

Artificial Neural Networks

Physical Address for Prof. Ilias Tagkopoulos

Computer Science:

Office: 3063 Kemper Hall

Phone: (530) 752-4821

Fax: (530) 752-4767

Instructor: Ilias Tagkopoulos

iliast@ucdavis.edu

Genome and Biomedical Sciences Facility:

Office: 5313 GBSF

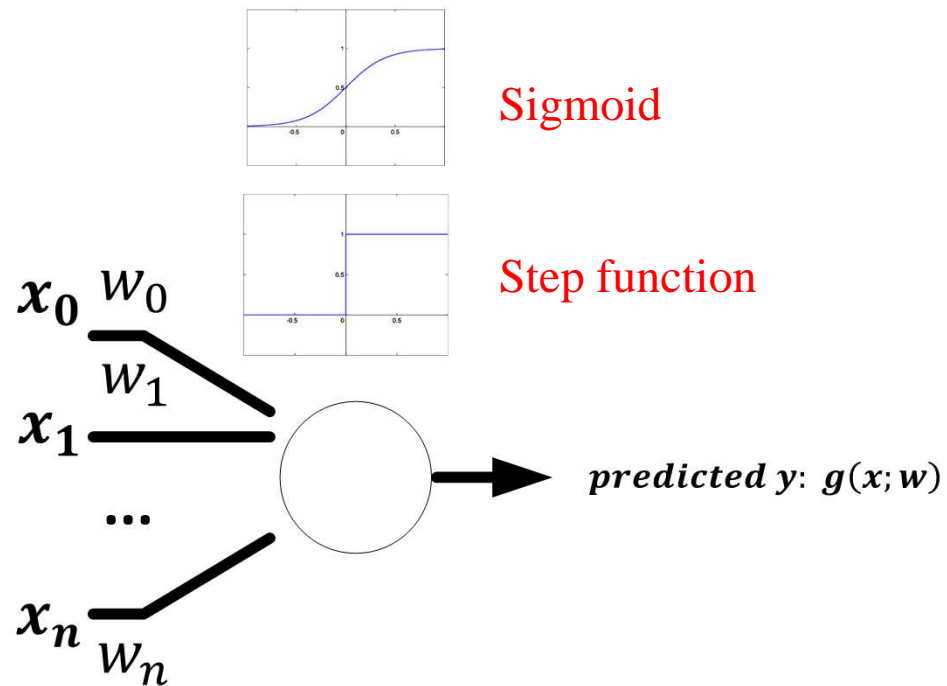
Phone: (530) 752-7707

Fax: (530) 754-9658

■ Building Boolean Functions

Lecture	Date	Topic	Comments
1	9/24/2015	Introduction	HW1 posted
2	9/29/2015	Linear Regression	
3	10/1/2015	Other Regression methods	
4	10/6/2015	Classification	HW1 due - HW2 posted
5	10/8/2015	Artificial Neural Networks	Project Topics
6	10/13/2015	Artificial Neural Networks	
7	10/15/2015	Support Vector Machines	Projects Assigned
8	10/20/2015	Support Vector Machines	
9	10/22/2015	Support Vector Machines	HW2 due - HW3 posted
10	10/27/2015	Midterm	
11	10/29/2015	Classification issues: Kernels, Overfitting, Regularization	
12	11/3/2015	Dimensionality Reduction	
13	11/5/2015	Reinforcement Learning	
14	11/10/2015	Decision support: Markov Decision Processes	
15	11/12/2015	Graphical Models - Naïve Bayes	
16	11/17/2015	Clustering: K-means - Hierarchical	HW3 due - HW4 posted
17	11/19/2015	Special topics: Deep Learning	
18	11/24/2015	Project Presentation I	Project Reports Due
19	11/26/2015	NO CLASS (Thanksgiving)	HW4 due
20	12/1/2015	Project Presentation II	
21	12/3/2015	Project Presentation III - Overview	

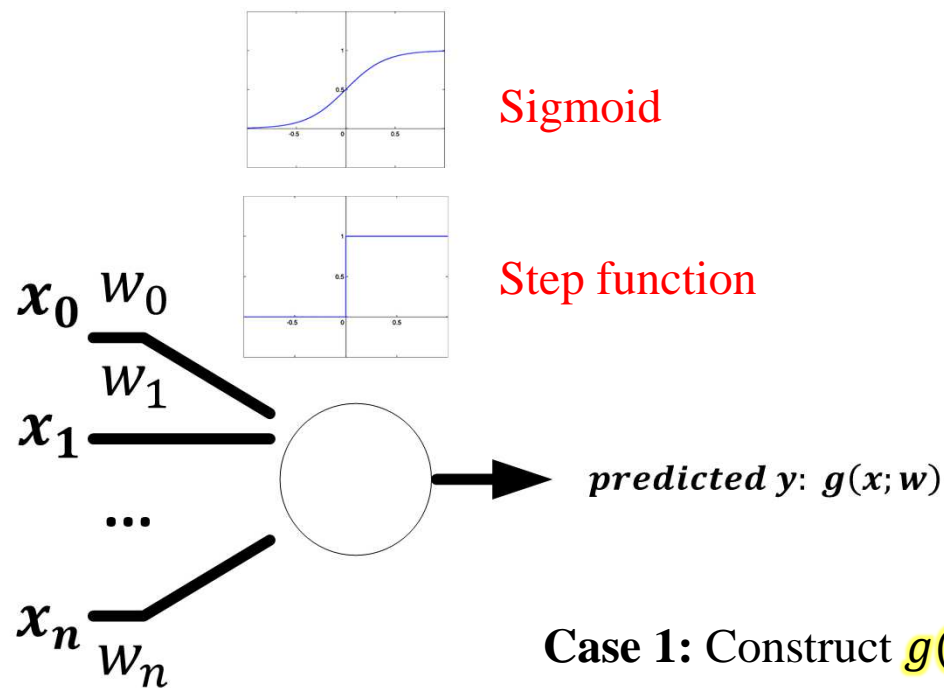
■ Building Boolean Functions



Picking the **right weights, function and thresholds** (whenever needed), we can represent any Boolean function.

For simplicity **assume** here that we use the **step function**.

■ Building Boolean Functions



Sigmoid

Step function

Picking the **right weights, function and thresholds** (whenever needed), we can represent any Boolean function.

For simplicity **assume** here that we use the **step function**.

$$g(x; w) = \begin{cases} 0 & \text{if } w^T x < 0 \\ 1 & \text{if } w^T x \geq 0 \end{cases}$$

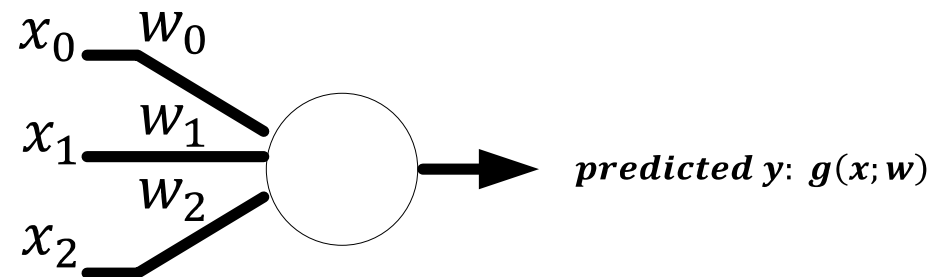
Case 1: Construct $g(x; w) = x_1 \bar{x}_2$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1 \\ -1 \\ \vdots \\ 0 \end{bmatrix}$$

Note: this solution is **not unique**, e.g.

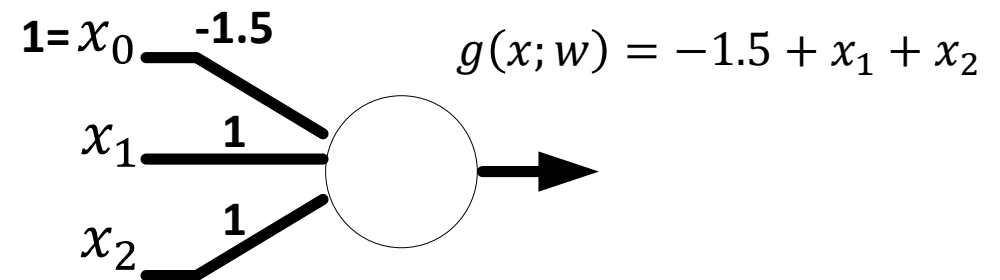
$$w = \begin{bmatrix} 105 \\ 200 \\ -100 \\ \vdots \\ 0 \end{bmatrix}$$

