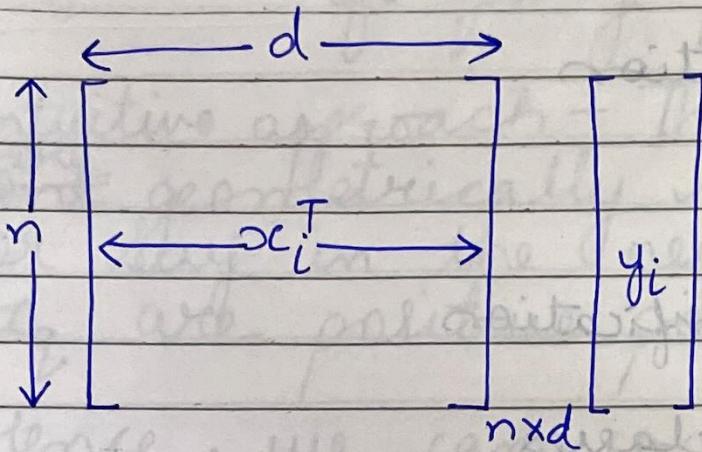


Video 2 : Notations of Data Matrix



By default, every data point is represented as column vector. So, transpose taken here : we want row vector of it. x_i is column vector to take transpose

$$D_n = \{ (x_i, y_i) \}_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\} \}$$

for binary classification

D is a dataset of pair of values, n -such pairs, where the first pair is a d -dimensional vector and second as either 0 or 1.

Video 3: Classification vs Regression

1) Binary classification

$$y_i \in \{0, 1\}$$

2) Multiclass classification

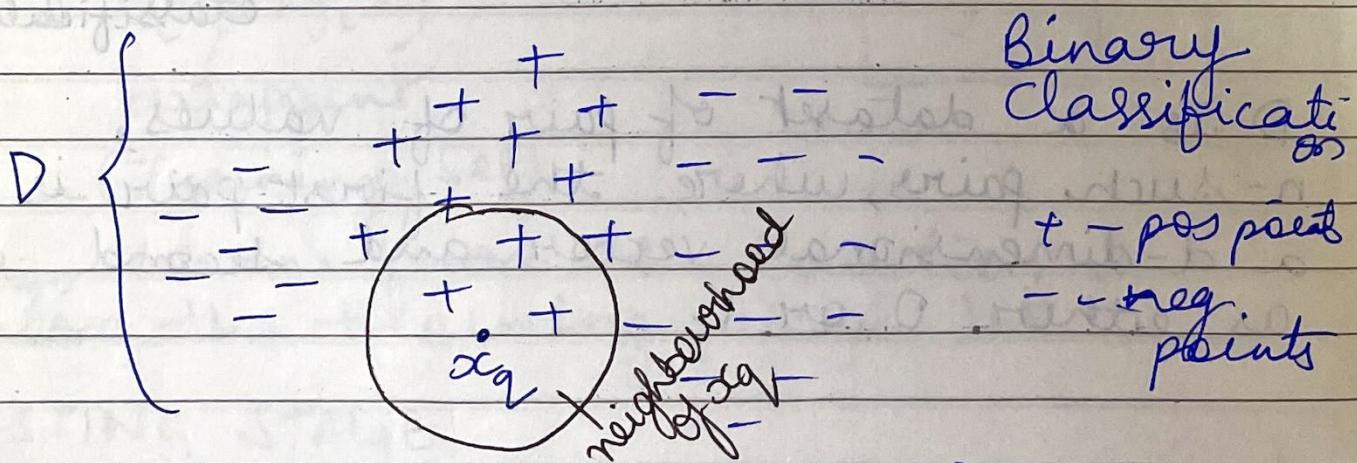
More than 2 classes.

3) Regression

$$y_i \in \mathbb{R} \text{ (real value/number)}$$

Video 4: K-Nearest Neighbours Geometric Intuition

2D Toy dataset



$$D = \{(x_i, y_i) \mid x_i \in \mathbb{R}^2, y_i \in \{0, 1\}\}$$

\mathbb{R} in \mathbb{R}^2 denotes 2D space.

Objective - To determine class label of query point.

Intuitive approach - The points that are geometrically close to x_q or say in the neighbourhood of x_q are positive points mostly.

Hence, we can conclude that x_q is also positive point.

KNN Algo.:-

1) Find K nearest points to x_q in D
Say $K=3$ & x_1, x_2 and x_3 are closest to x_q .

$$x_1 \rightarrow y_1$$

$$x_2 \rightarrow y_2$$

$$x_3 \rightarrow y_3$$

2) Decide class label of x_q using majority vote of y_i 's

$$\{y_1, y_2, y_3\} \quad y_q$$

$$+ + +$$

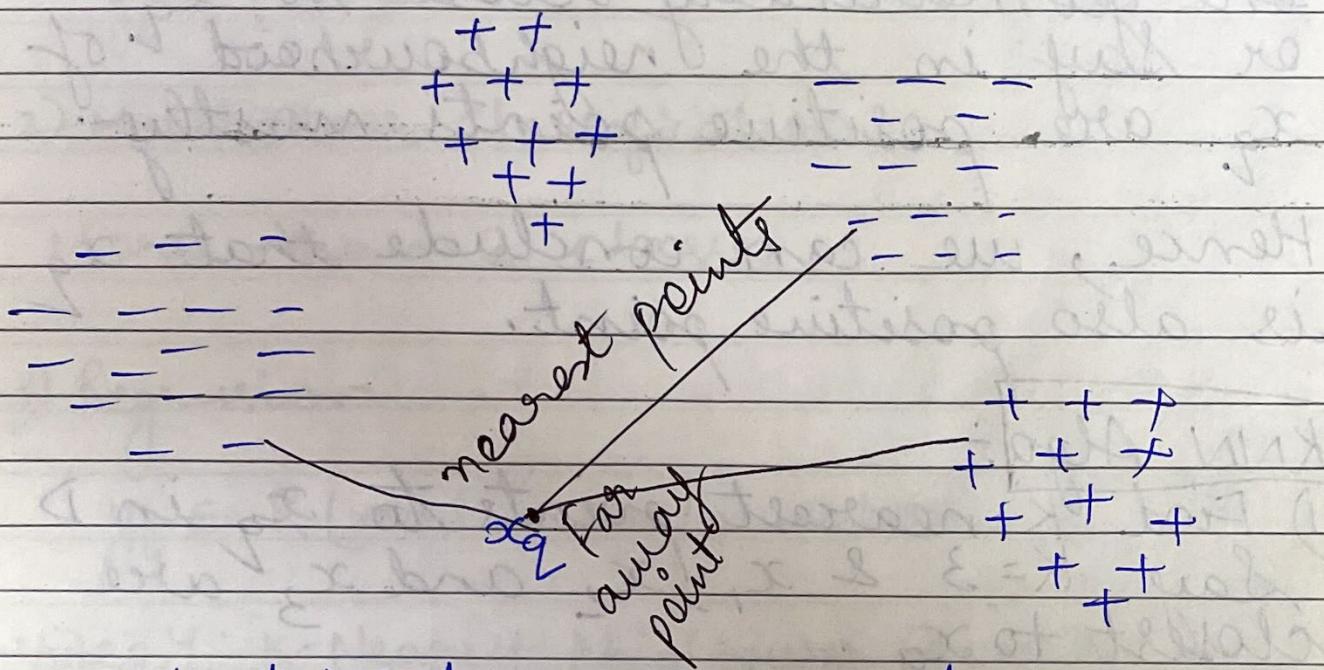
$$+ + -$$

$$- + -$$

Note - K should be an odd no.

Videos 5: Failure cases of KNN

- 1) query point is far away from clusters.



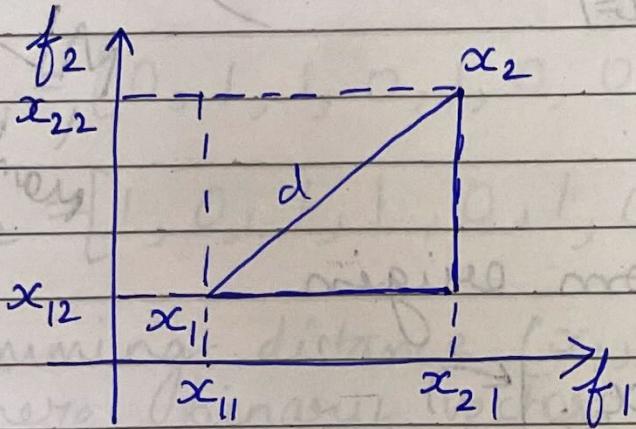
Output label may say point x_q as -ve
of majority vote but it doesn't make
much sense. It would be better if the
model could output as - Don't know

- 2) +ve and -ve points are randomly spread.

+ - + - + - + - +	No useful information
- + - + - + - + -	
+ - + - + - + - +	
- + - + - + - + - +	
+ - + - + - + - +	

Grouped data is useful but random noisy
data is not useful.

