

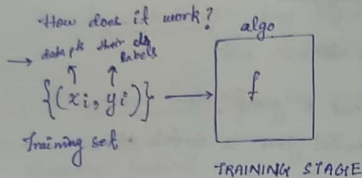
K-nearest neighbours

Intro to classification

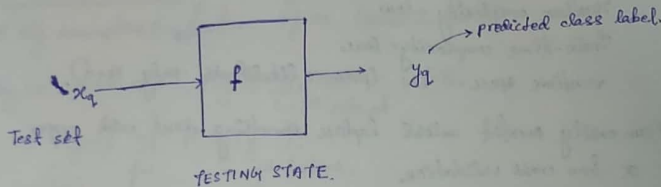
Given a new data-point, predict which class it belongs to.

$$y = f(x)$$

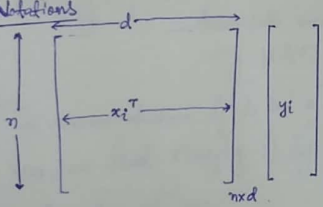
x = featurization of query data pt.
 y = predicted class label
 $f()$ = classifier.



Algo learns the function f from the training set data.



Notations



* By default, every data pt. is represented as col. vector, so, transpose taken here.

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

for binary classification.

D is a dataset of pair of values, n -such pairs, where the 1st of pair is a d -dimensional vector & 2nd is either 0 or 1.

binary classification $\rightarrow y_i \in \{0, 1\}$ only 2 classes out there.
 Ex \rightarrow Amazon food reviews. (true & -ve class).
 Multi-class \rightarrow more than 2 classes. Ex \rightarrow Iris dataset, MNIST, etc.

Regression $\rightarrow y_i \in \mathbb{R}$ real number. Ex \rightarrow predicting housing prices, etc.

K-nearest neighbour (Geometric intuition).

2-D toy dataset



Binary classification.

0 - blue datapoints

1 - red datapoints

⓪ - query datapoint

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

Objective

\rightarrow We want to determine class label of query point.

Intuitive Approach \rightarrow The points that are geometrically close to x_q or say in the neighbourhood of x_q are blue points mostly. Hence, we can conclude that x_q is also blue point. That's what KNN does.

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 .



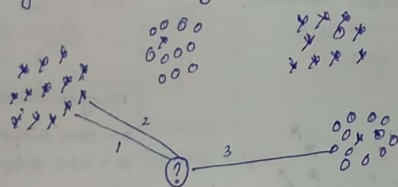
y_1, y_2, y_3

$\{y_1, y_2, y_3\}$

y_1	y_2	y_3
+	+	+
+	+	+
-	+	-
-	+	-

That's it

① query point is far away from clusters.



output label may say its
-ne because of majority
vote but it doesn't make
much sense.

It would be better if the
model could op - "don't know"

② (+)ve & (-)ve pts are randomly spread.

\rightarrow No useful information.

$d = \text{len of shortest line fm } x_1 \text{ to } x_2.$

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

Also called Euclidean distance.

Euclidean Distance :- If $x_1, x_2 \in \mathbb{R}^d$,

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

\downarrow
 L2 norm of $x_1 - x_2$

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

L1 norm of $x_1 - x_2$

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

$\underbrace{\hspace{10em}}_{L_p \text{ norm of } \underline{x_1 - x_2}}$

for $p=1$, \rightarrow manhattan
 $p=2$, \rightarrow euclidean.

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

In general, L_q norm of a vector x is given by. ($x \in \mathbb{R}^d$)

$$\|x\|_q = \left(\sum_{i=1}^d |x_i|^q \right)^{1/q}.$$

Hamming distance (b/w boolean vectors ~~only~~).

$x_1, x_2 \rightarrow$ boolean vectors \rightarrow Binary BOW.

$x_1, x_2 \rightarrow$ boolean vectors \rightarrow binary row.

$x_1 = [0, 1, 1, 0, 1, 0, \dots]$
 $x_2 = [1, 0, 1, 1, 0, 1, \dots]$

Hamming-dist(x_1, x_2) \Rightarrow # locations where binary vectors differ.

Also, hamming-distance could be used for measuring distances b/w strings.

cosine-distance and cosine similarity

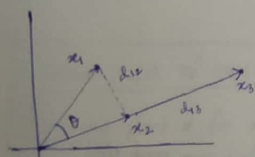
Intuitively, similarity $\uparrow \Rightarrow$ distance \downarrow
 & similarity $\downarrow \Rightarrow$ distance \uparrow .

