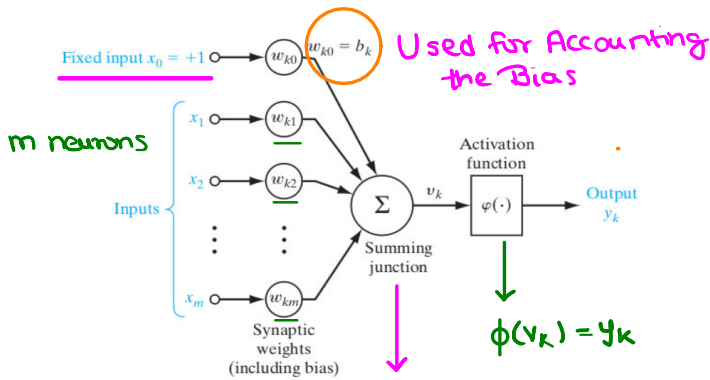


★ Neuron Model



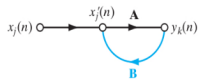
★ Adder Function

$$v_k = \sum_{j=0}^m w_{kj} x_j$$

$$\left. \begin{array}{l} w_{k0} = b_k \\ x_0 = +1 \end{array} \right\} \text{Bias accounting}$$

★ v_k we have w_{kj}

★ Feedback Loop



- Three Nodes are there $x_j(n)$, $x'_j(n)$ and $y_k(n)$
- Two black colored directed links
- One blue colored directed link
- Node $x'_j(n)$ has two input links
 - One from node $x_j(n)$
 - One from node $y_k(n)$

$$\begin{aligned} y_k(n) &= A[x'_j(n)] \\ x'_j(n) &= x_j(n) + \frac{B[y_k(n)]}{\text{feedbackoutput}} \quad (\text{Addition in Directed Graph}) \\ y_k(n) &= A[x_j(n) + \frac{B[y_k(n)]}{\text{feedbackoutput}}] \\ &= A[x_j(n)] + A \frac{B[y_k(n)]}{\text{feedbackoutput}} \\ &= A[x_j(n)] + AB[y_k(n)] \\ y_k(n) &= \frac{A}{1-AB} [x_j(n)] \end{aligned}$$

$$y_k(n) = \frac{A}{1-AB} [x_j(n)]$$

★ Gradient Descent

Algorithm

Step 1 Choose a maximum number of iterations M to be performed, an initial point $x^{(0)}$, two termination parameters ϵ_1, ϵ_2 , and set $k = 0$.

Step 2 Calculate $\nabla f(x^{(k)})$, the first derivative at the point $x^{(k)}$.

Step 3 If $\|\nabla f(x^{(k)})\| \leq \epsilon_1$, **Terminate**; (Slope almost zero)

Else if $k \geq M$; **Terminate**; (No of Iterations Exceeded)

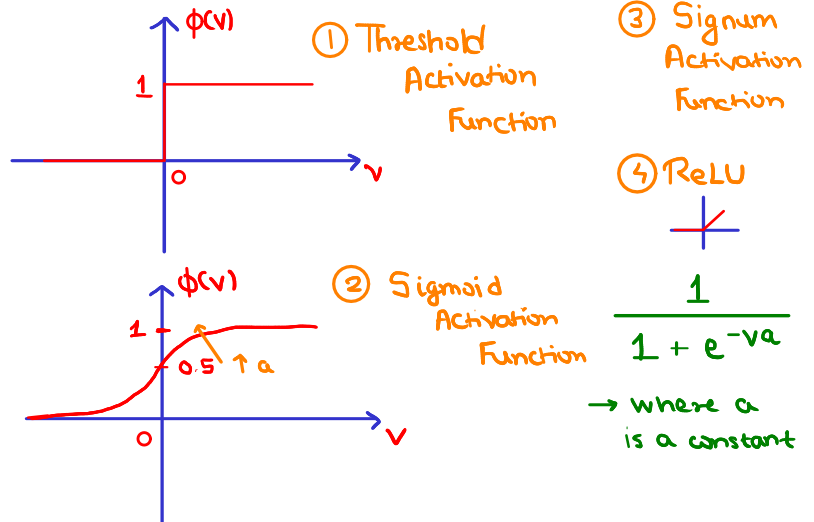
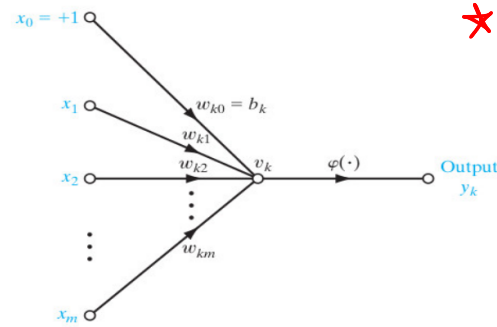
Else go to Step 4.

