×

-

(http://play.google.com/store/apps/details?id=com.analyticsvidhya.android)

(https://www.analyticsvidhya.com/datahack-
summit-2018/?utm_source=AVStickybanner2)

Analytics Vidhya
Learn everything about analytics

(https://www.analyticsvidhya.com/blog/)

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?
utm_source=AVTopBanner)

DEEP LEARNING (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/DEEP-LEARNING/)

PYTHON (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/PYTHON-2/)     R (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/R/)

# Understanding and coding Neural Networks From Scratch in Python and R

**SUNIL RAY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/AUTHOR/SUNIL-RAY/), MAY 29, 2017**

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?
utm_source=AVBlogbwBanner2)

## Introduction

You can learn and practice a concept in two ways:

- **Option 1:** You can learn the entire theory on a particular subject and then look for ways to apply those concepts. So, you read up how an entire algorithm works, the maths behind it, its assumptions, limitations and then you apply it. Robust but time taking approach.

- **Option 2:** Start with simple basics and develop an intuition on the subject. Next, pick a problem and start solving it. Learn the concepts while you are solving the problem. Keep tweaking and improving your understanding. So, you read up how to apply an algorithm – go out and apply it. Once you know how to apply it, try it around with different parameters, values, limits and develop an understanding of the algorithm.

I prefer Option 2 and take that approach to learning any new topic. I might not be able to tell you the entire math
(https://www.analyticsvidhya.com/datahack-
behind an algorithm, but I can tell you the intuition. I can tell you the best scenarios to apply an algorithm based
summit-2018/?utm_source=AVStickybanner2)
on my experiments and understanding.

In my interactions with people, I find that people don't take time to develop this intuition and hence they struggle to apply things in the right manner.

In this article, I will discuss the building block of a neural network from scratch and focus more on developing this intuition to apply Neural networks. We will code in both "Python" and "R". By end of this article, you will understand how Neural networks work, how do we initialize weigths and how do we update them using back-propagation.

Let's start.

## Table of Contents:

## Simple intuition behind neural networks

If you have been a developer or seen one work – you know how it is to search for bugs in a code. You would fire various test cases by varying the inputs or circumstances and look for the output. The change in output provides you a hint on where to look for the bug – which module to check, which lines to read. Once you find it, you make the changes and the exercise continues until you have the right code / application.

Neural networks work in very similar manner. It takes several input, processes it through multiple neurons from multiple hidden layers and returns the result using an output layer. This result estimation process is technically known as "Forward Propagation".

Next, we compare the result with actual output. The task is to make the output to neural network as close to actual (desired) output. Each of these neurons are contributing some error to final output. How do you reduce the error?

We try to minimize the value/ weight of neurons those are contributing more to the error and this happens while traveling back to the neurons of the neural network and finding where the error lies. This process is known as "Backward Propagation" (https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)
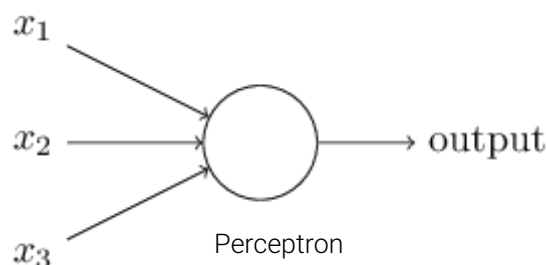
In order to reduce these number of iterations to minimize the error, the neural networks use a common algorithm known as "Gradient Descent", which helps to optimize the task quickly and efficiently.

That's it – this is how Neural network works! I know this is a very simple representation, but it would help you understand things in a simple manner.

## Multi Layer Perceptron and its basics

Just like atoms form the basics of any material on earth – the basic forming unit of a neural network is a perceptron. So, what is a perceptron?

A perceptron can be understood as anything that takes multiple inputs and produces one output. For example, look at the image below.



Perceptron

The above structure takes three inputs and produces one output. The next logical question is what is the relationship between input and output? Let us start with basic ways and build on to find more complex ways.

