

# Machine Learning

## Lect 1 : Introduction

Dr. PRERANA MUKHERJEE

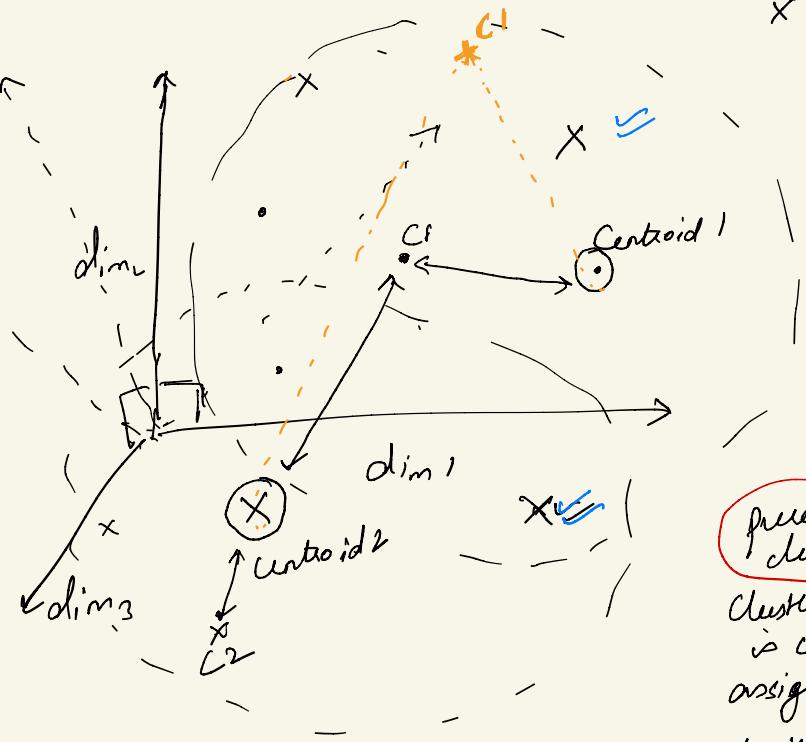


Training samples  $\downarrow$   $\downarrow$   $\downarrow$  features  
 $\rightarrow \mathcal{X}_1 = \{d_1, d_2, \dots, d_n\} \in \mathbb{R}^n$   
 $\rightarrow \mathcal{X}_2 = \{d_2, d_4, \dots, d_n\}$   
 $\rightarrow \mathcal{X}_3 = \{d_3, d_4, \dots, d_n\}$   
 $\vdots$   
 $\mathcal{X}_n = \{d_5, d_6, \dots, d_n\}$

wiki  
 {Distance metrics  
Unsupervised algorithms

k-means Input "k" = # clusters

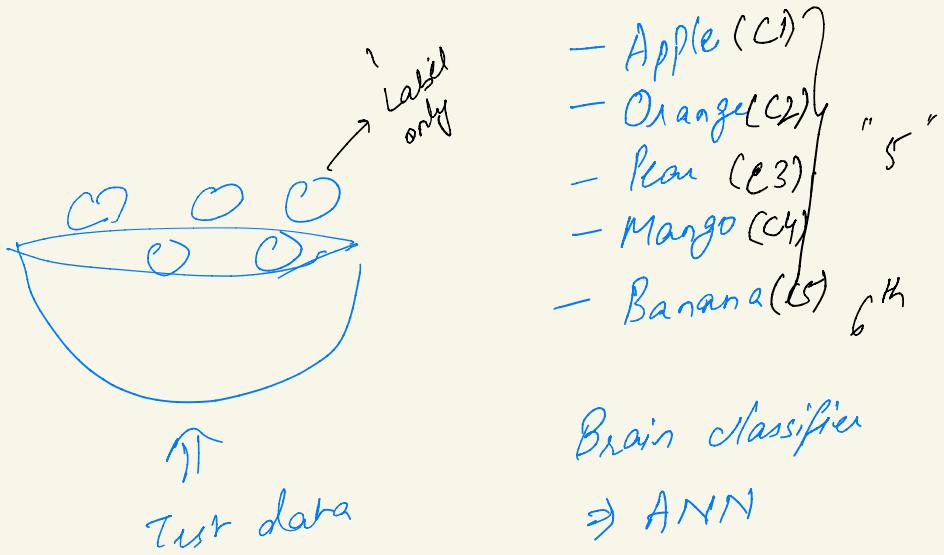
- Cluster 1
- Cluster 2



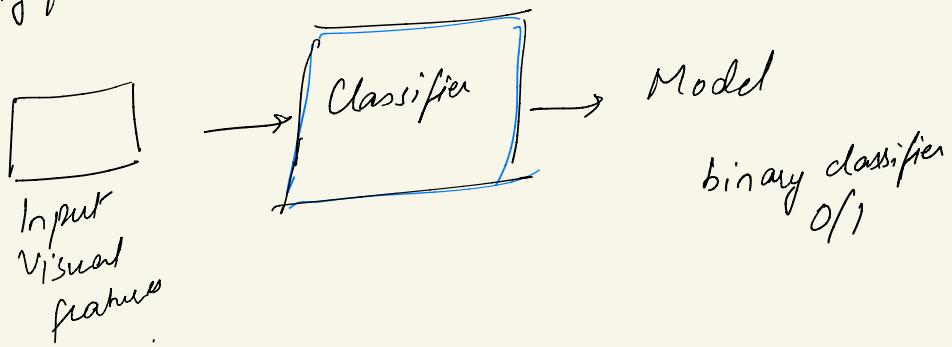
stopping criteria  $\Rightarrow \checkmark$

cluster 1 is correctly assigned  
 cluster 2 is correctly assigned

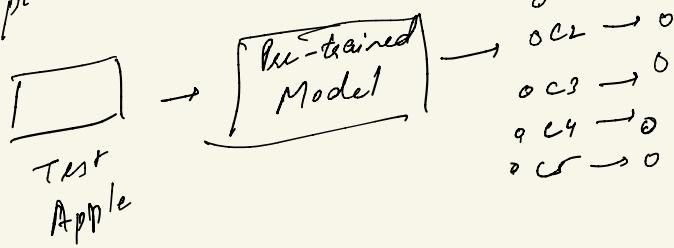
Supervised Setting  
Train a Classifier  $\Rightarrow$  SVM, kNN, Bayes, DNN, ...  
CNN, LSTM

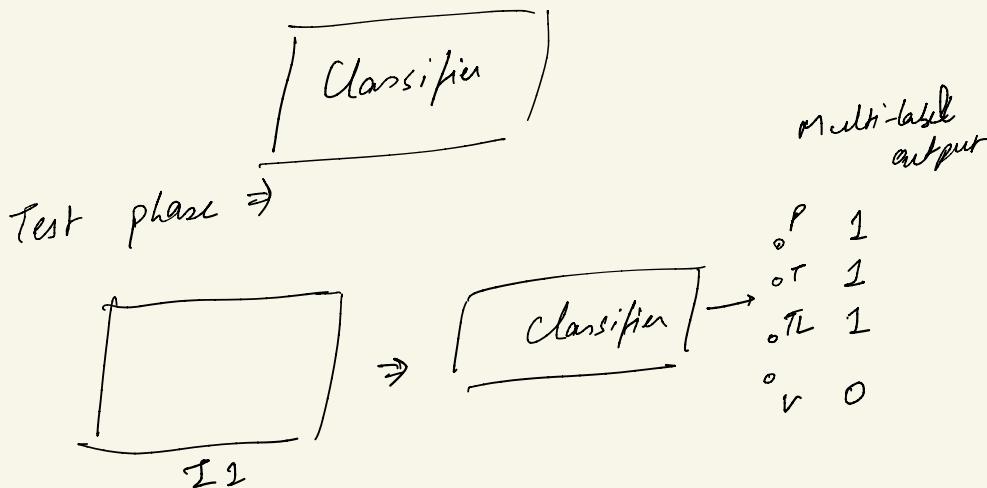
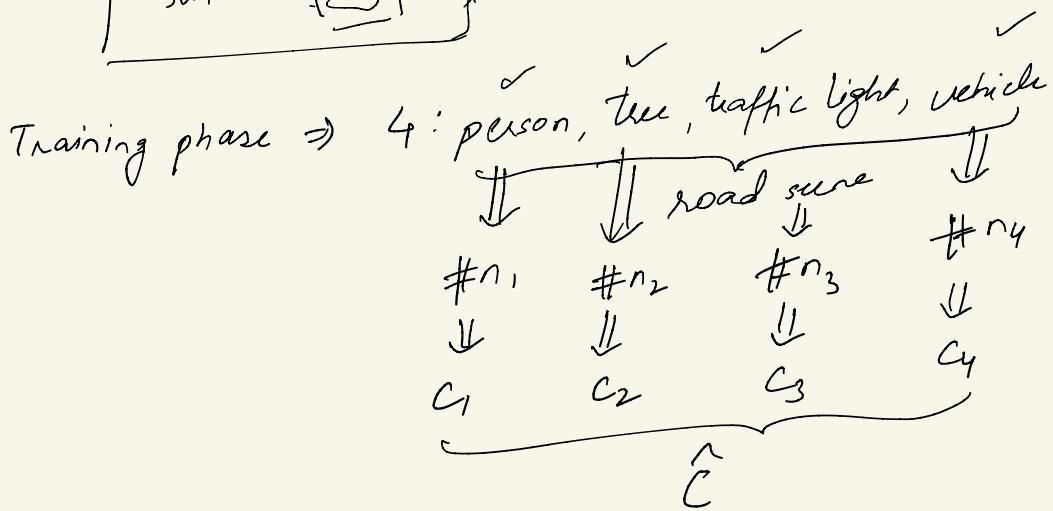
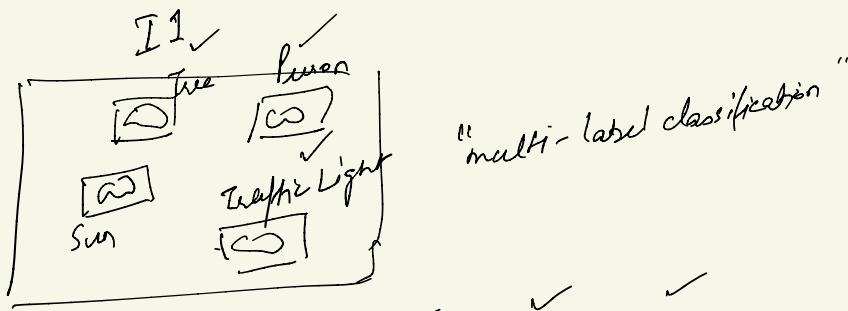


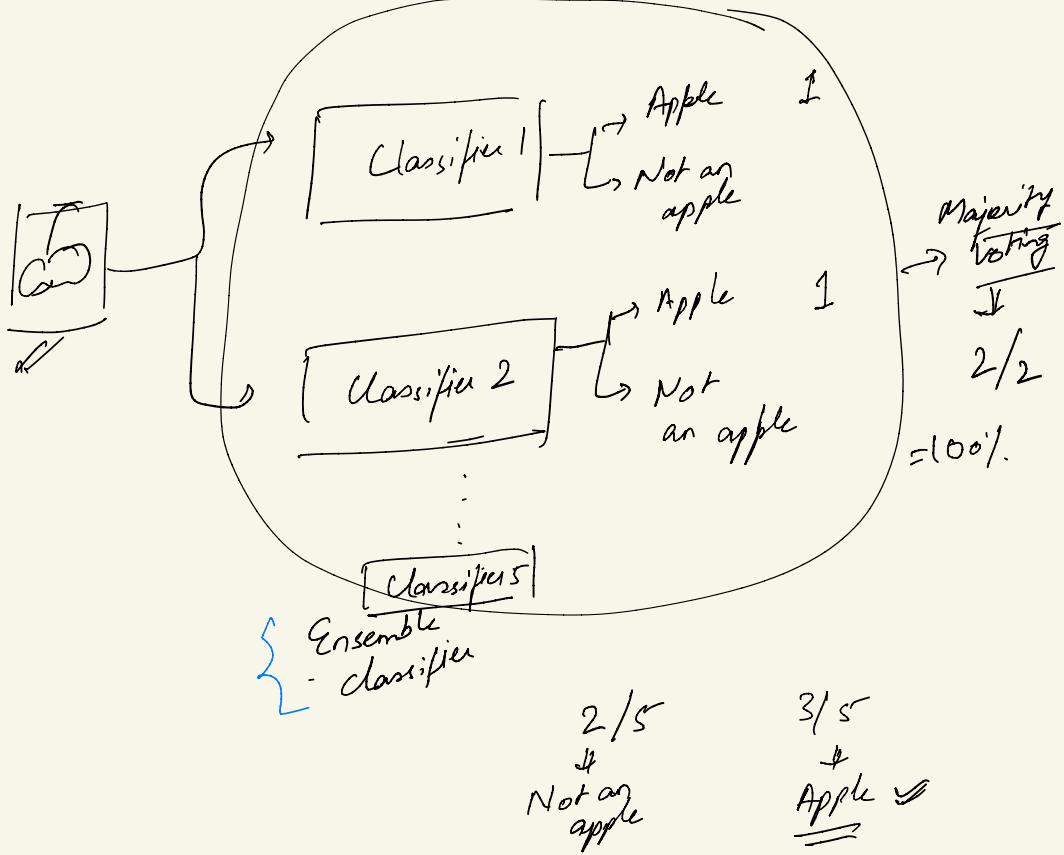
Training phase:



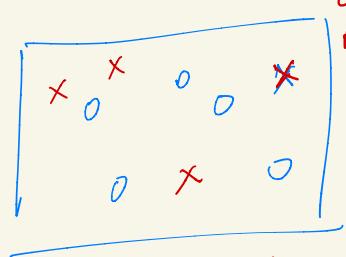
Test phase:







Semi supervised setting  $\Rightarrow$  EM, R.L.



Expectation maximization  $\Rightarrow$  Reinforcement learning

$\Rightarrow$  Regression  $\Rightarrow$

Linear +  
Non-linear

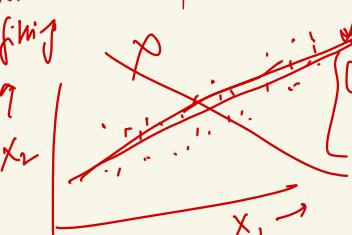
$\Rightarrow$  Classification

Unsupervised  
Supervised  
Semi-supervised



Curve  
fitting

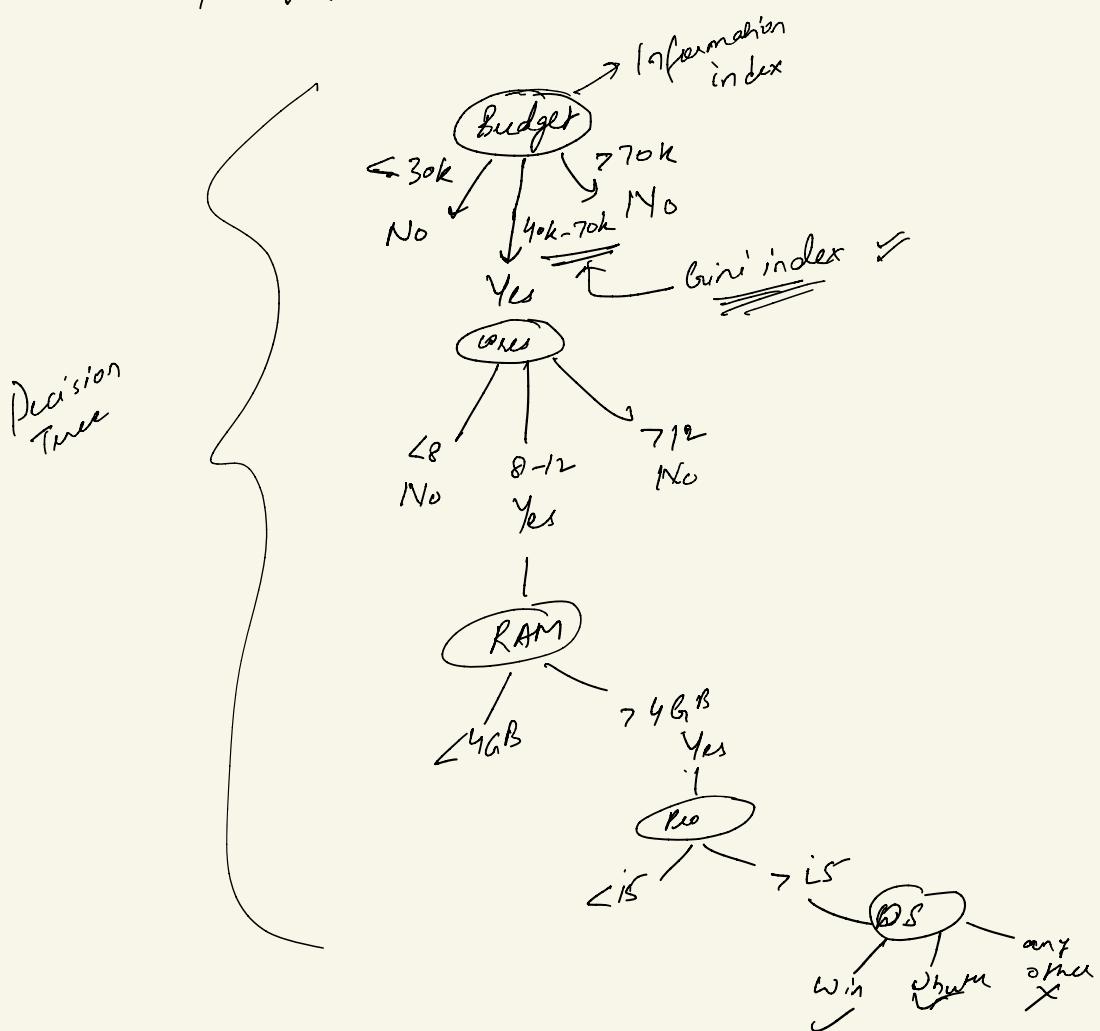
$y$   
 $x_2$



Curve fitting  
line fitting

## # Purchase of Laptop

- Budget ( $30k < \text{INR} < 70k$ )
- cores of processor (CPU) ( $8 < \# \text{cores} < 12$ )
- RAM ( $> 4GB$ )
- Processor generation ( $> i5$ )
- Operating System (Windows 10, Ubuntu)





## Machine Learning Lect -2

### Introduction to Neural Networks

DR. PRERANA MUKHERJEE

---

---

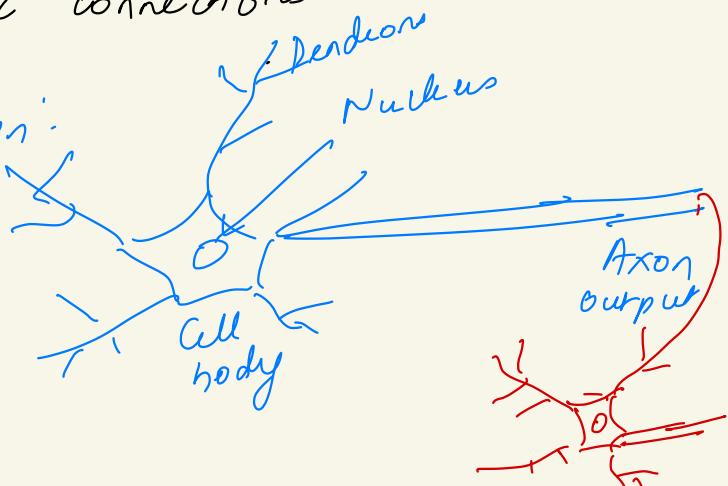
---

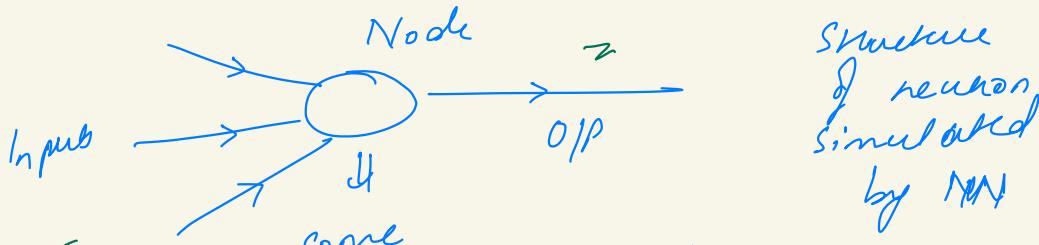


# Introduction

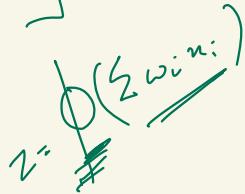
- Neural Networks
  - ↳ Inspired by how human brain works & maps the inputs to logical outputs.
- Brain is the most complex part of human body & consists of 10 billion + neurons
- Neurons are highly interconnected & <sup>act as</sup> individual units.
- Characteristics of neurons:
  - Massive parallelism
  - Highly interconnected (solve the complex problems of your real world)
  - Model distributed associative memory through weights in the synaptic connections

Diagram of neuron:





some simple computations  
weighted sum of inputs  
& then applies some squashing/  
activation function (could be  
sigmoid or other)



ANN: through synaptic connections O/P meets  
another neuron input using  
electric impulses through the  
synaptic layers.

