

* classification and Regression models K-NN

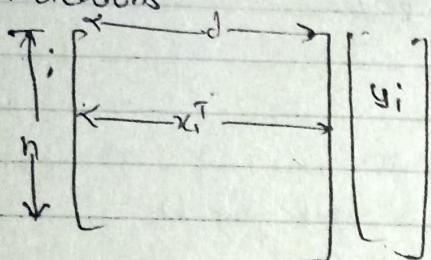
classification :- given a new data point , predict its category by classification model .

$$y = f(x)$$

category datapoint

classification is all about finding the funcn .

→ Notations



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

set of pair x_i, y_i

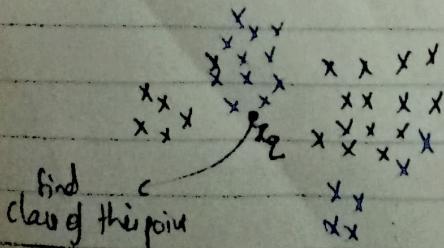
If y have only 2 values then where $x_i \in \mathbb{R}^d$ and $y_i \in \{0, 1\}$, it is called 2 class classif or binary classification

In MNIST , $y \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$:: 10-class classification
also called multi-class classification problem .

But what if y_i is know more pair to smaller fintest of classes and y_i can have any real value . In that case classification Not work . In that case we use Regression .

→ K-NN geometric intepretation

2D-toy dataset : binary classification



x^+ : +ve data point

x^- : -ve data point

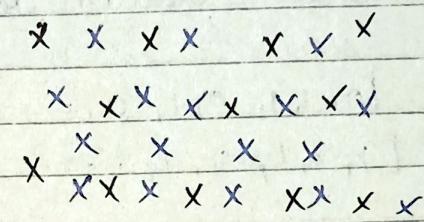
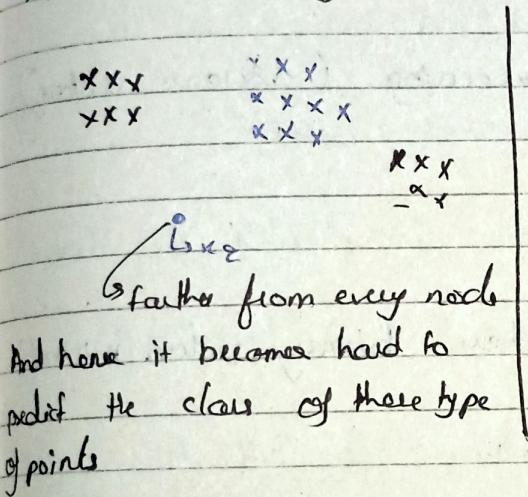
$$D = \{(x_i, y_i) \mid x_i \in \mathbb{R}^2, y_i \in \{0, 1\}\}$$

Take all the points which are close to x_2 :
 As we can see that, all the neighbours of x_2 are +ve we can
 hence conclude that x_2 is +ve.

Step 1: find K-nearest point to x_2 in D

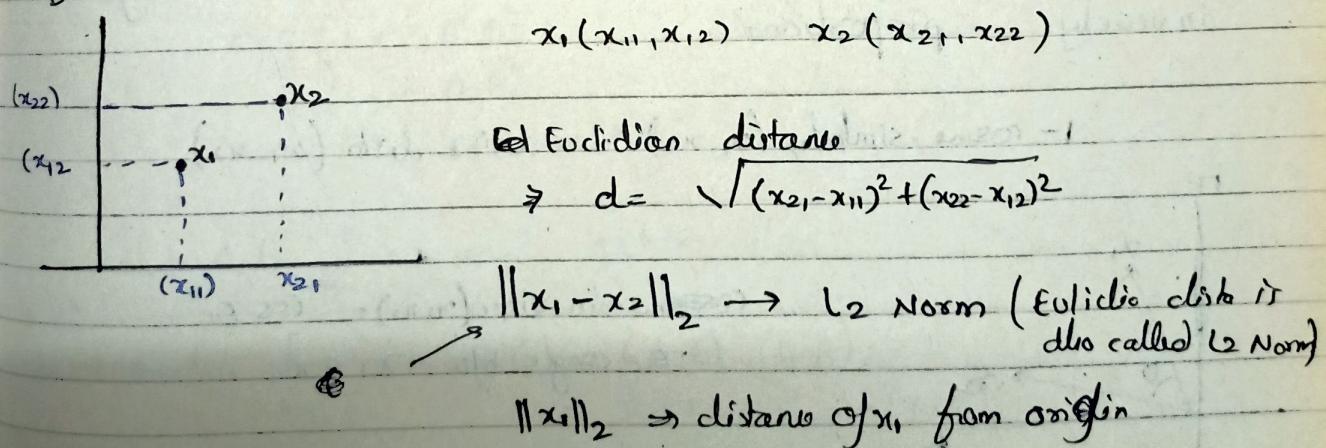
Step 2: choose the majority y_i is + among neighbours
 decide the y_i at y_2 (class of x_2)

→ Failure cases of K-NN



→ In these types of plots it is hard to work K-NN.

* Distance Measures :-

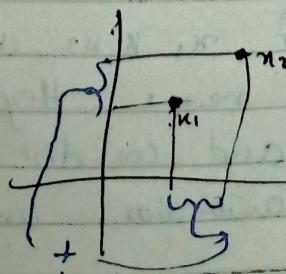


→ Manhattan's distance:

$$\sum_{i=1}^d |x_{1i} - x_{2i}|$$

This is L1 Norm. Hence written as

$$||x_1 - x_2||_1$$



L_p norms \rightarrow Minkowski distance.

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

If $p=2$, Minkowski dist. = Euclidean dist.

If $p=1$, Minkowski dist. = Manhattan dist.

$$L_p$$
 Norm of single vector: $\|x_i\|_p = \left(\sum_{i=1}^d |x_{is}|^p \right)^{1/p}$

\rightarrow Hamming distance: used in text processing (boolean vectors)

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

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

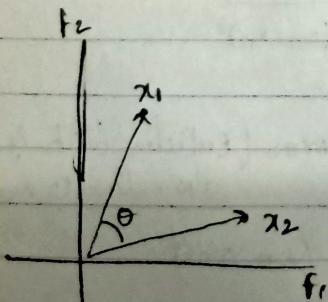
Hamming dist. $(x_1, x_2) = \# \text{loc where binary values mismatch}$

\rightarrow Cosine similarity & cosine-distance:-

a) similarity b) distance

If similarity increases, distance decreases. They are inversely proportional.

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



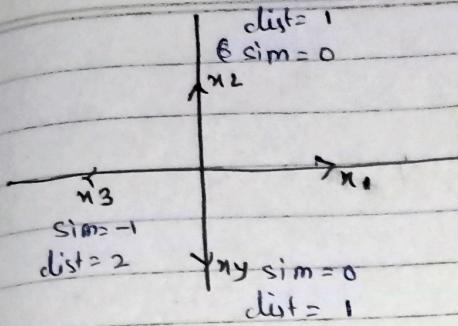
$$\text{cosine-sim.}(x_1, x_2) = \cos \theta$$

$$\therefore \theta = \text{angle b/w } x_1 \text{ & } x_2$$

If x_1 & x_2 are in same dir. then $\text{cos-sim.}(x_1, x_2) = 1$
that means they are very much similar.

$$\text{and cos-dist.} = 1 - \text{cos-sim.} = 1 - 1 = 0$$

Hence, cosine-sim. ~~will~~ give the angular distance.



$$\cos\text{-sim}(u_1, u_2) = \cos(\theta) = \frac{\mathbf{u}_1 \cdot \mathbf{u}_2}{\|\mathbf{u}_1\|_2 \|\mathbf{u}_2\|_2}$$

If \mathbf{u}_1 & \mathbf{u}_2 are unit vec : Then $\|\mathbf{u}_1\|_2 = \|\mathbf{u}_2\|_2 = 1$

