

Markov chain network training and conservation law approximations: Linking microscopic and macroscopic models for evolution

Roderick V.N. Melnik

M² NeT Lab, Wilfrid Laurier University, 75 University Avenue West, Waterloo, ON, Canada N2L 3C5

Abstract

In this paper, a general framework for the analysis of a connection between the training of artificial neural networks via the dynamics of Markov chains and the approximation of conservation law equations is proposed. This framework allows us to demonstrate an intrinsic link between microscopic and macroscopic models for evolution via the concept of perturbed generalized dynamic systems. The main result is exemplified with a number of illustrative examples where efficient numerical approximations follow directly from network-based computational models, viewed here as Markov chain approximations. Finally, stability and consistency conditions of such computational models are discussed.

© 2007 Elsevier Inc. All rights reserved.

Keywords: Training dynamics; Neural networks; Complex dynamic systems; Markov chain approximation method; Conservation law approximations

1. Introduction

Many concepts and tools used in information theory, control, theory of approximation, dynamical systems and artificial intelligence are closely linked. Among such tools are neural networks. The success in using neural networks in many application areas is well documented in the literature. This includes, but not limited to, aeronautics, pattern and speech recognition (in particular in the context of various classification problems) [6,10], optoelectronics, robotics and autonomous navigation, determining structure–property relationships in material science applications [35,33], constructing dynamic observers [3], identification and control of complex biotechnological cycles, industrial processes, and nonlinear systems in general [36,7], modelling complex dynamic systems and nonlinear phenomena such as hysteresis [16]. In many cases, a key to this success is kept by various associations of (computational) neural networks with *the atomistic approach* to molecular design aimed at determining some relationships between the properties of the structure (thermomechanical, electromagnetic, etc.) and the structure itself described at the microscopic (mesoscopic, and eventually macroscopic)

E-mail address: rmelnik@wlu.ca

level. Another set of tools and successful approaches to structure modelling is *the macroscopic approach* essentially based on the conservation law equations. In this contribution, we developed a general framework for establishing a link between neural network concepts and numerical approximations of conservation laws.

One of the major objectives in neural network theory is the study of constructive approaches to the design of effective learning algorithms [14]. We observe that models for learning of neural networks and models for the evolution of dynamic systems are closely connected. The idea of using neural networks as components in dynamic systems is now well established in the literature (e.g. [28]). Models for the dynamic interactions between evolution and learning have been used by a number of authors who use evolutionary algorithms for the optimization of neural structures. It is often argued that artificial neural networks (ANN) and evolutionary algorithms (EA) can be efficiently combined together. Indeed, if ANN is interpreted as a model for biological nervous systems and the learning process, then EA can be interpreted as a model for biological evolution itself and the process of adaptation of large scales [32]. In this case, ANN becomes a useful mathematical tool for linking deterministic and probabilistic approaches in the approximation theory. This idea was developed in [23,24] where the evolution was proposed to be modelled by a generalized equation with training/learning rules (algorithms) understood in a probabilistic sense. The problem of learning in neural network theory is formulated in terms of the minimization of an error function which is a function of adaptive parameters (weights and biases). In developing training/learning algorithms for ANN we have to deal with incomplete data, and an initial step of this procedure can be associated with the initializing of the weights in the network. Assigning these weights random values [11], leads us to a probabilistic framework [34]. Postprocessing of the available information is required for deriving deterministic models. For example, given a (training) set of input–output pairs $(u(i), y_i)$ ($i = 1, \dots, n$ for n training/example subjects), we can “smooth” the data. The task is to construct a map which provides a good generalization (i.e. for given u that does not belong to this map should provide a reasonable estimate/prediction of the unobservable output). This is reducible to finding the best (functional) approximation to multivariate empirical (possible noisy, sparse, with some unobservable regions) data. However, if some additional “smoothing” constraints are used, we solve not the original, but a regularized problem. This provides a link between standard (Tichonov-like) approaches to deterministic, but ill-posed, problems and statistical (e.g. the Bayesian) approaches [8]. In this contribution, we explore further this link between deterministic and probabilistic approaches in neural network theory by considering learning and evolution in an intrinsic dynamic connection.

The starting point of our discussion is a formalization of the learning/training process in terms of the minimization of an error function based on the neural network approach. In particular, starting with a minimal structure (no hidden layers), according to certain rules, new connections, neurons, and layers are added. Most commonly used are related to a (probabilistic) adaptation of the network weights, a procedure which is intrinsically coupled with the overall network size (number of neurons on each layer and the total number of layers). Since the size and topology of neural networks are closely connected with the required training time, in dealing with the training/learning problem we will have to construct *dynamic* models. Such models arise naturally in applications of ANN to optimal control problems in particular in those areas where system parameters are time varying [21]. The approach developed in this contribution has common features with the on-line evolution approach (previously discussed in the literature in the context of game theory and control problems [1]) in a sense that our procedure can be interpreted as an on-line evolution algorithm applied to a certain time-dependent equation. The type of this equation is motivated by the fact that network architectures can be established by approximating the dynamic programming equations [4]. In this paper, we also use a dynamic equation that corresponds to a certain neural network architecture and is derived from the framework of generalized dynamic systems (GDS) developed in [23,24].

In what follows, although we focus our attention on feedforward-based networks complemented by back-propagation-like learning numerical procedures, the underlying procedure is the same for other architectures of neural systems where we have to define the connections between layers, the parameters (e.g., initial weights) and some learning rules. We demonstrate that these connections can be established by combining “forward evolution” (using neural networks) and “backward evolution” (using Markov chains associated with the original process). This combination ensures the minimization of uncertainty in the following sense. By processing information in response to discrete and continuous inputs, the ANN allows us to model dynamic systems. Since without a correction mechanism uncertainty may increase dramatically over time, we use Markov chains

to construct an appropriate mechanism where the architecture of the ANN depends on the cone of macroscopic events [24]. Although our approach is completely different in principle from the classical Gelenbe approach (see, e.g. [2]), we note that in both cases one has to exploit the idea of neural-network random structures. Since the original time-dependent models are discretized in our approach, one can interpret the resulting scheme via interactions among neurons so it is possible to calculate probabilities of activation of network neurons in a way similar as it is done in the Gelenbe approach.

Finally, we note that when applying ANNs to complex dynamic systems, one of the most important issues is to stabilize the learning mechanism [15] (e.g., one has to avoid the unbounded growth of the adjustable parameter in ANN). We show in this paper that this non-trivial task can be linked directly to discretizing conservation law models, and that the ANN stability can be treated effectively through the stability of numerical approximations of conservation laws, rather than through the Lyapunov stability (e.g. [20]).

2. Systems theoretic framework for constructing neural networks and dynamic training rules

By applying neural networks to learning and identification of dynamic systems, we attempt to influence the behavior of these systems by some components built-in into the systems [28,27]. Hence, the procedures for controlling such dynamic systems using ANN should be connected in one way or another with approximations of system dynamics. Moreover, since control, in addition to the requirements of fast and accurate, should ensure stability and robustness of the system, the stability of control is closely connected in such cases with the stability of numerical approximations of models describing the systems dynamics.

