

## UNIT-2

### Artificial Neural Networks

- \* Neural Network learning methods, provide a robust approach to approximations real-valued, discrete-valued and vector-valued target functions.
- \* ANN learning is well matched for certain types of problems, like:
  - interpret complex real-world sensor data.
  - Speech recognition.
  - visual scenes.
  - Learning robot control strategies.
- \* The Back propagation algorithm is one of the efficient ANN algorithm, which is successful in many problems such as

- Learning to recognize handwritten characters.

- Learning to recognize face.

- Learning to recognize spoken word.

\* The Back propagation algorithm uses gradient descent to tune network parameters to best fit a training set of input-output pairs

Neurons and Biological Motivation

\* The study of ANN has been inspired by the observation that biological learning systems are built of very high complex, non-linear parallel interconnections of neurons.

\* The human body is made up of

billions of cells. Cells of the nervous system called nerve cells or neurons, are specialized to carry "messages" through an electrochemical process.

\* The information processing abilities of biological neural systems must follow from highly parallel processes, operations on representations that are distributed over many neurons.

One motivation for ANN is to capture this kind of highly parallel computation based on distributed representations.

- \* historically, two groups of researchers have worked with artificial neural networks.
- \* ANNs to study and model biological learning process.

- To obtain highly effective machine learning algorithms

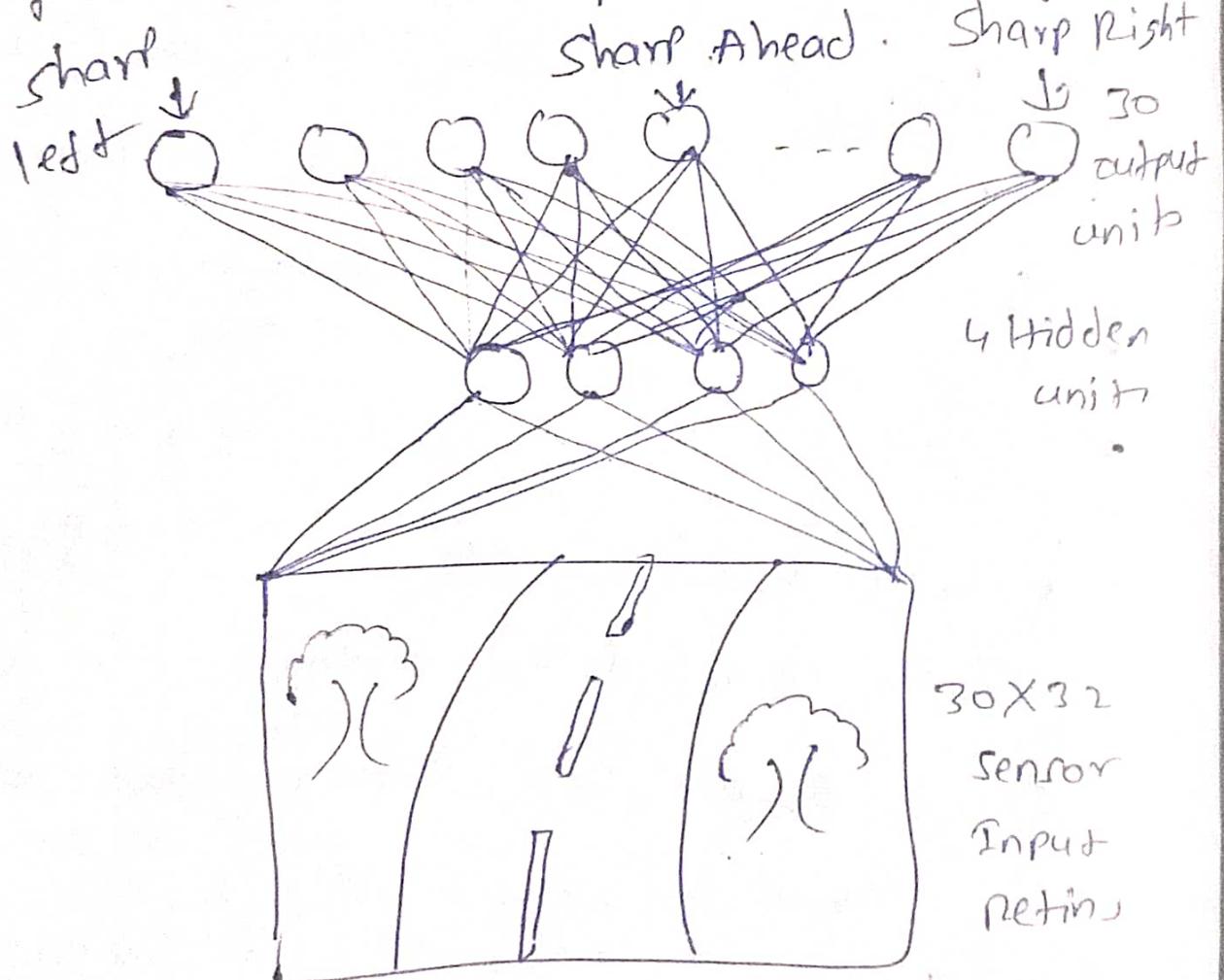
## Neural Network Representations

\* A prototypical example of ANN learning is AlvinN (autonomous land vehicle). In a Neural network is a perception system, which uses a learned ANN to steer an autonomous vehicle driving at normal speeds on public highways.

\* Input:- The input to the neural network is a  $30 \times 32$  grid of pixel intensities obtained from a forward-pointed camera mounted on the vehicle.

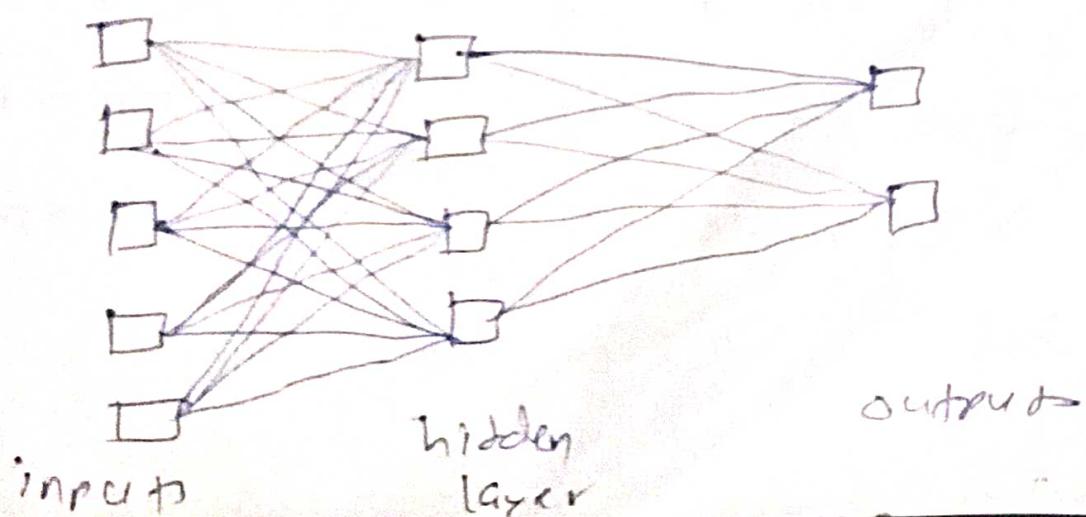
\* Output:- The network's output is the direction in which the vehicle is steered.

