
Applied Machine Learning with Big Data “EE 6973”



Topic:
Deep Learning

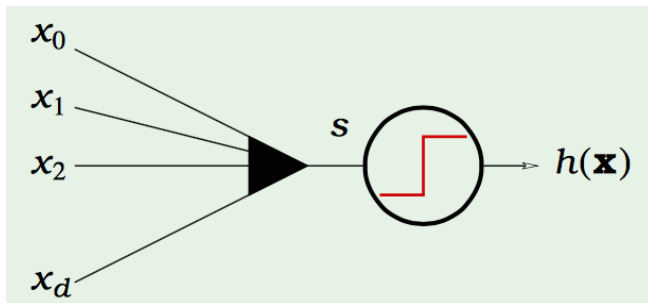
Paul Rad, Ph.D.

Chief Research Officer
UTSA Open Cloud Institute(OCI)
University of Texas at San Antonio

Review Models

Linear Classification

$$h(x) = \text{Sign} \left(\sum_{i=0}^n w_i x_i \right)$$



Sign function:

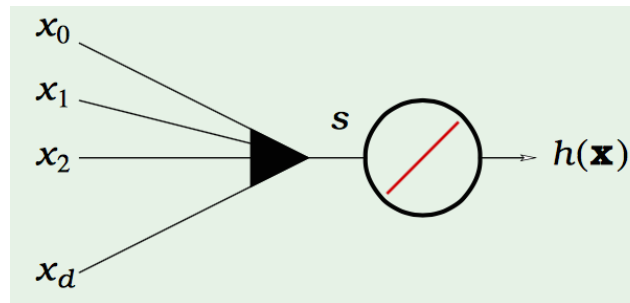
$$h(s) = 1 \quad s \geq 0$$

$$h(s) = 0 \quad s < 0$$

Hard Threshold : Certainty

Linear Regression

$$h(x) = \sum_{i=0}^n w_i x_i$$

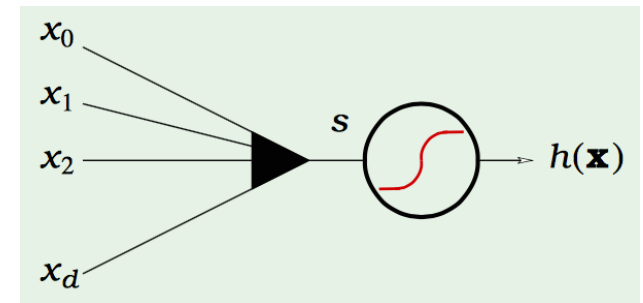


identity function:

$$h(s) = s$$

Logistic Regression

$$h(x) = \text{Sigmoid} \left(\sum_{i=0}^n w_i x_i \right)$$

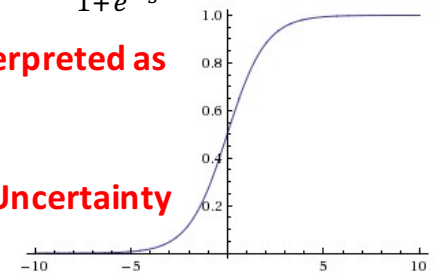


Sigmoid function:

$$h(s) = \frac{1}{1 + e^{-s}}$$

The output is interpreted as probability

Soft Threshold : Uncertainty



Probability Interpretation

$h(\mathbf{x}) = \text{Sigmoid}(\sum_{i=0}^n w_i x_i) = \theta(s)$ is interpreted as a probability

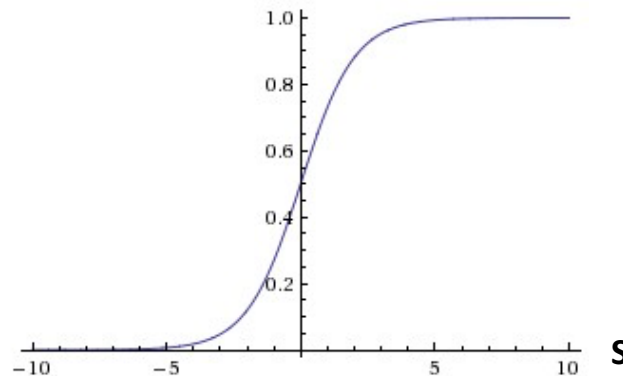
Example: Prediction of heart attacks

Input \mathbf{X} : x_1 =cholesterol level, x_2 =patient age, x_3 =patient weight, etc.

