

$M_1$ : 95% of errors are below 0.1  
 $M_2$ : 80% of errors are below 0.1

$\therefore M_1$  is better than  $M_2$

## Classification Algorithm in Various Situation

### 31.1 Introduction

All the algorithms are used and applicable in various real world cases.

Lets say, I have a dataset with very few positive points and lot of negative points

$D = \frac{+ve}{\uparrow}, \frac{-ve \text{ points}}{\downarrow}$

very few      lot of them

} Happens a lot in Medical Applications

- In 10k patients, 1k has cancer and rest 9k dont have cancer
- 1M people visit website but only 1000 people purchasing the product.

31.2

## Imbalanced Vs Balanced Dataset

2 class classification :

$$D_n \rightarrow n_1 \text{ +ve points}$$

$$\rightarrow n_2 \text{ -ve points}$$

$$n_1 + n_2 = n$$

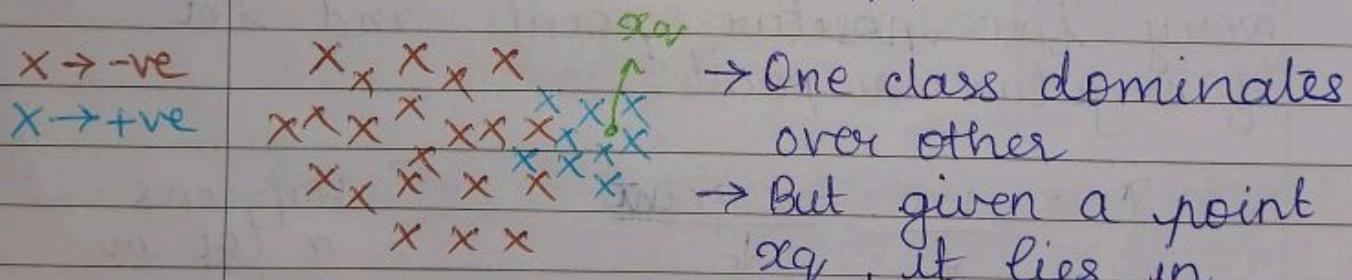
If  $n_1 \approx n_2 \rightarrow$  Balanced dataset

If  $n_1 \ll n_2 \text{ OR } n_2 \ll n_1 \rightarrow$  Imbalanced dataset

Moderately Imbalanced / Severely Imbalanced :

$$n_1 = 50; n_2 = 950.$$

So, dataset would appear like this



→ One class dominates over other

→ But given a point  $xq$ , it lies in

+ve cluster, hence we would classify it as positive point

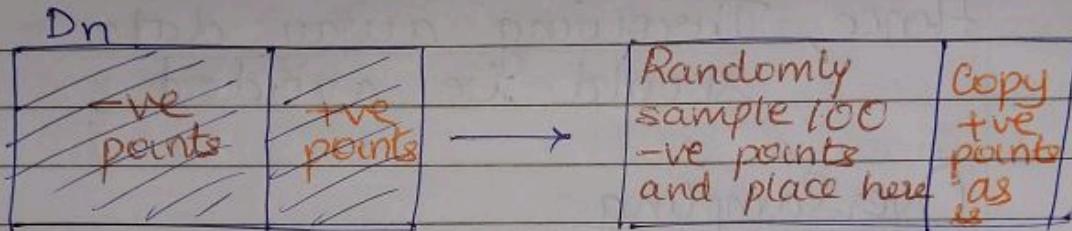
→ By seeing the number of  $n_2$ , we may wonder that any query point may have high chance to be -ve. But that's not always the case

How to work around Imbalanced Dataset problem?



### Solution 1: Undersampling

$$\begin{array}{l} n_1 = 100 \text{ (+ve)} \\ n_2 = 900 \text{ (-ve)} \end{array} \quad > D_n$$



Undersampling states that given  $D_n$ , use a sampling trick to create a new dataset called  $D'_n$ .

Now  $D'_n$  has 100 -ve and 100 +ve points which makes  $D'_n$  Balanced dataset.

Hence in  $D'_n$ ; we keep Minority points as is and sample Majority points.

### Problem with Undersampling

Originally in  $D_n$ , we have 1000 points

900 -ve      100 +ve

In  $D'_n$ , we have only 200 points

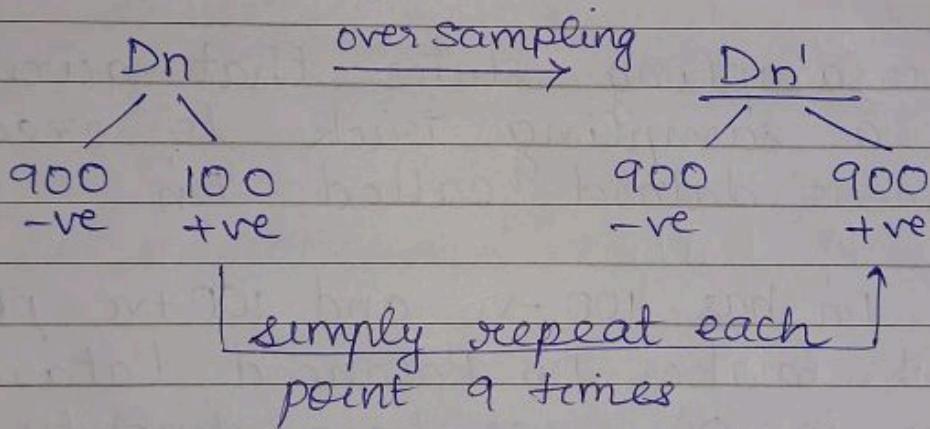
100 -ve      100 -ve

So, we are throwing 800 points from  $D_n$ .

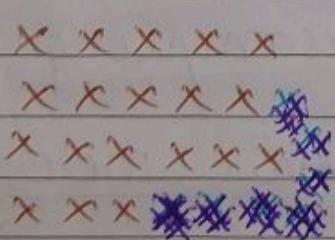
Since  $|D_n'| \ll |D_n|$ , the Model may not be accurate since it doesn't work on all the dataset points.

Hence, Throwing away data should be avoided.

### OverSampling



Geometrically :



X : +ve

X : -ve

X : oversampled/  
Repeated +ve points

Placing the minority points on top  
of existing

∴ OverSampling is placing more points from the minority class in the dataset.

Is there other way of oversampling?  
 → Instead of Repeating, take artificial or synthetic points within the minority region



$X$ : -ve point

$\textcolor{blue}{X}$ : +ve point

$X$ : extrapolating  
synthetic points  
in same region  
as minority (+ve)

One more variation of oversampling

Class Weight:

$D_n \rightarrow$  100 +ve points  
 $\rightarrow$  900 -ve points

I would give weights to points:

Weight (+ve) = 9; More weight to Minority  
 Weight (-ve) = 1; Less weight to Majority

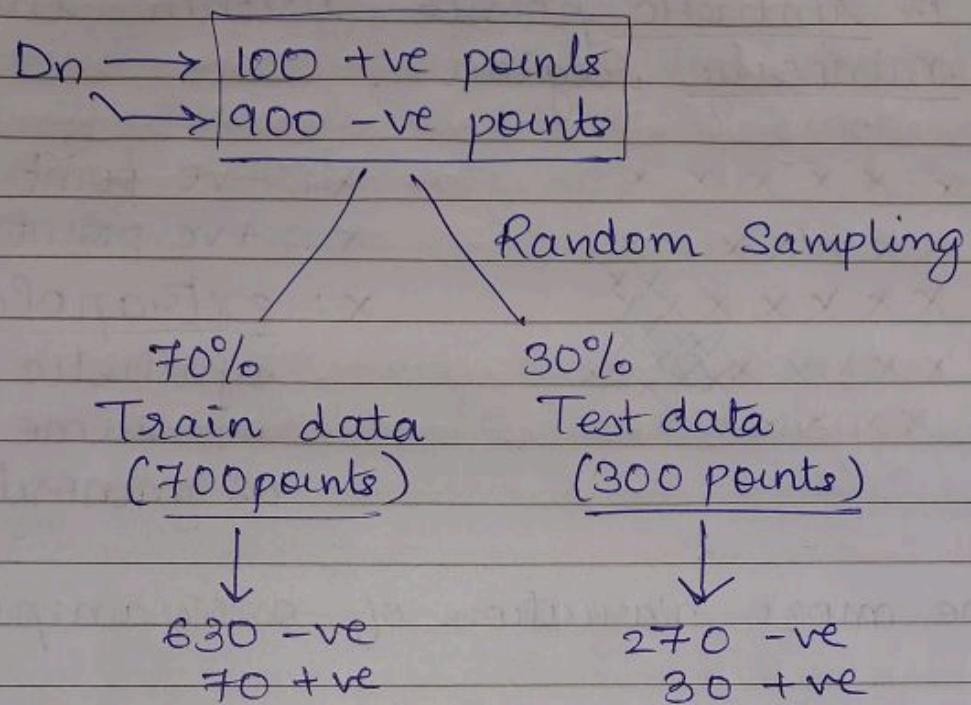
$X \times \times \times \rightarrow$  count it as 1 point

$\times \times \times \textcolor{red}{\circled{X}} \rightarrow$  count it as 9 points

$\times \times \times$

This trick is exactly same as Repeating technique of oversampling.

Other problem because of Imbalanced data :



We know Metric called Accuracy

Lets say my model  $f(x)$  says that every point is -ve

$f(x)$ : Every point is -ve {Dumb model}

$$\text{Accuracy} = \frac{\text{Test data}}{\frac{270}{30}} = 90\% \text{ (Very High)}$$

So, we would get very high Accuracy for DUMB ~~Accura~~ model.

Imbalanced dataset occur a lot of time in Real world.

### 31.3 Multi class Classification

Binary Classifica<sup>n</sup>:  $y_i \in \{0, 1\}$

Multi class Classifica<sup>n</sup>  $y_i \in \{0, 1, 2, \dots, 9\}$   
Eg: MNIST

K-NN:

$$D_n = \{(x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{1, 2, \dots, c\}\}$$

c-class classifier

Let's say for K-NN ; there are total c classes;

Given a query point  $x_q$ .

$$x_q : [0 | 6 | 1 | 0 | \dots | 0] \quad \begin{matrix} 1 & 2 & 3 & \dots & c \end{matrix}$$

6 nearest neighbor belong to class 2;  
Rest belong to class 1.

Assume its 7-NN;  
So, by majority vote,  
 $\hat{y}_q = 2$  (class 2)

### Probability Classifier

$$P(y_q = 2) = \frac{6}{7}$$

$$P(y_q = 3) = \frac{1}{7}$$

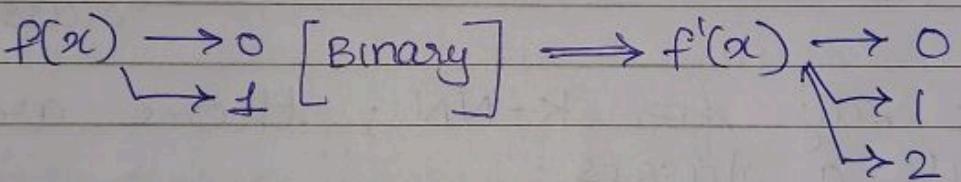
$$P(y_q = 1) = 0 = P(y_q = 4) = P(y_q = c)$$

Hence k-NN can easily be extended for multi-class classification.

There are classifiers like  
Logistic Regression → Binary classifier  
    └ Cannot do  
        • Multiclass classifier easily.

Q] Given a multiclass classification problem, can we convert it to a binary classification problem?

Imagine we are given



Let's say, we are given Multiclass prob  
 $y_i \in \{1, 2, 3, \dots, C\}$

We want to convert it to Binary prob.

② Similarly I can create Binary classifier for  $y_i = 2$  &  $y_i \neq 2$

③ Again, we will do for  $y_i = 3$  and  $y_i \neq 3$

So, if I have  $c$  classes, I will build  $c$  binary classifiers, I can build and train  $c$  models.

$\therefore$  Multi-class

Problem.  $\longrightarrow$   $c$ -Binary classifier  
( $c$ -classes) Problems

Such a technique  
is called

1 Vs Rest

$\rightarrow f_1(x) \rightarrow$  class 1 or not  
 $\rightarrow f_2(x) \rightarrow$  not 2 or not  
 $\vdots$   
 $\rightarrow f_c(x) \rightarrow$  not  $c$  or not

3.4 K-NN; Given a distance or similarity Matrix

During classification;  $x_i \in \mathbb{R}^d$   
 $\hookrightarrow$  Vector is not easy to frame

Let  $x_i$ : chemical compound  
(some paracetamol)

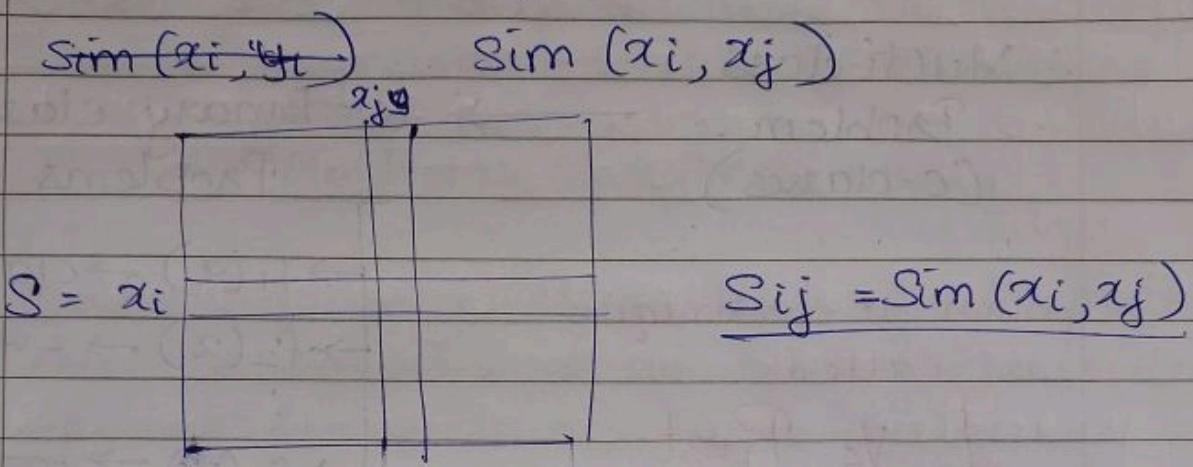
$y_i : 1 \rightarrow$  useful in curing fever.

It won't be easy to represent  $x_i$  in vector form.

A Pharmacist would be able to tell similarity between two chemical compound.

