

---

# 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

---

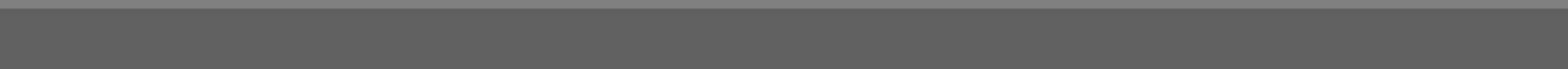
# Outline

---

Neural Network Model

Forward Propagation (Prediction)

Backpropagation (Learning)

- Parameters (Weights)
  - Training Data
  - Cost Function and Error
  - Learning Rate
- 

# Deep Learning Architecture

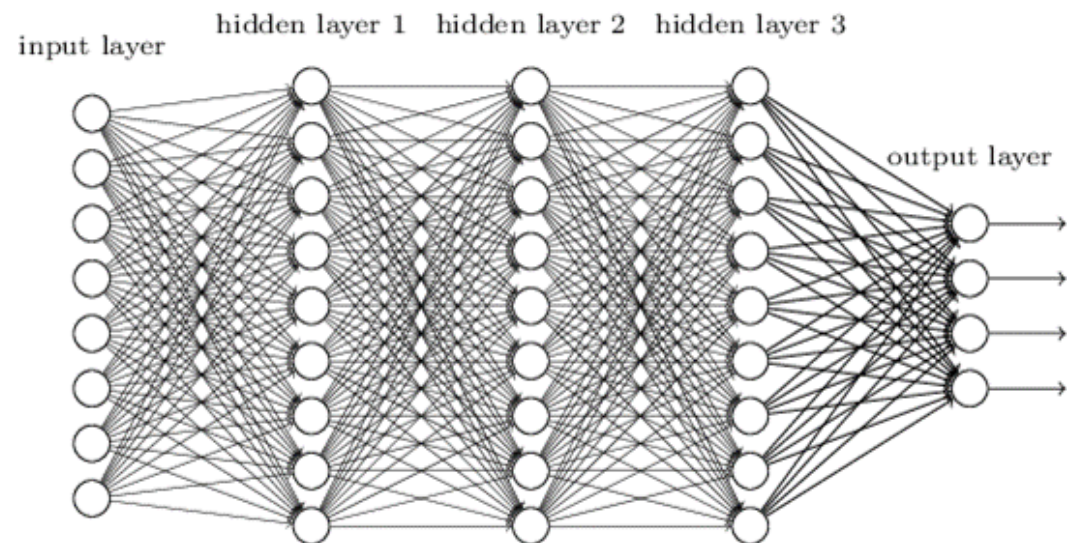
An Artificial Neural Network with one or more hidden layers = Deep Learning

They typically consist of many hundreds of simple processing units which are wired together in a complex communication network.

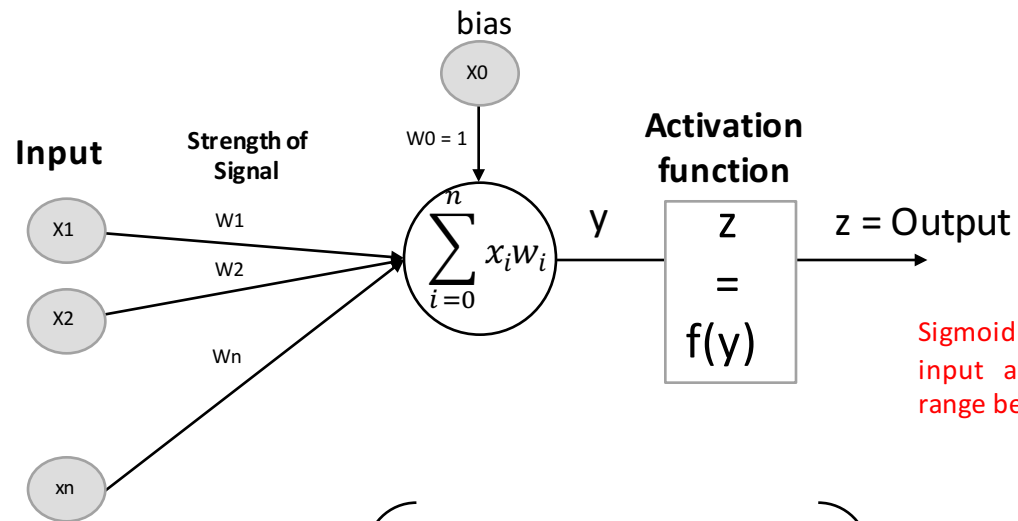
Each unit or node is a simplified model of a real neuron which fires (sends off a new signal) if it receives a sufficiently strong input signal from the other nodes to which it is connected.

