

SOFT COMPUTING
Subject Code : 18MIT24C
Prepared by Dr. N.Thenmozhi

UNIT-II: Back propagation Networks: Architecture of a Back-Propagation Network – Back propagation Learning- Effect of Tuning parameters of the Back Propagation Neural Network – Selection of various parameters in BPN.

TEXT BOOK

1. S.Rajasekaran & G.A.Vijayalakshmi Pai, “Neural Networks, Fuzzy logic, and Genetic Algorithms Synthesis and Applications, PHI, 2005.

REFERENCE BOOKS

1. James A. Freeman, David M.Skapura, “Neural Networks-Algorithms, Applications, and Programming Techniques”, Pearson Education.
2. Fredric M. Ham, Ivica Kostanic, “Principles of Neuro computing for science of Engineering”, TMCH.

2. Back propagation Networks

What is Backpropagation?

Back-propagation is the essence of neural net training. It is the method of fine-tuning the weights of a neural net based on the error rate obtained in the previous epoch (i.e., iteration). Proper tuning of the weights allows you to reduce error rates and to make the model reliable by increasing its generalization.

Backpropagation is a short form for "backward propagation of errors." It is a standard method of training artificial neural networks. This method helps to calculate the gradient of a loss function with respects to all the weights in the network.

Why We Need Backpropagation?

Most prominent advantages of Backpropagation are:

- Backpropagation is fast, simple and easy to program
- It has no parameters to tune apart from the numbers of input
- It is a flexible method as it does not require prior knowledge about the network
- It is a standard method that generally works well
- It does not need any special mention of the features of the function to be learned.

History of Backpropagation

- In 1961, the basics concept of continuous backpropagation were derived in the context of control theory by J. Kelly, Henry Arthur, and E. Bryson.
- In 1969, Bryson and Ho gave a multi-stage dynamic system optimization method.
- In 1974, Werbos stated the possibility of applying this principle in an artificial neural network.
- In 1982, Hopfield brought his idea of a neural network.
- In 1986, by the effort of David E. Rumelhart, Geoffrey E. Hinton, Ronald J. Williams, backpropagation gained recognition.
- In 1993, Wan was the first person to win an international pattern recognition contest with the help of the backpropagation method.

Disadvantages of using Backpropagation

- The actual performance of backpropagation on a specific problem is dependent on the input data.
- Backpropagation can be quite sensitive to noisy data
- You need to use the matrix-based approach for backpropagation instead of mini-Batch

2.1 Architecture of a Backpropagation Network

The Perceptron Model

The **perceptron** is a classification algorithm. Specifically, it works as a linear binary classifier. It was invented in the late 1950s by Frank Rosenblatt. The **perceptron** basically works as a threshold function — non-negative outputs are put into one class while negative ones are put into the other class.

Architecture of a Backpropagation network

The perceptron model

Rosenblatt's perceptron was introduced for linearly inseparable (non-linearly separable) problem.

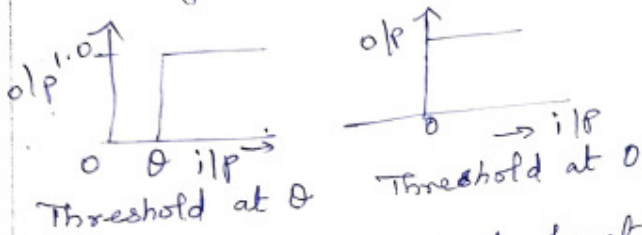
The initial approach to solve this, more than one perceptron, each setup identifying small linearly separable section of i/p's, and then combining their outputs into another perceptron. For problem was solved.

Each neuron in the structure takes the weighted sum of inputs, thresholds it and outputs either a one or zero. For the perceptron in the first layer, the i/p comes from the actual i/p of the problem. For 2nd layer the i/p are o/p of the first layer. The perceptrons of the 2nd layer do not know which of the real i/p's from the first layer were on/off.

It is impossible to strengthen the connections between active i/p's & strengthen the correct parts of the network. No actual i/p's are effectively

masked off from the o/p units of the intermediate layer.

