

Chapter 1

Introduction to Neural Networks

Let us start this dissertation by introducing the key components of an artificial neural network and discussing the way it can be configured for a specific application. Please note that this chapter is not meant to provide a comprehensive overview on neural networks, rather to investigate some aspects and concepts functional to the following chapters. For further reading, we refer the reader to, e.g., [5, 6, 9], from which we retrieved many of the informations provided in this chapter.

Throughout this work we confine the attention to the most-spread neural network paradigm - the *feedforward* neural network, presented in Section ?? - employing the well-known *backpropagation of error* as learning rule, derived in Section ?. Actually, in the numerical experiments we carried out and whose results will be discussed in Chapter ??, we mainly refer to a variant of backpropagation - the Levenberg-Marquardt algorithm [], shortly discussed in Section ?.

Before moving to the description of technical neural networks, let us provide a brief excursus on their biological counterparts. The goal is to highlight the basic features of the human nervous system, focusing on the working principles of neurons and the way informations are processed, thus to extract the key concepts which should be taken over into a mathematical, simplified representation.

1.1 Biological motivation

The information processing system of a vertebrate can coarsely be divided into the *central nervous system* (CNS) and the *peripheral nervous system* (PNS). The former consists of the *brain* and the *spinal cord*, while the latter mainly comprises the *nerves*, which transmit informations from all other parts of the body to the CNS (*sensory nerves*) and viceversa (*motor nerves*). When an output stimulus hits the sensory cells of an organ sense, these generate an electric signal, called *sensory signal*, which is transferred to the central nervous system via the *sensory nerves*. Within the CNS, informations are stored and managed to provide the muscles with a suitable *motor signal*, broadcast through the *motor nerves* and finally converted by the effectors into a system output [6].

Hence, both the central and peripheral nervous system are directly involved in the information processing workflow. At the cellular level, this is accomplished through a huge amount of modified cells called *neurons*. These processing elements continuously communicate each other by means of electric signals, traveling through a thick net of connections. For instance, in a human being each neuron is linked in average with $10^3 - 10^4$ other neurons. As detailed in the next paragraph, a neuron is characterized by a rather simple structure, specifically designed to rapidly collect input signals and generate an output pulse whenever the accumulated incoming signal exceeds a threshold - the *action potential*. In other terms, a neuron acts as a switch, establishing a typically nonlinear input-output mapping [9].

From a simplifying perspective, a neuron consists of three main components: the *dendrites*, the *nucleus* or *soma*, and the *axon*. Dendrites are tree-like networks of nerve fibers receiving input signals from many sources and conveying them directly to the nucleus of the neuron. Here, input signals are accumulated

and thresholded, as mentioned before. The possible output pulse is then broadcast to the cells contiguous the neuron through the axon - a unique, slender fiber constituting an extension of the soma and splitting in many branches at the opposite extremity [15]. To ease the electrical conduction of the signal, the axon is isolated through a myelin sheath which consists of Schwann cells (in the PNS) or oligodendrocytes (in the CNS). However, this insulating film is not continuous, rather presents gaps at regular intervals called *nodes of Ranvier*, which lets the signal be conducted in a saltatory way.

The signal coming from the axon of another neuron or from another cell is transferred to the dendrites of a neuron through a particular connection called *synapsis*¹. A synaptic may be either electrical or chemical. In the former, the presynaptic side, i.e. the sender axon, and the postsynaptic side, i.e. the receiver dendrite, are directly in contact, so that the potential can simply travel by electrical conduction. Conversely, a chemical synapsis consists of a synaptic *cleft*, physically separating the presynaptic side from the postsynaptic side. Then, to let the action potential reach the postsynaptic side, at the presynaptic side the electrical pulse is converted into a chemical signal. This is accomplished by releasing some chemical substances called *neurotransmitters*. These neurotransmitters then cross the cleft and bind to the receptors dislocated onto the membrane of the postsynaptic side, where the chemical signal is re-converted into an electrical potential. On the other hand, neurotransmitters do not simply broadcast the action potential. Indeed, we can distinguish between excitatory and inhibitory neurotransmitters, respectively amplifying or modulating the signal. Hence, the pulse outgoing a neuron is preprocessed within the synapsis before reaching the target cell. In other terms, a neuron gets in input many *weighted* signals, which should then be collected.

Different studies have unveiled the tight correlation between the synapses the neurons establish among each other, and the tasks a neural network can address [5]. That is, the set of interneuron connection strengths represent the information storage, i.e. the knowledge, of a neural network [9]. Knowledge is acquired through a *learning* or *training* process, entailing adjustments at the synaptic level to adapt to environmental situations. The adjustments may not only involve the modification of existing synapses, but also the creation of new synaptic connections. Hence, the nervous system is a distributed memory machine whose evolutionary structure is shaped by experience.

As mentioned above, a biological neural network acquaints itself with problems of a specific class through a learning procedure. During the learning, the network is exposed to a collection of situations, giving it the possibility to derive a set of tools which will let it provide reasonable solutions in similar circumstances. In other terms, the cognitive system should be able to *generalize*. Furthermore, after a successful training a neural network should also show a discrete level of *fault tolerance* against external errors, e.g. noisy inputs. It worths notice here that the nervous system is also naturally fault tolerant against *internal* errors. Indeed, in case a neuron or a (relatively small) group of neurons got damaged or died, the other processing nodes would take care of its tasks, so that the overall cognitive capabilities would be only slightly affected [9].

1.2 Artificial neural networks

Inspired by the biological information processing system discussed so far, an artificial neural network (ANN), usually simply referred to as "neural network", is a computational model capable to learn from observational data, i.e. by example, thus providing an alternative to the algorithmic programming paradigm [12]. Exactly as its original counterpart, it consists of a collection of processing units, called (artificial) neurons, and directed weighted synaptic connections between the neurons themselves. Data travel among neurons through the connections, following the direction imposed by the synapses themselves. Hence, an artificial neural network is an *oriented graph* to all intents and purposes, with the neurons as *nodes* and the synapses as oriented *edges*, whose weights are adjusted by means of a *training* process to configure the network for a specific application [15].

¹For the sake of completeness, we mention that there exist synapses directly connecting the axon of the sender neuron with either the soma or the axon of the receiver. Actually, a synapsis may also connect the axon of a neuron with the dendrite or soma of the same neuron (autosynapsis). However, for our purposes we can confine the attention to the axon-dendrite synapsis.

Formally, a neural network could be defined as follows [9].

Definition 1.1 (Neural network). *A neural network is a sorted triple $(\mathcal{N}, \mathcal{V}, w)$, where \mathcal{N} is the set of neurons, with cardinality $|\mathcal{N}|$, $\mathcal{V} = \{(i, j), 1 \leq i, j \leq |\mathcal{N}|\}$ is the set of connections, with (i, j) denoting the oriented connection linking the sending neuron i with the target neuron j , and $w : \mathcal{V} \rightarrow \mathbb{R}$ is the weight function, defining the weight $w_{i,j} = w((i, j))$ of the connection (i, j) . A weight may be either positive or negative, making the underlying connection either excitatory or inhibitory, respectively. By convention, $w_{i,j} = 0$ means that neurons i and j are not directly connected.*

In the following, we dive deeper into the structure and training of a neural network, starting by detailing the structure of an artificial neuron.