$\theta(s)$: probability of a heart attack

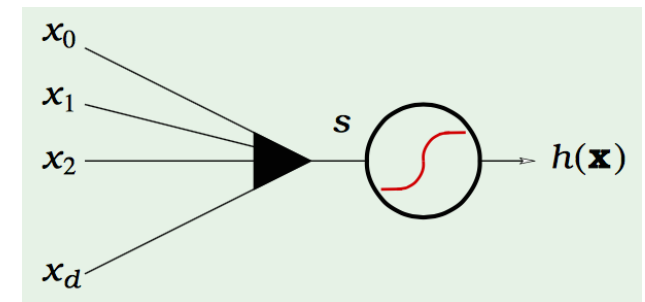
$S = \sum_{i=0}^n w_i x_i$ “risk score”

$$\theta(s) = \frac{1}{1+e^{-s}}$$



Logistic Regression

$$h(\mathbf{x}) = \text{Sigmoid}(\sum_{i=0}^n w_i x_i)$$



Sigmoid function

Outline

Neural Network Model

Error and Learning

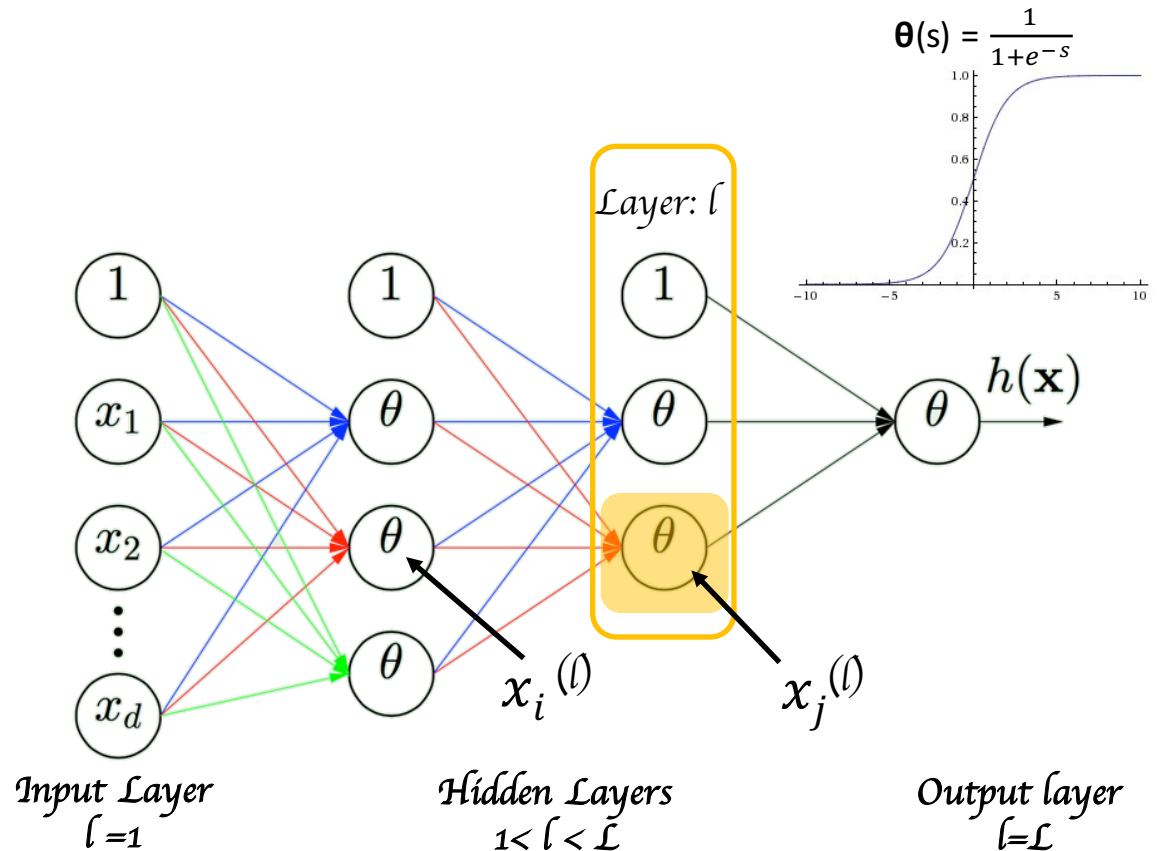
Backpropagation

Neural Network Model (forward propagation)

$$w_{ij}^{(l-1,l)} \left\{ \begin{array}{ll} 2 \leq l \leq L & \text{Layers} \\ 0 \leq i \leq d^{(l-1)} & \text{inputs} \\ 1 \leq j \leq d^{(l)} & \text{outputs} \end{array} \right.$$

