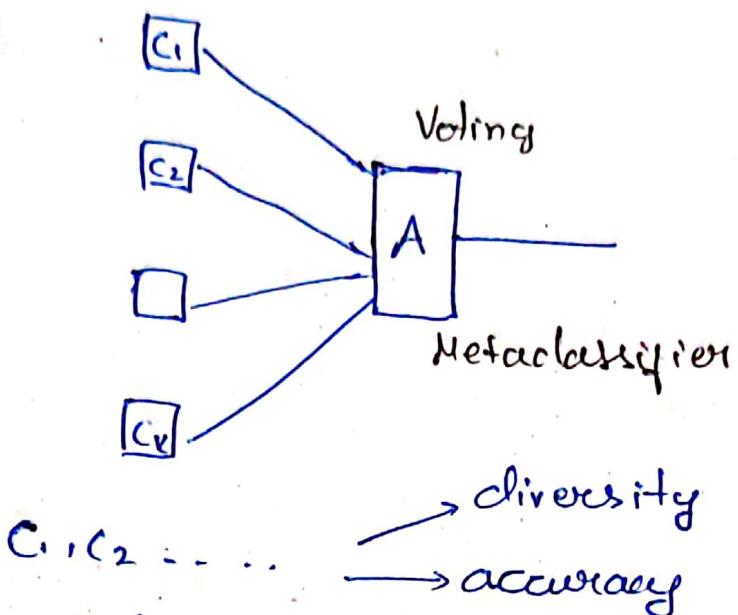


→ Ensemble Classifier :-



- 1) Some classifiers but trained on different training data sets.
- 2) Classifier is same, input features are different.
- 3) Nature of classifier is also different.

Metaclassifier → Inputs are the outputs of the 'k' classifiers

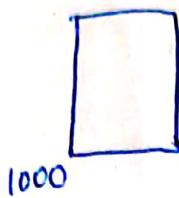
→ If is trained on them and final output is given.

* Training can be done in 2 ways :-

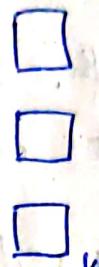
- 1) Train C_i and A together.
- 2) first train C_i s and then train ' A '.

→ Bagging : (Bootstrap Aggregation)

Bootstrap sample → getting multiple samples from a given single dataset.



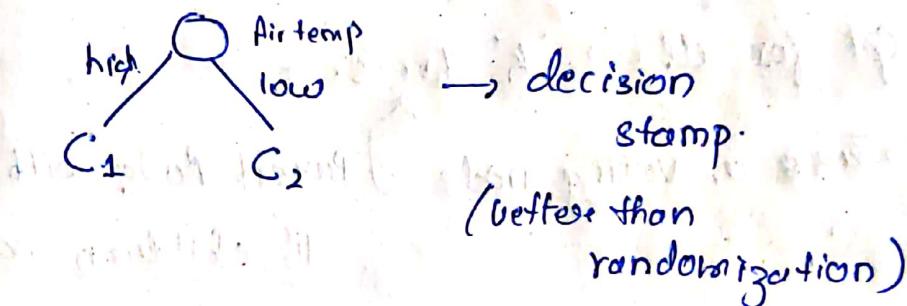
Sample randomly with replacement



Bootstrapped samples.

Voting

→ Base classifier :- weak classifier



→ Instead of training a large classifier, training very large weak classifiers and combining them is good enough.

Random Forest (Extension of above):-

→ Bootstrapped training set

→ random feature subsets.

1 — Height Bounded decision tree.

2
⋮
⋮
15

$x_1 \dots x_d$

↳ Randomly select the subset of features.

level 1
i.e randomly project 6-d data to 3-d etc.

→ On this randomized / projected data, create a decision tree.

level 2 :- Repeat the same, excluding the root attribute.

until required height-

∴ T_1 (Bootstrap set) $\rightarrow BC_1$

Get for all $\rightarrow k$ trees.

→ Take a voting node \uparrow Parent node with
 k children nodes.

→ Useful if input dimension is very high.

Random Forest \rightarrow Random Projections

+
Bagging

+
Bayes Classifier (on top).

Ensemble

↳ bias + variance ↓ | decrease

the variance by maintaining bias.

to 3-d etc.
a, create a

ling the

Bootstrap \longrightarrow Special case of Re-sampling
data.

Bagging, Random
Forest
↓
Resample
or
feature
subset
Resampling
on
features
also.

\longrightarrow Resampling
+
Training
+
Combining

* Adaptive Re-Sampling :-

↳ 1st Practice set \rightarrow Random

2nd " " \rightarrow where 1st failed

3rd " " \rightarrow where 1st & 2nd failed.

i.e. kth resample depends on its previous
resamples & classifiers.

\rightarrow 1st classifier is close a weak classifier that
is never good on all features.

* Add Boost :- (Adaptive boosting)

\rightarrow Boosting

x_1, x_2, \dots, x_n

w_1, w_2, \dots, w_n \rightarrow change with
iterations

i.e. weight of set where

if previous classifier fails it's increased
in next iteration.

based on the failure
set.

maintaining

Iteration 1 :- $w_1 = w_2 = \dots = w_n = \frac{1}{n}$



\rightarrow Samples with replacement.

choose a weak classifier as base classifier.

$\rightarrow C_1 \rightarrow \epsilon$ (error)

\hookrightarrow wrongly predicted.

Iteration 2 :-

$w_{i+1} = w_i e^{+\alpha}$ if x_i is misclassified by C_i

$= w_i e^{-\alpha}$ if x_i is correctly classified by C_i .

Repeat & combine $C_1 - C_k$.

Adaboost :-

Ensemble

x_1, x_2, \dots, x_n

$y_1, y_2, \dots, y_n \rightarrow$ jth iteration.

D_1, D_2, \dots

$D_n, \sum_{i=1}^n D_i = 1$

y_n

$\downarrow n$

Probabilities that they are selected.