The output is aiming for a target.

Deep neural network



# Single Layer Neural Networks with Nonlinear Math Model

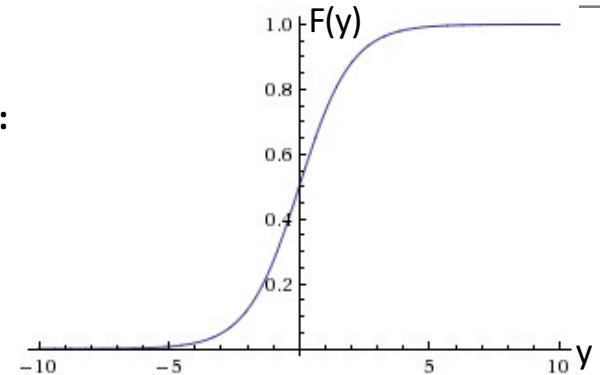


**Sigmoid function:**

$$f(y) = \frac{1}{1 + e^{-y}}$$

$$[0, 1]$$

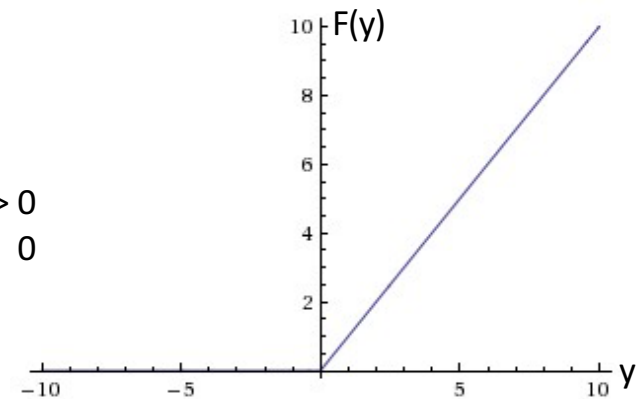
Sigmoid takes a real-valued input and squashes it to range between 0 and 1



$$\text{Output} = f \left( \begin{bmatrix} x_0 & x_1 & x_2 & \dots & x_n \end{bmatrix} \times \begin{bmatrix} 1 \\ w_1 \\ w_2 \\ \dots \\ w_n \end{bmatrix} \right)$$