1.2.1 Neuronal model

As its name may suggest, an artificial neuron represents a simplified model of a biological neuron, retaining its main features discussed in Section ???. To introduce the components of the model, let us consider the neuron j represented in Figure ???. Suppose that it is connected with m sending neurons s_1, \dots, s_m , and n receiving (target) neurons r_1, \dots, r_n . Denoting by $y_\Omega(t) \in \mathbb{R}$ the scalar output fired by a generic neuron Ω at time t , neuron j gets the weighted inputs $w_{s_k,j} \cdot y_{s_k}(t)$, $k = 1, \dots, m$, at time t , and sends out the output $y_j(t + \Delta t)$ to the target neurons r_1, \dots, r_n at time $t + \Delta t$. Please note that in the context of artificial neural networks the time is discretized by introducing the timestep Δt . This is clearly not plausible from a biological viewpoint; on the other hand, it dramatically eases the implementation. In the following, we will avoid to specify the dependence on time unless strictly necessary, thus to lighten the notation.

An artificial neuron j is completely characterized by three functions: the propagation function, the activation function, and the output function. These will be defined and detailed hereunder in the same order they get involved in the data flow.

Propagation function. The propagation function f_{prop} converts the vectorial input $\mathbf{p} = [y_{s_1}, \dots, y_{s_m}]^T \in \mathbb{R}^m$ into a scalar u_j often called *net input*, i.e.

$$u_j = f_{prop}(w_{s_1,j}, \dots, w_{s_m,j}, y_{s_1}, \dots, y_{s_m}). \quad (1.1)$$

A common choice for f_{prop} (used also in this work) is the weighted summer, adding up the scalar inputs multiplied by their respective weights:

$$f_{prop}(w_{s_1,j}, \dots, w_{s_m,j}, y_{s_1}, \dots, y_{s_m}) = \sum_{k=1}^m w_{s_k,j} \cdot y_{s_k}. \quad (1.2)$$

The function (1.2) provides a simple yet effective way of modeling the accumulation of different input electric signals within a biological neuron; this motivates its popularity.

Activation or transfer function. At each timestep, the *activation state* a_j , often shortly referred to as *activation*, quantifies at which extent neuron j is currently active or excited. It results from the activation function f_{act} , which combines the net input u_j with a threshold $\theta_j \in \mathbb{R}$ [9]:

$$a_j = f_{act}(u_j; \theta_j) = f_{act}\left(\sum_{k=1}^m w_{s_k,j} \cdot y_{s_k}; \theta_j\right), \quad (1.3)$$

where we have employed the weighted summer (1.2) as propagation function. From a biological perspective, the threshold θ_j is the analogous of the action potential mentioned in Section ???. Mathematically, it represents the point where the absolute value $|f'_{act}|$ of the derivative of the activation function is maximum. Then, the activation function reacts particularly sensitive when the net input u_j hits the threshold value θ_j [9].

Furthermore, noting that θ_j is a parameter of the network, one may like to adapt it through a training

process, exactly as can be done for the synaptic weights, as we shall see in Section ???. However, θ_j is currently incorporated in the activation function, making its runtime access somehow cumbersome. This is typically overcome by introducing a *bias neuron* in the network. A bias neuron is a continuously firing neuron, with constant output $y_b = 1$, which gets directly connected with neuron j , assigning the *bias weight* $w_{b,j} = -\theta_j$ to the connection. As can be deduced by Figure ??, θ_j is now treated as a synaptic weight, while the neuron threshold is set to zero. Moreover, the net input becomes

$$u_j = \sum_{k=1}^m w_{s_k,j} \cdot y_{s_k} - \theta_j, \quad (1.4)$$

i.e. the threshold is now included in the propagation function rather than in the activation function, which we can now express in the form

$$a_j = f_{act}\left(\sum_{k=1}^m w_{s_k,j} \cdot y_{s_k} - \theta_j\right). \quad (1.5)$$

Let us point out that this trick can be clearly applied to all neurons in the network which are characterized by a non-vanishing threshold: just connect the neuron with the bias, weighting the connection by the opposite of the threshold value. However, for ease of illustration in the following we shall avoid to include the bias neuron in any graphical representation of a neural network.

Conversely to the propagation function, there exist various choices for the activation function, as the Heaviside or binary function, which assumes only 0 or 1, according to whether the argument is negative or positive, respectively:

$$f_{act}(v) = \begin{cases} 0, & \text{if } v < 0, \\ 1, & \text{if } v \geq 0. \end{cases} \quad (1.6)$$

Neurons characterized by such an activation function are usually named McCulloch-Pitts neurons, after the seminal work of McCulloch and Pitts [6], and their employment is usually limited to single-layer perceptrons implementing boolean logic (see Section ??). In addition, note that (1.6) is discontinuous, with a vanishing derivative everywhere except that in the origin, thus not admissible for the backpropagation training algorithm presented in Section ??.

Among continuous activation maps, sigmoid functions have been widely used for the realization of artificial neural networks due to their graceful combination of linear and nonlinear behaviour [6]. Sigmoid functions are s-shaped and monotonically increasing, and assume values in a bounded interval, typically $[0, 1]$, as the logistic function,

$$f_{act}(v) = \frac{1}{1 + e^{-v/T}} \quad \text{with } T > 0, \quad (1.7)$$

or $[-1, 1]$, as the hyperbolic tangent,

$$f_{act}(v) = \frac{e^v - e^{-v}}{e^v + e^{-v}}. \quad (1.8)$$

Both functions are displayed in Figure ??. Note that the logistic function resemble the Heaviside function as T increases.

Output function. Finally, the output function f_{out} is in charge of calculating the scalar *output* $y_j \in \mathbb{R}$ based on the activation state a_j of the neuron:

$$y_j = f_{out}(a_j); \quad (1.9)$$

typically, f_{out} is the identity function, so that activation and output of a neuron coincides, i.e. $y_j = f_{out}(a_j) = a_j$. Let us point out that while the input $\mathbf{p} = [y_{s_1}, \dots, y_{s_m}]^T \in \mathbb{R}^m$ of the neuron is generally vectorial, i.e. $m > 1$, the output is scalar. The output y_j could then either be sent to other neurons, included the outputting neuron itself (autosynapsis), or constitute a component of the overall output vector of the network, as for the neurons in the output layer of a feedforward neural network (see Section ??).

It worths mention here that, as the activation function, also the output function is usually *globally* defined, i.e. all neurons in the network (or at least a group of neurons) are equipped with the same output function.

The neural model presented so far actually refers to the so called *computing* neuron, i.e. a neuron processing input informations to provide a response. However, in a neural network one may also distinguish *source* neurons, supplying the network with the respective components of the activation pattern (input vector) [6]. The role of source neurons will be clearer in the next section, where we will introduce the multilayer feedforward neural network. Here we just mention that such a neuron receives a scalar and unweighted input, which is simply forward to the connected neurons; no computations are performed.

1.2.2 Network topologies: the feedforward neural network

The way neurons are interconnected within a network defines the *topology* of the network itself, i.e. its design. In literature, many network architectures have been proposed, sometimes tailored to a specific application or task. In this section, we expand our investigation for the two most common network topologies: the feedforward and the recurrent neural network.