$$x_j^{(l)} = \theta \left(\sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l-1,l)} x_i^{(l-1)} \right)$$

$$2 \leq l \leq L, 1 \leq j \leq d^{(l)}$$



Derivative of Sigmoid Function

$$\frac{ds(x)}{dx} = \frac{1}{1 + e^{-x}}$$

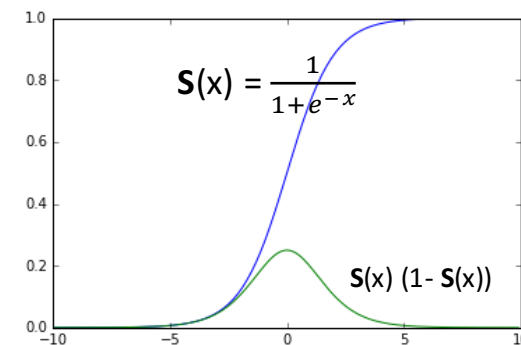
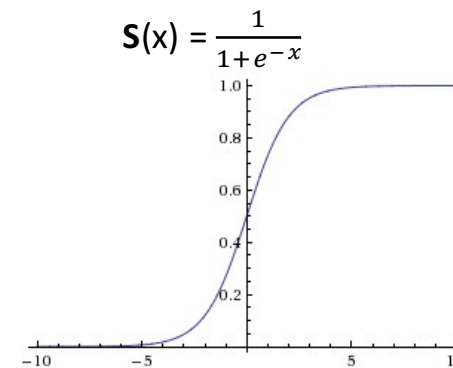
$$= \left(\frac{1}{1 + e^{-x}} \right)^2 \frac{d}{dx} (1 + e^{-x})$$

$$= \left(\frac{1}{1 + e^{-x}} \right)^2 e^{-x} (-1)$$

$$= \left(\frac{1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) (-e^{-x})$$

$$= \left(\frac{1}{1 + e^{-x}} \right) \left(\frac{-e^{-x}}{1 + e^{-x}} \right)$$

$$= s(x)(1 - s(x))$$



Error

$$x_j^{(l)} = \theta \left(\sum_{i=0}^d w_{ij}^{(l-1)} x_i^{(l-1)} \right) \quad h(x) = x_j^{(l)} \text{ when } l = \mathcal{L}$$

If we have all the weights $w = \{w_{ij}^{(l-1,l)}\}$ for each layer then we can determine $h(x)$.

Error on example (X_n, Y_n) is $e(h(X_n), Y_n) = e(w)$

We need to calculate Gradient: $\nabla e(\mathbf{w}): \frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$ for all i, j, l

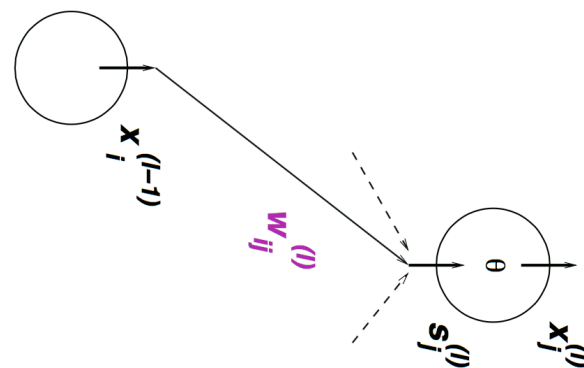
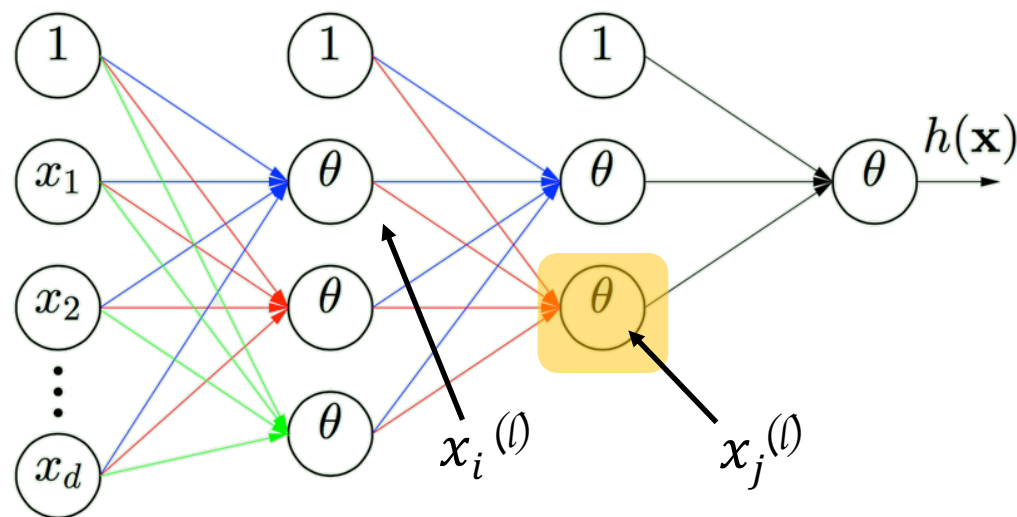
Calculating partial derivatives of $e(\mathcal{W})$

