# Neural networks and the backpropagation algorithm

## S. Marchand-Maillet

Computer Science – University of Geneva

# 1 introduction

Given a function  $\mathbf{f}: \mathbb{R}^D \to \mathbb{R}^d$ , we wish to approximate (regress) it from training data  $\mathfrak{X} \times \mathfrak{Y} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  (where  $\mathbf{y}_i = f(\mathbf{x}_i)$ ) using a mapping

$$\begin{array}{cccc} \varphi_{\theta} : & \mathbb{R}^{D} & \rightarrow & \mathbb{R}^{d} \\ & x_{i} & \mapsto & \varphi_{\theta}(x_{i}) \end{array}$$

parameterized with parameters  $\theta$  (using notation from statistics).

In machine learning terms, f is the mapping that underlies the phenomenon we wish to learn and is unknown. We can access annotated samples  $\mathcal{X} \times \mathcal{Y} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  and we wish to train a learner  $\boldsymbol{\varphi}_{\boldsymbol{\theta}}$  by adjusting its parameters  $\boldsymbol{\theta}$  to minimize a loss function  $\mathcal{L}(\mathcal{X} \times \mathcal{Y}; \boldsymbol{\theta})$ .

Here, we present  $\phi_{\theta}(\cdot)$  as a feedforward neural network taking D input values and producing d output values. We will further consider the practical case of the squared loss:

$$\mathcal{L}(\mathcal{X} \times \mathcal{Y}; \boldsymbol{\theta}) = \frac{1}{2} \frac{1}{N} \sum_{i=1}^{N} \| \boldsymbol{\Phi}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i}) - \boldsymbol{y}_{i} \|_{2}^{2} = \frac{1}{2N} \sum_{i=1}^{N} \sum_{k=1}^{d} (\boldsymbol{\Phi}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i})_{[k]} - \boldsymbol{y}_{i}_{[k]})^{2}$$
(1)

Note. When the label  $y_i$  is of dimension d = 1, the squared loss reduces to the classical

$$\mathcal{L}(\mathcal{X} \times \mathcal{Y}; \boldsymbol{\theta}) = \frac{1}{2N} \sum_{i=1}^{N} (\widehat{y}_i - y_i)^2 = \frac{1}{2N} \sum_{i=1}^{N} (\varphi_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - y_i)^2$$

A Neural Network model is a Directed Acyclic Graph (DAG) of units (referred to as "neurons") organized in layers (see figure 2 for an example). In the particular case of a feedforward network, conventionally, connections between neurons only exist from one layer to the next.

## 1.1 "Neuron" model and notation

A generic neuron is located at layer  $l \in [\![L]\!]$  in a network of L layers and takes its inputs from the outputs of other neurons. To ease the presentation and because this is the dominant "feedforward" model, we assume that the neuron takes its inputs from the output of the neurons in the previous layer l-1 (see figure 1)<sup>(a)</sup>.

A neuron (l,k) located at index  $k \in [N_l]$  of layer  $l \in [L]$  comprising  $N_l$  neurons:

- 1. receives its inputs from weighted connections with weights  $w_{ik}^l$
- 2. aggregates these inputs into a linear sum into its activation  $a_k^l = \sum_{j=0}^{N_{l-1}} w_{jk}^l \varphi_k^{l-1}$
- 3. outputs the result of its activation function  $\phi_k^l(\cdot)$  applied to its activation:  $\phi_k^l = \phi_k^l(a_k^l)$

Note that in this model, the weights are attached to the neuron itself. They form the adjustable parameters of the model

$$\theta = \{w_{ik}^l \parallel j \in [N_{l-1}], k \in [N_l], l \in [L]\}$$

The hyperparameters of the model concern its structure and choice of activation functions, i.e L > 0,  $\{N_l\}_{l=1}^L$  and  $\{\varphi_l^l\}_{l=1}^{N_l,L}$ .

