

## Classification Algorithm

How to deal with imbalanced data ? :

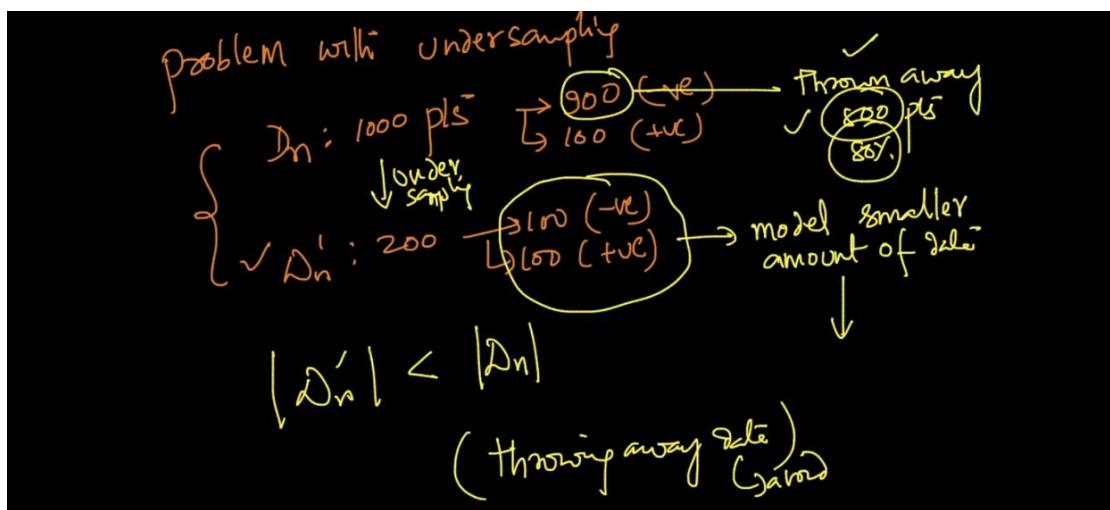
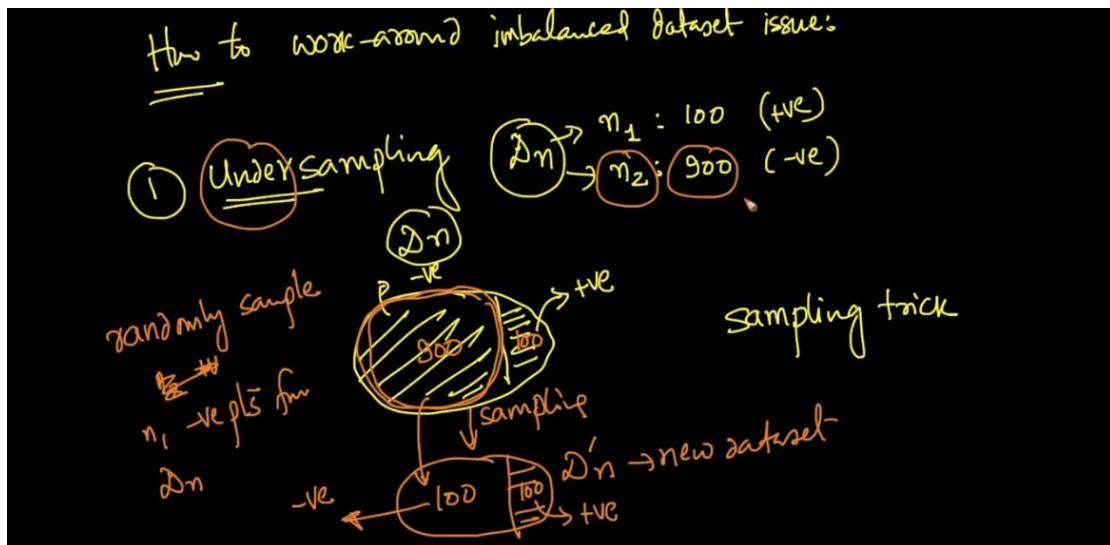
Imbalanced data means data is not proper for example lets say we have cancer data set 1000 values . In that 100 are + patient and 900 are - patient this is called not balanced .

There are different methods used to handle Imbalanced data :

1. Under Sampling :

In this method what we do is we can make data balance by removing 800 (-) points .

So we get now 100 (+) and 100 (-) . But this method is very bad because we are missing too much data so our model will not work well.

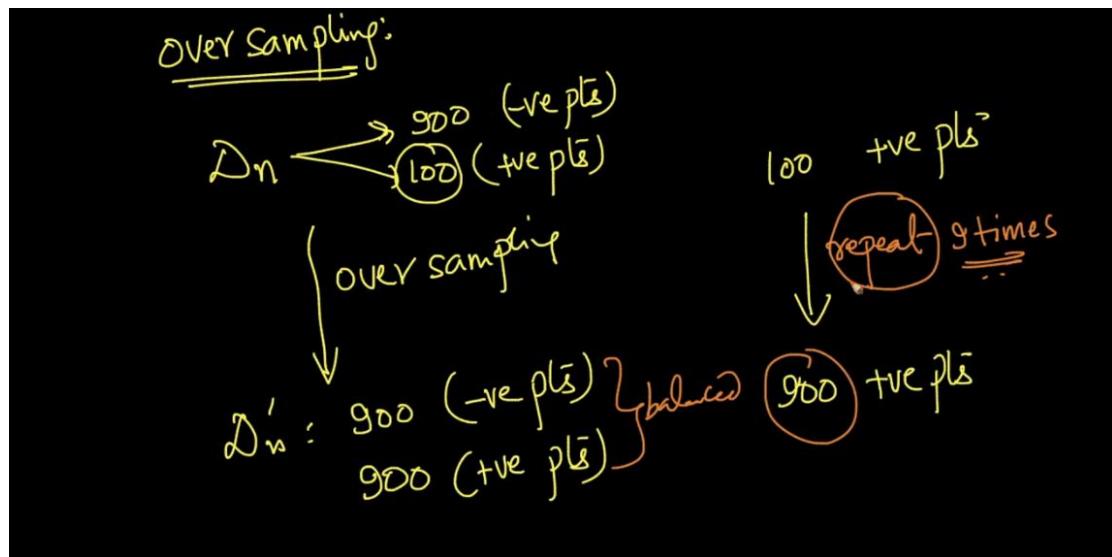


## 2. Over Sampling :

In this method we make data equal just by repeating same data for minority class .

In above example we repeat (+) points 8 times so we get + points 900 and - points 900 now our data is balanced .

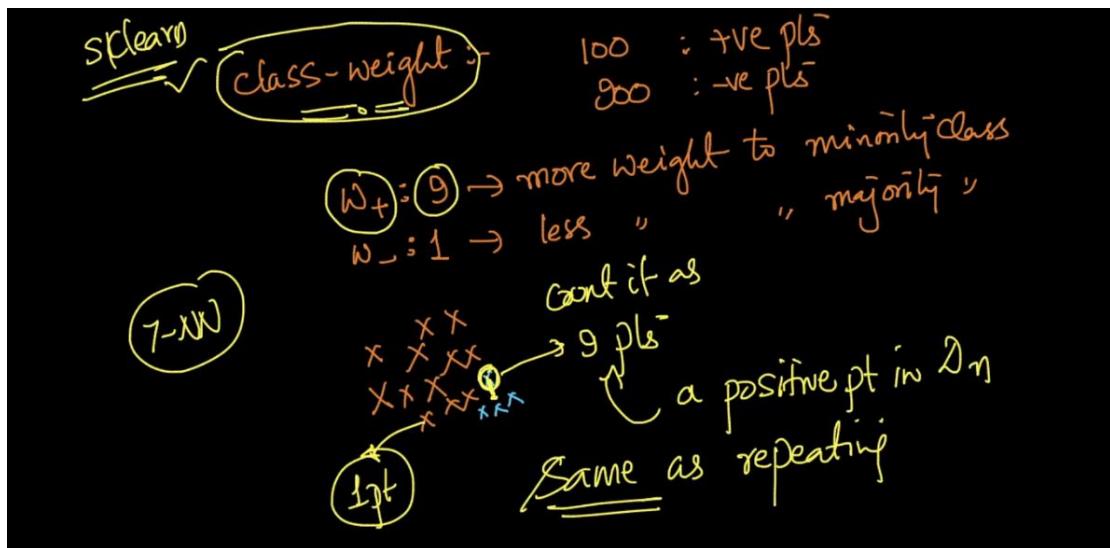
This method is too much better than previous because here we are using all data .



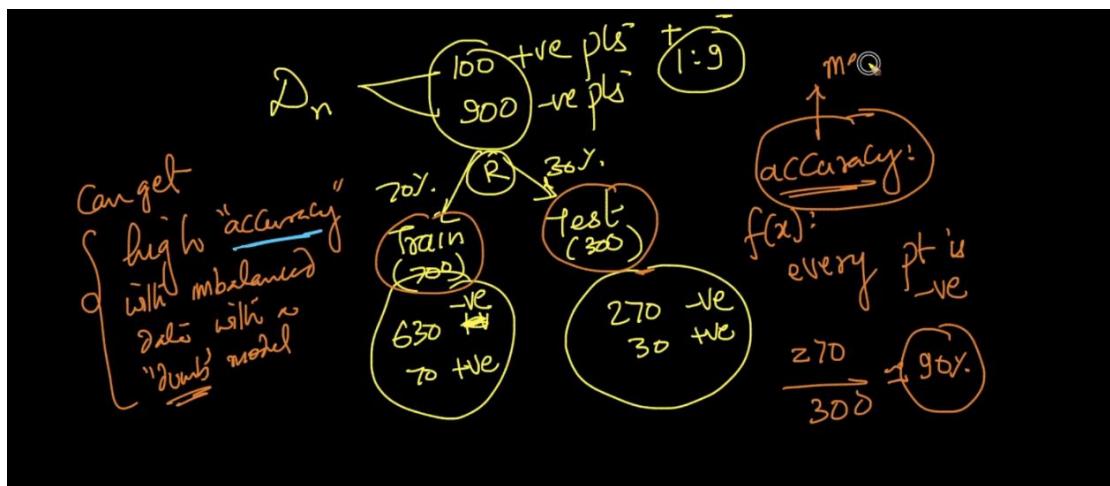
### 3. Class Weight :

In this method what we do is we will more weight to the minority class (+ class in our example) .

That means 1 + point = 9 (-) point in this way we can balance data its same like 2<sup>nd</sup> method .



Imbalanced data set occurs lots of time in real world problem .



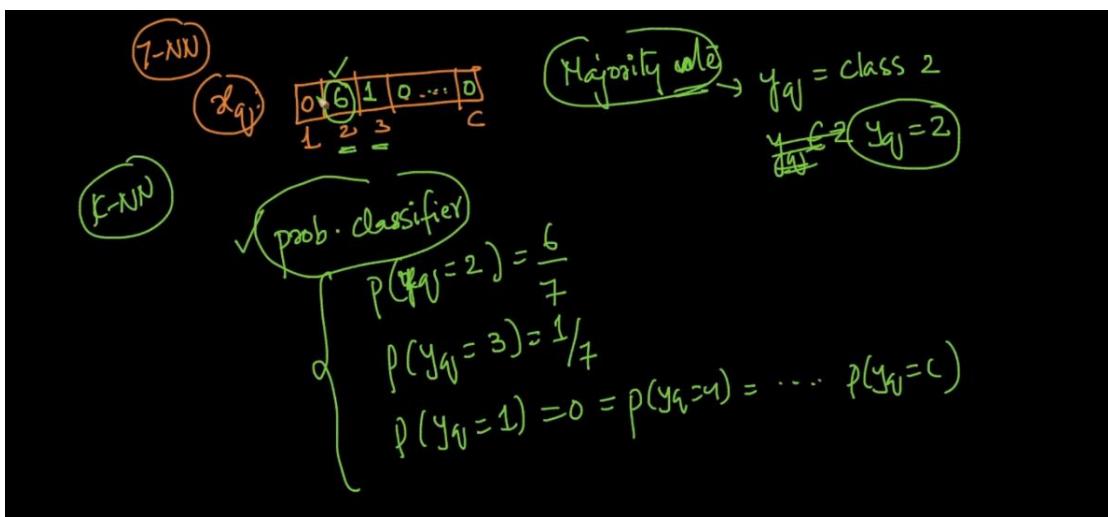
See above image we get 90% accuracy on test data its very wrong because of not balanced data .

## M class Classification :

Algorithm like KNN can very simple do multiple classification but logistic regression algorithm can do easily .

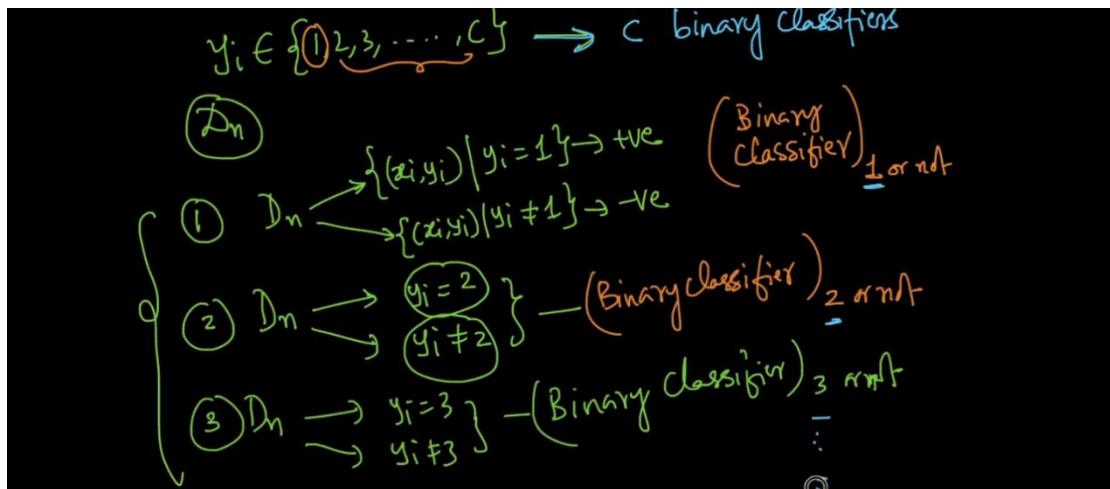
Lets say we have x class we are using KNN then we can easily predict lets say  $k = 7$  .

Then we simply see say how much nearest value to class 1 ,2 . .... x . and we take majority vote and we classify .



## Binary Classifier :

Not able to classify multiple class but we can solve that problem see below image .



We will make  $c$  classifier and will make classification . this will increase time and space complexity too much .

Similarity matrix means we create matrix for two variable and say how much they are similar .

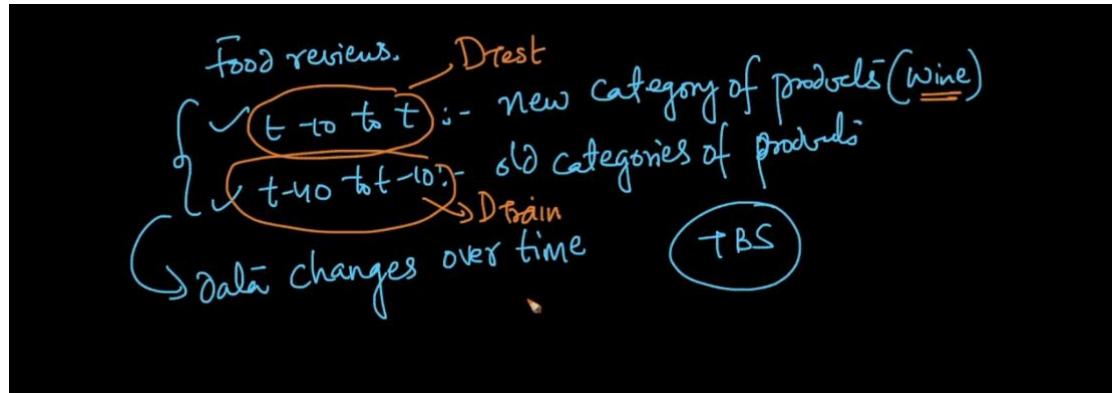
Ex. Drug1 and Drug2 if we go to the pharmacy and ask then how much they are similar we get values so that is nothing but similarity matrix .

