

CLASSIFICATION ALGORITHM IN VARIOUS CASES

[INTRODUCTION](#)

[IMBALANCED v/s BALANCED DATASET](#)

[k-NN, GIVEN A DISTANCE MATRIX OR SIMILARITY MATRIX](#)

[TRAIN AND TEST SET DIFFERENCES](#)

[IMPACT OF OUTLIERS](#)

[LOCAL OUTLIER FACTOR \(MEAN DISTANCE TO KNN\)](#)

[K - DISTANCE & REACHABILITY DISTANCE](#)

[LOCAL REACHABILITY DENSITY](#)

[LOCAL OUTLIER FACTOR](#)

[IMPACT OF COLUMN & SCALE STANDARDIZATION](#)

[INTERPRETABILITY](#)

[FEATURE IMPORTANCE AND FORWARD FEATURE SELECTION](#)

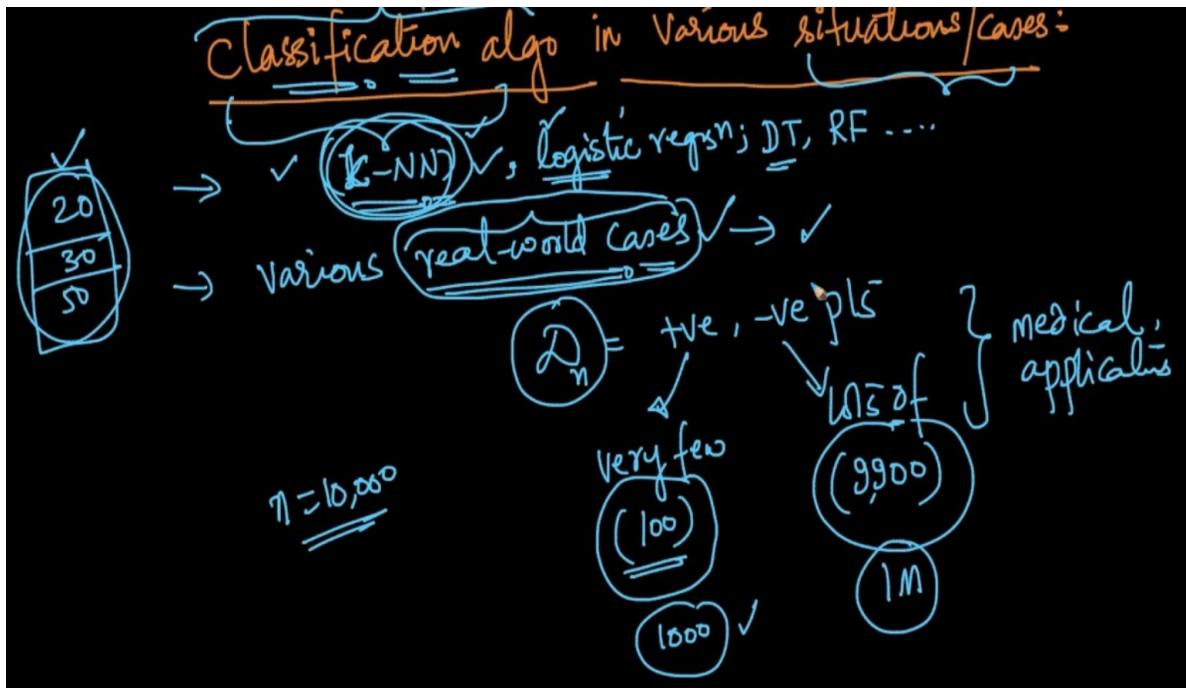
[HANDLING CATEGORICAL AND NUMERICAL VALUES](#)

[HANDLING MISSING VALUES BY IMPUTATION](#)

[CURSE OF DIMENSIONALITY](#)

CLASSIFICATION ALGORITHM IN VARIOUS CASES

INTRODUCTION

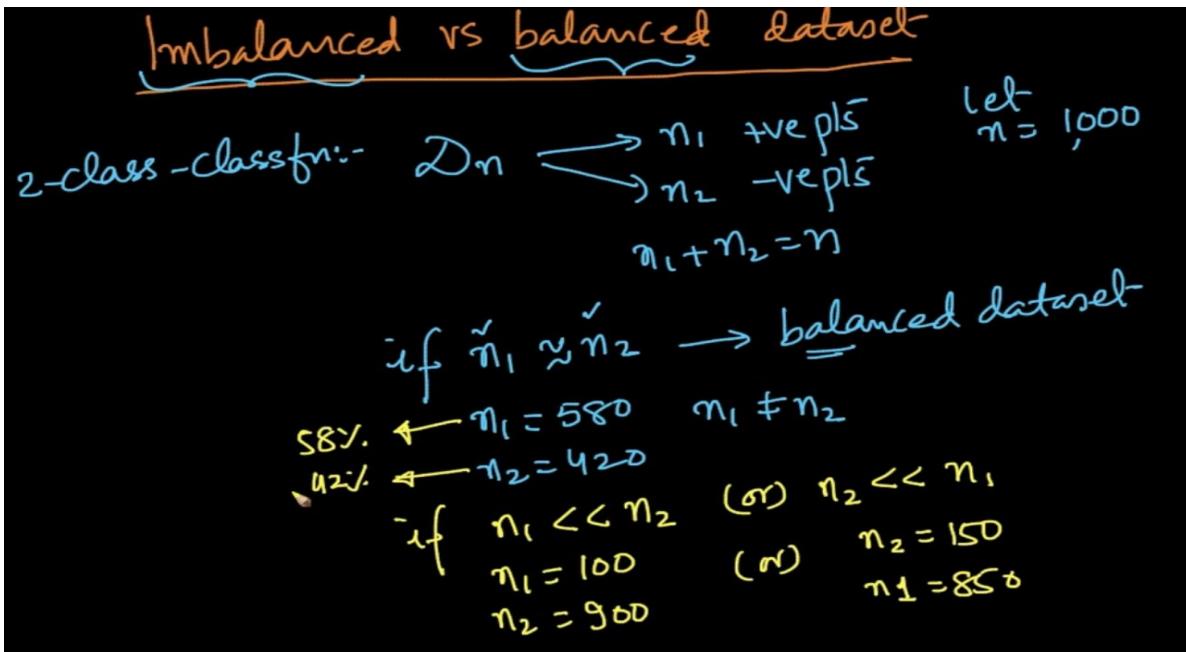


We'll know various classification Algorithms like k-NN, Logistic Regression, SVM etc but if we don't know how to apply that knowledge in different cases then it's not worth it.

Here, suppose we've medical dataset. In that there are only a few which are diagnosed with diseases from the total dataset. So we need to know how to modify our algorithm for this case

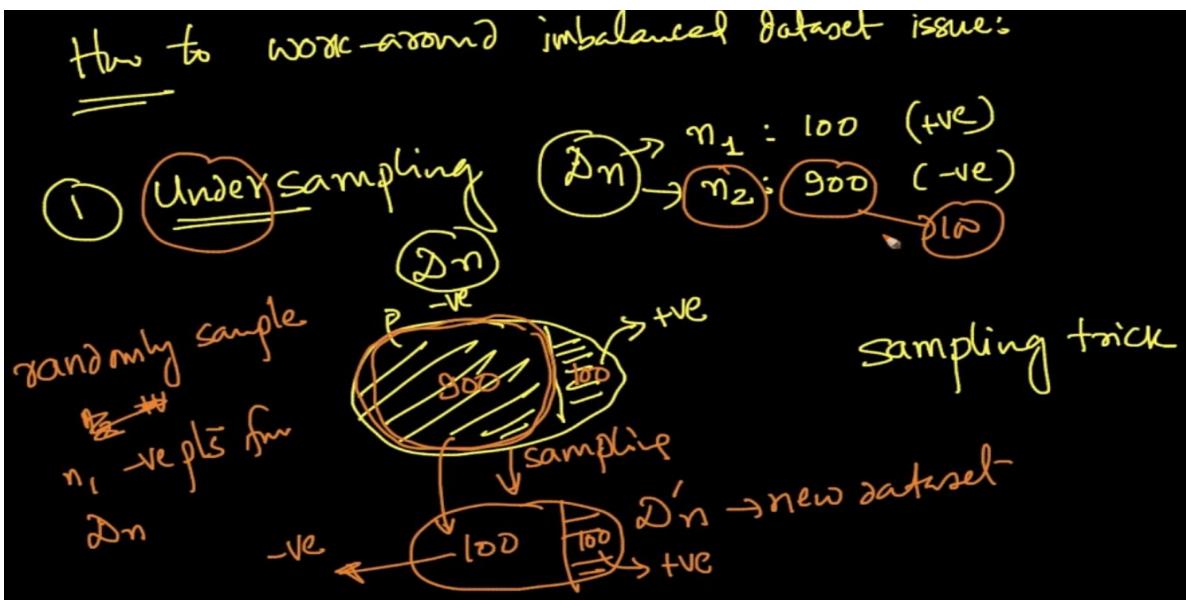
Knowing how to apply Algorithms in various cases is far more important then knowing the number of tricks/Methods you know

IMBALANCED v/s BALANCED DATASET

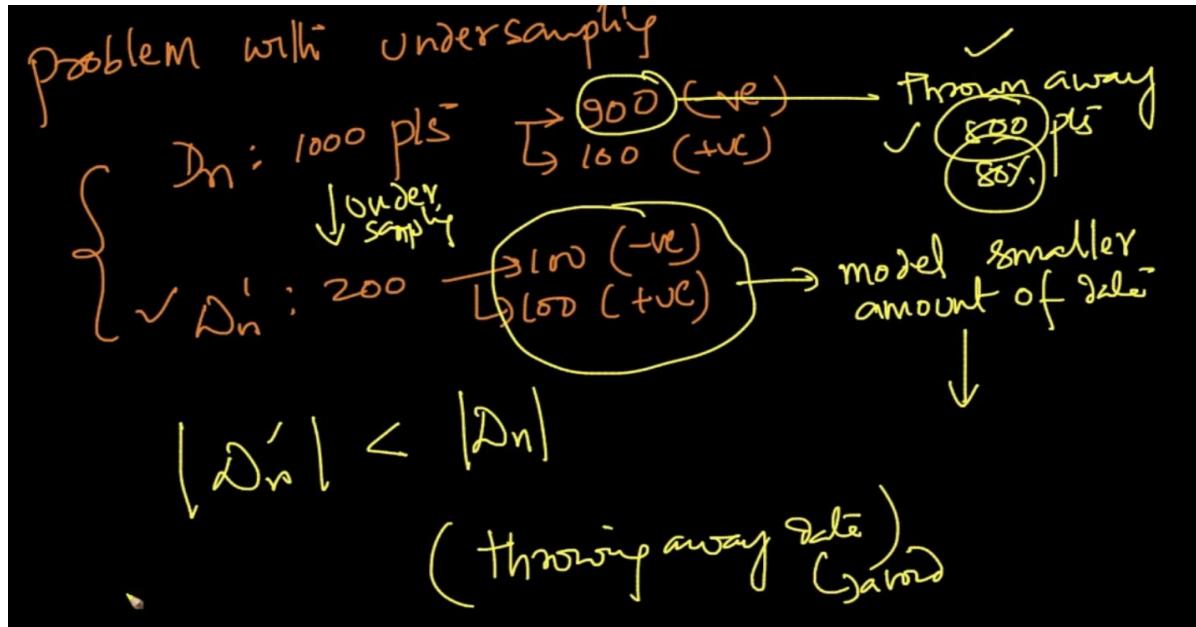


If there's a huge difference number of different class labels then it's imbalanced dataset

Here, $n = 1000$, $n_1 = 900$ (-ve pts), $n_2 = 100$ (+ve pts) so it's an imbalanced dataset

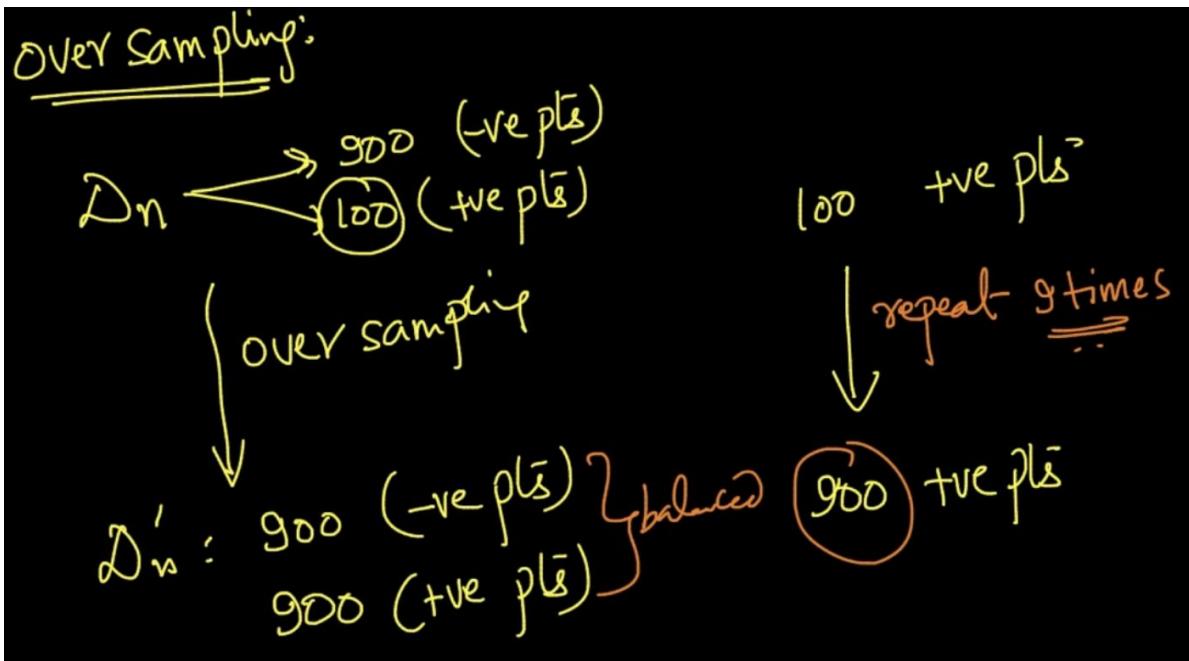


This is undersampling (1 method to handle imbalanced data) . We keep the minority class (n1) as it is and randomly take points as same as n1 (n2 is 900 we take 100 pts from that)



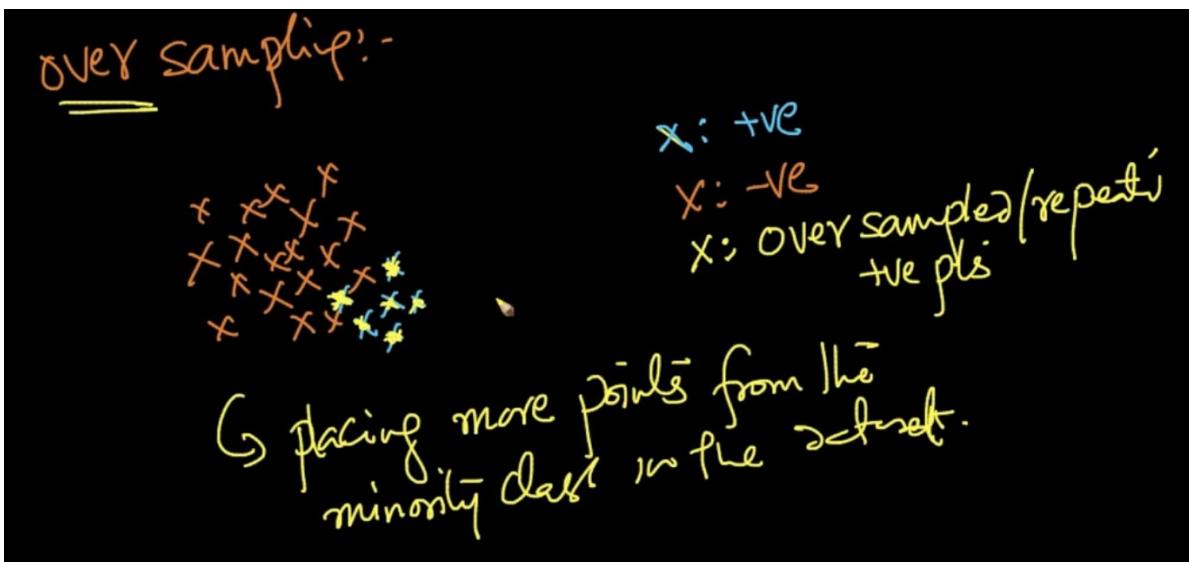
Problem with undersampling is we are throwing away data (80% here) which should not happen

OVERSAMPLING



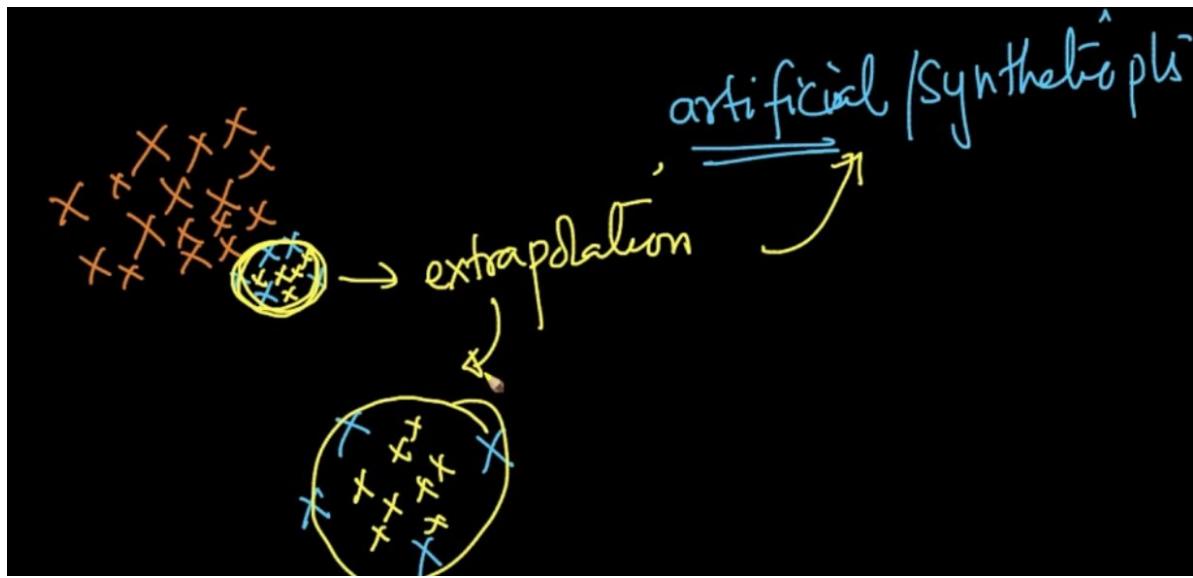
In this method, the minority class is repeated to level with majority class. Here, $n_1 = 100$

So we repeat it 9 times to balance it with majority class ($n_2 = 900$). Now both balanced