Step 4 Perform a unidirectional search to find $\alpha^{(k)}$ using ϵ_2 such that $f(x^{(k+1)}) = f(x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}))$ is minimum. One criterion for termination is when $|\nabla f(x^{(k+1)}) \cdot \nabla f(x^{(k)})| \leq \epsilon_2$.

Step 5 Is $\frac{\|x^{(k+1)} - x^{(k)}\|}{\|x^{(k)}\|} \leq \epsilon_1$? If yes, **Terminate**;

Else set $k = k + 1$ and go to Step 2.

★ Directed Graph Rep of the Network



★ A search direction d^t is a descent direction at point x^t if the condition $\nabla f(x^t) \cdot d^t \leq 0$ is satisfied

→ $-\nabla f(x^t)$ → Direction of max descent

★ Here t is the current time / iteration

★ Descent direction Condition:

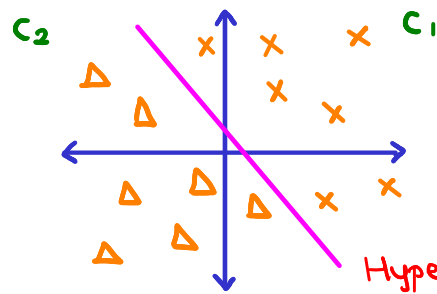
$$\rightarrow f(x^{t+1}) = f(x^t + \alpha \nabla f(x^t) \cdot d^t) < f(x^t)$$

★ Binary Search Type Algorithm

Interval halving method - algorithm

- Step 1 Given interval (a, b) , choose ϵ . Let $x_m = \frac{(a+b)}{2}$; $L = (b - a)$
- Step 2 Initialize $x_1 = a + \frac{1}{4}$; $x_2 = b - \frac{1}{4}$; Compute $f(x_1), f(x_2)$
- Step 3 If $f(x_1) < f(x_m)$ then $b = x_m$; $x_m = x_1$; Go to step 5; else go to step 4
- Step 4 If $f(x_2) < f(x_m)$ then $a = x_m$; $x_m = x_2$; Go to step 5; else $a = x_1$, $b = x_2$; go to step 5; (Very Important)
- Step 5 Calculate $L = (b - a)$. If $(|L| < \epsilon)$ terminate else go to step 2

Linearly Seperable data



$$W = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} \rightarrow \text{Coefficients}$$

$\eta \rightarrow$ Learning Rate

$$x = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} \rightarrow \text{Inputs}$$

★ Higher the value of the learning rate faster the opti sol arrival but less accuracy.

★ Can be generalized to higher dimmensions

$$\star w^T x > 0 \quad \forall x \in C_1$$

$$\star w^T x \leq 0 \quad \forall x \in C_2$$

Update Rule:

$$1) \quad \underline{w(n+1) = w(n)} \quad \text{if } \underline{w^T(n)x(n) > 0} \quad \& \quad \underline{x(n) \in C_1}$$

$$\underline{w(n+1) = w(n)} \quad \text{if } \underline{w^T(n)x(n) \leq 0} \quad \& \quad \underline{x(n) \in C_2}$$

$$2) \quad \underline{w(n+1) = w(n) - \eta(n)x(n)} \quad \text{if } \underline{w^T(n)x(n) > 0} \quad \& \quad \underline{x(n) \in C_2}$$

$$\underline{w(n+1) = w(n) + \eta(n)x(n)} \quad \text{if } \underline{w^T(n)x(n) \leq 0} \quad \& \quad \underline{x(n) \in C_1}$$

$$\text{Let } w = w(0) = 0$$

$$\& \eta(n) = 1$$

$$\Rightarrow x(n) \in C_1 \& w^T(n)x(n) \leq 0$$

$$\star w(n+1) = w(n) + x(n) \quad \because \eta(n) = 1$$

$$\Rightarrow \underline{w(n+1) = x(1) + x(2) + \dots + x(n)} \quad \dots \textcircled{1}$$

★ If C_1 & C_2 are linearly seperable
 $\exists w_0 \ni w_0^T(n)x(n) > 0$
 for $x(1), x(2), \dots, x(n) \in C_1$

$$\text{Let } \alpha = \min_{x(n) \in C_1} w_0^T(n)x(n)$$

$$\star w(k+1) = w(k) + x(k)$$

$$\Rightarrow \|w(k+1)\|^2 - \|w(k)\|^2 \leq \|x(k)\|^2$$

$$\text{Let } \beta = \max_{x(k) \in C_1} \|x(k)\|^2$$

$$\text{Then } \|w(n+1)\|^2 \leq n\beta$$

$$\star \text{Hence} \rightarrow w_0^T w(n+1) \geq n\alpha \quad (\# \text{Multiply both sides})$$

$$\star \|w_0\| \|w(n+1)\| \geq w_0^T w(n+1) \quad (\text{Cauchy Schwarz})$$

$$\Rightarrow \|w(n+1)\|^2 \geq \frac{n^2 \alpha^2}{\|w_0\|^2}$$

$$\Rightarrow n\beta \geq \frac{n^2 \alpha^2}{\|w_0\|^2}$$

$$\Rightarrow \underline{n_{\max} = \frac{\beta \|w_0\|^2}{\alpha^2}} \quad (\text{Max Number of iterations})$$