Recall that in most typical situations a nonlinear dynamic system can be described from a systems theoretic point of view by the following set of equations:

$$\begin{cases} x(k+1) = f[x(k), u(k)], & f(0, 0) = 0, \\ y(k) = h[x(k)], & h(0) = 0, \end{cases} \quad (2.1)$$

where $u(k) \in \mathbb{R}^l$, $y(k) \in \mathbb{R}^m$ and $x(k) \in \mathbb{R}^n$ are input, output, and state vectors, respectively, at discrete-time k , and mappings $f: \mathbb{R}^n \times \mathbb{R}^l \rightarrow \mathbb{R}^n$ and $h: \mathbb{R}^n \rightarrow \mathbb{R}^m$ are given functions, while control based on the state vector and control based only on input–output data need not be the same [28]. Model (2.1) provides adaptive means for implementing digital neuro-controlled systems and for further insight into the neural learning and adaptation. We note that instead of model (2.1), the dynamic system model can be determined through identification (e.g., with multi-layered neural networks) using the data obtained by a recurrence equation describing the output signal and taking into account the internal perturbations (which influence the output signal):

$$y(k+1) = \tilde{h}(x(k+1), y(k), u(k), e(k+1)). \quad (2.2)$$

New output $y(k+1)$, determined from the state vector $x(k+1)$ by using the given input data $[y(k), u(k)]$, will introduce error $e(k+1)$, representing internal perturbations at time $k+1$. For example, it can mimic the effects of synaptic time delay representing the “memory” of neurons which may differ from connection to connection. Mapping \tilde{h} in (2.2) can be viewed as a perturbed function h . In principle, this equation, which combines the information given by model (2.1), is sufficient to describe the training procedure. From a numerical point of view such a nonlinear-control-framework consideration [27] can be viewed as an approximation to the first equation in (2.1), describing the evolution of states of the system, supplemented by an approximation of the “constitutive” law (the second equation in (2.1)), describing some properties of the system, control rules, and/or requirements imposed by the user. Hence, if the original description of the system is set by using control-theoretic approach, model (2.2) can be viewed as an approximation to the Hamilton–Jacobi–Bellman-type (HJB) equation describing some “energetic” characteristic of the system such as cost or value function [23–25]. If the entire controlled process is modelled by the ANN, then the training/learning process for the neural network (using specific outputs with each of several inputs) can be associated with the dynamics of a certain conservation law model that describes the evolution of the nonlinear system. In this case, the inputs of ANN are associated with the initial conditions of the conservation law model. The description of the output–input training process will be understood here as a constructive backpropagation based on a

process-associated Markov chain approximation. This approximation will be linked to numerical approximations of conservation laws.

Following [27], we consider a neural network as a conveniently parameterized class of nonlinear maps. In the next paragraphs, we develop the methodology for linking this class of maps to conservation law approximations. In this context it is important to note that when a network is chosen to approximate a given mapping using IO data, the IO set is finite, whereas a conservation law defines an infinite set of input–output data. However, once such a conservation law model is approximated the process of neural network training and the approximate conservation laws can be described within the same general framework based on the concept of generalized dynamic systems [23,24]. In what follows, our consideration will be pertinent to discrete-time (even though the original system might be continuous in time), which reflects the fact that most complex systems are controlled by computers and therefore it is quite natural to consider approximations to complex dynamic systems as being discrete in time [29]. Moreover, any model, no matter how sophisticated it is, should account for uncertainties of the environment [23,24] which are easier to deal with in the discrete-time formulations [29]. Furthermore, the stability issues are much more transparent in the discrete space-time of events where dynamic systems evolve. Finally, in the context of neural networks the discrete-time consideration is the most natural way to proceed with the analysis because neural networks are proved to be universal approximators and any continuous function can be approximated arbitrarily well on a compact set with a (multilayer feed-forward, radial basis function, or another) neural network (e.g. [13,18,5,7]).

3. Approximating system Hamiltonians with neural networks

If the dynamic system under consideration is Hamiltonian, the problem of training neural networks with dynamic system rules governing the evolution of this system requires efficient procedures for approximating the Hamiltonian of such a system by a neural network. As a starting point, let T be a given set of times during which the network is trained, Σ is a state space of the dynamic system, U_T is a set of all permissible strategies for training, and X_T is the domain of definition of the system Hamiltonian assumed to be a compact Borel set [12]. Further constructions are based on the following fundamental property of applications of neural networks in the theory of approximation [13,18,5].

Theorem 3.1. *If $H \in \mathbb{L}^1(X_T)$, then for any arbitrarily small $\epsilon > 0$ there exists a network \tilde{H} such that:*

$$\|H - \tilde{H}\|_{\mathbb{L}^1(X_T)} < \epsilon. \quad (3.1)$$

For the description and control of complex dynamical systems via neural networks both feedforward (FNN) evolutionary networks (for the modelling of nonlinear input–output dynamics) and feedback evolutionary networks (for the implementation of training algorithms) are required. The link between those can be established on the basis of the Markov chain approximations as follows. First, recall that if H is a continuous function we can construct a uniform approximation (e.g., by using a FNN [13,18]). Note also that if smoothness of the function H is measured in terms of its Fourier representation an actual estimation of the network performance can also be obtained (e.g. [5]). Therefore, it is natural to associate a FNN model with the approximation \tilde{H} of the system Hamiltonian. For example, in the simplest case this approximation can be thought of as a linear combination with one layer of sigmoidal nodes or units in \mathbb{R}^n :

$$H_n(\mathbf{x}) = \sum_{i=1}^n \alpha_i \bar{x}_i(\mathbf{y}_i \cdot \mathbf{x} + \beta_i) + \alpha_0, \quad (3.2)$$

with $\mathbf{x}, \mathbf{y}_i \in \mathbb{R}^n$, $\mathbf{y}_i \cdot \mathbf{x}$ being their inner product, $\alpha_i, \beta_i \in \mathbb{R}$, and \bar{x}_i interpreted as the activation function of the hidden layer. The above approximation can be generalized to [34]:

$$H_n(\mathbf{x}) = x_0 \left(\sum_{i=1}^n \alpha_i \bar{x}_i(\mathbf{y}_i \cdot \mathbf{x} + \beta_i) + \alpha_0 \right), \quad (3.3)$$

where x_0 is understood as the activation function of the output and can be, for example, linear, sigmoidal, or a threshold unit. The point that we are making here is that such models alone may be of limited applicability

when applied in uncertain dynamic environments, because the quality of network performance is dependent strongly on its training capabilities with respect to the regularity of H and the structure of X_T . In considering computational models of system dynamics the structure of the network-approximator like (3.2) cannot be considered to be fixed, but rather it has to be adapted to the change of environment and internal perturbations. As a consequence, the most appropriate approximation in the framework of (3.1) can lead to a situation where activation function \bar{x}_t may not be continuous, subject to regularity of H . Furthermore, even if this function is continuous, in order to account for the topological structure of X_T each network layer has to be characterized by both feedforward and feedback operators.

Remark 3.1. In a special case, where $X_T \equiv \mathbb{R}^N$ and hence $H : \mathbb{R}^N \rightarrow \mathbb{R}^M$ ($N, M \in \mathbb{N}$), such operators may be interpreted through weight matrices [31].

The situation we described above is of the same nature as in control theory where a regularity balance between control and value functions has to be achieved [22,23]. In the context of our framework, this balance is achieved by simultaneous computational treatments of the “feedforward” and “backward” evolutionary processes, which can be effectively implemented by using along with feedforward approximators feedback networks such as block feedback networks (BFN). Below we explain this idea in detail.

Consider Hamiltonian H mapping X_T into \mathbb{R}^M and its computational equivalent in the form of a Turing computable function H_T such that:

$$H : X_T \rightarrow \mathbb{R}^M \quad \text{and} \quad H_T : \mathbb{N} \rightarrow \mathbb{N}. \tag{3.4}$$

Then, since BFN models have the same computing power as Turing machines (e.g. [31]), we conclude that there exists a BFN \tilde{H}_T such that for any input n from a subset of \mathbb{N} a finite number of network steps produces $H_T(n)$. Therefore, due to the Gödel numeration procedure (e.g. [9] and references therein), any Hamiltonian function $H : \mathbb{R}^N \rightarrow \mathbb{R}^M$ can be arbitrary well approximated by a network implementation of function H_T . Since both feedforward \tilde{H} and feedback \tilde{H}_T networks provide an approximation to H , the development of constructive algorithms for training requires further study of the connection between \tilde{H} and \tilde{H}_T . This connection between the concepts pertinent to feedforward and feedback networks is often overlooked. Examples where this connection becomes transparent are provided by associative memory networks or recurrent neural networks, where feedforward connections should be supplemented by some “memory” neurons or synaptic time delay, which may vary from connection to connection. A natural way to establish such a connection is to construct a Markov chain associated with the system evolution in such a way that it reflects the process of network training on the problem specific information.

