



Modern Machine Learning — Neural Networks & Backpropagation

Kenneth E. Barner
Department of Electrical and Computer Engineering
University of Delaware



Artificial neural networks are biologically-inspired networks of highly interconnected (simple) processing elements (neurons)

Optimization is performed through steps analogous to adjusting synaptic connections between neurons — adaptive feedback learning

Backpropagation is utilized in conjunction with an optimization technique, such as **gradient descent**

Optimization consists of two stages: **propagation** and **weight update**

The effect of an **input is propagated forward** through the network, the error determined, and the **error is propagated back** through the network

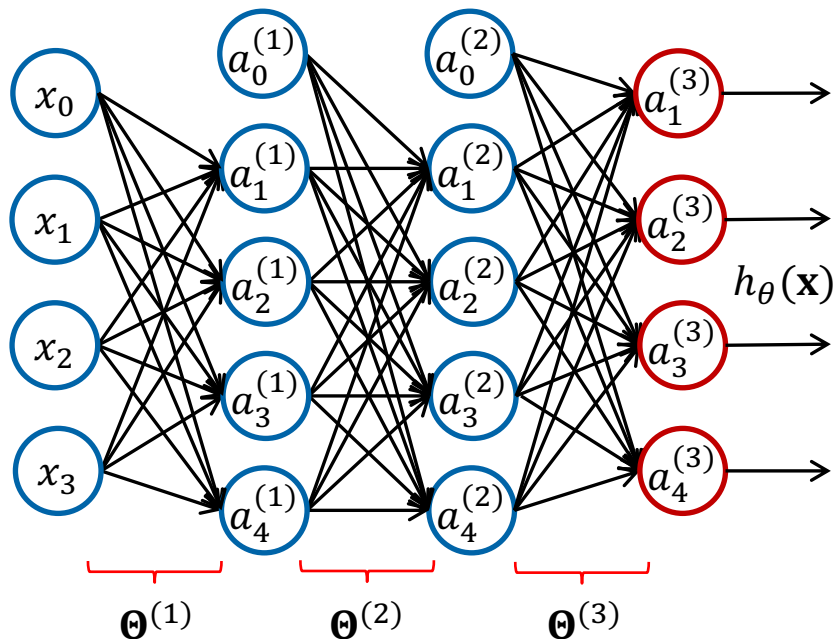
Weights are adjusted according to **gradient descent** and the back propagated error

Regularization, or shrinkage, can be applied to the optimization in order to minimize overfitting and reduce the number of non-zero weights



Methodology: Establish a cost function $J(\Theta)$ to be minimized as a function of Θ by gradient descent, called **backpropagation** in the NN setting

Observations: The gradient is easily derived using the chain rule because of the NN structure



Gradient Computation

The gradient can be determined via the chain rule by a forward and backward sweep over the network

Assume a neural network consisting of L layers, and a training example (\mathbf{x}, \mathbf{y}) consisting of observation \mathbf{x} and its corresponding label \mathbf{y}

The forward propagation equations are given by

$$\begin{aligned} \mathbf{z}^{(1)} &= \Theta^{(1)} \mathbf{x} \\ \mathbf{a}^{(1)} &= g(\mathbf{z}^{(1)}) \\ \mathbf{z}^{(2)} &= \Theta^{(2)} \mathbf{a}^{(1)} \\ \mathbf{a}^{(2)} &= g(\mathbf{z}^{(2)}) \\ &\vdots \\ \mathbf{z}^{(L)} &= \Theta^{(L)} \mathbf{a}^{(L-1)} \\ \mathbf{a}^{(L)} &= h_{\theta}(\mathbf{x}) = g(\mathbf{z}^{(L)}) \end{aligned}$$

where

$$g(\mathbf{z}) = \frac{1}{1 + e^{-z}}$$

is the logistic function



Define the cost function $J(\Theta)$ as the squared output error:

$$J(\Theta) = \frac{1}{2} \|\mathbf{y} - \mathbf{a}^{(L)}\|^2 = \frac{1}{2} \sum_i (a_i^L - y_i)^2$$

