

Feed-forward neuronet recognising digits

Batarin Egor, Beksultanova Aikun, Zubritsky Ivan

April 21, 2021

Abstract

This pdf is done to descript the neuronet architecture and give step-by-step description of learning algorithms.

1 Accepted notations

$x \in \mathbb{R}^p$ - input value of neuronet, $p = \text{lenght} \times \text{width}$ - number of image pixels. x consists of 0 - "white color" and 1 - "black color".

$y = y(x) \in \mathbb{R}^{10}$ - desired output value. If input x responds to a digit $i \in \{0, 1, \dots, 9\}$, then $y(x)$ consists of nine zeros, and one, which is in a i -th place.

L - number of the last (output) layer.

$n : \{1, \dots, L\} \rightarrow \mathbb{N}$ - number of neurons in the l -th layer, $n(1) = p$, $n(L) = 10$.

$$\omega^l = \begin{pmatrix} w_{11}^l & w_{12}^l & \dots & w_{1,n(l-1)}^l \\ w_{21}^l & w_{22}^l & \dots & w_{2,n(l-1)}^l \\ \dots & \dots & \dots & \dots \\ w_{j1}^l & w_{j2}^l & \dots & w_{j,n(l-1)}^l \\ \dots & \dots & \dots & \dots \\ w_{n(l),1}^l & w_{n(l),2}^l & \dots & w_{n(l),n(l-1)}^l \end{pmatrix}$$

w^l is weight matrix, which consist of w_{jk}^l - weight between j -th neuron in l layer and k -th neuron in $l - 1$ layer. $l = 2, \dots, L$.

$b^l = (b_1^l, b_2^l, \dots, b_{n(l)}^l)^T$ is a vector of biases, b_j^l - bias of j -th neuron in the l -th layer. $l = 2, \dots, L$.

$a^l = (a_1^l, a_2^l, \dots, a_{n(l)}^l)^T$ is a vector of activations, b_j^l - bias of j -th neuron in the l -th layer, $x = a^1 \in \mathbb{R}^p$ - input value, $a^L \in \mathbb{R}^{10}$ - output value, predicted by neuronet. $l = 1, \dots, L$.

We define weighted input vector $z^l = (z_1^l, z_2^l, \dots, z_{n(l)}^l)$, where $z_j^l = \sum_{k=1}^{n(l-1)} \omega_{jk}^l a_k^{l-1} + b_j^l$. It can be written in vectorized form $z^l = \omega^l a^{l-1} + b^l$. $l = 2, \dots, L$.

$C_{\text{quad}}(x) = \frac{1}{2} \|y(x) - a^L(x)\| = \frac{1}{2} \sqrt{\sum_{i=1}^{10} (y_i(x) - a_i^L(x))^2}$ is quadratic cost, refered to a single input value x .

We will also use cross-entropy cost $C_{\log}(x) = - \sum_{i=1}^{10} y_i(x) \ln a_i^L(x) + (1 - y_i(x)) \ln (1 - a_i^L(x))$

$\sigma(z) = \frac{1}{1+e^{-z}}$ - sigmoid activation function, $\sigma'(z) = \sigma(z) (1 - \sigma(z))$. This funtion links activation and weighted input: $a^l = \sigma(z^l)$, $l = 2, \dots, L$.

For any cost function we define $\nabla_a C = \left(\frac{\partial C}{\partial z_1^L}, \dots, \frac{\partial C}{\partial z_{10}^L} \right)^T$, $C = \frac{1}{\text{number of inputs}} \sum_x C(x)$ - average cost.

For two vectors $a = (a_1, a_2, \dots, a_n)^T$ and $b = (b_1, b_2, \dots, b_n)^T$ the Hadamard product $a \odot b$ is defined as

$$a \odot b = (a_1 b_1, a_2 b_2, \dots, a_n b_n)^T$$

We define an error δ_j^l of j -th neuron in l -th layes as $\delta_j^l = \frac{\partial C}{\partial z_j^l}$, $\delta^l = (\delta_1^l, \dots, \delta_{n(l)}^l)^T$. $l = 2, \dots, L$.

2 Neuronet architecture

The structure of neuronet is presented in the picture below:

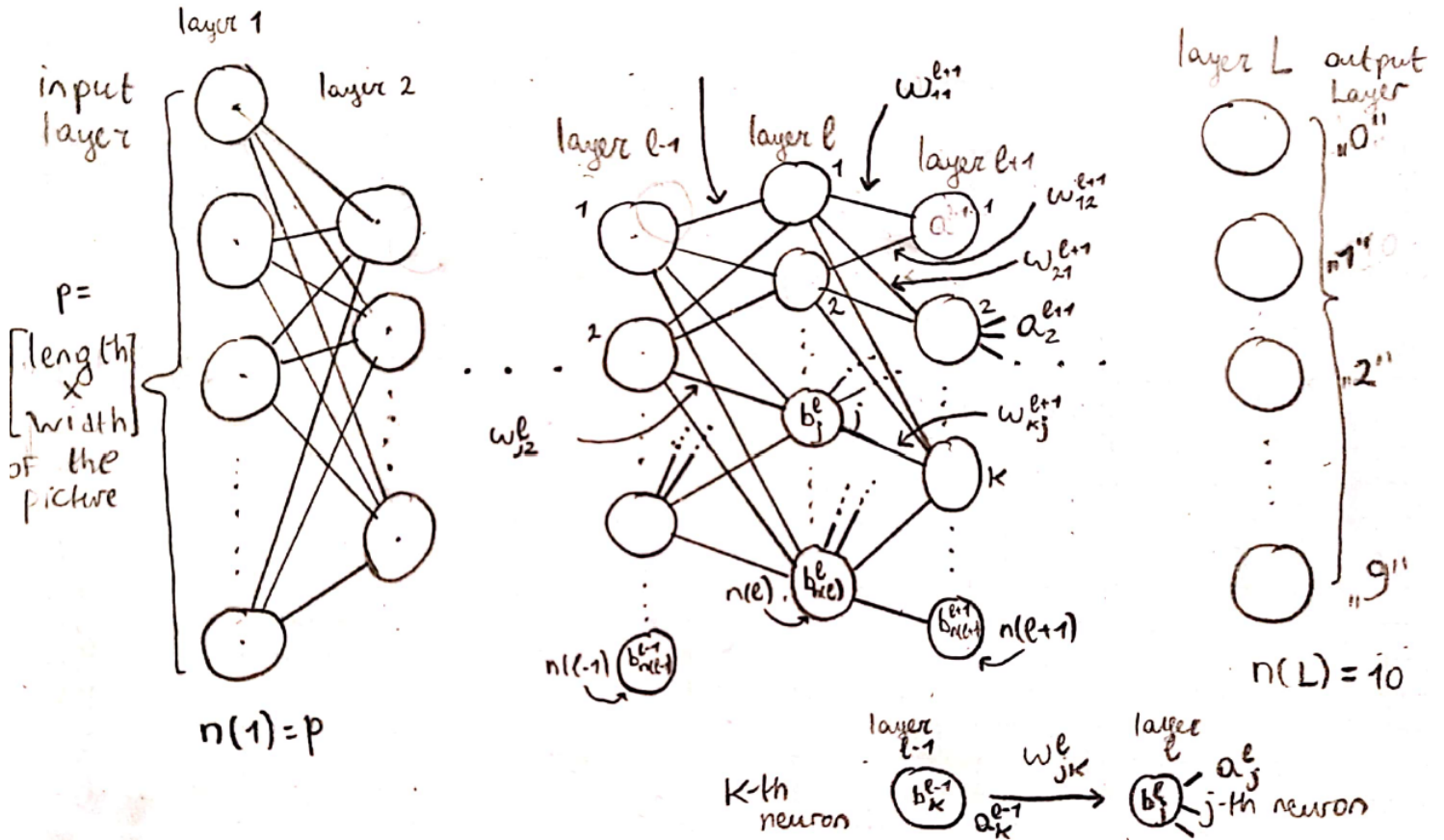


Figure 1: Neuronet structure

3 Backpropagation equations

Below four significant equations are presented, which will be used in backpropagation algorithm:

$$\delta^L = \nabla_a C \odot \sigma'(z^L) \quad (1)$$

$$\delta^l = \left((\omega^{l+1})^T \delta^{l+1} \right) \odot \sigma'(z^l) \quad (2)$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad (3)$$

$$\frac{\partial C}{\partial \omega_{jk}^l} = a_k^{l-1} \delta_j^l \quad (4)$$

4 Learning algorithm

1. **Set hyper-parameters:** learning rate η , number N of epoch of training, size m of mini-batch.
2. **Repeat for each of N -th epochs of training:**
 - (a) **Input a random mini-batch of m training examples x from training data.**
 - (b) **For each training example x :** Set the corresponding input activation $a^{x,1} = x$, and do the following:
 - i. **Feedforward:** For each $l = 2, 3, \dots, L$ compute $z^{x,l} = \omega^l a^{x,l-1} + b^l$ and $a^{x,l} = \sigma(z^{x,l})$
 - ii. **Output error $\delta_j^{x,L}$:** Compute the error vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$.
 - iii. **Backpropagate the error:** For each $l = L - 1, L - 2, \dots, 2$ compute $\delta^l = \left((\omega^{l+1})^T \delta^{l+1} \right) \odot \sigma'(z^l)$.
 - (c) **Gradient descent:** For each $l = L - 1, L - 2, \dots, 2$ update the weights: $\omega^l \rightarrow \omega^l - \frac{\eta}{m} \sum_x \delta^{x,l} (a^{x,l-1})^T$,
 $b^l \rightarrow b^l - \frac{\eta}{m} \sum_x \delta^{x,l}$.