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

Just as in Euclidean terms, points which have less distance b/w 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, $\cos\text{-sim}(x_1, x_2) = \cos \theta$ where θ = angle b/w x_1 & x_2 .



$$d_{12} < d_{13} \\ \text{but } \cos\text{-dist}_{13} < \cos\text{-dist}_{12}$$

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

If the problem deals with angular stuff, then taking euclidean distance doesn't make much sense and vice-versa.

* Also, $\cos \theta = \frac{x_1 \cdot x_2}{\|x_1\|_2 \|x_2\|_2}$ if x_1 & x_2 are unit vectors, $\cos \theta = x_1 \cdot x_2$

Relationship between euclidean & cosine distances
(See cosine similarity wikipedia for details.)

$$[\text{eucl-dist}(x_1, x_2)]^2 = 2 \times \cos\text{-dist}(x_1, x_2) \quad \text{if } x_1 \text{ \& } x_2 \text{ are unit vectors.}$$

$$= 2 \times (1 - \cos\text{-sim}(x_1, x_2))$$

Measuring effectiveness of KNN

① Split D into D_{train} ($n_1 \approx 70\%$) and D_{test} ($n_2 \approx 30\%$).
 $D_{\text{train}} \cap D_{\text{test}} = \phi$
 $D_{\text{train}} \cup D_{\text{test}} = D$

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

Algo:-

count $\leftarrow 0$

for each pt in D_{test} :

$x_q \leftarrow \text{pt}$

use D_{train} & KNN to determine y_q

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

count $\leftarrow \text{count} + 1$

In the end,

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

$$\text{Accuracy} = \frac{\text{count}}{n_2} \quad 0 \leq \text{accuracy} \leq 1$$

Accuracy is one of the measures. There are many others that'll come.

Time and Space complexity of KNN

Evaluation:- given x_q . Need to determine y_q .

① $\text{KNNpts} = []$
for each x_i in D_{train} :
- compute $d(x_i, x_q) \rightarrow d_i \rightarrow O(d)$
- keep the smallest k -distances (x_i, y_i, d_i) to KNN pts. $\rightarrow O(k)$

$k \ll n$, Hence, neglected
so, $O(n \times d)$ for above logic.

② cnt_pos = 0 ; cnt_neg = 0
for each x_i in KNNpts:
if y_i is +ve:
cnt_pos += 1
else
cnt_neg += 1
return (cnt_pos >= cnt_neg)
So, total $O(k)$.
So, overall time complexity = $O(n \times d) + O(k)$.
= $O(nd)$
if $d \ll n$, $O(n)$
time takes to compute label of a single query pt.

Space complexity - We need to calculate distance for each pt in D_{train} .

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

Space $= O(n \times d)$ \rightarrow Not good. Very high.

Limitation of KNN (simple implementation).

Consider amazon food reviews example.

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

With simple KNN,
Time $= O(nd)$
Space $= O(nd)$.

$n \approx 364K$
 $d \approx 100K$

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

② Time complexity - 36 Billion computations.

Amazon
services

Review $\xrightarrow{1ms}$ (+)ve / (-)ve label.
low latency.

this should be the case.
If latency \uparrow , bad user
experience.

for Trading & finance, as much as $1\mu s$ delay is permitted only.

Non-low latency systems \rightarrow in medicines recommendation, etc.

This is one of the main reason of not implementing KNN in production much
instead of it being so simple & intuitive.

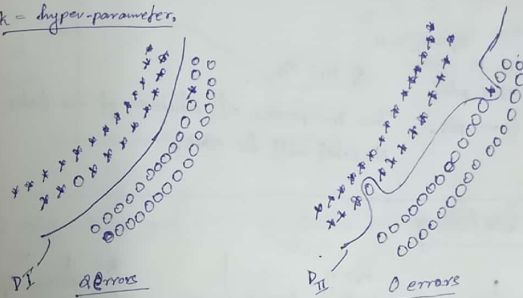
Some modifications to KNN later to be covered -

① Kd-tree

② LSH (Locality sensitive hashing).

Decision surface of KNN as k -changes.

k = hyper-parameter.



D_I & D_{II} are the decision surfaces.

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

\rightarrow because, of generalization of higher dimensions.

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

\rightarrow for KNN, decision surface could be found by following

- \rightarrow take each point in the space of the datapoints, surrounding them.
- \rightarrow find its k -nearest neighbours.
- \rightarrow find the majority vote & assign that label to that point.
- \rightarrow 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 & consider a point in close proximity of 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, but due to majority voting of remaining 4 nearest neighbours, that pt will be labelled the opposite of outlier, giving smooth surface as D_I .