V6: Distance Measures: Euclidean (L^2), Manhattan (L^1), Minkowski, Hamming



d = length of shortest line from x_1 to x_2

$$x_1 = (x_{11}, x_{12})$$

$$x_2 = (x_{21}, x_{22})$$

$$d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2}$$

For d -dimensions:-

1) Euclidean distance,

$$\text{If } x_1, x_2 \in \mathbb{R}^d$$

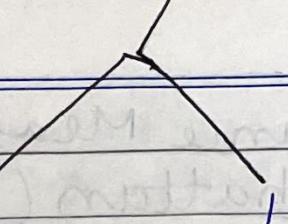
$$\|x_1 - x_2\|_2 = \left(\sum_{i=1}^d (x_{1i} - x_{2i})^2 \right)^{\frac{1}{2}}$$

L^2 norm of $x_1 - x_2$ vector

$$\text{Distance of } x \text{ from Origin} = \left(\sum_{i=1}^d x_{1i}^2 \right)^{\frac{1}{2}}$$

2) Manhattan distance

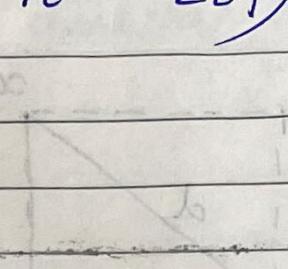
$$\|x_1 - x_2\|_1 = \left(\sum_{i=1}^d |x_{1i} - x_{2i}| \right)$$


 absolute value

L_1 norm of $x_1 - x_2$

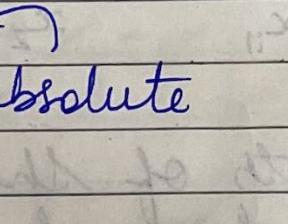
Distance of x from origin.

$$\|x_1\|_1 = \sum_{i=1}^d |x_{1i}|$$


 absolute

3) Minkowski distance

$$\|x_1 - x_2\|_p = \left(\sum_{i=1}^d |x_{1i} - x_{2i}|^p \right)^{1/p}$$


 absolute

L_p norm of $x_1 - x_2$

For $p \neq 0$ but any no.

For $p=1 \rightarrow$ manhattan

For $p=2 \rightarrow$ Euclidean

L_p norm operates on a single vector
whereas distances operate on
2 vectors.

Minkowski - L_p norm

Euclidean - L_2 "

Manhattan - L_1 "

Distances are b/w 2 pts
Norms are b/w vectors. Note

4) Hamming distance (b/w boolean vectors)

$x_1, x_2 \rightarrow$ boolean vectors \rightarrow like Binary BoW

$$x_1 = [0, 1, 1, 0, 1, 0, 0]$$

$$x_2 = [1, 0, 1, 1, 0, 1, 0]$$

Hamming distance $(x_1, x_2) \Rightarrow$ # of locations where binary vectors differ.

Also, Hamming distance could be used for measuring distance b/w strings.

V7: Cosine distance and Cosine Similarity

Intuitively, when similarity \uparrow distance \downarrow
& when similarity \downarrow distance \uparrow

So, by definition,

$$\text{cosine-distance}(x_1, x_2) = 1 - \text{cosine-similarity}(x_1, x_2)$$

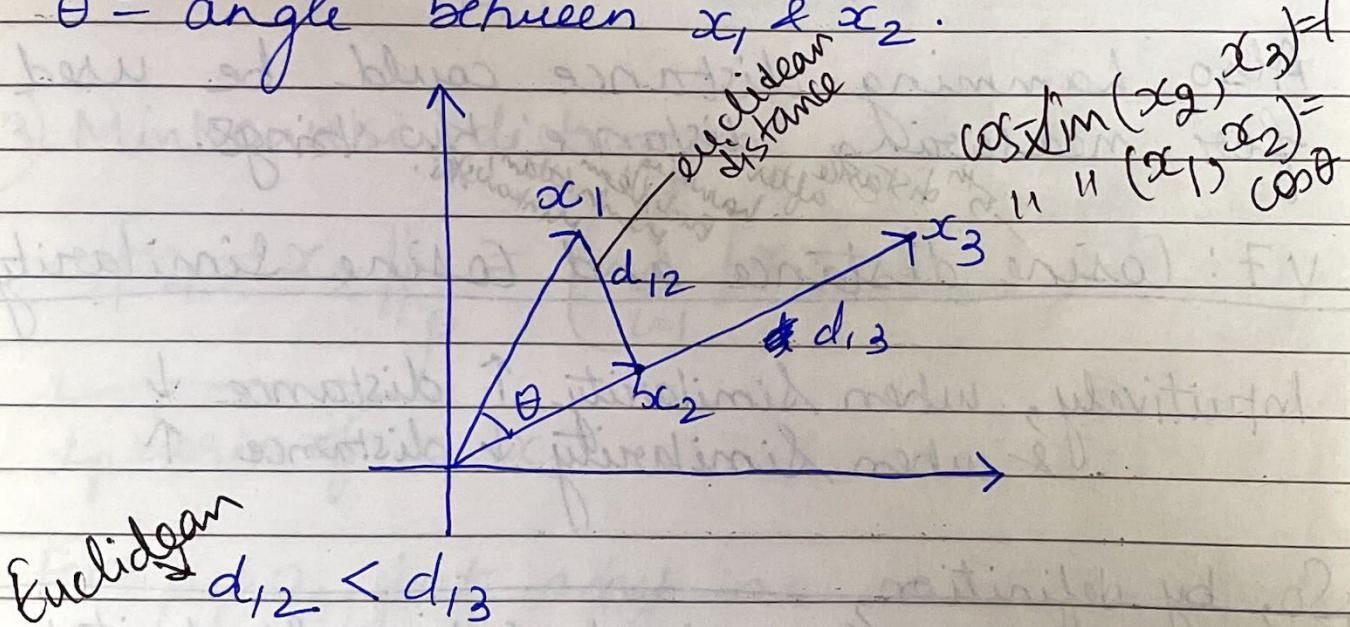
Just as in Euclidean terms, points which have less distance between them are similar. Similarly, here points that make less angle with each other are similar.

So, lesser the angle, more should be similarity.

And we know $\cos \theta \propto \frac{1}{\theta}$.

Hence, similarity $\propto \cos \theta$

So, $\text{cos-sim}(x_1, x_2) = \cos \theta$ where
 $\theta = \text{angle between } x_1 \text{ & } x_2$.



But $\cos dist_{13} < \cos dist_{12}$

So, depending upon problem specification, we may need either of the metrics as the distance measure.

If the problem deals with angular terms, then taking euclidean distance doesn't make

sense and vice-versa.

↳
norm

$$\cos \theta = \frac{\mathbf{x}_1 \cdot \mathbf{x}_2}{\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2}$$

default
value for any
norm is 2 if
not specified
if assume 2

∴ If \mathbf{x}_1 & \mathbf{x}_2 are unit vectors,

$$\cos \theta = \mathbf{x}_1 \cdot \mathbf{x}_2$$

Relationship between Euclidean and cosine distance (Refer. cosine-similarity wikipedia)

$$[\text{euc-dist}(\mathbf{x}_1, \mathbf{x}_2)]^2 = 2 \times \cos \text{dist}(\mathbf{x}_1, \mathbf{x}_2) \\ = 2 \times (1 - \cos \text{sim}(\mathbf{x}_1, \mathbf{x}_2))$$

If \mathbf{x}_1 & \mathbf{x}_2 are unit vectors

V8: Measuring Effectiveness of KNN

1) Split D_n into D_{Train} ($n_1 \approx 70\%$)

D_{Test} ($n_2 \approx 30\%$)

$$D_{\text{Train}} \cap D_{\text{Test}} = \emptyset - \text{null set}$$

$$D_{\text{Train}} \cup D_{\text{Test}} = D_n$$

2) Train the KNN model using D_{Train} and measure how well it does using D_{Test} .

Algo:

count = 0

for each point in D_{Test} :

x_q = point

use D_{Train} & KNN to find y_q

if $y_q = y_{\text{point}}$:

Count = count + 1

$$\text{Accuracy} = \frac{\text{count}}{n_2}$$

where,

count = # pts for which D_{Train} + KNN gave correct class label.

n_2 = # pts. in total D_{Test} .

Accuracy is one of the measures. There are many other measures.

V9: Test Time and Space Complexity of KNN

We need to evaluate:

Given x_q determine y_q .

Input: D_{Train} , K , $x_q \in \mathbb{R}^d$

Output: y_q

pseudocode :-

```
KNN pts = [] - List
for each  $x_i$  in  $D_{Train}$ :
    - compute  $d(x_i, x_g) \rightarrow d_i$ 
    - Keep the smallest K distances
    →  $(x_i, y_i, d_i)$  to KNN pts
        tuple
```

$O(d)$

$O(K)$

when K = small

say 1 it becomes $O(1)$

$K \ll n$. Hence neglected.