■ Building Boolean Functions



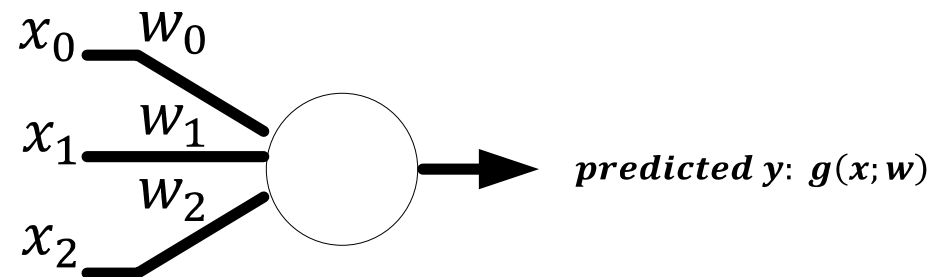
Case 2: Construct $g(x; w) = x_1 \text{ AND } x_2$

x_1	x_2	$y = g(x; w)$
0	0	0
0	1	0
1	0	0
1	1	1



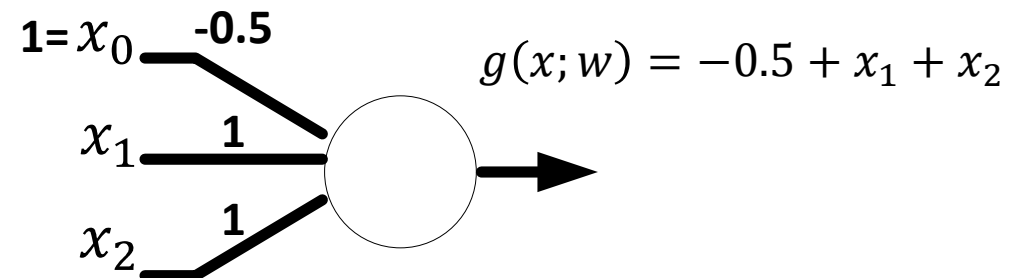
$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} -1.5 \\ 1 \\ 1 \end{bmatrix}$$

■ Building Boolean Functions



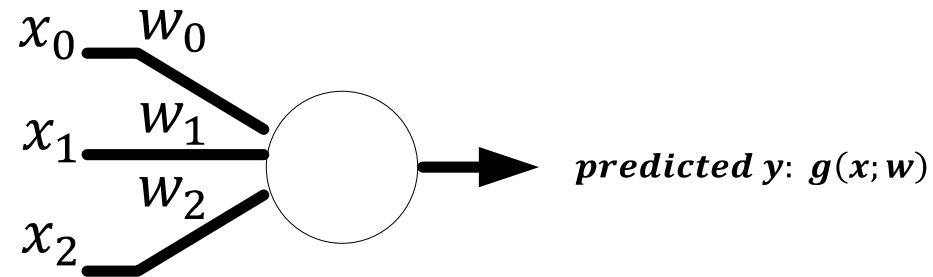
Case 3: Construct $g(x; w) = x_1 \text{ OR } x_2$

x_1	x_2	$y = g(x; w)$
0	0	0
0	1	1
1	0	1
1	1	1



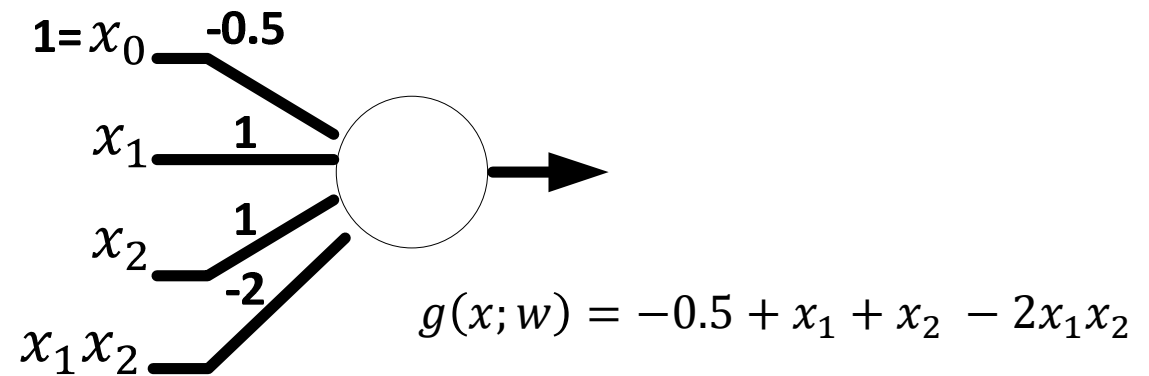
$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1 \\ 1 \end{bmatrix}$$

■ Building Boolean Functions



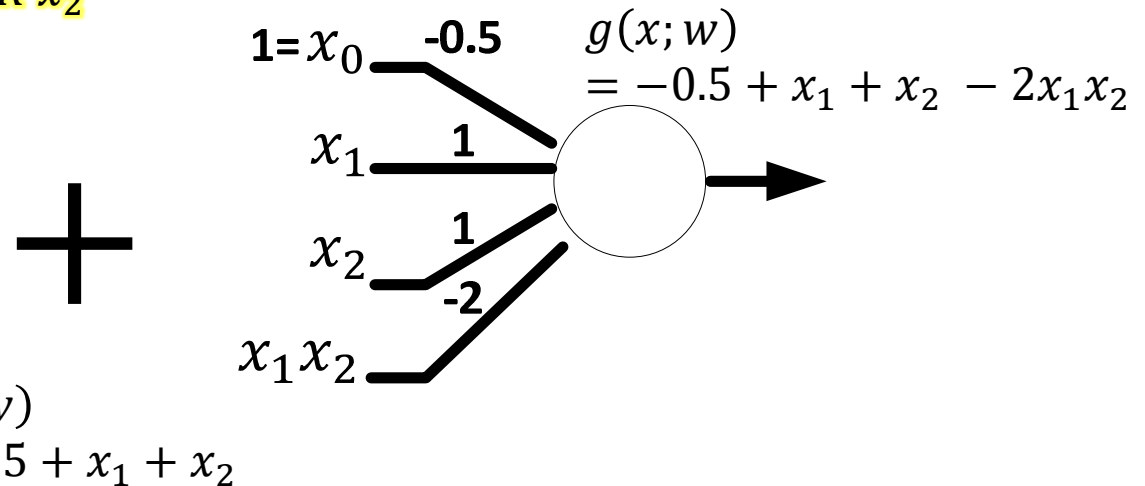
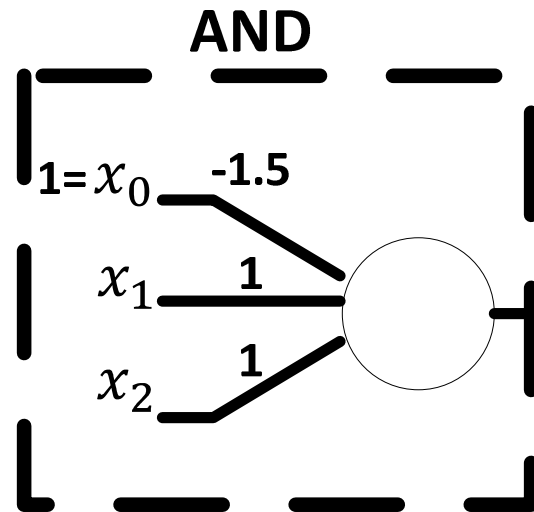
Case 3: Construct $g(x; w) = x_1 \text{ XOR } x_2$