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

Let, say $k = n_0$. $n_1 + n_2 = n$.

$n_1 =$ no of (+)ve pts

$n_2 =$ no of (-)ve pts.

$\& n_1 > n_2$

Then irrespective of whatever pt we take, its label will be (+)ve.

Overfitting & Underfitting

$k = 1$ x

→ Lazy model.

Checks which class has more data pts & assigns that label to new query pts.

→ Underfit model.

$k = 5$

→ tries to find a balance b/w both extremes

→ Less prone to noise. Makes more sensible decisions.

$k = 1$ x

→ this model overworks & get more complicated decision surface & as to ensure no errors.

→ It may take noise as data & classify new pts incorrectly.

Determining optimal value of k . (Need for cross validation).

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

overfit

n

underfit.

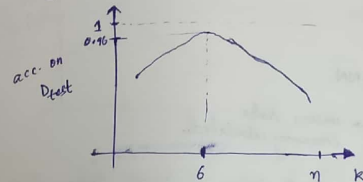
So, the optimal value will be somewhere middle.

Initial thoughts

trying out diff value of k & training on D_{train} and checking the accuracy on D_{test} .

	train data	acc on D_{test} .
$k=1$	D_{train}	0.78
$k=2$	"	0.82
$k=3$	"	0.85
⋮	⋮	⋮

A typical graph will look as follows.



So, we can conclude $k=6$ gives the best accuracy on D_{test} when using D_{train} as training data.

using D_{train} and 6-NN on Amazon reviews dataset, accuracy = 96%.

Small problems.

The sole objective of ML is to perform well prediction on unseen data.

We used the whole D_n here to calculate the model "function" f using D_{train} as training data.

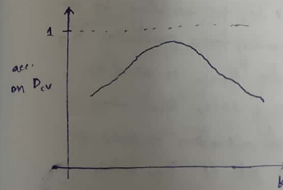
gives the Nearest Neighbours k value.

So, we can't say here the accuracy of model is 96% on unseen data. (Typical results ~ 80%)

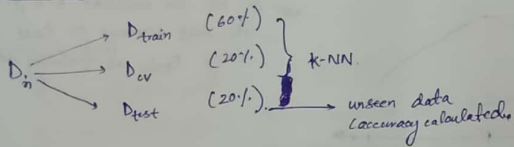
Cross validation

D_n splits into D_{train} (used to determine right 'k' value), D_{dev} (used to determine right 'k' value), and D_{test} (apply f to get acc. on unseen data).

Now, if D_{test} gives 73% accuracy on the optimal value of 'k' found on cross-validation, we can say the accuracy of model as 73% on unseen data.



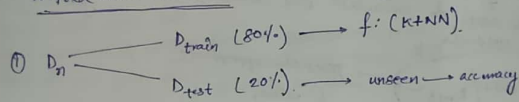
K'-fold cross validation



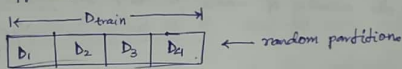
Problem:- we get only 60% of data as D_{train} .

* More the training data, better is the algorithm. So, we should somehow try to get whole 80% ($D_{test} + D_{cv}$) as training data.

K'-fold cross validation does this very smartly.



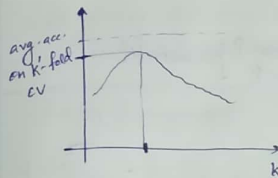
② Say $K'=4$, So, we divide D_{train} to 4 parts of equal size.



	Train	cv	acc. on cv
$K=1$	$D_1 D_2 D_3$	D_4	a_4
$K=1$	$D_1 D_2 D_4$	D_3	a_3
$K=1$	$D_1 D_3 D_4$	D_2	a_2
$K=1$	$D_2 D_3 D_4$	D_1	a_1
$K=2$	$D_1 D_2 D_3$	D_4	a'_4
$K=2$	$D_1 D_2 D_4$	D_3	a'_3
$K=2$	$D_1 D_3 D_4$	D_2	a'_2
$K=2$	$D_2 D_3 D_4$	D_1	a'_1

The points are plotted on Accuracy Vs K plot.

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



So, as we wanted, all the 80% of data used as training data.

But, what should be value of K?

→ $K'=10$ - Rule of Thumb.

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

Knowing Overfitting or Underfitting from given 'K'

* Given a value of 'K', how to determine whether it underfits or overfits?