Incorporate the 2 fundamental components  
of the biological neural nets.

- Nodes - Neurons
- Weights - synapses

Basic unit in NN: Perception

# Perceptron

Inputs :  $n$

$x_1$

$w_1$

$x_2$

$w_2$

$x_3$

$w_3$

$x_n$

$w_n$

$\in \mathbb{R}^n$

bias or  $w_0$

transfer function

$\sum$  |  $\phi$  |  $\int$   
weighted sum

$y (0 \text{ or } 1)$

Training Data  
 $\bar{x}$ ,  $y$ ,  
 $\bar{x}_1, y_1$ ,  
 $\bar{x}_2, y_2$ ,  
 $\vdots$ ,  
 $\bar{x}_m, y_m$ .  
 Training the  
 n/w means learning  
 the weights to have  
 good fit for  
 training data.

$$\phi(z) = z \quad (\text{identity})$$

$$\phi(z) = \begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases} \quad (\text{unit step})$$

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad (\text{sigmoid})$$

$$y = \sum_{i=1}^n w_i x_i + b \quad \text{or} \quad \underbrace{\sum_{i=0}^n w_i x_i}_{\text{bias}}$$

$\phi$  is identity function  $\Rightarrow$  same passed as output

Or

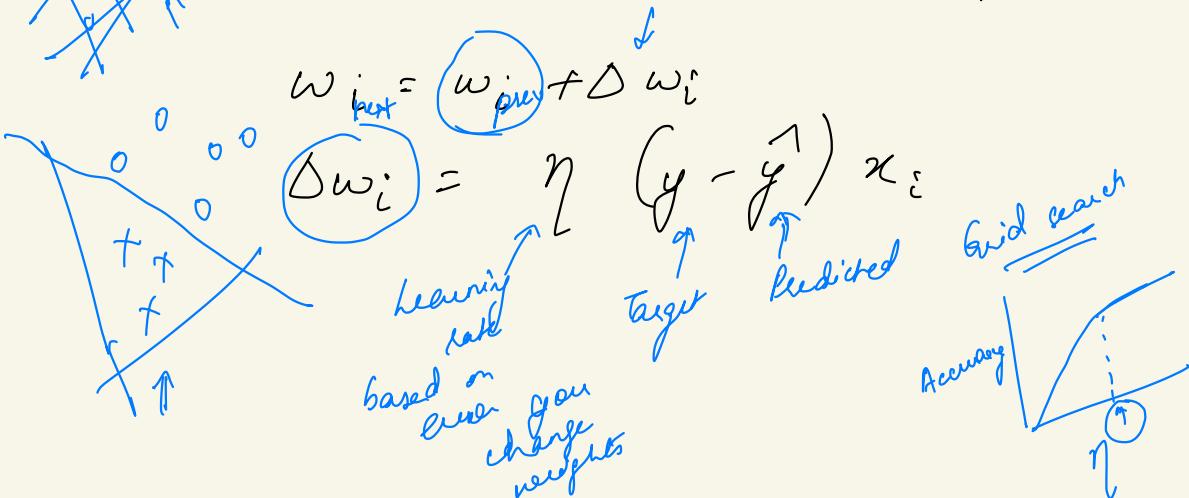
$$\phi_2(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases} \quad (\text{threshold function})$$

## Perceptron learning rules

- we feed 1 eg at a time & update the weights based on output.



if output = = output  
(predicted) (target)  $\Rightarrow$  no update  
in weight required.



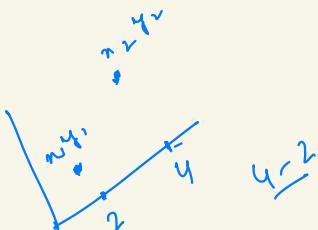
Shopping  $\Rightarrow$  exception learning convergence

Criteria if  $D$  is linearly separable.

(Training data)  $\Rightarrow$  linear decision surface

What if data is not linearly separable?

+/  
separate



Solution : Gradient descent  $\leftarrow$

can be used when w.r.t to particular parameter set / weight values you can define the error of the alg. Perform GD on error func. to find the optimal parameter set for which the error fn. is minimized

$$E = \frac{1}{2} \sum_d (y_d - \hat{y}_d)^2 \Rightarrow \text{is like a surface \& we want to find minima of this surface}$$

so that differentiation is easy



find partial derivative of error fn. w.r.t each weight and find the direction of the gradient & we go towards the -ve direction of gradient in order to go towards minima.



In certain cases, this error surface can be convex or quadratic & there will be single minima  $\Rightarrow$  thus gradient descent guarantees the minima.

Multilayer networks <sup>networks</sup> surface will be ill-behaved (non-convex)

there will be local minima

get stuck

$\phi_1(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$  step function  $\rightarrow$  is not differentiable  $\Rightarrow$  so gradient descent can't be done

$\phi_1(z) = z$  is differentiable  $\Rightarrow$  gradient descent can be done

Stochastic gradient descent

1 eg. from training set at a time

$$E = \frac{1}{2} (y_d - \hat{y}_d)^2$$

Training rule  
of G.D

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = n_{ij} (\hat{y}_j - y_j)$$

Single layer n/w capture linear decisions.  
linear separable data.

Capture non-linear functions  $\rightarrow$  multi-layer n/w's  
 $\rightarrow$  different n/w's connected with each other  
 $\rightarrow$  If linear units are connected

$$\phi_3(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

combination is again linear fn.

$$\hat{p}_3'(z) = \phi_3(z) \cdot (1 - \phi_3(z))$$

Assignment  
Leave it

$$E = \frac{1}{2} \sum_{d \in D} \left( y_d - \underbrace{\sigma(w \cdot z_d)}_{\hat{y}_d} \right)^2$$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_d \frac{\partial E}{\partial \hat{y}_d} \cdot \frac{\partial \hat{y}_d}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = \sum_d (y_d - \hat{y}_d) \cdot \underbrace{\sigma'(w \cdot z_d) \cdot z_d}_{\text{to represent complex functions we want non-linear units which are differentiable}}$$

$\frac{\partial}{\partial w_i} (\underbrace{y_d - \sigma(w \cdot z_d)}_{\hat{y}_d})$   
 eg sigmoid/  
 transfer function

Using (2) in (1)

$$\frac{\partial E}{\partial w_i} = \sum_d (y_d - \hat{y}_d) \text{nid}$$
$$\hat{y}_d = \frac{1}{1 + e^{-w_i^T \theta}}$$

$$\Delta w_i = -\eta \sum_d (y_d - \hat{y}_d) \hat{y}_d (1 - \hat{y}_d) \text{nid}$$

↑  
Training  
rule  
of sigmoid  
unit

↑  
single layered  
'NN' → for linearly  
separable function.

# Machine Learning Lect 3

## Gradient Descent

Dr. Neeran Narkhede

9 Aug '20





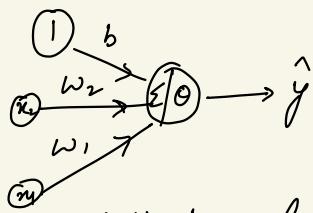
Ques AND function implementation  
using Single Layer Perceptron

Truth Table of AND			$x_1$	$x_2$	$t$
$x_1$	$x_2$	$t$	1	1	1
1	1	1	1	-1	-1
1	0	0	-1	1	-1
0	1	0	-1	-1	-1
0	0	0			

$\textcircled{1} \rightarrow -1$  (label)

$1 \rightarrow 1$  (Label)

Perception NW



$$\eta = \alpha = 1$$

some  
value  
b/w 0 to 1

Activation function  $\theta = 0$

$$\theta(y_{in}) = \begin{cases} 1 & y_{in} > 0 \\ 0 & y_{in} = 0 \\ -1 & y_{in} < 0 \end{cases}$$

Initialize weights = 0  
 $w_1, w_2, b = 0$

$$y_{in} = \sum w_i n_i + b$$

$$= b + w_1 n_1 + w_2 n_2 = 0 + 0(1) + 0(1)$$

$$= 0$$

$$\hat{y} = \phi(y_{in}) = 0$$

$\hat{y} \neq t \Rightarrow$  update weights

$$w_i(\text{new}) = w_i(\text{old}) + \Delta t \cdot n_i$$

$$\Delta \omega = \eta \left( t - \frac{v}{c} \right) n_i$$

$$\begin{aligned}w_{1\text{new}} &= w_{1\text{old}} + \alpha t x_1 \\&= 0 + 1 \cdot 1 \cdot 1 = 1\end{aligned}$$

$$\begin{aligned}w_{2\text{new}} &= w_{2\text{old}} + \alpha t x_2 \\&= 0 + 1 \cdot 1 \cdot 1 = 1\end{aligned}$$

$$\begin{aligned}b_{\text{new}} &= b_{\text{old}} + \alpha t \\&= 0 + 1 \cdot 1 = 1\end{aligned}$$

$$\begin{aligned}\text{i)} \quad y_{\text{in}} &= b + w_1 x_1 + w_2 x_2 \\&= 1 + (1)(1) + (1)(-1) \\&= 1 + 1 - 1 = 1\end{aligned}$$

$$\hat{y} = \Theta(y_{\text{in}}) = 1$$

$\hat{y} \neq t \Rightarrow$  update weights

$$\begin{aligned}w_{1\text{new}} &= w_{1\text{old}} + \alpha (t - \hat{y}) x_1 \\&= 1 + 1(-1 - 1) = 1 \\&= 1 - 2 = -1\end{aligned}$$

$$\begin{aligned}w_{2\text{new}} &= w_{2\text{old}} + \alpha (t - \hat{y}) x_2 \\&= 1 + 1(-1 - 1) = 1 \\&= 1 + 2 = 3\end{aligned}$$

$$\begin{aligned}b_{\text{new}} &= b_{\text{old}} + \alpha (t - \hat{y}) \\&= 1 + 1(-2) = -1\end{aligned}$$

$$\begin{aligned}\text{ii)} \quad y_{\text{in}} &= b + w_1 x_1 + w_2 x_2 \\&= -1 + (-1)(-1) + (3)(1) \\&= -1 + 1 + 3 = 3\end{aligned}$$

$$\hat{y} = \Theta(y_{\text{in}}) = 1$$

$\hat{y} \neq t \Rightarrow$  update weights

$$\omega_1 \text{ new} = \omega_1 \text{ old} + \eta (t - \hat{y}) n_1$$

$$\omega_1 \text{ new} = -1 + 1 (-1 - 1) 1$$

$$= -1 + 2 = 1$$

$$\omega_2 \text{ new} = \omega_2 \text{ old} + \eta (t - \hat{y}) n_2$$

$$= 3 + 1 (-1 - 1) 1$$

$$= 3 - 2 = 1$$

$$b \text{ new} = -1 + 1 (-2) = \frac{-1 - 2}{-3} = -3$$

(iv)  $\hat{y}_{\text{in}} = -3 + (1)(-1) + 1(-1)$   
 $= -3 - 1 - 1 = -5$

$\hat{y} = -1$   
 $\hat{y} = t \Rightarrow \text{no change in weight}$

Final weights are  
 $\omega_1 = 1 \quad \omega_2 = 1 \quad b = -3$

Decision boundary

$$b + \omega_1 x_1 + \omega_2 x_2 = 0$$

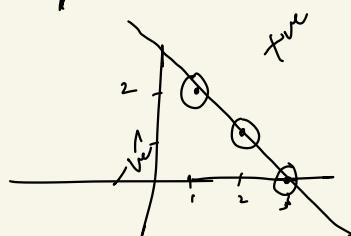
$$-3 + x_1 + x_2 = 0$$

$$x_1 + x_2 = 3$$

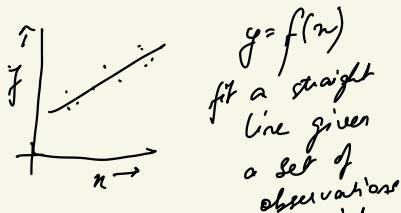
$$x_1 = 0 \quad x_2 = 3$$

$$x_1 = 1 \quad x_2 = 2$$

$$x_1 = 2 \quad x_2 = 1$$

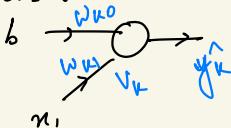


## Gradient Descent



linear regression  
 $y = f(n)$

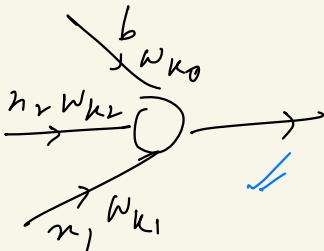
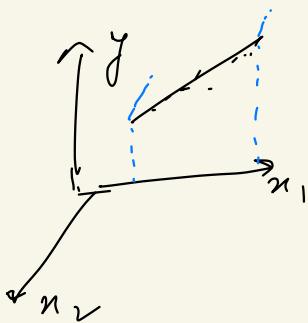
works for single neuron which linear activation unit  $v_k = \sum_{j=0}^n w_{kj} n_j + \hat{y}_k$



2 parameters: slope & intercept

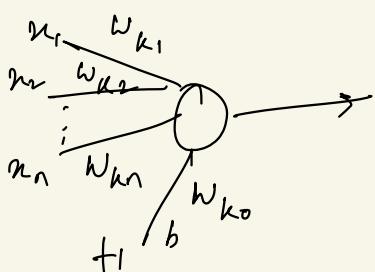
$$y = f(n_1, n_2)$$

line fitting → 3D line fitting



2 slopes along  $n_1, y$  plane  
components  $n_2, y$  plane  
 $w_{k1}$   
 $w_{k2}$

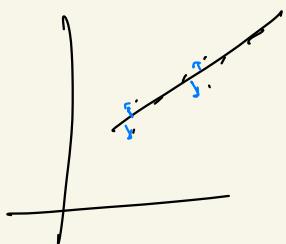
$$y = f(n_1, n_2, \dots, n_n) + b$$



$$y = \sum w_{kj} n_j$$

n-d straight line  
so many gradients along individual planes + bias

①



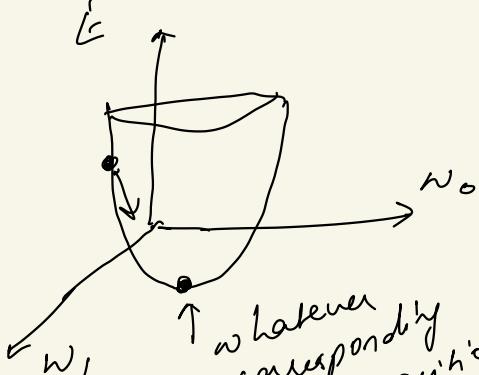
not passing  
exactly  
the origin  
points

Error values  
some positive  
& some  
negative

$\Rightarrow$  combined  
error

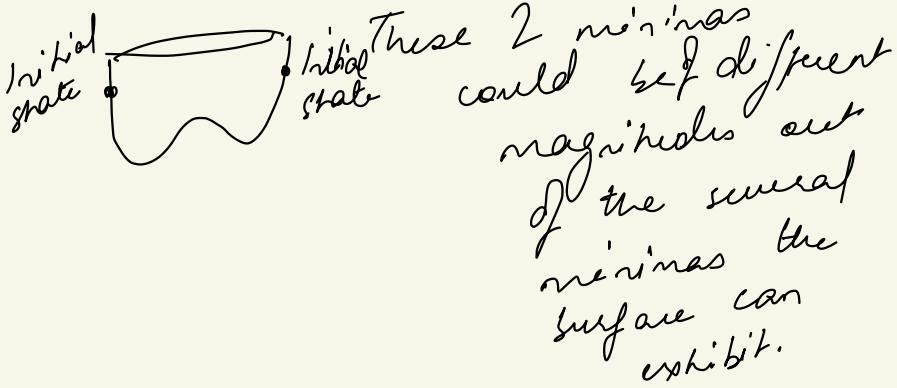
can be plot  
as a function  
of 2 variables  
 $w$ , &  $w_0$   
slope intercept  
natural tendency  
is to reach the  
minima (minimum  
error)

fit a curve which  
searching all over the plane  
gives min error for best combination  
of parameters  
( $w$ ,  $w_0$ )



whatever  
correspondingly  
position of  
 $w_1$  &  $w_0$  at this point  
(best solution)

②



Only 1 global minima rest are local minima.

System adopts & learns to reach best fitment & in the process it reaches global minima or gets trapped at local minima.

Gradient point up but the point moves opposite to the direction of the gradient.

