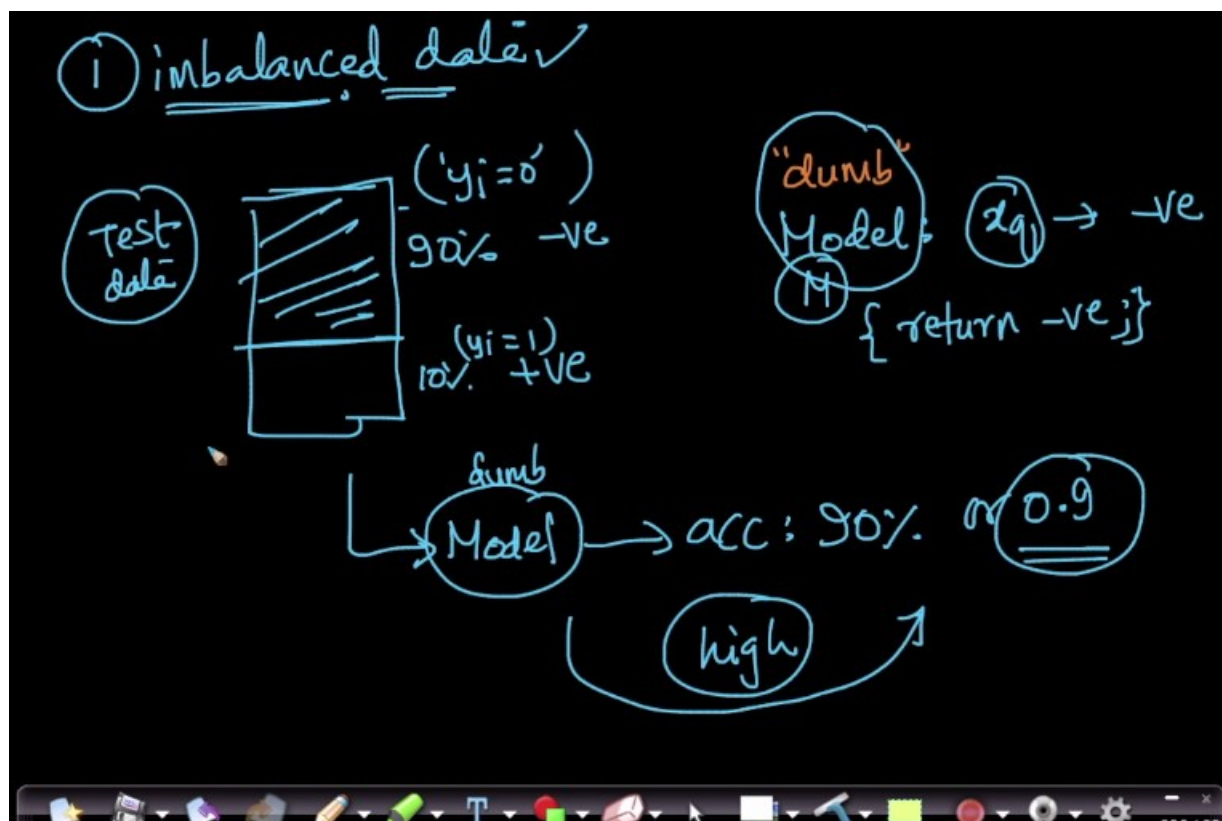 accuracy = Number of correctly classified points / Total number of points.
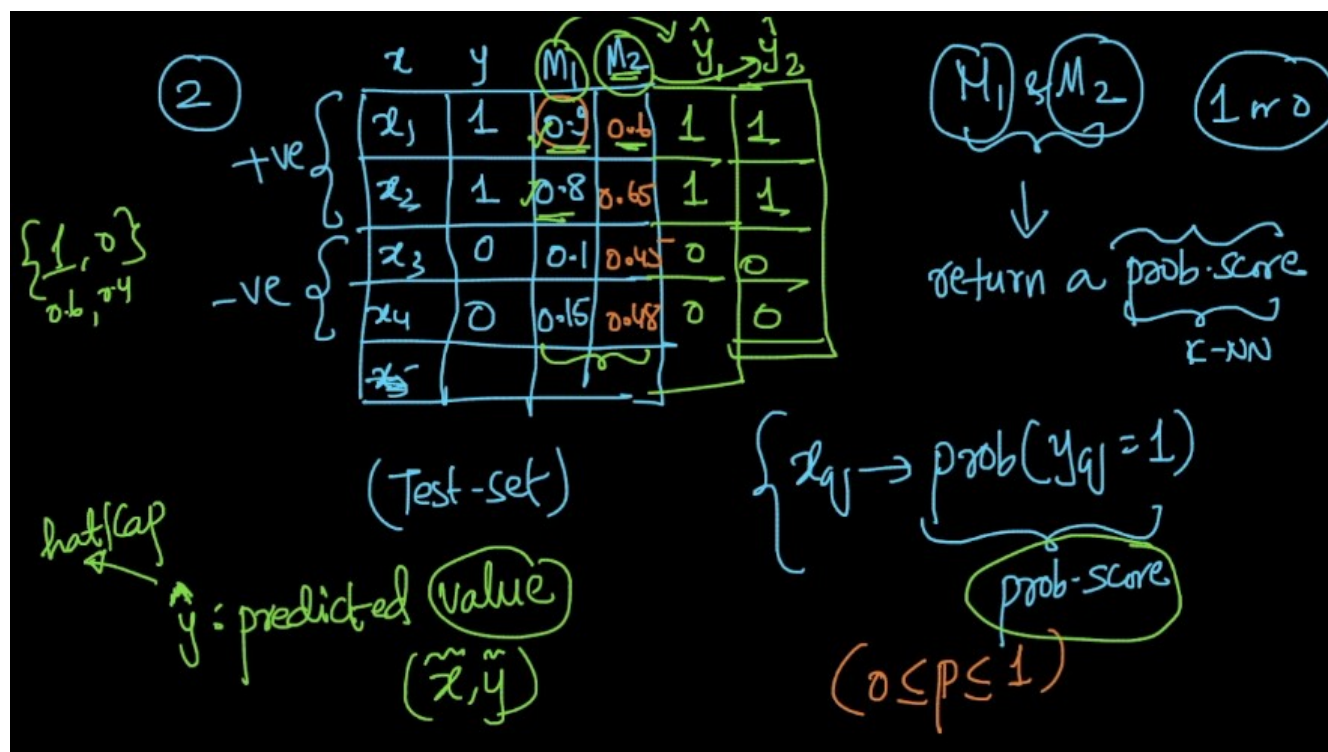