Also following the feedforward structure, we fix the following boundary conditions for the model:

$$\begin{cases} \varphi_k^0 = x_i \text{\tiny{[k]}} & N_0 = D & (\text{input}) \\ \varphi_\theta(x_i) = \varphi^L(x_i) & N_L = d & (\text{output}) \end{cases}$$

<sup>(</sup>a) A more generic formulation using the notation for graph adjacency is left as exercise.

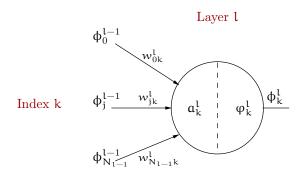


Figure 1: A generic neuron in a feedforward network, taking input from the previous layer

where the notation  $\phi_k^l(x_i)$  indicates the output of neuron (l,k) when the network takes  $x_i$  as input.  $a_k^l(x_i)$  is defined similarly.

Using variable homogenization, it is also customary to consider that  $\phi_0^1 = 1$  where a bias (hence  $w_0^1$ ) should be introduced in order to escape from the pure linear activation.

*Note.* The logistic neuron defines its activation functions as the logistic function:

$$\varphi_{k}^{l}(z) = \frac{1}{1 + e^{-z}} \quad \forall k \in [N_{l}] \ \forall l \in [L]$$

The present study is more general. It is left as exercise to model the logistic regression as a neural network.

As a result we have the following system of equations:

$$\begin{cases} \varphi_k^l = \varphi_k^l(\alpha_k^l) & \forall k \in \llbracket N_l \rrbracket & \forall l \in \llbracket L \rrbracket \\ \alpha_k^l = \sum_{j=0}^{N_{l-1}} w_{jk}^l \varphi_j^{l-1} & \forall k \in \llbracket N_l \rrbracket & \forall l \in \llbracket L \rrbracket \\ \varphi_k^0 = x_i [k] & \forall k \in \llbracket N_l \rrbracket & \forall l \in \llbracket L \rrbracket \\ \varphi_0^l = 1 & \forall k \in \llbracket D \rrbracket \\ \varphi_0^l = 1 & \forall l \in \llbracket L - 1 \rrbracket \\ \varphi_\theta(x_i) [k] = \varphi_k^L(x_i) & \forall k \in \llbracket d \rrbracket & \forall i \in \llbracket N \rrbracket \\ \mathcal{L}(\mathfrak{X} \times \mathfrak{Y}; \theta) = \frac{1}{2N} \sum_{i=1}^N \sum_{k=1}^d (\varphi_\theta(x_i) [k] - y_i [k])^2 \end{cases}$$

to adjust the parameters of the model  $\boldsymbol{\theta} = \left\{w_{jk}^{l} ~\|~ j \in [\![N_{l-1}]\!], k \in [\![N_{l}]\!], l \in [\![L]\!]\right\}$ 

Note. In this document, the assumption is that of a DAG. The connectivity may be more general than "feedforward" (e.g with "skip-connections" connecting layers that are not successive). Only the notation would change.

# 1.2 Training the network using the backpropagation algorithm

The aim of training is the minimization of the approximation (regression, prediction) error (the "loss"):

$$\theta^* = \operatorname*{argmin}_{\theta \in \Theta} \mathcal{L}(\mathfrak{X} \times \mathfrak{Y}; \theta)$$

As indicated, the parameters of the models are the weights  $\theta = \{w_{ik}^l\}_{j,k,l}$ .

The algorithm used for training is the gradient descent over the error function to adjust the parameters. In other words, the simple structure of the algorithm is as shown in algorithm 1.

The key element of the algorithm is therefore the computation of all values of

$$\frac{\partial \mathcal{L}}{\partial w_{jk}^{l}} \qquad \forall j \in [\![ N_{l-1} ]\!], k \in [\![ N_{l} ]\!], l \in [\![ L ]\!]$$

as the variation of the loss against the modification of weights  $w_{jk}^l$ .