Let us have  $p : n \cdot d$  observations

no of  
different  
ops

$$E = \sum_p E^p$$

$$E^p = \sum_{y_0} \left( t_0^p - y_0^p \right)^2 \quad - \textcircled{1}$$

$t_0^p$   $\uparrow$   
 Desired

$y_0^p$   $\uparrow$   
 Actual

Industrial output

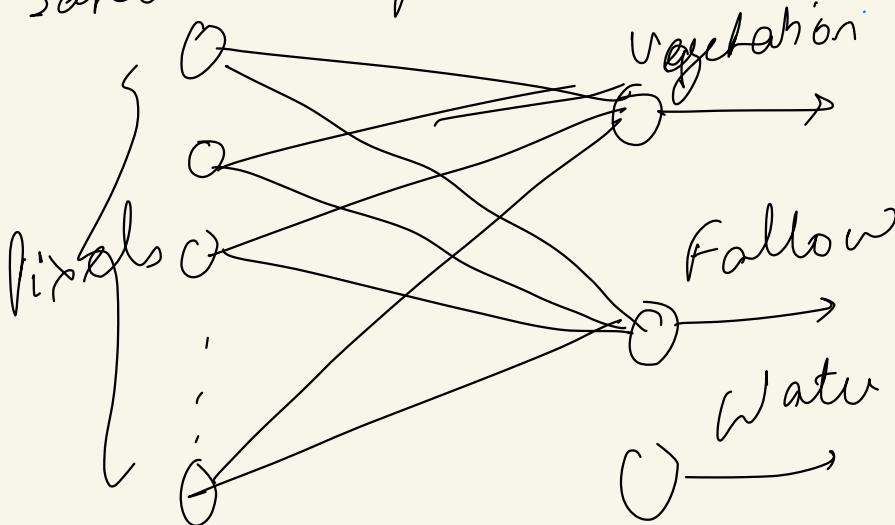
$$y_0 = f_1 (n_1, n_2, \dots, n_n)$$

$$y_1 = f_2 (n_1, n_2, \dots, n_n)$$

$$\vdots$$

$$y_n = f_n (n_1, n_2, \dots, n_n)$$

Satellite image



↗ semantic segmentation  
 ↗ image level vegetation

4

$$g_i = \frac{\partial E}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_p E^p$$

$$= \sum_p \frac{\partial E^p}{\partial w_{ij}}$$

Chain rule

$$\frac{\partial E}{\partial w_{oi}} = \frac{\partial E}{\partial y_o} \cdot \frac{\partial y_o}{\partial w_{oi}}$$

From ①

$$\frac{\partial E}{\partial y_o} = - (t_o - y_o) \dots \textcircled{2}$$

$$y_o = \sum_j w_{oj} x_j$$

⑤

$$\frac{\partial y_o}{\partial w_{oi}} = \frac{1}{\sum_j w_{oj}^2 n_j} = n_i - ③$$

$$\frac{\partial E}{\partial w_{oi}} = - \frac{(t_o - y_o) n_i}{w_{oi}}$$

particular weight input

Diagram illustrating the backpropagation error term:

$\Delta w_{oi} = f'(t_o - y_o) n_i$

$w_{oi} = w_{oi} + \Delta w_{oi}$

rate at which it is falling

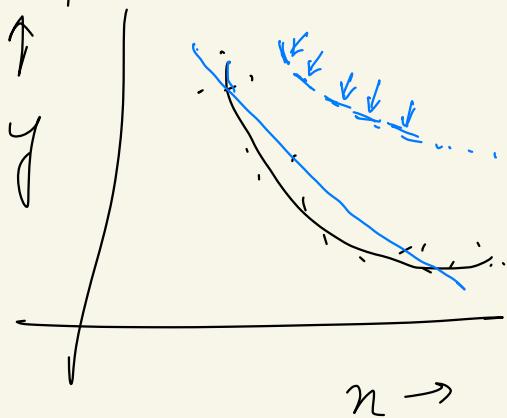
$\eta \rightarrow$  high if it's faster (right across the minima)  
else slower

Linear and straight line fitting

new neural netw activation  
linear fn.

oscillate to reach minima

$O^{-1}$



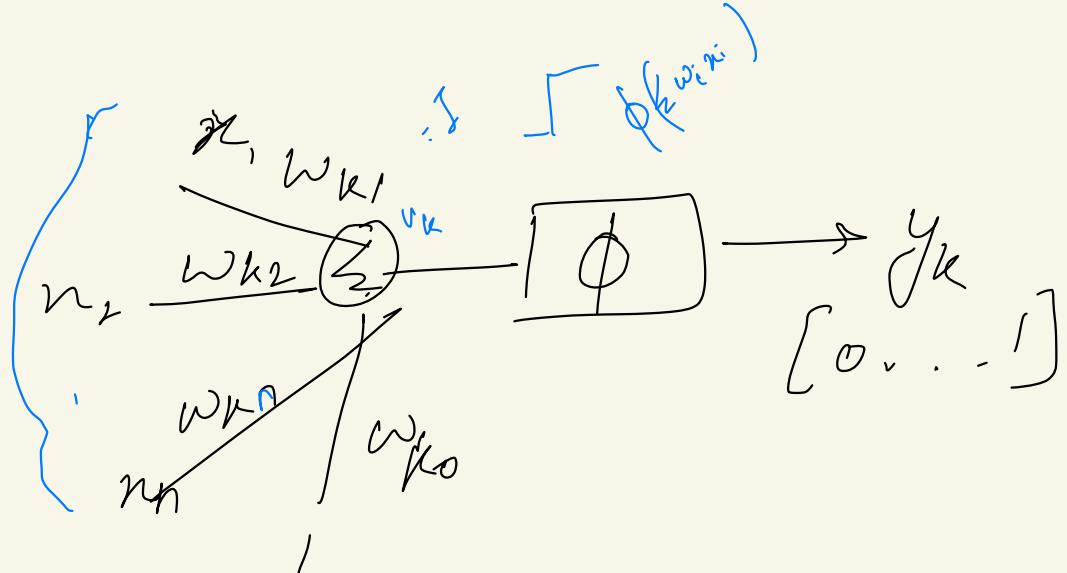
non - linear curve fitting

break it into piecewise linear components

↓  
complicated process

Soln: non - linear model

⑦



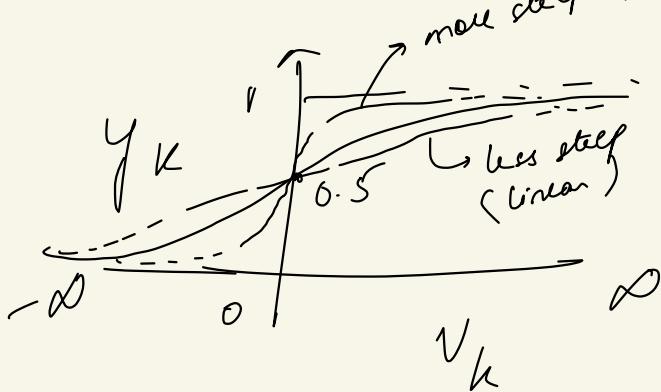
$$v_k = \sum_{j=0}^n j \cdot w_{kj}$$

more steep (step)

$$y_k = \phi(v_k)$$

$$y = f(v_k)$$

iff



i) monotonically increasing

ii) continuously differentiable

(McCulloch & Pitts fn.)

fn.  $\Rightarrow$  not continuous  
differentiable

Tunability : <sup>control</sup>  
shape

$$\phi(v_k) = \frac{1}{1 + e^{-av_k}} \quad (\text{sigmoid})$$

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

Quotient rule:

$$\sigma'(z) = \frac{(1+e^{-z}) \cdot 0 - (1)(-e^{-z})}{(1+e^{-z})^2}$$

$$= \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1 - 1 + e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{1 + e^{-z}}{(1+e^{-z})^2} - \frac{1}{(1+e^{-z})^2}$$

$$= \frac{1}{(1+e^{-z})} - \frac{1}{(1+e^{-z})^2}$$

$$= \left( \frac{1}{1+e^{-z}} \right) \left( 1 - \frac{1}{1+e^{-z}} \right)$$

$$= \sigma(z) (1 - \sigma(z))$$

⑨



# Machine Learning Lect 5<sup>48</sup>

## Multi Layer Neural Network

Dr. Partha Mukherjee



## # Multi layer Neural N/w

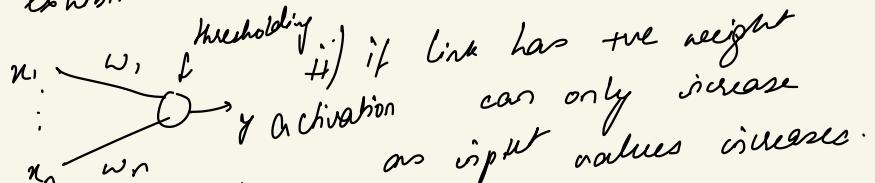
Single layer NN can represent the linear functions but in order to represent non-linear functions

we require multi layer NN which can be achieved by stacking perceptrons into different architectures

## # Feed forward NN.

Limitations in Perceptron:

i) They exhibit the monotonic behaviour



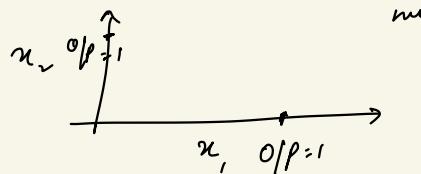
ii) If link has the weight & activation can only increase as input values increases.

(It cannot represent functions where input interactions can cancel each other.)

iii) Logic gates like AND can be represented as it is monotonic function but XOR cannot be represented by SLP as it is not a monotonic function O/P is 1 when  $x_1=1, x_2=1$

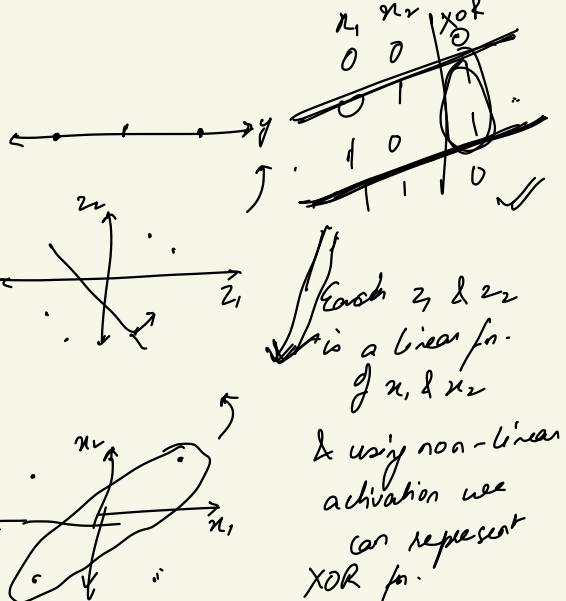
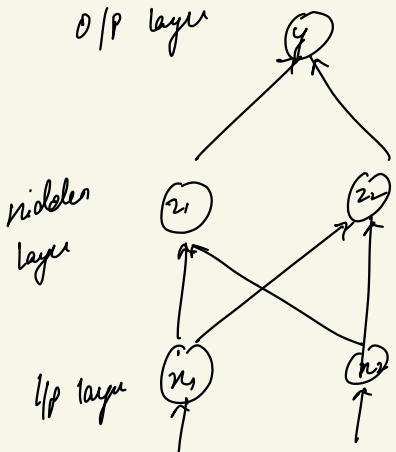
Can't represent functions where

input interactions can cancel each other



Solo: Multi layer NN.

iv) can represent only linear separable functions

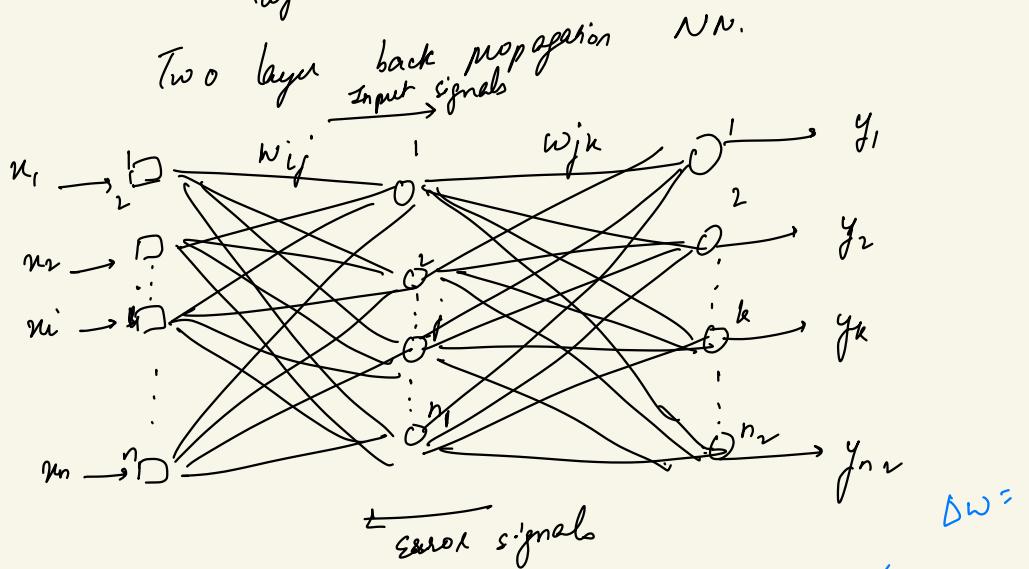
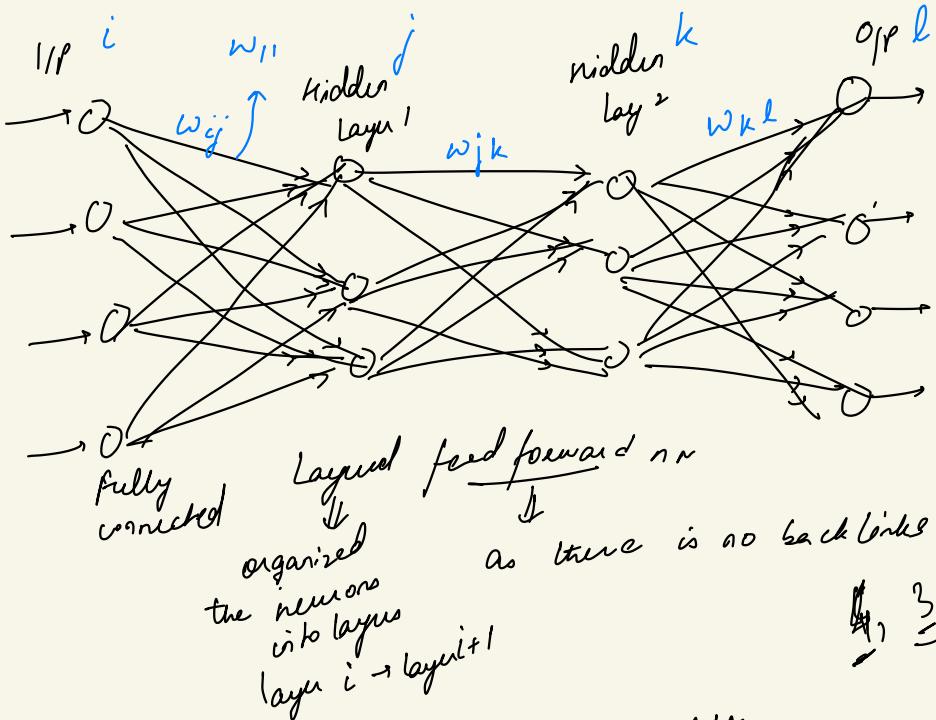


Why hidden units : because in training we can't observe them & using them we can represent many non-linear functions

### Power / Expressiveness of MLNN

- can represent interactions among inputs
- Two layer NNs can represent any Boolean function & continuous functions (within a tolerance) as long as the number of hidden units is sufficient & appropriate activation functions used
- Learning algorithm exists, but weaker guarantee than perceptron learning algorithms.

Three layer NN can represent all computable functions as they have good representational capability.

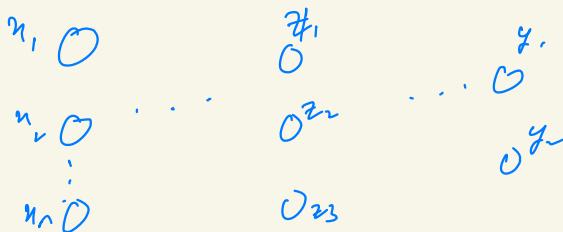


Perception : if there is change in opp (and)  
training rule you change the weights  
It can be achieved in SLP  
But there is problem in MLP??

If different b/w ideal & obtained opp

We do not know the ideal output in hidden units to compute the change in weights ( $w_{ij}$ )

Solution: Error observed at O/P layer is propagated backwards.



Error at  $y_1$  is backpropagated to  $z_1, z_2$  &  $z_3$

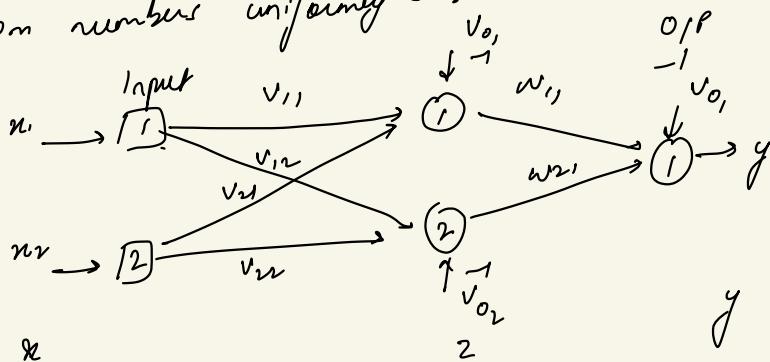
Error at  $y_2$  is also backpropagated to  $z_1, z_2$  &  $z_3$

& thus we get the notion of error & then we can update the weights

## # Backpropagation training algo.

Step 1 : Initialization

Set all weights & threshold levels of the network to random numbers uniformly distributed inside a small range.



## Step 2 Forward Computing

- Apply an input vector  $\bar{x}$  to input units
- Compute the activation / output at the hidden layer

$\Sigma \phi$

$$\bar{z}_j = \phi \left( \sum_i v_{ij} x_i \right)$$

- Compute the output vector  $\bar{y}$  at output layer

$$y_k = \phi \left( \sum_j w_{jk} z_j \right)$$

$\bar{y}$  is the result of computation.

Training examples

$\bar{x}_1$	$\bar{y}_1$	$\hat{y}_1$
$\bar{x}_2$	$\bar{y}_2$	$\hat{y}_2$
.	.	.
$\bar{x}_m$	$\bar{y}_m$	$\hat{y}_m$

error  $\Rightarrow$  based on  
this change  
weights

## Step 3 Learning for BP Net

- Update of weights in  $w$   
(between opp & hidden layers)
- delta rule
- Not applicable to updating  
 $v$  (b/w input & hidden)
- don't know the target  
values of hidden units

$$z_1, z_2 \dots z_p$$

Solution: Propagate error at opp units  
to hidden units to derive  
the update of weights  
 $v$  (again by delta rule)  
(error backpropagation learning)

- Error backpropagation can be continued downward if net has more weight than one hidden layer
- How to compute errors on hidden units?

For 1 output neuron, error fn. is

$$E = \frac{1}{2} (y - \hat{y})^2$$

For each unit  $j$ , the output  $o_j$  is defined as,

$$o_j = \phi(\text{net}_j) = \phi\left(\sum_{k=1}^n w_{kj} o_k\right)$$



The input net<sub>j</sub> to a neuron is the weighted sum of outputs  $o_k$  of previous  $n$  neurons.

