

Summary of the Back Propagation Algorithm

→ Note: Using sequential, or stochastic mode
In this mode, weight updation is
to be performed after the presentation of
each training example. An Epoch is
complete after one complete presentation
of the entire training set during the
learning process.

In this sense, sequential mode
operates by adjusting network weights
after the presentation of each training
example within the epoch.

Advantages of this mode:

- The algorithm is simple to implement.
- It provides effective solns to large & difficult problems.

Algorithm:

I] Initialization.

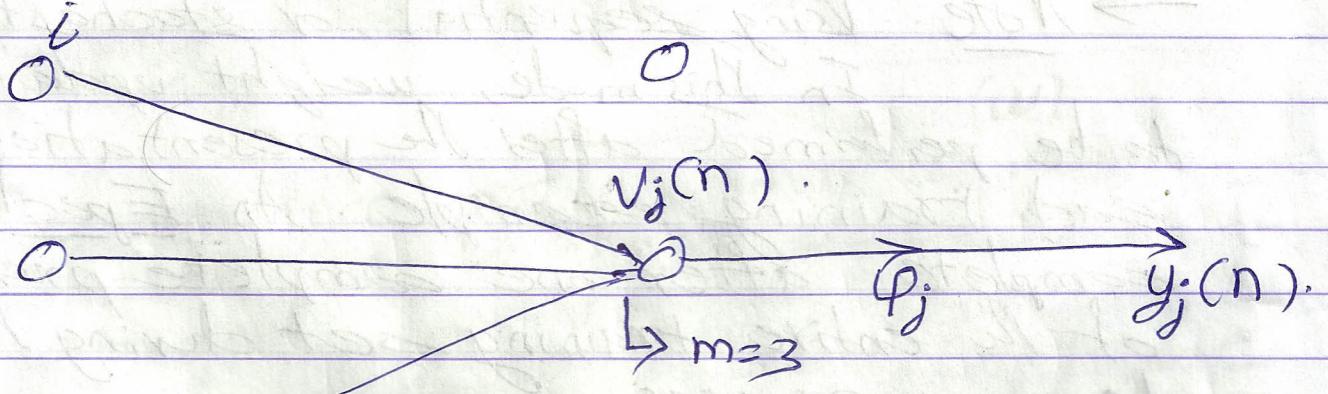
Assuming no prior information is available,
pick the synaptic weights & thresholds from
a uniform distribution whose mean is zero &
whose variance is chosen to make the
standard deviation of the induced local
fields of the neurons lie at the
transition betw the linear & saturated
parts of the sigmoid activation functn.

L (I)

For explanation, turn over.

Layer i
Previous

Layer j
current



where; $v_j(n)$: Sum of all ilps at neuron n in j^{th} layer.

m : No. of ilps to a given neuron

y_j : Activation functⁿ.

y_j : Olp of neuron j

The induced Local Field $v_j(n)$ produced at the ~~ilp~~ ilp of the activation functⁿ of neuron j is:

$$v_j(n) = \sum_{i=0}^m w_{ji}(n) y_i(n)$$

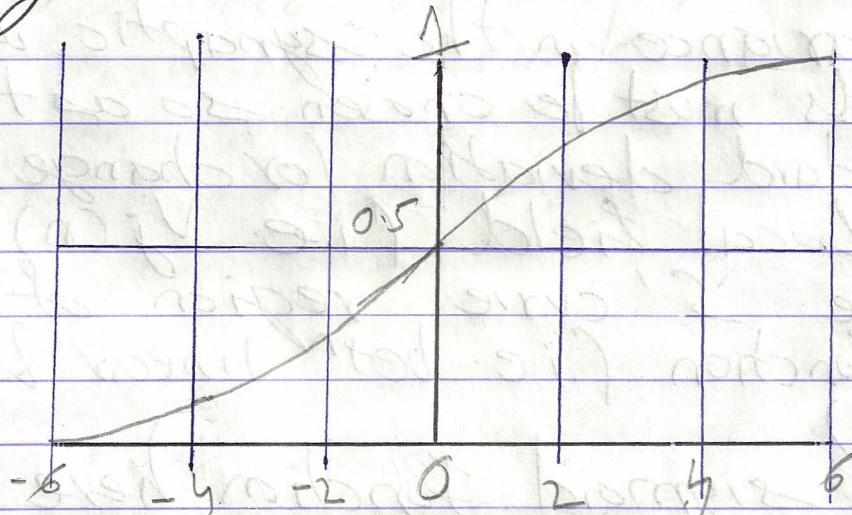
$w_{ji}(n)$: weight of connection from i^{th} to j^{th} neuron layers, to n^{th} neuron in j^{th} layer, i.e., current neuron.

