# Critique of backprop by Hinton et. Al

Based on the work on Artificial Neural Networks that had made some progress, Hinton et. al. have constructed an algorithm for adjusting the errors/difference as compared to the actual outputs and thereby reflect them to the weights of the hidden units.

The network can be as below. There are a number of neurons in each layer. While the neurons between two consecutive layers are connected, the neurons run parallel (i.e., not connected) in the same layer. The last layer is supposed to produce the desired output. But, how can that occur? This is easily possible by taking the squared error of the actual and predicted output.
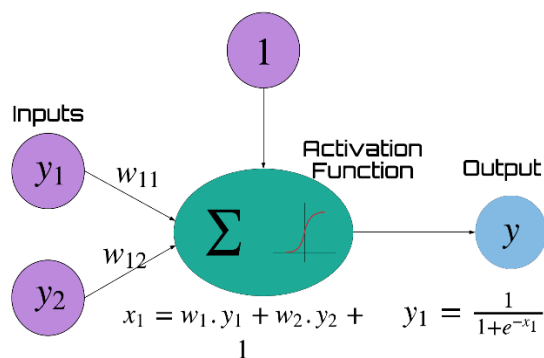
Before we dive into the back-propagation algorithm, let's understand the forward propagation in terms of equations. Let's assume the input to a layer is given $x_j$ is the total sum of the outputs of the previous layer $y_i$ multiplied by their respective weights $w_{ji}$. They can be expressed as below:

$$x_j = \sum_i y_i . w_{ji}$$

An additional bias of **1** can be added to this. Lastly, the neuron gets activated by a non-linear function as below:
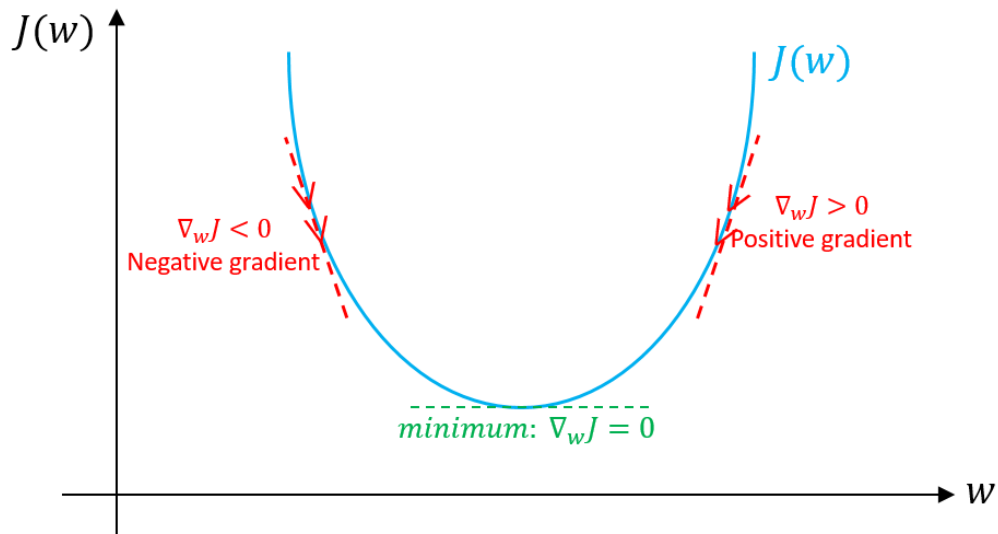
$$y_j = \frac{1}{1 + e^{x_j}}$$

The pictorial representation can be as below:



This operation recursively occurs to all the neurons till the output unit. The output unit is compared with the actual output using squared error as below:

$$E = \frac{1}{2} \sum_c \sum_j (y_{j,c} - d_{j,c})^2$$

Now let's move to the next phase, i.e., back pass. We can accomplish the reduction of error by reflecting a reduction in the error to all the weights of the network. Hence, we perform gradient descent where you find you effectively find the slope of the $Error\ Function$.

To get an intuition od the gradient descent, let's consider the picture above (Please ignore the notations as they are different). Suppose the slope(gradient) is negative and we subtract the gradient from $w$, we are actually adding some value to $w$ ($\because -(-) = +$)and hence, $w$ moves towards the minima, while if the gradient is positive, we still subtract from the gradient but this time, we are effectively subtracting as the gradient is positive. Hence, this methodology of reducing the error so that we reach the minima is called gradient descent. This above procedure is performed for all the variables as follows:

**$1^{st}$ the output:**

$$\frac{\partial E}{\partial y_j} = y_j - d_j$$

and since $y_j$ is dependent of $x_j$ as:

$$y_j = \frac{1}{1 + e^{x_j}}$$

We reflect the same changes to $x_j$ as:

$$\frac{\partial E}{\partial x_j} = \frac{\partial E}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_j}$$

$$= \frac{\partial E}{\partial y_j} * y_j(1 - y_j)$$

Now the last set of variables left are $w_{ji}$ which are in turn dependent on $x_j$ in the back-tracking order as:

$$x_j = \sum_i y_i \cdot w_{ji}$$

$$\Rightarrow \frac{\partial x_j}{\partial w_{ji}} = y_i$$

And

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial x_j} \cdot \frac{\partial x_j}{\partial w_{ji}}$$

$$= \frac{\partial E}{\partial x_j} \cdot y_i$$

Now moving a layer backward, we calculate $y_i$ as:

$$\frac{\partial x_j}{\partial y_i} = w_{ji}$$

$$\& \therefore \frac{\partial E}{\partial y_i} = \frac{\partial E}{\partial x_j} \cdot \frac{\partial x_j}{\partial y_i}$$

$$\Rightarrow \frac{\partial E}{\partial y_i} = \sum_j \frac{\partial E}{\partial x_j} \cdot w_{ji}$$

Now that we have calculated all the gradients, we must now reduce the errors of the respective variables by subtracting the gradients from them. Additionally, we could use a dampening factor incase the gradient values are too high. Let the factor be $\epsilon$ lying between $0$ and $1$. This procedure is apparently quiet slow and hence a momentum factor $\alpha$ can be utilized along with the weights. These two combined are shown below:

$$\Delta w(t) = \alpha . \Delta w(t-1) - \epsilon . \frac{\partial E}{\partial w(t)}$$

Where $t$ is the current sweep/iteration and $w$ can be any of the weights.

The other aspects mentioned in the paper include:

1. Random initialization of the weights. If the weights were to be initialized with the same value, then all would get updated with the same values, leading to all the weights learning the same values. Hence, random initialization ensures that the symmetry is broken.
2. It must also be noted that when only sufficient number of neurons are used, the network might get caught up in a local minima.

This paper has presented a breakthrough in terms of the learning that an artificial network has achieved although this does not exemplify the mechanism or learning in a biological neural network!