Find derivative of the error:

~~Backpropagation training rule~~

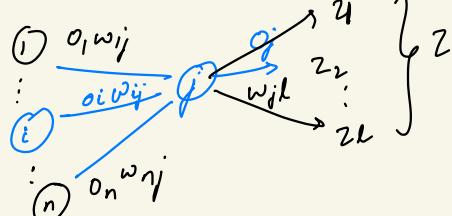
$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ij}}$$

- ①

$$\sum \phi'(j) \rightarrow o_j$$

$$\frac{\partial \text{net}_j}{\partial w_{ij}} = w_{i1} o_1 + w_{i2} o_2 + \dots + w_{in} o_n$$

$$\frac{\partial \text{net}_j}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \left( \sum_{k=1}^n w_{kj} o_k \right) = o_i \quad - \textcircled{1}$$



$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial}{\partial \text{net}_j} \phi(\text{net}_j) = \phi(\text{net}_j)(1 - \phi(\text{net}_j)) \quad - \textcircled{2}$$

Depends on form activation fn. (let's assume sigmoid) ✓

$$\frac{\partial E(o_j)}{\partial o_j} = \frac{\partial E(1, 2, \dots, l)}{\partial o_j}$$

$$\frac{\partial E(o_j)}{\partial o_j} = \sum \left( \frac{\partial E}{\partial \text{net}_l} \cdot \frac{\partial \text{net}_l}{\partial o_j} \cdot w_{jl} \right) \quad - \textcircled{3}$$

$$\frac{\partial E}{\partial w_{ij}} = \left( S_j \right) \cdot o_i$$

Errors computed at 2<sup>nd</sup> layer which is backpropagated at j

$$S_j = \frac{\partial E}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j}$$

$\left\{ \begin{array}{l} (o_j - t_j) o_j (1 - o_j) \Rightarrow \text{if it is an output} \\ \leq (\sum w_{il}) o_j (1 - o_j) \Rightarrow \text{if there is some intermediate unit} \end{array} \right.$

- ✓ Perception training rule
  - ✓ Gradient training rule
  - ✓ Backpropagation training rule
- NN ↗

# Machine Learning

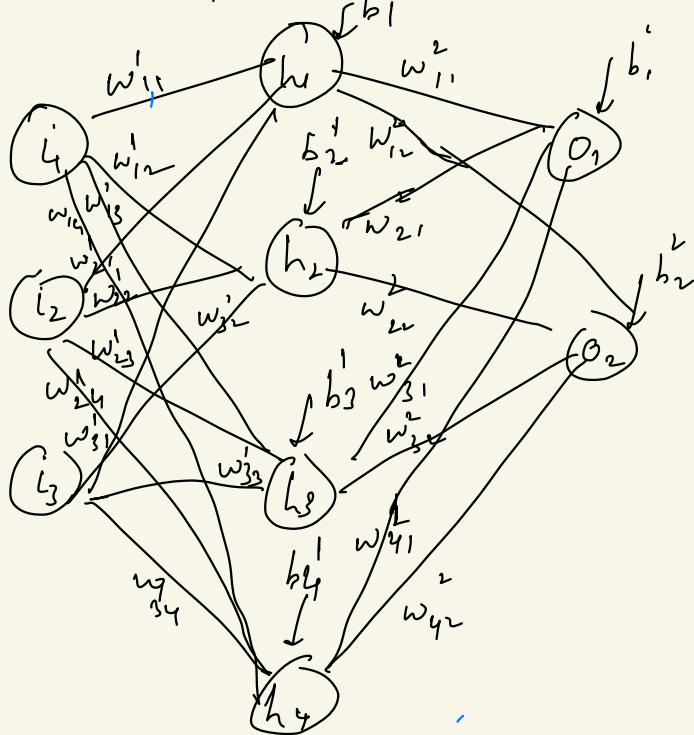
Lect . 6

Dr. Purna  
Mukherjee



# Multi-layer Perception

XOR problem

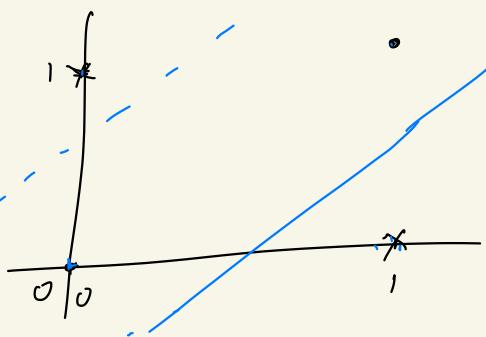


$$h_n = \sigma \left( \sum_{k=1}^3 i_k w_{kn} + b_n \right)$$

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

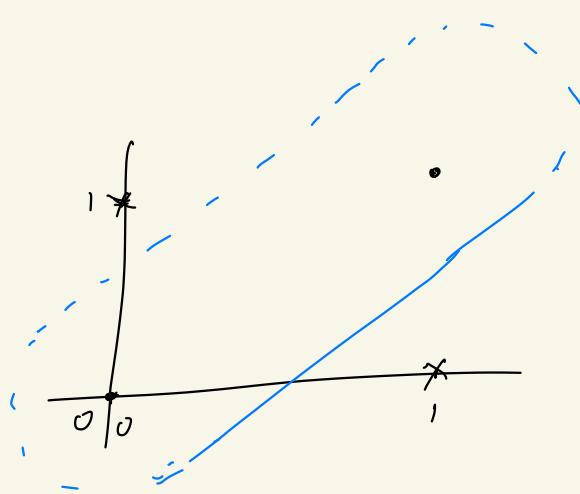
Universal approximation theorem

$n_1$	$n_2$	O/P	$y$
0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	0

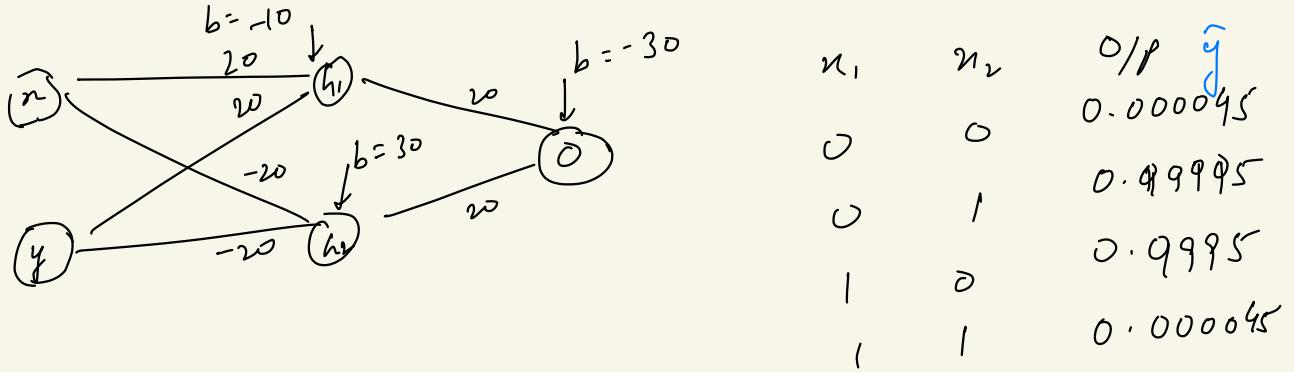


How to separate this with a linear boundary

↓  
Perception can't solve  
this  
but MLP can!



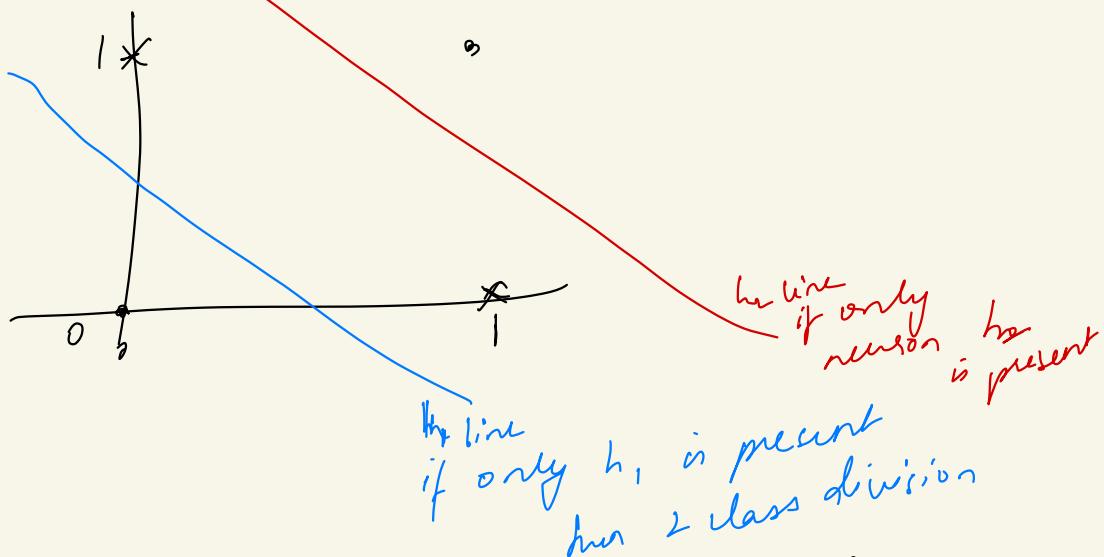
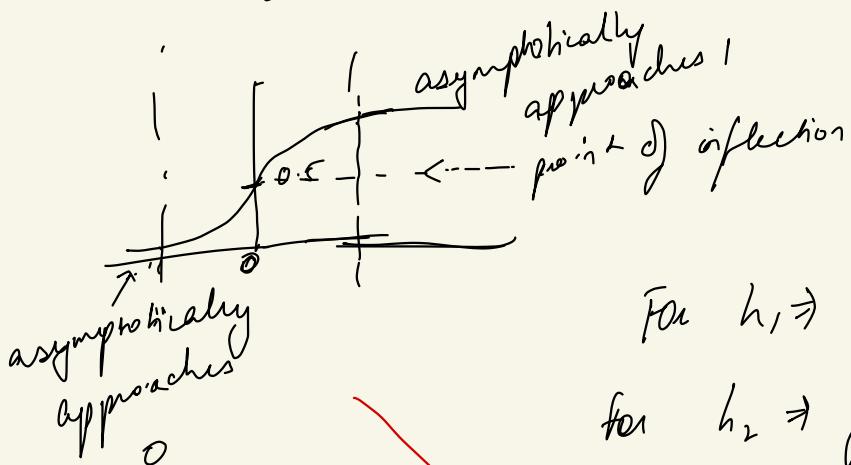
⇒ Non linear separation



$$h_1 = \text{sigmoid} \left( 20x + 20y - 10 \right) = 0.5$$

$$h_2 = \text{sigmoid} \left( -20x - 20y + 30 \right) = 0.5$$

$$o = 20h_1 + 20h_2 - 30$$



With only these given set of weights you can learn such complex XOR functions.

But how to learn these weights?

↓

Sol: Backpropagation

$$\begin{array}{lll} \tau(20 \times 0 + 20 \times 0 - 10) \approx 0 & \tau(-20 \times 0 - 20 \times 0 + 30) \approx 1 & \tau(20 \times 0 + 20 \times 1 - 30) \approx 0 \\ \tau(20 \times 1 + 20 \times 1 - 10) \approx 1 & \tau(-20 \times 1 - 20 \times 1 + 30) \approx 0 & \tau(20 \times 1 + 20 \times 0 - 30) \approx 0 \\ \tau(20 \times 0 + 20 \times 1 - 10) \approx 1 & \tau(-20 \times 0 - 20 \times 1 + 30) \approx 1 & \tau(20 \times 1 + 20 \times 1 - 30) \approx 1 \\ \tau(20 \times 1 + 20 \times 0 - 10) \approx 1 & \tau(-20 \times 1 - 20 \times 0 + 30) \approx 1 & \tau(20 \times 1 + 20 \times 1 - 30) \approx 1 \end{array}$$

# Machine Learning

Lect 8

Dr. Pearson  
Mukherjee



SVM : Support Vector Machines : hyperplane based classifiers

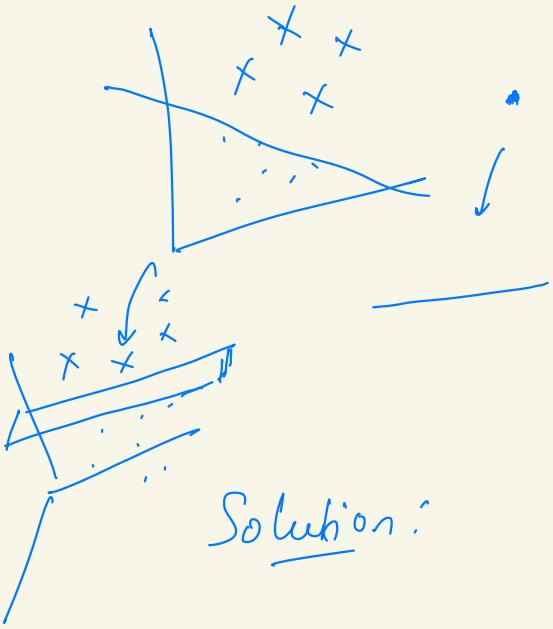
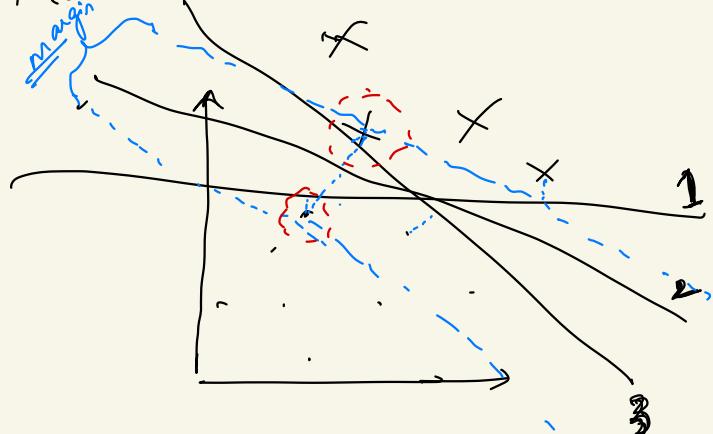
Vapnik Chervonenkis

Twin SVM  
MCML: Minimum complexity Machine

[VC dimension]

↓  
complexity of machine

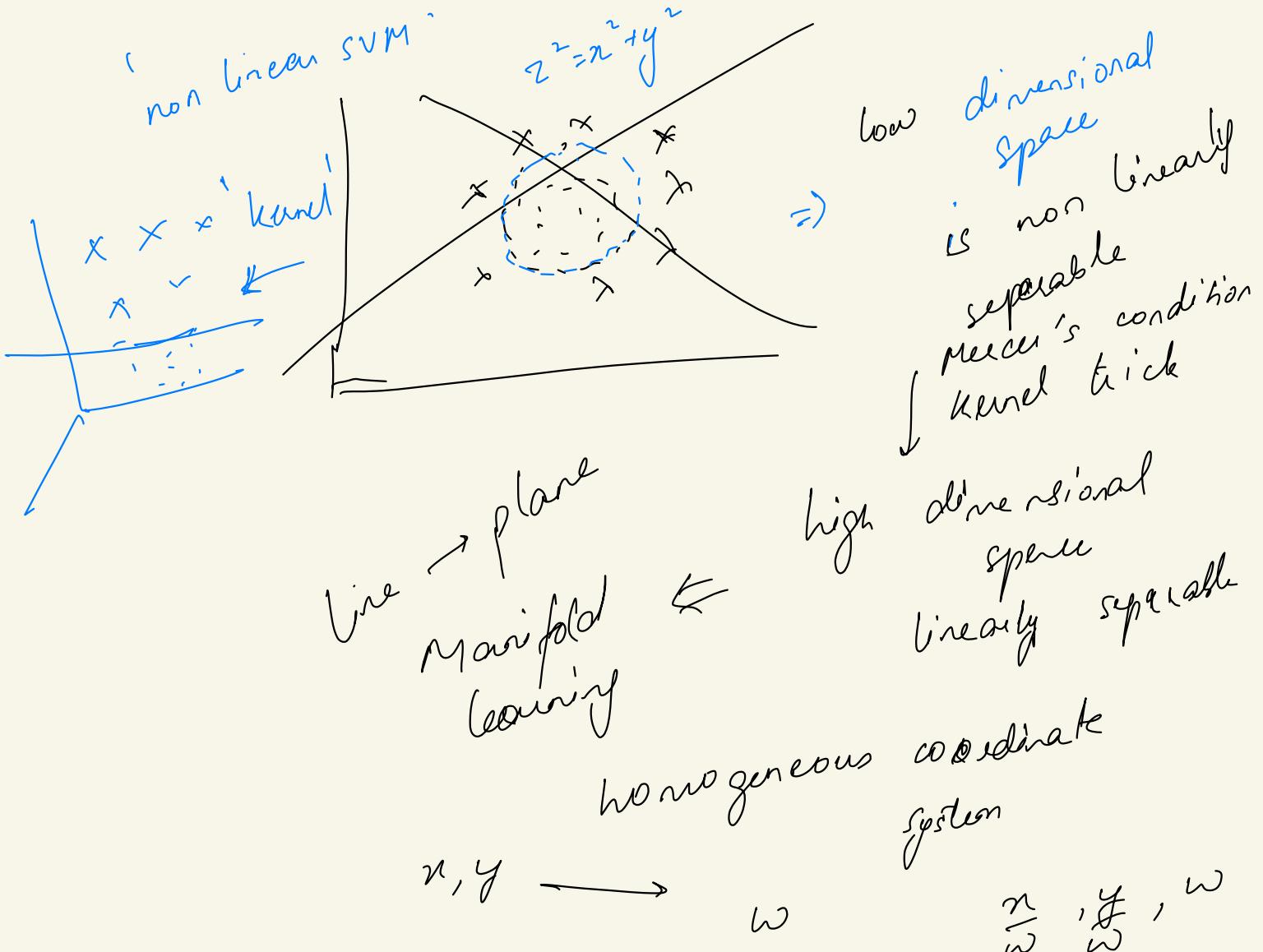
$$\approx \begin{cases} O(n) \\ O(n) \end{cases}$$



1, 2, 3 valid hyperplanes  
1, 2, 3 is a good  
hyperplane!

Solution: Find that hyperplane which has maximum margin!

Support vectors

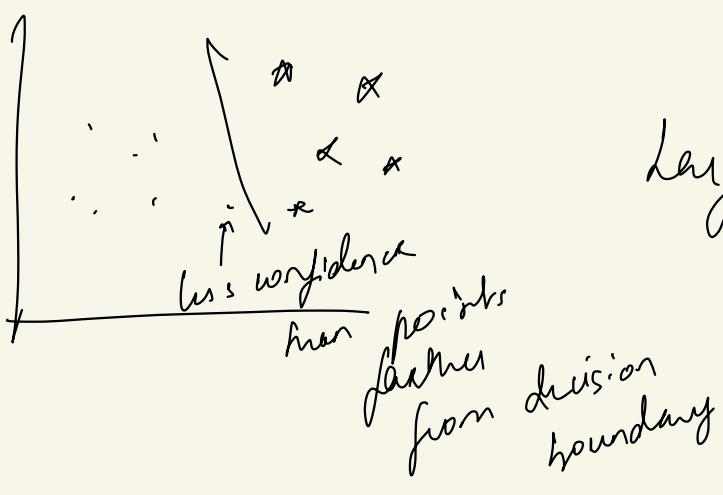
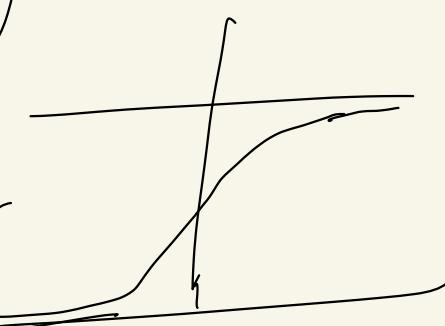


Logistic Regression:

$$P(y=1/n) = h(w) = \sigma(\beta^T n)$$

Predict 1 when  $h(n) > 0.5$

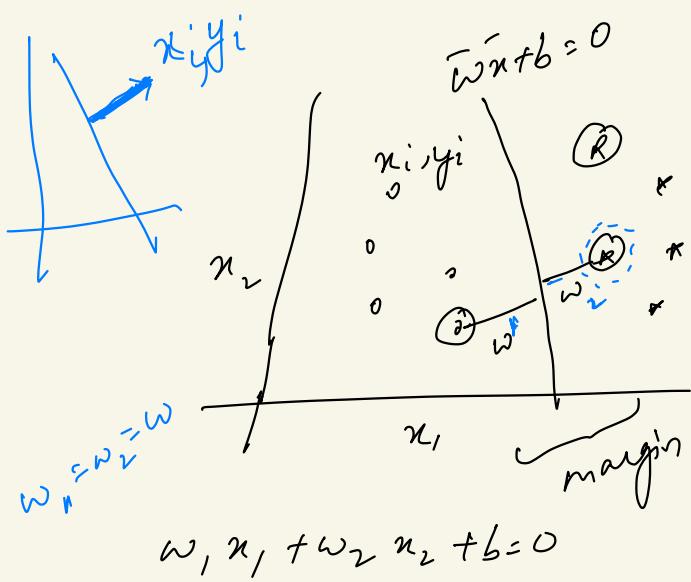
0 when  $h(n) < 0.5$



large no of features  
use Bayesian Learning  
we look at all

possible classifiers  
weigh them by their

aposteriori probability & combined to get the answer  
It is not computationally efficient.