x_1	x_2	$y = g(x; w)$
0	0	0
0	1	1
1	0	1
1	1	0

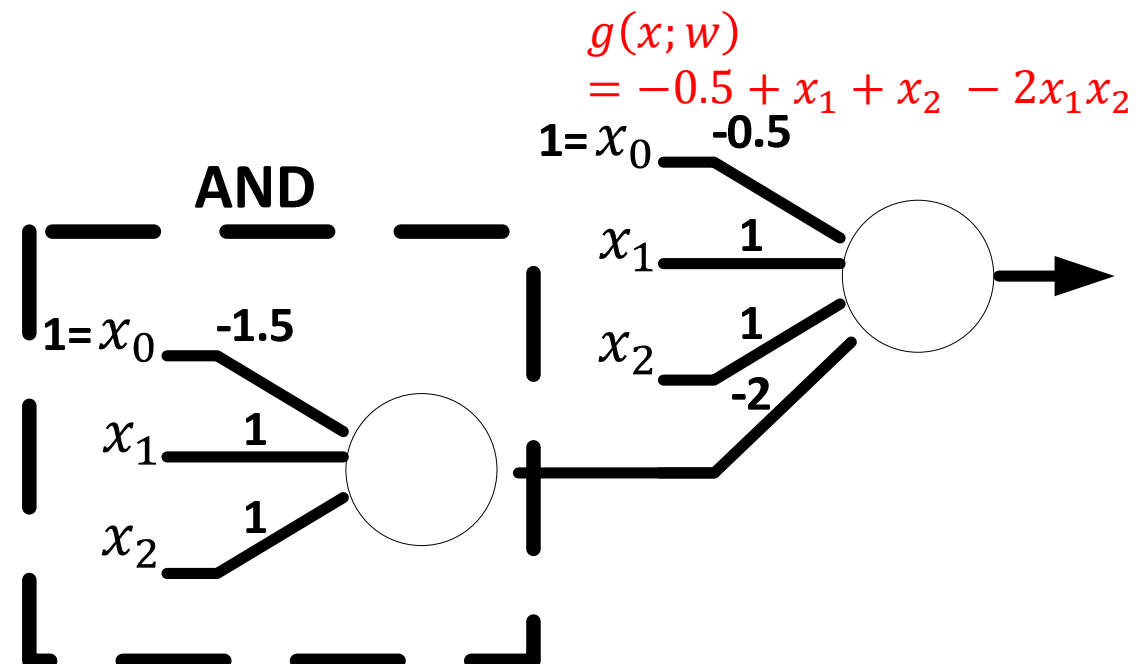


3-layer ANN encoding for an XOR

Case 3: Construct $g(x; w) = x_1 \text{ XOR } x_2$



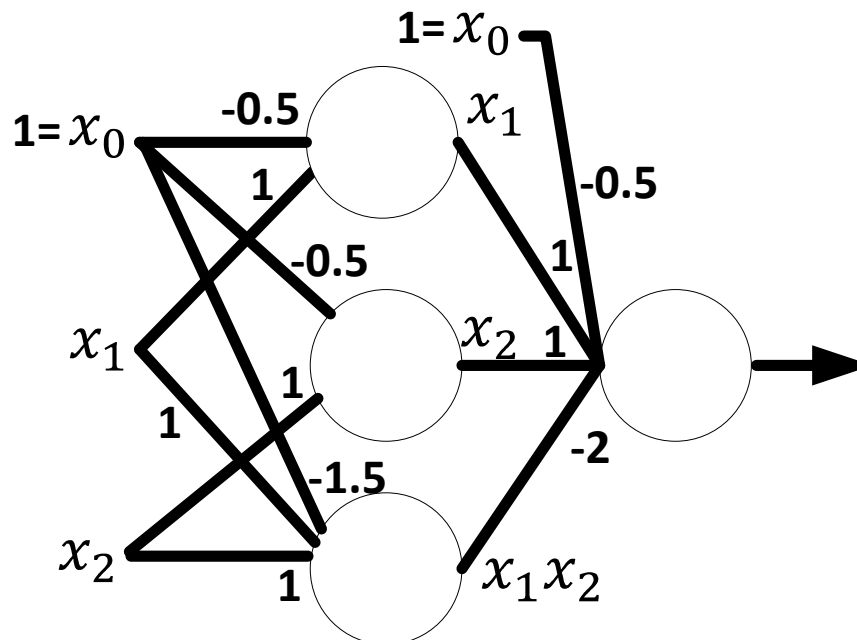
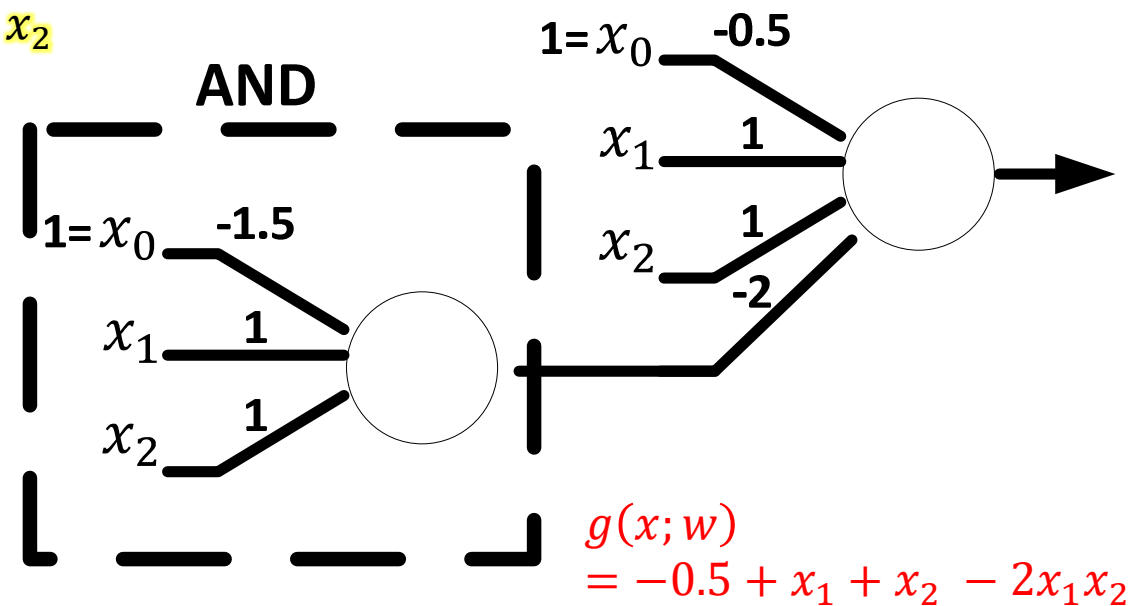
x_1	x_2	$y = g(x; w)$
0	0	0
0	1	1
1	0	1
1	1	0



3-layer ANN encoding for an XOR

Case 3: Construct $g(x; w) = x_1 \text{ XOR } x_2$

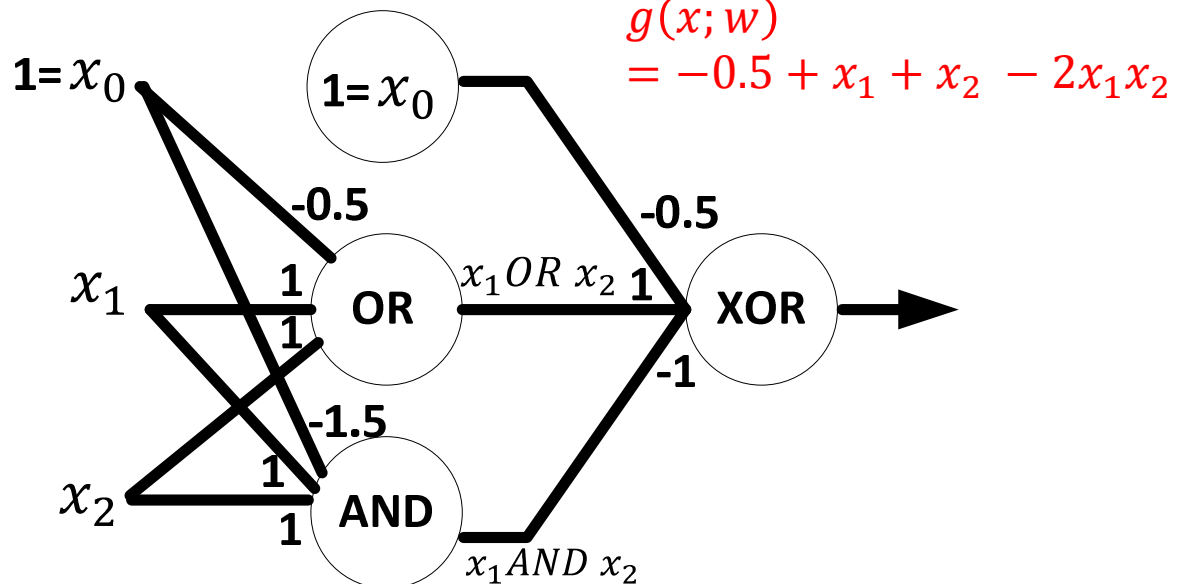
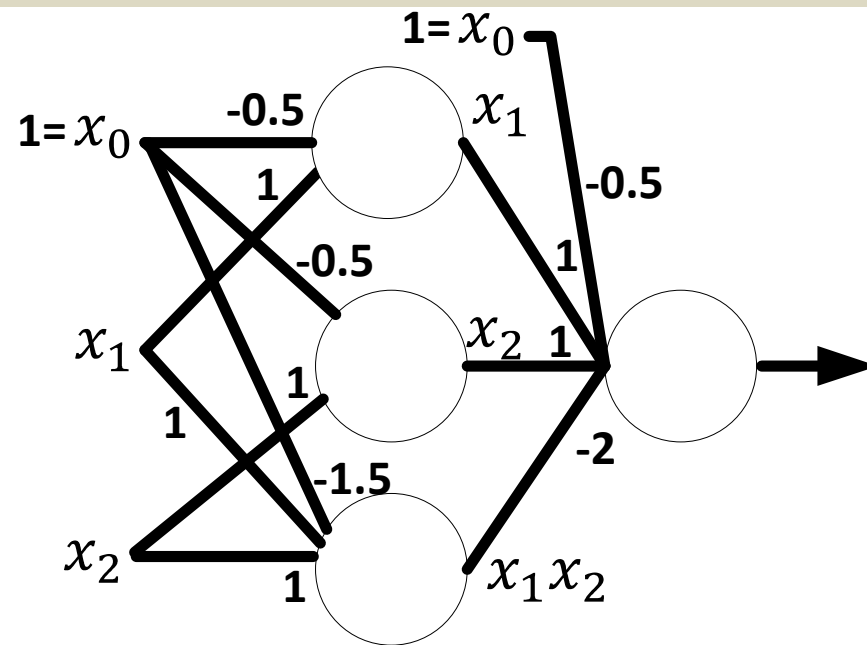
x_1	x_2	$y = g(x; w)$
0	0	0
0	1	1
1	0	1
1	1	0