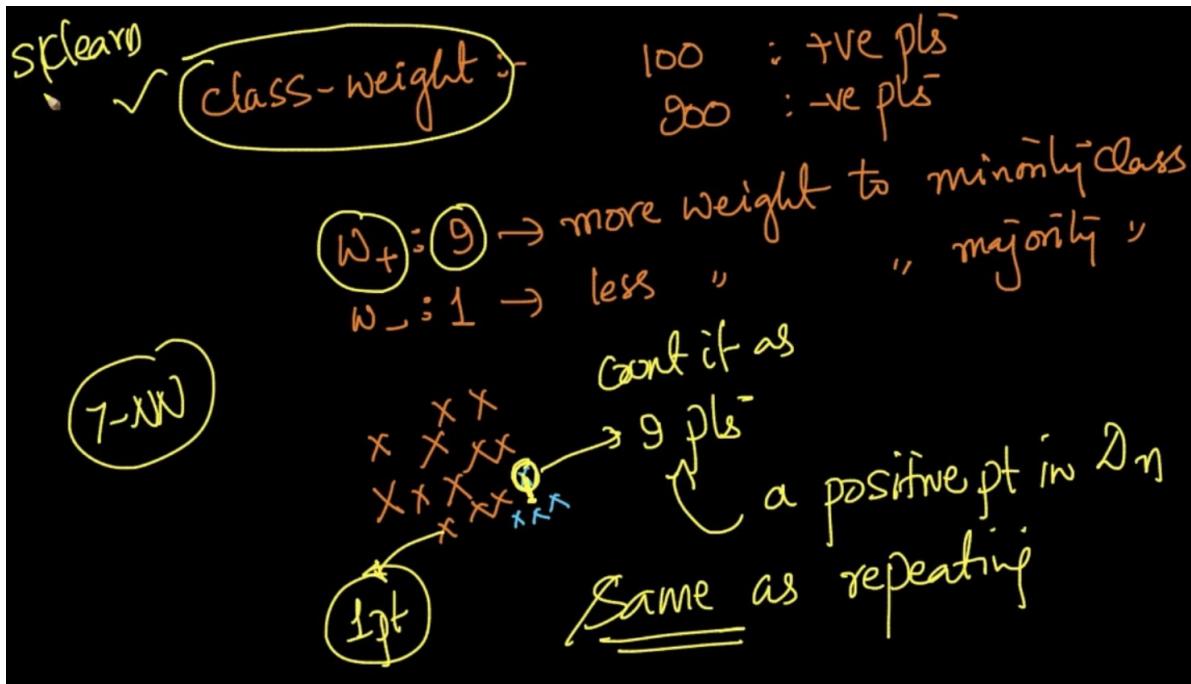
Here in the minority class, we are placing the same points 9 times on top of minority class points

EXTRAPOLATION/ ARTIFICIAL POINTS CREATION



For the minority class region we are artificially creating points on it to balance the dataset

CLASS WEIGHT



More weight is being given to minority class . So, if we are using a pt from minority class We count it as 9 point. So if we are using 7-NN and a point is near it all the nearest neighbors will be from that point. Same as repeating

MULTI CLASS CLASSIFICATION

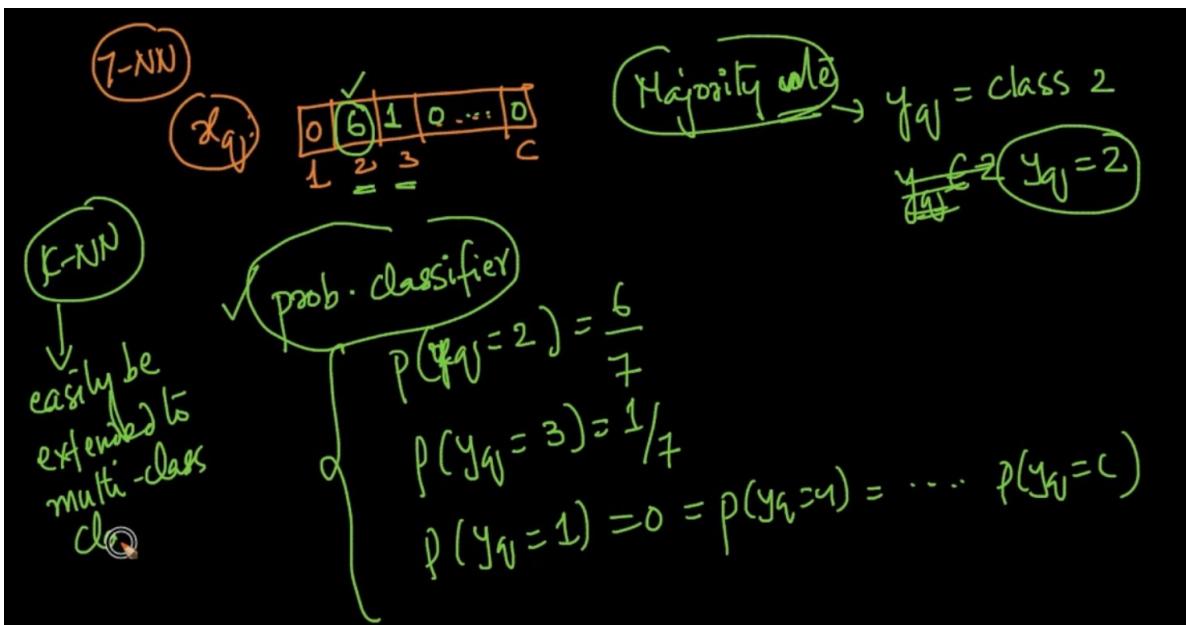
Multi-class classifier:

✓ binary classifier : $y_i \in \{0, 1\}$

✓ Multi-class-classifier: - $y_i \in \{0, 1, 2, \dots, 9\}$ → MNIST
 (10-classes) \square

$$KNN = \mathcal{D}_n = \left\{ (x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \overbrace{\{1, 2, 3, \dots\}}^{C-class} \right\}$$

There's a c number of classes in our dataset. How does k-NN work in that?



If we've x_q and running 7-NN and we are getting 6 values in class 2 and 1 in class 3

Then by majority vote $y_q = 2$ or if probability classifier then $P(x_q = 2) = 6/7$. Same can be done for other classes as well. So k-NN can easily be extended to multi-class classification

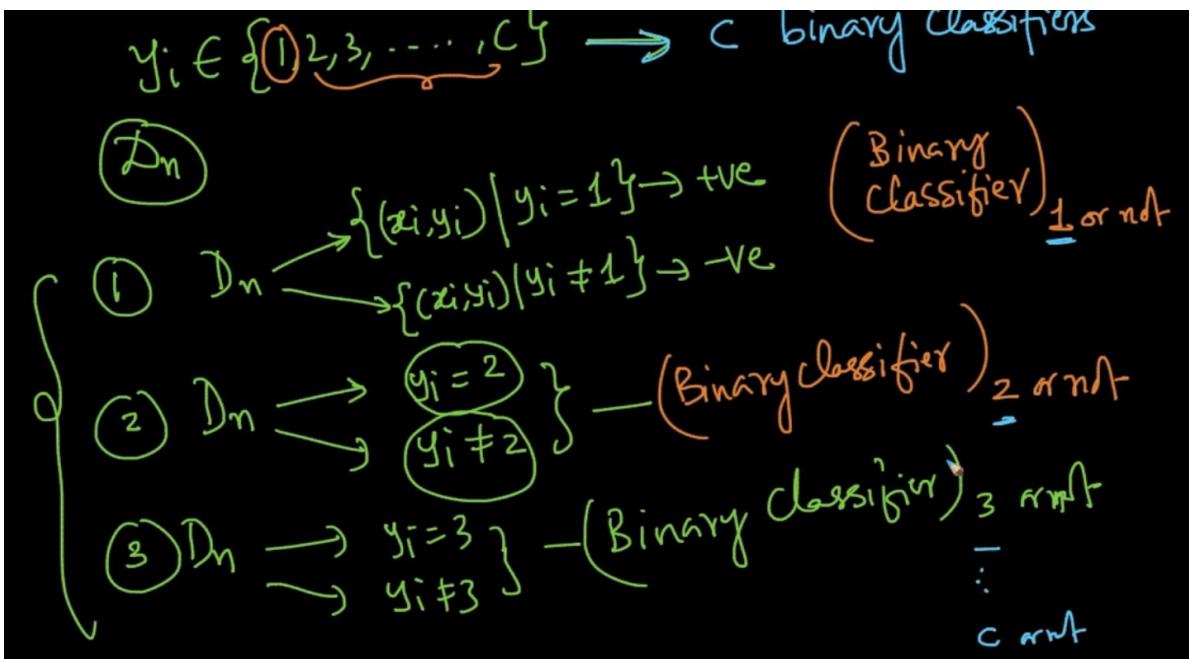
✓ Logistic regresn $\xrightarrow{\text{binary classifier}}$ $\xrightarrow[\text{cannot go}]{} \text{multiclass-classfn. as easily as k-NN}$

(*) Given a multi-class classfn prob;
can we convert it into a binary classfn. prob

$$f(x) \xrightarrow{0} \text{binary}$$

$$f'(x) \begin{cases} 0 \\ 1 \\ 2 \\ \vdots \\ c-1 \end{cases} \left\{ \begin{array}{l} \text{c classes} \\ \text{Binary classifier} \end{array} \right.$$

We can't do multi class classification on Logistic regression . So can we do anything to convert multiclass to a binary class $f'(x) \rightarrow f(x)$



Here we are converting multi class into binary classification by breaking the dataset into 2 parts . For e.g : $y_i \in \{1, 2, 3, 4, \dots, c\}$. D_n is broken into two parts 1st is $y_i = 1$,2nd is $y_i \neq 1$ Same for other classes . This is called 1 vs Rest method. We are training c- models here as No. of classes are c

k-NN, GIVEN A DISTANCE MATRIX OR SIMILARITY MATRIX

k-NN, given a distance measure

Classfn :- $x_i \in \mathbb{R}^d$ → don't get x_i 's as vec^ts
 ↳ vector is not easy

y_i ← x_i → useful to treatment

x_i → Chemical compound in Pharma → numeric vect^r (may not be easy)
 ↓ Paracetamol → [] [] [] []

$\text{Sim } (x_i, x_j) = \begin{cases} \text{Pharmacist} \\ \text{Chemist} \end{cases}$

Sometimes, x cannot be converted into a numeric vector. For e.g a chemical compound but similarity between 2 chemical compounds can be calculated

$\text{Sim } (x_i, x_j) =$

$S = x_i \begin{array}{|c|c|c|c|} \hline & x_j & & \\ \hline & & & \\ \hline \end{array} n \times n$

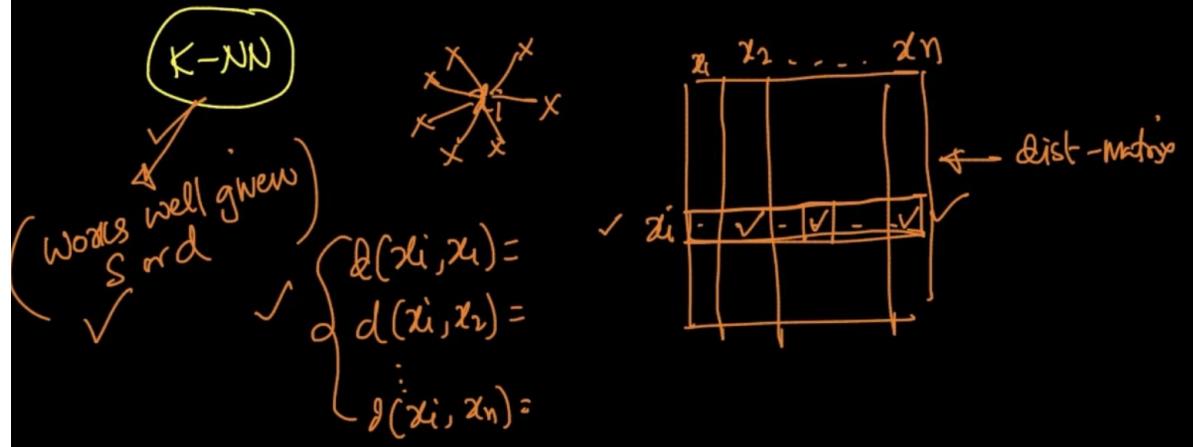
n data-points
 ↳ 10K

$S_{ij} = \text{Sim } (x_i, x_j)$

$\text{Dist} : d_{ij} = \frac{1}{S_{ij}}$

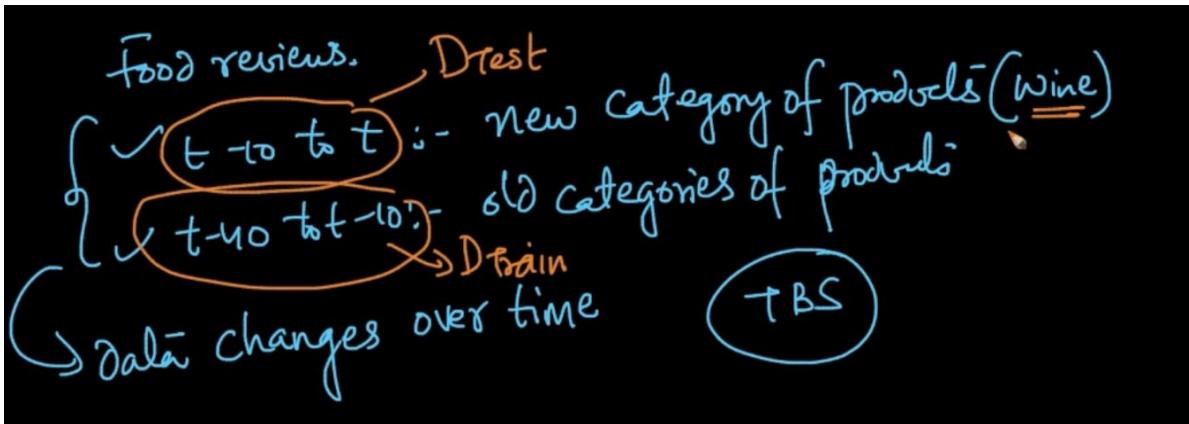
A similarity matrix S of all the chemical compound is created. S_{ij} tells us about the similarity between two chemicals x_i, x_j . Distance matrix can also be created.

Given Sim-matrix (Σ) or dist-matrix (d)
✓ instead of $x_i \in \mathbb{R}^d$ explicitly \leftarrow Logistic Regn

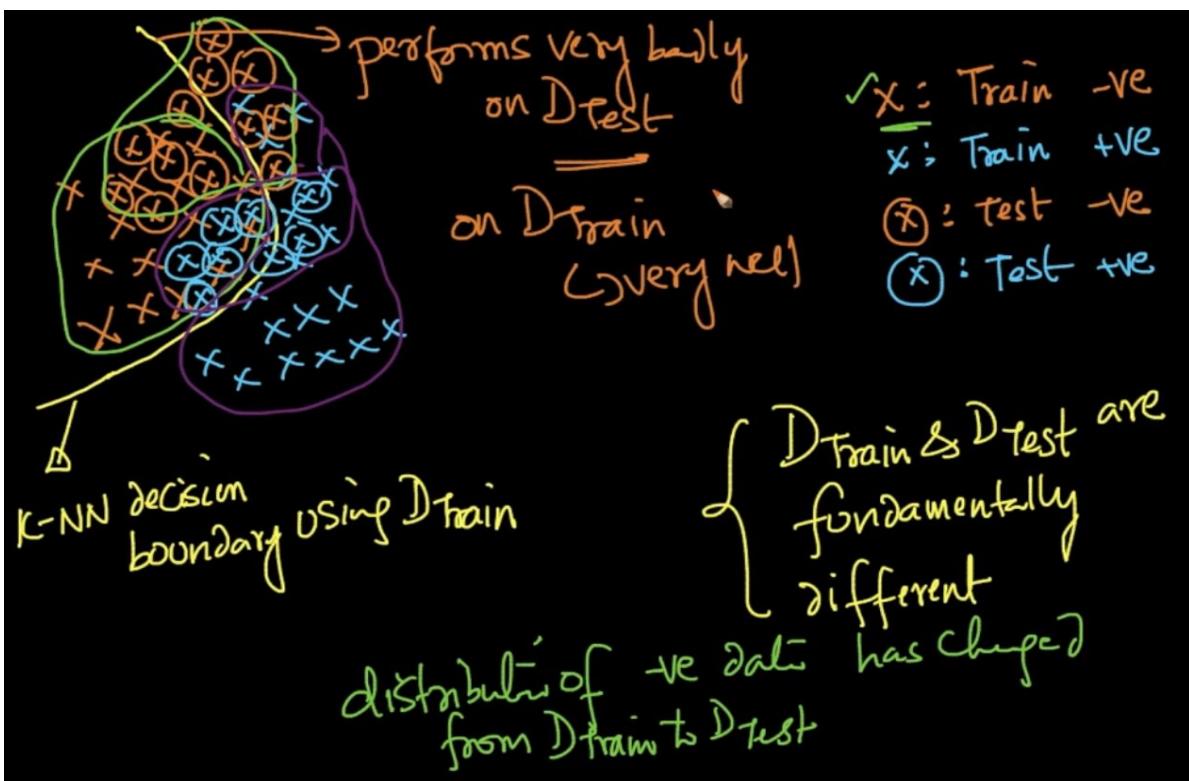


Given distance matrix for all x . We get the 3 nearest neighbors if $k=3$ by using distance matrix. So, k-NN can work awesomely with Similarity matrix or distance matrix

TRAIN AND TEST SET DIFFERENCES

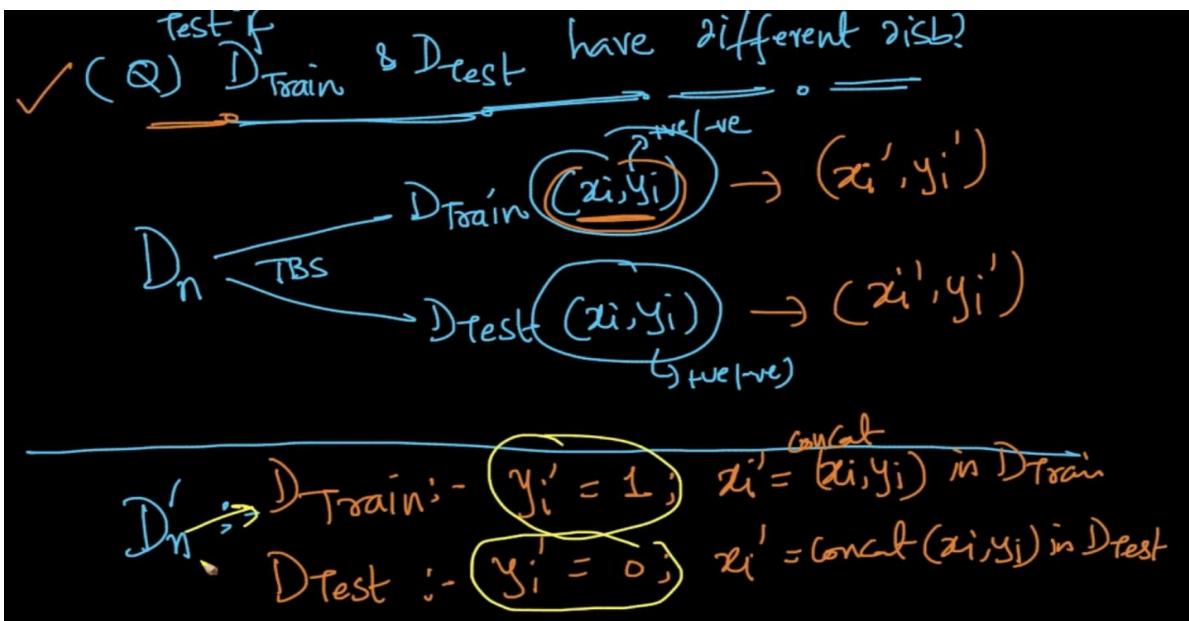


In TBS data can be changed overtime as seen i.e new products can be added or old products removed in Amazon Food Reviews Dataset

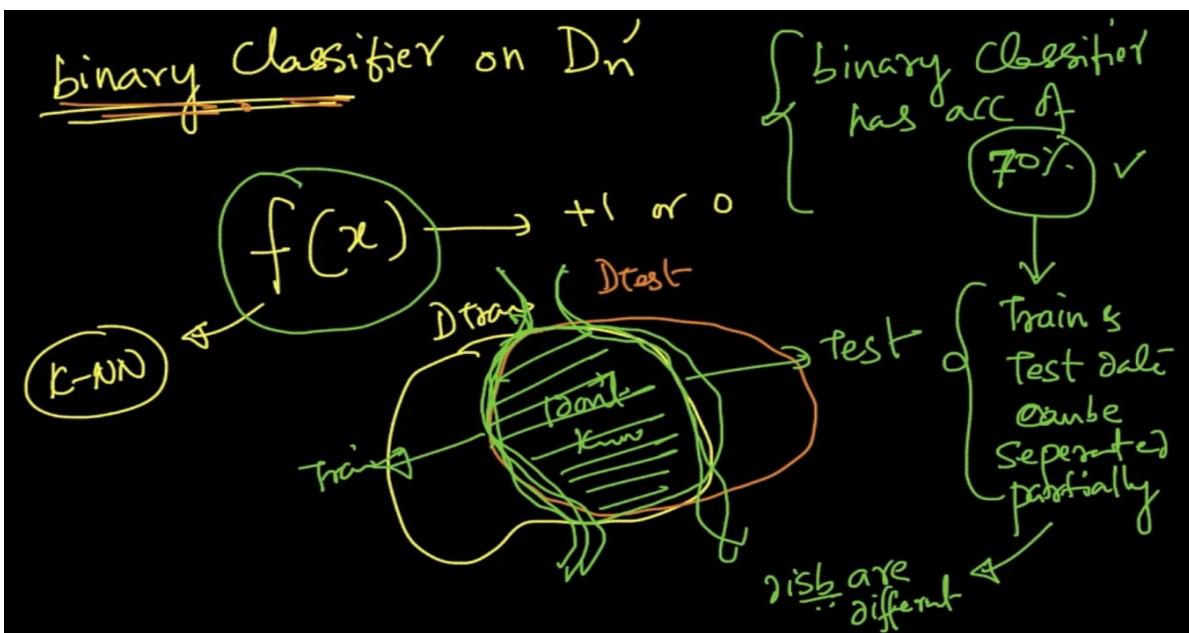


Our decision boundary from k-NN performs very well on Train data but with a change of data overtime our distributions can change . So k-NN will perform badly in Test data. So our CV error will be low but Test error is very high.

