

## COSC 522 – Machine Learning

### Backpropagation (BP) and Multi-Layer Perceptron (MLP)

Hairong Qi, Gonzalez Family Professor  
Electrical Engineering and Computer Science  
University of Tennessee, Knoxville  
<https://www.eecs.utk.edu/people/hairong-qi/>  
Email: hqi@utk.edu

# Questions

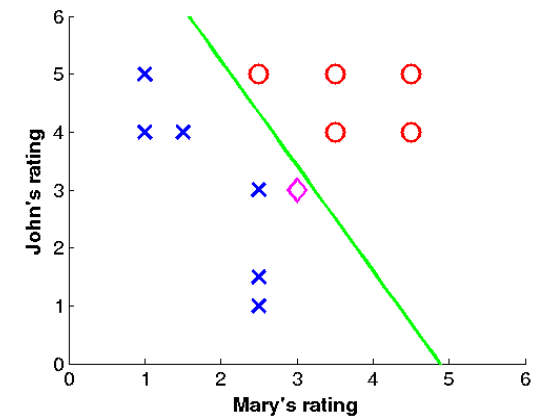
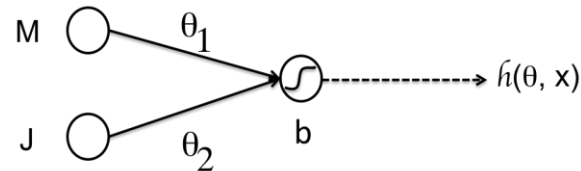
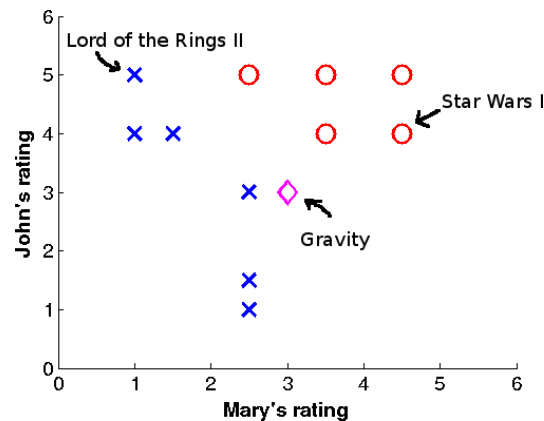
- Limitations of perceptron
- Why go deeper?
- MLP structure
- MLP cost function and optimization method (BP)
- The importance of the threshold function
- Relationship between BPNN and MPP
- Various aspects of practical improvements of BPNN

# Limitations of Perceptron

- The output only has two values (1 or 0)
- Can only classify samples which are linearly separable (straight line or straight plane)
- Single layer: can only train AND, OR, NOT
- Can't train a network functions like XOR

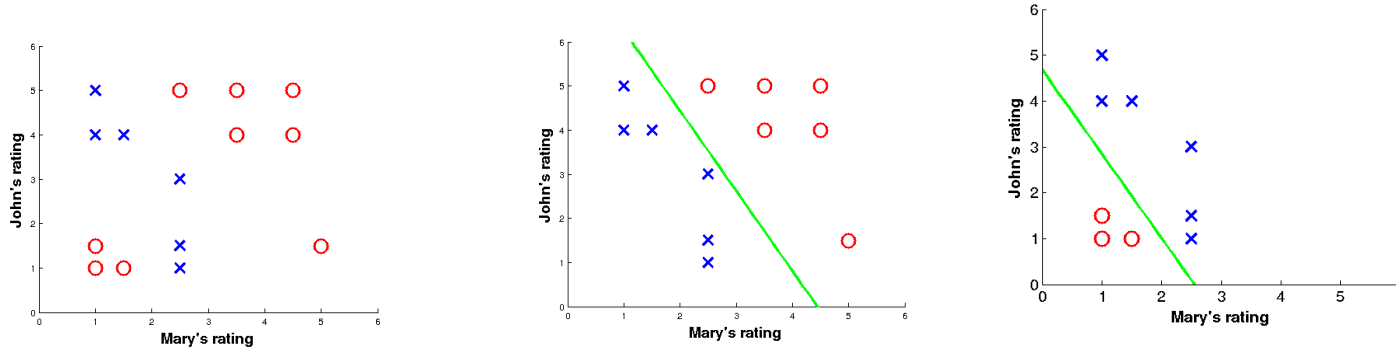
# Why deeper?