### Train and Test Data :

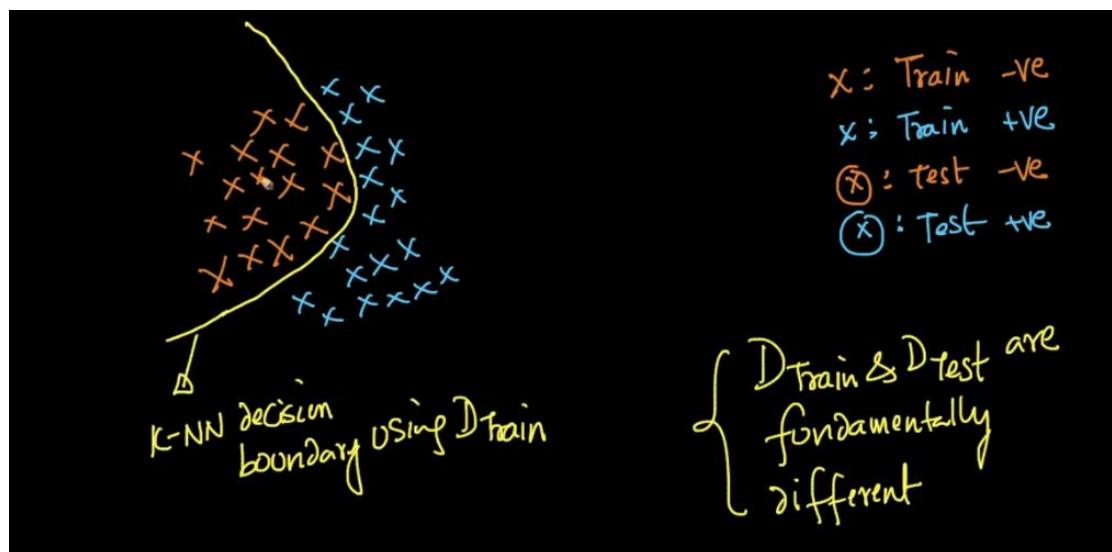
Lets take food review example as we have time based splitting so we have train data for today-60 days . and test data for last 10 days .

Now in last 10 day amazon add new category bear . so review for this product will be very different as

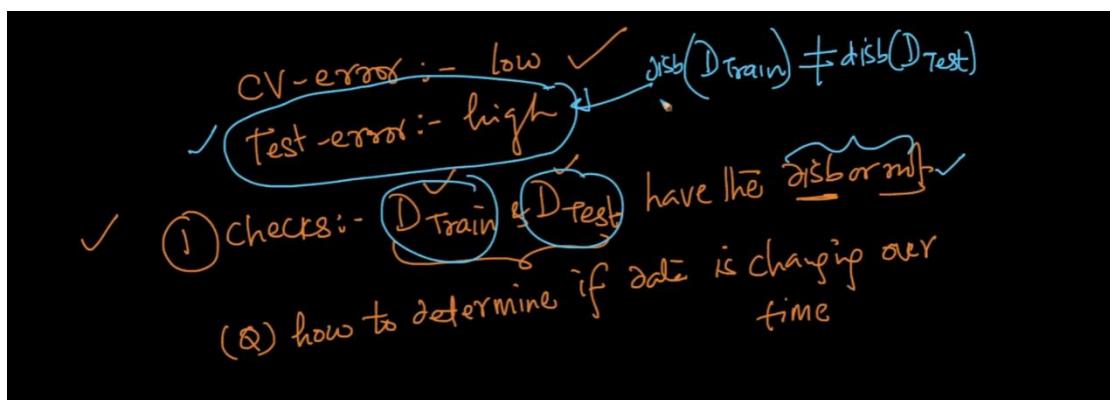
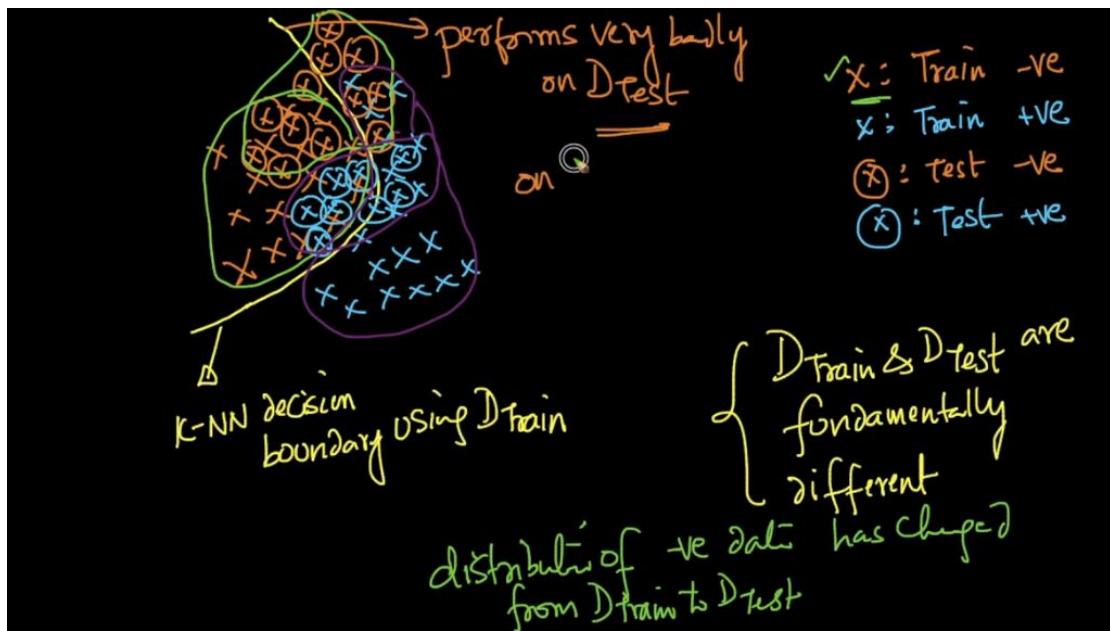
it is alcoholic food . but in train data we don't have this data ..



Because of this our whole distribution of data will change from train to test data .

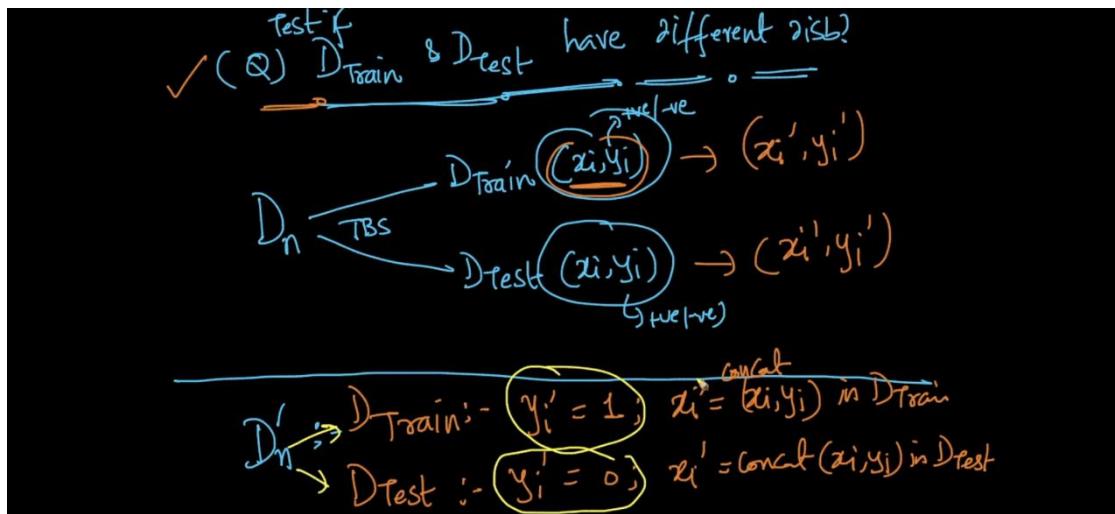


See below image now that yellow line work only well for train data but if we test same model on new data it will give very high error .

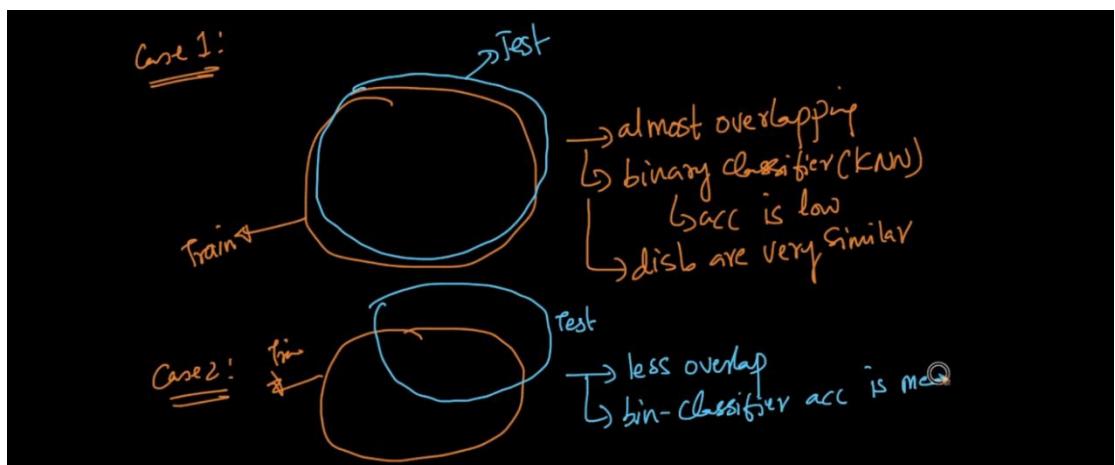


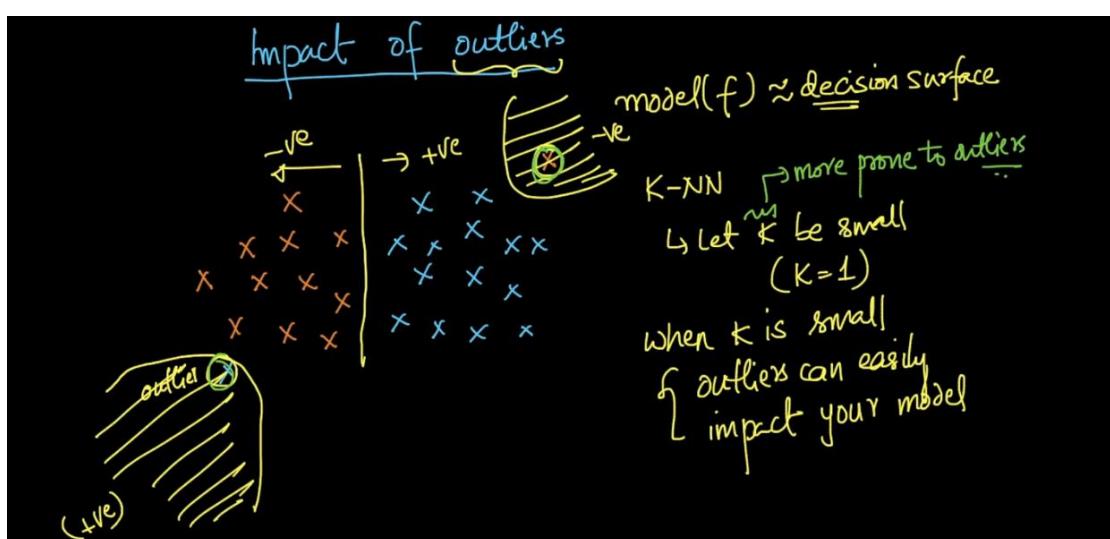
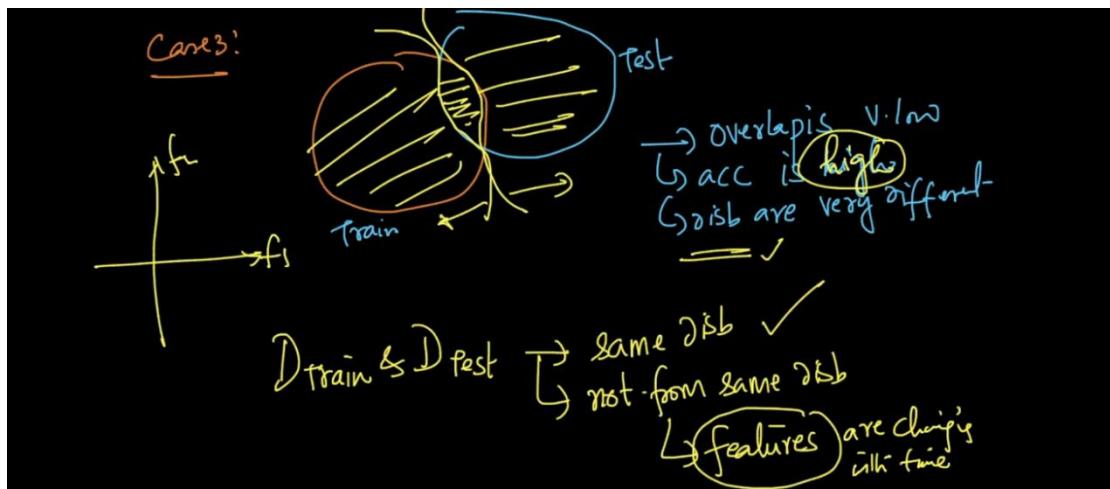
If distribution changes very speed then no model work very well so we need to check first whether both train and test follows same distribution or not .

So we create new data set  $D'$  with all train data as label =1 and all test data label = 0 .

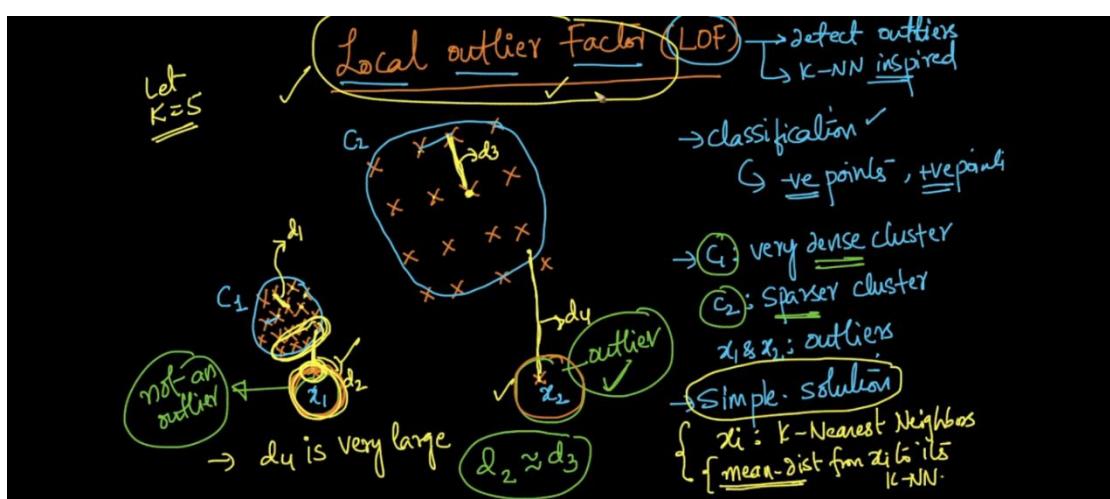


Then we can apply binary classification on new data to check whether they are from same dist or different .