& m is the total number of ilps (here, 3).

→ The synaptic weight w_{ji} equals the bias b_j , applied to neuron j. Hence, funcⁿ \rightarrow $y_j(n)$ appearing at the olp of neuron j at iteration n is $\rightarrow y_j(n)$

Now, to understand the initialization statement.

a] Sigmoid Function:



As seen above, a sigmoid function is a mathematical function having a 'S' shaped curve (sigmoid curve).

Often, the sigmoid function refers to the special case of the Logistic function (shown above) & defn by:

$$\sigma(I) = \frac{1}{1+e^{-I}}. \quad \text{--- (a)}$$

Now, referring to statement (I), we have:

→ synaptic weights: $w \rightarrow w_{ji}$

→ thresholds

→ Induced Local Field: $V_j(n)$

↳ This is produced at the i/p of the activation functⁿ, i.e. Φ_j , to produce y_j

$$\therefore y_j(n) = \Phi_j(V_j(n))$$

Considering these points, we understand the following:

(w) The variance in the synaptic weights & thresholds must be chosen so as to make the standard deviation (or change) in the induced local field (i.e. $V_j(n)$) lie within the 'S' curve region of the sigmoid function (i.e. between linear & saturated regions)

The sigmoid function here is the activation function (φ_j)

$\therefore \varphi_j$ must lie within the 'S'-shaped curve.

Conclusion: Synaptic weight & threshold values must be chosen so as to ensure V_j (the locally induced field) when applied to φ_j must lie between, or within, the 'S'-shaped curve.

II] Presentations of Training Examples:

- Present the network with an epoch of training examples. For each example set perform the sequence of forward & backward computations described below in points III & IV.
- Present the examples in an epoch in a random order each epoch!

III] Forward Computation

- Let a training example in the epoch be denoted by $(x(n), d(n))$.
where:
 $x(n) \rightarrow$ i/p vector applied to i/p nodes
 $d(n) \rightarrow$ desired o/p presented to the o/p layer.
- Compute the induced local fields & functⁿ signs by proceeding forward through the n/w layer-by-layer
- ⇒ Now, recall eqⁿ(5):

$$v_j(n) = \sum_{i=0}^m w_{ji}(n) y_i(n)$$

→ Thus the induced local field for a neuron (v_j) j in the Layer L is:

$$v_j^{(L)} = \sum_{i=0}^{m_o(L)} w_{ji}(n) y_i^{(L-1)}$$

where,

(L)

$v_j^{(n)}$ → Induced Local field of j^{th} , i.e., current neuron in layer L .

(L) $m \rightarrow$ Total no. of ips.

$w_{ji}(n)$ → Weight of connection from i^{th} to j^{th} neuron, i.e., from neuron to

$y_i^{(L-1)}$ previous layer to that in current layer

$y_i^{(n)}$ → o/p functⁿ signal of i^{th} neuron in the previous layer, i.e. $(L-1)$ layer.

Note: For $i=0$, i.e. for ip layer:

$$y_i^{(L-1)} = y_0^{(L-1)} = +1$$

$$\& w_{ji}^{(L)} = w_{j0}^{(L)} = b_j^{(L)}(n)$$

\hookrightarrow i.e., bias app'd to current neuron.

