

# Model Representation I

Let's examine how we will represent a hypothesis function using neural networks. At a very simple level, neurons are basically computational units that take inputs (**dendrites**) as electrical inputs (called "spikes") that are channeled to outputs (**axons**). In our model, our dendrites are like the input features  $\mathbf{x}_1 \dots \mathbf{x}_n$ , and the output is the result of our hypothesis

function. In this model our  $\mathbf{x}_0$  input node is sometimes called the "bias unit." It is always equal to 1. In neural networks, we use the same logistic function as in classification,  $\frac{1}{1+e^{-\theta^T x}}$ , yet we sometimes call it a sigmoid (logistic) **activation** function. In this situation, our "theta" parameters are sometimes called "weights".

Visually, a simplistic representation looks like:

$$\begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \rightarrow [ \quad ] \rightarrow \mathbf{h}_\theta(\mathbf{x})$$

Our input nodes (layer 1), also known as the "input layer", go into another node (layer 2), which finally outputs the hypothesis function, known as the "output layer".

We can have intermediate layers of nodes between the input and output layers called the "hidden layers."

In this example, we label these intermediate or "hidden" layer nodes  $\mathbf{a}_0^2 \dots \mathbf{a}_n^2$  and call them "activation units."

$\mathbf{a}_i^j$  = "activation" of unit i in layer j

$\Theta^j$  = matrix of weights controlling function mapping from layer j to layer j+1

If we had one hidden layer, it would look like:

$$\begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{a}_1^2 \\ \mathbf{a}_2^2 \\ \mathbf{a}_3^2 \end{bmatrix} \rightarrow \mathbf{h}_\theta(\mathbf{x})$$

The values for each of the "activation" nodes is obtained as follows:

$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_{\theta}(x) = a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

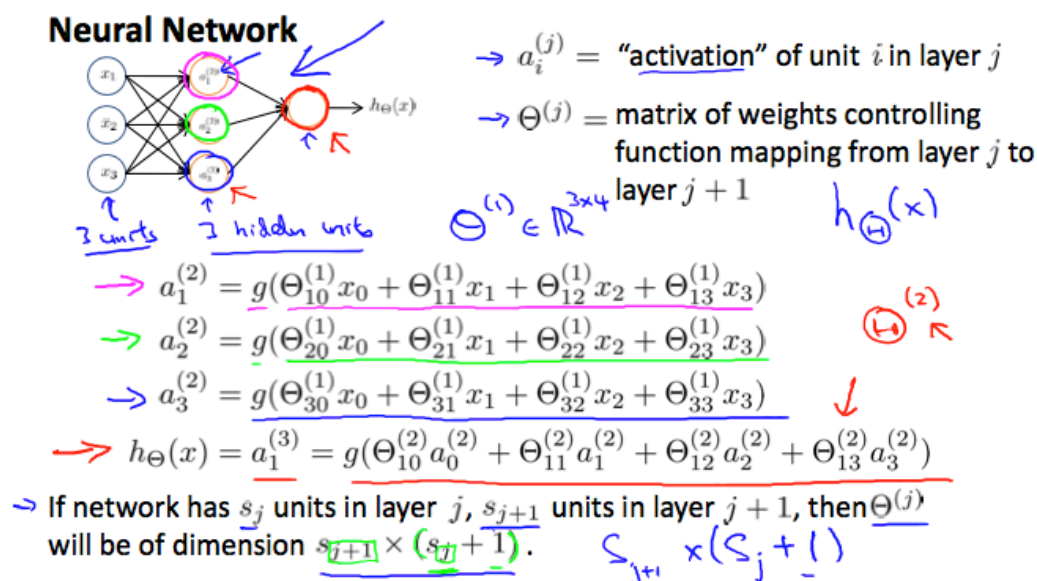
This is saying that we compute our activation nodes by using a  $3 \times 4$  matrix of parameters. We apply each row of the parameters to our inputs to obtain the value for one activation node. Our hypothesis output is the logistic function applied to the sum of the values of our activation nodes, which have been multiplied by yet another parameter matrix  $\theta^{(2)}$  containing the weights for our second layer of nodes.

Each layer gets its own matrix of weights,  $\theta^{(j)}$ .

The dimensions of these matrices of weights is determined as follows:

If network has  $s_j$  units in layer  $j$  and  $s_{j+1}$  units in layer  $j+1$ , then  $\theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .

The  $+1$  comes from the addition in  $\theta^{(j)}$  of the "bias nodes,"  $x_0$  and  $a_0^{(j)}$ . In other words the output nodes will not include the bias nodes while the inputs will. The following image summarizes our model representation:



Example: If layer 1 has 2 input nodes and layer 2 has 4 activation nodes. Dimension of  $\Theta^{(1)}$

is going to be  $4 \times 3$  where  $s_j = 2$  and  $s_{j+1} = 4$ , so  $s_{j+1} \times s_j = 4 \times 3$ .