Movie name	Mary's rating	John's rating	I like?
Lord of the Rings II	1	5	No
...	...	...	...
Star Wars I	4.5	4	Yes
Gravity	3	3	?

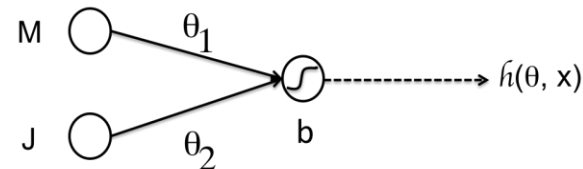
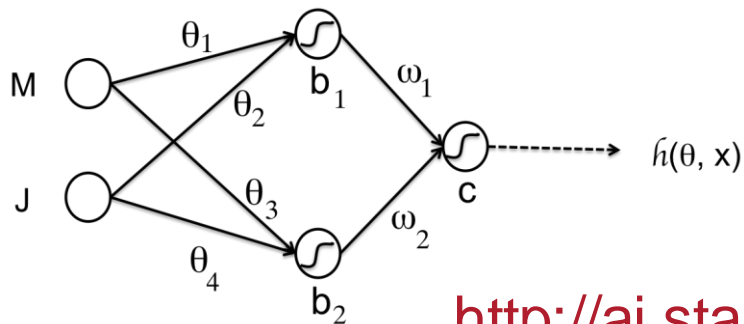