Below, I have discussed three ways of creating input output relationships:

1. **By directly combining the input and computing the output** based on a threshold value. for eg: Take x1=0, x2=1, x3=1 and setting a threshold =0. So, if x1+x2+x3>0, the output is 1 otherwise 0. You can see that in this case, the perceptron calculates the output as 1.
2. **Next, let us add weights to the inputs.** Weights give importance to an input. For example, you assign w1=2, w2=3 and w3=4 to x1, x2 and x3 respectively. To compute the output, we will multiply input with respective weights and compare with threshold value as w1*x1 + w2*x2 + w3*x3 > threshold. These weights assign more importance to x3 in comparison to x1 and x2.
3. **Next, let us add bias:** Each perceptron also has a bias which can be thought of as how much flexible the perceptron is. It is somehow similar to the constant *b* of a linear function *y = ax + b. It allows us to move the line up and down to fit the prediction with the data better. Without b the line will always goes through*

*the origin (0, 0) and you may get a poorer fit.* For example, a perceptron may have two inputs, in that case, it requires three weights. One for each input and one for the bias. Now linear representation of input will look like, w1*x1 + w2*x2 + w3*x3 + 1*b.

But, all of this is still linear which is what perceptrons used to be. But that was not as much fun. So, people thought of evolving a perceptron to what is now called as artificial neuron. A neuron applies non-linear transformations (activation function) to the inputs and biases.
(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

## What is an activation function?

Activation Function takes the sum of weighted input (w1*x1 + w2*x2 + w3*x3 + 1*b) as an argument and return the output of the neuron.

$$a = f(\sum_{i=0}^{N} w_i x_i)$$

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28091327/eq1-neuron.png)In above equation, we have represented 1 as x0 and b as w0.

The activation function is mostly used to make a non-linear transformation which allows us to fit nonlinear hypotheses or to estimate the complex functions. There are multiple activation functions, like: "Sigmoid", "Tanh", ReLu and many other.

## Forward Propagation, Back Propagation and Epochs

Till now, we have computed the output and this process is known as "Forward Propagation". But what if the estimated output is far away from the actual output (high error). In the neural network what we do, we update the biases and weights based on the error. This weight and bias updating process is known as "Back Propagation".

Back-propagation (BP) algorithms work by determining the loss (or error) at the output and then propagating it back into the network. The weights are updated to minimize the error resulting from each neuron. The first step in minimizing the error is to determine the gradient (Derivatives) of each node w.r.t. the final output. To get a mathematical perspective of the Backward propagation, refer below section.

This one round of forward and back propagation iteration is known as one training iteration aka "Epoch".

## Multi-layer perceptron

Now, let's move on to next part of Multi-Layer Perceptron. So far, we have seen just a single layer consisting of 3 input nodes i.e x1, x2 and x3 and an output layer consisting of a single neuron. But, for practical purposes, the single-layer network can do only so much. An MLP consists of multiple layers called Hidden Layers stacked in

between the Input Layer and the Output Layer as shown below.



(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/26094834/Screen-Shot-2017-05-26-at-9.47.51-AM.png)

The image above shows just a single hidden layer in green but in practice can contain multiple hidden layers. Another point to remember in case of an MLP is that all the layers are fully connected i.e every node in a layer(except the input and the output layer) is connected to every node in the previous layer and the following layer.

Let's move on to the next topic which is training algorithm for a neural network (to minimize the error). Here, we will look at most common training algorithm known as Gradient descent (https://www.analyticsvidhya.com/blog/2017/03/introduction-to-gradient-descent-algorithm-along-its-variants/).

## Full Batch Gradient Descent and Stochastic Gradient Descent

Both variants of Gradient Descent perform the same work of updating the weights of the MLP by using the same updating algorithm but the difference lies in the number of training samples used to update the weights and biases.

Full Batch Gradient Descent Algorithm as the name implies uses all the training data points to update each of the weights once whereas Stochastic Gradient uses 1 or more(sample) but never the entire training data to update the weights once.