#### Algorithm 1 The base structure of the backpropagation algorithm

```
1: procedure Backpropagation(\mathcal{X},\mathcal{Y})

2: w_{jk}^{l,(0)} \leftarrow \text{random}() \ \forall j \in \llbracket N_{l-1} \rrbracket, k \in \llbracket N_{l} \rrbracket, l \in \llbracket L \rrbracket \triangleright Initialize random weights

3: for iteration t = 1 until t = T do \triangleright or until convergence

4: for i \in \llbracket N \rrbracket do \triangleright Gradient iteration for data \mathbf{x}_i

5: for j \in \llbracket N_{l-1} \rrbracket, k \in \llbracket N_{l} \rrbracket, l \in \llbracket L \rrbracket do

6: w_{jk}^{l,(t)} \leftarrow w_{jk}^{l,(t-1)} - \frac{1}{N} \eta \frac{\partial \mathcal{L}}{\partial w_{jk}^{l}} (\mathbf{x}_i) \triangleright \eta > 0 is the gradient step (learning rate)
```

The name "backpropagation" for this algorithm comes from the realization that when the network has more than one layer (L>1), there is no direct access to  $\frac{\partial \mathcal{L}}{\partial w_{jk}^1}$ . The chain rule should be used and the gradient propagated back ("backpropagated") from the end of the network towards its input.

In the specific case of the squared loss, given  $p \in [N_{l-1}], r \in [N_l], m \in [L]$  we write

$$\frac{\partial \mathcal{L}}{\partial \mathcal{W}_{pr}^{m}} = \frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{d} \frac{\partial \varphi_{\theta}[\mathtt{k}]}{\partial \mathcal{W}_{pr}^{m}} (x_{i}) \left( \varphi_{\theta}(x_{i})_{[\mathtt{k}]} - y_{i}_{[\mathtt{k}]} \right)$$

where  $\varphi_{\theta}[k]$  is the  $k^{\text{th}}$  component function of  $\varphi_{\theta}(\cdot)$ . By definition,  $\varphi_{\theta}[k](x_i) = \varphi_k^L(x_i)$  so we can write

$$\frac{\partial \varphi_{\theta^{[k]}}}{\partial w_{pr}^{m}}(x_i) = \frac{\partial \varphi_k^L}{\partial w_{pr}^{m}}(x_i) \overset{\text{[1]}}{=} \left[ \frac{\partial \varphi_k^L}{\partial \alpha_k^L} \cdot \frac{\partial \alpha_k^L}{\partial w_{pr}^{m}} \right](x_i)$$

[1]: using the chain rule for derivation.

Now,

$$\frac{\partial \phi_k^l}{\partial w_{pr}^m} = \frac{\partial \phi_k^l}{\partial a_k^l} \cdot \frac{\partial a_k^l}{\partial w_{pr}^m} \tag{2}$$

$$\frac{\partial a_{\mathbf{k}}^{\mathbf{l}}}{\partial w_{\mathbf{pr}}^{\mathbf{m}}} = \sum_{\mathbf{j}=0}^{N_{\mathbf{l}-1}} w_{\mathbf{j}\mathbf{k}}^{\mathbf{l}} \frac{\partial \phi_{\mathbf{j}}^{\mathbf{l}-1}}{\partial w_{\mathbf{pr}}^{\mathbf{m}}}$$

$$\frac{\partial \phi_{\mathbf{k}}^{\mathbf{l}}}{\partial \phi_{\mathbf{k}}^{\mathbf{l}}} = \frac{\partial \phi_{\mathbf{k}}^{\mathbf{l}}}{\partial w_{\mathbf{pr}}^{\mathbf{m}}} \frac{\partial \phi_{\mathbf{k}}^{\mathbf{l}-1}}{\partial \phi_{\mathbf{k}}^{\mathbf{l}-1}}$$
(3)

$$\Rightarrow \frac{\partial \phi_{k}^{l}}{\partial w_{pr}^{m}} = \frac{\partial \phi_{k}^{l}}{\partial a_{k}^{l}} \cdot \sum_{j=0}^{N_{l-1}} w_{jk}^{l} \frac{\partial \phi_{j}^{l-1}}{\partial w_{pr}^{m}}$$