So, $O(n * d)$ for above algo.

```
count_pos = 0; count_neg = 0
for each  $x_i$  in KNN pts:
    if  $y_i$  is +ve:
        countpos += 1
    else
        count_neg += 1
```

$O(K)$
 $K =$
small

```
if count_pos > count_neg
    return  $y_g = 1$  ( $y_g \rightarrow +ve$ )
else
    return  $y_g = 0$ . ( $y_g \rightarrow -ve$ )
```

So, total $O(k)$

Overall time complexity

$$= O(nd) + O(1) + O(1)$$

$$= \boxed{O(nd)}$$

If $d \ll n$,

$$\boxed{O(n)}$$

Time taken to compute label of a single query pt.

Space Complexity:

We need to calculate distance for each point in D_{Train}

So, the entire D_{Train} should be present in RAM.

$$\text{Space} = \boxed{O(nd)}$$

V10: Limitation of KNN

We need to have a real time system with very low latency.

1) With Simple KNN

$$\text{Time} = O(nd)$$

$$\text{Space} = O(nd)$$

If $n = 364K$ — D_{Train}
 $d = 100K$

$nd = 36 \text{ GB RAM}$, means
min 64GB RAM needs to be installed.

2) Time Complexity 36 Billion computations

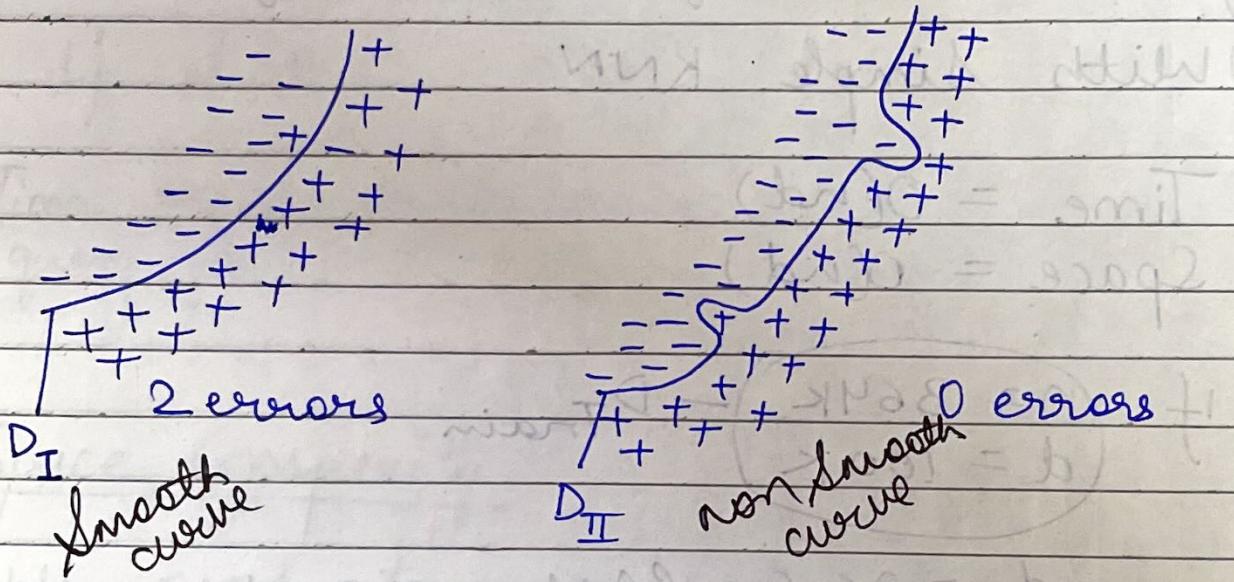
In Internet applications low latency requirement is there.

For trading and finance as much as 1μs delay is permitted only.

Non low latency systems → Medicine
This is one of the main reasons for not implementing KNN in productions and instead of it being so simple algo

VII: Decision Surface for KNN as K changes

K is a hyperparameter



D_I and D_{II} are the decision surfaces. We call them decision surfaces because of generalization of higher dimensions.

The surface decides whether a point is positive or negative depending upon which side of the decision surface the point lies.

But how to algorithmically come up with a decision surface given datapoints?

For KNN decision surface could be found by following:

- Take each point in the space of the datapoints surrounding them.
- Find its k nearest neighbours.
- Find the majority vote and assign that label to that point.
- Now draw a surface that separates both kind of points.

Let $k=1$

For the given pattern of datapoints distribution as above, if we take a consideration that a point is close to the outlier points, since we are considering 1 closest neighbour, the label of that point will also be same as of outlier point. Thus giving a skewed surface as D_{II} .

But if $k=5$, even if the sample point is close to outlier and due to majority voting of remaining 4 nearest neighbours that point will be labelled as opposite of outlier giving surface as smooth as D_I .

So, the smoothness of the decision surface increases as k increases.

If say $k = n = 1000$

$$n_1 + n_2 = n$$

n_1 = no. of +ve points 600

n_2 = " " -ve " 400

$$\text{and } n_1 > n_2$$

Then irrespective of whatever point we take, its label would be positive \therefore of imbalanced dataset.

When K becomes n x_k becomes majority class's label.

V12: Overfitting & Underfitting

I) Overworking/overfitting System

$$\text{If } k=1$$

This model overworks to get more complicated decision surface and has to ensure no errors.

It may take noise as data & classify new points incorrectly.
Overfit

2) Well fit / Robust System

If $k = 5$

It tries to find a balance b/w both extremes.

It is less prone to noise. Makes more sensible decisions.

Appropriately fit function

3) Lazy / Underfitting System

If $k = n$

It checks which class has more data points and assigns that label to the new query point.

Underfit

v13: Need for Cross Validation

$k = 1, 2, 3, \dots, n$

1

overfit

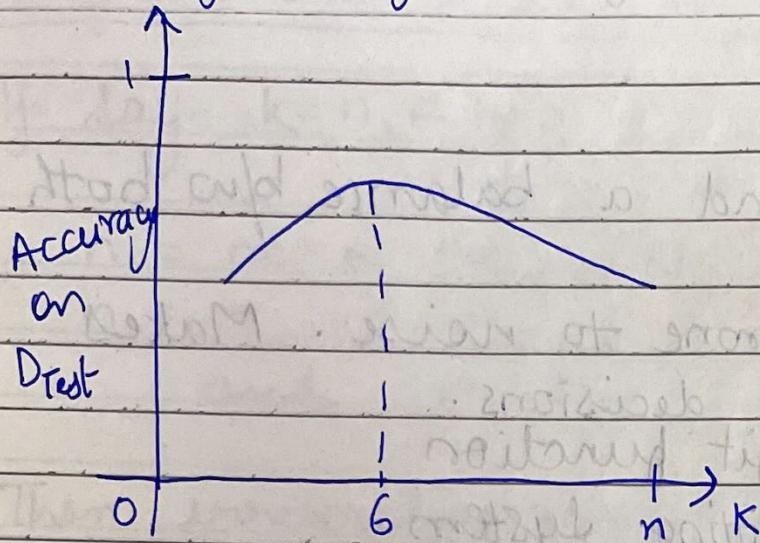
underfit

So the optimal value would be somewhere in the middle.

In order to find the right k , we must check the accuracy on D_{Test} and draw graph.
eg.

	Train data	Accuracy on D_{Test}
$K=1$	D_{Train}	0.78
$K=2$	"	0.82
:		

So, a typical graph will look as follows:-



So we conclude that $K=6$ gives the best accuracy on D_{Test} when using D_{Train} as training data.

Using D_{Train} and 6-NN on Amazon Fine Food Reviews accuracy = 96%.

But there is a small problem.
The sole objective of Machine learning is to perform well prediction on unseen data.
We need used the whole D_n here to calculate the model function.

D_{Train} gives the Nearest Neighbours.

D_{Test} gives appropriate value of K .

If a model works on future unseen data it has generalization

Find

$D_{Train} = NN$

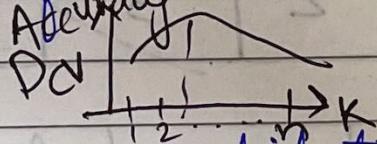
$D_{CV} = \text{Right K value}$

$D_{Test} = \text{Accuracy after applying the function}$

Cross Validation

Divide D_n into 3 parts (D_{Train} to find nearest neighbour) (D_{CV} to be used to determine right K value) (D_{Test} to apply function to get accuracy on unseen data.)

Now, if D_{Test} gives 93% accuracy on the optimal value of K found on D_{Cross} validation, we can say the accuracy of model is 93% on unseen data.



generalization accuracy
generalization error

V14: K-fold Cross Validation

Let's suppose we divide our D_n dataset into 3 parts D_{Train} , D_{CV} and D_{Test} into 60%, 20% and 20% division.

