# Formal Derivation of Mesh Neural Networks with Their Forward-Only Gradient Propagation

Federico A. Galatolo, Mario G. C. A. Cimino, Gigliola Vaglini

## This is a preprint. Please cite using:

# Formal derivation of Mesh Neural Networks with their Forward-Only gradient Propagation

Federico A. Galatolo · Mario G.C.A. Cimino · Gigliola Vaglini

**Abstract** This paper proposes the Mesh Neural Network (MNN), a novel architecture which allows neurons to be connected in any topology, to efficiently route information. In MNNs, information is propagated between neurons throughout a state transition function. State and error gradients are then directly computed from state updates without backward computation. The MNN architecture and the error propagation schema is formalized and derived in tensor algebra. The proposed computational model can fully supply a gradient descent process, and is potentially suitable for very large scale sparse NNs, due to its expressivity and training efficiency, with respect to NNs based on back-propagation and computational graphs.

**Keywords** Artificial Neural Networks · Gradients Computation · Supervised Learning · Deep Learning

## 1 Introduction and background

A huge amount of research has been made during the last years on a variety of applications of Artificial Neural Networks (ANNs). As a consequence, many ANNs architectures have been developed, generating surrogate models from different types of big data, such as image, audio, video, text, time series, and so on. With ANNs, the underlying relationships among data can be approximated with little knowledge of the system to be modelled. In spite of this success, ANNs are computational models vaguely inspired to biological brains, and require relevant computation and management with respect to the biological counterpart.

Specifically, Deep Learning is achieving good levels of performance, via architectures composed of several layers. The Deep Learning research is mostly based on gradient-based optimization methods and on the well-known *backpropagation* (BP) algorithm. In essence, BP includes a forward and backward layer-wise computation of the loss function with respect to the neurons weights. Actually, BP is not biologically plausible. Moreover, convergence problems, such as vanishing and exploding

Department of Information Engineering, University of Pisa, 56122 Pisa, Italy E-mail: federico.galatolo@ing.unipi.it, mario.cimino@unipi.it, gigliola.vaglini@unipi.it

gradients, occur when using many layers. Finally, BP can be very unstable when dealing with recurrent networks and can be ineffective to exploit long-lasting relationships [1]. In the last decade, an increasing number of alternative strategies have proposed to simplify the ANN training. A first strategy consists in removing the backward computation by deriving a forward only computation. A reference work for this approach is [2]. Specifically, the proposed method improves the efficiency of Jacobian matrix computation, for fully or partially connected ANNs. An interesting advantage of this approach is that it can train arbitrarily connected ANNs, and not just Multi-Layer Perceptron (MLP)-based architectures. Indeed, ANNs with connections across layers are much more powerful than MLPs. A more recent research in which the Jacobian matrix is calculated only in the forward computation was made by Guo *et al.* [3]. In general, to remove the backward computation is not costless: an additional calculation in the forward computation must be considered. However, the forward-only computation is more parallelizable than traditional forward and backward computations, as the dataset is large and the number of hidden neurons increases. A different strategy is presented in [4], in which the training method is based on a different principle called information bottleneck, which does not require backpropagation. In general, a performance comparison with BP is difficult, since performance can heavily depend on the minibatch size. The minibatch size is usually a constant that is based on available GPU memory. On the other side, a quantity of interest is the learning convergence, which is unknown for either BP or other methods. Since the backward computation is removed for such approaches, they are more suitable for parallel computation. Another type of strategy is proposed by Jaderberg [5]: a model for predicting gradient, called synthetic gradient, is calculated in place of true backpropagation error gradients. With such synthetic gradients, layers can be independently updated, removing forward and update locking.

According to this research trend, this paper formally introduces recent advances leading to a novel, arbitrarily connected, ANNs architecture, in which error gradients are computed throughout a state transition function without backward computation. The paper is organized as follows. In Section 2, the fundamentals of the problem are defined. A formal derivation of the proposed architecture is presented in Section 3. Section 4 covers the implementation and experimental aspects. Section 5 is devoted to conclusions and future work.

## 2 Problem statement

An Artificial Neurons Layer (ANL) with $n_i$ inputs and $n_o$ outputs can be described by its *layer weights matrix* $\mathbf{W} \in \mathbb{R}^{n_i \times n_o}$ and activation function $\hat{\varphi}(\mathbf{x}) : \mathbb{R}^{n_o} \to \mathbb{R}^{n_o}$. Let us consider activation functions for which it holds that $\hat{\varphi}(\mathbf{x})_i = \varphi(x_i)$ (where $\varphi(x) : \mathbb{R} \to \mathbb{R}$). Each column $\mathbf{W}_{*,i}$ of $\mathbf{W}$ represents the weights vector from the inputs to the $i$-th perceptron, in which biases are represented as weights of fictitious inputs that always produce the constant value 1. Given the input vector $\mathbf{x} \in \mathbb{R}^{n_i}$, the output vector $\mathbf{y} \in \mathbb{R}^{n_o}$ of the ANL is $\mathbf{y} = \varphi(\mathbf{x}\mathbf{W})$. In multilayer neural networks, or MLPs, ANLs are stacked, i.e., the ANL$_i$ is fed by the output of the ANL$_{i-1}$: each set of weights connecting the $i$-th layer is represented by a different matrix $\mathbf{W}_i$, and the input/output layers are considered as special topological elements with respect to the hidden layers.

In the popular BP training algorithm, the gradients of the weights are iteratively computed exploiting a propagation rule between layers [6, 7]. Let us consider a generic error function $E(\mathbf{y}, \overline{\mathbf{y}}) : \mathbb{R}^{n \times 2} \to \mathbb{R}$ that computes the error between a network output $\mathbf{y}$ and a desired one $\overline{\mathbf{y}}$, and a generic error function with respect to the $o$-th output $y_o$ $E_o(y_o, \overline{y_o}) : \mathbb{R}^2 \to \mathbb{R}$. Let us assume that $E(\mathbf{y}, \overline{\mathbf{y}})$ is a composition of $E_o(y_o, \overline{y_o})$ for every output unit. Considering an MLP with $n_l$ layers, the objective of the BP algorithm is to compute the gradients of every output error $\frac{\partial E(y_o, \overline{y_o})}{\partial p_i}$ with respect to every parameter $p_i$. Such gradients can be used by a *Stochastic Gradient Descent* (SGD) algorithm to train the MLP [8]. Let $net_{i,o}$ be the $o$-th output of the $i$-th hidden layer. Applying the chain rule for differentiating composite functions to $\frac{\partial E(y_o, \overline{y_o})}{\partial p_i}$, the corresponding error gradient is:

$$\frac{\partial E(y_o, \overline{y_o})}{\partial p_i} = \frac{\partial E(y_o, \overline{y_o})}{\partial y_o} \frac{\partial y_o}{\partial p_i} = \frac{\partial E(net_{L-1,o}, \overline{y_o})}{\partial net_{L-1,o}} \frac{\partial net_{L-1,o}}{\partial p_i}. \tag{1}$$