→ graphical methods.

Some terminologies:-

$$\text{Accuracy} = \frac{\text{\#pts correctly classified}}{\text{Total \#pts.}}$$

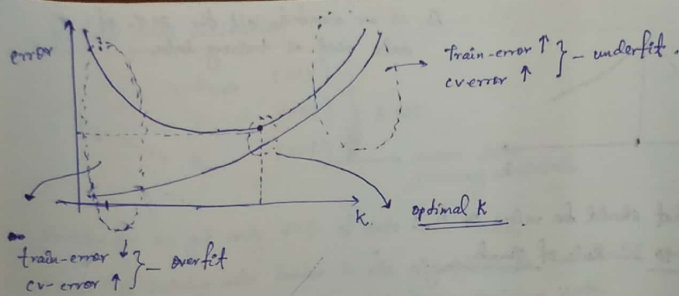
$$\text{error} = 1 - \text{Accuracy}$$

Training error:-

Usually we use D_{train} to supply the NN pts. and use D_{cv} to calculate accuracy against different values of K.

The corresponding loss of this accuracy is called cross validation loss.

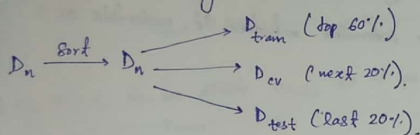
Instead of using D_{cv} pts, if we use D_{train} pts itself to calculate accuracy, the corresponding loss would be training loss.



Time-based splitting

→ works only if we have timestamp data.

① sort D_n in ascending order of time.



In datasets like Amazon food reviews, a time-based splitting is better than random splitting due to its use case.

* The model has to be trained on the old data we have. (no other option). But when deployed, it has to operate on new data.

With time, products and their reviews change.

So, we have to design the model in such a way that it takes older (or train) data.

more older data & gets good predictions on less older (more recent) data.

This can be achieved via by splitting in the above mentioned way.

Random splitting won't give this ability to the model.

So, whenever time data available and behaviour/data changes over time, time based splitting preferable.

K-NN for regression

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

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

① given x_1 , find k-nearest neighbours $(x_1, y_1), (x_2, y_2), \dots$

② $y_1 = \text{mean}(y_i)_{i=1}^k$
or
 $y_2 = \text{median}(y_i)_{i=1}^k$

weighted K-NN

The closer pts among the KNN pts are given more weight than farther pts.

Now, while voting, these weights are added for each class & majority value wins.

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

→ there can be many more as well.

Voronoi diagram

→ follow similar procedure as while trying to draw decision boundary. & take $k=1$ (one nearest neighbour). Colour each region with different colour.

Efficient implementation of K-NN

Simple implementation $\left. \begin{array}{l} \text{K-NN} - O(n) \text{ time} \\ O(n) \text{ space} \end{array} \right\} \text{ if } k \text{ \& } d \text{ are small.}$

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

We can improve time complexity using data-structure called kd-tree.

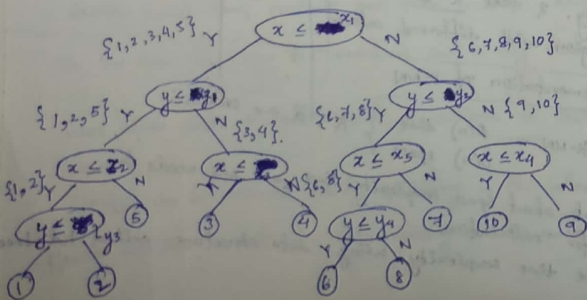
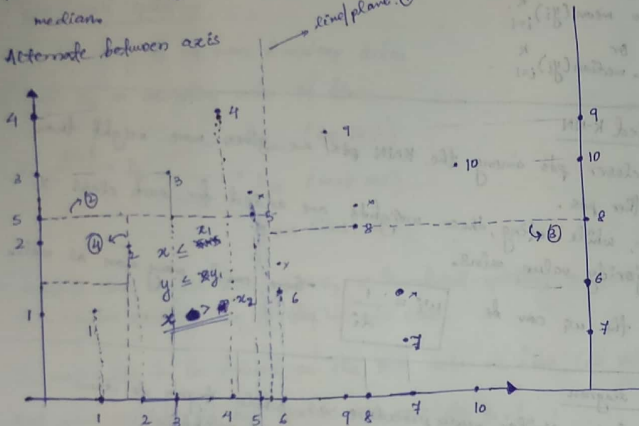
k-d tree — makes time complexity $O(\log n)$.
 — very popularly used in computer graphics.

