

Machine learning

Lecture 4 - Neural Networks

*Lecturer: Haim Permuter**Scribe: Nave Algarici*

Throughout this lecture we introduce Neural Networks, starting from a single neuron, and ending with the Backpropagation method. Most of the material for this lecture was inspired by the work of Michael Nielsen [1], and the UFLDL Stanford University Wiki [2].

I. NEURAL NETWORKS

Neural networks are limited imitations of how our own brains work. They've had a big recent resurgence because of advances in computer hardware. There is evidence that the brain uses only one "learning algorithm" for all its different functions. At a very simple level, neurons are basically computational units that take input (dendrites) as electrical input (called "spikes") that are channeled to outputs (axons).

Neural networks are typically organized in layers. Layers are made up of a number of interconnected 'nodes' which contain an 'activation function'. Patterns are presented to the network via the 'input layer', which communicates to one or more 'hidden layers' where the actual processing is done via a system of weighted 'connections'. The hidden layers then link to an 'output layer', which is the output of the network.

Neural Networks can be applied to many problems, such as: function approximation, classification, data processing, etc.

We will start by looking at a single neuron, define its model, and combine neurons to a complete network.

A. Single Neuron Model

A neuron takes k inputs, x_1, x_2, \dots, x_k , a set of weights w_1, w_2, \dots, w_k corresponding to the inputs, a bias b and an activation function $\sigma(z)$ and produces a single output y :

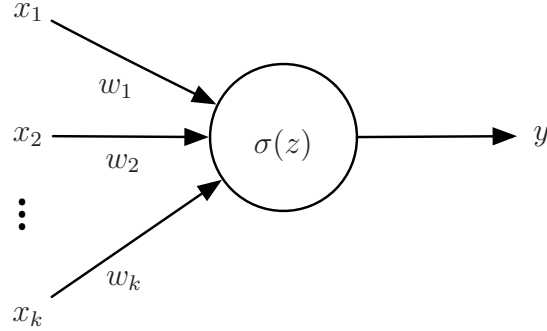


Figure 1. Scheme of a single neuron.

where the output of the neuron is determined by:

$$z = \sum_{i=1}^k w_i x_i + b, \quad (1)$$

$$y = \sigma(z). \quad (2)$$

There are many different options for the choice of the activation function. A few of them are:

- $\sigma(z) = \frac{1}{1+e^{-z}}$ - sigmoid function
- $\sigma(z) = \text{sign}(z)$ - sign function
- $\sigma(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ - hyperbolic tangent
- $\sigma(z) = \max(0, z)$ - rectified linear unit (RLU)

B. Neural Network Model

As we mentioned before, a Neural Network is organized in layers, where the first layer contains the inputs of the network, the last layer is the output of the network, and the layers in between are called hidden layers. Each layer gets its inputs from the layer before, and passes its outputs to the next. We call this step *forward propagation*.

