

# Back Propagation

Alon Filler, Nadav Menirav

August 2025

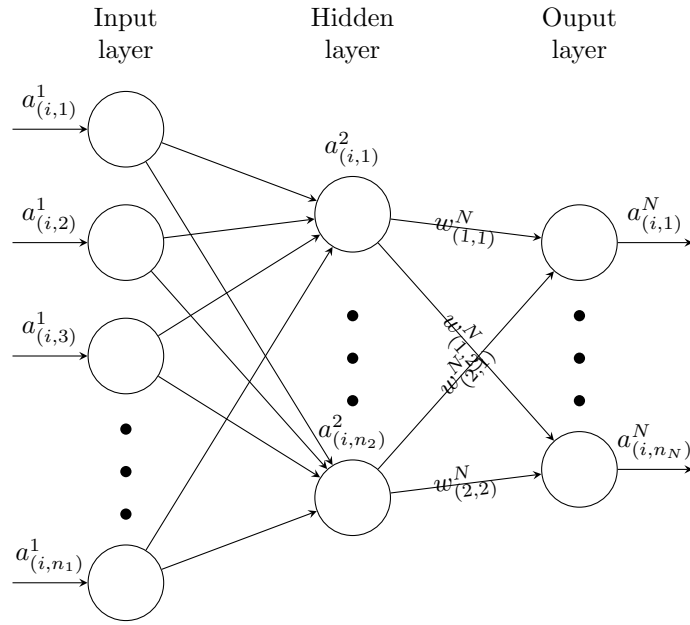
## Abstract

Back Propagation is one of many algorithms used to determine the gradient of a neural network's cost function. Its approach differs from other algorithms, in that it computes the derivatives for the last layer, and uses the results to compute the derivatives of earlier layers recursively. In this paper, we are going to derive the formula for the derivatives of each of the parameters (weights and biases) in a general neural network.

## Notations

Layers	$= N$
Neurons per layer	$= n_l$
Training data count	$= D$
Entry in training data	$= i \in \{1, \dots, D\}$
Neuron in layer $l$	$= j \in \{1, \dots, n_l\}$
Activation of neuron before non-linearity	$= \sigma^{-1}(a_{(i,j)}^l) = \left( \sum_{k=1}^{n_{l-1}} w_{(k,j)}^l a_{(i,k)}^{l-1} \right) + b_j^l$
Activation of neuron	$= a_{(i,j)}^l = \sigma \left( \left( \sum_{k=1}^{n_{l-1}} w_{(k,j)}^l a_{(i,k)}^{l-1} \right) + b_j^l \right)$
Input	$= a_{(i,j)}^1$
Output	$= a_{(i,j)}^N$

For the ease of the reader, here is an illustration of the notations for the  $i$  input:



## Calculating Gradient

Deriving the cost

$$\begin{aligned}
 cost(X) &= \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} (prediction_{(i,j)} - expected_{(i,j)})^2 \right) \\
 &= \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} (a^N_{(i,j)} - expected_{(i,j)})^2 \right)
 \end{aligned}$$

$$\begin{aligned}
\partial (cost(X))_{b_m^l} &= \partial \left( \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} (a_{(i,j)}^N - expected_{(i,j)})^2 \right) \right)_{b_m^l} \\
&= \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} \partial \left( (a_{(i,j)}^N - expected_{(i,j)})^2 \right)_{b_m^l} \right) \\
&= \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} 2(a_{(i,j)}^N - expected_{(i,j)}) \partial \left( a_{(i,j)}^N - expected_{(i,j)} \right)_{b_m^l} \right) \\
&= \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} 2(a_{(i,j)}^N - expected_{(i,j)}) \partial \left( a_{(i,j)}^N \right)_{b_m^l} \right)
\end{aligned}$$

The partial derivative was calculated below.

$$\partial (cost(X))_{w_{(m,t)}^l} = \frac{1}{D} \sum_{i=1}^D \left( \sum_{j=1}^{n_N} 2(a_{(i,j)}^N - expected_{(i,j)}) \partial \left( a_{(i,j)}^N \right)_{w_{(m,t)}^l} \right)$$

This partial derivative was also calculated below.

