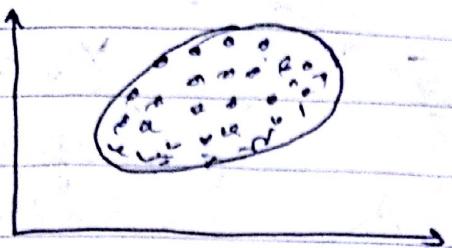


17-OCT-2014

classmate

Date _____
Page _____

PRINCIPAL COMPONENT ANALYSIS (PCA)



1. Best 0-dim representation.

$$\text{Minimising } J(x_0) = \sum_{n=1}^N \|x_0 - x_n\|^2$$

$$= \sum_{n=1}^N \| (x_0 - m) - (x_n - m) \|^2$$

$$= \sum_{k=1}^n \|x_0 - m\|^2 - 2(x_0 - m) \sum_{k=1}^n (x_k - m)$$

$$+ \sum_{n=1}^N \|x_n - m\|^2$$

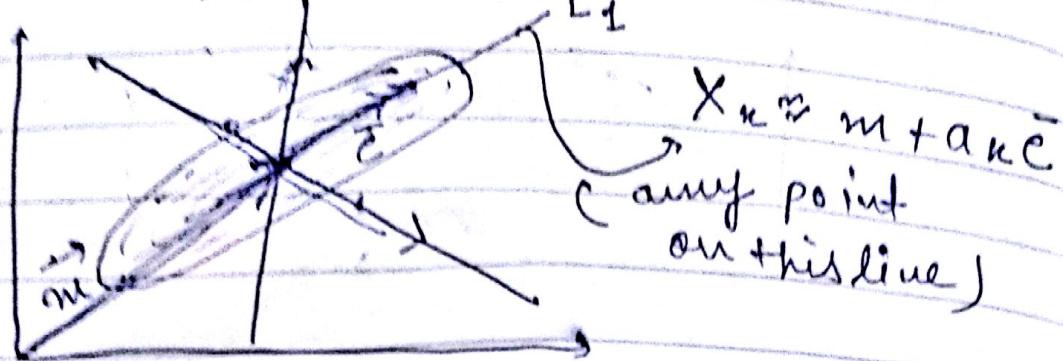
→ should '0' add up to 0

find an ' x_0 ' that will minimise above expression.

Here 'm' is assumed to be the mean
So in order to minimise above we can only change x_0 and it is clearly seen that

$\underline{x_0 = m}$ is the mini. condition.

Q. Best 1-D representation:-



$$x_m = a_n \vec{e} + m$$

d dimensional
scalar

(same \vec{e})

In above we know (\vec{e} & m) so for every point we will store ' a_k ' (scalar).

$$J(a_1, \dots, a_n, \vec{e}) = \sum_{k=1}^N \|c_{\text{mtake}} - x_k\|^2$$

measured

$$= \sum_{k=1}^N \|a_n \vec{e} - (x_k - m)\|^2$$

\vec{e} = unit vector

$$(1) \quad J(\cdot) = \sum_{k=1}^N \|a_n \vec{e}\|^2 - 2 \sum_{n=1}^N a_n \vec{e}^T (x_k - m) + \sum_{n=1}^N \|x_k - m\|^2$$

Now we want to minimise above by finding a_n & \vec{e} .

So 2 questions

(1) - what is the best \vec{c} .

(2) - what are the best axis.

Let us find best axis:-

$$\frac{\partial J}{\partial a_k} = 0.$$

[Here we are differentiating only w.r.t a_k [$k \rightarrow$ particular]]

So

$$2a_k - 2e^T(x_n - m) = 0$$

So

$$\overrightarrow{a_k} = e^T(\overrightarrow{x_n} - \overrightarrow{m}) = (\overrightarrow{x_n} - \overrightarrow{m})^T e$$

→ simply the linear projection of $(\overrightarrow{x_n} - \overrightarrow{m})$ over \overrightarrow{e} . [Dot product]

Now we will find Best \vec{e} :-

From (1)* The term $e^T(x_n - m)$ will become a_n for minimality (ideality)

So

$$J(\cdot) = - \sum_{n=1}^N a_n^2 + \sum_{n=1}^N \|x_n - m\|^2$$

Since we have to find $\frac{\partial J}{\partial e}$ expand a_k .

$$J(\cdot) = - \sum_{n=1}^N e^T (x_n - m)(x_n - m)^T e$$

$$\left| - \sum_{k=1}^N \|x_k - m\|^2 \right| \rightarrow \text{when } m \text{ will Diff. This will go to 0.}$$

$$J(\cdot) = -e^T \left[\sum_{n=1}^N (x_n - m)(x_n - m)^T \right] e$$

↓
Scatter Matrix

{ scaled version of
Covariance matrix }

$$\text{Max } J(e) = -e^T S e$$

$$\# \text{Maximise } J(e) = e^T S e$$

We are having here a constraint

Maximization problem

subject to constraint $\|e\|_2$

$$e^T e = 1$$

So

$$J(e) = e^T S e + \lambda e^T e$$

$$\frac{\partial J}{\partial e} = S e - \lambda e = 0$$

$$\rightarrow S e = \lambda e$$

So perfect 'e' would be a eigen vector, but which one

so that one which corresponds to Maximum eigen value.

(Best \vec{e}' = Eigen vector of C' with Max. eigen value)

→ Projection That will minimize the error, $J()$ will be giving the same thing (solu) if we have found the \vec{e}' , it's that Maximises the Inter class scattering.

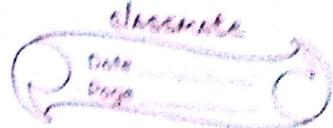
→ So our Low dimensional Feature vectors can be found by

$$E = \begin{bmatrix} | & | & | \\ e_1 & e_2 & \dots & e_{10} \\ | & | & \dots & | \end{bmatrix} \rightarrow \begin{array}{l} \text{Take} \\ \text{Top} \\ \text{such} \\ \text{eigen} \\ \text{vectors} \end{array}$$

$$\underline{x}_k = E^T \begin{bmatrix} x_{n-M} \\ \vdots \\ x_1 \end{bmatrix}_{d \times 1} \quad \text{HDFV}$$

So how many Eigen vectors should we choose, I may $\sum \lambda = 95\%$ of Total

So Take λ of those



Remember

$$\overline{X_{\text{new } k}} = E^T (X_{\text{old } k} - \bar{m})$$

Eigen faces:- (Samples) - $(N^2 \times 1)$ vector
 $N \times N$ image

Face :- $f_1, f_2, f_3, \dots, f_M$

$$\text{Mean Face} = \frac{1}{M} \sum_{i=1}^M f_i = \Psi \quad (\text{average face})$$

$$\overline{\phi_i = f_i - \Psi} = \underbrace{(x_n - \bar{m})}_{\text{of a sample}}$$

So let us define

$$A = [\phi_1 \ \phi_2 \ \dots \ \phi_M] \quad [\text{Data matrix of mean faces}]$$

So the Scatter Matrix

$$S = \sum_{n=1}^N \phi_i \phi_i^T = \sum (x_n - \bar{m})(x_n - \bar{m})^T$$

$$S = A A^T$$

Now find eigen vectors of S

$$Se = \lambda e$$

Now The problem here is

$$S = AAT = [10,000 \times 10,000]$$

matrix

for a 100×100

Suppose ~~suppose~~ image

$$\begin{matrix} S \\ 10K \times 10K \end{matrix} = \begin{matrix} A \\ 10K \times 100 \end{matrix} \times \begin{matrix} AT \\ 100 \times 10K \end{matrix}$$

So we have in S the rank of the matrix S is 100

{ rest all eigenvalues are 0 }

and 100×100 The number of samples. { images }

→ Here we can deal with

$$B = S_1 = A^T A \quad [100 \times 100]$$

so let its eigen vectors be

$$\underbrace{(A^T A)}_{100 \times 100} v = \lambda v \quad \text{--- (2)}$$

Premultiply above with λ

$$[A \ A^T] A v = \lambda A v$$

$$S \boxed{Av} = \lambda \boxed{Av}$$

$\hookrightarrow 10k \times 1$

vector

So our original $\boxed{e = Av}$

Face Recognition becomes Too easy
(solvable) because of above 3 lines

$$\boxed{e = Av}$$

{Came got a paper out of it}

→ Cuba Project, The community soln'

