

Model Representation II

To re-iterate, the following is an example of a neural network:

$$\mathbf{a}_1^{(2)} = g(\Theta_{10}^{(1)} \mathbf{x}_0 + \Theta_{11}^{(1)} \mathbf{x}_1 + \Theta_{12}^{(1)} \mathbf{x}_2 + \Theta_{13}^{(1)} \mathbf{x}_3)$$

$$\mathbf{a}_2^{(2)} = g(\Theta_{20}^{(1)} \mathbf{x}_0 + \Theta_{21}^{(1)} \mathbf{x}_1 + \Theta_{22}^{(1)} \mathbf{x}_2 + \Theta_{23}^{(1)} \mathbf{x}_3)$$

$$\mathbf{a}_3^{(2)} = g(\Theta_{30}^{(1)} \mathbf{x}_0 + \Theta_{31}^{(1)} \mathbf{x}_1 + \Theta_{32}^{(1)} \mathbf{x}_2 + \Theta_{33}^{(1)} \mathbf{x}_3)$$

$$\mathbf{h}_{\Theta}(\mathbf{x}) = \mathbf{a}_1^{(3)} = g(\Theta_{10}^{(2)} \mathbf{a}_0^{(2)} + \Theta_{11}^{(2)} \mathbf{a}_1^{(2)} + \Theta_{12}^{(2)} \mathbf{a}_2^{(2)} + \Theta_{13}^{(2)} \mathbf{a}_3^{(2)})$$

In this section we'll do a vectorized implementation of the above functions. We're going to define a new variable $\mathbf{z}_k^{(j)}$ that encompasses the parameters inside our g function. In our previous example if we replaced by the variable z for all the parameters we would get:

$$\mathbf{a}_1^{(2)} = g(\mathbf{z}_1^{(2)})$$

$$\mathbf{a}_2^{(2)} = g(\mathbf{z}_2^{(2)})$$

$$\mathbf{a}_3^{(2)} = g(\mathbf{z}_3^{(2)})$$

In other words, for layer $j=2$ and node k , the variable z will be:

$$\mathbf{z}_k^{(j)} = \Theta_{k,0}^{(1)} \mathbf{x}_0 + \Theta_{k,1}^{(1)} \mathbf{x}_1 + \dots + \Theta_{k,n}^{(1)} \mathbf{x}_n$$

The vector representation of \mathbf{x} and $\mathbf{z}^{(j)}$ is:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \dots \\ \mathbf{x}_n \end{bmatrix} \quad \mathbf{z}^{(j)} = \begin{bmatrix} \mathbf{z}_1^{(j)} \\ \mathbf{z}_2^{(j)} \\ \dots \\ \mathbf{z}_n^{(j)} \end{bmatrix}$$

Setting $\mathbf{x} = \mathbf{a}^{(1)}$, we can rewrite the equation as:

$$\mathbf{z}^{(j)} = \Theta^{(j-1)} \mathbf{a}^{(j-1)}$$

We are multiplying our matrix $\Theta^{(j-1)}$ with dimensions $\mathbf{s}_j \times (\mathbf{n}+1)$ (where \mathbf{s}_j is the number of our activation nodes) by our vector $\mathbf{a}^{(j-1)}$ with height $(\mathbf{n}+1)$. This gives us our vector $\mathbf{z}^{(j)}$

with height s_j . Now we can get a vector of our activation nodes for layer j as follows:

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

Where our function g can be applied element-wise to our vector $\mathbf{z}^{(j)}$.

We can then add a bias unit (equal to 1) to layer j after we have computed $\mathbf{a}^{(j)}$. This will be element $\mathbf{a}_0^{(j)}$ and will be equal to 1. To compute our final hypothesis, let's first compute another z vector:

$$\mathbf{z}^{(j+1)} = \Theta^{(j)} \mathbf{a}^{(j)}$$

We get this final z vector by multiplying the next theta matrix after $\Theta^{(j-1)}$ with the values of all the activation nodes we just got. This last theta matrix $\Theta^{(j)}$ will have only **one row** which is multiplied by one column $\mathbf{a}^{(j)}$ so that our result is a single number. We then get our final result with:

$$\mathbf{h}_{\Theta}(\mathbf{x}) = \mathbf{a}^{(j+1)} = \mathbf{g}(\mathbf{z}^{(j+1)})$$

Notice that in this **last step**, between layer j and layer j+1, we are doing **exactly the same thing** as we did in logistic regression. Adding all these intermediate layers in neural networks allows us to more elegantly produce interesting and more complex non-linear hypotheses.