But Hence

Sometimes, we don't get  $x_i$  as vector  
We get similarity matrices



Once we receive similarity matrices,  
we can convert that into distance  
matrices

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

Q Given the similarity matrix (s)  
OR distance Matrix (d) instead  
of  $x_i \in \mathbb{R}^d$ : can K-NN work?

Ans: Yes.

$x_i$	$x_1$	$x_2$	$\dots$	$x_n$
$x_i$	✓	✓	✓	✓

← dist  
matrix

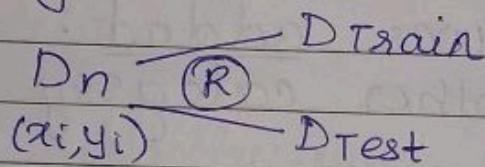
$$\begin{aligned} d(x_i, x_1) &= \\ d(x_i, x_2) &= \\ d(x_i, x_3) &= \\ \vdots & \\ d(x_i, x_n) &= \end{aligned}$$

} kNN will pick up  
the Neighbours  
Nearest to  $x_i$  based  
on the distances

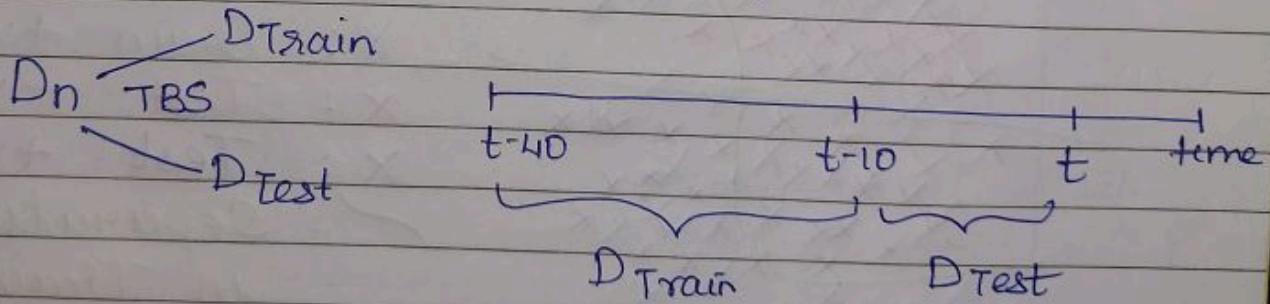
31.5

## Train and Test Set Difference

Given  $D_n$ , we have seen random splitting



The problem happens when we have time based slicing (TBS)



Eq: Amazon food Reviews

$D_{\text{Train}}$  &  $D_{\text{Test}}$  could be very different.

Eq: In food reviews, data might have changed in past.

## Food Reviews

~~t-10 to t~~

$\rightarrow D_{Test}$

$(t-10) \text{ to } (t)$  :- new category of products

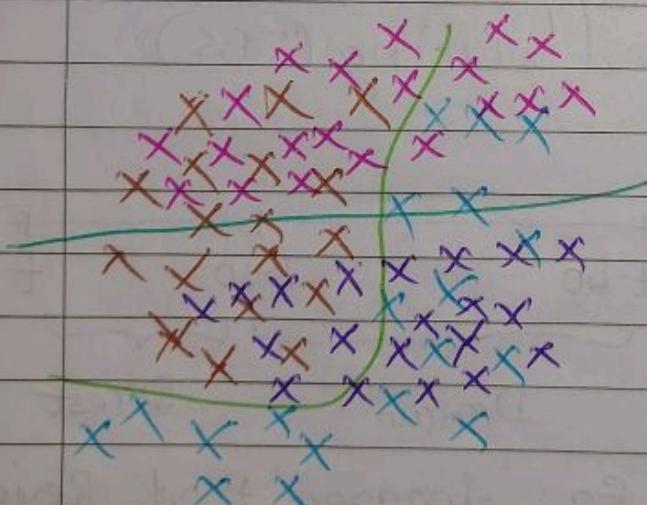
$(t-40) \text{ to } (t+10)$  : old category of products

$D_{Train}$

$\rightarrow$  Data changes Over Time

Ex: In  $(t-10) \text{ to } (t)$ ; new category of wine was added.  
And some other category was removed.

## Difference between Training & Test data



$x$ : Train -ve

$x$ : Train +ve

$x$ : Test -ve

$x$ : Test +ve

✓ : Separation  
in Train data

~ : Separation  
in Test data

$D_{Train}$  &  $D_{Test}$  are fundamentally different.

Distribution of -ve data has changed from  $D_{Train}$  and  $D_{Test}$ .

So if I use  $\checkmark$  line on Test data, it performs pathetically

$\therefore$  C.V. Error = Low  
But Test Error = high

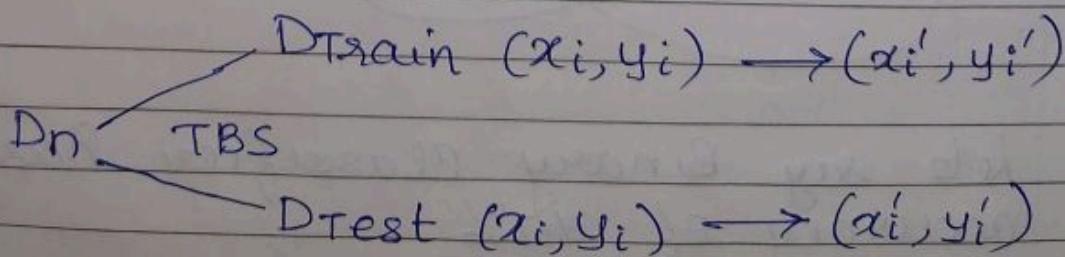
One check we need to do:

Do  $D_{Train}$  and  $D_{Test}$  have same distribution?

If the distributions are changing then we cannot use the same Model.

Q] How to determine if data is changing over time?

i.e. How to test if  $D_{Train}$  and  $D_{Test}$  have diff. distribu<sup>n</sup>?



For each  $D_{Train}$  and  $D_{Test}$  point, we will create a new point  $(x'_i, y'_i)$

I will create a new dataset  $D_n'$  such that

$D_n' := D_{Train} : - y_i' = 1 ; x_i' = (x_i, y_i) \text{ in } D_{Train}$

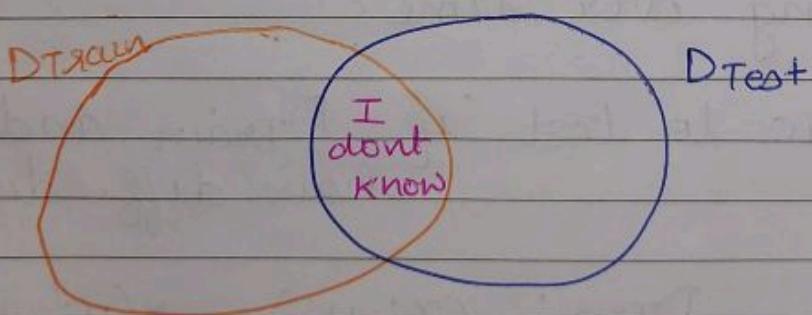
$D_{Test} : - y_i' = 0 ; x_i' = (x_i, y_i) \text{ concat in } D_{Test}$

Now, build a binary classifier in  $D_n'$

we can use any  
Binary classifier

Eg : KNN

Binary classifier will be able  
to separate  $D_{Train}$  and  $D_{Test}$   
to some extent

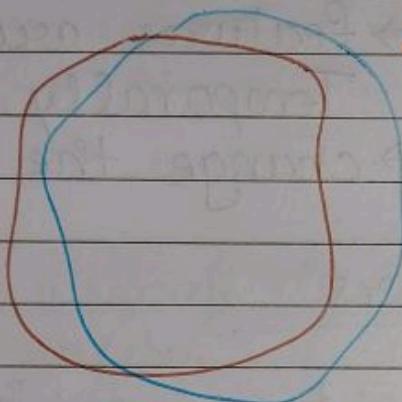


Let's say Binary classifier has  
accuracy of 70%

Means 30% of points are in "I don't know"  
region

70% Accuracy → Train and Test data  
can be separated  
partially

If we can separate them partially  
it means their distributions are  
different.



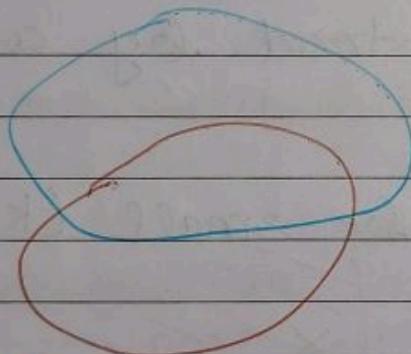
Almost Overlapping

We would get

LOW ACCURACY

It means the  
DISTRIBUTIONS are

VERY SIMILAR



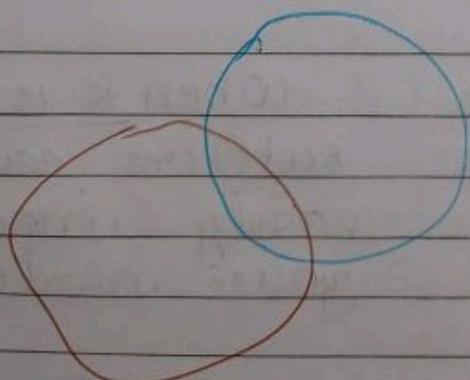
Less Overlap .

Binary classifier

accuracy is Medium

Distributions are

NOT VERY SIMILAR



Overlap is Very Low  
Accuracy is HIGH

distributions are

VERY DIFFERENT

Hence

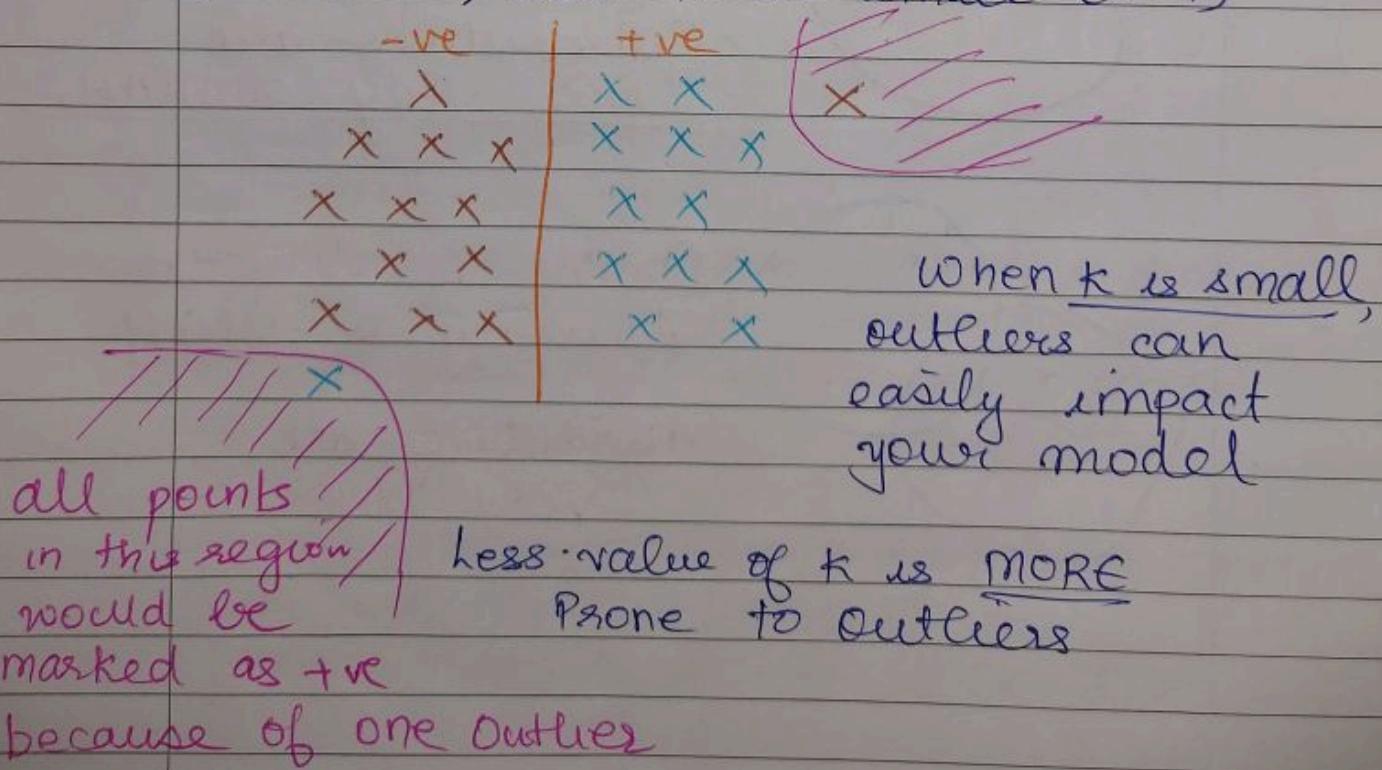
- $D_{Train}$  &  $D_{Test} \rightarrow$  same distributions ✓
- $\rightarrow$  not same distributions
- $\rightarrow$  features are changing with time
- $\rightarrow$  features are not Temporally stable
- $\Rightarrow$  change the features

### 3I-7

### Impact of Outliers

A model is a function which can be easily understood by a decision surface.

For k-NN, let  $k$  be small ( $k=1$ )



Let's say we are doing 10-fold CV.

	<u>Accuracy</u>
$K=1$	97%
$K=2$	97%
$K=3$	97%
$K=4$	97%
$K=5$	97%
$K=6$	95%
$K=7$	92%

→ we should choose largest value of  $K$  with Maximum Accuracy because lower values of  $K$  are more prone to outliers.

Other way to handle outliers:

Imagine if we can detect outliers and remove them → that's great. Then we don't have to worry about them. But how do we detect them?  
 → Visualizing the data (possible in 2D but not in higher dimension)

Technique to detect outliers  
 → Local Outlier Factor.

### 31.7 Local Outlier Factor (Technique to detect outliers in data)

→ Inspired by ideas like k-NN.