★ Algorithm

Initialization Set $w(0) = 0$; Perform following computations for $n = 1, 2, \dots$

Activation At time step n , provide the input vector $x(n)$ and desired response $d(n) \rightarrow$ **desired output**

Response $\text{sgn}(w^T(n)x(n))$ Output is $\{-1, +1\} \rightarrow y(n)$

Adaptation $w(n+1) = w(n) + \eta[d(n) - y(n)]x(n)$

Where

$$d(n) = \begin{cases} +1 & \text{if } x(n) \in C_1 \\ -1 & \text{if } x(n) \in C_2 \end{cases}$$

Iterate Increment n and go to activation step

★ **Example:**

$$\underline{\text{If } w^T(n)x(n) > 0}$$

$$\Rightarrow y(n) = +1$$

$$\underline{\text{if } x(n) \in C_1} \Rightarrow d(n) = 1 \rightarrow w(n+1) = w(n)$$

$$\underline{\text{if } x(n) \in C_2} \Rightarrow d(n) = -1 \rightarrow w(n+1) = w(n) - 2\eta x(n)$$

★ Batch Algorithm

- Compute: $\mathbf{w}^T(n)\mathbf{x}(n)$
- Treat the above quantity as the objective function
- With the modification $\mathbf{w}^T(n)\mathbf{x}(n)d(n)$
- For one $\mathbf{x}(n)$ the above objective function is used:
- For many $\mathbf{x}(n)$'s we have:

$$J(\mathbf{w}) = \sum_{\mathbf{x}(n) \in \mathcal{H}} \left(-(\mathbf{w}^T(n)\mathbf{x}(n)d(n)) \right)$$

- The above objective function should be **minimized**

→ Objective Cost Function

- We have to minimize or maximize a given objective function
- Perceptron rule: $\mathbf{w}^T(n)\mathbf{x}(n) > 0 \Rightarrow \mathbf{x}(n) \in \mathcal{C}_1$
- $\mathbf{w}^T(n)\mathbf{x}(n)$ quantity for \mathcal{C}_1 is positive, $d(n) = +1$. Decrease it by multiplying it -1
- Perceptron rule: $\mathbf{w}^T(n)\mathbf{x}(n) \leq 0 \Rightarrow \mathbf{x}(n) \in \mathcal{C}_2$
- $\mathbf{w}^T(n)\mathbf{x}(n)$ quantity for \mathcal{C}_2 is negative, $d(n) = -1$. Decrease it by multiplying it -1
- That is for any $\mathbf{x}(n)$, the quantity $-(\mathbf{w}^T(n)\mathbf{x}(n)d(n))$ to be **minimized**

→ Intuition

- Compute direction: $\nabla J(\mathbf{w}) = \sum_{\mathbf{x}(n) \in \mathcal{H}} (-\mathbf{x}(n)d(n))$
- Update $\mathbf{w}(n+1) = \mathbf{w}(n) - \eta(n) \nabla J(\mathbf{w})$
- That is $\mathbf{w}(n+1) = \mathbf{w}(n) - \eta(n) \sum_{\mathbf{x}(n) \in \mathcal{H}} (-\mathbf{x}(n)d(n))$

→ Gradient Descent Rule

★ MultiLayer Perceptron

- error at each neuron: $e_j(n) = d_j(n) - y_j(n)$
- Instantaneous error energy $\mathcal{E}_j(n) = \frac{1}{2} e_j^2(n)$
- Avg error $\mathcal{E}_{av}(N) = \frac{1}{2N} \sum_{n=1}^N \sum_{j \in \mathcal{C}} e_j^2(n)$ → Here N = no. of training examples

- $\mathcal{E}_{av} \rightarrow \mathcal{E}(n)$
- $\mathcal{E}(n) \rightarrow e_j(n)$
- $e_j(n) \rightarrow y_j(n)$
- $y_j(n) \rightarrow v_j(n)$
- $v_j(n) \rightarrow w_{ji}(n)$

$$\begin{aligned} \mathcal{E}(n) &= \frac{1}{2} \sum_{j \in \mathcal{C}} e_j^2(n) \\ e_j(n) &= d_j(n) - y_j(n) \\ y_j(n) &= \phi(v_j(n)) \\ \sum_{i=0}^m w_{ji}(n) y_i(n) &= v_j(n) \end{aligned}$$