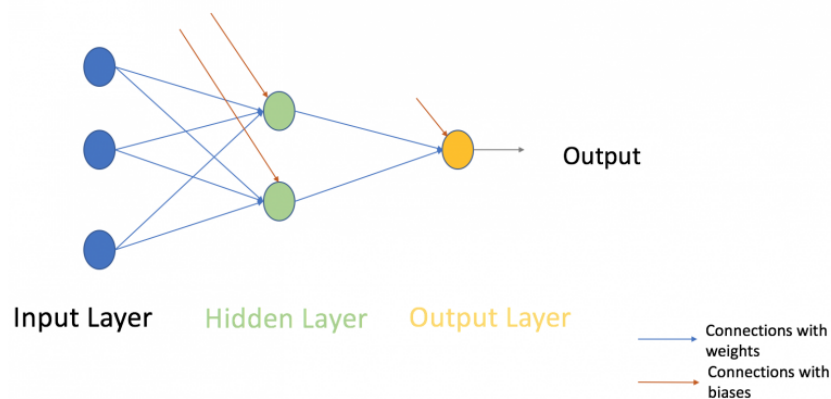Let us understand this with a simple example of a dataset of 10 data points with two weights w1 and w2.

Full Batch: You use 10 data points (entire training data) and calculate the change in w1 (Δw1) and change in w2(Δw2) and update w1 and w2.

SGD: You use 1st data point and calculate the change in w1 (Δw1) and change in w2(Δw2) and update w1 and w2. Next, when you use 2nd data point, you will work on the updated weights

For a more in-depth explanation of both the methods, you can have a look at this article (https://www.analyticsvidhya.com/blog/2017/03/introduction-to-gradient-descent-algorithm-along-its-variants/).

(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

## Steps involved in Neural Network methodology



(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/26094834/Screen-Shot-2017-05-26-at-9.47.51-AM.png)

Let's look at the step by step building methodology of Neural Network (MLP with one hidden layer, similar to above-shown architecture). At the output layer, we have only one neuron as we are solving a binary classification problem (predict 0 or 1). We could also have two neurons for predicting each of both classes.

First look at the broad steps:

0.) We take input and output

- X as an input matrix
- y as an output matrix

1.) We initialize weights and biases with random values (This is one time initiation. In the next iteration, we will use updated weights, and biases). Let us define:

- wh as weight matrix to the hidden layer
- bh as bias matrix to the hidden layer
- wout as weight matrix to the output layer
- bout as bias matrix to the output layer

2.) We take matrix dot product of input and weights assigned to edges between the input and hidden layer then add biases of the hidden layer neurons to respective inputs, this is known as linear transformation:

$$hidden\_layer\_input = matrix\_dot\_product(X, wh) + bh$$

3) Perform non-linear transformation using an activation function (Sigmoid). Sigmoid will return the output as $1/(1 + exp(-x))$.

(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

$$hiddenlayer\_activations = sigmoid(hidden\_layer\_input)$$

4.) Perform a linear transformation on hidden layer activation (take matrix dot product with weights and add a bias of the output layer neuron) then apply an activation function (again used sigmoid, but you can use any other activation function depending upon your task) to predict the output

$$output\_layer\_input = matrix\_dot\_product\ (hiddenlayer\_activations * wout\ ) + bout$$
$$output = sigmoid(output\_layer\_input)$$

All above steps are known as "Forward Propagation"

5.) Compare prediction with actual output and calculate the gradient of error (Actual – Predicted). Error is the mean square loss = $((Y-t)^2)/2$

$$E = y - output$$

6.) Compute the slope/ gradient of hidden and output layer neurons ( To compute the slope, we calculate the derivatives of non-linear activations x at each layer for each neuron). Gradient of sigmoid can be returned as $x * (1 - x)$.

$$slope\_output\_layer = derivatives\_sigmoid(output)$$
$$slope\_hidden\_layer = derivatives\_sigmoid(hiddenlayer\_activations)$$

7.) Compute change factor(delta) at output layer, dependent on the gradient of error multiplied by the slope of output layer activation

$$d\_output = E * slope\_output\_layer$$

8.) At this step, the error will propagate back into the network which means error at hidden layer. For this, we will take the dot product of output layer delta with weight parameters of edges between the hidden and output layer (wout.T).