Feedforward neural network. In a feedforward neural network, also called *perceptron* [14], neurons are arranged into *layers*, with one *input layer* of M_I source neurons, K *hidden layers*, each one consisting of H_k computing neurons, $k = 1, \dots, K$, and an *output layer* of M_O computing neurons. As characteristic property, neurons in a layer can only be connected with neurons in the next layer towards the output layer. Then, an *activation pattern* $\mathbf{p} \in \mathbb{R}^{M_I}$, supplied to the network through the source nodes in the first layer, provides the input signal for the neurons in the first hidden layer. For each hidden layer, its output signal gives the input pattern for the following layer. In this way, informations travel towards the last layer of the network, i.e. the output layer, whose outputs constitute the components of the overall output $\mathbf{q} \in \mathbb{R}^{M_O}$ of the network, response to the input pattern \mathbf{p}^2 . Hence, a feedforward network establish a map between the *input space* \mathbb{R}^{M_I} and the *output space* \mathbb{R}^{M_O} . This makes this network architecture particularly suitable for, e.g., classification and continuous function approximation.

Feedforward networks can be classified according to the number of hidden neurons they present, or, equivalently, the number of layers of trainable weights. Single-layer perceptrons (SLPs) just consist of the input and output layer; no hidden layers. Note that the layer which the name refers to is the output layer; the input layer is not accounted for since it does not perform any calculation. Despite their quite simple structure, the range of application of single-layer perceptrons is rather limited. Indeed, consider a binary input vector, supplied to an SLP provide with a unique output neuron, equipped with a binary activation function. The network acts then as a classifier, splitting the input space, i.e. the unit hypercube, by means of a hyperplane. Therefore, only *linearly separable* data can be properly represented [9]. On the other hand, the share of linearly separable sets decreases as the space dimension increases, making single-layer perceptrons seldom attractive.

Conversely, multi-layer perceptrons (MLPs), providing at least one hidden layer, are universal function approximators, as stated by Cybenko [1, 2]. In detail:

- (i) a multi-layer perceptron with one layer of *hidden neurons* and differentiable activation functions can approximate any *continuous* function [2];
- (ii) a multi-layer perceptron with two layers of *hidden neurons* and differentiable activation functions can approximate *any function* [1].

Therefore, in many practical applications there is no reason to employ MLPs with more than two hidden layers. Considering again the binary classifier discussed above, one can represent any convex polygon by

²Please note that while the output of a single neuron is denoted with the letter y , we use the letter \mathbf{q} (bolded) to indicate the overall output of the network. Clearly, for the j -th output neuron the output y_j coincides with the correspondent entry of \mathbf{q} , i.e. $q_j = y_j$, $j = 1, \dots, M_O$.

adding one hidden layer, and any arbitrary set by adding two hidden layers; further increasing the number of hidden neurons would not improve the representation capability of the network. However, we should point out that (i) and (ii) do not give any practical advice on both the number of hidden neurons and the number of samples required to teach the network.

An instance of a three-layer (i.e. two hidden layer plus the output layer) feedforward network is offered in Figure ???. In this case, we have $M_I = 3$ input neurons (denoted with the letter i), $H_1 = H_2 = 6$ hidden neurons (letter h), and $M_O = 4$ output neurons (letter o). In particular, we point out that it represents an instance of a *completely linked* perceptron, since each neuron is directly connected with all neurons in the following layer.

Finally, let us just mention that, although we have previously stated that in a feedforward neural network a synapses can only connect pairs of neurons in contiguous layers, recent years have seen the development of different variants. For instance, *shortcut connections* skip one or more layers, while *lateral connections* takes place between neurons within the same layer. However, throughout this work we shall be faithful to the original definition of the perceptron.

Recurrent neural network. In recurrent networks any neuron can bind with any other neuron, but aut-synapses are forbidden, i.e. the output of a neuron can be input into the same neuron at the next time step. If each neuron is connected with all other neurons, then the network is said *completely linked*. As a consequence, one can not distinguish neither input or output neurons any more: neurons are all equivalent. Then, the input of the network is represented by the initial *network state*, which is the set of activation states for all neurons in the network. Similarly, the overall network output is given by the final network state. So, communication between a recurrent neural network and the surrounding environment takes place through the states of the neurons. Examples of recurrent networks are the Hopfield networks [], inspired by the behaviour of electrically charged particles within a magnetic field, and the self-organizing maps by Kohonen [], highly suitable for cluster detection.

As mentioned in the introductory chapter, in this work we refer to neural networks for the approximation of the unknown map ϕ between the parameters μ of a parametrized partial differential equation, and the coefficients α of the corresponding reduced basis solution. To accomplish this task, we rely on a collection of samples $\{\mu_i, \alpha_i\}_{i=1}^{N_{tr}}$ generated through a high-fidelity model. Although a detailed explanation will be provided later in Chapter ??, what is worth notice here is that we are concerned with a *continuous function approximation* problem. Then, motivated by what previously said in the section, a multi-layer feedforward neural network turns out as the most suitable network architecture for our purposes.

We are now left to investigate the way the weights of a perceptron can be *trained* to meet our purposes.

1.2.3 Training a multilayer feedforward neural network

As widely highlighted so far, the primary characteristic of a neural network lies in the capability to *learn* from the environment, storing the acquired knowledge through the network internal parameters, i.e. the synaptic and bias weights. Learning is accomplished through a training process, during which the network is exposed to a collection of examples, called *training patterns*. According to some performance measure, the weights are then adjusted by means of a well-defined set of rules. Therefore, a learning procedure is an *algorithm*, typically iterative, modifying the neural network parameters to make it knowledgeable of the specific environment it operates in [6]. Specifically, this entails that after a successful training, the neural network has to provide reasonable responses for unknown problems of the same class the network was acquainted with during training. This property is known as *generalization* [9].

Training algorithms can be firstly classified based on the nature of the training set, i.e. the set of training patterns. We can then distinguish three *learning paradigms*.

Supervised learning. The training set consists of a collection of *input patterns* (i.e. input vectors) $\{p_i\}_{i=1}^{N_{tr}}$, and corresponding desired responses $\{t_i\}_{i=1}^{N_{tr}}$, called *teaching inputs*. Then, t_i is the output the neural network should desirably provide when it gets fed with p_i . As we shall see, any training procedure aims to minimize

(in some appropriate norm) the discrepancy between the *desired* output t_i and the *actual* output q_i given by the network as response to p_i .

Unsupervised learning. Although supervised learning is a simple and intuitive paradigm, there exist many tasks which require a different approach to be tackled. Consider for instance a cluster detection problem. Due to lack of prior knowledge, rather than telling the neural network how it should behave in certain situations, one would like the neurons to *independently* identify rules to group items. Therefore, a training pattern just consist of an input pattern. Since no desired output is provided, such a pattern is referred to as *unlabeled*, as opposed to the *labeled* examples involved in the supervised learning paradigm.

Reinforcement learning. Reinforcement learning is the most plausible paradigm from a biological viewpoint. After the completion of a series of input patterns, the neural network is supplied with a boolean or a real saying whether the network is wrong or right. In the former case, the *feedback* or *reward* may also indicate to which extent the network is wrong [9]. Conversely to the supervised and unsupervised paradigms, reinforcement learning focuses on finding a balance between exploration (of uncharted territory) and exploitation (of current knowledge). Hence, this paradigm particularly suits problems involving a trade-off between a long-term and a short-term reward [8].