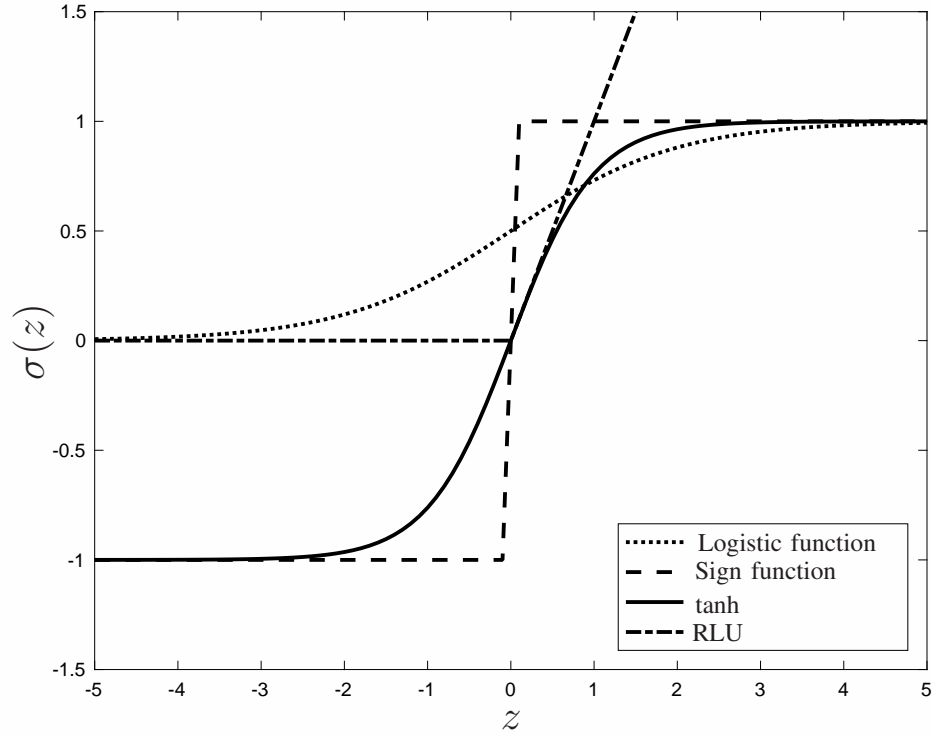


Figure 2. Plot of the activation functions mentioned above.

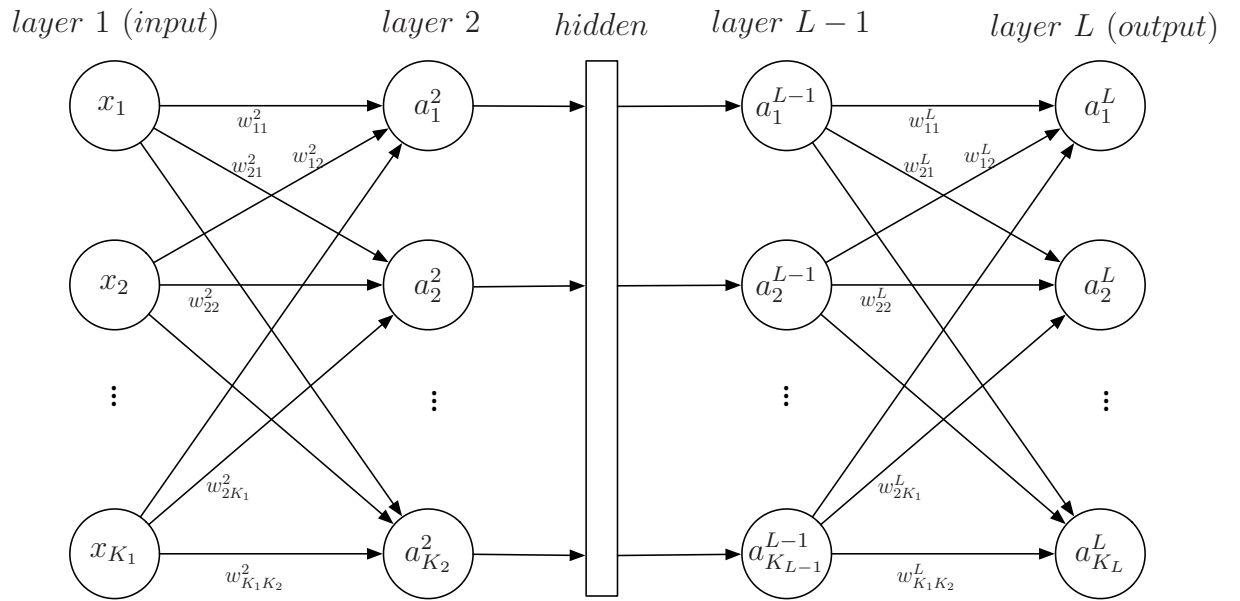


Figure 3. Structure of a general Neural Network

where:

- w_{jk}^l - the weight for the connection from the k^{th} neuron in the $(l-1)^{th}$ layer to the j^{th} neuron in the l^{th} layer.
- b_j^l - the coefficient we add to the j^{th} neuron in the l^{th} layer.
- K_l - the number of neurons in the l^{th} layer.
- $z_j^l = \sum_{k=1}^{K_l} w_{jk}^l a_k^{l-1} + b_j^l$
- $a_j^l = \sigma(z_j^l)$

By organizing our parameters in matrices and using matrix-vector operations, we can take advantage of fast linear algebra routines to quickly perform calculations in our network.