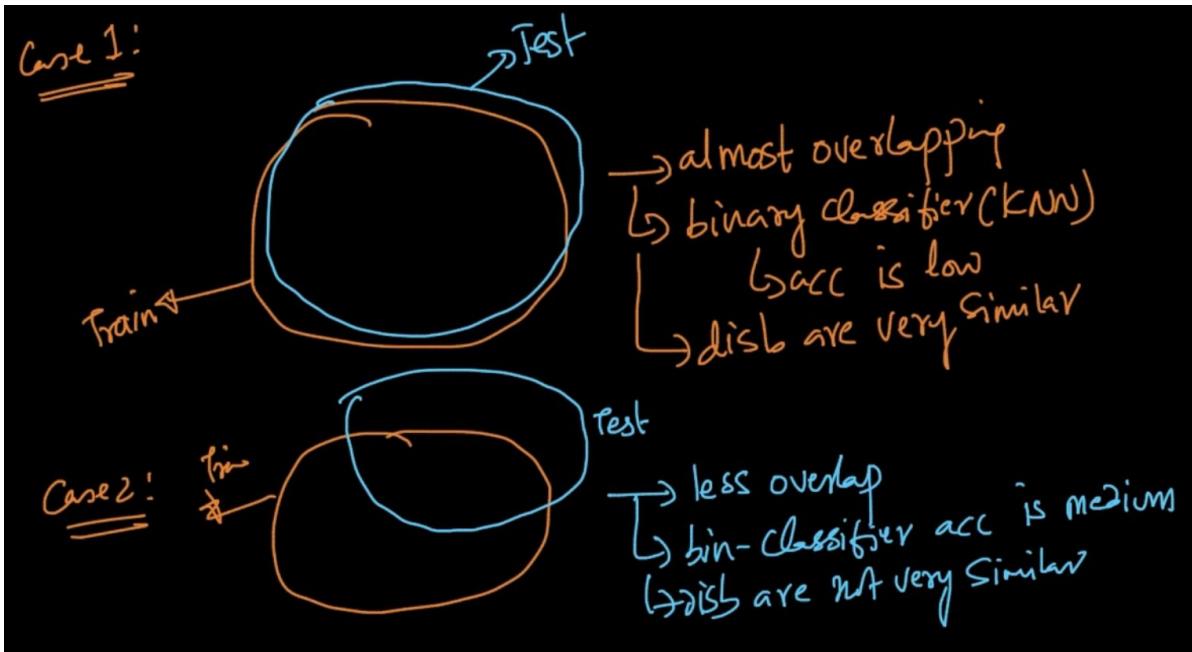
So how to determine if data is changing overtime ?



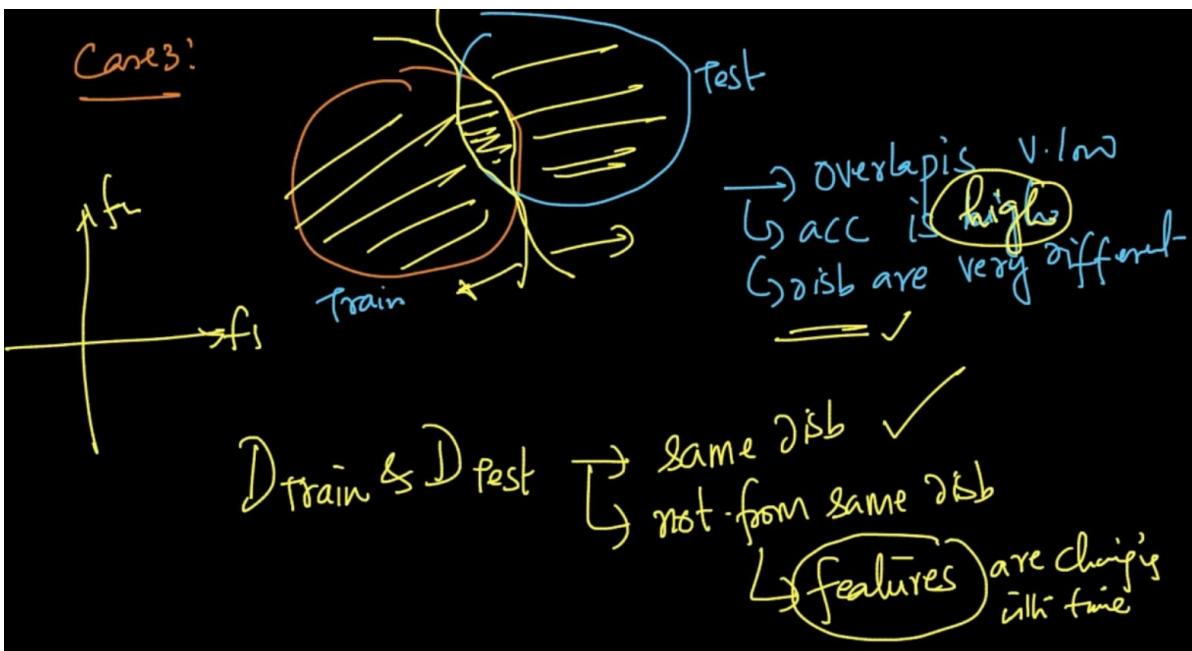
So , we make a new dataset D'_n and Train data $x'_i = \text{concat}(x_i, y_i)$ of previous Train data D_{train} and $y'_i = 1$. Same for test data except $y'_i = 0$



So we are doing binary classification on D'_n and check accuracy. Here it is 70% means it Train and Test data can be separated partially which implies distributions are different. Since the overlap isn't much binary classification performs good here

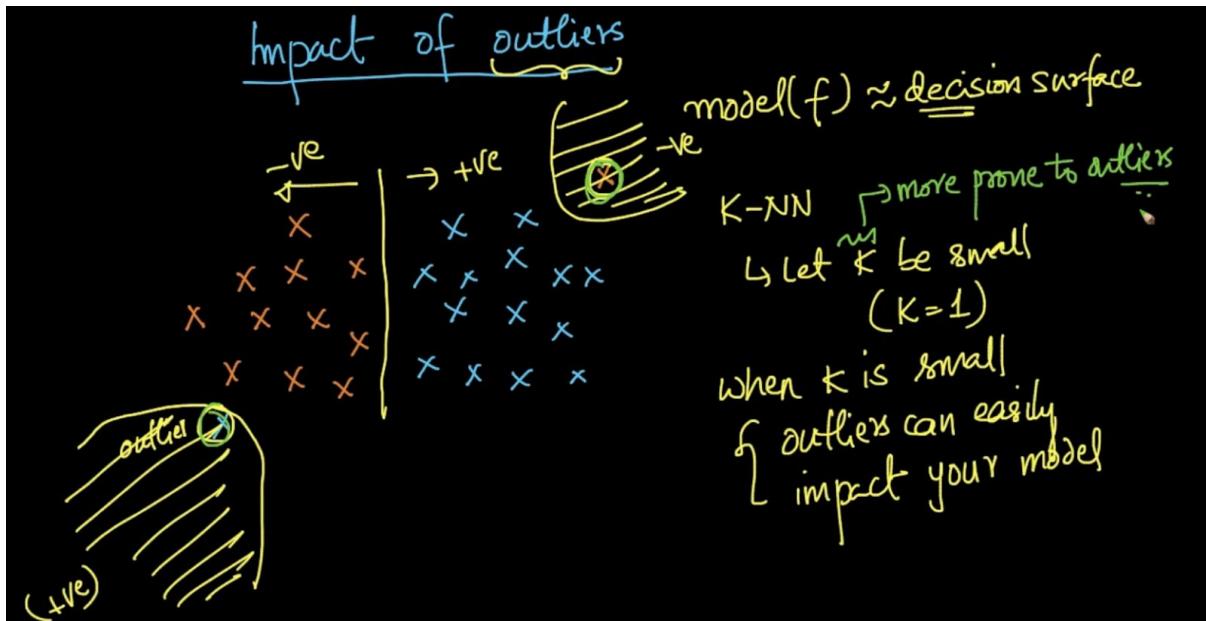


Case 1: Since they are overlapping they've the same distributions therefore they can't be separated by binary classification i.e its accuracy will be low

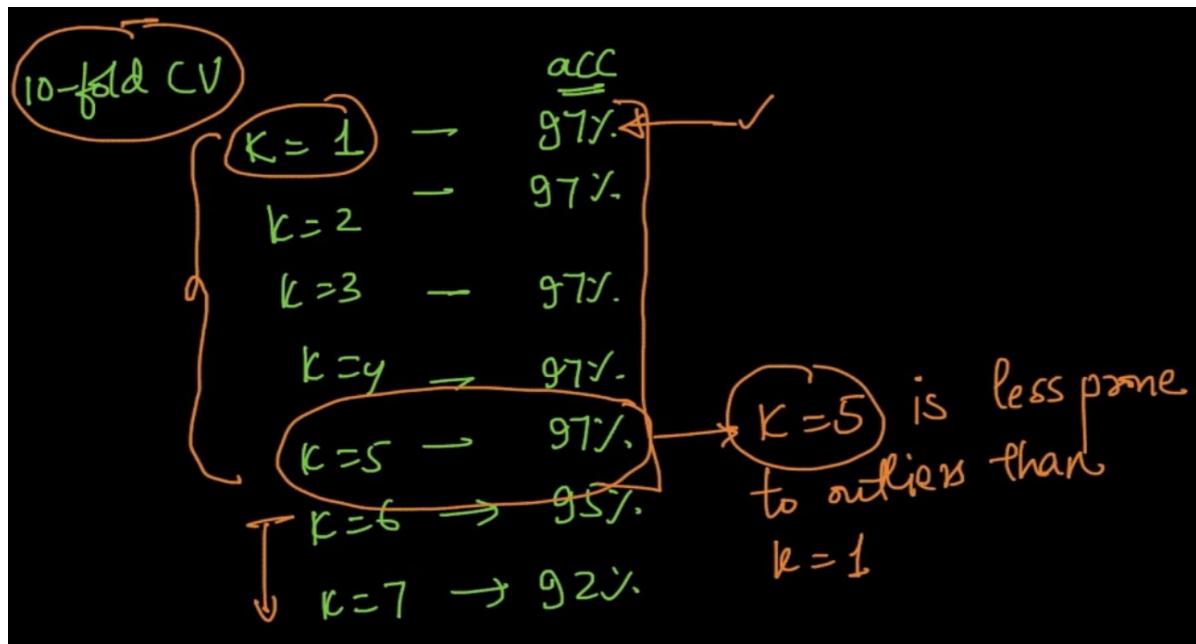


Case 3: Overlap is very low so binary classifier's accuracy will be very high . It means that features are changing over time or they are temporally stable. To tackle this we build/change/design new features

IMPACT OF OUTLIERS

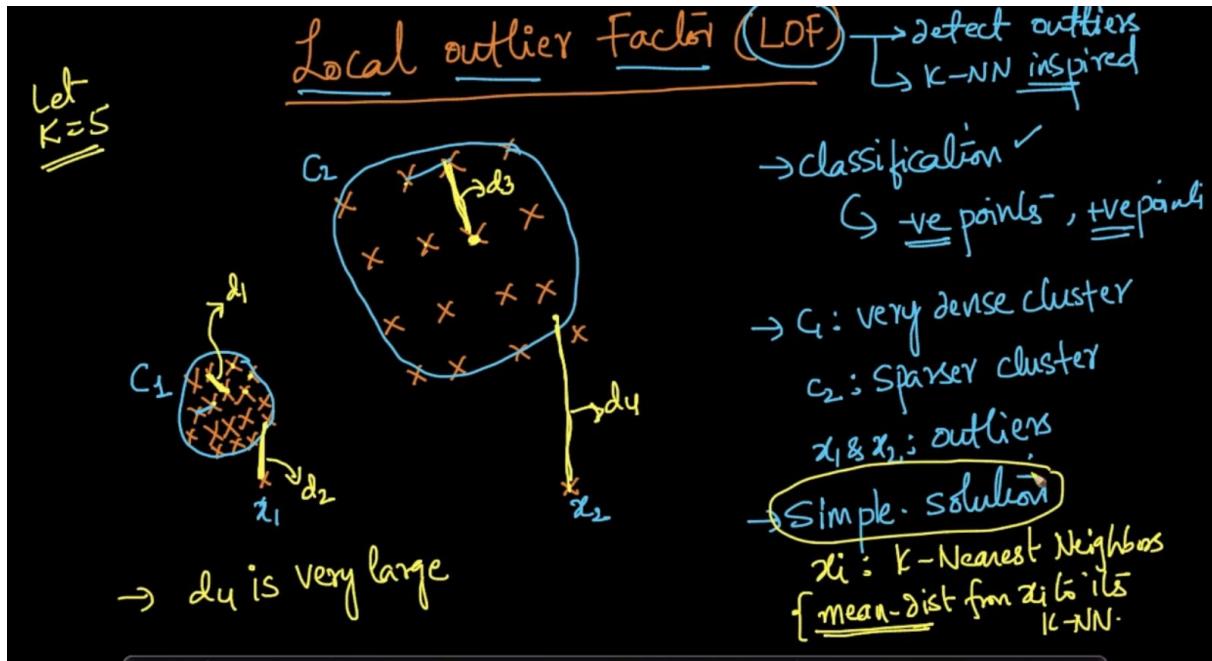


If $k=1$ and outlier is present then it'll mark all the query point with its class label which is wrong



In k-Fold CV if we are getting same accuracies on different k 's we'll choose a high k since they are less prone to errors. Here, $k = 5$

LOCAL OUTLIER FACTOR (MEAN DISTANCE TO KNN)



Mean distances of all data points x_i 's are calculated.

- Simple Solution
- 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

So, x_2 & x_4 are outliers but with this solution only x_2 will be considered an outlier as $d_2 \approx d_3$ means if we remove outlier x_1 with mean dist d_2 then all the other points i.e of sparse cluster C_2 has to be removed

k - DISTANCE & REACHABILITY DISTANCE

$\checkmark \quad \underline{k\text{-distance}}(x_i) = \frac{\text{distance to the } k^{\text{th}} \text{ Nearest Neighbor of } x_i \text{ from } x_i}{\text{Set of all points that belong to the } k\text{-NN of } x_i}$

$N_{k=5}(x_i) = \{x_1, x_2, x_3\}$
 $N_{k=5}(x_i) = \text{Neighborhood of } x_i$

$5\text{-distance}(x_i) = d_5$
 $1\text{-distance}(x_i) = d_1$

5-Distance (d_5) > 1-distance (d_1). Neighborhood is also explained pretty well

$\checkmark \quad \left\{ \begin{array}{l} \text{reachability-distance}(x_i, x_j) \\ = \max \left(k\text{-distance}(x_j), \text{dist}(x_i, x_j) \right) \end{array} \right\}$

$\max(d_5, d_4) = d_5$

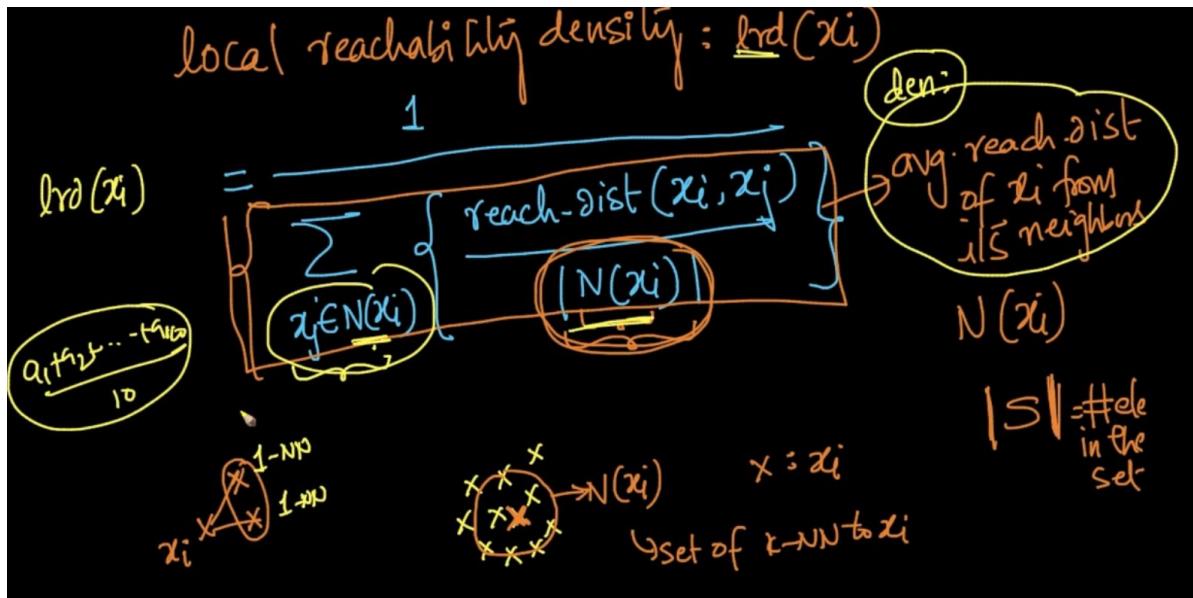
$\max(d_5, d) = d$

$\text{dist. of } k^{\text{th NN}} \text{ of } x_j \text{ from } x_i$
 \downarrow
 $\text{actual dist. b/w } x_i \text{ & } x_j$