Goal: margin width is highest  
classifier has highest generalization  
ability  
points closest to margin: Support vectors

Distance of closest +ve =  
distance of closest -ve  
minimum of 2 support vectors

Functional margin: dist of  $n_i, y_i$   
 $n_i, y_i$  wrt  $w, b$  from decision boundary

Larger functional

margin more

confidence

to

classify

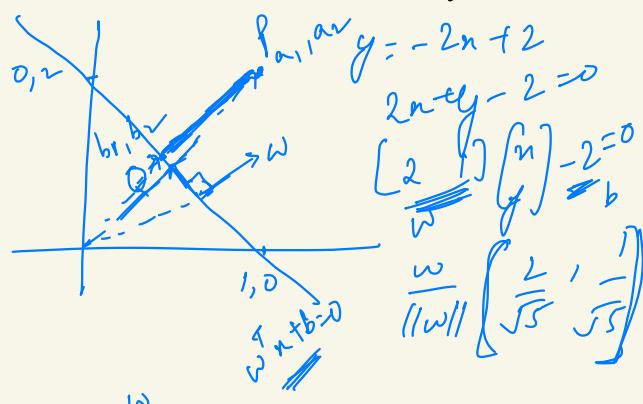
and also depends

on coefficients

of line

(Scaling factors)

$$\text{Line: } 2x + 3y + 1 = 0$$



$$\gamma^i = \underbrace{y_i}_{n_i^T w + b} \underbrace{(w^T n_i + b)}_{\text{Normal to the line}} = \min_{i=1}^m \gamma^i$$

Geometric margin: to get rid  
of scaling factors

unit vector:  $\frac{w}{\|w\|}$

$$= \left( \frac{2}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right)$$

$$P = Q + \gamma \frac{w}{\|w\|}$$

$$(a_1, a_2) = (b_1, b_2) + \gamma \frac{w}{\|w\|}$$

$$w^T ((a_1, a_2) - \gamma \frac{w}{\|w\|}) + b = 0$$

$$\gamma = \frac{w^T (a_1, a_2) + b}{\|w\|} = 0$$

$$y \leftarrow \frac{y}{\|w\|} (w^T x_i + b) \rightarrow \frac{\|w\|}{\|w\|} (w^T x_i + b) + \frac{b}{\|w\|}$$

Let  $\|w\| = 1$  (after normalization)

$$y_i \rightarrow y = y \left( \frac{w^T (x_1, x_2)}{\|w\|} + \frac{b}{\|w\|} \right)$$

# Maximize Margin width

→ Maximize  $\frac{y}{\|w\|}$  subject to

$$\rightarrow \frac{y}{\|w\|} \geq 1 \text{ for } i=1, 2, \dots, m$$

→ Scale so that  $y = 1$  has geometric margin

→ Maximizing  $\frac{1}{\|w\|}$  is same as minimizing  $\|w\| = w \cdot w$

→ minimize  $w \cdot w$  subject to constraints

for all  $x_i, y_i, i=1, \dots, m$ :

$$w^T x_i + b \geq 1 \text{ if } y_i = 1$$

$$w^T x_i + b \leq -1 \text{ if } y_i = -1$$

If  $(w, b)$  is decision boundary

$\frac{y}{\|w\|}$  is the geometric margin

margin

$$P = Q + \frac{y w}{\|w\|}$$

$$w^T w = \|w\|^2$$

$$w^T P = w^T Q + y \frac{w^T w}{\|w\|}$$

$$w^T P = w^T Q + y \|w\|$$

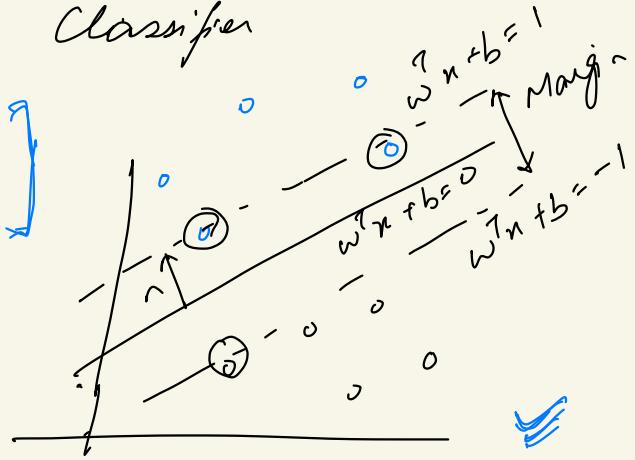
## # Large Margin Linear Classifier

Formulation:

$$\text{minimize}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

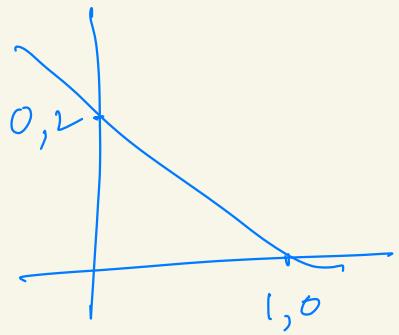
such that

$$y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad \forall i$$



Primal form

- Optimization problem with convex quadratic objectives & m inequality linear constraints
- Can be solved using QP
- Lagrange duality to get the optimization problem's dual form.
- Allow us to use kernels to get optimal margin classifiers to work efficiently in very high dimensional spaces.
- Allow us to derive an efficient algorithm for solving the above optimization problem that will typically do much better than generic QP software.



$$x + \frac{y}{2} = 1$$

$$2x + y = 2$$

$$2x + y - 2 = 0$$

$$w = \begin{bmatrix} 2 & 1 \end{bmatrix}$$

$$\frac{w}{\|w\|} = \left[ \frac{2}{\sqrt{5}}, \frac{1}{\sqrt{5}} \right]$$

# Machine Learning

Lect 9

Dr. Pranav  
Mukherjee



## # Solving the Optimization Problem

$$\text{minimize } \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 \quad (\text{linear inequality constraints})$$

Summary  
Classification

→ Optimization problem with convex quadratic objectives and linear constraints

→ Can be solved using QP.

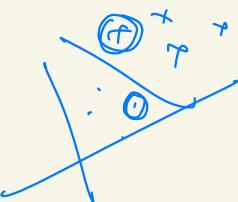
→ Lagrange duality to get the optimization problem's dual form:

- Allow us to use kernels to get optimal margin classifiers to work efficiently in very high dimensional space

- Allow us to derive an efficient algorithm for solving the above optimization problem that will typically do much better than generic QP software.

## # Lagrangian Duality

The primal Problem:  $\min_w f(w)$



Parameters: to find values of  $w$  to minimize  $f(w)$ ,  
 s.t.  $\begin{cases} g_i(w) \leq 0, i=1, \dots, k \\ h_i(w) = 0, i=1, \dots, l \end{cases}$   
 ('k' linear inequality constraints)  
 ('l' linear equality constraints)

The generalized Lagrangian:

$$L(w, \lambda, \beta) = f(w) + \sum_{i=1}^k \lambda_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

the  $\lambda_i$ 's ( $\lambda_i \geq 0$ ) and  $\beta$ 's are called the

Lagrange multipliers.

Lemma:

$$\max_{\lambda, \beta, \lambda_i \geq 0} L(w, \lambda, \beta) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise} \end{cases}$$

A re-written Primal:

$$\min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$$

$$L(w, \alpha, \beta)$$

The Primal Problem  $P^* = \min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$

The Dual Problem:  $d^* = \max_{\alpha, \beta, \alpha_i \geq 0} \min_w L(w, \alpha, \beta)$

Theorem (weak duality):

$$d^* = \max_{\alpha, \beta, \alpha_i \geq 0} \min_w L(w, \alpha, \beta) \leq$$

$$\min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta) = P^*$$

Theorem (strong duality):

iff there exist a saddle point of  $L(w, \alpha, \beta)$   
we have  $d^* = P^*$   $\Downarrow$  optimal value of  
both primal & dual forms

# KKT conditions

if there exists some saddle point of  $L$ , then  
it satisfies the following Karush Kuhn Tucker (KKT)

conditions:

$$\frac{\partial}{\partial w_i} L(w, \alpha, \beta) = 0, i=1, \dots, k$$

$$\frac{\partial}{\partial \beta_i} L(w, \alpha, \beta) = 0, i=1, \dots, l$$

$$\frac{\partial}{\partial \alpha_i} g_i(w) = 0, i=1, \dots, m$$

$$g_i(w) \leq 0, i=1, \dots, m$$

$$\alpha_i \geq 0, i=1, \dots, m$$

Theorem: If  $w^*, \alpha^*, \beta^*$  satisfy the KKT condition, then it  
is also a solution to the primal and the dual problems.

## # Support Vectors

In SVM, we just have  $g_i$ 's (inequality constraint)  
no  $h_i$ 's (equality constraints)

- Only a few  $\alpha_i$ 's can be nonzero  $\Rightarrow$  it's are S.V.
- Call the training data points whose  $\alpha_i$ 's are nonzero the support vectors.

$$\alpha_i g_i(\omega) = 0, i = 1, \dots, m$$

if  $\alpha_i > 0$  then  $g(\omega) = 0$

## # Solving the Optimization problem

Quadratic programming  
with linear  
constraints

$$\text{minimize } \frac{1}{2} \|\omega\|^2$$

$$\text{st. } y_i(\omega^T n_i + b) \geq 1$$

$$\text{minimize } L_f(\omega, b, \alpha) = \frac{1}{2} \|\omega\|^2 - \sum_{i=1}^m \alpha_i (y_i(\omega^T n_i + b) - 1)$$

$$\text{st. } \alpha_i \geq 0 \Leftrightarrow \text{S.V.}$$

①

minimize  
w.r.t  $\omega$  &  $b$   
for fixed  $\alpha$

$$\frac{\partial L_f}{\partial \omega} = 0 \Rightarrow \boxed{\omega = \sum_{i=1}^m \alpha_i y_i n_i} \rightarrow ②$$

$$\frac{\partial L_f}{\partial b} = 0 \Rightarrow \boxed{\sum_{i=1}^m \alpha_i y_i = 0} \rightarrow ③$$

Put ② in ①

$$L_f(\omega, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (n_i^T n_j) - \sum_{i=1}^m \alpha_i y_i \rightarrow ④$$

Put (3) in (4)

$$L_p(w, b, \lambda) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (w_i^T w_j)$$

The Dual Problem is

Now we have the following dual opt problem:

$$\max_{\alpha} J(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (w_i^T w_j)$$

$$\text{s.t. } \alpha_i \geq 0, \quad i=1, \dots, k$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

This is a quadratic programming problem:  
 $\rightarrow$  A global maximum of  $\alpha_i$  can always be found.

This dual opt. problem is much easier  
 to solve than primal formulation bce  
 constraints are simpler

## # SVM

Once we have the Lagrange multipliers  $\{\alpha_i\}$   
 we can reconstruct the parameter vector  $w$   
 as a weighted combination of the training

example:

$$w = \sum_{i=1}^m \alpha_i y_i w_i$$

$$w = \sum_{i \in SV} \underline{\alpha_i y_i w_i}$$

$\alpha_i \geq 0$

$w_i$ 's non zero  
 for  $i \in SV$

vis  
3

For testing with new data  $\underline{z}$

labels - compute  $\underline{w}^T \underline{z} + b = \sum_{i \in SV} \alpha_i y_i (\underline{n}_i^T \underline{z}) + b \quad \text{--- (1)}$

and classify  $\underline{z}$  as class 1 if sum is positive and class 2 otherwise

~~Labels~~ Note:  $w$  need not be formed explicitly we can just use (1) as we take dot product of support vectors with  $\underline{z}$

over classifier  
over classifier

## # Solving the Opt. Problem

→ The discrimination function is:

$$g(\underline{n}) = \underline{w}^T \underline{n} + b = \sum_{i \in SV} \alpha_i y_i \underline{n}^T \underline{n} + b$$

→ It relies on a dot product b/w test point  $\underline{n}$  and support vectors  $\underline{n} \Rightarrow$  we get a scalar

→ Solving the opt. problem involved in computing dot products  $\underline{n}_i^T \underline{n}_j$  b/w all pairs of training points

→ The optimal  $w$  is a linear combination of a small no of data points

Multilabel classifier

CL 1    CL 2    CL 3

OVO

CL 1    vs

CL 2

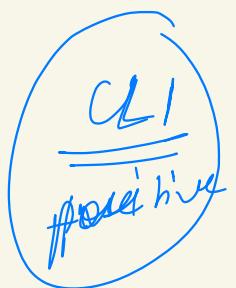
CL 1    vs

CL 3

binary classifier

OVA

=



vs

CL 2    CL 3  
negative

kNN

# Machine learning

Lect 10

De : Sumeet  
Mukherjee



## # Non linear SVM and Kernel functions

$$\phi: \mathcal{X} \rightarrow \phi(\mathcal{X})$$

original feature space can be mapped  
into new feature space where  
it is linearly separable  $\rightarrow$  higher dimensional

But what about computational cost ??

Linear  
SVMs

$$\text{linear } : \sum L_i d_j y_i y_j n_i \cdot n_j = 0$$

$\uparrow$  d-dimensional  
 $\uparrow$  d-dimensional  
we need to  
find pairwise  
dot products

if we have m training examples  
we need to do  $m^2$  computations  
like this

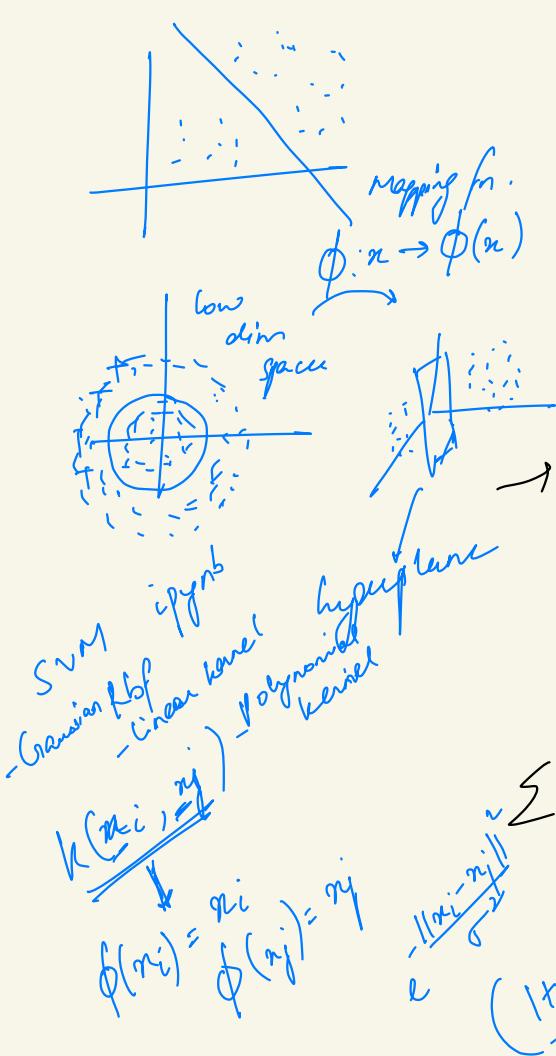
computation time  $O(d^2)$

If transform to high  
dimensional space  $D \gg d$

$$\sum L_i d_j y_i y_j \phi(n_i) \cdot \phi(n_j)$$

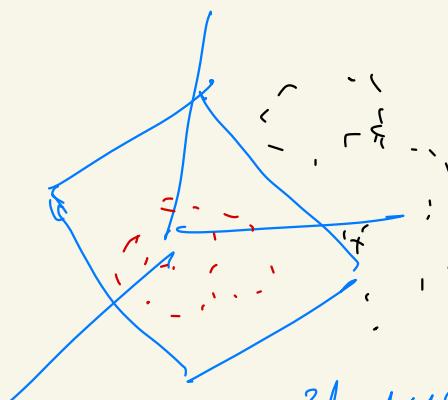
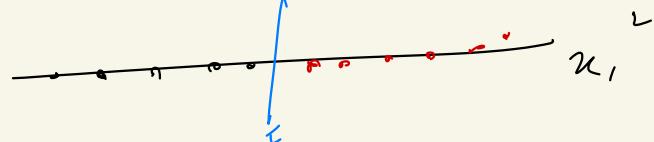
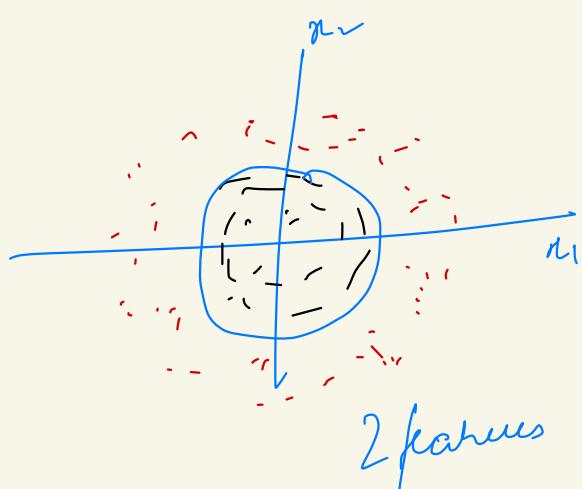
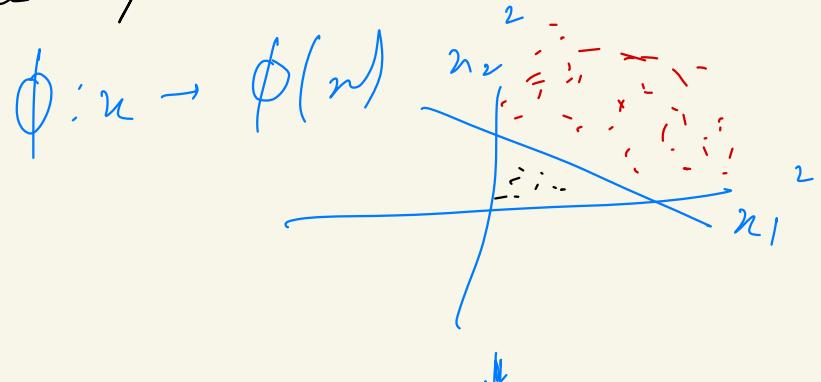
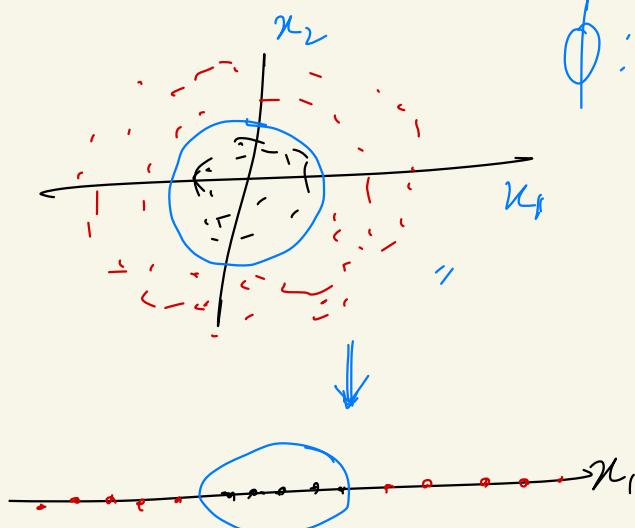
$O(D^2)$

$\uparrow$  more computational cost



Soln: Use kernel function to reduce computational costs

## # Visualizing Feature Space



3 features  
linearly separable

In certain feature transformation where after transformation, SVM can be solved efficiently.