With the following particular cases:

$$\frac{\partial a_{r}^{m}}{\partial w_{pr}^{m}} = \frac{\partial}{\partial w_{pr}^{m}} \left[ \sum_{j=0}^{N_{m-1}} w_{jr}^{m} \varphi_{j}^{m-1} \right] = \varphi_{p}^{m-1}$$

and for p = 0,

$$\frac{\partial \alpha_r^m}{\partial w_{0r}^m} = \varphi_0^{m-1} = 1 \hspace{1cm} \forall m \in \llbracket L \rrbracket$$

or, if  $k \neq r$  or 1 < l < m,  $w_{pr}^{m}$  is not a involved in the computation of  $a_{k}^{l}$  and

$$\frac{\partial a_k^l}{\partial w_{pr}^m} = \frac{\partial}{\partial w_{pr}^m} \left[ \sum_{j=0}^{N_{l-1}} w_{jk}^l \phi_j^{l-1} \right] = 0 \tag{4}$$

finally, if  $\mathfrak{m} = 1$  and  $\mathfrak{p} \in [\![ D ]\!]$ ,

$$\frac{\partial a_{\mathbf{r}}^1}{\partial w_{\mathbf{pr}}^1}(\mathbf{x}_{\mathbf{i}}) = \frac{\partial}{\partial w_{\mathbf{pr}}^1} \left[ \sum_{j=0}^{N_0} w_{j\mathbf{r}}^1 \phi_j^0 \right] (\mathbf{x}_{\mathbf{i}}) = \phi_{\mathbf{p}}^0(\mathbf{x}_{\mathbf{i}}) = \mathbf{x}_{\mathbf{i}}[\mathbf{p}]$$

This set of equations (in particular (3) and (2)) creates a recursive chain of computation for derivatives, starting from layer L backpropagating values towards the first (input) layer. Note that equation (4) confirms that the backtracking of derivation of  $\frac{\partial a_k^l}{\partial w_{pr}^m}$  concern only paths in the network containing  $w_{pr}^m$  and stopping at layer m (equation (4)). We can therefore define the recursive algorithm 2 where the highlighted functions ensure the recursive backpropagation.

Clearly, the weights  $\{w_{jk}^l\}_{jkl}$  may be stored in 3D arrays. Similarly, the values for  $\{a_k^l(x_i)\}_{ikl}$  and  $\{\phi_k^l(x_i)\}_{ikl}$  may be stored using matrices (for every  $x_i$ ) or 3D arrays. Algebraic operations (addition, multiplication and more) can be defined over these so-called "tensors" and accelerated using GPUs (or TPUs).

```
Algorithm 2 The recursive backpropagation algorithm
```

```
1: function \frac{\partial \Phi_{k}^{l}}{\partial w_{pr}^{p}}(\mathbf{x}_{i})

2: return \frac{\partial \Phi_{k}^{l}}{\partial a_{k}^{l}} \cdot \frac{\partial a_{k}^{l}}{\partial w_{pr}^{m}}(\mathbf{x}_{i})
                                                                                                                                                                                                \triangleright \ \mathrm{Use \ chain \ rule \ since} \ \varphi_k^l = \phi_k^l(\alpha_k^l)
  3:
  4: function \frac{\partial \varphi_k^1}{\partial a_k^1}(x_i)
                                                                                                                                                                                                                                    ▶ Activation function
                Depends on l , k and the definition of \varphi_k^l(z)
                e.g if \varphi_k^l(z) = z return 1
e.g if \varphi_k^l(z) = z^2 return 2a_k^l(\mathbf{x}_i)
e.g if \varphi_k^l(z) = \frac{1}{1+e^{-z}} return \varphi(a_k^l(\mathbf{x}_i))(1-\varphi(a_k^l(\mathbf{x}_i)))
  6:
  7:
  9:
10: function \frac{\partial a_k^1}{\partial w_{pr}^m}(x_i)
                                                                                                                                                                                                                                                         ▶ Activation
                if l > m then return \sum_{j=0}^{N_{l-1}} w_{jk}^l \frac{\partial \phi_j^{l-1}}{\partial w_{pr}^m}(x_i)
11:
                if l = m then
12:
                         if k = r then return \varphi_p^{m-1}(x_i)
13:
                         else return 0
14:
                if l = 1 then return x_{i[p]}
15:
16:
17: function \frac{\partial \mathcal{L}}{\partial w_{jk}^1}(\mathbf{x}_i)
                \mathbf{return} \, \sum_{r=1}^{d} \frac{\frac{\partial \phi_{r}^{L}}{\partial w_{ik}^{L}}}{\frac{\partial \phi_{r}^{L}}{\partial w_{ik}^{L}}}(x_{i}) \left( \phi_{r}^{L}(x_{i}) - y_{ik}^{L} \right)
18:
19:
         procedure Backpropagation(\mathfrak{X},\mathfrak{Y})
20:
                \begin{aligned} w_{jk}^{l,(0)} \leftarrow \mathsf{random}() \ \forall j \in [\![N_{l-1}]\!], k \in [\![N_l]\!], l \in [\![L]\!] \\ \mathbf{for} \ \mathrm{iteration} \ t = 1 \ \mathrm{until} \ t = T \ \mathbf{do} \end{aligned}
                                                                                                                                                                                                                     ▶ Initialize random weights
21:
                                                                                                                                                                                                                                  ▷ or until convergence
22:
                         for i \in [N] do
23:
                                Forward pass computes and stores values for a_k^l(x_i) and \phi_k^l(x_i) \forall k, l
                                                                                                                                                                                                                                   \triangleright Forward pass for x_i
24:
                                for j \in [N_{l-1}], k \in [N_l], l \in [L] do w_{jk}^{l,(t)} \leftarrow w_{jk}^{l,(t-1)} - \frac{1}{N} \eta \frac{\partial \mathcal{L}}{\partial w_{jk}^{l}}(\mathbf{x}_i)
25:
                                                                                                                                                                                \triangleright \eta > 0 is the gradient step (learning rate)
26:
```