$$Error\_at\_hidden\_layer = matrix\_dot\_product(d\_output, wout.Transpose)$$

9.) Compute change factor(delta) at hidden layer, multiply the error at hidden layer with slope of hidden layer activation

$$d\_hiddenlayer = Error\_at\_hidden\_layer * slope\_hidden\_layer$$

10.) Update weights at the output and hidden layer: The weights in the network can be updated from the errors calculated for training example(s).

$$wout = wout + matrix\_dot\_product(hiddenlayer\_activations.Transpose, d\_output)*learning\_rate$$

$$wh = wh + matrix\_dot\_product(X.Transpose,d\_hiddenlayer)*learning\_rate$$

learning_rate: The amount that weights are updated is controlled by a configuration parameter called the learning rate)

11.) Update biases at the output and hidden layer: The biases in the network can be updated from the aggregated errors at that neuron.

- bias at output_layer =bias at output_layer + sum of delta of output_layer at row-wise * learning_rate
- bias at hidden_layer =bias at hidden_layer + sum of delta of output_layer at row-wise * learning_rate

$$bh = bh + sum(d\_hiddenlayer, axis=0) * learning\_rate$$
$$bout = bout + sum(d\_output, axis=0)*learning\_rate$$

Steps from 5 to 11 are known as "Backward Propagation"

One forward and backward propagation iteration is considered as one training cycle. As I mentioned earlier, When do we train second time then update weights and biases are used for forward propagation.

Above, we have updated the weight and biases for hidden and output layer and we have used full batch gradient descent algorithm.


## Visualization of steps for Neural Network methodology

We will repeat the above steps and visualize the input, weights, biases, output, error matrix to understand working methodology of Neural Network (MLP).

Note:

- For good visualization images, I have rounded decimal positions at 2 or3 positions.
- Yellow filled cells represent current active cell
- Orange cell represents the input used to populate values of current cell

Step 0: Read input and output

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | | | | | | | | | | | | | | | | 1 | |
| 1 | 0 | 1 | 1 | | | | | | | | | | | | | | | | 1 | |
| 0 | 1 | 0 | 1 | | | | | | | | | | | | | | | | 0 | |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/26185640/0.0NN.jpg)

Step 1: Initialize weights and biases with random values (There are methods to initialize weights and biases but for now initialize with random values)
(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | | | | | | | 0.30 | 0.69 | | 1 | |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | | | | | | | 0.25 | | | 1 | |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | | | | | | | 0.23 | | | 0 | |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28120922/Screen-Shot-2017-05-28-at-12.06.49-PM.png)

Step 2: Calculate hidden layer input:
hidden_layer_input= matrix_dot_product(X,wh) + bh

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | | | | 0.30 | 0.69 | | 1 | |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | | | | 0.25 | | | 1 | |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | | | | 0.23 | | | 0 | |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28120936/Screen-Shot-2017-05-28-at-12.07.23-PM.png)Step 3: Perform non-linear transformation on hidden linear input
*hiddenlayer_activations = sigmoid(hidden_layer_input)*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | | 1 | |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | | 1 | |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | | 0 | |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28120948/Screen-Shot-2017-05-28-at-12.07.32-PM.png)Step 4: Perform linear and non-linear transformation of hidden layer activation at output layer

output_layer_input = matrix_dot_product (hiddenlayer_activations * wout ) + bout
*output = sigmoid(output_layer_input)*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | 0.79 | 1 | |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121001/Screen-Shot-2017-05-28-at-12.07.41-PM.png)Step 5: Calculate gradient of Error(E) at output layer

(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

*E = y-output*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121014/Screen-Shot-2017-05-28-at-12.07.48-PM.png)Step 6: Compute slope at output and hidden layer