Kernel function :  $n_a \quad n_b$   
 $\phi(n_a) \quad \phi(n_b)$

$$\frac{\phi(n_a) \cdot \phi(n_b)}{[ ]}$$

$$K(n_a, n_b) = \phi(n_a) \cdot \phi(n_b)$$

also called kernel trick

may be easy to compute if

is a fn. of  $n_a$  &  $n_b$  without expanding  $\phi(n_a)$  or  $\phi(n_b)$

What is Kernel: i) Original input attributes is mapped to a new set of input features via feature mapping  $\phi$   
 ii) Since the algo can be written in terms of the scalar product, we replace  $n_a \cdot n_b$  with  $\phi(n_a) \cdot \phi(n_b)$   
 iii) For certain  $\phi$ , there is a simple operation on 2 vectors in the low-dim space that can be used to compute the scalar product of their 2 images in the high dim space  
 $k(n_a, n_b) = \phi(n_a) \cdot \phi(n_b)$   
 Let the kernel do the work rather than do the scalar product in high dim space

# Non linear SVM → With this mapping, our discriminant function is now:  

$$g(n) = w^T \phi(n) + b$$

$$= \sum_{i \in SV} \alpha_i \boxed{\phi(n_i)^T \phi(n)} + b$$

$$\begin{matrix} w^T n \\ w^T \phi(n) \\ w^T \phi(n) \end{matrix}$$

- In dual form  
 $\rightarrow$  We only use the dot product of feature vectors in both - the training & test
- $\rightarrow$  A kernel function is defined as a function that corresponds to a dot product of two feature vectors in some expanded feature space:
- $$k(n_a, n_b) = \underline{\phi(n_a) \cdot \phi(n_b)}$$

This  $k(n_a, n_b)$  may be very inexpensive to compute even if  $\phi(n_a)$  may be extremely high dimensional

### Kernel Example

2 dimensional vectors  $\bar{n} = [n_1, n_2]$

$$\text{let } k(n_i, n_j) = (1 + n_i \cdot n_j)^2$$

$$k(n_i, n_j) = (1 + n_i \cdot n_j)^2$$

$$= 1 + n_i^2 n_j^2 + 2 n_i \cdot n_j, n_{i2} n_{j2} + n_{i2}^2 n_{j2}^2$$

$$+ 2 n_{i1} n_{j1} + 2 n_{i2} n_{j2}$$

$$= \begin{bmatrix} 1 & n_{i1}^2 & \sqrt{2} n_{i1} n_{i2} & n_{i2}^2 & \sqrt{2} n_{i1} \sqrt{2} n_{i2} \\ 1 & n_{j1}^2 & \sqrt{2} n_{j1} n_{j2} & n_{j2}^2 & \sqrt{2} n_{j1} \sqrt{2} n_{j2} \end{bmatrix}$$

$$= \phi(n_i) \cdot \phi(n_j)$$

where  $\phi(n) = \begin{bmatrix} 1 & n_i^2 & 2n_i n_j & n_j^2 & \sqrt{2}n_i & \sqrt{2}n_j \end{bmatrix}$

$\left[ 1 + n_i n_j \right] \quad \left[ \begin{matrix} n_i & n_j \\ n_j & n_i \end{matrix} \right]$

$$(a+b+c)^2 = a^2 + b^2 + c^2 + 2ab + 2ac + 2bc$$

$$= (f n_{ii}^2 n_{jj}^2 + n_{ij}^2 n_{ji}^2 + 2 n_{ii} n_{ji}) + 2 n_{ij} n_{ji} + 2 n_{ij} n_{ji}$$

$n_{ij} n_{ji}$

Dot product

$$a \cdot b = |a| |b| \cos \theta$$

or

$$a \cdot b = \underbrace{a_x b_x + a_y b_y}_{\text{Multiplying the } n \text{-components of 2 vectors}} + \underbrace{a_z b_z}_{\text{then } z\text{-components}}$$

Scalar quantity

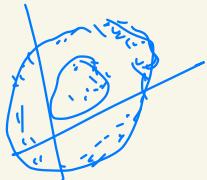
$$n_i = \begin{bmatrix} n_{ii} & n_{i2} \\ n_{j1} & n_{j2} \end{bmatrix}$$

$$n_i \cdot n_j = n_{ii} n_{jj} + n_{i2} n_{j2}$$

$$(1 + n_i \cdot n_j)^2 = (1 + n_{ii} n_{jj} + n_{i2} n_{j2})^2 \Rightarrow (a+b+c)^2$$

- ## # Commonly used kernel functions
- linear SVM corresponds to linear kernel:  
 $k(n_i, n_j) = n_i \cdot n_j$  where  $\phi(n_i) = n_i$  &  $\phi(n_j) = n_j$
  - Polynomial of degree  $P$ :  
 $k(n_i, n_j) = (1 + n_i \cdot n_j)^P$  Identifying mapping
  - Gaussian (radial basis function):  
 $k(n_i, x_j) = e^{-\frac{\|n_i - x_j\|^2}{2\sigma^2}}$
  - Sigmoid:

$$k(n_i, n_j) = \tanh(\beta_0 n_i \cdot n_j + \beta_1)$$

 In general, functions that satisfy Mercer's condition can be kernel functions

Eg: Text classification we can use kernels based on similarity of words of text

## # Kernel Functions

- kernel function can be thought of as a similarity measure b/w the input objects
- Not all similarity measure can be used as kernel function
- Mercer's condition states that any positive semi-definite kernel  $k(n, g)$  ie
  $\sum k(n_i, n_j) c_i c_j \geq 0$  for any real no.
 

Generally this similarity measure is symmetric

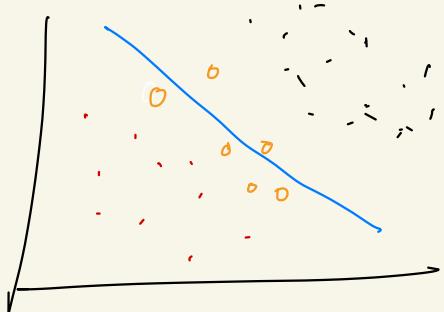
$$k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$$

↓  
is symmetric  
+ all the eigen values  
are even

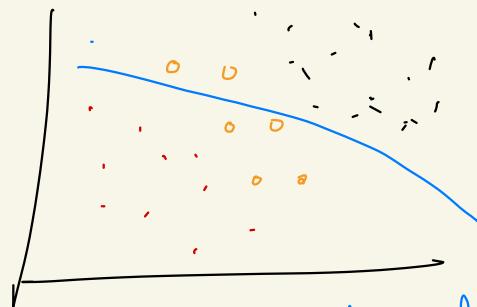
If you can write similarity b/w points as a matrix & this matrix is symmetric & semi definite then a kernel function will exist

→ Such functions can be expressed as a dot product in high dimensional space

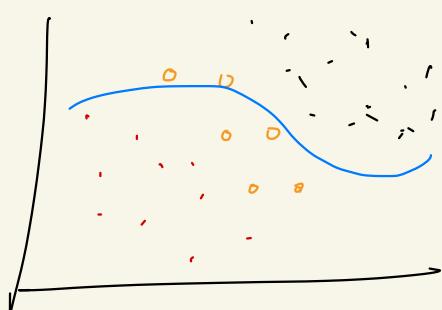
### SVM examples



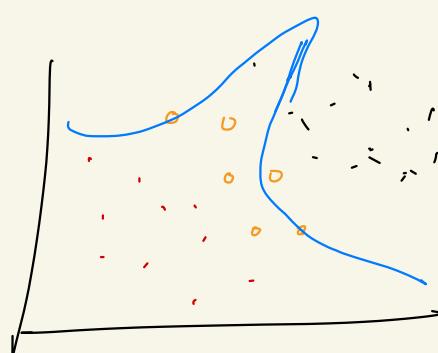
↓  
linear kernel  
with noise



↓ 2nd order  
polynomial



↓  
4th order  
polynomial



↓  
8th order  
polynomial

Non linear SVM optimization  
 Formulation (Lagrangian Dual Problem)

$$\text{maximize}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \quad \text{K}(x_i, x_j)$$

s.t.

$$0 \leq \alpha_i \leq C \quad \rightarrow \text{hyperparameter}$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

$C=0.001$   $\forall K$  condition

The solution of the discriminant function is any test point  $x$

$$= g(x) = \sum_{i \in SV} \alpha_i K(x_i, x) + b$$

$\downarrow$   $b = \phi(x_i) \cdot \underline{\phi(x)}$

$\phi(x_i) \cdot \phi(x)$

# Performance:

- SVM work very well in practice
- We must choose the kernel function & its parameters ( $C$ ) appropriately
- They can be expensive in time & space for big dataset
- Computation of the maximum margin hyperplane depends on the square of number of training cases
- We need to store all support vectors
- The kernel trick can also be used to do PCA in a much higher-dimensional

Space, thus giving a non linear version of PCA  
in the original space

## # Multiclass classification

SVMs can only handle 2-class outputs

Learn N-SVMs

- SVM1 learns Class 1 vs Rest

- SVM2 learns Class 2 vs Rest

SVM N learns Class N vs Rest

Then to predict off for a new input, just  
predict with each SVM & find out which  
one put the prediction the furthest  
into the positive region

majority voting

Class<sup>1</sup> Class<sup>2</sup> Class<sup>3</sup> ... Class<sup>N</sup>

OVO  
one vs one

— SVM<sub>1</sub> Class 1 vs Rest  
Class 2 vs Rest

—

—

OvA  
one vs all

Class<sup>1</sup> ————— negative  
R

sk-learn SVM predict  
svm.classify

pymlab lib-svm library

# Machine learning

Lect 11

Dr. Haran

Mukherjee



## # SVM Numerical

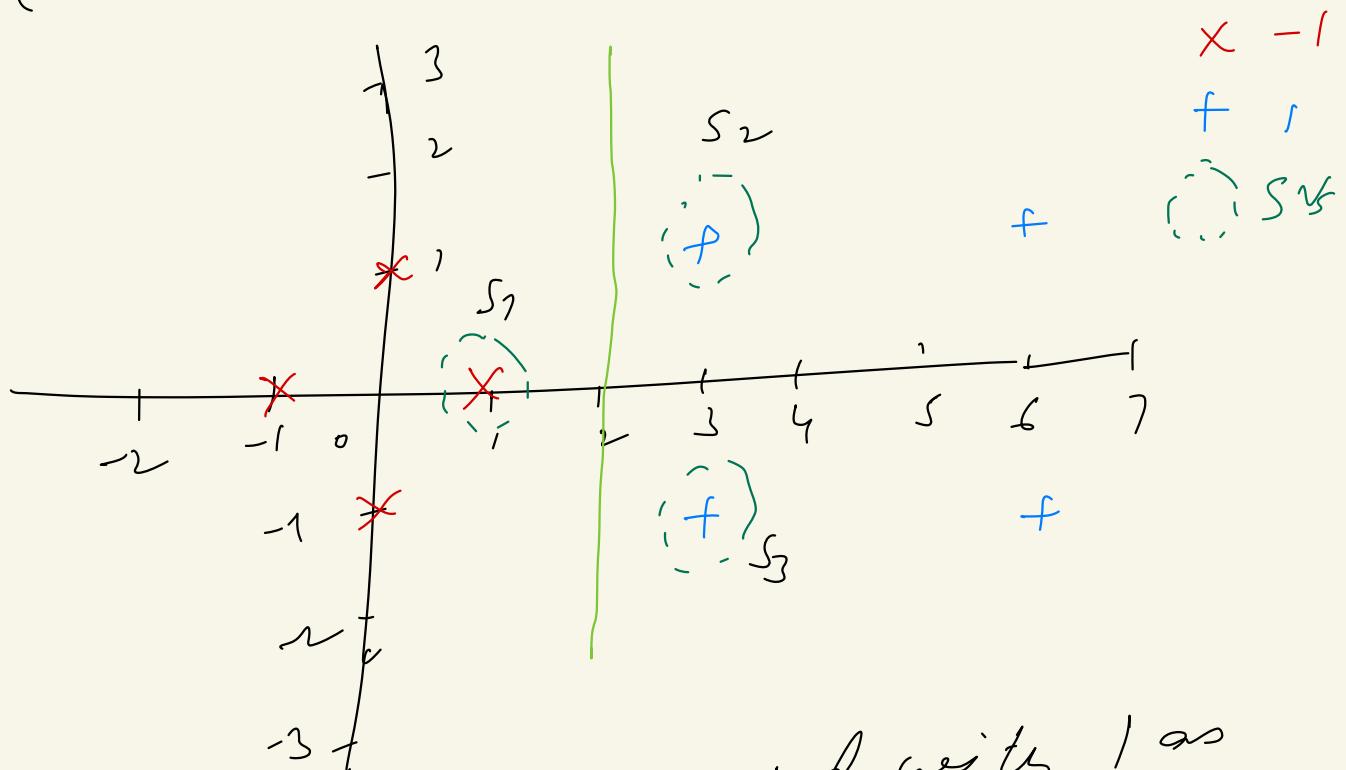
Suppose we are given the following data points

$$\left\{ \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 0 \end{pmatrix}, \begin{pmatrix} 6 \\ -1 \end{pmatrix} \right\}$$

five labeled +1

$$\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix} \right\}$$

unlabeled -1



Each vector is augmented with 1 as bias input

$$S_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \tilde{S}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

$$S_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \quad \tilde{S}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \quad S_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix} \quad \tilde{S}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

$$\mathcal{L}_1 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} = -1$$

$$\cancel{\mathcal{L}_1} \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} = +1$$

$$\mathcal{L}_1 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} = +1$$

$$\mathcal{L}_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\mathcal{L}_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = +1$$

$$\mathcal{L}_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = +1$$

$$\begin{aligned} \alpha_1(1+0+1) + \alpha_2(3+0+1) \\ + \alpha_3(3+0+1) = -1 \end{aligned}$$

$$\begin{aligned} \alpha_1(3+0+1) + \alpha_2(9+1+1) \\ + \alpha_3(9-1+1) = 1 \end{aligned}$$

$$\begin{aligned} \alpha_1(3+0+1) + \alpha_2(9-1+1) + \\ \alpha_3(9+1+1) = 1 \end{aligned}$$

$$2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1$$

$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1$$

$$\begin{aligned} \alpha_1 = -3.5 & \quad \alpha_2 = 0.75 \\ \alpha_3 = 0.75 \end{aligned}$$

$$\begin{aligned} \tilde{\omega} &= \sum_i \alpha_i s_i \\ &= -3.5 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \end{aligned}$$

$$= \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix} ] \rightarrow \tilde{\omega}$$

sinc

$$y = \omega x + b$$

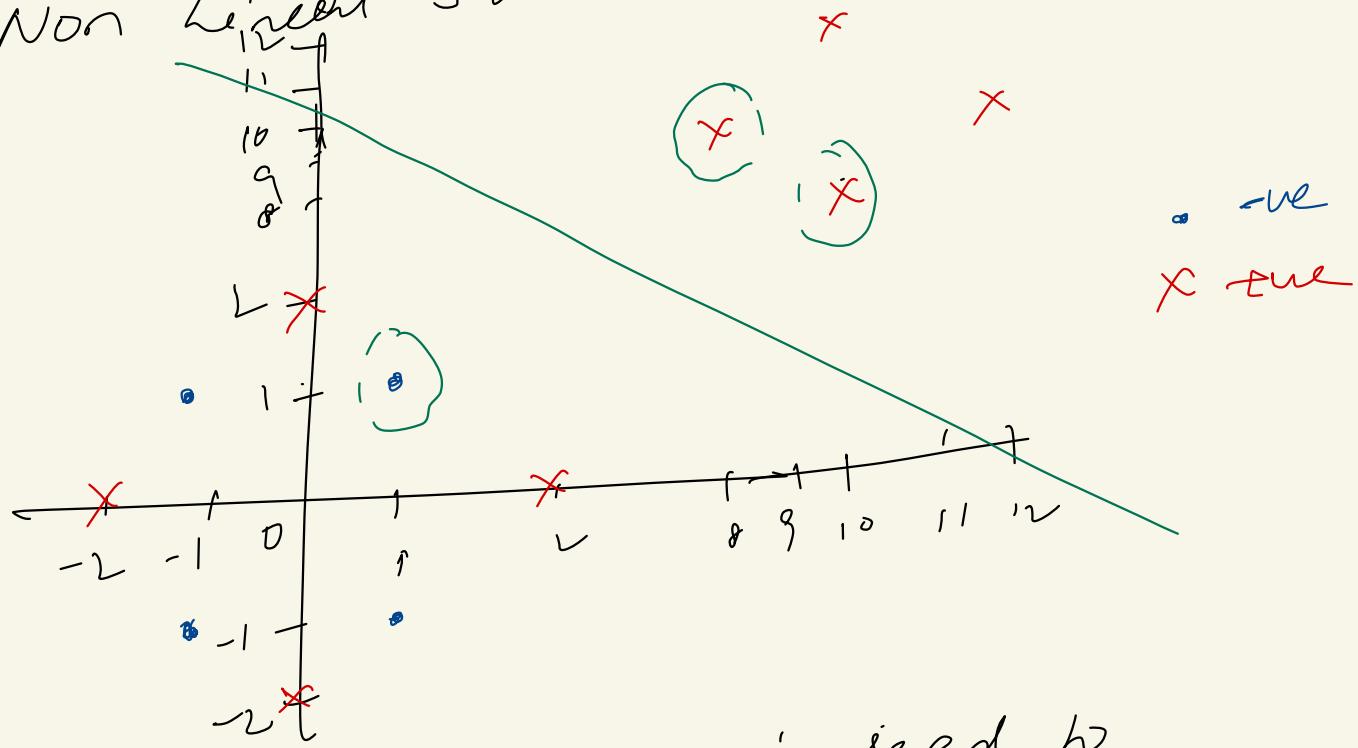
$$\omega = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad b = -2$$

R -ve class points  $\rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ -1 \end{pmatrix}$

the +ve class points  $\rightarrow \begin{pmatrix} 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 5 \\ 1 \end{pmatrix}, \begin{pmatrix} 5 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 0 \end{pmatrix}$

find the generality function

## # Non Linear SVM



$\phi \rightarrow$  Mapping function is reqd to transform these data to a new feature space where a separating hyperplane can be found.

$$\phi(n_1, n_2) = \begin{cases} \begin{pmatrix} 6 - n_1 + (n_1 - n_2)^2 \\ 6 - n_2 + (n_1 - n_2)^2 \end{pmatrix} & \text{if } n_1 + n_2 \geq 0 \\ \begin{pmatrix} n_1 \\ n_2 \end{pmatrix} & \text{otherwise} \end{cases}$$

$\rightarrow$  Blue class vectors are :  $\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix},$