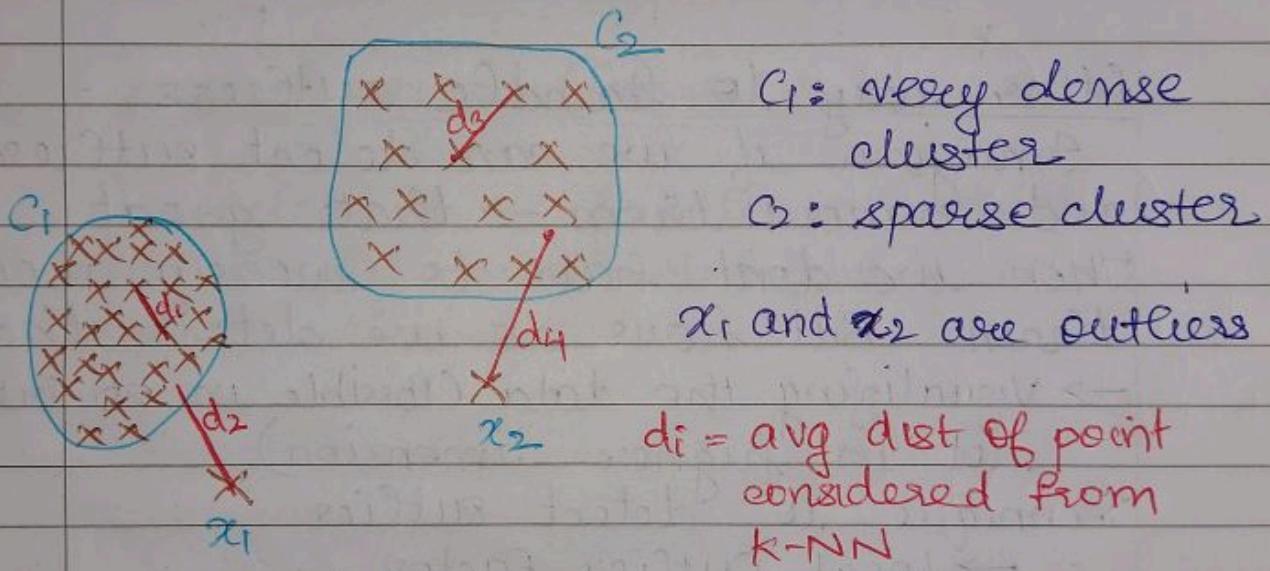
Let's take a classification setting

→ Take all -ve points. Find outliers.

→  $\rightarrow$  +ve  $\rightarrow$

→ Remove them.

Consider group of -ve points and their cluster



#### Simple selection

→ Take each point  $x_i$

→ Take its k-Nearest Neighbours  
 Let's  $k=5$

→ Find avg distance to all 5 points

→ Sort all the points as per the distances.

→  $d_4$  is large. Hence its outlier.

Remove that point.

i.e. If Avg dist is high, then the point is outlier

## Problem with this approach

Now  $x_2$  will be declared as an outlier.  
 But  $x_1$  will NEVER be declared as outlier  
 Because  $d_2 \approx d_3$ ; it is close to dense cluster.

If we have to declare  $x_2$  as outlier,  
 then all points with distance  $d_2$   
 would be considered as outlier.

We have to compute local density

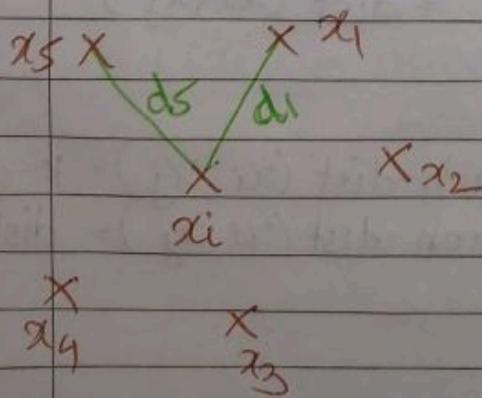
The current method fails when we have more than one cluster.

One of the clusters is dense and other is sparse.

## 31.8 K-distance

K-distance ( $x_i$ ) = distance of the  $k^{th}$  Nearest Neighbour of  $x_i$  from  $x_i$

$N(A)$  = Neighbour of  $A$  = set of all point belong to k-NN of  $x_i$



$d_5$  = 5-distance.

$d_1$  = 1-distance

$$N(x_i) = \{x_1, x_2, x_3, x_4, x_5\}_{k=5}$$

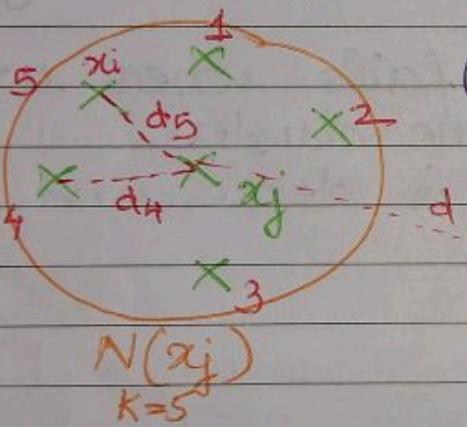
### 31.9. Reachability distance (A, B)

Reachability Distance ( $x_i, x_j$ ) =  $\max \left[ k\text{-distance} (x_j), \underbrace{\text{dist} (x_i, x_j)}_{\substack{\text{Distance of} \\ \text{kth NN of} \\ x_j \text{ from } x_j}} \right]$

$\substack{\text{Actual} \\ \text{distance} \\ \text{between} \\ x_i \text{ and } x_j}$

Geometrically,

Assume  $x_j$  and  $k=5$



① if  $x_i \in N(x_j)$   
then.

$$\text{Reach-dist} (x_i, x_j) = \max (d_5, d_4)$$

$= d_5 \quad \left\{ \begin{array}{l} \text{Because } d_5 \\ \text{is greater} \\ \text{than } d_4 \end{array} \right\}$

$$= k\text{-dist} (x_j)$$

② else

$\{ x_i \notin N(x_j) \Rightarrow x_i \text{ is not in Neighbourhood} \}$

$$\begin{aligned} \text{Reach-dist} (x_i, x_j) &= \max (d_5, d) \\ &= \frac{d}{\text{dist} (x_i, x_j)} \end{aligned}$$

To sum it up:

if  $x_i \in N(x_j)$  then  $\text{Reach-dist} (x_i, x_j) = k\text{-dist} (x_j)$   
else.  $\text{Reach-dist} (x_i, x_j) = \text{dist} (x_i, x_j)$

31.10

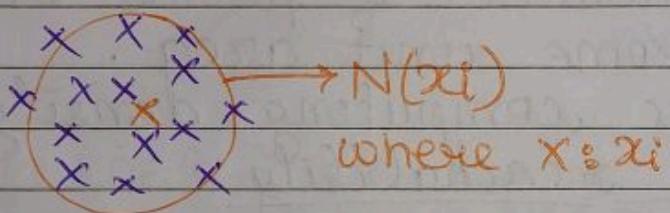
Local Reachability Density : lrd( $x_i$ )

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

↓  
All points  
belonging  
to  $N(x_i)$

Number of  
elements in  
Neighbourhood ( $x_i$ )

This formula looks like Average.



$$N(x_i) = \text{Set of } k\text{-NN to } x_i \\ \therefore |N(x_i)| = K$$

But Not always ; because 2 points may have same distance from  $x_i$

So, we can think of LRD as  
[Average Reachability distance of  
 $x_i$  from its neighbourhood]

lrd( $x_i$ ) = Inverse of Average  
Reachability distance  
of  $x_i$  from its  $k$  Neighbourhood

We can also write LRD as

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

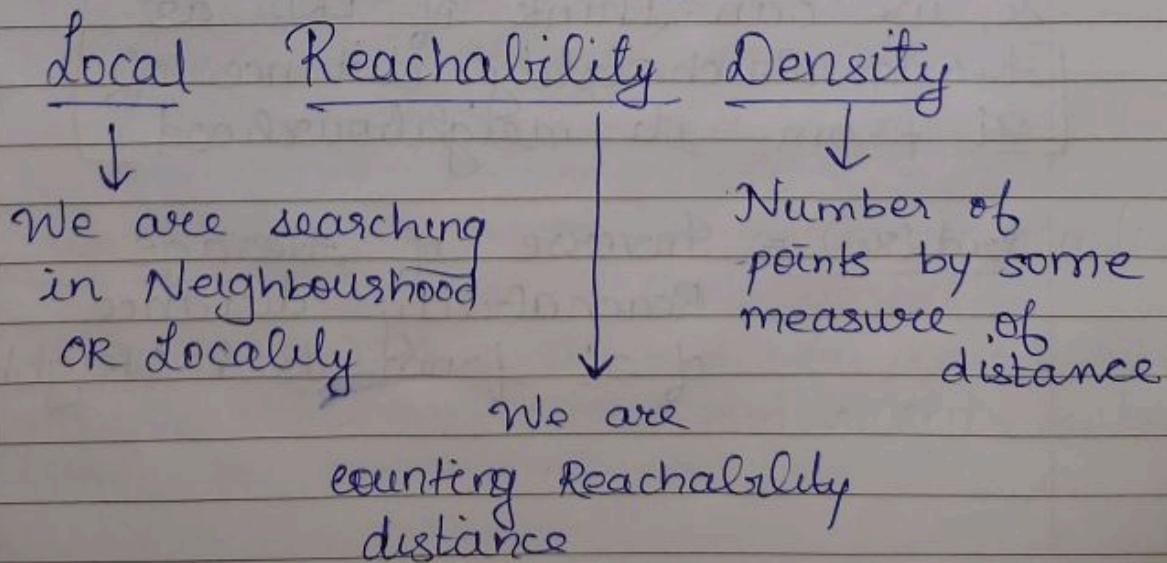
↓ some measure  
 of distance

Intuitively, if we think about it, I am measuring density.  
 Density  $\rightarrow$  number of points in some unit area.

Here we are computing density using local reachability

In denominator, I am counting the sum of reach-dist in the locality of point  $x_i$ .

Hence we are calling it as



31.11

Local Outlier Factor ( $x_i$ )

$$\text{LOF}(x_i) = \frac{\sum_{x_j \in N(x_i)} \text{ld}(x_j)}{|N(x_i)|} * \frac{1}{\text{ld}(x_i)}$$

Average ld  
of points in  
Neighbourhood  
of  $x_i$

If  $\text{LOF}(x_i)$  is large when.  
Avg. component is  $\uparrow$  OR  $\text{ld}(x_i)$  is  $\downarrow$

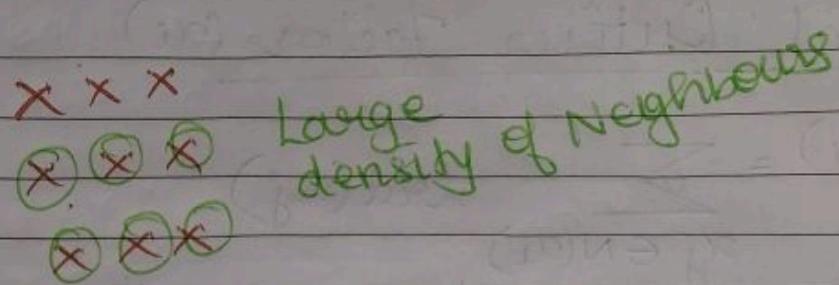
(Density of points in Neighbourhood is large) (My  $\overset{(x_i)}{\text{out}}$  density is small)

$\text{LOF}(x_i)$  is small when.

Avg Comp is  $\downarrow$  OR  $\text{ld}(x_i)$  is  $\uparrow$   
 density around  $x_i$  is Large .

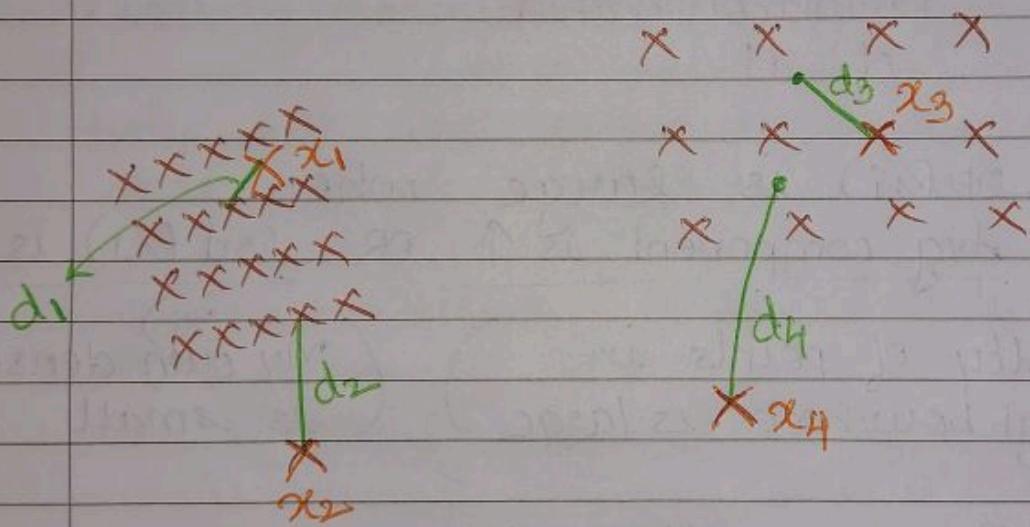
If  $\text{LOF}(x_i)$  is large.  $\rightarrow$  outlier  
 is small  $\Rightarrow$  outlier

[density ( $\text{ld}$ ) of  $x_i$  is small compared to its neighbours]



$\rightarrow$  small density compared to Neighbours

Other example:



$x_1$ : point in dense cluster

$x_2$ : point slightly away from dense cluster

$x_3$ : point in sparse cluster

$x_4$ : point far away from sparse cluster

Let's take  $x_1$

All of its neighbours are very close to each other

lets say  $K=5$ .

Assume avg distance to 5-NN =  $d_1$ .  
i.e. Average Reachability dist =  $d_1$

Avg Reach-dist of  $x_1 \rightarrow d_1$   
 $x_2 \rightarrow d_2$   
 $x_3 \rightarrow d_3$   
 $x_4 \rightarrow d_4$

Intuitively, we understand that  
 $d_1 < d_2 < d_3 < d_4$

We know  $\text{erd} = \frac{1}{\text{avg Reach-dist}}$  { Intuitively }

Since  $d_1 < d_2 < d_3 < d_4$

$$\therefore \frac{1}{d_1} > \frac{1}{d_2} > \frac{1}{d_3} > \frac{1}{d_4}$$

$\downarrow$        $\downarrow$        $\downarrow$        $\downarrow$   
 $\text{erd}(x_1)$      $\text{erd}(x_2)$      $\text{erd}(x_3)$      $\text{erd}(x_4)$

Now  $\text{erd}(x_1)$  would be very similar to  $\text{erd}$  of all points close to  $x_1$ .