The derivative $\frac{\partial E(net_{n_l-1,o}, \overline{y_o})}{\partial net_{n_l-1,o}}$ depends on the error function and is known. In the derivative $\frac{\partial net_{n_l-1,o}}{\partial p_i}$, each parameter of a layer influences the output values of all the subsequent layers. Hence, in order to compute $\frac{\partial net_{n_l-1,o}}{\partial p_i}$, the chain rule is applied up to the term $\frac{\partial net_{i,o}}{\partial p_i}$. For this purpose, the BP algorithm iteratively applies the chain rule on each layer in reverse order for efficiently computing the partial derivatives with respect to all parameters. More formally, given the output of the $l$-th layer, $\mathbf{net_l} = \varphi(\mathbf{net_{l-1}}\mathbf{W}_l)$, let us say its $o$-th element $t_{l,o} = (\mathbf{net_{l-1}}\mathbf{W}_l)_o$. The chain rule is applied to $\varphi(t_{l,o})$, and in order to compute the term $\frac{\partial \varphi(t_{l,o})}{\partial t_{l,o}}$, $\mathbf{t_l}$ needs to be saved for each layer.

To train ANNs without a layered topology, the approach commonly used is the automatic differentiation on *computational graphs* (CGs) [9], in which computations are represented in a graph. In essence, for each operation (e.g., matrix multiplication, element-wise sum, etc.) the inputs $x_0, x_1, \cdots, x_{n-1}$ and the output $\mathbf{y}$ are represented as incoming and outgoing edges of a graph, respectively. For each edge $\frac{\partial y}{\partial x_i}$ is computed. For a given ANN, the operations to compute its output $y_o$ and the error $E(y_o, \overline{y_o})$ are then represented as a CG. Let us consider, a "factoring path", i.e., a path between two nodes in which the derivatives $\frac{\partial y}{\partial x_i}$ encountered on the traversed edges are all multiplied together. Then, the partial derivative of the error function with respect to a parameter, i.e., $\frac{\partial E(y_o, \overline{y_o})}{\partial p_i}$, is the sum of all the reverse factoring paths from $E(y_o, \overline{y_o})$ to $p_i$, i.e., the paths belonging to the set $\mathcal{P}_i$:

$$\frac{\partial E(y_o, \overline{y_o})}{\partial p_i} = \sum_{p \in \mathcal{P}_i} \prod_{(x,y) \in p} \frac{\partial y}{\partial x}. \tag{2}$$

A CG representation is a general formalism to represent all network topologies, such as feedforward, recurrent, convolutional, residual, and so on. To train arbitrarily connected ANNs topologies is very important, because ANNs with connections across layers are much more powerful than classical MLP architectures. However, a CG increases the space complexity with respect to a corresponding MLP-based representation (where an MLP representation is possible). Indeed, the underlying data structure needs to store both the graph topology and the partial derivatives $\frac{\partial y}{\partial x_i}$ of each edge. Moreover, it results in a higher time complexity, because all the reverse factoring paths have to be found.

In the next section, a novel ANNs representation is introduced, which is capable of training arbitrarily connected neural networks and, as a consequence, ANNs with reduced number of neurons and good generalization capabilities. The interesting properties of the training algorithm is the lack of a backpropagated computation, and an iteration without need of memory relationships than the one with the previous step. Hence, the proposed method is much simpler than traditional forward and backward procedure. Indeed, the training iteration can be described by three matrix operations. Due to the possibility of training unstructured ANNs, the proposed architectural model is called Mesh Neural Network (MNN).

## 3 Formal derivation of a Mesh Neural Network

3.1 Structure, activation and state of an MNN

The proposed MNN is based on a matrix representation that is not a transfer matrix, but it is an *adjacency matrix* (AM), i.e., a square matrix representing the ANN as a finite graph. The elements of the AM indicate whether pairs of vertices are adjacent or not in the graph, by means of a non-zero or zero weight, respectively.

More formally, an AM $\mathbf{A}$ is a matrix in which each element $A_{i,j}$ represents the weight from the node $i$ to the node $j$. For example MLPs are a subset of the representable topologies with AMs: since in MLPs only connections between layers are possible, their AMs are block matrices. Figure 1 shows an MLP topology with the corresponding AM. Here, each $\mathbf{W}_i$ is the weights matrix of the $i$-th layer and occupies a corresponding block in the AM.
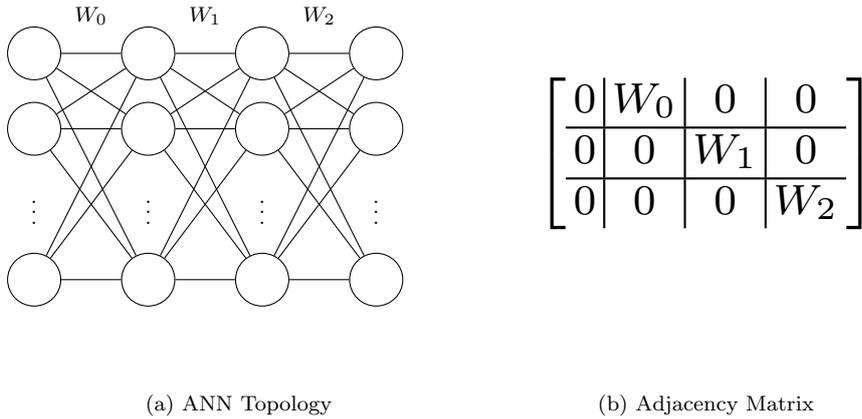


$$\begin{bmatrix} 0 & W_0 & 0 & 0 \\ 0 & 0 & W_1 & 0 \\ 0 & 0 & 0 & W_2 \end{bmatrix}$$

(a) ANN Topology                              (b) Adjacency Matrix

**Fig. 1** An MLP and its adjacency matrix

An example of unstructured topology and its corresponding AM is shown in Figure 2.



