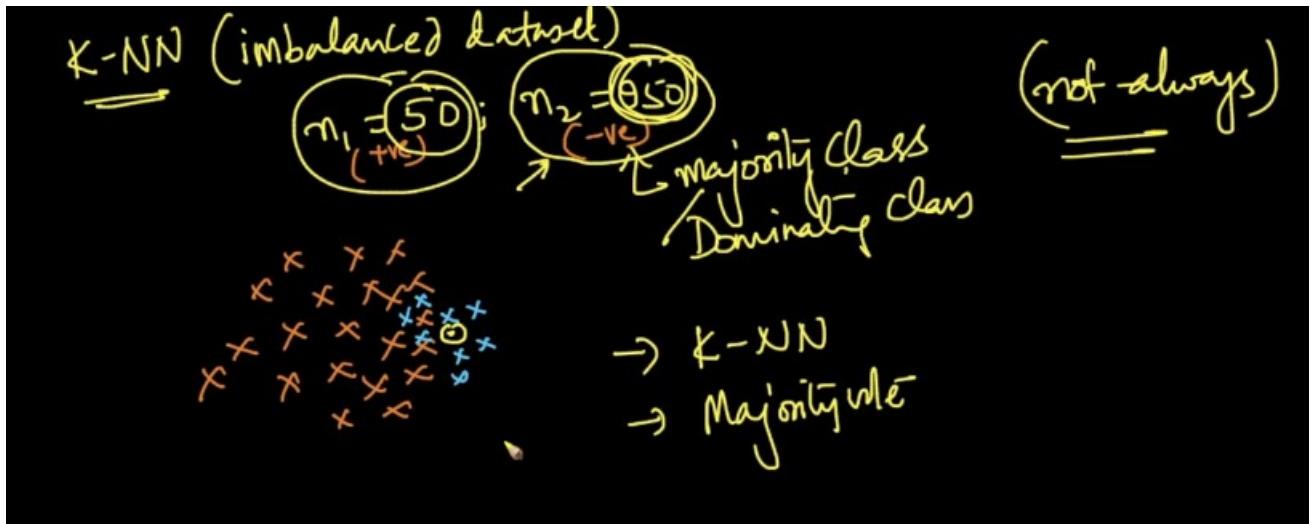


## NOTES

Balanced vs. Imbalanced data set:

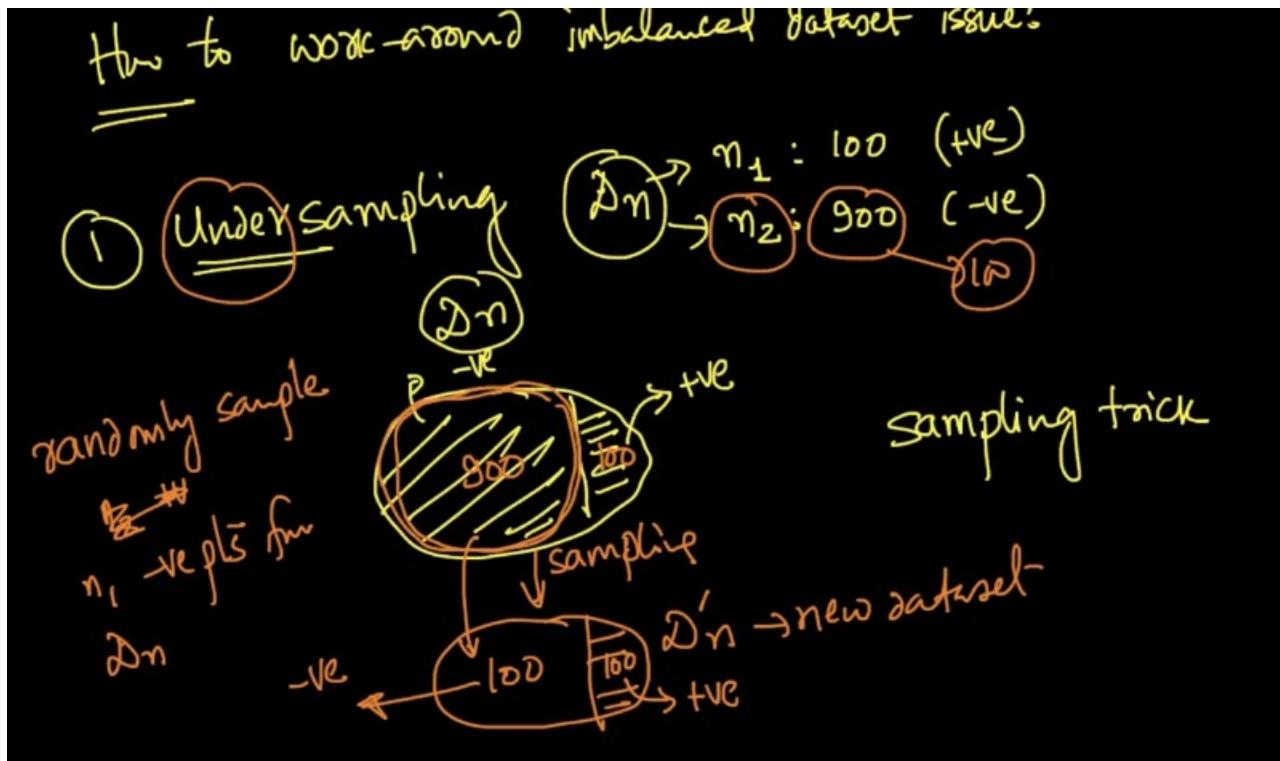
Imbalanced data-sets:

In case of imbalanced data sets majority class will be dominate over the minority class. In case of KNN the majority will get high priority. This may occur.



Resolving imbalanced data sets:

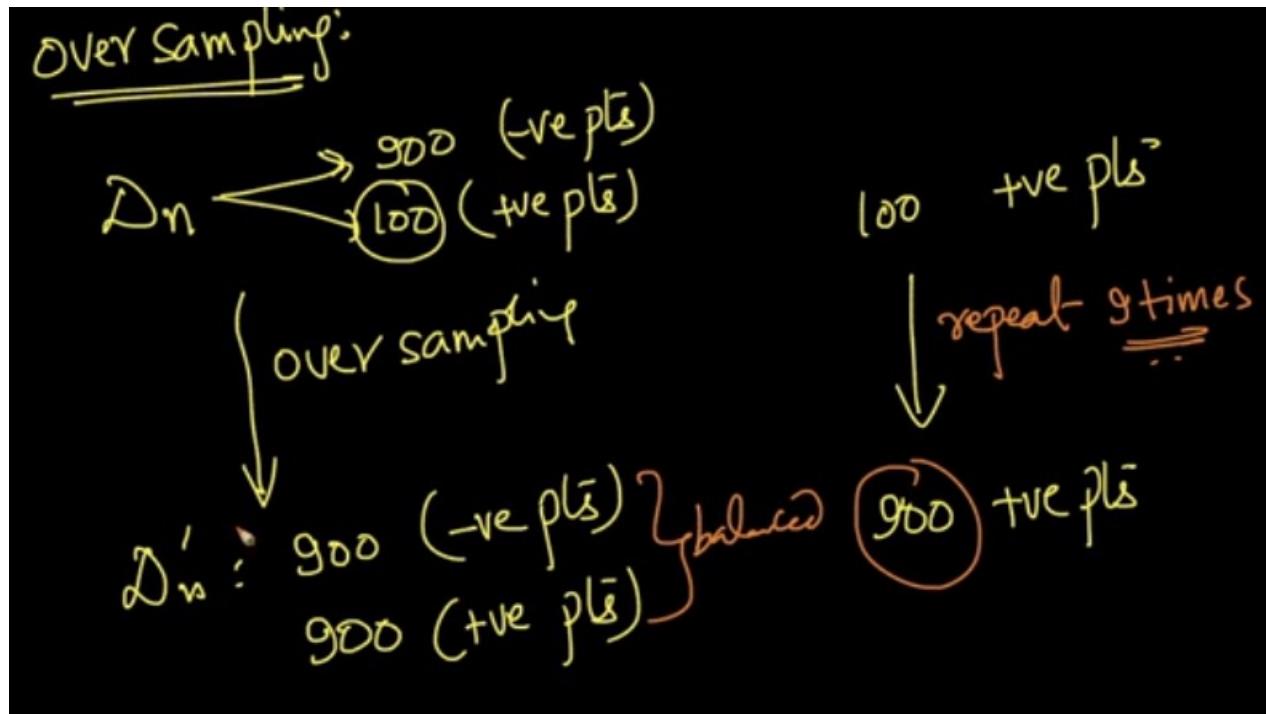
1. Under sampling.
2. Over sampling.



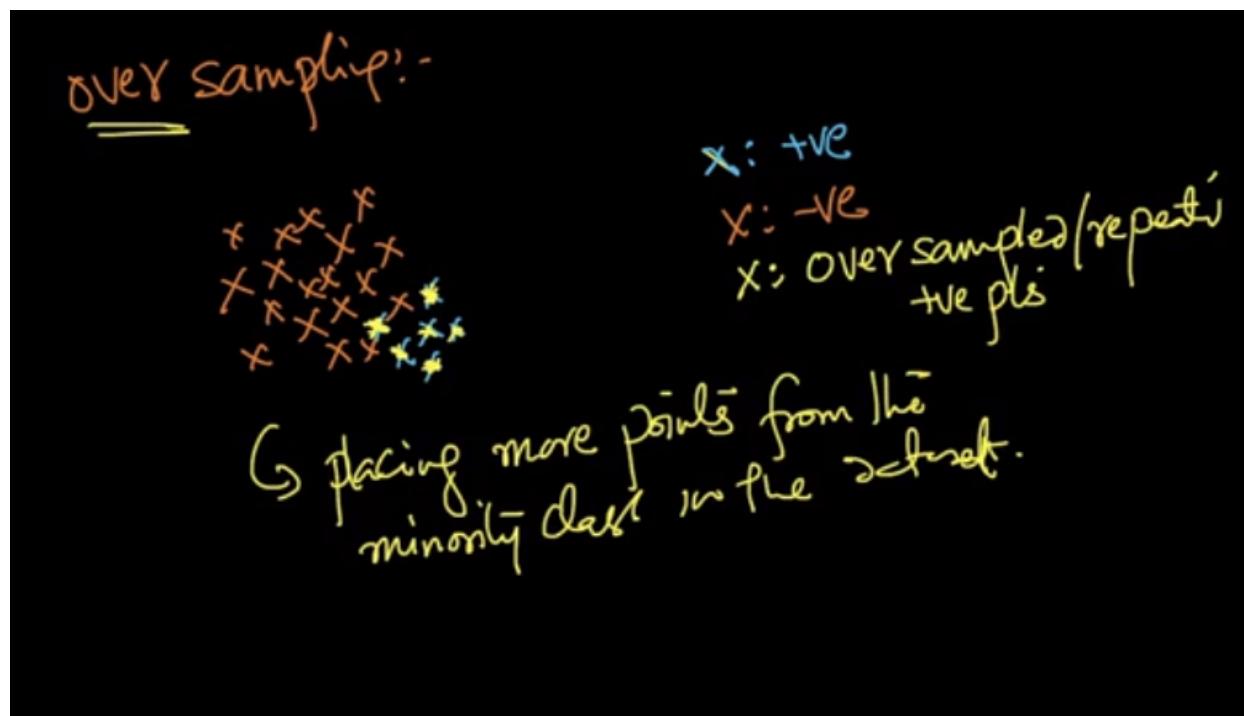
The sampling must be done **randomly**.

There is data loss in case of under sampling. We can overcome this situation by Oversampling the data.

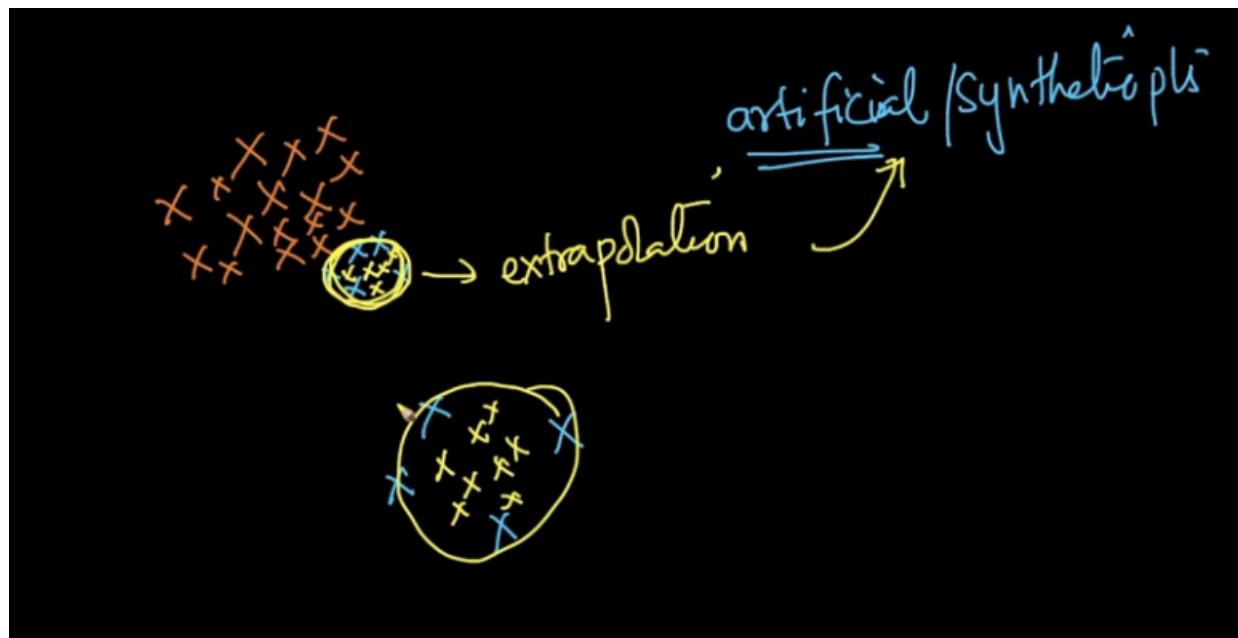
Oversampling is taking the same data point repeatedly.



The data points lie on top of each other.

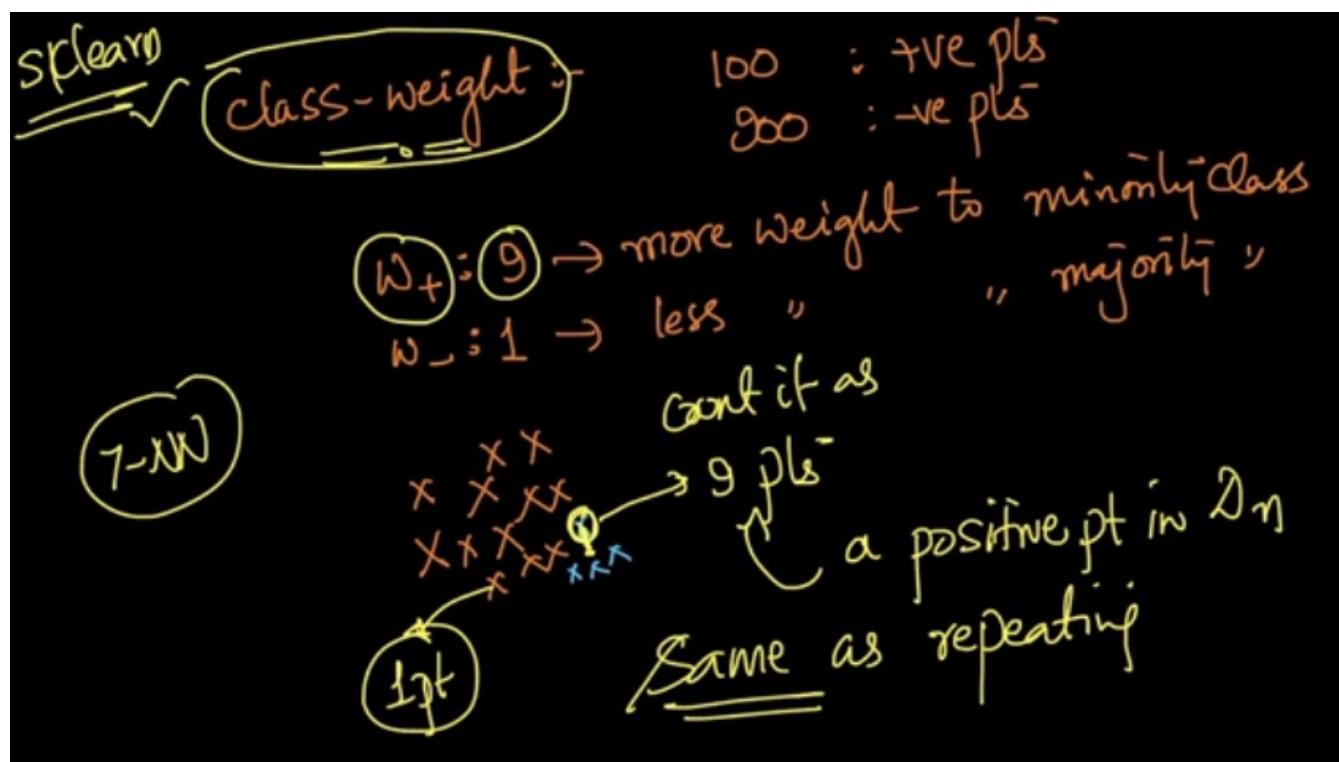


Making synthetic points is the other alternative to oversampling. There artificial points are created with in the same class points.

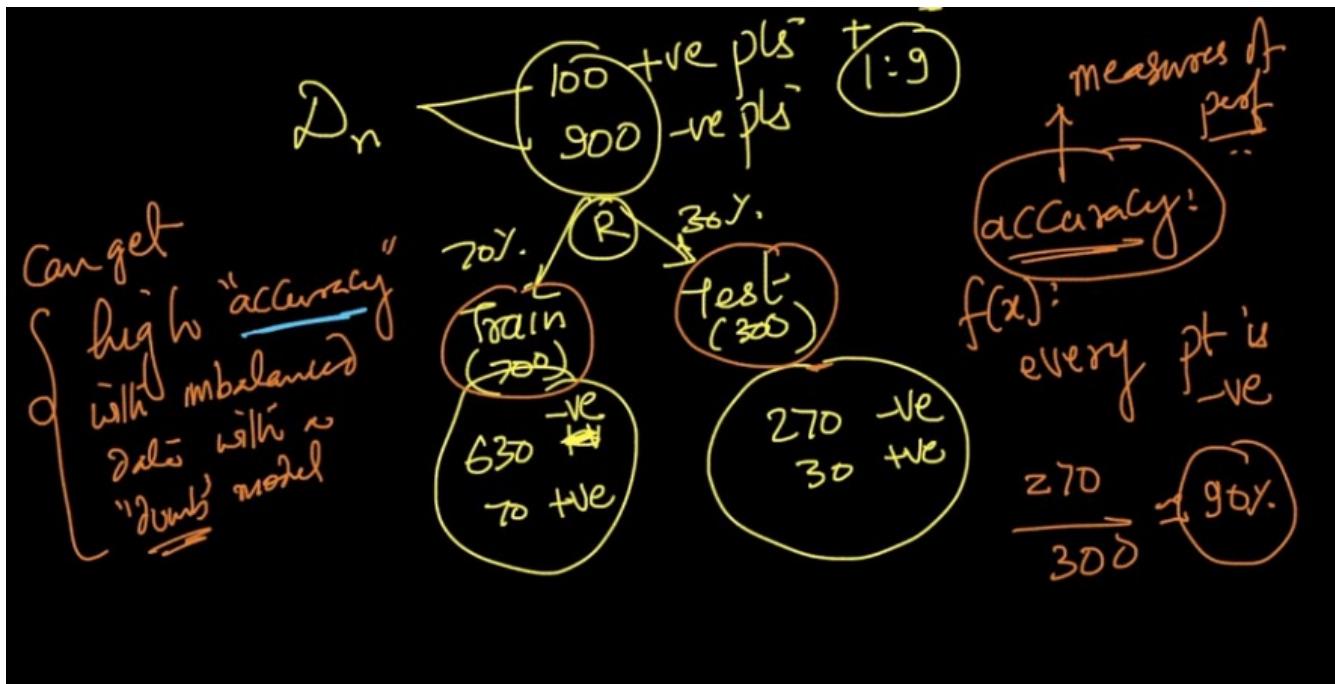


Implementation of class – weight for oversampling:

The more weight is given to the minority class the less weight to the majority class. This is same as repeating the points.



In case of Imbalanced data sets we can get high accuracy with a dumb model.