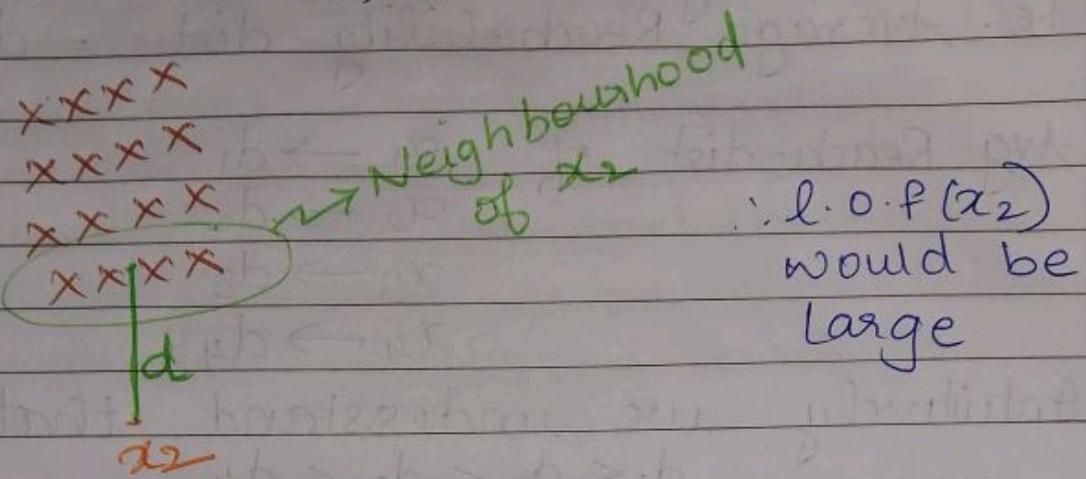
$$\text{erd}(x_1) \approx \text{erd}(x_i \in N(x_1))$$

Density around  $x_1$  would be roughly same as point ~~count~~ density around  $x_1$ .

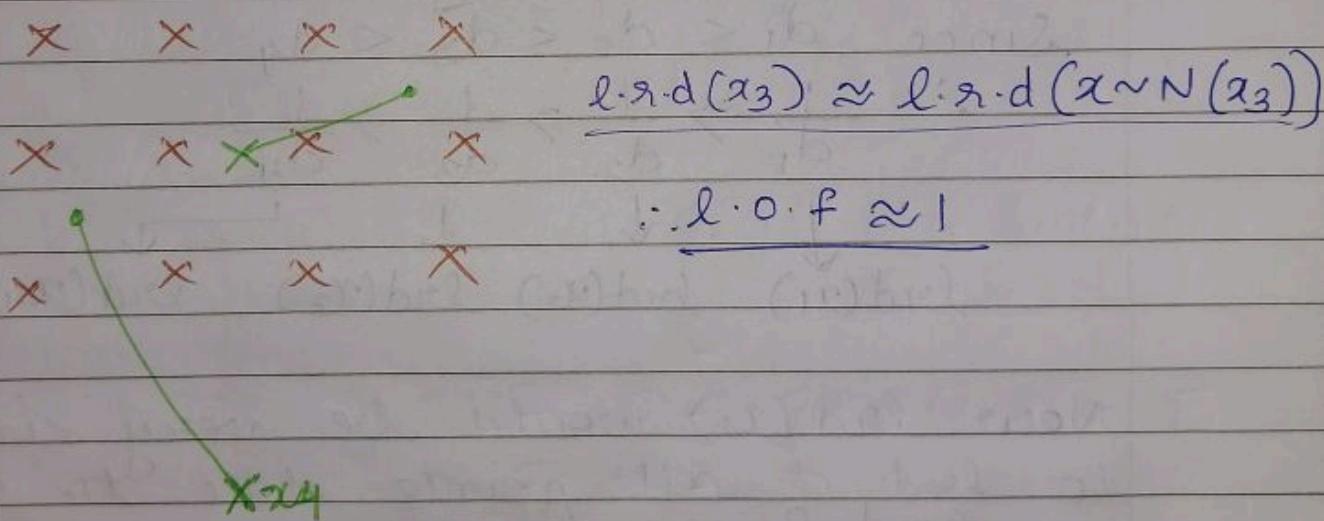
$\therefore \text{lod}(x_1)$  would be small  
 $\Rightarrow$  close to 1

About  $x_2$ :

$$\text{l.r.d.}(x_2) \leq \text{l.r.d.}(\{x_i \in N(x_2)\})$$



About  $x_3$ :



About  $x_4$

$$\text{l.r.d.}(x_4) \ll \text{l.r.d.}(\{x \in N(x_4)\})$$

∴ l.o.f( $x_4$ ) is large

- ∴  $L.O.F(x_1 \in c_1) \rightarrow$  small
- $L.O.F(x_2) \rightarrow$  Large  $\rightarrow$  Outlier
- $L.O.F(x_3 \in c_2) \rightarrow$  small
- $L.O.F(x_4) \rightarrow$  Large  $\Rightarrow$  Outlier

$\therefore$  Large L.O.F  $\Rightarrow$  Outlier

Hence, how do we determine Outliers ?

- $\rightarrow$  For each point  $x_i$ , compute  $L.O.F(x_i)$
- $\rightarrow$  Pick points with highest L.O.F.
- Most of them would be outliers
- $\Rightarrow$  It is very subjective to decide the value of L.O.F to be considered as outlier (on where to set the threshold)
- $\Rightarrow$  We would require Domain Knowledge

### Disadvantages of L.O.F

1. There is no interpretability of L.O.F value means, we know that if value is large, then its bad. But how bad it is
2. ~~to~~ L.O.F value could lie anywhere from 0 to  $\infty$ . Imagine we have a modified L.O.F which gives values from 0 to 1. This could act as a probability of a point which could be a outlier. This is called Local Outlier Probability.