$$\nabla e(\mathbf{w}): \frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}} \text{ for all } i, j, l$$

$$\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}$$

\swarrow \searrow
 $\frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} = \delta_j^{(l)}$ $x_i^{(l-1)}$

$$\frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} = \delta_j^{(l)}$$



For the final layer $\delta_j^{(l)}$

$$\frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} = \delta_j^{(l)} \quad \text{For the final layer } l = \mathcal{L} \text{ and } j=1$$

$$e(\mathbf{w}) = (x_1^{(\mathcal{L})} - y)^2$$

$$x_1^{(\mathcal{L})} = \Theta(s_1^{(\mathcal{L})})$$

$$\Theta'(s) = \Theta(x) (1 - \Theta(x))$$

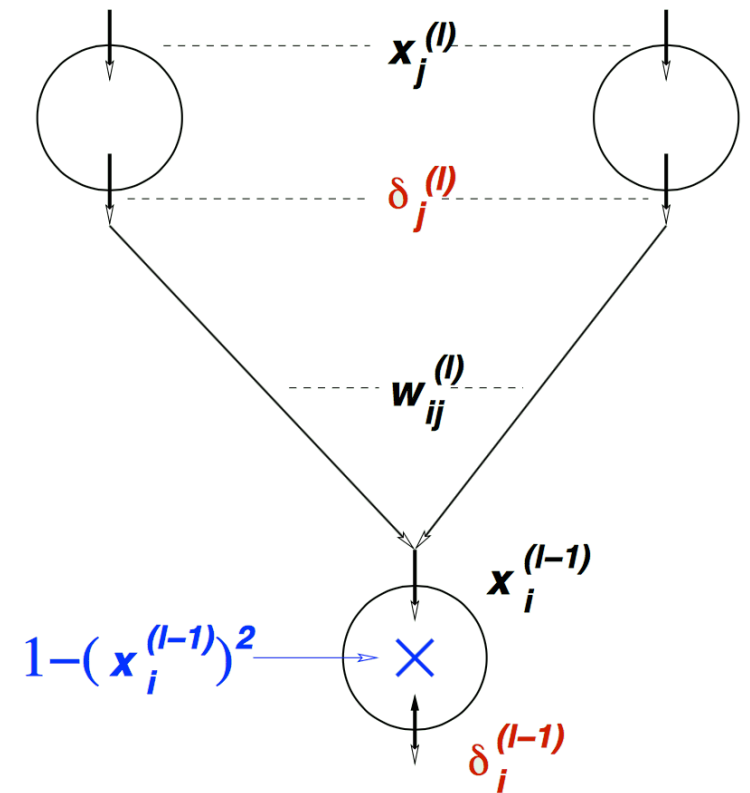
$$\frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} = \delta_j^{(l)}$$

$$\frac{\partial (\Theta(s_1^{(\mathcal{L})}) - y)^2}{\partial s_j^{(l)}}$$

Back propagation of $\delta_j^{(l)}$

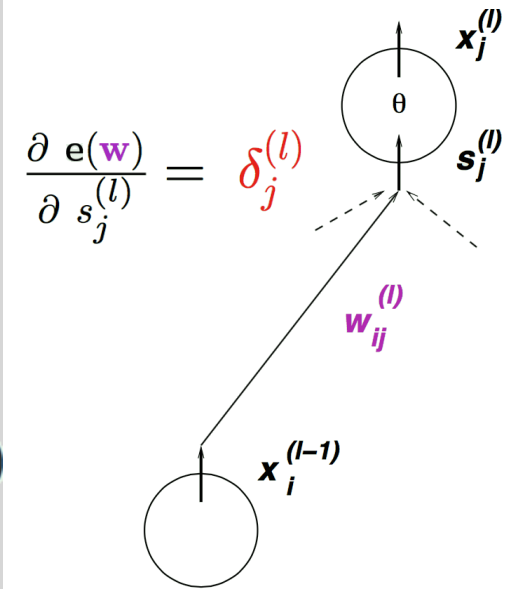
$$\begin{aligned}
 \delta_i^{(l-1)} &= \frac{\partial e(\mathbf{w})}{\partial s_i^{(l-1)}} \\
 &= \sum_{j=1}^{d^{(l)}} \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}} \\
 &= \sum_{j=1}^{d^{(l)}} \delta_j^{(l)} \times w_{ij}^{(l)} \times \theta'(s_i^{(l-1)})
 \end{aligned}$$