## Sub derivatives of the cost

$$\partial \left( a_{(i,j)}^N \right)_{a_{(i,m)}^{N-1}} = \partial \left( \left( \sum_{k=1}^{n_{N-1}} w_{(k,j)}^N a_{(i,k)}^{N-1} \right) + b_j^N \right)_{a_{(i,m)}^{N-1}} = w_{(m,j)}^N$$

$$\partial \left( a_{(i,j)}^N \right)_{w_{(m,j)}^N} = \partial \left( \left( \sum_{k=1}^{n_{N-1}} w_{(k,j)}^N a_{(i,k)}^{N-1} \right) + b_j^N \right)_{w_{(m,j)}^N} = a_{(i,m)}^{N-1}$$

$$\partial \left( a_{(i,j)}^N \right)_{b_j^N} = \partial \left( \left( \sum_{k=1}^{n_{N-1}} w_{(k,j)}^N a_{(i,k)}^{N-1} \right) + b_j^N \right)_{b_j^N} = 1$$

$$\frac{\partial a_{(i,t)}^l}{\partial w_{(m,t)}^l} = \partial \left( a_{(i,t)}^l \right)_{w_{(m,t)}^l} = a_{(i,t)}^l \cdot \left( 1 - a_{(i,t)}^l \right) \cdot a_{(i,m)}^{l-1}$$

$$\begin{aligned} \frac{\partial a_{(i,t)}^{l+1}}{\partial b_m^l} &= \partial \left( a_{(i,t)}^{l+1} \right)_{b_m^l} = \partial \sigma \left( \left( \sum_{k=1}^{n_l} w_{(k,t)}^{l+1} a_{(i,k)}^l \right) + b_t^{l+1} \right)_{b_m^l} \\ &= / : \partial (\sigma(f(x)))_{x_i} = \sigma(f(x)) (1 - \sigma(f(x))) \partial (f(x))_{x_i} \\ &= a_{(i,t)}^{l+1} \cdot \left( 1 - a_{(i,t)}^{l+1} \right) \cdot \partial \left( \sum_{k=1}^{n_l} w_{(k,t)}^{l+1} a_{(i,k)}^l \right)_{b_m^l} \\ &= a_{(i,t)}^{l+1} \cdot \left( 1 - a_{(i,t)}^{l+1} \right) \cdot w_{(m,t)}^{l+1} \cdot \partial \left( a_{(i,m)}^l \right)_{b_m^l} \\ &= a_{(i,t)}^{l+1} \cdot \left( 1 - a_{(i,t)}^{l+1} \right) \cdot w_{(m,t)}^{l+1} \cdot a_{(i,m)}^l \cdot \left( 1 - a_{(i,m)}^l \right) \partial \left( \left( \sum_{k=1}^{n_{l-1}} w_{(k,m)}^l a_{(i,k)}^{l-1} \right) + b_m^l \right)_{b_m^l} \\ &= a_{(i,t)}^{l+1} \cdot \left( 1 - a_{(i,t)}^{l+1} \right) \cdot w_{(m,t)}^{l+1} \cdot a_{(i,m)}^l \cdot \left( 1 - a_{(i,m)}^l \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial a_{(i,k)}^{l+1}}{\partial a_{(i,t)}^l} &= \partial \left( a_{(i,k)}^{l+1} \right)_{\partial a_{(i,t)}^l} = a_{(i,k)}^{l+1} \cdot \left( 1 - a_{(i,k)}^{l+1} \right) \cdot \partial \left( \sigma^{-1} \left( a_{(i,k)}^{l+1} \right) \right)_{a_{(i,t)}^l} \\ &= a_{(i,k)}^{l+1} \cdot \left( 1 - a_{(i,k)}^{l+1} \right) \cdot \partial \left( \left( \sum_{j=1}^{n_l} w_{(j,k)}^{l+1} a_{(i,j)}^l \right) + b_k^{l+1} \right)_{a_{(i,t)}^l} \\ &= a_{(i,k)}^{l+1} \cdot \left( 1 - a_{(i,k)}^{l+1} \right) \cdot \partial \left( w_{(t,k)}^{l+1} a_{(i,t)}^l \right)_{a_{(i,t)}^l} \\ &= a_{(i,k)}^{l+1} \cdot \left( 1 - a_{(i,k)}^{l+1} \right) \cdot w_{(t,k)}^{l+1} \end{aligned}$$

$$\begin{aligned}
\frac{\partial a_{(i,j)}^N}{\partial a_{(i,t)}^{l+1}} &= \partial \left( a_{(i,j)}^N \right)_{a_{(i,t)}^{l+1}} = \sum_{k=1}^{n_{l+2}} \frac{\partial a_{(i,j)}^N}{\partial a_{(i,k)}^{l+2}} \cdot \frac{\partial a_{(i,k)}^{l+2}}{\partial a_{(i,t)}^{l+1}} \\
\partial \left( a_{(i,j)}^N \right)_{b_m^l} &= \sum_{t=1}^{n_{l+1}} \frac{\partial a_{(i,j)}^N}{\partial a_{(i,t)}^{l+1}} \cdot \frac{\partial a_{(i,t)}^{l+1}}{\partial b_m^l} \\
\partial \left( a_{(i,j)}^N \right)_{w_{(m,t)}^l} &= \frac{\partial a_{(i,j)}^N}{\partial a_{(i,t)}^l} \cdot \frac{\partial a_{(i,t)}^l}{\partial w_{(m,t)}^l}
\end{aligned}$$