Surviving Peak oil
Agriculture

↳ Urban Gardens

↳ Sustainable practices

↳ Land Distributions

↳ Organic Farming

↳ Roof Top Gardens

{Bicycles distributed,

∴ No oil }

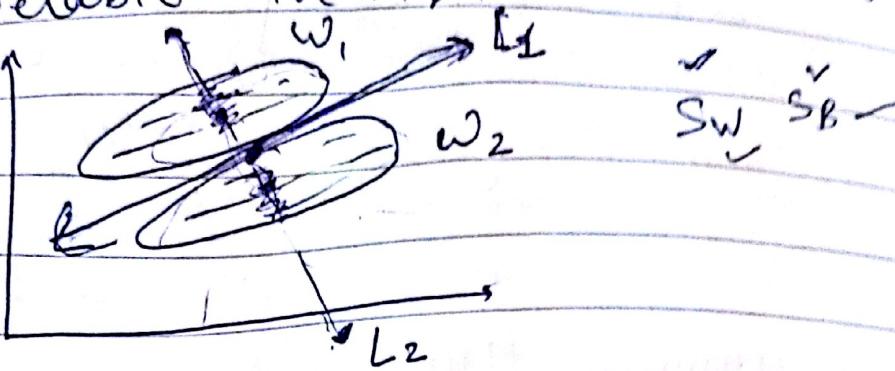
28-OCT-2014

END-SEM-

CLASSMATE

Date _____
Page _____

- If we use PCA to reduce the dimensionality of samples, then they might not be separable in Lower Dim. if they were separable in HD.



However in LD there exist a line (line L_2) which avails the info. of separability just that PCA might not find it

for L_1 , $J = \frac{1}{2} \text{trace}(\mathbf{S}_w) - \mathbf{x}^T \mathbf{x}$ is minimum
 but for L_2 is maximum (higher error), so now our idea would be to reduce to LD in order to have higher separability in LD.

- To have separability in LD we would want

HD	\bar{m}_1	\bar{m}_2	S_1	S_2
LD	\tilde{m}_1	\tilde{m}_2	\tilde{S}_1	\tilde{S}_2

$$\|\tilde{m}_1 - \tilde{m}_2\| \rightarrow \text{maximise}$$

\tilde{S}_1 and \tilde{S}_2 should be as low as possible (in class variance).

$$\frac{S_B}{S_W}$$

$$S_W = \sum S_i$$

classmate
Date _____
Page _____

so our criterion function should look like

$$J(\mathbf{w}) = \frac{\|\mathbf{w}_1 - \mathbf{w}_2\|^2}{(\mathbf{w}_1 + \mathbf{w}_2)^T} \quad [\text{want to maximise it}]$$

S_W = within class scatter (HD)

\tilde{S}_W = " (LD)

$$S_i = \sum_{k=1}^{n_i} (\mathbf{x}_k - \mathbf{m}_i) (\mathbf{x}_k - \mathbf{m}_i)^T$$

$$S_W = \sum_{i=1}^C S_i \quad S_T = S_W + S_B$$

If S_W is the scatter Matrix in HD

Then $\boxed{S_W = \mathbf{W}^T S_W \mathbf{W}}$ is scatt. in LD

where \mathbf{w} (vector in line)

so we want to minimise S_W

When going to HD we can think of $\|\mathbf{m}_1 - \mathbf{m}_2\|$ to be replaced by

scatter of means and we want to maxi. it.

$$S_B = \sum_{i=1}^C n_i (\mathbf{m}_i - \mathbf{m}) (\mathbf{m}_i - \mathbf{m})^T$$

[Between class]

for 2 class $S_B = (m_1 - m_2)(m_1 - m_2)^T$
else $S_B = \sum_{n=1}^N N_n(m_n - \bar{m})(m_n - \bar{m})^T$

 $\hat{S_B} = \|w^T S_B w\|$

we want to Maximise

$$J = \frac{S_B}{S_W}, \text{ or effectively}$$

we want to Maximise Total scatter, $S_T = S_B + S_W$

\rightarrow probably the result of PCA

206 (Bishop)

So our $J = \frac{S_B}{S_W}$

Find $(S_W^{-1} S_B) \rightarrow$ eigen vector
 \rightarrow line L_2 will be given by it.

so above thing is

LDA (Linear discriminant analysis)

or MDA (multi DA)

\hookrightarrow proposed by Fischer

$\therefore = FDA | FLDA | FMDA -$

\Rightarrow So we can apply this to face Recog.
Fisher faces vs eigenfaces.
(works better)

Consider S_B ($d \times d$) Matrix so \rightarrow

$\begin{bmatrix} 0 & - & - & - \\ - & 0 & - & - \\ - & - & 0 & - \\ - & - & - & 0 \end{bmatrix}$ all rows are gR,
so Rank of the mat^{ix}
is '1'.

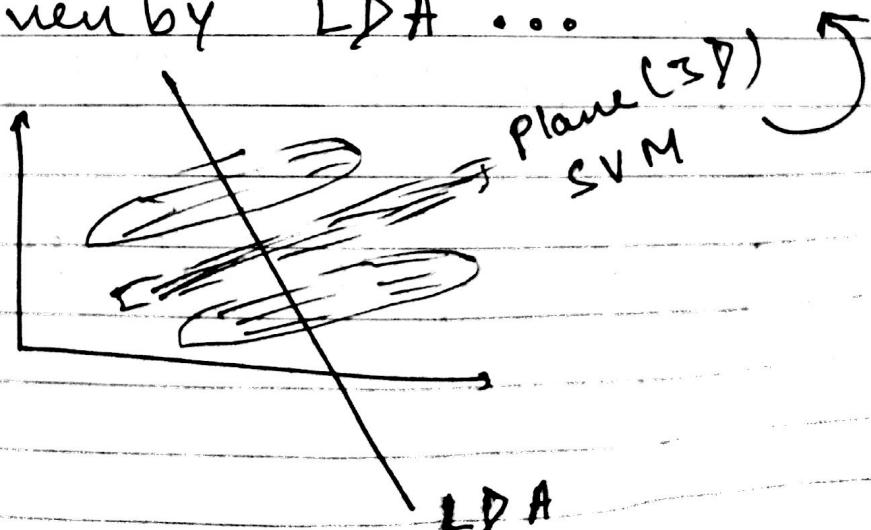
if $A = B \times C$ or $\sum B C$

The Rank of A is not greater than
B or C. ?? (Doubt +)

↳ if we have C classes in all
we can have utmost $C-1$ non-
zero eigenvalues of S_B

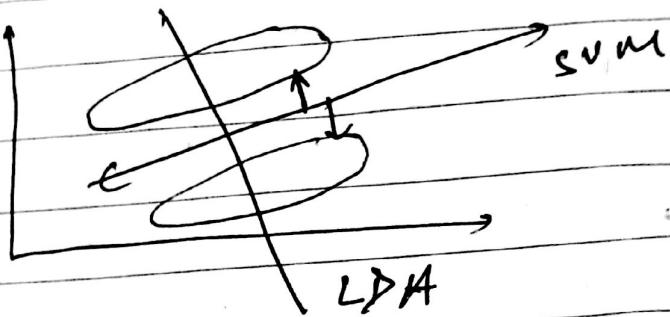
\leftarrow SVM \rightarrow

... we want to find the Hyperplane
that "MAXIMALLY" separates
the sample, that will be a HP that
is perpendicular to line (2D)
given by LDA ...



It is somewhat similar to 'Perceptron Learning' but it finds Hyperplane that maxi. sep. the sample.