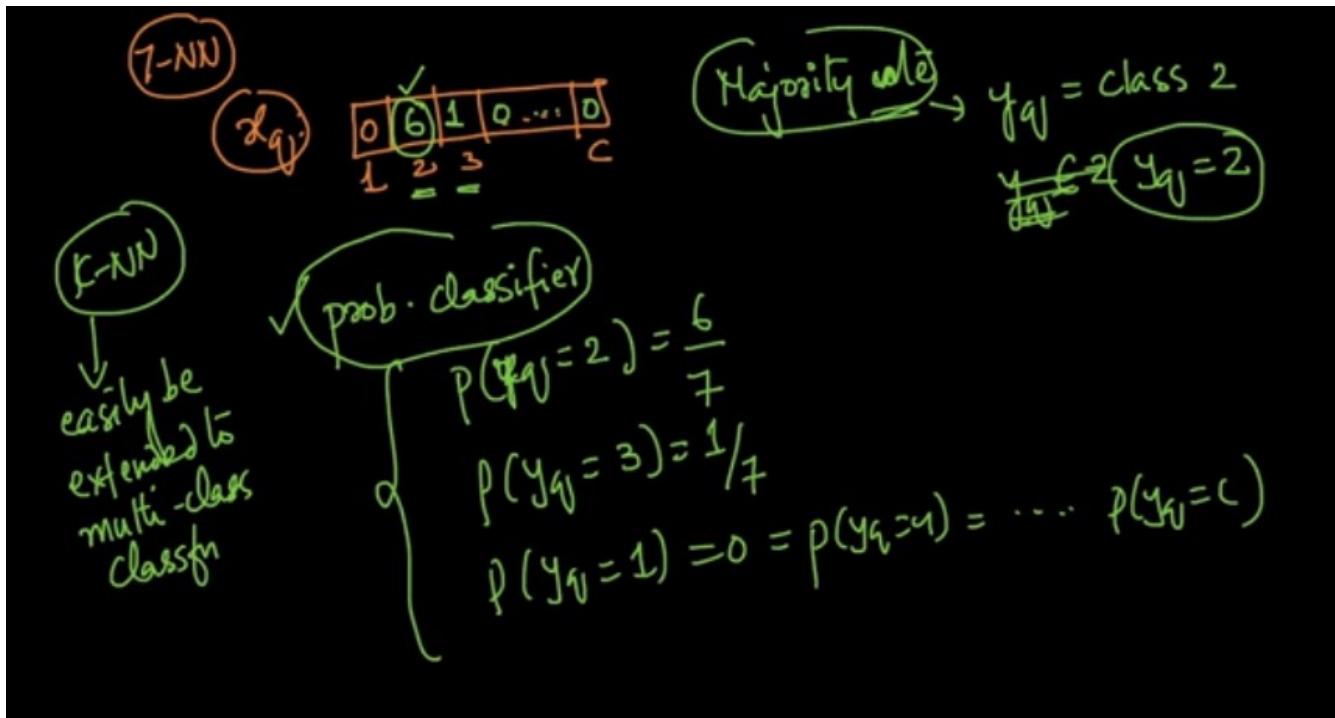


Trying oversampling is the good practical strategy.

Multi – class classification:

Example:

K – nearest neighbors: KNN can be extended to multi – class classification by using probability scores.



Logistic regression can only do binary classification, that can be extended to multi – class by using soft - max classification.

$\checkmark \text{ Logistic regresn } \rightarrow \text{ binary classifier}$

cannot do multiclass-classfn. as easily as K-NN

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

$f(x) \rightarrow \begin{cases} 0 \\ 1 \end{cases}$  (binary)

$f'(x) \begin{cases} 0 \\ 1 \\ 2 \\ \vdots \\ c-1 \end{cases}$  c-classes

Breaking the data – set to small parts with each class. If we have ‘c’ classes we can break into ‘c’ binary classification problems.

$y_i \in \{0, 1, 2, \dots, c\} \rightarrow c \text{ binary classifiers}$

$D_n$

①  $D_n \rightarrow \{(x_i, y_i) | y_i = 1\} \rightarrow \text{tve} \quad (\text{Binary classifier})_1 \text{ or not}$

②  $D_n \rightarrow \{(x_i, y_i) | y_i \neq 1\} \rightarrow \text{-ve}$

③  $D_n \rightarrow \{(x_i, y_i) | y_i = 2\} \rightarrow \text{Binary classifier}_2 \text{ or not}$

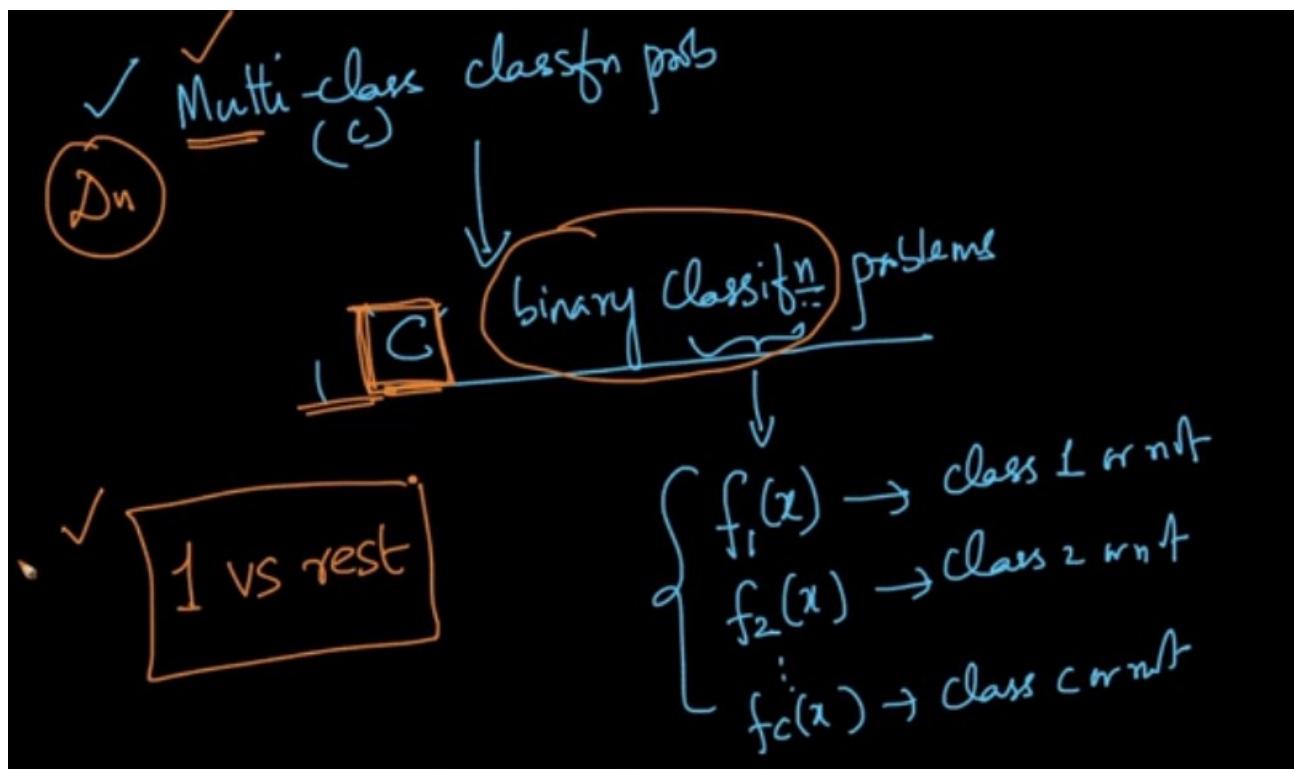
④  $D_n \rightarrow \{(x_i, y_i) | y_i \neq 2\} \rightarrow \text{-}$

⑤  $D_n \rightarrow \{(x_i, y_i) | y_i = 3\} \rightarrow \text{Binary classifier}_3 \text{ or not}$

⋮

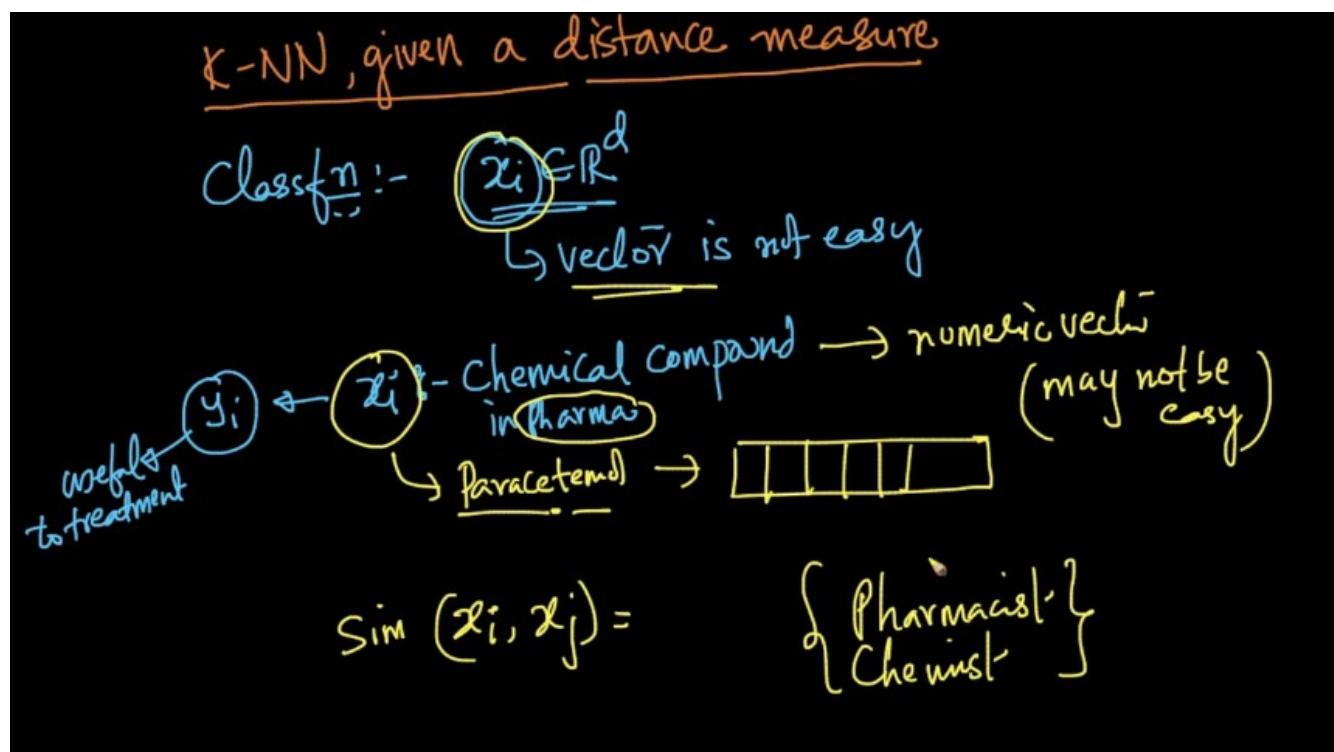
c or not

We must train 'c' classifiers in this case.  
This solution is called one vs. rest solution.



Distance similarity matrix:

Example: Chemical compounds cannot be easily expressed in the form of numeric values.



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

$n$  data points  
Lok

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

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

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

Distance matrix is calculated by taking inverse of the similarity value of each feature.

Given Sim-matrix ( $S$ ) or dist-matrix ( $d$ )  
✓ instead of  $x_i \in \mathbb{R}^d$  explicitly ← Logistic regn

K-NN  
works well given  
 $S$  or  $d$



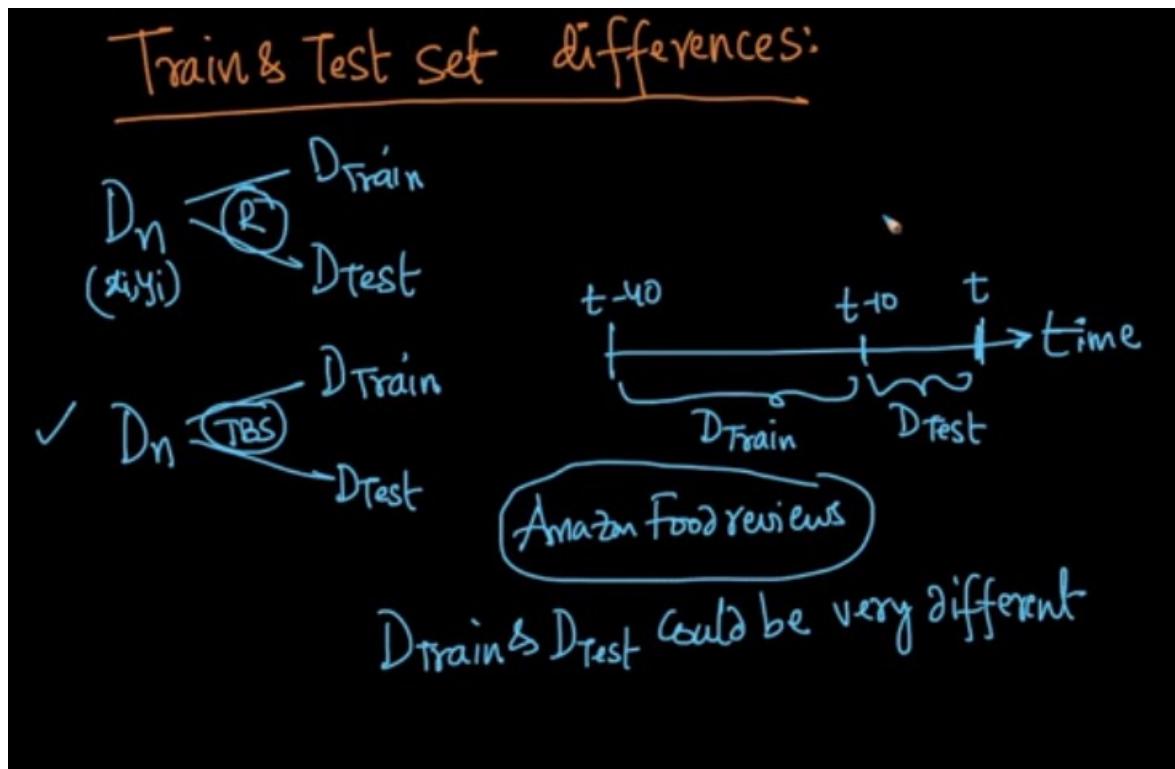
$$\begin{cases} d(x_i, x_1) = \\ d(x_i, x_2) = \\ \vdots \\ d(x_i, x_n) = \end{cases}$$

$$\begin{matrix} & x_1 & x_2 & \dots & x_n \\ x_i & \begin{array}{|c|c|c|c|} \hline & \checkmark & - & \checkmark & - \\ \hline & - & \checkmark & - & \checkmark \\ \hline & \checkmark & - & \checkmark & - \\ \hline & - & \checkmark & - & \checkmark \\ \hline & \vdots & & & \vdots \\ \hline \end{array} & \end{matrix} \leftarrow \text{dist-matrix}$$

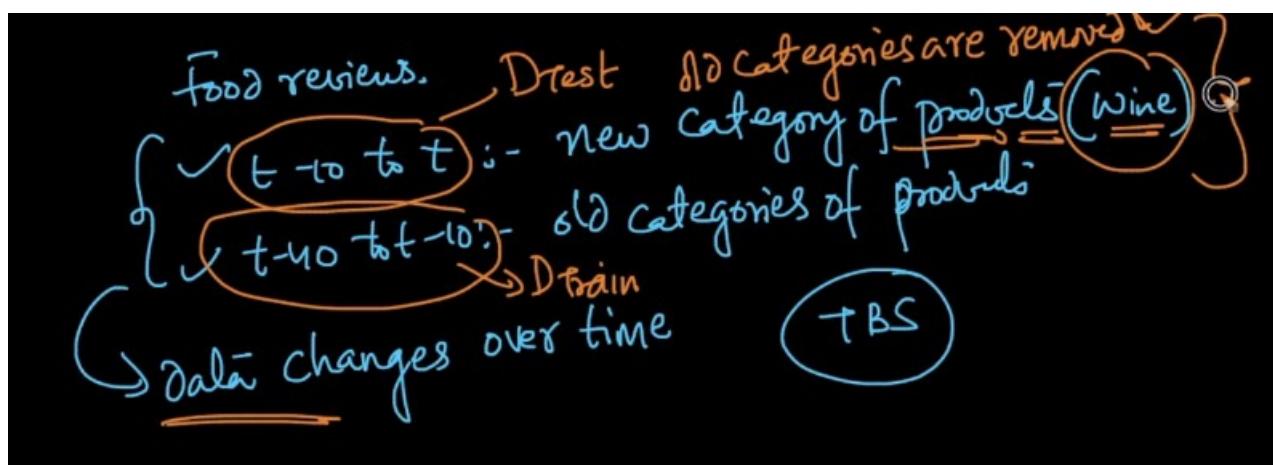
Train – test set differences:

1. Random splitting.

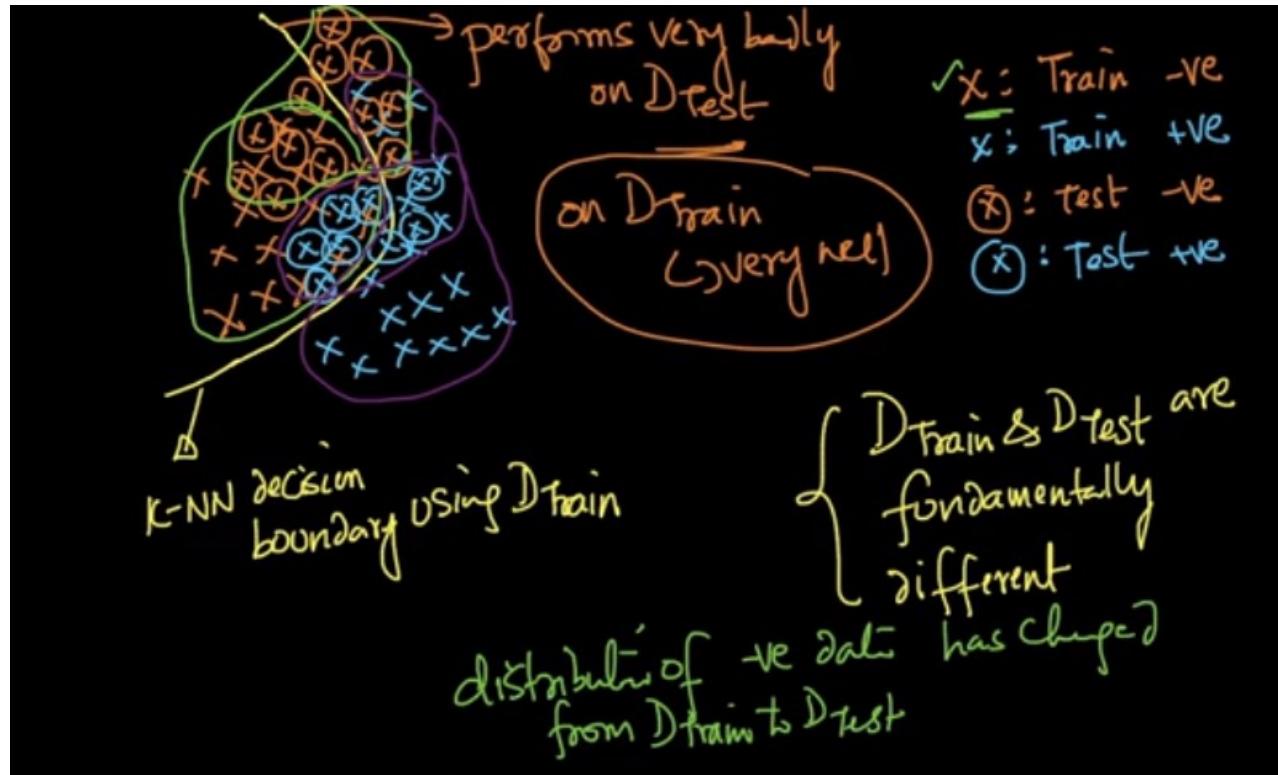
2. Time – based splitting.



The underlying data changes over time then choosing time – based slicing is the good choice of splitting the data.

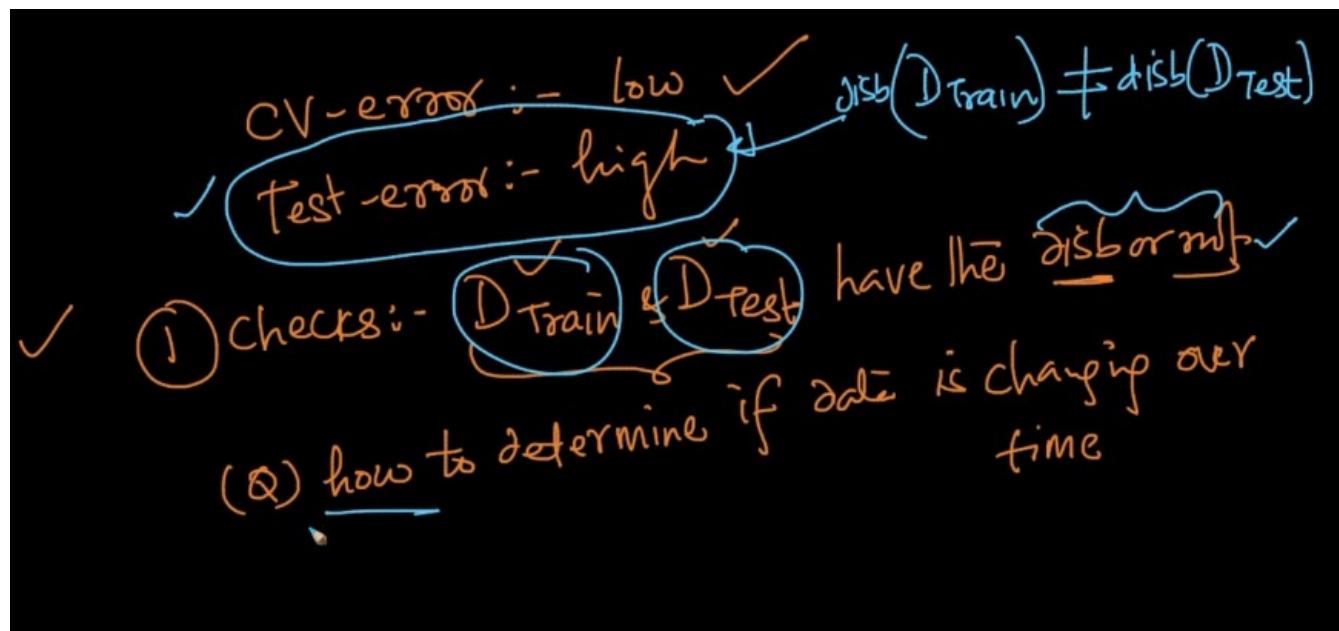


There will be miss match of the train and test data sets. The distribution of train and test data-sets changes.