$\begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}$  no change since

$$\sqrt{n_1^2 + n_2^2} < 2 \text{ ff vectors}$$

→ Red class vectors  $\begin{pmatrix} 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} -2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -2 \end{pmatrix}$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} 2 \\ 0 \end{pmatrix} = \begin{pmatrix} 6 - 2 + (2-0)^2 \\ 6 - 0 + (2-0)^2 \end{pmatrix} = \begin{pmatrix} 8 \\ 10 \end{pmatrix}$$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} 0 \\ 2 \end{pmatrix} = \begin{pmatrix} 10 \\ 8 \end{pmatrix}$$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} -2 \\ 0 \end{pmatrix} = \begin{pmatrix} 12 \\ 10 \end{pmatrix}$$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} 0 \\ -2 \end{pmatrix} = \begin{pmatrix} 10 \\ 12 \end{pmatrix}$$

$$s_1 = \begin{pmatrix} 8 \\ 10 \end{pmatrix} \quad s_2 = \begin{pmatrix} 10 \\ 8 \end{pmatrix} \quad s_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\tilde{s}_1 = \begin{pmatrix} 8 \\ 10 \\ 1 \end{pmatrix} \quad \tilde{s}_2 = \begin{pmatrix} 10 \\ 8 \\ 1 \end{pmatrix} \quad \tilde{s}_3 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$d_1 \tilde{s}_1 \tilde{s}_1 + d_2 \tilde{s}_2 \tilde{s}_1 + d_3 \tilde{s}_3 \tilde{s}_1 \stackrel{\sim}{=} +1$$

$$d_1 \tilde{s}_2 \tilde{s}_2 + d_2 \tilde{s}_2 \tilde{s}_2 + d_3 \tilde{s}_2 \tilde{s}_2 \stackrel{\sim}{=} +1$$

$$d_1 \tilde{s}_3 \tilde{s}_1 + d_2 \tilde{s}_3 \tilde{s}_2 + d_3 \tilde{s}_3 \tilde{s}_3 \stackrel{\sim}{=} -1$$

$$165\omega_1 + 161\omega_2 + 19\omega_3 = +1$$

$$164\omega_1 + 165\omega_2 + 19\omega_3 = +1$$

$$19\omega_1 + 19\omega_2 + 3\omega_3 = -1$$

$$\omega_1 = \omega_2 = 0.859 \quad \omega_3 = -1.4219$$

$$\tilde{\omega} = \sum_i \omega_i s_i$$

$$\tilde{\omega} = \omega_1 \begin{pmatrix} 8 \\ 10 \\ 1 \end{pmatrix} + \omega_2 \begin{pmatrix} 10 \\ 8 \\ 1 \end{pmatrix} + \omega_3 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 0.1243 \\ 0.1243 \\ -1.2501 \end{pmatrix}$$

$$y = \omega^n + b$$

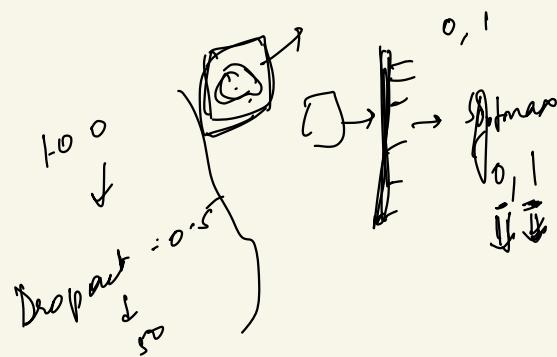
$$\omega = \begin{pmatrix} 0.1243 & 0.1243 \\ 0.1243 & 0.1243 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$b = -1.2501 / 0.1243 \\ = -10.057$$

Test point  $n_1, n_2 = -1, 2$  belongs to red point class five

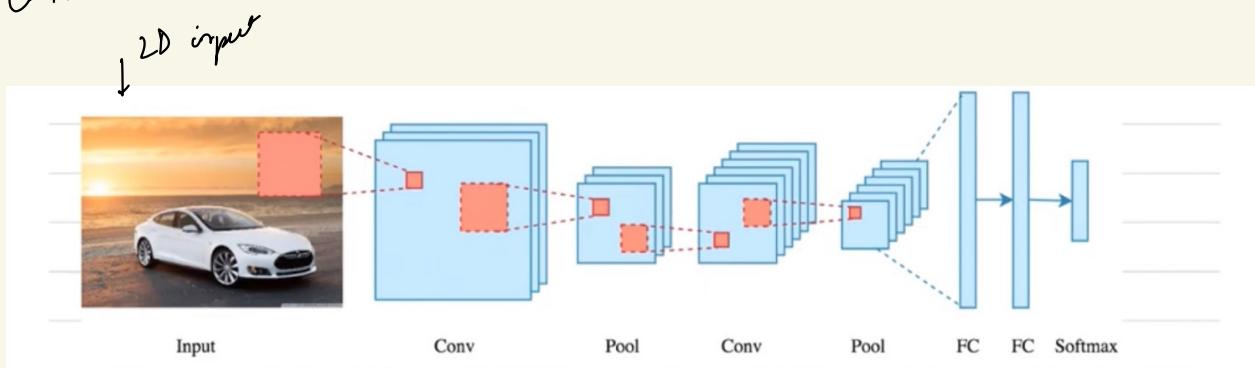
$$\phi(n_1, n_2) = \begin{pmatrix} 16 \\ 13 \end{pmatrix}$$

$$w \phi(n) = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 16 \\ 13 \end{pmatrix} = 29710$$



## # Introduction to Deep Learning (CNNs)

Convolutional Neural Network



Correlation  
Convolution  
Input  
1D  
input  
 $h(x)$   
f $\theta^0$

CNNs over ANNs - i) No context learnt (Neighborhood of pixels)  
ii) More parameters to learn in ANNs

no flipping  $\downarrow$  flip kernel 18°  
Correlation / Convolution  
In case of CNNs  $\rightarrow$  it is actually correlation only  
Conv (filter size, stride, padding)  $\rightarrow$  pooling

① Filter size

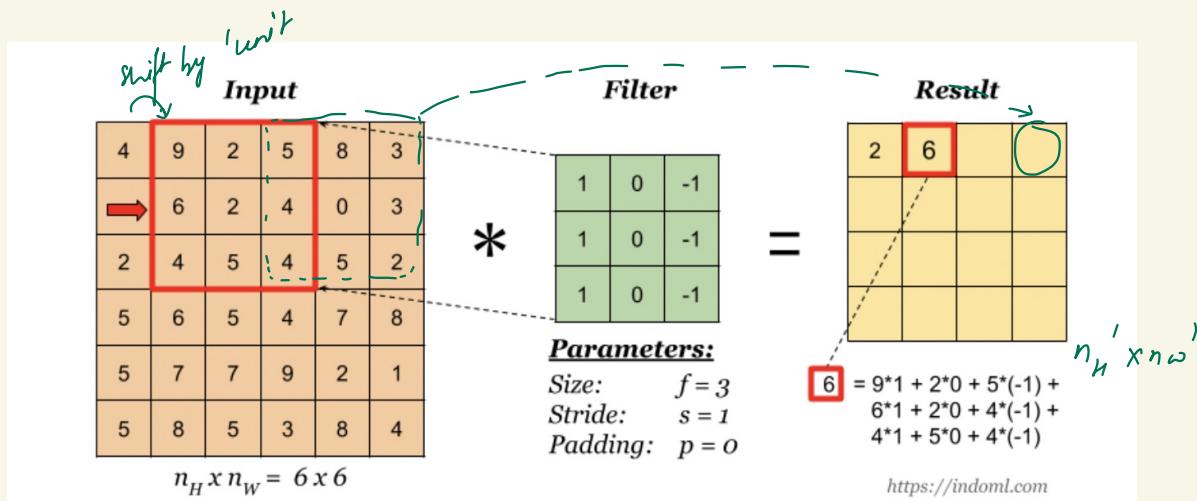
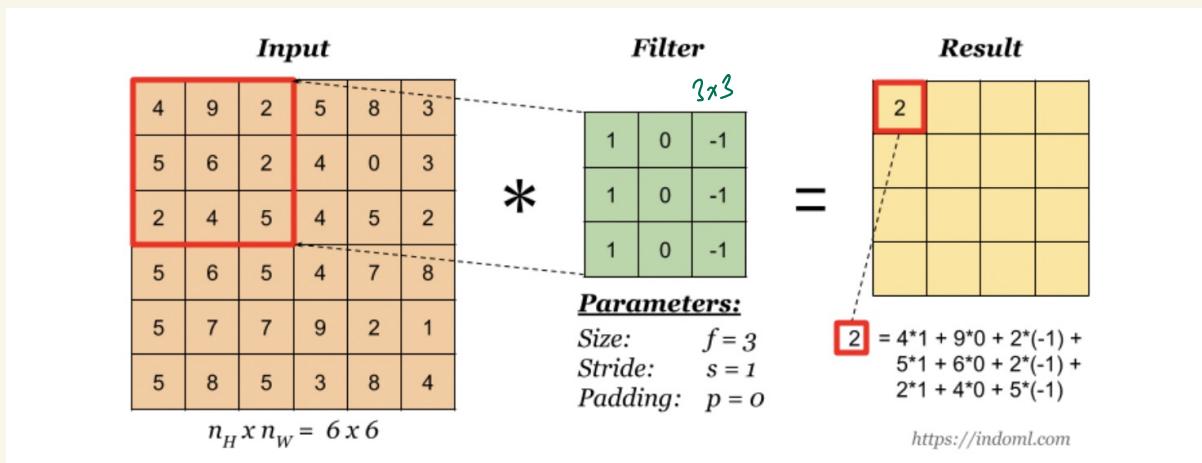
1	0	1
0	0	0
-1	0	-1

$3 \times 3$

2	3	4	2	1
6	9	2	1	6
2	8	7	2	7
1	1	2	6	1
3	3	5	3	2

$8 \times 8$

## ② Convolution



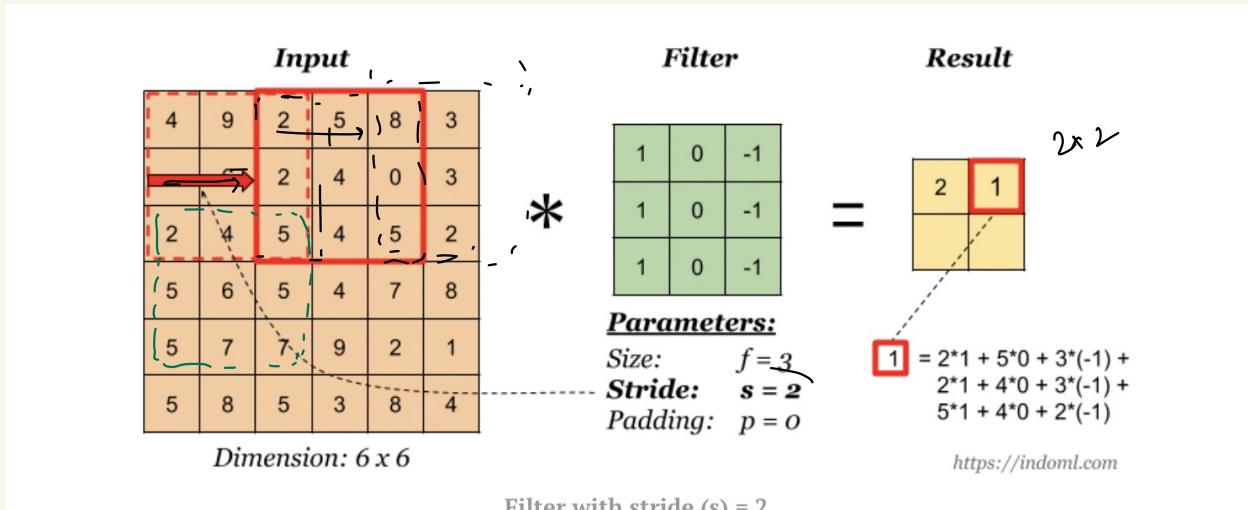
$$n_H' = n_H - f + 1 = \\ = 6 - 3 + 1 = 4$$

stride

$$w = 5-1, p=0$$

$$n_H' = \frac{n_H - f + 1}{s} + 1$$

$$n_H' = \left\lfloor \frac{n_H + 2p - f + 1}{s} \right\rfloor$$



Filter with stride ( $s$ ) = 2

Padding

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} \quad 4 \times 4 \rightarrow 6 \times 6 \rightarrow 8 \times 8$$

$p = 1 \quad p = 2$

since  $p=1$  increase image size by  $2^{units}$

$$n_H' = n_H + 2p - f + 1$$

Convolution reduces size, padding = valid  $\downarrow$  same  $\uparrow$   $\frac{p}{s} = \frac{f-1}{2}$   
 reduces size of image  $\downarrow$  keeping original size of image same  $= \frac{3-1}{2} = 1$   
 $6 \times 6 \rightarrow 6 \times 6$

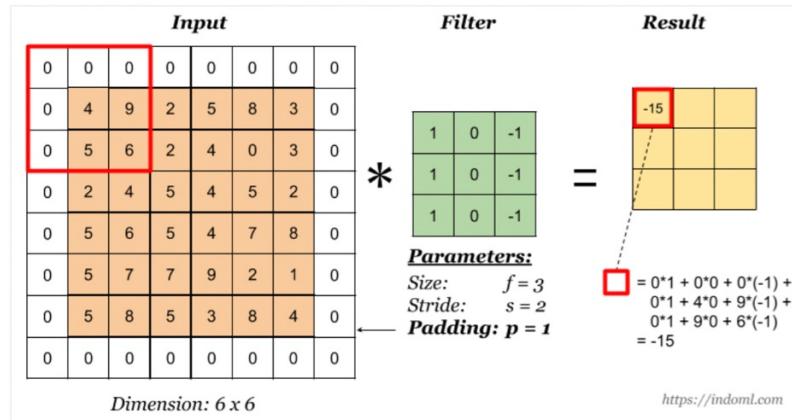
$$n_H' = \left\lfloor \frac{n_H + 2p - f}{s} + 1 \right\rfloor = \underline{\underline{n_H}}$$

$$= R = \frac{f-1}{2}$$

# Padding

Padding has the following benefits:

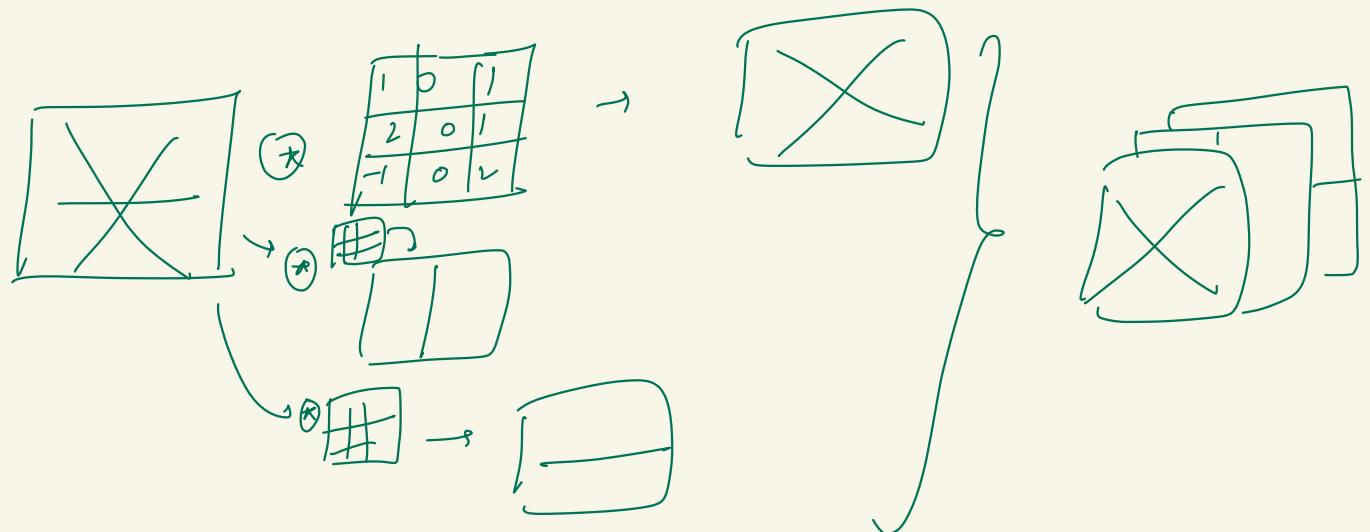
1. It allows us to use a CONV layer without necessarily shrinking the height and width of the volumes. This is important for building deeper networks, since otherwise the height/width would shrink as we go to deeper layers.
2. It helps us keep more of the information at the border of an image. Without padding, very few values at the next layer would be affected by pixels as the edges of an image.

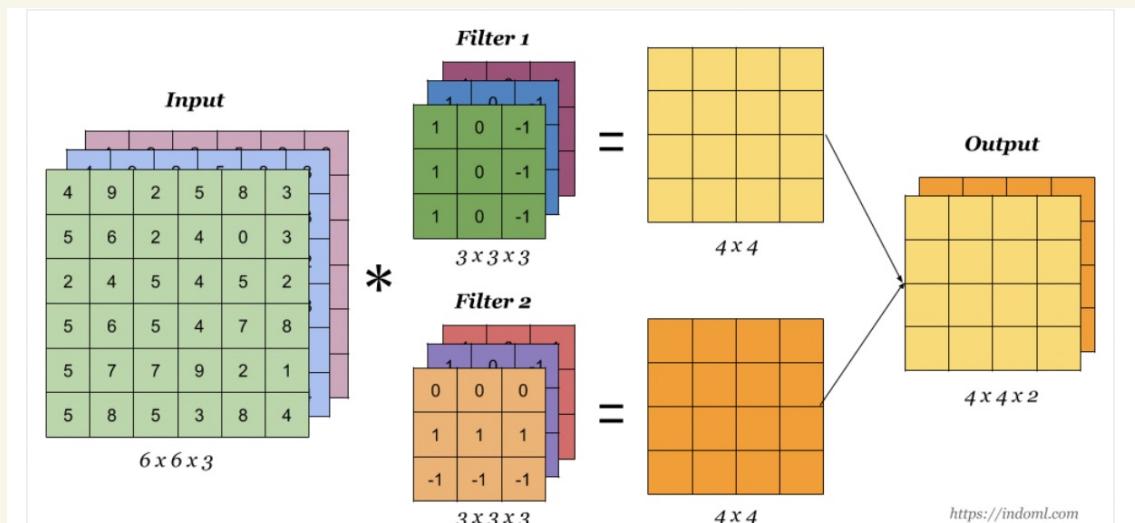
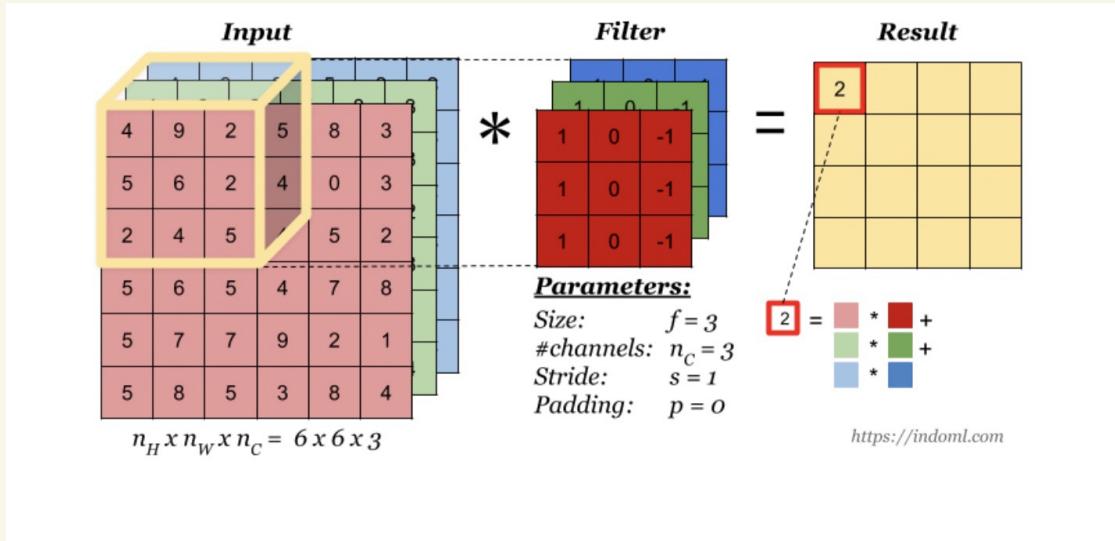


Notice the the dimension of the result has changed due to padding. See the following section on how to calculate output dimension.