Problems of Accuracy as the measure:

Accuracy is only done on the test data. Accuracy fails in case of imbalanced data.
**In case of imbalanced data-set we should never use accuracy as the measure, because a DUMB model can give high accuracy.**

① imbalanced data ✓

Test data

$(y_i = 0)$
$90\%$ -ve

$(y_i = 1)$
$10\%$ +ve

"dumb" Model; $x_q \to$ -ve

M { return -ve }

dumb Model → acc : 90% or $\underline{0.9}$

high ↗

In case of the models output is a probability score, we assume a threshold to be and make the labels above that to be one class, below that is other class.(threshold ~ 0.5).



②

$\{1, 0\}$
$0.6, 0.4$

+ve $\begin{cases} \end{cases}$

-ve $\begin{cases} \end{cases}$

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

(Test-set)

hat/cap
$\hat{y}$ : predicted value
$(\tilde{x}, \tilde{y})$

$M_1$ & $M_2$    (1 ir 0)

↓

return a prob·score
K-NN

$\{ x_q \to prob(y_q = 1)$
prob-score
$(0 \le p \le 1)$

The predicted class labels are exactly the same for the two models.(**m**1 and **m**2), though the models are different.



These are two major drawbacks for accuracy as the measure of the model.
Confusion – matrix:
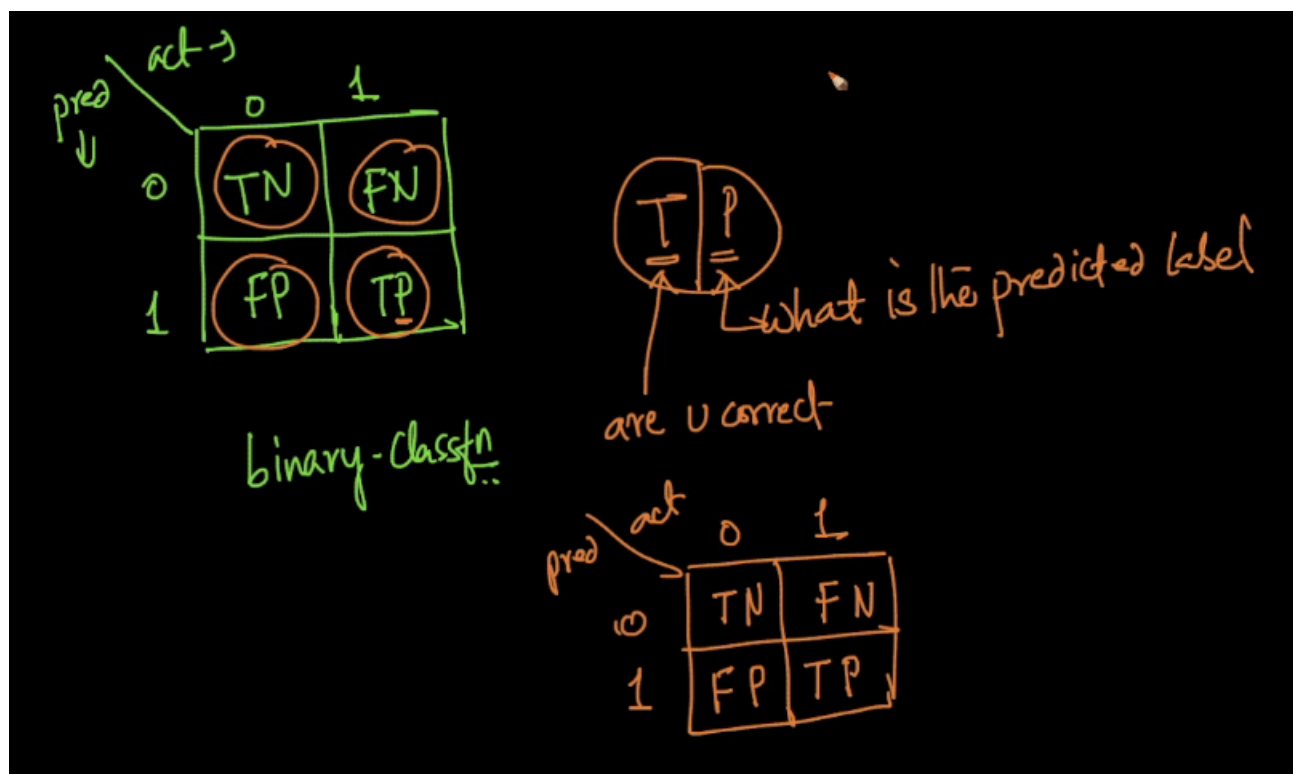It cannot take probability scores. They only take binary values.