$S_1 =$

x_1, x_2, \dots, x_n

C_1

(weak classifier)

$C_1 \xrightarrow{\text{train on}} S_1$

$\frac{1}{N}$

ement.

use classifiers.

ed.

classified by
 C_i

rectly classified
by C_i



weak classifier \rightarrow slightly better than
random classifier.

Now use C_1 on $\mathcal{D} \rightarrow$ predict final output.

Let $\epsilon_1 = \sum_j D_j^0 + j$ if x_j is misclassified
else compare with y_j .

$$\text{Let } \alpha_1 = \frac{1}{2} \ln \left(\frac{1 - \epsilon_1}{\epsilon_1} \right)$$

\hookrightarrow so α is +ve (as weak classifier).
(as $\epsilon_1 < 1/2$).

$$S_2 = D_1^1 \cdot D_2^1 \cdots D_n^1$$

if x_j is correctly classified by C_1

$$D_j^1 = D_j^0 e^{-\alpha_1}$$

else if misclassified.

$$D_j^1 = D_j^0 e^{+\alpha_1}$$

By proof

of convergence

Construct $C_2 \xrightarrow{\text{train on}} S_2$.

$$\epsilon_2 = \sum_j D_j^1$$

+ mis-classified

$$\Rightarrow \alpha_2 = \frac{1}{2} \ln \left(\frac{1 - \epsilon_2}{\epsilon_2} \right)$$

\hookrightarrow use for next iteration.

$$C_1 \quad C_2 \quad \dots \quad C_T$$

$$\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_T$$

\rightarrow all +ve.

$$\left(\alpha_i = \frac{1}{2} \ln \frac{1 - \epsilon_i}{\epsilon_i} \right)$$

C_1
(weak
classifier)

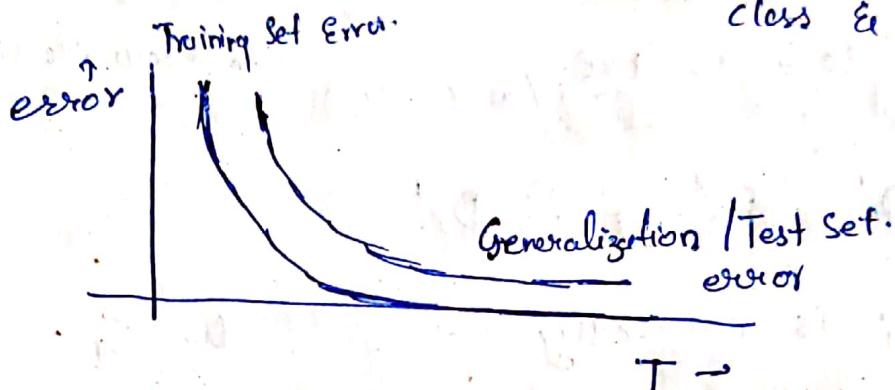
Final Classifier :-

$$C = \sum_{i=1}^T \alpha_i C_i$$

(C_i is the class output values), i.e. lower ' ϵ ', higher weight for classifier.

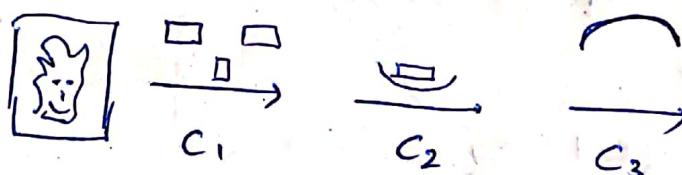
↳ AdaBoost Algorithm.

Add all α 's that predict particular class & compare.



↳ Refer the Paper on Websize for mathematical proof:

Viola-Jones Face recognition :-

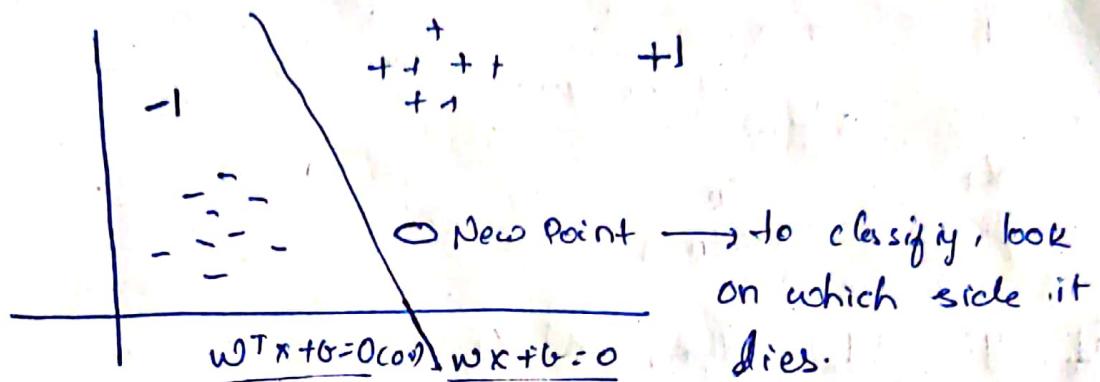


↳ combine these classifiers.

- * All the classifiers are decorrelated.
i.e no classifier agrees on more than 50% of the examples.
- * New classifiers need not be considered if similar to previous ones.

Support Vector Machine :-

Linear Discriminants.



How to draw the discriminant?

