

ANOPP SIR

Books -

Pattern classification - Duda Hartle stock

Neural Networks - Simon Haykin

Read - (Prerequisite)

Prob. distribution func., multivariate densities (joint & conditional)

Vectors, spaces, subspaces, Eigen values

& vectors

Assignments - 24

Project - 20

Mid(2) - 30

Final - 20

Others - 6

Tasks in AI

- sensing
- Recognition
- Reasoning
- Acting

Main focus on Recognition -

- ↳ Where am I
- ↳ What do I see around
- ↳ Who is this person
- Identification of a pattern as a member of a category we already know or are familiar with
(Like image → text)

- A pattern is the opposite of a chaos ; it is an entity vaguely defined, that could be given a name.

→ Intra class variability

- T T T

Inter class ~ Variability (similarity)

- O D 8 B Z Z

characters that look similar

* In recog. we categorize in class.

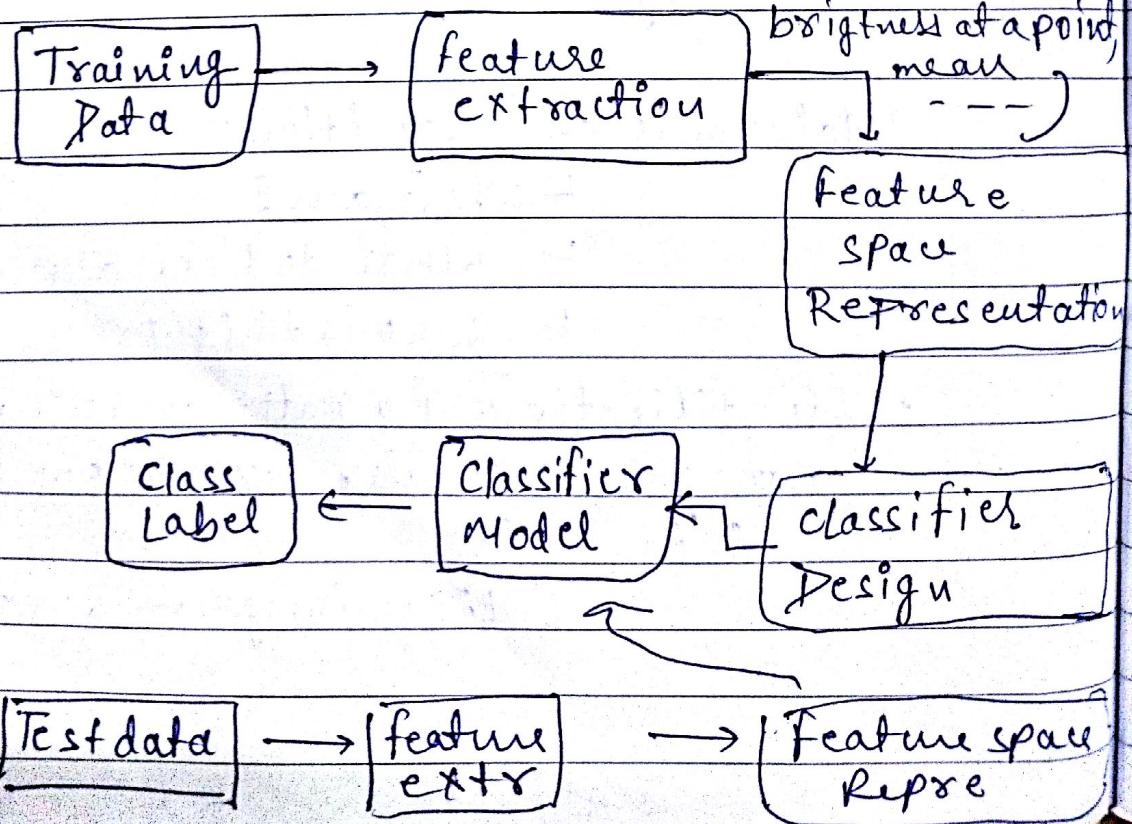
COGNITION - formation of new classes

RECOGNITION - known classes

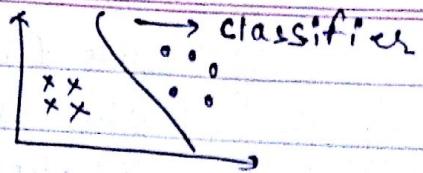
PATTERN RECOGNITION PROCESS :-

(basis of image can be light,

brightness at a point,
mean
---)



Like feature space -



08-08-2014

Note:- 22nd Aug. class test

~~Model~~

signal

(can be email,
speech, image--)

Representation
(feature space) $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ (n dimensional
features)

+++ | --
+++ | -- learning
(Training)

{figure out what
separates the classes}

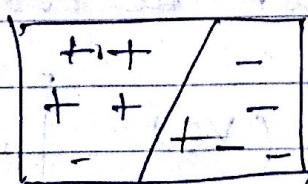
** The process of
Raw data \rightarrow
Representa.
is the soul of your
system

testing

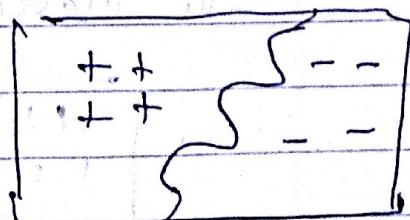
(classification)

I). LEARNING -

How the separator Look Like



or



or any other

We use the concept of Generalisation
 [the generalised curve will have few mistakes but simple representation)

→ 1. Generative Model * (You can generate a new model
 $M(\text{Apple}) \text{ or } M(\text{orange})$)
 Probability of fruit being orange
 being apple

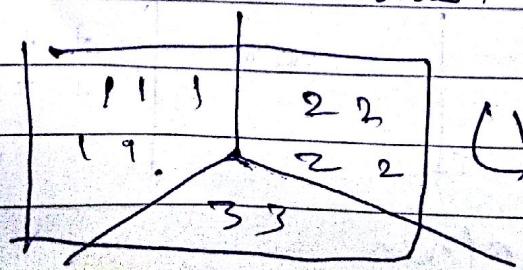
(Note here that our universe is limited to Apple + Oranges)

2. Discriminative Model

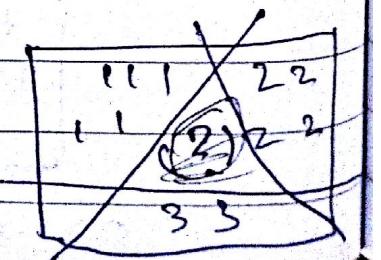
Here we are only interested in discriminating the models (what is the classifier) rather than being generative.

