

19/03/2020

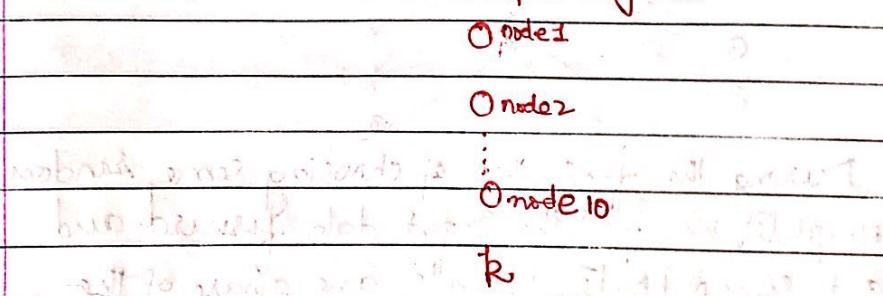
## The Backpropagation Algorithm.

The algorithm starts by constructing a network with the required input units (dependent on the task at hand), desired number of hidden layers, desired number of nodes in each hidden layer, and the output units.

The feedforward step is performed first (as done in the last class). For each input sample, the output is noted. Since the output may have more than one node, the error associated with the input sample  $i$  can be written as

$$E_i(\mathbf{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

You may visualize this as



Say, we were solving the digit recognition problem, say we input a **1** and the output is

$$\begin{bmatrix} 0.98 \\ 0.02 \\ 0.03 \\ 0.15 \\ 0.05 \\ 0.01 \\ 0.70 \\ 0.20 \\ 0.30 \\ 0.05 \end{bmatrix}$$

$O_k$

and the desired output is

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_{10} \end{bmatrix} = t_K$$

So, the error of this sample is

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k=1}^K (t_k - o_k)^2$$

Reoutput

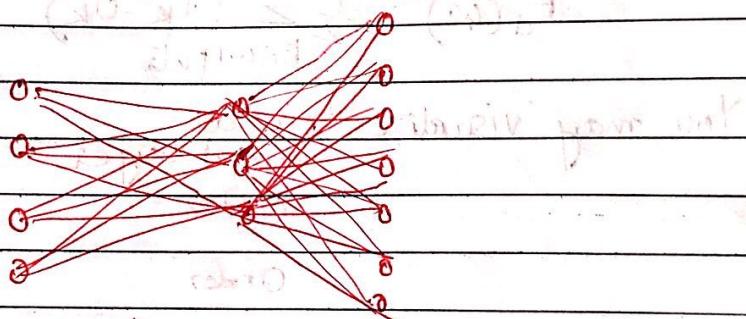
If we were to find the total error due to ALL the input samples, it would be

$$E(\vec{w}) = \frac{1}{2} \sum_{d=1}^D \sum_{k=1}^K (t_{kd} - o_{kd})^2$$

where the inner summation is for each sample and outer summation is over ALL D samples.

Intuition:

Say a network has an input layer, one hidden layer and an output layer.



During the first step of choosing some random weights, we feed the input data forward and get some outputs. We also are aware of the target values that we desire.

So in essence we are aware of the error at the output layer, when a particular sample is input. We desire to somehow update the weight matrix that is between

Input layer and the hidden layer, as well as the weight matrix between the hidden layer and the output layer.

In the lab question given to you, I had given both these matrices to you.

We shall again apply the stochastic gradient rule.  
(We did for training a linear unit).

For each training example  $d$ , every weight,  $w_{ji}$  is updated by the following  $\Delta w_{ji}$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} \quad (\text{Negative to the steepest ascent})$$

where  $E_d$  is the error on the training sample  $d$ ,

$$E_d(\vec{w}) = \frac{1}{2} \sum (t_k - o_k)^2$$

Outputs

Some notations that you need to know are in your book. Diagrammatically they are as follows —

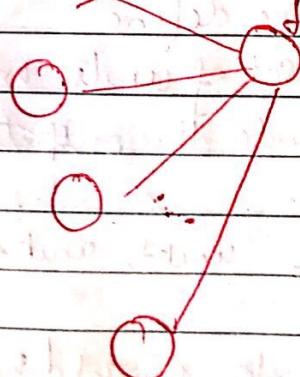
node  $j$



This node  $j$  may be anywhere, in one of the hidden layers, or in the output layer.