(a) ANN Topology

$$\begin{bmatrix} 0 & 0 & w_{0,2} & w_{0,3} & 0 & 0 & 0 & 0 & w_{0,8} & 0 \\ 0 & 0 & 0 & w_{1,3} & 0 & w_{1,5} & w_{1,6} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & w_{2,5} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & w_{3,9} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & w_{4,7} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & w_{5,9} \\ 0 & 0 & 0 & 0 & w_{6,5} & 0 & 0 & 0 & w_{6,8} & 0 \\ 0 & 0 & w_{7,2} & 0 & 0 & 0 & 0 & 0 & w_{7,8} & w_{7,9} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$
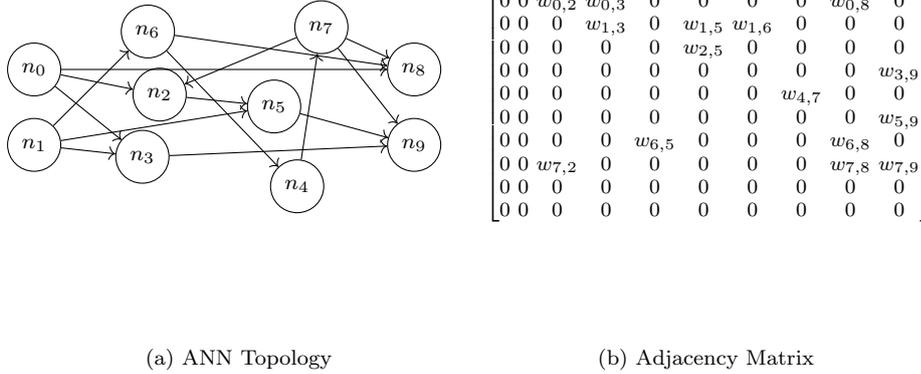
(b) Adjacency Matrix

**Fig. 2** An unstructured ANN and its adjacency matrix

A generic MNN topology with $n$ neurons is represented by a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$. It is worth noting that this representation does not include the topological distinction between input, hidden and output neurons. Let $n_i, n_o$, and $n_o$ be the number of input, hidden and output neurons. Since all neurons are identified by a position in the matrix, a good convention (hereinafter called "iho positioning convention") to distinguish the three sets without loss of generality is to assign them a positioning: to consider the first $n_i$ elements as input neurons, the subsequent $n_h$ elements as hidden neurons, and the last $n_o$ elements as output neurons.

Let the *state* be $\mathbf{s_i} \in \mathbb{R}^n$ the output value of each neuron in the MNN at the $i$-th instant of time. The output of an MNN is provided along a temporal sequence, whose length depends on the distances between input and output neurons. This allows an MNN to exhibit temporal dynamic behavior. Let us recall that: (i) $A_{i,j}$ represents the weight from neuron $i$ to neuron $j$; (ii) the $h$-th neuron output is computed as $\varphi(\sum_{k=0}^{N} w_{k,h} x_k)$; (iii) biases are represented as weights of fictitious inputs that always produce the constant value 1. Hence, given an initial state $\mathbf{s_0}$, which is set to the input value for input neurons and to zero for the other neurons, the next state is calculated as:

$$\mathbf{s_n} = \hat{\varphi}(\mathbf{s_{n-1}} \mathbf{A}) \tag{3}$$

At each time tick, the state transition of each neuron can influence the outputs values of all adjacent neurons. For subsequent ticks, the initial piece of information contained in $\mathbf{s_0}$ can traverse subsequent neurons and can influence their states, up to the output neurons.

3.2 Derivation of state and error gradients

In this section, the error derivative $\frac{\partial E(y,\overline{y})}{\partial p_i}$ for every parameter $p_i$ of an MNN are formally determined. It can be observed from Equation (3) that the unique parameter is $\mathbf{A}$. Let us assume an MNN with $n$ neurons, of which $n_i$ input neurons and $n_o$ output neurons positioned in the matrix according to the iho ordering convention. Let be the MNN processed for $t$ states. The $o$-th output value is then $y_o = s_{t-1,o} = \hat{\varphi}(\mathbf{s}_{t-2}\mathbf{A})_o$ where $o \in \{n - n_o, \cdots, n - 1\}$. Recalling the chain rule:

$$\frac{\partial E(y_o, \overline{y_o})}{\partial p_i} = \frac{\partial E(y_o, \overline{y_o})}{\partial y_o}\frac{\partial s_{t-1,o}}{\partial p_i}. \tag{4}$$

Let us consider a generic state $\mathbf{s_n} = \hat{\varphi}(\mathbf{t_n})$ where $\mathbf{t_n} = \mathbf{s_{n-1}}\mathbf{A}$. According to the chain rule, the derivative for a generic output $o$ is:

$$\frac{\partial s_{n,o}}{\partial A_{i,j}} = \frac{\partial \varphi(t_{n,o})}{\partial t_{n,o}}\frac{\partial t_{n,o}}{\partial A_{i,j}} = \frac{\partial \varphi(t_{n,o})}{\partial t_{n,o}}\frac{\partial(\mathbf{s_{n-1}}\mathbf{A})_o}{\partial A_{i,j}} \tag{5}$$

where $(\mathbf{s_{n-1}}\mathbf{A})_o$ is:

$$(\mathbf{s_{n-1}}\mathbf{A})_o = \sum_{k=0}^{N} s_{n-1,k}A_{k,o} \tag{6}$$

Let us distinguish two cases in Equation (6): (i) if $o = j$, one of the $A_{k,o}$ is $A_{i,j}$; (ii) if $o \neq j$, all the $A_{k,o}$ are constant with respect to $A_{i,j}$. Let us consider the case $o = j$. For linearity of differentiation:

$$\frac{\partial(\mathbf{s_{n-1}}\mathbf{A})_j}{\partial A_{i,j}} = \frac{\partial(\sum_{k=0}^{N} s_{n-1,k}A_{k,j})}{\partial A_{i,j}} = \sum_{k=0}^{N}\frac{\partial(s_{n-1,k}A_{k,j})}{\partial A_{i,j}} \tag{7}$$

In the partial derivatives $\frac{\partial(s_{n-1,k}A_{k,j})}{\partial A_{i,j}}$, all the $s_{n-1,k}$ elements depend on $A_{i,j}$. Moreover, in the case $k \neq i$, the matrix elements $A_{k,j}$ are constants with respect to $A_{i,j}$. Let us distinguish in Equation (7) the term with $k = i$:

$$\sum_{k=0}^{N}\frac{\partial(s_{n-1,k}A_{k,j})}{\partial A_{i,j}} = \sum_{k=0,\ k\neq i}^{N}\frac{\partial(s_{n-1,k}A_{k,j})}{\partial A_{i,j}} + \frac{\partial(s_{n-1,i}A_{i,j})}{\partial A_{i,j}} \tag{8}$$

Since $A_{k,j}$ is a constant, the first term of Equation (8) is:

$$\sum_{k=0,\ k\neq j}^{N}\frac{\partial(s_{n-1,k}A_{k,j})}{\partial A_{i,j}} = \sum_{k=0,\ k\neq j}^{N}\frac{\partial s_{n-1,k}}{\partial A_{i,j}}A_{k,j} \tag{9}$$

By applying the product rule to the second term of Equation (8):

$$\frac{\partial(s_{n-1,i}A_{i,j})}{\partial A_{i,j}} = \frac{\partial s_{n-1,i}}{\partial A_{i,j}}A_{i,j} + \frac{\partial A_{i,j}}{\partial A_{i,j}}s_{n-1,i} = \frac{\partial s_{n-1,i}}{\partial A_{i,j}}A_{i,j} + s_{n-1,i} \tag{10}$$

