

Prof. Dr. Jens Honore Walther  
Dr. Georgios Arampatzis  
ETH Zentrum, CLT  
CH-8092 Zürich

## Set 10

Issued: 21.05.2021

In this exercise, you will learn about the workhorse behind neural networks training, i.e. the Backpropagation Algorithm, and get hands-on experience with PyTorch within JupyterLab.

We will deal with the **multilayer perceptron** (MLP, also called feed-forward neural network). We assume that our neural network has  $L$  layers, each with dimensions  $\{n_0, \dots, n_L\}$  (defined by the user), where  $n_0$  denotes the dimension of the input, and  $n_L$  the dimension of the output. The input  $\mathbf{x} = (x_1, \dots, x_{n_0})$  is mapped to the first hidden layer via

$$a_j^1 = \sum_{i=1}^{n_0} W_{ij}^1 x_i, \quad \forall j \in \{1, \dots, n_1\}, \quad \text{and} \quad z_j^1 = \varphi_1(a_j^1), \quad (1)$$

where  $\varphi_1$  denotes the activation function (non-linearity), and  $z_j^1$  the output of the neuron  $j$  in the first layer. In similar fashion, the equations connecting consecutive hidden layers are given as such:

$$a_j^k = \sum_{i=1}^{n_{k-1}} W_{ij}^k z_i^{k-1}, \quad \forall j \in \{1, \dots, n_k\}, \quad \text{and} \quad z_j^k = \varphi_k(a_j^k). \quad (2)$$

where  $k$  denotes the  $k$ -th layer of the neural network. Finally, for the output layer we have

$$a_j^L = \sum_{i=1}^{n_{L-1}} W_{ij}^L z_i^{L-1}, \quad \forall j \in \{1, \dots, n_L\}, \quad \text{and} \quad y_j = \varphi_L(a_j^L). \quad (3)$$

In vector-matrix notation with  $\mathbf{w} = \{W^1, \dots, W^L\}$  this reads

$$\mathbf{y}(\mathbf{x}, \mathbf{w}) = \varphi_L \left( W^L \varphi_{L-1} \left( W^{L-1} \varphi_{L-2} \left( \dots W^2 \varphi_1 \left( W^1 \mathbf{x} \right) \right) \right) \right) \quad (4)$$

To train the network, the most commonly employed technique is **supervised learning** based on **Backpropagation**, where **training data**  $D = \{(\mathbf{x}_n, \hat{\mathbf{y}}_n)\}_{n=1, \dots, N}$  is used to learn a mapping from the input to the output which minimizes the mean squared error

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N |\hat{\mathbf{y}}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w})|^2 \quad (5)$$

In the following, we are going to explain how backpropagation works, and derive its mathematical expression.

## Question 1: Shedding light on Backpropagation

In this exercise, we are going to explain how the Backpropagation Algorithm works. Backpropagation is the workhorse of state-of-the-art deep neural network architectures. The name is a shortcut for **Backpropagation of Errors**. Backpropagation is a method that allows us to compute the **gradient** of the loss function with respect to the **weights** of the neural network. This gradient is expressed mathematically as  $\nabla_{\mathbf{w}} E$ , where  $\mathbf{w}$  denotes the weights and  $E$  represents the loss function. For simplicity, we will derive the backpropagation rules based on the loss of a single training sample  $E_n = |\hat{y}_n - y(\mathbf{x}_n, \mathbf{y})|^2$ . For the gradient, we need to compute the partial derivative  $\partial E_n / \partial W_{ij}^k$  for every weight  $W_{ij}^k$  of the neural network. Using the gradient, we can apply a gradient descent update for each weight, in order to move towards a direction that will minimize the loss.

The first step of the backpropagation algorithm is the **forward pass** (or **forward propagation**). In the forward pass, we compute the values  $a_j^k$  and  $z_j^k$  of each node  $j$  at each layer  $k$  of the neural network, as well as the final output  $y(\mathbf{x}_n, \mathbf{w})$ , according to equations (1-3), using a fixed input  $\mathbf{x}_n$ . Then, we need perform the **backward pass**, for which you will derive the formulas.

As a first step, we use the chain rule to rewrite the error  $\partial E_n / \partial W_{ij}^k$  as

$$\frac{\partial E_n}{\partial W_{ij}^k} = \frac{\partial E_n}{\partial a_j^k} \frac{\partial a_j^k}{\partial W_{ij}^k}, \quad (6)$$

where  $a_j^k$  is the output of node  $j$  at layer  $k$  before applying the activation function.

- Find a formula for the second factor in equation 6, i.e.  $\partial a_j^k / \partial W_{ij}^k$ , using equation (2).
- Next, we require an expression for the error gradient  $\delta_j^k$ . Again, we can use the chain rule to write

$$\delta_j^k = \frac{\partial E_n}{\partial a_j^k} = \sum_{i=1}^{n_{k+1}} \frac{\partial E_n}{\partial a_i^{k+1}} \frac{\partial a_i^{k+1}}{\partial a_j^k} = \sum_{i=1}^{n_{k+1}} \delta_i^{k+1} \frac{\partial a_i^{k+1}}{\partial a_j^k}, \quad (7)$$

where we decompose  $\delta_j^k$  to contributions made up of the error gradients of the next layer  $\delta_i^{k+1}$ . Your job is to find an expression for the term  $\partial a_i^{k+1} / \partial a_j^k$ , and then recover the error gradient  $\delta_j^k$ .

**Hint:** use the fact that  $z_j^k = \varphi^k(a_j^k)$  and  $a_i^{k+1}$  is a weighted sum of  $z_{j'}^k$  for  $j' \in \{1, \dots, n_k\}$ . The outcome includes the derivative of the activation function  $\varphi'_k$ .

- Now, plug the expressions you found for  $\delta_j^k$  and  $\partial a_j^k / \partial W_{ij}^k$  in (6) to write the loss function gradient with respect to the weights ( $\partial E_n / \partial W_{ij}^k$ ). Now assume you are using a gradient descent optimizer with step-size  $\eta$  to update the weights of the neural network. Formulate the weight update equation in terms of the weight update  $\Delta W_{ij}^k$ .
- Optional* There exist many activation functions that are commonly used in practice (ReLU, tanh, sigmoid, ...). Find one that satisfies the following expression:

$$f'(x) = f(x) \left( 1 - f(x) \right) \quad (8)$$

where  $f$  denotes the activation function.

- e) *Optional* You build a deep neural network with **tanh** activation functions, and observe that training converges slowly. After debugging, you find out that most **tanh** activation functions **saturate**. Justify this behavior based on the backpropagation algorithm and the form of the nonlinearity. Which activation function would you propose to alleviate the problem?