$i$  (say)

$x_{ji}$

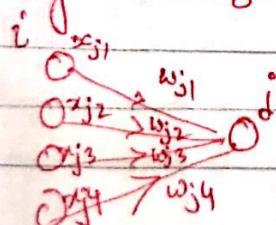


$x_{ji}$  =  $i^{\text{th}}$  input to node  $j$

$w_{ji}$  = weight from  $i^{\text{th}}$  input to node  $j$

$$\text{net}_j = \sum_i w_{ji} x_{ji}$$

You may visualize  $\text{net}_j$  as follows —



So  $\text{net}_j$  is the weighted sum of inputs to node  $j$  from the previous layer

$o_j$  = output computed by unit  $j$

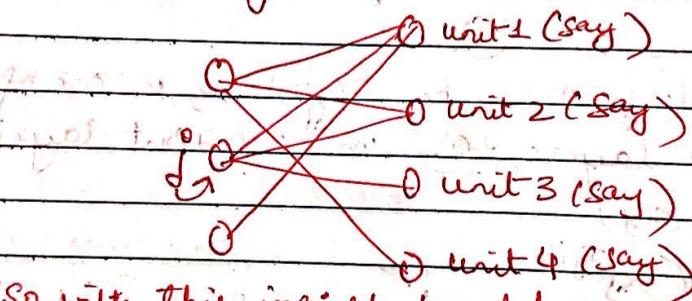
which is in effect

$$\sigma(\text{net}_j) = \sigma\left(\sum_i w_{ji} x_{ji}\right)$$

$t_j$  = target output for unit  $j$

$\tau$  = sigmoid function (any squashing function)

Now, though not mentioned earlier, you should now understand that it is not necessary that <sup>all</sup> nodes of one layer are connected to all nodes of the next layer. Depending on the need and architecture, one layer may be partially connected to the next layer.



So, with this insight, we define

$\text{Downstream}(j) =$  the set of units whose immediate inputs include the outputs of unit  $j$

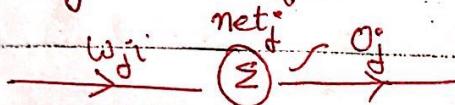
$$\text{So } \text{Downstream}(j) = \{\text{unit 1, unit 2, unit 3}\}$$

In a fully connected network,  $j$  would have been connected to EACH node in the next layer.

To implement the stochastic gradient descent rule i.e  $\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$

we need the value of  $\frac{\partial E_d}{\partial w_{ji}}$

In a very basic form



Using chain Rule

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \cdot x_{ji} \quad \leftarrow \text{To find } -\textcircled{1}$$

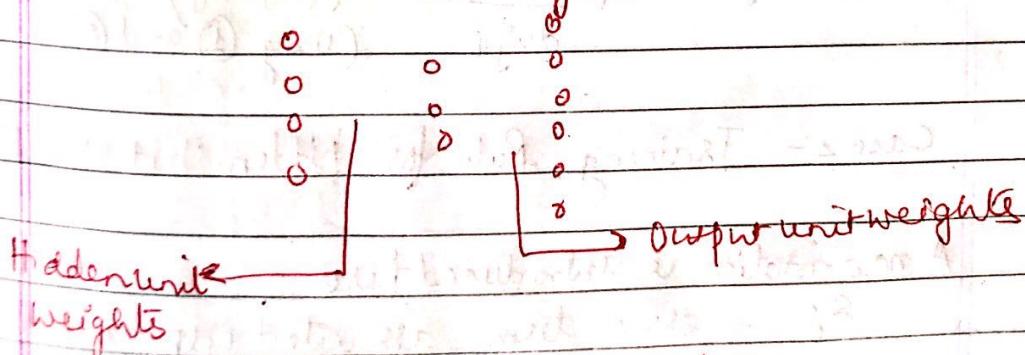
This value is essential for updating all the weights.

For a three layer architecture

There are 2 sets of weights

- Output unit weights

- Hidden unit weights



In a higher number of hidden layers, we will have more hidden unit weight matrices.

The training rule for Output Unit weights & the training rule of Hidden Unit weights is different. So we shall handle it with 2 cases.