In case of multiple - class classification:

We can draw a matrix for the predicted and actual values.

If the model is sensible, the number of values along the principal diagonal must be more than the off diagonal elements.



Important to remember the confusion matrix.

Various confusion matrix measurements:

$$TPR = \frac{TP}{P}$$ ✓

$$TNR = \frac{TN}{N}$$

$$FPR = \frac{FP}{N}$$

$$FNR = \frac{FN}{P}$$

pred ↓  act →

|        | 0  | 1   |
|--------|----|-----|
| 0      | TN | FN  |
| 1      | FP | TP  |

(N) (P)

$$N + P = n$$

---

pred ↓  act →

|     | 0       | 1      |
|-----|---------|--------|
| 0   | 850 (TN)| 6 (FN) |
| 1   | 50 (FP) | 84 (TP)|

900 = N      P = 100

Test :— 900 −ve ⎫ im balanced
        100 +ve ⎬

Model

↑ TPR = 94%        FPR = $\frac{50}{800}$ ↓

↑ TNR = $\frac{850}{900}$        FNR = 6% ↓

Cases of confusion matrix:



In case of imbalanced data-sets, confusion matrix helps in making good inference from the model.
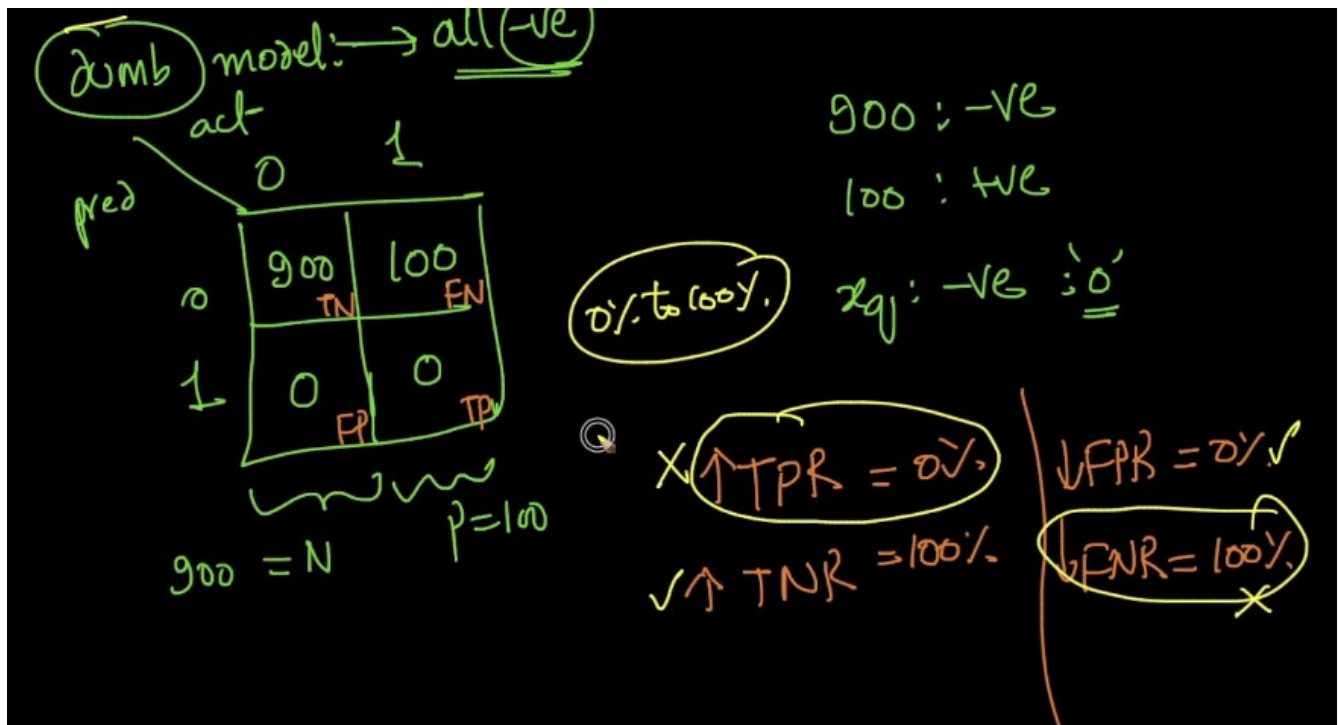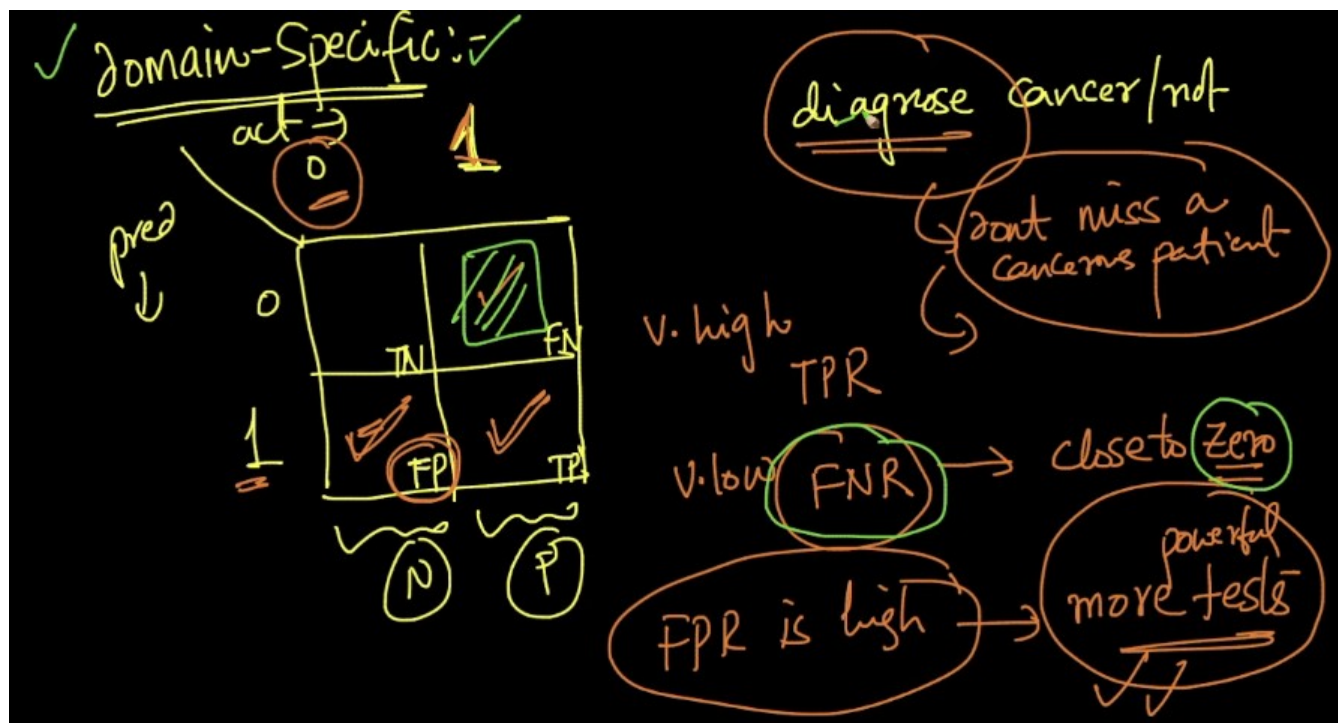The importance of the various measures of the confusion matrix is more domain specific.

Precision, recall and F1 – score:



Precision: Precision only computes the positive class predicted rate.
Recall: Recall only computes on actual positive class values rate.
These are only used, If we want to measure positive class performance.
F1- score:
F1 – score is the harmonic mean of the precision and recall.

Precision and recall are more interpret able, than f1 – score.

Receiver Operation characteristic curve and accuracy curve:
We use thresholds as the base in AUC and ROC curve. The values must be sorted in decreasing order of class scores.



We can have 'n' thresholds for 'n' points in the data set.

For each threshold we can compute the true positive and false positive rate.



Next, step is to draw the plot of TPR vs. FPR. This curve is called ROC curve.
The blue line breaks the total plot is broken into two halves.

**AUC is the area of the curve under the ROC curve. ROC curve can only be used for the binary classification tasks.**

The values of AUC will lie between 0 and 1. The higher the better is the accuracy.

Properties of AUC:

1. In case of imbalanced data  the AUC can he high.

2. AUC does not care of the actual value of the accuracy, It cares only the ordering of the class scores.



AUC of several models can be the same.

If the model is random then the AUC will be same as the diagonal line, i.e, 0.5.

