

Derivation of a Forward and Backpropagation Pass on a Feedforward Neural Network

Raul Morales Delgado

Contents

1	Objective	2
2	Scope	2
3	Introduction	2
4	The Feedforward Neural Network (FFNN)	3
5	Analysis of the Forward Pass	5
6	Defining the Loss Function and its Gradient	6
7	Gradient Descent and the Learning Rate	7
8	Analysis of the Backpropagation Pass	8
9	Conclusions	13
	Appendix A	14
	Appendix B	15

1 Objective

The objective of this tutorial is to demonstrate, through thorough calculations, how to derive a forward pass and a back-propagation pass through a feedforward neural network (FFNN) of n inputs, one hidden layer with m nodes, and p outputs. In this sense, this tutorial aims to be a general case application for any FFNN with a single hidden layer.

2 Scope

In this tutorial, the architecture of a FFNN with a single hidden layer will be presented. This architecture considers n inputs, m nodes in the hidden layer, and p outputs. Also, it includes the use of biases for the hidden and output layers, as well as generic activation functions for each node of these layers.

For the back-propagation process, a variation of the Mean Squared Error (MSE) function will be used as the loss function in order to set the iterative process to aim to minimize the error of the network in each iteration.

Since there is plenty documentation, articles and tutorials about how to code and use a NN, and how to tweak their hyperparameters for better learning and resulting accuracy, these topics will not be covered here. This tutorial will only focus on the mathematical processes involved in a simple FFNN.

3 Introduction

Feedforward neural networks (FFNN) are among the simplest neural networks and are the core upon which much more complex neural networks are built. A few days ago (from when I started writing this tutorial), I was taking a course on Natural Language Processing (NLP) and, as a way to explain RNN and LSTM networks, FFNN were revisited in a detailed manner, such that the student would follow along with the derivations that allow one to obtain the results of a forward pass and of a backpropagation pass.

However, throughout the process, there were some inconsistencies in the material which did not allow me to correctly follow the derivation. Because of this, I decided to do my very own, which is the core of this tutorial. Although there are plenty examples out there of the results of this process, they tend to be extremely summarized — they present the network and they leap through the equations — such that the underlying math is not clear if you want to follow along. Furthermore, these examples are usually limited to a single node in the output layer (i.e. for binary classification, e.g. `dog` / `no_dog`) and very few neurons in the hidden layer (three in the course I took).

For these reasons, I decided to derive the general case of a feedforward neural network with n input nodes, a single hidden layer with m nodes, and p output nodes. The derivation will consist on the mathematical procedures involved in one forward pass through the network and one backpropagation pass. To execute the backpropagation process, a loss function is considered, as previously explained.

Finally, please bear in mind that I am not a Computer Scientist, let alone a Mathematician. While there might be some lack of rigorosity in the definitions or derivations, I trust that the process has been performed adequately and that the results are correct. If you find any mistake, please do not hesitate to email me or submit a GitHub issue.