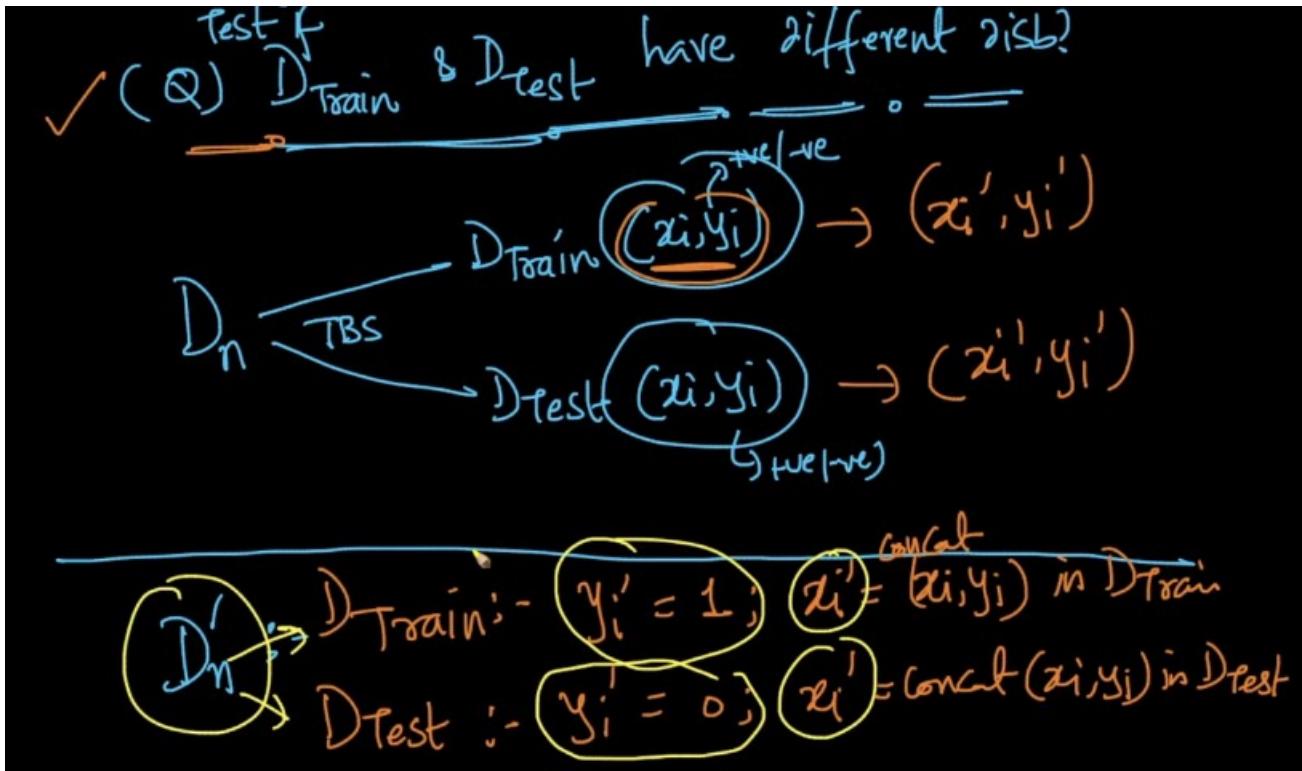
The cross-validation error will be low but the test accuracy will be low.

The distribution of the train and test must be checked.



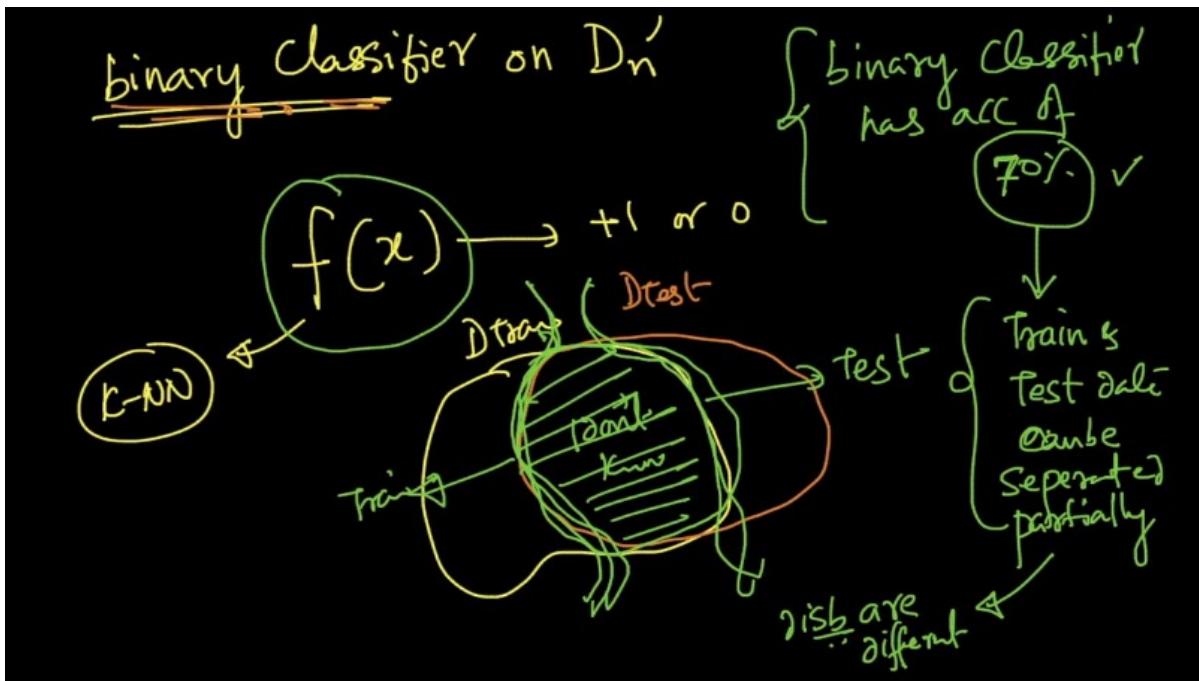
Q) Determining the train and test distributions.

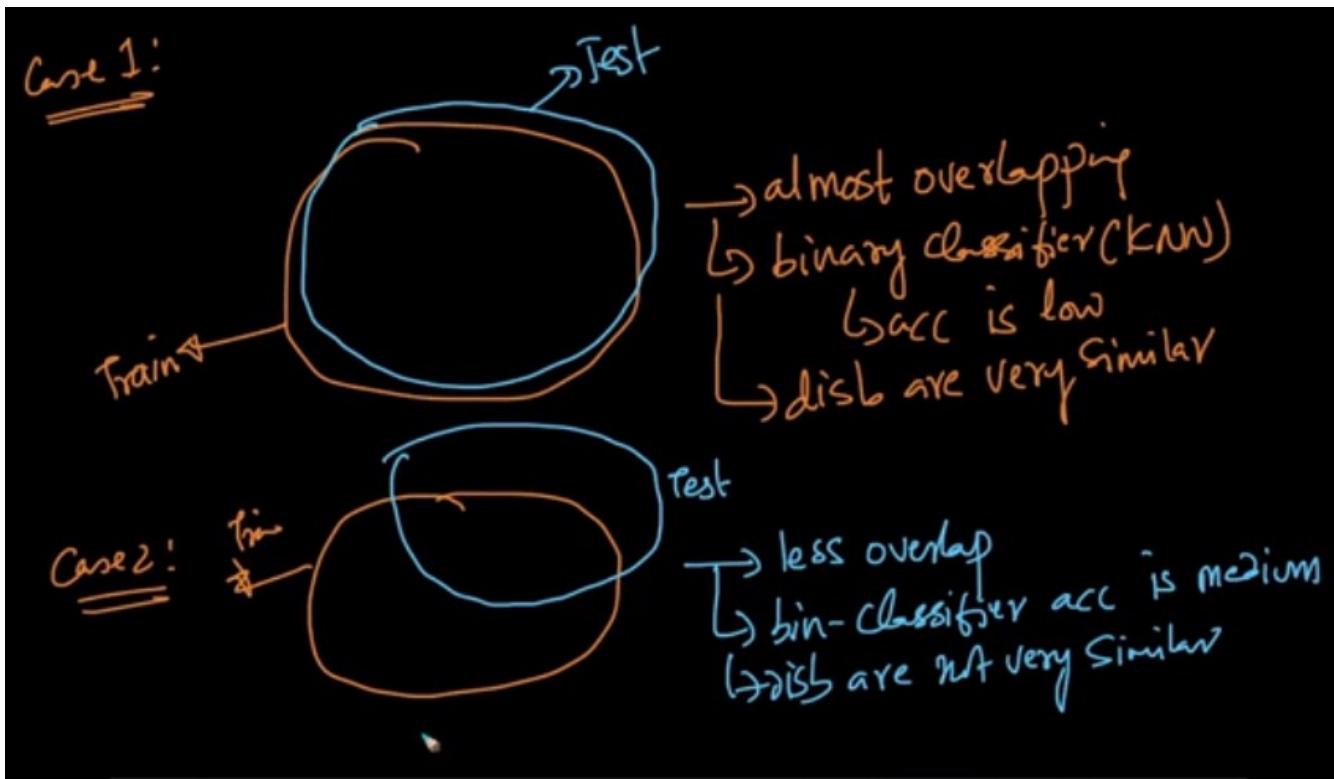
We can use a classifier to determine the distribution of train and test data of the new data-set and the old data-set.



Build the binary classifier on the new data set.

The accuracy of the model gives the differences of the train and test data-sets.

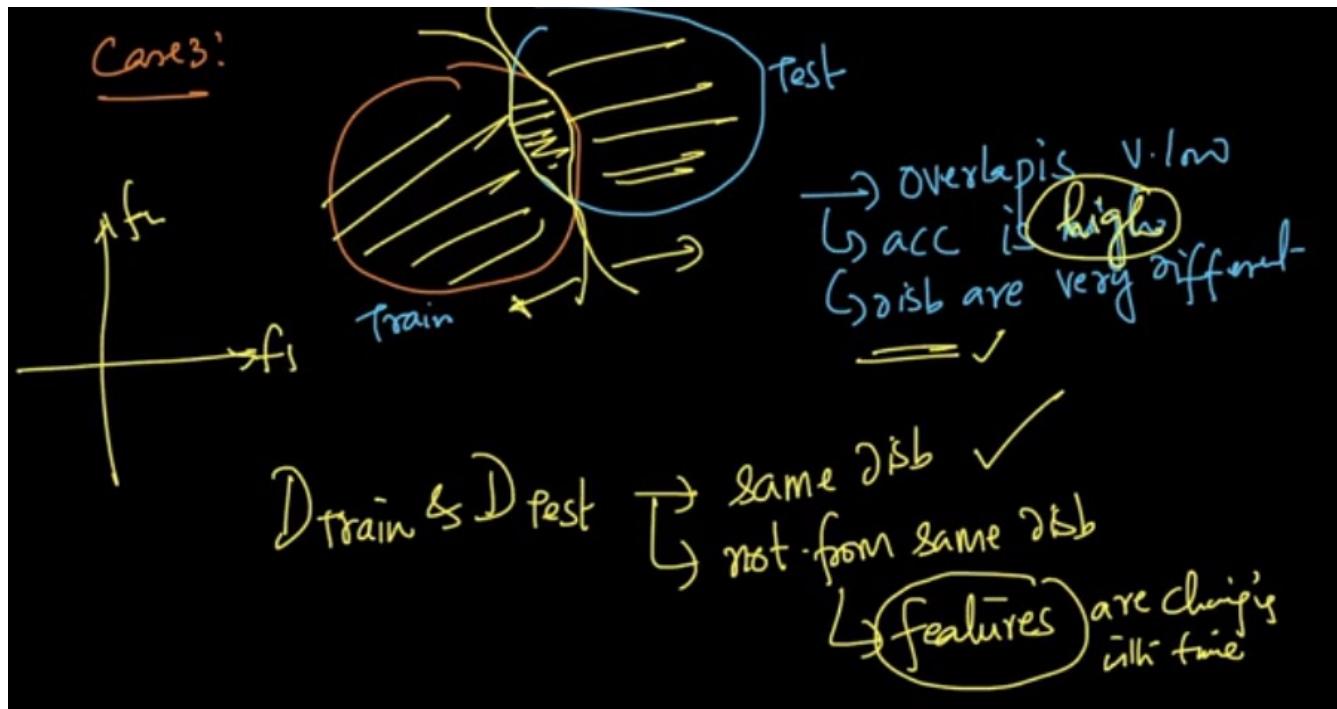




We can use any **binary classifier**.

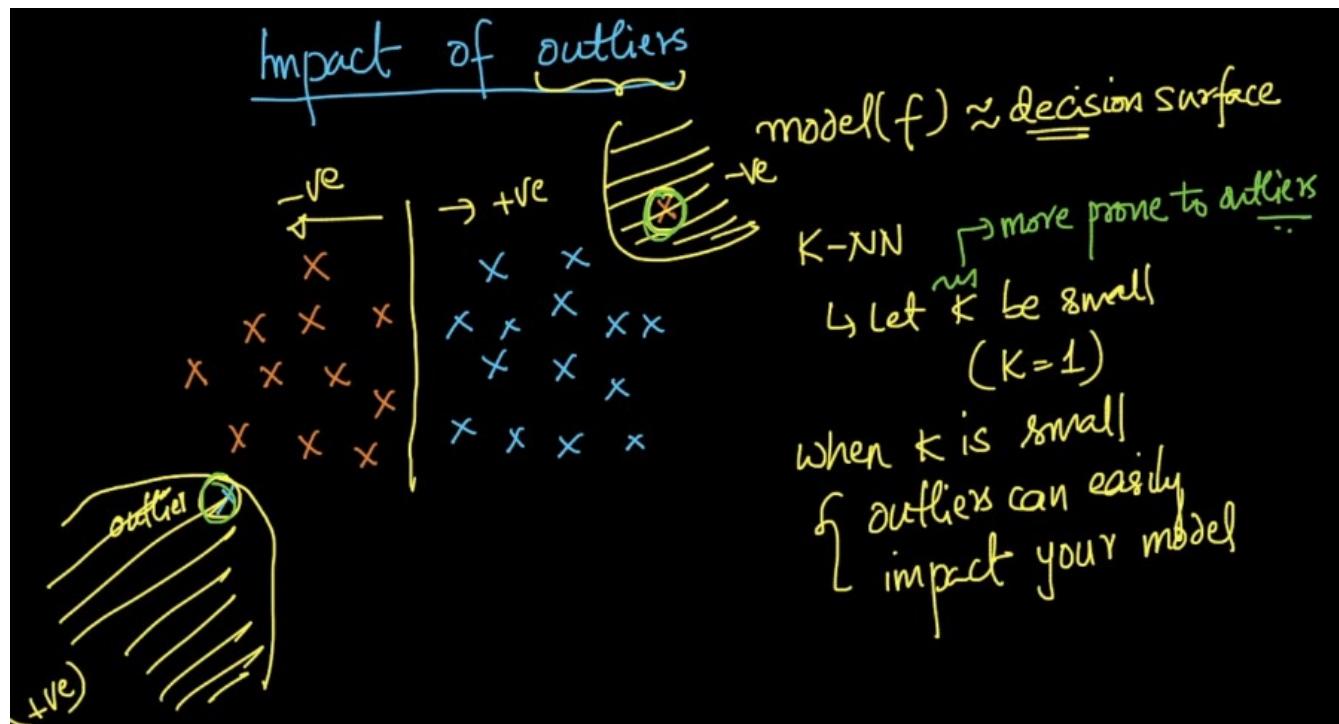
This accuracy will be low as the distributions are very different.

If the data is not coming from the same distribution then feature space is changing, building the new feature will be the way.

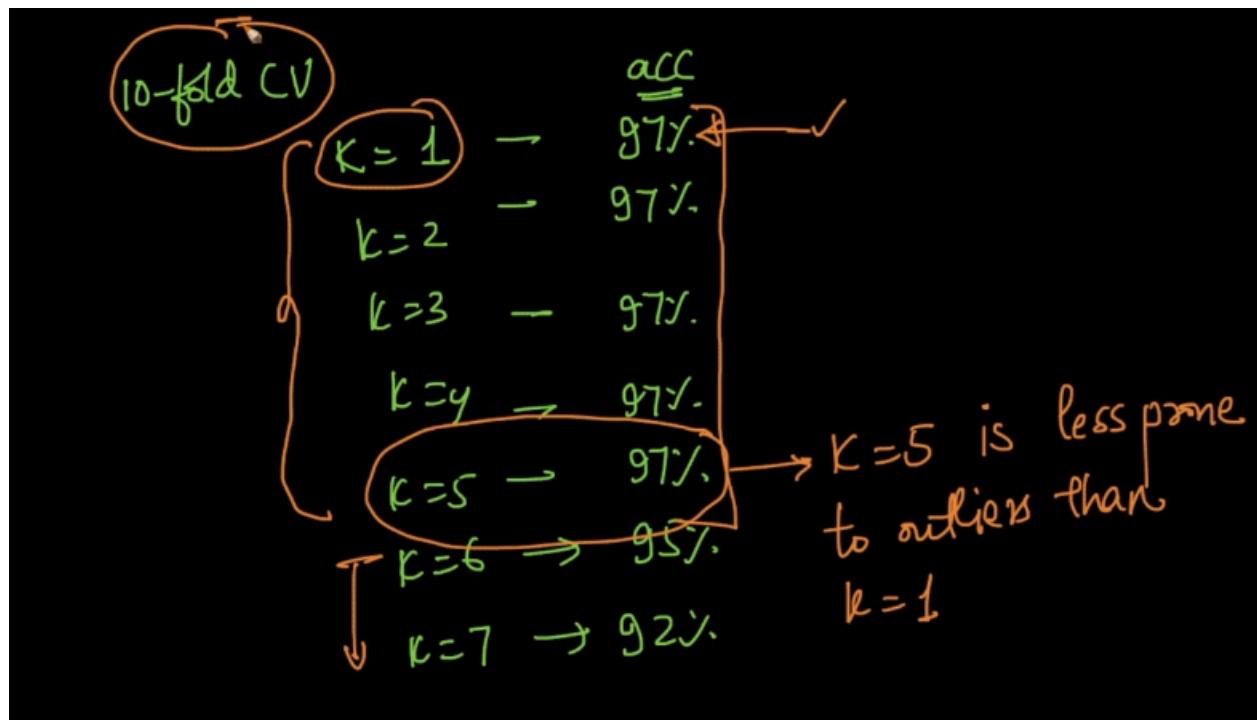


## Impact of outliers:

In case of KNN, when K is small outliers can easily impact the model.



It is always to pick the greater value of L, Since the model is less prone to outliers.



Local outlier – factor is used to detect outliers.

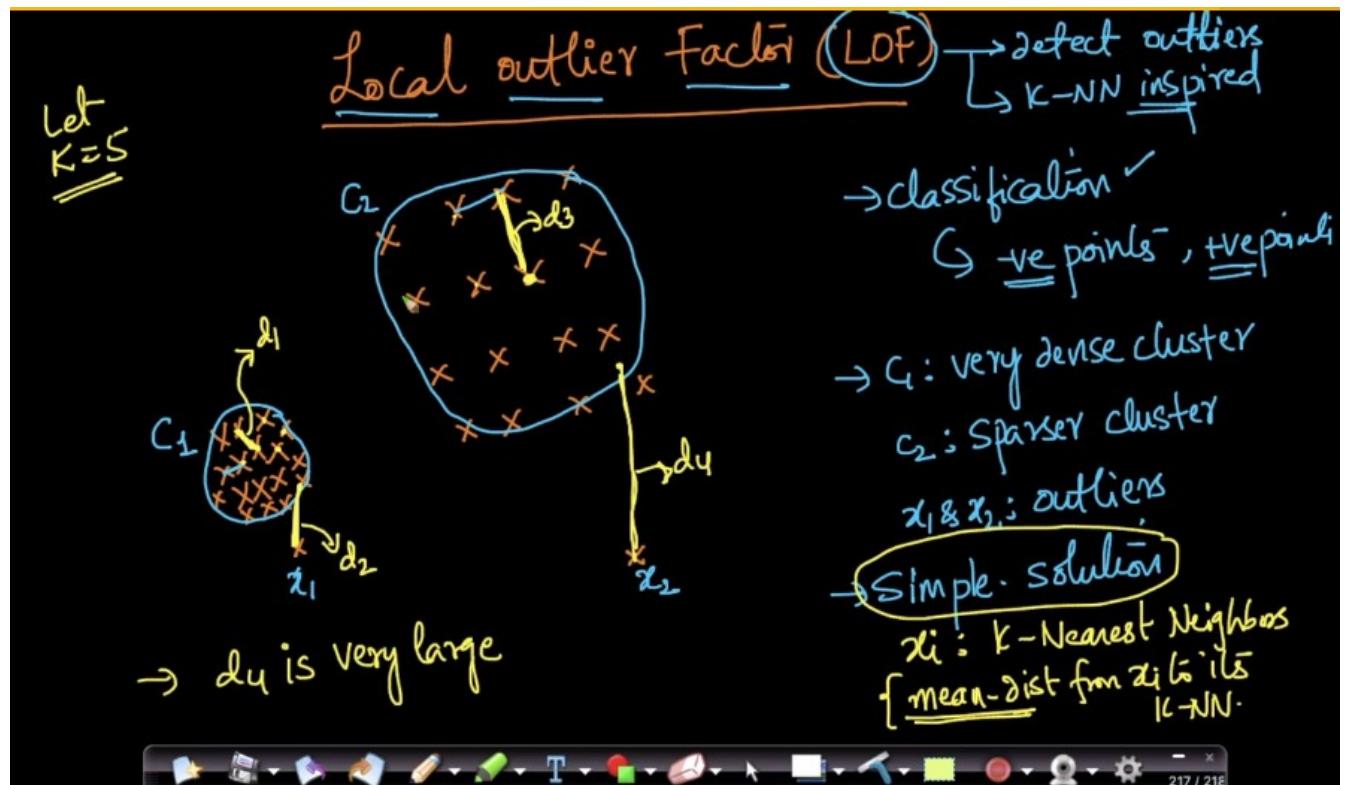
Local – outlier factor:

Mean distance to KNN:

It is being inspired by KNN.

