

Training of Neural Network (This Lecture And Next Week)

Notation (e chapter 7)

Input Layer: $L=0$

Hidden Layers: $1 \leq L \leq L-1$

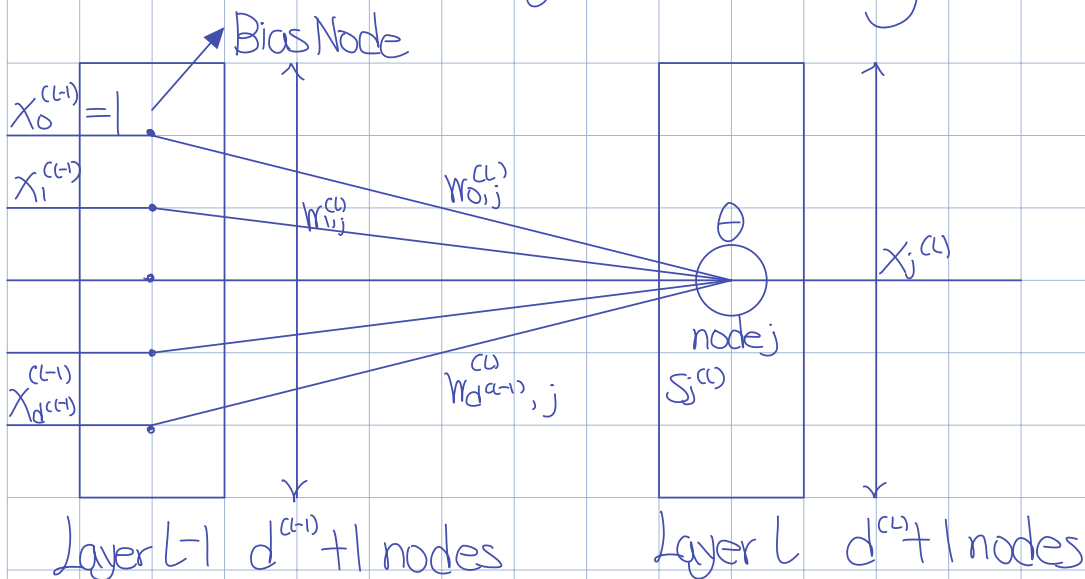
Output Layer: $L=L$

} L layers neural network

$w_{i,j}^{(L)} \equiv$ weight connecting node i in layer $(L-1)$ to node j in layer L

Not Power, But Superscript

$d^{(L)} \equiv$ # of nodes in layer L (not counting bias node)



$x_j^{(L)} \equiv$ output from node j in layer L

$s_j^{(L)} \equiv$ input into node j in layer L

$$s_j^{(L)} = w_{0,j}^{(L)} + \sum_{i=1}^{d^{(L-1)}} w_{i,j}^{(L)} x_i^{(L-1)} \quad (*)$$

$$x_j^{(L)} = \theta(s_j^{(L)})$$

$\Theta(s) \equiv$ Activation Function

$$\Theta(s) = \tanh(s) \\ = \text{ReLU}(s)$$

Vector Notation

$$S^{(l)} = \begin{bmatrix} S_1^{(l)} \\ \vdots \\ S_{d^{(l)}}^{(l)} \end{bmatrix} \quad \Theta(S^{(l)}) = \begin{bmatrix} \Theta(S_1^{(l)}) \\ \vdots \\ \Theta(S_{d^{(l)}}^{(l)}) \end{bmatrix}$$

$$X^{(l)} = \begin{bmatrix} 1 \\ \Theta(S^{(l)}) \end{bmatrix}$$

Weight Matrix

$$W^{(l)} = \left\{ W_{i,j}^{(l)} \right\}_{\substack{0 \leq i \leq d^{(l-1)} \\ 1 \leq j \leq d^{(l)}}}$$

Bias Node From Layer $(l-1)$ to Nodes In Layer (l)

$$W^{(l)} = \begin{bmatrix} W_{0,1}^{(l)} & W_{0,2}^{(l)} & \dots & W_{0,d^{(l)}}^{(l)} \\ W_{1,1}^{(l)} & W_{1,2}^{(l)} & \dots & W_{1,d^{(l)}}^{(l)} \\ \vdots & \vdots & \ddots & \vdots \\ W_{d^{(l-1)},1}^{(l)} & W_{d^{(l-1)},2}^{(l)} & \dots & W_{d^{(l-1)},d^{(l)}}^{(l)} \end{bmatrix}$$

j th column Incident to node j in Layer l

In Vector Form

$$S^{(l)} = W^{(l)T} X^{(l-1)}$$

Given, Input Vector

$$X^{(0)} = [1 \ x_1 \ x_2 \ \dots \ x_d]^T$$

$$X^{(0)} \xrightarrow{W^{(1)}} S^{(1)} \longrightarrow X^{(1)}$$

$$X^{(1)} \xrightarrow{W^{(2)}} S^{(2)} \longrightarrow X^{(2)}$$

$$X^{(L-1)} \xrightarrow{W^{(L)}} S^{(L)} \longrightarrow X^{(L)} \text{ (Final Output)}$$

