

# Lecture Notes 02: Feed-forward networks

**Tudor Berariu**  
tudor.berariu@gmail.com



October 12, 2016

## 1 Feed forward networks

*Feed forward networks* (also called *multi-layer perceptrons*) represent the fundamental neural architecture. Its computational graph is a directed acyclic graph that processes information from some input  $\mathbf{x}$  to an output  $\mathbf{y}$ . The information flows without any feedback loops.

$$\mathbf{y} = f_{\text{net}}(\mathbf{x}, \boldsymbol{\theta})$$

Feed forward networks represent a composition of functions that are usually called *layers*.

$$\mathbf{y} = f_L(f_{L-1}(\dots f_1(\mathbf{x}, \boldsymbol{\theta}_1) \dots, \boldsymbol{\theta}_{L-1}), \boldsymbol{\theta}_L)$$

## 2 Cost functions

### 2.1 Regression

In regression tasks the data set comprises of pairs of inputs and correct outputs which are real values. The usual cost function is the Mean Squared Error.

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_n^N \left| \mathbf{y}(\mathbf{x}^{(n)}, \boldsymbol{\theta}) - \mathbf{t}^{(n)} \right|_2^2 = \sum_n^N \sum_k^K \left( y_k(\mathbf{x}^{(n)}, \boldsymbol{\theta}) - t_k^{(n)} \right)^2$$

The derivative of the cost function w.r.t.  $\mathbf{y}$  is:

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \mathbf{y}} = \sum_n \mathbf{y}(\mathbf{x}^{(n)}, \boldsymbol{\theta}) - \mathbf{t}^{(n)}$$

### 2.2 Classification

For classification tasks the targets are one-hot encoded vectors with the single 1 value corresponding to the correct class.

One way to compute a classification model is to train a neural network to compute a posterior probability distribution over the classes:

$$y_k(\mathbf{x}, \boldsymbol{\theta}) \approx P(C_k | \mathbf{x})$$

Given a data set, our goal is to find the parameters  $\boldsymbol{\theta}^*$  that maximize the probability of the examples being in the correct class.

$$P(\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}^{(n)}, \boldsymbol{\theta})^{t_k^{(n)}}$$

Maximizing  $P(\boldsymbol{\theta})$  is equivalent to minimizing its negative logarithm.

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} -\log(P(\boldsymbol{\theta})) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} -\sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log(y_k(\mathbf{x}^{(n)}, \boldsymbol{\theta}))$$

The derivative of the cross-entropy function w.r.t the outputs is given by this formula:

$$\delta_{Lk} = \frac{\partial \mathcal{L}}{\partial y_k} = -\frac{t_k}{y_k} = \begin{cases} 0 & , \text{if } t_k = 0 \\ -1/y_k & , \text{if } t_k = 1 \end{cases} \quad (1)$$

### 3 Layers

The classic feed-forward networks alternate fully-connected linear layers with non-linear transfer functions (e.g. logistic, hyperbolic tangent). For regression problems the last layer does not need to go through a squash function. For classification tasks, the last layer is usually a softmax one, forcing the network to approximate a probability distribution.

#### 3.0.1 The fully connected layer

**Calculul ieirilor** A fully-connected layer computes a projection of an input vector  $\mathbf{x} \in \mathbb{R}^D$  into an output space  $\mathbb{R}^K$ . Formula 2 describes the computation for a single output unit  $y_k$ .

$$y_k = \sum_{i=1}^D \theta_{ki} x_i + b_k, \quad \forall k \in \{1, \dots, K\} \quad (2)$$

Formula 2 can be written in matrix form as in Formula 3.

$$\mathbf{y} = \boldsymbol{\Theta} \mathbf{x} + \mathbf{b} \quad (3)$$

The matrix  $\boldsymbol{\Theta} \in \mathbb{R}^{K \times D}$  and the vector  $\mathbf{b} \in \mathbb{R}^K$  are the parameters of the layer.

**Error backpropagation** In the error backpropagation phase, the partial derivatives of the loss function with respect to the inputs  $\vec{x}$  are being computed given  $\boldsymbol{\delta}_y = \frac{\partial \mathcal{L}}{\partial \mathbf{y}}$ . The computation follows Formula 4 or, the equivalent matrix expression in Formula 5.

$$\frac{\partial \mathcal{L}}{\partial x_i} = \sum_{k=1}^K \frac{\partial \mathcal{L}}{\partial y_k} \cdot \frac{\partial y_k}{\partial x_i} = \sum_{k=1}^K \delta_{y_k} \theta_{ki} = \boldsymbol{\delta}_y^T \boldsymbol{\theta}_i \quad (4)$$

$$\frac{\partial E}{\partial \mathbf{x}} = \boldsymbol{\theta}^T \boldsymbol{\delta}_y \quad (5)$$

In a similar fashion the gradient of the loss function with respect to the parameters  $\mathbf{x}\boldsymbol{\Theta}$  is  $\mathbf{b}$ :

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Theta}} = \boldsymbol{\delta}_y \mathbf{x}^T \quad (6)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \boldsymbol{\delta}_y \quad (7)$$

#### 3.0.2 Tanh

A **TanH** layer applies the hyperbolic tangent function element wise on an input vector.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (8)$$

$$\tanh'(x) = 1 - \tanh(x)^2 \quad (9)$$

**Computing outputs**

$$y_i = \tanh(x_i) \quad \forall i \quad (10)$$

**Backpropagating gradients** Working on Formula 9, the vector  $\delta_x$  can be computed using element-wise multiplication  $\odot$ :

$$\delta_x = (\mathbf{1} - \mathbf{y} \odot \mathbf{y}) \odot \delta_y \quad (11)$$

A TanH layer has no parameters.