Detect outliers:

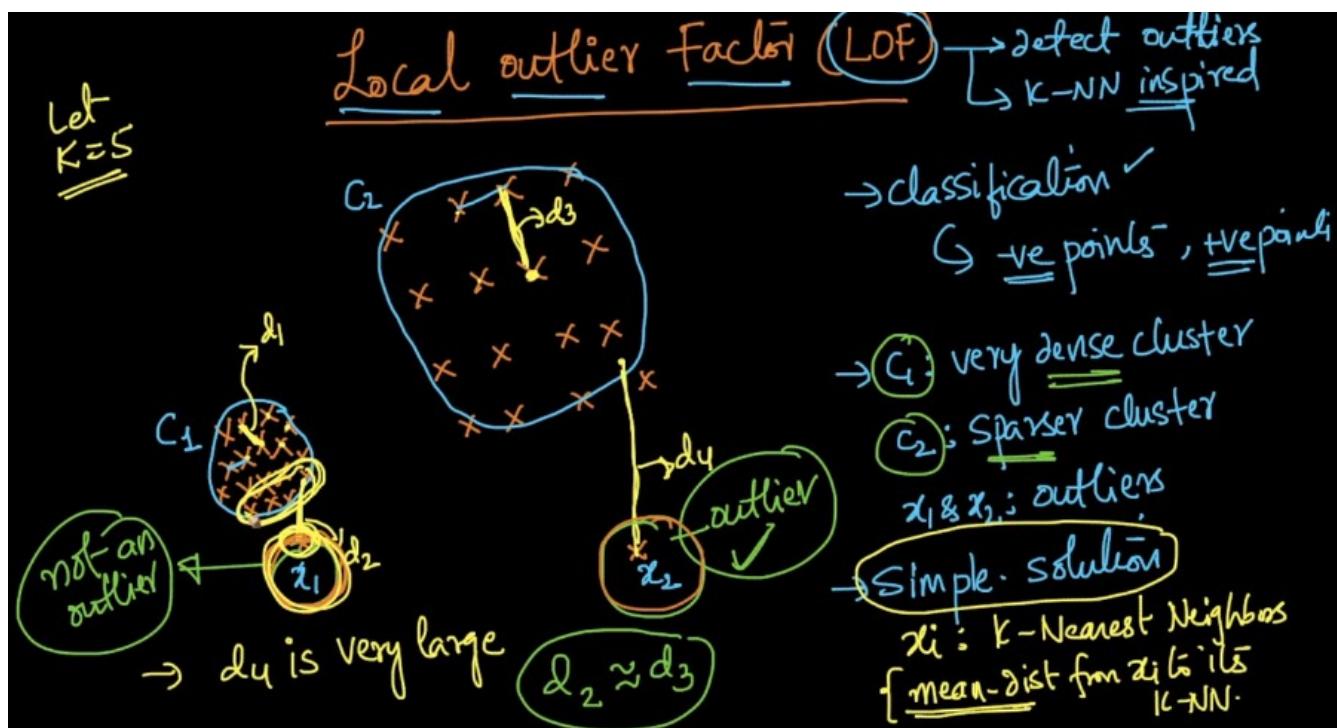
Consider the two clusters of points C1 and C2, dense and sparse respectively. Both X1 and X2 are outliers.



- 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

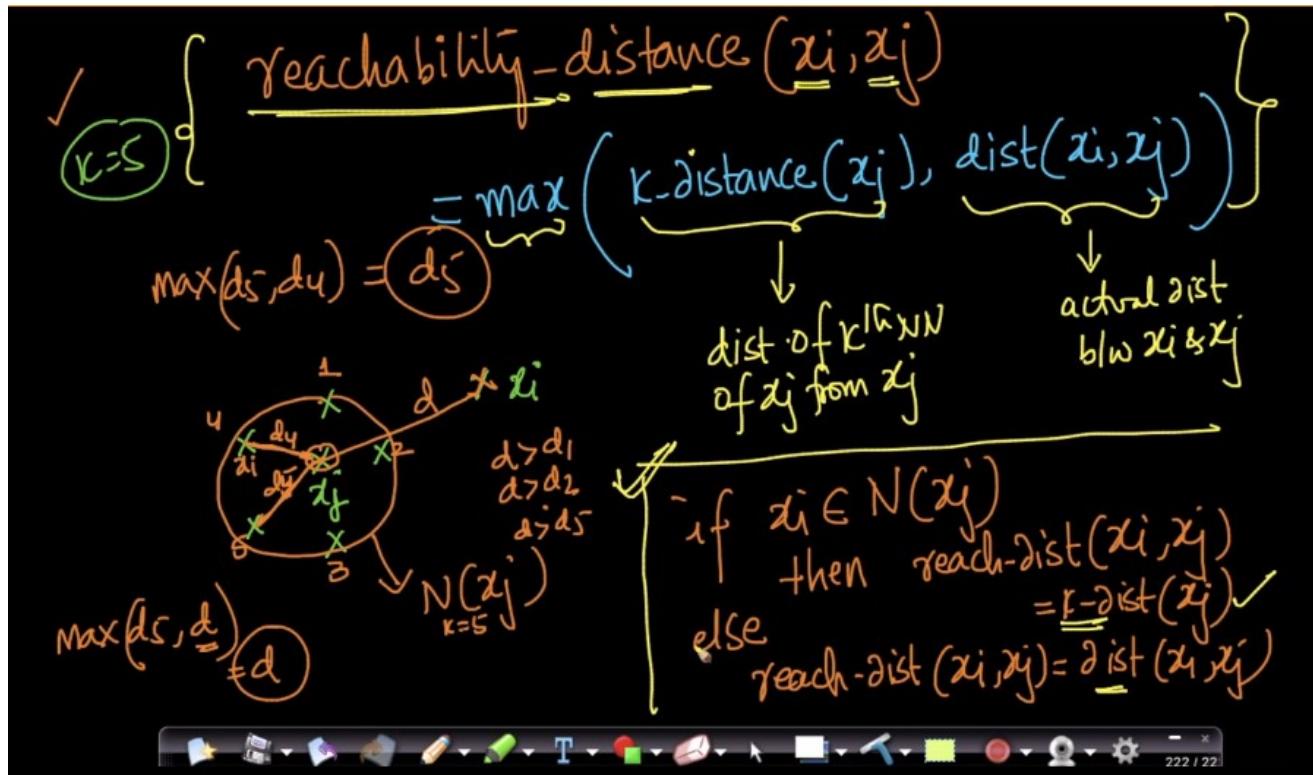
By this process we cannot conclude  $X_2$  as an outlier. Though it is an outlier to the dense cluster.

Therefore this situation leads to computing local density.

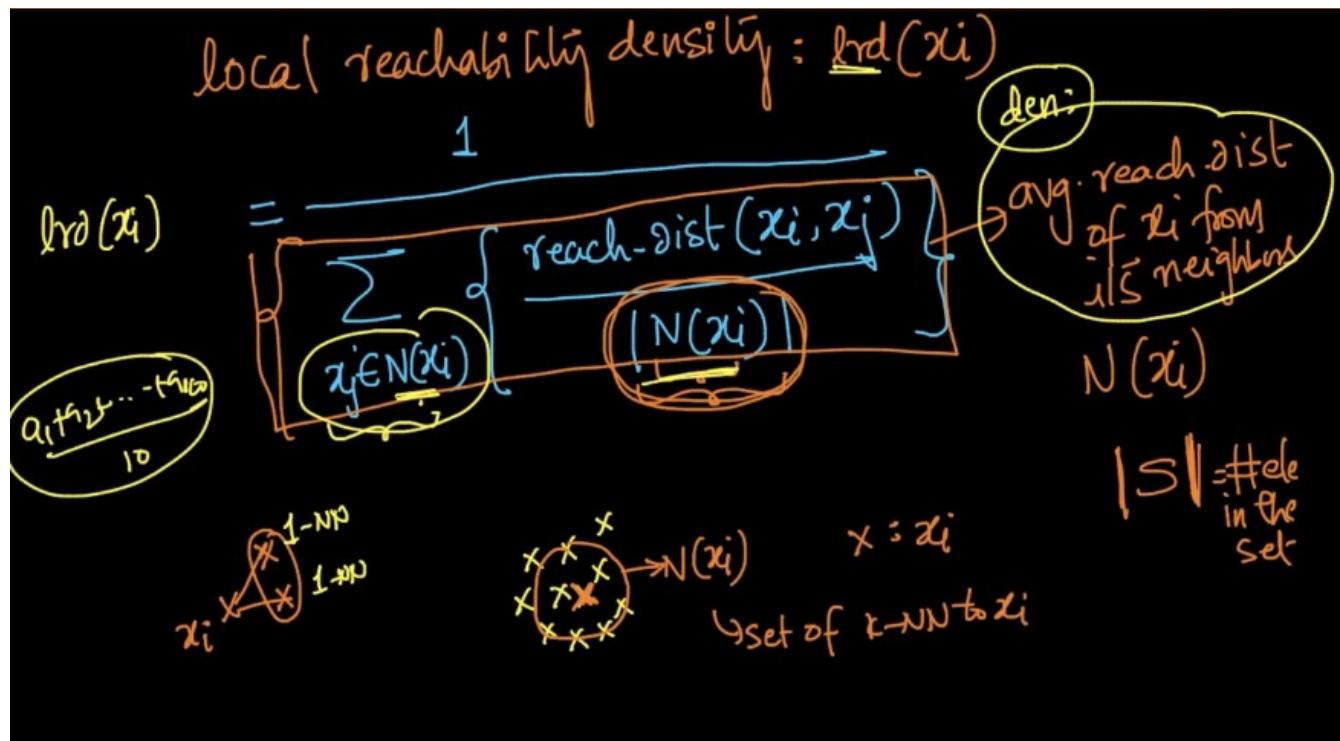


Reachability – Distance(A, B):

The reachability distance of two points  $x_1$  and  $x_j$  is given as follows:



Local reach-ability – density(A):



eqn - english

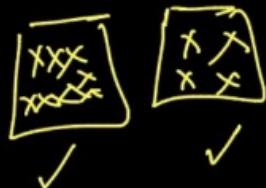
$\{ \text{lrd}(x_i) = \text{inverse of avg. reach-dist of } x_i \text{ from its neighbors} \}$

$\{ \text{reach-dist}(x_i) = \begin{cases} \text{if } \dots \\ \text{else } = \dots \end{cases} \}$

Reach-ability distance lead to local reach-ability – density.

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

local-reachability-density



Local outlier Factor:

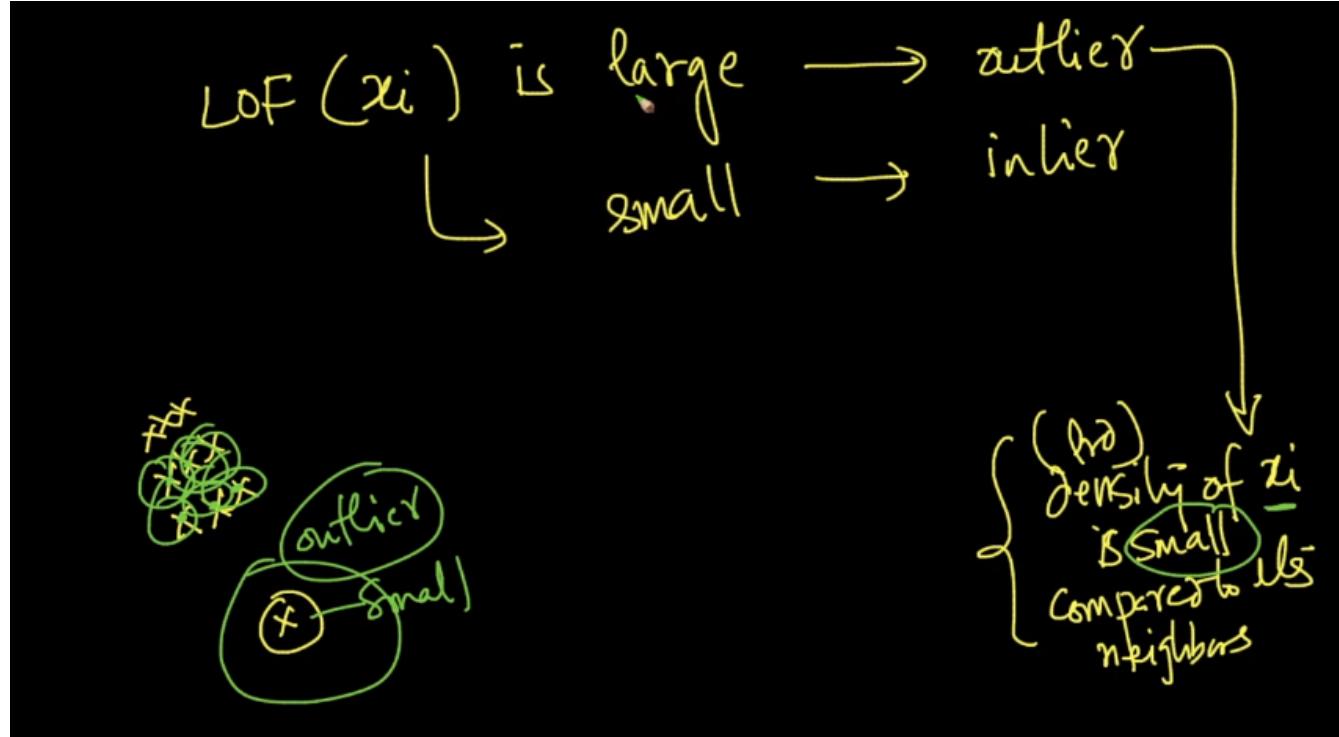
Local outlier factor ( $x_i$ )

$$= \frac{\sum_{j \in N(x_i)} \text{ld}(x_j)}{|N(x_i)|}$$

Large  $\Rightarrow$   $\text{ld}(x_i)$  is small

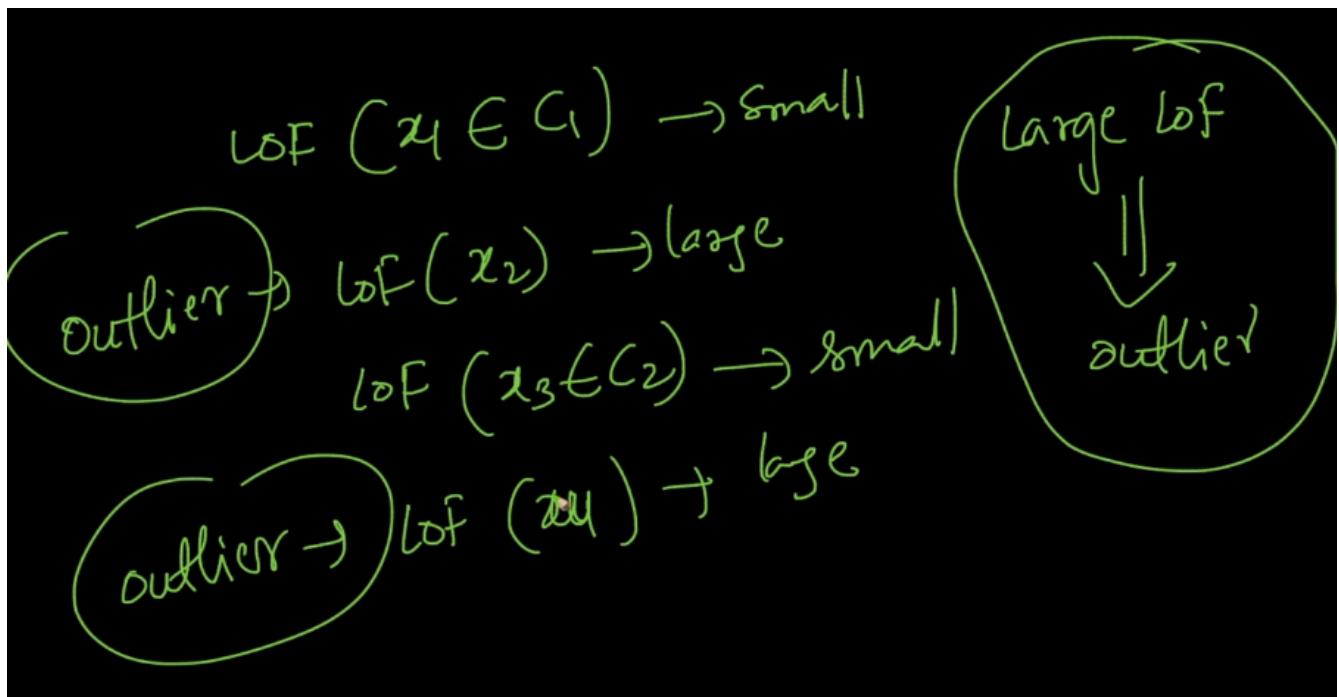
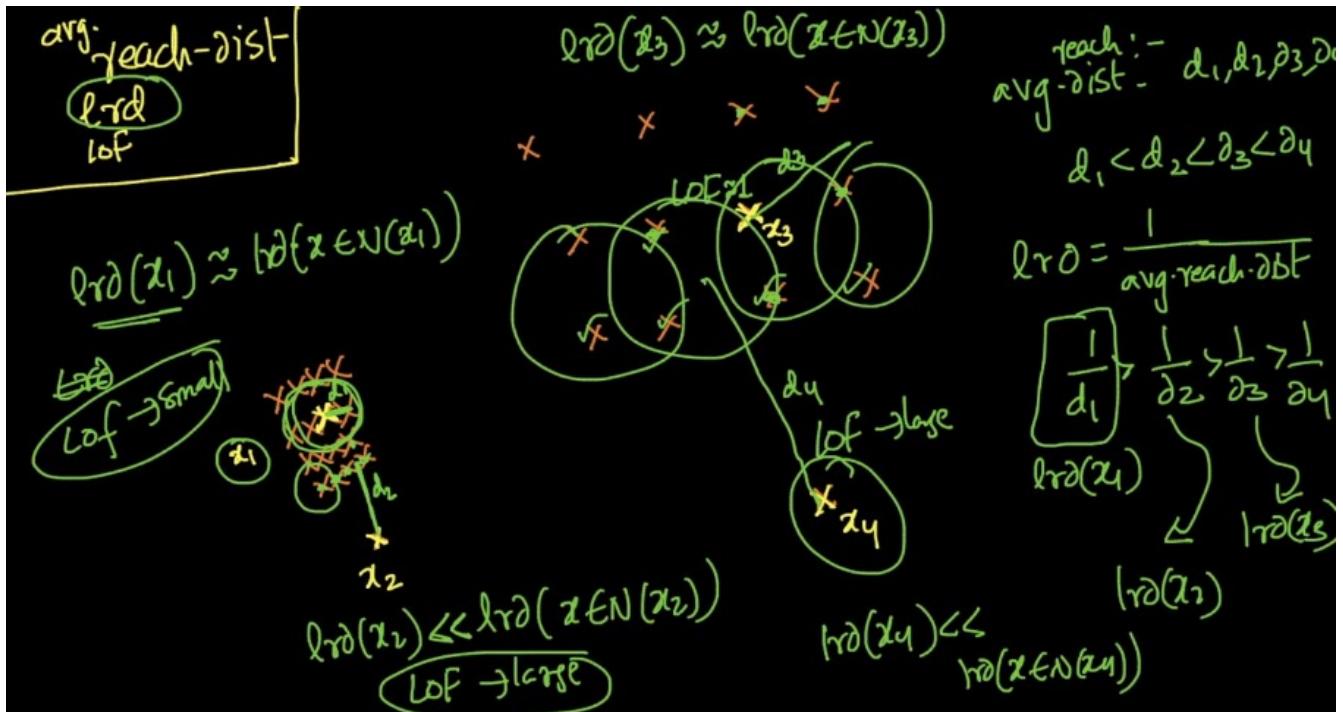
$\text{LOF}(x_i)$  is large  $\Rightarrow$   $x_i$  is an outlier

(avg.  $\text{ld}$  of pts in the neighborhood of  $x_i$ )



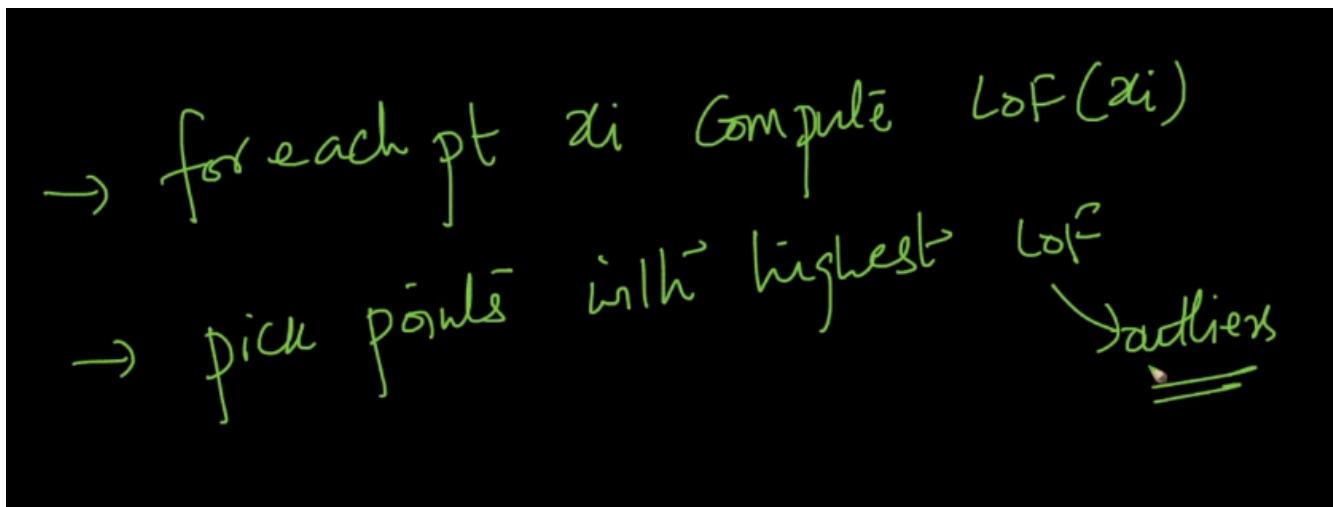
Example:

Explanation of local outlier factor with example.