→ $X \in \mathbb{R}^d$ i.e d-dimensional vectors
i.e (each point).

$$X = [x_1 \dots x_d]$$

→ class level $\in [-1, +1]$ → assumption

(only 2 classes).
 $S = [(x_1, y_1) \dots (x_n, y_n)]$

$$w = [w_1 \dots w_d]$$

$$w \cdot X = w_1 x_1 + w_2 x_2 \dots w_d x_d$$

$$\text{line } w \cdot X + b = 0$$

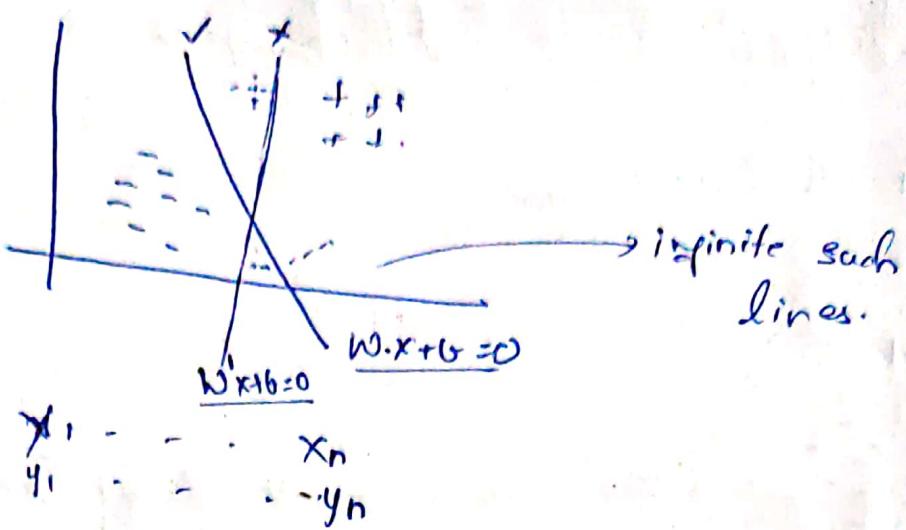
$$(or) w^T \cdot X + b = 0$$

⇒ Notations.

+ve side → non-origin side.

$$w \cdot z + b > 0 \text{ if } y(z) = +1$$

$$w \cdot z + b < 0 \text{ if } y(z) = -1$$



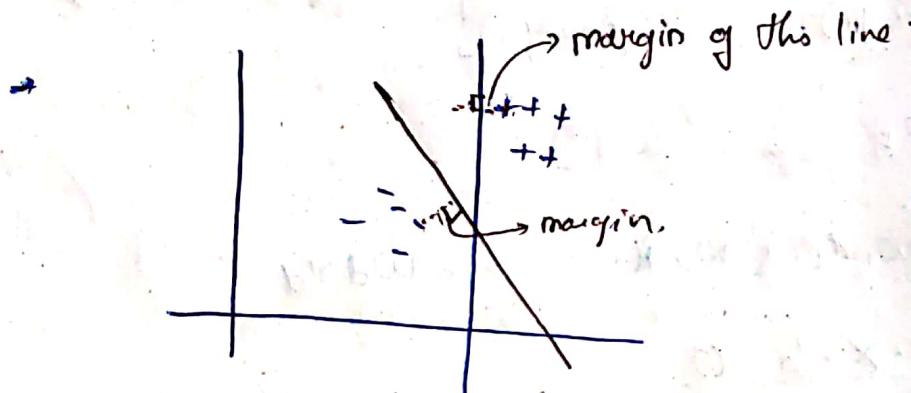
if $y_i = +1$ then $w \cdot x_i + b > 0$

if $y_i = -1$ then $w \cdot x_i + b < 0$.

$\forall i$

for all i , $y_i \cdot (w \cdot x_i + b) \geq 0 \cdot \forall i$

→ Prefer more equidistantly placed line.



Margin of the line → \perp distance from that classifying line to the closest point.

→ central line has more margin than the others.

Choose the one having highest Margin.

→ See missed notes:

Primal Problem:-

$$\text{Minimize } \frac{1}{2} w^T w$$

$$\text{such that } y_i(w \cdot x_i + b) \geq 1 \quad \forall i.$$

Lagrange Multiplier:-

$$(\text{Minimize}) L = \frac{1}{2} w^T w - \sum_{i=1}^n \lambda_i (y_i(w \cdot x_i + b) - 1)$$

$$(L) \leftarrow \text{s.t. } \lambda_i \geq 0 \quad \forall i$$

$$\frac{\partial L}{\partial w} = 0, \quad \frac{\partial L}{\partial b} = 0,$$

$$w = \sum \lambda_i y_i x_i$$

some solution as
above (optimization
problem)

↓ direction factor.
and $\sum \lambda_i y_i = 0$ scaling factor.

$$L = \sum \lambda_i - \frac{1}{2} \leq \sum \lambda_i x_j y_i y_j x_i^T x_j$$

$$L = \Lambda U^T - \frac{1}{2} \Lambda H \Lambda^T$$

$$\Lambda = [\lambda_1, \dots, \lambda_n], \quad U = [1, \dots, 1]^T$$

$$H = [y_i^T y_j, x_i^T x_j]_{n \times n}$$

↳ Solving the Dual Problem.

→ finding $\underline{\lambda}$'s to optimize the problem.

↓ optimizing variable. (w & b can be derived from λ).

if $\lambda_i = 0$, then corresponding x_i will not determine 'w'.

↳ has a very sparse solution.

$\lambda_i > 0 \rightarrow x_i$ are called support vectors.

Most of the $\lambda_i = 0 \rightarrow$ Boundary solutions.

only few $\underline{\lambda_i > 0} \rightarrow$ interior solutions.

$$w = \sum \lambda_i y_i x_i$$

$$\lambda_i \geq 0.$$

Primal Problem \leftarrow Some solutions \rightarrow Dual Problem.

$$\text{Minimize } \frac{1}{2} w^T w$$

$$L = \frac{1}{2} w^T w - \epsilon \sum y_i ($$

$$\text{s.t. } y_i (w \cdot x_i + b) \geq 1$$

$$(w \cdot x_i + b) - 1)$$

$$\begin{cases} \lambda_i = 0 \\ \lambda_i > 0 \end{cases}$$

Karush-Kuhn-Tucker Condition:- (KKT)

In primal problem:-

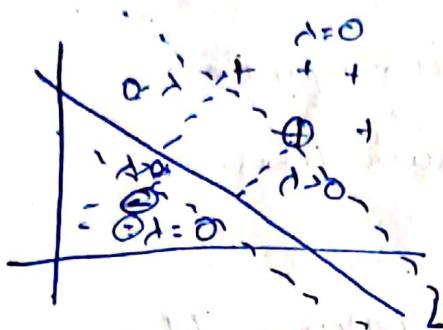
$$y_i (w \cdot x_i + b) \geq 1 \leftarrow \lambda_i = 0$$

$$y_i (w \cdot x_i + b) = 1 \leftarrow \lambda_i > 0$$

i.e. the boundary & interior area in dual invented in primal problem.