$\cos\theta = \mathbf{u}_1 \cdot \mathbf{u}_2$ if \mathbf{u}_1 & \mathbf{u}_2 are unit vector.

what is Relationship b/w Euclidean dis and ~~cosine sim~~

$$\|\mathbf{A} - \mathbf{B}\|^2 = (\mathbf{A} - \mathbf{B})^T (\mathbf{A} - \mathbf{B}) = \|\mathbf{A}\|^2 + \|\mathbf{B}\|^2 - 2\mathbf{A}^T \mathbf{B}$$

If \mathbf{A} & \mathbf{B} are unit vector the $\|\mathbf{A}\|^2 = \|\mathbf{B}\|^2 = 1$

$$\begin{aligned}\|\mathbf{A} - \mathbf{B}\|^2 &= 2 - 2\mathbf{A}^T \mathbf{B} \\ &= 2(1 - \mathbf{A}^T \mathbf{B}) \\ &= 2(1 - \mathbf{A} \cdot \mathbf{B}) \\ &= 2(1 - \cos\text{-sim}(\mathbf{A}, \mathbf{B}))\end{aligned}$$

∴ If \mathbf{A}, \mathbf{B} are unit vector then:

$$\begin{aligned}\text{sqar of Eucli dis } \|\mathbf{A} - \mathbf{B}\|^2 &= (\mathbf{A} - \mathbf{B})^T (\mathbf{A} - \mathbf{B}) \\ &= 2(1 - \cos\text{-sim}(\mathbf{A}, \mathbf{B})) \\ &= 2 * \text{cos-clst}(\mathbf{A}, \mathbf{B})\end{aligned}$$

* Measuring how good K-NN is.

* Time & space complexity:

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

O/P: y_q

KNN pts = []

$O(n) \rightarrow$ for each x_i in Train :

$O(d) \rightarrow$ - compute $d(x_i, x_q) \rightarrow d_i$

- Keep the smallest K -distances (x_i, y_i, d_i)
put it in knn_{pts}

$O(n \cdot d)$ ↗ dimension / # of features

$c_{nt} = 0$; $c_{nt_neg} = 0$;

for each x_i in knn_{pts}:

if y is +ve

c_{nt_pos} ++

else

c_{nt_neg} ++;

If $c_{nt_pos} > c_{nt_neg}$

then $y_q = 1$.

else

then $y_q = 2$

Time : $O(n \cdot d)$

space x_q : ~~Train~~ D_{Train} should be in memory
 $\hookrightarrow O(n \cdot d)$

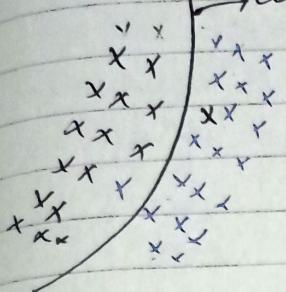
Extra space is not req. but to run K-NN we have to store D_{Train} .

\therefore Time = Space = $O(n \cdot d)$

* Limitation of K-NN: ~~•~~ High time & space complexity.
There are some other implementations of K-NN,
like K-d tree & LSH. We will learn them later.

→ How to pick right K in K-NN :- K in K-NN is known as hyper-parameter.

decision surface. (which categorizes the data)



How draw decision surface :- take many query points and check their output and using them draw a decision surface.

as K increase smoothness of decision surface increase.

what if $K=n$ i.e. max $K = \text{total no. of points}$.

Suppose there are n_1 positive points in dataset and n_2 negative points.

Here if $n_1 > n_2$ then we get positive class for every test class.

else we get negative class.

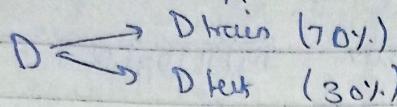
* overfitting & underfitting :-

overfitting :- when we try to fit every point; extremely non-smooth surface. It may also try to fit outliers.

underfit :- We are really not fitting all data points. There are still chance of improvement.

* Need for ~~c~~ cross-validation:

How to determine K

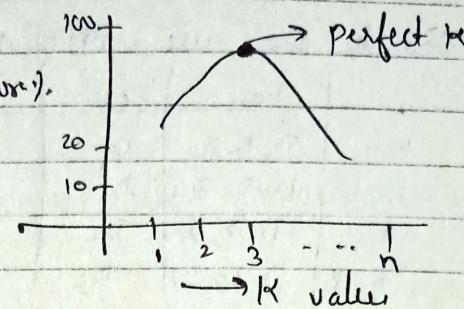


check for several values of K and check. accu.

Acc

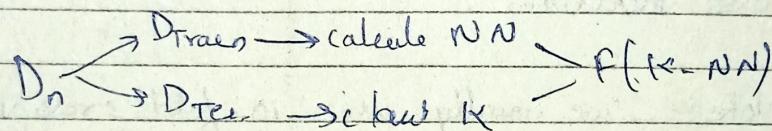
$K = 1$	0.78
$K = 2$	0.79
$K = 3$	0.84
:	

Accuracy.

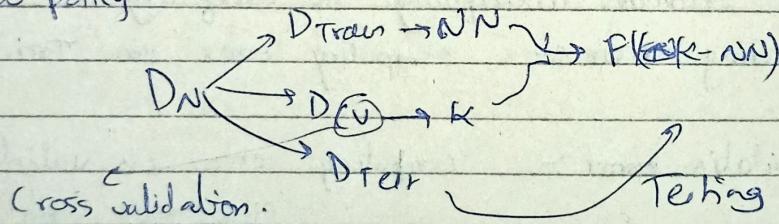


Now, one thing to be noted is we are using D_{test} to calculate K and now we are testing on same D_{test} so, results are not reliable.

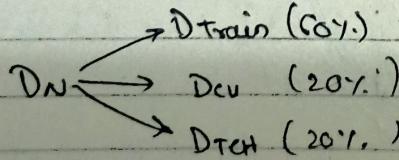
~~How~~ previously



New policy



* K -fold cross-validation

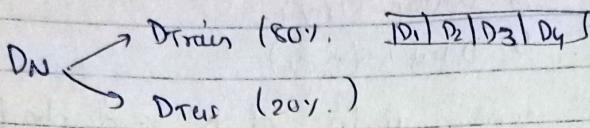


There is a problem we are only using 80% of data as training. We can't do anything for 20% test but can we somehow use $D_{\text{train}} + D_{\text{CV}} = 80\%$ of data. That's where the K -fold cross-validation is used.

Note: (K) -fold : , (K) NN

Their roles not be same.

What we do is we divide the D_{train} (80%) data into 4 equal parts



One by one, we ~~pick~~ one of D_i as cross valid.

	Train	CV	Acc
$K=1$	D_1, D_2, D_3	D_4	x
$K=2$	D_2, D_3, D_4	D_1	y
$K=3$	D_1, D_3, D_4	D_2	z
$K=4$	D_1, D_2, D_3	D_4	e

→ Take average and this is the value of 1-NN same with 2-NN, 3-NN
4-NN, and check which NN have max accuracy.

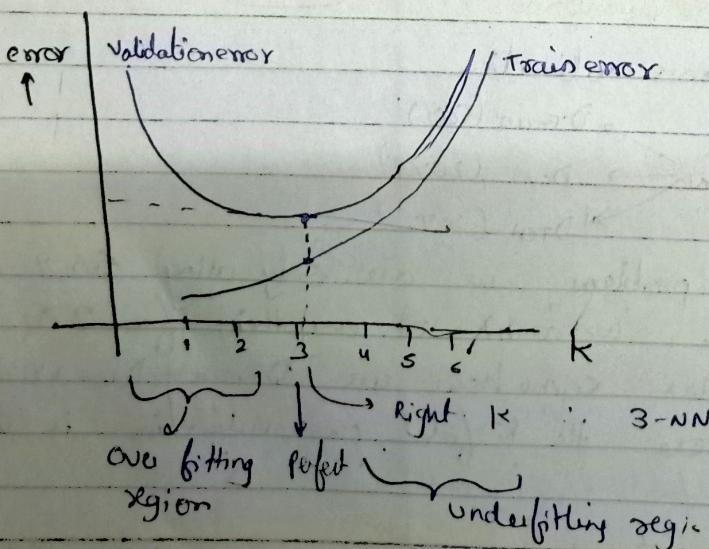
We are divided D_{train} into 4 folds, then it is called 4-fold cross validation.