So Now we do not want to separate means but to separate the close samples, Maximize the margin



→ Vapnik - Chervonenkis (VC dimension — can we quantify generalisation)
"Huge Breakthrough"

Bound on expected Loss

$$R(\alpha) \leq R_{\text{train}}(\alpha) + \sqrt{\frac{f(n)}{N}}$$

classify

$R \rightarrow$ Risk (error in last)

$\alpha \rightarrow$ a solution

None this is a Hoeffding thing to say
that my future error is bounded
-ed, factor = $\sqrt{f(h)/N}$

$$f(h) = ?$$

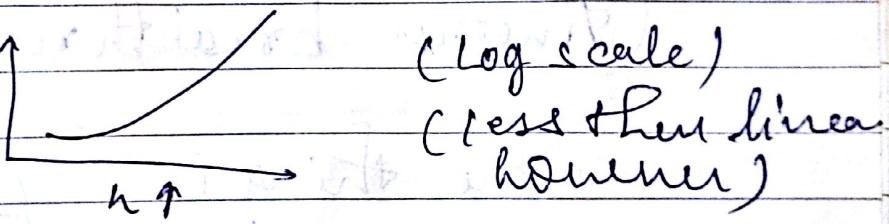
'h' is the VC dimension

$$f(h) = h + h \log(2N) - h \log h - c$$

$$= h \left[1 + \log\left(\frac{2N}{h}\right) \right] - c$$

overall as $h \uparrow f(h) \uparrow$

in fact \uparrow (Log scale)



→ for minimizing $R(x)$ minimize

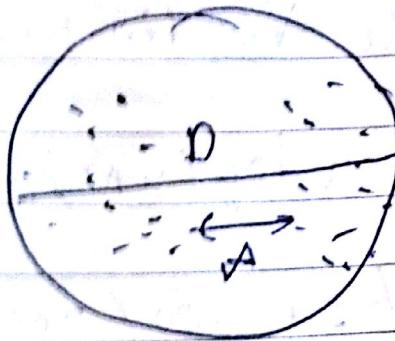
$R_{\text{train}}(x) \rightarrow 0$ (if we can separate perfectly)

So we need to minimize $\sqrt{\frac{f(h)}{N}}$

$\Rightarrow \min(f(h)) \Rightarrow \min. \text{VC dim}(h)$

Marginal length

as $\begin{bmatrix} A & P & h \downarrow \end{bmatrix}$



relative margin
 $= A/D$

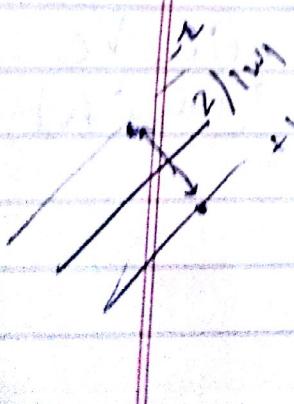
$V(C-D)$, $h \leq \min\{d\}$, $\left[\frac{D^2}{\sqrt{\Delta^2}} \right] + 1$

(Another breakthrough)

so the thing is $\left(\frac{D}{\sqrt{\Delta}} \right) \rightarrow \cancel{\text{---}}$
does not contain 'd'

so for minimizing $\frac{D}{\sqrt{\Delta}} \Rightarrow$ Maximize

$\frac{\sqrt{\Delta}}{D}$, { curse of Higher dimensions has been removed
it does not matter how big 'd' is }



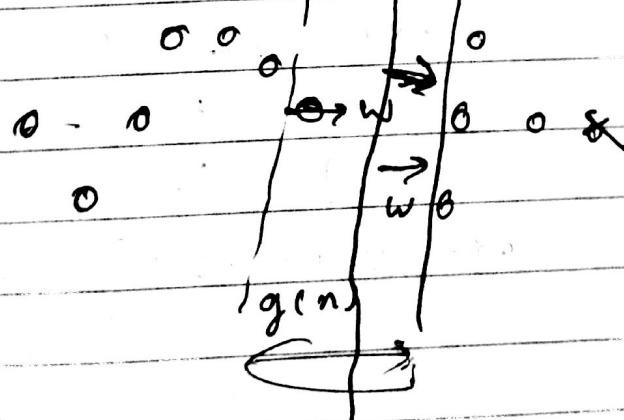
They said this in 1970',
 knew how to use it, nobody
 \rightarrow 1980's
 \rightarrow 1980's \rightarrow (SVM)

\rightarrow so if D is fixed Maximising it
 will minimise 'h'

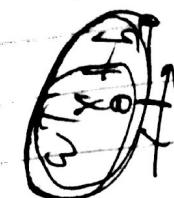
31-10-14

- $w^T x + b = \pm \epsilon$ (margin), but what about
 Higher Dimensions — How to define
 margin?

- fix $g(x) \geq 1$ and minimise $\|w\|$



$$\text{minimise } \|w\| \Rightarrow \|w\|^2 = \frac{1}{2} w^T w$$



$$\|w\|$$

so we have $\min(\|w\|)$ or $\min(w^T w)$
subject to $d_i(w^T x_i + b) \geq 1 + \epsilon_i$

#1 It's a QP Problem
 $\min_w \phi(w) = \frac{1}{2} w^T w$
 subject to: $d_i(w^T x_i + b) - 1 \geq \epsilon_i$

QP solver exists \therefore SVM exist.

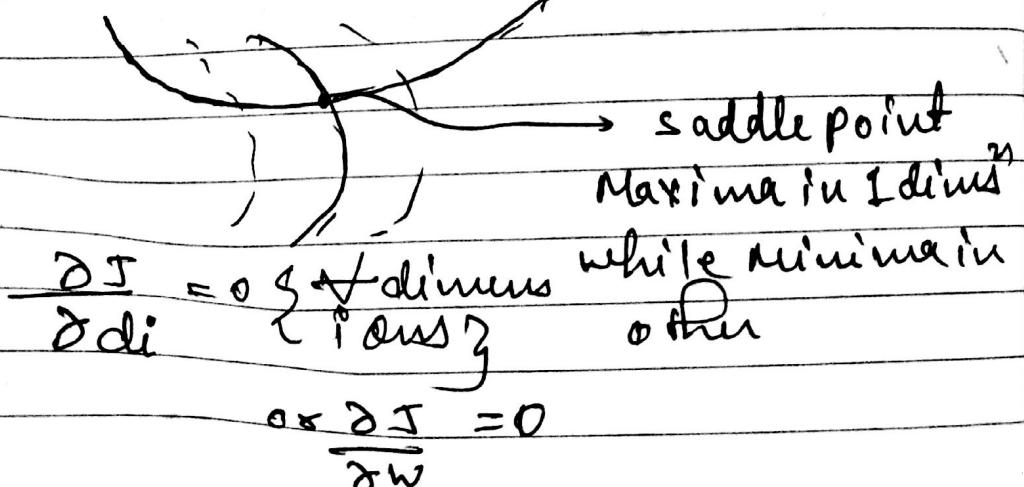
equivalently Lagrangian form -

#2 $J(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i d_i(w^T x_i + b)$
 $\min_w + \sum_{i=1}^N \alpha_i$
 subject to $\alpha_i \geq 0$

we want maximum α_i 's with $w \neq 0$

→ so minimize J with respect to w and b , and max wrt to α] — because of Lagrangian form

J → Horse saddle



{ Max flow and Min cut are duals of each other
 solve one \Rightarrow other } - good when one is easy.



Tough: Primal \rightarrow Easy: Dual

= so find dual of $J(w, b, \alpha)$

At the optimum -

$$① - \frac{\partial J}{\partial w} = 0, \quad ② \frac{\partial J}{\partial b} = 0$$

$$a) - \frac{\partial J}{\partial w} = w - \sum_{i=1}^n \alpha_i d_i x_i$$

$$\text{or } w = \sum_{i=1}^n \alpha_i d_i x_i \rightarrow (1)$$

and $d_i = \pm 1$ so w_0 is basically a LC of Training examples, only question find α 's.

$$b) - \frac{\partial J}{\partial b} = 0, \Rightarrow \sum_{i=1}^n \alpha_i d_i = 0 \rightarrow (2)$$

α_i 's are Lagrange Multiplier.

KKT conditions:-

[Karush - Kuhn - Tucker]

$$(3) = \boxed{\alpha_i [d_i (w_0^T x_i + b_0) - 1] = 0}$$

(for all 'i') — at optimal point

KKT :-

$$\underbrace{\alpha_i}_{(B)} \underbrace{\left[d_i (w_0^T x_i + b_0) - 1 \right]}_{(A)} = 0$$

either (A) = 0 $\rightarrow d_i (w_0^T x_i + b_0) = 1$
∴ They are basically the points on the
boundary $\| = 1$ support vectors

~~(A)~~ or (B) = 0, $\alpha_i = 0$ means The far
away samples from the boundary
have $\alpha_i = 0$.