- $z^l = [z_1^l, z_2^l, \dots, z_{K_l}^l]^T$
- $b^l = [b_1^l, b_2^l, \dots, b_{K_l}^l]^T$
- $w^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \dots & w_{1K_{l-1}}^l \\ w_{21}^l & w_{22}^l & \dots & w_{2K_{l-1}}^l \\ \vdots & \vdots & \ddots & \vdots \\ w_{K_l1}^l & w_{K_l2}^l & \dots & w_{K_lK_{l-1}}^l \end{bmatrix}$
- $a^l = [a_1^l, a_2^l, \dots, a_{K_l}^l]^T$
- $z^l = w^l a^{l-1} + b^l$
- $a^l = \sigma(z^l)$
- $\sigma([x_1, \dots, x_n]^T) = [\sigma(x_1), \dots, \sigma(x_n)]^T$

Example 1 (XOR function) Consider a simple two input XOR function:

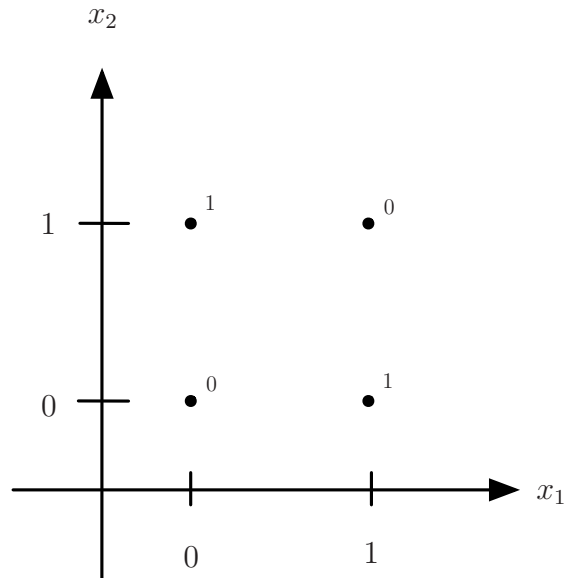


Figure 4. XOR function - illustration of output for binary inputs x_1, x_2

It is easy to see, that the XOR function cannot be approximated using a linear function (there is no line that could separate the two groups of answers). We now will show, that it is possible to approximate the XOR function using a Neural Network. We will build a Network with three inputs, x_1, x_2 , and the third set to '1', a hidden layer of two neurons, and a single neuron output layer. All biases are set to zero. Set the weights of the network as such:

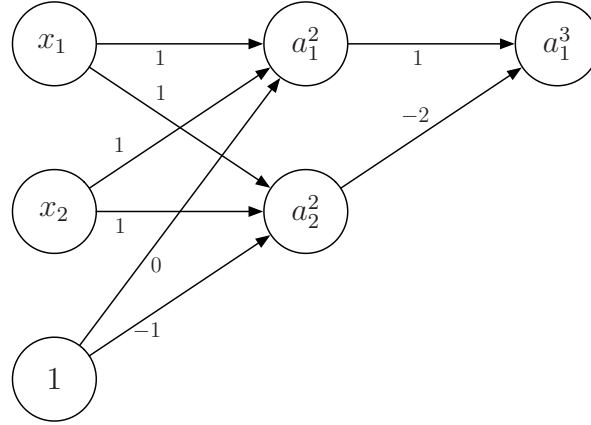


Figure 5. Example: XOR function Neural Network

For this example, the activation function will be the RLU function.