4 The Feedforward Neural Network (FFNN)

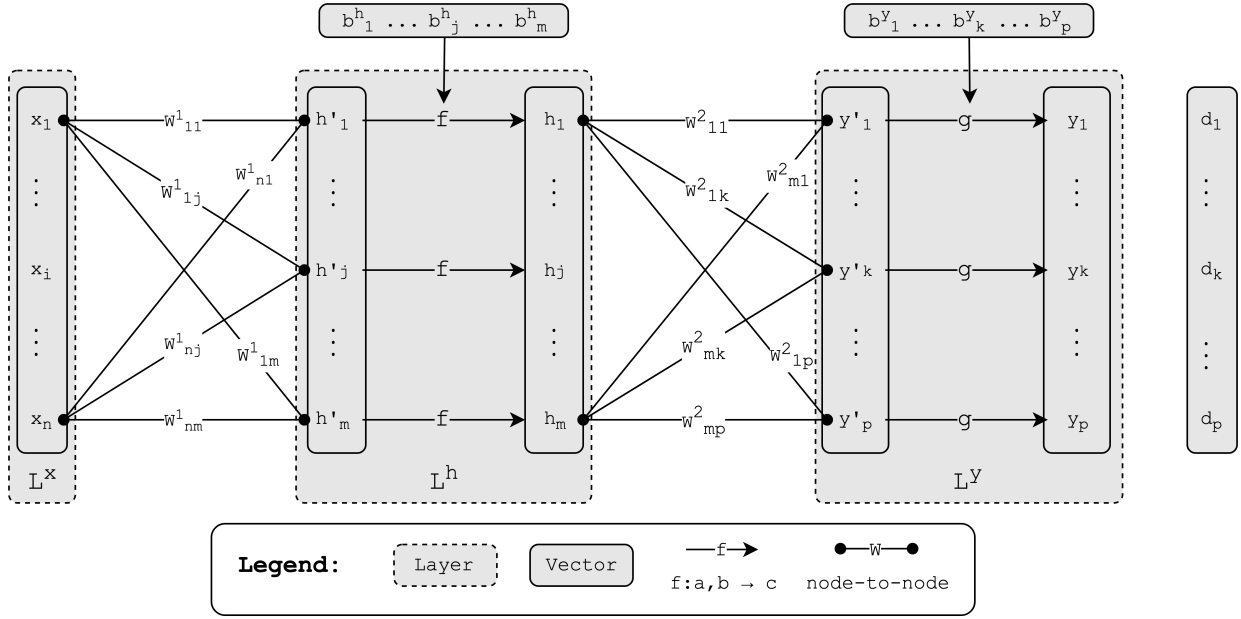


Figure 1: Diagram of a FFNN with a single hidden layer.

For the purpose of this tutorial, let's consider the FFNN of Figure 1. In this neural network, from left to right, the following can be defined:

- L^x, L^h, L^y Input, hidden and output layer, respectively.
- \mathbf{x} : Input vector of n components, where x_i represents the i^{th} arbitrary component of the vector, such that:

$$\mathbf{x} = [x_1, \dots, x_i, \dots, x_n]$$

- \mathbf{W}^1 : Matrix of weights 1, with size $n \times m$, which connects L^x to L^h , where:

$$\mathbf{W}^1 = \begin{bmatrix} W_{1,1}^1 & \dots & W_{1,j}^1 & \dots & W_{1,m}^1 \\ \vdots & \ddots & \vdots & & \vdots \\ W_{i,1}^1 & \dots & W_{i,j}^1 & \dots & W_{i,m}^1 \\ \vdots & & \vdots & \ddots & \vdots \\ W_{n,1}^1 & \dots & W_{n,j}^1 & \dots & W_{n,m}^1 \end{bmatrix}$$

- \mathbf{h}' : Input vector to L^h of m components, where h'_j represents the j^{th} arbitrary component of the vector, such that:

$$\mathbf{h}' = [h'_1, \dots, h'_j, \dots, h'_m]$$

- \mathbf{b}^h : Bias vector of L^h of m components, where b_j^h represents the j^{th} arbitrary component of the vector, such that:

$$\mathbf{b}^h = [b_1^h, \dots, b_j^h, \dots, b_m^h]$$

- f : Activation function of the hidden layer, L^h , such that $f : a, b \rightarrow c$, where a and b are inputs of the function, and c its output.
- \mathbf{h} : Output vector to L^h of m components, where h_j represents the j^{th} arbitrary component of the vector, such that:

$$\mathbf{h} = [h_1, \dots, h_j, \dots, h_m]$$

- \mathbf{W}^2 : Matrix of weights 2, with size $m \times p$, which connects L^h to L^y , where:

$$\mathbf{W}^2 = \begin{bmatrix} W_{1,1}^2 & \dots & W_{1,k}^2 & \dots & W_{1,p}^2 \\ \vdots & \ddots & \vdots & & \vdots \\ W_{j,1}^2 & \dots & W_{j,k}^2 & \dots & W_{j,p}^2 \\ \vdots & & \vdots & \ddots & \vdots \\ W_{m,1}^2 & \dots & W_{m,k}^2 & \dots & W_{m,p}^2 \end{bmatrix}$$