*Slope_output_layer= derivatives_sigmoid(output)*

*Slope_hidden_layer = derivatives_sigmoid(hiddenlayer_activations)*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

| Slope hidden layer | | |
|---|---|---|
| 0.15 | 0.12 | 0.19 |
| 0.08 | 0.11 | 0.14 |
| 0.15 | 0.14 | 0.17 |

| Slope Output |
|---|
| 0.17 |
| 0.16 |
| 0.17 |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121030/Screen-Shot-2017-05-28-at-12.07.57-PM.png)Step 7: Compute delta at output layer

*d_output = E * slope_output_layer*lr*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

| Slope hidden layer | | |
|---|---|---|
| 0.15 | 0.12 | 0.19 |
| 0.08 | 0.11 | 0.14 |
| 0.15 | 0.14 | 0.17 |

| Slope Output |
|---|
| 0.17 |
| 0.16 |
| 0.17 |

| E |
|---|
| 0.21 |
| 0.20 |
| -0.79 |

| delta output |
|---|
| 0.04 |
| 0.03 |
| -0.13 |

(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121043/Screen-Shot-2017-05-28-at-12.08.06-PM.png)

Step 8: Calculate Error at hidden layer

Error_at_hidden_layer = matrix_dot_product(d_output, wout.Transpose)

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

| Slope hidden layer | | | error at hidden layer | | |
|---|---|---|---|---|---|
| 0.15 | 0.12 | 0.19 | 0.010 | 0.009 | 0.008 |
| 0.08 | 0.11 | 0.14 | 0.010 | 0.008 | 0.008 |
| 0.15 | 0.14 | 0.17 | -0.039 | -0.033 | -0.031 |

| Slope Output |
|---|
| 0.17 |
| 0.16 |
| 0.17 |

| E |
|---|
| 0.21 |
| 0.20 |
| -0.79 |

| delta output |
|---|
| 0.04 |
| 0.03 |
| -0.13 |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121057/Screen-Shot-2017-05-28-at-12.08.14-PM.png)Step 9: Compute delta at hidden layer

*d_hiddenlayer = Error_at_hidden_layer * slope_hidden_layer*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|------|------|------|------|------|------|------|------|------|------|------|------|------|------|--------|---|------|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.30 | 0.69 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.52 | | | | | | | | | | | | | | |

| Slope hidden layer | | | error at hidden layer | | |
|------|------|------|--------|--------|--------|
| 0.15 | 0.12 | 0.19 | 0.010 | 0.009 | 0.008 |
| 0.08 | 0.11 | 0.14 | 0.010 | 0.008 | 0.008 |
| 0.15 | 0.14 | 0.17 | -0.039 | -0.033 | -0.031 |

| Slope Output |
|------|
| 0.17 |
| 0.16 |
| 0.17 |

| E |
|------|
| 0.21 |
| 0.20 |
| -0.79 |

(https://www.analyticsvidhya.com/datahack-summit-2018/?utm_source=AVStickybanner2)

| delta hidden layer | | |
|--------|--------|--------|
| 0.002 | 0.001 | 0.002 |
| 0.001 | 0.001 | 0.001 |
| -0.006 | -0.005 | -0.005 |

| delta output |
|------|
| 0.04 |
| 0.03 |
| -0.13 |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121113/Screen-Shot-2017-05-28-at-12.08.20-PM.png)

Step 10: Update weight at both output and hidden layer

wout = wout + matrix_dot_product(hiddenlayer_activations.Transpose, d_output)*learning_rate
wh =  wh+ matrix_dot_product(X.Transpose,d_hiddenlayer)*learning_rate

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|------|------|------|------|------|------|------|------|------|------|------|------|------|------|--------|---|------|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.29 | 0.69 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.51 | | | | | | | | | | | | | | |

| Slope hidden layer | | | error at hidden layer | | |
|------|------|------|--------|--------|--------|
| 0.15 | 0.12 | 0.19 | 0.010 | 0.009 | 0.008 |
| 0.08 | 0.11 | 0.14 | 0.010 | 0.008 | 0.008 |
| 0.15 | 0.14 | 0.17 | -0.039 | -0.033 | -0.031 |

| Slope Output |
|------|
| 0.17 |
| 0.16 |
| 0.17 |

| E |
|------|
| 0.21 |
| 0.20 |
| -0.79 |

| Learning Rate | 0.1 |
|------|------|

| delta hidden layer | | |
|--------|--------|--------|
| 0.002 | 0.001 | 0.002 |
| 0.001 | 0.001 | 0.001 |
| -0.006 | -0.005 | -0.005 |

| delta output |
|------|
| 0.035 |
| 0.033 |
| -0.131 |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121126/Screen-Shot-2017-05-28-at-12.08.30-PM.png)