**ReLU function:**

$$f(y) = \begin{cases} y & y \geq 0 \\ 0 & y < 0 \end{cases}$$

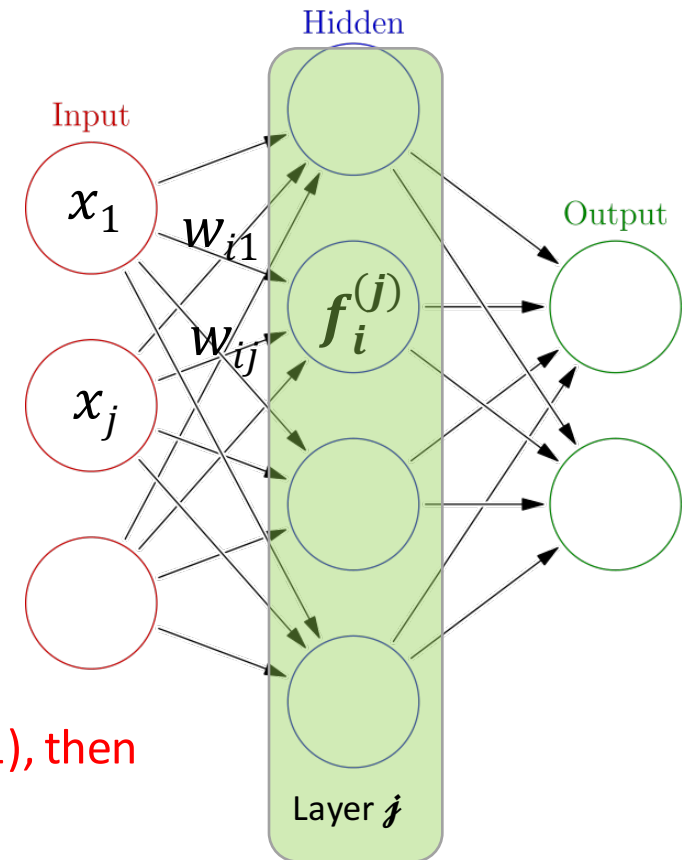


# Multi Layer NN Nonlinear Math Model

$f_i^{(j)}$  = “activation function” of unit  $i$  in layer  $j$

$W^{(j)}$  = matrix of weights controlling function mapping from layer  $j$  to layer  $j+1$

$$y_i = f_i \left( \sum_{j=0}^n w_{ij} x_j \right)$$



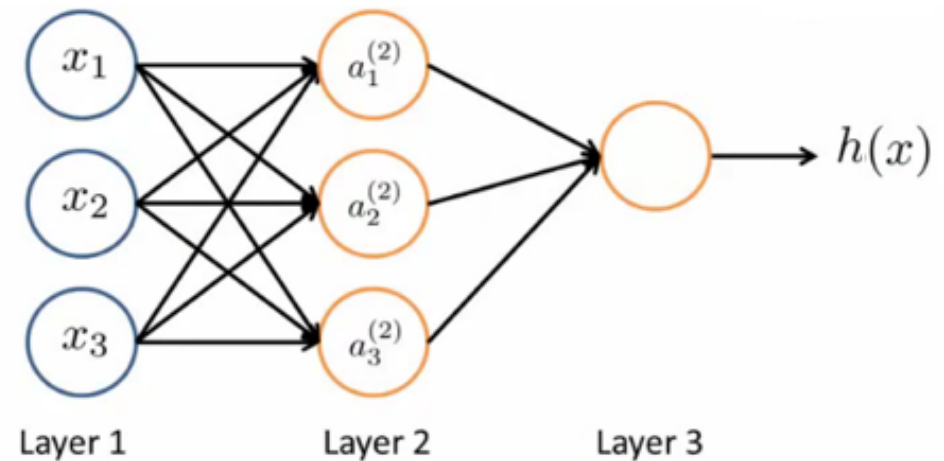
If network has  $N$  units in layer  $j$ , and  $M$  units in layer  $(j-1)$ , then  $W^{(j)}$  will be of dimension of  $(M+1) \times N$

