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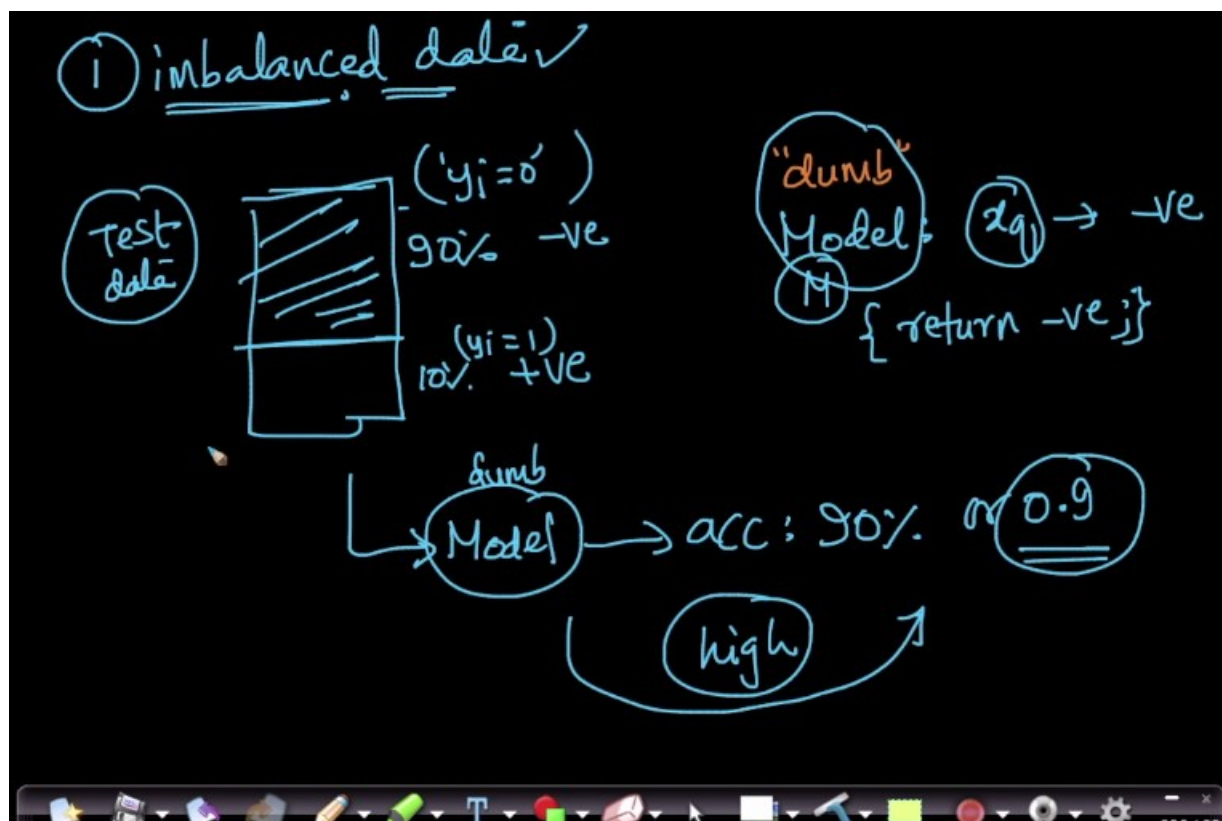


Problems of Accuracy as the measure:

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 ✓



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).