Note: We can combine models also like we can use Pro. of generative models to design a discriminative ones.

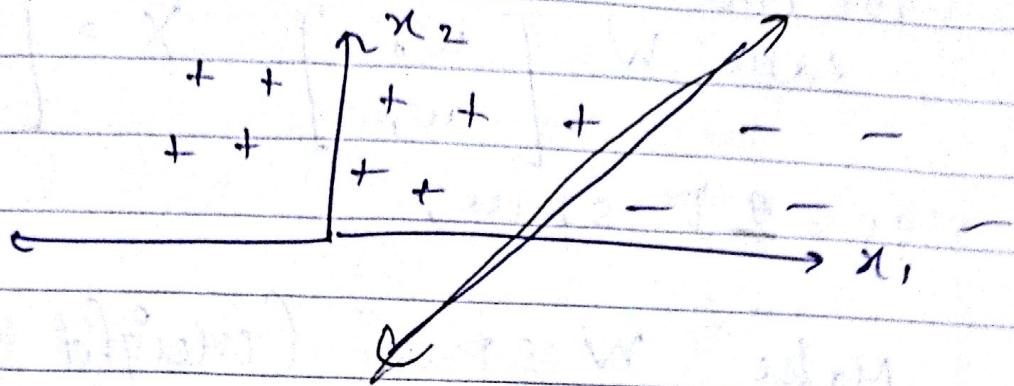
The classifier must be good



(vs)



eg - In 2D space



$$\text{Line} = w_1 x_1 + w_2 x_2 + w_0 = 0$$

$$\therefore W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad \text{and given a new sample}$$

$$X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

(added 1 for convenience)

We need to put point in eqn, for samples, compute \therefore find W such that

$$W^T X > 0 \quad \text{if point is +ve}$$

$$W^T X < 0 \quad \text{if point is -ve}$$

Find W ?

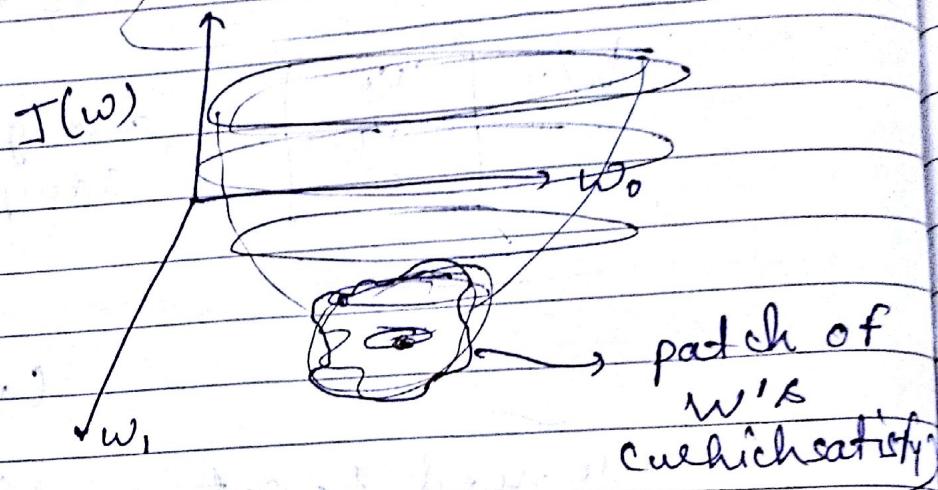
W = classifier

X = feature vector

for find W (best)
take $W = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$ $X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$

for 2D space,

Make W space (weight space)



We have a function $J(w)$ such that when its value \downarrow it comes nearer to patch

$$\text{Make } a = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix}$$

$a^T y = 0$, For samples $y_0 = 1$ (if +ve class)

$y_0 = +1$ (if -ve class
and negate all)

Now by doing this

$a^T y > 0$ for all samples (y)

Now we need to find $J(w)$ such that its min. \rightarrow Patch

We randomly choose ' a '

and make set \mathcal{A} where the values are misclassified

(based on that)

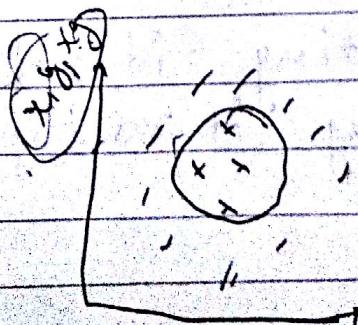
$$J(w) = \sum_{y_i \in \mathcal{A}} -a^T y_i$$

If value of $J(w)$ is very high \rightarrow Not a good classifier
so if $J(w) = 0 \rightarrow$ Yes it is

** So after randomly choosing an ' a ' and computing ' $J(w)$ ' move in the direction of gradient descent to reach good ' a '

Note that ' a ' and ' y ' are vectors

'Vector Algebra', differentiation etc

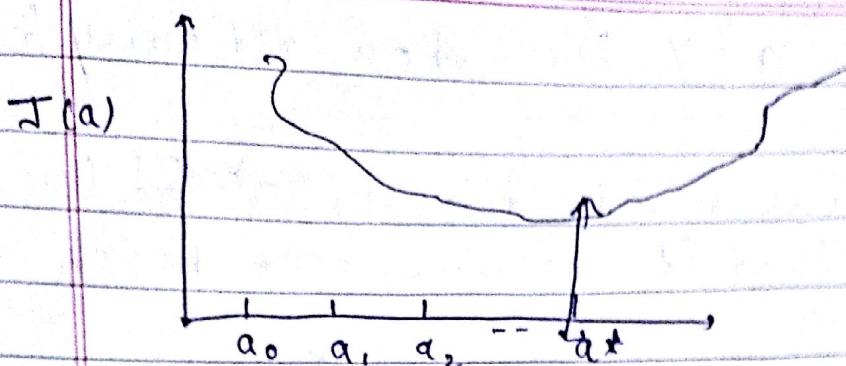


12th Aug - 2014

classmate

Date _____

Page _____



We need to move in the direction of a for which slope is maximum. (It's opp. dir)
[it is a vector as it gives direction]

$\nabla J(a)$ is pointing towards left
we need to move in descent dir.