The term $\frac{\partial s_{n-1,i}}{\partial A_{i,j}} A_{i,j}$ can be integrated in the summation of Formula (9):

$$\sum_{k=0}^{N} \frac{\partial (s_{n-1,k} A_{k,j})}{\partial A_{i,j}} = \sum_{k=0}^{N} \frac{\partial s_{n-1,k}}{\partial A_{i,j}} A_{k,j} + s_{n-1,i} \qquad (11)$$

Similarly, considering the case $o \neq j$ in Equation (6), the $A_{k,o}$ elements are constant with respect to $A_{i,j}$, leading to:

$$\frac{\partial (\mathbf{s_{n-1}} \mathbf{A})_o}{\partial A_{i,j}} = \frac{\partial (\sum\limits_{k=0}^{N} s_{n-1,k} A_{k,o})}{\partial A_{i,j}} = \sum_{k=0}^{N} \frac{\partial s_{n-1,k}}{\partial A_{i,j}} A_{k,o} \qquad (12)$$

Hence, Equation (5) can be formulated as follows:

$$\frac{\partial s_{n,o}}{\partial A_{i,j}} = \begin{cases} \frac{\partial \varphi(t_{n,o})}{\partial t_{n,o}} \left( \sum\limits_{k=0}^{N} \frac{\partial s_{n-1,k}}{\partial A_{i,j}} A_{k,j} + s_{n-1,i} \right) & \text{if } o = j \\ \frac{\partial \varphi(t_{n,o})}{\partial t_{n,o}} \left( \sum\limits_{k=0}^{N} \frac{\partial s_{n-1,k}}{\partial A_{i,j}} A_{k,o} \right) & \text{if } o \neq j \end{cases} \qquad (13)$$

As a result, Equation (13) determines a very efficient algorithm for computing the partial derivative of the MNN state, which is, in turn, essential for applying an SGD-based training. In three terms: (i) the partial derivatives of the activation function $\frac{\partial \varphi(t_{n,o})}{\partial t_{n,o}}$, (ii) the previous states $s_{n-1,k}$, and (iii) the partial derivatives previous state $\frac{\partial s_{n-1,k}}{\partial A_{i,j}}$. Consequently, it is possible to compute both the next states $s_{n,o}$ and the next state partial derivatives $\frac{\partial s_{n,o}}{\partial A_{i,j}}$, concurrently and in the same iteration step. Moreover, an iteration does not need to store any intermediate values except for those of the current state, which can then be overwritten in the next iteration. Since the error gradient can be directly calculated from state gradient, Equation (4) results in a simplified iterative method without any memory dependency than the one with the previous step.

Operations in Equation (13) can be performed with scalars, vectors, and matrices, and then can be reformulated so as to be efficiently performed with tensors. In the next section, Equation (13) and the error gradient propagation schema are formalized and derived by tensor algebra.

### 3.2.1 Tensor Algebra formulation of the error gradient

Let us denote by $\nabla_{\mathbf{A}} \mathbf{s_n} \in R^{N \times N \times N}$ the tensor of the partial derivatives $\frac{\partial s_{n,o}}{\partial A_{i,j}}$

$$(\nabla_{\mathbf{A}} \mathbf{s_n})_{i,j,o} = \frac{\partial s_{n,o}}{\partial A_{i,j}} \qquad (14)$$

and by $\nabla_{\mathbf{x}} \varphi(\mathbf{x})$ the tensor of partial derivatives $\frac{\partial \varphi(x_i)}{\partial x_i}$

$$(\nabla_{\mathbf{x}} \varphi(\mathbf{x}))_i = \frac{\partial \varphi(x_i)}{\partial x_i} \qquad (15)$$

and by $\tilde{\mathsf{S}}_n \in R^{N \times N \times N}$ a tensor such that:

$$\tilde{\mathsf{S}}_{i,j,o} = \begin{cases} s_{n,i} & \text{if } o = j \\ 0 & \text{otherwise} \end{cases} \tag{16}$$

Hence, it is possible to formulate Equation 13 as:

$$\nabla_{\mathbf{A}}\mathbf{s_n} = \nabla_{\mathbf{t_n}}\varphi(\mathbf{t_n}) \odot (\nabla_{\mathbf{A}}\mathbf{s_{n-1}}\mathbf{A} + \tilde{\mathsf{S}}_n) \tag{17}$$

where the symbol $\odot$ denotes the Hadamard product.

As a result, the error gradient Forward-Only Propagation (FOP) algorithm of an MNN can be formulated in terms of the following steps, i.e., initialization, state derivatives forward propagation, and error derivative computation:

---

$\nabla_{\mathbf{A}}\mathbf{s} \leftarrow 0$
**for** $i$ **in** $\{1, 2, \cdots, n_t - 1\}$ **do**
    $\mathbf{s}_{0:n_i} \leftarrow x$
    $t \leftarrow \mathbf{s}\,\mathbf{A}$
    $\nabla_{\mathbf{A}}\mathbf{s} \leftarrow \nabla_{\mathbf{t}}\varphi(\mathbf{t}) \odot (\nabla_{\mathbf{A}}\mathbf{s}\,\mathbf{A} + \tilde{\mathsf{S}})$
    $\mathbf{s} \leftarrow \varphi(\mathbf{t})$
**end**
$\mathbf{y} \leftarrow \mathbf{s}[n - n_h : n]$
$\nabla_{\mathbf{A}}E(\mathbf{y}, \overline{\mathbf{y}}) \leftarrow \nabla_{\mathbf{y}}E(\mathbf{y}, \overline{\mathbf{y}}) \odot \nabla_{\mathbf{A}}\mathbf{s}$

**Algorithm 1** FOP algorithm for the error gradient of an MNN

---

where $n_t$ is the number of timesteps needed to the input to traverse the network and provide a sufficiently accurate output. In Recurrent Neural Networks (RNNs) a careful consideration is required to determine the value of $n_t$, because any recurrent connection results in a potentially undefined number of loops. However, a relevant advantage of MNN with respect to RNN based on back-propagation is that an MNN does not need to save the prior steps determined by a loop. In RNNs a bounded-history approximation strategy is used to simplify the computation and provide an adequate approximation to the true gradient: relevant information is saved in the fixed number of timesteps $n_t$ and any information older than that is forgotten. According to this strategy, in Backpropagation Through Time [10], a backward pass through the most recent $n_t$ time steps is performed at each time the network is run through an additional time step. In contrast, in MNN the lack of an error backpropagation sensibly reduces the impact of $n_t$: it should be large enough to capture the temporal structure of the problem to model. Thus, after $n_t$ timesteps the computation is simply truncated to take the output value. It is worth noting that already in the training phase weights are adjusted according to the specified $n_t$. Consequently, recurrent connections are adequately weakened when producing noise on the error, reducing the impact of the recurrent computation.In conclusion, a sufficiently large $n_t$ results in an adequate approximation to the true gradient, and it is not a sensitive parameter of the network.