Clearly, the choice of the learning paradigm is task-dependent. In particular, function approximation (i.e. what we are interested in) perfectly fits the *supervised learning* paradigm. Indeed, consider the nonlinear unknown function f ,

$$\begin{aligned} f : \mathbb{R}^{M_I} &\rightarrow \mathbb{R}^{M_O} \\ x &\mapsto y = f(x), \end{aligned}$$

and a set of labeled examples $\{x_i, y_i\}_{i=1}^{N_{tr}}$. The task implies to approximate f over a domain $V \subset \mathbb{R}^{M_I}$ up to a user-defined tolerance ϵ , i.e.

$$\|F(x) - f(x)\| < \epsilon \quad \forall x \in V,$$

where $F : \mathbb{R}^{M_I} \rightarrow \mathbb{R}^{M_O}$ is the actual input-output map established by the neural network, and $\|\cdot\|$ is some suitable norm on \mathbb{R}^{M_O} . Surely, as necessary condition the neural system must satisfy to well approximate f for each input in the domain V , the system should provide accurate predictions for the input patterns, i.e.

$$F(x_i) \approx y_i, \quad \forall i = 1, \dots, N_{tr}.$$

Then, we could train the network through a supervised learning algorithm, employing $\{x_i\}_{i=1}^{N_{tr}}$ as input patterns and $\{y_i\}_{i=1}^{N_{tr}}$ as teaching inputs. That is,

$$p_i = x_i \text{ and } t_i = y_i, \quad \forall i = 1, \dots, N_{tr}.$$

As defined in the first part of the section, a training algorithm involves:

- (a) a set of well-defined rules to modify the synaptic and bias weights;
- (b) a *performance function*, quantifying the current level of knowledge of the surrounding environment.

Regarding (a), a plethora of weight updating techniques have been proposed in literature, sometimes tailored on specific applications. Nevertheless, most of them relies on the well-known Hebbian rule, proposed by Donal O. Hebb in 1949 [7]. Inspired by neurobiological considerations, the rule can be stated in a two-steps fashion [6]:

- (i) if two neurons on either side of a synapse (connection) are activated simultaneously (i.e. synchronously) then the strength of that synapse is selectively increased;
- (ii) if two neurons on either side of a synapse are activated asynchronously, then that synapse is selectively weakened (or eliminated).

Actually, (ii) was not included in the original statement; nevertheless, it provides a natural extension to [6]. In mathematical terms, consider the synopsis between a sending neuron i and a target neuron j . Then, at the t -th iteration (also called *epoch*) of the training procedure, the weight $w_{i,j}(t)$ of the connection (i, j) is modified by the quantity

$$\Delta w_{i,j}(t) \sim \eta y_i(t) a_j(t), \quad (1.10)$$

where $\eta > 0$ is the *learning rate*, and we have exploited the fact that $y_i = a_i$ since the output function is usually represented as the identity. Hence, at the subsequent iteration $t + 1$ the synaptic weight is simply given by

$$w_{i,j}(t + 1) = w_{i,j}(t) + \Delta w_{i,j}(t). \quad (1.11)$$

Many of the supervised learning rules proposed in literature, included the backpropagation of error derived later, can be recast in the following form, which turns out as a generalization of the Hebbian rule (1.10) [9]:

$$\Delta w_{i,j} = \eta h(y_i, w_{i,j}) g(a_j, t_j), \quad (1.12)$$

with the functions h and g dependent on the specific learning rule, and t_j the j -th component of the teaching input \mathbf{t} . Note that all variables involved in (1.12) are supposed to be evaluated at time t , i.e. the correction $\Delta w_{i,j}$ is time-dependent. In addition, let us remark that (1.12) represents a *local* and *interactive* mechanism, since it involves both and only the neurons at the end-points of the synapse.

The second ingredient required to define a training algorithm is the performance or error function E . This function somehow measures the discrepancy between the neural network knowledge of the surrounding environment, and the actual state of the environment itself. In other terms, the larger the performance function, the farer the neural network representation of the world is from the actual reality, i.e. the farer we are from the application goal. Therefore, every learning rule aims to *minimize* the performance E as much as possible. For this purpose, E should be intended as a scalar function of the free parameters, i.e. the weights, of the network, namely

$$E = E(\mathbf{w}) \in \mathbb{R}. \quad (1.13)$$

Recalling the notation and assumptions introduced in Definiton 1.1, here $\mathbf{w} \in \mathbb{R}^{|\mathcal{V}|}$ is a vector collecting the weights $\{w_{i,j} = w((i, j))\}_{(i,j) \in \mathcal{V}}$, with \mathcal{V} the set of admissible connections in the network³. Thus, Equation (1.13) implies that the point over the error surface reached at the end of a successful training process provides the *optimal* configuration \mathbf{w}_{opt} for the network.

The steps a generic supervised training procedure should perform are listed by Algorithms ?? and ?? for online and offline learning, respectively. *Online* learning means that the weights are updated after the exposition of the network to each training pattern. In other terms, each epoch involves only one training pattern. Conversely, in an *offline* learning procedure the modifications are based on the entire training set, i.e. the weights are adjusted only after the network has been fed with all input patterns and the corresponding errors have been accumulated. Therefore, we should distinguish between the *specific error* $E_{\mathbf{p}}(\mathbf{w})$, specific to the activation pattern \mathbf{p} , and the *total error* $E(\mathbf{w})$ accounting for all specific errors, namely

$$E(\mathbf{w}) = \sum_{\mathbf{p} \in P} E_{\mathbf{p}}(\mathbf{w}), \quad (1.14)$$

with P the training set. For instance, we could think of the specific error as the Mean Squared Error (MSE):

$$E_{\mathbf{p}}(\mathbf{w}) = \frac{1}{M_O} \sum_{j=1}^{M_O} (t_{\mathbf{p},j} - q_{\mathbf{p},j})^2, \quad (1.15)$$

³Please note that while in Definiton 1.1 \mathcal{V} denoted the set of all *possible* connections, here \mathcal{V} disregards those connections which are not consistent with the network topology in use. For instance, a feedforward neural network can not be endowed with connections oriented towards the input layer, then such connections will not be included in \mathcal{V} . In this way, we reduce the size of \mathbf{w} - which is a practical advantage.

where we have added the subscript \mathbf{p} to the components of the teaching input \mathbf{t} and the actual output \mathbf{q} to remark they refer to the input pattern \mathbf{p} . Accordingly, we could provide the following definition for the accumulated MSE:

$$E(\mathbf{w}) = \sum_{\mathbf{p} \in P} E_{\mathbf{p}}(\mathbf{w}) = \sum_{\mathbf{p} \in P} \frac{1}{M_O} \sum_{j=1}^{M_O} (t_{\mathbf{p},j} - q_{\mathbf{p},j})^2. \quad (1.16)$$

So far, we have not yet discussed how the update $\Delta w_{i,j}$ for the weight $w_{i,j}$ can be practically computed at each iteration. That is, we have still to rigourously define the functions h and g involved in the generalized Hebbian rule (1.12). Then, let us recall that any given operation carried out by the neural network can be thought as a point over the error surface $E(\mathbf{w})$. Therefore, to increase the performance of the network, we need to iteratively move toward a (local) minimum of the surface [6]. For this purpose, we may employ a *steepest descent* technique, thus following the direction of the anti-gradient, i.e.

$$\Delta w_{i,j} = -\eta \frac{\partial E(\mathbf{w})}{\partial w_{i,j}}, \quad (1.17)$$

$\eta > 0$ being the learning rate, modulating the size of the step; its role will be clearer later in the chapter. Among the others, the *backpropagation of error* [11] is surely the most-known supervised, gradient-based training procedure for a multi-layer feedforward neural network whose neurons are equipped with a *semi-linear*, i.e. continuous and differentiable, activation function⁴. The derivation of the backpropagation algorithm is provided in the following.