The 2 states of neuron being on or off, do not give us any indication of the scale by which we have to adjust the weights.



The hard-hitting threshold functions remove the information that is needed if the network is to successfully learn. Hence the network is unable to determine which of the i/p weights should be increased and which one should not and so, it is unable to work to produce a better solution next time.

Using the step function as the thresholding process is to adjust it slightly and to use a slightly different nonlinearity.

2. The solution

If we smoothen the threshold function, it more or less turns on or off as before but has a sloping region in the middle that will give us some information on the inputs.

We will be able to determine when it need to strengthen or weaken the relevant weight

Now, the new will be able to learn as required. A couple of possibilities for the new thresholding function are given as



Combining perceptrons to solve XOR problem.

1. If the value of the output = 1, then if the i/p exceeds the value of the threshold a lot. Else it is = 0 if the i/p is far less than the threshold.

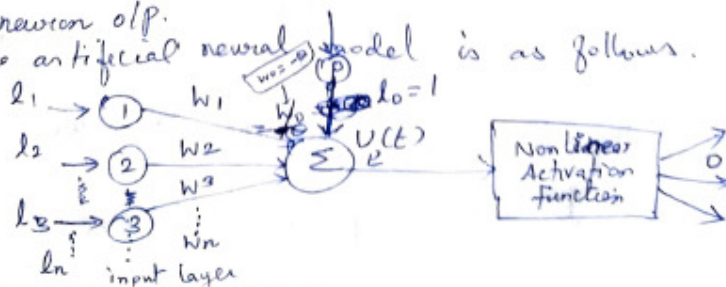
2. When i/p & threshold are almost same, the o/p of the neuron will have a value between 0 & 1. (i.e. o/p to the neuron can be related to i/p).

⇒ Based on the biological neuron, an artificial neuron receives much i/p representing the o/p of the other neurons.

⇒ Each i/p is multiplied by the corresponding weight analogous to synaptic strengths.

⇒ All of these weighted i/p's are then summed up & passed through an activation function to determine the neuron o/p.

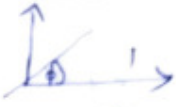


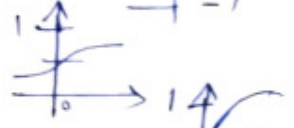


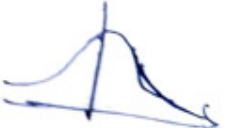
∴ The artificial neural model is as follows.



$$u(t) = w_1 I_1 + \dots + w_n I_n$$

$$u = \{W\} \{I\}$$

Typical nonlinear activation operators

Type	Equation	Function Form
Linear	$O = gI$ $g = \tan \phi$	
Piecewise Linear	$O = \begin{cases} 1 & \text{if } m > 1 \\ gI & \text{if } m < 1 \\ -1 & \text{if } m < -1 \end{cases}$	
Hard Limiter	$O = \text{sign}[I]$	
Unipolar Sigmoidal	$O = \frac{1}{(1 + \exp(-\lambda I))}$	
Bipolar Sigmoidal	$O = \tanh[\lambda I]$	
Unipolar multipole	$O = \frac{1}{2} \left[1 + \frac{1}{m} \sum_{m=1}^m \tanh(g^m (I - W_0^m)) \right]$	
Radial Basis Function (RBF)	$O = \exp(-I)$ $I = \left[\frac{-\sum_{i=1}^N (W_i(t) - X_i(t))^2}{2\sigma^2} \right]$	

Considering threshold θ , the relative input to the neuron is given by.

$$u(t) = w_1 I_1 + w_2 I_2 + \dots + w_n I_n - \theta$$

$$= \sum_{i=0}^n w_i I_i \quad \text{where } w_0 = -\theta, I_0 = 1$$

the o/s using the nonlinear transfer function 'f' is given by

$$O = f(u)$$

Even if the sigmoidal fn is differentiated, it gives continuous val of the output.