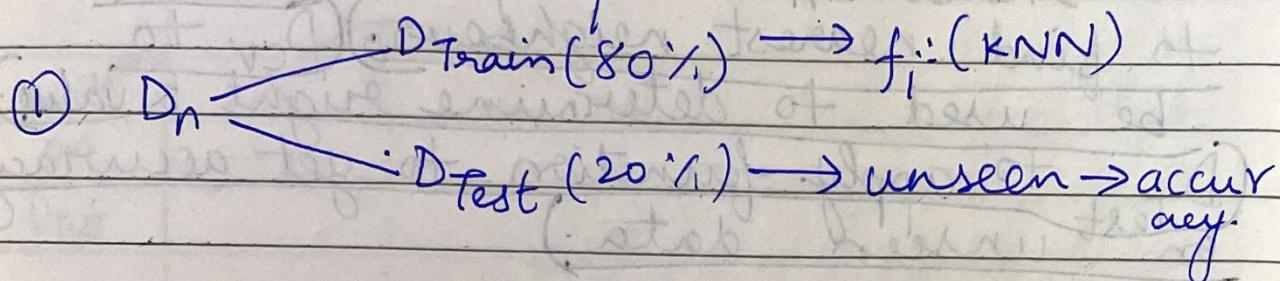
D_{Test} should be unseen data and accuracy must be calculated on that.

But the main problem is that we get only 60% of our data as train data.

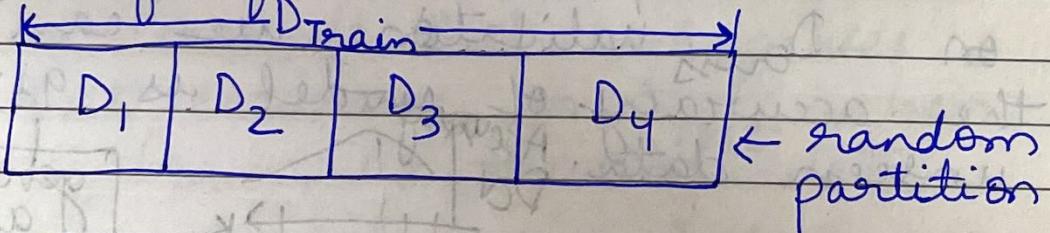
So there must be a solution as more the training data, better is the algorithm. So, we should somehow try to get whole 80% ~~D_{Test}~~ D_{Train} + D_{CV} as

training data different from test KNN

Now we use K-Fold Cross Validation to solve the above problem.



② Say $K' = 4$, so we divide D_{Train} to 4 parts of equal size:

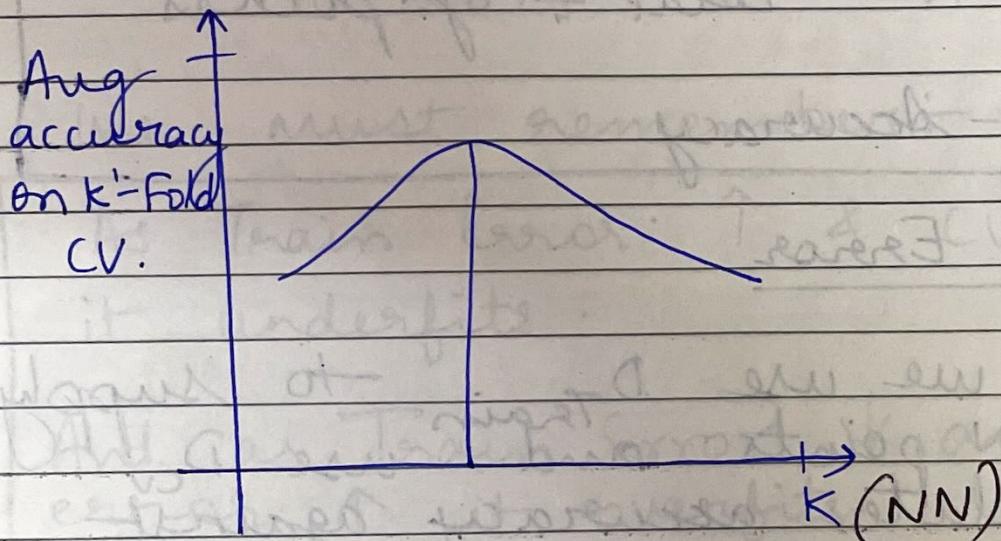


$K=1$	Train	cv	Acc. on cv
D_1, D_2, D_3	D_4	a_4	average $a_K = 1$
D_1, D_2, D_4	D_3	a_3	
D_1, D_3, D_4	D_2	a_2	
D_2, D_3, D_4	D_1	a_1	

$K=2$	Train	cv	Acc. on cv
D_1, D_2, D_3	D_4	a'_4	average $a_K = 2$
D_1, D_2, D_4	D_3	a'_3	
D_1, D_3, D_4	D_2	a'_2	
D_2, D_3, D_4	D_1	a'_1	

These average accuracy score for each k' must be plotted on Accuracy vs. k plot.

So, for arriving on conclusion of accuracy for a particular ' k ', all the train is used as training data. Also, more cross-validation is done as we get 4 D_{cv} sets.



But what should be value of k'

$k' = 10$ Rule of Thumb

10 fold cross validation

Time taken to compute optimal ' k ' using k' fold cross validation increases k' times as compared to simple cross validation set.

V16: How to determine overfitting and underfitting?

Given a value of 'K', how to determine whether it underfits or overfits? We use Graphical method for that.

$$\text{Accuracy} = \frac{\text{Total \# of points correctly classified}}{\text{Total \# of points}}$$

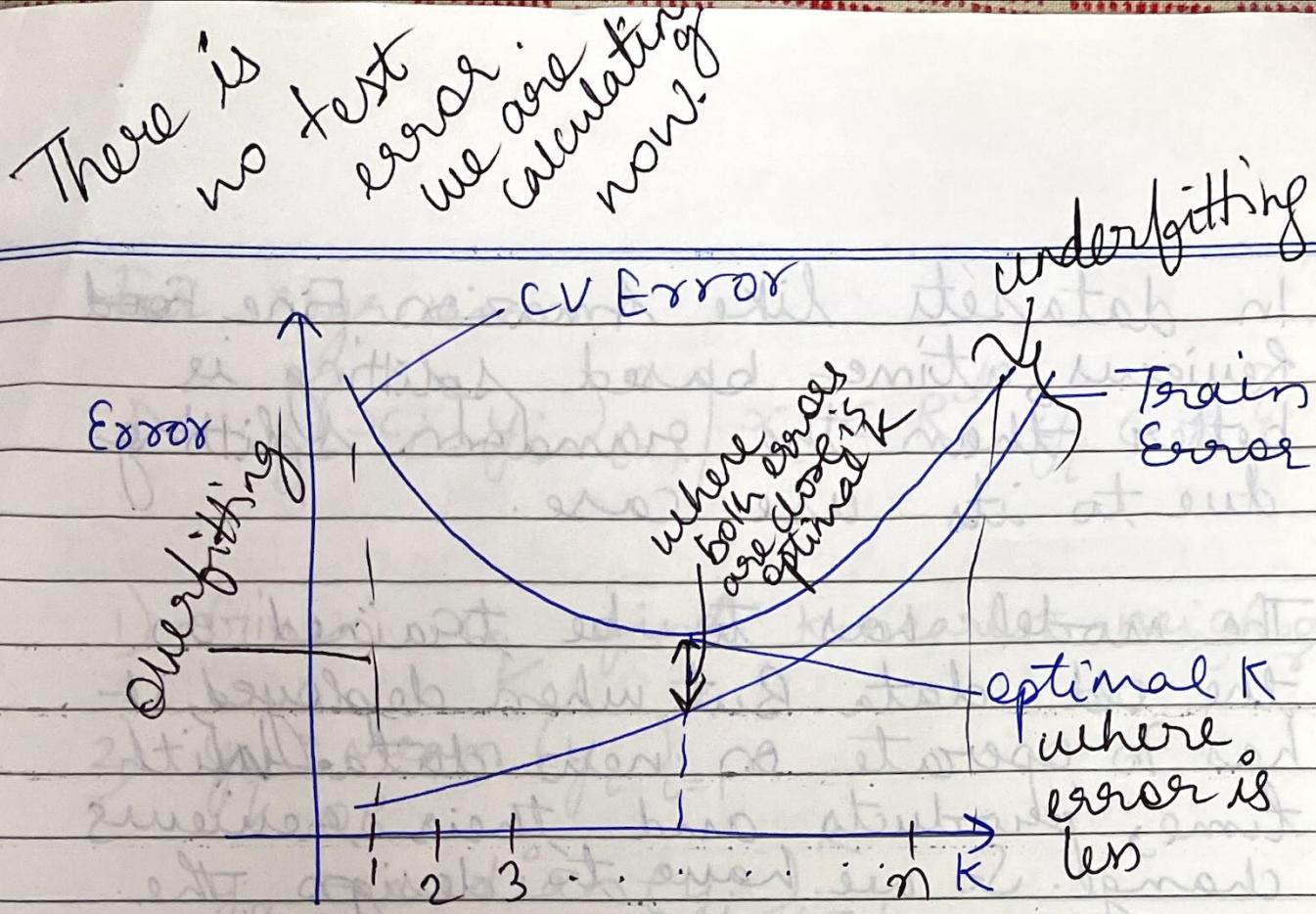
$$\text{Error} = 1 - \text{Accuracy}$$

Training Error

Usually, we use D_{Train} to supply the D_{CV} points and use D_{CV} to calculate accuracy against different values of K.