\* The ANN is trained to mimic the observed steering commands of a human driving the vehicle for approximately 5 minutes. ALVINN has used its learned networks to successfully drive at speeds up to 70 miles per hour and for distances of 90 miles on public highways.



- \* The ALVINN system uses Backpropagation to learn to steer an autonomous vehicle (photo at top) driving at speeds up to 70 miles per hour.
- \* The diagram shows how the image of a forward-mounted camera is mapped to 960 neural network inputs which are fed forward to 4 hidden units connected to 20 output units

\* Each ANN is composed of a collection of perceptrons grouped in layers. A typical structure is shown below



Note:- The three layers: input, output and hidden layer several hidden layers can be placed between the input and output layers

### Problems for Neural network learning

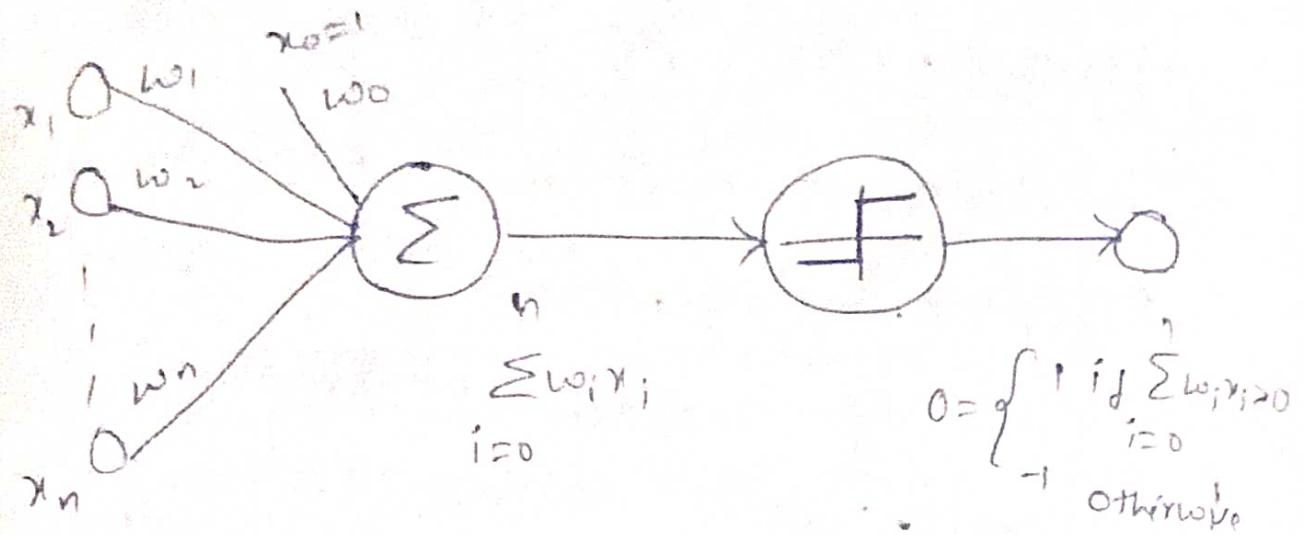
- \* ANN learning is well suited to problems in which the training data corresponds to noisy, complex sensor data. It is also applicable to problems for which more symbolic representations are used.
- \* The back propagation (BP) algorithm is the most commonly used ANN learning technique. It is appropriate for problems with the characteristics:-

① Instances are represented by many attribute-value pairs.

- 2) The target function output may be discrete-valued, real-valued, or a vector of several real- or discrete-valued attributes.
- 3) The training examples may contain errors.
- 4) Long training time are acceptable.
- 5) Fast evaluation of the learned target function may be required.
- 6) The ability of humans to understand the learned target functions is not important.

## Perceptron

\* one type of ANN system is based on a unit called a perceptron.



### A Perception

- \* An ANN consist of perceptrons; each of the perceptrons receives input, processes input and delivers a single output

- \* A perceptron take a vector of real-valued inputs, then calculates linear combination of these inputs after that it returns 1 if the result is greater than some threshold and -1 otherwise as its

$$O(x_0, x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

- \* where each  $w_i$  is a real-valued constant, or weight, that determines the contribution of input  $x_i$  to the perceptron output.

### Representational Power of perceptron

- \* we can view the perceptron as representing a hyperplane decision surface in the n-dimensional space of instances. (point).

- \* The perceptron outputs a +1 for instances lying on one side of the hyperplane and outputs a -1 for instances lying on the other side.

- \* Examples can be separated by a single perceptron are called linearly separated set of examples.



- \* A single perceptron can be used to represent many boolean functions.
  - \* For example, if we assume boolean values of 1 (true) and -1 (false), then one way to use a two-input perceptron to implement the AND function is to set the weights  $w_0 = -0.5$ , and  $w_1 = w_2 = 0.5$ .

- \* The perceptron can be made to represent the OR function instead by altering the threshold to  $w_0 = -0.3$ .
- \* Unfortunately, however, some boolean functions cannot be represented by a single perceptron, such as the XOR function. These are called linearly non separable training examples.

- \* A single perceptron can be used to represent many boolean functions.

AND function (linearly separable):-

Trainings examples

$x_1$	$x_2$	output
0	0	-1
0	1	-1
1	0	-1
1	1	1

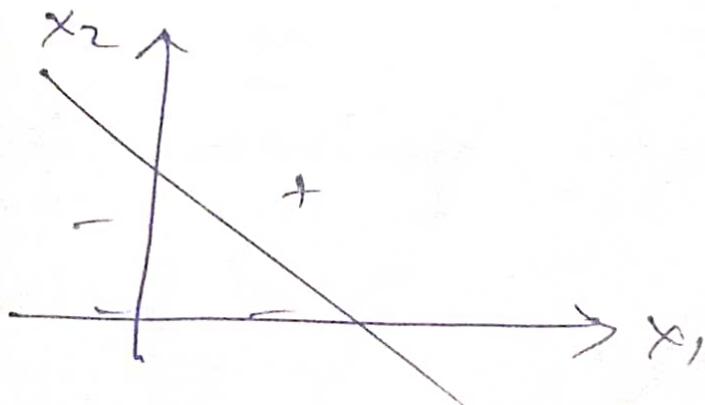
test results

$x_1$	$x_2$	$w_i x_i$	output
0	0	-0.8	-1
0	1	-0.3	-1
1	0	-0.2	-1
1	1	0.2	1

Decision hyper plane:-

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

$$-0.8 + 0.5 x_1 + 0.5 x_2 = 0$$



$$-0.8 + 0.5 x_1 + 0.5 x_2 = 0$$

## OR function (Linearly Separable)

The two-input perceptron can implement the OR function when we set the weights:  $w_0 = -0.3w_1$ ,

### Training examples

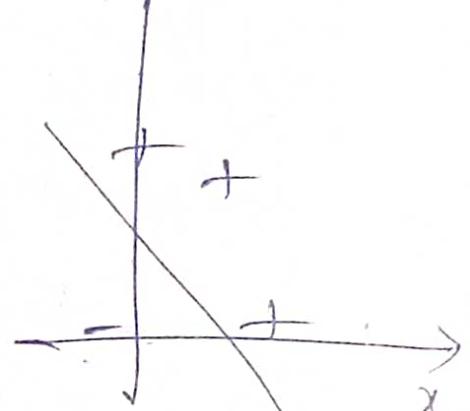
$x_1$	$x_2$	Output
0	0	-1
0	1	1
1	0	1
1	1	1