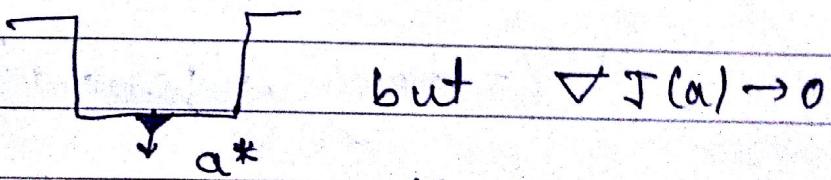
$$\underline{a^{(k+1)} = a^{(k)} - \eta \nabla J(a)_k}$$

η = Learning rate
choose fast you want
to learn)

Loop above eqn until $|\eta \nabla J(a)| < \theta$

$$\theta \approx 0.0 - 1$$

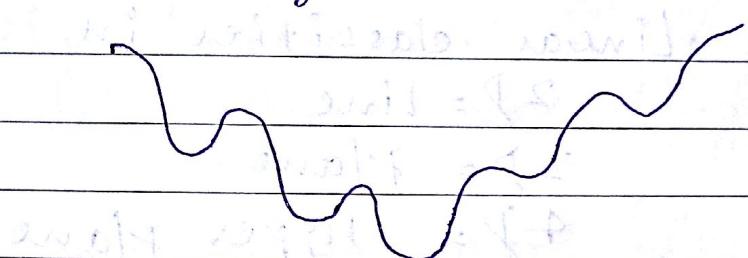
* $J(a)$ doesn't work always as



in the starting itself, however
Practically it rarely happens.

- if η is constant throughout we might oscillate hence η_k is used; learning is fast in the beginning and reduces as we approach the solution.

- Problem might arise



we might get stuck in a local minimum.

Try out differentiable $J(a)$ functions

1. Perceptron

$$J_p(a) = \sum_{y_i \in Y_{mc}} - (a^T y)$$

$$\nabla J_p(a) = \sum_{y_i \in Y_{mc}} - y$$

Dif.. wrt 'a' [Read Vector algebra]

So algorithm \rightarrow

```

     $a_0 = \text{init}$ ,  $k=0$ 
    do  $k = k+1$ 
         $a(k) = a(k-1) + \sum_{y_i \in Y_{mc}} y$ 
    until  $Y_{mc} = \emptyset$ 
  
```

"Batch Perceptron Algorithm"

2. Single Sample Perceptron Algorithm:-

$$[a_{k+1} = a_k + y^{k+1}]$$

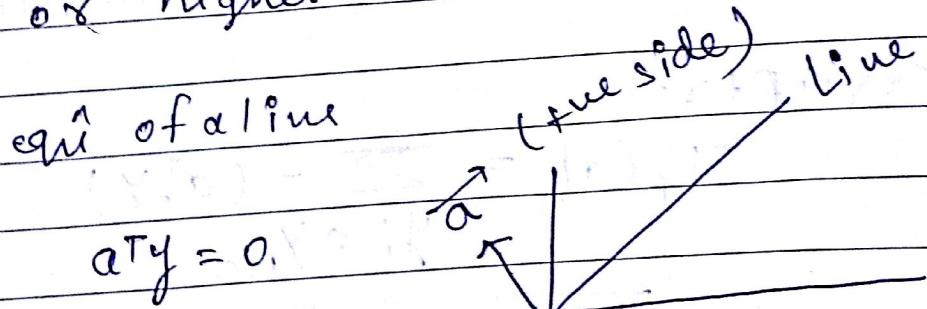
where y^{k+1} = any one misclassified sample [It is not power]

→ Linear classifier in
 $2D = \text{Line}$

$3D = \text{Plane}$

$4D = \text{Hyper Plane}$

and we will use Hyper planes
 as a generalised plane for lower
 or higher dimensions also

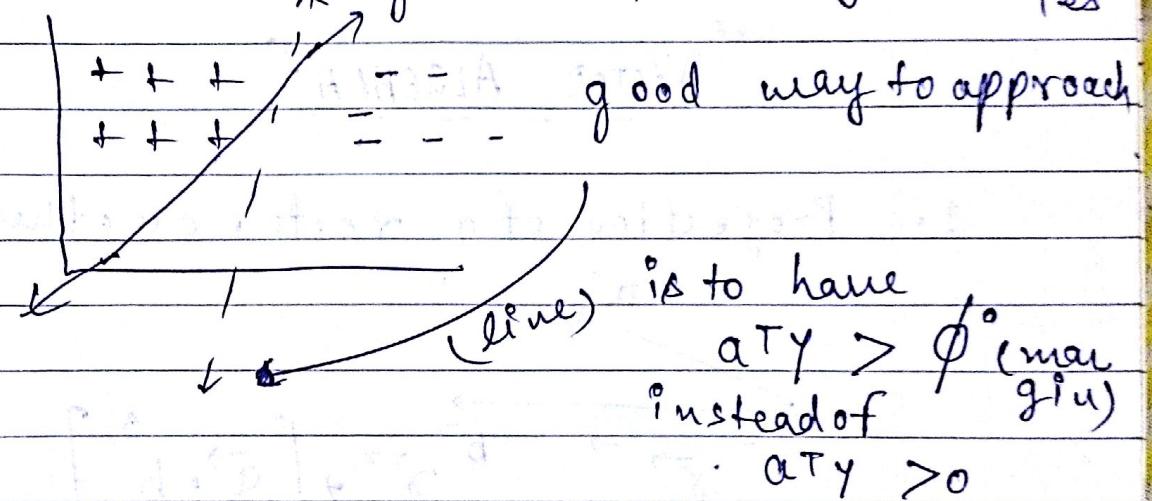


can be thought as \vec{a} passing
 through origin and $\vec{1}$ to line

[Note that for above method we have
 assumed that such an \vec{a} exist
 which linearly separates the
 samples..]