# Forward Propagation to Calculate $h(x)$

$$a_1^{(2)} = f\left(\sum_{j=0}^n w_{1j} x_j\right) =$$

$$a_2^{(2)} = f\left(\sum_{j=0}^n w_{2j} x_j\right) =$$

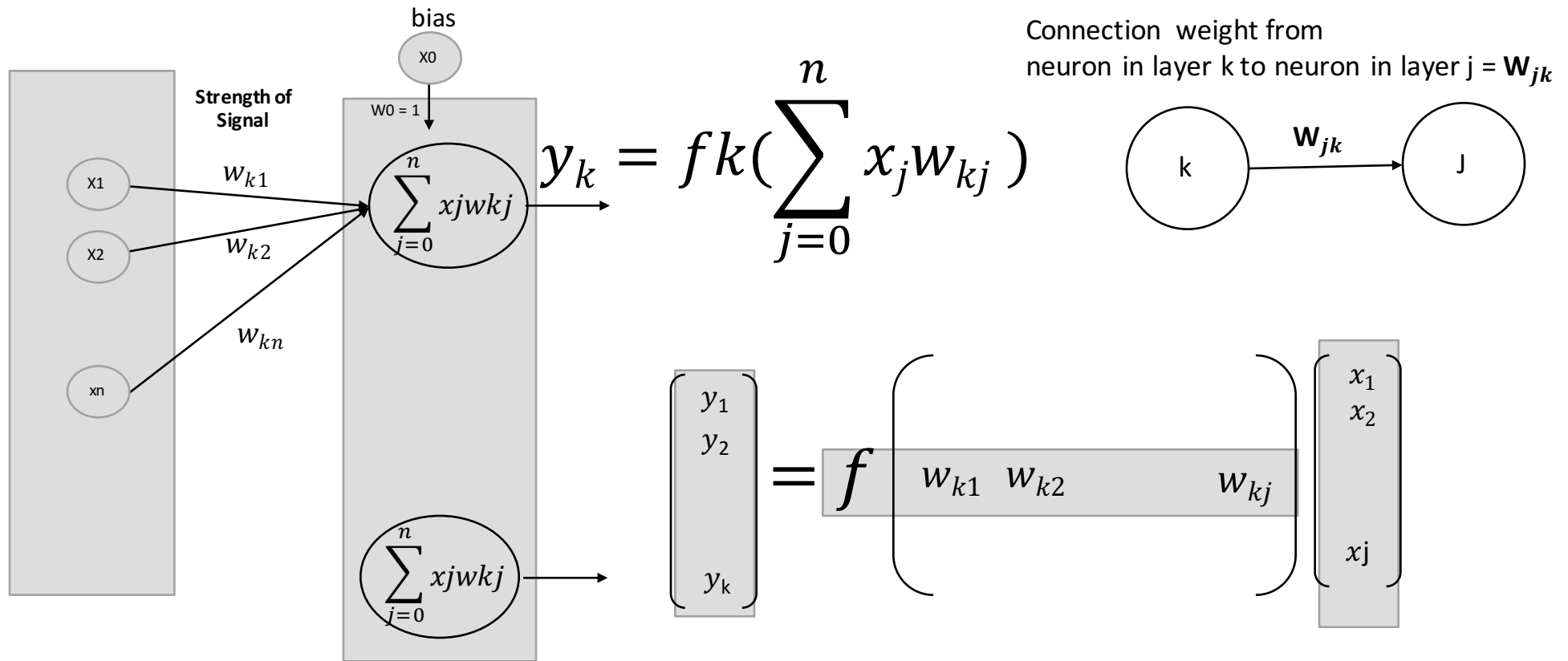
$$a_3^{(2)} = f\left(\sum_{j=0}^n w_{3j} x_j\right) =$$



$$a_i^{(2)} = f\left(\sum_{j=0}^n w_{ij} x_j\right)$$

$$h(x) = f(w_3 a_3^{(2)} + w_2 a_2^{(2)} + w_1 a_1^{(2)} + a_0)$$

# Multi Layer NN Nonlinear Math Model



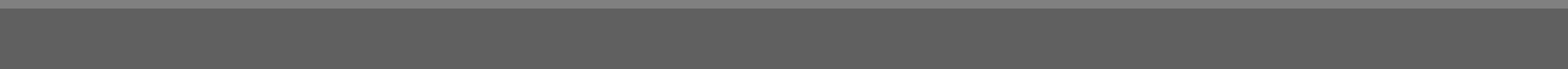
# Cost Function

---

For logistic regression,

$$E(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left( 1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right)$$

Compare to linear regression:

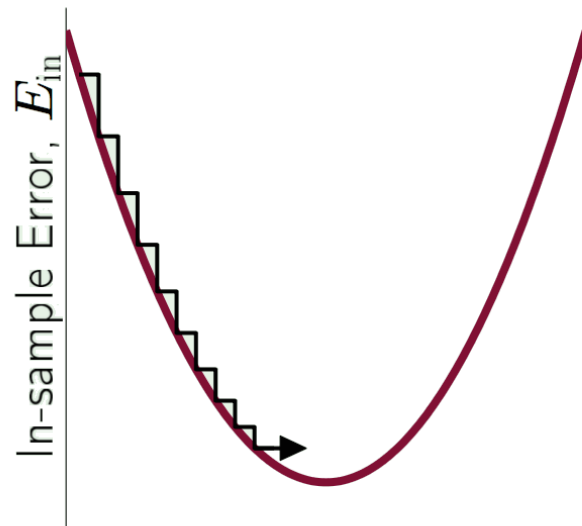
$$E(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2$$




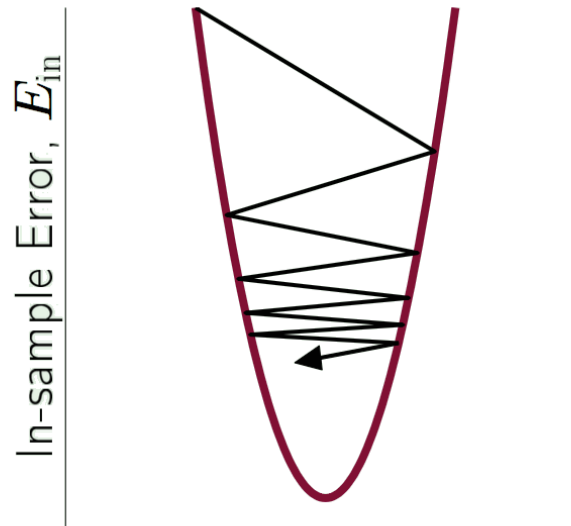
# Backpropagation Algorithm (Learning)

- 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  $\nabla e(\mathbf{w}): \frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$  for all  $i, j, l$
- 6:   Update the weights:  $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \nabla e(\mathbf{w}): \frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$  for all  $i, j, l$
- 7:   Iterate to the next step until it is time to stop
- 8: Return the final weights  $w_{ij}^{(l)}$