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:
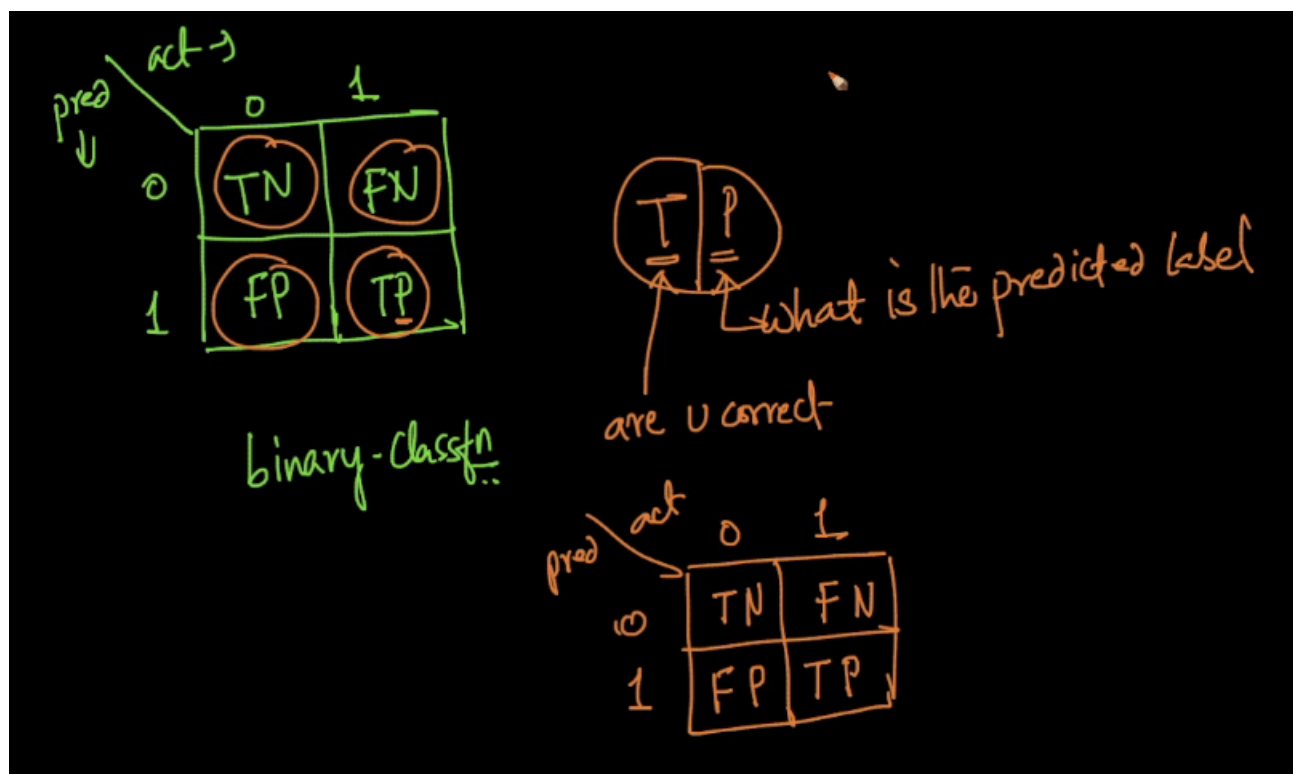It cannot take probability scores. They only take binary values.

In case of multiple - class classification:

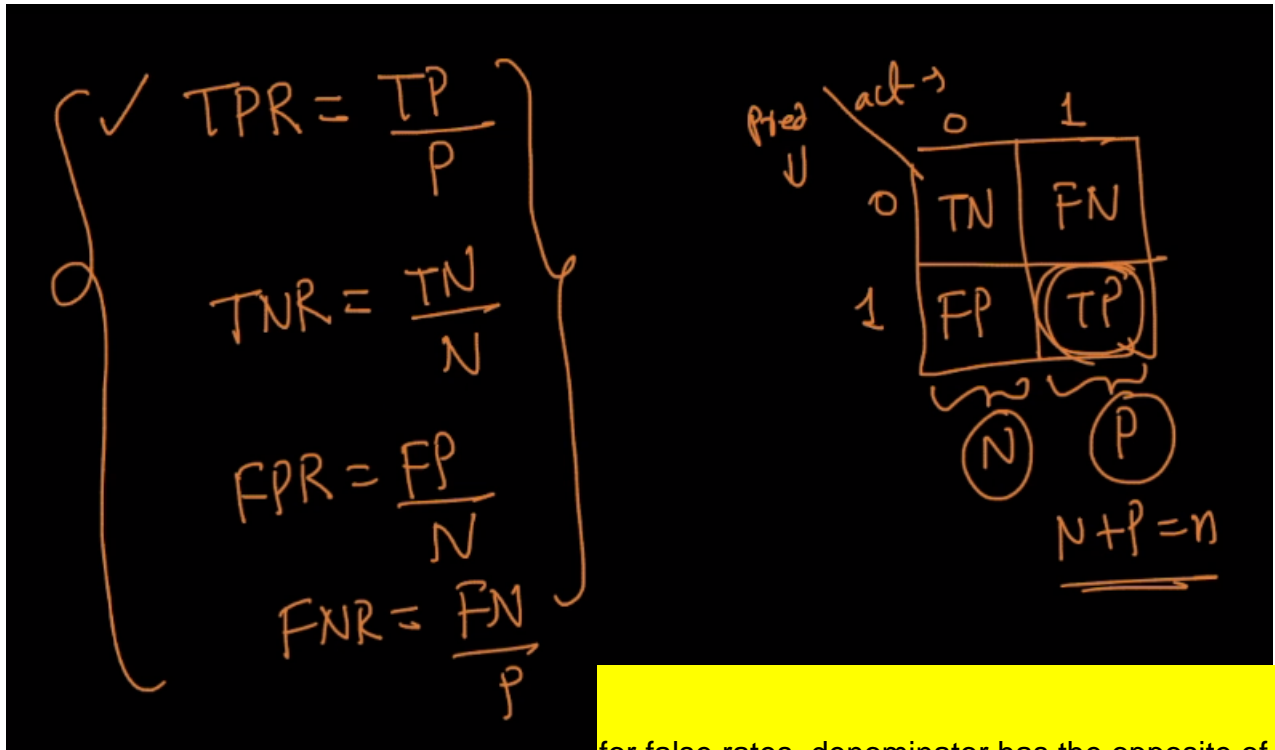We can draw a matrix for the predicted and actual values.