Step 11: Update biases at both output and hidden layer

*bh = bh + sum(d_hiddenlayer, axis=0) * learning_rate*
*bout = bout + sum(d_output, axis=0)*learning_rate*

| X | | | | wh | | | bh | | | hidden_layer_input | | | hidden_layer_activations | | | wout | bout | output | y | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0.42 | 0.88 | 0.55 | 0.46 | 0.72 | 0.08 | 1.48 | 1.78 | 1.10 | 0.81 | 0.86 | 0.75 | 0.29 | 0.68 | 0.79 | 1 | 0.21 |
| 1 | 0 | 1 | 1 | 0.10 | 0.73 | 0.68 | | | | 2.40 | 1.89 | 1.61 | 0.92 | 0.87 | 0.83 | 0.25 | | 0.80 | 1 | 0.20 |
| 0 | 1 | 0 | 1 | 0.60 | 0.18 | 0.47 | | | | 1.48 | 1.56 | 1.27 | 0.81 | 0.83 | 0.78 | 0.23 | | 0.79 | 0 | -0.79 |
| | | | | 0.92 | 0.11 | 0.51 | | | | | | | | | | | | | | |

| Slope hidden layer | | | error at hidden layer | | | Slope Output | E |
|---|---|---|---|---|---|---|---|
| 0.15 | 0.12 | 0.19 | 0.010 | 0.009 | 0.008 | 0.17 | 0.21 |
| 0.08 | 0.11 | 0.14 | 0.010 | 0.008 | 0.008 | 0.16 | 0.20 |
| 0.15 | 0.14 | 0.17 | -0.039 | -0.033 | -0.031 | 0.17 | -0.79 |

(https://www.analyticsvidhya.com/datahack-
summit-2018/?utm_source=AVStickybanner2)

Learning Rate   0.1

| delta hidden layer | | | delta output |
|---|---|---|---|
| 0.002 | 0.001 | 0.002 | 0.035 |
| 0.001 | 0.001 | 0.001 | 0.033 |
| -0.006 | -0.005 | -0.005 | -0.131 |

(https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2017/05/28121142/Screen-Shot-
2017-05-28-at-12.08.39-PM.png)

Above, you can see that there is still a good error not close to actual target value because we have completed
only one training iteration. If we will train model multiple times then it will be a very close actual outcome. I have
completed thousands iteration and my result is close to actual target values ([[ 0.98032096] [ 0.96845624] [
0.04532167]]).

# Implementing NN using Numpy (Python)

```
import numpy as np

#Input array
X=np.array([[1,0,1,0],[1,0,1,1],[0,1,0,1]])

#Output
y=np.array([[1],[1],[0]])

#Sigmoid Function
def sigmoid (x):
return 1/(1 + np.exp(-x))

#Derivative of Sigmoid Function
def derivatives_sigmoid(x):
return x * (1 - x)

#Variable initialization
epoch=5000 #Setting training iterations
lr=0.1 #Setting learning rate
```

```
inputlayer_neurons = X.shape[1] #number of features in data set
hiddenlayer_neurons = 3 #number of hidden layers neurons
output_neurons = 1 #number of neurons at output layer

 #weight and bias initialization
wh=np.random.uniform(size=(inputlayer_neurons,hiddenlayer_neurons))
bh=np.random.uniform(size=(1,hiddenlayer_neurons))
wout=np.random.uniform(size=(hiddenlayer_neurons,output_neurons))
bout=np.random.uniform(size=(1,output_neurons))

 for i in range(epoch):

     #Forward Propogation
     hidden_layer_input1=np.dot(X,wh)
     hidden_layer_input=hidden_layer_input1 + bh
     hiddenlayer_activations = sigmoid(hidden_layer_input)
     output_layer_input1=np.dot(hiddenlayer_activations,wout)
     output_layer_input= output_layer_input1+ bout
     output = sigmoid(output_layer_input)

     #Backpropagation
     E = y-output
     slope_output_layer = derivatives_sigmoid(output)
     slope_hidden_layer = derivatives_sigmoid(hiddenlayer_activations)
     d_output = E * slope_output_layer
     Error_at_hidden_layer = d_output.dot(wout.T)
     d_hiddenlayer = Error_at_hidden_layer * slope_hidden_layer
     wout += hiddenlayer_activations.T.dot(d_output) *lr
     bout += np.sum(d_output, axis=0,keepdims=True) *lr
     wh += X.T.dot(d_hiddenlayer) *lr
     bh += np.sum(d_hiddenlayer, axis=0,keepdims=True) *lr

 print output
```