■ 3-layer ANN encoding for an XOR

Case 3: Construct $g(x; w) = x_1 \text{ XOR } x_2$

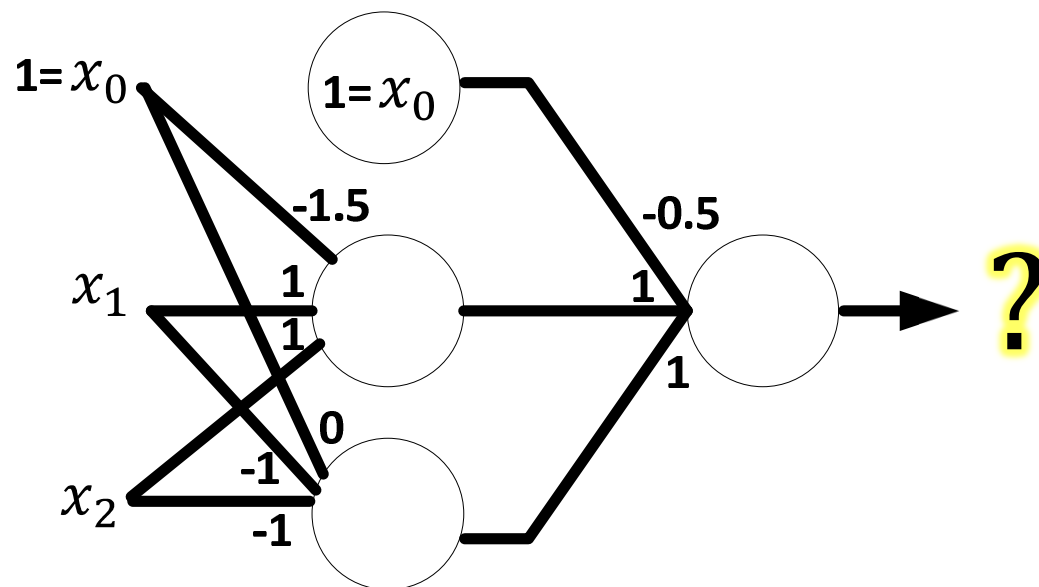
x_1	x_2	$y = g(x; w)$
0	0	0
0	1	1
1	0	1
1	1	0



■ 3-layer ANN encoding for an ?

Case 4: What is this gate? $g(x; w) = x_1 ? x_2$

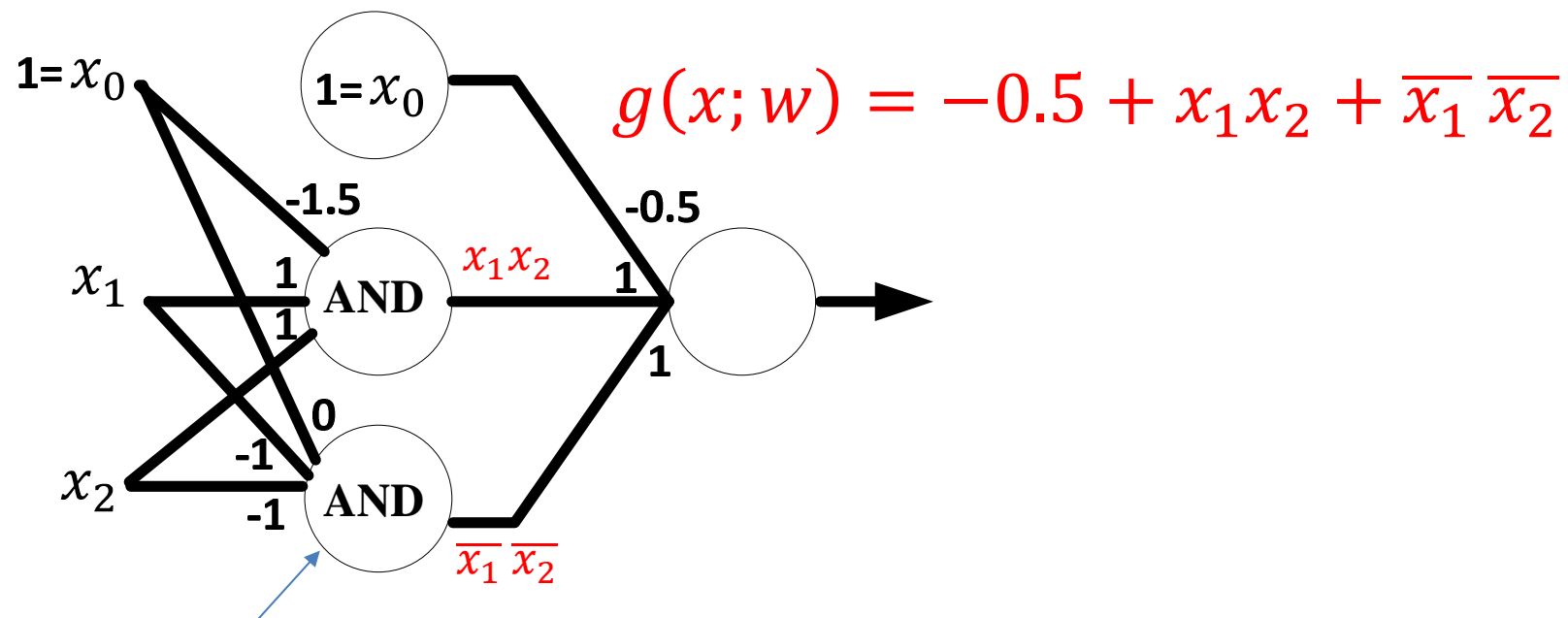
x_1	x_2	$y = g(x; w)$
0	0	?
0	1	?
1	0	?
1	1	?



■ 3-layer ANN encoding for an XOR

Case 4: What is this gate? $g(x; w) = x_1 \text{ XNOR } x_2 = \text{NOT}(x_1 \text{ XOR } x_2)$

x_1	x_2	$y = g(x; w)$
0	0	1
0	1	0
1	0	0
1	1	1



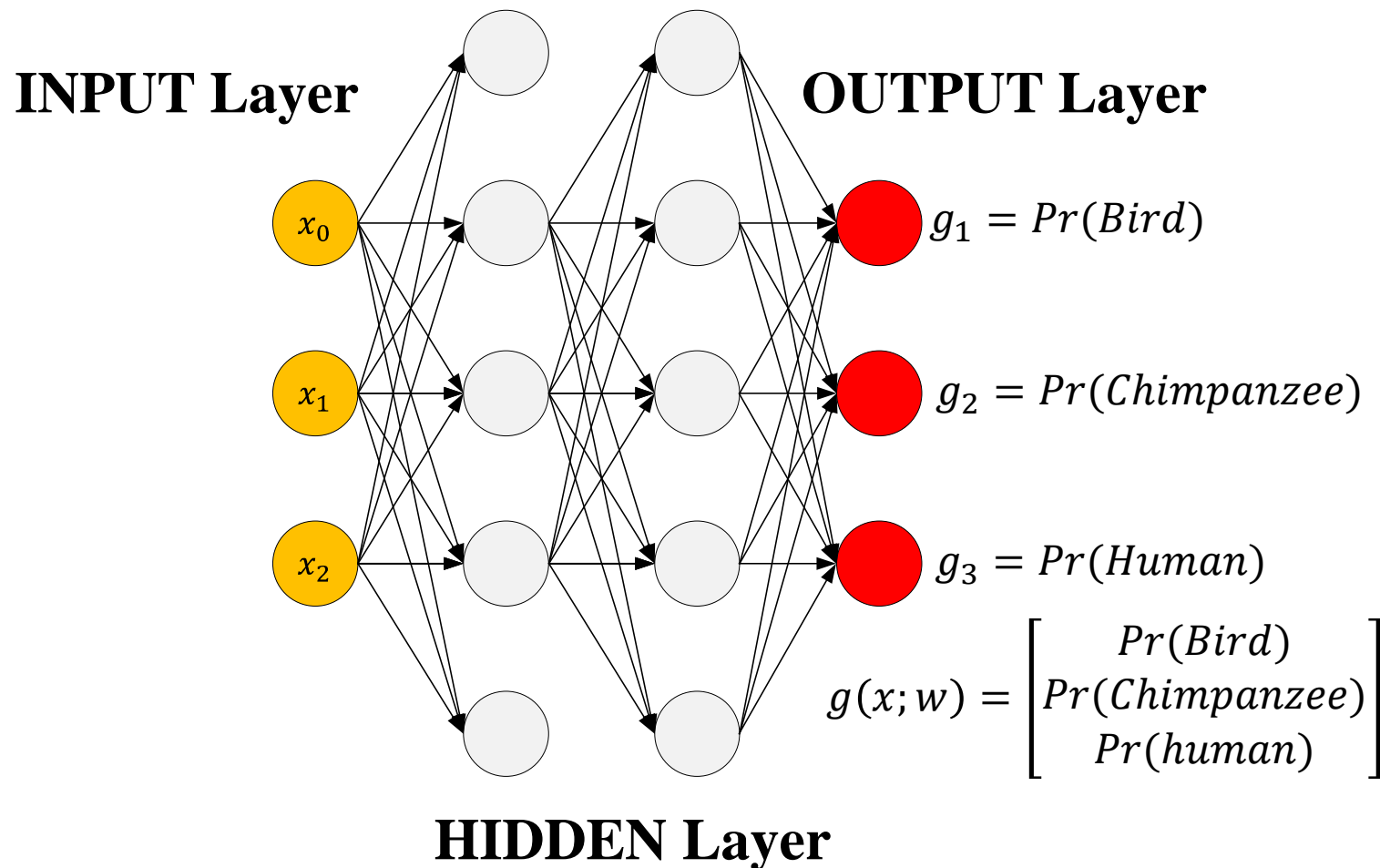
Actually a NOR, not an AND, on the original inputs