Now, we are in a position to formulate the problem of training in terms of approximating the system Hamiltonian in such a way that the necessity of applying both forward and backward dynamic rules is transparent. For this purpose we consider a relatively simple case, where we assume that:

$$\bar{x}_t : \mathbb{R} \rightarrow \Sigma \tag{3.5}$$

is a sigmoidal function, meaning that:

$$\lim_{t \rightarrow -\infty} \bar{x}_t = 0, \quad \lim_{t \rightarrow +\infty} \bar{x}_t = 1. \tag{3.6}$$

The function \bar{x}_t here is the activation function for a neural network defined by its neurons as the following mapping:

$$\bar{x}_t \circ \mu : T \otimes \Sigma \otimes U_T \rightarrow \Sigma, \tag{3.7}$$

where $\mu : T \otimes \Sigma \otimes U_T \rightarrow \mathbb{R}$ is known as the decision maker (DM) function [24] such that it has to be adjusted (trained) so that the neural network (3.7) leads to an approximation of H . If the network (3.7) is trained with some dynamic rules by using a dynamic system whose Hamiltonian is H , then the process of approximation of function H in terms of (3.1) can be seen as the construction of a training strategy for a new network \tilde{H}_n^ϵ depending on the arbitrarily small positive parameter ϵ and arbitrarily large number of sigmoidal nodes of the associated network n (e.g. (3.2)) so that:

$$\|H - \tilde{H}_n^\epsilon\|_{L^1(X_T)} \rightarrow \min, \tag{3.8}$$

where $H \in \mathbb{L}^1(X_T)$. The main difficulty in the solution of problem (3.8) stems from a priori unknown character of dependency of the network on parameters which determine the function μ . In the most general setting, the problem of constructing the mapping μ is intrinsically connected with the definition of dynamic rules in singular stochastic control problems [23], and in interpreting (3.8) one expects that:

$$\text{if } \epsilon \rightarrow 0^+ \text{ and } n \rightarrow \infty \text{ then } \tilde{H}_n^\epsilon \rightarrow H. \quad (3.9)$$

However, for the constructive solution of (3.8) by using neural networks we need additional information. For example, if FNNs are applied to the solution of the problem, we need information on a coupling rule between ϵ , n , and the topology specified by X_T . This coupling rule will determine conditions for the system stability. If BFNs are applied to the approximation of H on an arbitrary set X_T , we need additional information on the network dimension and architecture [31]. This information can be made available only in a *sequential manner*, and the appropriate tool for the analysis of the network performance in such situations is a family of Discrete Markovian Decision Processes (DMDP). In this case a model for the training process, written in terms of the decision maker function μ , depends on Markov chain parameters with possible discontinuities that are dependent on values of the sigmoidal function \bar{x}_t . If the functional dependency of μ is chosen a priori, this may lead to the underestimation of possible irregularities of the function H which can happen if, e.g., smoothness of H is measured in terms of its Fourier representation [5,26]. In the next section, we aim at constructing a model that gives an approximation to μ with respect to some learning rules determined from the Markovian property of the process (\bar{x}_t, μ) .

4. Modelling dynamics of network training

From a constructive approximation point of view, the training process of neural networks describing dynamic systems evolution can be seen as the solution to the following problem. We have to construct a model (possibly, a hierarchy of models [22]) for the decision maker function μ in such a way that a neural network approximates the Hamiltonian of a dynamic system evolution, while the dynamics of this evolution is described by the process of network training.

Any two sets of input–output data in the training process can be represented by the mathematical model of a dynamic system that couples two space-time events of the system evolution, e_n and e_{n+1} , by a function of the perturbed velocity v_ϵ and the system Hamiltonian or its approximation H :

$$e_{n+1} = H(v_\epsilon, e_n), \quad n = 0, 1, \dots \quad (4.1)$$

The perturbed velocity is introduced to account for the changing environment and/or internal perturbations [24]. Then, if we specify a sequence of events (e_0, e_1, \dots) by temporal evolution and formalize the dynamics of the system by a discrete-time model, we obtain the following two equations [24]:

$$\begin{cases} x_{t+1} = H_1(v_1, x_t), \\ h_{\tau+1} = H_0(v_0, h_\tau), \end{cases} \quad (4.2)$$

where H_1 is an approximation to H and $v_1 = \lim_{\epsilon \rightarrow 0^+} v_\epsilon$, H_0 is an operator for sequential corrections of such an approximation needed due to system external/internal perturbations and v_0 is the function characterizing the rate of such perturbations. The dynamics of the system described by (4.2) operates on two coupled temporal scales, t and τ with spatial trajectories described by x_t and h_τ , respectively. In the simplest cases, the trajectory x_t may be approximated by \bar{x}_t defined in the previous section. Since perturbations (internal and/or external) are an intrinsic part of any dynamic system, we associate the class of neural networks with built-in training capability (determined by models for the decision maker function μ) as Infinite Length Perturbed Markov Chains (ILPMC) [24]. Since the form of the functional dependency of μ cannot be a priori fixed for the entire dynamic path, the weights and the network structure may change which leads to a situation much more complicated than traditionally dealt with. Some algorithms dealing with such a situation are known (e.g., cascade-correlation, pruning algorithms, etc.). However, our approach is different from the previously developed since it is based on the Markovian character of the process (x_t, μ) , associated with both system dynamics and the process of training, and the concept of generalized dynamic systems (with appropriate assumptions on dynamic stochastic rules [24]). In our approach, by using conditional probabilities of the

Markov chain, the velocity function between two macroscopic events in the evolution of the generalized dynamic system is introduced as a measure of changes which take place on the microscopic level with respect to the macroscopic behavior of the system. This velocity function provides a link between microscopic and macroscopic models for the evolution of dynamic systems. Similar to model (3.7), models (4.1) and (4.2) are interpreted in terms of the limit $\epsilon \rightarrow 0^+$ and $n \rightarrow \infty$, and the connection between them can be analyzed when learning rules for the system evolution in response to perturbations are introduced. In particular, in the limit (3.9) both sequences (x_0, x_1, \dots) and (h_0, h_1, \dots) merge. As an important part of the system dynamics, perturbations can be formalized from the very beginning of the modelling process as a decision making process with limited information available.

Remark 4.1. Since perturbed and unperturbed models might give rise to qualitatively distinct types of descriptions of system behavior for any arbitrary $\epsilon > 0$, the perturbation parameter alone cannot be an appropriate characteristic of the model uncertainty [24], and this fact results in two informational sequences in model (4.2) considered on two different time scales.

The learning rules are introduced into the model by the equation written in terms of function μ (Section 3). In particular, following [24,25], we arrive the following proposition.

Proposition 1. *If $H \in \mathbb{L}^1(X_T)$, then the process of training/learning can be modelled by the following equation:*

$$(1 + v_1) \left[\frac{\partial \mu}{\partial x} + \frac{1}{v_1} \left(\frac{\partial \mu}{\partial t} + f_0 \right) \right] = 0, \tag{4.3}$$

where, in the context of neural networks, v_1 is the velocity of information transmission between neurons, and f_0 is a training/learning goal defined by a priori knowledge/assumptions on X_T and the function H .

This process can be approximated by an appropriately constructed Markov chain associated with the evolution of the dynamic system under consideration. Since function μ is required to be adapted to a priori given information, the practical implementation of training procedures may lead to approximations of the network-activator functions by a piecewise-deterministic stochastic process. Such non-diffusion stochastic models have been previously studied in theory of DMDP and in theoretical physics (see [24] and references therein). Such models are multiscale models for the evolution of dynamic systems described by (4.2) and, as pointed out in [24], can be interpreted by a system of coupled differential equations on two scales:

$$\begin{cases} \dot{h}(\tau) = v_0(\tau, h, \mu), \\ \frac{\partial \mu}{\partial t} + v_1(t, x, \mu) \frac{\partial \mu}{\partial x} = 0. \end{cases} \tag{4.4}$$

We conclude this section with the following observation.

Remark 4.2. Both parts of the perturbed velocity functions v_0 and v_1 inherit their dependency on the decision-maker function μ . If two events between which GDS evolution has to be studied are specified (e.g., two sets of input–output data are entered), then a pair of functions $(h(\tau), \mu(t, x))$ gives the solution to the training problem under consideration.

Next, we describe a procedure that allows us to approximate this pair.