The activation function $f(x)$ is chosen as a nonlinear function to emulate the nonlinear behaviour of conduction current mechanism in a biological neuron. Artificial neuron is not intended to be a xerox copy of the biological neurons, many forms of nonlinear functions are used in various engineering applications. From the above neurons with sigmoidal fns ~~are preferred~~ resemble to biologic calculation. Hard limiter & radial basis fns are also equally popular.

3. Single Layer Artificial Neural Network

Single Layer ^{feed}forward neural network consisting of an i/p layer to receive the i/p & an o/p layer to o/p the vector respectively.

i/p layer \rightarrow n neurons & o/p layer \rightarrow m neurons
 Indicate the weight of the synapse connecting ith i/p neuron to the jth o/p neurons as w_{ij} . The i/p's of the i/p layer & the corresponding o/p's of the o/p layer are given as

$$I_i = \begin{Bmatrix} I_{i1} \\ I_{i2} \\ \vdots \\ I_{in} \end{Bmatrix}_{n \times 1} \quad O_o = \begin{Bmatrix} O_{o1} \\ O_{o2} \\ \vdots \\ O_{om} \end{Bmatrix}_{m \times 1}$$

Linear transfer function for the neurons in the i/p layer and the unipolar sigmoidal function for the neurons in the o/p layer.

$$\{O_o\} = \{I_i\} \in \text{linear transfer function.}$$

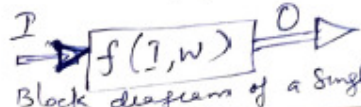
$$I_{oj} = w_{1j} I_{i1} + w_{2j} I_{i2} + \dots + w_{nj} I_{in}$$

Hence, the i/p to the o/p layer can be given as

$$\{I_o\} = [W]^T \{O_o\} = [W]^T \{I_i\}$$

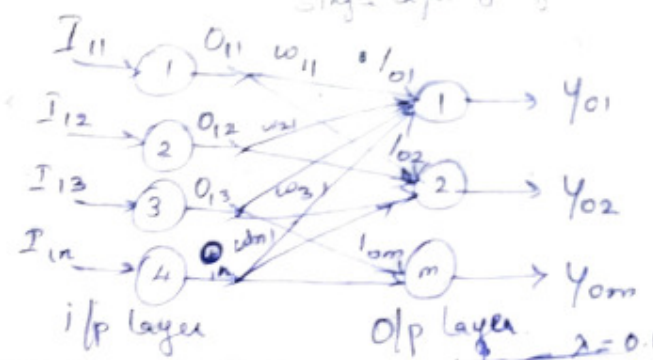
$m \times 1 \quad m \times n \quad n \times 1$

\rightarrow weight or connection matrix.



Block diagram of a single layer feedforward neural network.

single layer feedforward network



Unipolar Sigmoidal \rightarrow $\lambda = 0.125$
 & Slope of this fn is $\lambda = 1$ for various values of λ .
 squashed-S slope

Neurons in the o/p layer, the o/p is given by

$$O_{ok} = \frac{1}{(1 + e^{-\lambda I_{ok}})} \rightarrow \text{non linear activation fn.}$$

$$\{O_o\} = f\{W I\}$$

λ is known as sigmoidal gain.
 Each activation value is in turn a scalar product of the i/p with respect to weight vectors.

