

Neural Network with Layers

Let $n \in \mathbb{N}$ be the count of layers (do not count the input layer) in a given fully-connected feed-forward neural network. Then one can describe the network through an activation function $\sigma: \mathbb{R} \rightarrow [0, 1]$, weight matrices $(W_k)_{k=1}^n$ and bias vectors $(B_k)_{k=1}^n$ with the following properties for all $k \in \mathbb{N}$ with $k \leq n$.

$$W_k \in \mathbb{R}^{p_k \times q_k} \quad B_k \in \mathbb{R}^{p_k} \quad p_k, q_k \in \mathbb{N}$$

p_k is the number of neurons and q_k the number of inputs in the k . layer. The dimensions have to fulfill a certain consistency condition such that the matrix-vector-product can be defined between each layer. For all $k \in \mathbb{N}$ with $k < n$ it holds

$$p_k = q_{k+1}$$

Then the output of the k . layer with respect to its input is described by the following function. We define $\bar{\sigma}$ to be the element-wise application of σ to a vector.

$$f[W_k, B_k]: [0, 1]^{q_k} \rightarrow [0, 1]^{p_k} \quad f[W_k, B_k](x) := \bar{\sigma}(W_k x + B_k)$$

The overall output f of the complete network is then described by composition of these functions.

$$f[(W_k, B_k)_{k=1}^n]: [0, 1]^{q_1} \rightarrow [0, 1]^{p_n}$$

$$f[(W_k, B_k)_{k=1}^n] := f[W_n, B_n] \circ \dots \circ f[W_1, B_1] = \bigcirc_{k=1}^n f[W_k, B_k]$$

Error Function

Given $m \in \mathbb{N}$ labeled training samples $(x_k, y_k)_{k=1}^m$ with $x_k \in [0, 1]^{q_1}$ and $y_k \in [0, 1]^{p_n}$ for all $k \in \mathbb{N}$ with $k \leq m$ we can define the total error of the given neural network by the mean-squared-error $\varepsilon[(x_i, y_i)_{i=1}^m]$.

$$\varepsilon[(x_i, y_i)_{i=1}^m]: \bigtimes_{k=1}^n (\mathbb{R}^{p_k \times q_k} \times \mathbb{R}^{p_k}) \rightarrow \mathbb{R}$$

$$\varepsilon[(x_i, y_i)_{i=1}^m]((W_k, B_k)_{k=1}^n) := \frac{1}{m} \sum_{i=1}^m \|y_i - f[(W_k, B_k)_{k=1}^n](x_i)\|^2$$

Gradients of the Error Function

The gradients of the error function can be computed by using the abstract index notation. This process is error-prone and cumbersome and will only be shown for the example gradient $\nabla_{W_n} \varepsilon[(x_i, y_i)_{i=1}^m]$. First, we start by reformulating the error function with the abstract index notation.

$$\varepsilon[(x_i, y_i)_{i=1}^m]((W_k, B_k)_{k=1}^n)$$

$$\begin{aligned}
&= \frac{1}{m} \sum_{i=1}^m \|y_i - f[(W_k, B_k)_{k=1}^n](x_i)\|^2 \\
&= \frac{1}{m} \sum_{i=1}^m (y_i - f[(W_k, B_k)_{k=1}^n](x_i))^T (y_i - f[(W_k, B_k)_{k=1}^n](x_i)) \\
&\equiv \frac{1}{m} \sum_{i=1}^m (y_{i\alpha} - f_\alpha[(W_k, B_k)_{k=1}^n](x_i)) (y_i^\alpha - f^\alpha[(W_k, B_k)_{k=1}^n](x_i))
\end{aligned}$$

Now we compute the partial derivate of a certain component of the gradient matrix.

$$\begin{aligned}
&[\nabla_{W_n} \varepsilon[(x_i, y_i)_{i=1}^m] ((W_k, B_k)_{k=1}^n)]^T \\
&\equiv \partial_{W_n^{\mu_\nu}} \frac{1}{m} \sum_{i=1}^m (y_{i\alpha} - f_\alpha[(W_k, B_k)_{k=1}^n](x_i)) (y_i^\alpha - f^\alpha[(W_k, B_k)_{k=1}^n](x_i)) \\
&= -\frac{2}{m} \sum_{i=1}^m (y_{i\alpha} - f_\alpha[(W_k, B_k)_{k=1}^n](x_i)) \partial_{W_n^{\mu_\nu}} f^\alpha[(W_k, B_k)_{k=1}^n](x_i) \\
&\quad \partial_{W_n^{\mu_\nu}} f^\alpha[(W_k, B_k)_{k=1}^n](x_i) \\
&= \partial_{W_n^{\mu_\nu}} \sigma(W_n^{\alpha_\beta} f^\beta[(W_k, B_k)_{k=1}^{n-1}](x_i) + B_n^\alpha) \\
&= \sigma'(W_n^{\alpha_\beta} f^\beta[(W_k, B_k)_{k=1}^{n-1}](x_i) + B_n^\alpha) \delta_\mu^\alpha \delta_\beta^\nu f^\beta[(W_k, B_k)_{k=1}^{n-1}](x_i) \\
&= \sigma'(W_n^{\alpha_\beta} f^\beta[(W_k, B_k)_{k=1}^{n-1}](x_i) + B_n^\alpha) \delta_\mu^\alpha f^\nu[(W_k, B_k)_{k=1}^{n-1}](x_i)
\end{aligned}$$