- By this method all α^* 's are moving towards α^* (Perfectly right)
 $\alpha = \text{constant, may not be 1}$
 ie- Magnitude of difference keeps on decreasing..

⇒ Note that by this algorithm we might find a line just separating the samples



or you can use different J 's

$$J_{\text{smooth}} = \sum_{y \in Y} (\alpha^T y)^2 \quad \begin{cases} \text{smooth} \\ \text{fn, } \end{cases}$$

→ slightly preferable.

or use

$$(a^T y - b)^2 \quad \text{will take care of slope handling issues of } (a^T y)^2$$

so we have =)

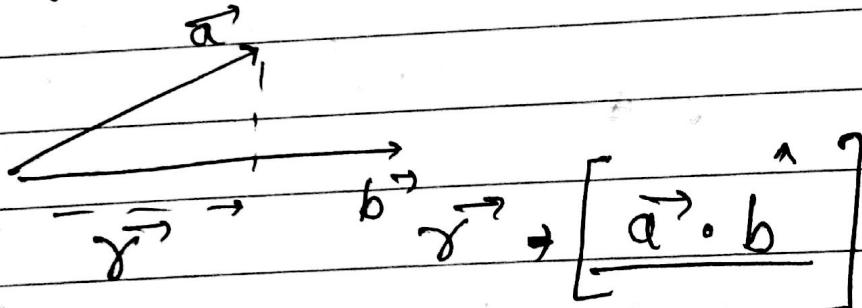
$$J_{\text{reg}}(a) = \frac{1}{2} \sum_{y \in Y_{\text{MC}}} (a^T y - b)^2$$

$$\nabla J_{\text{reg}}(a) = \sum_{y \in Y_{\text{MC}}} \frac{a^T y - b}{\|y\|^2} \cdot y$$

(Relaxation algorithm)

"VECTOR ALGEBRA"

1. Projection of a vector on other



2. A. $\hat{f(u)}$ is linear if

$$\therefore f(\alpha a) = \alpha f(a)$$

$$\bullet f(u+y) = f(u) + f(y)$$

So

$$f(\alpha u_1 + \beta u_2) = \alpha f(u_1) + \beta f(u_2)$$

- for making the classifier, select some subset of labeled samples to make the classifier and use the rest for checking the authenticity of the classifier (pt's accuracy)

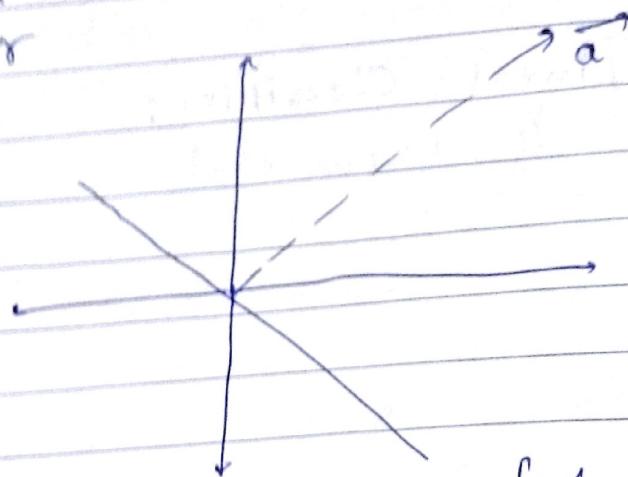
3. Linear Classifier:-

You know that.

19-Aug-2014

Note: $b = \mathbf{a}^T \mathbf{y}_i$; "b" is the perpendicular distance between \vec{a} and \mathbf{y}_i

Very IMP. \Rightarrow If \vec{a} is a solution vector
Then $\alpha \vec{a}$ is also a solution vector



so when you are doing $\{\forall_i \mathbf{a}^T \mathbf{y}_i \geq b\}$
b can be any number
you can have a solution for

$\forall_i \mathbf{a}^T \mathbf{y}_i \geq 1$, and multiply
solution \vec{x} with 'b' to have

$b\vec{x}$ as another sol'n satisfying
 $\forall_i \mathbf{a}^T \mathbf{y}_i \geq b$, [we will find]

a solution always because if \vec{x} is a
sol'n then $\beta\vec{x}$ is also]

But for $[\forall_i \mathbf{a}^T \mathbf{y}_i = b]$, probably
we can't find a sol'n.

To prove that $\forall i; a^T y_i = b_i$ may not give a soln write it as.

(each row
is a
+ training
solution)

$$\left[\begin{array}{c} y_{10} \\ | \\ | \\ | \\ y_{no} \end{array} \right] - \left[\begin{array}{c} y_{11} \\ | \\ | \\ | \\ - \end{array} \right] - \left[\begin{array}{c} y_{id} \\ | \\ | \\ | \\ - \end{array} \right] = \left[\begin{array}{c} a_0 \\ | \\ | \\ | \\ a_n \end{array} \right] = \left[\begin{array}{c} b_1 \\ | \\ | \\ | \\ b_n \end{array} \right]$$

$$Y a = b$$

Now if we have $\forall i; a^T y_i = b_i$,
you can find the solution as \vec{a}

by

$$\boxed{Y a = b}$$

Y = Augment feature vector Matrix
Now if

$n > d + 1$, overdetermined \Rightarrow no soln

$n = d + 1$, unique soln

$n < d + 1$, \Leftrightarrow solution,
underdetermined

\Rightarrow So basically $[Y]$ has to be a
square matrix $[n = d + 1]$

Rows in Y has to be linearly independent
for unique soln.

$$\boxed{a = Y^{-1} b}$$

So if we want to find $\{\forall i; a^T y_i = b\}$
in 2D space

$$\Rightarrow d = 3, n = [3+1]$$

i.e. It will produce a limit on
number of points (samples).