Note :- we usually use 10-fold cross validation.

→ How to determine underfitting and overfitting.

Training error :- computing error on Train data points.

Validation error :- computing error on validation data,



- if training error and cross-validation error both are high then it signifies : Underfitting
- if train error is low and cv error is high then it is overfitting

* Time-base splitting

$D_n \rightarrow D_{\text{train}}$
 \downarrow
 $\rightarrow D_{\text{cv}}$
 \downarrow
 $\rightarrow D_{\text{test}}$

Till now we are randomly splitting the data.

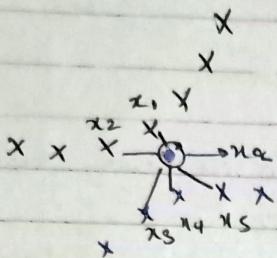
But random splitting is not the only way to split.
 like Time-based splitting.

In Time-base splitting :

x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}	x_{19}	x_{20}	x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}	x_{29}	x_{30}	x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}	x_{39}	x_{40}	x_{41}	x_{42}	x_{43}	x_{44}	x_{45}	x_{46}	x_{47}	x_{48}	x_{49}	x_{50}	x_{51}	x_{52}	x_{53}	x_{54}	x_{55}	x_{56}	x_{57}	x_{58}	x_{59}	x_{60}	x_{61}	x_{62}	x_{63}	x_{64}	x_{65}	x_{66}	x_{67}	x_{68}	x_{69}	x_{70}	x_{71}	x_{72}	x_{73}	x_{74}	x_{75}	x_{76}	x_{77}	x_{78}	x_{79}	x_{80}	x_{81}	x_{82}	x_{83}	x_{84}	x_{85}	x_{86}	x_{87}	x_{88}	x_{89}	x_{90}	x_{91}	x_{92}	x_{93}	x_{94}	x_{95}	x_{96}	x_{97}	x_{98}	x_{99}	x_{100}	x_{101}	x_{102}	x_{103}	x_{104}	x_{105}	x_{106}	x_{107}	x_{108}	x_{109}	x_{110}	x_{111}	x_{112}	x_{113}	x_{114}	x_{115}	x_{116}	x_{117}	x_{118}	x_{119}	x_{120}	x_{121}	x_{122}	x_{123}	x_{124}	x_{125}	x_{126}	x_{127}	x_{128}	x_{129}	x_{130}	x_{131}	x_{132}	x_{133}	x_{134}	x_{135}	x_{136}	x_{137}	x_{138}	x_{139}	x_{140}	x_{141}	x_{142}	x_{143}	x_{144}	x_{145}	x_{146}	x_{147}	x_{148}	x_{149}	x_{150}	x_{151}	x_{152}	x_{153}	x_{154}	x_{155}	x_{156}	x_{157}	x_{158}	x_{159}	x_{160}	x_{161}	x_{162}	x_{163}	x_{164}	x_{165}	x_{166}	x_{167}	x_{168}	x_{169}	x_{170}	x_{171}	x_{172}	x_{173}	x_{174}	x_{175}	x_{176}	x_{177}	x_{178}	x_{179}	x_{180}	x_{181}	x_{182}	x_{183}	x_{184}	x_{185}	x_{186}	x_{187}	x_{188}	x_{189}	x_{190}	x_{191}	x_{192}	x_{193}	x_{194}	x_{195}	x_{196}	x_{197}	x_{198}	x_{199}	x_{200}	x_{201}	x_{202}	x_{203}	x_{204}	x_{205}	x_{206}	x_{207}	x_{208}	x_{209}	x_{210}	x_{211}	x_{212}	x_{213}	x_{214}	x_{215}	x_{216}	x_{217}	x_{218}	x_{219}	x_{220}	x_{221}	x_{222}	x_{223}	x_{224}	x_{225}	x_{226}	x_{227}	x_{228}	x_{229}	x_{230}	x_{231}	x_{232}	x_{233}	x_{234}	x_{235}	x_{236}	x_{237}	x_{238}	x_{239}	x_{240}	x_{241}	x_{242}	x_{243}	x_{244}	x_{245}	x_{246}	x_{247}	x_{248}	x_{249}	x_{250}	x_{251}	x_{252}	x_{253}	x_{254}	x_{255}	x_{256}	x_{257}	x_{258}	x_{259}	x_{260}	x_{261}	x_{262}	x_{263}	x_{264}	x_{265}	x_{266}	x_{267}	x_{268}	x_{269}	x_{270}	x_{271}	x_{272}	x_{273}	x_{274}	x_{275}	x_{276}	x_{277}	x_{278}	x_{279}	x_{280}	x_{281}	x_{282}	x_{283}	x_{284}	x_{285}	x_{286}	x_{287}	x_{288}	x_{289}	x_{290}	x_{291}	x_{292}	x_{293}	x_{294}	x_{295}	x_{296}	x_{297}	x_{298}	x_{299}	x_{300}	x_{301}	x_{302}	x_{303}	x_{304}	x_{305}	x_{306}	x_{307}	x_{308}	x_{309}	x_{310}	x_{311}	x_{312}	x_{313}	x_{314}	x_{315}	x_{316}	x_{317}	x_{318}	x_{319}	x_{320}	x_{321}	x_{322}	x_{323}	x_{324}	x_{325}	x_{326}	x_{327}	x_{328}	x_{329}	x_{330}	x_{331}	x_{332}	x_{333}	x_{334}	x_{335}	x_{336}	x_{337}	x_{338}	x_{339}	x_{340}	x_{341}	x_{342}	x_{343}	x_{344}	x_{345}	x_{346}	x_{347}	x_{348}	x_{349}	x_{350}	x_{351}	x_{352}	x_{353}	x_{354}	x_{355}	x_{356}	x_{357}	x_{358}	x_{359}	x_{360}	x_{361}	x_{362}	x_{363}	x_{364}	x_{365}	x_{366}	x_{367}	x_{368}	x_{369}	x_{370}	x_{371}	x_{372}	x_{373}	x_{374}	x_{375}	x_{376}	x_{377}	x_{378}	x_{379}	x_{380}	x_{381}	x_{382}	x_{383}	x_{384}	x_{385}	x_{386}	x_{387}	x_{388}	x_{389}	x_{390}	x_{391}	x_{392}	x_{393}	x_{394}	x_{395}	x_{396}	x_{397}	x_{398}	x_{399}	x_{400}	x_{401}	x_{402}	x_{403}	x_{404}	x_{405}	x_{406}	x_{407}	x_{408}	x_{409}	x_{410}	x_{411}	x_{412}	x_{413}	x_{414}	x_{415}	x_{416}	x_{417}	x_{418}	x_{419}	x_{420}	x_{421}	x_{422}	x_{423}	x_{424}	x_{425}	x_{426}	x_{427}	x_{428}	x_{429}	x_{430}	x_{431}	x_{432}	x_{433}	x_{434}	x_{435}	x_{436}	x_{437}	x_{438}	x_{439}	x_{440}	x_{441}	x_{442}	x_{443}	x_{444}	x_{445}	x_{446}	x_{447}	x_{448}	x_{449}	x_{450}	x_{451}	x_{452}	x_{453}	x_{454}	x_{455}	x_{456}	x_{457}	x_{458}	x_{459}	x_{460}	x_{461}	x_{462}	x_{463}	x_{464}	x_{465}	x_{466}	