### Decision hyperplane

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

$$-0.3 + 0.5x_1 + 0.5x_2 = 0$$

$x_2$



$$-0.3 + 0.5x_1 + 0.5x_2 = 0$$

### Test results

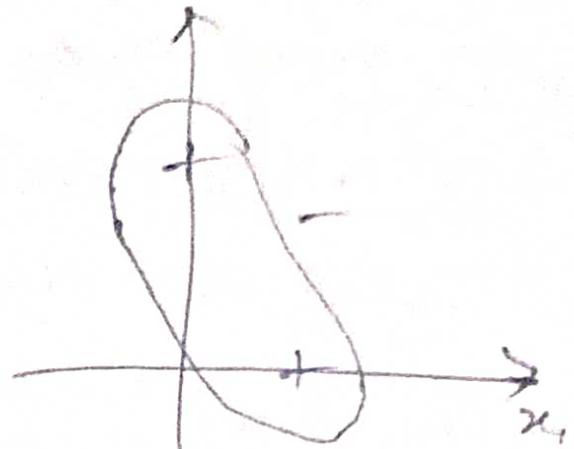
$x_1$	$x_2$	$w_1 x_1$	Output
0	0	-0.7	-1
0	1	0.2	-1
1	0	0.2	-1
1	1	0.7	1

## XOR function (Non linearly Separable)

\* It is impossible to implement the XOR function by a single exception.

## Training Example

$x_1$	$x_2$	Output
0	0	-1
0	1	1
1	0	1
1	1	-1



A two-layer network of perceptrons can represent XOR function

$$\boxed{x_1 \oplus x_2 = x_1 \bar{x}_2 + \bar{x}_1 x_2}$$

## Perceptron Training Rule

- \* They are important to ANNs because they provide the basis for learning networks of many units.
- \* Several algorithms are known to solve this learning problem.

Here we consider two:-

- \* The perceptron rule.
- \* The delta rule

## The Perceptron Rule:-

- \* One way to learn an acceptable weight vector is to begin with random weights, then iteratively apply the perceptron to each training example, modifying the perceptron weights whenever it misclassifies an example.
- \* The process is repeated until the perceptron classifies all training examples correctly.

\*  $w_i \leftarrow w_i + \Delta w_i$

when  $\Delta w_i = \eta(t - o) x_i$

$t \rightarrow$  is the target output for the current training example

$o \rightarrow$  is the output generated by the perceptron.

$\eta \rightarrow$  is positive constant called learning rate

Example:-

$$x_i = 0.8$$

$$\eta = 0.1$$

$$t = 1$$

$$o = -1$$

$$\Delta w_i = \eta(t-o)x_i$$

$$= 0.1(1 - (-1))0.8$$

$$= 0.16$$

In this example we calculated the Perceptron.

Gradient Descent and the Delta Rule:

- \* The perceptron rule finds a successful weight vector when the training examples are linearly separable, it can fail to converge if the examples are not linearly separable.

- \* The second training rule, called the delta rule, is designed to overcome this difficulty.

\* The key idea of rule is to use gradient descent to search the space of possible weight vector to find the weights that best fit the training examples.

\* The delta training rule is best understood by considering the task of training an unthresholded perceptron. That is, a linear unit for which the output  $o$  is given by

$$o(\vec{x}) = \vec{w} \cdot \vec{x}$$

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{id}$$

Update weight as

$$w_i \leftarrow w_i + \Delta w_i$$

## Gradient - Descent Algorithm:-

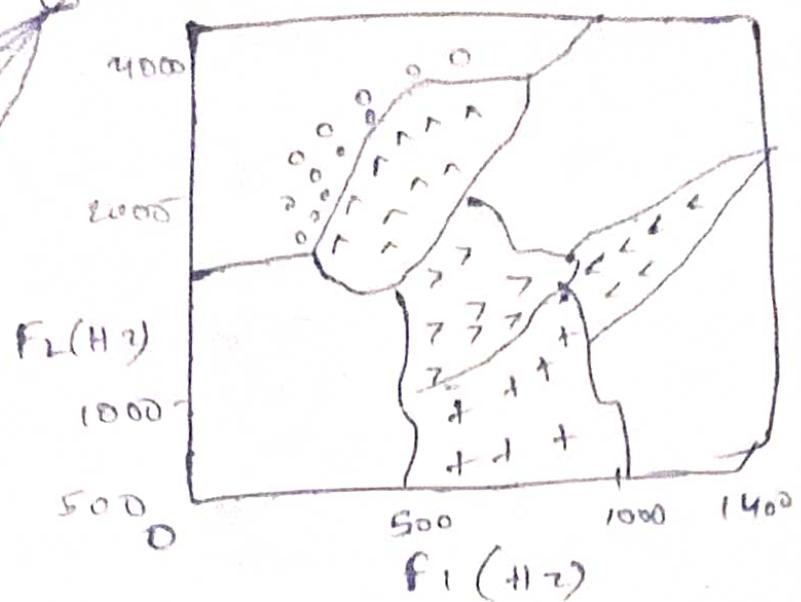
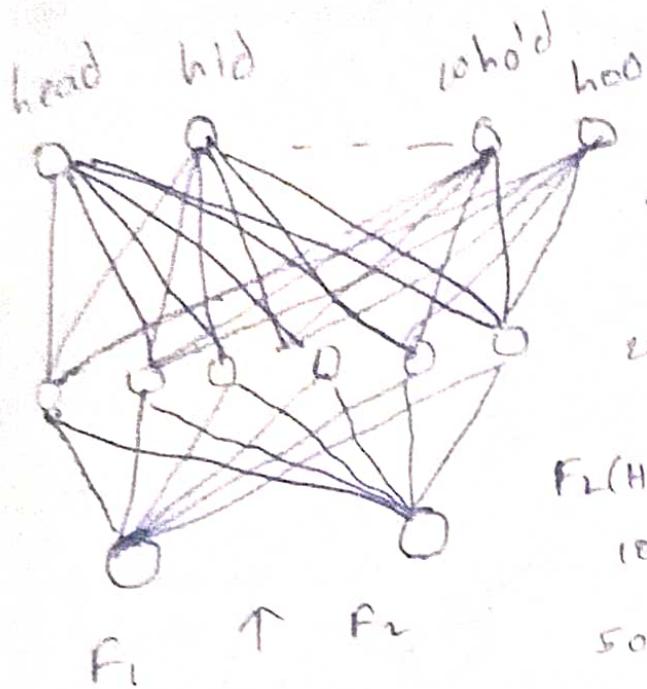
Each training example is a pair of the form  $(\vec{x}, t)$ , where  $\vec{x}$  is the vector of input values, and  $t$  is the target value.  $\eta$  is the learning rate (e.g., .05).