- \mathbf{y}' : Input vector to L^y of p components, where y'_k represents the k^{th} arbitrary component of the vector, such that:

$$\mathbf{y}' = [y'_1, \dots, y'_k, \dots, y'_p]$$

- \mathbf{b}^y : Bias vector of L^y of p components, where b_k^y represents the k^{th} arbitrary component of the vector, such that:

$$\mathbf{b}^y = [b_1^y, \dots, b_k^y, \dots, b_p^y]$$

- g : Activation function of the output layer, L^y , such that $g : a, b \rightarrow c$, where a and b are inputs of the function, and c its output.
- \mathbf{y} : Output vector to L^y of p components, where y_k represents the k^{th} arbitrary component of the vector, such that:

$$\mathbf{y} = [y_1, \dots, y_k, \dots, y_p]$$

- \mathbf{d} : Target vector (vector of *true* values) of p components, where d_k represents the k^{th} arbitrary component of the vector, such that:

$$\mathbf{d} = [d_1, \dots, d_k, \dots, d_p]$$

A forward pass in this FFNN can be defined as the set of signals that start at the input layer, L^x , and finish in the output layer, L^y ; they go from left to right. Inversely, a backpropagation pass starts at L^y and finish at L^x . In this FFNN, all nodes are interconnected (some of them, though, have not been connected in the diagram to not clutter it).

Also, for the purpose of this tutorial, vectors will be considered as 1-dimensional matrices. For instance, \mathbf{x} is a matrix of 1 row and n columns.

5 Analysis of the Forward Pass

Given the input vector \mathbf{x} and the weight matrix \mathbf{W}^1 , the input vector to the hidden layer, \mathbf{h}' , can be described by:

$$\mathbf{h}' = \mathbf{x}\mathbf{W}^1 \quad (1)$$

In the hidden layer, L^h , an activation function f , whose input parameters are vectors \mathbf{h}' and \mathbf{b}^h , is used to yield the output vector \mathbf{h} , such that $f : h'_j, b_j^h \rightarrow h_j$. Therefore:

$$\mathbf{h} = f(\mathbf{h}' + \mathbf{b}^h) \quad (2)$$

A good summary about activation functions can be found in this [here](#). Now, expanding the matrix multiplication of Eq.(1):

$$\mathbf{h}' = [x_1, \dots, x_i, \dots, x_n] \begin{bmatrix} W_{1,1}^1 & \dots & W_{1,j}^1 & \dots & W_{1,m}^1 \\ \vdots & \ddots & \vdots & & \vdots \\ W_{i,1}^1 & \dots & W_{i,j}^1 & \dots & W_{i,m}^1 \\ \vdots & & \vdots & \ddots & \vdots \\ W_{n,1}^1 & \dots & W_{n,j}^1 & \dots & W_{n,m}^1 \end{bmatrix};$$

And evaluating the result for the j^{th} component of \mathbf{h}' :

$$h'_j = x_1 \cdot W_{1,j}^1 + \dots + x_i \cdot W_{i,j}^1 + \dots + x_n \cdot W_{n,j}^1;$$

$$\therefore h'_j = \sum_i^n x_i \cdot W_{ij}^1. \quad (3)$$

Consequently, the j^{th} component in \mathbf{h} can be described by:

$$h_j = f\left(\sum_i^n x_i \cdot W_{ij}^1 + b_j^h\right). \quad (4)$$

If we define ψ_j as:

$$\psi_j = \sum_i^n x_i \cdot W_{ij}^1 + b_j^h \quad (5)$$

Then, Eq.(4) can be rewritten as:

$$h_j = f(\psi_j). \quad (6)$$

Moving forward in the FFNN, the previous process can be repeated between the hidden layer, L^h , and the output layer, L^y . Therefore,

$$\mathbf{y}' = \mathbf{h}\mathbf{W}^2 \quad (7)$$

And, in the output layer L^y , the activation function g , whose input parameters are \mathbf{y}' and \mathbf{b}^y , is used to yield the output vector \mathbf{y} , such that $g : y'_k, b_k^y \rightarrow y_k$, and:

$$\mathbf{y} = g(\mathbf{y}' + \mathbf{b}^y). \quad (8)$$

Consequently, the evaluated k^{th} component of \mathbf{y} can be described by:

$$y_k = g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right). \quad (9)$$

If we define ω_k as:

$$\omega_k = \sum_j^m h_j \cdot W_{jk}^2 + b_k^y \quad (10)$$

Then, Eq.(9) can be rewritten as:

$$y_k = g(\omega_k). \quad (11)$$

Finally, Eq.(4) can be plugged into Eq.(9) and the expression rearranged for improved readability:

$$y_k = g \left(b_k^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{jk}^2 \right). \quad (12)$$