→ set of closest points.

i.e. only the closest points have non-zero ' α ' and only they determine ' w '.



both sides.
→ keep expanding the line, till
H touches a point.

$$w = \sum \alpha_i y_i x_i \quad \rightarrow \text{Support Vectors.}$$

\hookrightarrow closest points.

$$\therefore w = \sum \alpha_i y_i x_i$$

s.t. x_i is a S.V. How to find b ?

We know for S.V.,

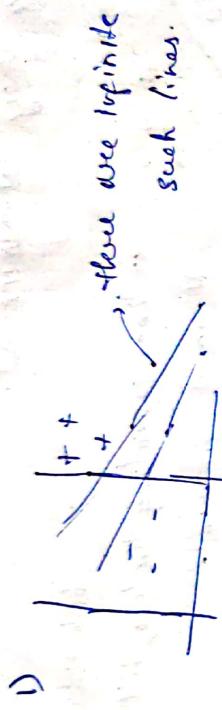
$$y_i(w \cdot x_i + b) = 1 \quad \rightarrow \text{all known except } b.$$

\hookrightarrow obtain 'b' from here.

steps :-

- 1) Prepare the $n \times n$ matrix.
- 2) Input that to a standard dual P solver.
- 3) Collect $\alpha_i > 0$.
- 4) find w & b .
- 5) New point can be classified using the line.

Assumptions :-



→ Only if there is a closed margin b/w the classes.

or -

+	+	+	→ No line.
-	+	-	$\rightarrow y_i(w \cdot x_i + b) \geq 1$
-	-	-	

for all.

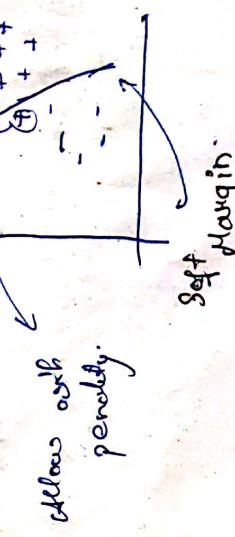
so to ignore the overlapping point i.e.

Consider $|y_i(w \cdot x_i + b)| \geq 1 - \epsilon_{j,i}$ for all.

Minimize $\frac{1}{2} w^T w + C$ and $\sum_{i=1}^n \epsilon_{j,i}^2$.

→ Allow Ep pencilize relaxation.

$\epsilon_{j,i}$ is different for various points.
 $\epsilon_{j,i}(x_i)$



Final dual P :-

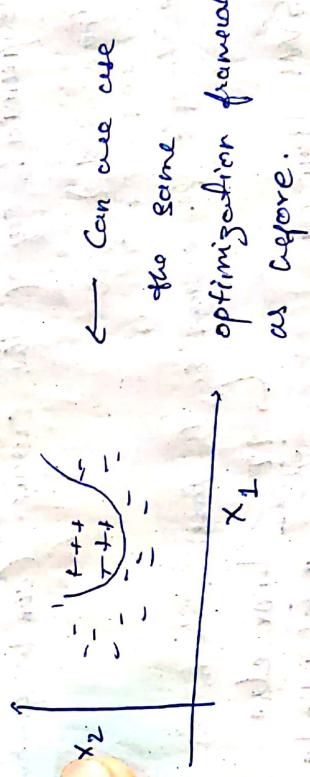
$$L = \lambda U^T - \frac{1}{2} \alpha H \alpha^T$$
$$0 \leq \alpha_i \leq C$$

↳ only change is an upper limit
for ' α '.

$C \Rightarrow$ Generalization constant .

→ No mathematical way to choose optimal ' C '.
→ trial & error method.

⇒ Curve Boundary :-



$$2x_1^2 + 3x_2^2 + 5x_1x_2 + 7 = 0 \rightarrow \text{example, curve}$$

→ Map it to a higher dimensional space i.e.

$$\begin{cases} x_1' \\ x_2' \\ x_3' \end{cases} \rightarrow \begin{cases} x_1' = x_1^2 \\ x_2' = x_2^2 \\ x_3' = x_1x_2 \end{cases}$$

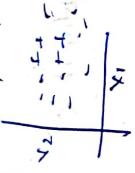
→ map each point accordingly ✓

Advantage?

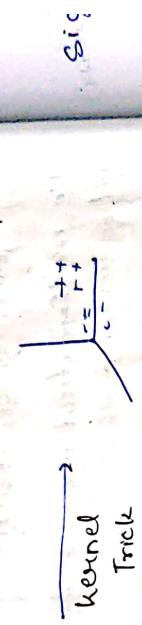
$$\text{the newt. curve exp} \Rightarrow 2x_1' + 3x_2' + 5x_3' + 7 = 0$$

Linear Equation.
i.e. the boundary is a line.

Input Space



Feature Space



- Introduce an axis for every term in the polynomial eq of the curve. (Can be greater than 3-d.)
- Since the curve is its possible form.

SVM Training:

$$\text{Min } L = \lambda u^T - \frac{1}{2} \|u\|^2$$

Subject to $\lambda = [\lambda_1 \dots \lambda_n] \geq 0$

$$H = [y_i x_i' \cdot x_p']_{n \times n} \rightarrow \text{dot product.}$$

SMV Testing :-

$$\text{Sign} (w \cdot x_i' + b) = w = \sum_{i=1}^n \alpha_i y_i x_i'$$

$\rightarrow \alpha_i = 0$ if x_i is not a support vector.