So for converting to dual one

$$J(w, b, \alpha)$$

$$= \sum \alpha_i + \frac{1}{2} w^T w - w^T \sum \alpha_i d_i x_i$$

$$- b \sum \alpha_i d_i$$

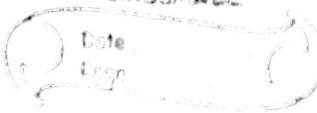
$\left\{ \begin{array}{l} \text{simply re} \\ \text{written the} \\ \text{previous J} \end{array} \right\}$

At the optimum

#3

$$J(\cdot) = \sum_{i=1}^N \alpha_i - \frac{1}{2} w_0^T w_0$$

$\left\{ \begin{array}{l} \text{will be} \\ \text{satisfi} \\ -d_3 \end{array} \right\}$



Now opening ' w_0 '

$$J(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i d_i x_i^T \sum_{j=1}^N \alpha_j d_j x_j$$

so we have

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j x_i^T x_j$$

{ Now this thing have to be Maximised
wrt to α $\because w$ & b are gone }

subject to $[\alpha_i \geq 0 \text{ and } \sum_{i=1}^N \alpha_i d_i = 0]$

→ Dual form

This is also an optimisation problem (QP)

$$Q(\alpha) \rightarrow [\text{QP solution}] \rightarrow (\alpha_i)$$

Now we need to find $g(w) = w_0^T x + b_0$

$$w_0 = \sum_{i=1}^N \alpha_i d_i x_i$$

for any support vector present in class

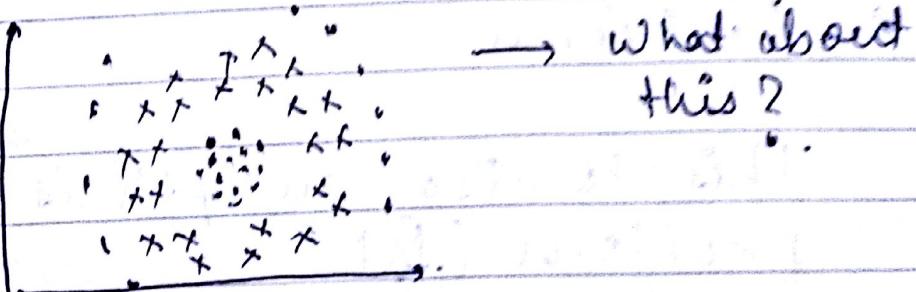
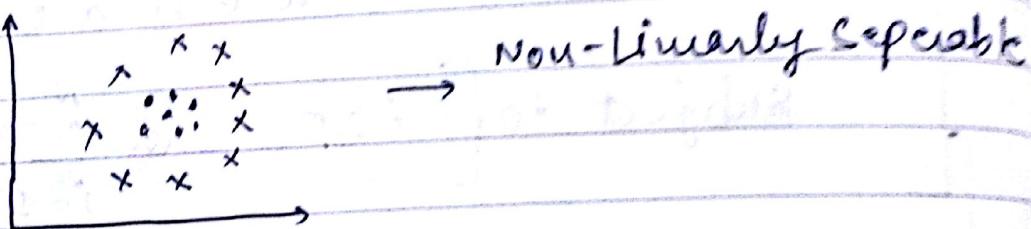
for finding b_0 , Take a support vector
 $(\alpha_i \neq 0)$ Then use KKT

$$d_i (w_0^T x_i + b_0) = 1$$

$$b_0 = 1 - w_0^T X_{c+}$$

{ It's a guaranteed optimal sol }

Note →



for Noisy samples:

are not Linearly Separable but we should allow some $R(\alpha)$

Train

so Introduce slack variables $\varepsilon_i \geq 0$

$$\text{so } \left\{ \begin{array}{l} \phi(w) = \frac{1}{2} w^T w \\ d_i(w^T x_i + b) \geq 1 - \varepsilon_i \end{array} \right\} \forall i$$

(by allowing $\varepsilon_i \geq 1$ we can handle mixing up of samples) — so we have a training error.

So minimise $\sum_{i=1}^N I(\varepsilon_i \geq 1)$
(0 or 1)

or minimise $\sum_{i=1}^N \varepsilon_i$ — additional penalties

So minimise

$$\# \bar{\phi}(w) = \frac{1}{2} w^T w + C \sum_{i=1}^N \varepsilon_i$$

subject to: $d_i(w^T x_i + b) \geq 1 - \varepsilon_i$
 $\varepsilon_i \geq 0, \forall i$

It's dual form →

#6

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j d_i d_j x_i^T x_j$$

subject to $0 \leq \alpha_i \leq C \quad \forall i$ and

$$\sum \alpha_i d_i = 0$$

\rightarrow [additional constraint]

$$Q(\alpha) \rightarrow [\text{QP solution}] \rightarrow \alpha^*$$

\rightarrow The parameter 'C' controls Relative error between VC dimension & training error.

STS \rightarrow

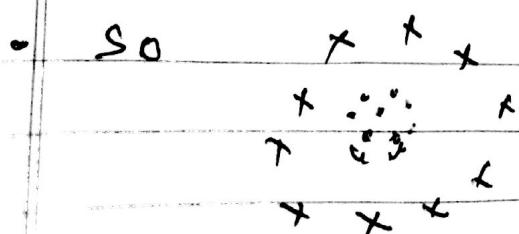
- ITER (International Thermonuclear Reactor)

2nd NOV 2014

Non-Linear Boundaries

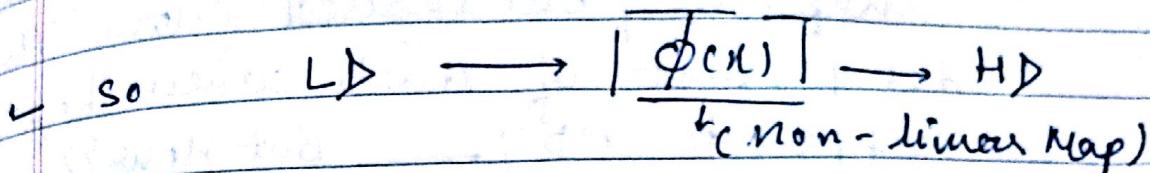
L Good way would be to map them to a ~~lower~~ Dimension so that higher

samples are linearly separable



add an additional dimension 3rd

$$x_3 = x_1^2 + x_2^2$$



We need not worry about 3 dimension because VC Dimension does not depend on it. $\phi(x)$ has to be Non-linear.

so during Training Phase: \rightarrow [TRAINING]

~~given~~ $\phi(x) = \sum x_i - \frac{1}{2} \sum \sum x_i x_j d_i d_j$
 $\phi(x_i) \cdot \phi(x_j)$

subject to

$$0 \leq x_i \leq C \quad \& \quad \sum x_i d_i = 0$$

{ same methods but do $x_i \rightarrow \phi(x_i)$ }

[TESTING PHASE]

$$\text{Label} = \text{sign}(w_0 \cdot \phi(x_{\text{test}}) + b_0)$$

$$w_0 = \sum_{i=1}^N x_i d_i \phi(x_i)$$

$$\therefore \text{Label} = \text{sign} \left(\sum_{i=1}^N (x_i d_i \phi(x_i)) \cdot \phi(x_{\text{test}}) + b_0 \right)$$

dot prod
of sample
and
vector

The thing that is required here is only $\phi(A) \cdot \phi(B)$ so we need not know mapping [but should know the dot product of $A \cdot B$ when they are mapped to H^D] — But How??