The corresponding loss of this accuracy is called Cross Validation Loss.

Instead of using D_{CV} points if we use D_{Train} points itself to calculate accuracy, the corresponding loss would be training loss.



We must remember:-

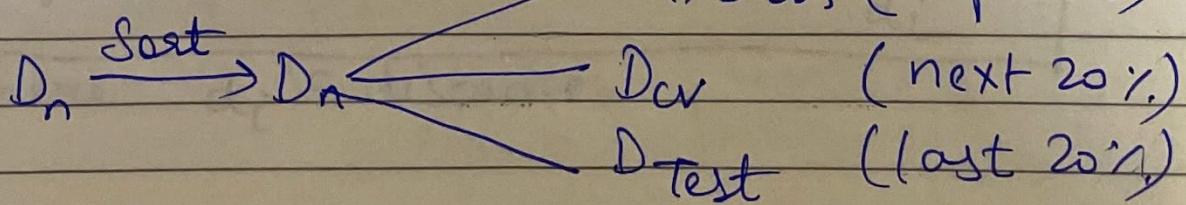
As Train Error ↑ & CV error ↑
it underfits.

And as Train Error ↓ and CV error ↑ it overfits.

V17: Time based Splitting

It works only, if we have timestamp data.

1) Sort D_n in increasing order of time.



In datasets like Amazon Fine Food Reviews, a time based splitting is better than the random splitting due to its use case.

The model has to be trained on the old data. But when deployed, it has to operate on new data. As with time, products and their reviews change. So we have to design the model in such a way that it takes more older data and gets good predictions on less older data.

This can be achieved by splitting in the above mentioned way. Merely doing random splitting won't give this ability to the model.

V18: k-NN for regression

For 2-class classification

$$D = \{(x_i, y_i)\}_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\}\}$$

For regression

$$D = \{(x_i, y_i)\}_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \mathbb{R}\}$$

1) Given x_q , find k -nearest neighbors

$$2) y_q = \text{mean}(y_i)_{i=1}^k$$

$$y_q = \text{median}(y_i)_{i=1}^k$$

The median formula is less prone to outliers.

V19: Weighted KNN

The closer points among the KNN points are given more weight than further points while

~~Now~~ ~~while~~ ~~voting~~, these weights are added for each class, and majority value wins. ~~in majority voting strategy.~~

$$\boxed{\text{One of the } w_i \text{ can be } w_i = \frac{1}{d_i}}$$

where d_i = distance.

V20: Voronoi Diagram

It follows similar procedure while trying to draw decision boundary and take $k=1$ (one nearest neighbour). It colours each region with different colours. In each cell there is only 1 nearest neighbour.

V21: Binary Search Tree

~~Suppose we have numbers:-~~

V22: How to build a Kd Tree

It is a simple implementation.

$$\begin{aligned} \text{KNN Time Complexity} &= O(n) \\ \text{Space} &= O(n) \end{aligned}$$

We can't do much about space complexity since all points made to be present for calculation.

We can improve time complexity using data structure called Kd tree.

It makes time complexity = $O(\lg n)$

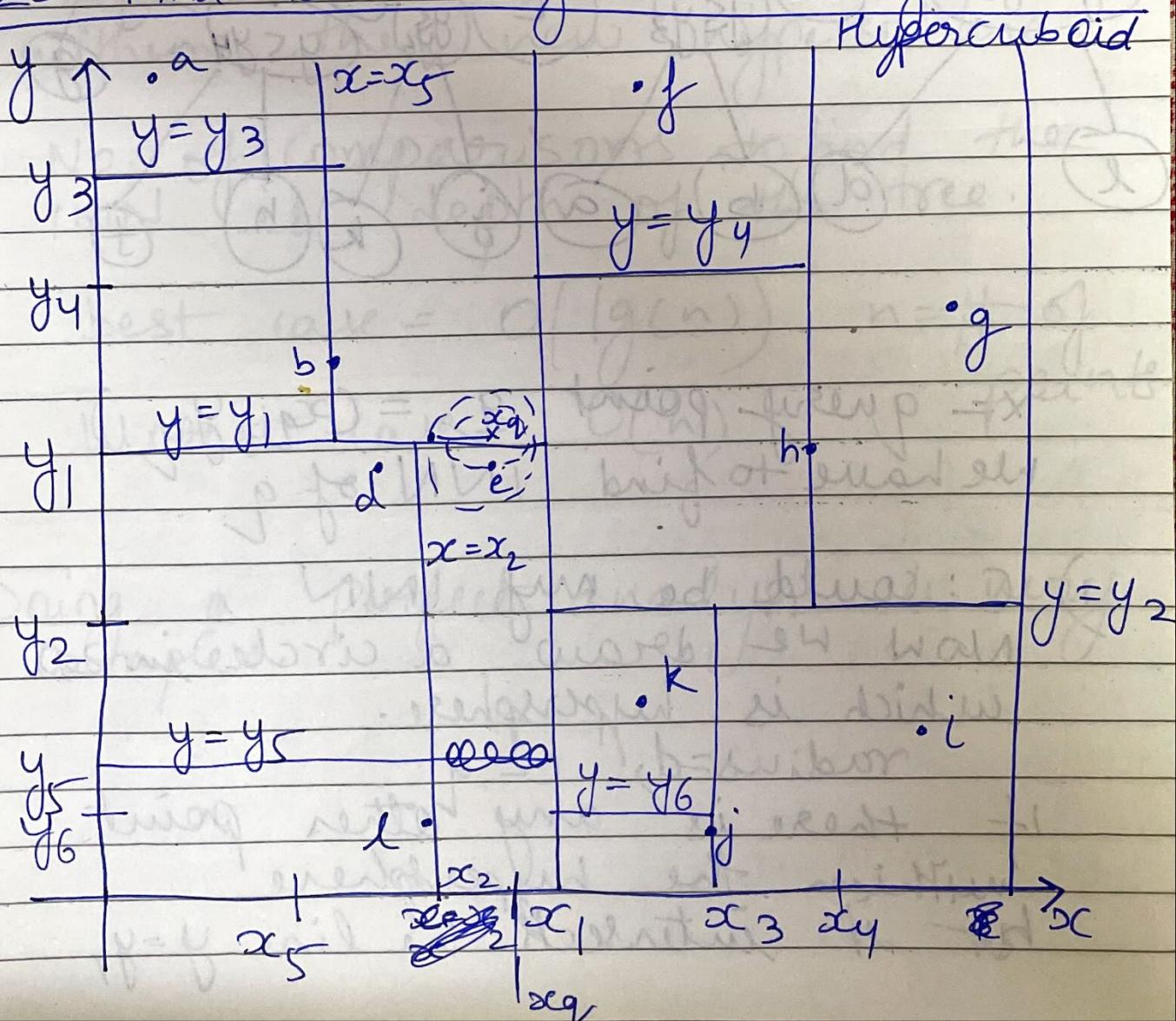
It is very popular and is used in Computer Graphics.