x_{467}	x_{468}	x_{469}	x_{470}	x_{471}	x_{472}	x_{473}	x_{474}	x_{475}	x_{476}	x_{477}	x_{478}	x_{479}	x_{480}	x_{481}	x_{482}	x_{483}	x_{484}	x_{485}	x_{486}	x_{487}	x_{488}	x_{489}	x_{490}	x_{491}	x_{492}	x_{493}	x_{494}	x_{495}	x_{496}	x_{497}	x_{498}	x_{499}	x_{500}	x_{501}	x_{502}	x_{503}	x_{504}	x_{505}	x_{506}	x_{507}	x_{508}	x_{509}	x_{510}	x_{511}	x_{512}	x_{513}	x_{514}	x_{515}	x_{516}	x_{517}	x_{518}	x_{519}	x_{520}	x_{521}	x_{522}	x_{523}	x_{524}	x_{525}	x_{526}	x_{527}	x_{528}	x_{529}	x_{530}	x_{531}	x_{532}	x_{533}	x_{534}	x_{535}	x_{536}	x_{537}	x_{538}	x_{539}	x_{540}	x_{541}	x_{542}	x_{543}	x_{544}	x_{545}	x_{546}	x_{547}	x_{548}	x_{549}	x_{550}	x_{551}	x_{552}	x_{553}	x_{554}	x_{555}	x_{556}	x_{557}	x_{558}	x_{559}	x_{560}	x_{561}	x_{562}	x_{563}	x_{564}	x_{565}	x_{566}	x_{567}	x_{568}	x_{569}	x_{570}	x_{571}	x_{572}	x_{573}	x_{574}	x_{575}	x_{576}	x_{577}	x_{578}	x_{579}	x_{580}	x_{581}	x_{582}	x_{583}	x_{584}	x_{585}	x_{586}	x_{587}	x_{588}	x_{589}	x_{590}	x_{591}	x_{592}	x_{593}	x_{594}	x_{595}	x_{596}	x_{597}	x_{598}	x_{599}	x_{600}	x_{601}	x_{602}	x_{603}	x_{604}	x_{605}	x_{606}	x_{607}	x_{608}	x_{609}	x_{610}	x_{611}	x_{612}	x_{613}	x_{614}	x_{615}	x_{616}	x_{617}	x_{618}	x_{619}	x_{620}	x_{621}	x_{622}	x_{623}	x_{624}	x_{625}	x_{626}	x_{627}	x_{628}	x_{629}	x_{630}	x_{631}	x_{632}	x_{633}	x_{634}	x_{635}	x_{636}	x_{637}	x_{638}	x_{639}	x_{640}	x_{641}	x_{642}	x_{643}	x_{644}	x_{645}	x_{646}	x_{647}	x_{648}	x_{649}	x_{650}	x_{651}	x_{652}	x_{653}	x_{654}	x_{655}	x_{656}	x_{657}	x_{658}	x_{659}	x_{660}	x_{661}	x_{662}	x_{663}	x_{664}	x_{665}	x_{666}	x_{667}	x_{668}	x_{669}	x_{670}	x_{671}	x_{672}	x_{673}	x_{674}	x_{675}	x_{676}	x_{677}	x_{678}	x_{679}	x_{680}	x_{681}	x_{682}	x_{683}	x_{684}	x_{685}	x_{686}	x_{687}	x_{688}	x_{689}	x_{690}	x_{691}	x_{692}	x_{693}	x_{694}	x_{695}	x_{696}	x_{697}	x_{698}	x_{699}	x_{700}	x_{701}	x_{702}	x_{703}	x_{704}	x_{705}	x_{706}	x_{707}	x_{708}	x_{709}	x_{710}	x_{711}	x_{712}	x_{713}	x_{714}	x_{715}	x_{716}	x_{717}	x_{718}	x_{719}	x_{720}	x_{721}	x_{722}	x_{723}	x_{724}	x_{725}	x_{726}	x_{727}	x_{728}	x_{729}	x_{730}	x_{731}	x_{732}	x_{733}	x_{734}	x_{735}	x_{736}	x_{737}	x_{738}	x_{739}	x_{740}	x_{741}	x_{742}	x_{743}	x_{744}	x_{745}	x_{746}	x_{747}	x_{748}	x_{749}	x_{750}	x_{751}	x_{752}	x_{753}	x_{754}	x_{755}	x_{756}	x_{757}	x_{758}	x_{759}	x_{760}	x_{761}	x_{762}	x_{763}	x_{764}	x_{765}	x_{766}	x_{767}	x_{768}	x_{769}	x_{770}	x_{771}	x_{772}	x_{773}	x_{774}	x_{775}	x_{776}	x_{777}	x_{778}	x_{779}	x_{780}	x_{781}	x_{782}	x_{783}	x_{784}	x_{785}	x_{786}	x_{787}	x_{788}	x_{789}	x_{790}	x_{791}	x_{792}	x_{793}	x_{794}	x_{795}	x_{796}	x_{797}	x_{798}	x_{799}	x_{800}	x_{801}	x_{802}	x_{803}	x_{804}	x_{805}	x_{806}	x_{807}	x_{808}	x_{809}	x_{810}	x_{811}	x_{812}	x_{813}	x_{814}	x_{815}	x_{816}	x_{817}	x_{818}	x_{819}	x_{820}	x_{821}	x_{822}	x_{823}	x_{824}	x_{825}	x_{826}	x_{827}	x_{828}	x_{829}	x_{830}	x_{831}	x_{832}	x_{833}	x_{834}	x_{835}	x_{836}	x_{837}	x_{838}	x_{839}	x_{840}	x_{841}	x_{842}	x_{843}	x_{844}	x_{845}	x_{846}	x_{847}	x_{848}	x_{849}	x_{850}	x_{851}	x_{852}	x_{853}	x_{854}	x_{855}	x_{856}	x_{857}	x_{858}	x_{859}	x_{860}	x_{861}	x_{862}	x_{863}	x_{864}	x_{865}	x_{866}	x_{867}	x_{868}	x_{869}	x_{870}	x_{871}	x_{872}	x_{873}	x_{874}	x_{875}	x_{876}	x_{877}	x_{878}	x_{879}	x_{880}	x_{881}	x_{882}	x_{883}	x_{884}	x_{885}	x_{886}	x_{887}	x_{888}	x_{889}	x_{890}	x_{891}	x_{892}	x_{893}	x_{894}	x_{895}	x_{896}	x_{897}	x_{898}	x_{899}	x_{900}	x_{901}	x_{902}	x_{903}	x_{904}	x_{905}	x_{906}	x_{907}	x_{908}	x_{909}	x_{910}	x_{911}	x_{912}	x_{913}	x_{914}	x_{915}	x_{916}	x_{917}	x_{918}	x_{919}	x_{920}	x_{921}	x_{922}	x_{923}	x_{924}	x_{925}	x_{926}	x_{927}	x_{928}	x_{929}	x_{930}	x_{931} </th

→ Weighted K-NN

Till now we have seen simple K-NN, there is also weighted K-NN.



neighbor of $x_q \rightarrow x_1, x_2, x_3, x_4, x_5$

$$x_1, y_1, d_1 = 0.1$$

$$x_2, y_2, d_2 = 0.2$$

$$x_3, y_3, d_3 = 1.0$$

$$x_4, y_4, d_4 = 2.0$$

$$x_5, y_5, d_5 = 4.0$$

As we can see that x_1 is very close to x_q , should not we give more weightage to x_1 , because there are more chance that x_q have same class label as x_1 .

In normal K-NN, we are treating each neighbour equally. But logically we should give more importance to closer neighbour.