Check and Read Book for

$$\boxed{Y_a = b}, \{n = d+1\}$$

**

But in reality

$[n >> d+1]$, so we need to
look into this, so instead of having

$$Y_a = b \quad \dots \dots (1)$$

Take

$$e = Y_a - b \quad \dots \dots (2)$$

and minimise the error 'e'

$$(e = Y_a - b)$$

$$\min (||e||^2)$$

so we need to minimise

$$||e||^2 = \sum_{i=1}^n (a^T y_i - b_i)^2$$

so

$$J_e(a) = ||e||^2 = \sum_{i=1}^n (a^T y_i - b_i)^2$$

$\{$ This is over all the samples $\}$

So we need to minimize $J_s(a)$ and find a unique solution.

at the minimum of a function its gradient is '0'.

$$\text{so } \nabla J_s(a) = 0 \quad \begin{cases} \text{Intuitively it's a min. point} \\ \text{as map is } \partial \end{cases}$$

$$\nabla J_s(a) = 2 \sum_{i=1}^n (a^T y_i - b_i) y_i = \vec{0}$$

Here $\vec{0} = \vec{0}$ as it is (constant) y_i

so it says at minimum we have



$$2Y^T(Ya - b) = \vec{0}$$

$$Y^T(Ya - b) = \vec{0}$$

$$Y^T Y a - Y^T b = 0$$

$$Y^T Y a = Y^T b$$

$(d \times n) \quad (n \times d)$

So : \rightarrow

$$a = (Y^T Y)^{-1} Y^T b$$

and $Y^T Y$ = square Matrix ($d \times d$)
solution for ($n > d + 1$)

Good Thing →

when $n = d$, $[a = Y^{-1} b]$

and $n > d$, $[a = \boxed{(Y^t Y)^{-1} Y^t b}]$
 $\downarrow +$
 Y^t (pseudo inverse)

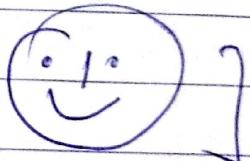
So: →

$$\boxed{\begin{array}{c} a = Y^+ \cdot b \\ (d \times 1) \quad (d \times n) \quad (n \times 1) \end{array}} \quad \text{unique sol'n}$$

(least square

Y^+ is called the left pseudo inverse
 as $[y^t y = I]$ (sol'n)
 and not $[yy^t \neq I]$

— [too Good sol'n



$$\Rightarrow \left\{ \begin{array}{l} Y_1 = [\text{ones}(\text{size}(X_1, 1), 1) \ X_1]; \\ Y_2 = -1 * [\text{ones}(\text{size}(X_2, 1), 1) \ X_2]; \end{array} \right.$$

$$Y = [Y_1' \ Y_2']'$$

$$b = \text{ones}(\text{size}(Y, 1), 1);$$

$$[a = (Y' * Y) * Y' * 1 * b]$$

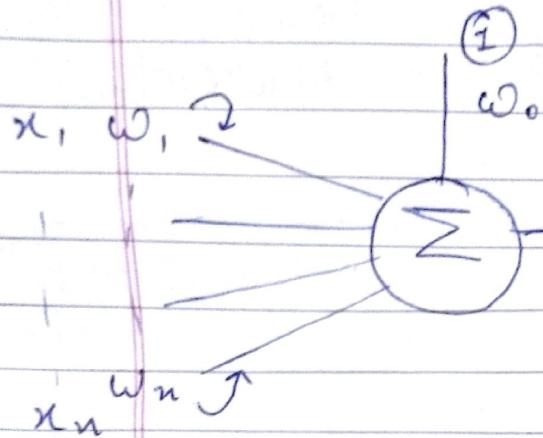
→ How come scientists know size of universe, distance to a distant galaxy, or even know that there is another galaxy (because probably light has not reached there if send from here) — find out...

2/Sept/2014

Generative Models:-

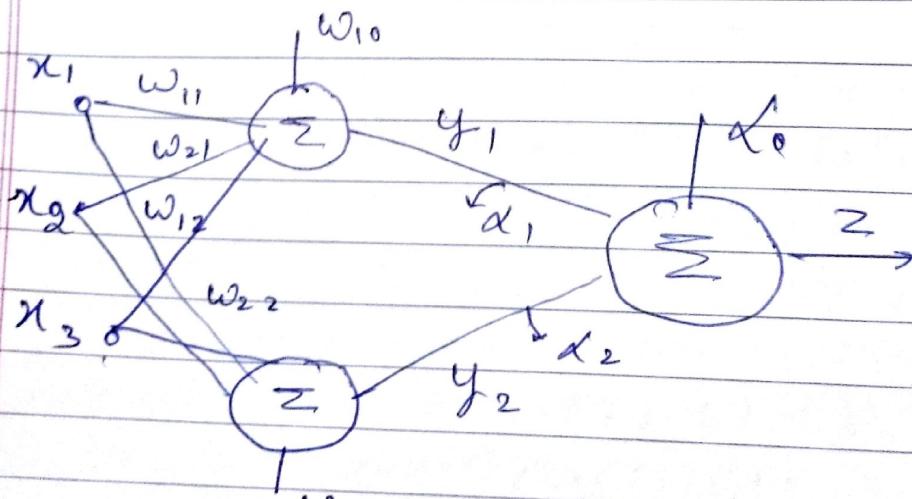
• Perceptron

$$\sum w_i x_i + w_0$$



Resembles something like brain.
(neurons)

(IT is a model of
a Linear classifier)



$$y_1 = \sum_{i=1}^d w_{1i} x_i + w_{10}$$

$$y_2 = \sum_{i=1}^d w_{2i} x_i + w_{20}$$

$$z = \alpha_1 y_1 + \alpha_2 y_2 + \alpha_0$$

So above thing when expanded

$$\mathbf{f}_x \rightarrow \mathbf{d}$$

$$\sum_{i=1}^n w_{new; i} x_i + w_{new_0}$$

(same eqn as LC but with different -t weights)

- Honey bee

- small brain

(few 1000 neurons)

with such a small number of them,
they are able to do very-very
complex tasks.