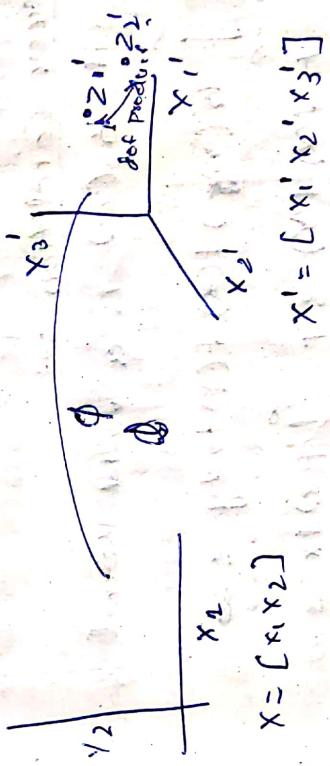
$\rightarrow \alpha_i > 0$ if x_i is a SV.

$$N = \sum_{x_i' \in \text{SV}} \alpha_i y_i x_i'$$

$$\text{Sign} (\sum_{x_i' \in \text{SV}} \alpha_i y_i x_i' \cdot x_k' + b)$$

↳ dot product.

so the only operation in both training & testing is dot product only.



$$x = [x_1 \ x_2]$$

$$x' = [x_1' \ x_2' \ x_3']$$

$$x_1' = x_1^2, \quad x_2' = x_2^2, \quad x_3' = x_1 x_2$$

$\rightarrow x' = \phi(x) \rightarrow$ high dimension / infinite dimensional vector.

$$z_1' = \phi(z_1), \quad z_2' = \phi(z_2)$$

$$z_1' \cdot z_2' = \phi(z_1) \cdot \phi(z_2) \rightarrow ?$$

$$(\text{Kernel trick}) = \text{Is } (z_1, z_2) \text{ directly on } \mathbb{B}_{1,2}$$

$$= b(z_1, z_2)$$

✓
 → Combination of dot product and ϕ .
 (inner)

$$H = \begin{bmatrix} y_i y_j x_i' x_j' \end{bmatrix}_{i,j} \quad \text{if } 'y' \text{ is a valid kernel, then 'H' is positive definite.}$$

↓
 scalar.

Sign ($C \in \mathbb{R}^{n \times n}$) is (x_i, x_k) .

H is a positive definite
 $x^T H x \geq 0$ for all x .

Mercer Condition:

If $K(\cdot, \cdot)$ is a valid kernel function,

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} K(x, y) f(x) f(y) dxdy < \infty$$

and

$$\int_{-\infty}^{\infty} |f(x)|^2 dx < \infty$$

→
 functional analysis. (square integrable / L_2 function)

Linear kernel:

$$K(x, y) = x \cdot y \quad (\text{dot product})$$

Quadratic kernel:

$$K(x, y) = x \cdot x + y \cdot y + x \cdot y$$

Polynomial kernel:

$$K(x, y) = (1 + x \cdot y)^d$$

Radical Basis Kernel (RBF) / Gaussian kernel:

$$K(x, y) = e^{-\|x-y\|^2}$$

→ distance b/w

the points.

↳ Some of the valid kernels:

→ Similarity can be defined between 2 objects (not only points) like graphs, complex objects, etc., using kernel functions.
↳ Classifier.

↳ Required to find? → Best kernel.

(Research → β -kernels)

- * If K_1, K_2 are valid kernels, then
→ $\alpha K_1(x, y) + (1-\alpha) K_2(x, y) \rightarrow$ is also valid
↳ (Composite kernels?)

Ex:- $X = ATGCTCAT$

$T = AACGTCCAT$

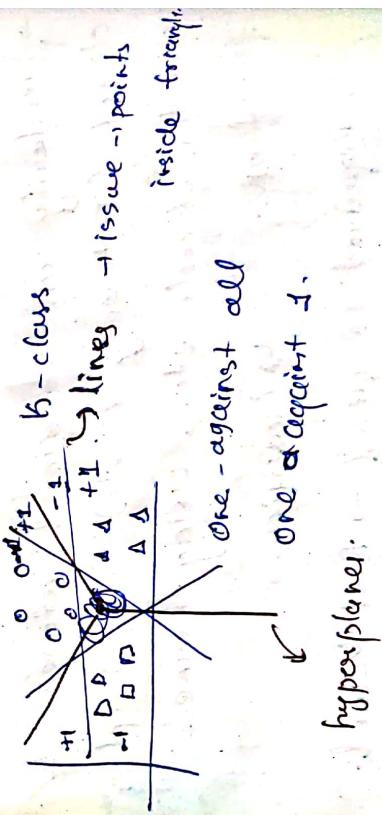
$in(X, Y) = \text{No. of common sub-strings}$
between X, Y

→ All possible anis for each possible
Sub-string : No. of its occurrences = its
value on that anis. etc..

∴ if a kernel is designed, no need of an
input vector

$\hookrightarrow (2\text{-class SVM})$
(all here).

Multiclass SVM :-

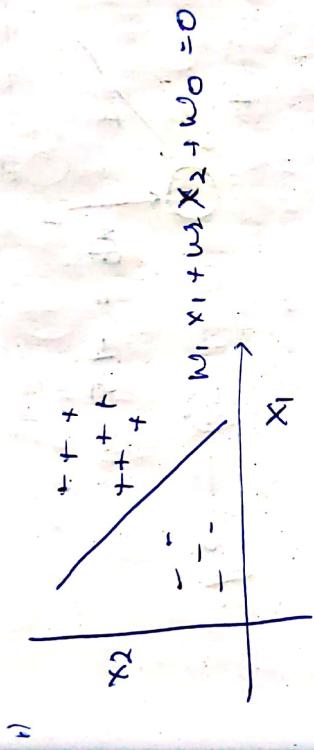


DAGs:-
→ Like a tree
of hyperplanes.