$\left\{ \begin{array}{l} \text{if } x_i \in N(x_j) \\ \quad \text{then reach-dist}(x_i, x_j) \\ \quad = 1\text{-dist}(x_i, x_j) \checkmark \\ \text{else} \\ \quad \text{reach-dist}(x_i, x_j) = \text{dist}(x_i, x_j) \end{array} \right.$

Check the if - else condition to understand it

LOCAL REACHABILITY DENSITY

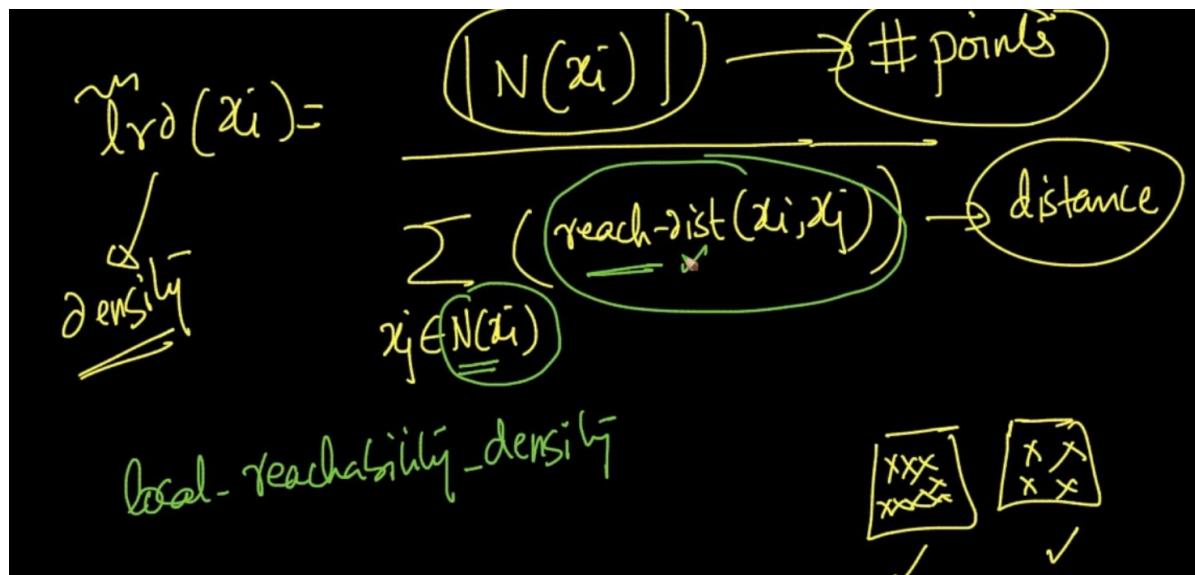


$\text{Lrd}(x_i)$ is inverse of summation of reachability distance of all the points belonging to the neighborhood of a point x_i divided by $|N(x_i)|$ i.e number of points in Neighborhood which is basically k but note that two points can also have the same NN .

eqn - english

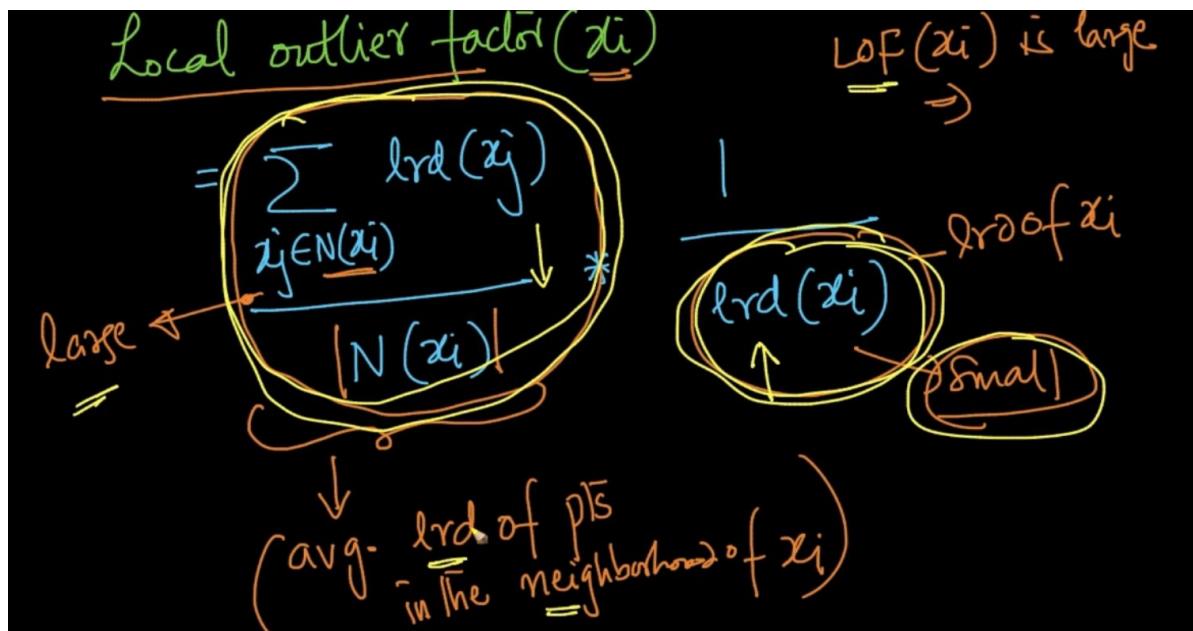
$$\left\{ \begin{array}{l} \text{Lrd}(x_i) = \text{inverse of avg. reach-dist of } x_i \\ \text{from i's neighbors} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{reach-dist}(x_i) = \text{if } \\ \text{else } = \end{array} \right.$$

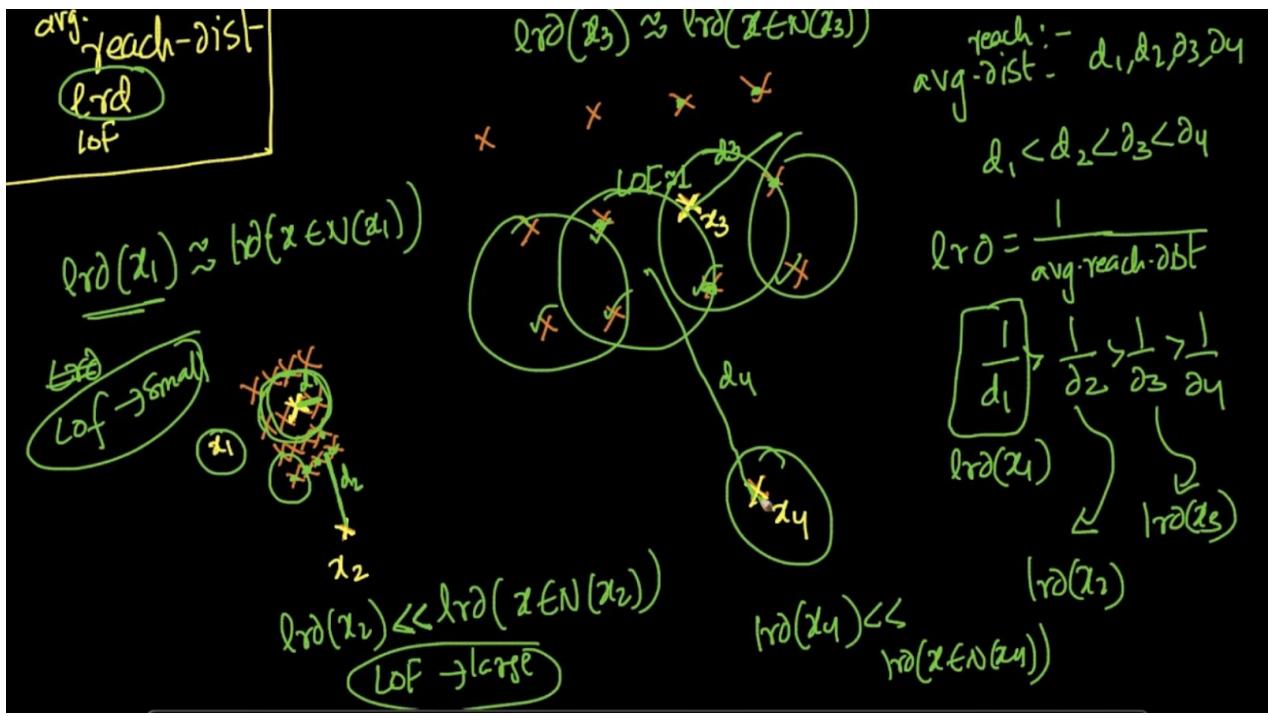


$\text{lrd}(x_i)$ is basically a density as it is No. of points/ Summation of reachability distance
Of all the neighborhood points

LOCAL OUTLIER FACTOR



$\text{LOF}(x_i)$ is large if the right lrd is small and left part is large. It simply means that density of points around x_i is small but density of its nearest neighbors is huge



Avg reach distance is calculated but it ain't enough and so is lrd calculation. So LOF is used.

For x_1 : Density is great so lrd is also great and $\text{lrd}(x_1) \approx \text{lrd}(x \in N(x_1))$ since they are very close and when we apply that to LOF formula above $\text{LOF}(x_1)$ will be small

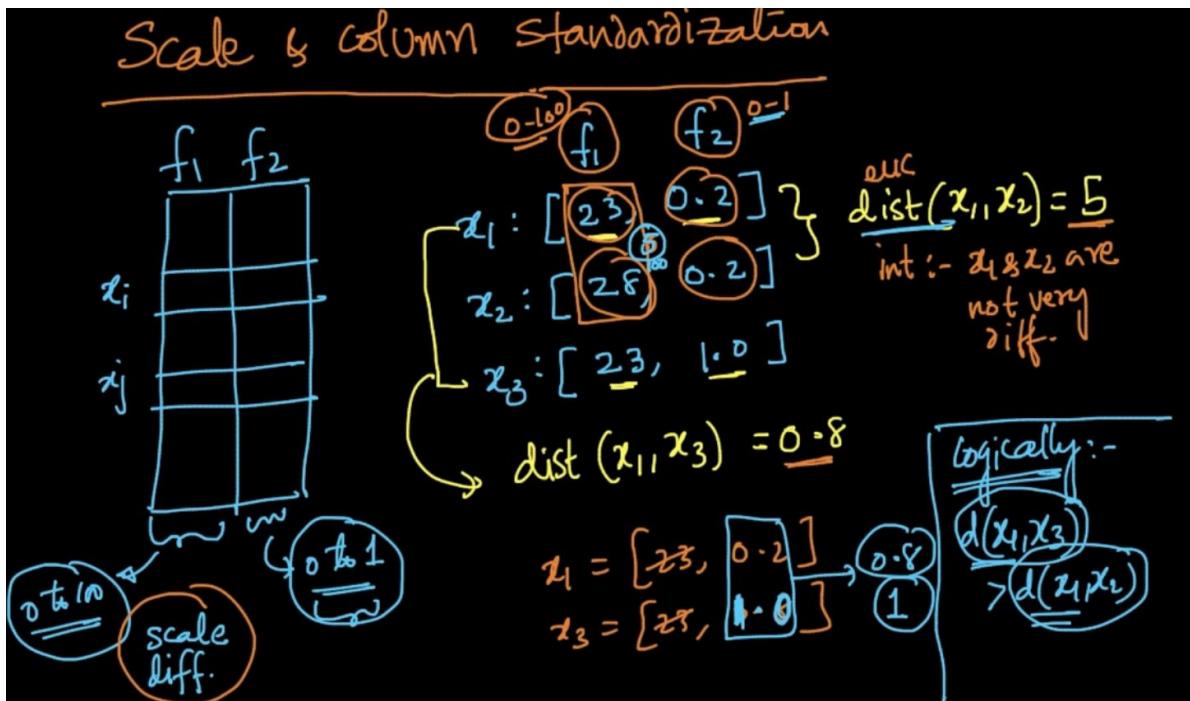
For x_2 : Density is none and therefore $\text{lrd}(x_2) \ll \text{lrd}(N(x_2))$ so $\text{LOF}(x_2)$ is large.

Same concept for other points x_3 & x_4 . So, we can observe that if LOF is large then it's an outlier

→ for each pt x_i compute $\text{LOF}(x_i)$
 → pick points with highest LOF → outliers

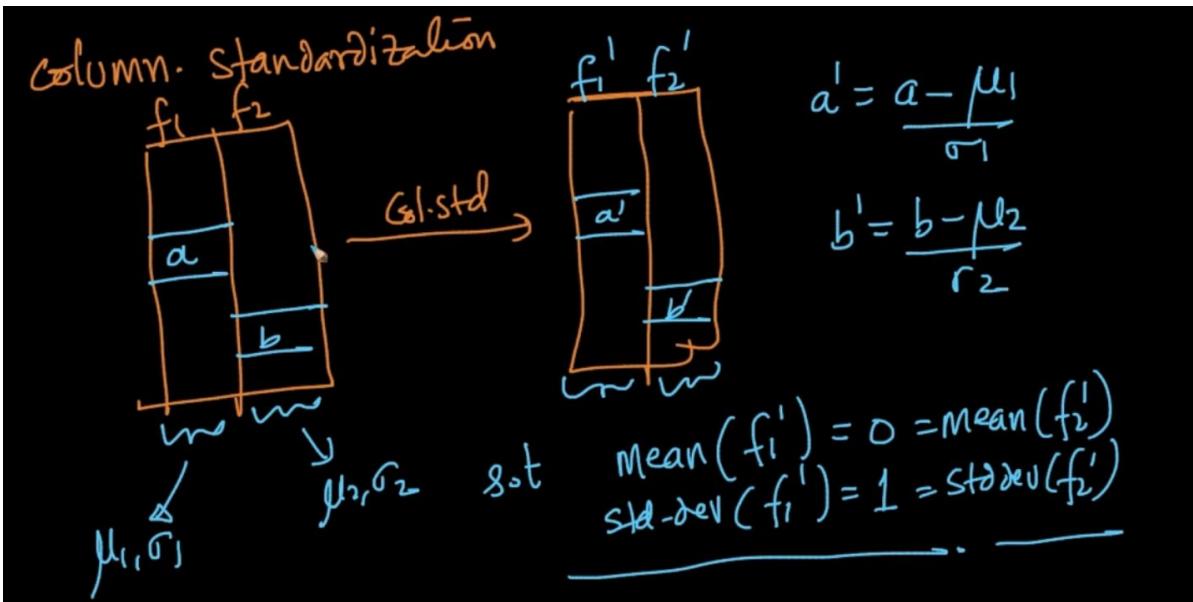
After doing it remove some percentage of outliers.

IMPACT OF COLUMN & SCALE STANDARDIZATION

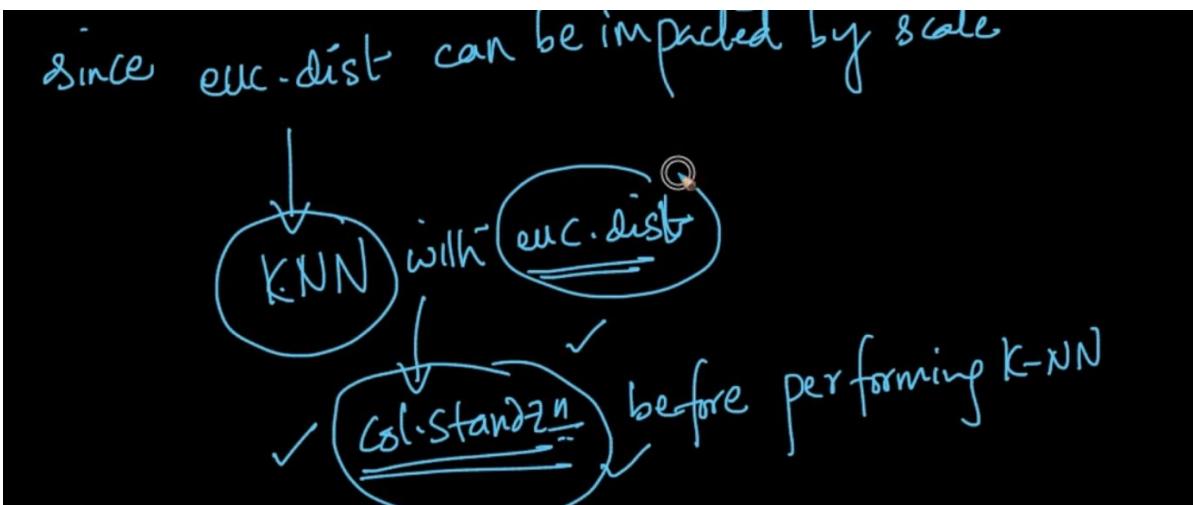


We've 2 features f_1, f_2 and their Euclidean distances are calculated. $\text{euc-dist}(x_1, x_2) = 5$ and $\text{euc-dist}(x_1, x_3) = 0.8$ but distance of x_1, x_2 is coming from f_1 scaling from 0-100 . So distance 5 isn't much of a difference because it is just 5% of variability but distance of x_1, x_3 is coming from f_2 which is scaling from 0-1 so distance 0.8 denotes 80% variability so logically $d(x_1, x_3) > d(x_1, x_2)$.So Euclidean distance is impacted by scales

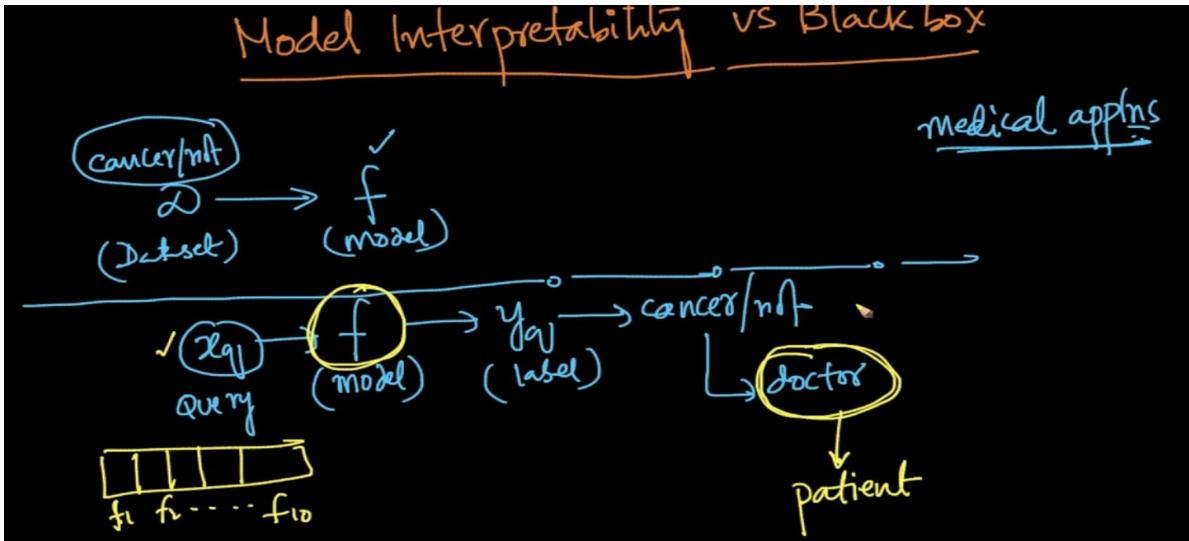
So how do we tackle this?