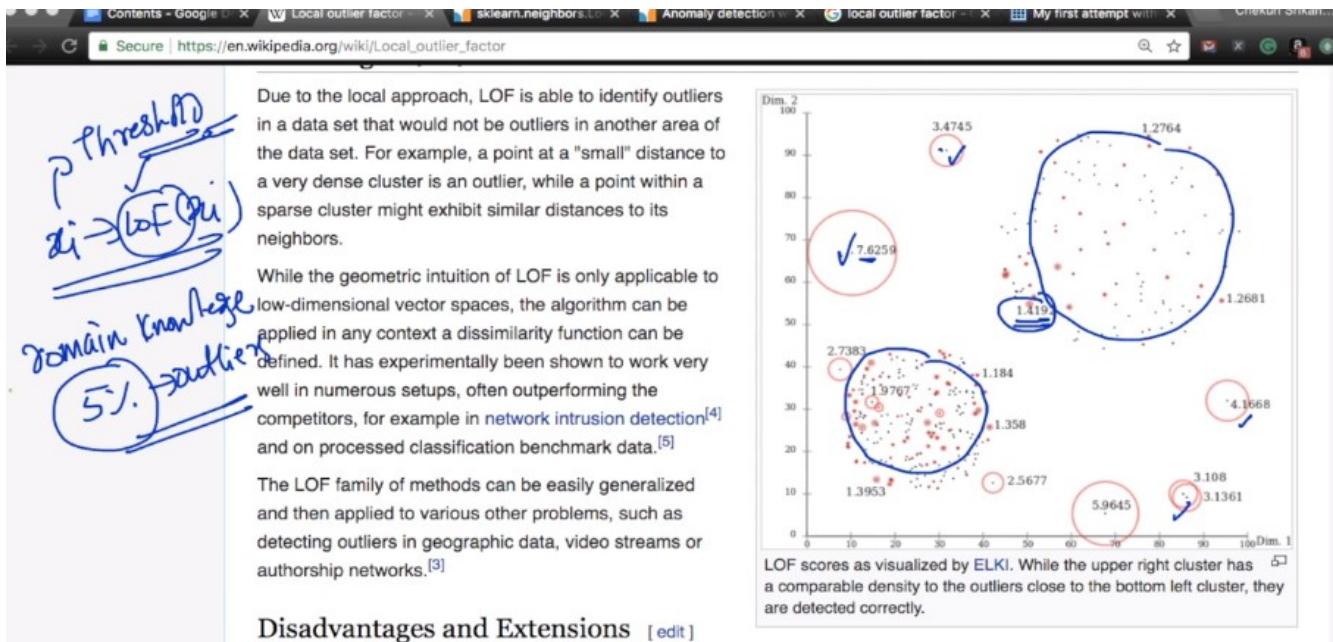


The larger the LOF is the outlier.

The threshold of LOF values depends on the application.

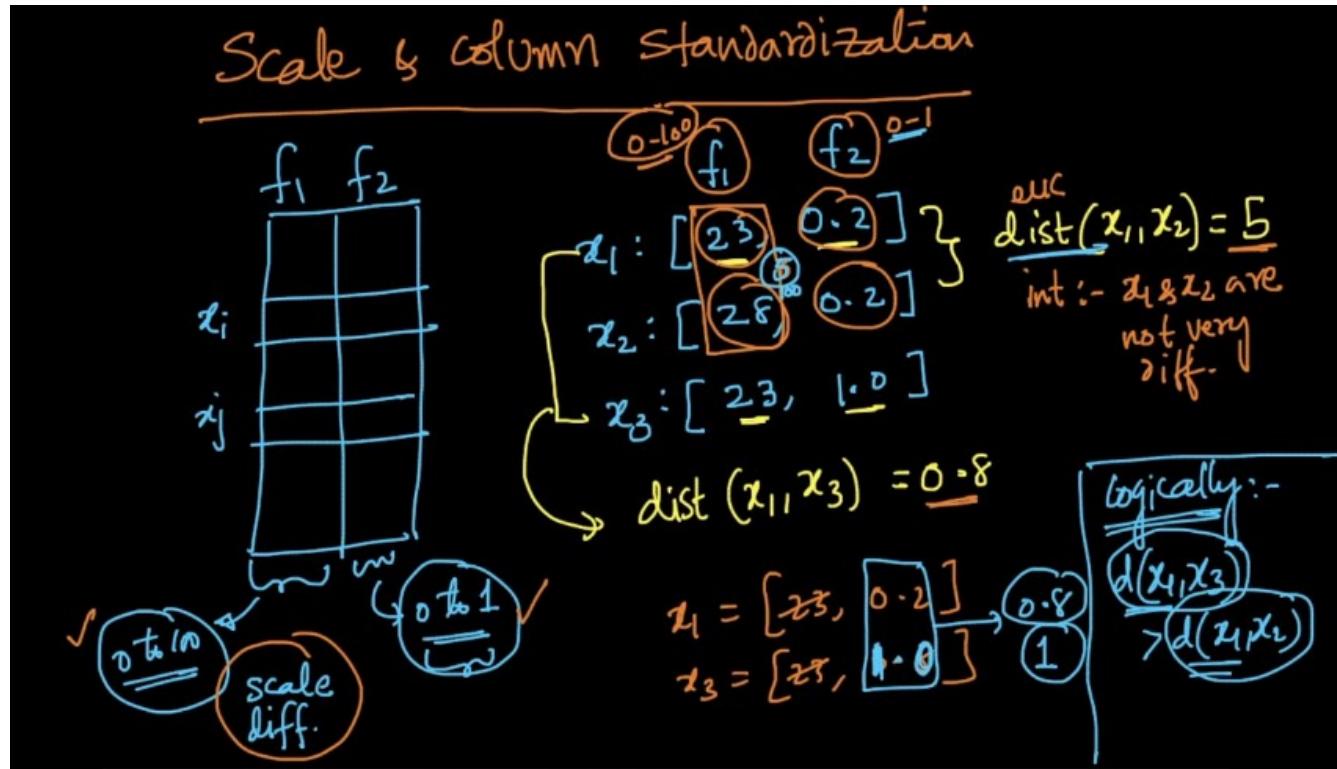


From the domain knowledge if we know the percentage of the points that will be outliers, then we can pick right LOF distance value and remove the points that are greater than the distance chosen.

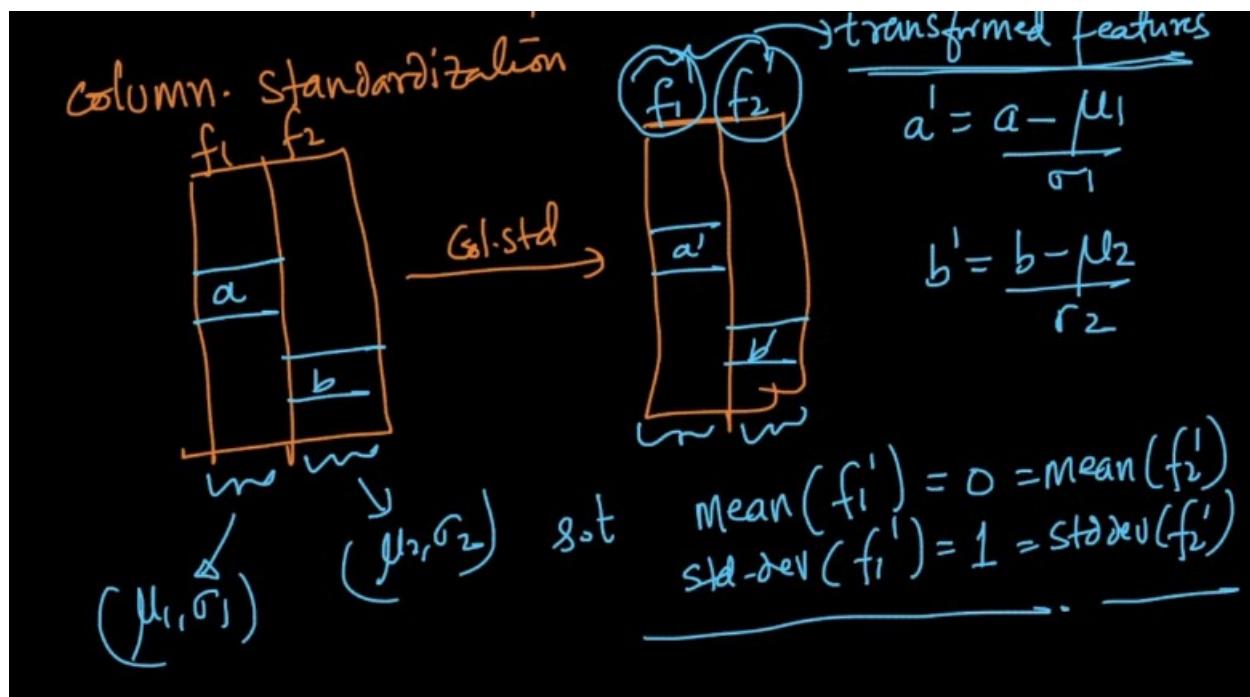


Impact of scale and column standardization:

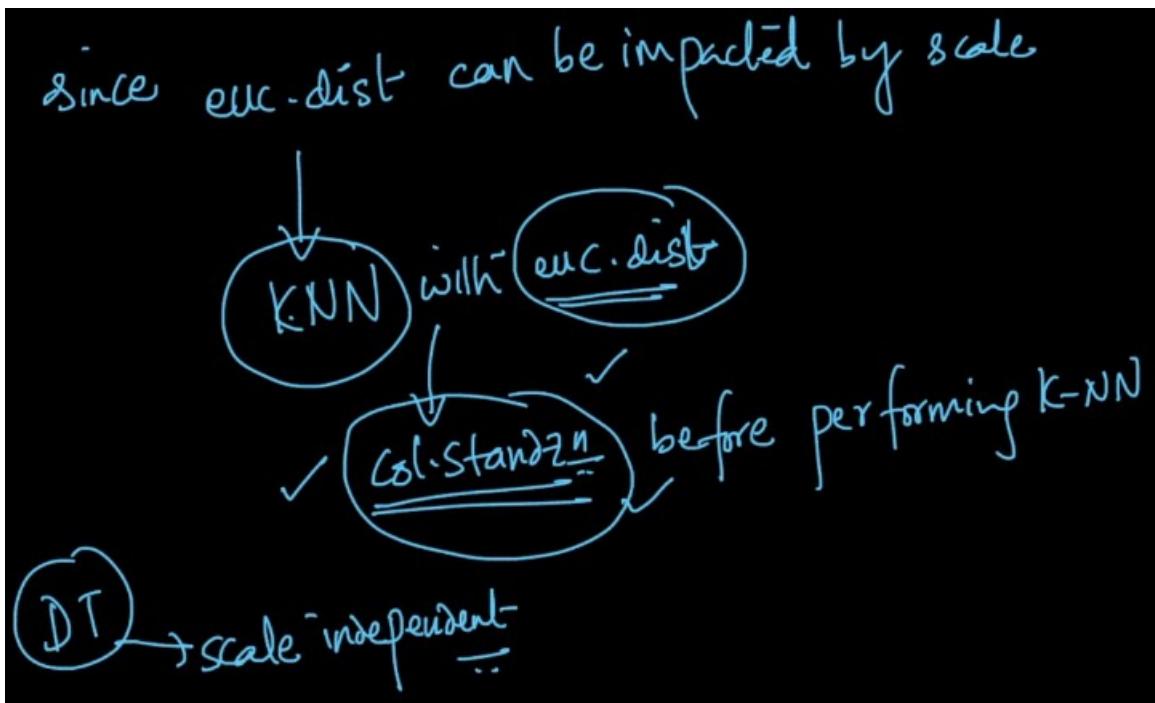
The values of the attributes will be different in scale compared one another. The distance between the points will be different when the scale of the attributes are different.



Column standardization solves the scale difference of the attributes.

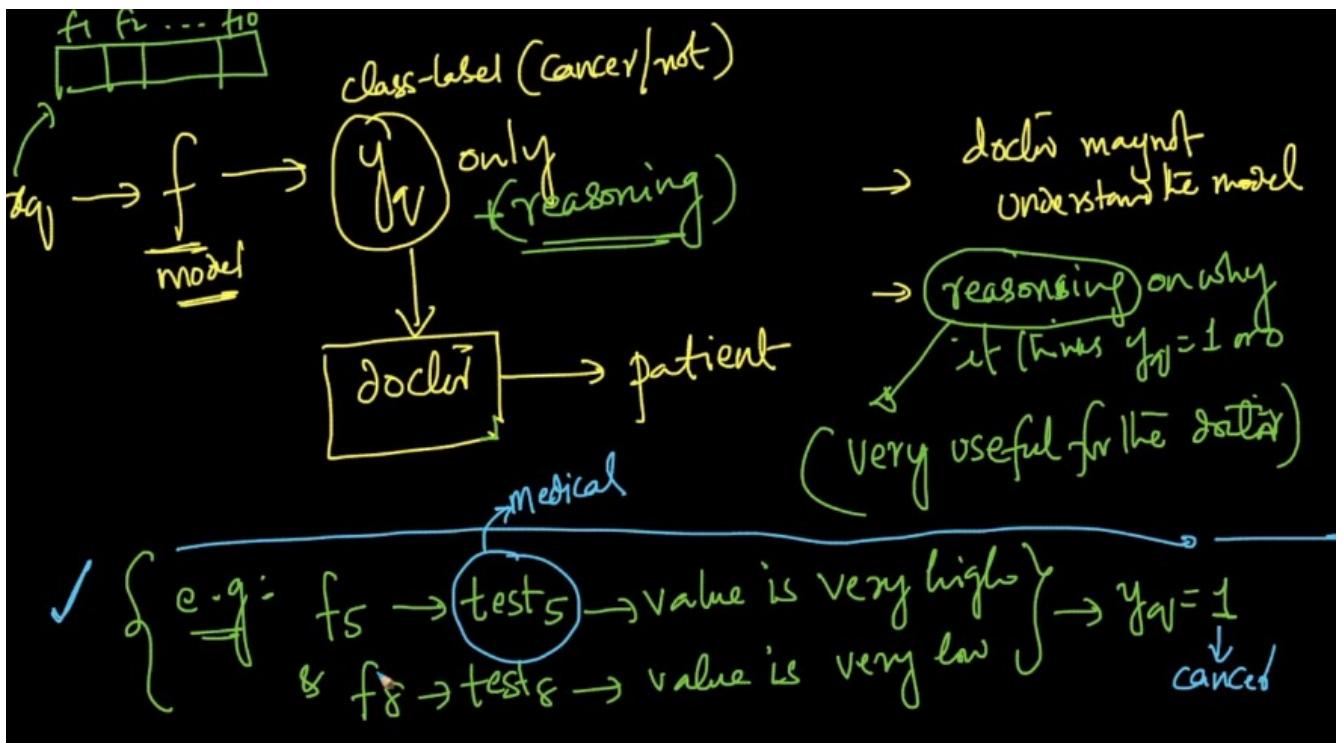


The distance measures will be impacted by scale. So, we need column standardization.



Interpretability:

If we provide reasoning of the prediction of the model, we can better understand the results.



Models that cannot give the reasoning is called black – box model.

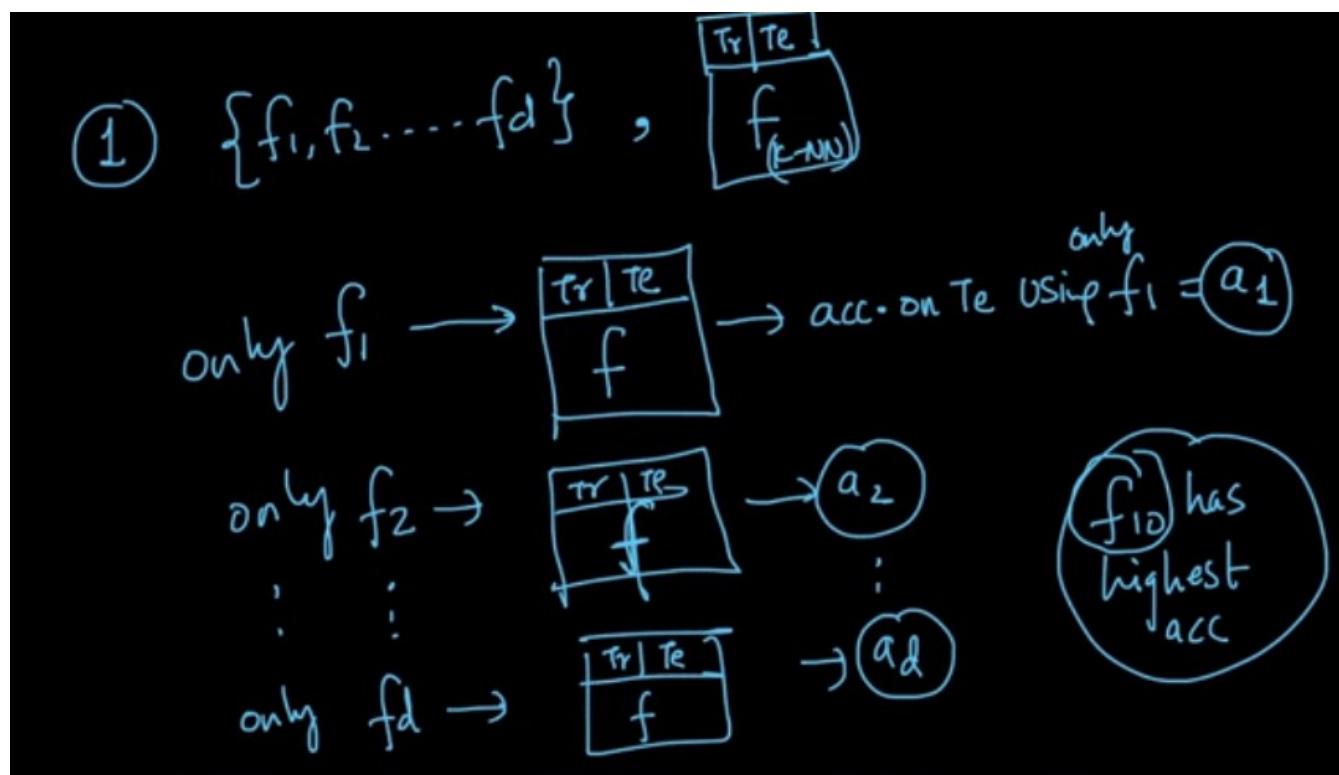
The interpretability plays important role in reasoning the results of the model.  
KNN is interpretable when d is small and k is small.

Feature importance and forward feature importance:

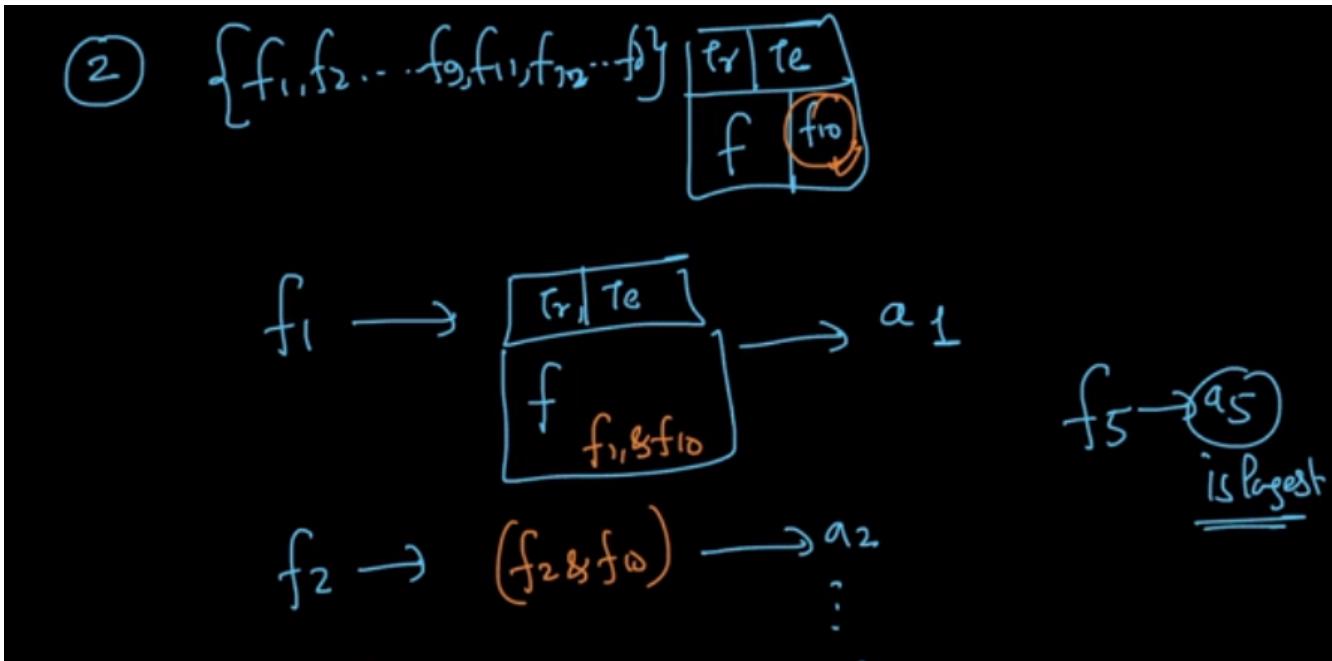
Forward feature selection:

Picking the top features in the data set is called forward feature selection.