5. Network-based computational models for dynamics as Markov chain approximations

Although the perturbed system dynamics x_t^ϵ might not be governed by the Markovian property (and the function x_t^ϵ might not be sigmoidal), the pair of functions $(h(\tau), \mu(t, x))$ does possess the Markovian property [24]. Therefore, the key idea here is based on the construction of a Markov chain approximation simultaneously with an approximation of the system dynamics (which depends on Markov chain parameters). This allows us to guarantee system stability and to derive stability conditions in explicit form. The Markov chain will play in our constructions the role of a “training/learning” rule for the system dynamics considered in the case where the perturbed system’s velocity is replaced by its approximation (function v_1) in the macroscopic (or decision maker’s) frame of reference.

More precisely, we will approximate the pair of functions $(h(\tau), \mu(t, x))$ (which describes the process of GDS evolution and possesses the Markovian property) by a pair of discrete functions:

$$(h(\tau), \mu(t, x)) \rightarrow (\xi_n^{\tau h}, \mu_n^{\tau h}), \tag{5.1}$$

where $\xi_n^{\tau h}$ is an associated (with the microscopic frame of reference) Markov chain state.

First, we consider model (4.1) describing the evolution of dynamic systems in discrete space-time of events. Let an elementary space-time cell be c_{ij} and the complete spatio-temporal region, where the systems dynamics is studied, be \overline{G} such that two events:

$$e_j, e_{j+1} \in c_{ij} \equiv [x_i, x_{i+1}] \otimes [t^j, t^{j+1}] \subset \overline{G} \tag{5.2}$$

of system evolution are governed by the process (x_t, μ_t) . Then we specify these events by two pairs of discrete functions:

$$e_j = (\xi_j^{\tau h}, \mu(x_i, t^j)), e_{j+1} = (\xi_{j+1}^{\tau h}, \mu(x_{i+1}, t^{j+1})), \tag{5.3}$$

where $\xi_j^{\tau h} = x_i^j$ and $\xi_{j+1}^{\tau h} = x_{i+1}^{j+1}$ are states of the associated Markov chain.

Next, we consider the perturbed generalized dynamic system whose training/learning dynamics can be described by the following equation, obtained from (4.3) or (4.4) under appropriate assumptions:

$$\frac{\partial \mu}{\partial t} + v_1(t, x, \mu) \frac{\partial \mu}{\partial x} = \tilde{f}_0(t, x, \mu) \tag{5.4}$$

with the approximation of the initial condition in the DM-time scale:

$$\mu(x, t)|_{t=t_0} = \delta(\epsilon), \tag{5.5}$$

where $\delta(\epsilon)$ is a set of initial conditions with ϵ dependent on the approximation of the function v_0 in the coupled system (4.2) (or (4.4)).

To preserve basic macroscopic features of the system, the values of jumps $\Delta \xi_j^{\tau h} = \xi_{j+1}^{\tau h} - \xi_j^{\tau h}$ of this chain should be subordinated to the corresponding approximation of system-environment boundaries. This can be done via establishing the connection between the GDS and perturbed GDS.

Definition 5.1. Let (5.2) and (5.3) be two subsequent macroscopic events of GDS evolution that are taking place with probability 1. Then the GDS velocity function between the macroscopic events e_j and e_{j+1} can be defined in an elementary space-time cell $c_{ij} \subset \overline{G}$ as:

$$v(t, x) = \lim_{\tau \rightarrow 0} \frac{E^{\tau h | (x_i, \mu^i)} \Delta \xi_j^{\tau h}}{\tau}, \tag{5.6}$$

where the numerator under the limit in (5.6) is the velocity of the Markov Chain (v_{MC}) between two subsequent macroscopic events.

We observe the following remark.

Remark 5.1. If $\lim_{\epsilon \rightarrow 0^+} v_\epsilon = v_1$ then:

$$\lim_{\epsilon \rightarrow 0^+} \epsilon \rightarrow 0^+, n \rightarrow \infty v(t, x) = v_1, \tag{5.7}$$

and together with (4.1) this defines an Infinite Length Unperturbed Markov Chain (ILUMC).

Since perturbations cannot be ignored, Eq. (5.7) can be fulfilled only approximately, which leads to the consideration of an ILPMC, the concept already mentioned in Section 4.

Before formulating consistency conditions for the Markov chain associated with the dynamics of training rule (5.4) and (5.5) we recall [24] with the following definition.

Definition 5.2. The state space in the initial moment of observation in the macroscopic frame of reference is defined as:

$$\Xi(i; 0) = \{x_i, i = 0, 1, 2, \dots, N; N = 2n, n = \lceil (T - t_0)/h \rceil\}, \tag{5.8}$$

and subsequently the cone of macroscopic events of system evolution is defined by a set of macroscopic events as a mapping from the minimal resolution set for the identification of all macroscopic events relevant to the system evolution in the limit $n \rightarrow \infty$ to $\Xi(i; j)$, where

$$\Xi(i; j) = \{(x_i, t_j), i = \overline{k, 2n - k}, j = k, k = \overline{0, n}\}. \tag{5.9}$$

Now we are in a position to formulate.

Proposition 2. *The consistency conditions (local and global, respectively) of the Markov chain $\xi_n^{\tau h}$, $n < \infty$ with the Markov process $(h(\tau), \mu(t, x))$, defined by the mathematical model of GDS evolution (5.4) and (5.5), are:*

$$E_n^{\tau h(x_i, \mu^j)} \Delta \xi_j^{\tau h} = v_1(x_i, t^j, \mu^j) \tau + o(h + \tau), \tag{5.10}$$

$$\text{cov}_n^{\tau h(x_i, \mu^j)} \Delta \xi_j^{\tau h} = o(h + \tau). \tag{5.11}$$

Using these ideas and following [24], a simple approximation to (5.4) and (5.5) can be constructed and in the next sections we generalize this approximation to a large class of schemes for approximating the dynamics of training rules in neural network applications. Let us first highlight the main steps for the construction of the scheme proposed in [24]:

- In the cone of macroscopic events we introduced a floating grid:

$$\omega_{\tau h}^\Delta = \{(x_i, t_j^{\tau j-1}), i = \overline{k, 2n - k}, j = k, k = \overline{0, n}\}, \tag{5.12}$$

where $t_j^{\tau j-1} = t^{j-1} + \tau_{j-1}$ when $j > 1$, $t_j^{\tau j-1} = t^0 + \tau$ when $j = 1$, and $t_j^{\tau j-1} = t^0$ when $j = 0$, where we constructed the following discrete scheme (based on upwind approximations with flux limiters [30]):

$$\begin{aligned} d_i^{j+1} = & d_i^j \left\{ 1 - \frac{\tau}{h} [|v| + v^- \gamma_4 - v^+ \gamma_1] \right\} + \frac{\tau}{h} d_{i-1}^j \{ [v^+ (1 + \gamma_2) + v^- \gamma_4] \} + \frac{\tau}{h} d_{i+1}^j \{ [v^- (1 - \gamma_3) - v^+ \gamma_1] \} \\ & + \frac{\tau}{h} d_{i-2}^j \{ [-v^+ \gamma_2] \} + \frac{\tau}{h} d_{i+2}^j \{ v^- \gamma_3 \}, \end{aligned} \tag{5.13}$$

where d is a discrete function approximating function μ on grid (5.12).