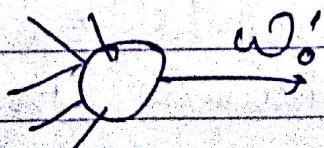
- our brain

- billion neurons.

- Trillion connections.

- So since it is found that every 100 neurons can down to, we are doing a linear combination of linear things \Rightarrow we get linear

To introduce non-linearity



w' + output =

$$y = f(w')$$

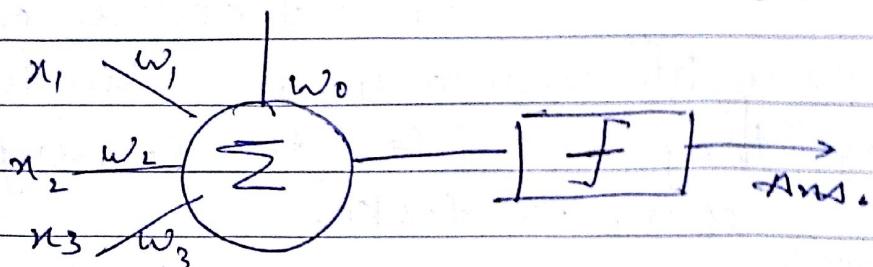
So we can have a variety of functions at the terminal.

Simplest Non-linearity.

= step function



(Thresholding)



(Simplest NL Model)

This small Non-linearity can solve many complex problems and there is a theory that shows that a (2 Layer or 3 Layer) Neural

Networks, can model any complex function [classifiers].

Theoretically

→ All problems of world can be solved using NN, but we have to train it...

"NAV LAB"

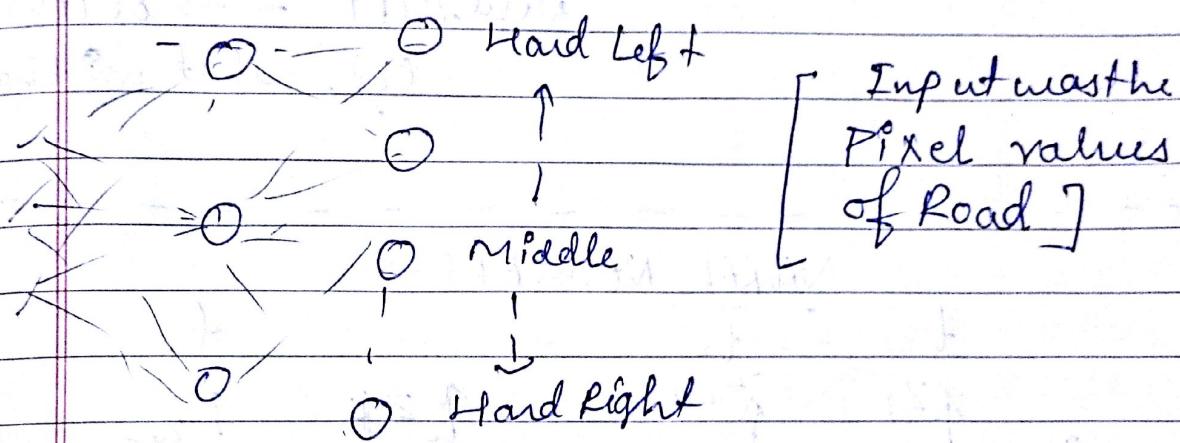
classmate

Date _____

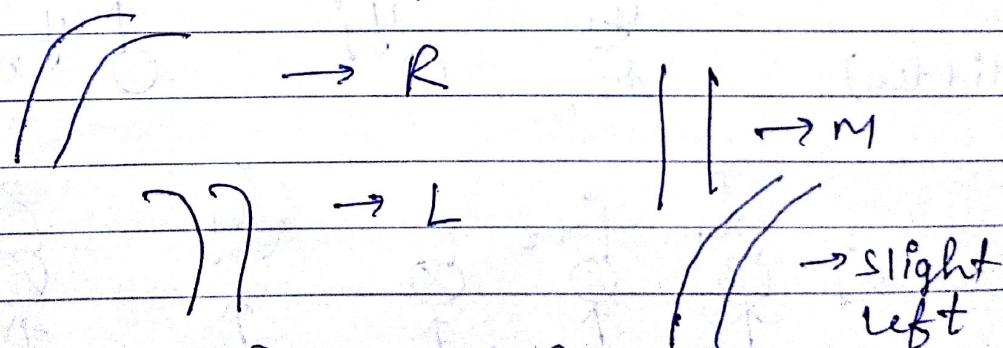
Page _____

- CMV, self driving car
 - ↳ Train the car how to drive rather than creating features, and doing a lot of things.

↳ Parameters → Accelerate
Brake
Turn Left, Right



- a large set of Input - output



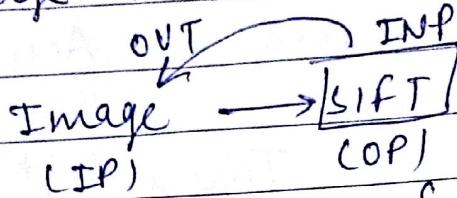
(so drive the car for sometime
NN will train and within 2 hrs
it will be trained to do so)

→ weights will be determined by
P & eff.