**If the model gives the AUC value between 0 and 0.5 then we just swap the class values to get the good model.**

(4) Model $M$ :            $AUC(M) = 0.2$

worse than random

$M$   $\begin{bmatrix} 0 \to 1 \\ 1 \to 0 \end{bmatrix}$   Swap

1

0.2

$1 - 0.2 = 0.8$

0

Swapping

$\begin{cases} y_i = 0 \to 1 \\ \tilde{y}_i = 1 \to 0 \end{cases}$

1

AUC
0.5 to 1 → ✓

$0.\bar{5}$ → random

0.0 to 0.5

Swap the class label
$AUC = 1 - 0.2$

---

Log – Loss:
This model is penalizing small deviations in a probability score. We want the log – loss to be as small as possible. Here we use actual probability score unlike other interpretations.
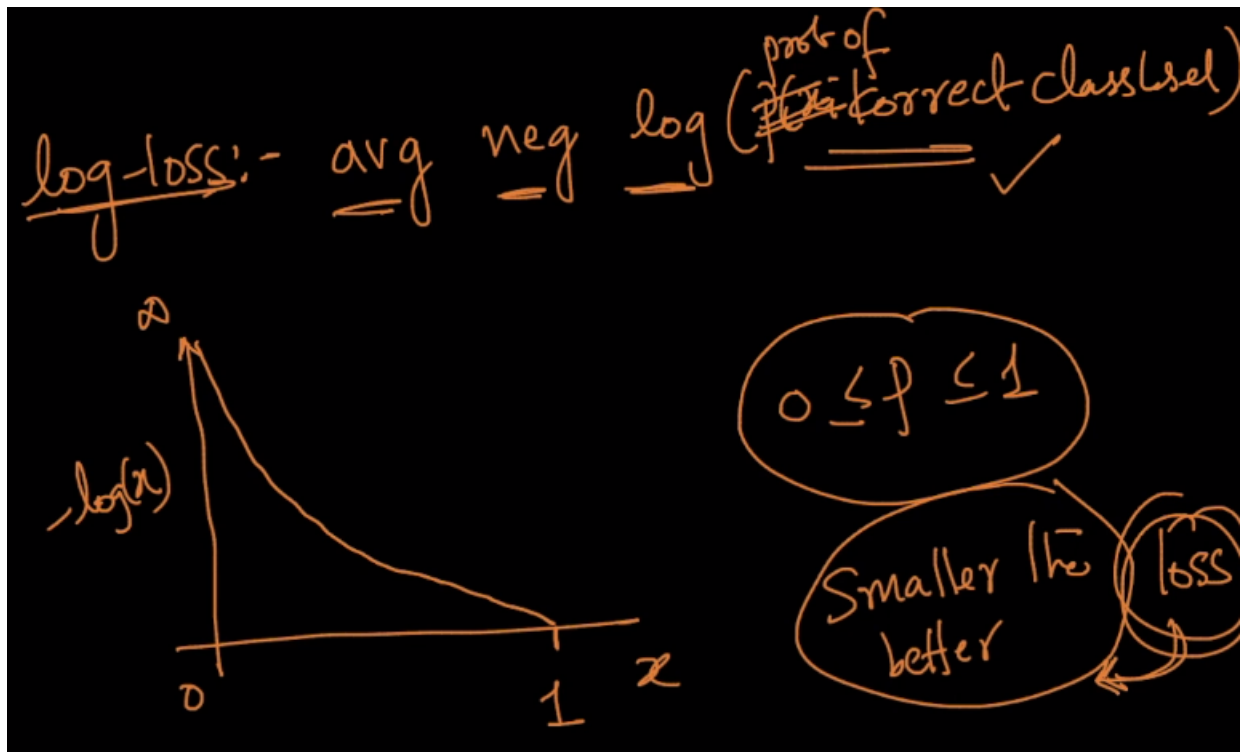
---



Log – loss :- prob-scores → model     $\hat{y}_i = p_i$

as small   0 to ∞

Binary-classfn :

Penalizing for small deviation in prob-score

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

Test-set of $n$-pts :-

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

avg

The smaller the better is the model.

$$\text{log-loss:-}\quad \text{avg}\quad \text{neg}\quad \log\left(\text{prob of correct class(label)}\right)\checkmark$$

$-\log(x)$

$0 \leq p \leq 1$

Smaller the loss better

Log loss function for multi − class classification:

$$\text{Multi-class}\quad \text{log loss:-}$$

$z_q \rightarrow \boxed{P_1\ P_2\ \cdots\ P_c}$

$$\left\{ -\frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{C} y_{ij}\ \log(p_{ij}) \right.$$
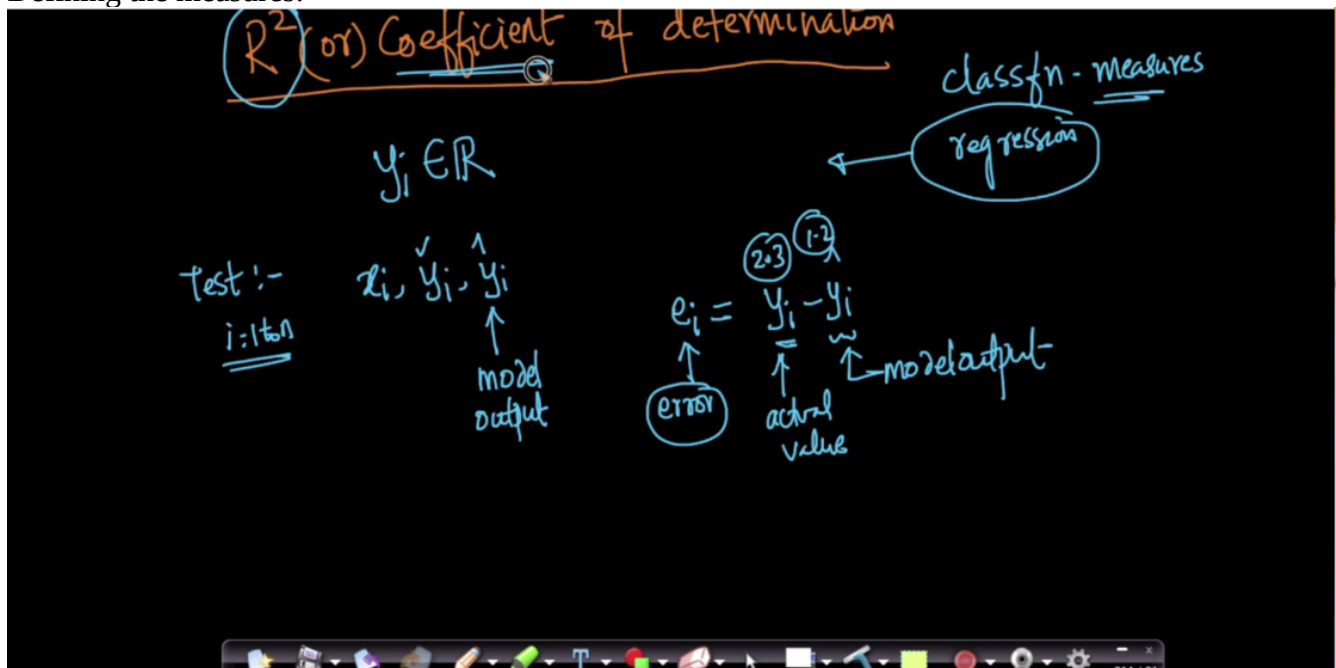
$\rightarrow$ prob that $x_i \in$ class $j$

$y_{ij} = 1$ if $x_i \in$ class $j$

$0$ o/w

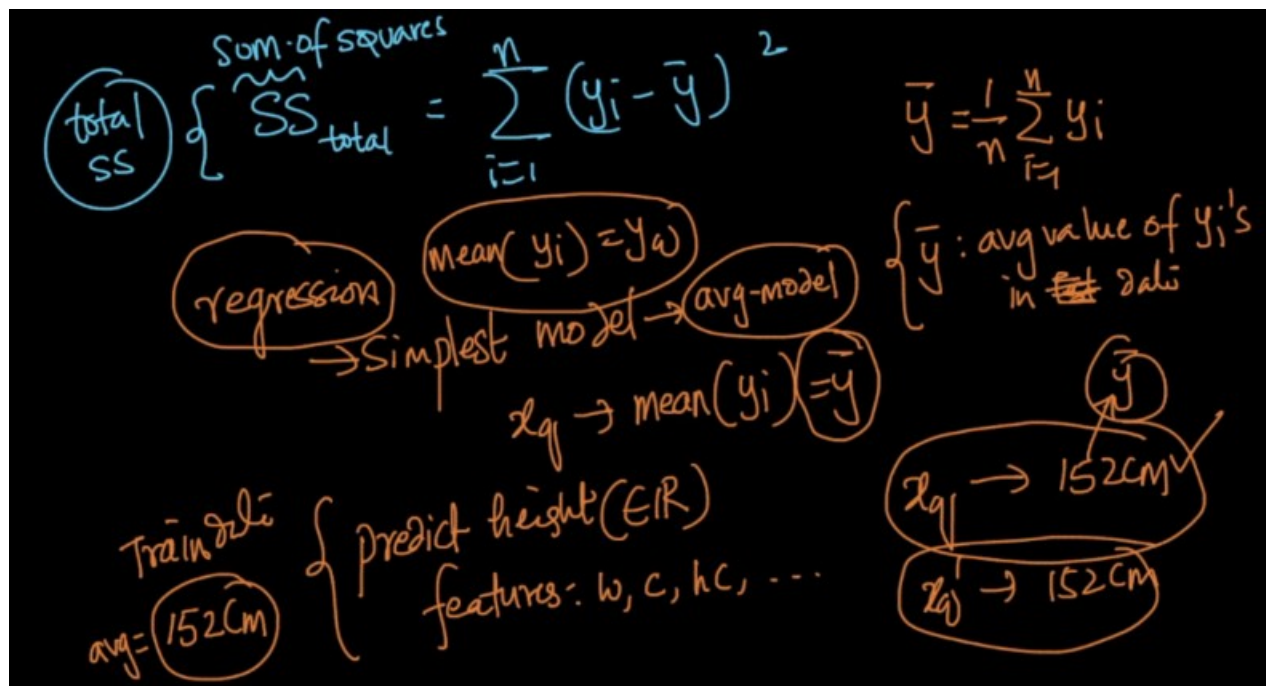The disadvantage of log − loss is we cannot interpret.