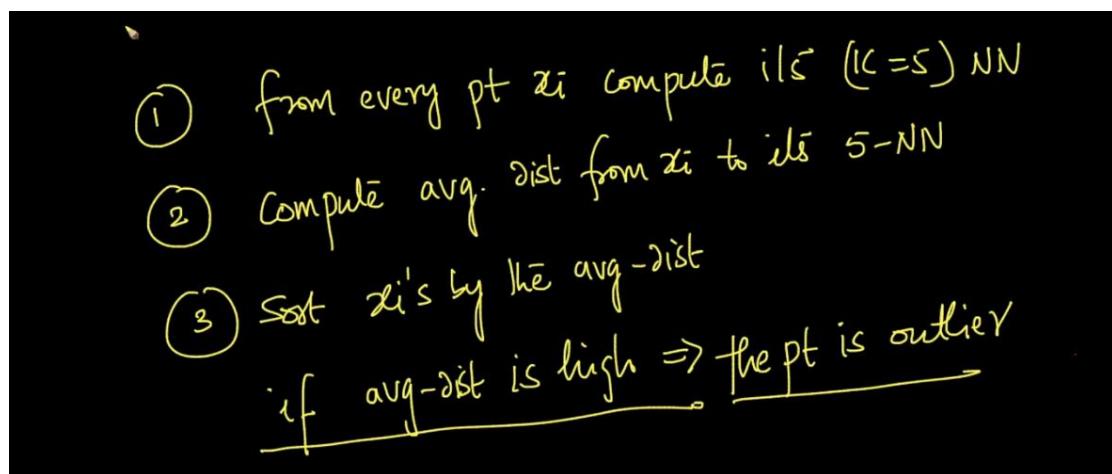
Always chose high K because if we chose  $k=1$  and there is outlier then output may be affected .

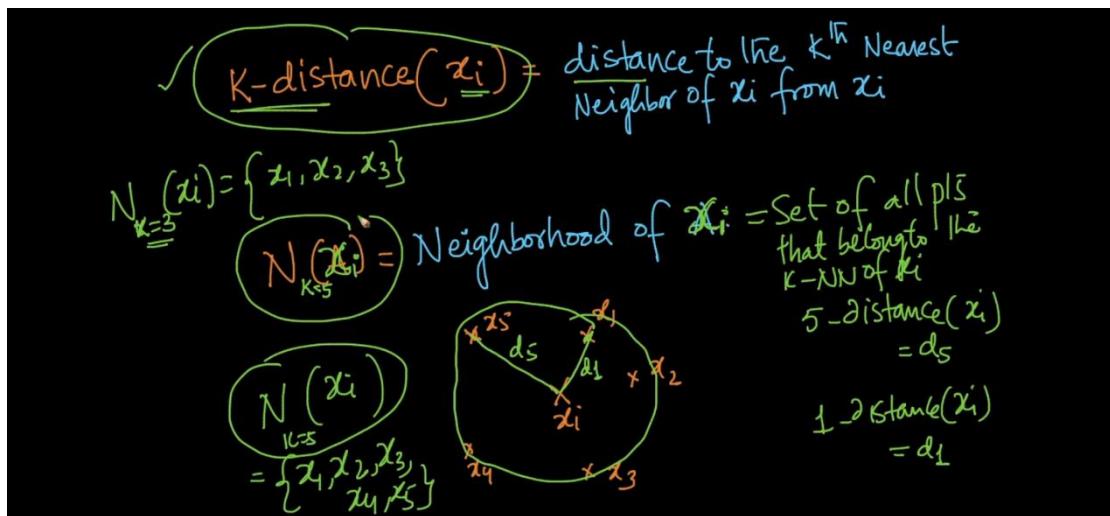


See above image LOF :

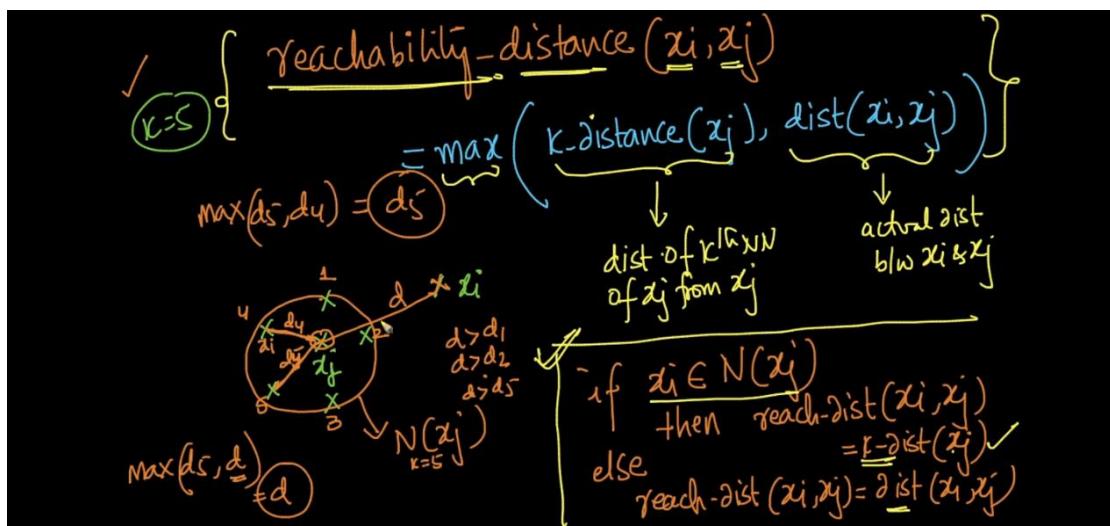
Lets say we have  $k = 5$  so for each point we will find 5 NN and will take average of that point that is nothing but distance  $d_1$ . similarly we find  $d_2$  for cluster  $c_2$  and  $d_3, d_4$  for 2 outliers . (Average distance for 5 neighbour)

Now we sort all  $d_1 \dots d_4$  . in this method problem is  $x_1$  not consider as out l. because  $d_2$  and  $d_3$  are almost same so if we say  $x_1$  is out l. then whole  $c_2$  will be out l. hence this method will not work well .

- 
- The image shows handwritten notes for the LOF algorithm. It consists of three numbered steps:
- ① from every pt  $x_i$  compute its ( $k=5$ ) NN
  - ② compute avg. dist from  $x_i$  to its 5-NN
  - ③ Sort  $x_i$ 's by the avg-dist  
if avg-dist is high  $\Rightarrow$  the pt is outlier



## Reachability Distance :



See above image . we will take max distance from formula .

local reachability density:  $\text{lrd}(x_i)$

$$\text{lrd}(x_i) = \frac{1}{|N(x_i)|} \sum_{x_j \in N(x_i)} \text{reach-dist}(x_i, x_j)$$

den: avg. reach dist of  $x_i$  from its neighbors

$|S| = \# \text{ele in the set}$

$x_i$

$N(x_i)$

$x = x_i$

set of  $x \rightarrow \text{NN to } x_i$

$\text{lrd}(x_i) = \frac{|N(x_i)| \rightarrow \# \text{points}}{\sum_{x_j \in N(x_i)} (\text{reach-dist}(x_i, x_j)) \rightarrow \text{distance}}$

density

local-reachability

LOF :

Local outlier factor ( $x_i$ )

$$= \frac{\sum_{x_j \in N(x_i)} \text{lrd}(x_j)}{|N(x_i)|} *$$

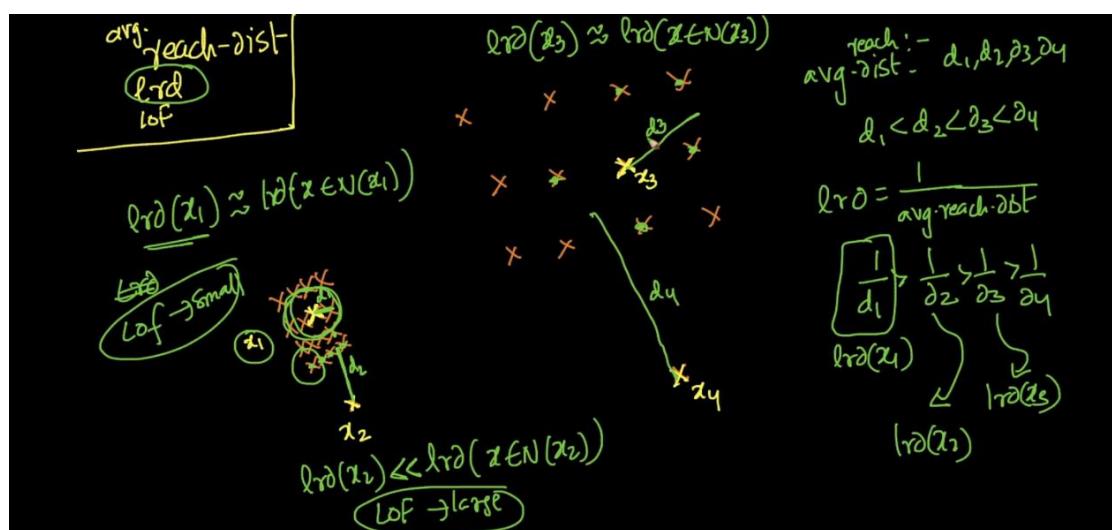
$\text{lrd}(x_i)$  is large  $\Rightarrow$  small

large

avg. lrd of pts in the neighbourhood of  $x_i$

See above image LOF formula . LOF will be high for any point if any one of the term is high in formula .

If both term are small then only LOF will be small .



See in this image we have 2 outliers. now we say point x2 and x4 are outliers .

Lets take 1 point from both cluster x1,x3 . If we find avg reach distance which is nothing but avg of all neighbour points . we get d1,d2,d3,,d4 .

Now LRD of d1 is high and density of NN around x1 is also high so LOF will be small .

Now LRD of  $x_2$  is much less than  $x_1$  and density of its NN is high see in fig . So high + low = Large LOF .

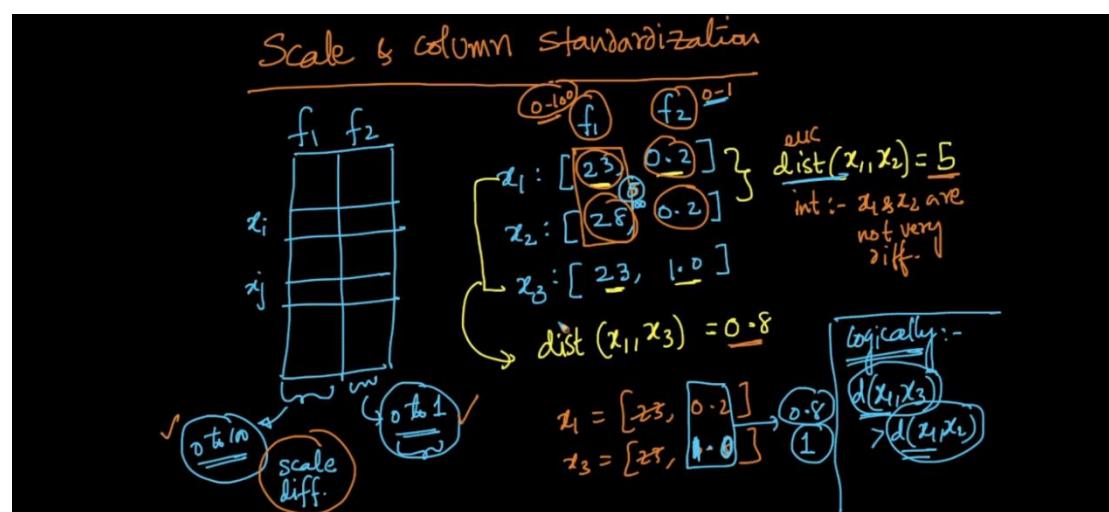
Same LRD of  $x_3$  is small and density of its NN is small (Not too dense) so LOF will be Low .

LRD of  $x_4$  is much less as compared to its NN density . and density of its NN is high so LOF is high .

So we can consider points as outlier . whose LOF is very high .

But there is no scientific logic how much high LOF is bad we need to take that decision as per our data ..

### Column Standardization :



Now see above image lets say we have feature 1 and feature 2  $f_1$  and  $f_2$  and both are from different scale means one lies from 1 to 100 and  $f_2$  lies from 0 to 1 .

If we measure E distance as per that we found distance between  $x_1$  and  $x_2$  is too much greater than  $x_1$  and  $x_3$  .

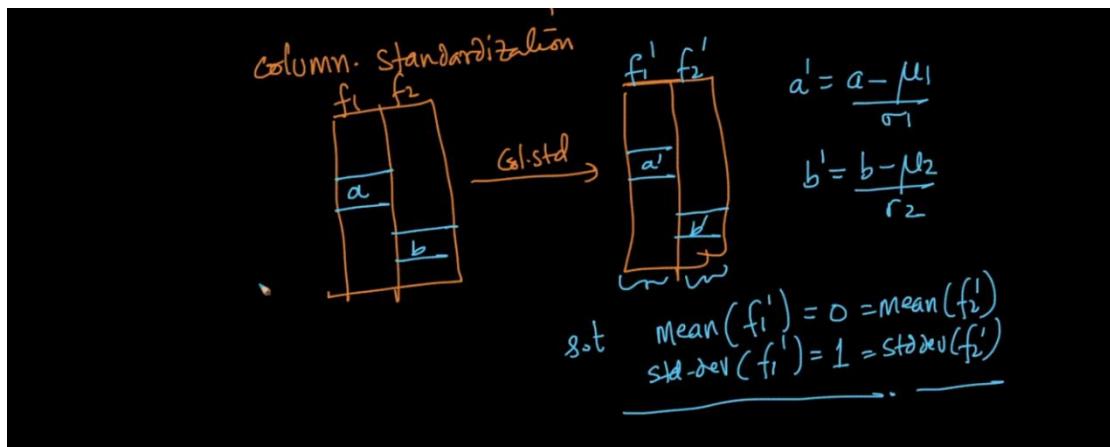
But logically this is wrong this issue because of different scale and due to this we get wrong NN in KNN .

If we see distance between  $x_1$  and  $x_2 = 5$  but range of  $f_1$  and  $f_2$  lies from 1 to 100 so as compare to that is very small .

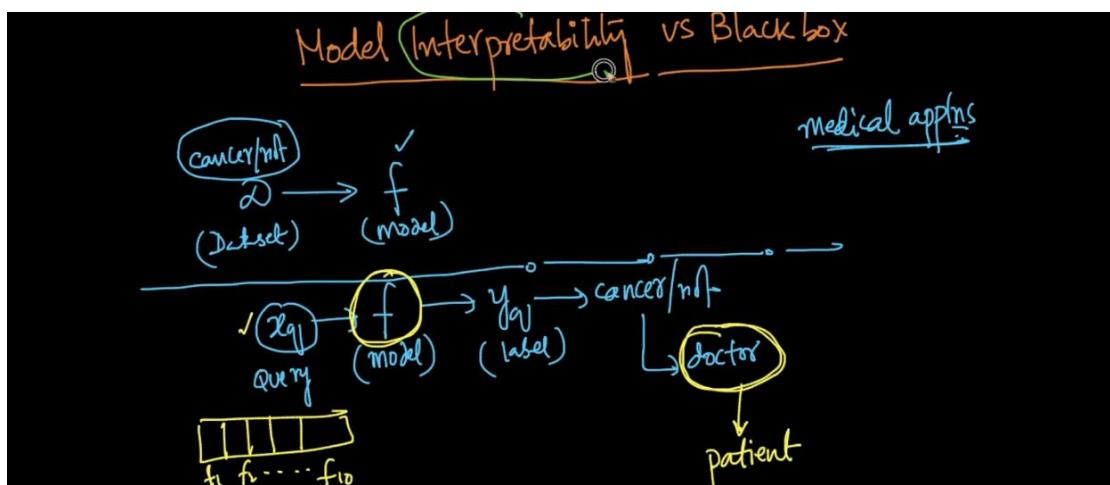
Other side distance between  $x_1$  and  $x_3$  is 0.8 and max distance is 1 so it is too high ..