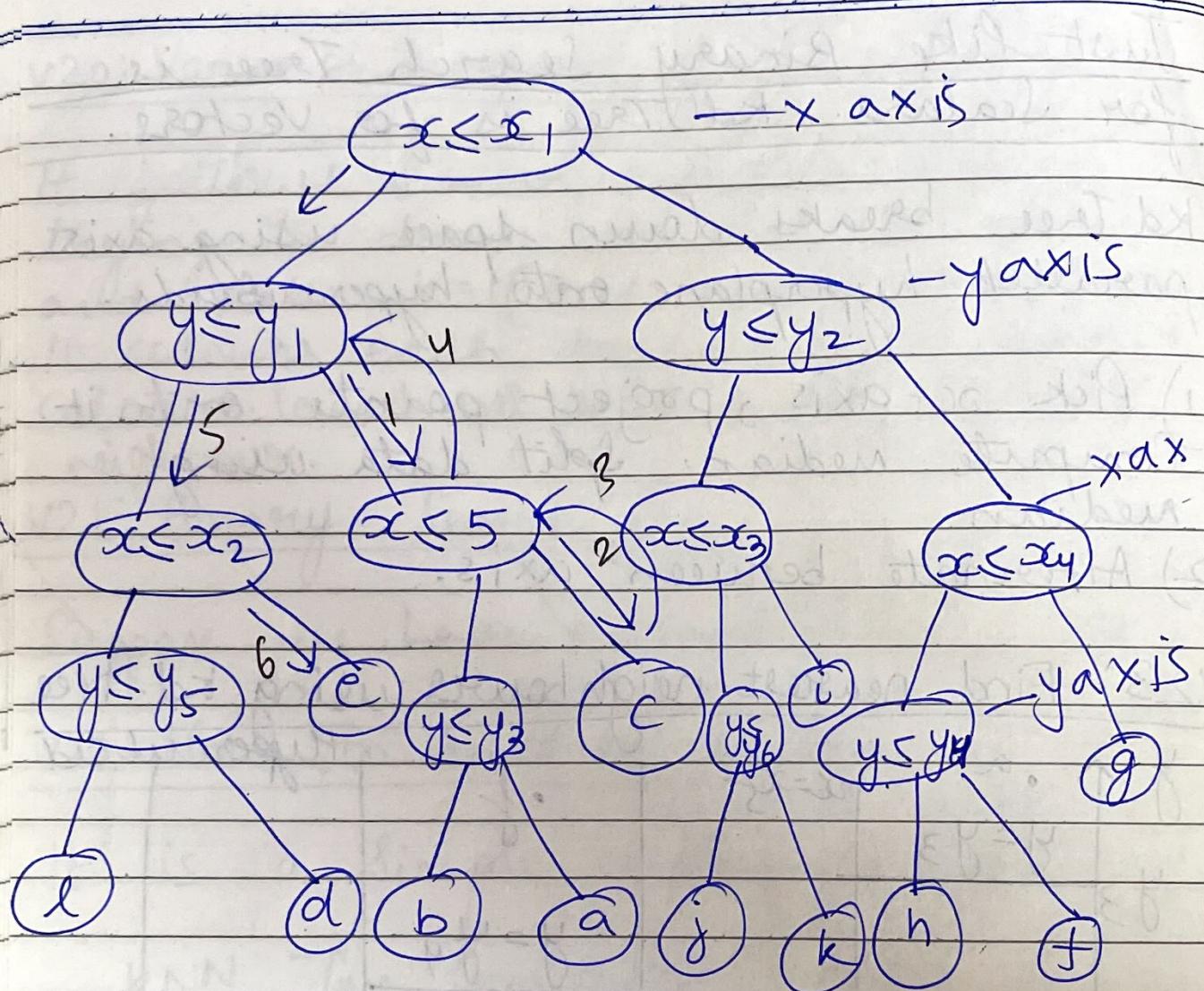
Just like Binary Search Tree is for Scalars. kdTree is for Vectors.

KdTree breaks down space using axis parallel hyperplane onto hypercuboids.

- 1) Pick x -axis, project points onto it. Compute median. Split data using median.
- 2) Alternate between axis.

v23: Find nearest neighbours using kd-tree





$x = \text{query point } q = (x_q, y_q)$
 we have to find INN of q

1) C: could be my INN
 Now we draw a circle in 2D
 which is hypersphere.
 $\text{radius} = d, C = q$

If there is any other point
 within the hypersphere
 but it intersects the line $y = y_q$,

So we track back so we reach another condition of on $y = y_1$ and we check the other condition and go down there

2) So we say e could be INN is So we discard the previous hypersphere and draw another with len radius.

And there is no other point within the new hypersphere.

No. of comparisons to find the INN is depth of the tree.

best case = $O(\lg(n))$ $n = \# \text{ of points}$

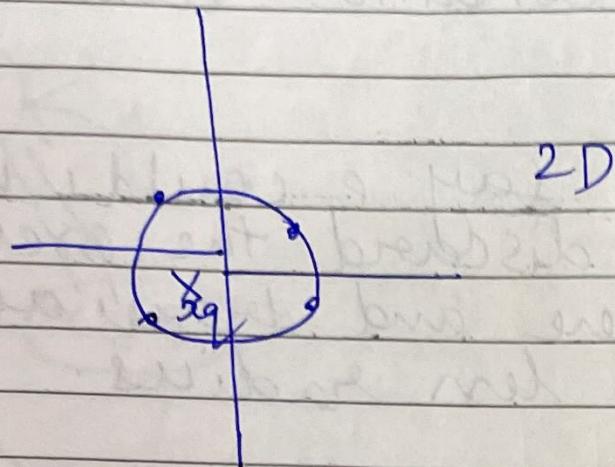
worst case = $O(n)$

Doing a perfect analysis is very complex.

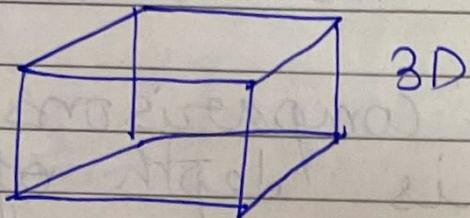
V24: Limitations of kd-tree

KRND

Q.



I have to check 4 regions.



I have to check 8 regions.

In nD I have to check $\Theta 2^d$
Disadvantages:-

So, (1) when d is not small.
 $d = 10; 2^d = 1024.$

$d = 20; 2^d = 10^6$ (Close to million)

As $d \uparrow$, worst case, # of adj. cells $2^d \uparrow$

when d is not small
So, $n = 1 \text{ million}$; $d = 20$.

So time complexity for $\text{INN} = O(2^d \lg n)$

$$O(2^d * \lg n)$$

In Brute Force $O(2^d \lg n)$

$$O(n)$$

So, since we have many more dimensions,
kd tree becomes useless.

② $O(\lg n)$ when d is small.

when data is uniformly distributed.

But non uniformly distributed in real world. So $O(\lg n)$ is no more valid.

It becomes $O(n)$ in non-uniform distribution.

Kd tree was not invented for KNN but for computer graphic since they require 2D or 3D data.

V26: Hashing vs. LSH

LSH = locality Sensitive Hashing.

Normal hashing

	0	1	2	3	4	5	6
Unordered array A ↴	2	1	5	4	6	7	5

We can store a in hashtable which is dictionary in python.

In normal hashing; there is hash function as $\text{hash}(x)$ which returns memory location, basically returns a no. which is interpreted as address.

If empty then data can be stored there, else a linked list like structure formed & its pointer is held by this memory location, in addition to data which it already had.

So hashing, we can go to memory location corresponding to an i/p.

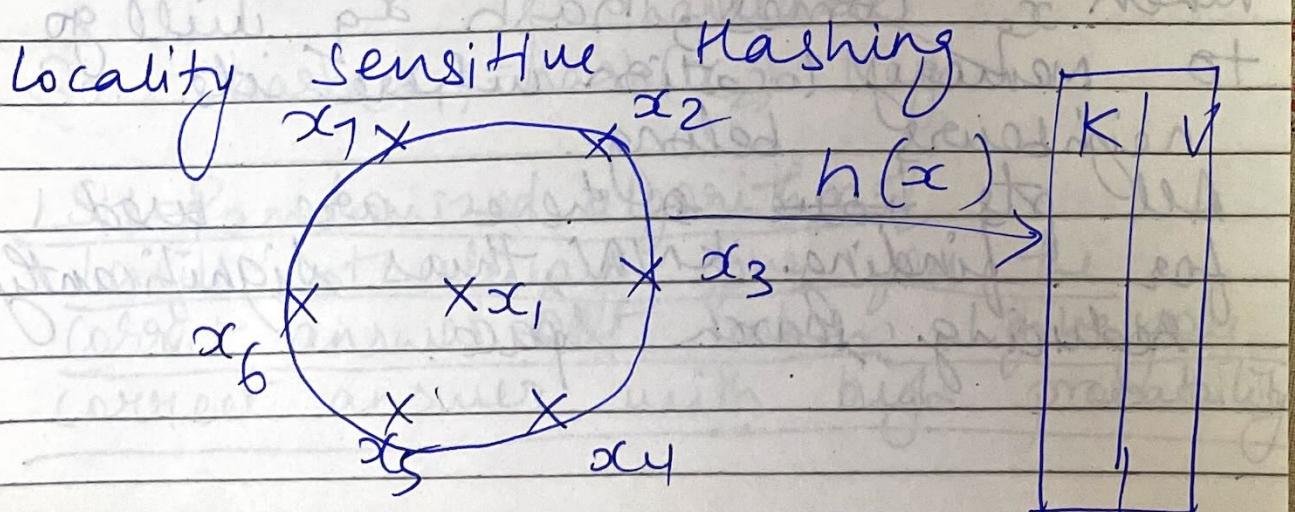
- Time complexity = $O(1)$.

Say an array is

0	1	2	3	4	5	6
21	15	16	17	15	18	

We want to store the indexes of each no. using hash table

x	hash(x)	loc.
2	0	
1	1	
5	2	5
6	3	
7	4	
8	6	
9		



~~go to each hashtable (total L tables)
& take union of all the sections
to which~~

It tries to achieve — all points that belong to the same neighbourhood hash to the same memory location.

So, here the $\text{hash}(x)$ takes a pt. x as IP and returns memory location.