So we bring the features to the same scale by making their mean 0 and std-dev 1

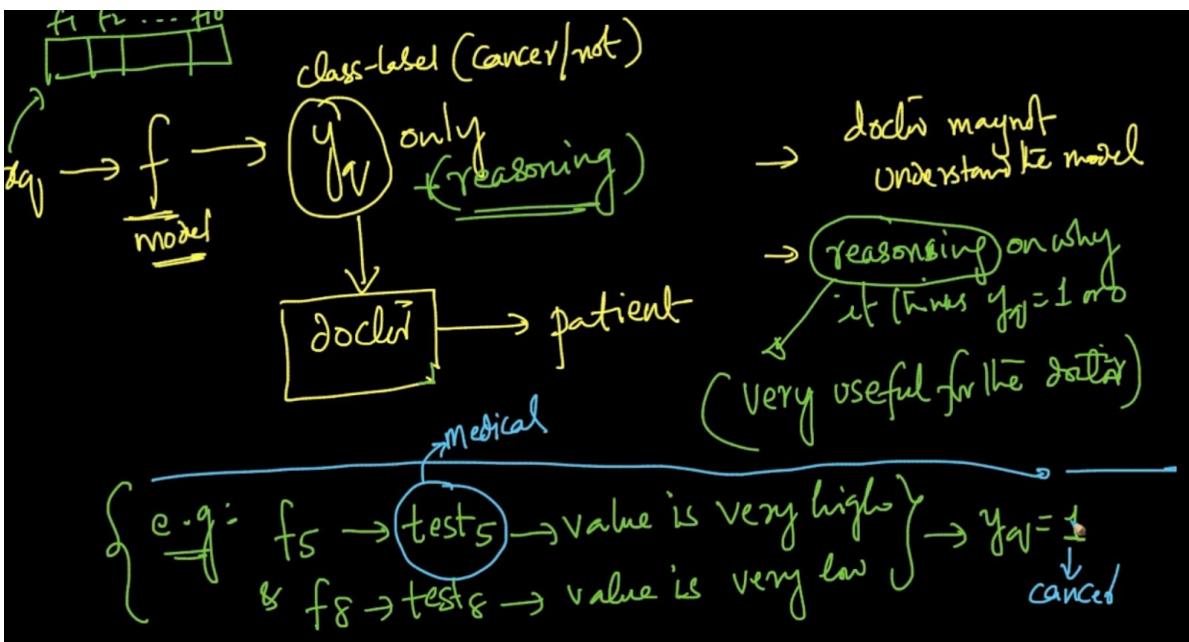


INTERPRETABILITY



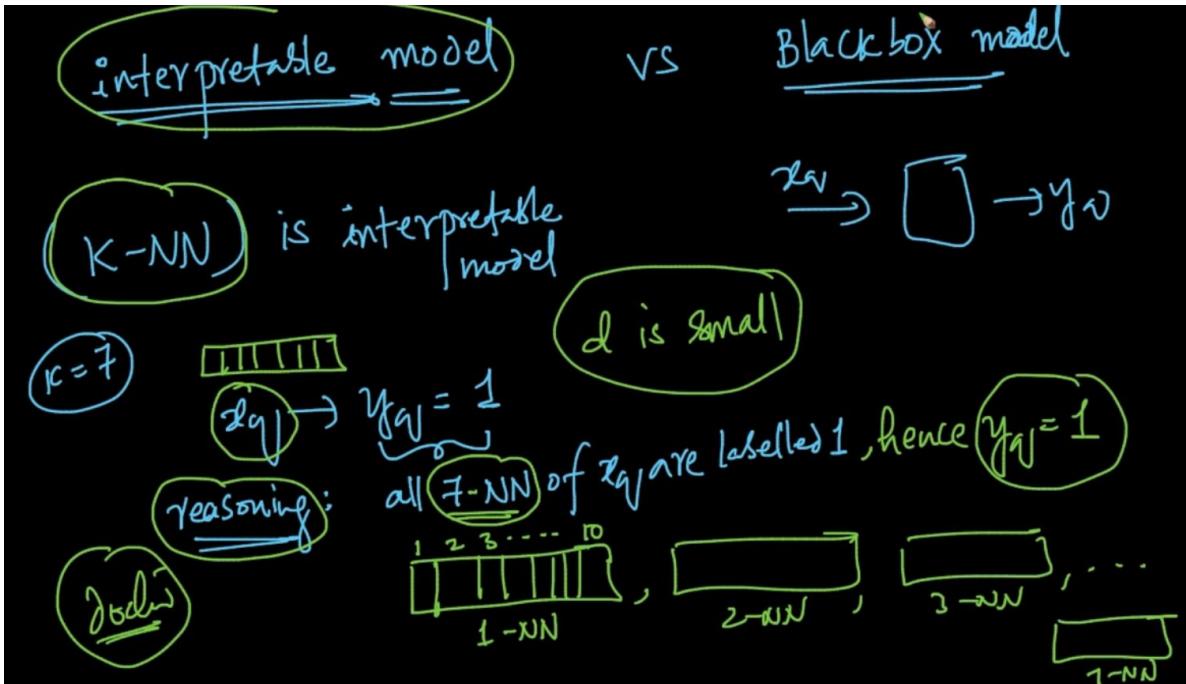
Suppose we've model for cancer diagnosis and we train it to predict cancer in a patient.

If our model just gives prediction label y_q only the Doctor can't be sure because he isn't an ML Engineer so he can't be sure with just y_q



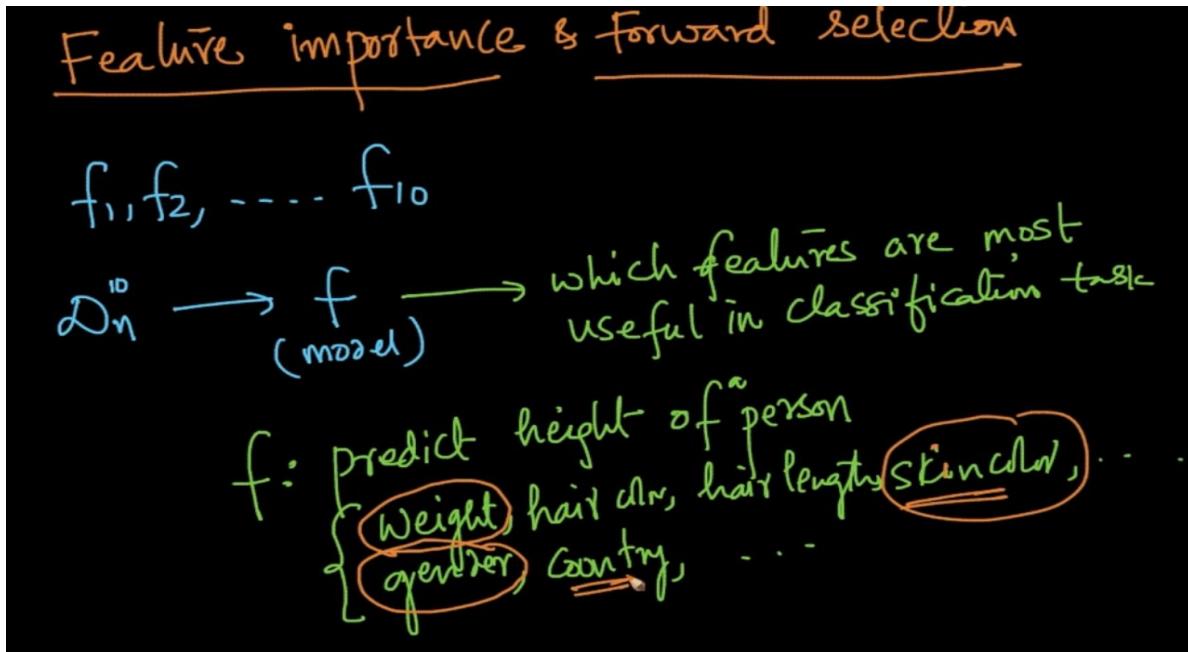
So if the model also give reasoning like above where some medical tests which are features of dataset have some values like high/low that's why $y_q = 1$ i.e cancer

If model can give explanation like this it's called interpretable model

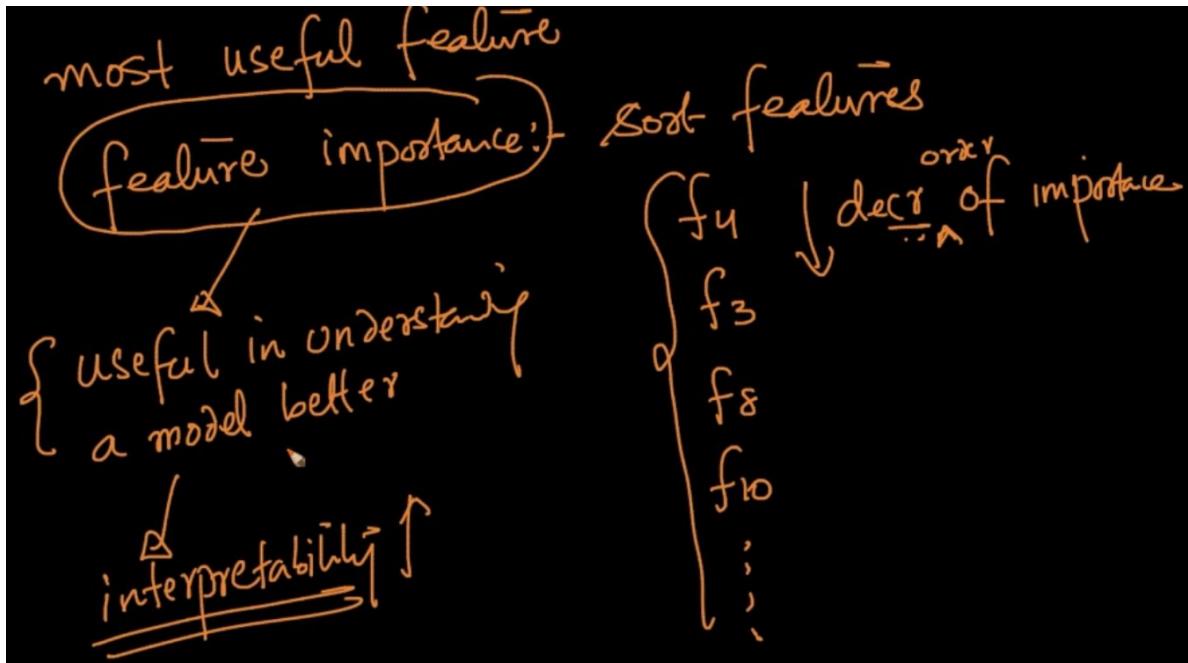


If we are running k-NN with 7-NN and we are getting $y_q = 1$ then we can interpret the model because our NN have medical tests as dimensions but d is low and our model is giving reasoning that all the neighbors have very similar results on tests and they all have cancer i.e $y_q = 1$. So, the doctor can understand our model because d is also small .

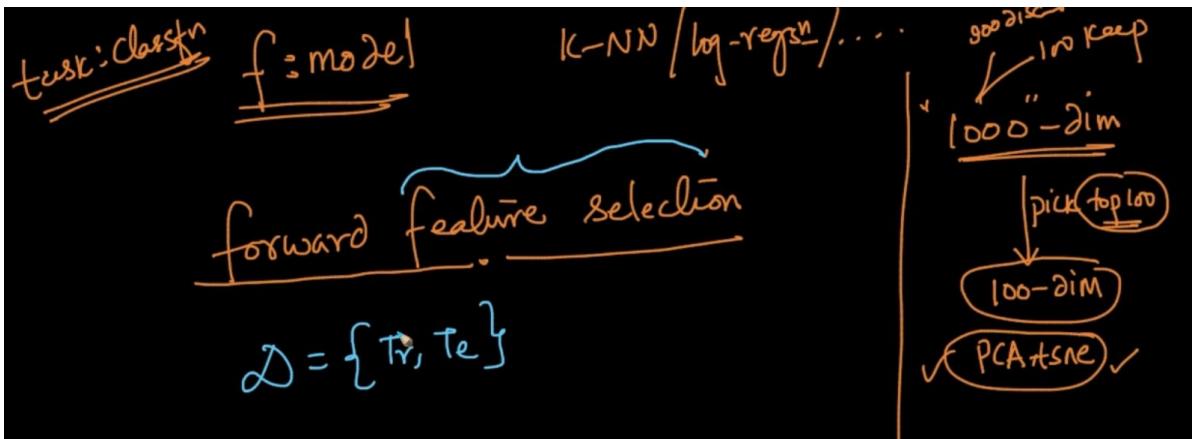
FEATURE IMPORTANCE AND FORWARD FEATURE SELECTION



For predicting heights we need to get the most important features

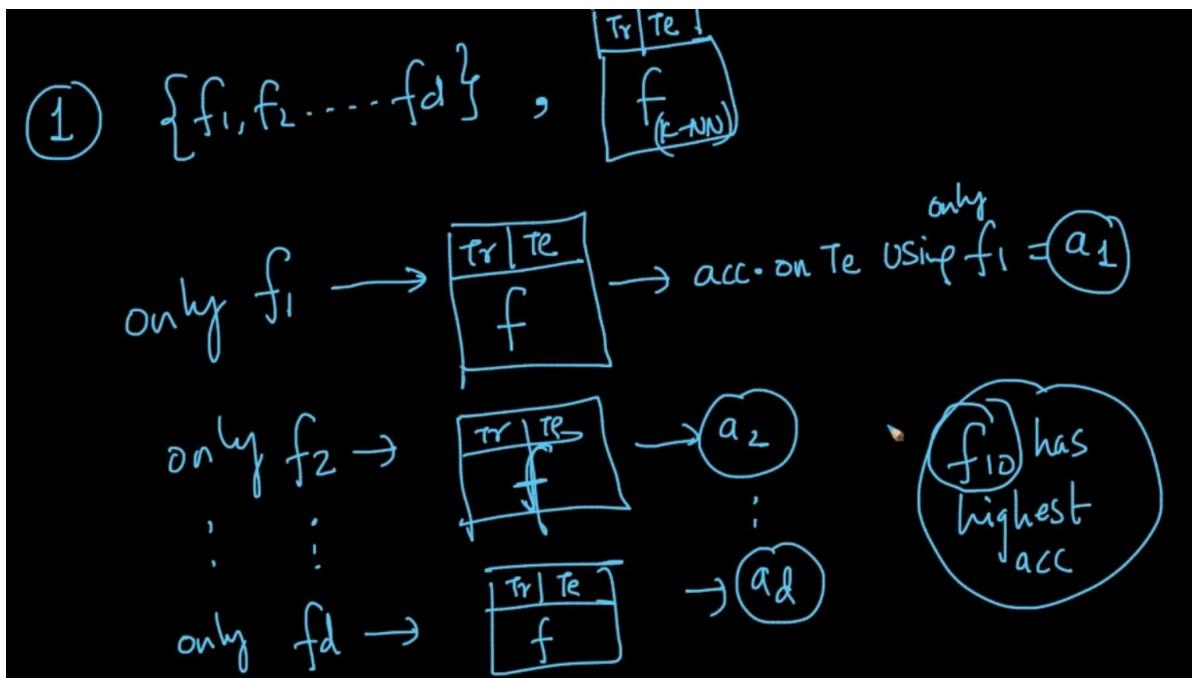


If we know the most important features it means that we can understand model better which means more interpretability

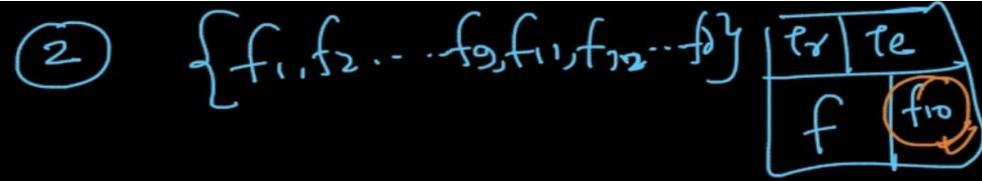


So a method called forward feature selection is used to pick the top 100 dimensions

Note - PCA, tsne are dimensionality reduction techniques but they maintain the distance instead of giving us the best features



We've d-features . We take each feature till we test our model on every feature with their corresponding Train and test data and check accuracy. Suppose feature f_{10} has the highest accuracy then we'll select that



$f_1 \rightarrow$

f	f_{10}
f	f_1, f_{10}

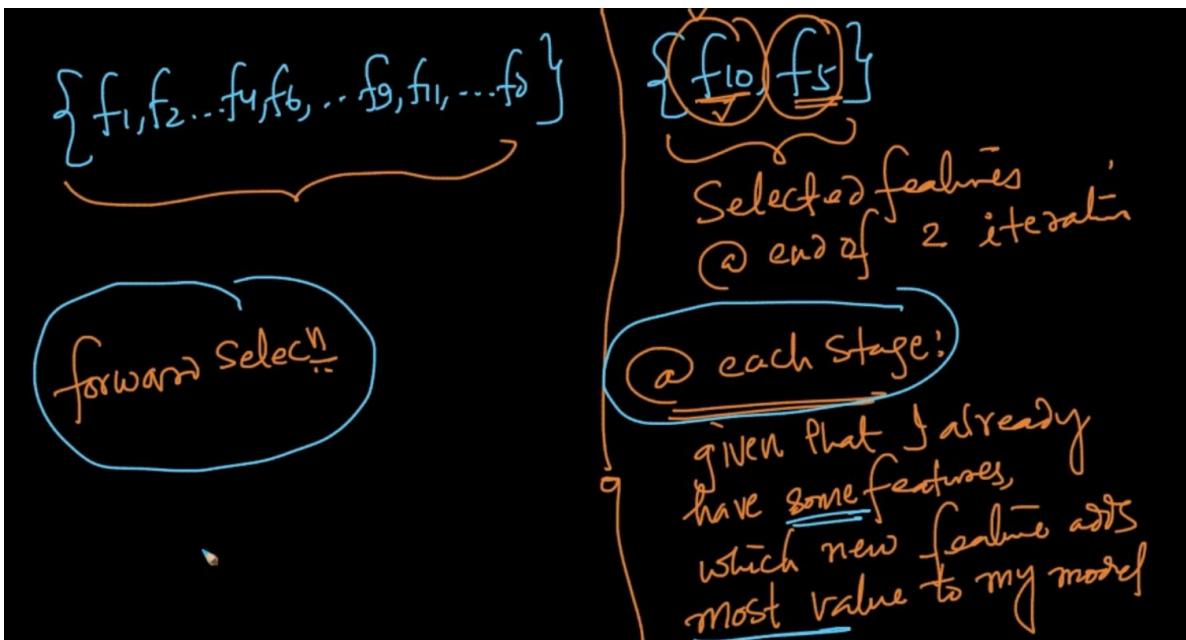
 $\rightarrow a_1$

$f_2 \rightarrow (f_2 \& f_{10}) \rightarrow a_2$

\vdots

a_d

Feature f_{10} is the most important feature. Now, we'll select it and along with it train our models and get accuracies taking both the features. Now let's say f_5 has got the highest accuracy. Therefore our features f_5, f_{10} are our selected features



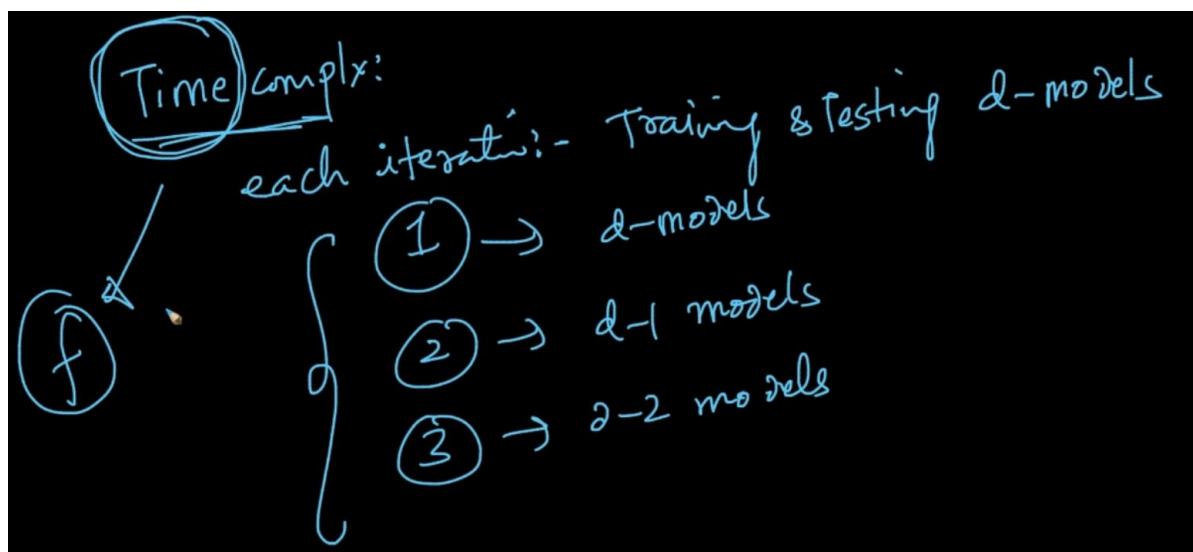
So we keep repeating this till we get our number of desired features. Let's say we wanted our top 100 features from 1000 features

Backward Selⁿ

$$\{ f_1, f_2, \dots, f_d \}$$

remove a feature in each iteration which results in lowest drop in accuracy

In Backward selection we remove a feature in each iteration which results in lowest drop of accuracy



Time complexity is very large obviously but it can be done in any model or algorithm

HANDLING CATEGORICAL AND NUMERICAL VALUES

e.g.