<http://ai.stanford.edu/~quocle/tutorial2.pdf>

# Why deeper?

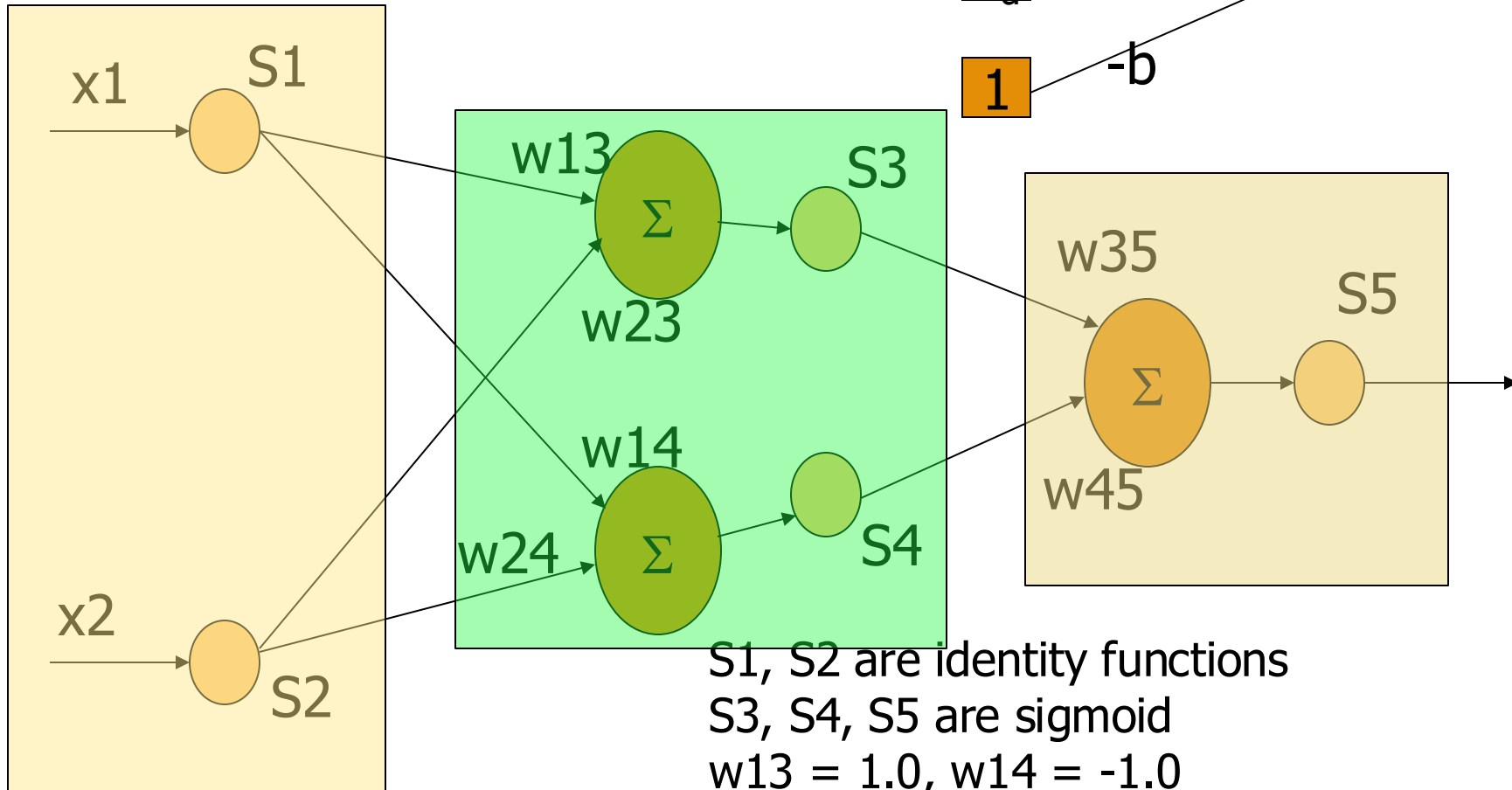
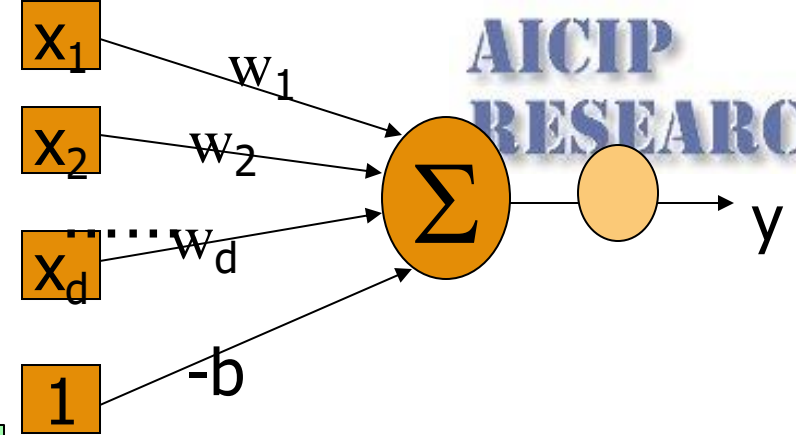


Movie name	Output by decision function $h_1$	Output by decision function $h_2$	Susan likes?
Lord of the Rings II	$h_1(x^{(1)})$	$h_2(x^{(2)})$	No
...	...	...	...
Star Wars I	$h_1(x^{(n)})$	$h_2(x^{(n)})$	Yes
Gravity	$h_1(x^{(n+1)})$	$h_2(x^{(n+1)})$	?