→ Assuming the use of a sigmoid functⁿ, the o/p of neuron j in layer L is:

$$y_j^{(L)} = \psi_j(v_j(n))$$

→ If neuron j is in the first hidden layer, i.e. $L=1$, set:

$$y_j^{(1)} = x_j(n)$$

→ i.e., d/p functⁿ slg of px neuron in previous (in this case input) Layer is the same as the d/p vector $x(n)$.

→ Now, if neuron is in the d/p layer, i.e. $L=L$ ($L \rightarrow$ depth of the n/w)

$$y_j^{(L)} = o_j(n)$$

→ Now, finally, compute the error slg:

$$e_j(n) = d_j(n) - o_j(n)$$

II] Backward Computation.

→ Involves computing the local gradient, i.e., δ 's of the n/w.

→ Recall eq's (14) & the final BPN eqⁿ (25):

$$\delta_j(n) = e_j(n) \varphi'_j(v_j(n))$$

→ For a neuron j in output layer L :

$$\delta_j^{(L)}(n) = e_j^{(L)}(n) \varphi'_j(v_j^{(L)}(n))$$

Δ → For a neuron j in hidden layer $\geq L$:

$$\delta_j^{(L)}(n) = \varphi'_j(v_j^{(L)}(n)) \sum_k \delta_k^{(L+1)}(n) w_{kj}^{(L+1)}(n).$$

$k \Rightarrow$ neuron k in the next $(L+1)$ layer.

⇒ Readjust synaptic weights (w) of the neurons in Layer L according to the generalised delta rule

$$w_{ji}^{(L)}(n+1) = w_{ji}^{(L)}(n) + \alpha [w_{ji}^{(L)}(n-1)] + \eta \delta_j^{(L)}(n) y_i^{(L-1)}(n)$$

where:

$\eta \rightarrow$ learning rate parameter.

$\alpha \rightarrow$ Momentum Constant.

↳ (Explanations to follow) ↳

Hence, the functional signal, $y_j(n)$ appearing at the o/p of neuron j at iteration n is:

$$y_j(n) = \varphi_j(v_j(n)).$$

\Rightarrow Activation Function:-

\rightarrow The computation of the s (Eq (25)) for each neuron of the multilayer perceptron netw requires knowledge of the derivative of the activation functn $\varphi(\cdot)$ associated with that neuron.

\rightarrow For this derivative to exist, $\varphi(\cdot)$ must be continuous.

\Rightarrow In basic terms, differentiability is the only requirement that an activation function must satisfy :-

\rightarrow Logistic Function.

Referring to eqn @ of the sigmoid function the logistic function is defined as:

$$s = \frac{1}{1+e^{-x}}.$$

For function $\varphi_j(v_j(n))$, i.e., the activation function over the locally induced field, we have

$$\varphi_j(v_j(n)) = \frac{1}{1 + \exp(-av_j(n))} \quad a > 0 \quad -\infty < v_j(n) < \infty$$

According to this non-linearity, the amplitude (value) of the output lies inside the range of $0 \leq y_j \leq 1$.

Differentiating eqⁿ ① w.r.t. $v_j(n)$, we get:

$$\varphi'_j(v_j(n)) = \frac{a \exp(-av_j(n))}{[1 + \exp(-av_j(n))]^2}$$

↳ explanation: $\frac{d}{dx} \frac{1}{u} = \frac{1}{u^2} \cdot \frac{du}{dx}$

$\delta \frac{du}{dx} = eu \cdot \frac{du}{dx}$

→ with $y_j(n) = \varphi_j(v_j(n))$, we may eliminate the exponential term $\exp(-av_j(n))$. & so express the derivative as

$$\varphi'_j(v_j(n)) = ay_j(n)[1 - y_j(n)] \quad \text{②} \quad \text{L(?)}$$

Now, remember eqⁿ ⑯, i.e.

$$\delta_j(n) = e_j(n) \varphi'_j(v_j(n))$$

Now, consider a neuron j located in the output layer, i.e., $y_j(n) = o_j(n)$

Hence, the local gradient for this neuron will be:

$$\delta_j(n) = e_j(n) \varphi'_j(v_j(n))$$

Substituting eqⁿ(2) & eqⁿ(1) up there:

$$\delta_j(n) = \alpha y_j(n) [1 - y_j(n)] [d_j(n) - y_j(n)]$$

As neuron in the output layer, $y_j(n) = o_j(n)$

$$\therefore \delta_j(n) = \alpha [d_j(n) - o_j(n)] o_j(n) [1 - o_j(n)]$$

$d_j \rightarrow$ Desired Response.