So to avoid this type of problems we need standardization of col . See below image .

So each time when we are applying KNN need to do 1<sup>st</sup> is standardization of data .



## Model Interpret and Black Box :



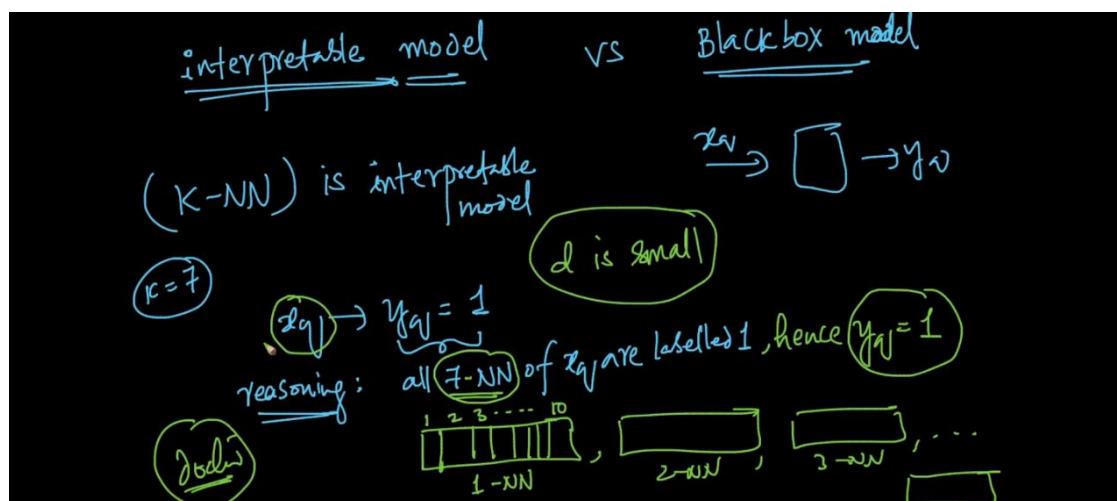
As we know that we have a data we train it and create a model ,when new data come we just use trained model to predict output .

Lets say we trained model on cancer data set now new patient come and our model predict he has a cancer . but patient and doctor both need reason why he has a cancer on what basis our model is telling that because doctor don't know any maths behind that we need a reason .

That is nothing but interpret model not black box .

Means it will give perfect reason and will tell I found f5 (test 5) result is very high and f8 is very low and because of this I say this patient has a cancer .

Then doctor can easily check test 5 and 8 and he will understand whether our model predict correct or wrong .

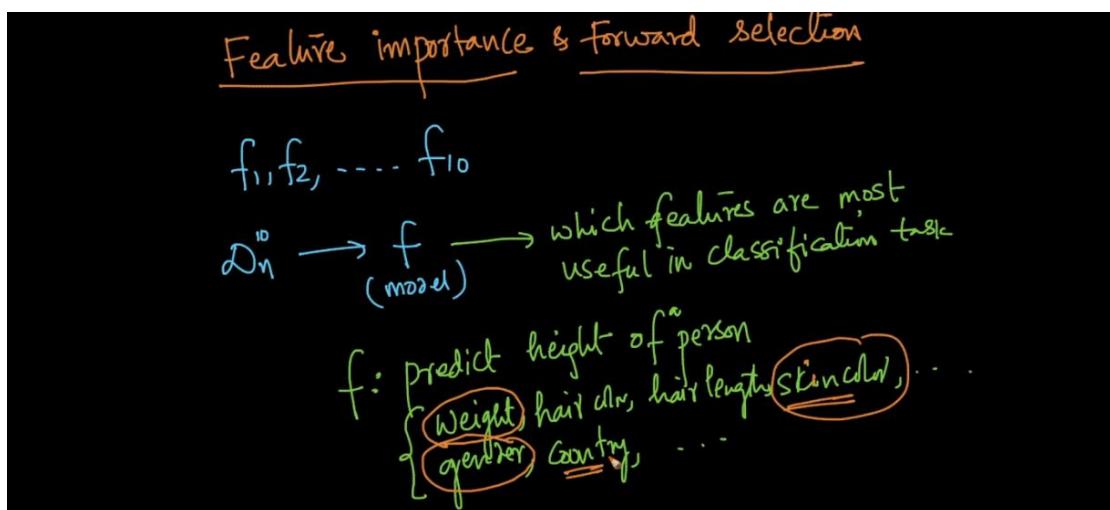


KNN is interpret model if d (dimension is small in our case here d = 10 ) so we can say it is interpret model ..

What KNN will do in this case is it will find lets say we decide k = 5 , so it will find 5 NN who has same

condition like new patient and as per that KNN will say Cancer or Not.

## Feature Importance and Selection :



In real time we can't say we have 10 features only we get 1000 sometimes more than that so using all of the feature will decrease model performance .

So now what we can do is we will try to remove non importance feature lets say from 1000 to 100 .

One way is PCA TSNE but they didn't care about classification one very powerful method is Features forward selection .

What we do here is we take f<sub>1</sub> and trained our whole model and will check accuracy on test lets say we get accuracy a<sub>1</sub> .

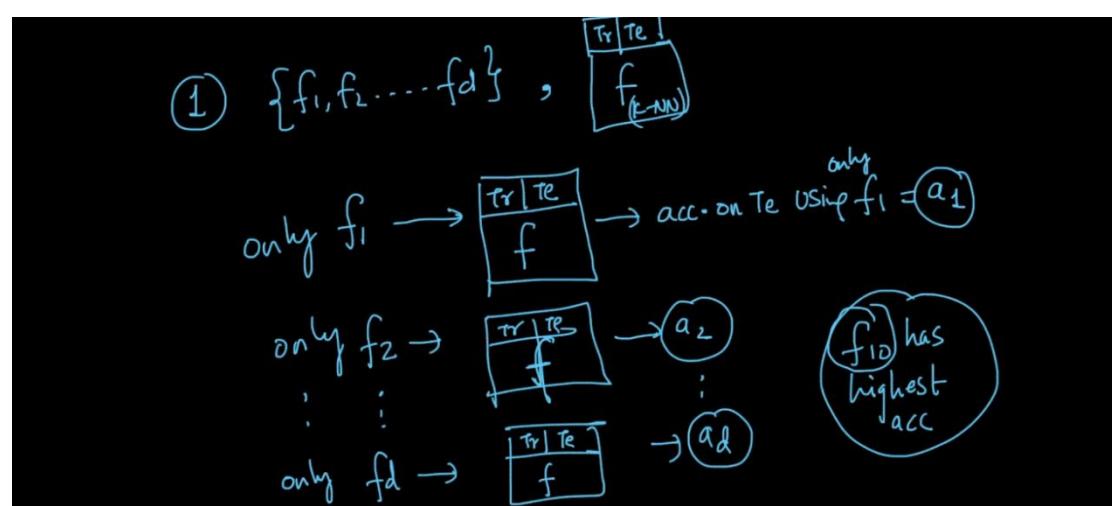
Similarly we repeat up to d features (1000) now we have 1000 accuracy we will check which one is highest then we select that feature . lets say f10 has very high accuracy so we select f10 .

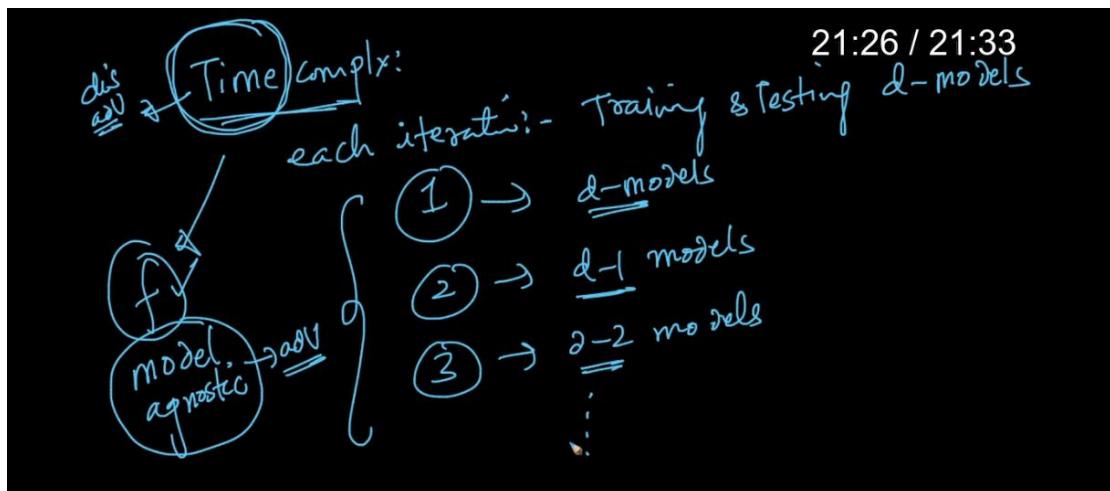
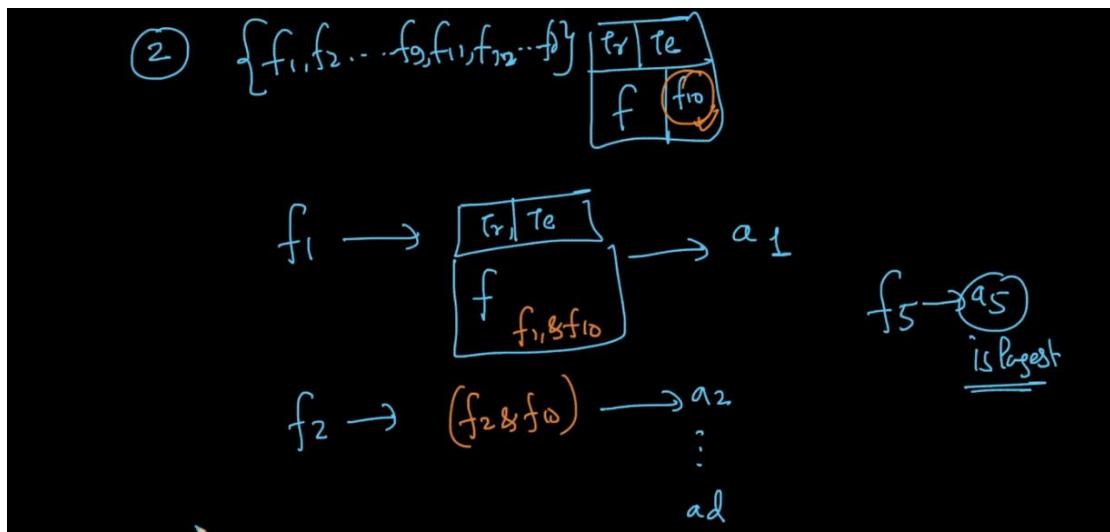
Then we take f10 + f1,f10+f2 ..... f10+ f d. then we chose highest accuracy feature lets we get

F10 and F5 .. similarly we repeat this process until we get 100 feature in this way we find best 100 features .

This method works on any any model .

Backward selection is exactly opposite to it we will remove 1,1 feature and will check how much effect on accuracy but forward is best .





Only the problem is it will increase time complexity very very high but still its good because its one time effort after that we get best features .

### How To deal with Categorical Data :

Lets say we category data for hair,color,country etc . and we want to predict height .

One way is we can convert feature into one hot . but the problem is lets say I have country feature and almost we know 200 unique name for country .

So one hot vector will be size of 200 and in that only one values will be 1 and 199 will be 0 this is like BOW problem sparse .

hair color = { Bl, Br, Red, G, Gray }

① Give a number ↗ number have an inherent order ↗ ordinal

Red > Bl X

② One-hot encoding: hair color Bl Br R G Gr

binary vector ↗  $x_i : \text{Bl} \rightarrow [0 \ 1 \ 0 \ 0 \ 0]^\top$

$x_j : G \rightarrow [0 \ 0 \ 0 \ 1 \ 0]^\top *$

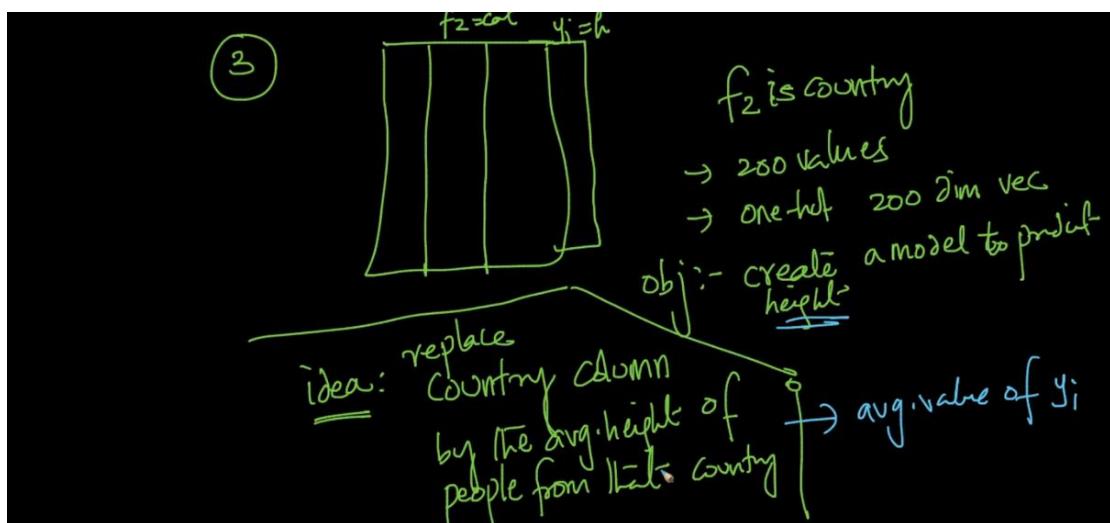
one-hot encoding:- bin-vector of the size of number of distinct elements-

country →  $\approx 200 \rightarrow$

$x_i : [0 \ 0 \ 0 \ \dots \ 1 \ 0 \ 0 \ \dots \ 0]^\top$  sparse

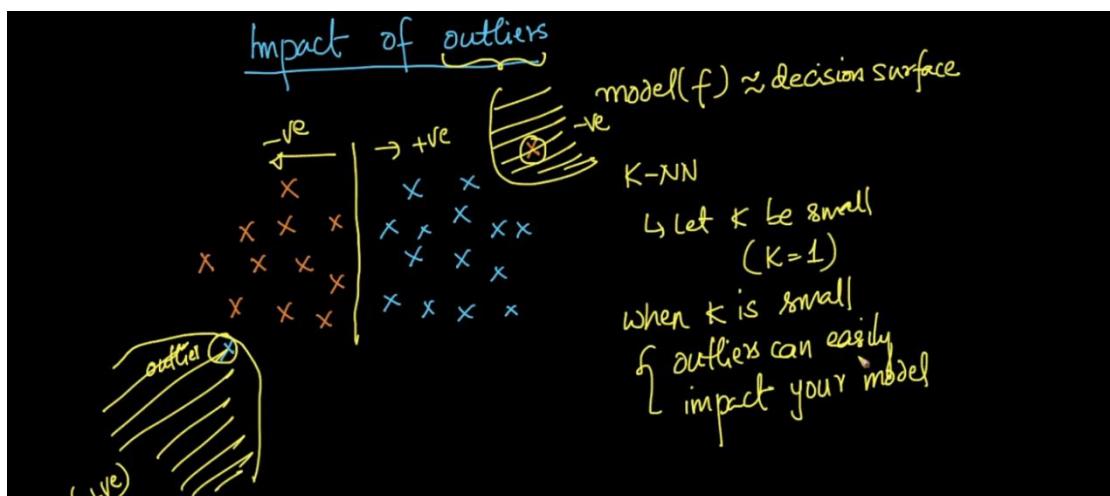
BOW ← one-hot encoding ↗ if the number of distinct values for a Categorical feature is large then one hot encoding can create Sparse & large vector

Another method is replace column with its mean value . like take height of people from India And find mean of that .