Using the **chain rule** we can express the error for layer L as

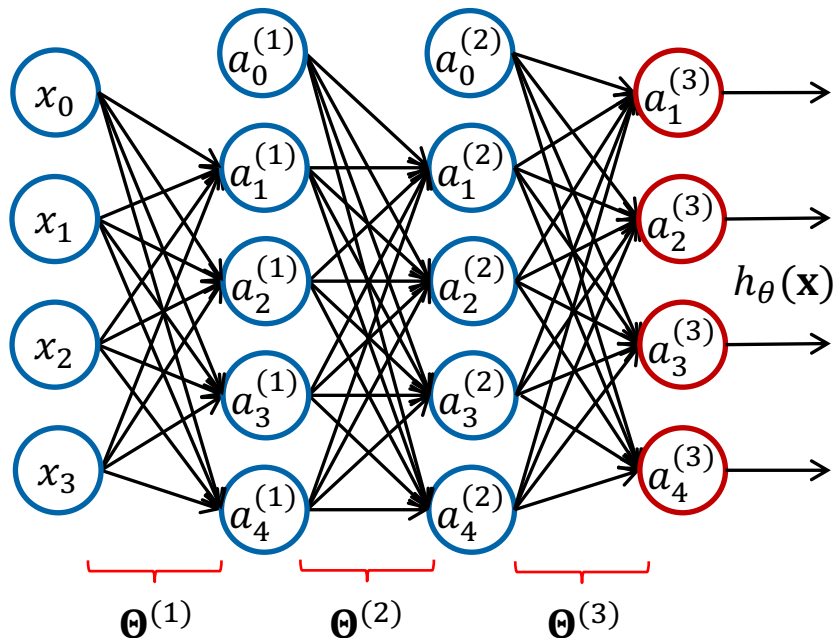
$$\frac{\partial J(\Theta)}{\partial \theta_{ij}^{(L)}} = \frac{\partial J(\Theta)}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \frac{\partial z_i^{(L)}}{\partial \theta_{ij}^{(L)}}$$

Observations:

$\frac{\partial J(\Theta)}{\partial a_i^{(L)}}$ — cost function derivative with respect to NN output

$\frac{\partial a_i^{(L)}}{\partial z_i^{(L)}}$ — NN output derivative with respect to output node inputs

$\frac{\partial z_i^{(L)}}{\partial \theta_{ij}^{(L)}}$ — output node input derivative with respect to weights (connecting outputs from prior layer nodes)





To determine

$$\frac{\partial J(\Theta)}{\partial \theta_{ij}^{(L)}} = \frac{\partial J(\Theta)}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \frac{\partial z_i^{(L)}}{\partial \theta_{ij}^{(L)}}$$

define the error

$$\delta_i^{(L)} = \frac{\partial J(\Theta)}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}}$$

and solve for each derivative component

Since

$$J(\Theta) = \frac{1}{2} \sum_j (a_j^L - y_j)^2$$

we have

$$\frac{\partial J(\Theta)}{\partial a_i^{(L)}} = (a_i^{(L)} - y_i)$$

Similarly,

$$\mathbf{a}^{(L)} = g(\mathbf{z}^{(L)}) \Rightarrow \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = g'(z_i^{(L)})$$

Note: backpropagation relies on the fact that the activation function, $g(\cdot)$, is differentiable

Therefore

$$\delta_i^{(L)} = \frac{\partial J(\Theta)}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = (a_i^{(L)} - y_i) g'(z_i^{(L)})$$

Finally, $z_i^{(L)} = \sum_k \theta_{ik}^{(L)} a_k^{(L-1)} \Rightarrow$

$$\frac{\partial z_i^{(L)}}{\partial \theta_{ij}^{(L)}} = \frac{\partial}{\partial \theta_{ij}^{(L)}} \left(\sum_k \theta_{ik}^{(L)} a_k^{(L-1)} \right) = a_j^{(L-1)}$$

Putting all the components together,

$$\begin{aligned} \frac{\partial J(\Theta)}{\partial \theta_{ij}^{(L)}} &= \frac{\partial J(\Theta)}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} \frac{\partial z_i^{(L)}}{\partial \theta_{ij}^{(L)}} \\ &= \delta_i^{(L)} a_j^{(L-1)} \end{aligned}$$