- It is easy to see that the sum of all coefficients near the unknown function on the right-hand side (5.13) gives the unity. This fact allows us associate these coefficients with transition probabilities of a Markov chain, provided those coefficients are nonnegative:

$$\begin{cases} 1 - \frac{\tau}{h} (|v| + v^- \gamma_4 - v^+ \gamma_1) \geq 0, & \gamma_2 \leq 0, & \gamma_3 \geq 0, \\ v^+ (1 + \gamma_2) + v^- \gamma_4 \geq 0, & v^- (\gamma_3 - 1) + v^+ \gamma_1 \leq 0. \end{cases} \tag{5.14}$$

- The consistency conditions (local and global, respectively) of the Markov chain (defined by time-transitions of the discrete scheme (5.13)) with the process $(h(\tau), \mu(t, x))$ (defined by the model (5.4) and (5.5)) couple flux limiters of the scheme:

$$\tau [v^- (1 - \gamma_4 + \gamma_3) - v^+ (1 + \gamma_1 - \gamma_2) - v] = o(\tau + h), \tag{5.15}$$

$$\tau \{ h [v^+ (1 - \gamma_1 - 3\gamma_2) + v^- (1 + \gamma_4 + 3\gamma_3)] - \tau v_{MC}^2 \} = o(\tau + h), \tag{5.16}$$

where the Markov chain velocity in this case is determined as $v_{MC} = v^- (1 - \gamma_4 + \gamma_3) - v^+ (1 + \gamma_1 - \gamma_2)$, and o is the Landau symbol. Using the idea of probabilistic characteristics the term $o(\tau + h)$ in (5.15) and (5.16) can be eliminated. The nonnegativeness of covariance leads to an additional stability condition:

$$\frac{\tau}{h} \leq \frac{v^+ (1 - \gamma_1 - 3\gamma_2) + v^- (1 + \gamma_4 + 3\gamma_3)}{[v^- (1 + \gamma_3 - \gamma_4) - v^+ (1 + \gamma_1 - \gamma_2)]^2}, \tag{5.17}$$

which can be satisfied under appropriate choice of the flux limiters. As a result, we have arrived to the stable Markov chain approximation of the training process.

- More precisely, if the interpolation interval τ is such that conditions (5.14) and (5.17) are satisfied and the transition probabilities of the Markov chain $(\xi_n^{\tau h}, n < \infty)$ are taken in the form:

$$p^{\tau h} [x_k^j, x_i^{j+1} | d(x_k^j, t^j)] = \begin{cases} 1 - \frac{\tau}{h} [v | + v^- \gamma_4 - v^+ \gamma_1], & k = i, \\ \frac{\tau}{h} [v^+ (1 + \gamma_2) + v^- \gamma_4], & k = i - 1, \\ \frac{\tau}{h} [v^- (1 - \gamma_3) - v^+ \gamma_1], & k = i + 1, \\ -\frac{\tau}{h} (v^+ \gamma_2), & k = i - 2, \\ \frac{\tau}{h} (v^- \gamma_3), & k = i + 2, \\ 0, & \text{otherwise,} \end{cases} \tag{5.18}$$

$\forall j = \overline{0, n - 1}$ and $i = \overline{j, N - j}$ ($\gamma_2 = 0$ for $i = j$ and $\gamma_3 = 0$ for $i = N - j$), then the Markov chain approximation of the process $(h(\tau), \mu(x, t))$ is stable, and

$$d(x_i^{j+1}, t^{j+1}) = \sum_k p^{\tau h} [x_k^j, x_i^{j+1} | d(x_k^j, t^j)] d(x_k^j, t^j). \tag{5.19}$$

A numerical procedure like (5.19) is an explicit (evolution forward) stabilization procedure where the DM-function is a stabilizing factor subject to the velocity of the system. Moreover, when $n \rightarrow \infty$ the velocity of the Markov chain converges to the velocity of the process in the sense of the Markov theorem (e.g. [24]). This puts the approach used here on a rigorous theoretical basis. In the next section we show that this approach can be generalized to an important class of conservation-law-based models for the training process of neural networks.

6. Network-based probabilistic approach to conservation law approximations

The network-based methodology developed in the previous sections can be applied to conservation law approximations. In what follows, we demonstrate this on a number of practical examples where we show that many efficient finite difference schemes for the approximation of conservation law equations follow directly from our technique.

Consider a model analogous to that of training/learning dynamics for perturbed GDS, but with $\mu \rightarrow u$, $v_1 \rightarrow a$, and $\tilde{f}_0 \equiv 0$:

$$\frac{\partial u}{\partial t} + a(t, x, u) \frac{\partial u}{\partial x} = 0, \tag{6.1}$$

where u is some physical quantity of interest, and a is the velocity related to the transport of that quantity. With the Cauchy initial condition:

$$u(x, 0) = u_0(x), \quad x \in \mathbb{R}, \tag{6.2}$$

this equation gives a model of a physical system without dissipation in the form easily amenable to a nonlinear (quasi-linear) conservation law. Indeed if $a(u) = F'(u)$, then model (6.1) and (6.2) is equivalent to

$$\frac{\partial u}{\partial t} + \frac{\partial F}{\partial x} = 0, \quad u(x, 0) = u_0(x), \quad x \in \mathbb{R}. \tag{6.3}$$

This model and model (5.4) and (5.5) describing the dynamics of network training belong to the same class of mathematical models, and such models provide building blocks for many mathematics applications in science and engineering. One of the sources of difficulties for numerical solutions of such problems is that in many practical situations the solution to (6.3) might not be continuous, and one has to deal with possible discontinuities across certain curve $x = \sigma(t)$. In other words, if:

$$\lim_{x \rightarrow \sigma(t)^-} u(x, t) = u_{\text{left}}, \quad \lim_{x \rightarrow \sigma(t)^+} u(x, t) = u_{\text{right}} \quad \text{and} \quad u_{\text{left}} \neq u_{\text{right}}, \tag{6.4}$$

one should be able to select a physical meaningful solution from a set of all possible solutions (since two or more classical characteristics in this case might pass the same point). In particular, considering the Cauchy problem (6.1) and (6.2), we note that its solution should be understood in the generalized sense, where both the Rankine–Hugoniot condition [17]:

$$[u] \frac{d\sigma}{dt} = [f], \tag{6.5}$$

where $[\cdot]$ denotes a jump across $x = \sigma(t)$ (i.e. $[u] = u_{\text{right}} - u_{\text{left}}$), and the entropy condition:

$$a(u_{\text{right}}) < \frac{d\sigma}{dt} < a(u_{\text{left}}), \tag{6.6}$$

meaning that the entropy increases as material crosses discontinuity, should be satisfied. Under these circumstances the development of constructive procedures for efficient numerical schemes for the solution of (6.1) and (6.2) becomes a very important task in theory and practice of conservation laws, as well as in a number of related fields. The stability of such schemes is at the heart of the success in achieving this task. Although viscosity solutions could provide some insight into these problems, they would not provide details on the solution stability. Hyperbolic partial differential equation (PDE) models are models for information propagations while probabilistic approaches imply diffusion. However, since such models need to be solved numerically, we will demonstrate that the stability conditions, derived in our approach in a straightforward manner, coincide with the stability conditions typical for numerical approximations of hyperbolic PDEs.

6.1. Neural networks in approximating conservation law models

There is an intrinsic analogy in the information flow pattern for neural networks used in applications and information flows represented by stencils of numerical approximations of conservation laws (e.g., Fig. 1). Consider first some basic classical schemes for numerical approximations of (6.1). We start from the forward centered Euler scheme:

$$u_i^{j+1} = u_i^j - \frac{\lambda}{2} (u_{i+1}^j - u_{i-1}^j). \tag{6.7}$$

The flow of information for these schemes is given in Fig. 1a and it can be seen that in the terminology of neural networks the stencil for this approximation (as well as for an improved approximation $u_i^j \rightarrow \frac{1}{2}(u_{i+1}^j + u_{i-1}^j)$) does not contain any hidden layers. The situation becomes more involved for the Leap-Frog scheme (see Fig. 1b), where several layers of network architecture are coupled:

$$u_i^{j+1} = u_i^{j-1} - \frac{\tau}{h} (F_{i+1}^j - F_{i-1}^j). \tag{6.8}$$

This coupling is sequential in nature, and hence our consideration is pertinent to such networks architectures where the layers are *sequentially* linked. From a numerical point of view such architectures can be represented by the predictor–corrector-type schemes. One of the most effective practical tools in numerical approximation of conservation laws that reflects this architecture is the Lax–Wendroff scheme:

$$u_i^{j+1} = u_i^j - \lambda (F_{i+1/2}^{j+1/2} - F_{i-1/2}^{j+1/2}), \quad \text{where } F_{i+1/2}^{j+1/2} = \frac{1}{2} a u_{i+1}^j + \frac{\lambda}{2} a^2 (u_{i+1}^j - u_i^j). \tag{6.9}$$