Putting all together we get the result for the complete gradient. Please note that I have transposed the matrix to build the real gradient. We define $\bar{\sigma}'$ to be the element-wise application of σ' to a vector.

$$\begin{aligned}
&\nabla_{W_n} \varepsilon[(x_i, y_i)_{i=1}^m] ((W_k, B_k)_{k=1}^n) \\
&= -\frac{2}{m} \sum_{i=1}^m \{ (y_i - f[(W_k, B_k)_{k=1}^n](x_i)) \odot \bar{\sigma}'(W_n f[(W_k, B_k)_{k=1}^{n-1}](x_i) + B_n) \} \\
&\quad \cdot \{ f[(W_k, B_k)_{k=1}^{n-1}](x_i) \}^T
\end{aligned}$$

The other gradients can be computed analogously. If we apply this process iteratively to every weight matrix and bias vector then we get recursive formulas to compute every gradient. The algorithm is given below.

Forward Feed Algorithm

For convenience we define the following for a given input $x \in [0, 1]^{q_1}$.

$$f_0 := x$$

For all $k \in \mathbb{N}$ with $k \leq n$ we compute the following recursive formulas.

$$t_k := W_k f_{k-1} + B_k$$

$$f_k := \bar{\sigma}(t_k)$$

t_k describes the net input vector and f_k the output vector of the k . layer. Therefore f_n is the output of the complete neural network.

Backward Feed Algorithm

Calculate a first temporary vector.

$$z_n := (y - f_n) \odot \bar{\sigma}'(t_n)$$

For all $k \in \mathbb{N}$ with $k < n$ compute the following vectors.

$$z_k := (W_{k+1}^T z_{k+1}) \odot \bar{\sigma}'(t_k)$$

Get the gradients from the temporary vectors.

$$\nabla_{W_k} e = -2z_k \cdot f_{k-1}^T$$

$$\nabla_{B_k} e = -2z_k$$

Use the gradients for an optimization algorithm like the steepest descent method with learning rate η .

$$W_k \leftarrow W_k - \eta \nabla_{W_k} e$$

$$B_k \leftarrow B_k - \eta \nabla_{B_k} e$$

Results

labeled training sample (x, y)

$$x := \begin{pmatrix} i_1 \\ i_2 \end{pmatrix} = \begin{pmatrix} 0.05 \\ 0.10 \end{pmatrix} \quad y := \begin{pmatrix} o_1 \\ o_2 \end{pmatrix} = \begin{pmatrix} 0.01 \\ 0.99 \end{pmatrix}$$

initial weight matrices W_1 and W_2

$$W_1^{(0)} := \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix} = \begin{pmatrix} 0.15 & 0.20 \\ 0.25 & 0.30 \end{pmatrix} \quad W_2^{(0)} := \begin{pmatrix} w_5 & w_6 \\ w_7 & w_8 \end{pmatrix} = \begin{pmatrix} 0.40 & 0.45 \\ 0.50 & 0.55 \end{pmatrix}$$

initial bias vectors B_1 and B_2

$$B_1 := \begin{pmatrix} b_1 \\ b_1 \end{pmatrix} = \begin{pmatrix} 0.35 \\ 0.35 \end{pmatrix} \quad B_2 := \begin{pmatrix} b_2 \\ b_2 \end{pmatrix} = \begin{pmatrix} 0.60 \\ 0.60 \end{pmatrix}$$

net input to first layer t_1

$$t_1 = W_1 x + B_1 = \begin{pmatrix} 0.3775 \\ 0.3925 \end{pmatrix}$$

output of first layer f_1

$$f_1 = \sigma(t_1) = \begin{pmatrix} 0.593270 \\ 0.596884 \end{pmatrix}$$

net input to second layer t_2

$$t_2 = W_2 f_1 + B_2 = \begin{pmatrix} 1.10591 \\ 1.22492 \end{pmatrix}$$

output of second layer f_2

$$f_2 = \sigma(t_2) = \begin{pmatrix} 0.751365 \\ 0.772928 \end{pmatrix}$$

total error e_0

$$e_0 = \frac{1}{2} \|y - f_2\|^2 = 0.298371$$

weight matrices after one backpropagation step $W_1^{(1)}, W_2^{(1)}$ without changing bias

$$W_1^{(0)} = \begin{pmatrix} 0.149781 & 0.199561 \\ 0.249751 & 0.299502 \end{pmatrix} \quad W_2^{(0)} = \begin{pmatrix} 0.358916 & 0.408666 \\ 0.511301 & 0.56137 \end{pmatrix}$$

total error after one backpropagation step and a second forward feed:

$$e_1 = 0.291028$$

$$\implies e_1 < e_0$$