- Initialize each  $w_i$  to some small random value
- until the termination condition is met, DO:
  - Initialize each  $\Delta w_i$  to zero.
  - For each  $(\vec{x}, t)$  in training examples, DO
    - \* Input the instance  $\vec{x}$  to the unit and compute the output  $O$
    - \* For each linear unit weight  $w_i$ , DO
$$\Delta w_i \leftarrow \Delta w_i + \eta(t-O)x_i$$

- for each ~~by~~ linear unit weight  
 $w_i, D_i$   
 $\omega_i \leftarrow \omega_i + \Delta \omega_i$

## Multilayer Networks

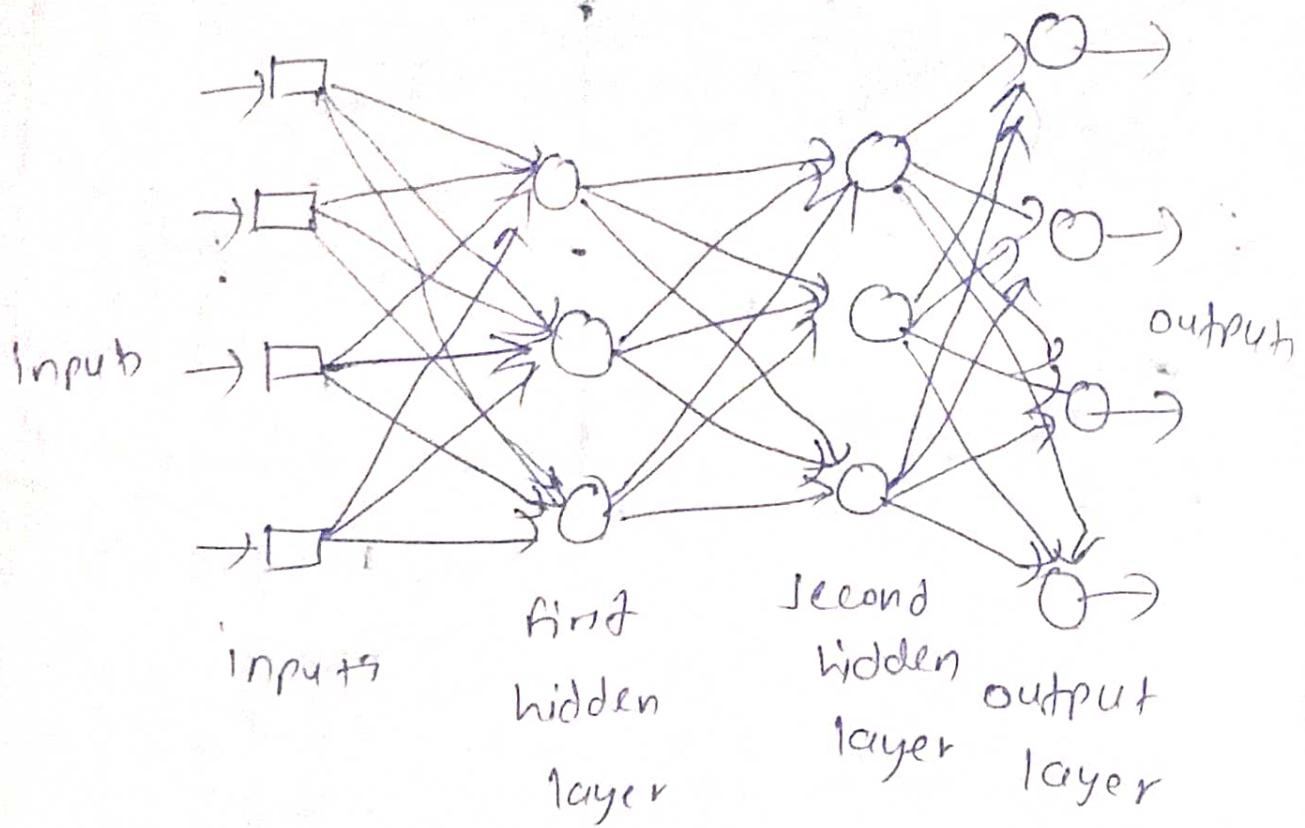
- \* The kind of multilayer networks learned by the Backpropagation algorithm are capable of expressing a rich variety of nonlinear decision surfaces
- \* for example, speech recognition task involves distinguishing among 10 possible vowels, all spoken in the context of "h-d". ("hid", "had", "head", "hood", etc).



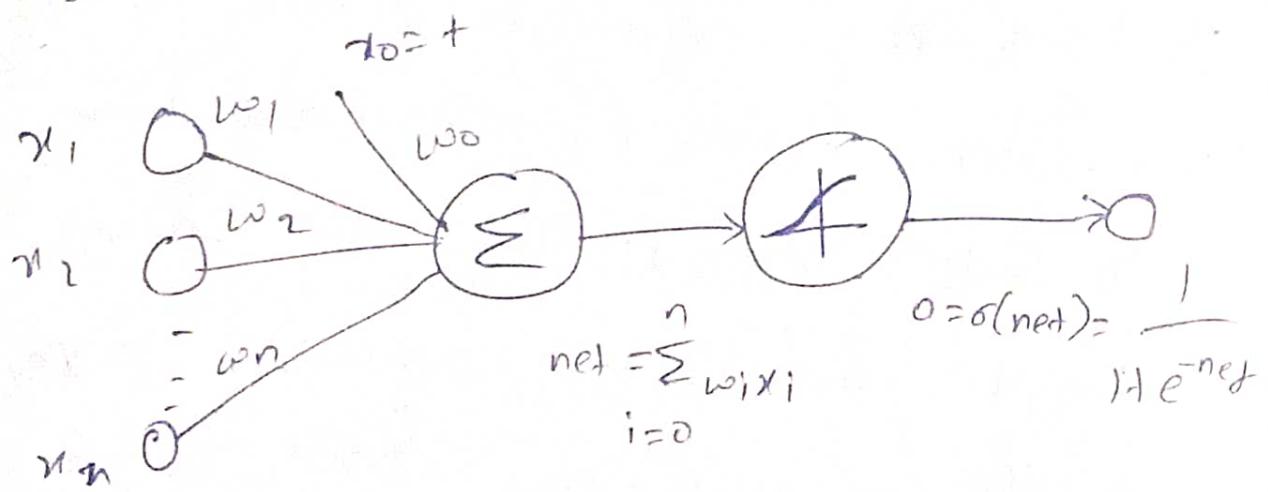
o head  
 + hid  
 < hid  
 ^ hood  
 > whold

The network input consists of two parameters,  $f_1$  and  $f_2$  obtained from a spectral analysis of the sound. The LO network output corresponds to the 10 possible vowel sounds. The plot on the right illustrates the highly nonlinear decision surface.

# multilayer network



A Differentiable Threshold unit :-



$\sigma(x)$  is the sigmoid function

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

\* like the perceptron, the sigmoid units first computes a linear combination of its inputs, then applies a threshold to the result.