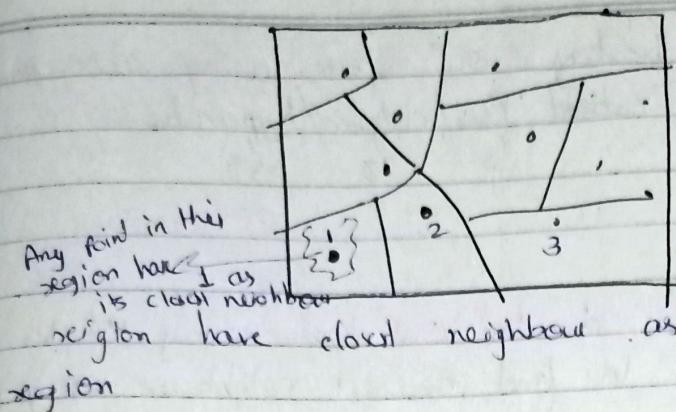
Hence, we can give weights to each point in neighbour. One way is assign weight $w_i = \frac{1}{d_i}$ if distance increase d_i or if distance is less then decrease w_i .

Now, to choose y_q for x_q .

Sum up the ~~positive~~ weights of each category and choose $y_q = \text{category with max weight}$.

There are many other weighing scheme other than $w_i = \frac{1}{d_i}$.

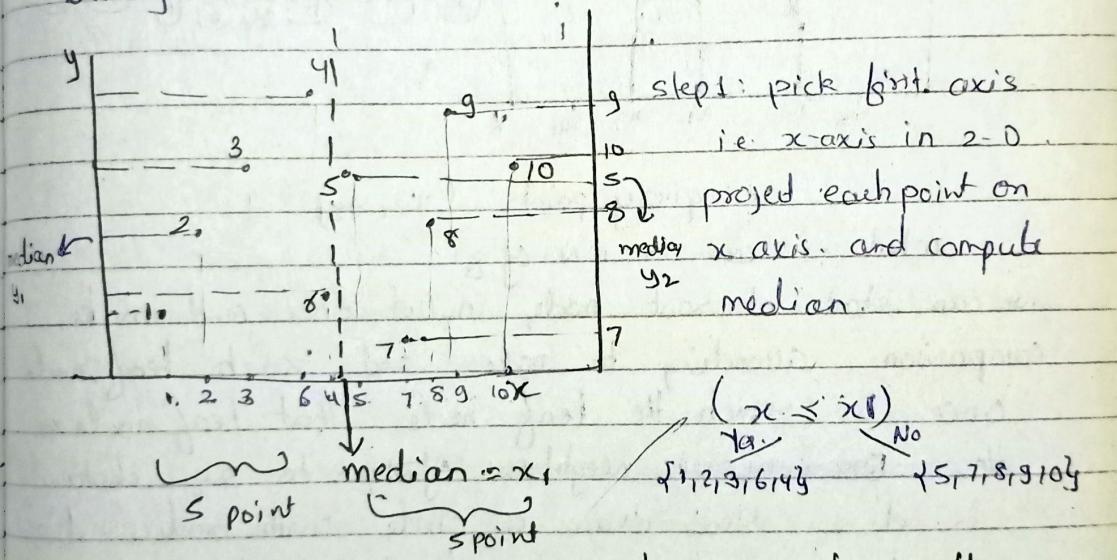
→ Voronoi diagram:- Voronoi diagram is a partitioning of a plane into regions based on distance to points in a specific subset of the plane.



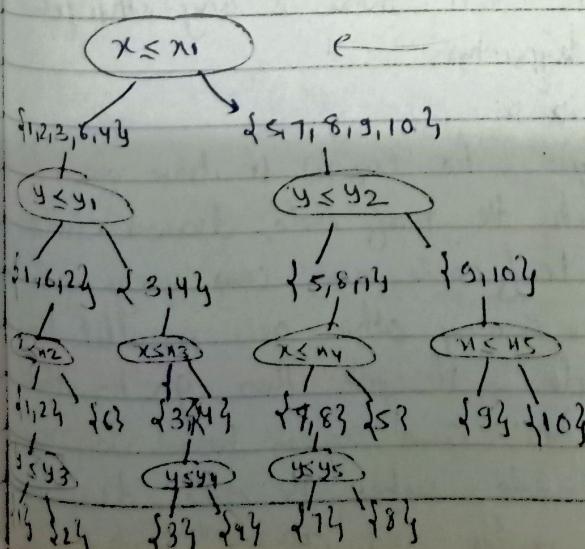
for $K=1$ int-k-nn
Voronoi diagram
divide the plane
in regions s.t
any point in that
region have closest neighbor as
the point of that
region

* K-d-tree:- As we have discussed earlier that Time complexity in $K\text{-NN} = O(nd) \times S: O(nd)$ we also discussed that we have other implementations of $K\text{-NN}$. One of them is Kd-tree which performs better.

Binary search tree can be converted to Kd tree:



step2: Now change the axis and perform some operation on child Node.

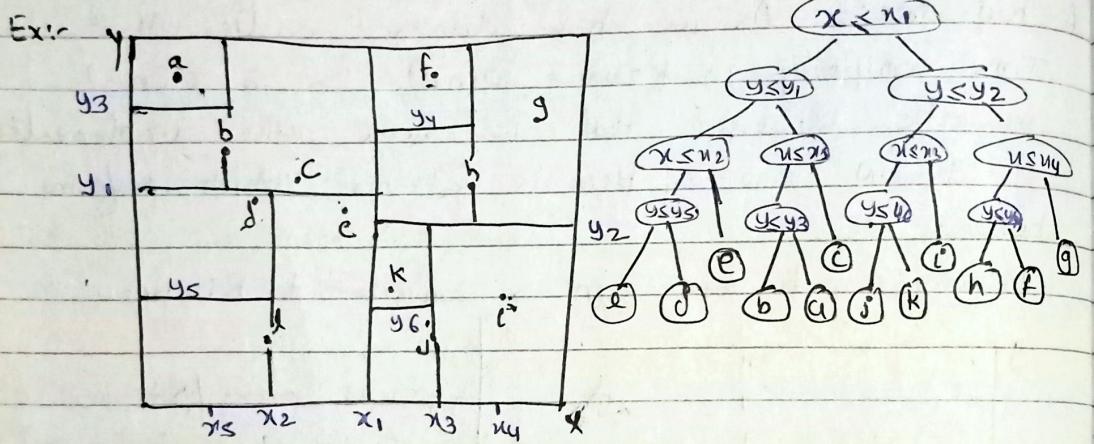


Now again use x-axis.

Actually we are breaking our space using axis-parallel plane into ~~cuboid~~ box/cuboid/hypercube
 2D - 3D - >3D

~~to search~~

Now, that we have divided our space using Kd-tree strategy we learn how to find nearest neighbours using Kd-tree.



x : query point (x_2, y_2)

what is ~~the~~ 1-NN of q .

we can start at root node in Kd-tree and make comparison according to node and search leaf node. Once we reach the leaf node that leaf node is ~~the~~ ~~one~~ nearest neighbour of x_2 . to x_2 distance is d . Now draw a circle with radius = d and centre as (x_2, y_2) . If there is any other point inside the circle / hypershape.

2D - >2D

Now to do it algorithmically:- for (x_2, y_2) if break down the tree once we reach the leaf node, draw a circle. a min dist b/w leaf node and ~~the~~ query node. Now, check if there is any other axis parallel line is intersecting the circle. If yes then go to that axis parallel node in tree and again break down until we reach a node where the circle drawn don't intersect other nodes over parallel axis.

$$T: O(\log(n)) \quad \{ \text{best case}\}$$

$$O(n) \quad \{ \text{worst case}\}$$

In K-NN

$$T: O(K \log n) \quad \{ \text{best case}\}$$

$$\therefore O(K \cdot n) \quad \{ \text{worst case}\}$$

with d-dim.