- so give the Black Box that performs something a name "Kernel"
Here our kernel (K) gives a scalar when you are inputting it two vectors i.e -

$$s = K(x_i, x_j) \\ = \phi(x_i) \cdot \phi(x_j)$$

so every where

$$\phi(x_i) \cdot \phi(x_j) = K(x_i, x_j)$$

We have delayed ourselves
Now see [MAGIC] :-

1. A simple Quad. Kernel.

let us find $K\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right)$

(Suppose)
where $\phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix}$

and here $K(X, Y) =$

$$(X \cdot Y)^2$$

$$(x_1 x_3)^2 + 2x_1 x_3 x_4 + (x_2 x_4)^2$$

also here $K(x, \beta)$ is computationally easier to compute by doing $(\alpha \cdot \beta)^2$ rather than

$$\alpha \rightarrow \phi(\alpha) \rightarrow \square \quad \beta \rightarrow \phi(\beta) \rightarrow \square \quad \text{dot} \quad \dots$$

original

$$2D \quad K(x, y) = (x \cdot y)^2 \Rightarrow \phi(x) = (4D) \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix}$$

$$2D \quad K(x, y) = (x \cdot y)^3 \Rightarrow \phi(x) = (4D) \begin{bmatrix} x_1^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_2^3 \end{bmatrix}$$

$$3D \quad K(x, y) = (x \cdot y)^3 \Rightarrow \phi(x) = (10D) \begin{bmatrix} x_1^3 \\ x_2^3 \\ x_3^3 \\ x_1^2 x_2 \\ x_2^2 x_1 \\ x_1^2 x_3 \\ x_2^2 x_3 \\ x_3^2 x_1 \\ x_2^2 x_3 \\ x_3^2 x_2 \\ x_1 x_2 x_3 \end{bmatrix}$$

KERNELS \Rightarrow

A) \hookrightarrow Adding to kernels gives you a new kernel.

$$K(x, y) = K_1(x, y) + K_2(x, y)$$

$$\therefore \phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix} \rightarrow \begin{bmatrix} [] \\ [] \end{bmatrix} \cdot []$$

also note $K(x, \beta)$ is computationally easier to compute by doing $(\alpha \cdot \beta)^2$ rather than
 $\alpha \rightarrow \phi(\alpha) \rightarrow \square \quad \square \text{ dot } \square \dots$
 $\beta \rightarrow \phi(\beta) \rightarrow \square \quad \square \text{ dot } \square \dots$

original
2D $K(x, y) = (x \cdot y)^2 \Rightarrow \phi(x) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$

2D $K(x, y) = (x \cdot y)^3 \Rightarrow \phi(x) = \begin{bmatrix} x_1^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_2^3 \end{bmatrix}$

3D $K(x, y) = (x \cdot y)^3 \Rightarrow \phi(x) = \begin{bmatrix} x_1^3 \\ x_2^3 \\ x_3^3 \\ x_1^2 x_2 \\ x_2^2 x_1 \\ x_1^2 x_3 \\ x_3^2 x_1 \\ x_2^2 x_3 \\ x_3^2 x_2 \\ x_3^2 x_1 x_2 x_3 \end{bmatrix}$

KERNELS \Rightarrow

A) \hookrightarrow Adding to kernels gives you a new kernel.

$$K(x, y) = K_1(x, y) + K_2(x, y)$$

$$\therefore \phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix} \xrightarrow{\text{dot product}} \begin{bmatrix} [] \\ [] \end{bmatrix} \cdot \begin{bmatrix} [] \\ [] \end{bmatrix}$$

$$\Rightarrow k(x, y) = \phi(x) \cdot \phi(y)$$

$$= \phi_1(x) \cdot \phi_1(y) + \phi_2(x) \cdot \phi_2(y)$$

also $(x \cdot y)^2, (x \cdot y)^3, \dots, (x \cdot y)^p$ are all valid kernels.

~~A Polynomial kernel~~

$$k_p(x, y) = (1 + x \cdot y)^p$$

$$= 1 + x \cdot y + (x \cdot y)^2 + (x \cdot y)^3 + \dots + (x \cdot y)^p$$

RHS is mapping to extremely HD which have all x_1, x_2, \dots many terms.

~~MERGER'S THEOREM~~

- Using kernels, we avoid explicit mapping with ϕ .
- So we do not have to know what ϕ is as long as we are sure that there exists one.

$$\text{so } k(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \cdot \phi_i(y)$$

$$\lambda_i > 0$$

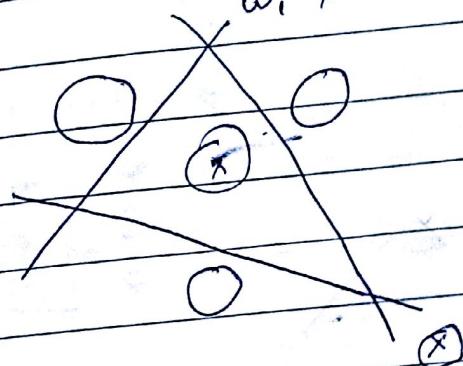
classmate
Date _____
Page _____

Iff $K(X, Y)$ satisfies above things.
→ there exist ϕ .

7-Nov-2014

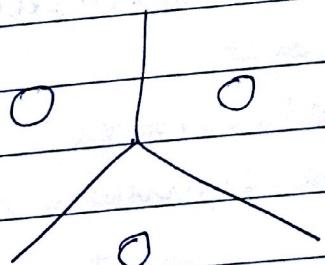
NON-LINEAR MULTICLASS CLASSIFIERS:

1.



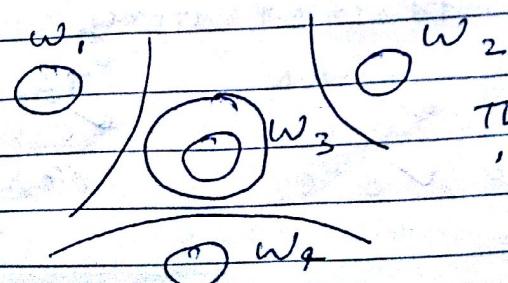
(n) for n classes.

2.



pairwise classifiers
(${}^n C_2$) for n classes.

→ for disputes we can have Confidence measure approach that will give us probability of a sample belonging to a class, but what about.

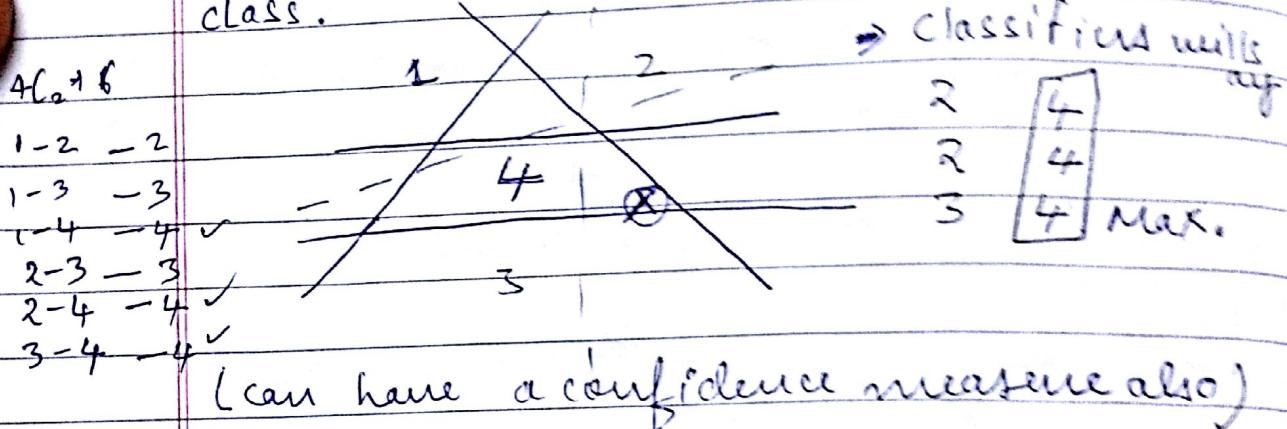


This may possible
"SVM", $w_i \neq w_j$,

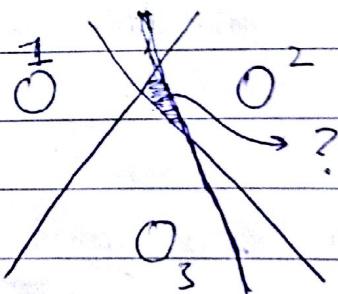
Libraries:-

- ↳ LibSVM
- ↳ SVMLight
- ↳ LibLinear.

- In case of pairwise classifiers find answer for each pair, take the maximum likelihood class.