The next section is devoted to the Python implementation and the evaluation of the proposed MNN.

## 4 Implementation and experimental studies

The MNN model has been developed, tested and publicly released on the Github platform, to make possible the initial roll-out of the approach, and to foster its application on various research environments. The implementation is based on *numpy*[11], a widespread package for tensor algebra in Python. The interested reader is referred to [12] for further implementation details.

The correctness of the symbolic derivatives is a critical aspect of the proposed network. To ensure it, in addition to the symbolic differentiation (SD), another implementation has been generated, in which gradients are calculated via automatic differentiation (AD) [13]. AD transforms a target function into a large graph of symbolic differentiation at elementary operation level, which are highly parallelizable [14]. This computational graph can efficiently manage orders of magnitude of gradients, providing highly accurate numerical values. Nevertheless, the AD-based system can be used for testing purposes only, since it is based on a back-propagated gradient error that has been criticized in the premise of this research work. To develop an efficient and coherent implementation of the proposed approach, the symbolic derivatives are then fundamental. To empirically evaluate the functional equivalence of the SD-based and the AD-based networks, the absolute differences between their corresponding output values have been computed over 100 tests. The two networks have been equipped with 5 input, 10 hidden and 3 output nodes. In each test, the two comparative networks have been set with (the same) random weights, and fed with a batch of (the same) 10 random inputs. As a result, the 95% confidence intervals of both the state and the gradient absolute differences are very low: $0.00024 \pm 0.000047$ and $0.000011 \pm 0.0000053$, respectively. The source code of the numerical test code has been publicly released [12].

### 4.1 Synthetic problems

In order to investigate the capabilities of the MNN model the dataset generator of *scikit-learn*[15] has been used to produce five types of two-dimensional dataset well-known in the literature (Figures 4, 5, 6): (a) *Moons*: a two-classes dataset made by two interleaving circles; (b) *Circles*: a two-class dataset made by concentric circles; (c) *Spirals*, which is considered as a good evaluation of training algorithms [2]; (d) *Single Blobs*: a three-class dataset made by isotropic Gaussian blobs with standard deviation 1.0, 2.5, 0.5; (e) *Double Blobs*: a three-class dataset made by two groups of isotropic Gaussian blobs with standard deviation 1.0.

Each dataset is made by 1,000 objects, balanced classes, and contains 10% of noise. Finally, a dataset from UCI Machine Learning Repository has been used, known as Iris [16]. Iris contains three classes of Iris plants. Each class consists of 50 objects characterised by 4 numeric features which describe, respectively, sepal length, sepal width, petal length and petal width. Class Iris Setosa is linearly separable from the other two. However, class Iris Versicolor and Iris Viginica are not separable from each other.

The MNN topology represented in Figure 3 has been used. Specifically, two output units have been assigned for the two classes datasets, and three output units for the three classes datasets. On the other side, three inputs units have been used: two inputs for the $(x, y)$ features of the dataset, and one input for the

bias input (constantly set to 1). 5 hidden units have been used. The Network has been evaluated for 3 time ticks. The ReLU activation function has been used for all units. Finally, the cross-entropy loss has been used as error function. For the experiments using the *Iris* dataset, it has been used an MNN with 5 input units (4 features and 1 bias), 10 hidden units and 3 output units (one for each class).
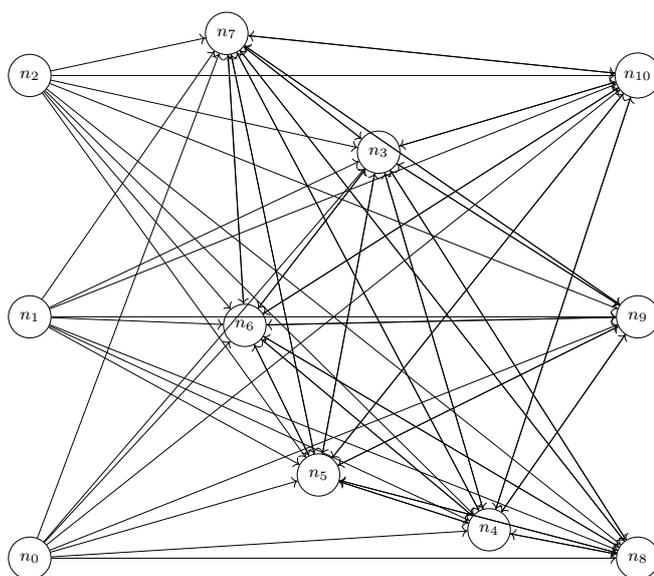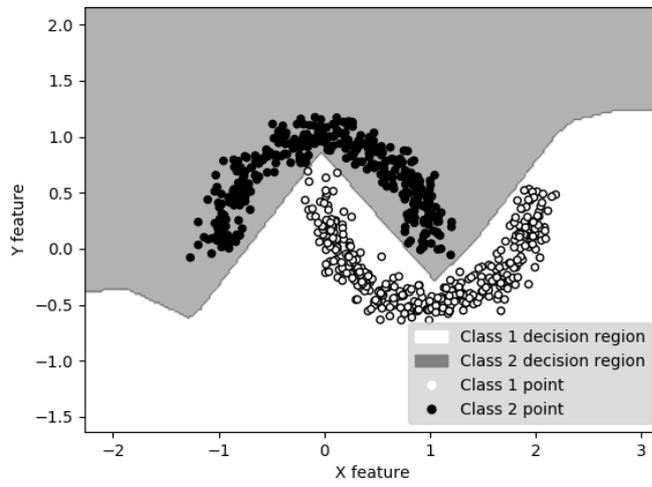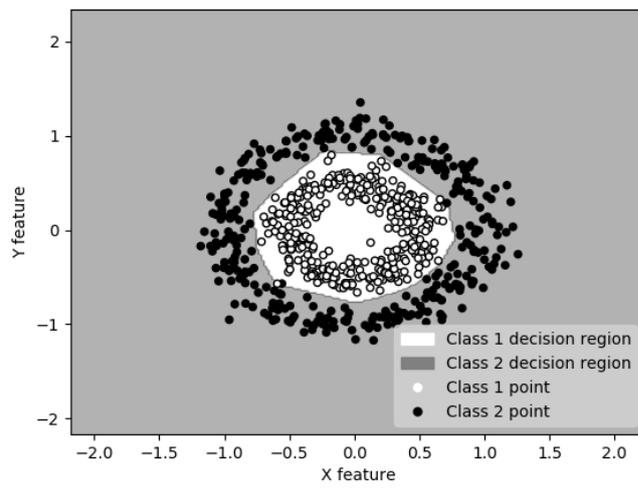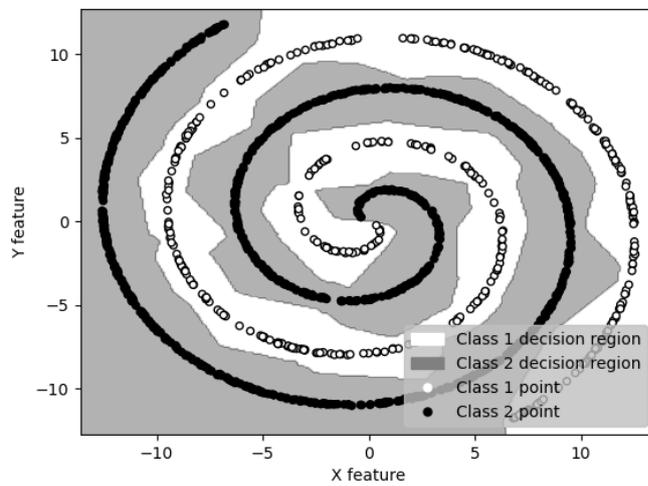