Else we can convert by domain knowledge .

There is no any fix logic for that we need to check as per our data . also we can experiment all method and check which one is best .



Always chose high k because out l with small K will disturb our model .

## Missing Data Handle :

✓ Missing Value imputation

data corrupted collection error

featurize it

① (imputation)

mean ✓  
median  
most freq (mode)

	$f_1$	$f_2$	$f_3$	$f_4$	$y$
$x_1$	✓	✓	✓	✓	
$x_2$	✓	✓			
$x_3$	✓	✓		✓	

MISSING  
NaN  
null  
-1

$x_i : f_3 \rightarrow \text{missing}$

all the non missing values for  $f_3$

1<sup>st</sup> is we can simply take mean of that column and will replace empty value with that mean value .

## Mean Imputation using Class label .

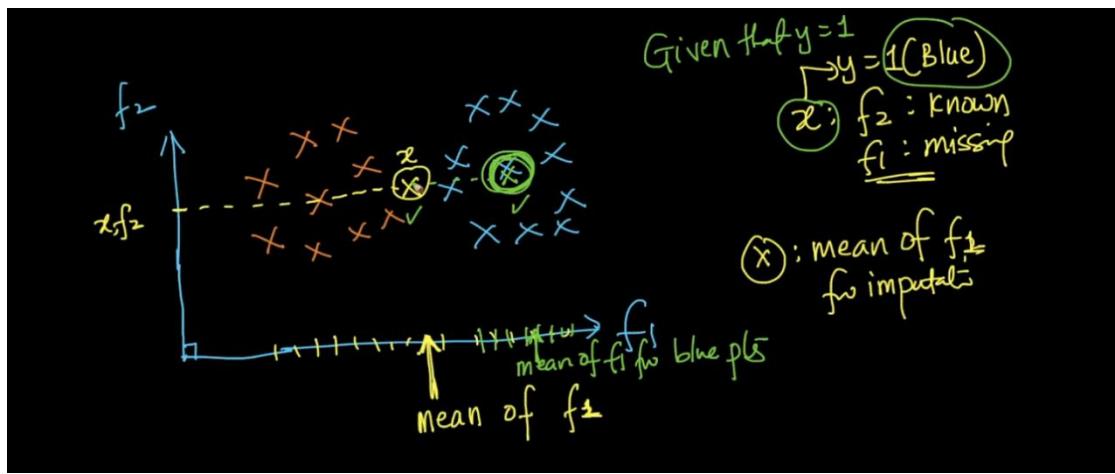
classification → 0, 1

impute based on class label

$\hookrightarrow x_i \rightarrow (f_3 : \text{missing})$

$\hookrightarrow y_i \rightarrow 1 \text{ (class label)}$

{ mean-imputation based on class label }      { mean of  $f_3$  for all pls whose class label is 1 }



Now instead of taking mean of whole column like f1 in above image we take mean of that class label area here class is blue so we take mean of blue area and that will very well than taking mean of whole column .

We took mean of blue because we already know that point belongs to that class.

new missing value features:-

	$f_1$	$f_2$	$f_3$	$f_4$	$f_1'$	$f_2'$	$f_3'$	$f_4'$
$x_1$	✓	✓	✗	✓				✓
$x_2$	✓	✗	✓	✓				✓

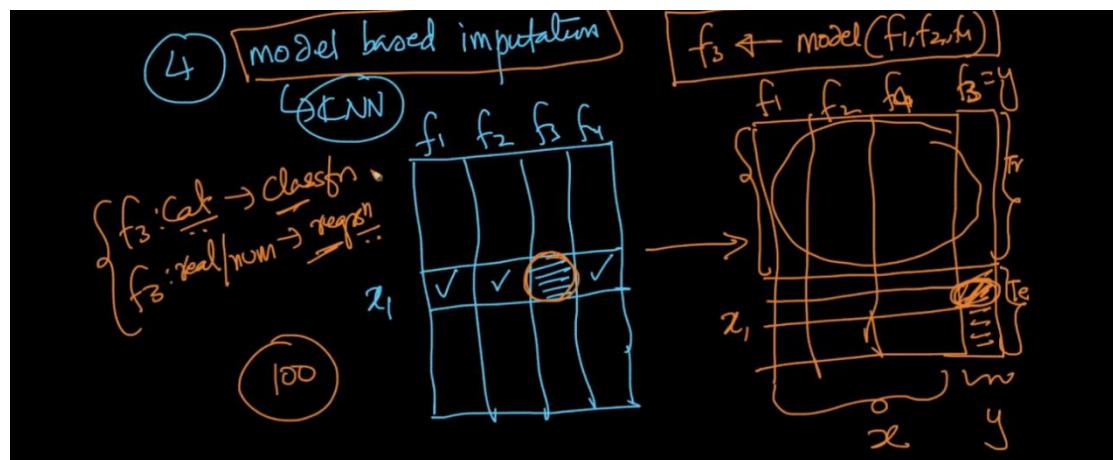
	$x_1$	$x_2$	$x_1'$	$x_2'$
$f_1$	✓	✓	✓	✓
$f_2$	✓	✗	✓	✓
$f_3$	✗	✓	✓	✓
$f_4$	✓	✓	✓	✓
$f_1'$	✓	✓	✓	✓
$f_2'$	✓	✓	✓	✓
$f_3'$	✓	✓	✓	✓
$f_4'$	✓	✓	✓	✓

Q) missing values → { Could be source of information } → imputation

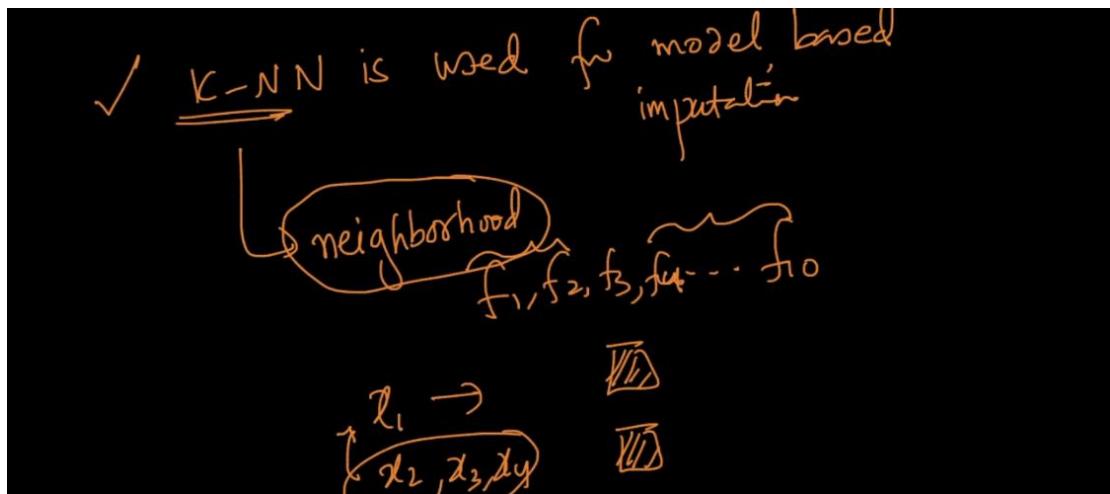
In above method we fill miss value by any standard method but after that we create another feature which keep details of missing value .

If u see above image for  $x_1$   $f_3$  is missing so we store binary value like 0 0 1 0

Means  $f_3$  is missing this record some times helps a lot .



In this method whatever data has missing we make it as test data and other data as a train we implement any algorithm and then that model will predict missing value .



We can use KNN for the same .

Means lets say we 10 features and  $f_3$  is missing then we will find all nearest point to  $x_3$  and if all points  $x_2, x_1, x_4$  is similar to  $x_3$  then it must similar in feature  $f_3$  also so we put value of that features in place of  $f_3$  .

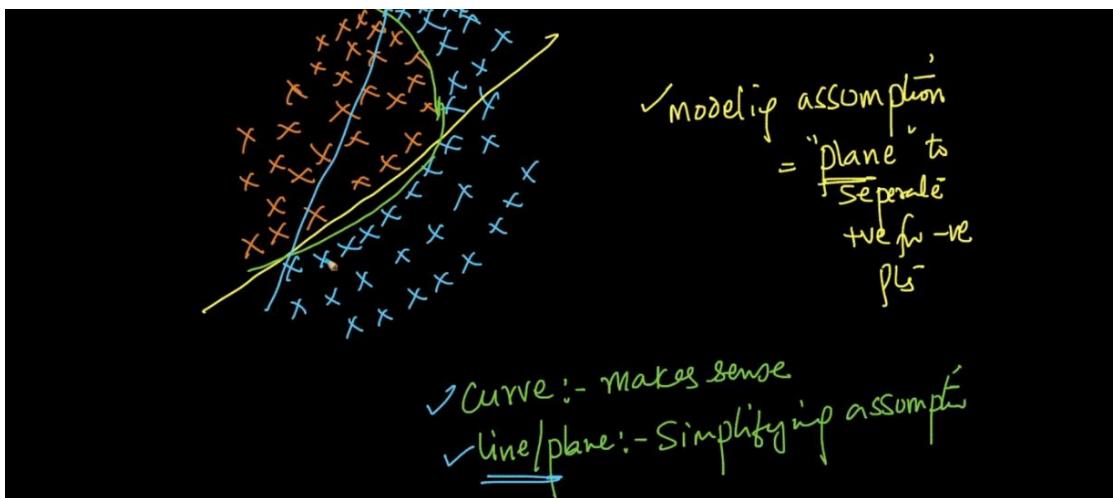
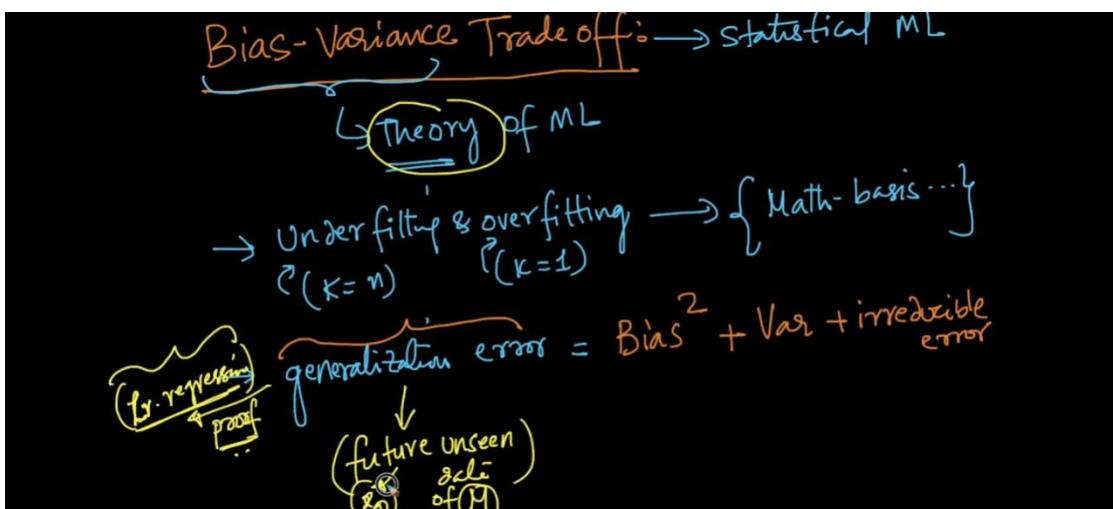
Curse of Dimension :

As dimension (features increase) no of data points also increase very fast to make a good model .

If no of data point are fixed but dimension increase then our model will not perform very well .

Need to watch video again .

## Bias Variance Trade Off :



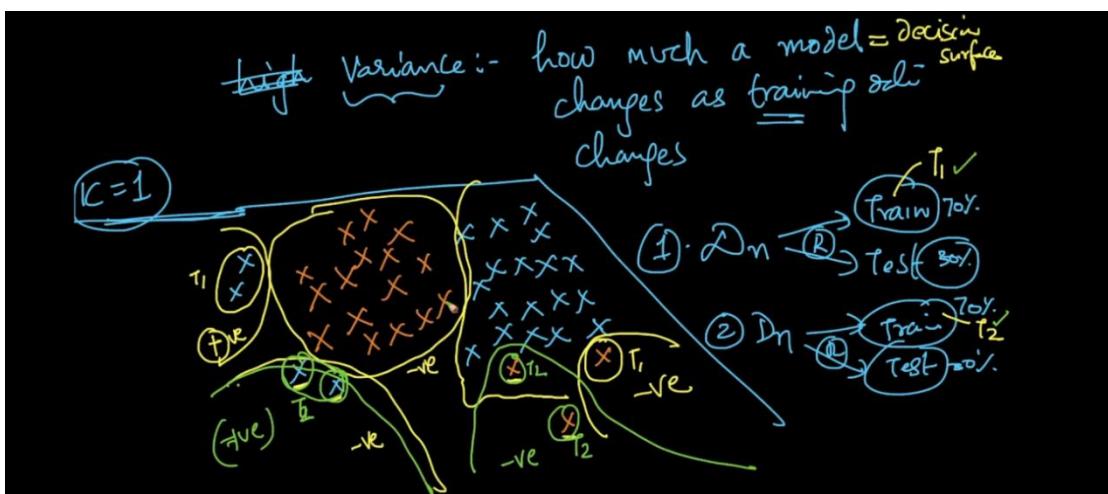
$$\downarrow \text{Gen. error} = \underbrace{\text{Bias}^2}_{\text{errors due to Simplifying assumptions}} + \underbrace{\text{Var}}_{\text{high bias} \Rightarrow \text{Underfitting}} + \underbrace{\text{irreducible error}}_{\text{error that you cannot further fr a given model}}$$

High bias is nothing but under fit model .  
It occurs just taking simplifying assumption .

Lets say we have a data set + and - points now in that 80 % points are - and only 20% points are positive so we blindly predict any new point as - because of its high probability in data . We are not considering our 20% data and this is very worst called under fit .

It will occur when  $k = n$  means we take k value which = no of data points .

High Variance :



This is nothing but over fit problem ..

Also this occurs because of change in training data .

Lets say we have data D we split into train and test randomly .