Multi-class ANN

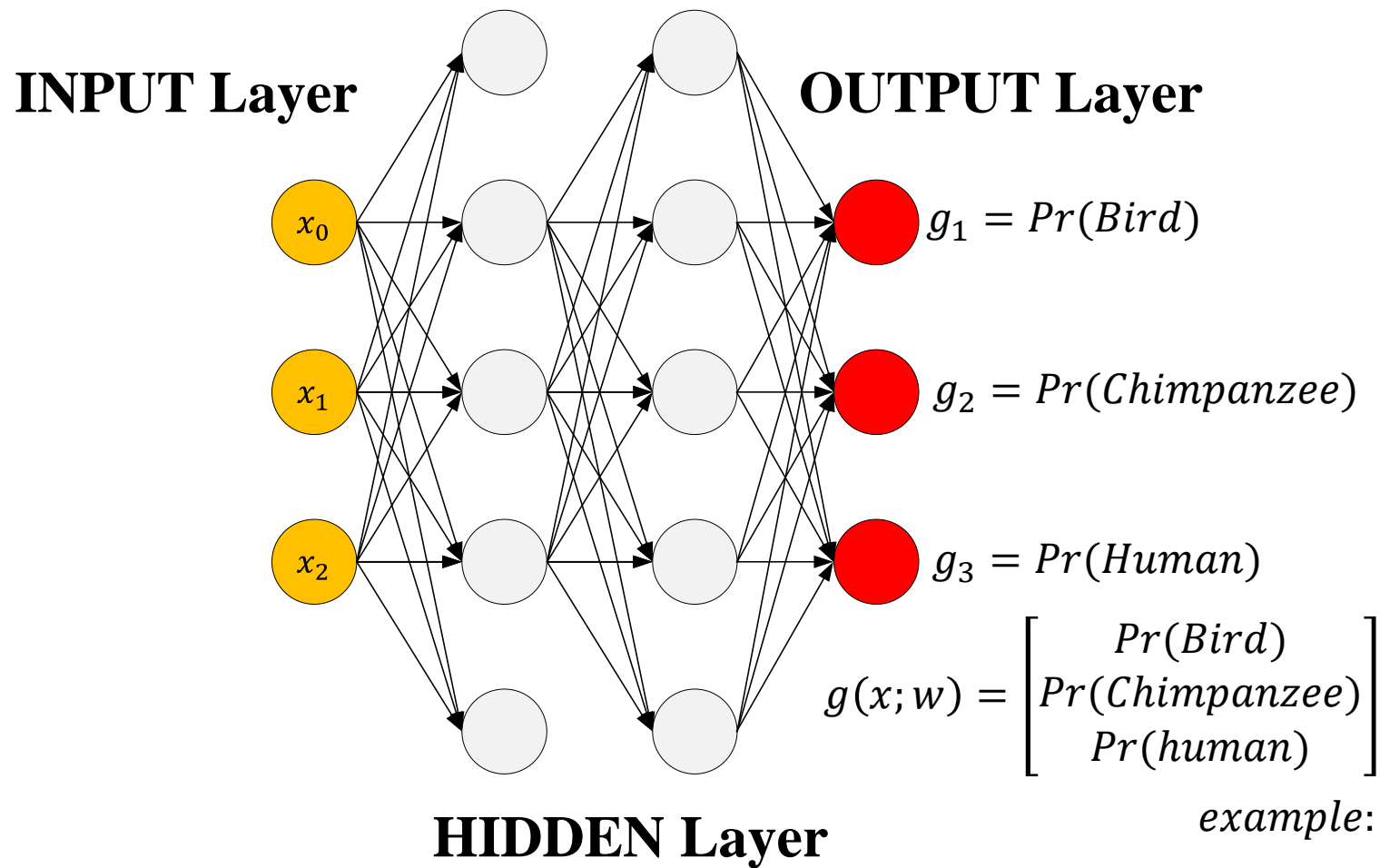
■ Multi-class classification

- We can train the ANN to **classify samples into multiple classes** (mutual exclusive or not)
- For example: **classify images** to {bird, chimpanzee, human}



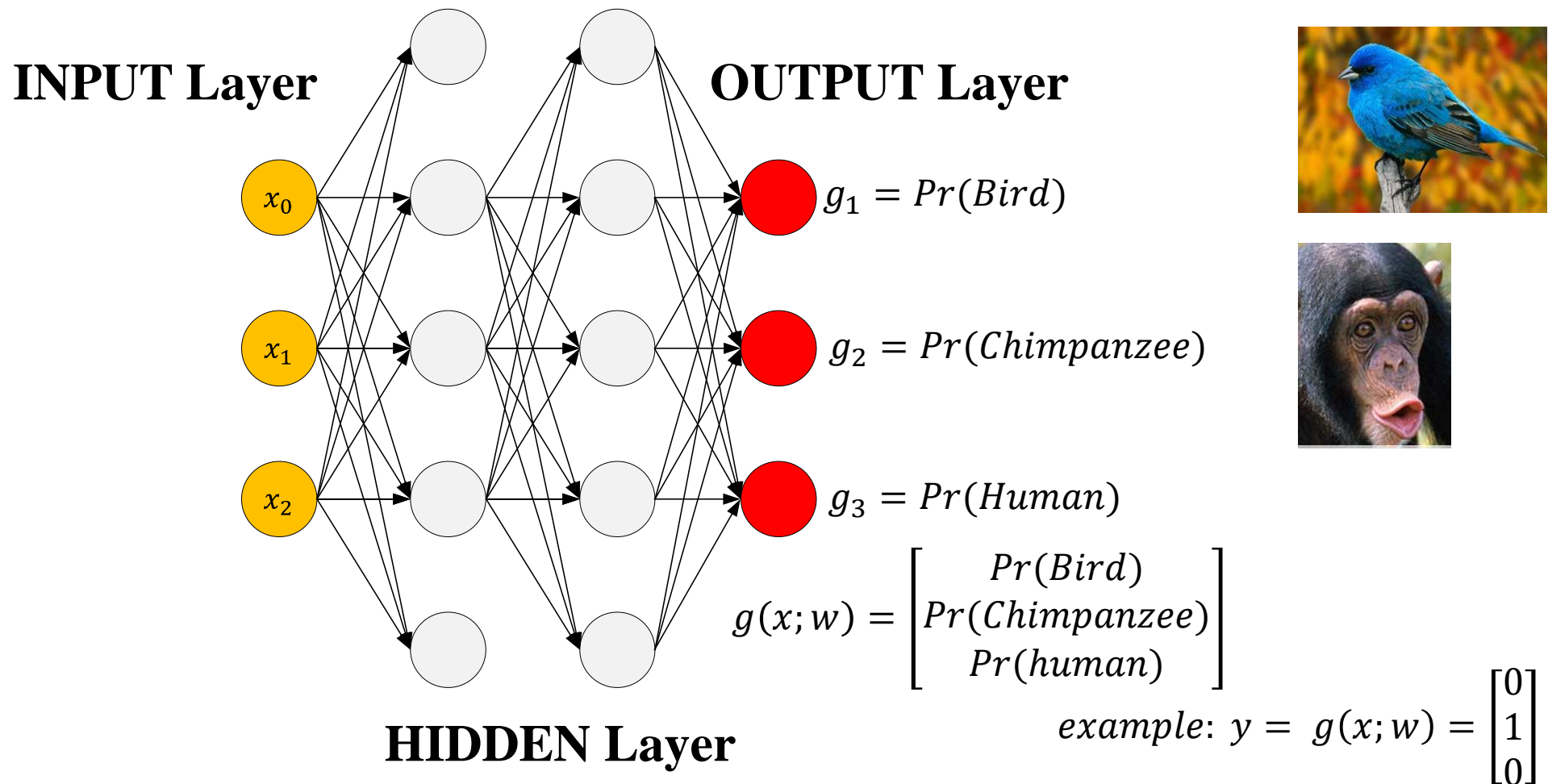
■ Multi-class classification

- We can train the ANN to classify samples into multiple classes (mutual exclusive or not)
- For example: classify images to {bird, chimpanzee, human}