Let's take the input $[x_1, x_2] = [0, 0]$ and insert it to the network:

$$a_1^2 = \max(0, 1 * x_1 + 1 * x_2 + 0 * 1) = 0$$

$$a_2^2 = \max(0, 1 * x_1 + 1 * x_2 + -1 * 1) = 0$$

$$a_1^3 = \max(0, 1 * a_1^2 - 2 * a_2^2) = 0$$

The same goes for the other options for the inputs: $[0, 1]$, $[1, 0]$, $[1, 1]$:

x_1	x_2	a_1^2	a_2^2	a_1^3
0	0	0	0	0
0	1	1	0	1
1	0	1	0	1
1	1	2	1	0

We can see that the output a_1^3 fits the XOR function perfectly. Note that the output of each layer is not dependant of previous layers, only it's own inputs and weights.

C. Cost Function

To quantify how well we're achieving this goal we define a cost function:

$$C(w, b) = \frac{1}{2N} \sum_{i=1}^N \|a(x_i) - y_i\|^2, \quad (3)$$

where x are the input vectors and y are their corresponding labels, both are determined. So the cost function changes its value depending on the weights w and biases b . The cost function given in (3) is the quadratic cost function.

Our main goal is to minimize the cost function, so that the output from the network will be close to the desired output as possible. To minimize the cost function, we will use a method called *Gradient Descent*.

Gradient descent is an iterative optimization algorithm. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point.

D. Backpropagation

Backpropagation is about understanding how changing the weights and biases in a network changes the cost function. Ultimately, this means computing the partial derivatives $\frac{\partial C}{\partial w_{jk}^l}$ and $\frac{\partial C}{\partial b_j^l}$. But to compute those, we first introduce an intermediate quantity, δ_j^l , which we call the error in the j^{th} neuron in the l^{th} layer.

Backpropagation will give us a procedure to compute the error δ_j^l , and then will relate δ_j^l to $\frac{\partial C}{\partial w_{jk}^l}$ and $\frac{\partial C}{\partial b_j^l}$.

Our goal is to minimize C as a function of w and b . To train our neural network, we initialize each parameter w_{jk}^l and each b_j^l to a small random value near zero, and then apply an optimization algorithm such as batch gradient descent. Since C is a non-convex function, gradient descent is susceptible to local optima. However, in practice gradient descent usually works fairly well. Note that it is important to initialize the parameters randomly, rather than to all 0's. If all the parameters start off at identical values, then all the hidden layer units will end up learning the same function of the input. The random initialization serves the purpose of *symmetry breaking*.

One iteration of gradient descent updates the parameters w, b as follows:

$$w_{jk}^l = w_{jk}^l - \alpha \frac{\partial C}{\partial w_{jk}^l}, \quad (4)$$

$$b_j^l = b_j^l - \alpha \frac{\partial C}{\partial b_j^l}, \quad (5)$$

where α is the learning rate. The key step is computing the partial derivatives above. We will now describe the backpropagation algorithm, which gives an efficient way to compute these partial derivatives.

We define the error δ_j^l of neuron j in layer l by

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}, \quad (6)$$

starting with the L^{th} layer, we get

$$\delta_j^L = \frac{\partial C}{\partial z_j^L}. \quad (7)$$

Applying the chain rule, we can re-express the partial derivative above in terms of partial derivatives with respect to the output activations

$$\delta_j^L = \sum_k \frac{\partial C}{\partial a_k^L} \frac{\partial a_k^L}{\partial z_j^L}, \quad (8)$$

where the sum is over all neurons k in the output layer. Of course, the output activation a_k^L of the k^{th} neuron depends only on the input weight z_j^L for the j^{th} neuron when $k = j$. And so $\frac{\partial a_k^L}{\partial z_j^L}$ vanishes when $k \neq j$. As a result we can simplify the previous equation to

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L}. \quad (9)$$

Recalling that $a_j^L = \sigma(z_j^L)$, the second term on the right can be written as $\sigma'(z_j^L)$, and the equation becomes

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L). \quad (10)$$

We can rewrite the equation in a matrix-based form, as

$$\delta^L = \nabla_a C \odot \sigma'(z^L). \quad (11)$$

Here, $\nabla_a C$ is defined to be a vector whose components are the partial derivatives $\frac{\partial C}{\partial a_j^L}$.