Now we trained model on that data with  $k = 1$  if you see some outliers - but second time I trained model because of random data split we get different train set so in that same point becomes + this is the effect of High variance .

$$\text{Gen-error} \downarrow = \underbrace{\text{Bias} \downarrow}_{\text{no Underfit}} + \underbrace{\text{Var} \downarrow}_{\text{no Overfit}} + \overbrace{\text{irreducible error}}$$

no Underfit

no Overfit

$K=1 \rightarrow \text{Bias} \downarrow$        $\text{Var} \uparrow$

$K \uparrow \rightarrow \text{Bias} \uparrow \text{slightly}$        $\text{Var} \downarrow \text{dramatically}$

$K \rightarrow n \rightarrow \text{Bias} \uparrow$

$$\text{Gen-error} = \text{Bias}^2 + \text{Var} + \text{Irr-error}$$

Let  $\bar{x} = 1; \mu = 5; n = 10$

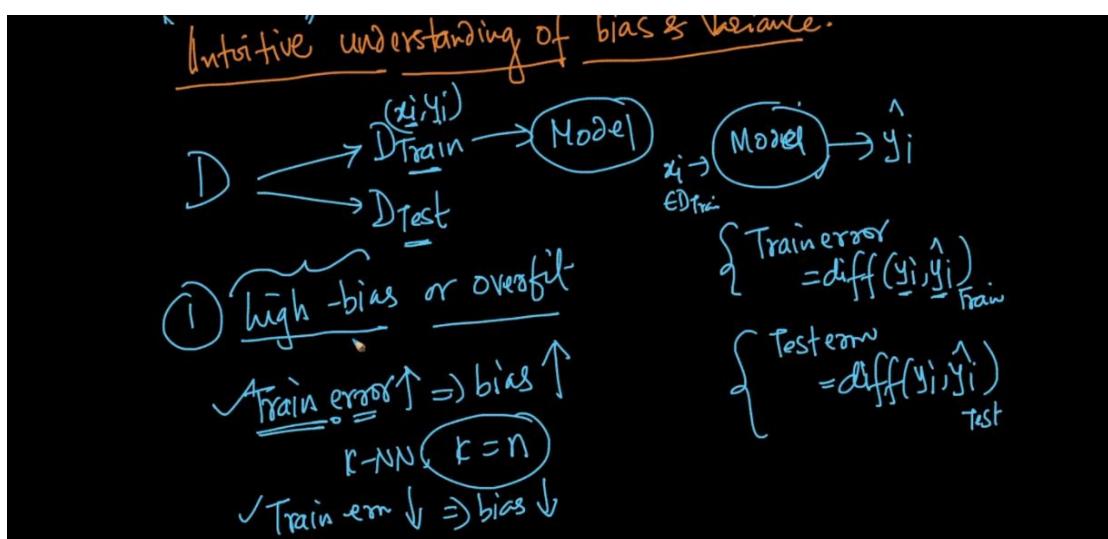
	$10 + 100$	$+ 3 \rightarrow 113 \xrightarrow{\text{overfit}}$
	$12 + 10$	$\rightarrow 25 \checkmark$
	$100 + 2$	$+ 3 \rightarrow 105 \xrightarrow{\text{underfit}}$

In this way we need to select proper K to avoid high bias high variance effect .

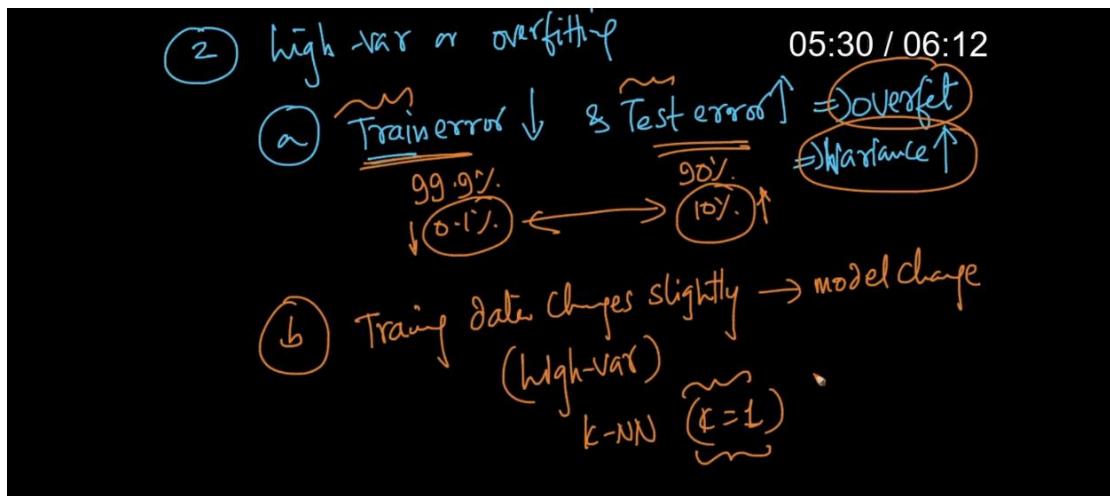
How can we say our model is high bias ?

Very simple if model on train data gives high error means it is high bias .

Means we are not considering 20% of + points just we are predicting each point as - because of its high probability  $k = n$  ..

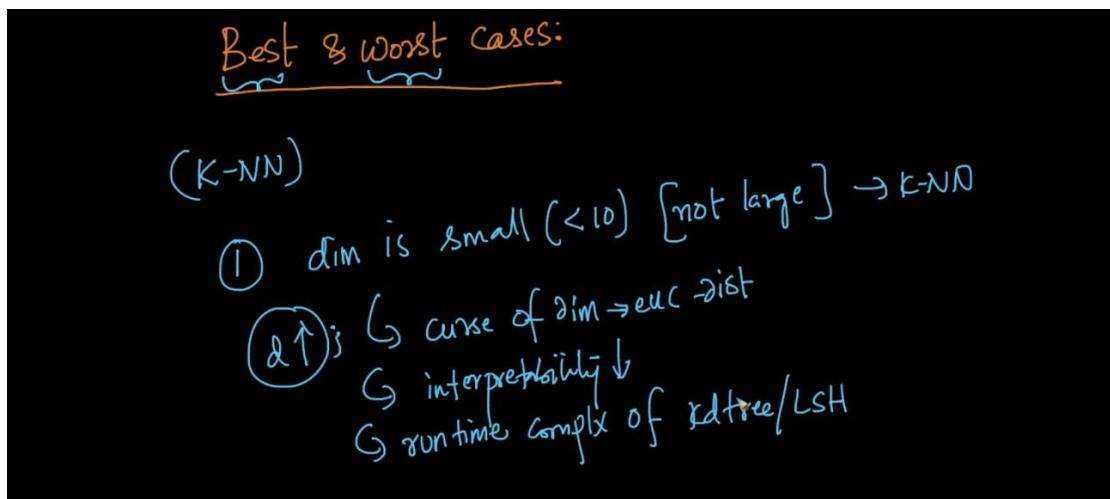


We can say model is high variance when our model train error is very less and test error is very high also called over fit .



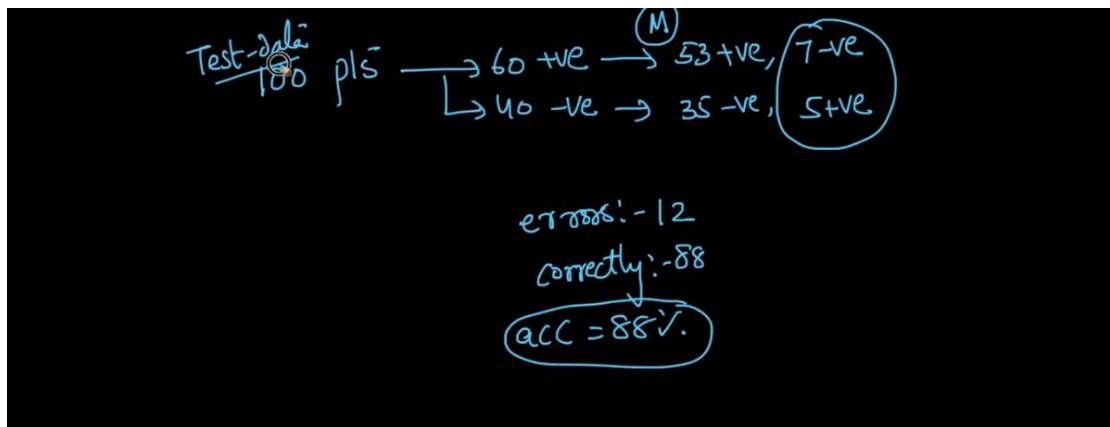
### Best and Worst Case KNN :

When dimension are less then KNN is good algorithm .



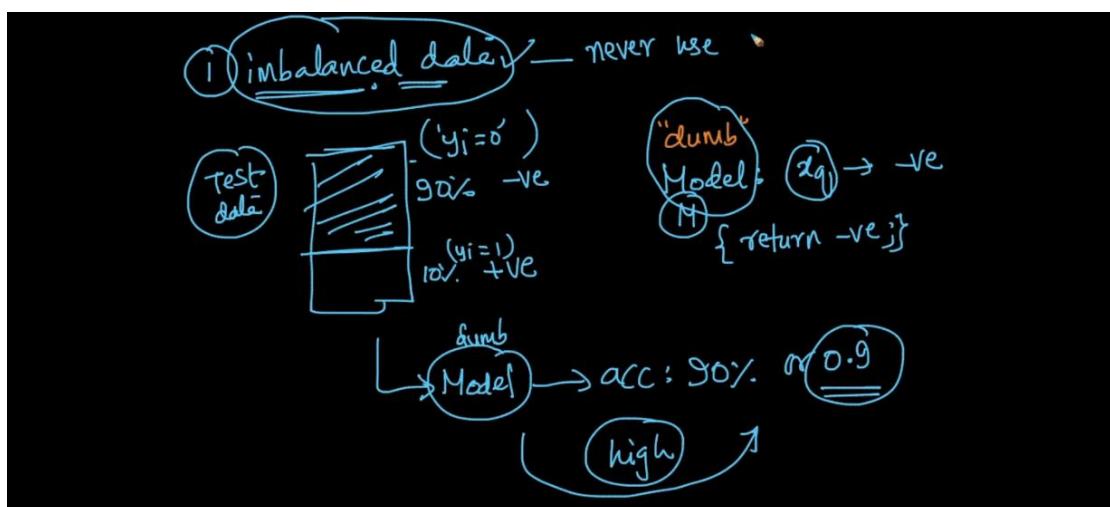
In low latency application KNN not used its time complexity is high .

## Accuracy :



As our model predicted 88 correct points out of 100 accuracy is 88 %.

This is not a good measure when data is not balanced .



(2)

$x_1$	$y$	$M_1$	$M_2$	$\hat{y}_1$	$\hat{y}_2$
$x_1$	1	0.9	0.4	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
$\bar{x}$					

+ve {  
 $\begin{cases} 1, 0 \\ 0.9, 0.4 \end{cases}$

-ve {  
 $\begin{cases} 1, 0 \\ 0.8, 0.65 \\ 0, 0.45 \\ 0, 0.48 \end{cases}$

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

$\hat{y}_1 \rightarrow \text{prob}(y_1 = 1)$   
 $\hat{y}_2 \rightarrow \text{prob}(y_2 = 1)$   
 $(0 \leq p \leq 1)$

$M_1 \sqcup M_2$   
 $1 \text{ or } 0$   
return a prob-score  
K-NN

See  $m_1$  predict low prob for  $x_3$  and  $x_4$  means prob of being 1 is very less and its correct .

So compare both model as per prob score we can easily say  $m_1$  is too better than  $m_2$  but if we use accuracy matric then it will predict accuracy of both model is same hence both models are same but this wrong .

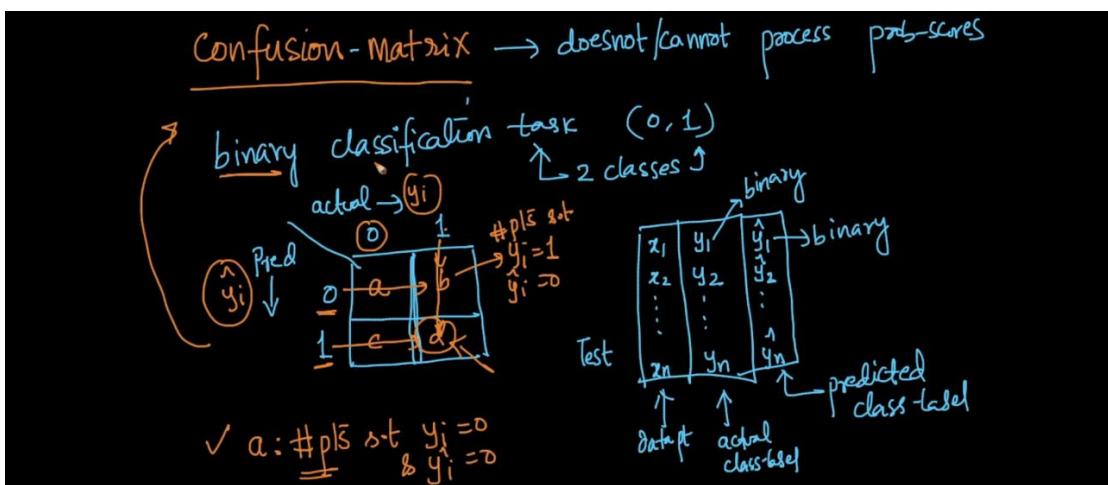
→ predicted class-labels are exactly the same  
 $M_1 \sqcup M_2$

→  $M_1$  is better than  $M_2$   
by looking at prob-scores.

accuracy → cannot use prob-scores

$\hat{y}_1, \hat{y}_2 \rightarrow M_1, M_2 \text{ have same acc}$

## Confusion Matrix :



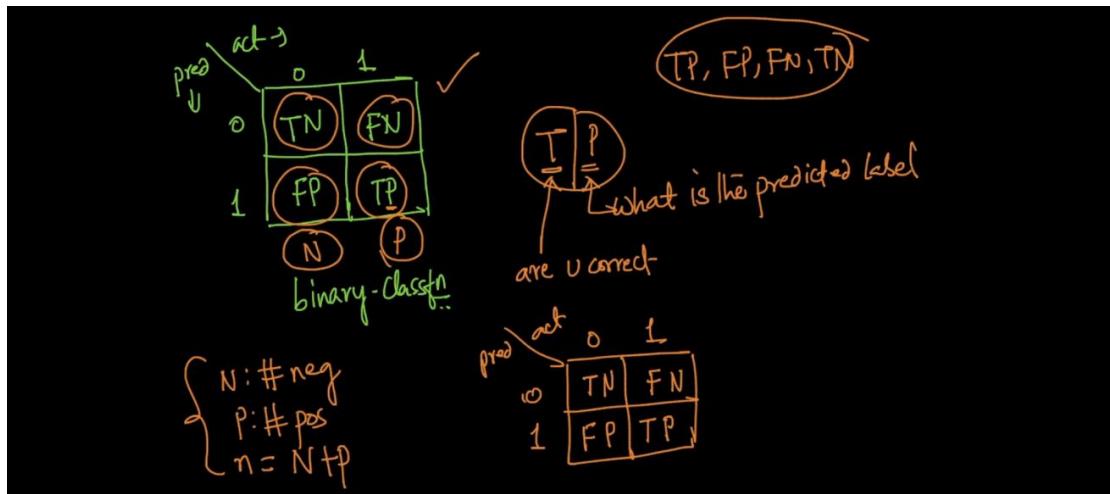
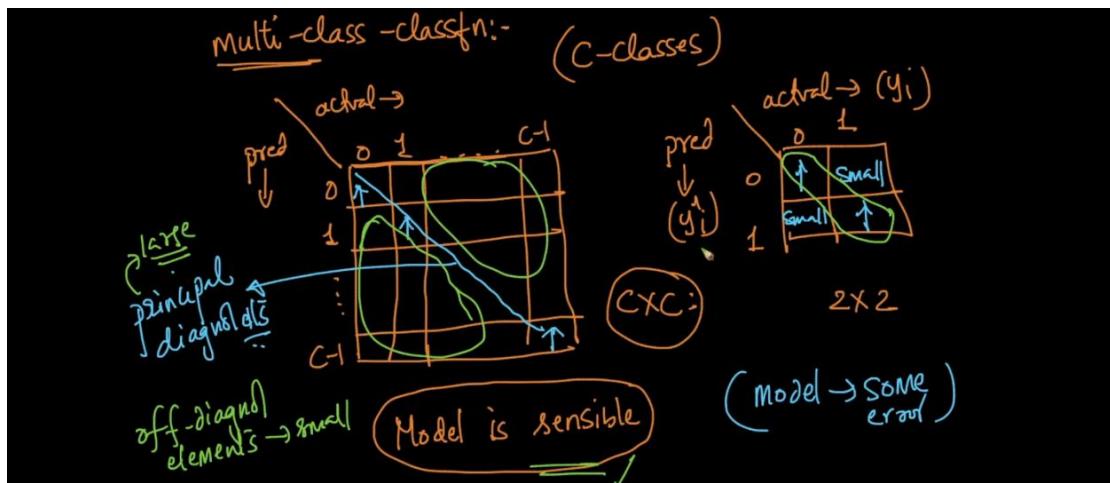
See actual value is 0 and predicted is also 0 means our model did great job . That is True Negative TN .