## 1.3 Example

#### 1.3.1 Quadratic regression

Given  $\mathfrak{X} \times \mathfrak{Y}$ , a quadratic regression considers the squared loss and a learner of the quadratic family:

$$\phi_{\theta}(\mathbf{x}_{i})_{[k]} = (\boldsymbol{\omega}_{2k}^{\mathsf{T}} \mathbf{x}_{i})^{2} + \boldsymbol{\omega}_{1k}^{\mathsf{T}} \mathbf{x}_{i} + \boldsymbol{\omega}_{0k}^{\mathsf{T}} \mathbf{1}_{\mathsf{D}}$$

$$\tag{5}$$

 $\phi_{\theta}$  may be represented with a 3-layer feedforward network with  $N_1=2, \ \phi_1^1(z)=z, \ \phi_2^1(z)=z^2, \ \phi_2^2(z)=z,$  as represented in figure 2.

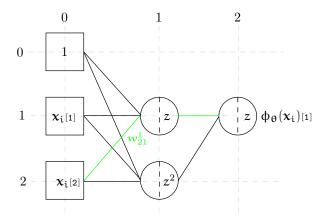


Figure 2: A neural network performing quadratic regression (D = 2, d = 1). Square nodes simply output their values. The green path shows the chain rule used for computing  $\frac{\partial \mathcal{L}}{\partial w_{21}^1}$ 

Following our notation, we obtain for the case  $D=2,\,d=1$ :

$$\phi_{\theta}(\mathbf{x}_{i}) = w_{11}^{2} (w_{01}^{1} + w_{11}^{1} \mathbf{x}_{i}_{[1]} + w_{21}^{1} \mathbf{x}_{i}_{[2]}) + w_{21}^{2} (w_{02}^{1} + w_{12}^{1} \mathbf{x}_{i}_{[1]} + w_{22}^{1} \mathbf{x}_{i}_{[2]})^{2}$$

$$= w_{11}^{2} \phi_{1}^{1} (a_{1}^{1}) + w_{21}^{2} \phi_{2}^{1} (a_{2}^{1})$$
(6)