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:

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

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

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

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

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
        100 +ve   balanced

Model

↑ TPR = 94%.

↑ TNR = $\frac{850}{900}$

FPR = $\frac{50}{800}$ ↓

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.
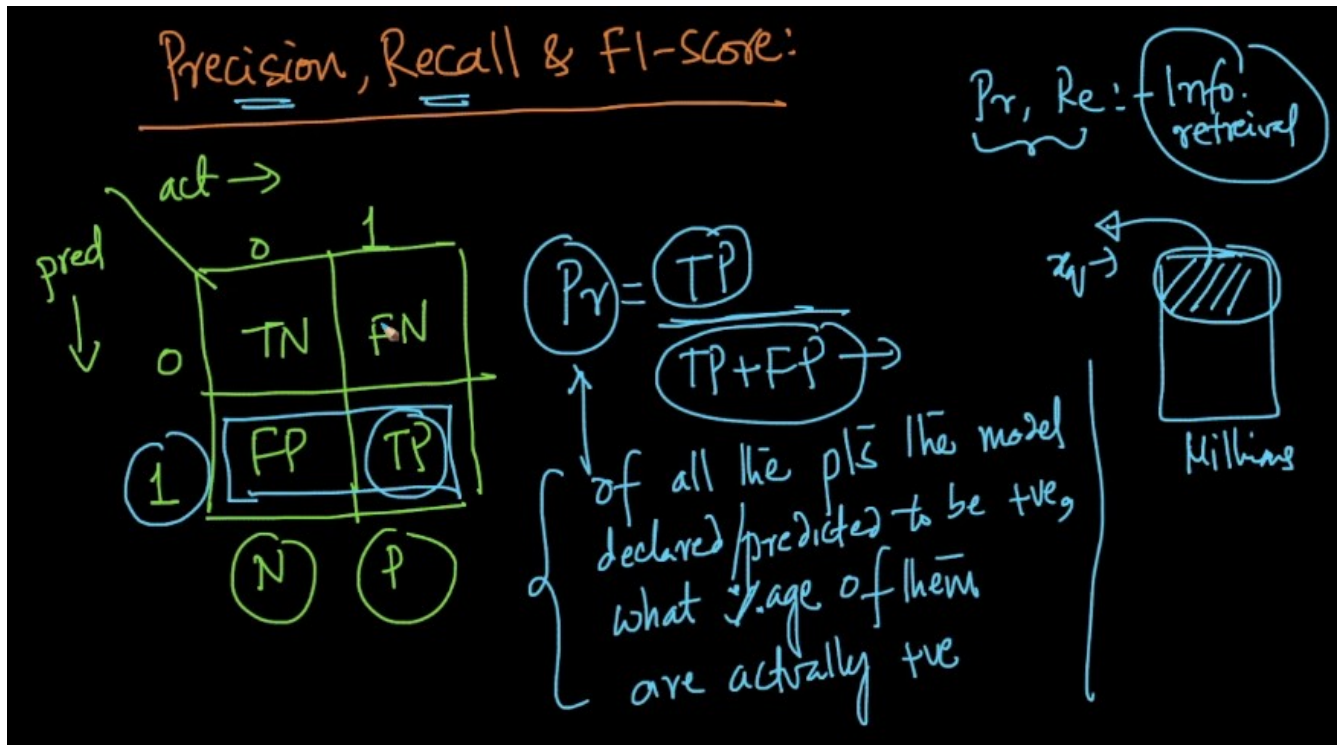
Precision, recall and F1 – score:



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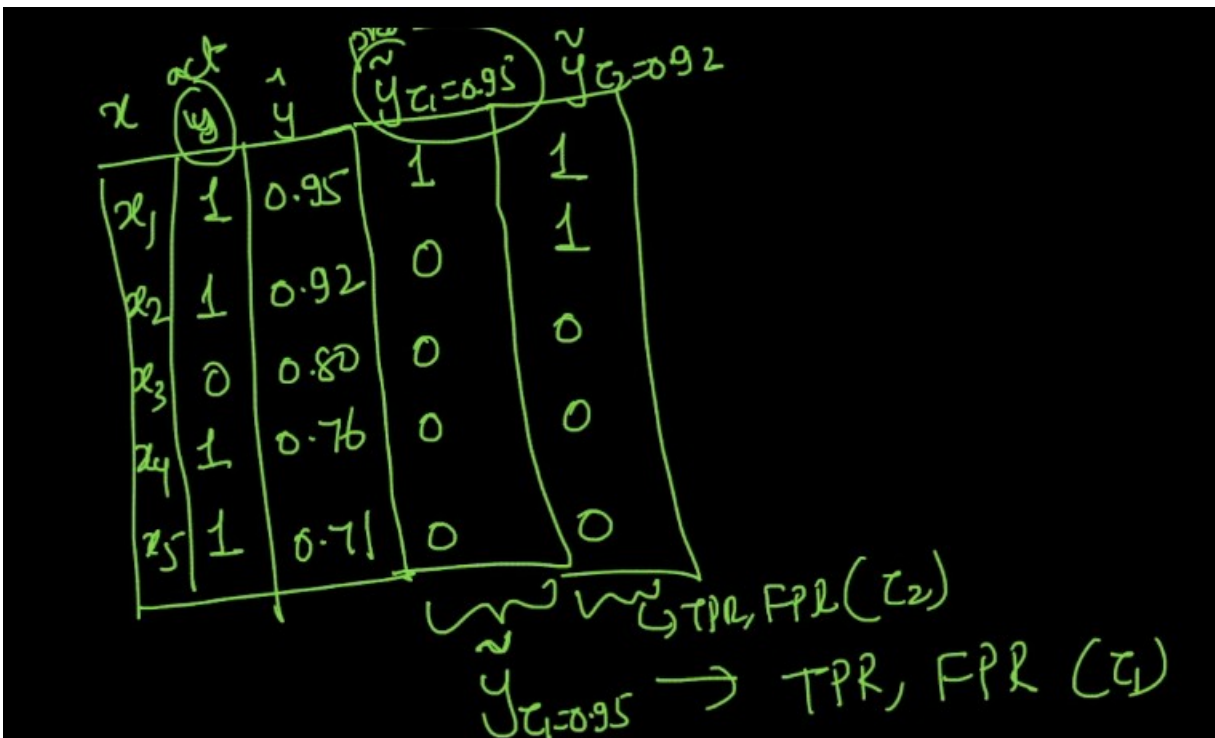
Receiver Operation characteristic curve and accuracy curve:
We use thresholds as the base in AUC and ROC curve. The values must be sorted in decreasing order of class scores.



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:

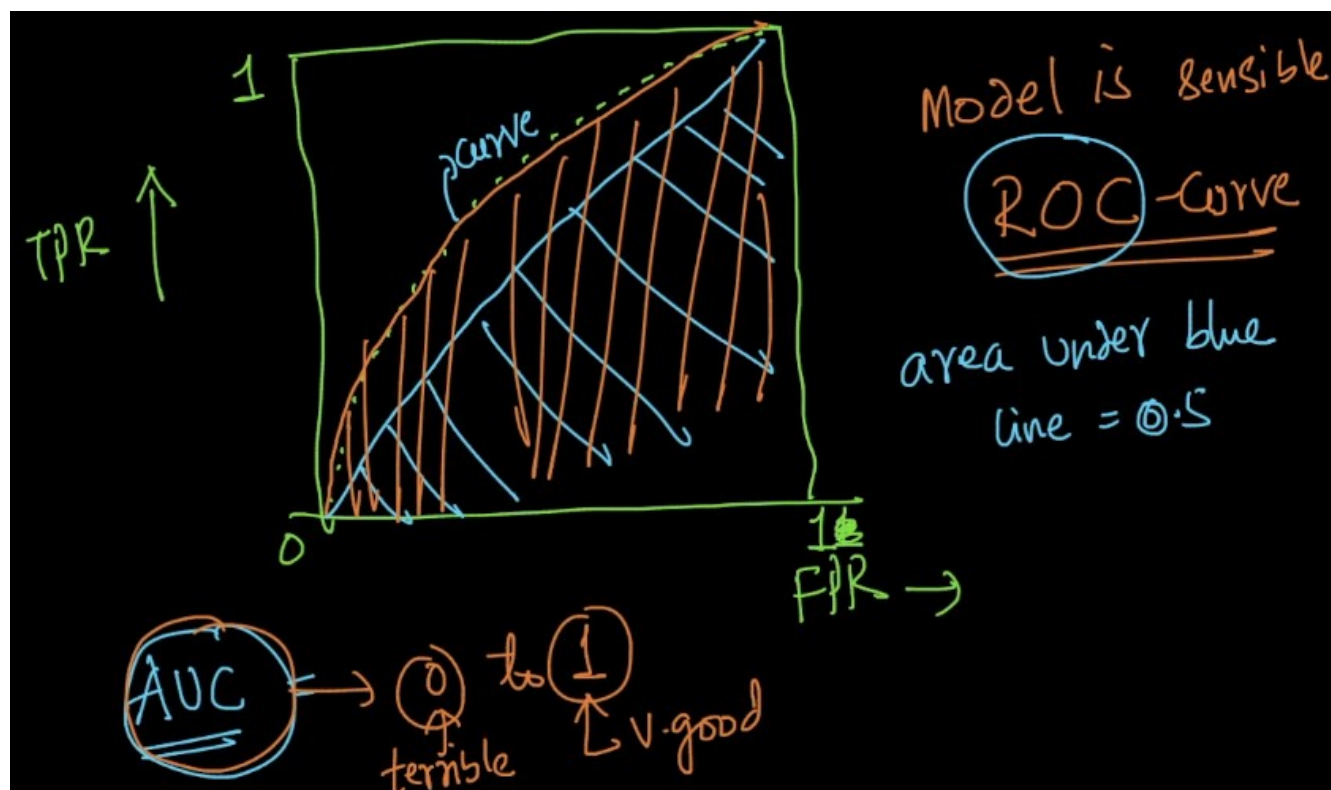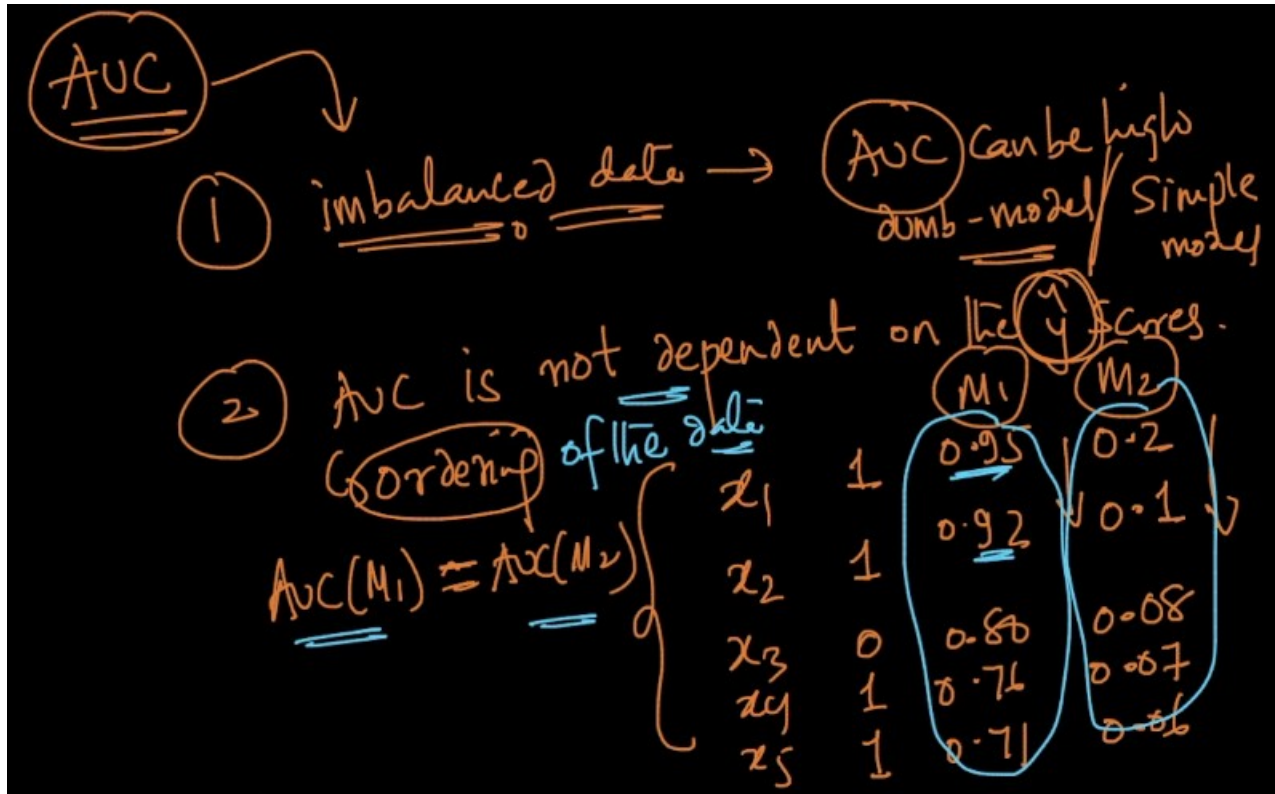1. In case of imbalanced data the AUC can he high.