**Figure 3** MNN topology used in experiments

The Adaptive Moment Estimation (Adam) [17] has been used to compute adaptive learning rates for each parameter of the gradient descent optimization algorithms, carried out with batch method. A learning rate of 0.001 has been set. The training has been carried on for 1000 epochs.
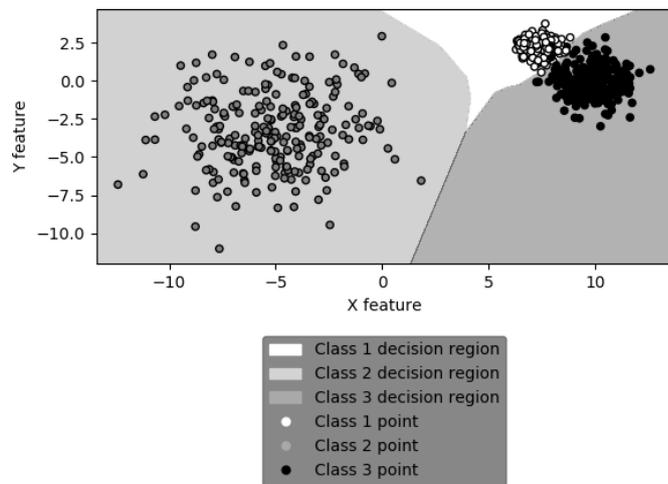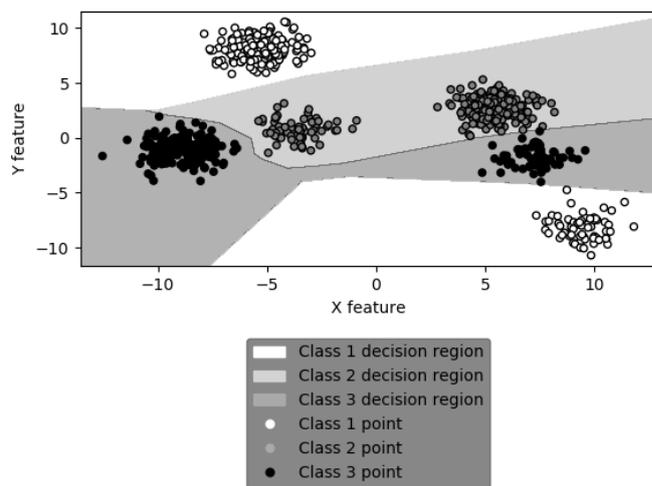
(a) Moons



(b) Circles



(c) Spirals

**Figure 4** Two-classes datasets and related decision regions

(a) Single Blobs



(b) Double Blobs

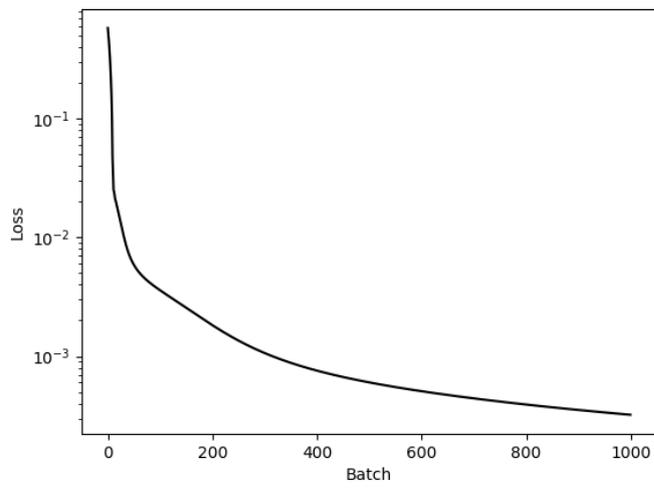**Figure 5** Three-classes datasets and related decision regions

The dataset has been partitioned into 70% and 30% for training and testing sets, respectively. Figures 4, 5, and 6 show with different gray levels the resulting partitioning of the input domain made by the MNN. Here, the generalization capabilities of the network are apparent. As a result, the MNN achieved the 100% accuracy for all datasets. In terms of complexity, the number of nodes of the MNN are $3 + 5 + 2 = 10$ and $3 + 5 + 3 = 11$ for 2 and 3 class datasets, respectively. The

corresponding number of parameters (weights) is $10 \cdot 10 = 100$ and $11 \cdot 11 = 121$, respectively. The interested reader is referred to [12] for a color animation of the MNN partitioning for each iteration. Table 1 shows the accuracy of the Spiral model generated by an MNN for increasing hidden neurons. It is interesting that, with 15 hidden neurons the problem is successfully modeled. Moreover, for a lower number of neurons, up to 7, the accuracy decreases gradually, in contrast to MLP and other approaches proposed in [2].
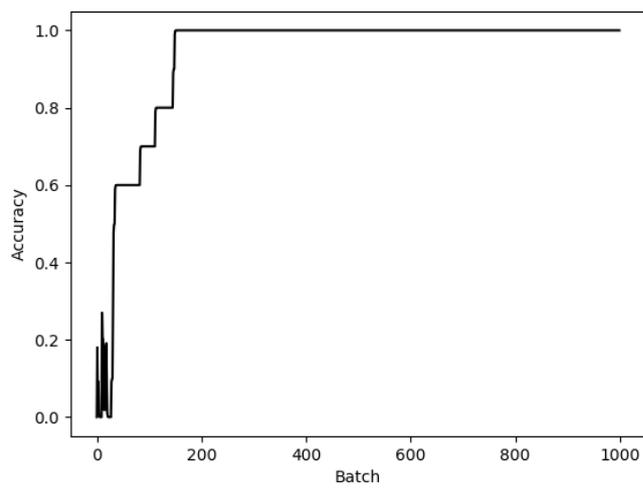
| Hidden Neurons | Accuracy |
|---|---|
| 5 | $0.75 \pm 0.079$ |
| 7 | $0.95 \pm 0.029$ |
| 10 | $0.94 \pm 0.039$ |
| 13 | $0.95 \pm 0.026$ |
| 15 | $0.99 \pm 0.011$ |

**Table 1** Accuracy of the Spiral model generated by an MNN for increasing hidden neurons

Figure 6(a) and Figure 6(b) show the training loss and the training accuracy over time for the *Iris* dataset. It is worth to note the convergence capabilities of the network. As a result, the MNN achieved $97.00\% \pm 1.62\%$ accuracy over 10 runs with a $3\sigma$ confidence interval. In terms of complexity, the number of nodes of the MNN is $5 + 10 + 3 = 18$. The corresponding number of parameters (weights) is $18 \cdot 18 = 324$.