\therefore Sigmoidal fn is

$$f(I) = \frac{1}{(1 + e^{-\lambda I})}$$

& $f'(I) = \lambda f(I)(1 - f(I))$

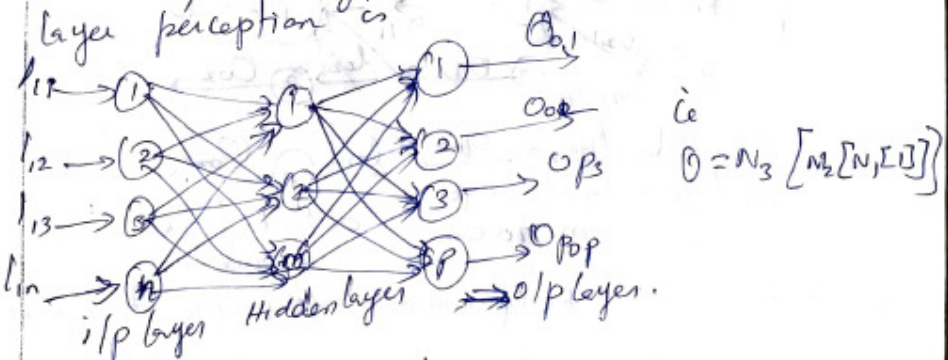
4. Model for multilayer perceptron

Adapted perceptrons are arranged in layers and are called multilayer perceptron. This model has 3 layers: an i/p layer & o/p layer & a layer in between not connected directly to the i/p or the o/p is called hidden layer.

For the perceptrons in the i/p layer, use linear transfer function & for the hidden & o/p layer we use sigmoidal or squashed-S functions.

The i/p layer serves to distribute the values they receive to the next layer & so, ~~it~~ does not perform a weighted sum or threshold. Because we have modified the single layer perceptron by changing the nonlinearity from a step function to sigmoid function and added a hidden layer, so we are forced to change the ^{learning} rules as well.

So we now should be able to learn to recognize more complex things. The i/p - o/p mapping of multi layer perceptron is



N_1, N_2, N_3 are nonlinear mapping provided by i/p, hidden & o/p layer respectively

multilayer perceptron provides
no increase in computational power
over single layer.

i/p layer represents raw information that is
fed into the net. Hidden layer is deter-
mined by the activities of the neurons
in the i/p layer & connecting weights between
i/p & hidden units.

the activity of o/p units depends on the
activity in the hidden layer & the weight
between the hidden & o/p layer.

Neurons in the hidden layer are free to
construct their own representations of the
o/p.



Backpropagation Learning

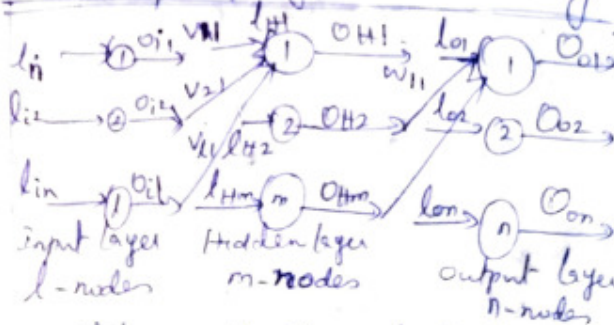


fig. multilayer feedforward backpropagation network

consider the problem \Rightarrow an n set of 'l' i/p
corresponding n set of 'n' o/p data

- n set of i/p and output data

NO	Input			output		
1	I_1	I_2	$\dots I_j$	O_1	O_2	$\dots O_n$
2	0.3	0.4	$\dots 0.8$	0.1	0.56	0.82
\vdots						
n set						

Input Layer Computation

consider linear activation function, the o/p of the i/p layer is i/p of i/p layer. ($g = \tan \phi = 1$)

Taking one set of data

$$\{O\}_l = \{I\}_l$$

$$l \times 1 = l \times 1$$

The hidden layer is connected by synapses to i/p neurons. & V_{ij} = weight of the arc between i th i/p neuron to j th hidden neuron

The i/p of hidden neuron = weighted sum of the o/p of the i/p neurons to get I_{Hp}
i.e. p th hidden neuron as

$$I_{Hp} = V_{p1}O_{i1} + V_{p2}O_{i2} + \dots + V_{pl}O_{il}$$

defined in weighted matrix $\{V\}_{l \times m}$ ($p = 1, 2, 3, \dots, m$)

We can get an i/p to the hidden neuron as

$$\{I\}_H = \{V\}^T \{O\}_l$$

$$m \times 1 \quad m \times l \quad l \times 1$$

Hidden Layer Computation

$$\text{sigmoidal fr. } \theta_{Hp} = \frac{1}{1 + e^{-\lambda(I_{Hp} - \theta_{Hp})}}$$

i.e. o/p of p th hidden neuron

O_{Hp} - o/p of the p th hidden neuron
 I_{Hp} - i/p
 θ_{Hp} - threshold of the p th neuron.