### 31.12 Impact of Scale and Column Standardization

Lets take a dataset with 2 features

	$f_1$	$f_2$
$x_i$	10	10
$x_j$	10	100
	~~~~~	~~~~~

Lets say values  
 values lie between  
 between 0 to 1  
 0 to 100

This is called scale difference between  
 $f_1$  and  $f_2$ .

Lets say  $x_1 = [23, 0.2]$   $\text{dist}(x_1, x_2) = 5$   
 $x_2 = [28, 0.2]$   
 $x_3 = [23, 1.0]$

$$\text{dist}(x_1, x_2) = 5$$

$$\text{dist}(x_1, x_3) = 0.8$$

Now there is a problem.

The distance varies because of  
 scale difference between  $f_1$  and  $f_2$

If  $f_1$  could have value from 0 to 100 so  $x_1$  and  $x_2$  are not very different from each other. Because out of the possible '100' differences, they have difference of only 5.

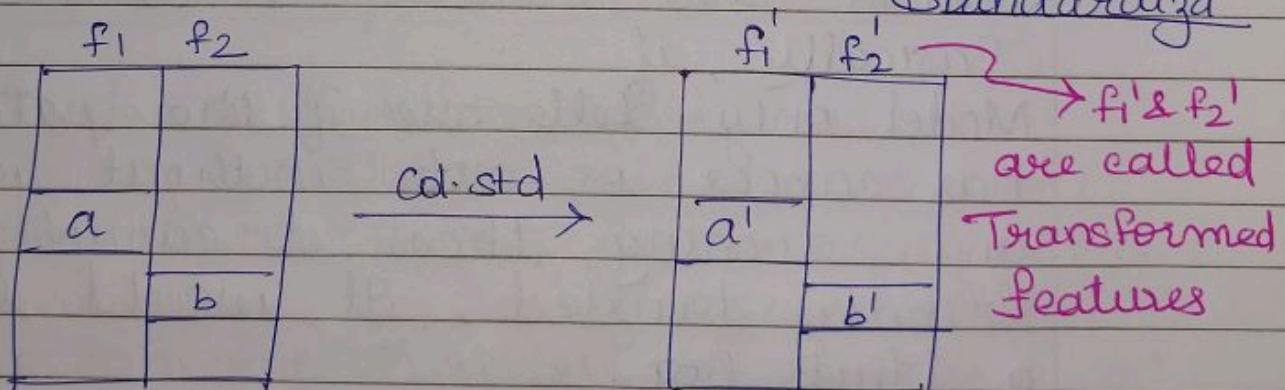
Between  $x_1$  and  $x_3$ , max diff possible is 1. Because  $f_2$  ranges from 0 to 1. Now  $\text{diff}(x_1, x_3) = \text{dist}(x_1, x_3) = 0.8$  so,  $x_1$  and  $x_3$  differ almost 80% of max allowable difference.

∴ Thinking logically :

$$\underline{d(x_1, x_3) > d(x_1, x_2)}$$

Way to solve this issue  $\Rightarrow$  Column

Standardizn



such that  $\text{Mean}(f_1') = 0 = \text{Mean}(f_2')$   
and  $\text{Std-dev}(f_1') = 1 = \text{Std-dev}(f_2')$

$$\boxed{a' = \frac{a - \mu_1}{\sigma_1} \quad \& \quad b' = \frac{b - \mu_2}{\sigma_2}}$$

where  $\mu_1, \sigma_1 \Rightarrow f_1$  { Mean & Std Dev }  
 $\mu_2, \sigma_2 \Rightarrow f_2$  { }

Since Euclidean dist can be impacted by differences in scale, and

Euclidean dist is used in k-NN then

[it is important to do column standardization before k-NN]

There are some techniques like Decision Trees which are scale independent.

31-13

## Interpretability

Consider a medical application where the model tells if the patient has cancer or not.

Basically, if

Model only tells me if the patient has cancer or not without giving any reasoning, then it cannot be blindly trusted. It would be a black box

$$x_q \rightarrow f \rightarrow y_q$$

Doctor has to be very sure.

If model gives only class label, then it is insufficient ~~now~~ without reasoning

Model should give reasoning on  
why  $y_q = 1/0$   
It would be very useful to doctor.

Let's say  $x_q$  has 10 features  
and model tells that

$f_5$  corresponds to test 5 and value is high  
 $f_8$  → test 8 → Low.

& hence, the patient has cancer ( $y_q = 1$ )  
Doctor understands these tests and those  
are medical test.

A Model which gives an explanation  
are called Interpretable Model.

Model which gives no explanation  
are Black Box model.

Is K-NN interpretable model or not?

Depends on Domain Model.

Let's say  $K = 7$

Now  $x_q \rightarrow y_q = 1$ ; Reasoning: all 7 NN  
of  $x_q$  are labelled as 1  
We can give 7 vectors of each NN

K-NN is interpretable when  $d$  is small  
and  $K$  is small.

This changes from domain to domain  
→ True in Medical Domain

31.14

Feature Importance & Forward Selection

Lets say we have 10 features  
 $f_1, f_2 \dots f_{10}$

$D_n^{10} \rightarrow f$  (model)  $\rightarrow$  if the model can tell which features are more useful.

Ex: If model  $f$ : predict height and features are  $f_1$ : <sup>weight</sup> height of person  
 $f_2$ : hair colour  
 $f_3$ : gender  
 $f_4$ : country  
 $f_5$ : skin colour

useful feature

Feature Importance  $\rightarrow$  sort features by important for classifica<sup>n</sup> task in order of importance.

- $\rightarrow$  Useful in understanding model better
- $\rightarrow$  Interpretability increases

k-NN doesn't give us feature importance by default.

Forward Feature Selection

Given large dimension, people try to reduce dimension.

1000-dim  $\rightarrow$  100 dim

One way is PCA  
But PCA doesn't return the useful or important features.

One technique to get that is forward feature selection

$$D = \{\text{Training}, \text{Test}\}$$

- ① We have all features  $f_1, f_2, \dots, f_d$ .
- ② Assume our model is K-NN
- ③ I will train model only using  $f_1$

only  $f_1 \rightarrow$ 

Tr.	Test
model	

 → Accuracy on Test data using only  $f_1 = a_1$

only  $f_2 \rightarrow$ 

Tr.	Test
Model	

 →  $= a_2$

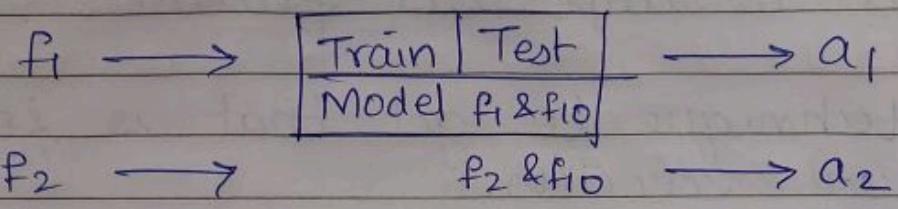
only  $f_d \quad \dots \quad = a_d$

- ④ Check which of them has highest accuracy.  
say  $f_{10}$  has highest Accuracy.

- ⑤ Now my model becomes

$\{f_1, f_2, \dots, f_9, f_{10}, \dots, f_d\}$ <span style="color: red;">Rest of feature</span>	Train   Test Model   $f_{10}$
	<span style="color: red;">→ successful feature</span>

⑥ Will train a model using  $f_1$  and  $f_{10}$



Now, say  $a_5$  is largest.  
So, will finalize  $f_{10}$  and  $f_5$  are best features.

⑦ At end of 2 Iterations ;

$\{f_1, f_2 \dots f_4, f_6 \dots f_d\}$

Remaining Features

$\{f_{10}, f_5\}$

Selected Features

One thing which people argue is  
→ In first iteration, we got the accuracies for features as  
 $f_{10} > f_6 > f_8 \dots$

So why didn't we take  $f_6$  as second best feature.

It is because if  $f_6$  and  $f_{10}$  are exactly same/similar feature (in error) then they would have same/similar accuracy

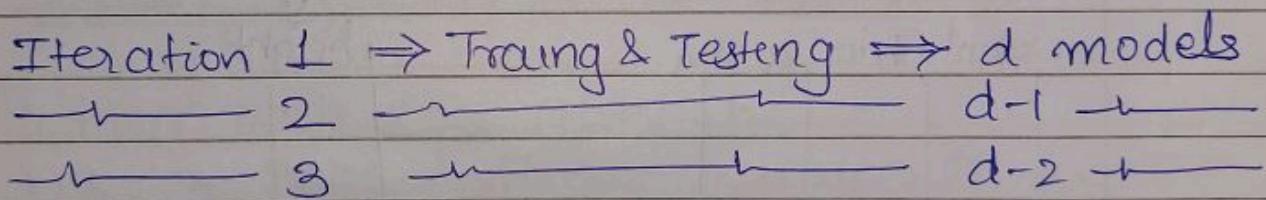
At each stage, we ask that given the best feature, which new feature adds the most value (highest accuracy).

Backward Selec<sup>n</sup>:

It takes all features and then removes feature in every iteration which adds least value.

Considering Time Complexity:

At each iteration, we are training and testing ' $d$ ' models.



Hence a LOT of time complexity is there. We should use it only when we have time.

Forward/ Backward Selection ~~to~~ techniques are Model Agnostic  $\Rightarrow$  means they dont care which model it is.

Disadv: Lot of time

## 31.15 Categorical and Ordinal features

We remember Iris dataset which had 4 features: PL, PW, SL, SW

## Real Valued Numbers

and we are given datapoint  $x_i \in \mathbb{R}^d$

But in real world, data is not always real values.

Eq: we are given dataset and asked to predict height

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

- ① weight  $\Rightarrow$  numerical  $\Rightarrow$  we can take value  
as is

② Hair colour  $\Rightarrow$  non-numerical  
 $\Rightarrow \{$  Black, Brown, Golden, ...  $\}$

① Give a number to each colour

② something like  
 $\{$  Black, Brown, Red, Golden, Gray  $\}$

1      2      3      4      5

Numbers are ordinal ; means they can be compared

Something like  $3 > 1$  which would imply Red  $>$  Black  $\Rightarrow$  Absurd.

One solution to it : One hot encoding

I would create an array for hair colour

$x_i : [0 \ 1 \ 0 \ 0 \ 0]$  if  $x_i : \text{Brown}$ .  
 Black Brown Red Golden Gray

So one hot encoding is the binary Vector of the size of Number of distinct elements

Problem with this approach :

Consider the feature country ; there are total 200 countries

$\therefore$  Vector would be of size 200

country vector :  $\boxed{\quad} \mid \dots \mid \boxed{\quad}$   
 $c_1 \ c_2 \ \dots \ c_{200}$ .

Also, the vector would be very sparse.  
 As the person would belong to only 1 country.

If the number of distinct values for a categorical feature is large then One hot encoding can create sparse and large vector.

(3)

One hack is to replace country vector with Avg height of people in the country

because our objective is to create a model to predict the height of a person.

Replace the country with Avg height

We are taking Mean value of

height for each category.

Its Mean Replacement technique

Another way to convert categorical feature to domain feature is through domain knowledge.

Eg: Country  $\xrightarrow{\text{Replace}}$  distance from India  
 (categorical) (numeric).

(A) Let's say we know that distance from India is related to height.

Ordinal Feature

Let's say a feature is

f: How do you rate Indian food?  
 {v.Good, Good, Avg, Bad, v.Bad}

↓  
 Categorical Feature

There is Logical ordering among features

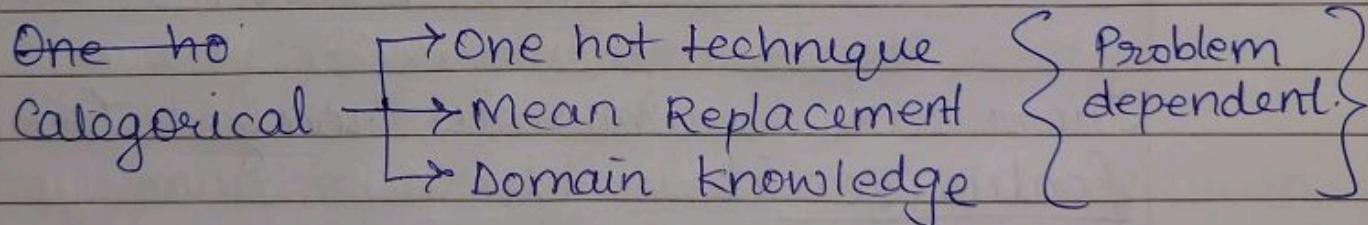
So, we can convert to Numeric  
Like: V.Good = 5 . . . . . , V.Bad = 1

Problem with this approach

$f_i$	$f_i'$	$f_i''$	
V.Good	5	10	} We can give any numeric values as long as they are <u>ordinal</u>
Good	4	6	
Avg	3	3	
Bad	2	2	
V.Bad	1	1	

The assignment of values depend on domain knowledge.

So we have following techniques



31/16

Missing Values Imputation

Let's assume I have 4 features

	$f_1$	$f_2$	$f_3$	$f_4$	$y$	
$x_i$	✓	✓	✗	✓		→ Missing value NaN, null, -1

And value is not present for one feature.

It may be due to

- Data corruption ↗ Very
- Collection error ↗ Common

There are numerous ways to featureize it :

① Imputation: 3 ways to do it

- Mean replacement
- Median replacement
- Most frequent value (mode)

Strategy: Take all non-missing value for the feature ; Take Mean, Median, mode .

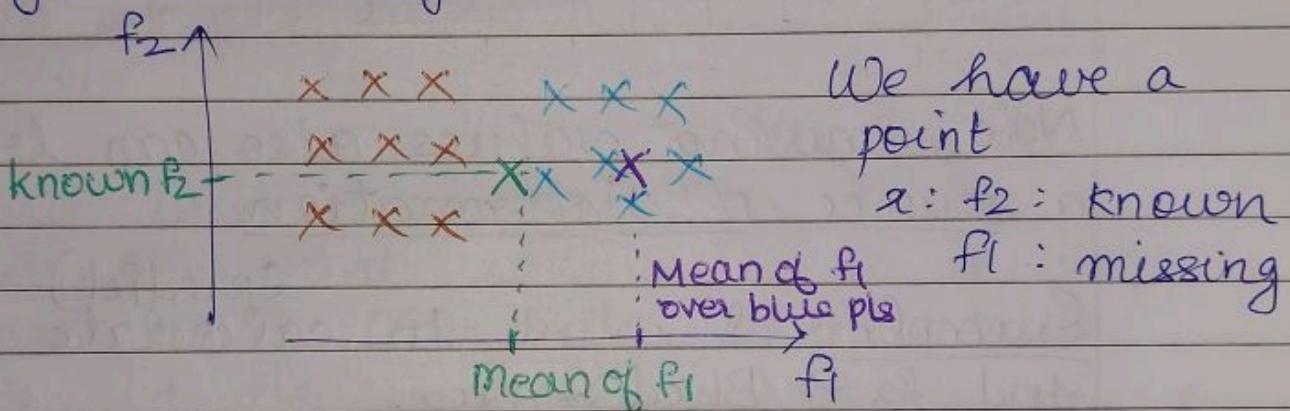
sklearn. → Imputer Method

while doing a classification task  
where the output is 0, 1

→ We impute based on class label.  
→  $x_i \rightarrow f_3$ : Missing and  $y_i = 1$

then instead of computing mean of all  $f_3$ ,  
compute the average of those  $f_3$  for  
which  $y_i$  is 1

Geometrically



lets say I calculated  $f_1$  as the mean of all  $f_1$  values, then my point  $x$  would lie over  $\times$  as shown

Now, we know that  $y = f(x) = 1$  (class)  
So, we would calculate mean of only those points which has  $y_i = 1$ .

so, my point  $x$  would lie in  $\otimes$  as shown

Logically  $\times$  is correct placement.  
given  $y = 1$

### ③ New Missing Value feature

for given dataset

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

$f_3$  is missing  
 $f_2$  is missing

Now, missing values also can be a source of information

Example: I want to calculate height  
And  $f_3 \rightarrow$  Blood Group.

If my feature  $f_3$  is missing for a particular region, it would mean the people are below poverty line and they don't have ~~expenses~~ money to conduct the blood test

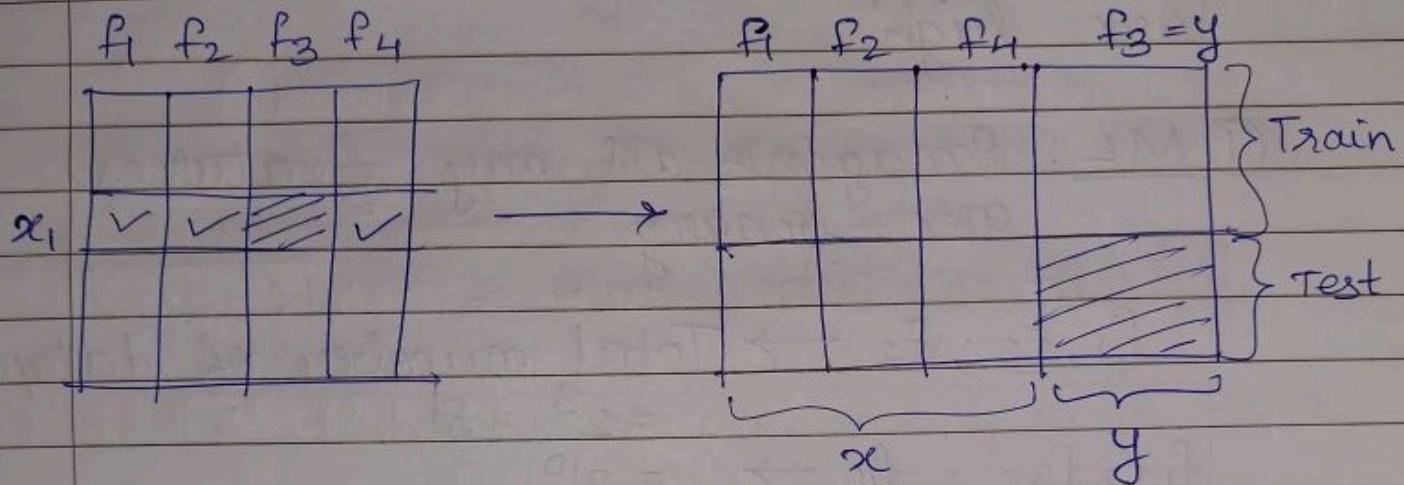
What we can do is add another set of features ( $f'_1, f'_2, f'_3, \dots$ ) which would indicate if the feature is missing

	$f_1$	$f_2$	$f_3$	$f_4$	$f'_1$	$f'_2$	$f'_3$	$f'_4$
$x_1$	✓	✓	✓	✓	0	0	1	0
$x_2$	✓	✓	✓	✓	0	1	0	0

→ Indicates if the feature is missing

#### ④ Model Based Imputation

We can take any model (say KNN)



We will take  $f_3 = y$

All the points where  $f_3$  is missing would be Test Data. Rest points would be Train Data. We would calculate the  $f_3 = y$  <sup>for Test data</sup> based on Train data.

K-NN is used a lot of time for Model Based imputation

Because K-NN has a feature of Neighbourhood.

For given point  $x_1$ ;  $f_3$  is missing

So, if  $x_2, x_3, x_4$  are in neighbourhood, their features are similar to  $x_1$ .  
 $\therefore$  We get idea of  $f_3$  for  $x_1$

## Curse of dimensionality

Explains the multiple phenomena that happens when dimensions are high.

- ① ML: Imagine all my features are binary

$$f_1, f_3; f_2 \rightarrow \text{Total number of datapoints} \\ = 2^3 = 8$$

$$f_1, f_2, f_3 \dots f_{10} \rightarrow = 2^{10}$$

As dimensionality  $\uparrow$ , the number of datapoints to perform good classification/regression for a good model increases exponentially