Just like BST used for numbers (scalars), when that concept extended to vectors, its called k-d tree.

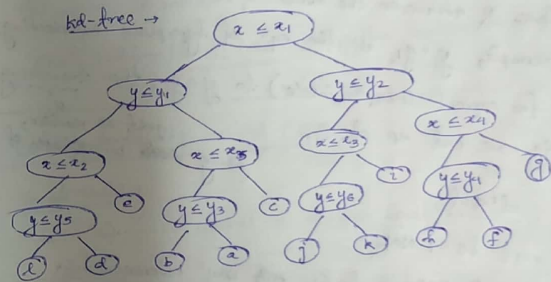
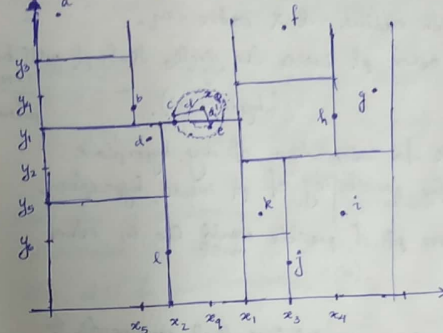
k-d tree — breaks down space using axis parallel hyperplane into hypercuboids.

① pick x-axis, project pts onto it, compute median, split data using median.

② Alternate between axis



Finding nearest neighbours using k-d tree.



x_q = query point.

We'll find 1-NN of x_q in the following, which could be extended to get K-NN.

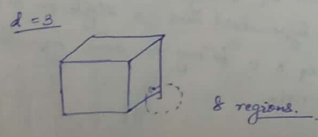
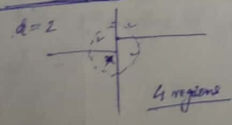
- ① Parse the tree for x_q (query pt). This will give the hypercuboid to which x_q belongs to.
- (Each hypercuboid has a single pt in it). Since x_q & c belong to same hypercuboid, c potentially could be the 1-NN.

- ② Join line joining x_q & c . let distance $= d$.
- ③ Draw a hypersphere with radius $= d$ & centre $= x_q$.
If there will be any other pt inside the circle that pt will be 1-NN (probably) (hypersphere).
- ④ We algorithmically check for intersection of any hyperplane with the hypersphere for possibility of pt inside hypersphere. This is because, any closer pt if possible must lie in other hypercube.
- ⑤ Track back to the condition in kd-tree that corresponds to that hyper-plane & traverse the other side.
Here, in this example, $y = y_1$ is the eqⁿ of hyperplane intersecting. Tracking back from current locatⁿ (i.e. 'c') to $y \leq y_1$, then traversing the other side, we end up at 'E'.
just as earlier, using ~~same~~ co-ordinates of x_q .
- ⑥ $x_q - c = d'$. $d' \leq d$
Hence 'c' can't be 1-NN.
go to ② & repeat till no further 'new' intersections.

Extending 1-NN to KNN,
 Best case $= O(\log n)$.
 worst case $= O(n)$.
 Best case $= O(K \log n)$.
 worst case $= O(K \times n)$.

All these if d is 'small'

Limitations of kd-tree



So, in worst-case, for d -dim data, hypersphere intersection may result in looking upto 2^d regions \rightarrow exponential.

~~if d is large~~, time complexity $= O(2^d \log n)$.
 If d is large, time complexity $= O(n \log n) \rightarrow$ worse than brute-force.
 If $2^d \geq n$, time complexity $= O(n \log n)$.

- ② Even when d is small, $O(\log n)$ holds when data uniformly distributed in space.
Else, it'll be more towards $O(n)$.
- * The whole kd-tree was designed for computer graphics (2-D & 3-D) for which it works pretty well.
Read wiki - /k-d_tree - python implementation, variations ball tree, etc.

Locality Sensitive Hashing (LSH)

Normal Hashing

hash(x) \rightarrow returns a memory locatⁿ { basically returns a no. which is interpreted as address }
 i/p no.
 If empty, then data can be stored there, else a linked-list like structure formed & its ptr is held by this memory locⁿ (in addition to data which it already had).

* So, by hashing, we can go to a memory locatⁿ corresponding to an i/p in $O(1)$.

Say, an array \rightarrow we want to store the indexes of each number using hash table.

	0	1	2	3	4	5	6
	2	1	5	6	7	5	8
x	hash(x) location						
2	0						
1	1						
5	2						
6	3						
7	4						
8	5						

Space: $O(nd)$. dimension (can be large).