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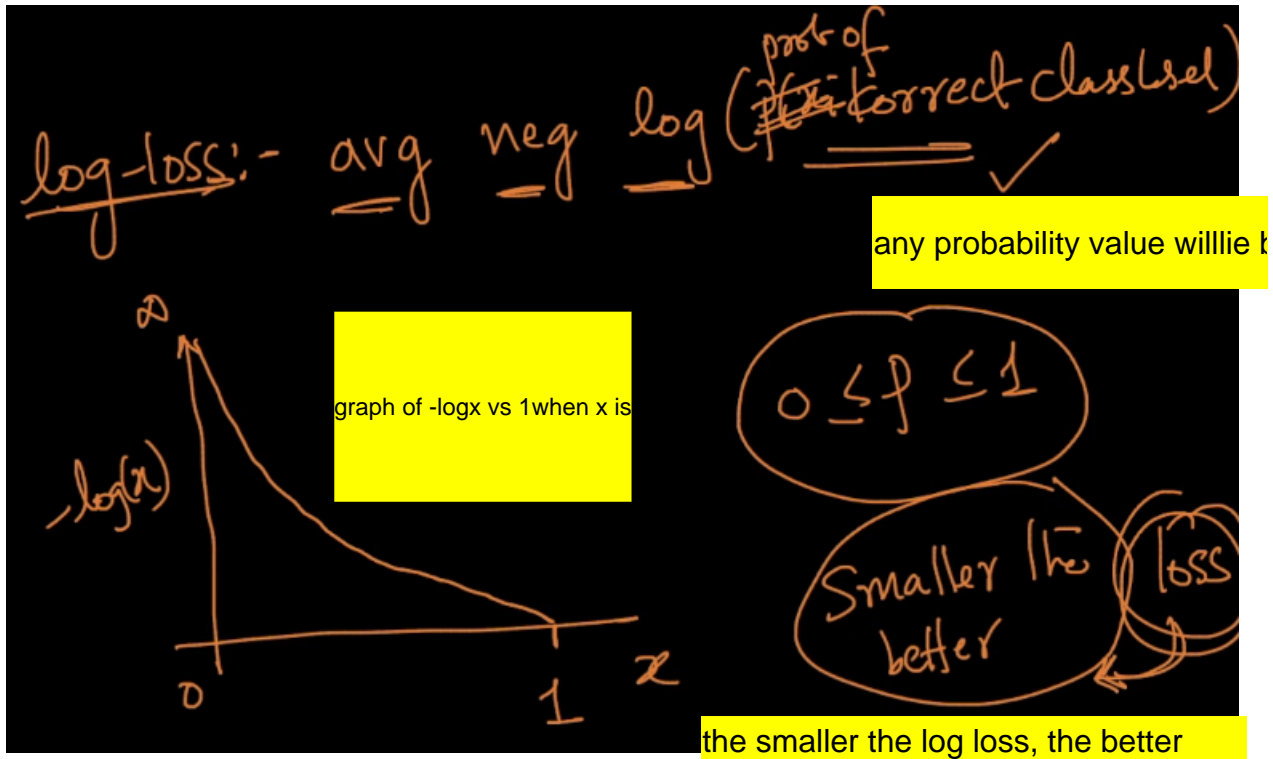
If the model is random then the AUC will be same as the diagonal line, i.e, 0.5.

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

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.

The smaller the better is the model.

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

$-\log(x)$ graph with x-axis labeled $x$ from $0$ to $1$, y-axis labeled $\infty$.