## Hughes Phenomenon

If 'n'  $\rightarrow$  size of dataset is fixed  
 Performance decreases [as Dimensionality increases].

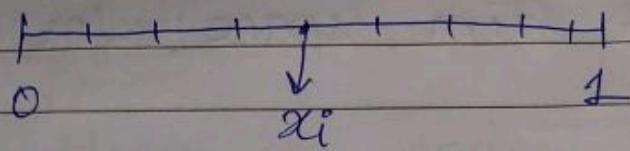
- ② Distance functions (esp. Euclidean dist)

The intuition we have about distances in 3D is not valid for high dimensional space.

171

171

Lets say in 1-D world: n-random points



If I take point  $x_i$

$$\text{Dist-Min}_{(x_i)} = \min_{x_j} \left\{ \text{dist}(x_i, x_j) \mid \text{where } x_i \neq x_j \right\}$$

It is distance of nearest point to  $x_i$

$$\text{Dist-Max}_{(x_i)} = \max_{x_j} \left\{ \text{dist}(x_j, x_i) \mid \text{where } x_j \neq x_i \right\}$$

Typically

$$\frac{\text{Dist-Max}(x_i) - \text{Dist-Min}(x_i)}{\text{Dist-Min}(x_i)} > 0 \quad \text{when } d=1, 2, 3$$

Now, when  $d$  increases significantly, the above term reaches near to 0.

$$\lim_{d \rightarrow \infty} \left( \frac{\text{distMax}(x_i) - \text{distMin}(x_i)}{\text{distMin}(x_i)} \right) \rightarrow 0$$

It means since the above terms becomes 0,

$$\text{distMax}(x_i) \approx \text{distMin}(x_i)$$

This is an alarming thing.

In high dim space, for n-random points

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



Every pair of points are equally distant from each other.

i.e.  $\text{dist}(x_i, x_j) \approx \text{dist}(x_i, x_k)$

In this scenario,  
k-NN does not make sense.  
Hence k-NN doesn't work well in  
high dimensional space.

Solution: Instead of Euclidean dist,  
use other distances like cosine  
similarity.

Hence, people prefer cosine  
similarity in high dimensional space  
problems like Text classification

Twist/caveat:

High dim & Dense → Impact of dimensionality  
is high

High dim & Sparse → Impact of dimensionality  
is lower  
non uniform/random  
spread of data

The entire concept is with assumption that points are distributed randomly and uniformly

### ③ Overfitting and Underfitting

As  $d \uparrow$ , the overfitting also increases.  
We would learn this during Linear Regression

- Solutions →
    - ① Forward feature selection  
Pick the most useful subset of features
    - ② Use dimensionality reduc<sup>n</sup>. like PCA, tSNE  
They do not use class label and are not classifica<sup>n</sup> oriented.
    - ③ Some people also change distance measures → cosine similarity instead of Euclidean distance.  
→ Sparse representation instead of dense representation.
- Bow ←  
W2V ←

31-18

## Bias Variance Tradeoff

We talked out underfitting & overfitting  
 $(K=1)$   $(K=n)$

How do we mathematically analyze underfitting & overfitting?

→ Bias Variance Tradeoff.

→ Most used idea in statistical ML.

We know generalization error as the error on future unseen data

$$\text{Generaliza}^n = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

↓  
 (future unseen data of M)

We would learn  
 Math proof of  
 same in linear  
 Regression

Irreducible error → error cannot be reduced further

Bias error → error due to simplifying assumptions.

High Bias  $\Rightarrow$  Underfitting

What is error due to simplifying assumption?



If my model assumption is the we use plane to separate the +ve from -ve points

The plane is not exact/perfect fit and it gives some errors

Such errors are called Bias error.

We can use a  $\sigma$  curve; if my model limit itself to plane, we get Bias errors.

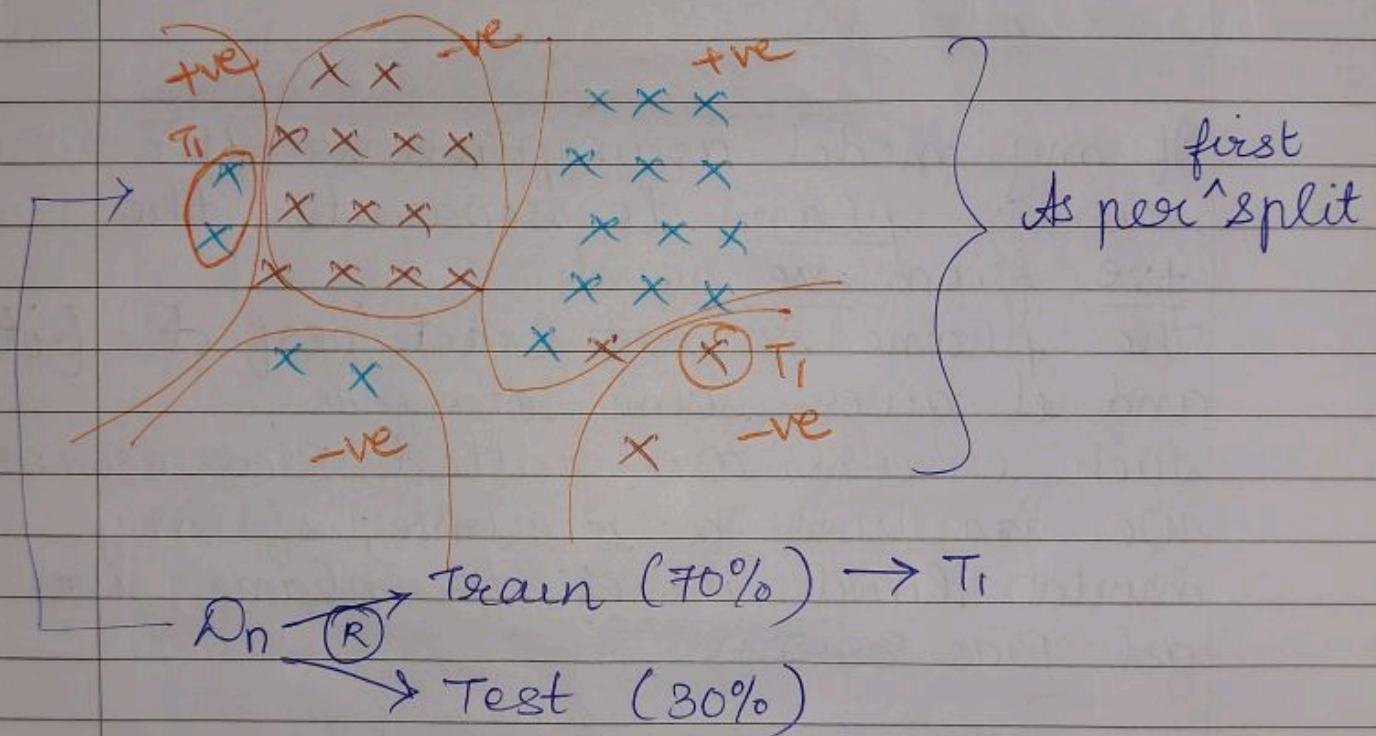
Let's say in K-NN, if  $K=n$  and 80% of points are -ve  
20% of points are +ve

So, for every query point  $x_q$ ,  $y_q = -ve$   
Hence, dominant class becomes the class label of  $x_q$ . Bias error would be high

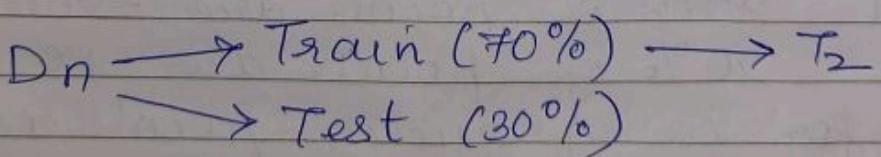
Hence, high Bias would be Underfitting

High Variance: means the difference in model OR how much a model changes as training data changes

Imagine I have following total data out of which I will select some data as training and other as test data.



I split  $D_n$  one more time



Model is nothing but decision surface

Now when one of the point from  $T_1$  is replaced by  $T_2$ , the decision surface changes.

When  $K=1$ , small change in dataset results in very diff model (= decision surface).

Small change in dataset → Large change in model → High Variance Model.

↓

Overfit Model

as  $K \uparrow$ ; Variance ↓

So if we want to :

$$\text{Generaliz}^n \downarrow = \text{Bias}^2 \downarrow + \text{Var} \downarrow + \text{Irreducible error} .$$

$\underbrace{\hspace{2cm}}$      $\underbrace{\hspace{2cm}}$

No                  No  
underfit      overfit

$K=1$ ;	Bias ↓	Var ↑
$\text{As } K \uparrow$ ;	Bias ↑ (slightly)	Var ↓ (drastically)
$K=n$ ;	Bias ↑ (Very High)	Var ↓ (Very Low)

Let's say

$K=1$ ;	10 + 100 + 3	= 113	overfit
$K=5$ ;	12 + 10 + 3	= 25 ✓	
$K=n$ ;	100 + 2 + 3	= 105	underfit

So, we want to design a model that balances both Bias & Variance

∴ Bias Variance Tradeoff  $\Rightarrow$  Theme of ML

→ Overfit Vs Underfit

### Intuitive Understand of Bias and Variance.

We have a dataset

D  $\xrightarrow{D_{\text{Train}} \rightarrow \underline{\text{model}}$   
 $D \xrightarrow{D_{\text{Test}}}$

Case 1: High Bias OR Underfit

Using the "Model", I can compute two types of errors.

Train Error: difference between  $y_i$  and  $\hat{y}_i$  on train data

Train Error =  $\text{diff}(y_i, \hat{y}_i)_{\text{Train}}$

Test Error:  $\text{diff}(y_i, \hat{y}_i)_{\text{Test}}$

If Train Error  $\uparrow$  then Bias  $\uparrow$ .

## ② High Variance OR Overfit

If Train Error  $\downarrow$  and Test Error  $\uparrow$   
 $\rightarrow$  Overfit  
 $\rightarrow$  Variance  $\uparrow$ .

If Train data changes slightly  
and there is high model change.  
 $\rightarrow$  High Variance

## 31.21 Best and Worst case of Algo

For any ML algo, it is imp to understand Best and Worst Case.

### ① Regarding K-NN

① when dimensionality is small ( $< 10$ )  
 $\Rightarrow$  K-NN is a good algo.

② when dimensionality is large  
 $\Rightarrow$  We get into curse of dimensionality  
 $\Rightarrow$  Interpretability reduces  
 $\Rightarrow$  Runtime complexity of LSH/K-D tree  $\uparrow$

### ② Low Latency System:

$\Rightarrow$  System where we want fast search results

e.g.: Google Search, Amazon products

K-NN should not be used in Low Latency Systems

③ find right distance measure

→ then K-NN is good. (e.g.: cosine similarity)

For text, cosine dist works well.

For Genome data, Euclidean dist works

If someone gives Similarity Matrix or dist Matrix → K-NN

We have to be careful while using K-NN.

## 32.1 NAIVE BAYES ALGO

Simplistic or unsophisticated algo using Bayes' Theorem

Imagine I have a datapoint  $x$  with  $n$  features and  $k$  possible outcomes or classes.

$$x = \langle x_1, x_2, x_3, \dots, x_n \rangle \Rightarrow d = n$$

$$y_i \in \{c_1, c_2, \dots, c_k\}$$

$P(c_k|x)$  = Probability of a class label given <sup>data</sup> point  $x$ . (variable)

Using simple Bayes Theorem.

$$\frac{P(c_k|x)}{\text{Posterior}} = \frac{P(c_k) \cdot P(x|c_k)}{P(x) \cdot \text{evidence}}$$

We know that

$$P(c_k) \cdot P(x|c_k) = P(x \cap c_k)$$

Now lets say I calculate the probability for each of the classes

$$P(c_1|x)$$

$$P(c_2|x)$$

$$P(c_3|x) \rightarrow \begin{array}{l} \text{say, its} \\ \text{Largest} \end{array}; \text{ we will assign } x \text{ to } c_3$$

$$P(c_k|x)$$

$P(C_1|x)$   
 $P(C_2|x)$   
 $P(C_3|x)$   
 $\vdots$   
 $P(C_k|x)$ 
} All of them have same denominator, hence we need to consider the numerator.

also the denominator doesn't contain the term  $C_k$  (class) which is the evaluation criteria. So, we can ignore denominator.

So we can say

$$P(C_k) \cdot P(x|C_k) = P(x \cap C_k) = P(C_k, x)$$

Another way of writing

Since denominator is same for all  
 ~~$P(C_k|x) \propto P(C_k, x)$~~

This is called Joint Probability Model.

We know  $x = \langle x_1, x_2, \dots, x_n \rangle$

$$\begin{aligned} \therefore P(C_k, x) &= P(C_k, \langle x_1, x_2, \dots, x_n \rangle) \\ &= P(\langle x_1, x_2, \dots, x_n \rangle, C_k). \end{aligned}$$

We know that  $P(A|B) = P(A \cap B) / P(B)$

$$P(A \cap B) = P(A|B) \cdot P(B)$$

... as per definition of cond. probab.

∴ Similarly

$$P(\underbrace{x_1, x_2, \dots, x_n}_{A}, c_k) = P(x_1 | x_2, x_3, \dots, x_n, c_k) * P(x_2, x_3, \dots, x_n, c_k)$$

keeping the first term as is, lets break the second term

$$= P(x_1 | x_2, x_3, \dots, x_n, c_k) \rightarrow \text{FIRST TERM}$$

$$* P(x_2 | x_3, \dots, x_n, c_k) * P(x_3, x_4, \dots, x_n, c_k)$$

SECOND Term.

So, if we keep on breaking further

$$= P(x_1 | x_2, x_3, \dots, x_n, c_k)$$

$$* P(x_2 | x_3, \dots, x_n, c_k)$$

\* ...

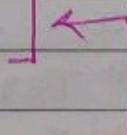
$$* P(x_{n-1} | x_n, c_k)$$

$$* P(x_n | c_k)$$

$$* P(c_k)$$

$$\rightarrow P(c_k, x_1, \dots, x_n)$$

$$P(c_k | x) \propto P(c_k, x)$$



Computing the probability calculation is very complex.

Eg :

$$P(\text{weight} = 180 | \text{hairColor} = \text{Black}, \text{hairLength} = 5\text{cm}, \text{skinColour} = \text{Brown}, \text{eyeColour} = \text{Black})$$

It is difficult to compute probability with exact same features.

So, we go with Naive Assumption of conditional independence.

$$P(A|B,C) = P(A|C)$$

A is conditionally dependent on B.  
hence, we can remove B.

OR

A is conditionally independent of B.

$$P(x_i | x_{i+1}, \dots, x_n, c_k) = P(x_i | c_k)$$

$x_i$  is conditionally independent of  
 $x_{i+1}, x_{i+2}, \dots, x_n$ .

Hence, we are making Naive assumption that features are conditionally independent of each other.

$$\therefore P(c_k | x_1, x_2 \dots x_n) \propto P(c_k) \cdot \prod_{i=1}^n P(x_i | c_k)$$

So, we arrive at a fundamental term of exact probability.

$$P(C_k | x_1, x_2, \dots, x_n) = \frac{1}{Z} P(C_k) \cdot \prod_{i=1}^n P(x_i, C_k)$$

So, we will compute

$$P(C_1|x), P(C_2|x), P(C_3|x), \dots, P(C_k|x)$$

and compute max value out of it

This is called Maximum a Posteriori Rule

These are Posterior Probability.

### 32.2 Toy Example

We have a dataset D with features like

Outlook  $\in \{\text{Sunny, Overcast}\}$

Temperature  $\in \{\text{Hot, Mild, Cold}\}$

Humidity  $\in \{\text{High, Normal}\}$

Wind  $\in \{\text{Weak, Strong}\}$

Response  $\in \{\text{Yes, No}\}$

Response indicates whether its good time to play.

Its a binary classification Task with 4 categorical features.

So, we have 2 stages

Training (Learning)

Test

} Refer  
shatterline  
Blog for  
exact  
Dataset

For K-NN there is no training stage  
There is only Test/Evaluation Stage.

Building k-d tree / LSH we require  
training stage.

Training phase  $\rightarrow$  also called Learning  
phase.

We have to predict

$$P(\text{Response} = \text{Yes} | x_q) \quad \text{and} \quad P(\text{Response} = \text{No} | x_q)$$

The higher probability will be  $y_q$ .

We saw in Naive Bayes that

$$P(C | f_1, f_2, f_3, f_4) \propto P(f_1 | C) * P(f_2 | C) * P(f_3 | C) * P(f_4 | C) * P(C)$$

Referring the dataset (shatterline blog)

$$P(C = \text{Yes}) = \frac{9}{14}$$

Now lets consider the individual probabilities

$P(\text{Outlook} = \text{o} | \text{class} (\text{Yes}/\text{No}))$

<u>Outlook</u>	<u>Frequency</u>		<u>Probability in class</u>	
	<u>Yes</u>	<u>No</u>	<u>Yes</u>	<u>No</u>
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0/5