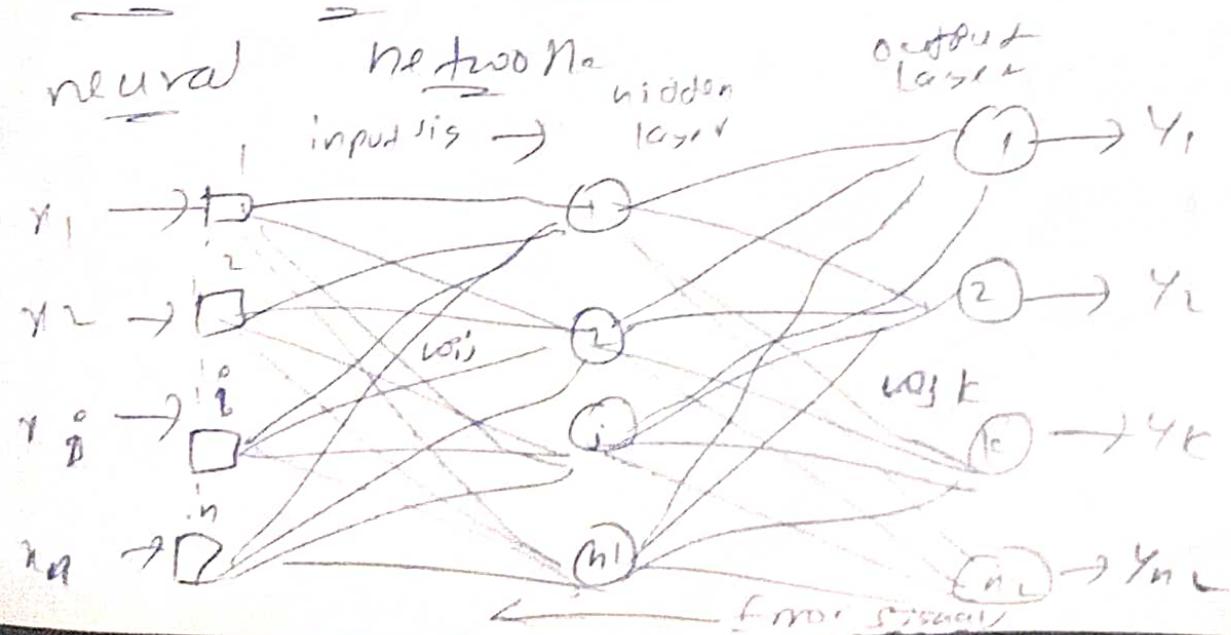
In the case of sigmoid units, however, the threshold output is a continuous function of its input.

The sigmoid function  $\sigma(x)$  is also called the logistic function.

Interesting properties:-

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1-\sigma(x))$$

Two-layer back-propagation

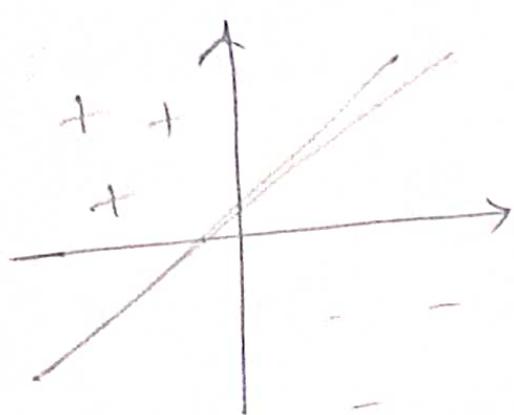


## The Back propagation (BP) Algorithm:

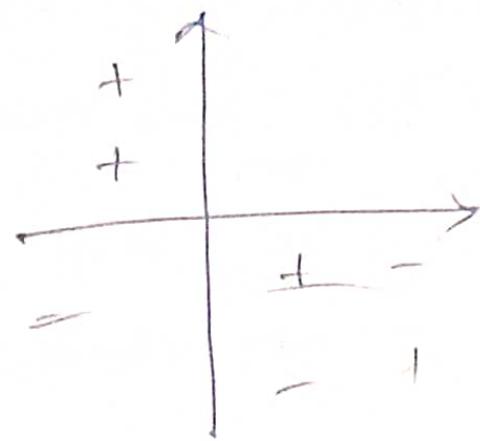
- \* The BP algorithm learns the weights for a multi-layer network, given a network with a fixed set of units and interconnections. It employs a gradient descent to attempt to minimize the squared error between the network output values and the target values for those outputs.
- \* The learning problem in Back Propagation is to search a large hypothesis space defined by all possible weight values for all the units in the network.
- \* In a multi-layer network the error surface can have multiple local minima.

\* The gradient descent is guaranteed only to converge toward some local minimum, but not necessarily to global minimum error.

### Gradient Descent and the Delta Rule



linearly separable

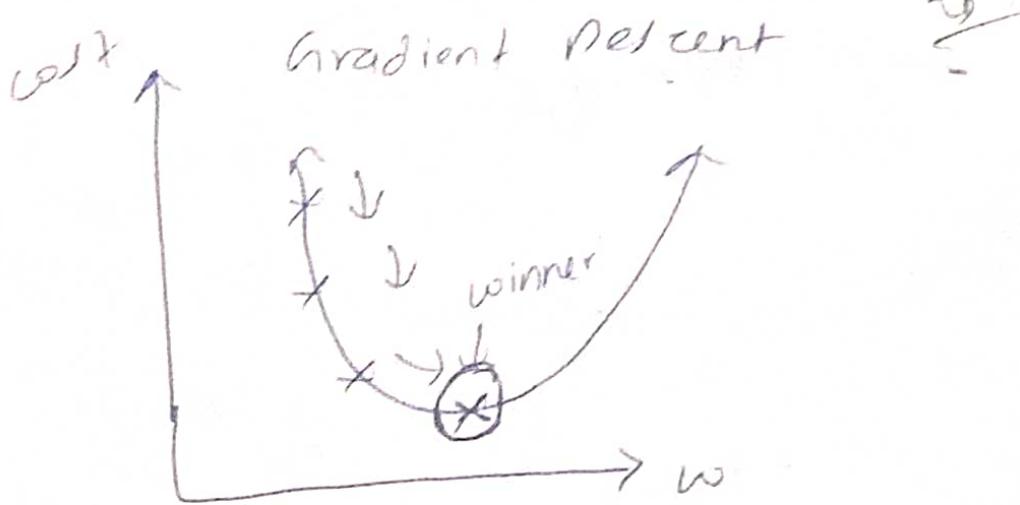
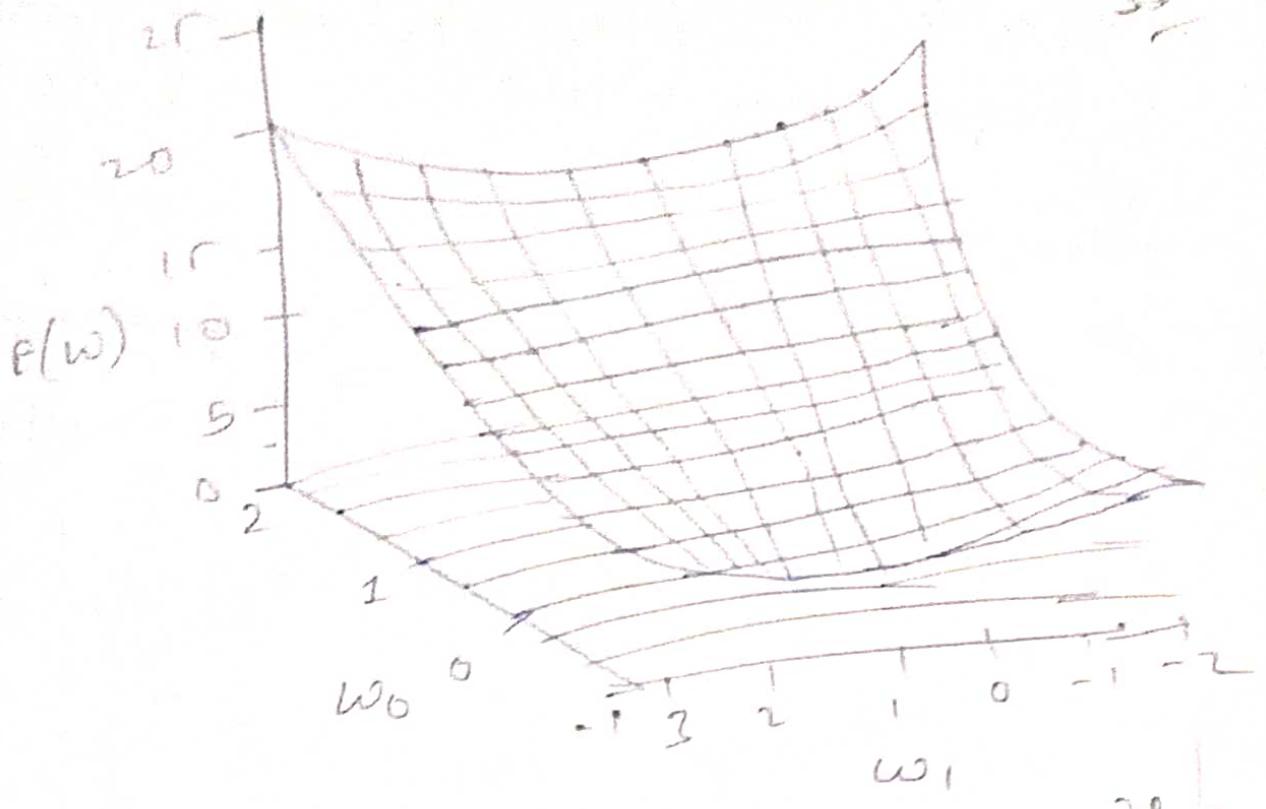