Some padding terminologies:

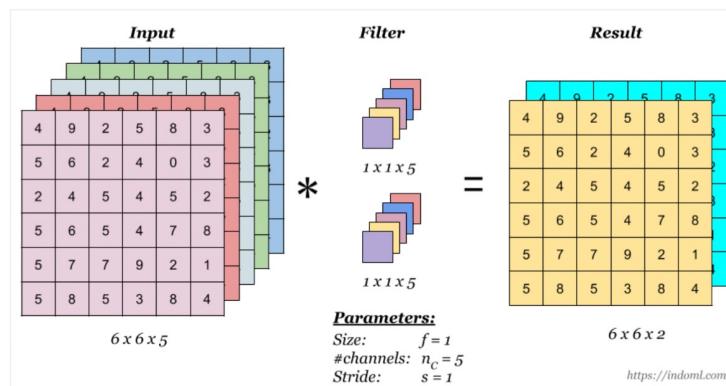
- “**valid**” padding: no padding
- “**same**” padding: padding so that the output dimension is the same as the input

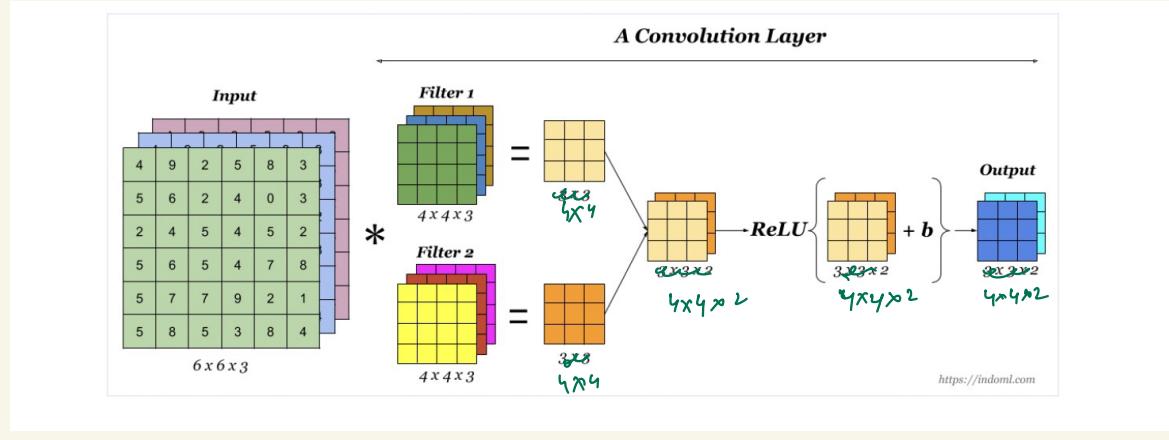




## 1 x 1 Convolution

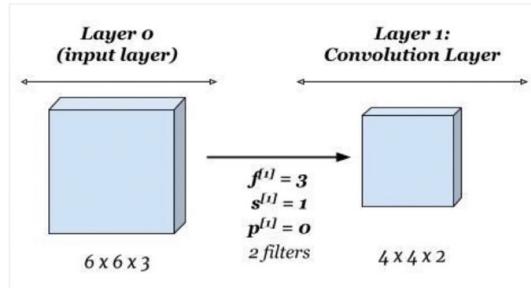
This is convolution with  $1 \times 1$  filter. The effect is to flatten or “merge” channels together, which can save computations later in the network:





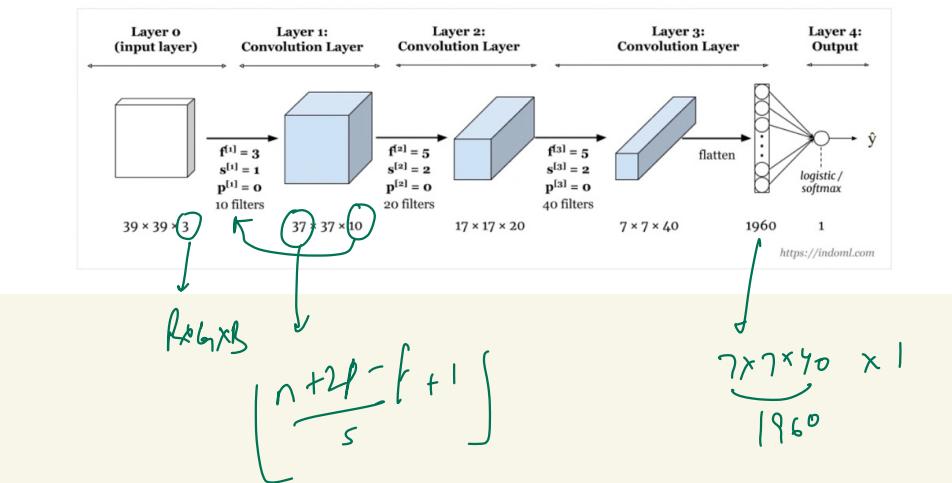
## Shorthand Representation

This simpler representation will be used from now on to represent one convolutional layer:

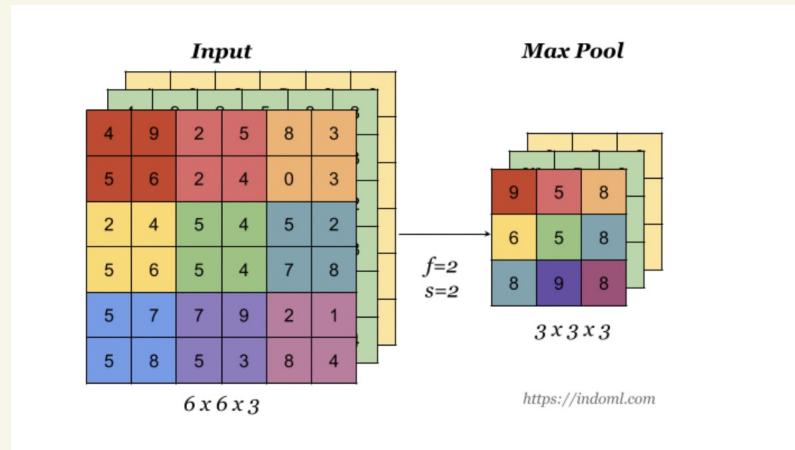
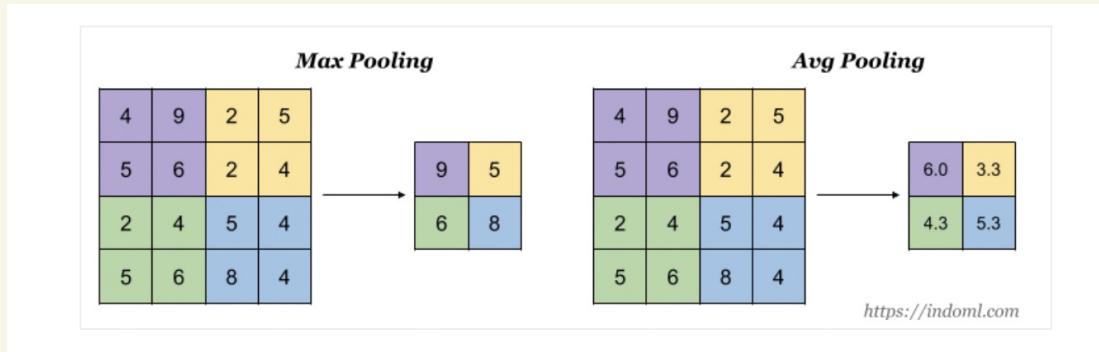


## Sample Complete Network

This is a sample network with three convolution layers. At the end of the network, the output of the convolution layer is flattened and is connected to a logistic regression or a softmax output layer.



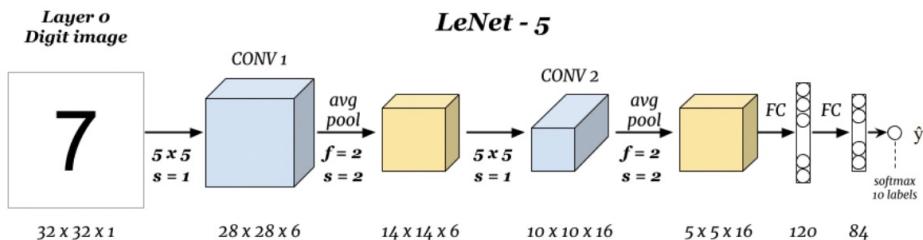
Pooling  
 ↗ Avg  
 ↗ Max (more used)  
 My parameters: size ( $f$ ) , stride(s) , type (max / avg)



## Well Known Architectures

### Classic Network: LeNet – 5

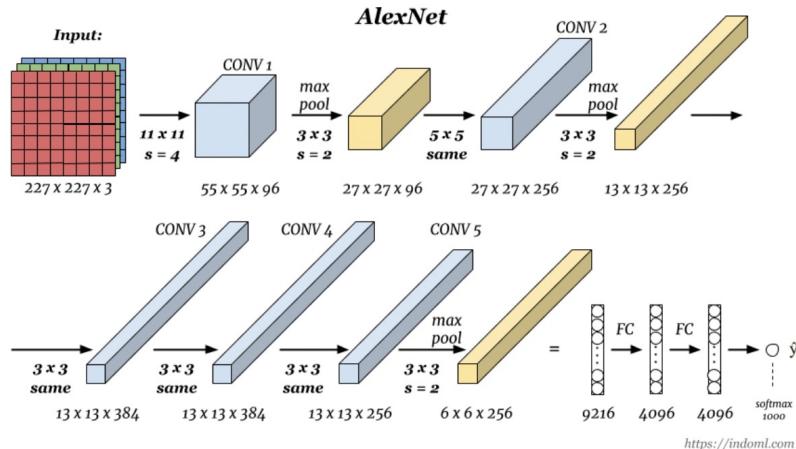
One example of classic networks is LeNet-5, from *Gradient-Based Learning Applied to Document Recognition* paper by Y. Lecun, L. Bottou, Y. Bengio and P. Haffner (1998):



- Number of parameters: ~ 60 thousands.

## Classic Network: AlexNet

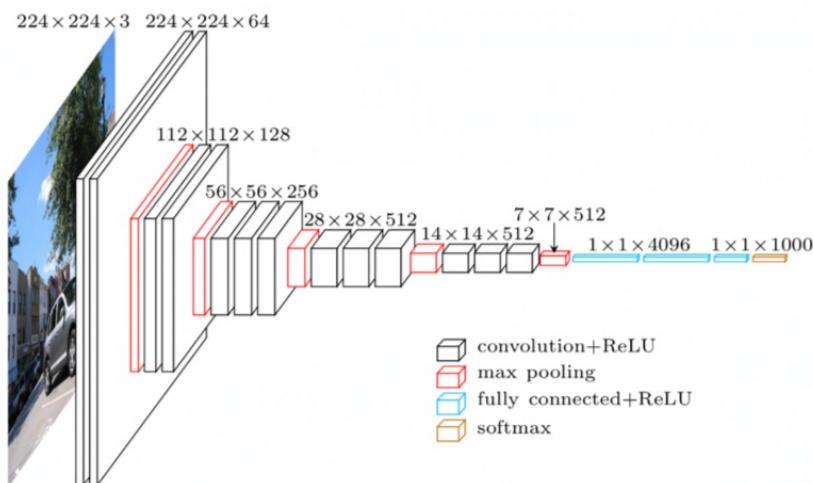
AlexNet is another classic CNN architecture from *ImageNet Classification with Deep Convolutional Neural Networks* paper by Alex Krizhevsky, Geoffrey Hinton, and Ilya Sutskever (2012).



- Number of parameters: ~ 60 millions

## Classic Network: VGG-16

VGG-16 from *Very Deep Convolutional Networks for Large-Scale Image Recognition* paper by Karen Simonyan and Andrew Zisserman (2014). The number 16 refers to the fact that the network has 16 trainable layers (i.e. layers that have weights).

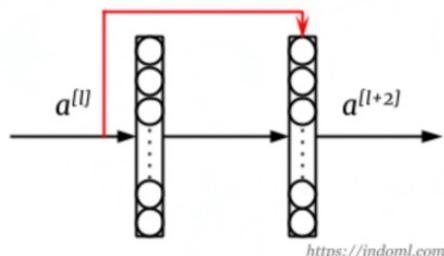


(image from [blog.heuritech.com](http://blog.heuritech.com))

- Number of parameters: ~ 138 millions.
- The strength is in the simplicity: the dimension is halved and the depth is increased on every step (or stack of layers)

## ResNet

The problem with deeper neural networks are they are harder to train and once the number of layers reach certain number, the training error starts to raise again. Deep networks are also harder to train due to exploding and vanishing gradients problem. ResNet (Residual Network), proposed by He et al in [\*Deep Residual Learning for Image Recognition paper\*](#) (2015), solves these problems by implementing skip connection where output from one layer is fed to layer deeper in the network:



In the image above, the skip connection is depicted by the red line. The activation  $a^{[l+2]}$  is then calculated as:

$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

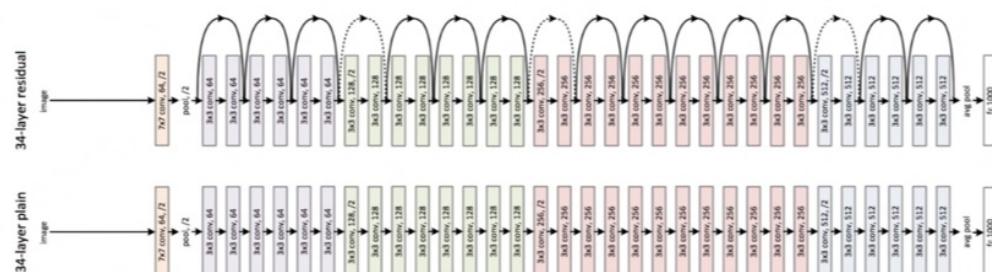
$$a^{[l+2]} = g^{[l+2]}(z^{[l+2]} + \color{red}{a^{[l]}})$$

The advantages of ResNets are:

- performance doesn't degrade with very deep network
  - cheaper to compute
  - ability to train very very deep network

ResNet works because:

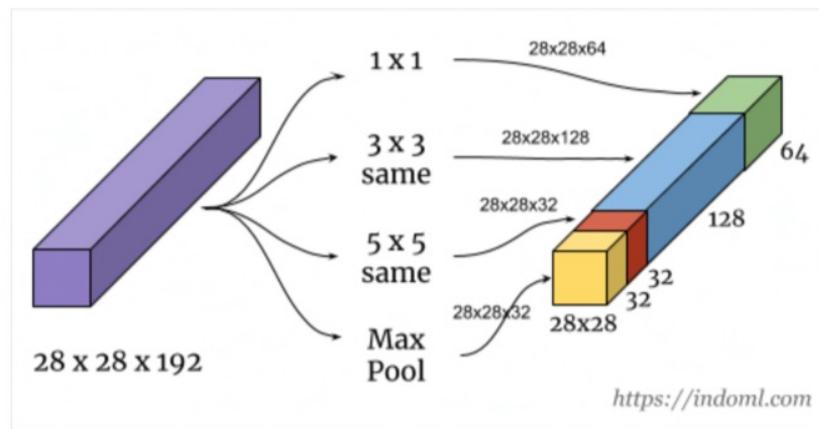
- identifying function is easy for residual block to learn
  - using a skip-connection helps the gradient to back-propagate and thus helps you to train deeper networks



Comparison of 34 layers ResNet with plain network (image from [euler.stat.yale.edu](http://euler.stat.yale.edu))

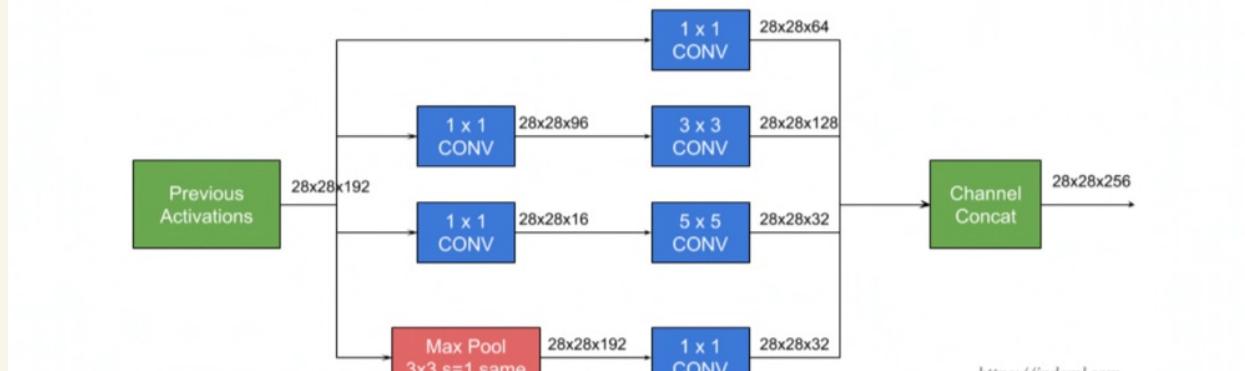
# Inception

The motivation of the inception network is, rather than requiring us to pick the filter size manually, let the network decide what is best to put in a layer. We give it choices and hopefully it will pick up what is best to use in that layer:

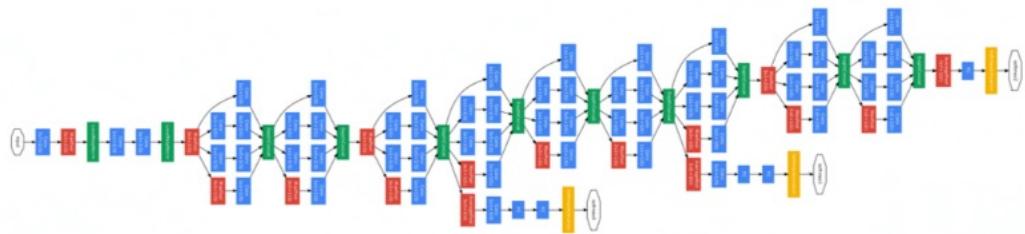


The problem with the above network is computation cost (e.g. for the  $5 \times 5$  filter only, the computation cost is  $(28 \times 28 \times 32) \times (5 \times 5 \times 192) = \sim 120$  millions).

Using  $1 \times 1$  convolution will reduce the computation to about 1/10 of that. With this idea, an inception module will look like this:



Below is an inception network called **GoogLeNet**, described in *Going Deeper with Convolutions paper* by Szegedy et all (2014), which has 9 inception modules:



GoogLeNet architecture (image [source](#))

## Data Augmentation

- Mirroring
- Random cropping
- Less common (perhaps due to their complexity):
  - Rotation
  - Shearing
  - Local warping
- Color shifting  $\rightarrow$  4