Rain	3	2	3/9	2/5
<u>total</u>	<u>total</u>			
	= 9	= 5		

Similarly count for rest of features.

In training phase, we compute all the likelihood probabilities

Time Complexity : Roughly  $O(ndc)$

num of ↓      ↓      ↓  
 data points    features    class  
 labels

Considering the small # of d and C.

Time Complexity :  $O(n)$

Space Complexity :

We need to save  $P(f_i | c)$  → 'd' space  
 and  $P(c)$  → 'c' space

∴ Space Compl :  $O(d * c)$

Overall Time Compl :  $O(nd)$  } Training  
 Space Compl :  $O(dc)$  } Phase.

lets say  $x_q = (\text{outlook} = \text{sunny},$   
 $\text{Temperature} = \text{cool},$   
 $\text{Humidity} = \text{High},$   
 $\text{Wind} = \text{Strong})$ .  
 $y_q = ?$

$$P(C = \text{Yes} | x_q) = P(\text{Sunny} | C = \text{Yes}) * P(\text{cool} | C = \text{Yes}) * P(\text{High} | C = \text{Yes}) * P(\text{Strong} | C = \text{Yes})$$

Calculated  
 in Training Phase  $\leftarrow \alpha = \frac{2/9 * 3/9 * 3/9 * 9/14}{0.0053}$

$P(C = \text{No} | x_q) \propto \frac{0.0025}{0.0205}$

} As per Aposterior Rule;  $y_q = \text{No}$

### Test Time Complexity - (Test data)

$$= O(d.c)$$

since we have to look up 'd' features for 'c' classes.

- \* Naive Bayes is much better as compared to K-NN in terms of space complexity.
- \* Naive Bayes is much more memory efficient at runtime

32.3

Naive Bayes on Text data

- Very popular technique for Text data
- First of techniques used for Email Spam filters.

Spam Filter :- email → spam  
 ↓  
 Not spam.

Review : +ve / -ve .

Task at hand :

	class
← text 1 →	1
← text 2 →	0
⋮	1
	0

Task :  $P(Y_q=1 | \text{text}_q)$  } compare which  
 $P(Y_q=0 | \text{text}_q)$  } is greater

Step we follow :

text → stop words }  
 Stemming } Preprocessing  
 n-grams }

↓  
 Binary BoW

- We just check if word is present or not
- don't count # occurrence

Now text  $\rightarrow$  represented  $\{w_1, w_2, \dots, w_d\}$

$$P(y=1 \mid \text{text}) \propto P(y=1 \mid \{w_1, w_2, \dots, w_d\})$$

$$= P(y=1)$$

$$\times P(w_1 \mid y=1)$$

$$\times P(w_2 \mid y=1)$$

\*

$$\times P(w_d \mid y=1)$$

$$\propto P(y=1) * \prod_{i=1}^d P(w_i \mid y=1)$$

$$\text{illy } P(y=0 \mid \text{text}) \propto P(y=0) * \prod_{i=1}^d P(w_i \mid y=0)$$

Now

$$P(y=1) = \frac{\# \text{ training points with } y=1}{\text{Total } \# \text{ of training points}}$$

$$P(y=0) = \frac{\# \text{ training points with } y=0}{\text{Total } \# \text{ of training points}}$$

$$P(w_i \mid y=1) = \frac{\# \text{ datapoints which contain word } w_i}{\# \text{ datapoints with } y=1 \text{ (Train)}}$$

We have divided Train data in 2 parts

$x_i$	$y_i$	Train data with $y=1$
$\text{text}_i$	1	
	0	Train data with $y=0$
	0	

Hence for Text Classification;

Spam Detection

Polarity of a review

? Naive

Naive Bayes is  
a very good  
baseline



Benchmark algo.

We usually compare the complexity of other algo with Naive Bayes.

Let's say:

Naive Bayes → 98% Accuracy

Logistics Regression →

G.B. D.T →

D.L →

Suppose.

Deep Learning gives 98.5% Accuracy  
& Naive Bayes → 98% Accuracy  
(Simplistic Algo)

So D.L. gives only small <sup>incremental</sup> accuracy over N.B.

## 32.4 Laplace Additive Smoothing

Coming to text processing

Training :- We have computed

$$P(y=1); P(y=0)$$

$$P(w_1|y=1); P(w_1|y=0)$$

$$P(w_2|y=1); P(w_2|y=0)$$

$$P(w_3|y=1); P(w_3|y=0)$$

:

$$P(w_m|y=1); P(w_m|y=0)$$

likelihoods

lets go to Test stage :

After all preprocessing,  
we get  $\text{textq} = (w_1, w_2, w_3, w')$

Lets say  $w'$  is not present in  
the set of words which we finalized  
after preprocessing

$w'$  is a completely new word and  
hence not used in training data at all  
OR if we are using bigrams or trigrams,  
then the sequence is not present  
~~in text~~ te

Now

$$P(y=1 \mid \text{text}_q) = P(y=1 \mid w_1, w_2, w_3, w')$$

$$= P(y=1) *$$

$$\{ P(w_1 \mid y=1) *$$

$$\begin{array}{l} \text{We have already } \\ \text{computed } \end{array} \left. \begin{array}{l} P(w_2 \mid y=1) * \\ P(w_3 \mid y=1) * \end{array} \right.$$

We need to  $\leftarrow P(w' \mid y=1)$   
find this

Challenge: How to get  $P(w' \mid y=1)$   
 $\& P(w' \mid y=0)$

→ One may say that why worry about one word. We can completely ignore it.

But ignoring it would imply

$$\frac{P(w' \mid y=1)}{P(w' \mid y=0)} = 1 \quad \text{which is incorrect}$$

Now

$$P(w' \mid y=1) = \frac{P(w', y=1)}{P(y=1)} \quad \xrightarrow{\text{Intersection}}$$

= # points such that  $w'$  occurs  
and  $y=1$

# points where  $y=1 \rightarrow n_1$

$$= \frac{0}{n_1} = 0 \Rightarrow \text{We cannot use this as it will make entire term as 0}$$

The workaround is

Laplace Smoothing (or Additive smoothing)  
(Not to be confused with Laplacian smoothing)

We have seen

$$P(w' | y=1) = \frac{0 + \alpha}{n_1 + \alpha K} \quad K = \# \text{ distinct values } w' \text{ can take}$$

So, here  $w'$  can be present  $\rightarrow 1$

or not present  $\rightarrow 0$

$\therefore K$  can take 2 values;  $K=2$

In a non text processing case,  
I want to take value of feature  $f_i$

$$P(f_i=a | y=1) = \frac{0 + \alpha}{n_1 + \alpha K}$$

↓

# of distinct values  
that  $f_i$  can take

$\alpha$  can be any numerical value.

Typically  $\alpha=1$ ; not always

Why is all this useful?

Let say  $n_1=100$ ; and  $K=2$  {since  $w'$  can be 1 or 0}

$$P(w' | y=1) = \frac{0 + \alpha}{100 + 2\alpha} = \frac{\alpha}{100 + 2\alpha}$$

$$\text{Case 1: } \alpha = 1; P(w' | y=1) = \frac{1}{102} \neq 0$$

Now, a question might arise is why to go through all such complex calculation.

Why not replace  $P(w' | y=1)$  with  $\epsilon = 0.001$  (very small value).

There is mathematical foundation behind this.

Case 2:  $\alpha$  is large;  $\alpha = 10000$

$$P(w' | y=1) = \frac{0 + 10K}{100 + 20K} \rightarrow \frac{10K}{20100} \approx \frac{1}{2}$$

When  $\alpha$  is large and we know that  $w' = 0 \quad \left. \begin{matrix} \\ \end{matrix} \right\} \text{Two possibilities}$   
 $w' = 1 \quad \left. \begin{matrix} \\ \end{matrix} \right\}$

$$P(w' | y=1) = P(w' \neq y=0) = \frac{1}{2}$$

We say that both possibilities are equally likely.

This is much better assumption where we say that  $w'$  is equally likely to appear in +ve class as in -ve class.

Hence, with Laplace Smoothing

$$P(w_i | y=1) = \frac{\text{# datapoints with } w_i \text{ and } y=1}{\text{# datapoints with } y=1} + \alpha$$

↓

for all words (Both irrespective of present in training data or not)

This is a smoothing operation. Also called as Negative Smoothing.

### Smoothing:

Consider we have a sparse dataset.

$$P(w_i | y=1) = \frac{2 + \alpha}{50 + K\alpha} \quad \text{where } K=2$$

$$\alpha=1 \Rightarrow P(w_i | y=1) = \frac{2+1}{50+2} = \frac{3}{52}$$

$$\alpha=10 \Rightarrow \frac{2+10}{50+20} = \frac{12}{70}$$

$$\alpha=100 \Rightarrow \frac{2+100}{50+200} = \frac{102}{250}$$

$$\alpha=1000 \Rightarrow \frac{2+1000}{50+2000} = \frac{1002}{2050} \approx \frac{1}{2}$$

As  $\alpha \uparrow$   
the  
Probab  
also  $\uparrow$



As  $\alpha \uparrow$ , we are moving the likelihood probabilities to uniform distribution

\* When  $n_1$  is small  $\Rightarrow$  denominator is small  
 $\Rightarrow$  we have less confidence in ratio

And when we have less confidence in ratio, we use  $\alpha$  and the probability tends to go to  $\frac{1}{2}$  as  $\alpha \uparrow$

↳ This is whole idea of Laplace smoothing.

It is called smoothing because we are smoothing the value towards uniform distribution.

Often times, we apply  $\alpha = 1$   
Also called 1 additive smoothing.  
This avoids 0 probability problem.

$\alpha$  is related to Bias Variance TradeOff

32.5

## Log Probabilities and Numerical Stability

We know

$$\begin{aligned}
 & P(y=1 | w_1, w_2, \dots, w_d) \\
 & = P(y=1) \\
 & * P(w_1 | y=1) \\
 & * P(w_2 | y=1) \\
 & : \\
 & * P(w_d | y=1)
 \end{aligned}
 \quad \left. \begin{array}{l} \text{all probab} \text{ } \cancel{\text{will}} \text{ lie 0-1} \\ \text{will lie 0-1} \end{array} \right\}$$

lets say  $d$  is large and all prob lie between 0 and 1.

Imagine  $d=4$ ; &  $P(y=1 | w_1, w_2, w_3, w_4)$

$$\begin{aligned}
 & = 0.2 \times 0.1 \times 0.2 \times 0.1 \\
 & = 0.0004
 \end{aligned}$$

The value is so small when  $d=4$ ;  
Imagine  $d=100$ ; the value would be very minuscule.

This results in numerical stability issues

In Python, double precision has only 16 significant values.

So, with further small values, python starts rounding.

This issue is called Numerical underflow

Reduction: Instead of using these values, we can use  $\log(\text{probability value})$ .

Why log?

$$\log(0.0004) = -2.23979$$

Instead of operating on small values, we can use log values.

Such values are called as Log Probabilities.

Also log is a monotonic function means as  $x \uparrow$ ;  $\log(x) \uparrow$ .

Hence we can compare logs of both probabilities ( $y=1$  &  $y=0$ ) to assign a class to  $y_i$ .

$$\log[P(y=1 | w_1, w_2, \dots, w_d)]$$

$$\rightarrow = \log(P(y=1)) + \sum_{i=1}^d \log[\cancel{\log} P(w_i | y=1)]$$

This is because  $\log(ab) = \log a + \log b$

Hence the product term in the equation gets converted to summation.

That's why log is used widely in math. It converts product to summation.

Similarly  $\log(a^b) = b \log a$ .

Hence log converts

$\log(ab) = \log(a) + \log(b)$   $\Rightarrow$  converts Mult  $\rightarrow$  Addition

$\log(a^b) = b \log(a)$   $\Rightarrow$  converts Exponential  $\rightarrow$  Multipl.

32.6

## Bias Variance Trade Off

We know

High Bias  $\rightarrow$  Underfitting

High Variance  $\rightarrow$  Overfitting

In Naive Bayes, there is one parameter alpha ( $\alpha$ ) which determines underfitting or overfitting

Case 1:  $\alpha = 0$  ;

$$P(w_i | y=1) = \frac{(\text{\# train points where } w_i \text{ occurs \& } y=1)}{(\text{\# train points where } y=1)}$$

Imagine total datapoints are 1000 and there is a word which occurs only twice (very rare)

$$P(w_i | y=1) = \frac{2}{1000}$$

This means we are giving a probability value (very small) to a very rarely occurring word.

- Means we are considering that outlier in calculation
- Means we are overfitting

Overfitting is related to High Variance ⇒ small change in Training data results in Dramatic Model changes.

In our training data, we had 2 texts containing the word  $w_i$ .

Now if we change the dataset slightly the probability changes dramatically

Since  $w_i$  occurs only in 2 out of 1000 cases,  
A small change in  $D_{Train}$   
i.e. Remove the 2 texts containing  $w_i$

$P(w_i | y=1)$  changes from  $\frac{2}{1000}$  to  $\frac{0}{1000}$

$\underbrace{\qquad\qquad\qquad}_{\text{Drastic change}} \text{from } \frac{2}{1000} \text{ to } 0$

Hence, when  $\alpha = 0$ ; small change in  $D_{\text{train}}$  implies large change in model  
 $\Rightarrow$  High Variance  
 $\Rightarrow$  Overfitting

Case 2:  $\alpha$  is very large; say  $\alpha = 10K$

$$P(w_i | y=1) = \frac{2 + 10K}{1000 + 2(10K)} \approx \frac{1}{2}$$

$\downarrow$   
 $K=2$

$$\therefore P(w_i | y=1) = P(w_i | y=0) \approx \frac{1}{2}$$

This happens for all words.  
So, a model cannot distinguish if word belongs to class 1 or 0.  
 $\Rightarrow$  This implies Underfitting