As it is obvious from Fig. 1c, this scheme includes “hidden” layer $j + 1/2$ where one has to determine fluxes before moving to the next network layer. In what follows we will show that this scheme and many other effective schemes for numerical approximations of conservation laws can be derived from the general probabilistic

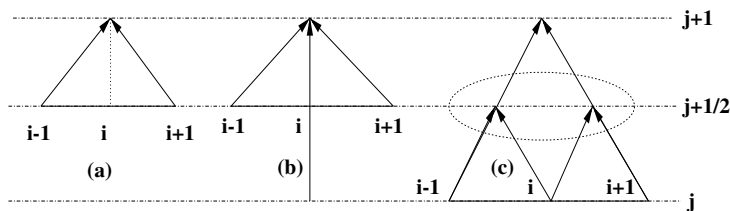


Fig. 1. Network architectures and information flow chart for numerical approximations of conservation laws: (a) forward centered Euler, (b) Leap-Frog and (c) Lax-Wendroff.

approach based on the association of the conservation law models with Markov chain training processes for some neural network architectures in a way similar to that described in Section 5. Moreover, the interpretation/association of the stencils of these schemes with neural network architectures will allow us to deal with the stability conditions in an efficient manner by using the consistency conditions of Markov chains with the original process.

In the spirit of Section 5, the architecture of neural networks is associated with the cone of macroscopic events. From the numerical point of view this is equivalent to the definition of the following points:

$$\begin{aligned}
 [(x_i, t_0^{\varepsilon_{i-1}}), i = 0, \dots, 2n] &\rightarrow [(x_i, t_1^{\varepsilon_0}), i = 1, \dots, 2n - 1] \rightarrow [(x_i, t_2^{\varepsilon_1}), i = 2, \dots, 2n - 2] \rightarrow \dots \\
 &\rightarrow [(x_i, t_{n-1}^{\varepsilon_{n-2}}), i = n - 1, n, n + 1] \rightarrow [(x_i, t_n^{\varepsilon_{n-1}}), i = n]
 \end{aligned}
 \tag{6.10}$$

in the cone of macroscopic events, where the approximation of the evolution of dynamic systems described by the conservation law takes place. An important point to emphasize is that in this general framework the stability conditions for the associated schemes can be obtained in the explicit form. Consider, for example, scheme (5.18) and (5.19) for the solution of homogeneous problem (5.4) and (5.5). First note that if the following condition:

$$v^+(1 - \gamma_1 - 3\gamma_2) + v^-(1 + \gamma_4 + 3\gamma_3) = [v^-(1 + \gamma_3 - \gamma_4) - v^+(1 + \gamma_1 - \gamma_2)]^2
 \tag{6.11}$$

is met, then the Courant–Friederichs–Lewy-type (CFL-type) stability condition (5.17) is satisfied [24]. Since either v^- or v^+ is zero while the other is non-zero, we consider two cases. Our aim in this example is to derive explicit expressions for the flux limiters. In its turn, as follows from our discussion in Section 5, this will define stability conditions of the associated scheme and will complete the definition of the transition probabilities of the associated Markov chain:

1. If $v^- = 0$ and $v^+ \neq 0$ then from (6.11) we have:

$$1 - \gamma_1 - 3\gamma_2 = (v^+)^2(1 + \gamma_1 - \gamma_2)^2.
 \tag{6.12}$$

This leads to a quadratic equation with respect to γ_1 :

$$v^+\gamma_1^2 + [2v^+(1 - \gamma_2) + 1]\gamma_1 + [v^+(1 - \gamma_2)^2 + 3\gamma_2 - 1] = 0.
 \tag{6.13}$$

In a particular case where $\gamma_2 = 0$ the solution to (6.13) can be easily found:

$$\gamma_1^{1,2} = \frac{-(2v^+ + 1) \pm \sqrt{8v^+ + 1}}{2v^+}.
 \tag{6.14}$$

Since γ_1 should be non-positive (see (5.14)) we take the sign “minus” in (6.14), leading to:

$$\gamma_1 = -1 - \frac{\sqrt{8v^+ + 1} + 1}{2v^+}.
 \tag{6.15}$$

A similar reasoning leads to the expression for γ_1 in the general case of (6.13):

$$\gamma_1^{1,2} = \frac{-[2v^+(1 - \gamma_2) + 1] \pm \sqrt{8v^+ + 1 - 16v^+\gamma_2}}{2v^+},
 \tag{6.16}$$

and finally to

$$\gamma_1 = -1 + \gamma_2 - \frac{\sqrt{8v^+ + 1 - 16v^+\gamma_2} + 1}{2v^+}.
 \tag{6.17}$$

2. Along the same vein we can obtain the expression for γ_4 in the case where $v^+ = 0$ and $v^- \neq 0$ as a solution to the equation:

$$1 + \gamma_4 + 3\gamma_3 = v^-(1 + \gamma_3 - \gamma_4)^2,
 \tag{6.18}$$

which can be written in a form easily amenable to its solution with respect to γ_4 :

$$v^-\gamma_4^2 - [2v^-(1 + \gamma_3) + 1]\gamma_4 + [v^-(1 + \gamma_3)^2 - 1 - 3\gamma_3] = 0.
 \tag{6.19}$$

The solution to the last equation is determined as:

$$\gamma_4 = 1 + \gamma_3 + \frac{\sqrt{8v^- + 1 + 16v^-\gamma_3 + 1}}{2v^-} \tag{6.20}$$

(the sign “plus” was taken due to the nonnegativeness of γ_4 , see (5.14)). In the case $\gamma_3 = 0$ (6.20) is reduced to the expression obtained in [24].

This probabilistic approach, based on the association between the conservation law approximations and the neural network training process, is amenable to the treatment of dissipative systems by using the framework of perturbed generalized dynamic systems developed in [24] (see also references therein). However, to make our basic concepts transparent to the reader, we concentrate in this paper on dynamic systems without dissipation. Since the representation (4.1) of dynamic system (and its consequences) from a numerical analysis point of view is an explicit scheme, we recall [30,19] that a family of explicit schemes for the solution of (6.1) and (6.2) can be written as:

$$u_j^{n+1} = u_j^n - \lambda \left(h_{j+1/2}^n - h_{j-1/2}^n \right), \tag{6.21}$$

where $h_{j+1/2}^n = h(u_j^n, u_{j+1}^n)$ is the numerical flux that can model hidden layers in the network architecture associated with the given scheme.

We propose the following general approximation of the flux at $x_j + 0.5\tau$ and $x_j - 0.5\tau$, respectively,

$$h_{j+1/2} = \sum_{k=j-1}^{j+2} b_k(\tau, h, v_k) u_k = b_{j-1} u_{j-1} + b_j u_j + b_{j+1} u_{j+1} + \tilde{b}_{j+2} u_{j+2}, \tag{6.22}$$

$$h_{j-1/2} = \sum_{k=j-1}^j b_k(\tau, h, v_k) u_{k-1} = \tilde{b}_{j-1} u_{j-2} + b_j u_{j-1} + b_{j+1} u_j + b_{j+2} u_{j+1}, \tag{6.23}$$

where all coefficients in (6.22) and (6.23) are velocity-dependent approximations of the flux limiters, for example, $b_{j-1} \equiv b_{j-1}(\tau, h, v_{j-1})$, $\tilde{b}_{j-1} \equiv \tilde{b}_{j-1}(\tau, h, v_{j-2})$, etc. Then, the general family of the explicit schemes can be written as follows:

$$\begin{aligned} u_j^{n+1} &= u_j^n - \lambda [-\tilde{b}_{j-1} u_{j-2} + (b_{j-1} - b_j) u_{j-1} + (b_j - b_{j+1}) u_j + (b_{j+1} - b_{j+2}) u_{j+1} + \tilde{b}_{j+2} u_{j+2}] \\ &= u_j^n - \lambda \left[-\tilde{b}_{j-1} u_{j-2} - \sum_{k=j-1}^{j+1} (b_k - b_{k+1}) u_k + \tilde{b}_{j+2} u_{j+2} \right]. \end{aligned} \tag{6.24}$$

Note that this technique is applicable to the case where $b_m \equiv b_m(\tau, h, \psi(v_{m-1}, v_m))$, e.g. $\psi(v_{m-1}, v_m) = (v_{m-1} + v_m)/2$.