Backpropagation of error

Let us consider the generic synapse (i, j) of a multi-layer feedforward neural network, connecting the *predecessor* neuron i with the *successor* neuron j . As mentioned before, suppose both i and j present a semi-linear activation function. By widely exploiting the chain rule, the backpropagation of error provides an operative formula for the evaluation of the antigradient of the error function $E(\mathbf{w})$ in an arbitrary point \mathbf{w} , thus allowing to compute the update $\Delta w_{i,j}$ for the weight $w_{i,j}$ according to (1.17). For this purpose, we shall need to distinguish whether the successor neuron j is either an output or an inner neuron. To improve the clarity of the following mathematical derivation, we shall decorate any variable concerning a neuron with the subscript \mathbf{p} , thus to explicitly specify the input pattern \mathbf{p} we refer to.

Let $S = \{s_1, \dots, s_m\}$ the set of m neurons sending their output to j through the synapses $\{(s_1, j), \dots, (s_m, j)\}$. Recalling the definition (1.14) of the accumulated error, via the chain rule we can formulate the derivative of $E(\mathbf{w})$ with respect to the weight $w_{i,j}$ as follows:

$$\frac{\partial E(\mathbf{w})}{\partial w_{i,j}} = \sum_{\mathbf{p} \in P} \frac{\partial E_{\mathbf{p}}(\mathbf{w})}{\partial w_{i,j}} = \sum_{\mathbf{p} \in P} \frac{\partial E_{\mathbf{p}}(\mathbf{w})}{\partial u_{\mathbf{p},j}} \frac{\partial u_{\mathbf{p},j}}{\partial w_{i,j}}. \quad (1.18)$$

Let us focus on the product

$$\frac{\partial E_{\mathbf{p}}(\mathbf{w})}{\partial u_{\mathbf{p},j}} \frac{\partial u_{\mathbf{p},j}}{\partial w_{i,j}} \quad (1.19)$$

occurring in Equation (1.18). Assuming the propagation function be represented as the weighted sum, so that the net input is given by (1.4), the second factor in (1.19) reads:

$$\frac{\partial u_{\mathbf{p},j}}{\partial w_{i,j}} = \frac{1}{\partial w_{i,j}} \left(\sum_{s \in S} w_{s,j} y_{\mathbf{p},s} - \theta_j \right) = y_{\mathbf{p},i}. \quad (1.20)$$

We then denote the opposite of the first term in (1.19) by $\delta_{\mathbf{p},j}$, i.e.

$$\delta_{\mathbf{p},j} = -\frac{\partial E_{\mathbf{p}}(\mathbf{w})}{\partial u_{\mathbf{p},j}}. \quad (1.21)$$

⁴Therefore, backpropagation cannot apply with a binary activation function.

$\delta_{\mathbf{p},j}$ is often referred to as the *sensitivity* of the specific error $E_{\mathbf{p}}$ with respect to neuron j . Plugging (1.21) and (1.20) into (1.18), then embedding the resulting equation into (1.17), the weight update can be cast in the form:

$$\Delta w_{i,j} = \eta \sum_{\mathbf{p} \in P} \delta_{\mathbf{p},j} y_{\mathbf{p},i}. \quad (1.22)$$

We now proceed to derive an operative formula for $\delta_{\mathbf{p},j}$. Consider the case j is an output neuron. Supposing a specific error of the form (1.15) and an identity output function, it follows:

$$\begin{aligned} \delta_{\mathbf{p},j} &= -\frac{\partial E_{\mathbf{p}}(\mathbf{w})}{\partial y_{\mathbf{p},j}} \frac{\partial y_{\mathbf{p},j}}{\partial u_{\mathbf{p},j}} \\ &= -\frac{1}{\partial y_{\mathbf{p},j}} \frac{1}{M_O} \sum_{k=1}^{M_O} (t_{\mathbf{p},k} - y_{\mathbf{p},k})^2 \frac{\partial a_{\mathbf{p},j}}{\partial u_{\mathbf{p},j}} \\ &= \frac{2}{M_O} (t_{\mathbf{p},j} - y_{\mathbf{p},j}) \frac{\partial f_{act}(u_{\mathbf{p},j})}{\partial u_{\mathbf{p},j}} \\ &= \frac{2}{M_O} (t_{\mathbf{p},j} - y_{\mathbf{p},j}) f'_{act}(u_{\mathbf{p},j}). \end{aligned} \quad (1.23)$$

It worths mention here that (1.23) implies the derivative f'_{act} of the activation function f_{act} with respect to its argument. This motivates the requirement of a differentiable transfer function.

On the other hand, Equation (1.23) does not apply when j lies within an hidden layer. In fact, no teaching input is provided for an hidden neuron. In that case, let us denote by $R = \{r_1, \dots, r_n\}$ the set of n neurons receveing the output generated by j , i.e. the *successors*. We then point out that the output of any neuron can directly affect only the neurons which receive the output itself, i.e. the successors [9]. Hence:

$$\begin{aligned} \delta_{\mathbf{p},j} &= \frac{\partial E_{\mathbf{p}}(\mathbf{w})}{\partial u_{\mathbf{p},j}} \\ &= \frac{\partial E_{\mathbf{p}}(u_{\mathbf{p},r_1}, \dots, u_{\mathbf{p},r_n})}{\partial u_{\mathbf{p},j}} \\ &= \frac{\partial E_{\mathbf{p}}(u_{\mathbf{p},r_1}, \dots, u_{\mathbf{p},r_n})}{\partial y_{\mathbf{p},j}} \frac{\partial y_{\mathbf{p},j}}{\partial u_{\mathbf{p},j}} \\ &= \sum_{k=1}^n \frac{\partial E_{\mathbf{p}}}{\partial u_{\mathbf{p},r_k}} \frac{\partial u_{\mathbf{p},r_k}}{\partial y_{\mathbf{p},j}} \frac{\partial y_{\mathbf{p},j}}{\partial u_{\mathbf{p},j}}. \end{aligned} \quad (1.24)$$

Applying the definition of sensitivity for neuron r_k , $k = 1, \dots, n$, under the same assumptions of Equation (1.23) we can further develop (1.24):

$$\begin{aligned} \delta_{\mathbf{p},j} &= \sum_{k=1}^n (-\delta_{\mathbf{p},r_k}) \frac{1}{\partial y_{\mathbf{p},j}} \left(\sum_{l=1}^{m_k} w_{s_l, r_k} y_{\mathbf{p},s_l} - \theta_{r_k} \right) \frac{\partial f_{act}(u_{\mathbf{p},j})}{\partial u_{\mathbf{p},j}} \\ &= - \sum_{k=1}^n \delta_{\mathbf{p},r_k} w_{j, r_k} f'_{act}(u_{\mathbf{p},j}). \end{aligned} \quad (1.25)$$

Here, we have supposed that r_k receives input signals from m_k neurons $\{s_1, \dots, s_{m_k}\}$.