On the other hand, for an arbitrary neuron j ,

$$\delta_j(n) = \varphi'_j(v_j(n)) \sum_k \delta_k(n) w_{kj}(n)$$

i.e. eqⁿ(2), or
main BPN eqⁿ.

Now from eqⁿ C.

$$\delta_j(n) = \alpha y_j(n)[1 - y_j(n)] \sum_k \delta_k(n) w_{kj}(n)$$

... neuron $j \geqslant$ hidden

→ Note from eqⁿ C that the derivative $y'_j(v_j(n))$ attains:

↳ Maximum Value at: $y_j(n) = 0.5$

Minimum Value at: $y_j(n) = 0$ or $y_j(n) = 1.0$
(0)

→ The amount of change in a synaptic weight of the network is proportional to the derivative $y'_j(v_j(n))$.

⇒ It follows that for a sigmoid activation function the synaptic weights are changed most for the neurons in the network where the function signals are in their midrange

⇒ It has been said that it is this feature of back-propagation learning that contributes to its stability as a learning algorithm.

V.] Iteration:

→ Iterate the forward & backward computations by presenting new epochs of training examples to the n/w until the stopping criterion is met.

Notes:

- ↳ As stated earlier, randomize the presentation of training examples
- ↳ Momentum & learning-rate parameters are typically adjusted usually decreased as the no. of iterations increase.

⇒ Delta Rule:

In machine learning, the delta rule is a gradient descent learning rule for updating the weights of the inputs to artificial neurons.

$\eta \rightarrow$ Learning Parameter.

↳ The smaller it is, the smaller the changes to the synaptic weights will be betⁿ iterations. This improved learning is obtained at the cost of a slower learning rate.

If η is too large, learning speeds up, but large changes to the synaptic weights make the network unstable.

\Rightarrow A simple method of increasing the rate of learning yet avoiding the danger of instability is to modify the delta rule of eqn ⑬ & ⑯ by including a momentum term.

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta \delta_j(n) y_i(n) \quad \text{--- (d)}$$

α \rightarrow Momentum, usually +ve.

\hookrightarrow controls the feedback loop acting around $\Delta w_{ji}(n)$.

Eq ⑬ ⑯

$$\Delta w_{ji}(n) = \eta \delta(n) y_i(n)$$

may be viewed as a special case of the delta rule, with $\alpha = 0$.

(d) is called the generalised delta rule.

\Rightarrow Two passes of computation.

\rightarrow In the forward pass, synaptic weights remain unaltered.

\rightarrow Function signals (y) are computed for each neuron.

$$y_j(n) = \varphi(v_j(n))$$

$v_j(n) \rightarrow$ Induced local field

$$v_j(n) = \sum_{i=0}^m w_{ji}(n) y_i(n).$$

\rightarrow For neuron j in first hidden layer.

$$y_i(n) = x_i(n)$$

$x_i(n) \rightarrow$ i/p vectors i^{th} element.

\rightarrow For neuron in the d/p layer.

$$y_j(n) = o_j(n).$$

$o_j(n) \rightarrow j^{th}$ element of d/p vector.

\rightarrow Compare $o_j(n)$ with $d_j(n)$

$d_j(n) \rightarrow$ desired response

\rightarrow By comparing $o_j(n)$ with $d_j(n)$, obtain error $e_j(n)$

→ Thus, forward phase begins at first hidden layer & terminates at d/p layer
First hidden layer → Present i/p vector
d/p layer → Compute error for each neuron of this layer.