<http://ai.stanford.edu/~quocle/tutorial2.pdf>

# XOR (3-layer NN)



S1, S2 are identity functions

S3, S4, S5 are sigmoid

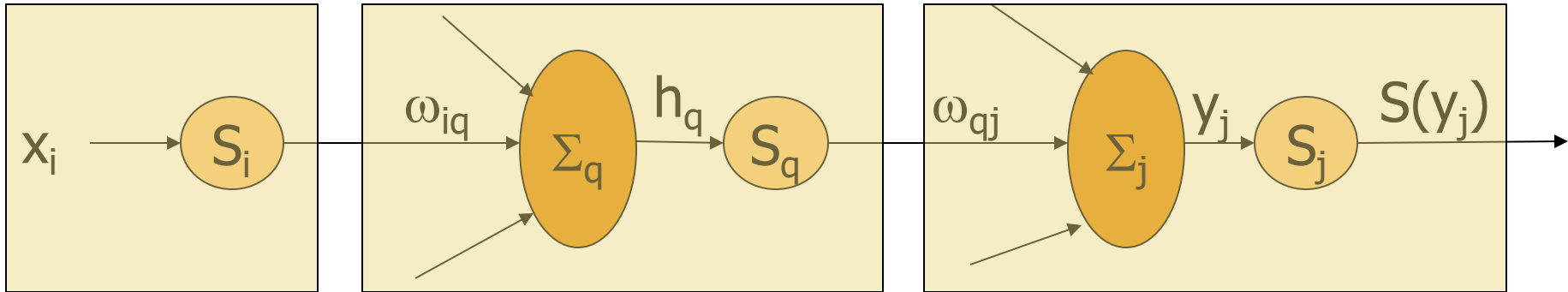
$w_{13} = 1.0$ ,  $w_{14} = -1.0$

$w_{24} = 1.0$ ,  $w_{23} = -1.0$

$w_{35} = 0.11$ ,  $w_{45} = -0.1$

The input takes on only  $-1$  and  $1$

# MLP – 3-Layer Network



$$E = \frac{1}{2} \sum_j \left( T_j - S(y_j) \right)^2$$

Choose a set of initial  $W_{st}$

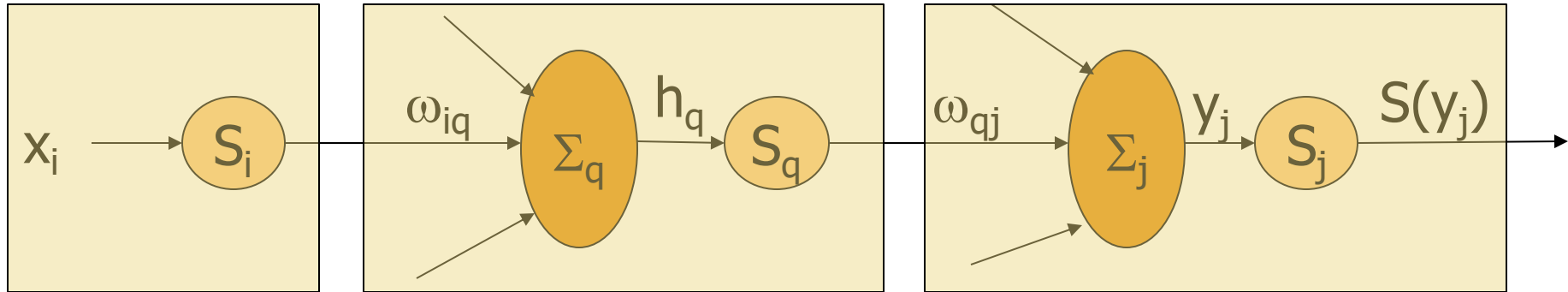
$$W_{st}^{k+1} = W_{st}^k - c^k \frac{\partial E^k}{\partial W_{st}^k}$$

$\omega_{st}$  is the weight connecting input  $s$  at neuron  $t$

The problem is essentially “how to choose weight  $\omega$  to minimize the error between the expected output and the actual output”

The basic idea behind BP is **gradient descent**

# Exercise

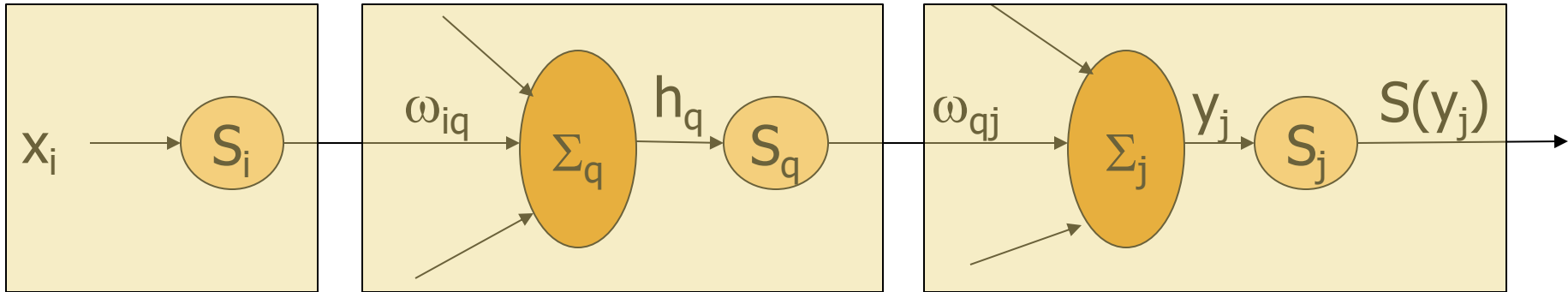


$$y_j = \mathring{a}_q s_q(h_q) w_{qj} \quad \mathbb{P} \quad \frac{\mathbb{P} y_j}{\mathbb{P} S_q} = w_{qj} \quad \text{and} \quad \frac{\mathbb{P} y_j}{\mathbb{P} w_{qj}} = s_q(h_q)$$