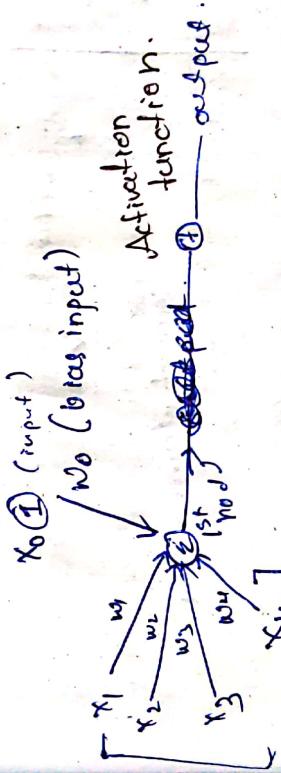
The discussed numerical methods previously have cubic complexity.

Optimizations to decrease complexity.

Eg:- Sequential Minimal Optimization (SMO)
etc...



→ Perceptron / Single Neuron :-



Input vector

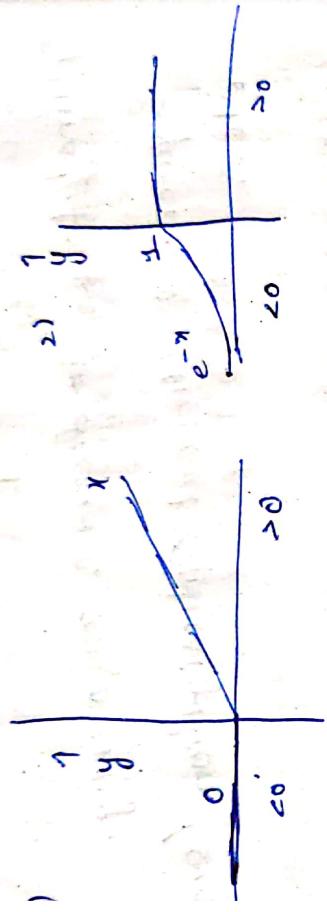
Ex coefficients.

1st node compute :- $\sum_{i=0}^n w_i x_i$, where
 $w_0 = 1$

$w_0 = \text{bias}$.

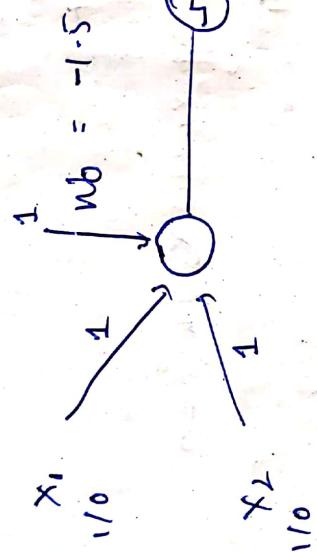
$$\text{output} = f(\sum_{i=0}^n w_i x_i)$$

Activation function example :-

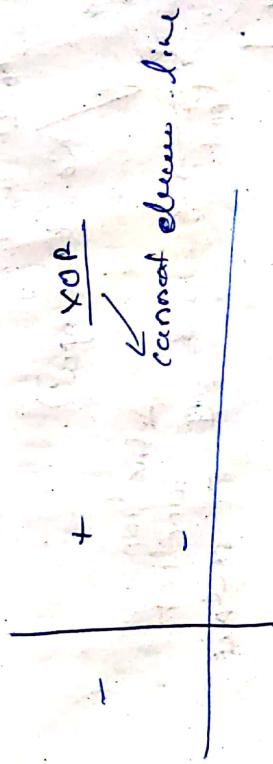


etc -:-

x_1 AND x_2 :-

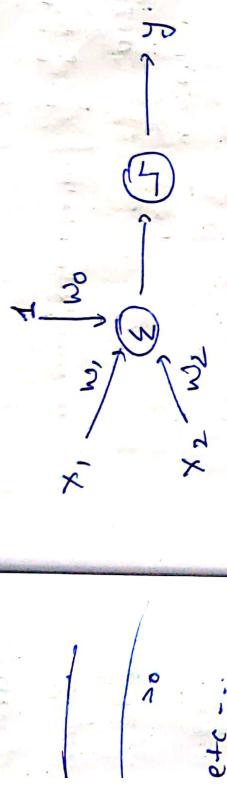


$y = x_1 \text{ AND } x_2$



\Rightarrow XOR cannot be generalized using this.

Perception Learning rule:-



→ Random initialization of weights.

Online - Training:-

- Give the inputs & using the present weights, compute 'y' & match it with the desired output.
- If they do not match, update weights else leave as it is.

$$w^{t+1} = w^t + \eta \Delta w$$

Learning rate ≈ 0.1

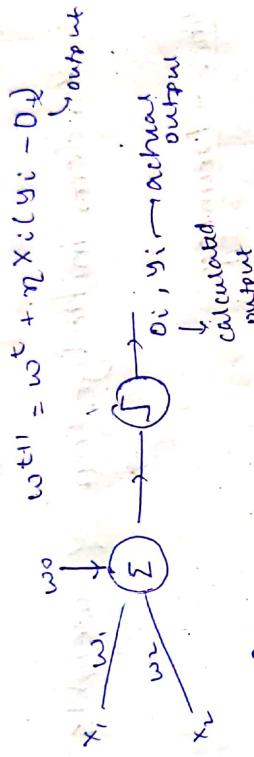
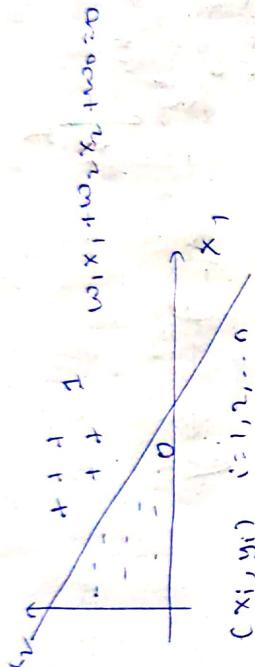
$$\Delta w = (t - y) X$$

↓ ↓
 predicted desired

Start with random weights. Continue the learning till it realizes the function.