## Implementing NN in R

```
# input matrix
X=matrix(c(1,0,1,0,1,0,1,1,0,1,0,1),nrow = 3, ncol=4,byrow = TRUE)

# output matrix
Y=matrix(c(1,1,0),byrow=FALSE)
```

```r
#sigmoid function
sigmoid<-function(x){
1/(1+exp(-x))
}

# derivative of sigmoid function
```

```r
derivatives_sigmoid<-function(x){
x*(1-x)
}

# variable initialization
epoch=5000
lr=0.1
inputlayer_neurons=ncol(X)
hiddenlayer_neurons=3
output_neurons=1

#weight and bias initialization
wh=matrix( rnorm(inputlayer_neurons*hiddenlayer_neurons,mean=0,sd=1), inputlayer_neurons, hiddenlaye
bias_in=runif(hiddenlayer_neurons)
bias_in_temp=rep(bias_in, nrow(X))
bh=matrix(bias_in_temp, nrow = nrow(X), byrow = FALSE)
wout=matrix( rnorm(hiddenlayer_neurons*output_neurons,mean=0,sd=1), hiddenlayer_neurons, output_neur

bias_out=runif(output_neurons)
bias_out_temp=rep(bias_out,nrow(X))
bout=matrix(bias_out_temp,nrow = nrow(X),byrow = FALSE)
# forward propagation
for(i in 1:epoch){

hidden_layer_input1= X%*%wh
hidden_layer_input=hidden_layer_input1+bh
hidden_layer_activations=sigmoid(hidden_layer_input)
output_layer_input1=hidden_layer_activations%*%wout
output_layer_input=output_layer_input1+bout
output= sigmoid(output_layer_input)

# Back Propagation

E=Y-output
slope_output_layer=derivatives_sigmoid(output)
slope_hidden_layer=derivatives_sigmoid(hidden_layer_activations)
```

```
d_output=E*slope_output_layer
Error_at_hidden_layer=d_output%*%t(wout)
d_hiddenlayer=Error_at_hidden_layer*slope_hidden_layer
wout= wout + (t(hidden_layer_activations)%*%d_output)*lr
bout= bout+rowSums(d_output)*lr
wh = wh +(t(X)%*%d_hiddenlayer)*lr
bh = bh + rowSums(d_hiddenlayer)*lr
}
output
```

## [Optional] Mathematical Perspective of Back Propagation Algorithm

Let $W_i$ be the weights between the input layer and the hidden layer. $W_h$ be the weights between the hidden layer and the output layer.

Now, $h=\sigma(u)=\sigma(W_iX)$, i.e h is a function of u and u is a function of $W_i$ and X. here we represent our function as $\sigma$

$Y=\sigma(u')=\sigma(W_hh)$, i.e Y is a function of u' and u' is a function of $W_h$ and h.

We will be constantly referencing the above equations to calculate partial derivatives.

We are primarily interested in finding two terms, $\partial E/\partial W_i$ and $\partial E/\partial W_h$ i.e change in Error on changing the weights between the input and the hidden layer and change in error on changing the weights between the hidden layer and the output layer.