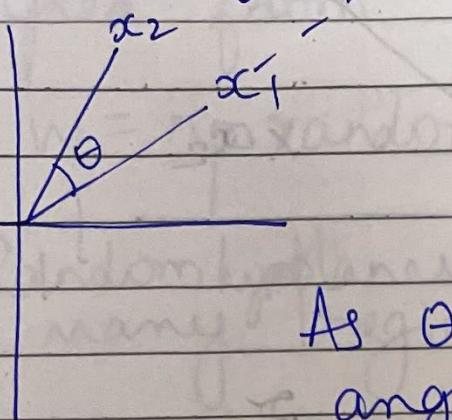
$$\text{hash}(x_1) = \text{hash}(x_2) = \text{hash}(x_3) = \dots \text{for ideal LSH.}$$

In the hash address, data pts are stored.

When x_q comes, $\text{hash } x_q$ will go to memory location where its neighbours belong.

All pts ~~have~~ there are used for finding KNN, thus significantly reducing search space.

V27 : LSH for cosine similarity

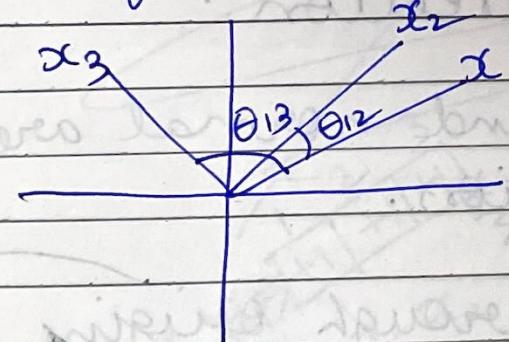


$$\cos\text{-sim}(x_1, x_2) = \cos\theta$$
$$\cos\text{-sim}(x_1, x_3) = \cos(\theta)$$
$$= 1$$

As $\theta \uparrow$, x_1 & x_2 are angularly separated \uparrow
distance \uparrow cosine \downarrow
similarity \downarrow

Cosine sim is all about angular similarity.

LSH for Cos-sim is:-



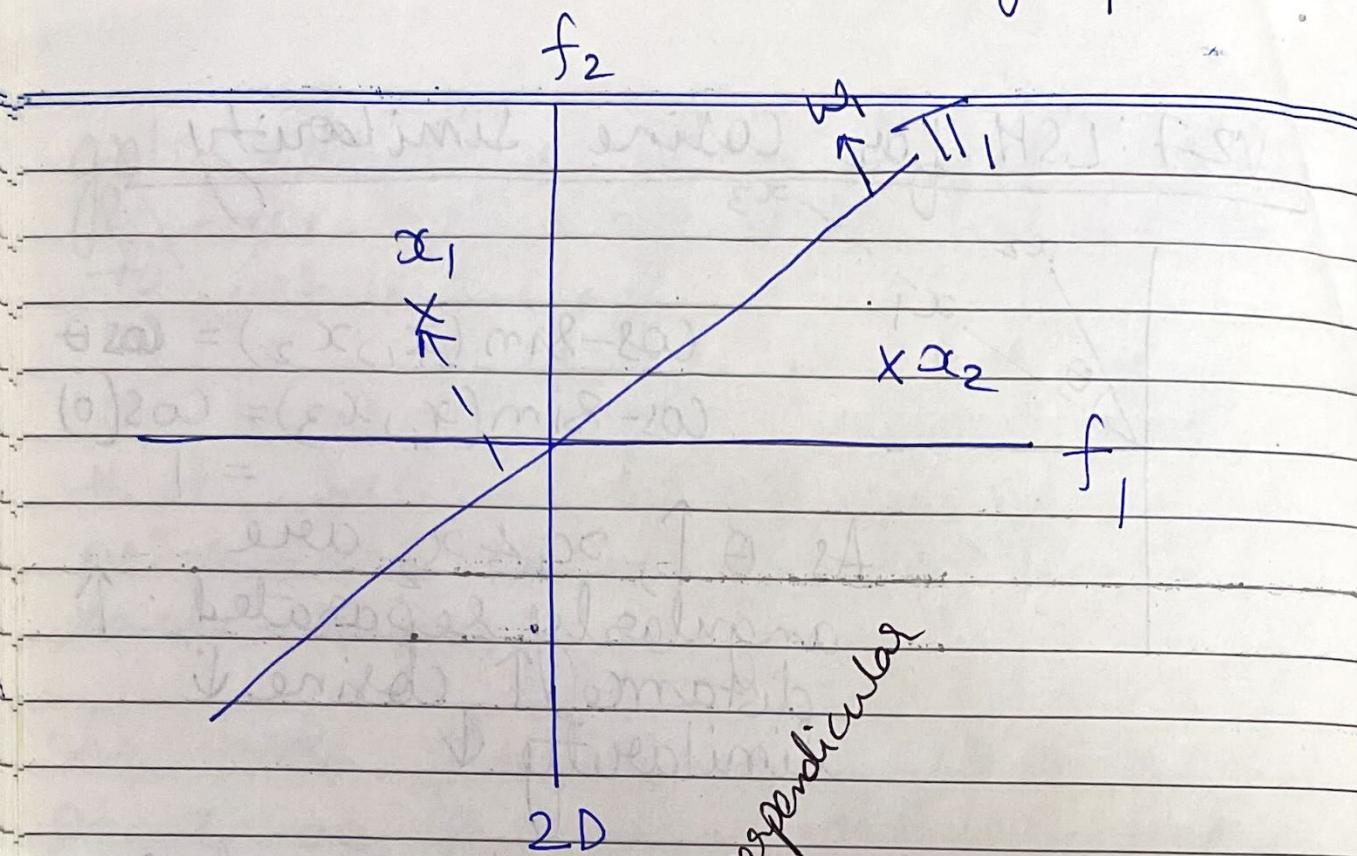
$\theta_{12} \approx 0 \Rightarrow x_1 \& x_2$
are similar
small
pls which are close,
go to same bucket in
hashtable.

So. $\text{hash}(x_1) = \text{hash}(x_2) = \text{key}$.

LSH = Randomized algorithm:

I will not always give you the correct answer but when I give you correct answer with high probability

Random Hyperplane.



w_1 is unit normal vector to π_1

$w^T x_1 \geq 0$ since x_1 and normal w , are in same direction.

$w^T x_2 \leq 0$ since x_2 and normal are in opposite direction.

Since plane passes through origin

$$w^T x = 0, 1, 2, 3, 4, \dots, d-1$$

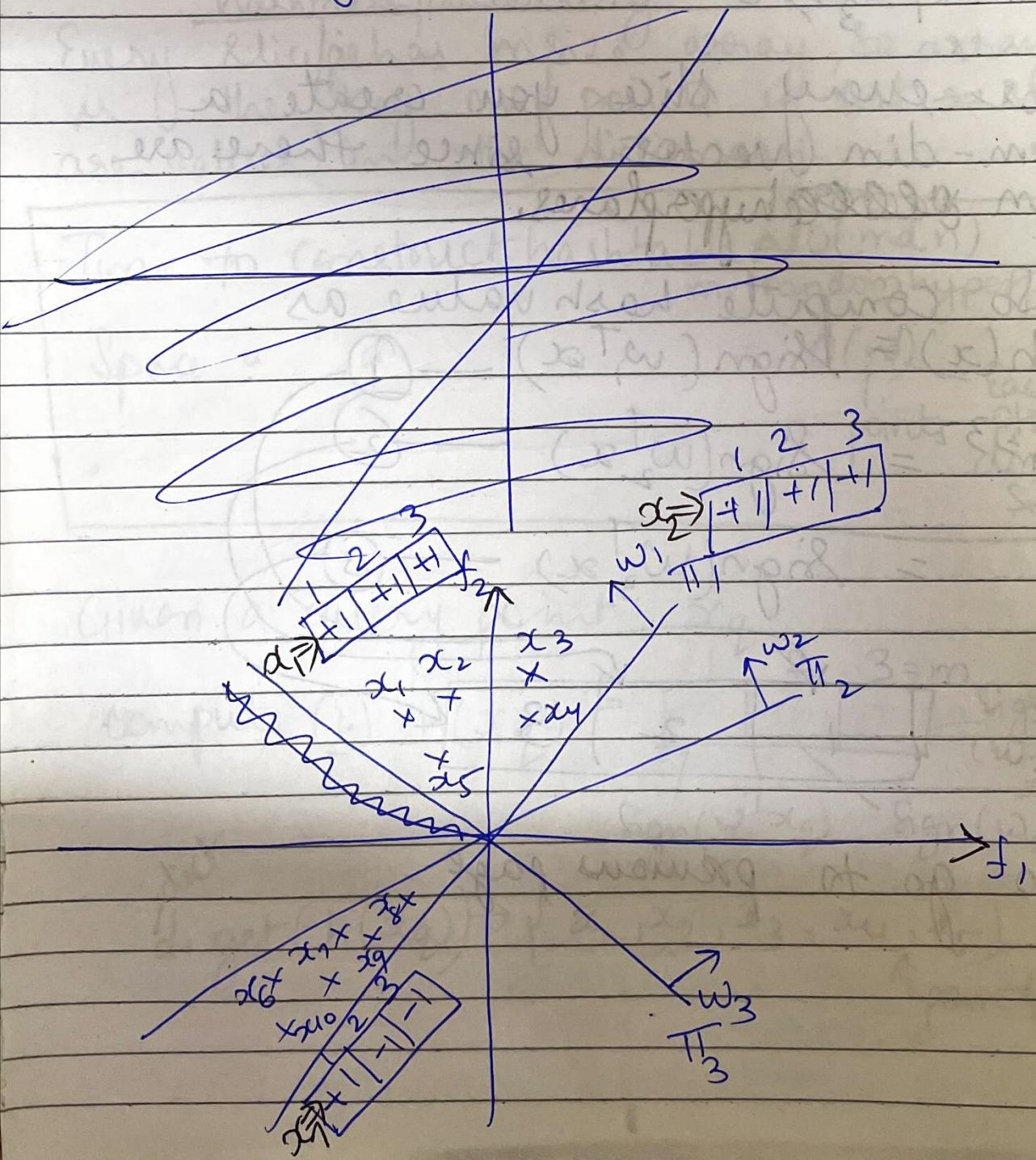
So, $w = \underbrace{[\quad | \quad | \quad | \quad | \dots | \quad]}_{d\text{-dim random vector}}$

(generate & dim vector
s.t.)