Let us now collect and summarize the results derived so far. At any iteration t of the backpropagation learning algorithm, the weight $w_{i,j}(t)$ of a generic connection (i, j) linking neuron i with neuron j is corrected by an additive quantity $\Delta w_{i,j}(t)$, i.e.

$$w_{i,j}(t+1) = w_{i,j}(t) + \Delta w_{i,j}(t).$$

Assume the accumulated mean squared error (1.16) as performance function. Understanding the dependence on time for the sake of clarity, the weight update $\Delta w_{i,j}$ reads:

$$\Delta w_{i,j} = \eta \sum_{\mathbf{p} \in P} y_{\mathbf{p},i} \delta_{\mathbf{p},j},$$

where $\eta > 0$, while $\delta_{p,j}$ is given by

$$\delta_{p,j} = \begin{cases} f'_{act}(u_{p,j}) \sum_{r \in R} \delta_{p,r} w_{j,r}, & \text{if } j \text{ inner neuron} , \\ \frac{2}{M_O} f'_{act}(u_{p,j}) (t_{p,j} - y_{p,j}), & \text{if } j \text{ output neuron} . \end{cases} \quad (1.26a)$$

$$(1.26b)$$

Some relevant remarks about the overall algorithm should be pointed out. First, observe that Equation (1.26a) defines $\delta_{p,j}$ for a hidden node j by relying on neurons in the following layer, whereas Equation (1.26b) only involves variables concerning the neuron (namely, $u_{p,j}$ and $y_{p,j}$) and the exact output $t_{p,j}$. Therefore, the coupled equations (1.26a)-(1.26b) implicitly set the order in which the weights must be adjusted: starting from the output layer, update all the connections ending in that layer, then move backwards to the preceeding layer. In this way, the error is *backpropagated* from the output down to the input, leaving traces in each layer of the network [9, 16].

The weight updating procedure detailed above corresponds to the offline version of the backpropagation of error, since it involves the total error E . On the other hand, the online algorithm readily comes from Equation (1.22): simply drop the summation over the elements of the training set P .

Although intuitive and very promising, backpropagation of error suffers of all those drawbacks that are peculiar to gradient-based techniques. For instance, we may get stuck in a local minimum, whose level is possibly far from the global minimum of the error surface E . Furthermore, since the step size dictated by the gradient method is given by the norm of the gradient itself, minima close to steepest gradients are likely to be missed due to a large step size. This motivates the introduction in (1.17) of the learning rate $\eta \in (0, 1]$, acting as a reducing factor and thus enabling a keener control on the descent. As suggested by Kriesel [9], in many applications reasonable values for η lies in the range $[0.01, 0.9]$. In particular, a time-dependent learning rate usually enables a more effective and more efficient training procedure. At the beginning of the process, when the network is far from the application goal, one often needs to span a large extent of the error surface, thus to identify a region of interest. Then, the learning rate should be large, i.e. close to 1, thus to speed up the exploration. However, as we approach the optimal configuration, we may want to progressively reduce the learning rate, then the step size, to fine-tune the weights of the neural network.

Nevertheless, selecting an appropriate value for η is still more an art than a science. Furthermore, for sigmoidal activations functions as the ones shown in Figure ??, the derivative is close to zero far from the origin. This results in the fact that it becomes very difficult to move neurons away from the limits of the activation, which could extremely extend the learning time [9]. Hence, different alternatives to backpropagation have been proposed in literature, either by modifying the original algorithm (as, e.g., the resilient backpropagation [13] or the quickpropagation [3]), or pursuing a different numerical approach to the optimization problem. The latter class of algorithms includes the Levenberg-Marquardt algorithm, which we shall extensively use in our numerical tests.

Levenberg-Marquardt algorithm

While backpropagation is a steepest descent algorithm, the Levenberg-Marquardt algorithm [10] is an approximation to the Newton's method [4]. As for backpropagation, the learning procedure is driven by the loss function $E = E(\mathbf{w})$, $\mathbf{w} \in |\mathcal{V}|$. Applying the Newton's method for the minimization of E , at each iteration the *search direction* $\Delta \mathbf{w}$ is found by solving the following linear system:

$$\nabla^2 E(\mathbf{w}) \Delta \mathbf{w} = -\nabla E(\mathbf{w}), \quad (1.27)$$

where $\nabla E(\mathbf{w})$ and $\nabla^2 E(\mathbf{w})$ denotes, respectively, the gradient vector and the Hessian matrix of E with respect to its argument \mathbf{w} . Assume now that the loss function is represented as the accumulated mean squared error,

$$E(\mathbf{w}) = \sum_{p \in P} \frac{1}{M_O} \sum_{j=1}^{M_O} (t_{p,j} - y_{p,j})^2 = \sum_{p \in P} \frac{1}{M_O} \sum_{j=1}^{M_O} e_{p,j}(\mathbf{w})^2; \quad (1.28)$$

with $e_{\mathbf{p},j}$ the j -th component of the error vector $\mathbf{e}_{\mathbf{p}} = \mathbf{t}_{\mathbf{p}} - \mathbf{y}_{\mathbf{p}}$ corresponding to the input pattern \mathbf{p} . Then, introducing the Jacobian $J_{\mathbf{p}}$ of the specific error vector $\mathbf{e}_{\mathbf{p}}$ with respect to \mathbf{w} , i.e.

$$J_{\mathbf{p}}(\mathbf{w}) = \begin{bmatrix} \frac{\partial e_{\mathbf{p},1}}{\partial w_1} & \frac{\partial e_{\mathbf{p},1}}{\partial w_2} & \cdots & \frac{\partial e_{\mathbf{p},1}}{\partial w_{|\mathcal{V}|}} \\ \frac{\partial e_{\mathbf{p},2}}{\partial w_1} & \frac{\partial e_{\mathbf{p},2}}{\partial w_2} & \cdots & \frac{\partial e_{\mathbf{p},2}}{\partial w_{|\mathcal{V}|}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial e_{\mathbf{p},M_O}}{\partial w_1} & \frac{\partial e_{\mathbf{p},M_O}}{\partial w_2} & \cdots & \frac{\partial e_{\mathbf{p},M_O}}{\partial w_{|\mathcal{V}|}} \end{bmatrix} \in \mathbb{R}^{M_O \times |\mathcal{V}|} \quad (1.29)$$

simple computations yield:

$$\nabla E(\mathbf{w}) = \sum_{\mathbf{p} \in P} \frac{2}{M_O} J_{\mathbf{p}}(\mathbf{w})^T \mathbf{e}_{\mathbf{p}} \in \mathbb{R}^{|\mathcal{V}|} \quad (1.30)$$

and

$$\nabla^2 E(\mathbf{w}) = \sum_{\mathbf{p} \in P} \frac{2}{M_O} [J_{\mathbf{p}}(\mathbf{w})^T J_{\mathbf{p}}(\mathbf{w}) + S(\mathbf{w})] \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}, \quad (1.31)$$

with

$$S(\mathbf{w}) = \sum_{\mathbf{p} \in P} \frac{2}{M_O} \sum_{j=1}^{M_O} e_{\mathbf{p},j} \nabla^2 e_{\mathbf{p},j}.$$

Assuming $S(\mathbf{w}) \approx 0$, inserting (1.30) and (1.31) into (1.27) we get the linear system

$$\left[\sum_{\mathbf{p} \in P} J_{\mathbf{p}}(\mathbf{w})^T J_{\mathbf{p}}(\mathbf{w}) \right] \Delta \mathbf{w} = - \sum_{\mathbf{p} \in P} J_{\mathbf{p}}(\mathbf{w})^T \mathbf{e}_{\mathbf{p}}(\mathbf{w}). \quad (1.32)$$

The Levenberg-Marquardt modification to the Newton's method reads [4, 10]:

$$\left[\sum_{\mathbf{p} \in P} J_{\mathbf{p}}(\mathbf{w})^T J_{\mathbf{p}}(\mathbf{w}) + \mu I \right] \Delta \mathbf{w} = - \sum_{\mathbf{p} \in P} J_{\mathbf{p}}(\mathbf{w})^T \mathbf{e}_{\mathbf{p}}(\mathbf{w}), \quad (1.33)$$

with $\mu \geq 0$ and I the identity matrix of size $|\mathcal{V}| \times |\mathcal{V}|$. Note that if $\mu = 0$ we recover the Newton's method (1.32), while for $\mu \gg 1$ the search direction $\Delta \mathbf{w}$ approaches the antigradient of E , i.e. we recover the backpropagation algorithm. Then, the Levenberg-Marquardt algorithm can be seen as an interpolation between the Newton's method and the steepest descent method, aiming to retain the advantages of both techniques.