	W	hair color	Country	...	$y = h$
numerical	100	Black	US	...	180

χ : numerical vector

$f(x) = y = \text{height}$

weight = numerical
↳ as-is

haircolor: {Bl, Br, Red, Golden, gray}

We are trying to predict heights of people and it's features X are {Weight. Hair color, country} . So, weight is numeric but haircolor isn't : {Black, Brown, Red,...} it's a categorical variable but we want our X should be a numerical vector so hair color should be converted to numeric instead of categorical.

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

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

Red > Bl χ

② One-hot encoding: haircolor

binary vector

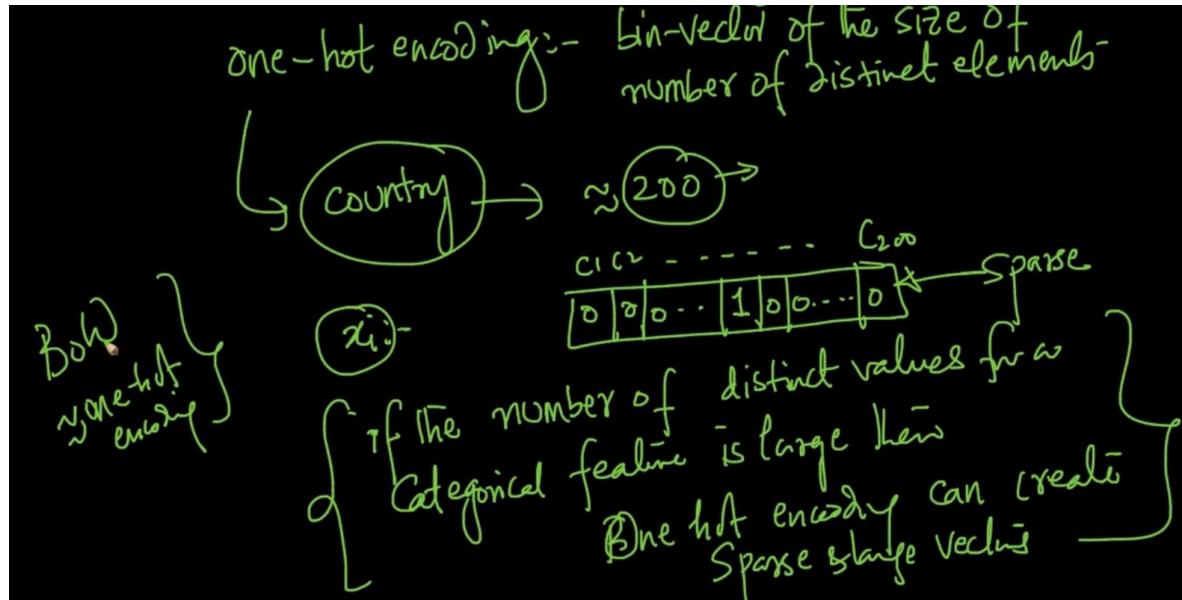
$x_i = (\text{Br}) \therefore \boxed{0|1|0|0|0}$

$x_i = G \therefore \boxed{0|0|0|1|0}$

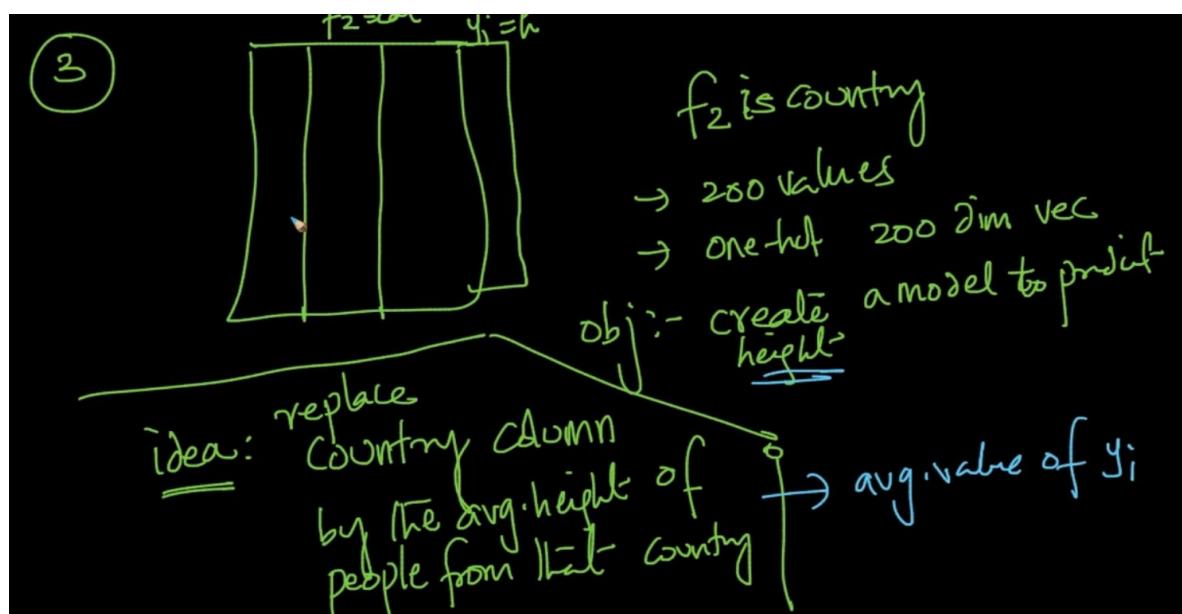
Bl	Br	Red	G	Gray

1) We are giving them a number but it'll create an inherent order i.e if Red => 3 and Black => 1 which will indicate Red > Bl

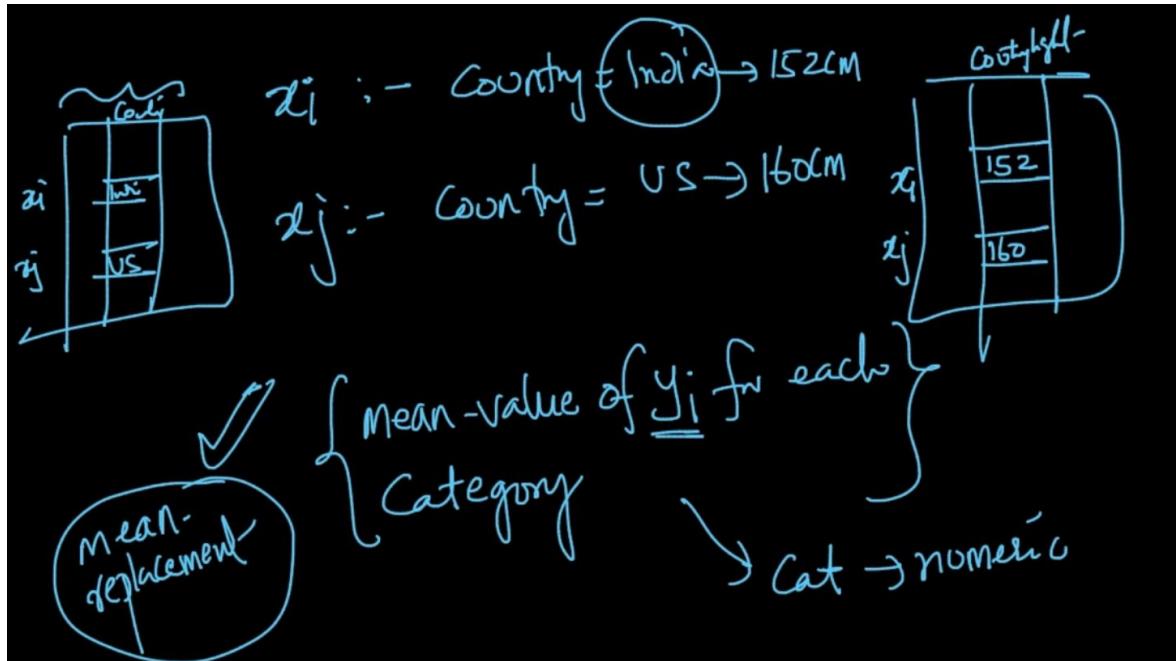
2) So we use one hot coding in which we create a d dimensional vector where d is number of distinct categories like color of hairs in haircolor feature in which we feed binary values
Here, no inherent order like Red > Black



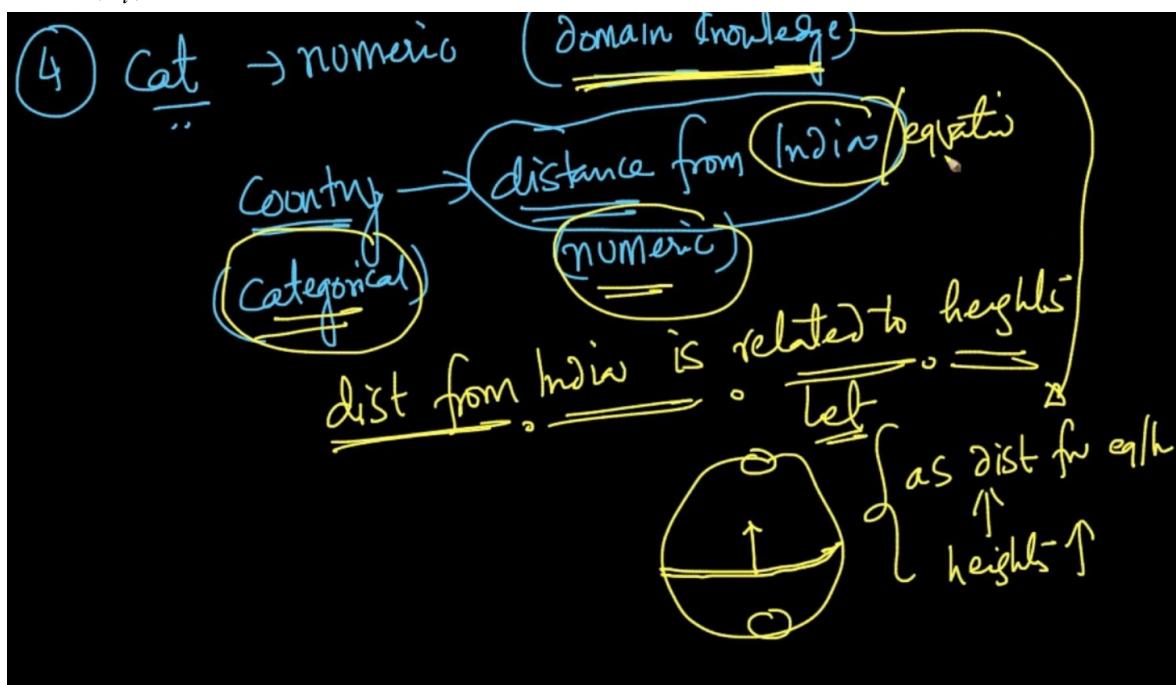
But there's a problem with one hot encoding i.e if the number of distinct values or d is large for ex : like countries where countries = 200 sparse & large vectors can be created



Suppose we've country feature . We'll replace it with average country height but note that it is very objective specific i.e related on predicted values



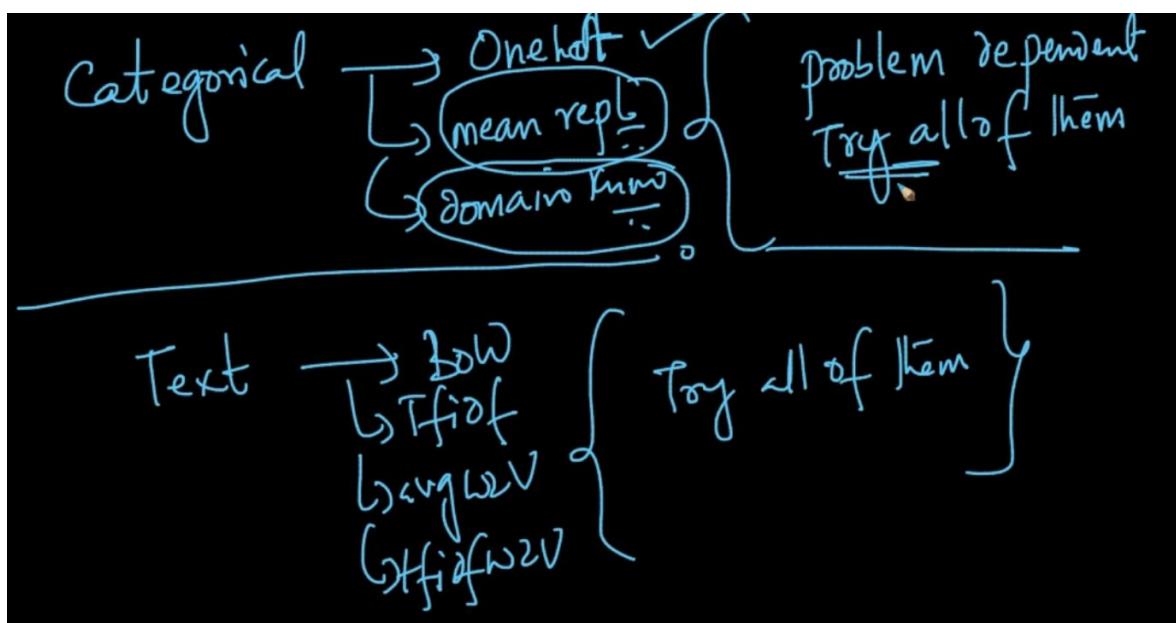
Replacing categorical variable Country with average/mean of y_i for each category i.e
Mean (y_i) for India is 152 cm so Country Height is 152 for Indians



Domain knowledge : If an anthropologist says that distance from India is related to heights as distance from India \uparrow heights \uparrow . So this is another method of converting categorical to numerical

f_1	f_1'	f_1''	
V-good	5	5	ordinal Cat-features
good	4	4	↓
avg	3	3	numeric
bad	2	2	
V.Bad	1	1	(?)

Ordinal categorical features like {V.Good, Good, avg, bad, v. bad} they can be converted into numeric by giving them numeric value based on their severity like 5 for V.Good and 1 for V.Bad but we there's no proper answer on how to give numeric values to them we can also give 10 => V.Good and 1 => V.Bad



So these are the methods for dealing with Categorical but there's no better because they are very problem dependent. So what do you do? You try all of them

HANDLING MISSING VALUES BY IMPUTATION

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_j	✓	✓	✓	✓	
x_i	✓	✓	✗	✓	
x_k	✓	✓	✗	✓	

missing
NaN
null
-1

$x_i : f_3 \rightarrow$ missing

{all the non missing values for f_3 }

Sometimes, in some features values can be missing or NaN. To deal with it we do imputation with the methods listed here

In our missing feature f_3 of x_i we can put the mean of all non missing values for f_3

classification → 0, 1

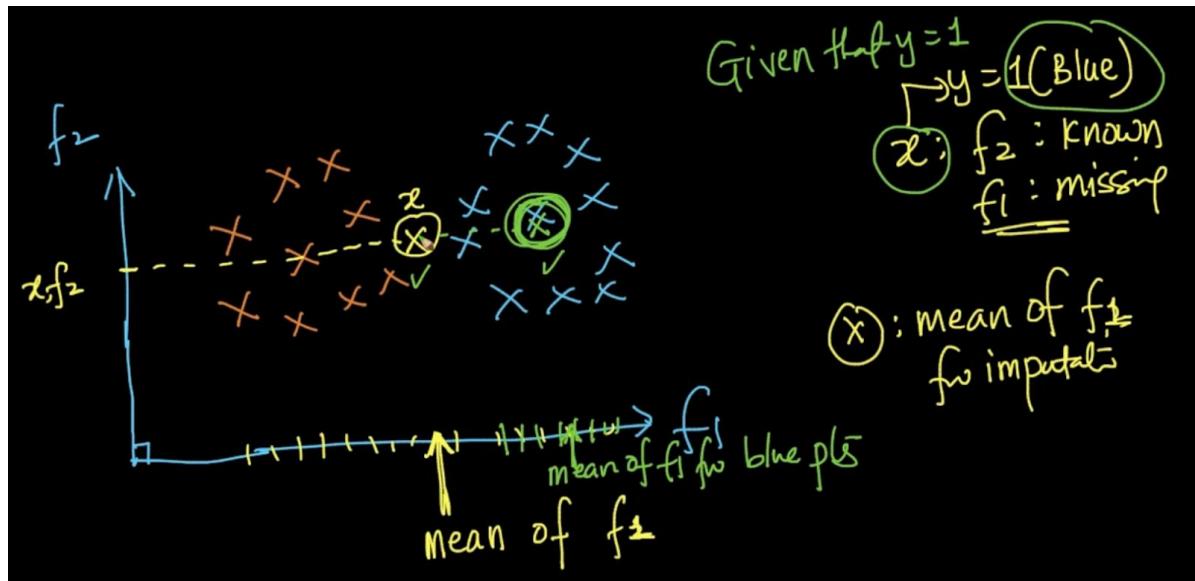
impute based on class label

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

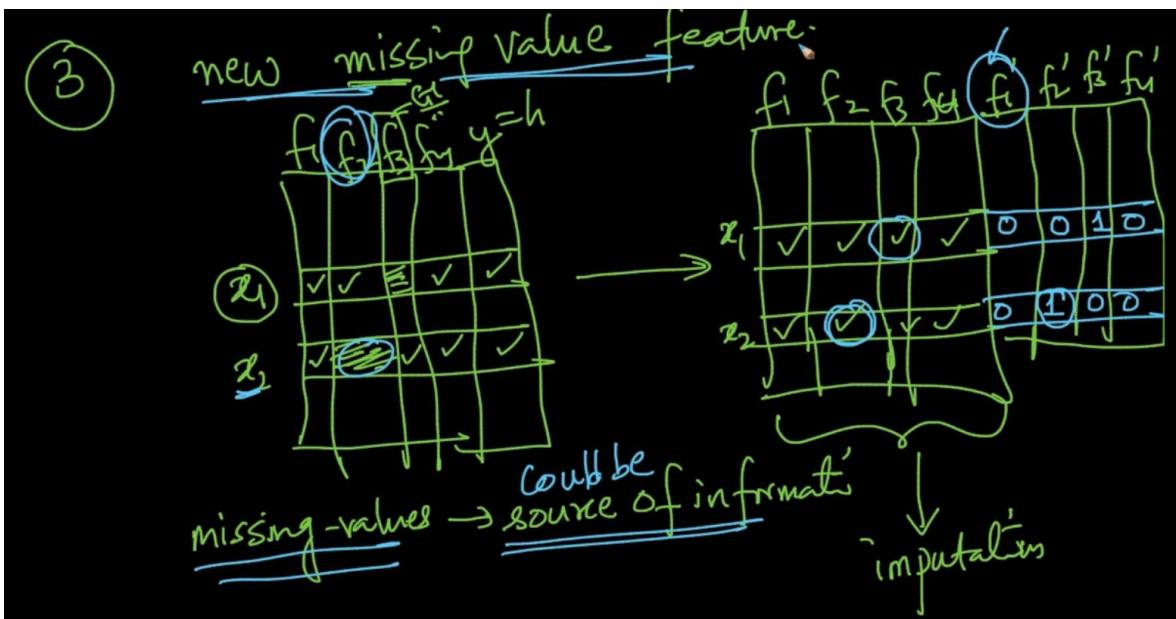
$\hookrightarrow y_i \rightarrow 1 (\text{class label})$