Considering a squared loss, one obtains

$$\frac{\partial \mathcal{L}}{\partial w_{jk}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial \phi_{\theta}}{\partial w_{jk}^{l}} (\mathbf{x}_{i}) (\phi_{\theta}(\mathbf{x}_{i}) - \mathbf{y}_{i})$$

and

$$\begin{split} \frac{\partial \varphi_{\theta}}{\partial w_{01}^1}(\mathbf{x}_{\mathfrak{i}}) &= w_{11}^2 - \frac{\partial \varphi_{\theta}}{\partial w_{11}^1}(\mathbf{x}_{\mathfrak{i}}) = w_{11}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{i}}_{\mathfrak{I}} - \frac{\partial \varphi_{\theta}}{\partial w_{21}^1}(\mathbf{x}_{\mathfrak{i}}) = w_{11}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} \\ \frac{\partial \varphi_{\theta}}{\partial w_{02}^1}(\mathbf{x}_{\mathfrak{i}}) &= 2w_{21}^2 (w_{02}^1 + w_{12}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}}) + w_{22}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} ) \\ \frac{\partial \varphi_{\theta}}{\partial w_{02}^1}(\mathbf{x}_{\mathfrak{i}}) &= 2w_{21}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}}(\mathbf{x}_{\mathfrak{i}}) + w_{12}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} + w_{22}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} ) \\ \frac{\partial \varphi_{\theta}}{\partial w_{22}^1}(\mathbf{x}_{\mathfrak{i}}) &= 2w_{21}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}}(\mathbf{x}_{\mathfrak{i}}) + w_{12}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} + w_{22}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} ) \\ \frac{\partial \varphi_{\theta}}{\partial w_{21}^2}(\mathbf{x}_{\mathfrak{i}}) &= 2w_{21}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}}(\mathbf{x}_{\mathfrak{i}}) + w_{12}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} + w_{22}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} ) \\ \frac{\partial \varphi_{\theta}}{\partial w_{21}^2}(\mathbf{x}_{\mathfrak{i}}) &= 2w_{21}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}}(\mathbf{x}_{\mathfrak{i}}) + w_{12}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} + w_{22}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} ) \\ \frac{\partial \varphi_{\theta}}{\partial w_{21}^2}(\mathbf{x}_{\mathfrak{i}}) &= 2w_{21}^2 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}}_{\mathfrak{I}} + w_{12}^1 \mathbf{x}_{\mathfrak{i}}_{\mathfrak{I}}_{\mathfrak{I}} \end{pmatrix}$$

where one can see that the development of the chain rule follows paths according to equations derived in section 1.2. For example (see figure 2):

$$\frac{\partial \phi_{\theta}}{\partial w_{21}^1}(\mathbf{x}_i) = \frac{\partial \phi_1^2}{\partial a_1^2} \cdot \frac{\partial a_1^2}{\partial \phi_1^1} \cdot \frac{\partial \phi_1^1}{\partial a_1^1} \cdot \frac{\partial a_1^1}{\partial w_{21}^1}(\mathbf{x}_i) = 1 \cdot w_{11}^2 \cdot 1 \cdot \mathbf{x}_{i[2]}$$

These are the only non-zero terms from the chain

$$\frac{\partial \varphi_{\theta}}{\partial w_{21}^1}(\boldsymbol{x}_i) = \frac{\partial \varphi_1^2}{\partial \alpha_1^2} \cdot \sum_{j=0}^{N_1} \left[ \frac{\partial \alpha_1^2}{\partial \varphi_j^1} \cdot \frac{\partial \varphi_j^1}{\partial \alpha_j^1} \cdot \frac{\partial \alpha_j^1}{\partial w_{21}^1} \right] (\boldsymbol{x}_i)$$

since  $\frac{\partial a_2^1}{\partial w_{21}^1} = 0$  by equation (4).

Note. Developing equation (5) for this case, one can see that the neural model is over-specified: many combinations of weights lead to the same solution. To fix a unique solution, one can e.g fix and not update (clamp)  $w_{02}^1$ ,  $w_{11}^2$  and  $w_{21}^2$ .