Non-linearly Separable

\* perceptron rule finds a successful weight vector when the training examples are linearly separable, but it can fail to converge if the examples are not linearly separable.

\* A second training rule, called the delta rule, is designed to overcome this difficulty.

- \* if the training examples are not linearly separable, the delta rule converges toward a best-fit approximation to the target concept.
- \* The key idea behind the delta rule is to use gradient descent to search the hypothesis space of possible weight vectors to find the weights that best fit the training examples.
- \* This rule is important because gradient descent provides the basis for the Back propagation algorithm, which can learn networks with many interconnected units.
- \* It is also important because gradient descent can serve as the basis for learning algorithms that must search hypothesis space containing many different types of continuously parameterized hypotheses.



$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$f_{\alpha_i} = w_{(\alpha_i, i)}$$

$$w_i \leftarrow w_i + \Delta w_i$$

where

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

learning rate

$$\begin{aligned}
 \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \\
 &= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\
 &= \frac{1}{2} \sum_{d \in D} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\
 &= \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\
 &= \sum_{d \in D} (t_d - o_d) (-x_{td})
 \end{aligned}$$

output calculated  
 by  
 $o = w_0 + w_1 x_1 + \dots + w_n x_n$

$$\Delta w_i = n \sum_{d \in D} (t_d - o_d) x_{id} \quad \Delta w_i = -n \frac{\partial E}{\partial w_i}$$

$$w_i \leftarrow w_i + \Delta w_i$$

## Back Propagation Algorithm

(training-example,  $\eta$ ,  $n_{in}$ ,  $n_{out}$ ,  $n_{hidden}$ )

- Each training example is a pair of the form  $(x, t)$ , where  $(x)$  is the vector of network input values, and  $(t)$  is the vector of target network output values.
- $\eta$  is the learning rate (e.g.: 0.05).
- $n_i$  is the number of network inputs.
- $n_{hidden}$  the number of units in the hidden layer, and
- $n_{out}$  the number of output units.
- The input from unit  $i$  into unit  $j$  is denoted  $x_{ji}$ , and the weight from unit  $i$  to unit  $j$  is denoted  $w_{ji}$

## Algorithm

- Create a feed-forward network with  $n_i$  inputs,  $n_{hidden}$  units, and  $n_{output}$  units.

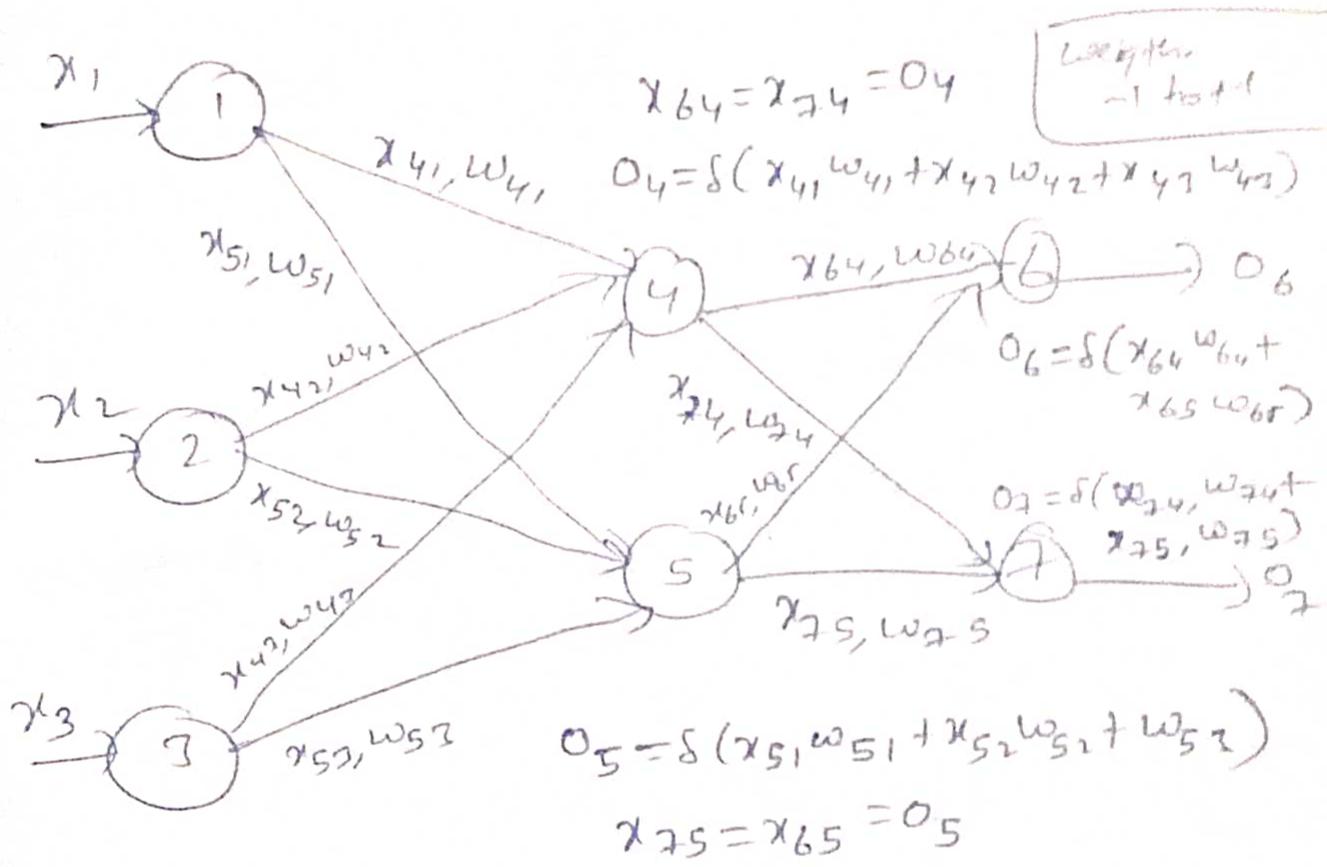
- Initialize all network weights to small random numbers.
  - until the termination condition is met, do
    - for each  $(x, t)$ , in training examples,
1. propagate the input forward through the network;
1. Input the instances  $x, t$  to the network and compute the output  $o_u$  of every unit  $u$  in the network.
  - propagate the errors backward through the network.
  2. for each network unit  $k$ , calculate its error term  $\delta_k$
- $$\delta_k \leftarrow o_k (1 - o_k) (t_k - o_k)$$
3. for each network unit  $h$ , calculate its error term  $\delta_h$
- $$\delta_h \leftarrow o_h (1 - o_h) \sum_{k \in \text{outputs}} w_{k,h} \delta_k$$