~~$$T: O((K \cdot d) \log n) \quad \{ \text{best case}\}$$~~
~~$$O(K \cdot d \cdot n) \quad \{ \text{worst case}\}$$~~

* → Limitation of kd-tree:

"In Kd we need to look at adjacent cells. To check if they are nearest."

b In 2D - we need to check 4 adjacent cells,

In 3D - 8

in nd - 2^d

Due to this Time complex becomes:

$$O(2^d * K * \log n) \quad \{ \text{In best case}\}$$

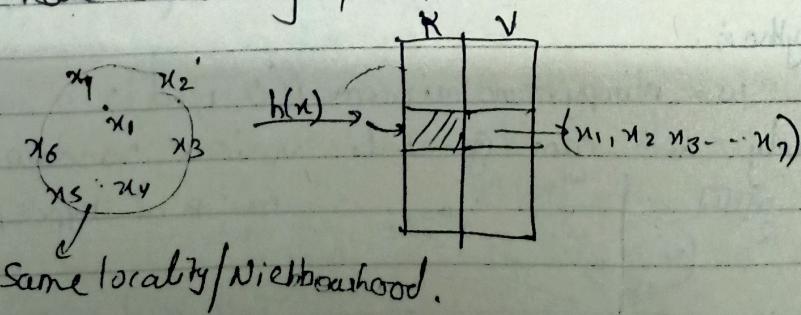
$$O(2^d * K * n) \quad \{ \text{In worst case}\}$$

K-d tree is ok for only small dimensions.
in bigger bidimen Time complex become very large.

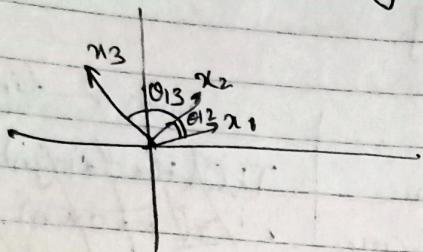
* Hashing & Locality sensitive hashing

LSH uses hashing to find Nearest Neighbours.

if we have bunch of points



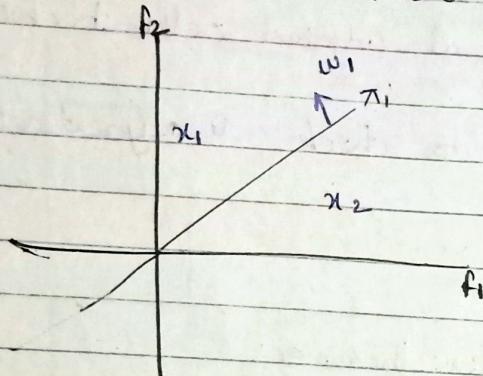
* LSH for cosine-similarity:



$\theta_{12} \approx 0$ hence x_1, x_2 are similar.

K	V
$b(n)$	x_1, x_2

LSH: It is a randomized algorithm.

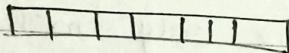


① break the whole space using Random-planes.

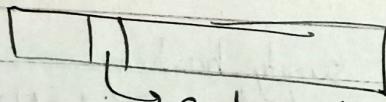
$w^T x_1 = +ve$ if x_1 is in some region
 $w^T x_2 = -ve$

, How we will generate Random-hyperplane.

w is a d-dim vector



If we can generate a random d-dim vector then it can be treated as Random-plane.

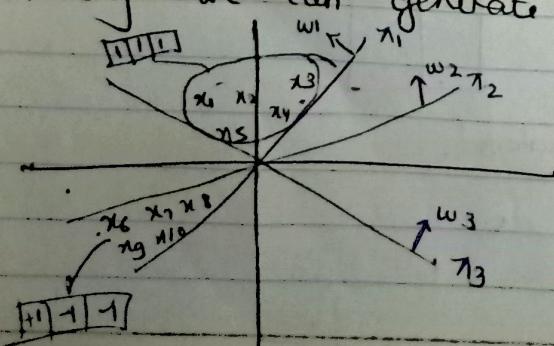


each point is randomly generated from a Normal dist (mean 0, $\sigma=1$)

in python:

$w = \text{numpy.random.normal}(0, 1, d)$;

similarly we can generate various random planes.



$m = \#$ of hyperplanes

Each point is now in a region.

$$x^T w_1 \geq 0 \rightarrow +1$$

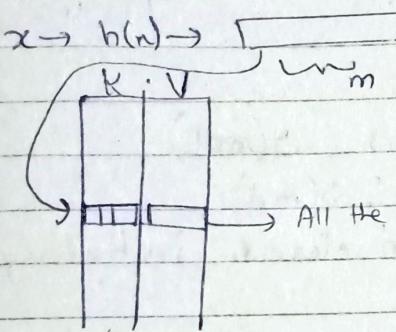
$$x^T w_2 \geq 0 \rightarrow +1$$

$$x^T w_3 \geq 0 \rightarrow +1$$

for every slice or region we create a m dim vector because there are m regions.

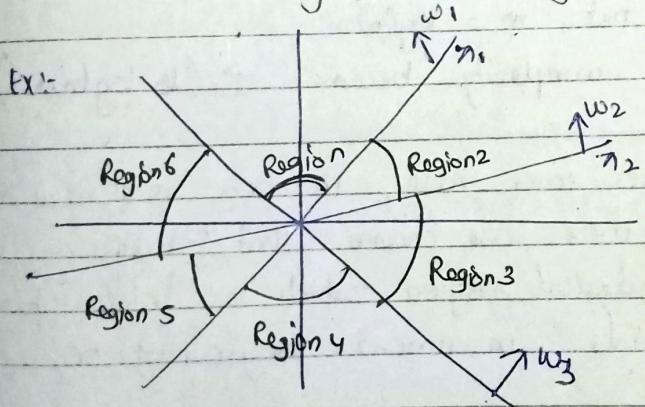
$$h(x) = \begin{bmatrix} \cdot & 1 & 2 & 3 \end{bmatrix} \rightarrow \text{sign}(w_3^T x) \quad \text{Sign}(w_2^T x) \quad \text{Sign}(w_1^T x)$$

Note



All the points ~~exist~~ in this region.

Key is not a single value, it is m-dim vector.



Each region have a m-dim vector

	w ₁	w ₂	w ₃
Region 1	+1	+1	+1
2	-1	+1	+1
3	-1	-1	+1
4	-1	-1	-1
5	+1	-1	-1
6	+1	+1	-1

These will act as key in hash table and all the points in the region will act as value in hash table.

How much Time & space to construct the hash table.

S : O(nd)

n points of d dim

* Every point has to be multiplied with region here
 $\sim O(mdn)$

- Now, imagine we are given a query point x_q .

Now step1: find $h(x_q) = \boxed{1} \downarrow \text{sign of } (w_1^T x_q), \text{ sign of } (w_2^T x_q)$

use $h(x_q)$ as key in hash-table and get the
get a list of data points in that region.

suppose we get n points from that region.

then we can select k nearest neighbor from these n
points.

∴ Time of one query $\Rightarrow O(\text{md})$

$x_q \rightarrow h(q) : O(\text{md})$

If there are n elements in that region,
 $O(n)$

typically we put $m = \log(n)$

hence Time complexity becomes $O(\log(n) \cdot md)$

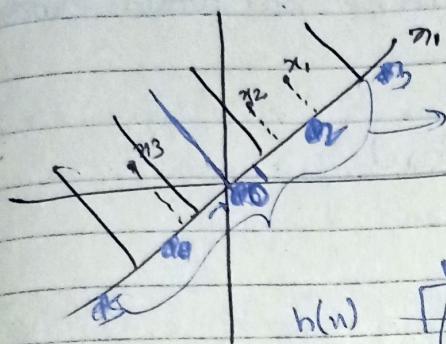
Note: There are corner cases here, if x_q is at near boundary of region then there are chances that there are some points in adjacent regions which are closer to x_q than the points in current region of x_q .