Time: $O(n \times m \times d)$ — training.

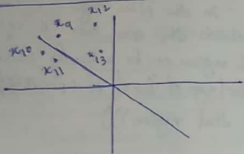
$O(m \times d + n' \times d)$ — Runtime n' — no of pts in each section.
small, typically.

So, $O(md)$.

$m = \log(n)$ — usually taken.

Hence, runtime $\rightarrow O(d \log n)$.

Problem with LSH



clearly x_{10} & x_{11} are nearest neigh. of x_9 . But due to separation to diff segments, x_9 thinks x_{12} & x_{13} are nearest.

How to overcome?

→ Instead of having a single set of m -hyperplanes, we'll have L -sets of m -hyperplane each, each having its own hash table. → Each set having different m hyperplane than other sets.

If two points are close, then they'll fall in the same section in atleast 1 of the hyperplane divisions.

So, for each point x_0 ,

① calculate which section it belongs to in each of the L sets.

hash₁(x_0) =

1	-1	1		-1
---	----	---	--	----

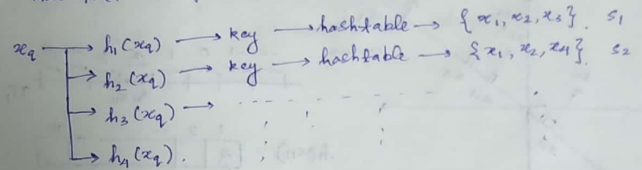
 → m -dimensional

② hash₂(x_0) =

1	-1	-1		1
---	----	----	--	---

goto each hash table (total L tables), & take union of all the sections to which x_0 belongs to.

Now find KNN from this set of pts.



KNN searched from $S_1 \cup S_2 \cup S_3 \dots S_p$

time complexity $\rightarrow O(mdL) \rightarrow \{ \text{Query pt} \}$

* As no of hyperplanes \uparrow , no of slices \uparrow , #pts per slice \downarrow .
 So, we need to be careful, because if #pts in a slice $< K$, then it can be problematic.

* But if #hyperplane \downarrow , no of slices \downarrow , #pts per slice \uparrow .

Hence, we may need to compare a lot of pts.

Thus, we need to find proper balance. b/w both extremes. usually $(m = \log(n))$

LSH for euclidean-distance (implementation).

→ similar to cosine-similarity with simple extension.

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

① For each hyperplane, take projections of all data pts on that hyperplane.

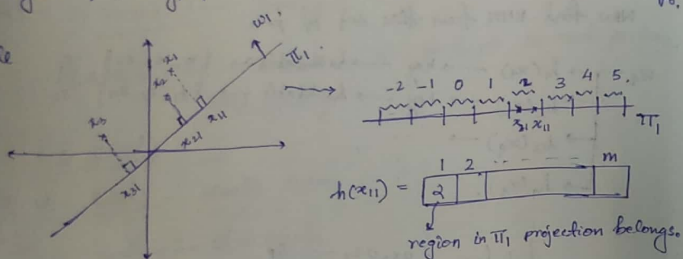
② each hyperplane is divided into say, 'a' regions.

③ The projections belonging to same region would be intuitively closer.

④ The hash(x_0) becomes this 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.

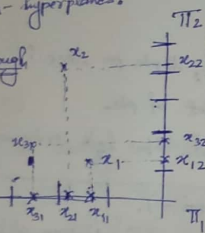
Example



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

Just like in cosine, here also we can have 1 -different sets of m -hyperplanes.

One catch though



even if $x_1x_2 < x_1x_3$ still x_1x_2 are in same region & x_1x_3 in different, for Π_1 .

But since planes generated randomly, Π_2 also may occur & x_1x_3 will be in same region there.

* LSH used extensively used in computer vision applications.

Probabilistic class label

say $K=7$ on KNN.

x_1 : 4 -ve pts & 3 +ve pts $\rightarrow y_1 = \bar{0}$ ve. } less informative,

x_2 : 7 -ve pts $\rightarrow y_2 = 4$ -ve

In case of x_2 we are more certain about the op label as compared to x_1 .

But the certainty doesn't reflect in the OP.

Quantifying the certainty \rightarrow giving probabilistic class label.

$$\begin{aligned} P(y_1 = -ve) &= 4/7 \\ P(y_2 = -ve) &= 7/7 \end{aligned} \quad \left. \begin{array}{l} \text{more informative} \\ \text{Makes models more interpretable.} \end{array} \right\}$$