Usually pairwise classifier is good
but it can also have situation like



so we can have confidence measure

for our first strategy

	TRAIN	TEST
TIME	$\checkmark o(c)$	$\checkmark o(c)$
SPACE	$\checkmark o(c)$	$o(c)$



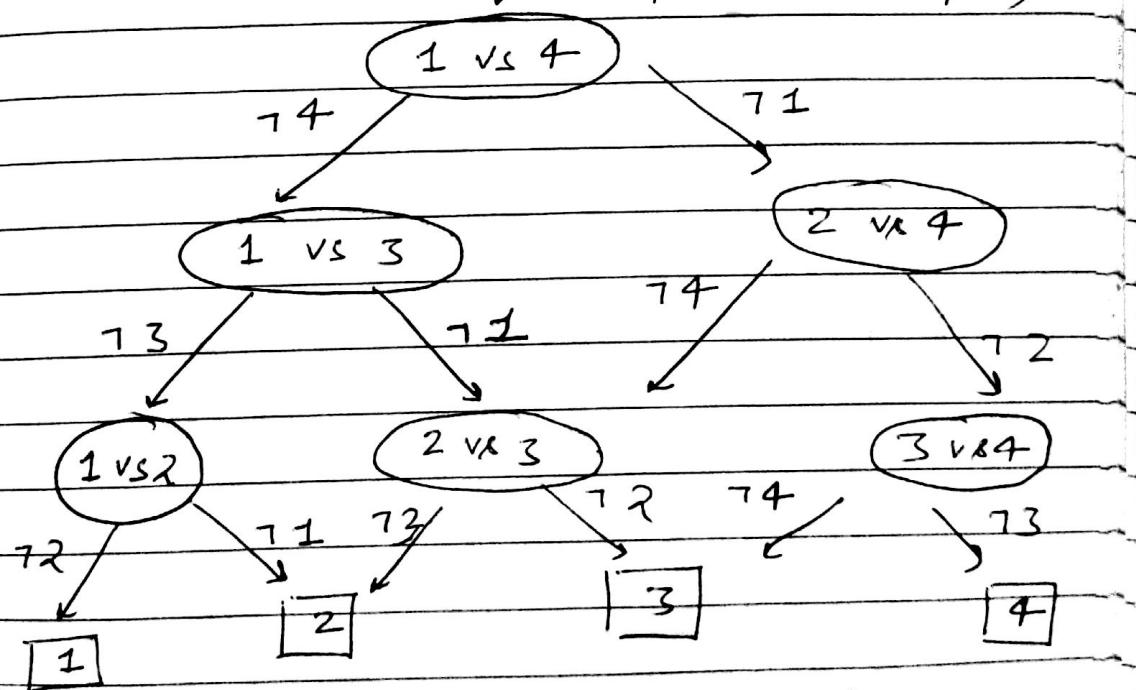
for second strategy

	TRAINING	TESTING
Time	$\checkmark O\left(\frac{C}{2}\right)$	$\checkmark O\left(\frac{C}{2}\right)$
space	$\checkmark O\left(\frac{C}{2}\right)$	$O\left(\frac{C}{2}\right)$

The ticked measures are quite a lot important.

We can also make decisions in pairwise classifiers like:-

↓ INPUT (Test sample)



+ { Time complexity is reduced }
 $\approx O(C)$

We have $DAG = \text{Directed acyclic graph}$

But since each edge is a decision so it is a $DDAG$ (Decision directed)

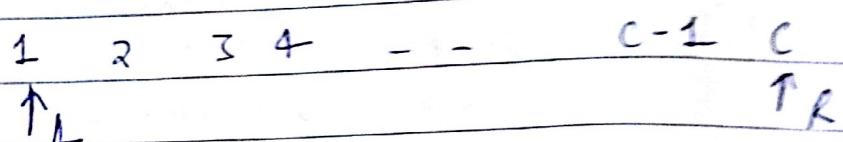
For DDAH

	TRAIN	TEST
Time	$O\left(\frac{C^2}{2}\right)$	$O(C)$
Space	$O\left(\frac{C^2}{2}\right)$	$O(C)$

→ Time complexity is greatly reduced for testing face from $O(C^2) \rightarrow O(C)$

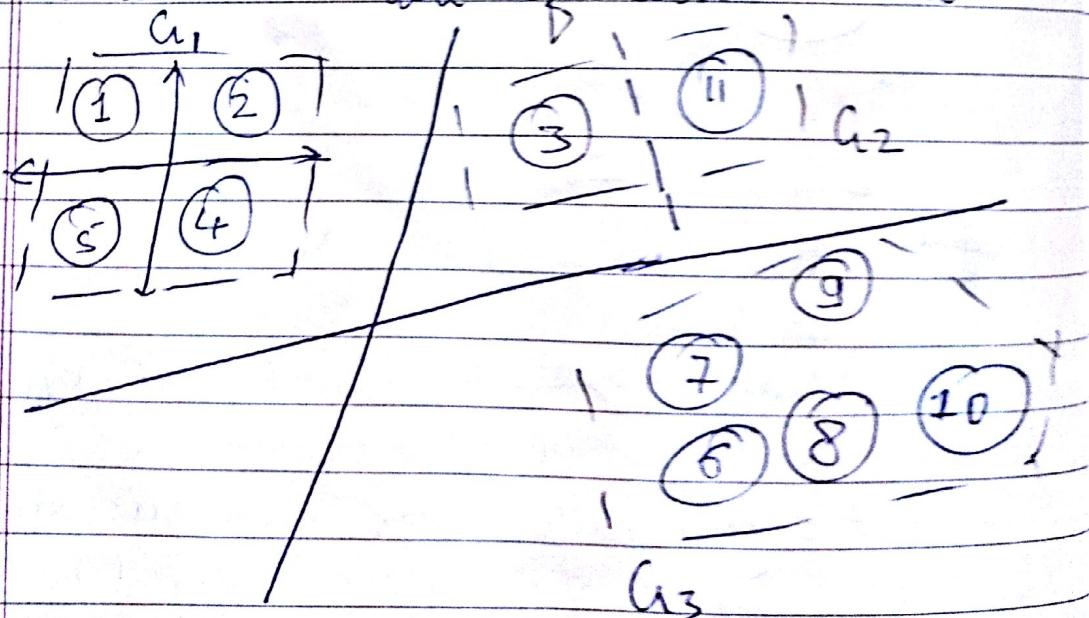
~~(new)~~

for implementing this

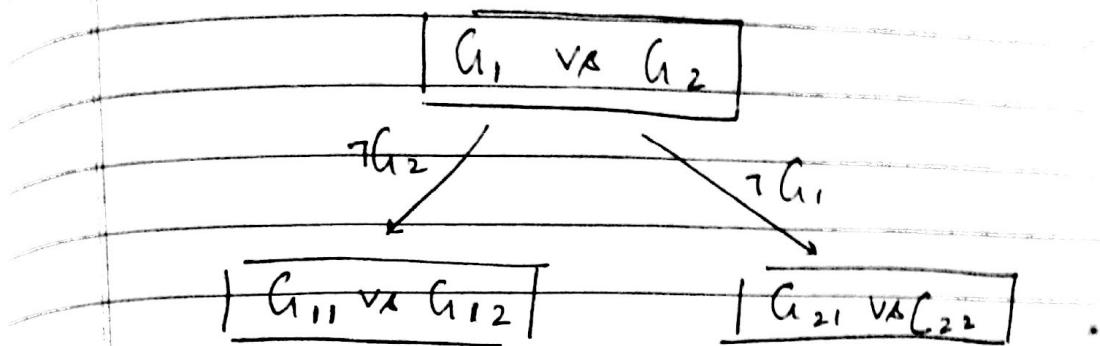


Move $L \rightarrow L+1$ or $R \rightarrow R-1$ based on decision of ($L \vee R$)

↳ Basically we are eliminating classes based on classifiers, so what if we can eliminate more of them at once.



Basically here we are creating groups and based on decisions we are going into ingroup classification.

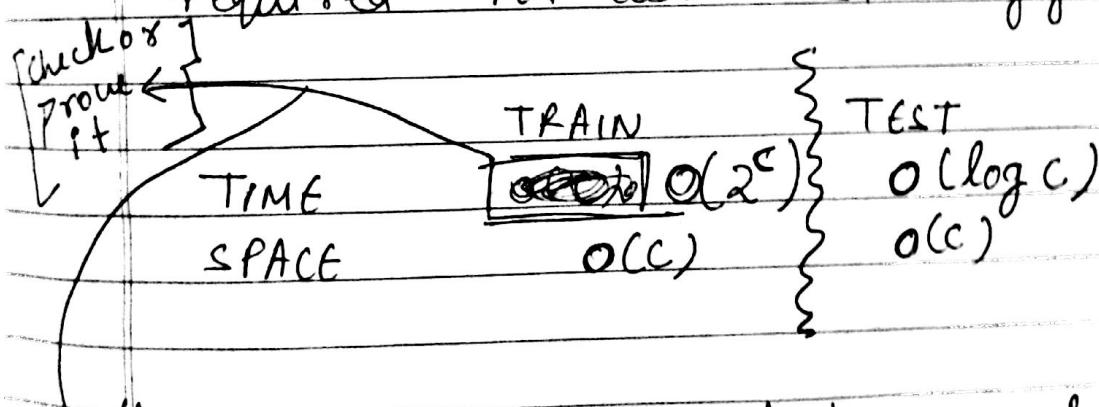


{ These are called BDT (Binary decision Trees) }



— But how to group ??