- Teaching a NN with more no. of layers is easier, than a 2 layer.

click out

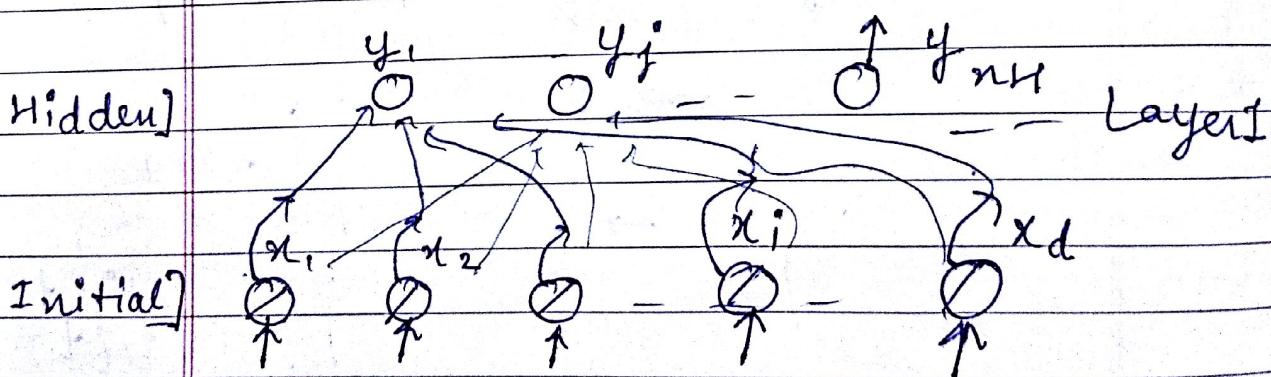
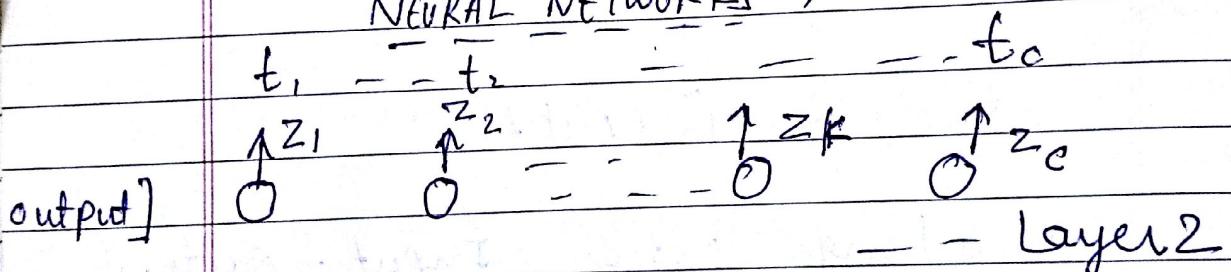
→ App. of NN in IP, [line filtering]
(Image Reconstruction from SIFT)



neural SIFT → OUTPUT?

(Q), what will happen

NEURAL NETWORKS →



'i' → Input layer $m \times n$ = number

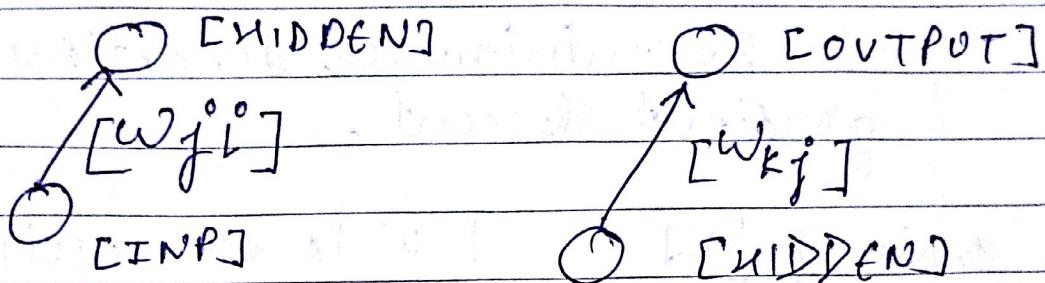
'j' → Hidden layer of nodes in
hidden layer

'k' → output layer

$nH \geq d$ or $nH < d$ both are possible.

'd' dimensional feature vector.

weight in general represented as →



"Target output is the classes"

Ideally $0, 1, 0 \dots 0$ should be the output but here we will talk about Probabilities,

Neural Networks Apps -

- Driving cars.
- function approximations.
- more appli other than classifications.

Sometimes we do have a biased input node.

Suppose for training t_1, \dots, t_c are values required. Then we will require

$$J = \sum_{k=1}^c (t_k - z_k)^2 \times \frac{1}{2}$$

to be minimised, direction of gradient descent.

$$\Delta w = -\eta \frac{\partial J}{\partial w} \quad \left. \begin{array}{l} w \text{ is the weight} \\ \text{vector} \end{array} \right\}$$

for hidden layer — [1 node]

$$\left. \begin{array}{l} \textcircled{1} \quad \text{net}_j = \sum_{i=1}^d w_{ji} x_i \\ \textcircled{2} \quad y_j = f(\text{net}_j) \\ \textcircled{3} \quad \text{net}_k = \sum_{j=1}^{n_H} w_{kj} \cdot y_j \\ \textcircled{4} \quad z_k = f(\text{net}_k) \end{array} \right\}$$

$$\frac{\partial J}{\partial w_{ij}} = \frac{\partial J}{\partial w_0} \cdot \frac{\partial w_0}{\partial w_{ij}} \quad \left. \begin{array}{l} \text{chain rule.} \\ \text{rule.} \end{array} \right\}$$

so J is a function of a function of a function - - -

a) If a weight between OUT \rightarrow HIDDEN
is changed only that node will
be effected.

$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial w_{kj}}$$

$$= \underbrace{-[t_k - z_k] \cdot f'(\text{net}_k) \cdot y_j}_{\delta_k}$$

$$\delta_k = (t_k - z_k) \cdot f'(\text{net}_k)$$

= (sensitivity of node) — we call

= [it captures the information
that how will the error
change if weights w_{nj} - (net_j)
is changed for that node]

$$\delta_k = -\frac{\partial J}{\partial \text{net}_k}$$

and $\Delta W_{kj} = - \frac{\partial J}{\partial z_k}$

$$\Delta W_{kj} = \eta \delta_k y_j \quad \text{--- Layer 2}$$

$$\delta_k = (t_k - z_k) \cdot f'(net_k)$$

so pass (forward pass) ones and then find Δw_{nj} , update the weights

$$w_{kj}^{(m+1)} = w_{kj}^{(m)} + \Delta w_{kj}$$

b) - $\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}}$

None if the inputs weight are changed
The output at both the forward nodes will change.