### Case 1 - Training Rule for Output Unit Weights

$$\frac{\partial E_d}{\partial w_{netj}} \quad \leftarrow \text{To find}$$

$$\frac{\partial E_d}{\partial w_{netj}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial w_{netj}} \quad \textcircled{2} \text{ (Chain Rule Again)}$$

$$= \frac{\partial}{\partial o_j} \left[ \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2 \right] * \frac{\partial}{\partial w_{netj}} [\sigma(w_{netj})]$$

~~this derivative will be non zero only when~~

$$t_j = o_j \quad \Rightarrow \quad = -(t_j - o_j) * \sigma'(1 - \sigma)$$

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

Substituting in the Stochastic Gradient Descent rule equation \*

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} = \eta(t_j - o_j) \sigma_j(1 - \sigma_j) \cdot x_{ji} \quad (\text{using } \textcircled{1} \text{ and } \textcircled{2})$$

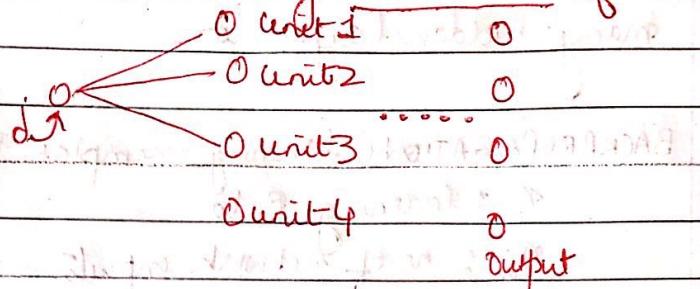
### Case 2 - Training Rule for hidden Unit Weights

A new notation is introduced here

$$\delta_i = \text{error term associated with unit } i \\ = -\frac{\partial E}{\partial w_{neti}}$$

$$\frac{\partial E_d}{\partial \text{net}_j} \leftarrow \text{To find}$$

Since we are in a hidden layer, each unit influences the set of units immediately downstream of unit  $j$ .



Since  $\text{net}_j$  effects  $E_d$  through downstream( $j$ )

Therefore

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{k \in \text{downstream}(j)} \frac{\partial E_d}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_j} \quad (\text{chain rule})$$

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_k -\delta_k \frac{\partial \text{net}_k}{\partial \text{net}_j} \quad (\text{Chain Rule + gain})$$

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_k -\delta_k w_{kj} \frac{\partial o_j}{\partial \text{net}_j} \quad \begin{cases} \delta_k = o_k(1-o_k) \\ \frac{\partial o_j}{\partial \text{net}_j} = \dots \\ \dots \\ = o_1 w_{1j} + o_2 w_{2j} \\ \dots \\ + o_j w_{kj} \end{cases}$$

$$\frac{\partial E_d}{\partial \text{net}_j} = -o_j(1-o_j) \sum_k \delta_k w_{kj}$$

Also, as defined

$$\delta_j = -\frac{\partial E_d}{\partial \text{net}_j} = o_j(1-o_j) \sum_{k \in \text{downstream}(j)} \delta_k w_{kj}$$

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

Having obtained the gradient descent weight update rule for output unit weights and for hidden unit weights, the complete algorithm can be written as follows (Pg 98, TM). This algo is for 2 layers of sigmoid units and can be easily extended to many hidden layers too.

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

$\eta$  : learning Rate

$n_{in}$  : no. of network inputs

$n_{hidden}$  : no. of units in hidden layer

$x_{ji}$  : input from unit  $i$  to unit  $j$

$w_{ji}$  : weight from unit  $i$  to unit  $j$

A. Create a FF network with  $n_{in}$  inputs,  $n_{hidden}$  hidden units and  $n_{out}$  output units

B. Initialize all network weights to small random numbers

C. Until termination do

- For each  $(\vec{x}, \vec{t})$  training example

Propagate the input forward through the network

1. Input  $\vec{x}$  and compute  $O_h$  of every unit  $h$

Propagate the errors backward

2. For each output unit  $k$ , calculate error term

$$\delta_k \leftarrow O_k(1-O_k)(t_k - O_k)$$

3. For each hidden unit  $h$ , calculate error term

$$\delta_h \leftarrow O_h(1-O_h) \sum_{k \in \text{Outputs}} w_{kh} \delta_k$$

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

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

where

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

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.

## Assignment.

Understand this much. In the last lab exercise involving only FF, the 2 matrices were provided to you. On implementing, you would have realized that the weights provided to you were optimum and that the actual output was same as target output for all the training samples.

I want you to do the following now,-

Take the input & the target output as given in the question. Randomly initialize both the weight matrices (say between -0.05 and +0.05)

Feed the inputs, observe the outputs

Check the error

(BP)

In case of error, backpropagate the error and update the weights

Repeat the FF and BP till you get a satisfactory model.

Brijeshwar  
19/03/2020