Each value in w are random values (nos.)
sampled from $N(0, 1)$

So $w = np.random.normal(0, 1, d)$

So Random hyperplanes can cut hyperspace
into many regions.



Each point gets assigned to each slice

So, $x^T w_1 \geq 0 \rightarrow +1$ if it is +ve we assign +1 i.e
 $x^T w_2 \geq 0 \rightarrow +1$ same direction of normal.
 $x^T w_3 \geq 0 \rightarrow +1$

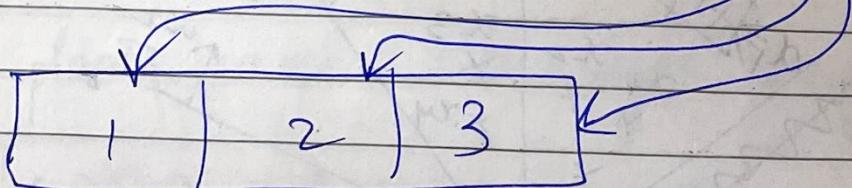
For every slice you create a m-dim vector since there are m ~~values~~ hyperplanes.

So compute hash value as

$$h(x) = \text{Sign}(w_1^T x) \quad \textcircled{1}$$

$$m=3 = \text{Sign}(w_2^T x) \quad \textcircled{2}$$

$$= \text{Sign}(w_3^T x) \quad \textcircled{3}$$



So go to previous page

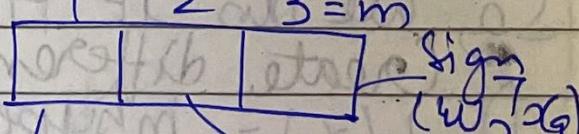
K	V
tuple S	+1 -1 -1 points
	+1 +1 -1 points

hashtable/dictionary
 Every slice has m size array, so array is the key and points in that slice are values in the dictionary.

Time to construct hashtable $\approx O(mdn)$
 $m = \text{random hyperplane}$

Space $n \cdot n = O(nd)$ Each pt's dimension

Given a query point x_q

compute (1) $h(x_q) =$ 

dict
 $d.\text{get}(h(x_q)) \rightarrow \{x_1, x_2, x_3, x_4, x_5\}$ points

So given a hash table
Time complexity for querying
 $sc_q = h(x_q) = O(md)$

If n' elements in slice then compare
 x_q So ~~total time~~ total time

$$O(md + n'd)$$

↑ If n' is small then
 $O(md)$

So typically m is set to be $\lg(n)$
So, Time $O(d * \lg n)$ for querying
complexity.

Problem if point lies in the adjacent slice, one could miss NN in cosine similarity as a measure of similarity.
So, we do it this way:-
1) Create different hash function for different types
 $h_1(x) = m$ hyperplanes (different)
2) $h_2(x) = m$ hyperplanes (different)

Different hyperplanes create different hashtable

Do repeatedly. Do it multiple times.

Your x_q will fall into same slice hence x_q and another point will fall in same slice.

After computing $h(x)$ of every set, ~~the~~ take set union of all the points, you get a big set.

Compare x_q with all pts. in big set.

Final Time Complexity to query
L Hashtables $O(mdL)$

L - hashtables (typically small)
 m = random hyperplanes
d - Each pts dimension.

$m = O(\lg n) \rightarrow$ No. of hyperplanes choice

As # of hyperplanes ↑
Slice ↑
pts per slice ↓

So we need to be careful : if # of pts in a slice $< k$, it can be problem.

Similar or closer points should go to same bucket ~~because it's~~ in hashtable

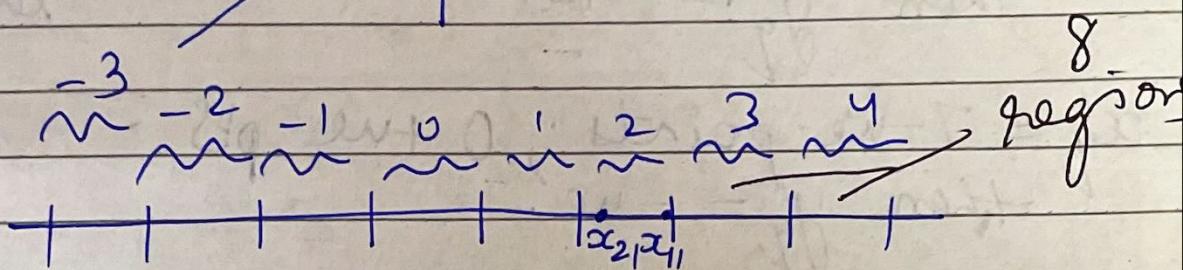
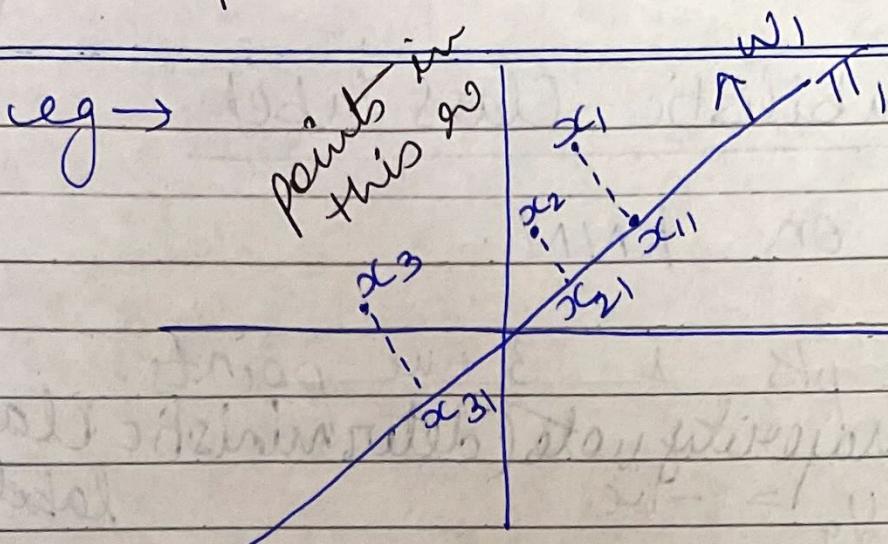
V28: LSH for Euclidean Distance

It is similar to cosine-similarity with simple extension.

As in previous case, here also space is divided by different hyperplanes.

- 1) For each hyperplane, take projections of all data points on that hyperplane.
- 2) Each hyperplane is divided into ~~say~~ regions.
- 3) The projections belonging to same region would be closer.
- 4) The hash (x_q) becomes the m-dimensional vector, where each component of the vector say 5th Component represents the region of 5th hyperplane where the projection of that point belongs to.

LSH is not perfect randomized algo.



$$\text{hash}(x_{11}) = \boxed{\begin{matrix} 1 & 2 \\ 2 & 1 & - & - & - & m \end{matrix}}$$

region in T_1 projection belongs

So, in euclidean distance $h(x)$ vector can have any numeric values unlike cosine similarity which only had $+1$ & -1 .

Just like cosine, here also we have L -different sets of m -hyperplanes.

Even if $x_1, x_2 \ll x_1, x_3$ still x_1 & x_2 are in same region & x_1 & x_3 in different after T_1 .

But since planes generated randomly T_2 also may occur & x_1 & x_3 will be in same region.

v29: Probabilistic Class Label

Say $k=7$ on KNN.

$x_g = 4$ -ve pts + 3 +ve points

Apply Majority vote (deterministic class labeling)
then $y_g = -ve$

$x'_g = 7$ -ve points; 0 +ve pts

then $y'_g = -ve$.

In case of x'_g we are more certain about the opposite label as compared to x_g .

But the certainty doesn't reflect in the o/p.

Quantifying the certainty giving probabilistic class label.

~~prob~~ $P(y_g = -ve) = 4/7$ } more informative

$P(y'_g = -ve) = 7/7$ } makes model more interpretable.