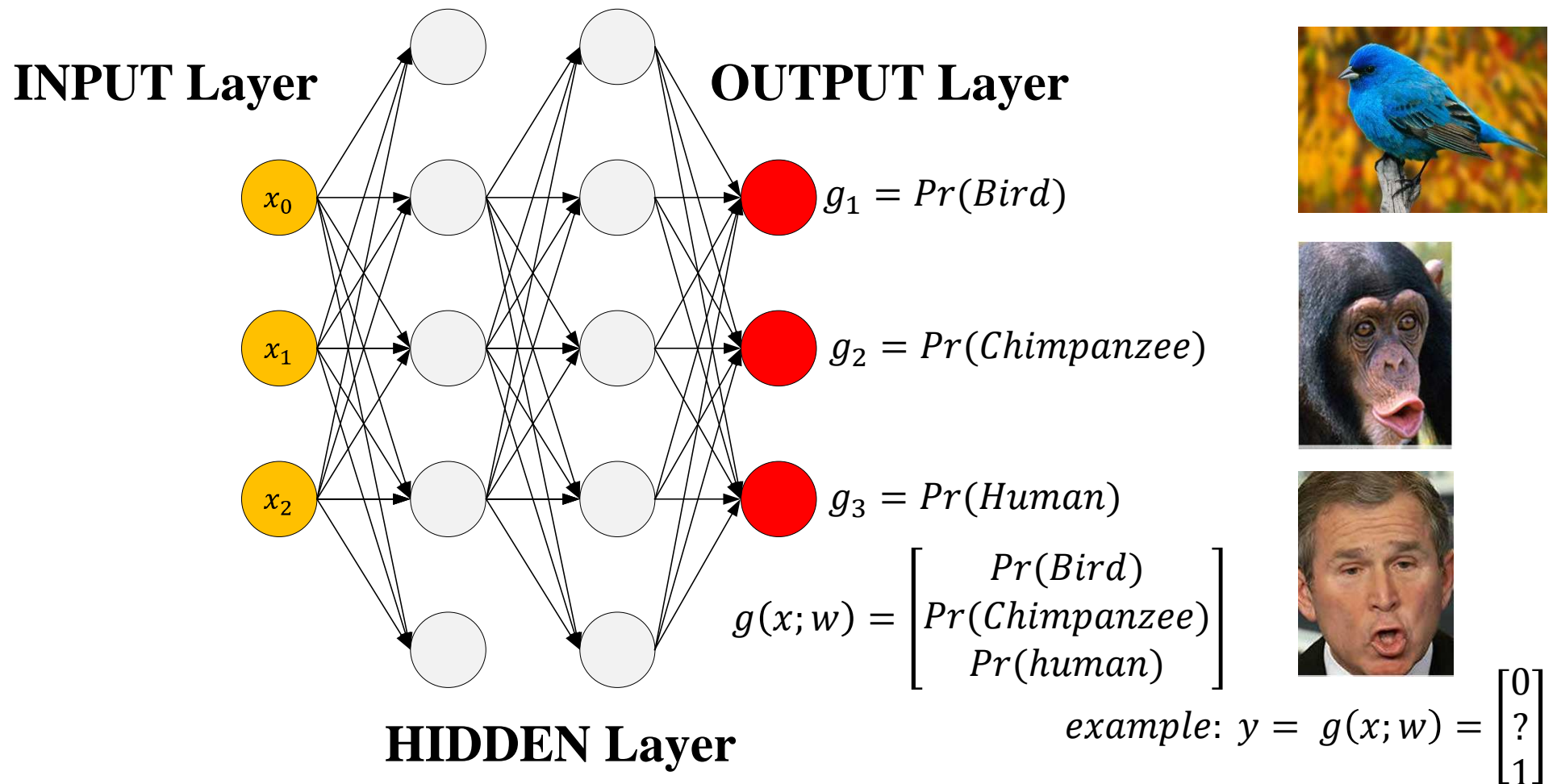
■ Multi-class classification

- We can train the ANN to classify samples into multiple classes (mutual exclusive or not)
- For example: classify images to {bird, chimpanzee, human}



■ Multi-class classification

- We can train the ANN to classify samples into multiple classes (mutual exclusive or not)
- For example: classify images to {bird, chimpanzee, human}





ANN training and backpropagation

■ Training the parameters

- If we use logistic regression as the **activation function**, then that the log likelihood for one NN node is (lecture 4, slide 11):

$$l(w) \triangleq \log p(D|w) = \sum_{i=1}^M (y^{(i)} \log g(x^{(i)}; w) + (1 - y^{(i)}) \log(1 - g(x^{(i)}; w)))$$

- So if we have **K nodes** then the log likelihood becomes:

$$l(w) = \sum_{k=1}^K \sum_{i=1}^M (y^{(i)}_k \log g(x^{(i)}; w)_k + (1 - y^{(i)}_k) \log(1 - g(x^{(i)}; w)_k))$$

- Similarly as before (logistic regression) we want to **maximize the log likelihood**
 - or conversely **minimize the negative log likelihood**, that can be thought as representing the **cost function**.

■ Back propagation

- To train the parameters, we can use any optimization method, for example **Gradient Descent** (or “ascent” in the case of the log-likelihood):

$$w_{ij} := w_{ij} + a \frac{\partial l(w)}{\partial w_{ij}}$$

with

$$l(w) = \sum_{k=1}^K \sum_{i=1}^M \left(y^{(i)}_k \log g(x^{(i)}; w)_k + (1 - y^{(i)}_k) \log (1 - g(x^{(i)}; w)_k) \right)$$

If we take the derivatives with respect to w , we see that:

$$\frac{\partial l(W)}{\partial w_{ij}^{(l-1)}} = -a_j^{(l-1)} \delta_i^{(l)}$$

Where $\delta_i^{(l)}$ is the error of the i^{th} neuron (node) in the l^{th} layer.

■ Some definitions

- To calculate the error, let's go through some definitions first:

$\alpha_i^{(l)}$ = Activation of unit i in layer l

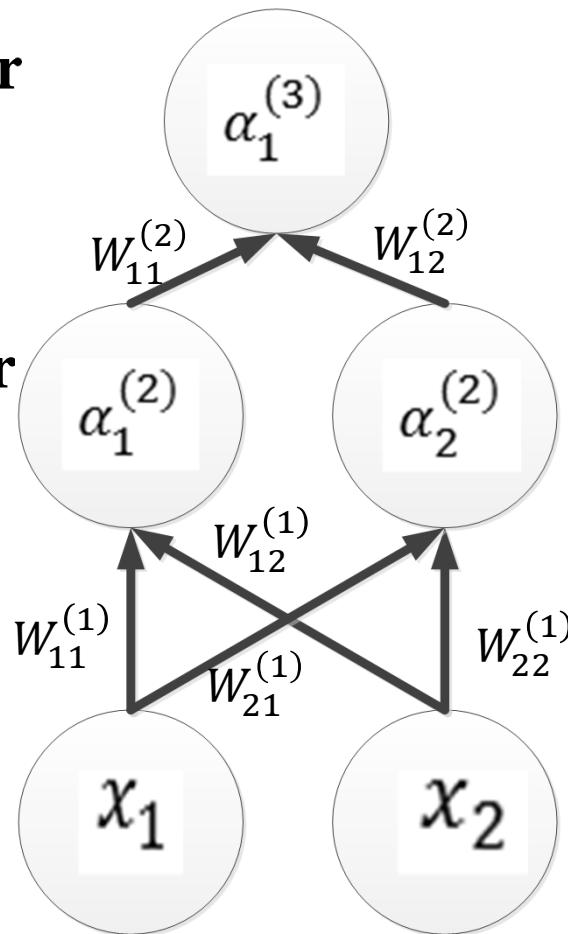
$W^{(l)}$ = weight matrix from layer l to l + 1

OUTPUT Layer

$$\alpha_1^{(3)} = g(W_{10}^{(2)}x_0 + W_{11}^{(2)}\alpha_1^{(2)} + W_{12}^{(2)}\alpha_2^{(2)})$$

HIDDEN Layer

INPUT Layer



$$\alpha_1^{(2)} = g(W_{10}^{(1)}x_0 + W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2)$$

$$\alpha_2^{(2)} = g(W_{20}^{(1)}x_0 + W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2)$$