A non zero threshold neuron = to an i/p always at -1
 becomes the connecting weight values.

The derivation \rightarrow o/p to the hidden neuron is

$$\{O\}_H = \left\{ \frac{1}{1 + e^{-\lambda(I_{Hp} - \theta_{Hp})}} \right\}$$

The weighted sum of the o/p of the hidden neurons
 to get I_{Oj} (Input of the j th output neuron)

$$\{I\}_O = \{W\}^T \{O\}_H$$

$n \times 1$ $n \times m$ $m \times 1$

o/p Layer Computation.

(1) \rightarrow

2.3 Training Algorithm

For training, BPN will use binary sigmoid activation function. The training of BPN will have the following three phases.

- **Phase 1** – Feed Forward Phase
- **Phase 2** – Back Propagation of error
- **Phase 3** – Updating of weights

All these steps will be concluded in the algorithm as follows

Step 1 – Initialize the following to start the training –

- Weights
- Learning rate α

For easy calculation and simplicity, take some small random values.

Step 2 – Continue step 3-11 when the stopping condition is not true.

Step 3 – Continue step 4-10 for every training pair.

Phase 1

Step 4 – Each input unit receives input signal x_i and sends it to the hidden unit for all $i = 1$ to n

Step 5 – Calculate the net input at the hidden unit using the following relation –

$$Q_{in_j} = b_{0j} + \sum_{i=1}^n x_i v_{ij}$$

Here b_{0j} is the bias on hidden unit, v_{ij} is the weight on j unit of the hidden layer coming from i unit of the input layer.

Now calculate the net output by applying the following activation function

$$Q_j = f(Q_{in_j})$$

Send these output signals of the hidden layer units to the output layer units.

Step 6 – Calculate the net input at the output layer unit using the following relation –

$$y_{in_k} = b_{0k} + \sum_{j=1}^p Q_j w_{jk}$$

Here b_{0k} is the bias on output unit, w_{jk} is the weight on k unit of the output layer coming from j unit of the hidden layer.

Calculate the net output by applying the following activation function

$$y_k = f(y_{in_k})$$

Phase 2

Step 7 – Compute the error correcting term, in correspondence with the target pattern received at each output unit, as follows –

$$\delta_k = (t_k - y_k) f'(y_{in_k})$$

On this basis, update the weight and bias as follows –

$$\Delta v_{jk} = \alpha \delta_k Q_{ij}$$

$$\Delta b_{0k} = \alpha \delta_k$$

Then, send δ_k back to the hidden layer.

Step 8 – Now each hidden unit will be the sum of its delta inputs from the output units.

$$\delta_{in_j} = \sum_{k=1}^m \delta_k w_{jk}$$

Error term can be calculated as follows –

$$\delta_j = \delta_{in_j} f'(Q_{in_j})$$

On this basis, update the weight and bias as follows –

$$\Delta w_{ij} = \alpha \delta_j x_i$$

$$\Delta b_{0j} = \alpha \delta_j$$

Phase 3

Step 9 – Each output unit ($y_k, k = 1 \text{ to } m$) updates the weight and bias as follows –

$$v_{jk}(\text{new}) = v_{jk}(\text{old}) + \Delta v_{jk} \quad v_{jk}(\text{new}) = v_{jk}(\text{old}) + \Delta v_{jk}$$

$$b_{0k}(\text{new}) = b_{0k}(\text{old}) + \Delta b_{0k} \quad b_{0k}(\text{new}) = b_{0k}(\text{old}) + \Delta b_{0k}$$

Step 10 – Each output unit ($z_j, j = 1 \text{ to } p$) updates the weight and bias as follows –

$$w_{ij}(\text{new}) = w_{ij}(\text{old}) + \Delta w_{ij} \quad w_{ij}(\text{new}) = w_{ij}(\text{old}) + \Delta w_{ij}$$

$$b_{0j}(\text{new}) = b_{0j}(\text{old}) + \Delta b_{0j} \quad b_{0j}(\text{new}) = b_{0j}(\text{old}) + \Delta b_{0j}$$

Step 11 – Check for the stopping condition, which may be either the number of epochs reached or the target output matches the actual output.