Note that this can be generalized to arbitrary level $l = 1, \dots, L$

$$\frac{\partial J(\Theta)}{\partial \theta_{ij}^{(l)}} = \delta_i^{(l)} a_j^{(l-1)}$$

where it is understood $a_j^{(0)} = x_j$ (input sample)



Consider the error term $\delta_i^{(L-1)}$, first without the chain rule:

$$\delta_i^{(L-1)} = \frac{\partial J(\Theta)}{\partial z_i^{(L-1)}}$$

To utilize the chain rule in this case, all paths between $z_i^{(L-1)}$ and the output terms must be considered

$$\delta_i^{(L-1)} = \sum_j \left(\frac{\partial J(\Theta)}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \right) \left(\frac{\partial z_j^{(L)}}{\partial a_i^{(L-1)}} \right) \left(\frac{\partial a_i^{(L-1)}}{\partial z_i^{(L-1)}} \right)$$

The first and third terms have previously determined results:

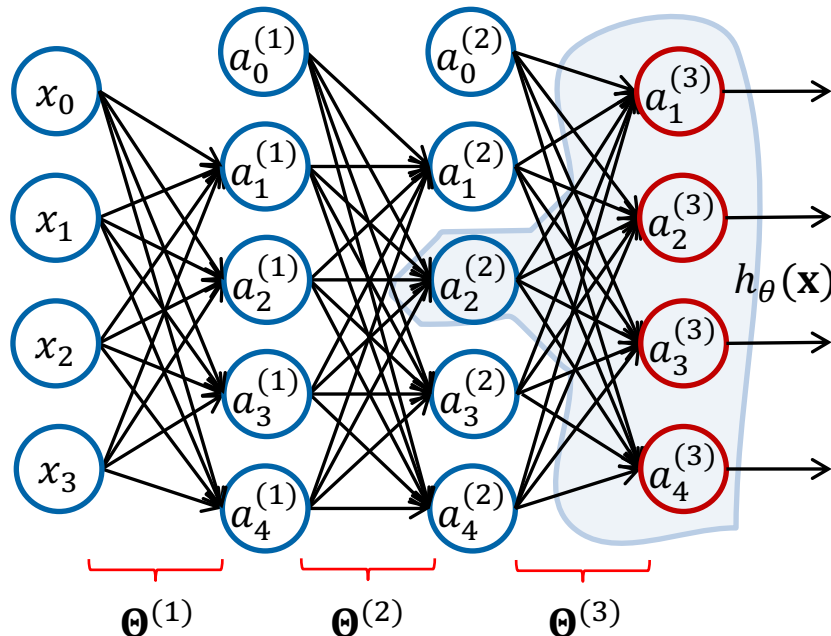
$$\delta_i^{(L-1)} = \sum_j \delta_j^{(L)} \left(\frac{\partial z_j^{(L)}}{\partial a_i^{(L-1)}} \right) g'(z_i^{(L-1)})$$

Now consider the middle term:

$$\frac{\partial z_j^{(L)}}{\partial a_i^{(L-1)}} = \frac{\partial}{\partial a_i^{(L-1)}} \left(\sum_n \theta_{jn}^{(L)} a_n^{(L-1)} \right) = \theta_{ji}^{(L)}$$

Thus

$$\delta_i^{(L-1)} = \left(\sum_j \delta_j^{(L)} \theta_{ji}^{(L)} \right) g'(z_i^{(L-1)})$$





The result

$$\delta_i^{(L-1)} = \left(\sum_j \delta_j^{(L)} \theta_{ji}^{(L)} \right) g' \left(z_i^{(L-1)} \right)$$

can be generalized as:

$$\delta_i^{(l)} = \left(\sum_j \theta_{ji}^{(l+1)} \delta_j^{(l+1)} \right) g' \left(z_i^{(l)} \right), \quad l = 1, \dots, L-1$$

With the result then utilized in the cost function derivative with respect to arbitrary weight $\theta_{ij}^{(l)}$:

$$\frac{\partial J(\Theta)}{\partial \theta_{ij}^{(l)}} = \delta_i^{(l)} a_j^{(l-1)}$$

Observation: We back propagate the error $\delta^{(l+1)}$ from layer $l+1$ to layer l

It is generally more convenient to express the operations in matrix form:

$$\delta^{(L)} = (\mathbf{a}^{(L)} - \mathbf{y}) \odot g'(\mathbf{z}^{(L)}),$$

and

$$\delta^{(l)} = \left((\Theta^{(l+1)})^T \delta^{(l+1)} \right) \odot g'(\mathbf{z}^{(l)})$$

for $l = 1, \dots, L-1$.