Log − loss is more useful for binary and as well as multi class classification tasks.

Accuracy measures for regression:
Defining the measures:



Steps in computing R^2:
Sum of squares:



The simplest model that can be constructed is reporting the mean of the data points for a query point xq. Here y(bar) is the average of all the data points in the data set.

Sum of squared of residuals:

This is called sum of square errors for each point that is being predicted by the model.

Defining R^2 term:



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

Case3 model: $SS_{res}$
└ same as a simple-mean-mod

Case 1: $SS_{res} = 0$ & $(e_i = 0)$ → $R^2 = 1$ (best-val)

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

Case 3:- $SS_{res} = SS_{tot}$ ; $R^2 = 1-1 = 0$ → Model same as S-M-Mode

When the model residue and the total sum of squared is same, then the model is same as simple squared model.

Another case:



Case4: $SS_{res} > SS_{tot}$

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

{ Model is worse than a simple- Mean-model

When one value is very large then R^2 can go wrong. The sum of residues can go for a toss in case of a large value as an outlier.
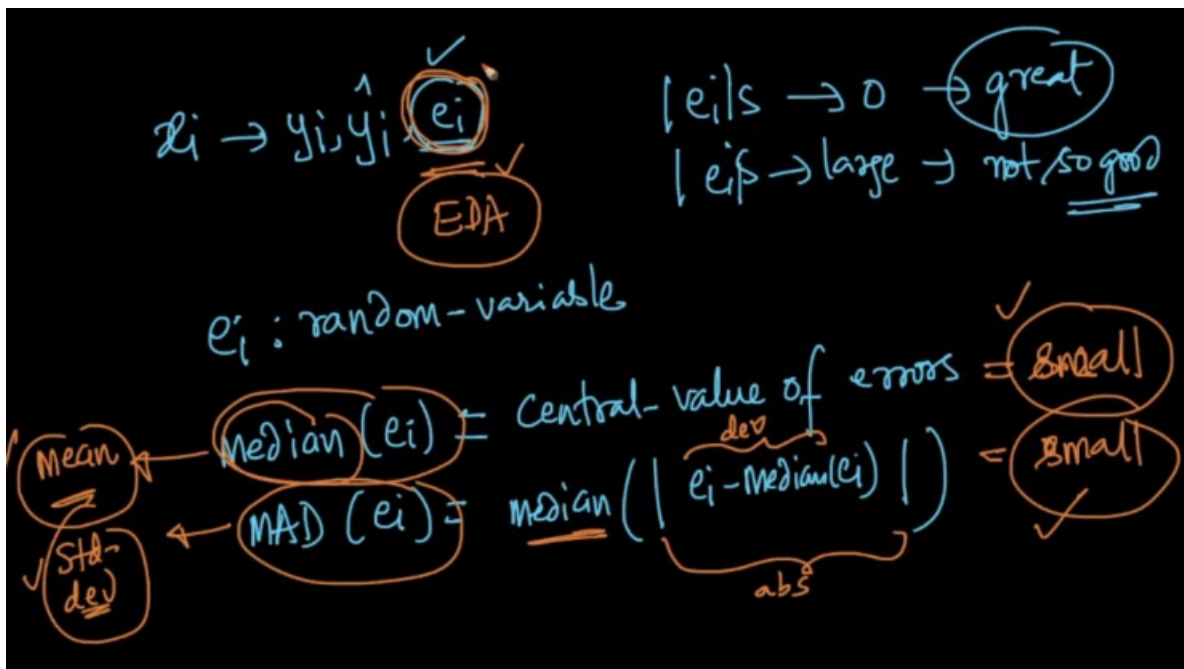
# Median Absolute deviation of errors

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

one $\boxed{e_3}$ is very large

$R^2$ is not very robust to outliers

$\boxed{MAD}$

That is why we use median absolute deviation as the measure, because it is prone from outliers.

$x_i \rightarrow y_i, \hat{y_i} \rightarrow \boxed{e_i}$

$\boxed{EDA}$

$|e_i|s \rightarrow 0 \rightarrow$ great

$|e_i|s \rightarrow$ large $\rightarrow$ not so good

$e_i$ : random-variable

$\boxed{mean}$ $\leftarrow$ $\boxed{median \,(e_i)} \rightarrow$ central-value of errors = small