{mean-imputation based on class label} \leftarrow {mean of f_3 for all pts whose class label is 1}

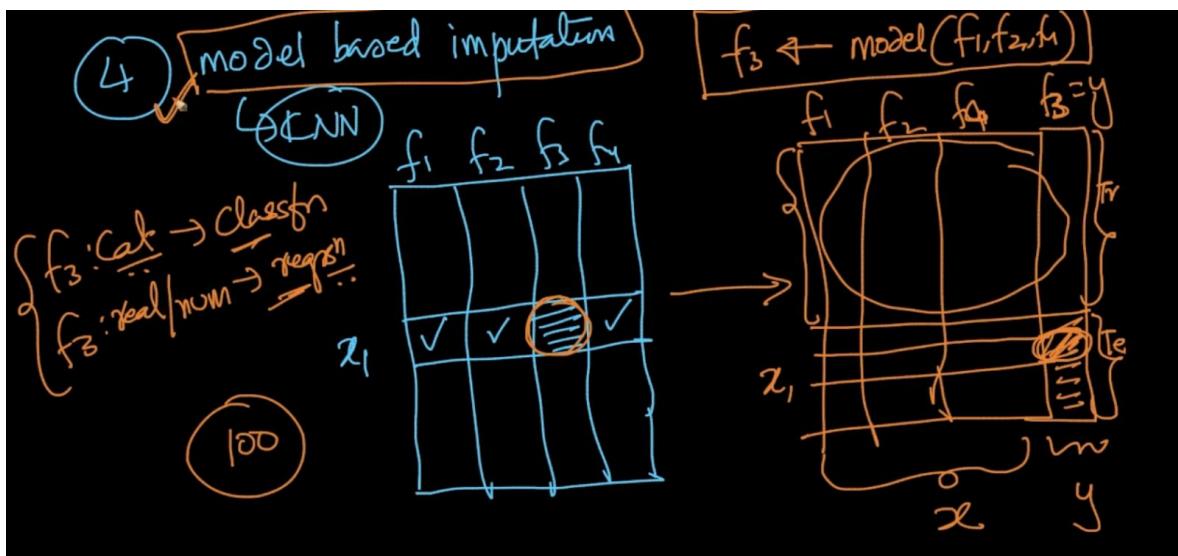
In classification , for f_3 of x_i whose class label $y_i = 1$ instead of taking mean of all the pts of f_3 for missing values we can just get the Mean of f_3 for points where $y = 1$



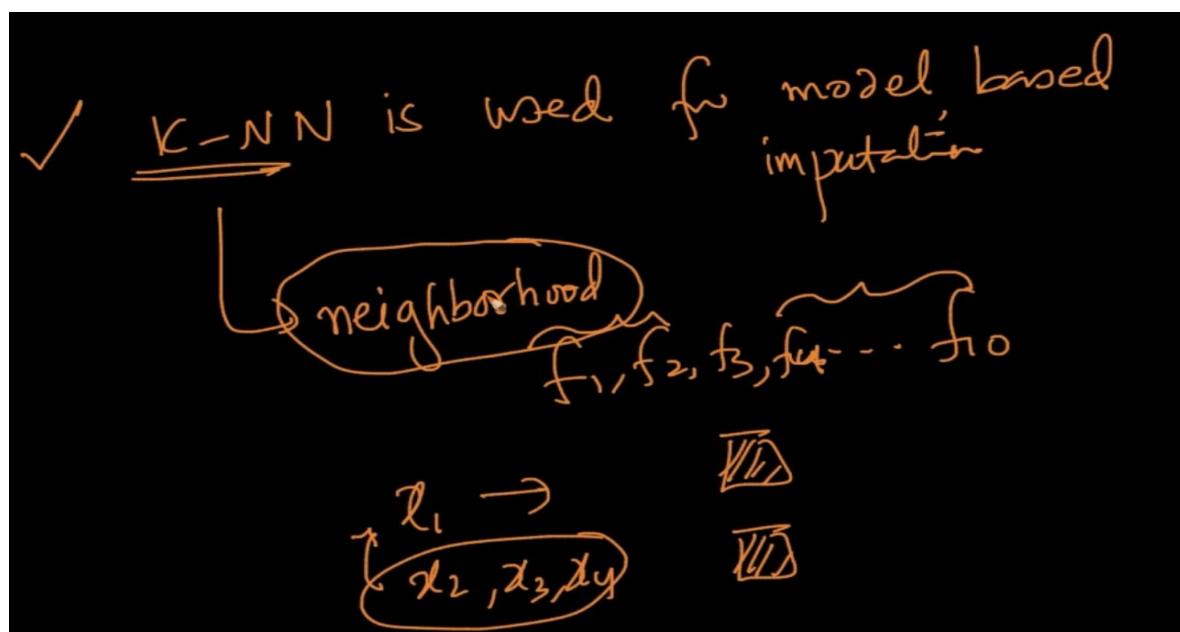
Here we want x whose feature f_1 is missing but $y = 1$. If we are taking the mean of all f_1 it'll lie in the yellow region but if we take the mean of f_1 where $y = 1$ then it'll be at a much better place as it seems because it's $y = 1$ as well



We add new features by adding a binary vector which is giving us info about the missing feature. Sometimes missing values could be the source of info



In model based imputation, e.g if f_3 has missing values . Then we can make it as a class label and use the other features to predict them. As you can see f_3 has been used as class label and in that all the missing values are in Test so we can predict them based on their Train data which has values of f_3



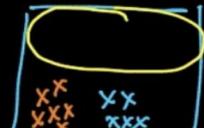
K-NN is used for model based imputation. Data point x_1 has missing value in feature f_3 . But x_2, x_3, x_4 are very similar to x_1 in other features so we can predict the missing value of f_3 as well

CURSE OF DIMENSIONALITY

Curse of dimensionality

→ d is high

↳ dim



① ML:

binary features

$\left\{ \begin{array}{l} f_1, f_2, f_3 \rightarrow \text{Total \# datapoints} = 2^3 = 8 \\ f_1, f_2, f_3, \dots, f_{10} \rightarrow \qquad \qquad \qquad = 2^{10} = 1024 \end{array} \right.$

$\left\{ \begin{array}{l} \text{as dim} \uparrow; \text{ The \# datapoints to perform good classfn (more) increase exponentially} \end{array} \right.$

Suppose we've three binary features/dimensions it'll require $2^3 = 8$ datapoints and if 10 features $2^{10} = 1024$ so as dimensions increase the number of datapoints to perform good classification increases exponentially

Highest phenomenon:-

$\text{lock} = n \rightarrow$ size of dataset is fixed

$\checkmark \text{perf} \downarrow \text{as dim} \uparrow$

Hughes phenomenon : Our performances \downarrow as dimensions \uparrow if your size of dataset is fixed

② Distance functions (euclidean dist)

Cure of Dim:- intuition of dist in 3D
is not valid in high dim spaces

1D-world → n-random pts

$$\left\{ \begin{array}{l} \text{dist_min}(x_i) = \min_{x_j \neq x_i} \left\{ \text{dist}^{\text{euc}}(x_i, x_j) \right\} \\ \text{dist_max}(x_i) = \max_{x_j \neq x_i} \left\{ \text{dist}^{\text{euc}}(x_i, x_j) \right\} \end{array} \right.$$

Here we've defined two terms 1) Dist_min(x_i) which gives us the minimum distant point or closest point from all points from x_i . 2) Dist_max(x_i) which gives us the maximum distant point or farthest point from all points from x_i

(1D)
(2D)
3D → $\left\{ \frac{\text{dist}_{\text{max}}(x_i) - \text{dist}_{\text{min}}(x_i)}{\text{dist}_{\text{min}}(x_i)} > 0 \right\}$ when $d=1, 2 \text{ or } 3$

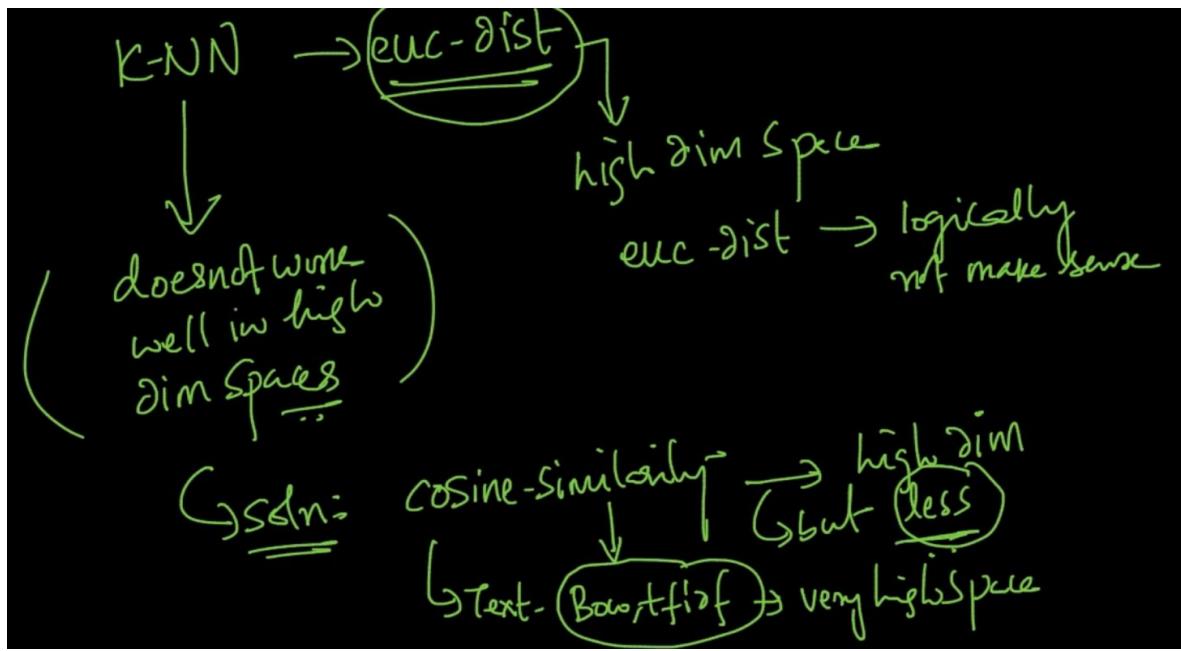
as $\dim \uparrow$

$$\lim_{d \rightarrow \infty} \left(\frac{\text{dist}_{\text{max}}(x_i) - \text{dist}_{\text{min}}(x_i)}{\text{dist}_{\text{min}}(x_i)} \right) \rightarrow 0$$

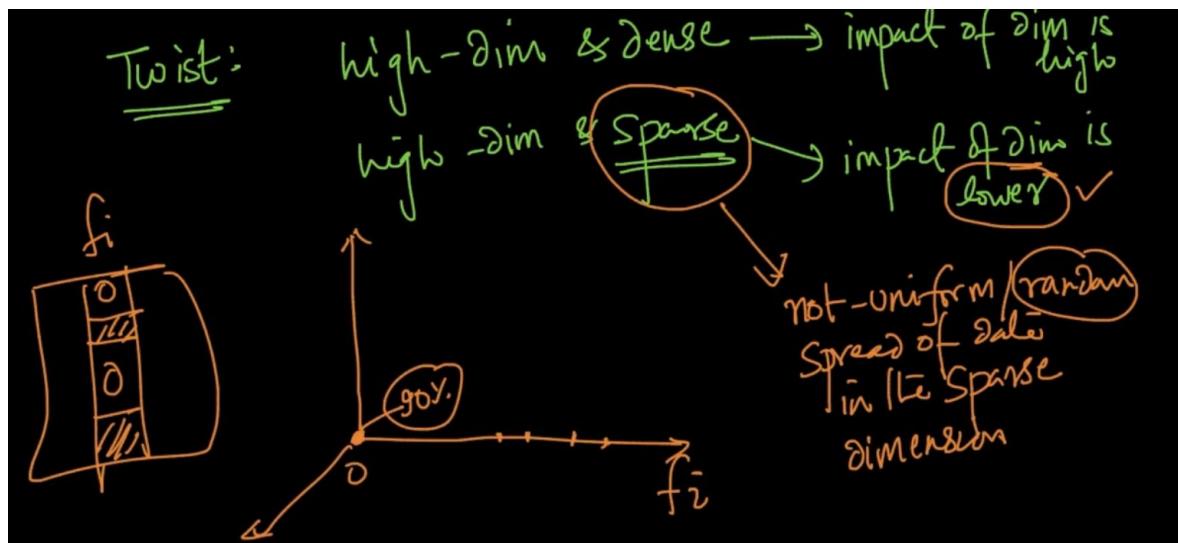
\Downarrow

$$(\text{dist}_{\text{max}}(x_i) \approx \text{dist}_{\text{min}}(x_i))$$

As dimensionality d increases the above ration tends to become 0 which means Distmax ≈ Distmin which is an alarming thing as it means every pair of points are equidistant from each other

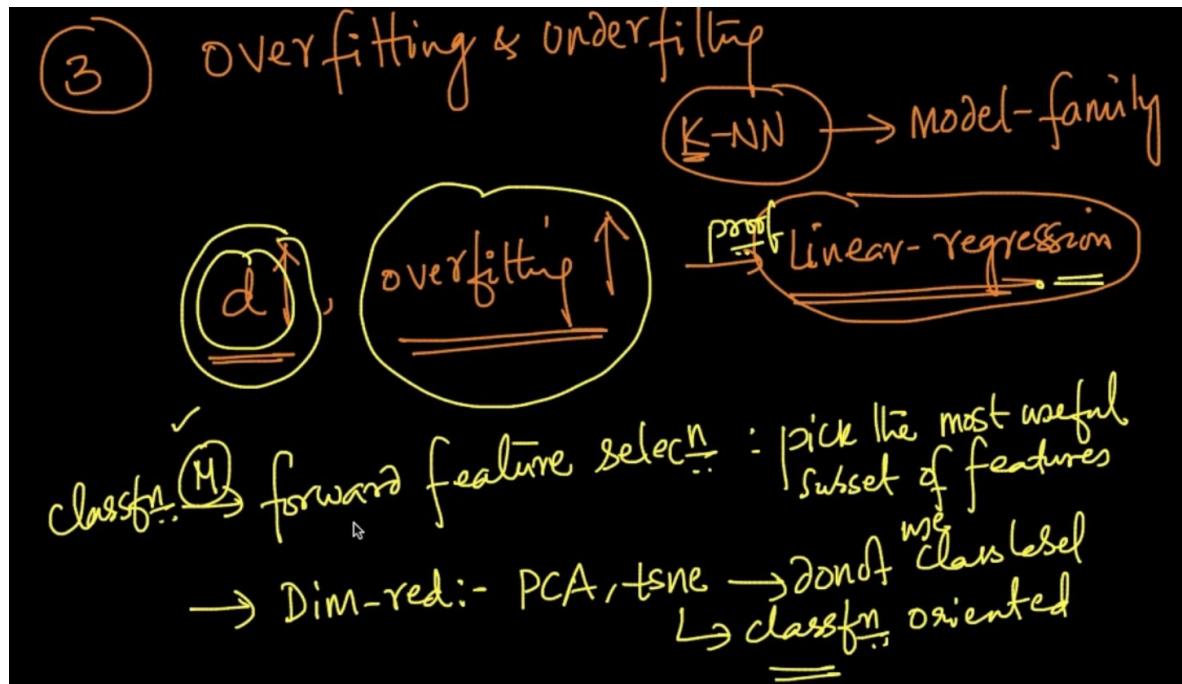


As distances don't make sense in high dimensions as $\text{Distmax} \approx \text{Distmin}$, k-NN don't work together very well so to tackle it we use Cosine similarity as it doesn't get affected by Curse of dimensionality and that's why we use it for text as it has very high dimensions,



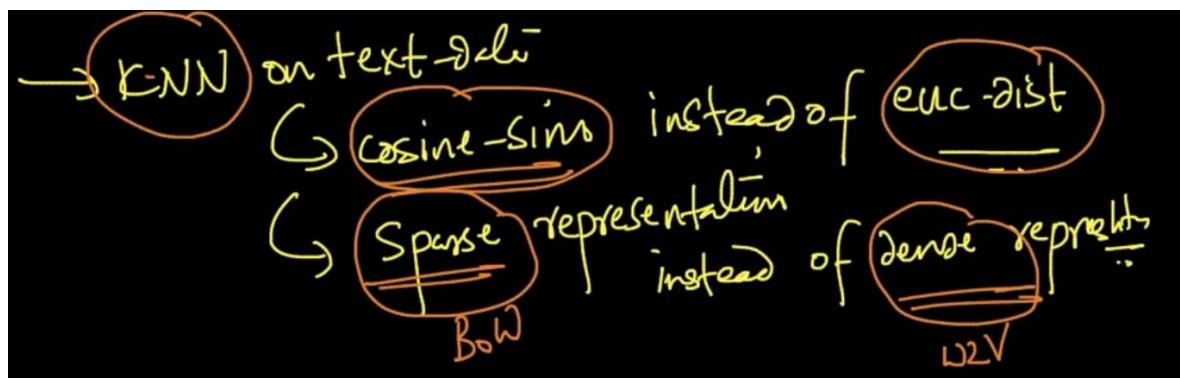
BOW, tf, idf are used since they are sparse

As dimensions become high & sparse the data get's non uniform as 90% of data is 0 so impact of Curse of dimensions become lower as whole intuition of it is dependent on random spread of data



As dimensionality increases overfitting also increases so we use forward feaure selection mostly to select the most important features hereby reducing dimensionality

Dimensionality reduction : PCA, t-SNE doesn't use class label in reducing the dimensions it is only related to maintain the spread /distance of points so not used for classification