$$\delta_i^{(l-1)} = x_i^{(l-1)} (1 - x_i^{(l-1)}) \sum_{j=1}^{d^{(l)}} w_{ij}^{(l)} \delta_j^{(l)}$$

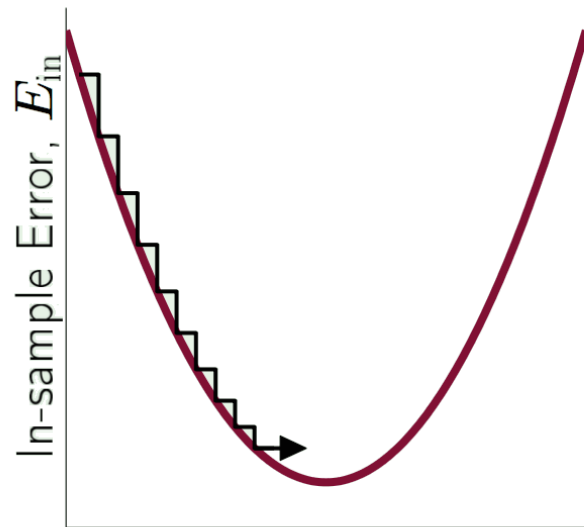


Backpropagation Algorithm

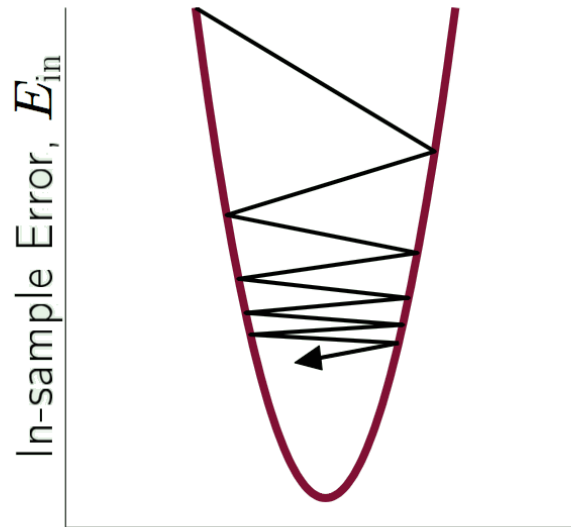
- 1: Initialize all weights $w_{ij}^{(l)}$ **at random**
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: Pick $n \in \{1, 2, \dots, N\}$
- 4: *Forward:* Compute all $x_j^{(l)}$
- 5: *Backward:* Compute all $\delta_j^{(l)}$
- 6: Update the weights: $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)}$
- 7: Iterate to the next step until it is time to stop
- 8: Return the final weights $w_{ij}^{(l)}$



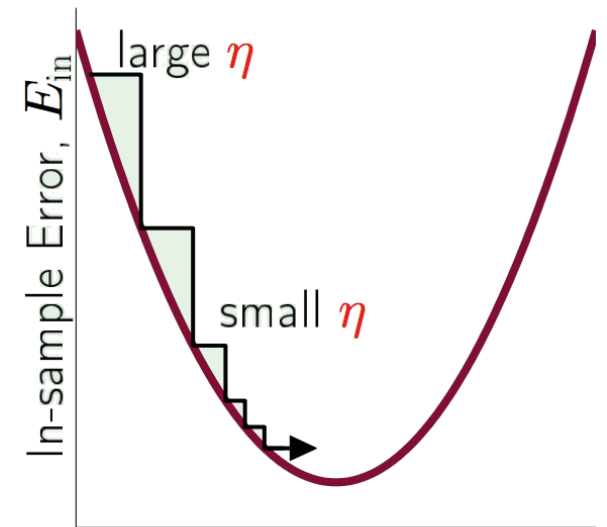
Learning Rate: Steps



η too small



η too large



variable η – just right

Learning Rate should increase with the slope