6.2. Special cases and examples

Now, we demonstrate our methodology on a number of examples and show that most efficient schemes for numerical approximations of conservation laws become just special cases of our general network-based approach.

Example 1. First of all, it is easy to confirm that by choosing the flux limiters in the forms:

$$b_{j-1} = 0, \quad b_{j-1} - b_j = -\frac{1}{2}a, \quad b_j - b_{j+1} = 0, \quad b_{j+1} - b_{j+2} = \frac{1}{2}a, \quad b_{j+2} = 0, \tag{6.25}$$

we obtain the forward centered Euler scheme for which the numerical flux is defined as:

$$h_{j+1/2} = \frac{1}{2}a(u_{j+1} + u_j). \tag{6.26}$$

Example 2. If we take now:

$$\begin{cases} b_{j-1} = 0, & -\lambda(b_{j-1} - b_j) = \frac{1}{2} + \frac{\lambda}{2}a, & 1 - \lambda(b_j - b_{j+1}) = 0, \\ -\lambda(b_{j+1} - b_{j+2}) = \frac{1}{2} - \frac{\lambda}{2}a, & b_{j+2} = 0, \end{cases} \quad (6.27)$$

we arrive at the following result:

$$b_{j+1} = \frac{1}{2}a - \frac{1}{2\lambda}, \quad b_j = \frac{1}{2}a - \frac{1}{2\lambda}. \quad (6.28)$$

Therefore, we confirm that under the above choice of the flux limiter coefficients, we obtain the Lax–Friedrichs scheme:

$$u_j^{n+1} = \frac{1}{2}(u_{j+1} + u_{j-1}) - \frac{\lambda a}{2}(u_{j+1} - u_{j-1}), \quad (6.29)$$

which can be written in the form (6.21) with:

$$h_{j+1/2} = \frac{1}{2} \left[a(u_{j+1} + u_j) - \frac{1}{\lambda}(u_{j+1} - u_j) \right]. \quad (6.30)$$

Example 3. By choosing the flux limiters in the general architecture (6.21)–(6.23) as:

$$b_{j-1} = b_{j+2} = 0, \quad b_{j-1} - b_j = -\frac{a}{2} - \frac{\lambda}{2}|a|, \quad b_j - b_{j+1} = |a|, \quad (6.31)$$

we also confirm that:

$$b_j = \frac{1}{2}a + \frac{|a|}{2}, \quad b_{j+1} = \frac{1}{2}a - \frac{|a|}{2}. \quad (6.32)$$

It can be seen that this leads to the (upwind) Euler uncentered scheme, representable in the general form (6.21) with:

$$h_{j+1/2} = \frac{1}{2} [a(u_{j+1} + u_j) - |a|(u - u_j)]. \quad (6.33)$$

Example 4. Taking,

$$b_{j-1} = 0, \quad b_j = \frac{1}{2}a + \frac{\lambda a^2}{2}, \quad b_{j+1} = \frac{1}{2}a - \frac{\lambda a^2}{2}, \quad b_{j+2} = 0, \quad (6.34)$$

in the general formulation (6.21)–(6.23) leads to the one of the most popular schemes for the numerical approximations of conservation law (6.3), the Lax–Wendroff scheme, discussed briefly earlier in this section, where

$$h_{j+1/2} = \frac{1}{2} [a(u_{j+1} + u_j) - \lambda a^2(u_{j+1} - u_j)]. \quad (6.35)$$

Example 5. Finally, we note that the family of schemes proposed in [24] and briefly reviewed in Section 5 is also a subset of the general representation given by (6.21)–(6.23). Indeed, if we take:

$$\begin{cases} b_{j-1} = v, & b_{j-1} - b_j = -v^+(1 + \gamma_2) - v^-\gamma_4, & b_j - b_{j+1} = |v| + v^-\gamma_4 - v^+\gamma_1, \\ b_{j+1} - b_{j+2} = -v^-(1 - \gamma_3) + v^+\gamma_1, & b_{j+2} = v^-\gamma_3 \end{cases} \quad (6.36)$$

(as before, γ_j are flux limiters of the scheme determined from the stability and consistency conditions), we can determine flux coefficient b_{j+1} in two different ways. First, “moving from the left” (see Fig. 2) it is easy to obtain the following chain of relationships:

$$b_j = v^+(1 + \gamma_2) + v^-\gamma_4, \quad b_{j+1} = b_j - |v| - v^-\gamma_4 + v^+\gamma_1 = v^+(1 + 2\gamma_2 + \gamma_1) - |v|. \quad (6.37)$$

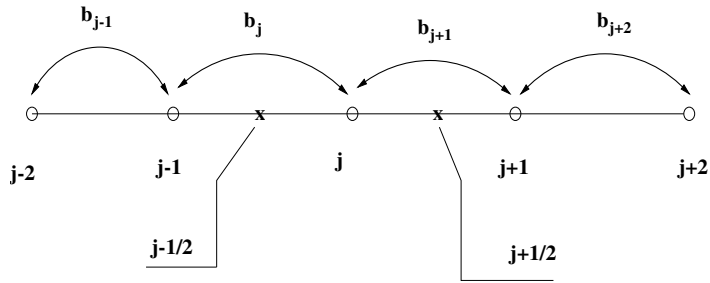


Fig. 2. Flux limiters in the probabilistic approach for conservation laws.

On the other hand, “moving from the right” (see Fig. 2), we get that:

$$b_{j+1} = v^- \gamma_3 - v^- (1 - \gamma_3) + v^+ \gamma_1 = v^- (2\gamma_3 - 1) + v^+ \gamma_1. \tag{6.38}$$

Therefore, by equating two expressions for b_{j+1} we arrive at the conclusion that:

$$v^+ (1 + 2\gamma_2) + v^- (1 - 2\gamma_3) = |v|, \tag{6.39}$$

which can be satisfied by setting $\gamma_2 = \gamma_3 = 0$. Having these flux limiters, γ_1 and γ_4 can be controlled by satisfying the stability conditions [24].

7. Stability and consistency requirements for network-based approximations of conservation laws

In the previous section we derived numerical approximations for conservation laws using the general network-based methodology. As we have discussed in Section 4 this methodology originates from the probabilistic foundations of the dynamics. Indeed, having approximate initial data (e.g., as a result of measurements) we attempt to extrapolate this data further in time. This procedure is a subject of a probabilistic error. Therefore, it is important to discuss further stability and consistency requirements for the schemes obtained earlier in this paper, in particular, for the general family of approximations (6.21)–(6.23).

The probabilistic approach to conservation laws allows us to explain a number of important phenomena related to numerical approximations of these models in a straightforward way. First, it is easy to check that the sum of all coefficients near the unknown function in every scheme discussed in the previous section is 1. However, this fact alone might not be sufficient to interpret those coefficients as transition probabilities. Indeed, the nonnegativity requirement plays a key role in the possibility of such an interpretation. From a numerical point of view this requirement leads to the stability conditions of the scheme. If we take the forward Euler/centered scheme this requirement leads to $\lambda a/2 \leq 0$, which can be satisfied only if $a < 0$ ($\lambda = \tau/h$). However, since we use here a forward scheme with a stencil depicted in Fig. 1, this should not come as a surprise, because for $a > 0$ the explicit scheme on this stencil is known to be unstable.

It is easy to check that coefficients in all schemes considered in Sections 5 and 6 will be nonnegative subject to the Courant–Friedrichs–Lewy-type stability conditions $|\lambda a| < 1$ (as one would expect for the models based on hyperbolic PDEs [30]) and therefore, all these schemes can be interpreted in a probabilistic sense. Naturally, however, that satisfying those stability conditions is subject to the appropriate choice of the flux limiters. For example, for the scheme (5.18) and (5.19) discussed in Section 5, such fluxes (γ_1 and γ_4) should be chosen as to satisfy the following conditions:

$$\frac{\tau}{h} (|v| + v^- \gamma_4 - v^+ \gamma_1) \leq 1, \quad \frac{1}{2} v^+ + v^- \gamma_4 \geq 0, \quad \frac{1}{2} v^- + v^+ \gamma_1 \leq 0. \tag{7.1}$$

Moreover, one has to satisfy the consistency conditions, which in the case of this scheme has the form:

$$\tau [v^- (3/2 - \gamma_4) + v^+ (3/2 + \gamma_1) - v] = o(\tau + h), \tag{7.2}$$

where, as before, o denotes the Landau symbol.