Write \Rightarrow

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \cdot \frac{\partial y_j}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}}$$

(x) \downarrow \downarrow
 $f'(net_j) \cdot x_i^*$

$$\frac{\partial J}{\partial y_j} = \frac{\partial}{\partial y_j} \left[\frac{1}{2} \sum_{n=1}^k (t_n - z_n)^2 \right]$$

$$\delta_j = - \sum_{k=1}^C (t_k - z_k) \underbrace{\frac{\partial z_k}{\partial y_j}}_{\frac{\partial z_n}{\partial \text{net}_n} \cdot \frac{\partial \text{net}_k}{\partial y_j}} \underbrace{f'(\text{net}_n) \cdot w_{kj}}$$

Hence $\frac{\partial J}{\partial y_j} = - \sum_{k=1}^C (t_k - z_k) \cdot f'(\text{net}_n) \cdot w_{kj}$

$$\frac{\partial J}{\partial y_j} = - \sum_{n=1}^C (t_n - z_n) \underbrace{f'(\text{net}_n) w_{nj}}_{\delta_k}$$

Hence finally -

$$\frac{\partial J}{\partial w_{ji}} = - \sum_{k=1}^C \delta_k w_{kj} \cdot f'(\text{net}_j) \cdot x_i$$

$$\boxed{-\delta_j = -\frac{\partial J}{\partial \text{net}_j}}$$

\hookrightarrow sensitivity of Hidden Node.

Always $\Delta w = \eta \text{ Gradient}$

classmate

Date _____

Page _____

$$\text{So } \Delta w_{ji} = \eta \delta_j x_i$$

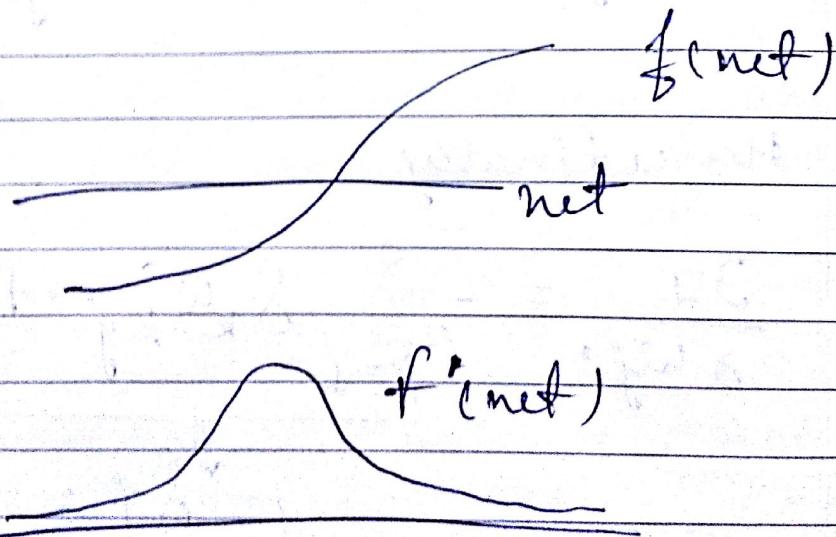
$$\delta_j = \sum_{k=1}^c \delta_{w_{kj}} \cdot f'(\text{net}_j)$$

layer 1

PROBLEMS:-

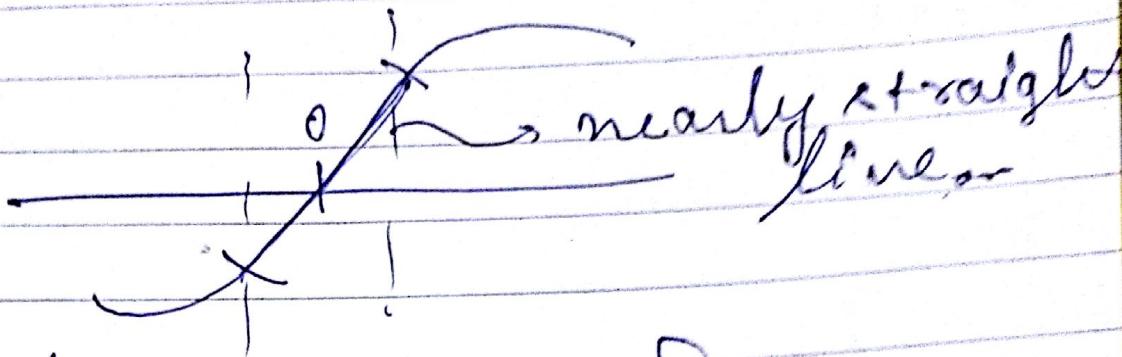
i) Activation function f :
↳ like step fun.

The function should
be differentiable for whole pro-
cess to work.

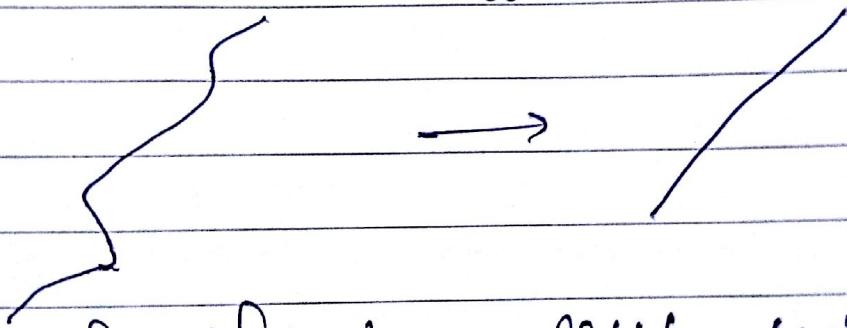


i) also $f'(\text{net})$ should not be
zero higher number of times
since update will be zero.

iii)- we would like to have functions like



If function is like this our curve will become less complex



for this do a little trick, in each update do

$$\text{weight decay}, \Rightarrow [w = 0.95 w]_{\text{New}}$$

iii)- t_n should be determined based on ^{update}

$f(\text{net } n) - \begin{cases} \text{Ranges, do not pick} \\ \text{their asymptotically higher values} \end{cases}$