Therefore, from the previous equation, it can be stated that:

$$y_k = y_k(x_i, W_{ij}^1, W_{jk}^2, b_j^h, b_k^y) \quad (13)$$

And when the learning process (training) is completed, \mathbf{W}^1 , \mathbf{W}^2 , \mathbf{b}^h and \mathbf{b}^y are constant. Consequently, for all components of \mathbf{y} :

$$\mathbf{y} = \mathbf{y}(\mathbf{x}). \quad (14)$$

In summary, the previous expression indicates that when the training is completed a set of optimized weights and bias matrices — namely, \mathbf{W}^1 , \mathbf{W}^2 , \mathbf{b}^h and \mathbf{b}^y for this FFNN — have been found such that the model will output a vector, \mathbf{y} , as a function of a given input vector, \mathbf{x} .

6 Defining the Loss Function and its Gradient

Deep Learning neural networks, like many other AI models, require an iterative process during the training phase to adjust their parameters, like weights and bias matrices. Particularly, in supervised learning, where a model is trained with a labeled dataset, this process allows to improve the model's ability to predict a set of features from a given input. To achieve this, the iterative process uses a loss function (or, more generally, a cost function) to compare the output vector of the neural network to the target vector — whose difference is called the *error*. In general, the minimization of the output of this function — the minimization of the error — is performed using gradient descent and it allows to: a) define the path that will allow to reduce the error by adequately updating the neural network's parameters, and, b) set the limit as to where the model has reached an adequate learning (e.g. an early stop callback for when the error reaches a certain value).

To accomplish its goal, the loss function operates in the backpropagation pass of the network, which happens once after every forward pass during the training phase. The backpropagation pass will be analyzed in detail in the following section, but for now it is important to highlight that this pass goes from right to left — it starts in the output layer L^y and finishes in the input layer L^x (more specifically, it finishes in the weight matrix *before* the input layer).

The loss function considered for this tutorial will be a variation of the Mean Squared Error (MSE). More information about this function can be found in this link. If the MSE is defined as the loss function to calculate the error, E , then:

$$E = \sum_k^p \frac{(d_k - y_k)^2}{2}. \quad (15)$$

From the previous equation and from Eq.(13), it can be stated that:

$$E = E(d_k, x_i, W_{ij}^1, W_{jk}^2, b_j^h, b_k^y). \quad (16)$$

During the backpropagation pass, which is when the loss function is applied to the neural network to calculate the error, the input vector \mathbf{x} and the output vector \mathbf{y} are constant — the error is analyzed for this state of the network. In consequence, the independent variables of E can be reduced to:

$$E = E(W_{ij}^1, W_{jk}^2, b_j^h, b_k^y). \quad (17)$$

To minimize the error of the network using gradient descent, the loss function must be differentiated with respect to its independent variables to find the gradient of the function, ∇E , which can be stated as:

$$\nabla E = \left(\frac{\partial E}{\partial W_{ij}^1}, \frac{\partial E}{\partial W_{jk}^2}, \frac{\partial E}{\partial b_j^h}, \frac{\partial E}{\partial b_k^y} \right). \quad (18)$$

7 Gradient Descent and the Learning Rate

During the iteration process of a neural network, the current iteration step can be defined as t , and the immediate next step as $t + 1$. In a single step t , a forward pass and a backpropagation pass are completed, and as mentioned earlier, the result of the backpropagation pass will be an updated set of the neural network's parameters that will be used in the following step. Furthermore, these parameters are updated such that the error, calculated with the loss function, is reduced by using a gradient descent algorithm — a type of optimization algorithm — to perform a small displacement in each of the independent variables values in the direction of the negative gradient — i.e. to slightly modify each of the parameters of the neural network. More information about gradient descent can be found here.