25/10/19

Perception



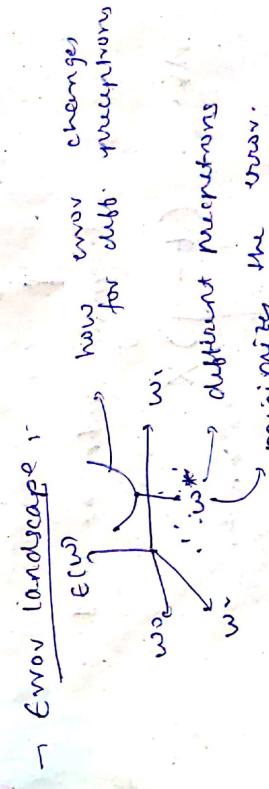
Delta Rule:

- assume a value of w_0, w_1, w_2, \dots & calculate the given output o_i for given input.

$$\rightarrow \text{error } e_i = \frac{1}{2} (y_i - o_i)^2$$

$$\text{Total Error } E(w) = \frac{1}{2} \sum_{i=1}^n (y_i - o_i)^2$$

error is a function of w



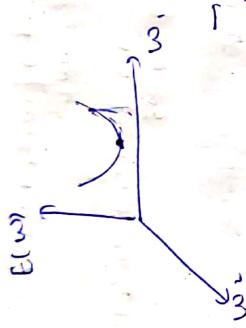
- if $E(w)$ has a parabolic shape then we can find w (min. error)

- we can find w^* by Newton-Raphson method
- Newton Raphson generalized to an algorithm called as

Gradient - Descent Optimization algo

2-min in a convex
otherwise 1

w' & w



→ gradient descent algo
takes us to the lowest point

→ direction of steepest descent

$$\nabla E = \left[\frac{\partial E(w)}{\partial w_1}, \frac{\partial E(w)}{\partial w_2}, \dots, \frac{\partial E(w)}{\partial w_d} \right]$$

→ negative direction
of that vector
corresponds
to direction
of tangent

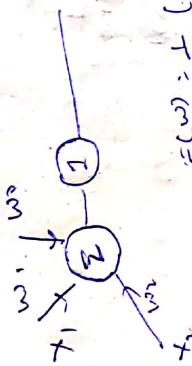
$$w^{t+1} = w^t + \eta (\nabla E)$$

↓
old vector
new vector

$$E(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - o_i)^2$$

$$O_i = \sum_{j=1}^n w_j x_{ji}$$

$$= w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$



$$\nabla E = \frac{1}{n} \sum_{i=1}^n \nabla E(w)$$

$$\nabla E = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^n w_j x_{ji} - y_i \right) \frac{\partial}{\partial w_j}$$

$$E(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - (\underbrace{w_0 + w_1 x_1 + w_2 x_2 + \dots}_o))^2$$

→ o_i → output of neurons

$$\nabla E(w) = \frac{1}{2} \sum_{i=1}^n \left[\frac{\partial}{\partial w_0} [w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots] \right. \\ \left. + \frac{\partial}{\partial w_1} [w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots] \right. \\ \left. + \frac{\partial}{\partial w_2} [w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots] \right]$$

$$\frac{\partial E}{\partial w_1} = \frac{1}{n} \sum (y_i - o_i) x_1$$

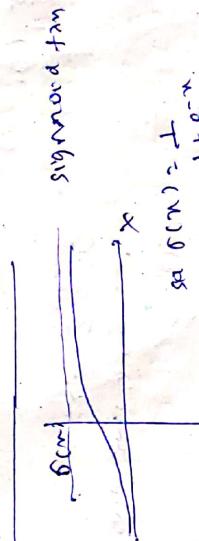
$$\frac{\partial E}{\partial w_2} = \frac{1}{n} \sum (y_i - o_i) x_2$$

x_{1i}, x_{2i}, \dots

$$\nabla E(w) = \frac{1}{n} \sum_{i=1}^n (y_i - o_i) x_i$$

$$w^{t+1} = w^t - \eta \frac{1}{n} \sum_{i=1}^n (y_i - o_i) x_i$$

↳ This is a batch algorithm unlike the previous one.
gradient descent algorithm.