For every feature we will make train, test and develop model to get the accuracy scores and the model with highest accuracy is chosen as the most favoring feature towards the accuracy.



The chosen features are again compared with other features and accuracy is calculated. This is called forward feature selection.



At every stage we will choose the feature that add most value to the model.

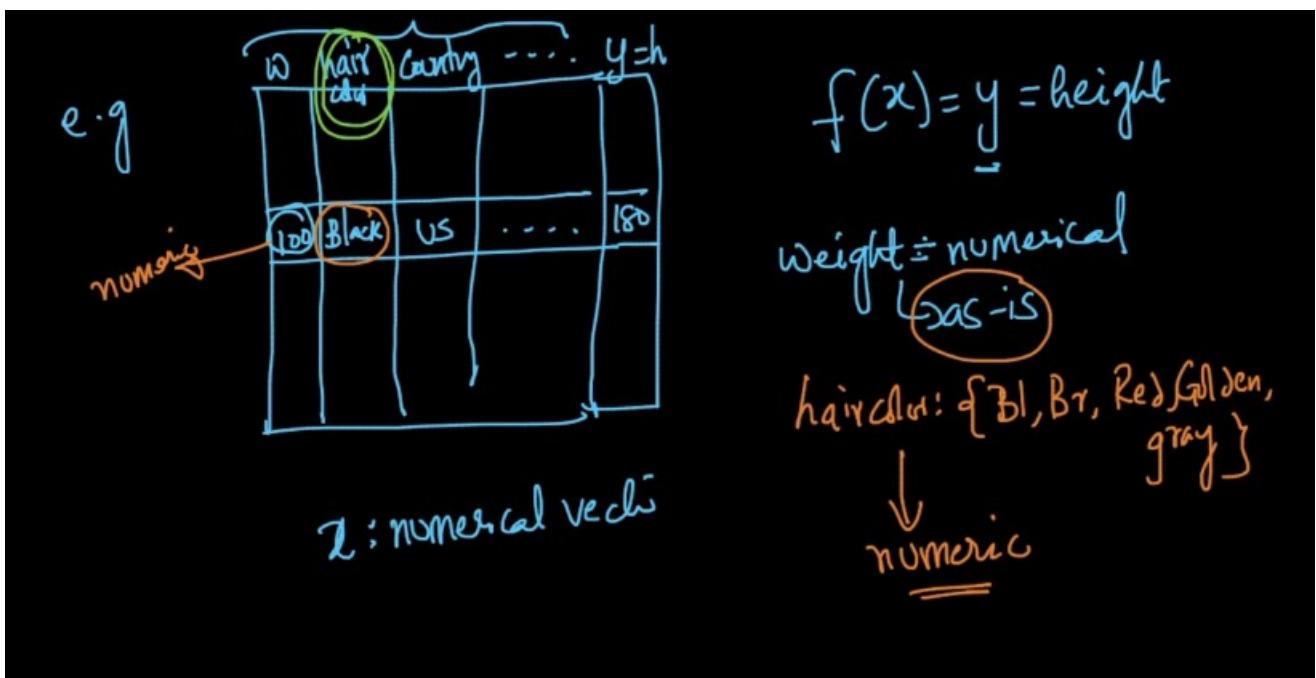
Backward feature selection:

In each iteration we will remove the feature which results in lowest drop in accuracy.

Handling categorical and numerical features:

Converting categorical features to numerical – Converting the categories to numerical will give an artificial constraint that changes the meaning of the categories.

To overcome this situation we will use one – hot encoding.



One-hot encoding makes the binary vector of each category of the size number of categories in the feature.

✓ hair color = {Bl, Br, Red, G, Gray}

① Give a number  
↳ numbers have an inherent order  
Red > Bl X

② One-hot encoding: hair color

Bl	Br	R	G	Gr

binary vector  
 $x_i: \text{Br} \rightarrow [0|1|0|0|0]$   
 $x_j: G \rightarrow [0|0|0|1|0]$

The problem with one-hot encoding, it increases the dimensions of the data.

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

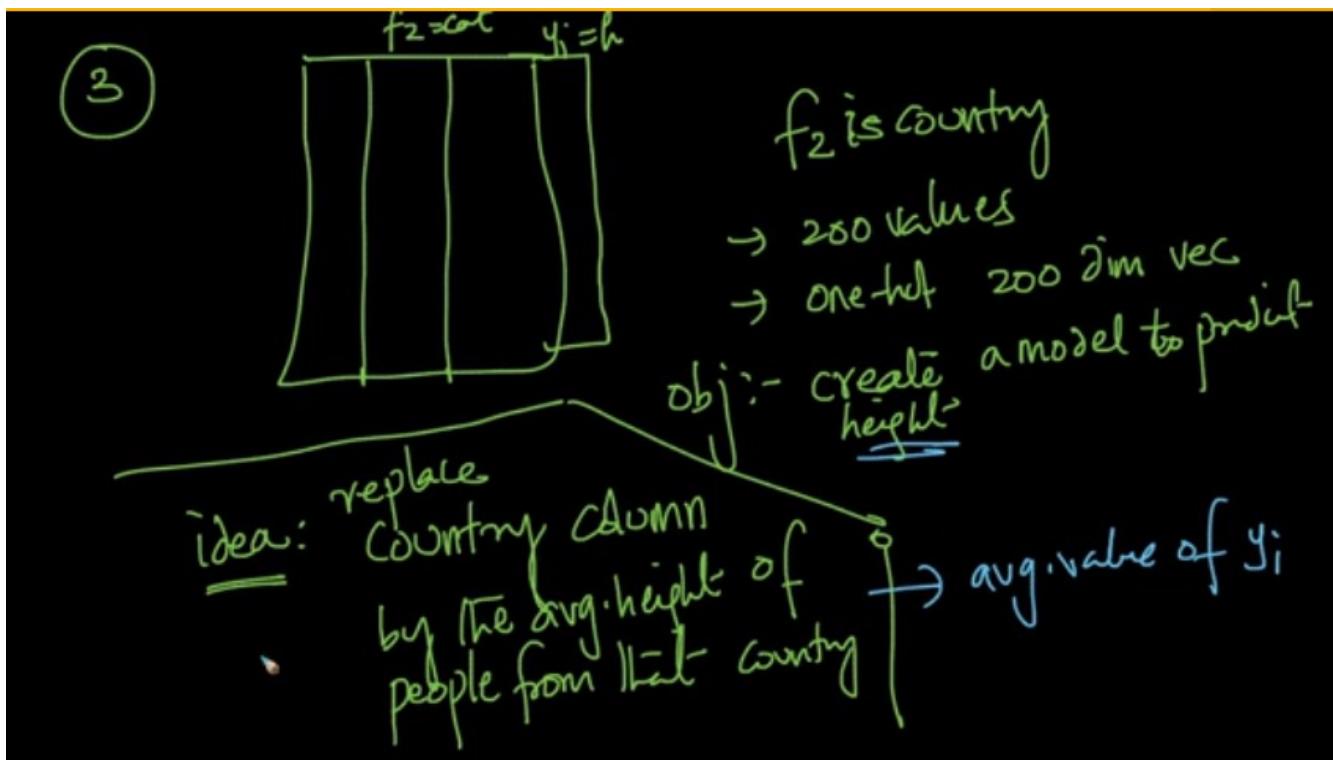
↳ counting  $\approx 200 \rightarrow$

$x_i:$   $c_1 c_2 \dots c_{200}$  sparse

Bow  
one-hot encoding

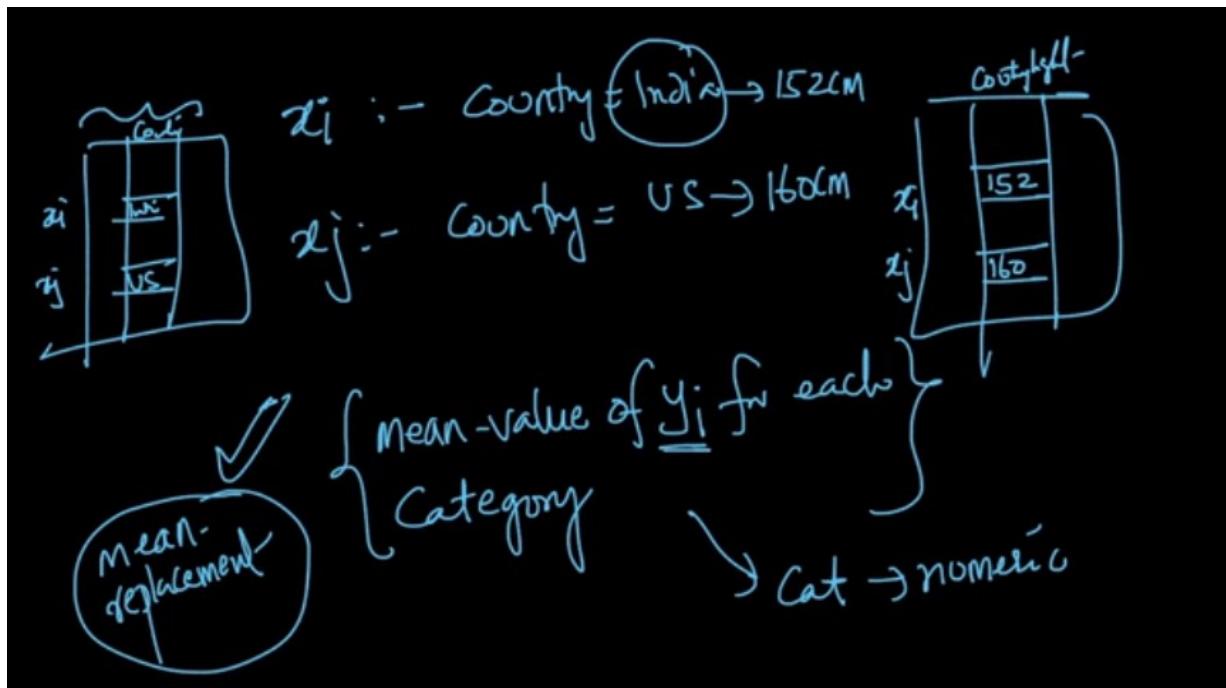
If the number of distinct values for a categorical feature is large then  
One-hot encoding can create sparse binary vectors

Other types of encoding:



Example:

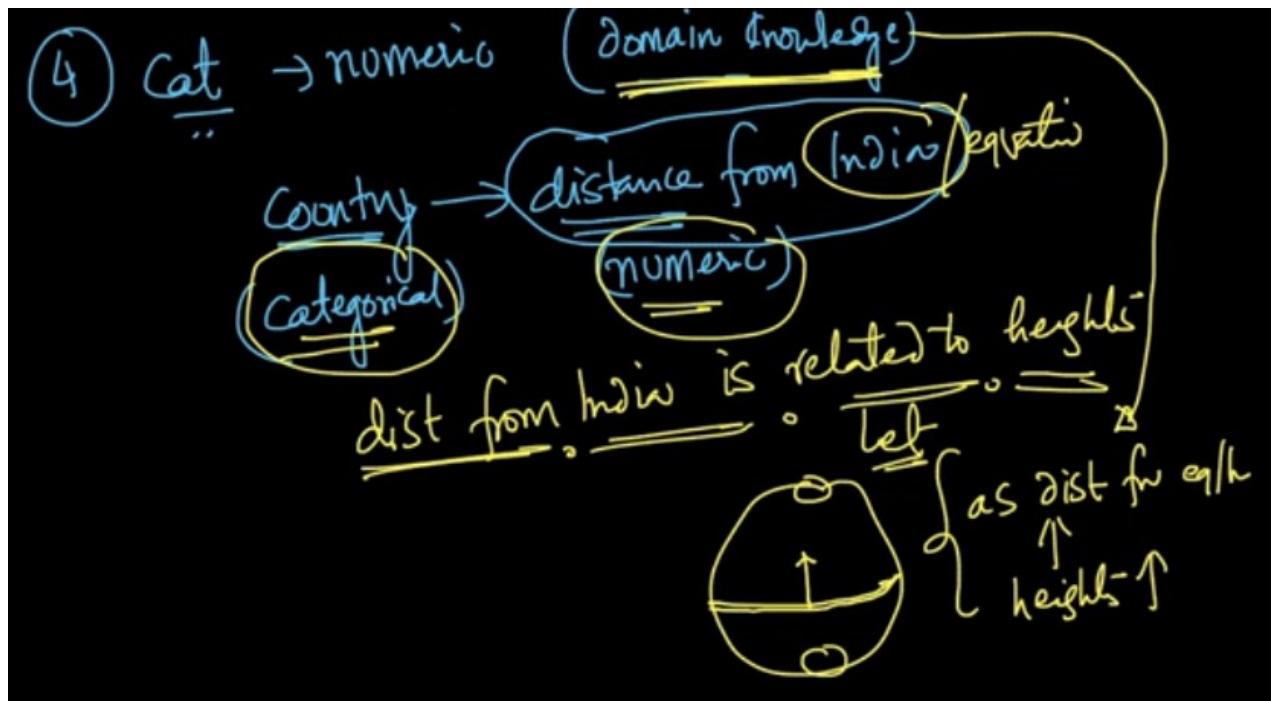
We can replace each category with the mean value of the predicting variable. This is called mean replacement of the categorical data.



This strategy is the problem specific approach.

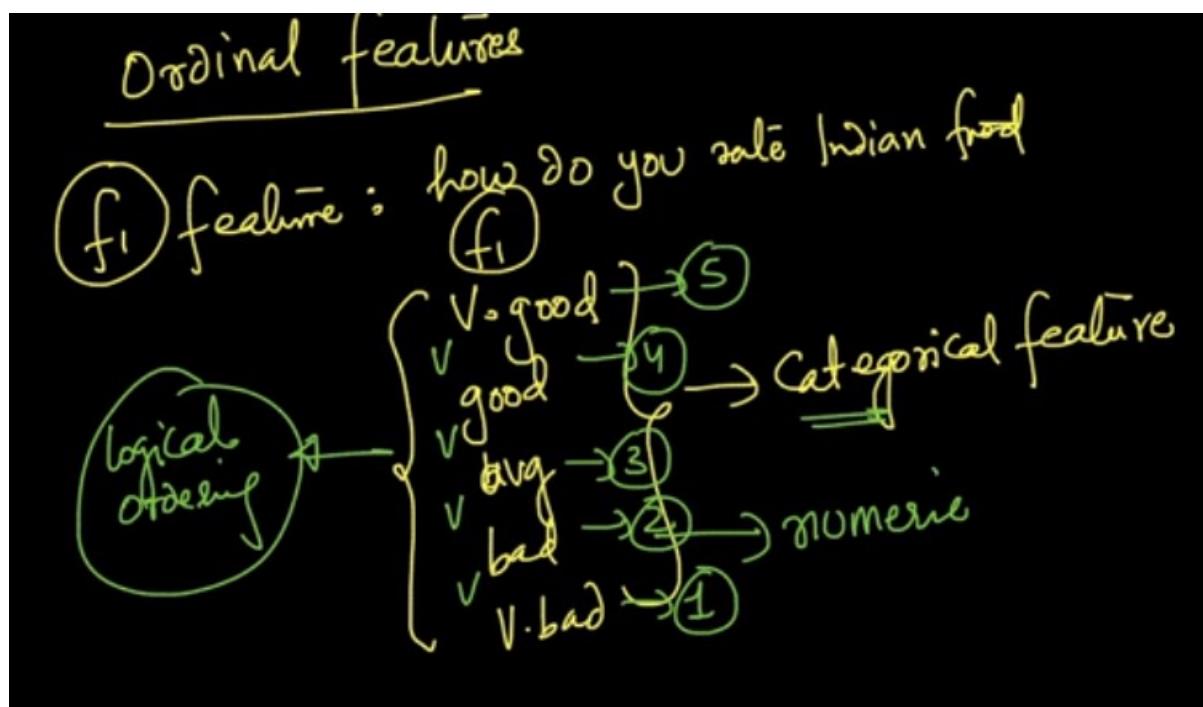
Another approach:

From domain knowledge we know that heights are related from one specific country.



Ordinal features:

Features that can be ordered in one specific order. If there is an order in the categories of the feature then we can convert them into numeric features by assigning the numbers in the same order.



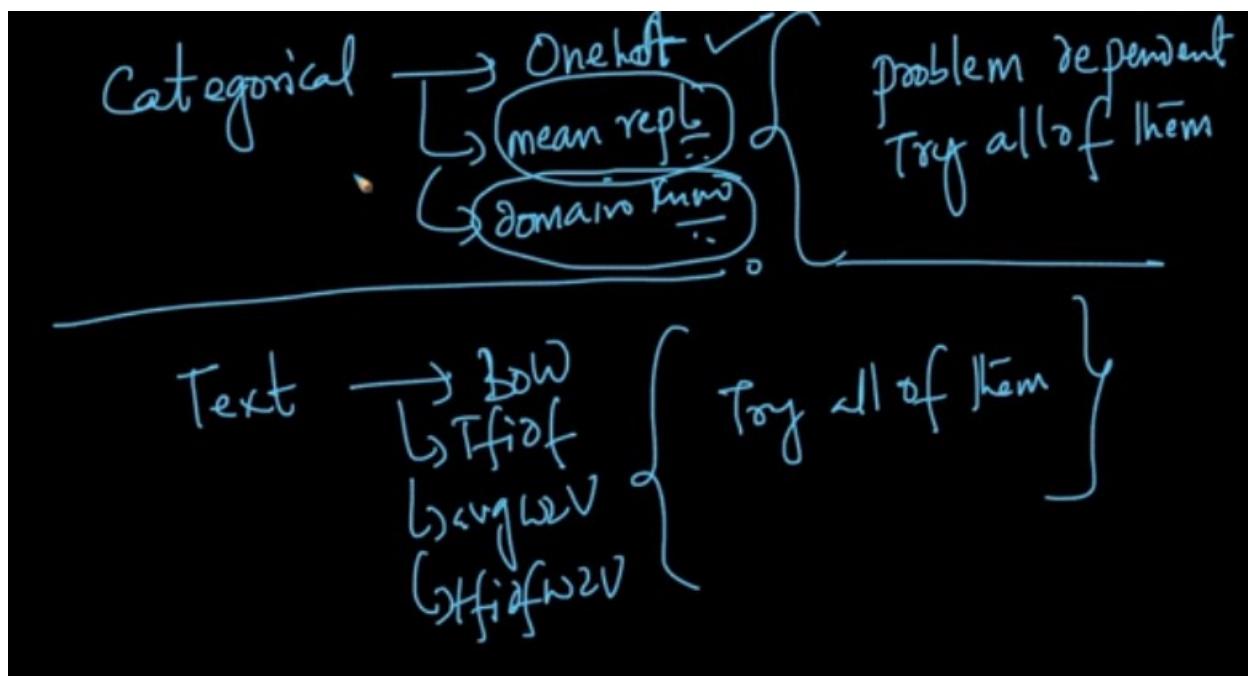
Ways of numeric encoding:

The order of giving numeric values is very problem specific.

$f_1$	$f_1'$	$f_1''$
V-good	5	6
good	4	5
avg	3	3
bad	2	2
V.bad	1	1

ordinal cat. features  
↓  
numeric  
(?)