To execute the small displacement on these variables, the gradient descent algorithm introduces a learning rate α and uses it to set the magnitude for the displacement at each step. For instance, for the FFNN of this tutorial, the displacements of the loss function's independent variables — ΔW_{ij}^1 , ΔW_{jk}^2 , Δb_j^h and Δb_k^y — can be defined as follows:

$$W_{ij@t+1}^1 - W_{ij@t}^1 = \Delta W_{ij}^1, \quad (19)$$

$$W_{jk@t+1}^2 - W_{jk@t}^2 = \Delta W_{jk}^2, \quad (20)$$

$$b_{j@t+1}^h - b_{j@t}^h = \Delta b_j^h, \quad (21)$$

$$b_{k@t+1}^y - b_{k@t}^y = \Delta b_k^y; \quad (22)$$

Where,

$$\Delta W_{ij}^1 = -\alpha \cdot \frac{\partial E}{\partial W_{ij}^1}, \quad (23)$$

$$\Delta W_{jk}^2 = -\alpha \cdot \frac{\partial E}{\partial W_{jk}^2}, \quad (24)$$

$$\Delta b_j^h = -\alpha \cdot \frac{\partial E}{\partial b_j^h}, \quad (25)$$

$$\Delta b_k^y = -\alpha \cdot \frac{\partial E}{\partial b_k^y}. \quad (26)$$

From the last four equations, it is important to highlight that $\alpha \in \mathbb{R} > 0$, and the negative sign in the expressions indicate that the displacement will be performed in the direction of the negative gradient. A simple explanation about how a gradient descent algorithm takes advantage of the negative gradient to minimize the error can be found in Appendix A.

Finally, when the small displacements are calculated and the network's parameters updated, the resulting error forms a monotonic sequences throughout the iteration process, such that:

$$E_{t-1} \geq E_t \geq E_{t+1} \geq \dots. \quad (27)$$

8 Analysis of the Backpropagation Pass

In simple terms, the backpropagation pass in iteration step t is used to minimize the error, such that an *optimized* set of parameters is calculated and used in the next step, $t+1$. The optimization of these values is defined by Eq.(19–22), and to solve those equations, Eq.(23–26) must be derived first. To keep a sequential order, this process will be carried out from the output layer towards the input layer, starting with the weights- and then the bias-related variables.

Derivation of ΔW_{jk}^2

To derive ΔW_{jk}^2 , Eq.(15) is plugged into Eq.(24), such that:

$$\Delta W_{jk}^2 = -\alpha \cdot \frac{\partial \left(\sum_k^p \frac{(d_k - y_k)^2}{2} \right)}{\partial W_{jk}^2}. \quad (28)$$

If Eq.(9) is plugged into the previous expression and the sum expanded:

$$\Delta W_{jk}^2 = -\alpha \cdot \left(\frac{\partial \left(\frac{(d_1 - g \left(\sum_j^m h_j \cdot W_{j1}^2 + b_1^y \right))^2}{2} \right)}{\partial W_{jk}^2} + \dots + \frac{\partial \left(\frac{(d_k - g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right))^2}{2} \right)}{\partial W_{jk}^2} + \dots + \frac{\partial \left(\frac{(d_p - g \left(\sum_j^m h_j \cdot W_{jp}^2 + b_p^y \right))^2}{2} \right)}{\partial W_{jk}^2} \right), \quad (29)$$

If each of the terms of the previous expression is differentiated with respect to W_{jk}^2 , only the middle term — the one who has the W_{jk}^2 component in the numerator — would remain, such that:

$$\Delta W_{jk}^2 = -\alpha \cdot \left(\frac{\partial \left(\frac{\left(d_k - g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right) \right)^2}{2} \right)}{\partial W_{jk}^2} \right),$$

And after differentiating by ∂W_{jk}^2 , the expression can be simplified back to its previous form:

$$\Delta W_{jk}^2 = \alpha \cdot (d_k - y_k) \cdot \frac{\partial y_k}{\partial W_{jk}^2}. \quad (30)$$

If δ_{W^2} is defined as:

$$\delta_{W^2} = \frac{\partial y_k}{\partial W_{jk}^2}, \quad (31)$$

Then, by substituting Eq.(9) into the previous expression:

$$\delta_{W^2} = \frac{\partial g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial W_{jk}^2}, \quad (32)$$

And by applying the chain rule and replacing with Eq.(10), the expression can be simplified to:

$$\begin{aligned} \delta_{W^2} &= \frac{\partial g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)} \cdot \frac{\partial \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial W_{jk}^2}; \\ \therefore \delta_{W^2} &= g'(\omega_k) \cdot h_j; \end{aligned} \quad (33)$$

Where:

$$g'(\omega_k) = \frac{\partial g(\omega_k)}{\partial \omega_k} = \frac{\partial g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}. \quad (34)$$

Note how the summation disappears in the chain derivation when differentiating with respect to W_{jk}^2 . This process is explained in Appendix B.