(a) Training loss over time



(b) Training accuracy over time

**Figure 6** Training convergence of MNNs with Iris dataset

4.2 Real-world problems

To investigate the effectiveness of the MNN architecture, some experiments have been carried out on two real-world problems used for benchmarking machine learning algorithms: MNIST [18] and Fashion-MNIST [19]. MNIST is a database of handwritten digits images, whereas Fashion-MNIST is a dataset of fashion article images. Both datasets are made by a training set of 60,000 examples, and a test set of 10,000 examples. Each example is a 28x28 image, with pixels in 0-255 grayscale values, associated with a class label of 10 possible classes. The task is to classify a given image into one of such 10 classes. Figure 7 shows some representative samples of the datasets. Both datasets contain samples ambiguous even for humans: MNIST and Fashion-MNIST have an average human performance of 98.29% [20] and 83.5% [19], respectively. Such datasets are widely used and deeply investigated: top-performing models, based on convolutional neural networks, achieve a classification accuracy higher than 99%, and have a layered structure made by feature extraction and classification. Feature extraction can be performed by alternating convolution and subsampling layers, whereas classification can be performed via dense layers, such as a fully connected feed forward (i.e, MLP-based) neural network. The purpose of this section is to use an MNN as a classification layer, to carry out a comparative analysis between MNN and MLP. Rather than providing the top performance, this solution simplifies the design of the classification layer for the sake of simplicity. Indeed, for a fair comparison it is essential to avoid complex architectures with many hyper-parameters, whose particular choices should be subject to in-depth discussion. Similarly, there are many choices for convolutional architectures, but using a general-purpose architecture with a high degree of automation reduces such choices.



(a) MNIST  (b) Fashion MNIST

**Figure 7** Representative samples of MNIST and Fashion MNIST datasets

With this premise, a Convolutional Auto-Encoder (CAE) is used for feature extraction, followed by an MNN or MLP based network for classification. The CAE is commonly used for unsupervised data encoding and noise reduction [21]. The following architecture is used in experiments. *Encoding*: a convolutional layer with a 3x3 kernel size and 16 channels, stride 3 and padding 1; a rectified linear unit (ReLU); a max pooling layer with 2x2 kernel size, stride 2; a convolutional layer with a 3x3 kernel size and 8 channels, stride 2 and padding 1; a ReLU; a max pooling layer with 2x2 kernel size, stride 2. *Decoding*: a transpose convolutional layer with 3x3 kernel size and 16 channels, stride 2; a ReLU; a transpose convolutional layer with 5x5 kernel size and 8 channels, stride 3, padding 1; a ReLU; a transpose convolutional layer with 2x2 kernel size and 1 channel, stride 2, padding 1; a hyperbolic tangent activation function. Overall, the CAE provides 32 features to the classification layer. Both MLP and MNN classification layers have been equipped with $n_i = 32$ input and $n_o = 10$ output neurons. The number $n_h$ of hidden neurons has been set accordingly, to have the same number of overall connections for the two comparative networks. Since both datasets are spatial, in the MNN network two recursion steps, i.e., $n_t = 3$ , are sufficient. The MNN network is statically pruned for better efficiency. Since the MNN model generalizes the other perceptron-based topologies, there are custom pruning that makes the MNN fully equivalent, for instance, to an RNN or to an MLP. However, for a significant comparison, such custom pruning is avoided, in favor of a randomly determined pruning.

Figure 8 represents the adjacency matrix of the MNN based classification layer. Here, $\mathbb{I}$, $\mathbb{H}$, and $\mathbb{O}$, represent the sets of indexes corresponding to the input, hidden and output neurons. Each block can then be characterized by a pair of related sets. In particular, white blocks have zero connections, whereas dotted blocks have a given percentage of randomly selected connections.
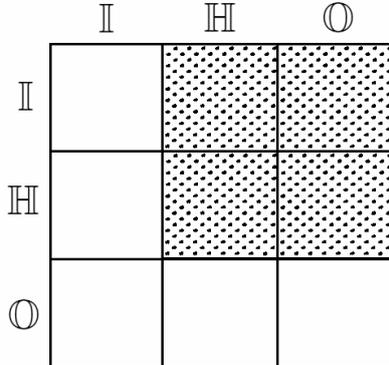


**Figure 8** Pruned adjacency matrix

Specifically, the white blocks represent the following connection types: all-to-input ($\mathbb{I} \rightarrow \mathbb{I}$, $\mathbb{H} \rightarrow \mathbb{I}$, $\mathbb{O} \rightarrow \mathbb{I}$), output-to-all ($\mathbb{O} \rightarrow \mathbb{I}$, $\mathbb{O} \rightarrow \mathbb{H}$, $\mathbb{O} \rightarrow \mathbb{O}$). The dotted blocks represent the following connection types: input-to-hidden ($\mathbb{I} \rightarrow \mathbb{H}$), input-to-output ($\mathbb{I} \rightarrow \mathbb{O}$), hidden-to-hidden ($\mathbb{H} \rightarrow \mathbb{H}$), and hidden-to-output ($\mathbb{H} \rightarrow \mathbb{O}$). As a consequence, assuming $n_h = 50$, the following connections and biases are available

in the dotted area: $(1-p) \cdot (n_i + n_h) \cdot (n_h + n_o) + (n_h + n_o)$, where $p$ is the pruning percentage, and the last term $(n_h + n_o)$ is the number of biases of hidden and output nodes. In order to have a similar number of connections, the MLP hidden neurons are made by two layers of $h_1$ and $h_2$ neurons. Hence, the total number of connections, considering also biases, is $(n_i \cdot h_1 + h_1) + (h_1 \cdot h_2 + h_2) + (h_2 \cdot n_o + n_o)$. The 95% confidence intervals achieved via the MNN and MLP based classifiers, calculated over 10 trials, are summarized in Table 2. It important to note that, for each trial, the percentage of randomly selected connections in MNN is completely renewed. As previously discussed, such classification rates are related to the features generated by the CAE layer. Consequently, the rates are not comparable with the top absolute performance of the literature. The significant result is that the MNN and the MLP based classifiers achieve very similar performance for increasing connections. For the sake of comparability, Table 2 shows only the settings with a very similar number of connections for the two networks. The best performance, of about 0.80 classification rate, is achieved via 1044-1047 connections. To evaluate the complexity of the classification task, it has been experimentally verified that for increasing number of connections (up to more than 3 thousand connections), both classifiers are not able to overcome the 0.8 classification rate.