* We can miss a NN in cosine as our distance measure.
To tackle this problem, we need to create multiple hash tables with different regions for each hash table and then we need to traverse through the list of each hash table and get a set of all possible NN and get k -NN.

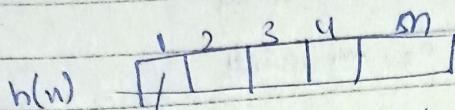
So final T: $O(1 \cdot md)$

No. of hash tables. \downarrow # of hyperplanes. \downarrow dimension/features.

→ LSH for Euclidean distance.



breaking the hyperplane in various parts and check how many points are projecting in that part.



part in which the projection of point lies in π_1 hyperplane.

$$h(x_1) = [2 | 1 | 1]$$

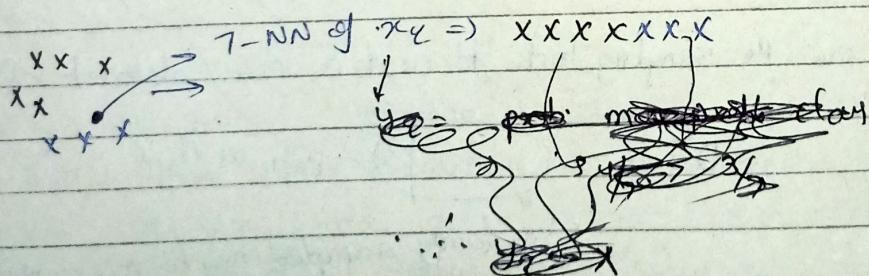
$$h(x_2) = [2 | 1 | 1]$$

$$h(x_3) = [4 | 1 | 1]$$

This is the only diff b/w LSH for Euclidean & LSH for cosine.

* LSH is not perfect.

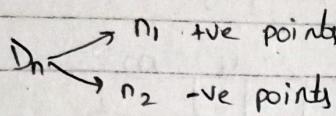
k-NN: probabilistic class label.



$y_2 = \text{circle with probability } p$ is +ve class with prob.
-ve " " " " 1 - p prob.

In this chapter we will learn how to apply classification with some real world cases. because in real world data may be in various forms.

→ Imbalanced vs balanced dataset



If $n_1 \approx n_2$ then D_n is balanced dataset

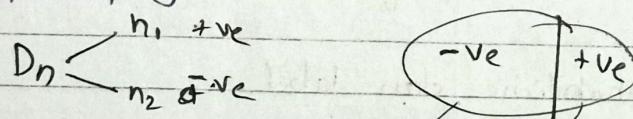
If $n_1 \ll n_2$ or $n_2 \ll n_1$, D_n is imbalanced dataset

So, how do we deal with imbalanced datasets.

Suppose we are using k-NN, we can get biased result. Because there are more chances that all the data points are near Neigh to from dominate class. and only few data points from other class.

So, How to deal with imbalanced data set

① Under sampling



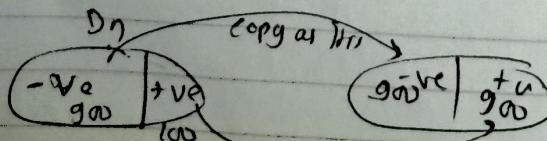
use the sampling trick to create a new dataset D'



randomly sampled neg in pair = positive points
Now both becomes equal.

But there is one problem with Under Sampling : we are not using all the data. we build our model with only small set of data. hence model may not work properly.

② Over sampling



repeat the +ve points until +ve become \approx -ve.

other way of doing oversampling is artificial/synthetic point we create new points using already points.

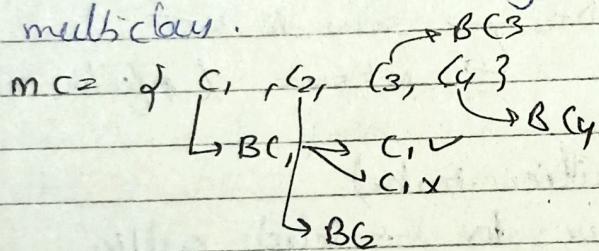
- Q) class weights? — Let we have 100 +ve, 900 -ve so we assign weights to the -ve class +ve points have weight 9, -ve point have 1

→ Multi-class vs binary classify.

All the classifications algo don't work on multiclass, like logistic regres works only on binary class

- Q) can we convert a multi-class problem to binary class:

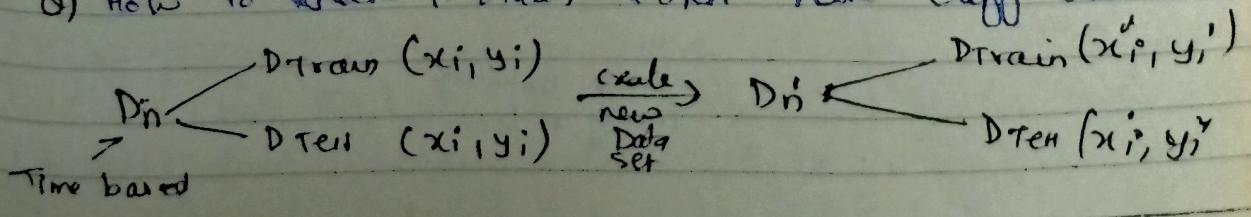
Ans:- for each class in multiclass problem we have to devise one binary classifier for each class of multiclass.



→ Test & Train data set diff =

When we are using ~~Time~~ based division of Data into Train D_{train} & Test D_{test} then there might be changes that Distribution of data may get changed and hence model perform bad. In Dynamic dataset where new point gets added time to time - we need to check if D_{train} & D_{test} are following the previous distribution or not.

- Q) How to check if D_{train} & D_{test} have diff distribution



for D training

$$(x'_i, y'_i) \Rightarrow x'_i = \text{concat}(x_i, y_i)$$
$$y'_i = 1$$

for D test

$$(x'_i, y'_i) \Rightarrow x'_i = \text{concat}(x_i, y_i)$$
$$y'_i = 0$$

Now build a Binary classifier on D'

If Dtrain & Dtest have different distribution then
Binary classifier have high accuracy, as it can separate well.

If Dtrain & Dtest have same distribution then BC have
low accuracy.

Dtrain, Dtest \rightarrow same dis ✓

\rightarrow not same (feature are not stable).

LOF (Local outlier factor)

It is a technique to detect outliers

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

$N(A)$ = Neighbourhood of A.

Reachability-distance (x_i, x_j) = $\max(K\text{-distance}(x_j), \text{dist}(x_i, x_j))$

Local reachability density : $lrd(x_i) =$

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

average reachability
dist of x_i from its Neighbours.

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

average ld of pt in the neighbourhood
of x_i

If $\text{LOF}(x_i)$ is large, if $\text{ld}(x_i) \downarrow$ or ld is large.

\downarrow it is small when

there are small points
around x_i , i.e. density
of x_i is small

\hookrightarrow this is large when
Neighbours of x_i are
darker.

so, simply: if neighbour average density is more but ld
density around x_i is small. Then x_i is outlier.

\$

If $\text{LOF}(x_i) \uparrow$ x_i is outlier

forward.

→ feature importance and ~~feature~~ selection.

Not all the features are important for classification
or regression. If we somehow sort the features
based on their importance. Then it becomes very
useful. Feature importance is useful to understand
the model better.

feature importance is a numerical score to each feature
+ K-NN doesn't provide feature importance (but
we can make changes in it to get feature import.)

forward feature selection is a technique to get importance
of feature importance.

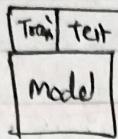
If we have very large dataset/features say 1000 it
becomes important to decrease the dimensionality.

Even though we can use PCA, t-SNE for dimensionality
reduction. But they don't care about classification task. They
just care about distance. Hence if we somehow know
the important features say 100 out of 1000 we can
discard 900 features.