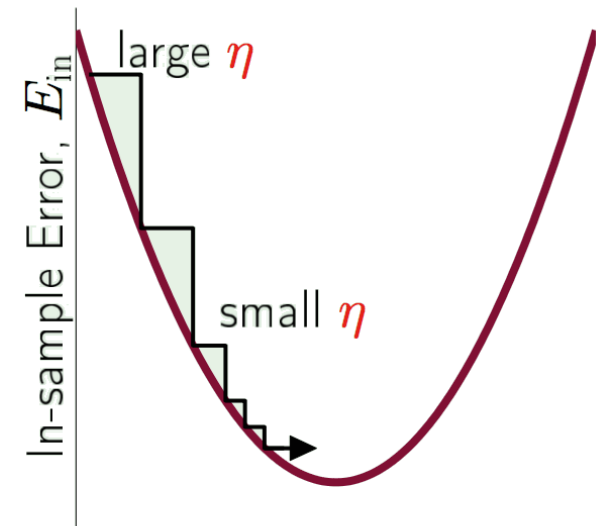
# Learning Rate: Steps



$\eta$  too small



$\eta$  too large



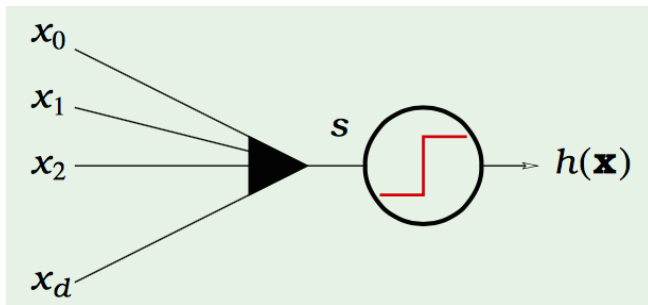
variable  $\eta$  – just right

Learning Rate should increase with the slope

# Review Models

## Linear Classification

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



**Sign function:**

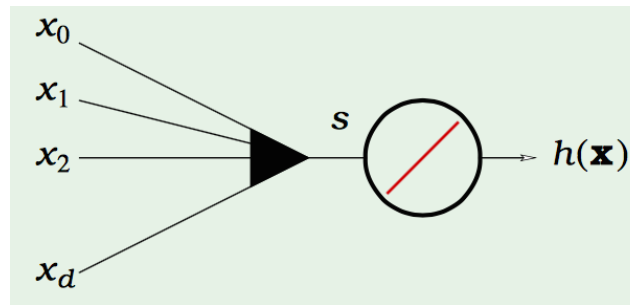
$$h(s) = 1 \quad s \geq 0$$

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

**Hard Threshold : Certainty**

## Linear Regression

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

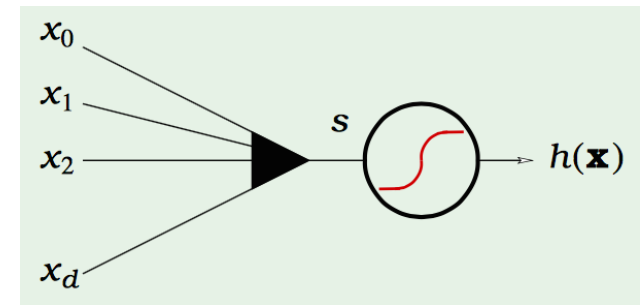


**identity function:**

$$h(s) = s$$

## Logistic Regression

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

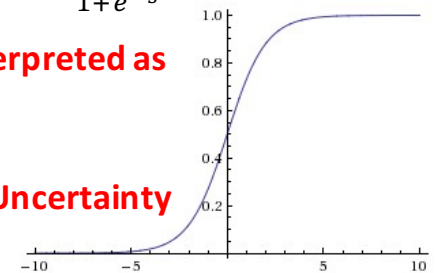


**Sigmoid function:**

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

**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

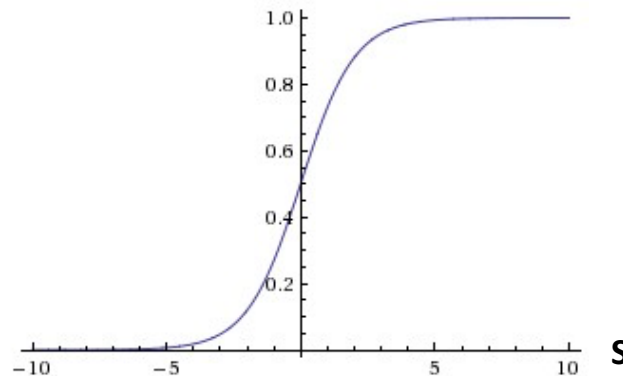
**Example: Prediction of heart attacks**

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

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

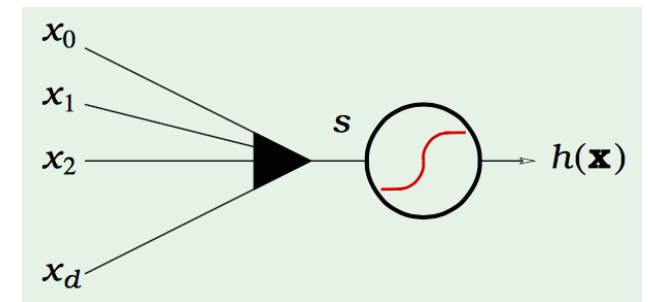
$S = \sum_{i=0}^n w_i x_i$  “risk score”