We use \odot to denote the elementwise product of the two vectors.

Next, we'll develop the equation for the error δ^l in terms of the error in the next layer, δ^{l+1} . To do this, we want to rewrite $\delta_j^l = \frac{\partial C}{\partial z_j^l}$ in terms of $\delta_k^{l+1} = \frac{\partial C}{\partial z_k^{l+1}}$. We can do this using the chain rule:

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} \quad (12)$$

$$= \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} \quad (13)$$

$$= \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1}, \quad (14)$$

where in the last line we have interchanged the two terms on the right-hand side, and substituted the definition of δ_k^{l+1} . To evaluate the first term on the last line, note that

$$z_k^{l+1} = \sum_j w_{kj}^{l+1} a_j^l + b_k^{l+1} = \sum_j w_{kj}^{l+1} \sigma(z_j^l) + b_k^{l+1}. \quad (15)$$

Differentiating, we obtain

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \sigma'(z_j^l). \quad (16)$$

Substituting back into (14) we obtain

$$\delta_j^l = \sum_k w_{kj}^{l+1} \delta_k^{l+1} \sigma'(z_j^l). \quad (17)$$

In a matrix-based form,

$$\delta^l = \left((w^{l+1})^T \delta^{l+1} \right) \odot \sigma'(z^l), \quad (18)$$

where $(w^{l+1})^T$ is the transpose of the weight matrix w^{l+1} for the $(l+1)^{th}$ layer.

By combining (18) with (11) we can compute the error δ^l for any layer in the network. We start by using (11) to compute δ^L , then apply Equation (18) to compute δ^{L-1} , then Equation (18) again to compute δ^{L-2} , and so on, all the way back through the network.

Now that we have the errors δ_j^l of all the layers of the network, we can compute the partial derivatives $\frac{\partial C}{\partial w_{jk}^l}$ and $\frac{\partial C}{\partial b_j^l}$ as a function of δ_j^l :

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} = \delta_j^l a_k^{l-1}, \quad (19)$$

$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} = \delta_j^l. \quad (20)$$

For each iteration we use (4) and (5) and compute the new values of the parameters.

To summarize the backpropagation algorithm:

- 1) Perform a feedforward pass, computing the activations for layers 2,3, and so on up to the output layer L .

- 2) For each output unit j in layer L (the output layer), set

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L). \quad (21)$$

- 3) For layers $l = L - 1, L - 1, \dots, 2$, for each node j in layer l , set

$$\delta_j^l = \sum_k w_{kj}^{l+1} \delta_k^{l+1} \sigma'(z_j^l). \quad (22)$$

- 4) Compute the desired partial derivatives, which are given as:

$$\frac{\partial C}{\partial w_{jk}^l} = \delta_j^l a_k^{l-1}, \quad (23)$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l. \quad (24)$$

- 5) Update the weights and biases of the network:

$$w_{jk}^l = w_{jk}^l - \alpha \frac{\partial C}{\partial w_{jk}^l}, \quad (25)$$

$$b_j^l = b_j^l - \alpha \frac{\partial C}{\partial b_j^l}, \quad (26)$$

REFERENCES

- [1] Michael A. Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2015.
- [2] Andrew Ng, Jiquan Ngiam, Chuan Yu Foo, Yifan Mai, Caroline Suen. *UFLDL Stanford Tutorial*. http://ufldl.stanford.edu/wiki/index.php/UFLDL_Tutorial.