| Architecture | Hidden Neurons | Pruning | Connections | MNIST | Fashion MNIST |
|:---:|:---:|:---:|:---:|:---:|:---:|
| MNN | 50 | 85% | 798 | $0.772 \pm 0.030$ | $0.762 \pm 0.005$ |
| MLP | 14+13 | - | 797 | $0.761 \pm 0.065$ | $0.750 \pm 0.026$ |
| MNN | 50 | 80% | 1044 | $0.796 \pm 0.013$ | $0.771 \pm 0.004$ |
| MLP | 17+17 | - | 1047 | $0.789 \pm 0.035$ | $0.786 \pm 0.013$ |

**Table 2** Testing classification rates of the MLP and MNN based networks, on MNIST and fashion-MNIST

## 5 Conclusions and future work

The purpose of this paper is to formally introduce recent advances leading to the MNNs, providing the key points to the reader.

Overall, the main advantages of the MNN model with the related FOP algorithm are: (i) the state partial derivatives can be computed along the forward propagation; (ii) the error gradient can be directly computed from state gradient; (iii) the state partial derivative update makes use only of short-lived variables, which can be overwritten at each state iteration; (iv) the state partial derivatives concern only one multidimensional parameter; (v) the overall gradient computation relies only on tensor multiplications, which can be easily distributed on parallel computing, thus potentially enabling large-scale sparse ANNs training [22].

In contrast, the BP-based family of algorithms is limited to layer-wise architectures, and needs to store all intermediate layer outputs, by comprising a forward and backward propagation through the network. On the other side, the CG-based gradient computation is not constrained in terms of network architecture, but it needs to store a large graph topology and the partial derivatives of each computation node, and it needs to compute all factoring paths for each parameter.

Due to its unconstrained structure, an interesting research perspective of MNNs is to adopt structural regularization techniques to dynamically drive the network topology. For small datasets the network topology is highly dense to exploit the available neurons, and then the adjacency matrix is highly dense. For large datasets, in general there are two strategies that can make the adjacency matrix sparse: a) offline pruning, i.e., to remove some types of connections according to some heuristics; b) online pruning, i.e., to remove iteratively some connections that do not contribute to model, in the training phase. The two strategies can be combined. In general, the possibility to have a sparse matrix depends on the problem complexity. Since the MNN needs large matrix operations, such strategies should also be supported by a framework implementation that efficiently exploits the hardware resources, e.g. via memory caching and highly parallel computation. However, the commercially available machine learning framework provide optimized libraries for back-propagated models. As a consequence, to test the MNN architecture on very large datasets, an optimized framework should be implemented on specific hardware. Such development is a long-term task and it is out of the scope of this paper, which focuses on the formal derivation of the technique and on pilot experimentation showing its potential application.

As a future work, in order to compare BP, CG and FOP according to a performance perspective, the scalability of each algorithm should be evaluated in terms of computational complexity. Moreover, a statistical performance evaluation should be carried out on benchmark problems, considering large-scale applications.

## Acknowledgements

## References

1. F. Galatolo, M. Cimino, and G. Vaglini, "Using stigmergy as a computational memory in the design of recurrent neural networks," *Proceedings of the 8th International Conference on Pattern Recognition Applications and Methods*, 2019.
2. B. M. Wilamowski and H. Yu, "Neural network learning without backpropagation," *IEEE Transactions on Neural Networks*, vol. 21, pp. 1793–1803, Nov 2010.
3. W. Guo, H. Huang, and T. Huang, "Complex-valued feedforward neural networks learning without backpropagation," in *Neural Information Processing* (D. Liu, S. Xie, Y. Li, D. Zhao, and E.-S. M. El-Alfy, eds.), (Cham), pp. 100–107, Springer International Publishing, 2017.
4. K. W.-D. Ma, J. P. Lewis, and W. B. Kleijn, "The hsic bottleneck: Deep learning without back-propagation," *ArXiv*, vol. ArXiv:1908.01580v3, 2019.
5. M. Jaderberg, W. M. Czarnecki, S. Osindero, O. Vinyals, A. Graves, D. Silver, and K. Kavukcuoglu, "Decoupled neural interfaces using synthetic gradients," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1627–1635, JMLR. org, 2017.

6.  J. M. Keller, D. Liu, and D. B. Fogel, *Multilayer Neural Networks and Backpropagation*, vol. Fundamentals of Computational Intelligence: Neural Networks, Fuzzy Systems, and Evolutionary Computation, p. 378. Wiley-IEEE Press, 2016.

7.  P. J. Werbos, *The roots of backpropagation: from ordered derivatives to neural networks and political forecasting*, vol. 1. John Wiley & Sons, 1994.

8.  S. Theodoridis, "Chapter 5 - stochastic gradient descent: The lms algorithm and its family," in *Machine Learning* (S. Theodoridis, ed.), pp. 161 – 231, Oxford: Academic Press, 2015.

9.  I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.

10. H. Tang and J. Glass, "On training recurrent networks with truncated backpropagation through time in speech recognition," in *2018 IEEE Spoken Language Technology Workshop (SLT)*, pp. 48–55, IEEE, 2018.

11. T. Oliphant, "NumPy: A guide to NumPy." USA: Trelgol Publishing, 2006–.

12. F. Galatolo, "https://github.com/galatolofederico/mesh-neural-networks," *GitHub repository*, 2019.

13. A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.

14. A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind, "Automatic differentiation in machine learning: a survey," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 5595–5637, 2017.

15. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

16. E. Anderson, "The Species Problem in Iris," *Annals of the Missouri Botanical Garden*, vol. 23, pp. 457–509.

17. D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *international conference on learning representations*, 2014.

18. Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

19. H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.

20. P. Simard, Y. LeCun, and J. Denker, "Efficient pattern recognition using a new transformation distance," in *Advances in Neural Information Processing Systems* (S. Hanson, J. Cowan, and C. Giles, eds.), vol. 5, pp. 50–58, Morgan-Kaufmann, 1993.

21. M. Chen, X. Shi, Y. Zhang, D. Wu, and M. Guizani, "Deep features learning for medical image analysis with convolutional autoencoder neural network," *IEEE Transactions on Big Data*, 2017.

22. D. Irony, S. Toledo, and A. Tiskin, "Communication lower bounds for distributed-memory matrix multiplication," *Journal of Parallel and Distributed Computing*, vol. 64, no. 9, pp. 1017–1026, 2004.