$\therefore$  When  $\alpha = \text{V. Large}$ :

$$P(y=1 | w_1, w_2 \dots w_d) = P(y=1) * \prod_{i=1}^d P(w_i | y=1)$$

$\approx \frac{1}{2}$

$$P(y=0 | w_1, w_2 \dots w_d) = P(y=0) * \prod_{i=1}^d P(w_i | y=0)$$

$\approx \frac{1}{2}$

~~\*\* V. Imp~~ Label is dependent on this term. Since val  $\rightarrow$  of both terms is  $\frac{1}{2}$

Comparing this scenario to K-NN  
 when  $K=n$ ; &  $n_1 \rightarrow +ve$      $\left. \begin{array}{l} n_1+n_2=n \\ n_2 \rightarrow +ve \end{array} \right\} n_1+n_2=n$

$\Rightarrow$  If  $n_1 > n_2$ ; Any  $x_q$  would be marked as +ve

$$\text{Similarly since } \prod_{i=1}^d P(w_i | y=1) = \prod_{i=1}^d P(w_i | y=0) \approx \frac{1}{2}$$

The overall probability is dependent on  $P(y=1)$  and  $P(y=0)$ .

which implies number of points with classlabel 1 and 0; whichever would be higher would have that label for any  $x_q$ ]

similar scenario

Hence when  $\alpha$  is very large

$\Rightarrow$  Underfitting

$\Rightarrow$  High Bias

Hence

Case 1:  $\alpha=0 \Rightarrow$  Overfitting  
 $\Rightarrow$  High Variance

Case 2:  $\alpha=1 \Rightarrow$  Underfitting  
 $\Rightarrow$  High Bias

So,

### & How to find Right $\alpha$

In K-NN ; we found right k using Simple cross validation OR 10-Fold  $\downarrow \downarrow \downarrow$ .

Similarly how we find

right  $\alpha$  : using simple C.V. OR 10-Fold C.V.

So,

since they control  
↑ Bias Variance Trade  
off

$\alpha$  in Naive Bayes  $\rightarrow$  Hyper Parameters

k in K-NN

$\downarrow$   
found using C.V.

## 32.7 Feature Importance & Interpretability

In K-NN we saw that we can obtain feature importance using forward selection

In Naive Bayes, we have likelihood ratio

+  $w_i, P(w_i | y=1)$  and  $P(w_i | y=0)$

We have two imp. & pieces of information

$w_i$	$P(w_i   y=1)$	$P(w_i   y=0)$
$w_1$	(Step 1).	
$w_2$	sort $w_i$ 's	
$w_3$	based on	
:	$P(w_i   y=1)$	
$w_n$	descending order	

(Step 1)

Returns top words with highest probability

→ Very informative words or features in determining a data point belongs to +ve class.

Similarly

$P(w_i | y=0) \rightarrow$  high; indicates that  $w_i$  is important word or feature in determining -ve class

Hence, feature importance is very trivial in Naive Bayes.

+ve class: find words ( $w_i$ ) with highest value of  $P(w_i | y=1)$

-ve class:  $\sim P(w_i | y=0)$

It can be determined directly from model

This is unlike K-NN where we used forward feature selection. Here, in Naive Bayes, we get it directly from Model.

What about Interpretability:

Given a point

$$x_q \rightarrow y_q = 1$$

↓

$$\{w_1, w_2, w_3, \dots, w_{10}\}$$

I am concluding  $y_q = 1$  because  $x_q$  contains words (say  $w_3, w_6, w_{10}$ ) which has high value of

$$P(w_3 | y=1);$$

$$P(w_6 | y=1);$$

$$P(w_{10} | y=1).$$

Similarly, I am also not concluding  $y_q = 0$ ; because  $x_q$  contains words (say  $w_8, w_{15}, w_{20}$ ) where

$$P(w_8 | y=0)$$

$$P(w_{15} | y=0)$$

$$P(w_{20} | y=0)$$

} is very small

Lets take one example:

In Amazon reviews; if we have a word named "phenomenal" or "great". the probability that review is +ve is very high.

Similarly ; words like "terrible", "not good", occurs with  $y_2 = 0$  very often.  
Hence, we declare the review to be negative.

32.8Imbalanced Data

What happens to Naive Bayes in terms of Imbalanced Data ?

$$P(y=1 | w_1, w_2, w_3 \dots w_d)$$

•

$$P(y=1 | w_1, w_2 \dots w_d) = P(y=1) * \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_1, w_2 \dots w_d) = P(y=0) * \prod_{i=1}^d P(w_i | y=0)$$

Let's say  $n \xrightarrow{\quad} n_1$  (+ve)  
 $\xrightarrow{\quad} n_2$  (-ve)

and assume  $n_1 \gg n_2$

90% of points in training data  $\rightarrow +ve$ .

$$\therefore \frac{n_1}{n} = 0.9 \text{ and } \frac{n_2}{n} = 0.1$$

$$\therefore P(y=1) = 0.9 \text{ and } P(y=0) = 0.1$$

Even if we assume that Product ( $\prod$ ) terms have same value.

$$\text{i.e. } \prod_{i=1}^d p(w_i | y=1) = \prod_{i=1}^d p(w_i | y=0)$$

then the probabilities solely depend on  $p(y=1)$  and  $p(y=0)$

The dominating class (which, in our case is  $p(y=1)$ ) gets the advantage which is incorrect

$p(y=1)$  is 9 times more than  $p(y=0)$

Solution :

① Standard technique of upsampling and downsampling

$$\text{i.e. } n_1 \approx n_2 ; \text{ Hence } p(y=1) = p(y=0) = \frac{1}{2}$$

② Drop  $p(y=1)$  and  $p(y=0)$

With upsampling and downsampling, we replace  $p(y=1)$  and  $p(y=0)$  with  $(\frac{1}{2})$

At the end we compare  $p(y=1 | w_1, \dots, w_d)$  with  $p(y=0 | w_1, \dots, w_d)$

So, replacing  $p(y=1) = p(y=0)$  with  $y_2$  doesn't matter. It is equivalent to dropping them.

③ Some scientist and researchers have modified Naive Bayes to account for class imbalance.  
 $\Rightarrow$  Not often used

### Other Problem with Imbalanced Data:

$$n_1 = 900 \text{ (+ve)} \rightarrow P(w_i | y=1)$$

Majority class =  $\frac{\text{Numerator}}{900}$  → Numerator Ranges from 0 to 900

$$n_2 = 100 \text{ (-ve)} \rightarrow P(w_i | y=0)$$

Minority class =  $\frac{\text{Numerator}}{100}$  → Numerator Ranges from 0 to 100  
 $\rightarrow$  Small value of Numerator

$$\underline{P(w_i | y=1)} \quad \text{Vs} \quad \underline{P(w_i | y=0)}$$

Large value in Numerator      Small value in Numerator

When we do Laplace Smoothing, it impacts Ratio with small Numerator more than large Numerator

Laplace smoothing impacts Minority Ratio more

Eq: Minority Class | Majority Class  
 (c-ve) | (c+ve)

$$P(w_i | y_i=1) = \frac{2}{100} = 2\% \quad \leftarrow \text{same val} \rightarrow P(w_i | y_i=0) = \frac{18}{900} = 2\%$$

Assume  $\alpha = 10$

$$\begin{aligned} &= \frac{2+10}{100+20} = \frac{12}{120} \\ &\approx 10\% \end{aligned} \quad \begin{aligned} &= \frac{18+10}{900+20} = \frac{28}{920} \\ &\approx 3.04\% \end{aligned}$$

Using same  $\alpha$ ;  
 $2\% \rightarrow 10\%$   
 because Num &  
 Denominator are  
 same

Using same  $\alpha$   
 $2\% \rightarrow 3\%$   
 because Numerator  
 & denominator  
 are large

Solution to this :

① Using std process of upsampling  
 & downsampling

② There are some hacks: I will  
 have one  $\alpha$  for +ve class &  
 other  $\alpha$  for -ve class.

32.9Outliers.How they are handled in Naive BayesWhat are Outliers

In Text Classification

$$\mathbf{x}_q = w_1 \ w_2 \ w_3 \ \dots \ w^T$$

and

$$w^T \notin \{w_1, w_2, w_3, \dots, w_m\}$$

these are set of words encountered in D<sub>Train</sub>

Laplace smoothing takes care of this.

$$P(w^T | y=1) = \frac{0 + \alpha}{n_i + 2\alpha}$$

Hence, if we get outlier at Test time, then Laplace smoothing takes care of the same

What happens to Outlier at Train data

$$\{w_1, w_2, w_3, \dots, w_m\} \leftarrow \text{set of words in } D_{Train}$$

Suppose there is a word

$w_8 \rightarrow$  Occurs few times in +ve and -ve class.

$\rightarrow$  outlier

## Hack/solutions for outliers in DTrain

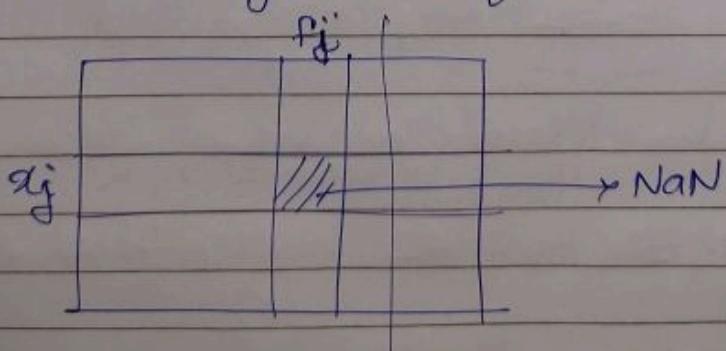
- ① If a word ( $w_j$ ) occurs ~~more~~<sup>fewer</sup> than  $\frac{10}{\downarrow}$  times, then just ignore that word (any number)
- ② and Remove that word from training set.  
 $w_j \notin \{w_1, w_2, \dots, w_m\}$   
 $w_j$  removed.
- ③ Using Laplace Smoothing with reasonable value of  $\alpha$ .

## 32.10 Missing Values

Case 1: Text Data:

There is no case of missing data as the word either occurs or doesn't.

Case 2: Categorical features



Now  $f_i \in \{a_1, a_2, a_3\} \rightarrow$  all <sup>possible</sup> values of  $f_i$  in Rest of dataset

We assume NaN as one of the categorical values of  $f_i$   
⇒ Imp. Hack.

$$\therefore f_i \in \{a_1, a_2, a_3, \underline{\text{NaN}}\}$$



### Case 3 : Numerical features

We used standard feature Imputations like Mean, Median, etc.

One such method is Gaussian Naïve Bayes. We would discuss it shortly.

## 32.14 Naive Bayes with Numerical Features

NB: Binary features } Very straight  
Categorical  $\rightarrow$  forward.

NB: Numerical / Real valued.

Imagine I have dataset

$f_1$	$f_2$	$f_3$	$f_j \dots f_d$	$y$	
$x_i$	$x_{i1}$	$x_{i2}$	$x_{i3}$	$x_{ij}$	$x_{id}$

Divide  
it in  
two class

$f_j \rightarrow$  real valued.

## Fitting NB in Real Valued scenario

$$P(y_i = 1 | x_{i1}, x_{i2}, \dots, x_{id})$$

Equal  
Remains  
Same

$$\propto [P(y_i = 1)] * \underbrace{\prod_{j=1}^d P(x_{ij} | y_i = 1)}$$

↓  
Class Prior  
Probab.  
 $= \frac{n_1}{n_1 + n_2}$

Probab of j<sup>th</sup>  
feature given  
class label

Remember, we have to find

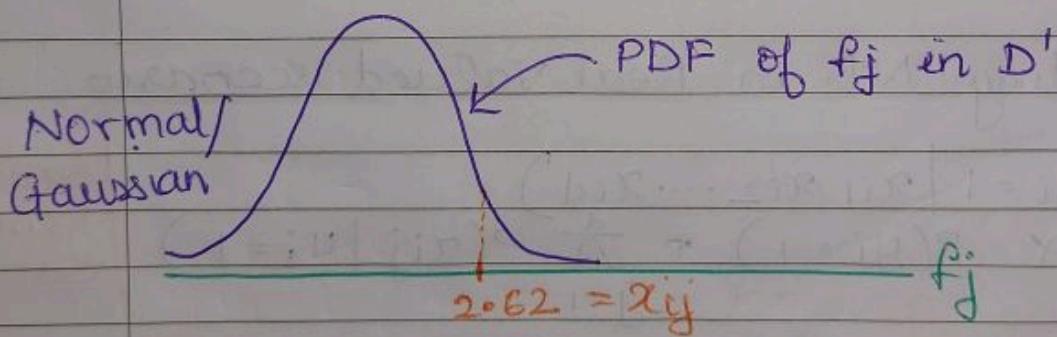
$$P(x_{ij} | y=1)$$

Since  $y=1$ ; Lets concentrate on  
+ve data in DTrain

$$\text{Let } D_{\text{Train}}^{+\text{ve}} = D'$$

$$\begin{aligned} P(x_{ij} | y=1) &= P(x_{ij} | D') \\ &= \text{Probability of } x_{ij} \\ &\text{occurring in } D' \end{aligned}$$

Imagine if I plot PDF of  
+ve Labeled data of  $f_j$



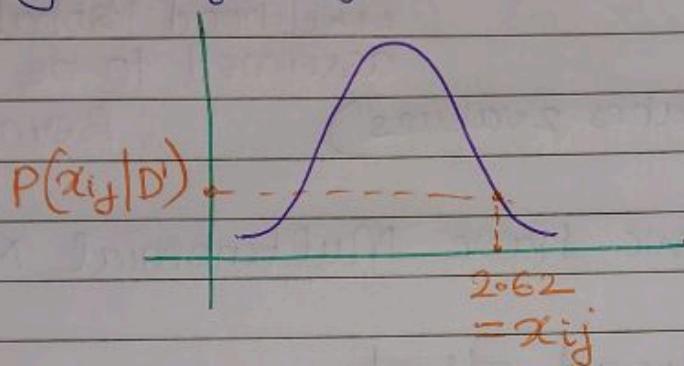
Lets assume, my curve is Gaussian

Lets say  $x_{ij} = 2.62$  which lies in PDF

We can obtain  $P(x_{ij} | D')$  from  
the PDF of  $f_j$  in  $D'$

Let's assume the PDF follows Gaussian dist with  $N(-1, 0.6)$

then I can get  $P(x_{ij} | D')$  by the height of  $x_{ij}$  in PDF



Similarly we can concentrate on dataset with  $y=0$  and get  $P(y=0 | P(x_{ij} | D''))$

Here, people assume  $f_j$  in  $D'$  follows Gaussian distribution with  $N(\mu_j, \sigma_j)$

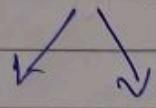
Also, they assume  $f_j$  in  $D''$  (-ve dataset) follows Gaussian dist with  $N(\mu_j^*, \sigma_j^*)$

When we make this assumption, such a model is called Gaussian Naive Bayes.

Likelihood Probabilities are Assumed to be Gaussian

Similarly, for words, we assume  
(text data)

$$P(w_i | y=1)$$



$\Rightarrow$  Bernoulli N.B.

(each word takes 2 values)

Likelihood Probab are  
assumed to be Binary/  
Bernoulli

Similarly we have Multinomial N.B.

Let's not forget that

Naïve Bayes makes fundamental  
assumption of Conditional Independence  
 $\Rightarrow$  All features are independent of  
each other

32.15

Can Naïve Bayes do Multi-class classification?

$\Rightarrow$  Yes

Just like we do

$$\left. \begin{array}{l} P(y_i = 1 | w_1, w_2, \dots, w_d) \\ P(y_i = 0 | w_1, w_2, \dots, w_d) \end{array} \right\}$$

We can also perform

$$\left. \begin{array}{l} P(y_i = 2 | w_1, w_2, \dots, w_d) \\ P(y_i = c-1 | w_1, w_2, \dots, w_d) \end{array} \right\}$$

Compare all  
and pick  
largest as  
class  
label

Naïve Bayes is designed for multiclass  
classification  $\Rightarrow$  It can easily do it

32-13

Similarity of Dist Matrix

Naive Bayes cannot handle similarity Matrix if given instead of feature values

Because we have to compute prob.  
and we have product term.  
and we need exact feature value.

$$P(y_i=1 | f_1, f_2, \dots) = P(y_i=1) \prod_{i=1}^d P(f_i | y_i=1)$$

↓  
exact feature  
value

& Naive Bayes doesn't used Distance

→ → is a probabilistic method  
where we need actual feature values

32-14

Large dimensionality

Naive Bayes used extensively in Text classification  
Hence, it can handle large dimensions

Only, we need to use log Probabilities

As dimensionality ↑, the Product term increases ; hence use log Probabilities

$$P(y=1 | w_1, w_2, \dots, w_d) = P(y=1) \prod_{i=1}^d P(w_i | y=1)$$

{  
Lot of terms.