Forward Propagation Algorithm

$$\text{Input } X^{(0)} = [1 \ x_1 \ x_2 \ \dots \ x_d]^T$$

for $l = 1, 2, \dots, L$, do

$$S^{(l)} = W^{(l)T} X^{(l-1)}$$

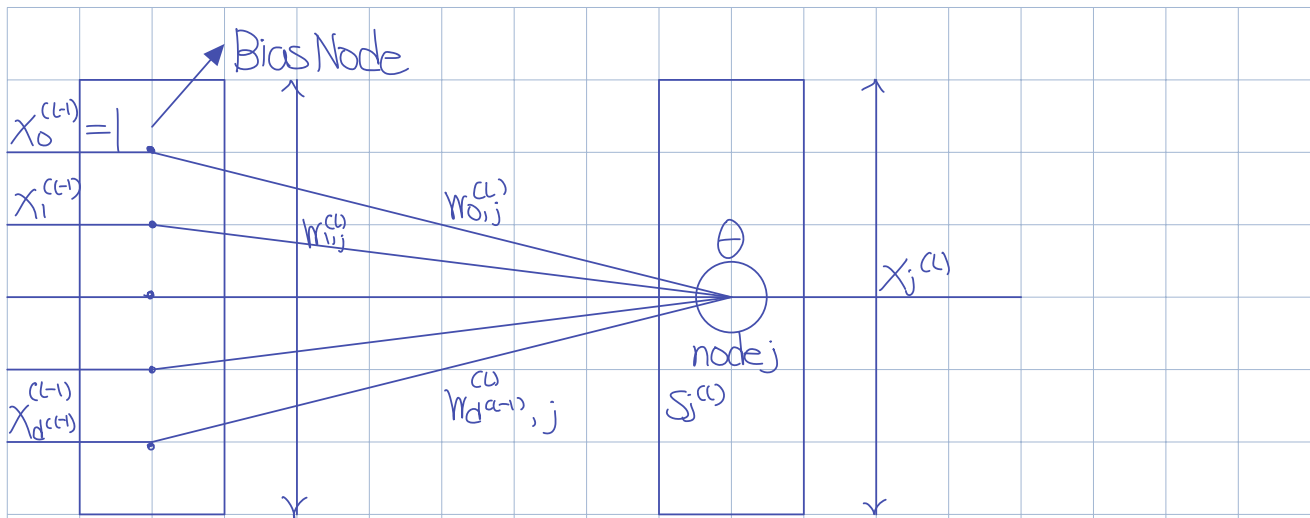
$$X^{(l)} = \Theta(S^{(l)})$$

end

output $X^{(L)}$

Computational Complexity of Forward Propagation

of computation needed in layer l :



$d^{(l-1)} + 1$ Computation for one node
 $(d^{(l-1)} + 1) d^{(l)}$ Computation for one layer

Operations for computing $s^{(l)}$

Activation Function $d^{(l)}$

$(d^{(l-1)} + 1) d^{(l)} + d^{(l)}$ for computing $x^{(l)}$

Total # of Computations

$$\sum_{l=1}^L (d^{(l-1)} + 1) d^{(l)} + \sum_{l=1}^L d^{(l)}$$

of edges = γ

of nodes = Q

\approx few millions

Linear

Summary So Far

Introduced Neural Network

Model Parameters

$$\Omega = \{W^{(1)}, W^{(2)}, \dots, W^{(L)}\}$$

$$\text{Input, } X^{(0)} = \begin{bmatrix} 1 & x_0 & x_1 & \dots & x_d \end{bmatrix}^T$$

$$\text{Output, } X^{(L)}$$

Loss Function

$$\text{Regression, } d^{(L)} = 1, g(X^{(L)}, y) = (X^{(L)} - y)^2$$

Classification,

$$\text{Output } X^{(L)} = \begin{bmatrix} \hat{p}_1 & \hat{p}_2 & \dots & \hat{p}_C \end{bmatrix}, d^{(L)} = C$$

$$\text{Loss } g(X^{(L)}, y) = -\log \hat{p}_y$$

Log Loss Function

Define $E_n(\Omega) \equiv$ Loss on training sample (x_n, y_n)

$$E_n(\Omega) = g(X_n^{(L)}, y_n) = (X_n^{(L)} - y_n)^2 \text{ Regression}$$

Average Loss

$$E_n(\Omega) = \frac{1}{N} \sum_{n=1}^N E_n(\Omega)$$

$$\Omega^* = \underset{\Omega}{\operatorname{argmin}} E_n(\Omega)$$

Gradient Descent Until Converge

Backpropagation