$\boxed{Std-dev}$ $\leftarrow \boxed{MAD \,(e_i)} = median \left( \,| e_i - median(e_i) | \, \right)$ = small
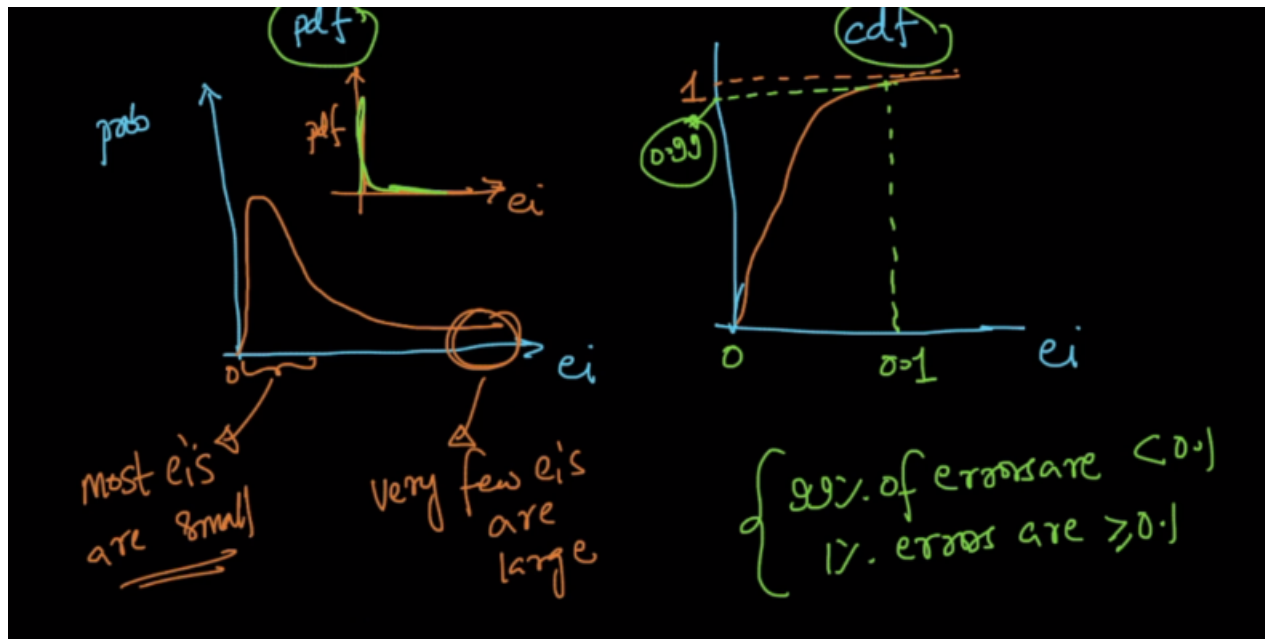
abs

We could use the
        Mean – standard deviation.
        Median – median absolute deviation. (robust).

For measuring the errors, we can use above methods to infer the distribution of errors(residuals).
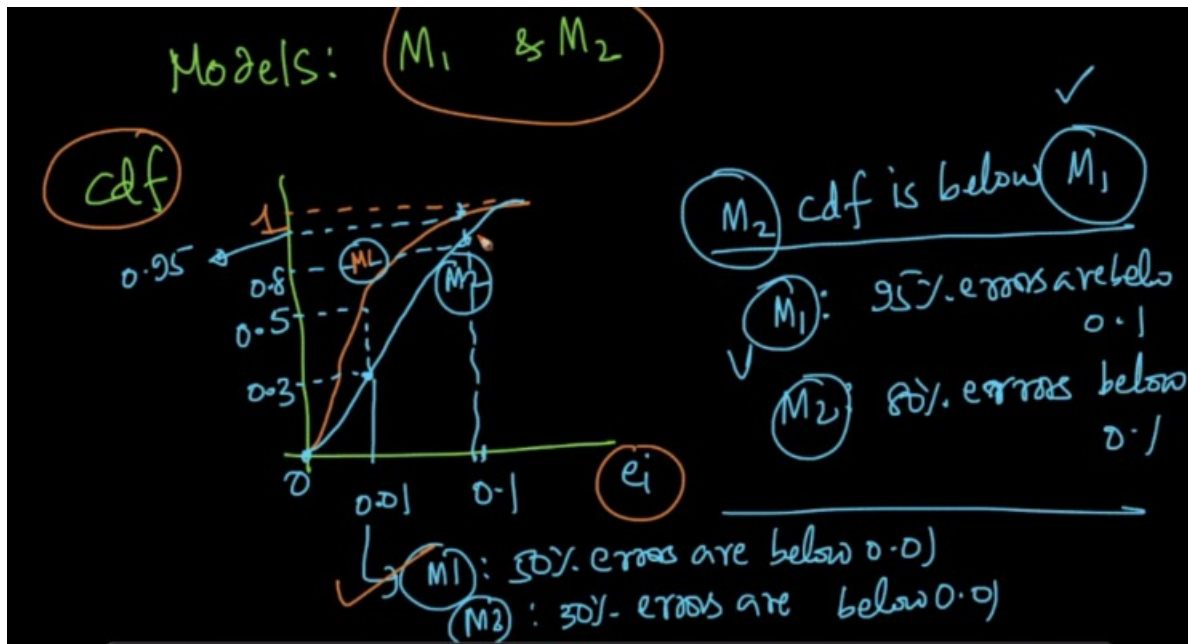
Distribution of errors:
We can use PDF and CDF for knowing the distributions of the errors.



If all the ei's are equal to zero then we performed best regressor.
Understanding the distribution of errors is important in case of regression.

Interpreting the errors of the two models:



By using CDF plots we can decide which model is better.