Actual value is 0 and model predict as 1 means model did a mistake called False Positive FP .

Actual value is 1 and model predict as 0 means model did a mistake called False Negative FN .

Actual value is 1 and model predict as 1 means model did a mistake called True Positive TP ..

So we need TP and TN values always high and FP and FN always low if so then only we can say our model working good .



$$\left\{ \begin{array}{l} \checkmark TPR = \frac{TP}{P} \\ TNR = \frac{TN}{N} \\ FPR = \frac{FP}{N} \\ FNR = \frac{FN}{P} \end{array} \right.$$

act  $\rightarrow$

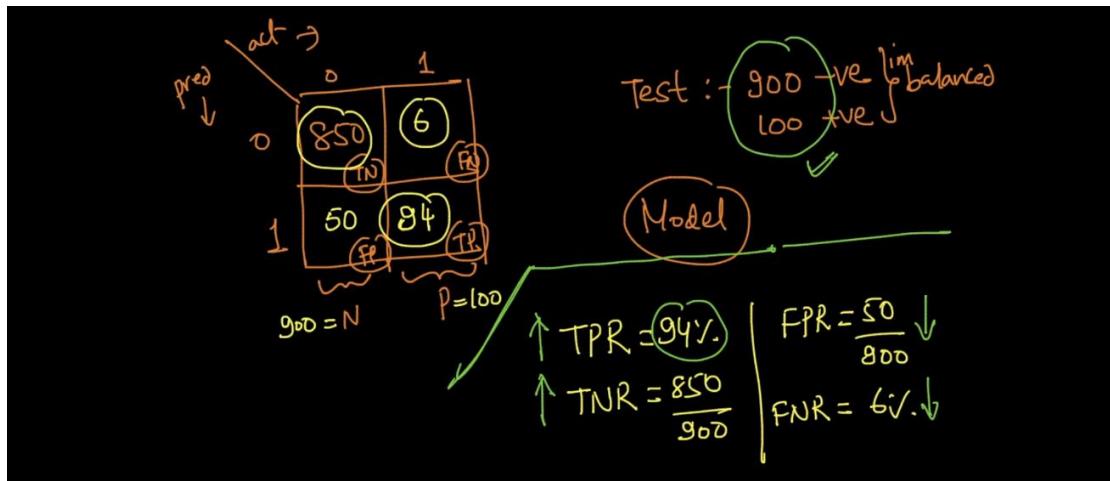
pred  $\downarrow$

TN FN

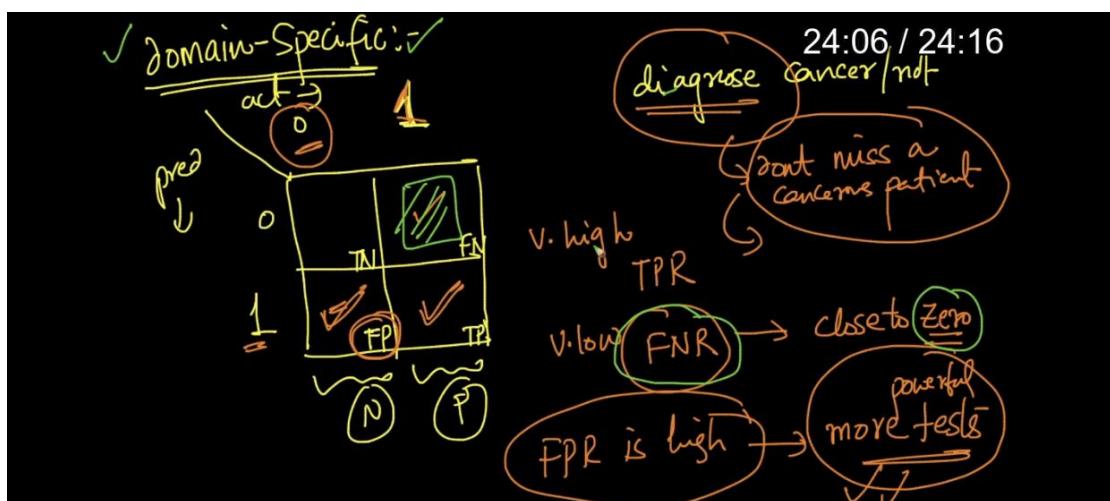
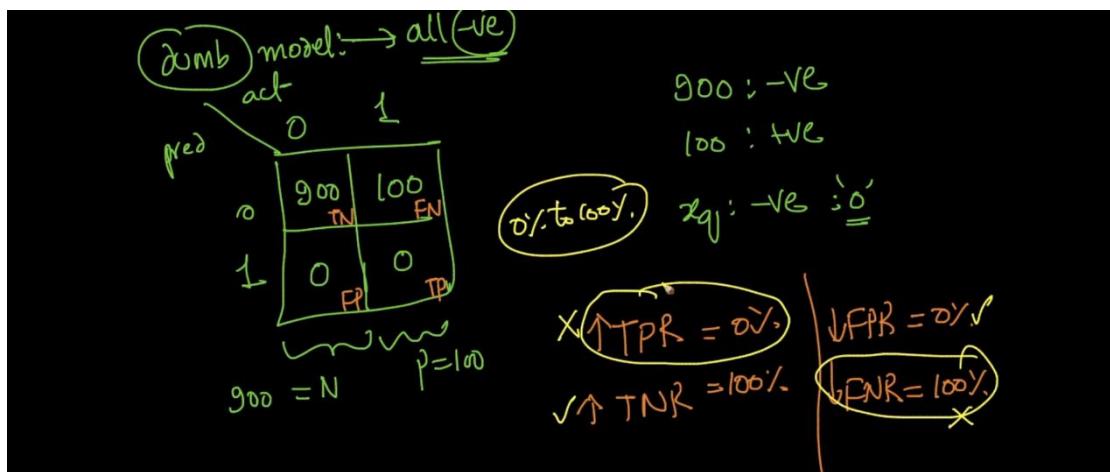
FP TP

N P

$N + P = n$



See below image by just seeing that rates we can say our model working something bad .

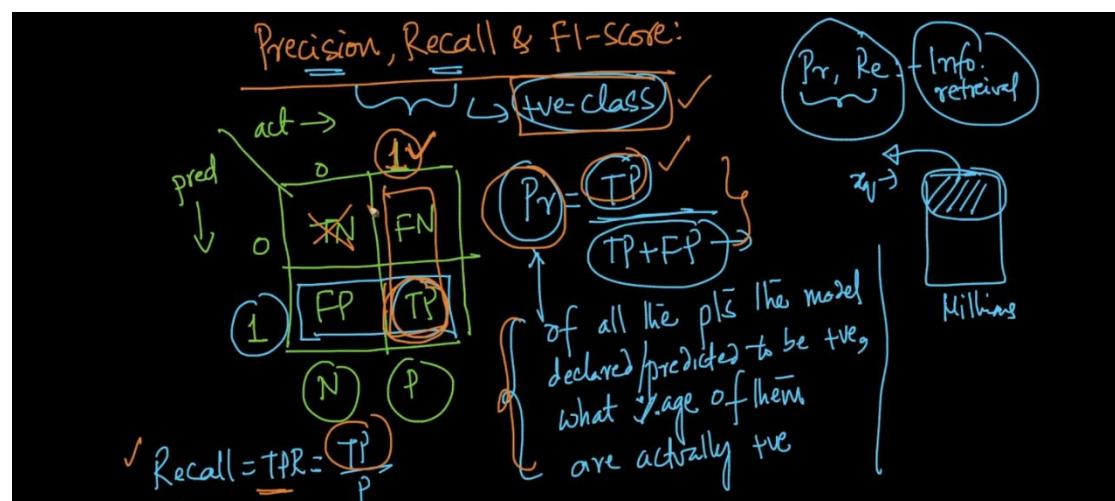


Above example of cancer or not . in this case we need TP rate very high means if patient has cancer then our model will predict yes he has cancer very sure ..

If our model predict wrong in this case then it would be very risky .

In this way confusion matrix helps a lot ..

Precision recall :



Precision or recall only cares about positive points .

Precision means of all points which model predicted as positive ,out of that how much are exactly positive .

Recall means of all actual positive points how much predicted as positive ..

F1 score cal by combing above 2 measures but its very easily understandable but some times it is useful in Kaggle competition .

We want P and R are always high and f1 also high .

Handwritten derivation of the F1-score formula:

Pr, Re → one measure interpretable

$\text{Pr} \uparrow \quad \text{Re} \uparrow$

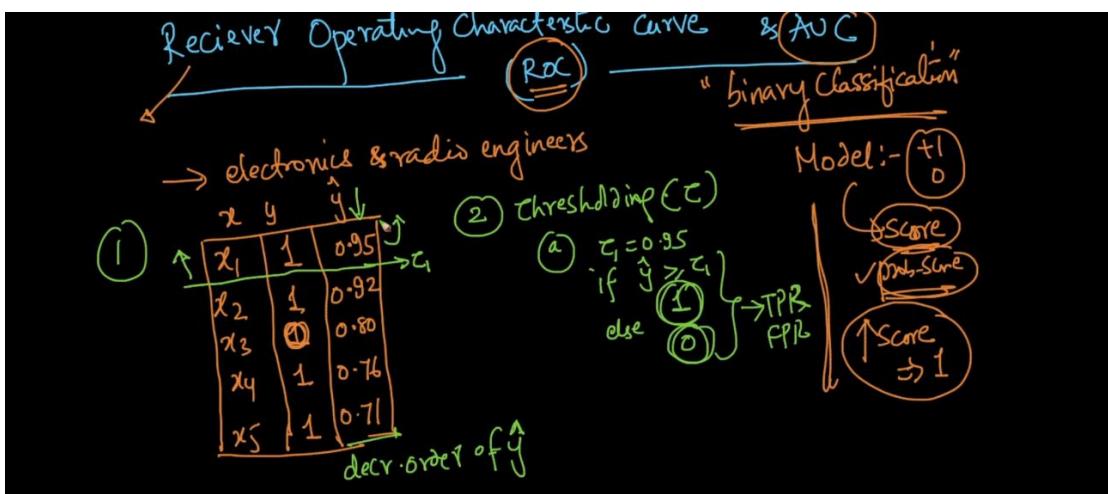
$(\text{Pr} - 1) \quad (\text{Re} - 1)$

$\rightarrow \text{F1-Score} = \left( \frac{2 * \text{Pr} * \text{Re}}{\text{Pr} + \text{Re}} \right)$

$\text{Pr} \uparrow \quad \text{Re} \uparrow$

$\text{F1-Score} \rightarrow 1$

Simple English



Lets say we have data x see above image and output y . so now our model predicted probability score of  $y = 1$  so we get value  $y_1 \dots y_5$  .

First step is we sort them in ascending order and pick first value as thresh hold called as tau 1 .

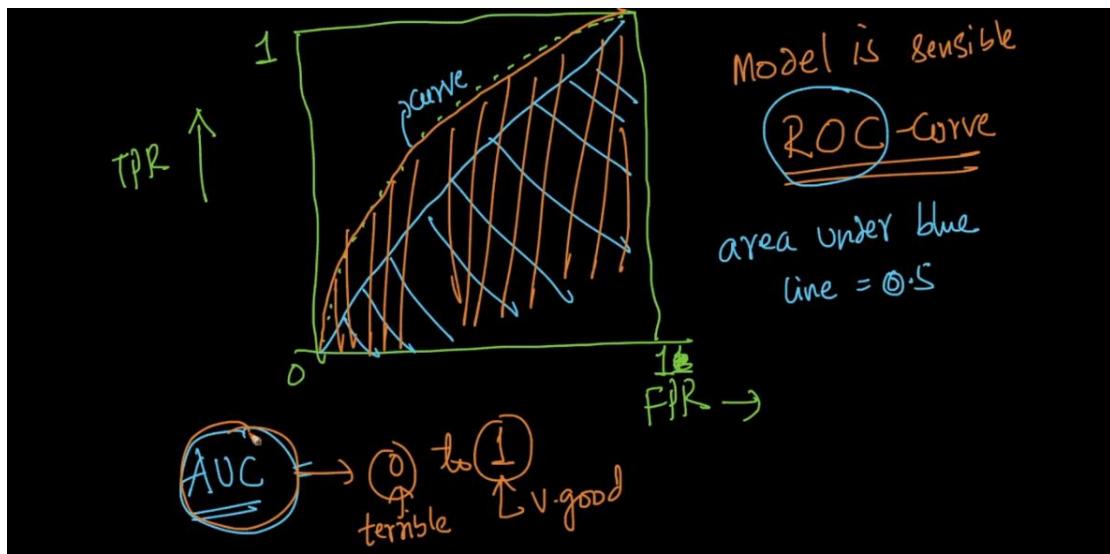
Then we compare  $y^{\wedge}$  with thresh hold if  $y > y^{\wedge}$  then it belong to class 1 .

We repeat same process until n points .

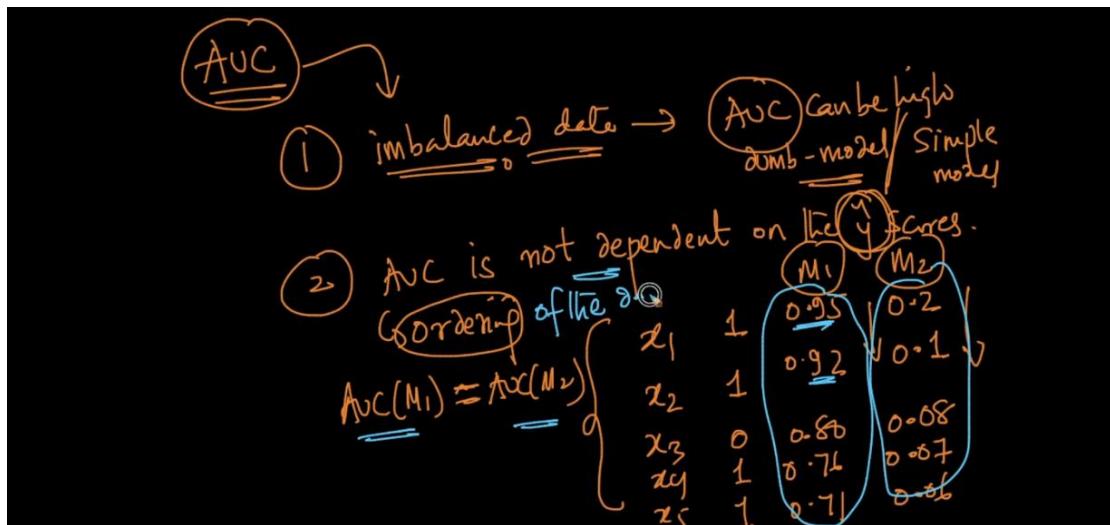
$x$	$y$	$y^{\wedge}$	$y_{T_1=0.95}$	$y_{T_2=0.92}$
$x_1$	1	0.95	1	1
$x_2$	1	0.92	0	0
$x_3$	0	0.80	0	0
$x_4$	1	0.76	0	0
$x_5$	1	0.71	0	0

$\underbrace{y_{T_1=0.95} \rightarrow TPR, FPL(T_1)}$   
 $\underbrace{y_{T_2=0.92} \rightarrow TPR, FPL(T_2)}$

For each  $\tau_1 \dots \tau_n$  we find TPR and FPR and plot a graph its called ROC .