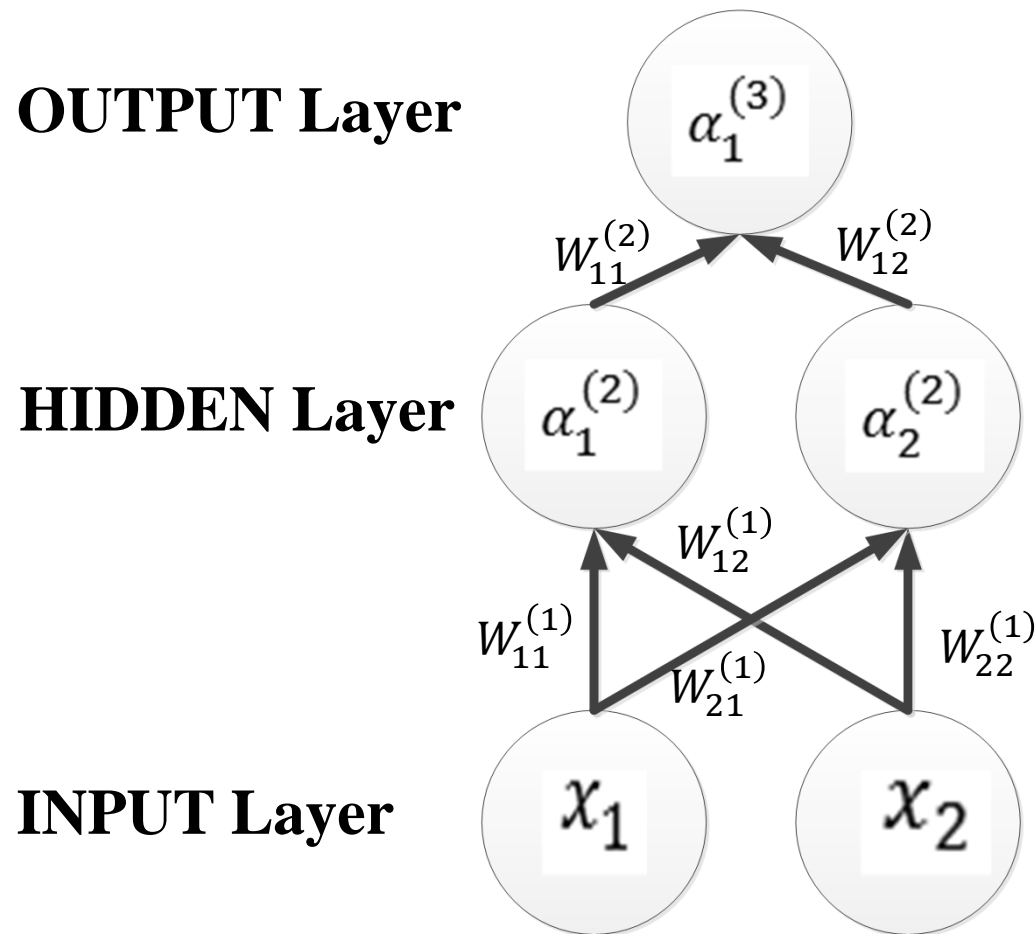
You can represent the sum in the () as **z** (since you didn't have enough variable definitions already...):

$$z_1^{(2)} = W_{10}^{(1)}x_0 + W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2$$

so that activation of unit 1 becomes:

$$\alpha_1^{(2)} = g(z_1^{(2)})$$

■ Artificial Neural Network



So output can be represented as

$$\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)})$$

With

$$\mathbf{a}^{(l)} = \begin{bmatrix} \alpha_0^{(l)} \\ \alpha_1^{(l)} \\ \alpha_2^{(l)} \end{bmatrix} \quad \mathbf{z}^{(l)} = \begin{bmatrix} z_1^{(l)} \\ z_2^{(l)} \end{bmatrix}$$

Each $\mathbf{a}^{(l)}$ can be **computed iteratively**, starting from $l = 1$ and ending at the output layer (here $l = 3$). This is called **forward-propagation**.

The final output is given by:

$$\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

■ Back-propagation

Let the error of **node i in layer l** be

$$\delta_i^{(l)} = \alpha_i^{(l)} - y_i$$

Then we can **calculate iteratively the errors for each layer** as:

$$\delta_i^{(l-1)} = (W^{(l-1)})^T \delta^{(l)} \cdot a^{(l-1)} \cdot (1 - a^{(l-1)})$$

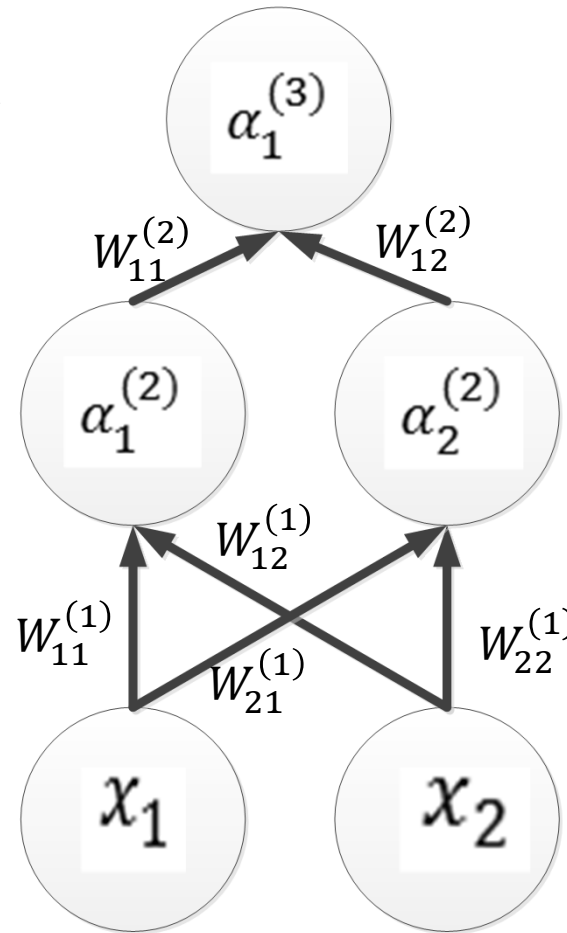
Where

$\delta^{(l)} = \begin{bmatrix} \delta_1^{(l)} \\ \vdots \\ \delta_n^{(l)} \end{bmatrix}$ is the vector of errors for the n nodes of the l^{th} layer.

OUTPUT Layer

HIDDEN Layer

INPUT Layer

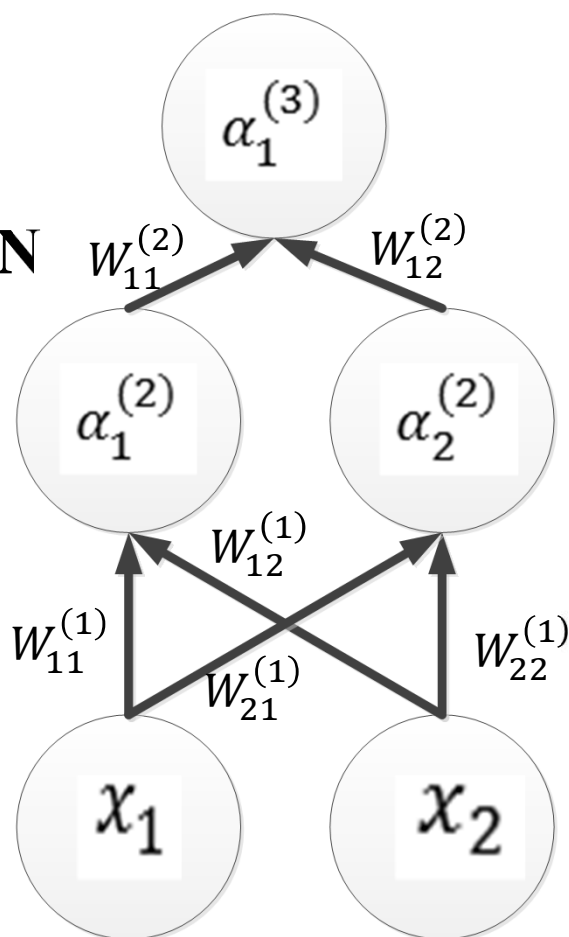


■ Back-propagation

OUTPUT Layer

HIDDEN Layer

INPUT Layer



It turns out that the **partial derivatives of the log-likelihood** with respect to the weights w are given by:

$$\frac{\partial l(W)}{\partial w_{ij}^{(l-1)}} = -a_j^{(l-1)} \delta_i^{(l)}$$

Which you now can calculate **recursively**:

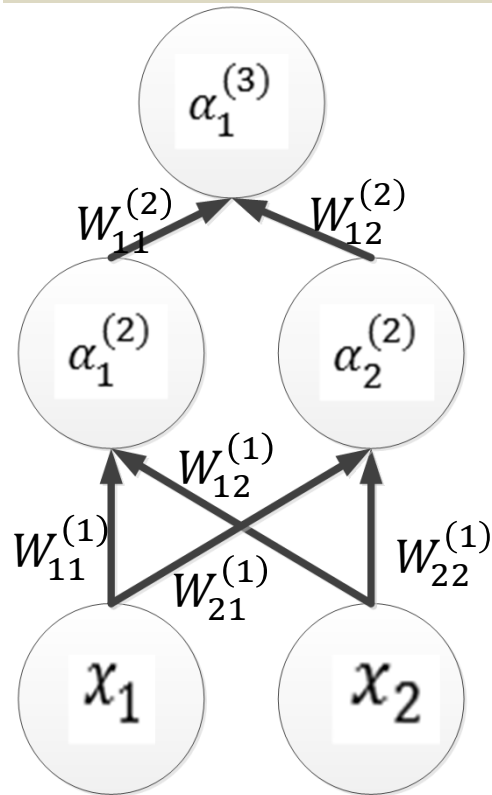
$$\delta^{(l-1)} = (W^{(l-1)})^T \delta^{(l)} \cdot a^{(l-1)} \cdot (1 - a^{(l-1)})$$

Element-by-element multiplication

Why you need the partial derivatives?