The Levenberg-Marquardt training algorithm proceeds as follows. At each epoch t of the training procedure, we solve the (potentially large) linear system (1.33). Whenever the step $\Delta \mathbf{w}(t)$ leads to a reduction in the performance function, i.e. $E(\mathbf{w}(t) + \Delta \mathbf{w}(t)) < E(\mathbf{w}(t))$, the parameter μ is reduced by a factor $\beta > 1$. Conversely, if $E(\mathbf{w}(t) + \Delta \mathbf{w}(t)) > E(\mathbf{w}(t))$ the parameter is multiplied by the same factor β . This reflects the idea that far from the actual minimum we should prefer the gradient method to the Newton's method, since the latter may diverge. Yet, once in a neighborhood of the minimum, we switch to the Newton's method so to exploit its faster convergence [10].

The key step in the algorithm is the computation of the Jacobian $J_{\mathbf{p}}(\mathbf{w})$ for each training vector \mathbf{p} . Suppose that the k -th element w_k of \mathbf{w} represents the weight $w_{i,j}$ of the connection (i, j) , for some i and j , with $1 \leq i, j \leq |\mathcal{N}|$. Then, the (h, k) -th entry of $J_{\mathbf{p}}(\mathbf{w})$ is given by:

$$\frac{\partial e_{\mathbf{p},h}}{\partial w_k} = \frac{\partial e_{\mathbf{p},h}}{\partial w_{i,j}}, \quad 1 \leq h \leq M_O, 1 \leq k \leq |\mathcal{V}|. \quad (1.34)$$

We recognize that the derivative on the right-hand side of Equation (1.34) is intimately related with the gradient of the performance function. Therefore, we can follow the very same steps performed to derive the backpropagation, namely:

$$\begin{aligned}
 \frac{\partial e_{\mathbf{p},h}}{\partial w_{i,j}} &= \frac{\partial e_{\mathbf{p},h}}{\partial u_{\mathbf{p},j}} \frac{\partial u_{\mathbf{p},j}}{\partial w_{i,j}} \\
 &= \frac{\partial e_{\mathbf{p},h}}{\partial y_{\mathbf{p},j}} \frac{\partial y_{\mathbf{p},j}}{\partial u_{\mathbf{p},j}} \frac{\partial u_{\mathbf{p},j}}{\partial w_{i,j}} \\
 &= \frac{\partial e_{\mathbf{p},h}}{\partial y_{\mathbf{p},j}} f'_{act}(u_{\mathbf{p},j}) y_{\mathbf{p},i} \\
 &= \delta_{\mathbf{p},h,j} y_{\mathbf{p},i},
 \end{aligned} \tag{1.35}$$

with

$$\delta_{\mathbf{p},h,j} := -\frac{\partial e_{\mathbf{p},h}}{\partial u_{\mathbf{p},j}} = -\frac{\partial e_{\mathbf{p},h}}{\partial y_{\mathbf{p},j}} f'_{act}(u_{\mathbf{p},j}). \tag{1.36}$$

For the computation of the derivative $\partial e_{\mathbf{p},h} / \partial y_{\mathbf{p},j}$ occurring in (1.36), further assume that within the set of neurons \mathcal{N} items are ordered such that the output neurons come first, i.e.

$$j \text{ output neuron} \Leftrightarrow 1 \leq j \leq M_O.$$

We can then distinguish three cases:

- (i) j output neuron, $j = h$: since $e_{\mathbf{p},h} = e_{\mathbf{p},j} = t_{\mathbf{p},j} - y_{\mathbf{p},j}$, then

$$\frac{\partial e_{\mathbf{p},h}}{\partial y_{\mathbf{p},j}} = -1; \tag{1.37}$$

- (ii) j output neuron, $j \neq h$: the output of an output neuron can not influence the signal fired by another output neuron, hence

$$\frac{\partial e_{\mathbf{p},h}}{\partial y_{\mathbf{p},j}} = 0; \tag{1.38}$$

- (iii) j inner neuron: letting R be the set of successors of j , similarly to (1.25) the chain rule yields

$$\frac{\partial e_{\mathbf{p},h}}{\partial y_{\mathbf{p},j}} = \sum_{r \in R} \frac{\partial e_{\mathbf{p},h}}{\partial u_{\mathbf{p},r}} \frac{\partial u_{\mathbf{p},r}}{\partial y_{\mathbf{p},j}} = - \sum_{r \in R} \delta_{\mathbf{p},h,r} w_{j,r}. \tag{1.39}$$

Ultimately, at any iteration of the learning algorithm the entries of the Jacobian matrix $J_{\mathbf{p}}$ are given by

$$\frac{\partial e_h}{\partial w_k} = \delta_{\mathbf{p},h,j} y_{\mathbf{p},i}, \quad 1 \leq h \leq M_O, 1 \leq k \leq |\mathcal{V}|, w_k = w_{i,j} \text{ for some } i \text{ and } j,$$

with

$$\delta_{\mathbf{p},h,j} = \begin{cases} f'_{act}(u_{\mathbf{p},j}) \sum_{r \in R} \delta_{\mathbf{p},h,r} w_{j,r}, & \text{if } j \text{ inner neuron} , \\ f'_{act}(u_{\mathbf{p},j}) \delta_{jh}^K, & \text{if } j \text{ output neuron} , \end{cases} \tag{1.40a}$$

$$\tag{1.40b}$$

where δ_{jh}^K is the Kronecker delta.

Let us finally remark that a trial and error approach is still required to find satisfactory values for μ and β [10]. Moreover, the dimension of the system (1.33) nonlinearly increases with the number of neurons in the network, making the Levenberg-Marquardt algorithm poorly efficient for large networks [4]. However, it is more efficient than backpropagation for networks with a few hundreds of connections, besides producing much more accurate results. We shall gain further insights into this topic in Chapter ??.