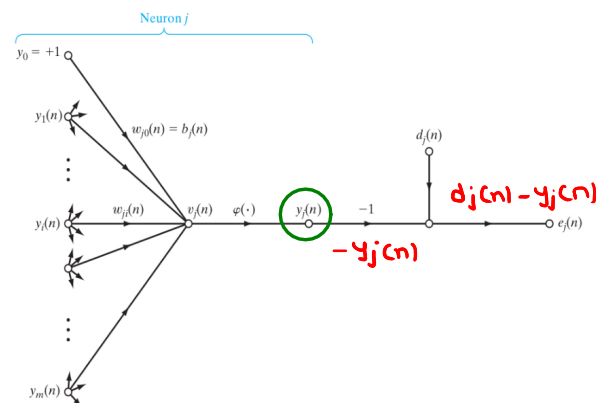
Dependencies used for gradient cal.

$$\frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} = \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \frac{\partial v_j(n)}{\partial w_{ji}(n)}$$

Chain Rule

\downarrow \downarrow \downarrow \downarrow
 $e_j(n)$ -1 $\phi_j'(v_j(n))$ $y_i(n)$

★ Back Propagation Algo Online



★ IF j is a hidden layer we don't know the desired output hence can't find $e_j(n)$

- $w_{ji}(n+1) = w_{ji}(n) - \eta \frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)}$ **★ Local gradient output neuron.**
- $\Delta w_{ji}(n) = -\eta \frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)}$
- $\Delta w_{ji}(n) = \eta \delta_j(n) y_i(n)$
- Where $\delta_j(n)$ is the local gradient at v_j (before the activation function)
- Local gradient is computed using chain rule as:

$$\begin{aligned} \star \delta_j(n) &= \frac{\partial \mathcal{E}(n)}{\partial v_j(n)} \\ &= \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \\ &= e_j(n) \times -1 \times \phi'_j(v_j(n)) \end{aligned}$$

$$(output) \delta_j(n) = \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)}$$

$$\begin{aligned} (hidden) \delta_j(n) &= \frac{\partial \mathcal{E}(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \\ &= \frac{\partial \mathcal{E}(n)}{\partial y_j(n)} \phi'_j(v_j(n)) \end{aligned}$$

★ When k is a hidden neuron, the error is computed by summing all the hidden neuron errors. $\Rightarrow \mathcal{E}(n) = \frac{1}{2} \sum_k e_k^2(n)$

$$\begin{aligned} \frac{\partial \mathcal{E}(n)}{\partial y_j(n)} &= \sum_k e_k(n) \frac{\partial e_k(n)}{\partial y_j(n)} \\ &= \sum_k e_k(n) \frac{\partial e_k(n)}{\partial v_k(n)} \frac{\partial v_k(n)}{\partial y_j(n)} \end{aligned}$$

$$\downarrow \quad \downarrow$$

$$\phi'_k(v_k(n)) \quad w_{kj}(n)$$

$$\star e_k(n) = d_k(n) - \phi_k(v_k(n))$$

$$\star v_k(n) = \sum_{j=0}^m w_{kj}(n) y_j(n)$$

- The complete derivative is:

$$\frac{\partial v_k(n)}{\partial y_j(n)} = w_{kj}(n)$$

$$\begin{aligned} \frac{\partial \mathcal{E}(n)}{\partial y_j(n)} &= \sum_k e_k(n) \frac{\partial e_k(n)}{\partial y_j(n)} \\ &= \sum_k e_k(n) \phi'_k(v_k(n)) w_{kj}(n) \\ &= \sum_k e_k(n) \phi'_k(v_k(n)) w_{kj}(n) \\ &= \sum_k \delta_k(n) w_{kj}(n) \end{aligned}$$

- For all the k^{th} neurons in the forward layer that connect to j^{th} neuron

★ Gradient Descent rule in Online MLP

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

\downarrow **Weight correction** \downarrow **Learning rate** \downarrow **Local gradient** \downarrow **Input signal at neuron j**

$$\bullet w_{ji}(n+1) = w_{ji}(n) - \eta \frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)}$$

$$\bullet \Delta w_{ji}(n) = -\eta \frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)}$$

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

$$\delta_j(n) = \begin{cases} a[d_j(n) - o_j(n)]o_j(n)[1 - o_j(n)] & \text{if } j \text{ is output neuron} \\ a y_j(n)[1 - y_j(n)] \sum_k \delta_k(n) w_{kj}(n) & \text{if } j \text{ is a hidden neuron} \end{cases}$$

→ for sigmoid function:

$$\phi_j(v_j(n)) = \frac{1}{1 + e^{-av_j(n)}}$$

$$\Rightarrow \phi'_j(v_j(n)) = a y_j(n) [1 - y_j(n)]$$

$$\# y_i(n) = \phi_j(v_j(n))$$