In this formulation, in this formulation, \odot denotes component-wise multiplication

In the particular case when $g(\mathbf{z})$ is the logistic function,

$$\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)}) = \frac{1}{1 + e^{-\mathbf{z}^{(l)}}},$$

we can express

$$g'(\mathbf{z}^{(l)}) = \mathbf{a}^{(l)} \odot (1 - \mathbf{a}^{(l)})$$



To implement backpropagation, suppose we have a training set of the form

$$\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}$$

The error is thus defined as

$$J(\Theta) = \sum_{k=1}^n \frac{1}{2} \|\mathbf{y}^{(k)} - \mathbf{a}^{(k)(L)}\|^2$$

and the derivative is

$$\frac{\partial J(\Theta)}{\partial \theta_{ij}^{(l)}} = \sum_{k=1}^n \delta_i^{(k)(l)} a_j^{(k)(l-1)},$$

Note that the superscript (k) denotes that the quantity was computed using the training sample pair: $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$



Indexing according to the training samples and utilizing vector representation:

$$\frac{\partial J(\Theta)}{\partial \Theta^{(l)}} = \sum_{k=1}^n \delta^{(k)(l)} (\mathbf{a}^{(k)(l-1)})^T$$

Putting everything together yields the following **backpropagation algorithm**, where step #4 is the standard **gradient descent update** (other update methodologies can be utilized, given the gradient)

Backpropagation algorithm:

Suppose we have a training set of the form $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}$

1. Feedforward: compute $\mathbf{z}^{(k)(l)} = \Theta^{(l)} \mathbf{a}^{(k)(l-1)}$ and $\mathbf{a}^{(k)(l)} = g(\mathbf{z}^{(k)(l)})$
2. Feedback: Compute $\delta^{(k)(L)} = (\mathbf{a}^{(k)(L)} - \mathbf{y}) \odot g'(\mathbf{z}^{(k)(L)})$ and
$$\delta^{(k)(l)} = \left((\Theta^{(l+1)})^T \delta^{(k)(l+1)} \right) \odot g'(\mathbf{z}^{(k)(l)})$$
3. Compute $J(\Theta) = \sum_{k=1}^n \frac{1}{2} \|\mathbf{y}^{(k)} - \mathbf{a}^{(k)(L)}\|^2$ and $\frac{\partial J(\Theta)}{\partial \Theta^{(l)}} = \sum_{k=1}^n \delta^{(k)(l)} (\mathbf{a}^{(k)(l-1)})^T$
4. Update the weights: $\Theta_{k+1}^{(l)} = \Theta_k^{(l)} - \alpha \frac{\partial J(\Theta)}{\partial \Theta^{(l)}}$



Observation: Regularization (shrinkage) can be included in the optimization to avoid over fitting and reduce the number of non-zero weights

For instance, including L_2 regularization yields:

$$J(\Theta) = \sum_{k=1}^n \frac{1}{2} \|\mathbf{y}^{(k)} - \mathbf{a}^{(k)(L)}\|^2 + \frac{\lambda}{2} \sum_l \sum_i \sum_j (\theta_{ij}^l)^2$$

and the derivative becomes:

$$\frac{\partial J(\Theta)}{\partial \Theta^{(l)}} = \sum_{k=1}^n \delta^{(k)(l)} (\mathbf{a}^{(k)(l-1)})^T + \lambda \Theta^{(l)}$$



Summary:

- Neural networks can be optimized utilizing **backpropagation** in conjunction **gradient descent**
- Optimization consists of two stages: **propagation** and **weight update**
- The effect of an **input is propagated forward** through the network, the error determined, and the **error is propagated back** through the network
- Weights are adjusted according to **gradient descent** and the back propagated error
- **Regularization**, or shrinkage, can be applied to the optimization in order to minimize overfitting and reduce the number of non-zero weights