⇒ Backward Pass starts at d/p layer

- Pass error δ_j leftward through the neurons layer by layer.
- Recursively compute δ (local gradient) for each neuron.

$$\delta_j(n) = e_j(n) \psi'_j(v_j(n)).$$

→ Using $\delta_j(n)$, compute change of synaptic weights in accordance with delta rule.

$$\Delta w_{ji}(n) = \eta \delta_j(n) y_i(n) \quad \text{when } d=0$$

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta \delta_j(n) y_i(n).$$

P.T.O. ...

→ For output layer neurons:

$$\delta_j(n) = e_j(n) \varphi'_j(v_j(n)).$$

$$\Delta w_{ji}(n) = \eta \delta_j(n) y_i(n).$$

→ For neurons in penultimate layer:

$$\delta_j(n) = \varphi'_j(v_j(n)) \sum_k \delta_k(n) w_{kj}(n).$$

accordingly calculate $\Delta w_{ji}(n)$

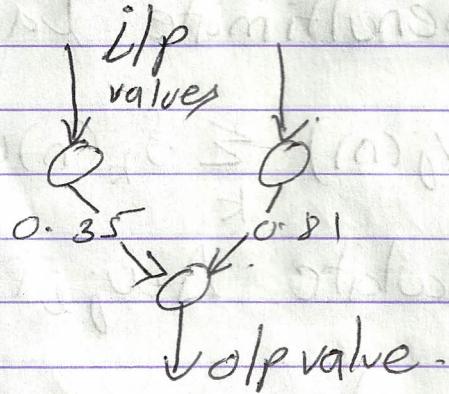
That is, we compute changes to the weights of connections feeding into current neuron.

→ Recursively compute layer by layer, propagating changes to all synaptic weights in the new.

Forward Propagation.

→ Set all weights to random values ranging from -1.0 to $+1.0$

Example: Consider a 2-layered perceptron.



$$\text{i/p's: } \begin{matrix} 0 & 1 \end{matrix} \rightarrow \text{o}^{\text{Targets}} \\ \begin{matrix} 1 & 1 \end{matrix} \rightarrow 1$$

Randomly set weights: 0.35 & 0.81 .
 $\eta \rightarrow \text{Learning Rate} \rightarrow 0.25$

a] o/p of i/p layer is same as i/p values

b] Activate neurons in 2nd layer. In this case, The o/p layer.

1 neuron: 2 i/p's.

i] For i/p {0 1}

$$\text{i/p 1} = 0.35 \times 0 = 0$$

$$\text{i/p 2} = 1 \times 0.81 = 0.81$$

$$\text{Adding} = 0 + 0.81 = 0.81$$

$$\text{target} = 0$$

$$\text{Subtracting from target} = -0.81$$

Value for changing wt. 1 = $0.25 \times 0 \times (-0.81) = 0$
Value for changing wt. 2 = $0.25 \times 1 \times (-0.81) = -0.2025$

change in wt 1 $\Rightarrow 0 + 0.35 = 0.35$ (no change).
change in wt 2 $\Rightarrow 0.81 + (-0.2025) = 0.6075$

2]

Now for ilp {1 1}

$$ilp 1 = 0.35 \times 1 = 0.35$$

$$ilp 2 = 0.6075 \times 1 = 0.6075$$

$$\text{Adding} = 0.9575.$$

$$\text{Target} = 1.$$

$$\text{Subtracting from target} = 1 - 0.9575 = 0.0425.$$

$$\text{Value for changing wt. 1} = 0.25 \times \cancel{0.35}^1 \times 0.0425 = 0.010625.$$

$$\text{Value for changing wt. 2} = 0.25 \times 1 \times 0.0425 = 0.010625.$$

$$\text{change in wt 1} \rightarrow 0.35 + 0.010625 = 0.360625$$

$$\text{change in wt 2} \rightarrow 0.6075 + 0.010625 = 0.618125.$$

#3] Compute net error:

$$(-0.81)^2 + (0.0425)^2 = 0.65790625.$$

Continue until 0 (\approx)

—