↳ How many classifiers would be required for above strategy ??



• We can partition given 'C' classes in groups in 2^C possible ways, so which partition should we choose?

TRAINING TIME $\approx O(2^C)$ since we do not know how to choose classifiers Before hand.

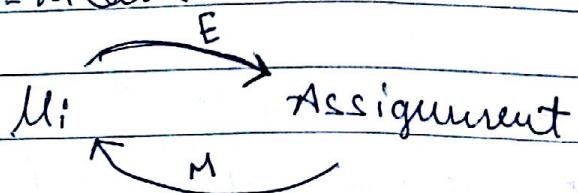
so like there is a trade off between training & testing time, good thing is that Training will become only one so in practice [like digit recognition Handwriting etc.] this can be preferred over others. { But training time is huge!! }

→ [1 class Missed on Clustering]

~~14-11-14~~

Expectation - Maximization Algorithm

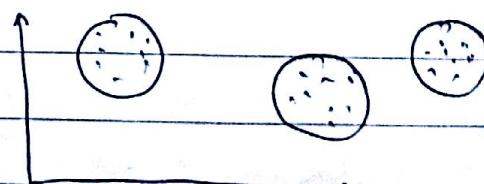
- Like K-Means



L_M = Maximising the likelihood of set of samples coming from μ_i

L_E = again think of μ_i

{ it is like a parameter estimation problem }



Like we will have cases like variances same, mean different

---.

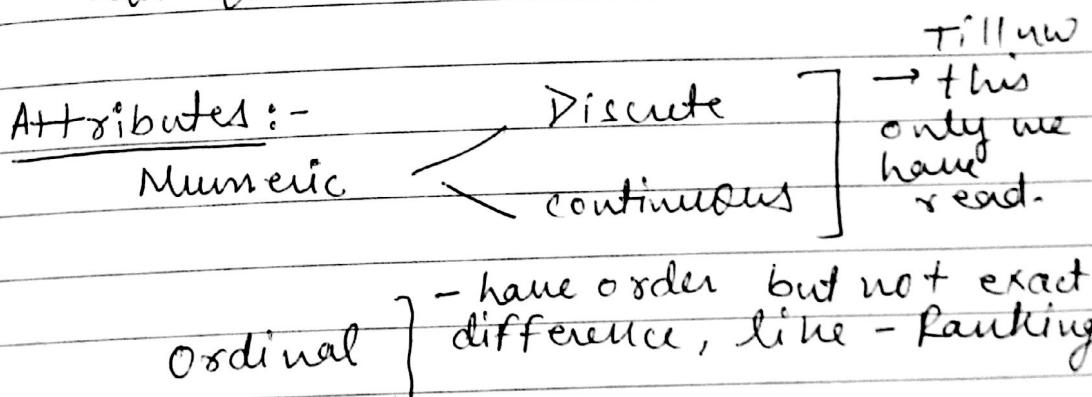
[see what happens for k-means]

Anneic version of K means \rightarrow ANN

- ↳ all the means are different
- ↳ $\Sigma = \text{a diagonal matrix}$
- ↳ inter class variance = same = 0
intra " = same }

.. DECISION TREES \rightarrow

- ↳ we have suppose data Like Language
Now we have to classify people based
on language -- attributes that are
not quantized but are labels.



Categorized/ Normal
{ color, language }
Male/ Female

Like consider classifying fruits :-
color - Green, Yellow, Red
size - small, med, Big
shape - Round, square, Long
Taste - sweet, sour, bitter -

So the problem of classifying by previous
In have to compute $x_1 - x_2$? ^{only}

[Those feature vector differences
can't be computed]

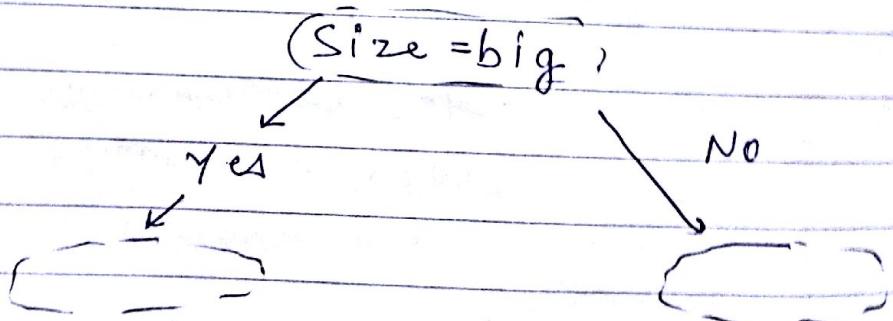
→ Probably we might use
Hamming Distance.

Fruits = Apples, Mangoes, Watermelon,
Banana, Grapes... .

Ques 1. A yellow, long fruit
↳ Banana,

2. A round, green, sweet fruit
↳ Watermelon.

→ Here we are eliminating the things
you can think of it like Decision
Trees.



at each point you can ask a
binary question.

we have what we call a generic decision tree (need not be binary)

ask questions till you reach leaf nodes for all fruits, In such trees you might have to ask same ques. multiple times.

- 1. How to decide what question to ask? — at each node (to have splitting).
- 2. When to stop?
- 3. How to label a Leaf-node?

Ans
.. we can define a measure of purity for splits & if more purity → go on..

Like define impurity measure and

$$i(N), i(N_L), i(N_R)$$

$$\text{so } \{i(N) - P_L i(N_L) - P_R i(N_R)\}$$

↓
(IMPURITY Redn)

We have suppose a node that splits into two so

$P(w_1)$ and $P(w_2)$ we have

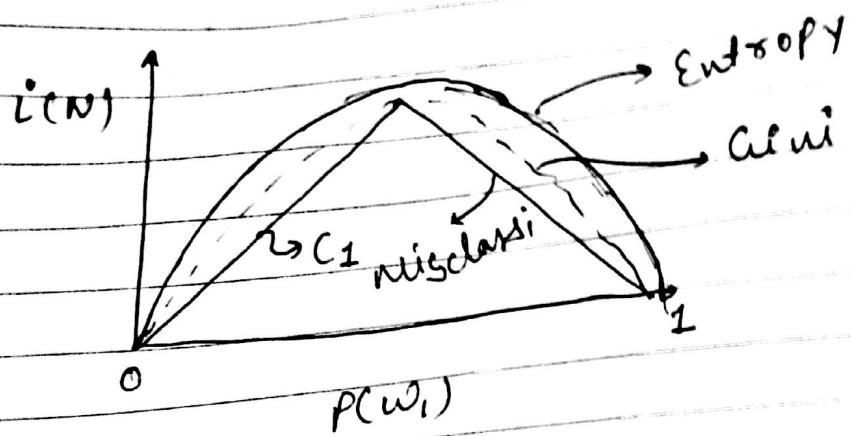
$$P(w_1) + P(w_2) = 1$$

$$0 \leq P(w_1) \leq 1$$

$$0 \leq P(w_2) \leq 1$$

also if $P(w_1) = P(w_2) = 0 \} 0 \text{ impurity}$

$P(w_1) = P(w_2) = 0.5 \} \text{ highest impurity}$



for C_1 -

1. Misclassification: $i(N) = \min(P(w_1), P(w_2), \dots)$

$$= 1 - \max_i P(w_i)$$

2. Variance / Gain

$$i(N) = \sum P(w_i) P(w_j) \quad i \neq j$$

Here

$$= P(w_1) \cdot P(w_2)$$

3. Entropy

$$i(N) = - \sum_i P(w_i) \log_2 P(w_i)$$

(Good question is that for which $i(N)$ is minimised)

{ Good One! } { Very Good }

or [Reduction of Impurity is Maximized]

→ Advantages:

- 1. Execution is very fast
- 2. If trained a human can directly analyze / understand

{ so we can explain why we classified or misclassified a sample
 (can be used in court) — sometimes very important }

In industry different decision trees might be available on cost.
 (can have a very complex training algo)