4. update each network weight  $w_{ji}$

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

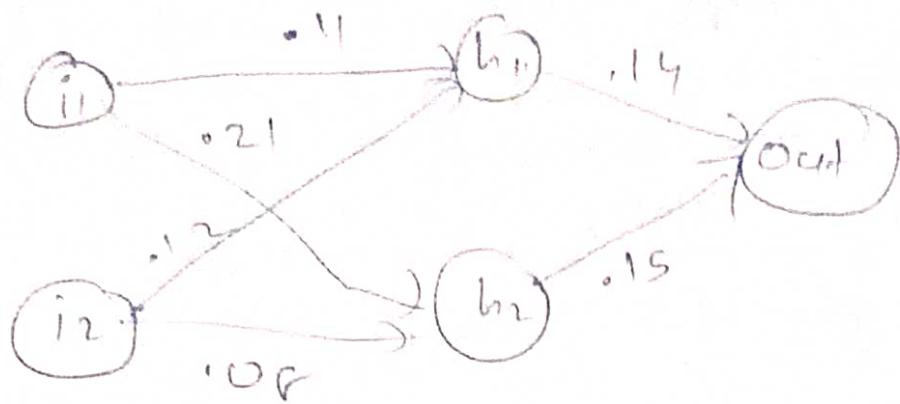
where

$$\Delta w_{ji} = \eta \delta_j \cdot x_{ji}$$



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

## Example



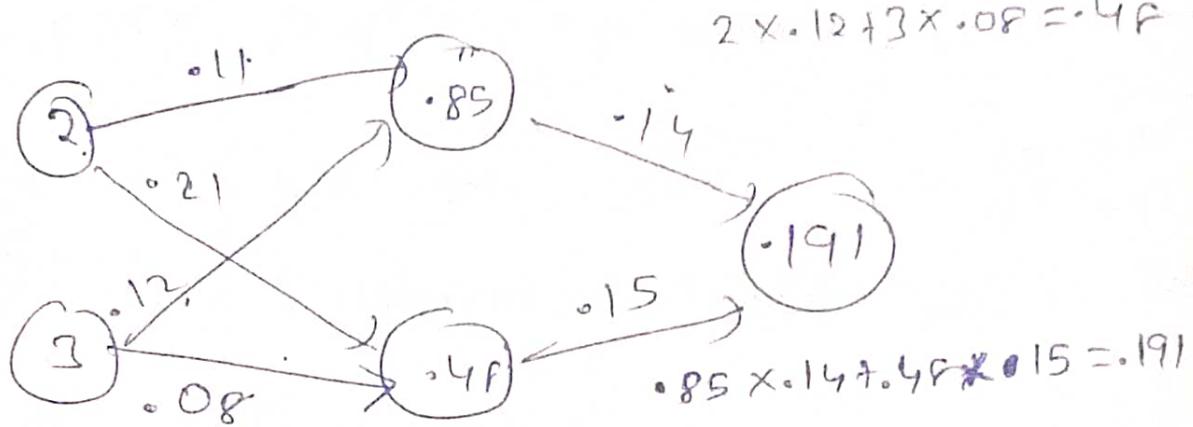
- \* Neural network training is about finding weights that minimize prediction error.
- \* we usually start our training with a set of randomly generated weights.
- \* Then, backpropagation is used to update the weights in an attempt to correctly map arbitrary inputs to outputs.
- \* Our initial weights will be as follows:  $w_1 = 0.11$ ,  $w_2 = 0.21$ ,  $w_3 = 0.12$ ,  $w_4 = 0.08$ ,  $w_5 = 0.14$  and  $w_6 = 0.15$ .

Dataset

our single sample is as follows

$$\text{input} = [2, 3] \text{ and output} = 1$$

Forward Pass



$$2 \times 0.11 + 3 \times 0.21 = 0.85$$

$$2 \times 0.12 + 3 \times 0.08 = 0.48$$

$$0.85 \times 0.14 + 0.48 \times 0.15 = 0.191$$

$$\begin{bmatrix} 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 0.11 & 0.12 \\ 0.21 & 0.08 \end{bmatrix} = \begin{bmatrix} 0.85 & 0.48 \\ 0.14 & 0.15 \end{bmatrix}$$

$$= \begin{bmatrix} 0.191 \end{bmatrix}$$

Calculating Error

$$\text{Error} = \frac{1}{2} (\text{Prediction} - \text{actual})^2$$

$$= \frac{1}{2} (0.191 - 1.0)^2 = 0.327$$

## Reducing Error

- Our main goal of the training is to reduce the error.
- Since actual output is constant, the only way to reduce the error is to change prediction value.
- The question now is, how to change prediction value?
- By decomposing prediction into its basic elements we can find that weights are the variable elements affecting prediction value.
- In other words, in order to change prediction value, we need to change weights values.

$$\text{Prediction} = \text{out}$$
$$\downarrow$$

$$\text{Prediction} = (h_1) w_1 + (h_2) w_2$$

$$\text{Prediction} = (i_1 w_1 + i_2 w_2) w_3 + (i_1 w_3 + i_2 w_4) w_5$$

to change prediction value we need  
to change weights.

## Backpropagation

- \* Backpropagation short for "backward propagation of errors", is a mechanism used to update the weights using gradient descent.
- \* It calculates the gradient of the error function with respect to the neural network's weights.
- \* The calculation proceeds backwards through the network

$$w_2 = w_2 - \alpha \left( \frac{\partial \text{Error}}{\partial w_2} \right) \text{with respect to weight}$$

↓      ↓      ↑  
new      old      learning  
weight    weight    rate

\*  $w_6 = w_6 - \alpha(h_2 \cdot A)$

\*  $w_5 = w_5 - \alpha(h_1 \cdot A)$

\*  $w_4 = w_4 - \alpha(i_2 \cdot \Delta w_6)$

$$w_3 = w_3 - \alpha(i_1 \cdot \Delta w_6)$$

$$w_2 = w_2 - \alpha(i_2 \cdot \Delta w_5)$$

$$w_1 = w_1 - \alpha(i_1 \cdot \Delta w_5)$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \alpha \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \begin{bmatrix} a_{11}, \Delta \\ a_{12}, \Delta \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \alpha \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix}$$

$$\begin{bmatrix} w_5 & w_6 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \begin{bmatrix} a_{11}, \Delta w_1, \Delta w_3 \\ a_{12}, \Delta w_5, \Delta w_6 \end{bmatrix}$$

### Backward Pass

$$\Delta = 0.91 - 1 = -0.809$$

$$\alpha = 0.05$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 0.85 \\ 0.48 \end{bmatrix}$$