$$h_q = \mathring{a}_i x_i w_{iq} \quad \mathbb{P} \quad \frac{\mathbb{P} h_q}{\mathbb{P} x_i} = w_{iq} \quad \text{and} \quad \frac{\mathbb{P} h_q}{\mathbb{P} w_{iq}} = x_i$$



# The Derivative – Chain Rule



$$D W_{qj} = - \frac{\partial E}{\partial W_{qj}} = - \frac{\partial E}{\partial S_j} \frac{\partial S_j}{\partial y_j} \frac{\partial y_j}{\partial W_{qj}}$$

$$= - (T_j - S_j) (S'_j) (S_q (h_q))$$

$$D W_{iq} = - \frac{\partial E}{\partial W_{iq}} = \left[ \sum_j \frac{\partial E}{\partial S_j} \frac{\partial S_j}{\partial y_j} \frac{\partial y_j}{\partial S_q} \right] \frac{\partial S_q}{\partial h_q} \frac{\partial h_q}{\partial W_{iq}}$$

$$= \left[ \sum_j (T_j - S_j) (S'_j) (W_{qj}) \right] (S'_q) (x_i)$$

# Threshold Function

- ◆ Traditional threshold function as proposed by McCulloch-Pitts is binary function
- ◆ The importance of differentiable
- ◆ A threshold-like but differentiable form for  $S$  (25 years)
- ◆ The sigmoid

$$S(x) = \frac{1}{1 + \exp(-x)}$$

# BP vs. MPP

$$\begin{aligned}
 E(\omega) &= \sum_{\mathbf{x}} [g_k(\mathbf{x}; \mathbf{w}) - T_k]^2 = \sum_{\mathbf{x} \in \omega_k} [g_k(\mathbf{x}; \mathbf{w}) - 1]^2 + \sum_{\mathbf{x} \notin \omega_k} [g_k(\mathbf{x}; \mathbf{w}) - 0]^2 \\
 &= n \left\{ \frac{n_k}{n} \frac{1}{n_k} \sum_{\mathbf{x} \in \omega_k} [g_k(\mathbf{x}; \mathbf{w}) - 1]^2 + \frac{n - n_k}{n} \frac{1}{n - n_k} \sum_{\mathbf{x} \notin \omega_k} [g_k(\mathbf{x}; \mathbf{w})]^2 \right\}
 \end{aligned}$$

$$\begin{aligned}
 \lim_{n \rightarrow \infty} \frac{1}{n} E(\mathbf{w}) &= P(\omega_k) \int [g_k(\mathbf{x}; \mathbf{w}) - 1]^2 p(\mathbf{x} | \omega_k) d\mathbf{x} + P(\omega_{i \neq k}) \int g_k^2(\mathbf{x}; \mathbf{w}) p(\mathbf{x} | \omega_{i \neq k}) d\mathbf{x} \\
 &= \int [g_k^2(\mathbf{x}; \mathbf{w}) - 2g_k(\mathbf{x}; \mathbf{w}) + 1] p(\mathbf{x}, \omega_k) d\mathbf{x} + \int g_k^2(\mathbf{x}; \mathbf{w}) p(\mathbf{x}, \omega_{i \neq k}) d\mathbf{x} \\
 &= \int g_k^2(\mathbf{x}; \mathbf{w}) p(\mathbf{x}) d\mathbf{x} - 2 \int g_k(\mathbf{x}; \mathbf{w}) p(\mathbf{x}, \omega_k) d\mathbf{x} + \int p(\mathbf{x}, \omega_k) d\mathbf{x} \\
 &= \int [g_k(\mathbf{x}; \mathbf{w}) - P(\omega_k | \mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x} + C
 \end{aligned}$$