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

if we use sigmoid for men

$$w^{t+1} = w^t - \eta \frac{1}{n} \left(\sum_{i=1}^n (y_i - o_i) x_i \right)$$

$$y_i = \sigma(w^t \cdot x_i)$$

$$x = [x_1 \ x_2 \ x_3 \ t] \quad | \quad x_1, t_1$$

$$w = [w_1 \ w_2 \ w_3 \ w_0] \quad | \quad x_2, t_2$$

$$s = \sum_{i=0}^3 w_i x_i \quad | \quad x_n, t_n$$

$$y = \sigma(s) \rightarrow \text{activation function.}$$

$$\sigma(s) = \frac{1}{1+e^{-s}}$$

$$\sigma(s) = \frac{1}{1+e^{-s}} \text{ sigmoid.}$$

$$\sigma'(s) = \sigma(s)(1-\sigma(s))$$

gradient descent

$$E(w) = \frac{1}{2} \sum_{i=1}^n (t_i - y_i)^2$$

$$w^{t+1} = w^t - \eta \nabla E$$

$$= w^t - \eta \sum_{i=1}^n (t_i - y_i) x_i$$

$$\begin{array}{c} + \\ - \\ - \\ - \end{array}$$



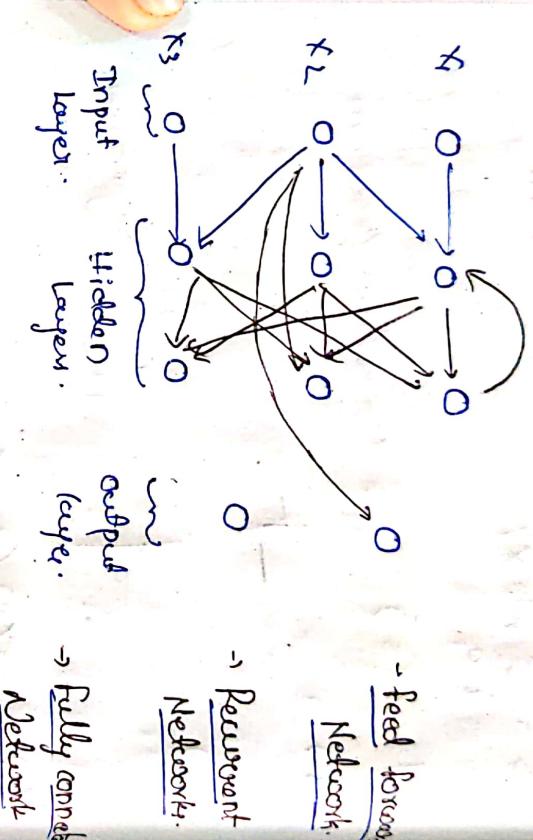
\rightarrow Single line cannot separate them.

But 2 lines can:

Similarly more than one perceptron can be used to distinguish such cases.

Multi Layered Perceptron (MLP) :-

Architecture , Training Algorithm.



In a feed forward Network :-

- 1) Output of a layer is the input to the next layer. Unidirectional flow.

In a Recurrent Network:

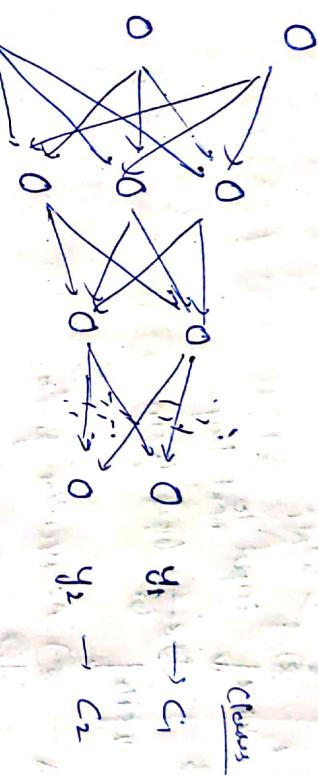
- 1) Data can flow backwards also.

→ Modifications in these networks can have data flow not only between immediate layers, but also further (or) previous layers.

In a Fully connected N-

→ Output of every node is input to every node of next layer.

\rightarrow MLP \rightarrow fully connected feed-forward network.



$$X = [x_1 \dots x_n] \quad ; \quad y = [y_1 \quad y_2]$$

\rightarrow Every node in a layer has a 'w' vector.

node \rightarrow w_i vector.

$$\text{Eq: } \text{Layer 2} \rightarrow \begin{cases} w_1 | n_1 \\ w_2 | n_2 \\ w_3 | n_3 \end{cases}$$

L Weight Matrix
for Layer 2.

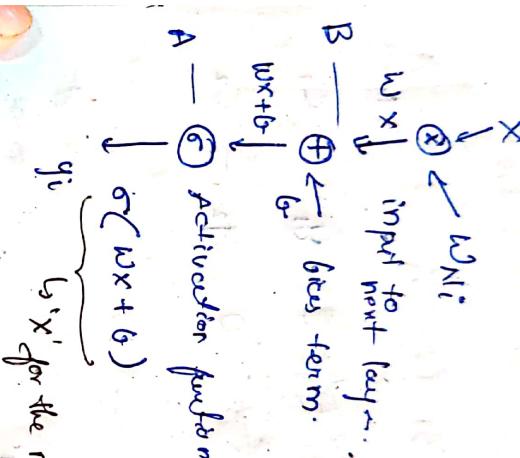
\rightarrow Now every node has a bias term.

Sgt max :-

$$P(c_i | x) = e^{y_i} / (e^{y_1} + e^{y_2})$$

Forwarded Propagation :-

Computation graph



gradient Descent :-

→ It was mathematically proved that even with back propagation.

one hidden layer, any mathematical function can be realized.

$$\frac{\partial y_i}{\partial x} = \frac{\partial y}{\partial A} \cdot \frac{\partial A}{\partial B} \cdot \frac{\partial B}{\partial x} \quad \rightarrow \text{By Chain Rule.}$$

$$\frac{\partial y}{\partial A} = \sigma'(s) = \sigma(s)(1 - \sigma(s))$$

$$\frac{\partial A}{\partial B} = \frac{\partial (wx+b)}{\partial B} = 1$$

$$\frac{\partial B}{\partial x} = \frac{\partial (wx)}{\partial x} = w.$$

$$w^{t+1} = w^t - \eta \nabla E.$$

$$\rightarrow \frac{1}{2} (t_i - y_i)^2$$

$$y = \sigma(\sum_{j=1}^n w_j x_j + b)$$

$$\Delta E = \sum_{i=1}^n \delta_i x_i$$

for output node :-

$$\delta_i = (t_i - y_i) \sigma'(s_i)$$

for hidden nodes:-

$$\delta_i = (\sum_j w_j \delta_j) \sigma'(s_i)$$

→ See slides for perfect derivation.

Back Propagation :- (Stochastic case) :-

→ Initialize all weights to small numbers.

→ for each training ex

compute network outputs.

→ for each output unit b :

$$\delta_b \leftarrow o_b (1 - o_b) (t_b - o_b)$$

→ for each hidden unit i

$$\delta_h \leftarrow o_h (1 - o_h) \sum_{k} w_{ki} \delta_k$$

→ for each weight

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

$$\text{where } \Delta w_{ji} = \eta \delta_j x_{ji}$$