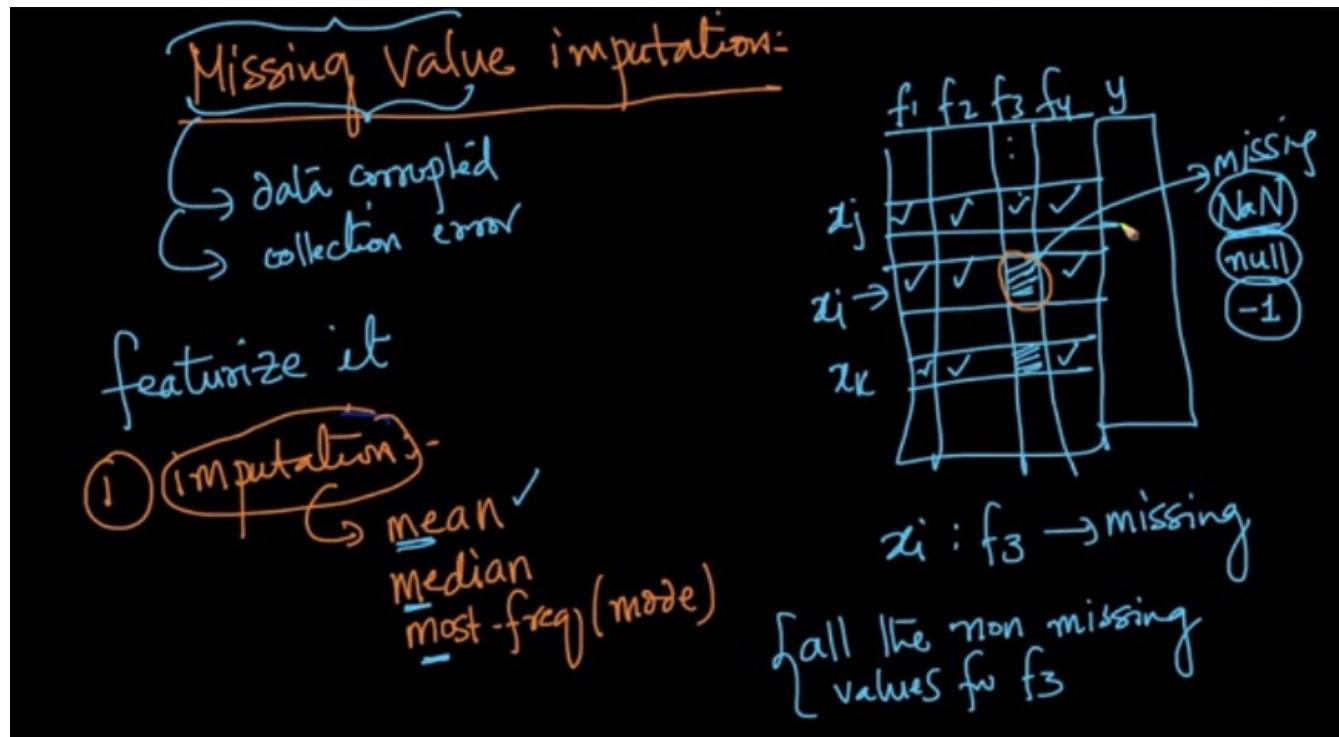
It is very problem specific to choose which will be the best for encoding the data.



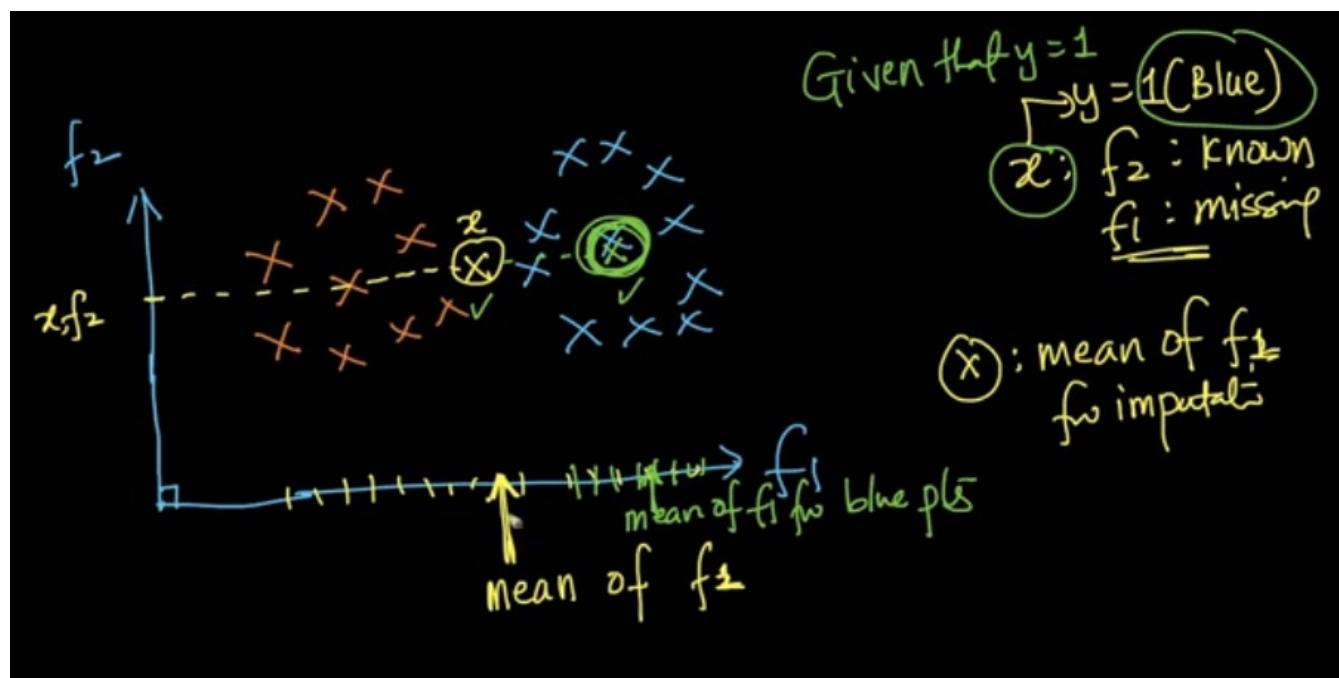
Handling missing values by imputation:

The missing values can be replaced by

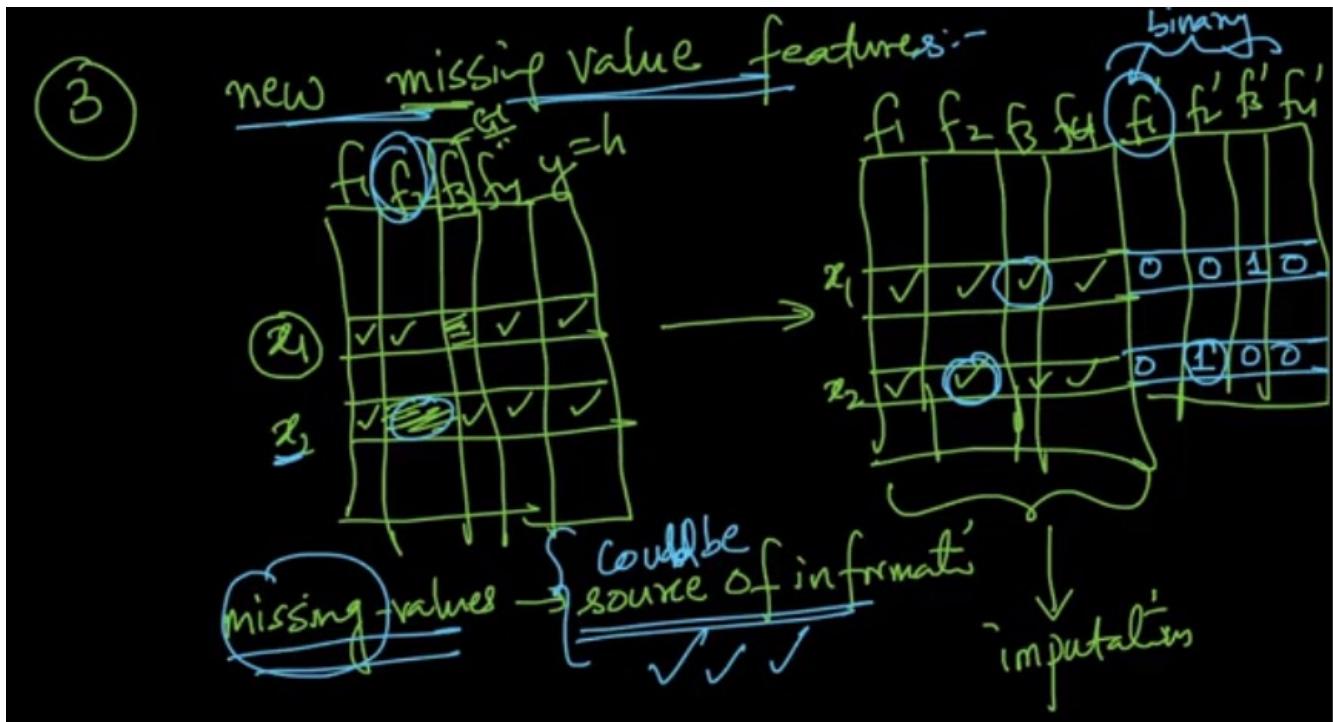
1. Mean
2. Median or
3. Mode.



Edge case of applying mean value.

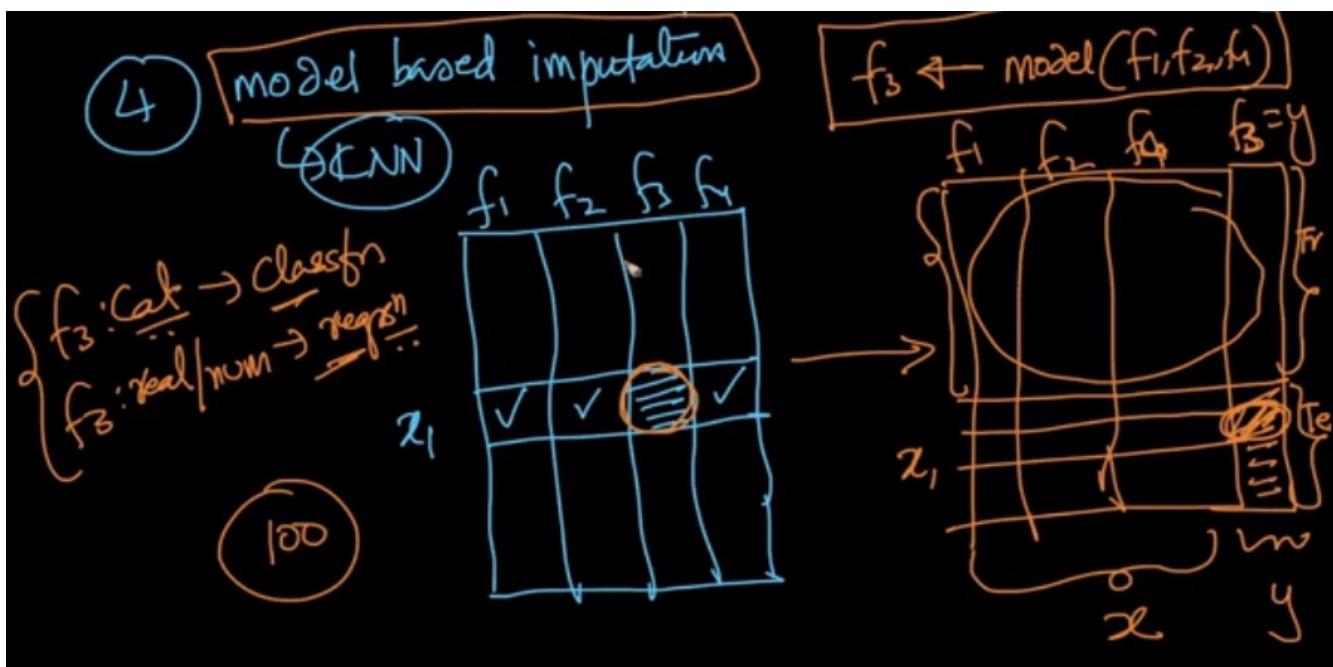


One more strategy of imputation is make a new features marking the features that are missing.



Model based imputation:

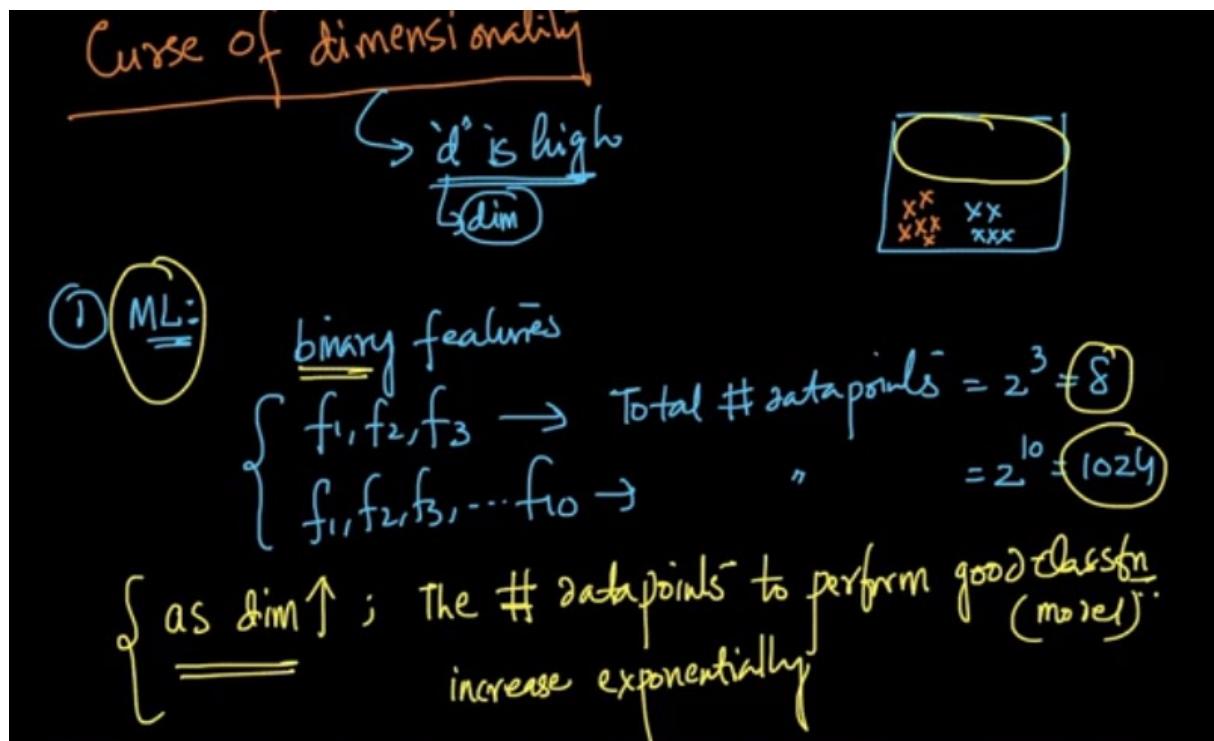
In this strategy we will make all the missing values as the test data and non – missing as training data and train KNN and predict the missing values in the test data.



KNN is most often used for model based imputation.

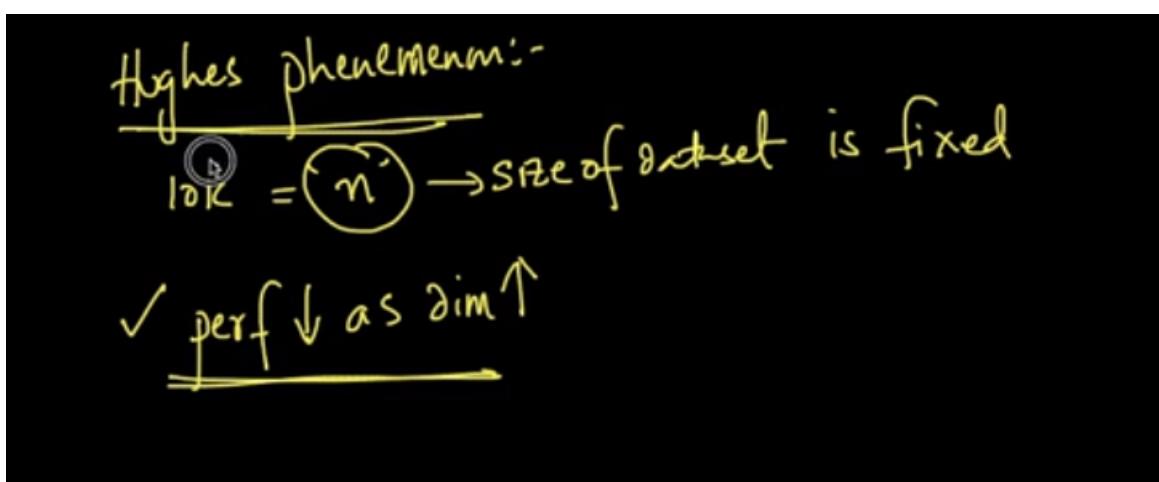
Curse of dimensionality:

As dimensionality increases then the number of data points to perform good classification increase exponentially.



Hughes phenomenon:

If the size of the data set is fixed, the performance of the model decreases as the dimensions increase.



Effect on distance functions if dimensions increase:

(2) Distance functions (euclidean dist)

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

1D-world  $\rightarrow$  n-random pts

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

(1)  
(2)  
3D  $\rightarrow \left\{ \frac{\text{distmax}(x_i) - \text{distmin}(x_i)}{\text{distmin}(x_i)} > 0 \right\}$  when  $d=1, 2 \text{ or } 3$  (heart icon)

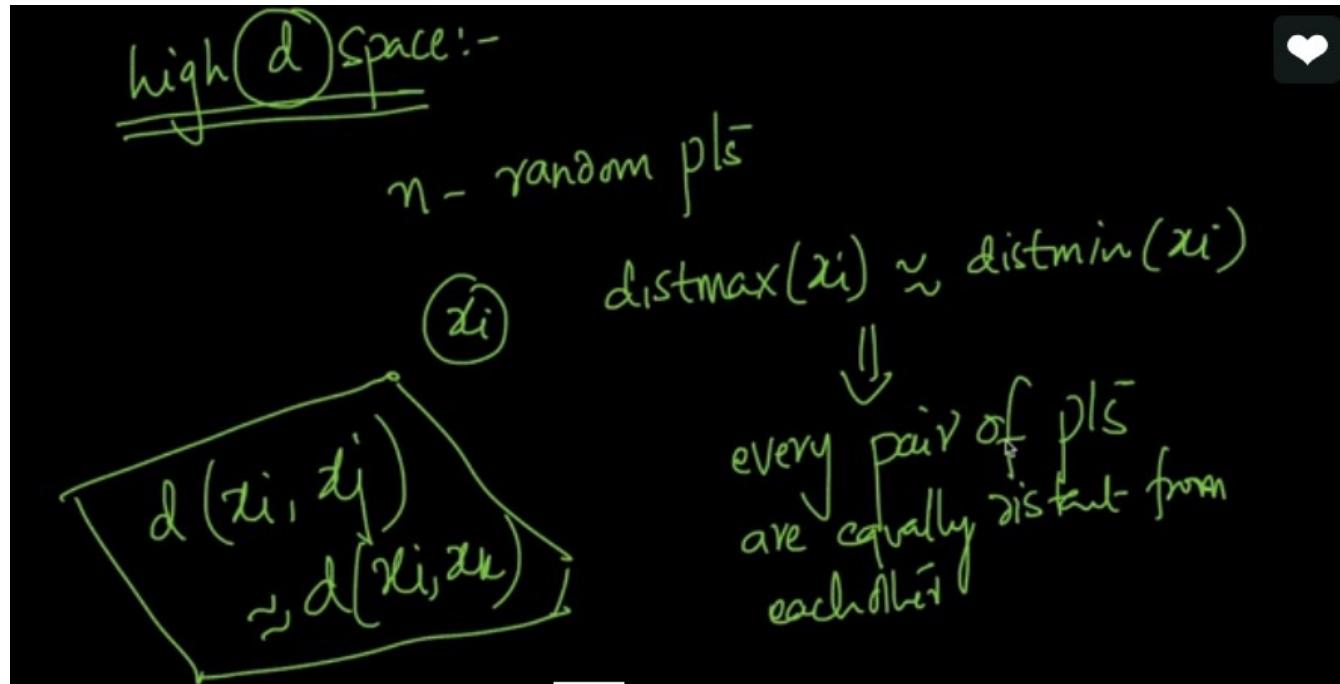
as  $\dim \uparrow$

$$\left[ \lim_{d \rightarrow \infty} \left( \frac{\text{distmax}(x_i) - \text{distmin}(x_i)}{\text{distmin}(x_i)} \right) \rightarrow 0 \right]$$

$\Downarrow$

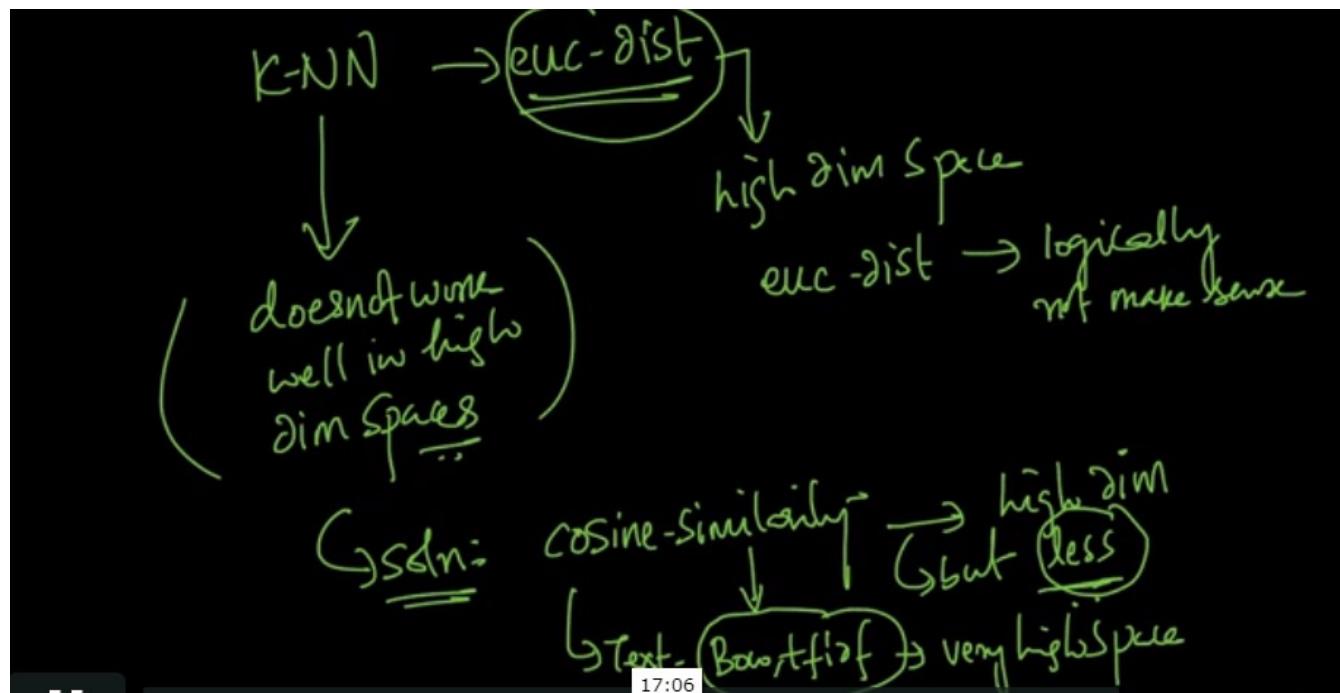
( $\text{distmax}(x_i) \approx \text{distmin}(x_i)$ )