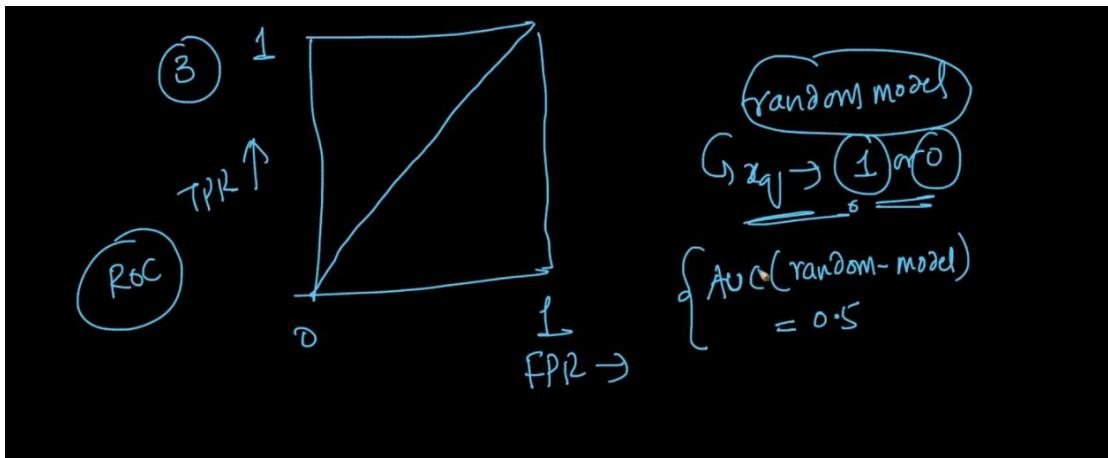
## Properties of AUC :



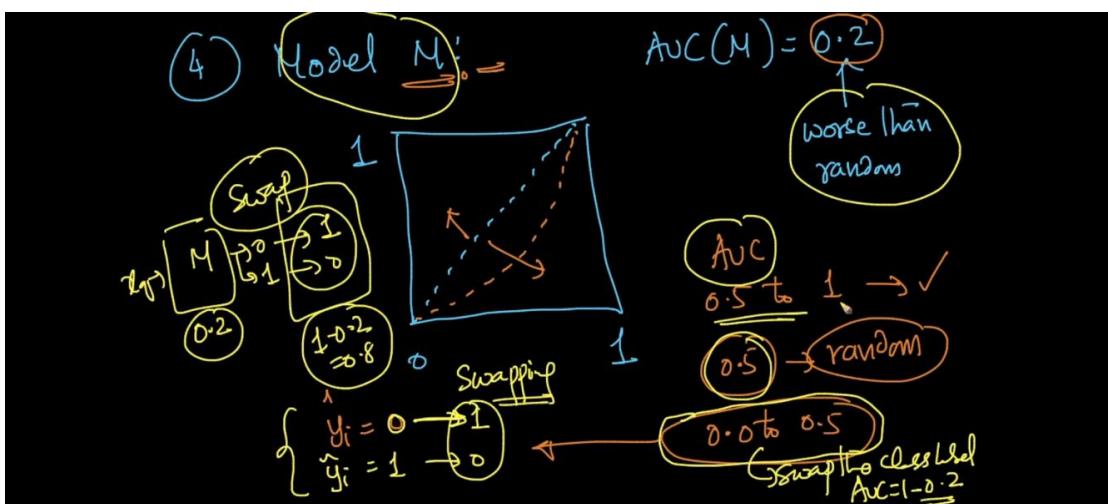
For not balanced data set AUC is always high .

AUC is not depend on probability score its just depend on ordering of score see above both model have same descending order then AUC for both will be same .

For random model AUC is 0.5



Means given any new point I am randomly telling its 0 or 1 for that AUC will be 0.5



If any model has AUC below 0.5 then just swipe output . means if it predict 0 then change it to 1 vice v .

## Log Loss :

<p><u>as small</u></p> <p><u>Binary-classfn:</u></p> <p>penalizing for small deviations in prob-score</p> <p>Test set of n-pls :-</p>	<p><u>Log-loss :-</u> prob-scores</p> <p><math>y_i = p_i</math></p>	<p>model</p>															
		<table border="1"> <thead> <tr> <th>x<sub>i</sub></th> <th>y<sub>i</sub></th> <th><math>\hat{y}_i = p_i</math></th> </tr> </thead> <tbody> <tr> <td>x<sub>1</sub></td> <td>1</td> <td>0.9</td> </tr> <tr> <td>x<sub>2</sub></td> <td>1</td> <td>0.6</td> </tr> <tr> <td>x<sub>3</sub></td> <td>0</td> <td>0.1</td> </tr> <tr> <td>x<sub>4</sub></td> <td>0</td> <td>0.4</td> </tr> </tbody> </table>	x <sub>i</sub>	y <sub>i</sub>	$\hat{y}_i = p_i$	x <sub>1</sub>	1	0.9	x <sub>2</sub>	1	0.6	x <sub>3</sub>	0	0.1	x <sub>4</sub>	0	0.4
x <sub>i</sub>	y <sub>i</sub>	$\hat{y}_i = p_i$															
x <sub>1</sub>	1	0.9															
x <sub>2</sub>	1	0.6															
x <sub>3</sub>	0	0.1															
x <sub>4</sub>	0	0.4															

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

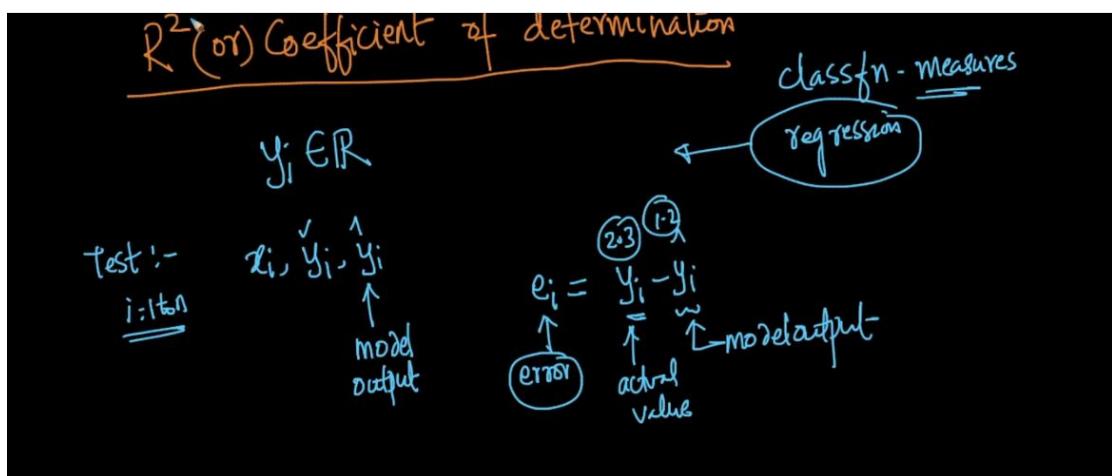
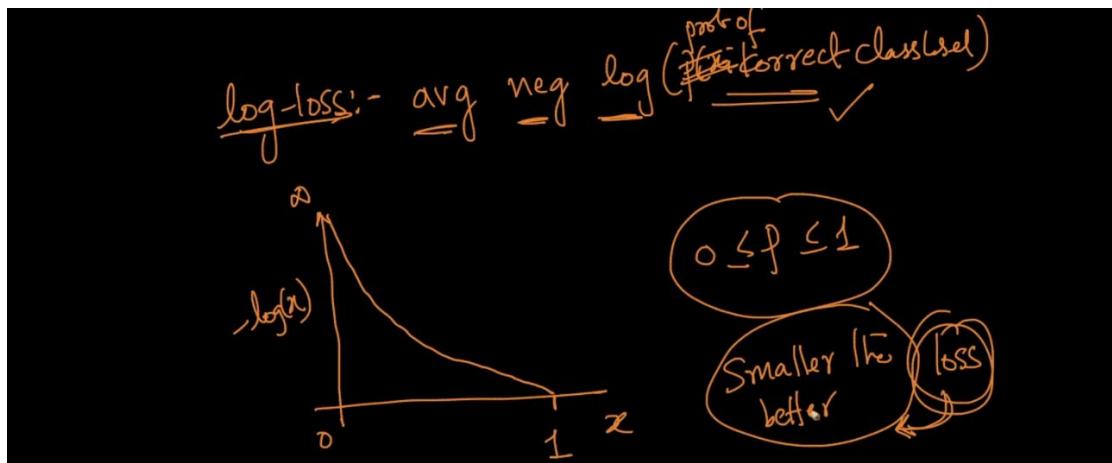
See above formula log loss use exactly all p values for  $y^{\wedge}$ .

Now see above table as error increases log loss value also start to increase .

Log loss can lie between 0 to inf .

Log loss close to 0 is always best means in that case there is very less loss or error in our model .

Log loss is not to easy interpretable that's the only problem now lets say some one say log loss of my model is 1 or 10 what can we say its very difficult because it lies between 0 to inf .. and best case is 0 .



Very simple model we can design using regression is mean model means for each new point we predict mean value of y col .

$$\text{total SS} \quad \left\{ \begin{array}{l} \text{sum of squares} \\ \sum_{i=1}^n (y_i - \bar{y})^2 \end{array} \right.$$

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

regression  $\rightarrow$  Simplest model  $\rightarrow$  avg-model  
 $x_{q1} \rightarrow \text{mean}(y_i) = \bar{y}$   
 Train. data  $\left\{ \begin{array}{l} \text{predict height (EIR)} \\ \text{features: w, c, hc, ...} \end{array} \right.$   
 avg = 152cm  $\rightarrow$   $\bar{y} \rightarrow 152\text{cm}$   
 $x_{q1} \rightarrow 152\text{cm}$   
 $y_q \rightarrow 152\text{cm}$

$\bar{Y}$  bar is mean of y col .

Sum if residual this is another term here we take difference between each y with predicted  $\hat{y}$  .

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

residue =  $e_i = \underline{\underline{y_i}} - \underline{\underline{\hat{y}_i}}$   
 actual  $\swarrow$  predicted  $\searrow$

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

Case 3  
Model:  $SS_{res}$   
↳ same as a simple-mean-model

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

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

Case 3:  $SS_{res} = SS_{tot}; R^2 = 1 - 1 = 0 \rightarrow \text{Model 1}$

R2 close to 1 means model is best R2 close to 0 means model is worst .

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

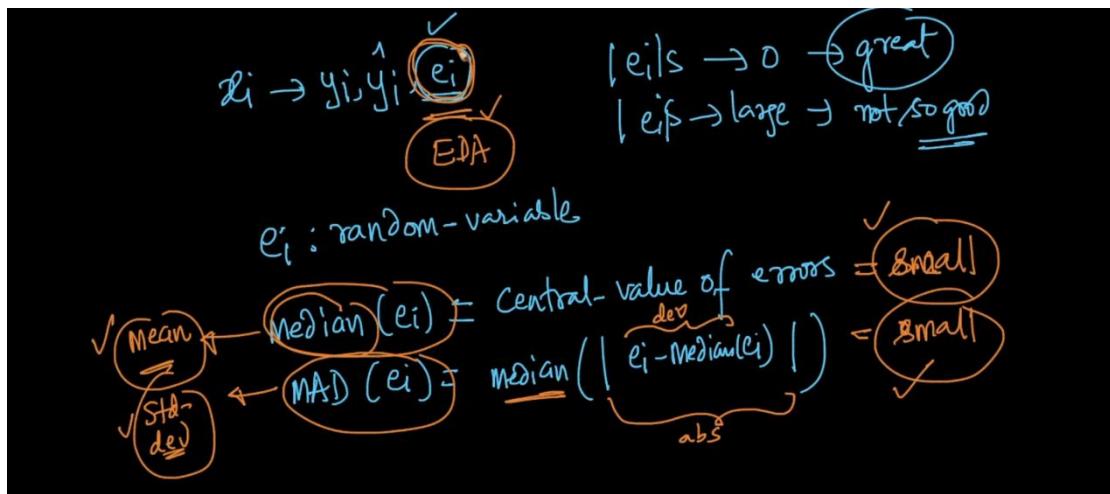
$$R^2 = 1 - (q > 1) = (-ve)$$

Model is worse than a simple-  
Model

Q

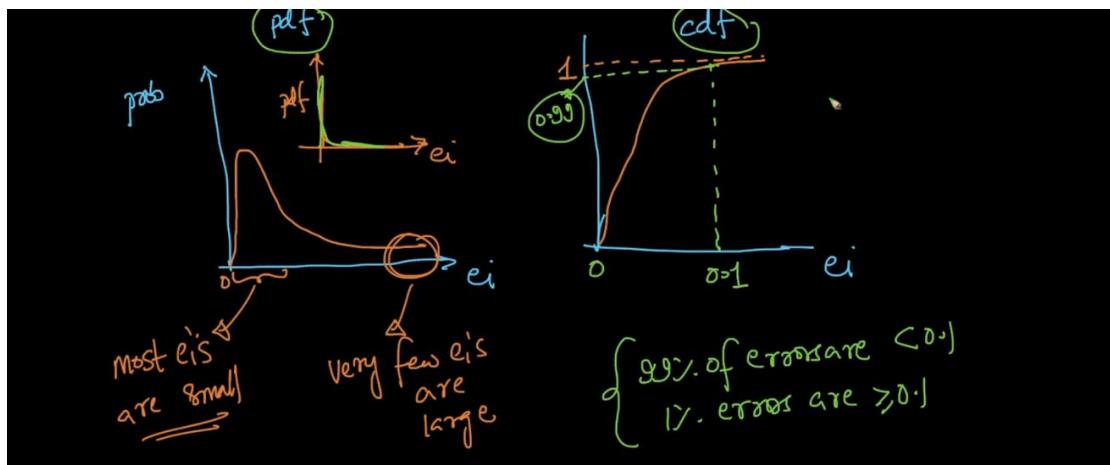
Now r2 is not robust to outliers. if lets say there is 1 outlier in data so only for that point we get very high error value and just because of that whole r2 vary .

So we use Median absolute deviation (MAD) for each error .



As we know that median and MAD not affected by outliers.

By using PDF and CDF we can very easily understand model performance on error .



From PDF we say that max points has error is low and very few points have error is high .

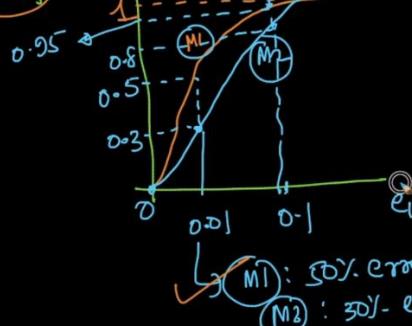
From CDF we can say 99% of value has error  $< 0.1$  and only 1 % have  $> 0.1$  error .

06:33 / 06:42

Models:  $M_1$  &  $M_2$



cdf



$M_2$  cdf is below  $M_1$

$M_1$ : 95% errors are below 0.1

$M_2$ : 80% errors below 0.1

$M_1$ : 50% errors are below 0.0

$M_2$ : 30% errors are below 0.0