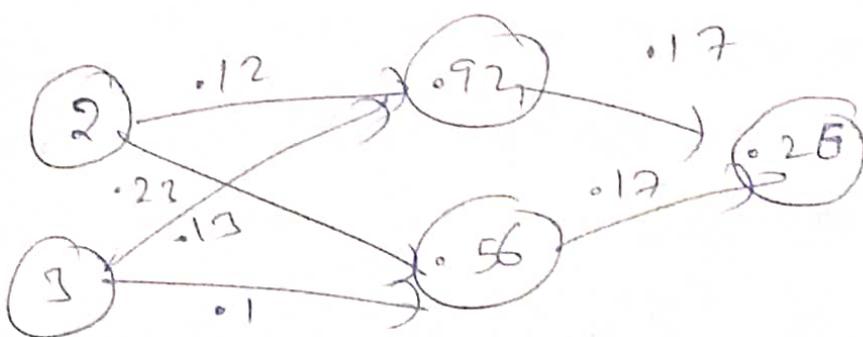
$$= \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - \begin{bmatrix} -0.074 \\ -0.019 \end{bmatrix} = \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .01 \end{bmatrix} = 0.05(-0.809)$$

$$\begin{bmatrix} \frac{2}{2} \\ \frac{2}{2} \end{bmatrix} \cdot \begin{bmatrix} 0.14 & 0.15 \end{bmatrix}$$

$$= \begin{bmatrix} 0.11 & 0.12 \\ 0.21 & 0.01 \end{bmatrix} - \begin{bmatrix} -0.011 & -0.012 \\ -0.017 & -0.014 \end{bmatrix}$$

$$= \begin{bmatrix} 0.12 & 0.13 \\ 0.23 & 0.10 \end{bmatrix}$$



$$\begin{bmatrix} 2 & 2 \end{bmatrix} \cdot \begin{bmatrix} 0.12 & 0.13 \\ 0.23 & 0.10 \end{bmatrix} = \begin{bmatrix} 0.92 & 0.56 \end{bmatrix}$$

$$= \begin{bmatrix} 0.17 \\ 0.12 \end{bmatrix}$$

$$= [0.26]$$

We notice that the prediction 0.26 is a little bit closer to actual output than the previously predicted one 0.191.

- we can repeat the same process of backward and forward pass until error is close or equal to zero.

### Gradient Descent

- \* Error is summed over all examples before updating weights
- \* Summing over multiple examples require more computation per weight update step
- \* Difficult when there are multiple local minimum

### stochastic Gradient Descent

- \* weights are updated examining each training example
- \* Less computation as individual weights are updated
- \* uses various  $E(\omega)$  rather than  $E(\omega)$ , hence handles multiple local minima in the rage.

## Unit 2 Continue

### Remarks on Back Propagation Algorithm:

There are 5 Remarks in it.

1. Convergence to local minima.
2. Representational power of feed forward network.
3. Hypothesis space search and inductive bias.
4. Hidden layer representation.
5. Generalisation, overfitting and stopping criterion.

#### ① Convergence to local minima:

All the neurons are interconnected to each other you can minimize the error using local minima but not global minima because it minimizes locally (within the net)

#### ② Representational power of feed forward network

How detailed we can represent the network depends on the width

and depth of the tree & types of functions used for building a tree are

1) Boolean.

2) continuous : gives us a small error.

3) Arbitrary (random variables)

③ Hypothesis Space Search and Inductive bias.

In this we move smoothly from one hypothesis to a nearby one. This is a form of hypothesis space search. Inductive Bias is to

the set of assumptions that the learner uses to predict output given inputs that it has not encountered.

④ Hidden layer Representation :-

Based on the given input values in some way that ~~computers~~ capture their relevant features, so that this hidden layer representation could be used by the output unit to compute the correct target values.

## ⑤ Generalization, overfitting and stopping criterion

The termination condition for the algorithm has been left unspecified. What is an appropriate condition for terminating the weight update loop.

① Continue training until the error  $E$  on the training examples falls below some predetermined threshold.

② Backpropagation is susceptible to overfitting the training examples at the cost of decreasing generalization accuracy over other unseen examples.

An Illustrative Example:- Face Recognition:

→ face Recognition is a category of biometric software that maps individual facial features and stores the data as a face print.

→ Machine learning and Image processing  
are backbones of Face Recognition.

### Image Processing

~~AI~~ Involves the process of computer vision.

Image Reading:- Reads any image between the range of 0-255.

for any image, there are 3 colors -

RGB.

- OpenCV:- python library to solve computer vision.

### Machine Learning

takes dataset as input and learn from that data.

- identifies the pattern (height, width, color etc).

- ML does 3 major functions in face Recognition.

1. Deriving the feature vector

2. Matching Algorithm,

3. Face Recognition Operations.

i) Face-detection

ii) Face Analysis

iii) Image to data conversion.

iv) Match finding.

## Face Recognition Software

Deepvision AI.

SenseTime

Amazon Rekognition.

facefirst

trueface

cognitec

## Applications

Hospital Security.

Airline Industry.

## Advanced Topics in Neural Networks:

3 advanced topics.

### 1. Alternative error functions:

Adding an extra penalty to  $f(w)$

can reduce the problem of overfitting.

- Each weight is multiplied with a constant in each iterations.

∴ we can redefine  $f$  as

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \mathcal{C}_d} (t_{kd} - o_{kd})^2 + \gamma \sum_i w_{ji}^2$$

2. Alternative error minimisation  
procedures

1. line search : optimized way to find the distance.
  2. conjugate gradient descent methods
- ↳
- Same as line search  
 but diff is error minimisation

3. Recurrent neural network (RNN) :-

ANN that can be applied to time series data and uses the output of n/w units at time  $t$  as input to other units at time  $t+1$

## Evaluating the Hypothesis:-

evaluating the accuracy of hypothesis is basic of ML.

for evaluating the accuracy of hypothesis, we mainly focus on 3 questions.

1. Hypothesis is accurate over limited sample of data, then what about additional data.

2. we have 2 hypothesis - one is better than other when used on sample of data.

Then will better hypothesis work good for general data also?

3. When we have have limited data, what is the best way to use this data for both learning and estimating hypothesis.

## Motivation

Sometimes, very precise hypothesis is required. In that case estimating the accuracy is very important.

Example: Medical Treatment.

Estimating hypothesis accuracy:-

Let us make some assumptions,

- There is some space of possible instances  $x$  over which many target functions

may be defined.

- different instances will have different frequencies.

- Based on frequency, each instance in  $x$  will have some unknown probability.

- Trainer will teach the machine about the training examples of target function.

Send a particular instance  $x_i$  along with target value  $f(x)$ .

Target function:  $f: x \rightarrow \{0, 1\}$

Classifies each instance into diff categories based on req. Now after classification, what is the probable error in this estimation we made.

Error : 2 types

- ① Sample Error.
- ② True Error.

Sample Error :-

The error rate of hypothesis over a sample of data.

$$\text{Sample error} = \text{error}_s(h) = \frac{1}{n} \sum_{\text{ses}} \delta(f(x), h(x))$$

$$\begin{aligned}\delta(f(x), h(x)) &= 1 \text{ if } f(x) \neq h(x) \\ &= 0 \text{ otherwise.}\end{aligned}$$

True error :-

The error rate of hypothesis over entire distribution P.

$$\text{True error} = \text{error}_D(h) = D[f(x) \neq h(x)]$$