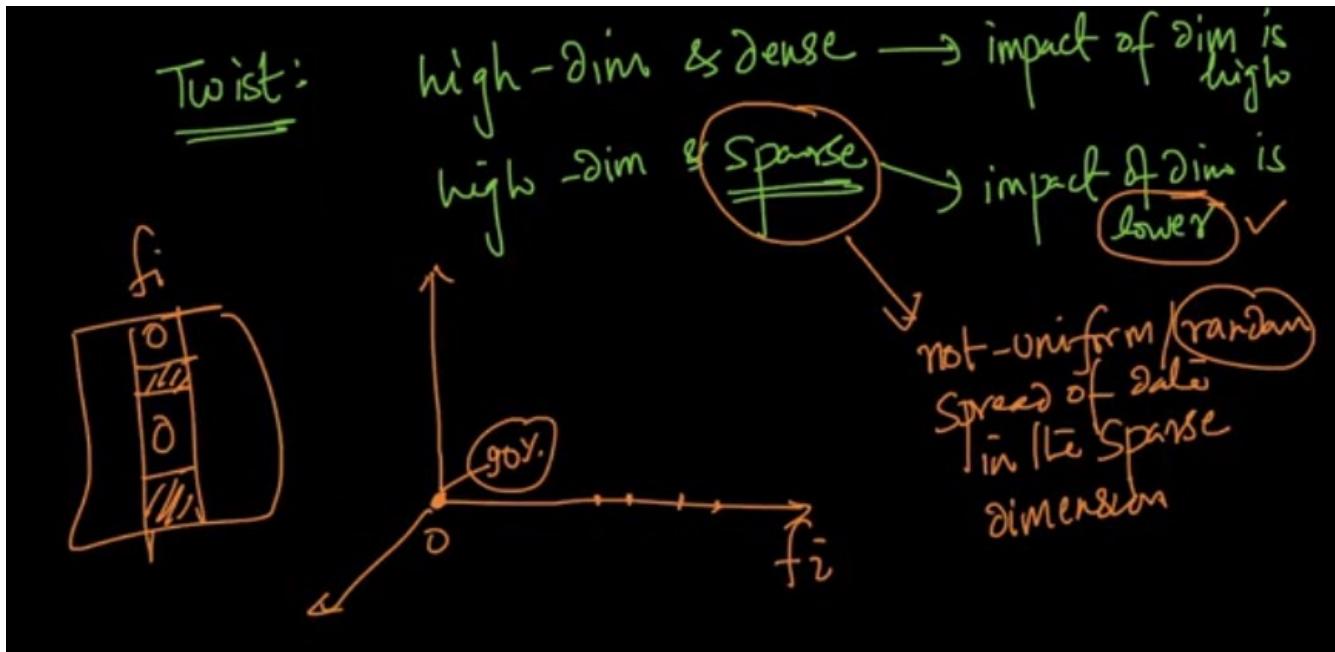
In a high dimensional space, if we take 'n' random points then its distance max is equal to distance min from a point. Every pair of points are equally distant from each other.



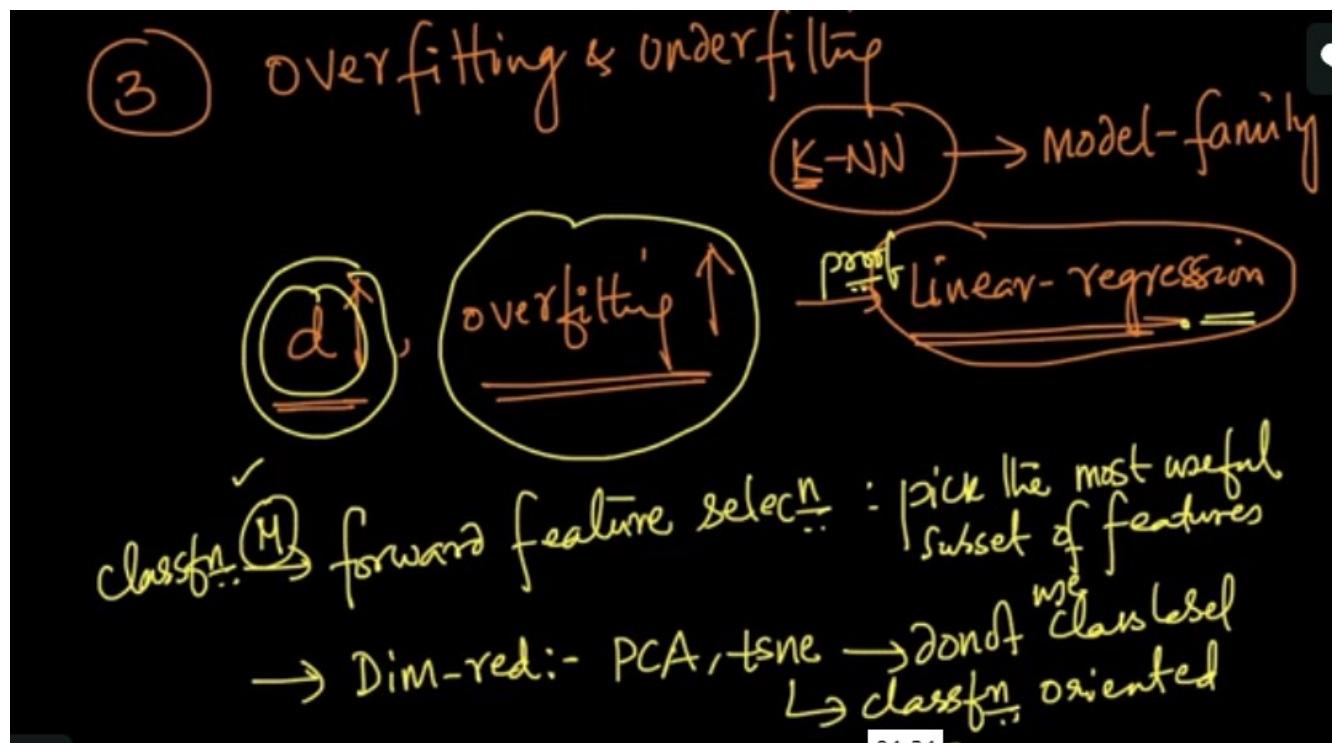
In high dimensional space euclidean distance does not logically make sense.

**KNN does not work very well in high dimensional data.**



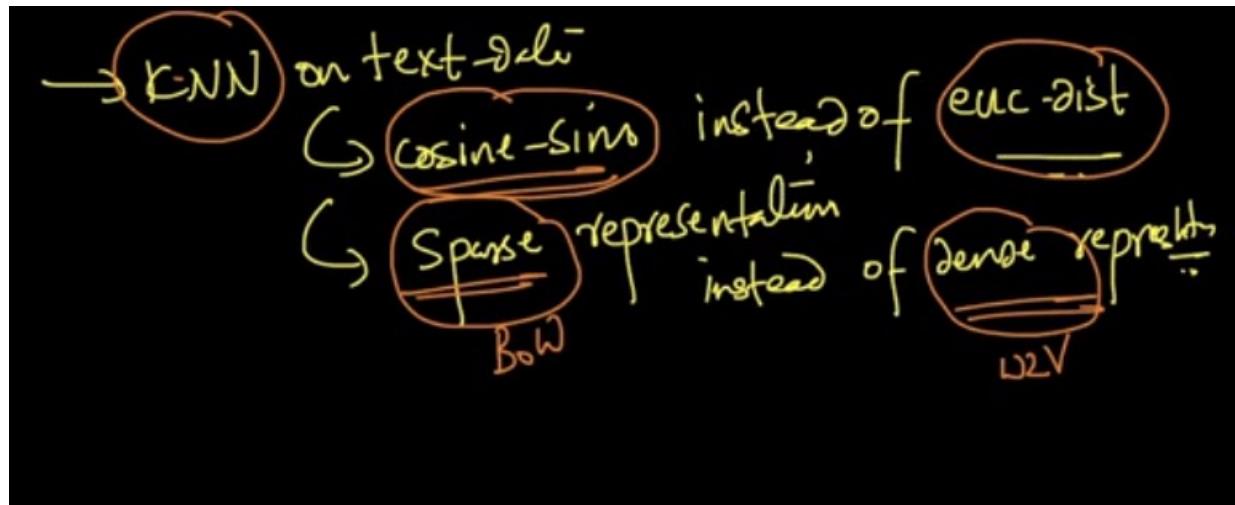


Overfitting and Underfitting:



If there are large dimensions we can use forward feature selection for classification.

Ways implementing KNN when dimensions are more.



Bias – Variance trade off:

Equation of generalization error:

The diagram shows handwritten notes on the Bias-Variance Trade Off and the equation of generalization error:

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

→ Underfitting & overfitting → {Math-basis ...}

$E_{\text{generalization}} = \underbrace{\text{Bias}^2}_{\substack{\text{(for regression)} \\ \text{proof}}} + \underbrace{\text{Var}}_{\substack{\text{(future unseen)} \\ \text{of } M}} + \underbrace{\text{irreducible error}}_{\substack{\text{for regression (proof)}}}$

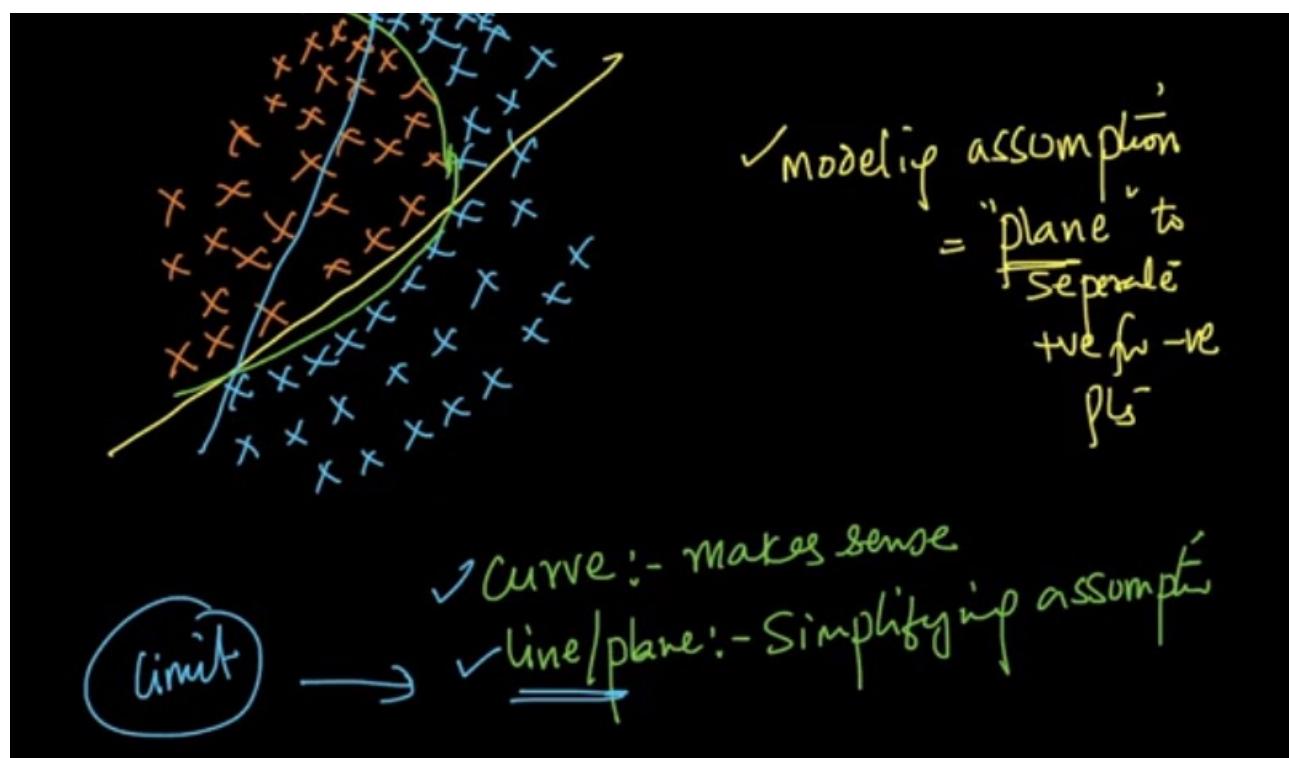
Irreducible error: This is the error that cannot be reduced further for a given model.

Bias error: This is due to simplified assumptions.

Variance:

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

high bias  $\Rightarrow$  underfitting

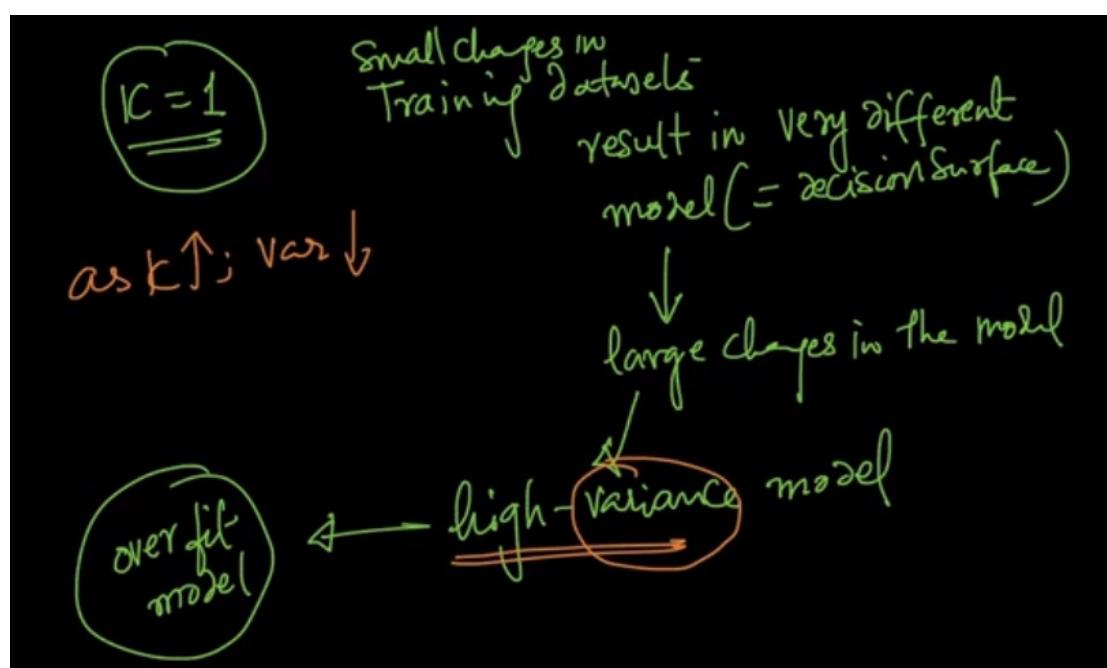
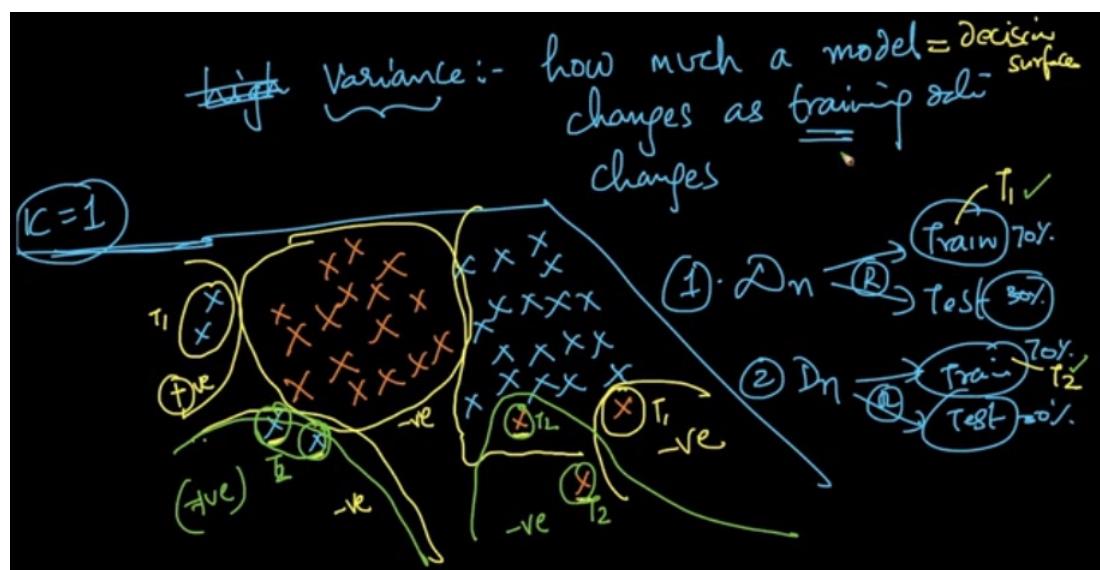


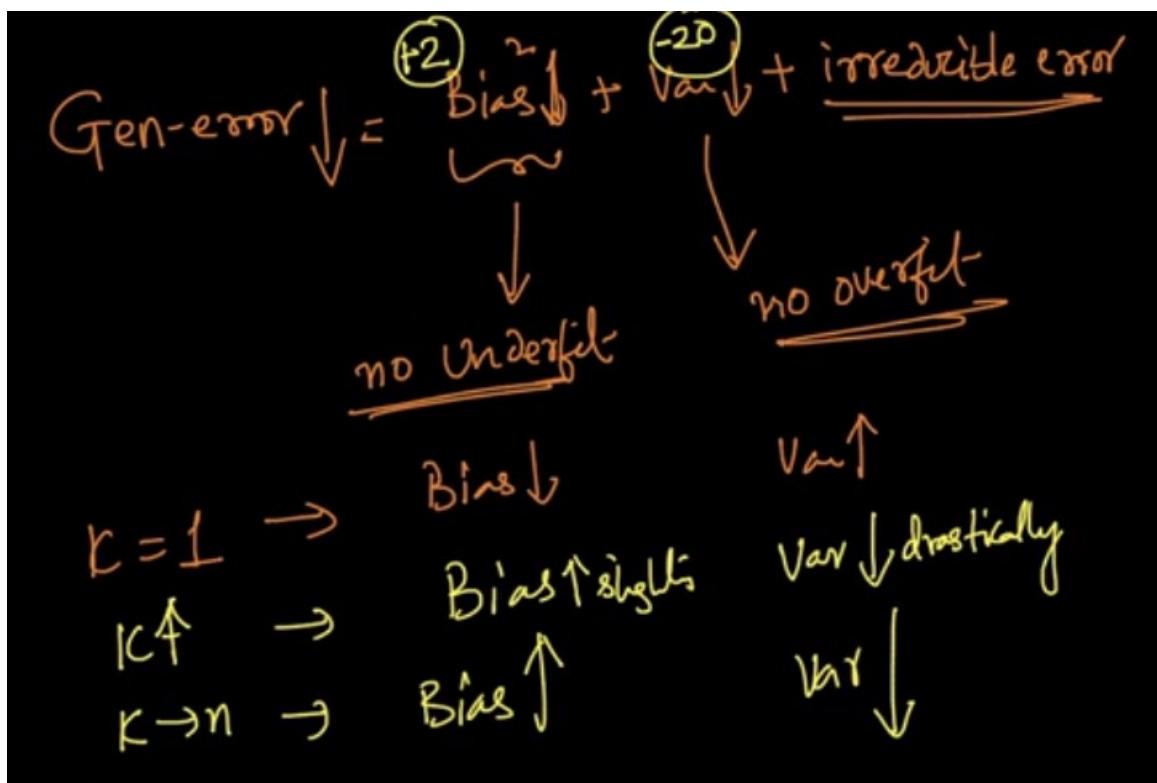
In KNN, when  $k = n$  the dominant class will be the class label of the query point.

This means we are underfitting the data-set.

① ✓ high bias  $\Rightarrow$  Underfitting  
↳ Simplifying assumptions

Variance: how much a model changes/decision surface as training data changes.





Balancing bias and variance will give the good model as output.

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

Let  $\begin{cases} k=1; \\ k=5; \\ k=n; \end{cases}$

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

Best and worst cases for an algorithm:

## Best & Worst Cases:

(K-NN)

- ①  $\text{dim}$  is small ( $< 10$ ) [not large]  $\rightarrow$  K-NN
- ②  $d \uparrow$  is
  - curse of  $\text{dim} \rightarrow \text{euc dist}$
  - $\hookrightarrow$  interpretability  $\downarrow$
  - $\hookrightarrow$  runtime complexity of kd-tree/LSH  $\uparrow$