Finally, replacing Eq.(33) into Eq.(31), and the latter into Eq.(30):

$$\Delta W_{jk}^2 = \alpha \cdot (d_k - y_k) \cdot g'(\omega_k) \cdot h_j. \quad (35)$$

Derivation of Δb_k^y

The derivation of Δb_k^y is very much like that of ΔW_{jk}^2 :

$$\Delta b_k^y = -\alpha \cdot \frac{\partial \left(\sum_k^p \frac{(d_k - y_k)^2}{2} \right)}{\partial b_k^y}, \quad (36)$$

After expanding the summation, because the differentiation is with respect to b_k^y , only the following remains:

$$\Delta b_k^y = -\alpha \cdot \left(\frac{\partial \left(\left(d_k - g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right) \right)^2 \right)}{\frac{2}{\partial b_k^y}} \right),$$

And after differentiating:

$$\Delta b_k^y = \alpha \cdot (d_k - y_k) \cdot \frac{\partial y_k}{\partial b_k^y}. \quad (37)$$

If δ_{b^y} is defined as:

$$\delta_{b^y} = \frac{\partial y_k}{\partial b_k^y}, \quad (38)$$

Then, by substituting Eq.(9) into the previous expression:

$$\delta_{b^y} = \frac{\partial g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial b_k^y}, \quad (39)$$

And by applying the chain rule and replacing with Eq.(10), the expression can be simplified to:

$$\delta_{b^y} = \frac{\partial g \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)} \cdot \frac{\partial \left(\sum_j^m h_j \cdot W_{jk}^2 + b_k^y \right)}{\partial b_k^y};$$

$$\therefore \delta_{b^y} = g'(\omega_k) \cdot (1) = g'(\omega_k); \quad (40)$$

Finally, replacing Eq.(40) into Eq.(38), and the latter into Eq.(37):

$$\Delta b_k^y = \alpha \cdot (d_k - y_k) \cdot g'(\omega_k). \quad (41)$$

Derivation of ΔW_{ij}^1

To derive ΔW_{ij}^1 , first, Eq.(15) is plugged into Eq.(23):

$$\Delta W_{ij}^1 = -\alpha \cdot \frac{\partial \left(\sum_k^p \frac{(d_k - y_k)^2}{2} \right)}{\partial W_{ij}^1}. \quad (42)$$

If Eq.(12) is plugged into the previous expression and the sum expanded:

$$\Delta W_{ij}^1 = -\alpha \cdot \left(\frac{\partial \left(\frac{\left(d_1 - g \left(b_1^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{j1}^2 \right) \right)^2}{2} \right)}{\partial W_{ij}^1} + \dots + \frac{\partial \left(\frac{\left(d_k - g \left(b_k^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{jk}^2 \right) \right)^2}{2} \right)}{\partial W_{ij}^1} + \dots + \frac{\partial \left(\frac{\left(d_p - g \left(b_p^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{jp}^2 \right) \right)^2}{2} \right)}{\partial W_{ij}^1} \right) \quad (43)$$

Because each of the components of the sum \sum_k^p has a W_{ij}^1 term nested inside, none of the components is zero after differentiation by ∂W_{ij}^1 ; thus, the summation remains. Consequently, proceeding with the differentiation:

$$\Delta W_{ij}^1 = \alpha \cdot \sum_k^p (d_k - y_k) \cdot \frac{\partial y_k}{\partial W_{ij}^1}. \quad (44)$$

If δ_{W^1} is defined as:

$$\delta_{W^1} = \frac{\partial y_k}{\partial W_{ij}^1}, \quad (45)$$

Then, by substituting Eq.(12) into the previous expression:

$$\delta_{W^1} = \frac{\partial g \left(b_k^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{jk}^2 \right)}{\partial W_{ij}^1}, \quad (46)$$