Generalized Delta Learning Rule

Delta rule works only for the output layer. On the other hand, generalized delta rule, also called as **back-propagation** rule, is a way of creating the desired values of the hidden layer.

Mathematical Formulation

For the activation function $y_k = f(y_{in,k})$ the derivation of net input on Hidden layer as well as on output layer can be given by

$$y_{in,k} = \sum_i z_i w_{ij} \quad y_{in,k} = \sum_i z_i w_{ij}$$

$$\text{And } y_{in,j} = \sum_i x_i v_{ij} \quad y_{in,j} = \sum_i x_i v_{ij}$$

Now the error which has to be minimized is

$$E = \frac{1}{2} \sum_k [t_k - y_k]^2 \quad E = \frac{1}{2} \sum_k [t_k - y_k]^2$$

By using the chain rule, we have

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} \left(\frac{1}{2} \sum_k [t_k - y_k]^2 \right) \quad \frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} \left(\frac{1}{2} \sum_k [t_k - y_k]^2 \right)$$

$$= \frac{\partial}{\partial w_{jk}} \left(\frac{1}{2} [t_k - f(y_{in,k})]^2 \right) = \frac{\partial}{\partial w_{jk}} \left(\frac{1}{2} [t_k - f(y_{in,k})]^2 \right)$$

$$= -[t_k - y_k] \frac{\partial}{\partial w_{jk}} f(y_{in,k}) = -[t_k - y_k] \frac{\partial}{\partial w_{jk}} f(y_{in,k})$$

$$= -[t_k - y_k] f'(y_{in,k}) \frac{\partial}{\partial w_{jk}} (y_{in,k}) = -[t_k - y_k] f'(y_{in,k}) \frac{\partial}{\partial w_{jk}} (y_{in,k})$$

$$= -[t_k - y_k] f'(y_{in,k}) z_j = -[t_k - y_k] f'(y_{in,k}) z_j$$

Now let us say $\delta_k = -[t_k - y_k] f'(y_{in,k})$ $\delta_k = -[t_k - y_k] f'(y_{in,k})$

The weights on connections to the hidden unit z_j can be given by –

$$\frac{\partial E}{\partial v_{ij}} = -\sum_k \delta_k \frac{\partial}{\partial v_{ij}} (y_{in,k}) \quad \frac{\partial E}{\partial v_{ij}} = -\sum_k \delta_k \frac{\partial}{\partial v_{ij}} (y_{in,k})$$

Putting the value of $y_{in,k}$ we will get the following

$$\delta_j = -\sum_k \delta_k w_{jk} f'(z_{in,j}) \quad \delta_j = -\sum_k \delta_k w_{jk} f'(z_{in,j})$$

Weight updating can be done as follows –

For the output unit

$$\begin{aligned}\Delta w_{jk} &= -\alpha \partial E \partial w_{jk} \\ \Delta w_{jk} &= -\alpha \partial E \partial w_{jk} \\ &= \alpha \delta_k z_j = \alpha \delta_k z_j\end{aligned}$$

For the hidden unit

$$\begin{aligned}\Delta v_{ij} &= -\alpha \partial E \partial v_{ij} \\ \Delta v_{ij} &= -\alpha \partial E \partial v_{ij} \\ &= \alpha \delta_j x_i\end{aligned}$$

2.4 Effect of Tuning parameters of the Back Propagation Neural Network

- Sigmoidal gain
- Threshold value

2.5 Selection of various parameters in BPN

- Number of hidden nodes
- Momentum coefficient
- Sigmoidal Gain
- Local Minima
- Learning Coefficient