Techniques mentioned above are used as impact of Curse of dimensionality is lower on them

[Curse of dimensionality - Wikipedia](#)

BIAS VARIANCE TRADEOFF

Bias-Variance Trade off → statistical ML
↳ Theory of ML

→ Under fitting & overfitting → {Math- basis...}
 $\hat{C}(K=n)$ $\hat{P}(K=1)$

(Ex. regression proof)
generalization error = $\underbrace{\text{Bias}^2}_{\substack{\downarrow \\ \text{future unseen} \\ x_0 \\ \text{data of } M}}$ + $\underbrace{\text{Var}}_{\substack{\downarrow \\ \text{scale}}} + \underbrace{\text{irreducible error}}$

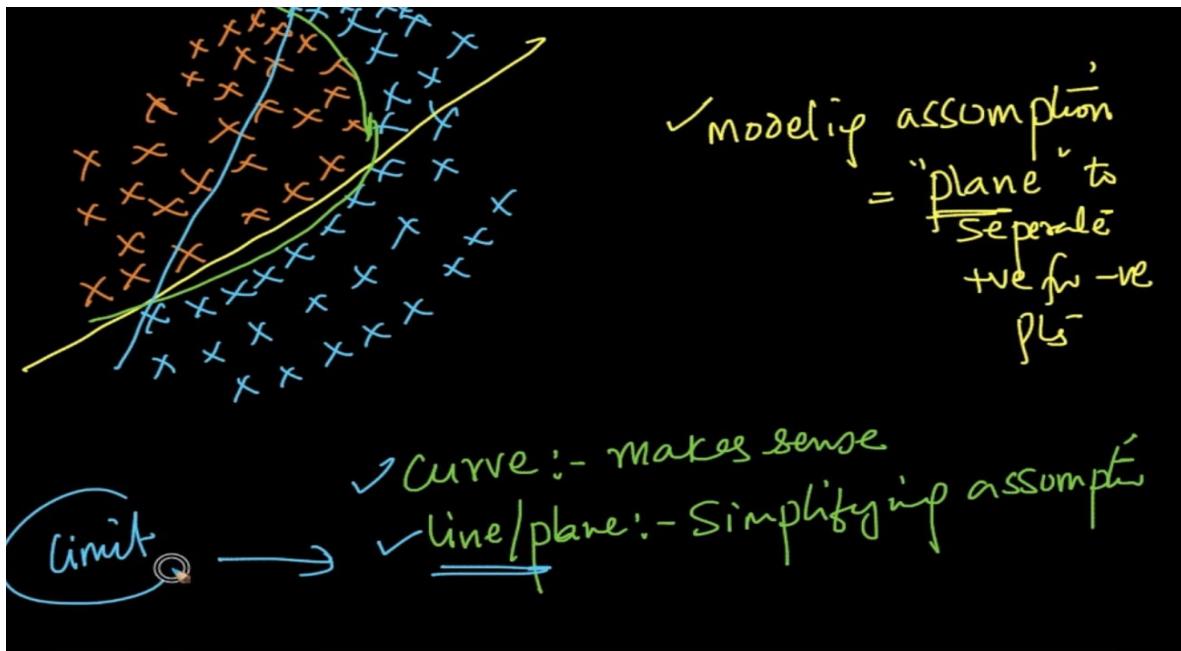
We understood the concept of underfitting & overfitting but now we are doing a Mathematical analysis of it. The answer of how to do it is Bias-Variance Tradeoff.

Generalization error = $Bias^2 + Variance + irreducible error$. It is that error which occurs on

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

unseen data

BIAS

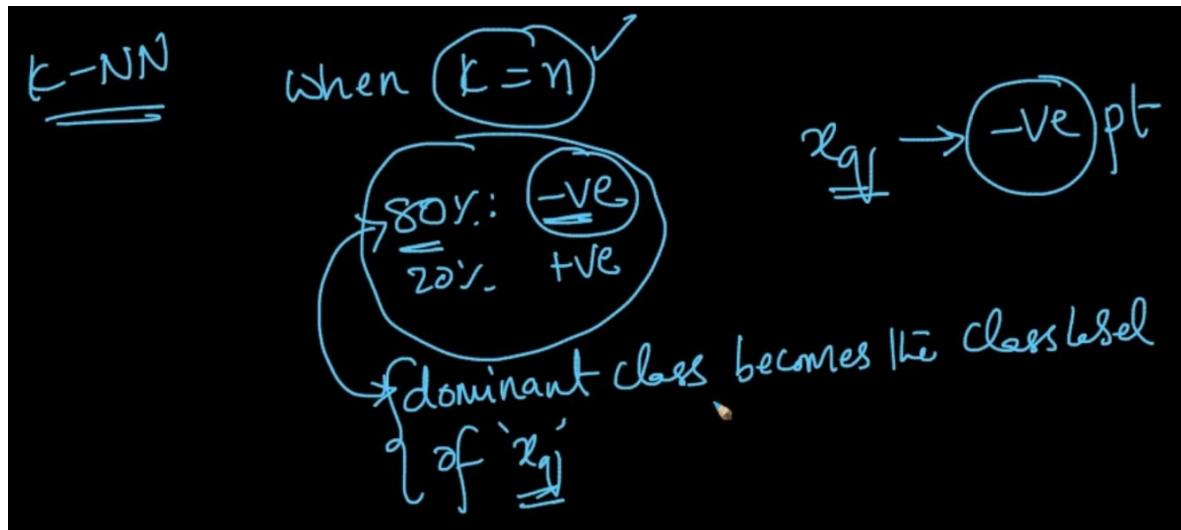


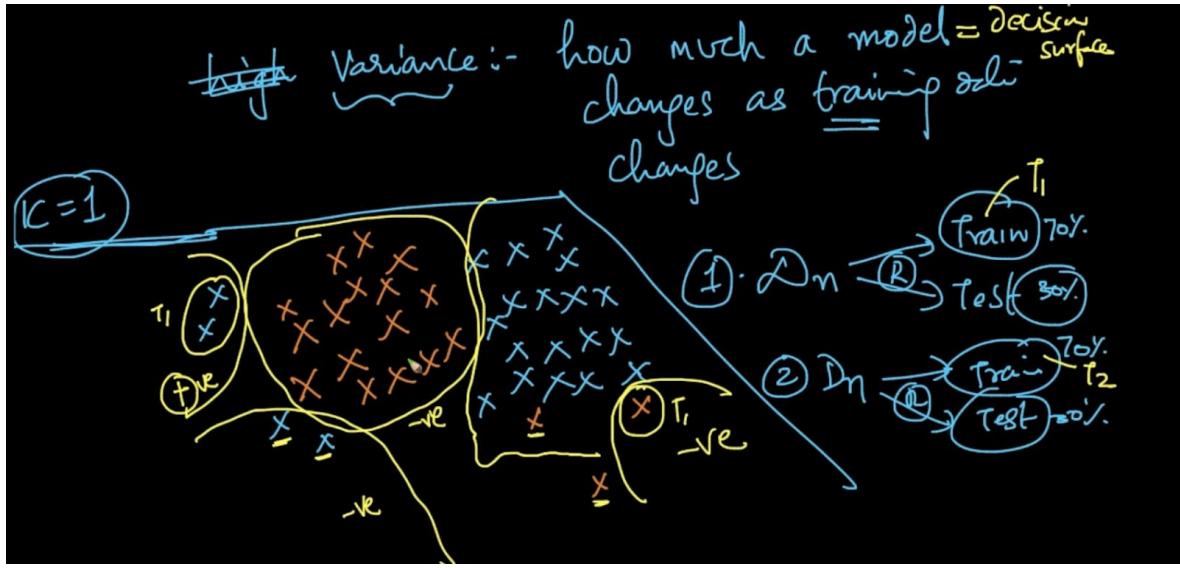
Bias error is error due to simplifying assumption .Suppose we are separating that data with a plane / line instead of curve (which can separate the data better) we are simplifying it and it leads to underfitting

In k-NN where $k = n$, if we make the dominant class as class label of every new query x_q we are doing a simplifying assumption which leads to high bias => underfitting

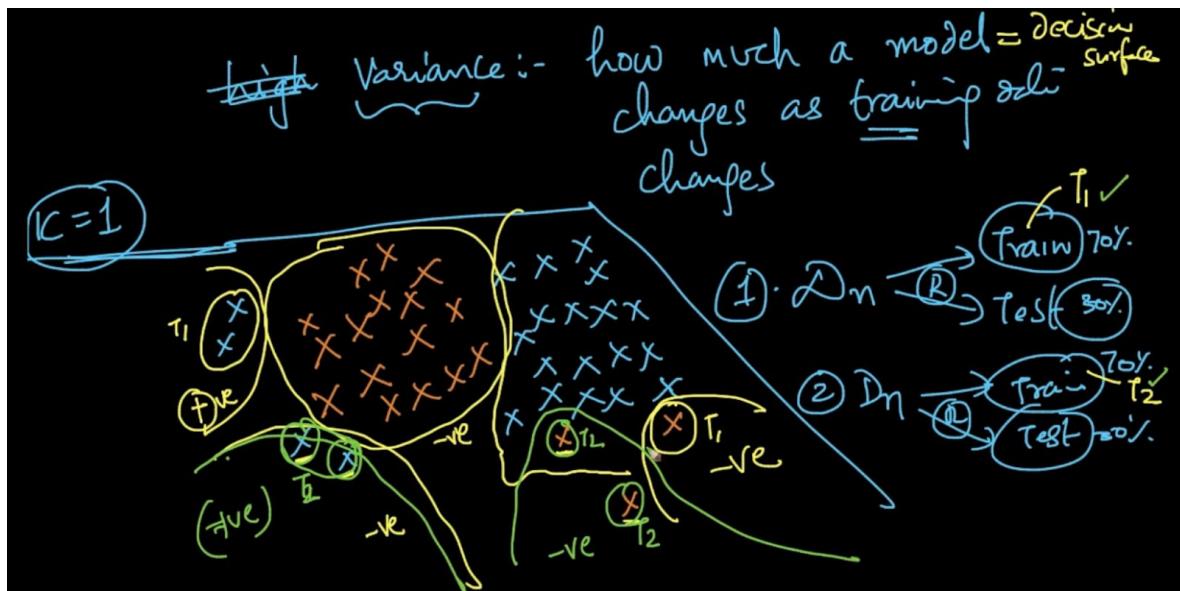
So whenever we've high bias we are underfitting as we are generalizing assumptions

VARIANCE

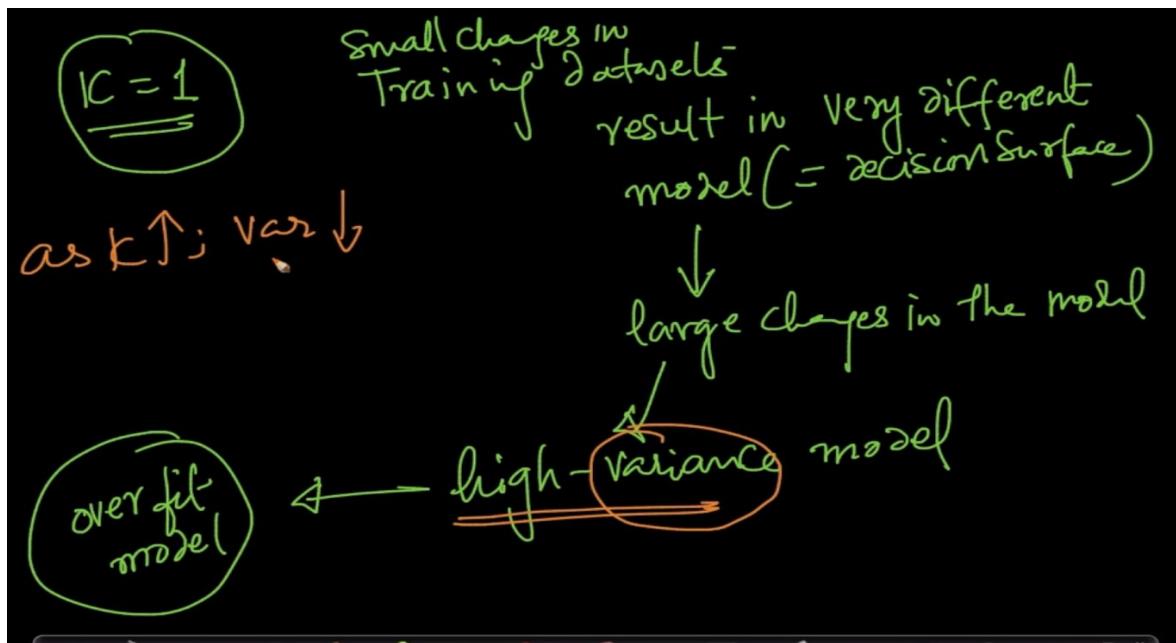




For $k = 1$, this is the decision surface with class labels on first random Data splitting where Train data T_1



As we do the 2nd Random split where Train data T_2 . Here, the green surface where it was first -ve in T_1 becomes +ve in T_2 . So, it's a variance



Small change in Training dataset result in very different model (= decision surface)

This change in decision surface/model is large which can lead to high variance which leads to overfit. As $k \uparrow$; $\text{var} \downarrow$

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

$k = 1 \rightarrow$ $\text{Bias} \downarrow$ $\text{Var} \uparrow$
 $k \uparrow \rightarrow$ $\text{Bias} \uparrow \text{slightly}$ $\text{Var} \downarrow \text{drastically}$
 $k \rightarrow n \rightarrow$ $\text{Bias} \uparrow$ $\text{Var} \downarrow$

So we want to balance bias error and variance error as seen that they vary on the values of their 'k'

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

let $k = 1; j$

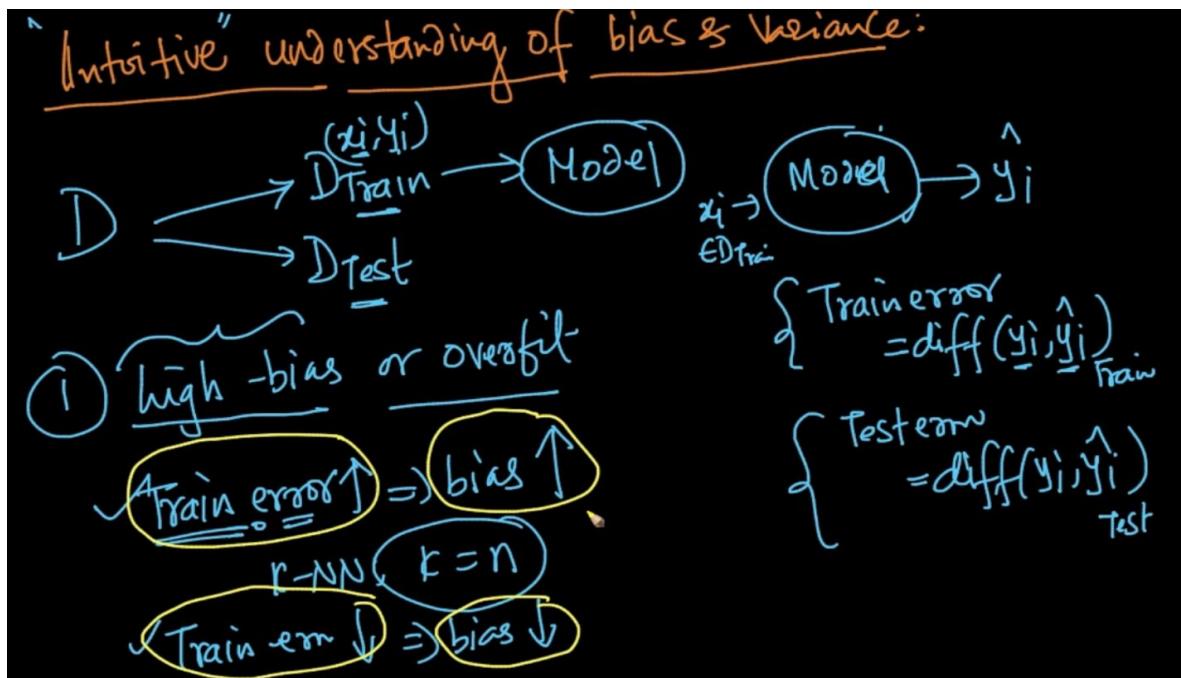
$10 + 100 + 3 \rightarrow 113 \rightarrow \text{overfit}$

$12 + 10 + 3 \rightarrow 25 \checkmark$

$100 + 2 + 3 \rightarrow 105 \text{ by underfit}$

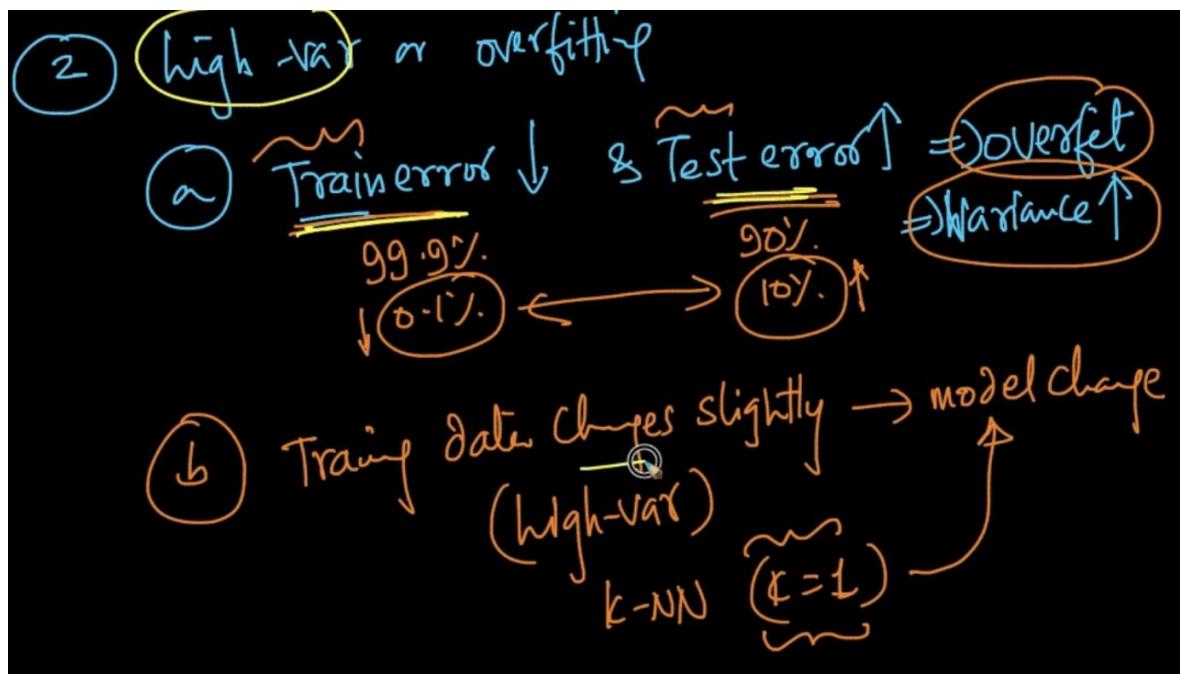
Here at $k = 5$, Generalization error is the lowest so we should be looking for finding the right balance

INTUITIVE UNDERSTANDING OF BIAS AND VARIANCE



Our model is given x_i and model gives \hat{y}_i and \hat{y}_i difference is tested with y_i which gives Train error and Test error.

High bias/ Underfit : A way to determine it is if Train error is high then bias is also high



High Variance / Overfit : if our Train error is low but Test error is high then we say that our model is overfit or has high variance