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



## Logistic Regression

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



Sigmoid function

# Derivative of Sigmoid Function

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

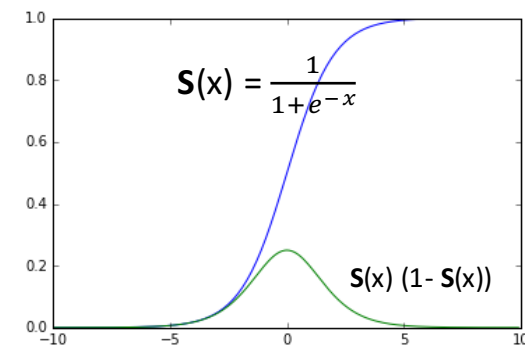
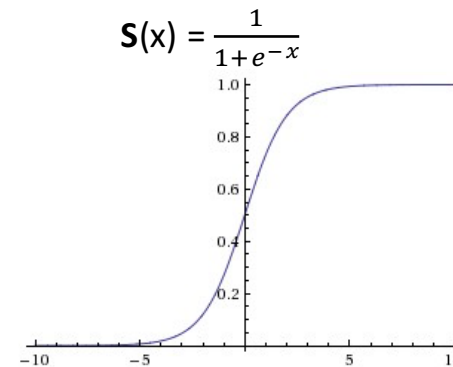
$$= \left( \frac{1}{1 + e^{-x}} \right)^2 \frac{d}{dx} (1 + e^{-x})$$

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

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

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

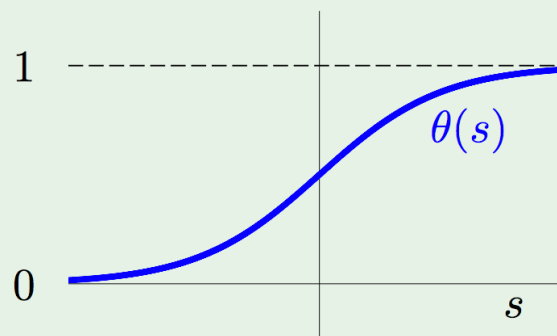
$$= s(x)(1 - s(x))$$



## Formula for likelihood

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute  $h(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x})$ , noting  $\theta(-s) = 1 - \theta(s)$



$$P(y \mid \mathbf{x}) = \theta(y \mathbf{w}^\top \mathbf{x})$$

Likelihood of  $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n)$$

Minimize

$$-\frac{1}{N} \ln \left( \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$= \frac{1}{N} \sum_{n=1}^N \ln \left( \frac{1}{\theta(y_n \mathbf{w}^\top \mathbf{x}_n)} \right)$$

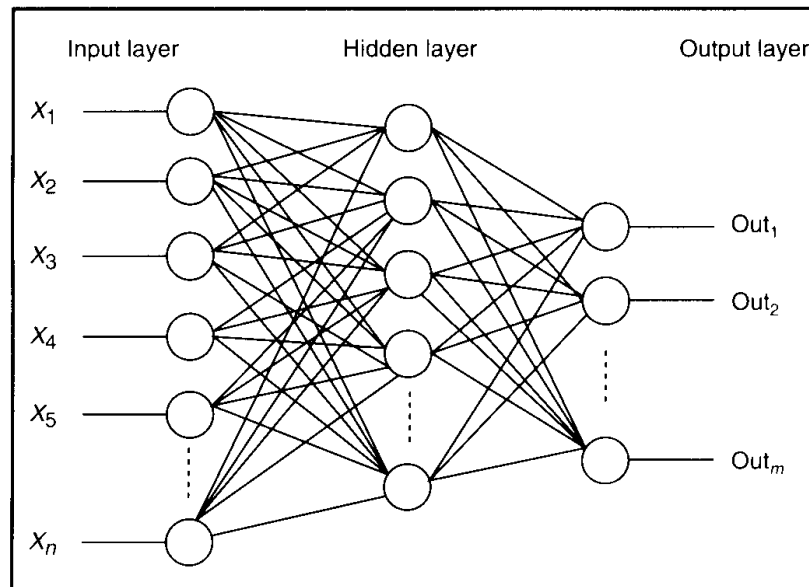
$$\left[ \theta(s) = \frac{1}{1 + e^{-s}} \right]$$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{\ln \left( 1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right)}_{e(h(\mathbf{x}_n), y_n)}$$

“cross-entropy” error

# Example: Multi output units

3



**0**  
**1**  
...  
**9**

$$h(x) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$h(x) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$h(x) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$