But to calculate both these partial derivatives, we will need to use the chain rule of partial differentiation since E is a function of Y and Y is a function of u' and u' is a function of $W_i$.

Let's put this property to good use and calculate the gradients.

$\partial E/\partial W_h = (\partial E/\partial Y).(\partial Y/\partial u').(\partial u'/\partial W_h)$, ........(1)

We know E is of the form $E=(Y-t)^2/2$.

So, $(\partial E/\partial Y)=(Y-t)$

Now, $\sigma$ is a sigmoid function and has an interesting differentiation of the form $\sigma(1-\sigma)$. I urge the readers to work this out on their side for verification.

So, $(\partial Y/\partial u')=\partial(\sigma(u')/\partial u'=\sigma(u')(1-\sigma(u'))$.

But, σ(u')=Y, So,

(∂Y/∂u')=Y(1-Y)

Now, ( ∂u'/∂W$_h$)= ∂( W$_h$h)/ ∂W$_h$ = h

Replacing the values in equation (1) we get,
(https://www.analyticsvidhya.com/datahack-
summit-2018/?utm_source=AVStickybanner2)
∂E/∂W$_h$ = (Y-t). Y(1-Y).h

So, now we have computed the gradient between the hidden layer and the ouput layer. It is time we calculate the gradient between the input layer and the hidden layer.

∂E/∂W$_i$ =(∂ E/∂ h). (∂h/∂u).( ∂u/∂W$_i$)

But, (∂ E/∂ h) = (∂E/∂Y).( ∂Y/∂u').( ∂u'/∂h). Replacing this value in the above equation we get,

∂E/∂W$_i$ =[(∂E/∂Y).( ∂Y/∂u').( ∂u'/∂h)]. (∂h/∂u).( ∂u/∂W$_i$)……………(2)

So, What was the benefit of first calculating the gradient between the hidden layer and the output layer?

As you can see in equation (2) we have already computed ∂E/∂Y and ∂Y/∂u' saving us space and computation time. We will come to know in a while why is this algorithm called the back propagation algorithm.

Let us compute the unknown derivatives in equation (2).

∂u'/∂h = ∂(W$_h$h)/ ∂h = W$_h$

∂h/∂u = ∂( σ(u)/ ∂u= σ(u)(1- σ(u))

But, σ(u)=h, So,

(∂Y/∂u)=h(1-h)

Now, ∂u/∂W$_i$ = ∂(W$_i$X)/ ∂W$_i$ = X

Replacing all these values in equation (2) we get,

∂E/∂W$_i$ = [(Y-t). Y(1-Y).W$_h$].h(1-h).X

So, now since we have calculated both the gradients, the weights can be updated as

W$_h$ = W$_h$ + η . ∂E/∂W$_h$

W$_i$ = W$_i$ + η . ∂E/∂W$_i$

Where η is the learning rate.

So coming back to the question: Why is this algorithm called Back Propagation Algorithm?

The reason is: If you notice the final form of $\partial E/\partial W_h$ and $\partial E/\partial W_i$, you will see the term (Y-t) i.e the output error, which is what we started with and then propagated this back to the input layer for weight updation.

(https://www.analyticsvidhya.com/datahack-
So, where does this mathematics fit into the code?
summit-2018/?utm_source=AVStickybanner2)

hiddenlayer_activations=h

E= Y-t

Slope_output_layer = Y(1-Y)

lr = η

slope_hidden_layer = h(1-h)

wout = $W_h$

Now, you can easily relate the code to the mathematics.


# End Notes:

This article is focused on the building a Neural Network from scratch and understanding its basic concepts. I hope now you understand the working of a neural network like how does forward and backward propagation work, optimization algorithms (Full Batch and Stochastic gradient descent), how to update weights and biases, visualization of each step in Excel and on top of that code in python and R.

Therefore, in my upcoming article, I'll explain the applications of using Neural Network in Python and solving real-life challenges related to:

1. Computer Vision
2. Speech
3. Natural Language Processing

I enjoyed writing this article and would love to learn from your feedback. Did you find this article useful? I would appreciate your suggestions/feedback. Please feel free to ask your questions through comments below.