To derive the previous expression, Eq.(4) is introduced to use the chain rule, such that:

$$\begin{aligned} \delta_{W^1} &= \frac{\partial g \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)}{\partial h_j} \cdot \frac{\partial h_j}{\partial W_{ij}^1} \\ \delta_{W^1} &= \frac{\partial g \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)}{\partial \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)} \cdot \frac{\partial \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)}{\partial h_j} \cdot \frac{\partial f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)} \cdot \frac{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial W_{ij}^1} \\ \therefore \delta_{W^1} &= g'(\omega_k) \cdot W_{jk}^2 \cdot f'(\psi_j) \cdot x_i, \end{aligned} \quad (47)$$

Where:

$$f'(\psi_j) = \frac{\partial f(\psi_j)}{\partial \psi_j} = \frac{\partial f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}. \quad (48)$$

Finally, replacing Eq.(47) into Eq.(45), and the latter into Eq.(44) and rearranging the expression:

$$\Delta W_{ij}^1 = \alpha \cdot f'(\psi_j) \cdot x_i \cdot \sum_k^p (d_k - y_k) \cdot g'(\omega_k) \cdot W_{jk}^2. \quad (49)$$

Derivation of Δb_j^h

The derivation of Δb_j^h is somewhat similar to that of ΔW_{ij}^1 . Plugging Eq.(15) into Eq.(25):

$$\Delta b_j^h = -\alpha \cdot \frac{\partial \left(\sum_k^p \frac{(d_k - y_k)^2}{2} \right)}{\partial b_j^h}, \quad (50)$$

And differentiating by ∂b_j^h :

$$\Delta b_j^h = \alpha \cdot \sum_k^p (d_k - y_k) \cdot \frac{\partial y_k}{\partial b_j^h}. \quad (51)$$

If δ_{b^h} is defined as:

$$\delta_{b^h} = \frac{\partial y_k}{\partial b_j^h}, \quad (52)$$

Then, by substituting Eq.(12) into the previous expression:

$$\delta_{b^h} = \frac{\partial g \left(b_k^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{jk}^2 \right)}{\partial b_j^h}, \quad (53)$$

To derive the previous expression, Eq.(4) is introduced to use the chain rule, such that:

$$\delta_{b^h} = \frac{\partial g \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)}{\partial h_j} \cdot \frac{\partial h_j}{\partial b_j^h} \cdot \frac{\partial f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)} \cdot \frac{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial b_j^h}$$

$$\delta_{b^h} = \frac{\partial g \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)}{\partial \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)} \cdot \frac{\partial \left(b_k^y + \sum_j^m h_j \cdot W_{jk}^2 \right)}{\partial h_j} \cdot \frac{\partial f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)} \cdot \frac{\partial \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right)}{\partial b_j^h}$$

$$\delta_{b^h} = g'(\omega_k) \cdot W_{jk}^2 \cdot f'(\psi_j) \cdot (1)$$

$$\therefore \delta_{b^h} = g'(\omega_k) \cdot W_{jk}^2 \cdot f'(\psi_j). \quad (54)$$

Finally, replacing Eq.(54) into Eq.(52), and the latter into Eq.(51) and rearranging the expression:

$$\Delta b_j^h = \alpha \cdot f'(\psi_j) \cdot \sum_k^p (d_k - y_k) \cdot g'(\omega_k) \cdot W_{jk}^2. \quad (55)$$

9 Conclusions

- In a forward pass throughout the FFNN, a output vector \mathbf{y} is generated, and its k^{th} component is described by Eq.(12):

$$y_k = g \left(b_k^y + \sum_j^m f \left(b_j^h + \sum_i^n x_i \cdot W_{ij}^1 \right) \cdot W_{jk}^2 \right)$$

In this equation, the sums \sum_j^m and \sum_i^n indicate that calculation of y_k takes into account all the weights and biases of the previous layers, which is what was proposed at the beginning of this FFNN — that previous layers must be fully interconnected.

- In the backpropagation pass, the effect of interconnecting layers is appreciated as the pass starts delving into the layers. For instance, the equations for ΔW_{jk}^2 and Δb_k^y (Eq.(35) and Eq.(41), respectively), depend only on the k^{th} component of \mathbf{y} , while the equations for ΔW_{ij}^1 and Δb_j^h (Eq.(49) and Eq.(55), respectively), depend on all the components of \mathbf{y} and only on the j^{th} component of \mathbf{h} . If there were more hidden layers, these dependencies would keep adding up until the last weights and bias matrices previous to the input layer.

Appendix A

This Appendix is under construction.

Appendix B

This Appendix is under construction.