→ we can also have a forest to decide

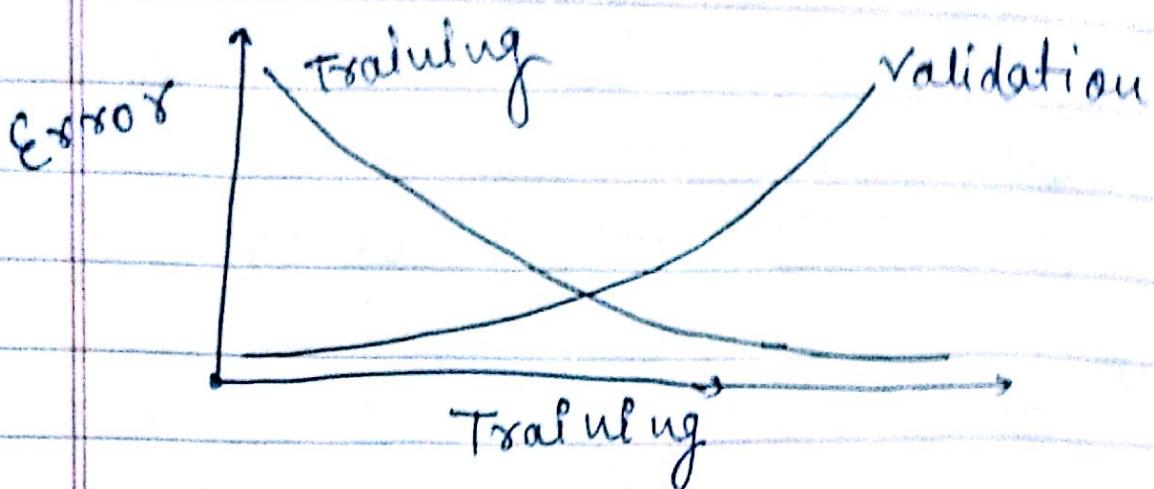
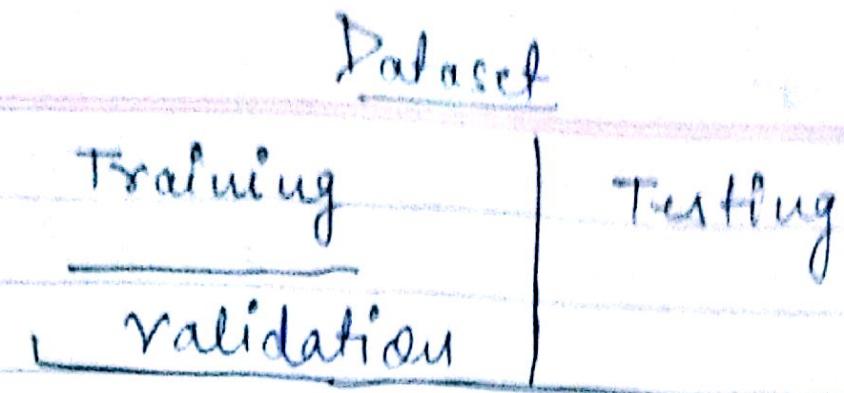
a forest = { set of Random Decision Trees }

In some trees, random noise is introduced.

can be used in deciding if credit card worthy or not, and at each node a complex question like if income > 10000 can be asked...

classmate

Date _____
Page _____



By → when to stop?

← when overfitting happens.

21-Nov-2014

classmate

Date _____
Page _____

Suppose we have features that are varying in dimensions [n, m length long]

for example if you try to compare spoken words - { have different lengths }

Problem - detail , detail
detaaill

(The length is different,
The position of phonem is
different --)

Gait Recognition { classify based on how you walk }

Also we can have a problem, where words are close in spoken format
debatable
detectable

one measure or distance spectable , frozen
edit distance objectable horizon

$$E[i, j] = \begin{cases} \Delta + E[i-1, j-1], & \text{if } i=j \\ \min \left[\text{either}(E[i, j-1], E[i-1, j]) + \Delta, \right. \\ \quad \quad \quad \left. \text{if } i \neq j \right] + 1 \end{cases}$$

→ Dynamic Programming
 {optimal substructure property
 --- }
 ↓

Litmus Test :-

$$R_1 \quad L_1 \quad R_2 = 0$$

$$--- b --- 0$$

If oracle tells you the value of L_1 then you can find soln of R_1 & R_2 (independently and that will be the global soln)

like spectable if t (oracle)
detectable $\downarrow t$
 ED_1 ED_2

→ Edit distance:-

insertion cost = $I(I)$

deletion cost = $D(I)$

substitution = $S(I)$

here we match prefixes.

$\boxed{S}(\textcircled{P})$ to $\rightarrow 2$
 \boxed{d} Nothing

or $(\textcircled{S})P$

$\downarrow d$

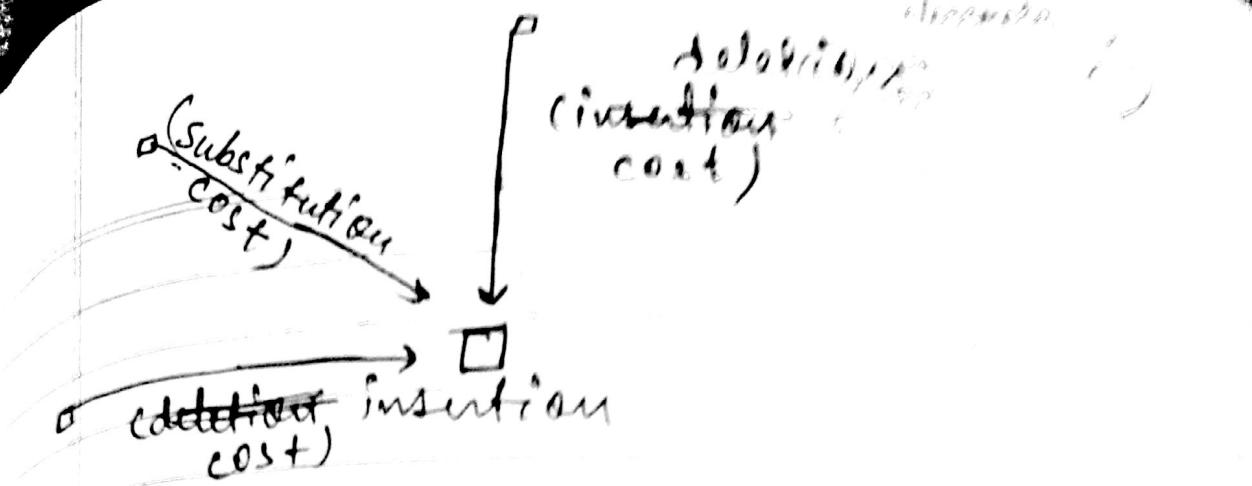
$(\textcircled{S}P) \rightarrow$ Map to

Nothing \rightarrow

$d \rightarrow$ inserted \rightarrow

:	\boxed{N}	$\boxed{\dots \cdot \cdot \cdot \cdot \cdot \cdot \cdot \cdot \cdot}$								
:	\boxed{d}	1	2	3	4	5	6	7	8	9
:	\boxed{e}	2	2							
:	\boxed{f}									

\rightarrow



$$E[i,j] = \min \left\{ \begin{array}{l} E[i-1, j-1] + \cancel{\text{substitution cost}} \\ E[i-1, j] + \cancel{I} \\ E[i, j-1] + \cancel{I} \end{array} \right.$$

To store the path have

a

Trace $(i, j) = \begin{cases} \nearrow & \text{Diagonal} \\ \uparrow & \text{Up} \\ \nwarrow & \text{Left} \end{cases}$

↳ So in case of speech matching
or

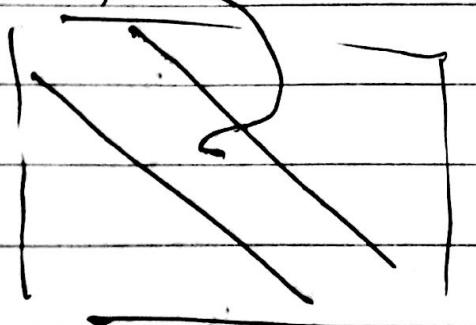
n to m Edit mapping
find edit distance & align

Then { detail } to { detail }

↳ the phonemes are aligned and
the costs will be different.

→ Previous method is very popular and is known as DTW (Dynamic Time Warping)

↳ Some optimisation by Sakoe chiba known as Sakoe chiba band



The above $E(i, j)$ can be used as a kernel function for SVM's.