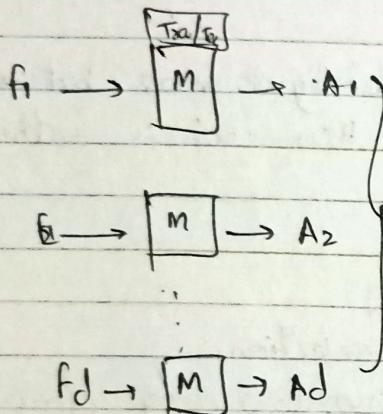
forward feature selection is used for the

$D_n \rightarrow D_{train}$
 $D_n \rightarrow D_{test}$

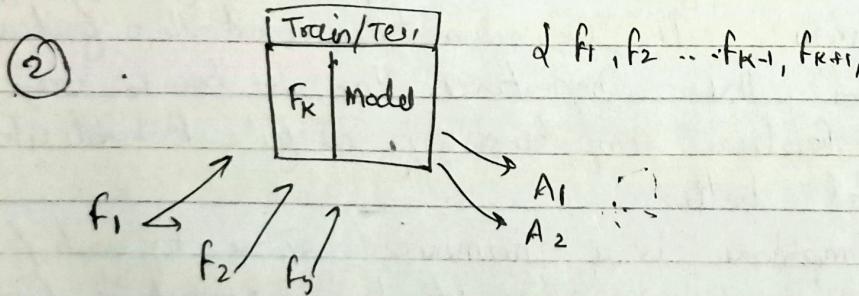
$$\textcircled{1} \quad \{ f_1, f_2, f_3, \dots, f_d \}$$



take only first feature and take Train & test data
only come to that feature. Now, suppose we get
some Accura - A_1 ,



Say f_k has high Accura
Then f_k is the most important feat.



Same step is repeated

Note:- One can argue that why we didn't selected Top 100 or K features ~~in~~ in the very first iteration because we are choosing only max accuracy feature, why not choose 2 or 3 or top k feature from 1st iter.
The Ans:- There may be chances that two feature are similar hence their accuracy is high and similar so we may end up selecting similar features.

In forward selection, we select one feature in each iteration. A similar concept is Backward set. In this we remove the least important feature at each iteration.

limitation - Time complexity ↑

Handling categorical & ordinal features.

Suppose we are given a data set of people where we have various features

Wt	hair	Country	height
1	Black	USA	160
2	Brown	UK	170
3	Red	Canada	155
4	Golden	Australia	180
5	Grey	Germany	165

→ weight is in No. so keep it as-is

→ {Black, Brown, Red, Golden, Grey}

If we can assign a number to them like {Black, Brown, Red, Golden} → {1, 2, 3, 4, 5}

Then there is a problem because numbers comes with an inherent orderliness i.e. $1 > 2, 5 > 4$, etc. But we can't say Black > Brown etc. so we cannot directly assign numbers.

→ One-hot encoding: a hair color: {Black, Brown, Red, Golden, Grey} → {1, 0, 0, 0, 0}

so creates 5-d vector

1	0	0	0	0
Black	Brown	Red	Golden	Grey

Let $x_1 = \text{Brown}$ Then its feature encoding → {0, 1, 0, 0, 0}

So one-hot encoding creates a binary vector of size of no. of unique values in that feature.

Hair color is a categorical feature.

Problem with one-hot encoding is sometimes there are lot of unique values in feature hence feature vector becomes larger.

(Ex.: country feature has around 200 values. here we get a 200 size vector.)

If the no. of distinct values for a categorical feature is high then one-hot encoding :: it can create sparse and large vectors because only one value is 1 and others are 0.

④ Solution for this is

$f_2 = \text{country}$

↳ 2D values

one hot-encoding

20V

idea: Replace "country col" by the average height of
first people from that country.

General idea:- Replace the categorical feature values to the average value of class label in each concept category.

But this don't work in every problem.

⑥ There is another way:- Using domain knowledge.

Ex:- Country → distance from a chosen country.
(categorical) (numerical)

→ ordinal features:

features: how do you rate Indian food.

V. good
good
avg
bad
V. bad

{ Eg. It is categorical, but one thing is
It have ordering. V-good > good
Sang > bad.

hence we can simply assign item no.

Handling missing values:-

① imputation

- ① mean replacement
- ② Median "
- ③ most freq value "

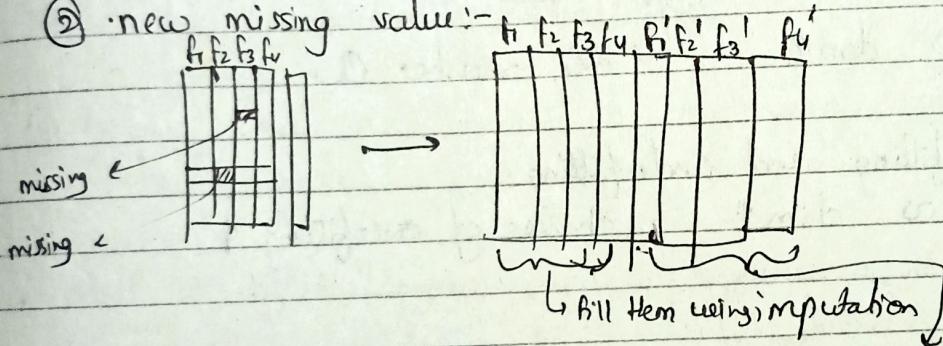
In classification we have to take care.

say $x_i : f_3$ is missing and $y_i = 1$

then replace $x_i : f_3$ with mean/median/mode of all the value have class label y_i

• we have to find mean/median/mode from same class in classification problem.

② new missing value:-



Put 0 if

missing else 0.

③ model based imputation

we can make a model to predict the value of missing values.

* Curse of dimensionality

it talks about the problems occur due to higher dimensionality.

① As dim \uparrow total no. of datapoints to perform good classification increase exponentially

• Hughes phenomenon:- $n \rightarrow$ size of dataset
if n is fixed, performance \downarrow as dim \uparrow

Higher phenomena:- on a fixed data set size, performance decrease as dim \uparrow .

② Distance funcⁿ (specially euclidean distance)
intuition of dist in 3D is not valid in higher dim.

$$\text{dist}_{\min}(x_i) = \min_{x_j \neq x_i} (\text{dist}(x_j, x_i))$$

$$\text{dist}_{\max}(x_i) = \max_{x_j \neq x_i} (\text{dist}(x_j, x_i))$$

$$\underbrace{\text{dist}_{\max}(x_i) - \text{dist}_{\min}(x_i)}_{\text{dist}_{\text{mean}}(x_i)} > 0 \quad \left. \right\} \text{for } 1D, 2D, 3D$$

as dim \uparrow , this reaches 0.

③ overfitting and underfitting:

as dim \uparrow , chances of overfitting \uparrow .

* Bias-variance tradeoff :-

mathematical basis to understand overfitting & underfitting

* Bias-variance tradeoff :-

generalization error: error that is made on future unseen
(future data error) error

$$\boxed{\text{generalization error} = \text{Bias}^2 + \text{Variance} + \text{irreducible error.}}$$

① irreducible error + error that cannot be reduced further.

② Bias error :- comes due to simplifying assumption.
high bias = underfitting.

③ Variance :- how much model changes as training data changes.

For example at k=1 in k-NN, model changes rapidly as train data changes.

high variance \Rightarrow overfitting.

so, our final goal is low generalisation error in order to attain this we should have \downarrow Bias \downarrow Variance.
i.e. \downarrow overfitting & underfitting.

\rightarrow Best & worst case of K-NN:-

- ① when dim is small (< 10) \rightarrow K-NN is good model.
- ② when dim \uparrow we get curse of dimensionality, run time increase,
- ③ KNN should not be used in low-latency applica. where we need fast result.
- ④ find the right distance measure:- K-NN is good