1.2.4 Practical considerations on the design of artificial neural networks

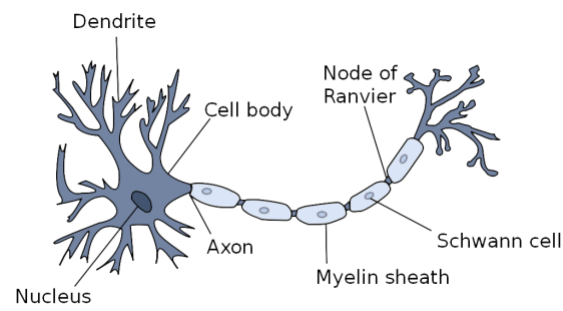


Figure 1.2.1

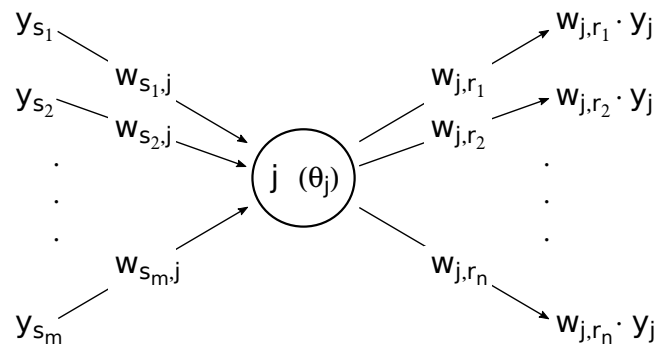


Figure 1.2.2

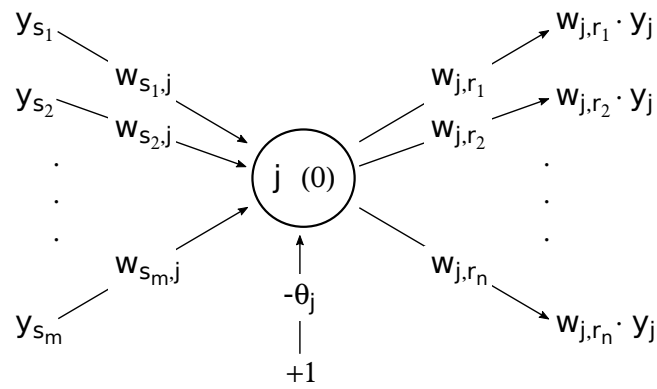
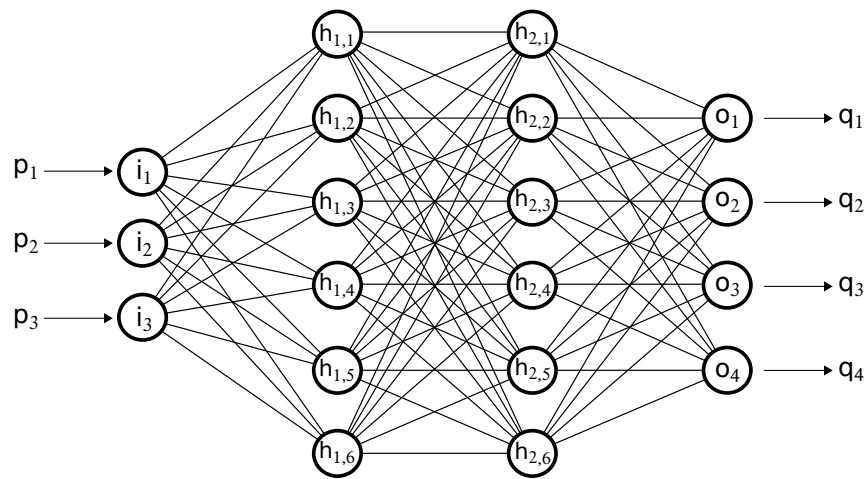


Figure 1.2.3

**Figure 1.2.4**

Chapter 2

Reduced Basis method for nonlinear elliptic PDEs

In this chapter, we derive a Reduced Basis (RB) approximation for a parametrized Boundary Value problem (BVP) involving a possibly nonlinear elliptic Partial Differential Equation (PDE). As a convenient test case, the nonlinear Poisson equation in one and two spatial dimensions is considered, with the nonlinearity stemming from a solution-dependent viscosity. In the one-dimensional case, the parameters may characterize the viscosity itself, the forcing term, or the boundary conditions. On the contrary, in the two-dimensional framework we are concerned with a parameter-free PDE defined on a domain undergoing geometrical transformations. In the RB context, when dealing with shape variations one needs to introduce a parametric transformation, mapping the whole domain to a parameter-independent reference configuration, e.g. a unit square [?]. For this aim, we refer to the boundary displacement-dependent transfinite map (BDD TM) proposed in [?], building a volumetric parametrization given the boundary parametrization of the domain.

For the sake of numerical efficiency, the RB procedure is usually carried out within an offline-online framework. The *offline* stage consists in generating the reduced basis out of an ensemble of *snapshots*, i.e. high-fidelity numerical solutions to the BVP for different realizations of the parameters. In this thesis, the well-known Proper Orthogonal Decomposition (POD) technique is used, relying on the Finite Element (FE) method for the computation of the snapshots. Then, given new values for the parameters, in the *online* phase one seeks an approximation to the high-fidelity solution in the reduced space, i.e. the linear space generated by the reduced basis. To retain the well-posedness of the problem, the system yielded by the FE method is projected onto the reduced space, thus enforcing the orthogonality to the reduced basis of the residual for each equation [?, ?]. Hence, the computational cost associated with the *resolution* of the *reduced* model is *independent* of the size of the original, large-scale model.

The whole procedure sketched so far is detailed in the following sections. Section ?? and ?? discuss the Finite Element method applied to the one- and two-dimensional Poisson equation, respectively. The Reduced Basis technique is described in Section ??, while in the final Section ?? we present an alternative, Neural Network-based approach for the computation of the reduced solution, which represents the actual novelty of the proposed reduced order modeling algorithm.

2.1 Finite Element method for one-dimensional Poisson

Bibliography

- [1] G. Cybenko. *Continuous valued neural networks with two hidden layers are sufficient*. Technical Report, Department of Computer Science, Tufts University, 1988.
- [2] G. Cybenko. *Approximation by superpositions of a sigmoidal function*. Mathematics of Control, Signals, and Systems, 2(4):303–314, 1989.
- [3] S. E. Fahlman. *An empirical study of learning speed in back-propagation networks*. Technical Report CMU-CS-88-162, CMU, 1988.
- [4] M. T. Hagan, M. B. Menhaj. *Training feedforward networks with the Marquardt algorithm*. IEEE Transactions on Neural Networks, 5(6):989–993, 1994.
- [5] M. T. Hagan, H. B. Demuth, M. H. Beale, O. De Jesús. *Neural Network Design, 2nd Edition*. eBook, 2014.
- [6] S. Haykin. *Neural Networks: A comprehensive foundation*. Prentice Hall, Upper Saddle River, NJ, 2004.
- [7] D. O. Hebb. *The organization of behaviour: A neuropsychological theory*. John Wiley & Sons, New York, NY, 1949.
- [8] L. P. Kaelbling, M. L. Littman, A. W. Moore. *Reinforcement Learning: A Survey*. Journal of Artificial Intelligence Research, 4:237–285, 1996.
- [9] D. Kriesel. *A Brief Introduction to Neural Networks*. http://www.dkriesel.com/en/science/neural_networks.
- [10] D. W. Marquardt. *An algorithm for least-squares estimation of nonlinear parameters*. Journal of the Society for Industrial and Applied Mathematics, 11(2):431–441, 1963.
- [11] J. L. McClelland, D. E. Rumelhart. *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*. MIT Press, Cambridge, UK, 1986.
- [12] M. A. Nielsen. *Neural Networkd and Deep Learning*. Determination Press, 2015.
- [13] M. Riedmiller, H. Braun. *A direct adaptive method for faster backpropagation learning: The rprop algorithm*. Neural Networks, IEEE International Conference on, 596–591, 1993.
- [14] F. Rosenblatt. *The perceptron: A probabilistic model for information storage and organization in the brain*. Psychological Review, 65:386–408, 1958.
- [15] C. Stergiou, D. Siganos. *Neural Networks*. https://www.doc.ic.ac.uk/~nd/surprise_96/journal/vol4/cs11/report.html#Introductiontoneuralnetworks.
- [16] B. Widrow, M. E. Hoff. *Adaptive switching circuits*. Proceedings WESCON, 96–104, 1960.