### 3.0.3 SoftMax

**Computing outputs** The SoftMax layer is used in classification problems where the goal is to compute one-hot encoded output representations. The outputs are interpreted as a posterior distribution over the set of classes. Given some input vector  $\mathbf{x} \in \mathbb{R}^K$ , the output vector  $\mathbf{y} \in \mathbb{R}^K$  is computed as in Formula 12.

$$y_k = \frac{e^{x_k}}{\sum_{k'=1}^N e^{x_{k'}}} \quad (12)$$

**Propagating error** Since the SoftMax layer has no parameters, in the backward phase only  $\delta_x$  needs to be computed as a function of  $\delta_y$ . Formula 13 describes this relation. For a mathematical proof, go to Section A.

$$\delta_{x_k} = y_k (\delta_{y_k} - Z) \quad (13)$$

where  $Z = \sum_{k'} \delta_{y_{k'}} y_{k'}$ .

## 4 The forward phase

---

**Algorithm 1** The forward phase

---

```

1: procedure FORWARD(net, x)                                ▷ Inputs: a FFN net, and an example x
2:   y0 ← x
3:   for l ← 1 . . . L do
4:     yl ← net.layers[l].forward(yl-1)                ▷ lth layer's input are (l - 1)th layers's outputs
5:   return yL                                              ▷ Return the last layer's outputs.
```

---

## 5 The backward phase

---

**Algorithm 2** Backpropagation of gradients through the network

---

```

1: procedure BACKPROPAGATE(net, x, δL)
2:   y0 ← x
3:   for l ← L . . . 1 do
4:     δl-1 ← net.layers[l].backward(yl-1, δl)        ▷ Gradients are accumulated internally  $\frac{\partial E}{\partial \theta_l}$ 
```

---

## 6 Training the network

Parameters are updated using stochastic gradient descent:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \frac{\partial \mathcal{L}(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}^{(t)}} \quad (14)$$

---

### Algorithm 3 Stochastic gradient descent

---

```

1: procedure SGD( $net, \mathcal{D} = \{(\mathbf{x}^{(n)}, \mathbf{t}^{(n)})\}_{1 \leq n \leq N}, \eta$ ) ▷  $net$ , the data set, the learning rate
2:   repeat
3:      $\mathcal{B} \leftarrow$  a mini-batch of size  $b$  from  $\mathcal{D}$ 
4:     for  $l \leftarrow 1 \dots L$  do
5:        $y_l \leftarrow net.layers[l].zeroGradients()$  ▷ Set gradients to zero
6:       for  $(\mathbf{x}, \mathbf{t}) \leftarrow \mathcal{B}$  do ▷ For each example in the mini-batch
7:          $\mathbf{y} \leftarrow forward(net, \mathbf{x})$  ▷ Compute the net's outputs
8:          $\boldsymbol{\delta}_L \leftarrow \frac{\partial E}{\partial \mathbf{y}}$  ▷ The gradient of the loss w.r.t. the outputs
9:          $backpropagate(net, \mathbf{x}, \boldsymbol{\delta}_L)$  ▷ The errors are backpropagated through the network
10:      for  $l \leftarrow 1 \dots L$  do
11:         $net.layers[l].updateParameters(\eta)$  ▷ The parameters are updated
12:   until convergence

```

---

## A Backpropagating errors through a SoftMax layer

Consider the **SoftMax** function which transforms an input vector  $\mathbf{x}$  into an output vector  $\mathbf{y}$  with the same dimension which represents a probability distribution (Formulae 15,16).

$$\mathbf{y} = softmax(\mathbf{x}) \quad (15)$$

$$y_k = \frac{e^{x_k}}{\sum_{k'} e^{x_{k'}}} \quad (16)$$

Given the gradient of the loss function  $\mathcal{L}$  w.r.t.  $\mathbf{y}$ , we need to find the expression of the gradient w.r.t. the inputs  $\mathbf{x}$ .

$$\boldsymbol{\delta}_y \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{y}} \quad (17)$$

$$\boldsymbol{\delta}_x \stackrel{not.}{=} \frac{\partial E}{\partial \mathbf{x}} \quad (18)$$

Formula 19 describes the computation of a single component of  $\boldsymbol{\delta}_x$ .

$$\delta_{xk} = \sum_{k'} \frac{\partial E}{\partial y_{k'}} \frac{\partial y_{k'}}{\partial x_k} = \sum_{k'} \delta_{y_{k'}} \frac{\partial y_{k'}}{\partial x_k} \quad (19)$$

$$\frac{\partial y_{k'}}{\partial x_k} = \frac{\mathbb{I}(k == k') e^{x_k} (\sum_{k''} e^{x_{k''}}) - e^{x_{k'}} e^{x_k}}{(\sum_{k''} e^{x_{k''}})^2} = \begin{cases} y_k - y_k y_{k'} & , \text{daca } k == k' \\ -y_k y_{k'} & , \text{daca } k \neq k' \end{cases} \quad (20)$$

By using Formula 20 in Formula 19:

$$\delta_{xk} = \sum_{k'} \delta_{y_{k'}} \frac{\partial y_{k'}}{\partial x_k} = y_k \left( \delta_{y_k} - \sum_{k'} \delta_{y_{k'}} y_{k'} \right) = y_k (\delta_{y_k} - Z) \quad (21)$$

where  $Z = \sum_{k'} \delta_{y_{k'}} y_{k'}$ .