$0 \le p \le 1$

Smaller the loss better

Log loss function for multi − class classification:

Multi-class log loss:-

$l_q \to$ [ $l_1 \ P_2 \ \cdots \ P_c$ ]

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

$P_{ij} \to$ prob that $x_i \in$ class $j$

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

$0$ o/w

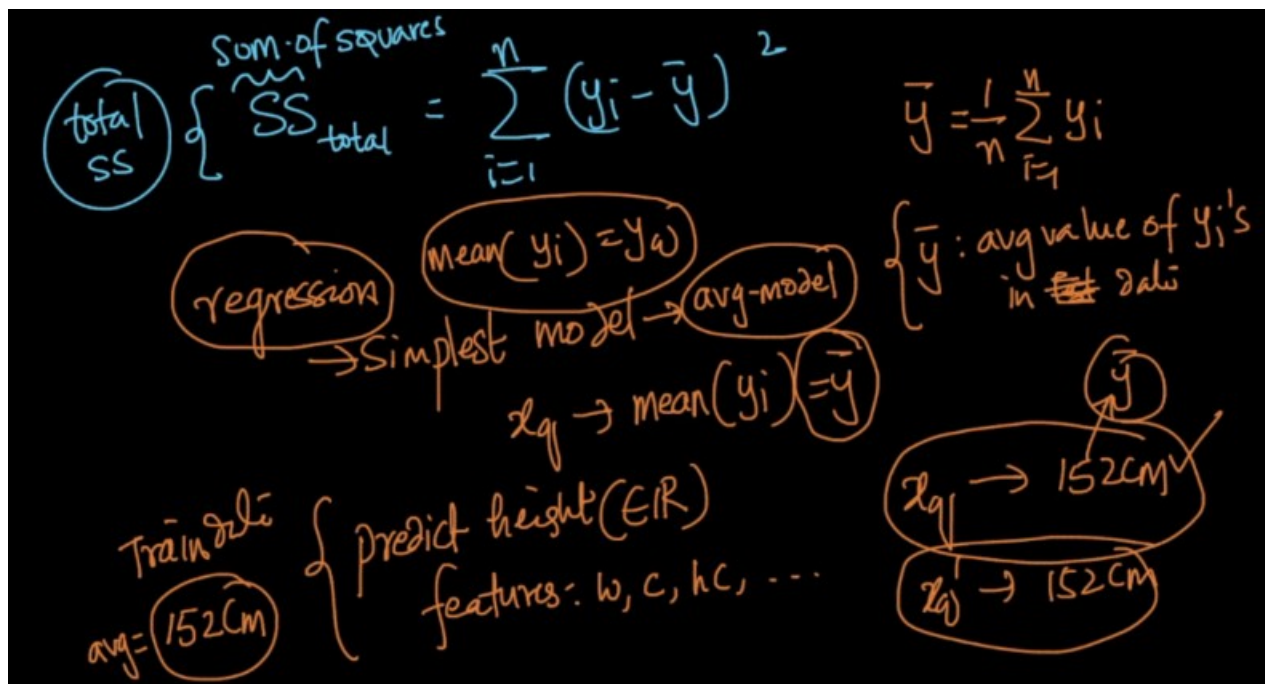The disadvantage of log − loss is we cannot interpret.

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. Herey(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:



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

Another case:



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.

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

$\boxed{EDA}$

$e_i$ : random-variable

$\boxed{\text{mean}}$

$\boxed{\text{median}}(e_i) =$ central-value of errors

$\boxed{\text{Std-dev}}$

$\boxed{MAD(e_i)} = \text{median}\left(\left| e_i - \text{median}(e_i) \right|\right)$