- **Gradient Descent or other optimization algorithms**

Back-propagation



For example, for our 3-layer network on the left we have:

$$\delta^{(3)} = \delta_1^{(3)} = \alpha_1^{(l)} - y_1$$

$$W^{(2)} = \begin{bmatrix} w_{10}^{(2)} \\ w_{11}^{(2)} \\ w_{12}^{(2)} \end{bmatrix} \quad W^{(1)} = \begin{bmatrix} w_{10}^{(1)} \\ w_{11}^{(1)} \\ w_{12}^{(1)} \end{bmatrix} \quad a^{(2)} = \begin{bmatrix} 1 \\ a_1^{(2)} \\ a_2^{(2)} \end{bmatrix}$$

So that for $\delta^{(2)}$ we have (following the equation from the previous slide:

$$\delta^{(2)} = \begin{bmatrix} \delta_0^{(2)} \\ \delta_1^{(2)} \\ \delta_2^{(l)} \end{bmatrix} = \begin{bmatrix} w_{10}^{(2)} \\ w_{11}^{(2)} \\ w_{12}^{(2)} \end{bmatrix}^T (\alpha_i^{(l)} - y_i)$$

Note: this is a vector/matrix representation. Follow the derivation that we did in class instead.

■ Backpropagation of a FF-ANN

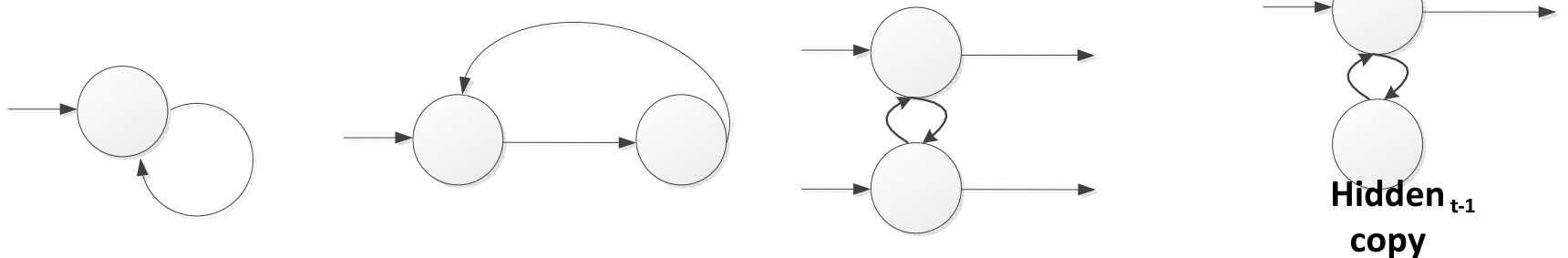
For each training sample $(x^{(k)}, y^{(k)})$:

- Set $a^{(1)} = x^{(k)}$
- Compute $a^{(l)}$ for all layers l (forward propagation)
- Compute the error in final layer $\delta_j^{(L)} = a_i^{(L)} - y^{(k)}$ and all hidden layers $\delta^{(l-1)} = (W^{(l-1)})^T \delta^{(l)} \cdot a^{(l-1)} \cdot (1 - a^{(l-1)})$
- Compute partial derivatives $\frac{\partial l(W)}{\partial w_{ij}^{(l-1)}} = -a_j^{(l-1)} \delta_i^{(l)}$
- Use the derivatives to update with a heuristic optimization method:

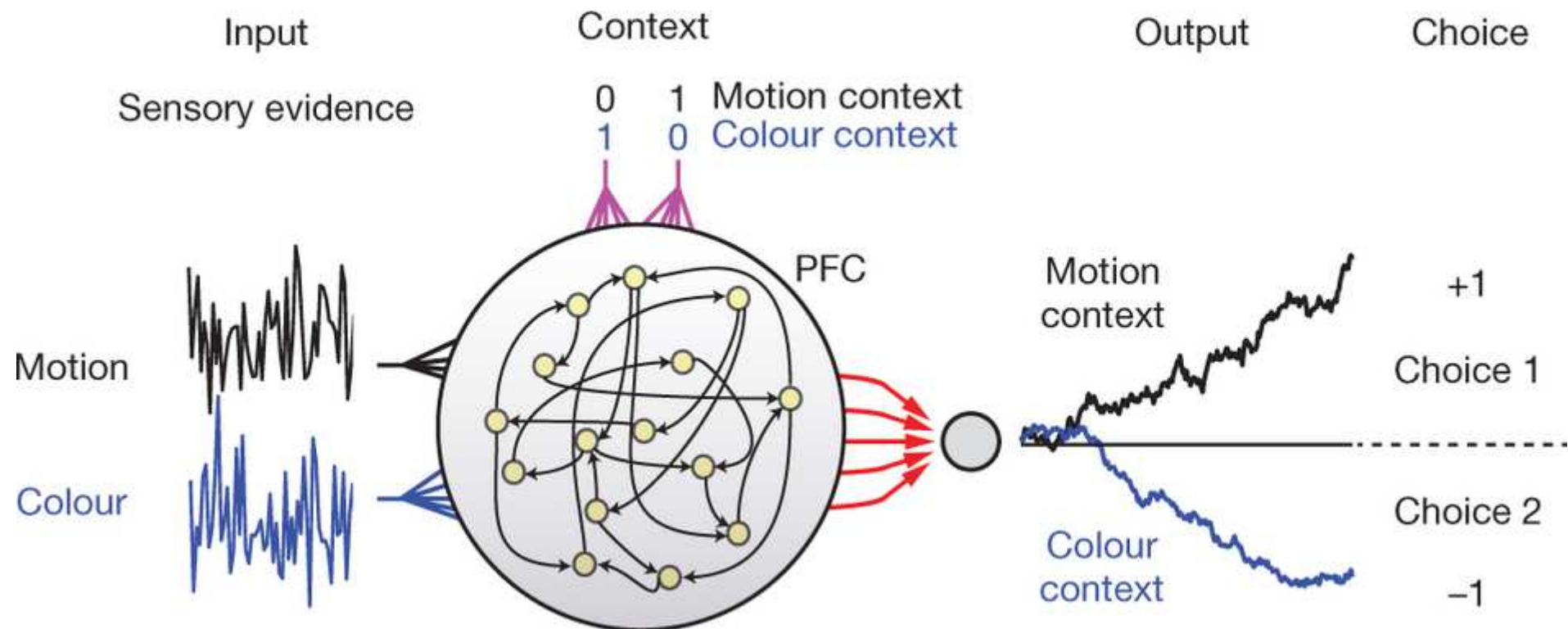
$$w_{ij}^{(l-1)} := w_{ij}^{(l-1)} - \alpha(a_j^{(l-1)} \delta_i^{(l)})$$

■ Recurrent Neural Networks

- Neural Networks can be split in two main categories:
 - **Feed-forward networks** (what we have done so far), when the ANN graph is acyclic
 - **Recurrent networks**, when the ANN graph contains cycles (feedback)
- There are many **other types** of networks that can fall into these categories (e.g. Radial Basis Function Networks, Hopfield Networks, long-short memory etc.)
- How can we train Recurrent networks?
 - Elman network, Backpropagation through time, real-time recurrent learning, etc.



Recurrent Neural Networks



Context-dependent computation by recurrent dynamics in prefrontal cortex

Valerio Mante^{1,†}, David Sussillo^{2,*}, Krishna V. Shenoy^{2,3} & William T. Newsome¹

78 | NATURE | VOL 503 | 7 NOVEMBER 2013

Figure 4 | A neural network model of input selection and integration. PFC is modelled as a network of recurrently connected, nonlinear, rate neurons that receive independent motion, colour and contextual inputs. The network is fully recurrently connected, and each unit receives both motion and colour inputs as well as two inputs that indicate context. At each time step, the sensory inputs are drawn from two normal distributions, the means of which correspond to the average strengths of the motion and colour evidence on a given trial. The contextual inputs take one of two values (0 or 1), which instruct the network to discriminate either the motion or the colour input. The network is read out by a single linear read-out, corresponding to a weighted sum over the responses of all neurons (red arrows). We trained the network (with back-propagation³⁵) to make a binary choice, that is, to generate an output of +1 at the end of the stimulus presentation if the relevant evidence pointed towards choice 1, or a -1 if it pointed towards choice 2. Before training, all synaptic strengths were randomly initialized.



End of Lecture 6