abs

$|e_i|s \rightarrow 0 \rightarrow \boxed{great}$

$|e_i|s \rightarrow$ large $\rightarrow$ not so good
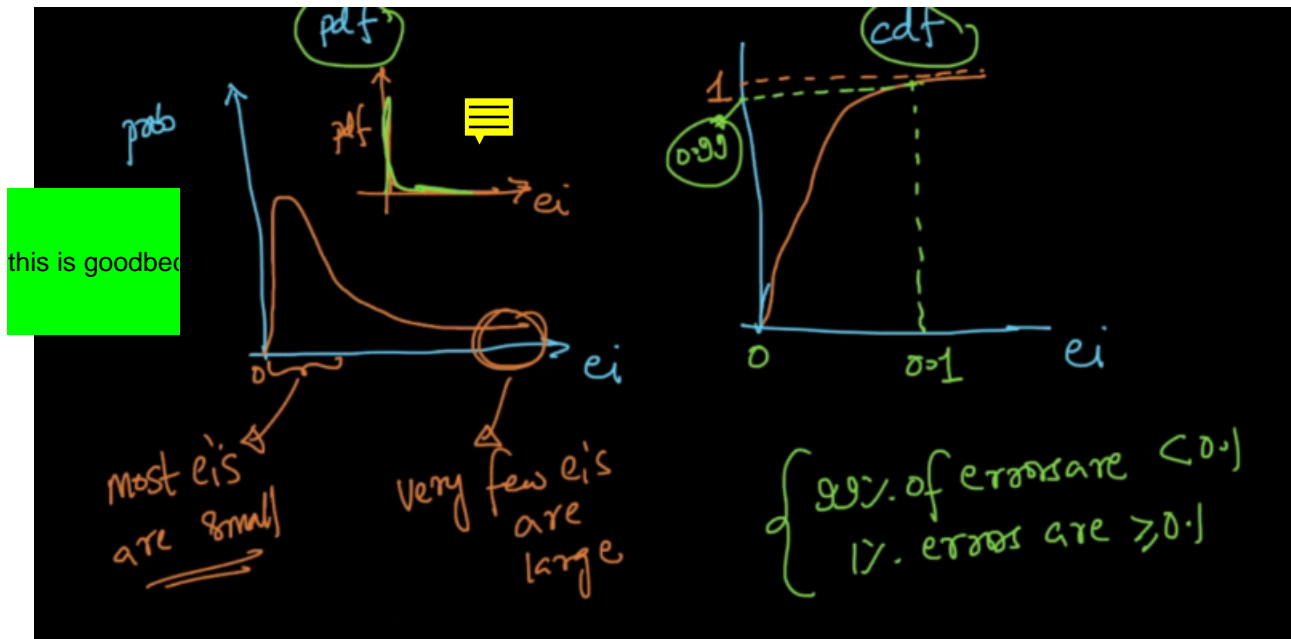
$\boxed{small}$

$= \boxed{small}$

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).
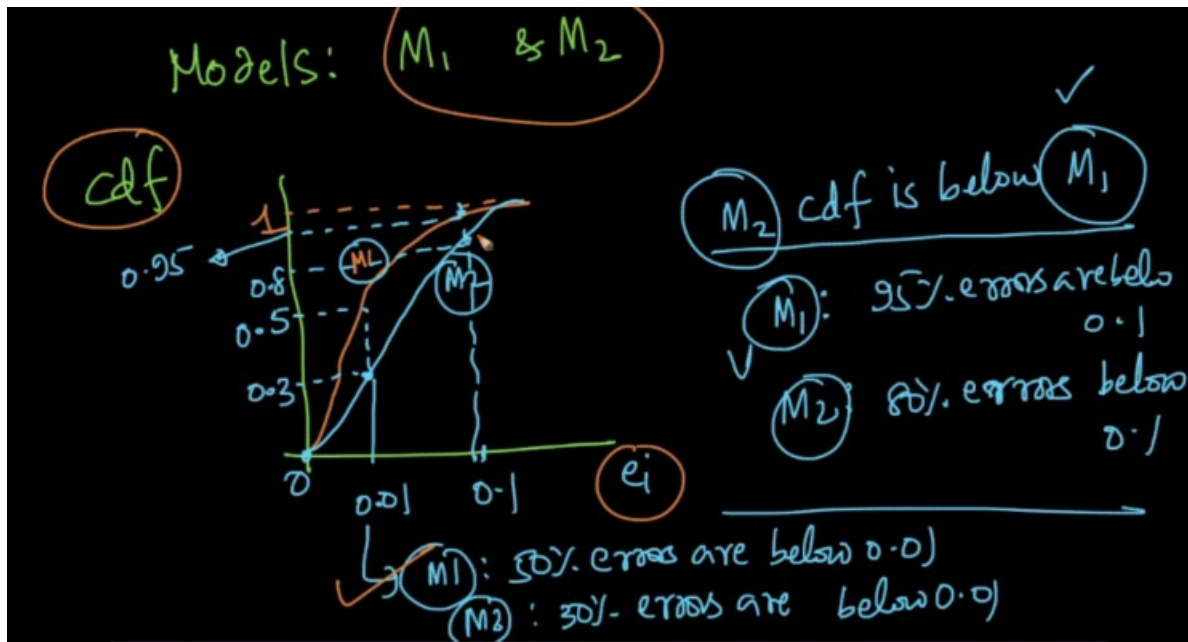
Distribution of errors:
We can use PDF and CDF for knowing the distributions of the errors.



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.