In our general case of the family (6.21)–(6.23) the nonnegativity of coefficients require:

$$\begin{cases} \lambda b_{j-1} \geq 0, & \lambda(b_j - b_{j-1}) \geq 0, & 1 - \lambda(b_{j+1} - b_j) \geq 0, \\ \lambda(b_{j+2} - b_{j+1}) \geq 0, & -\lambda b_{j+2} \geq 0. \end{cases} \tag{7.3}$$

This results in

$$\lambda(b_{j+1} - b_j) \leq 1, \quad \text{where } b_{j+1} \leq b_{j+2} \leq 0 \leq b_{j-1} \leq b_j \tag{7.4}$$

(see Fig. 3 for interpretation). Consistency conditions for the general scheme (6.21)–(6.23) can also be obtained by using the methodology developed in [24]. Indeed, we assume that each jump of the associated Markov chain is allowed within a region (e.g., rectangular) with characteristic lengths τ and h such that $\tau < h$ (if necessary such a region can be easily refined in a way similar, for example, to the finite element methodology). If the Markov chain is in state x we associate this state with the position of the Markov chain, i.e. assume that $\xi_j^{h\tau} = x$. At the next moment of time $\tau \rightarrow \tau + \Delta\tau$ the Markov chain is assumed to be in one of the following states: $x - 2h, x - h, x, x + h$, or $x + 2h$ (see Fig. 4, left). This can be extended to the case of any arbitrary number of states in a straightforward manner. The table of transition probabilities from the old to a new state of the Markov chain associated with our scheme (6.21)–(6.23) can be constructed as follows:

New state (x_i)	Probability of transition (p_i)
$x - h$	$\tau/h(b_j - b_{j-1})$
$x + h$	$\tau/h(b_{j+2} - b_{j+1})$
x	$1 - \tau/h(b_{j+1} - b_j)$
$x - 2h$	τ/hb_{j-1}
$x + 2h$	$-\tau/hb_{j+2}$

We can now easily relate this jump of the Markov chain to the transition to the next time layer in our numerical approximation of the conservation law. Therefore, we are in a position to derive the explicit form of the local consistency condition for the family of schemes (6.21)–(6.23). Since,

$$E_j^{\tau h}(x, d^j) = \sum_{i=1}^s x_i p_i, \tag{7.5}$$

where s is the number of states allowed in the next time layer (in our case $s = 5$), after some simplifications this condition can be given in the form:

$$\tau \sum_{k=j-1}^{j+2} b_k = o(\tau + h). \tag{7.6}$$

The sign is chosen to be positive due to the direction of computations (see Fig. 4, right).

For example, consider the family of schemes discussed in Section 5 (a subclass of the general family (6.21)–(6.23)). Note that in this case:

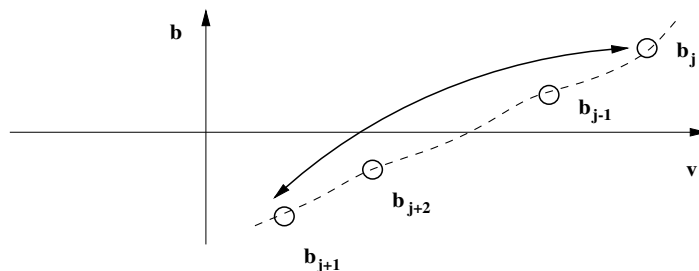


Fig. 3. A schematic representation of flux limiters as functions of velocity.

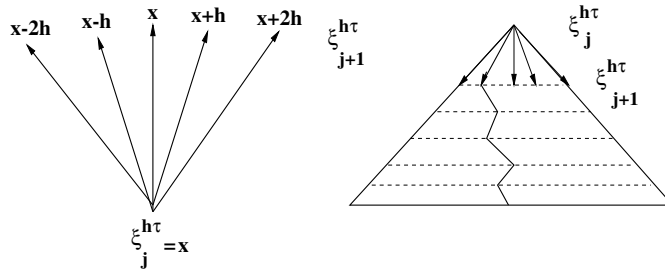


Fig. 4. Markov chain as a tool to combine feedforward and feedback neural networks.

$$\sum_{k=j-1}^{j+2} b_k = v^+ [2(1 + 2\gamma_2) + \gamma_1 + \gamma_2] + v^- (\gamma_3 + \gamma_4) - |v|. \tag{7.7}$$

Take, for example, the limiters as $\gamma_1 = \gamma_2 = -1/2$ and $\gamma_3 = \gamma_4 = 1/2$. This reduces the consistency condition to:

$$-v^+ + v^- - |v| = \sum_{k=j-1}^{j+2} b_k. \tag{7.8}$$

Note also that the resulting scheme in this case has the following form:

$$u_i^{j+1} = u_i^j + \frac{\tau}{h} |v| \left[-\frac{3}{2} u_i + \frac{1}{2} (u_{i-1} + u_{i+1}) \right] + \frac{\tau}{2h} [v^+ u_{i-2} + v^- u_{i+2}]. \tag{7.9}$$

It is easy to conclude that if $v > 0$ (and hence $-2v = \sum_{k=j-1}^{j+2} b_k$) scheme (7.9) is reducible to:

$$u_i^{j+1} = u_i^j + \frac{\tau}{2h} v [-3u_i^j + (u_{i+1} + u_{i-1}) + u_{i-2}], \tag{7.10}$$

whereas if $v < 0$ (and hence $\sum_{k=j-1}^{j+2} b_k = 0$), we have:

$$u_i^{j+1} = u_i^j - \frac{\tau}{2h} v [-3u_i + (u_{i+1} + u_{i-1}) + u_{i+2}]. \tag{7.11}$$

The global consistency condition involves the velocity of the associated Markov chain (see, e.g. (5.14) and further details in [24]). Recall that for the family of schemes constructed in [24], the Markov chain velocity (that coincides with the velocity of the dynamic system/process when $n \rightarrow \infty$) is determined in terms of the flux limiters as follows:

$$v_{MC} = v^- - v^+ = -v. \tag{7.12}$$

This result is independent of the sign of v . It shows that the Markov chain evolves in the direction opposite to the evolution of the dynamic system. This allows us to represent an explicit form of the global consistency condition for the family of schemes (6.21)–(6.23) in the form:

$$\tau [3h|v| - \tau v^2] = o(\tau + h). \tag{7.13}$$

In conclusion, we note that, as follows from (7.6), the Lax–Wendroff scheme’s consistency condition takes the form:

$$\tau (b_j + b_{j+1}) = o(\tau + h) \quad \text{or} \quad \tau a = o(\tau + h), \tag{7.14}$$

which leads to the CFL-type stability condition:

$$\frac{\tau a}{c(\tau + h)} \leq 1, \tag{7.15}$$

where $c = \lim_{\tau \rightarrow 0, h \rightarrow 0} \frac{2\tau v}{\tau + h}$ is a Landau constant participating in the definition of probabilistic characteristics of conservation laws [24]. Consistency conditions for other schemes considered in Section 6 can be obtained in a similar manner.

8. Conclusions

In this paper, a general framework for the analysis of a connection between the training of artificial neural networks via the dynamics of Markov chains and the approximation of conservation law models has been proposed. This framework allows us to demonstrate an intrinsic link between microscopic and macroscopic models for evolution via the concept of perturbed generalized dynamic systems. We have showed that mathematical models describing dynamics of network training can be treated effectively by using the concept of perturbed Markov chains associated with the original dynamic system. We have developed a general methodology allowing us to derive computational models for the dynamics of network training and to obtain constructive algorithms, as well as stability conditions for numerical approximations of conservation laws. We have demonstrated how this methodology can be applied to numerical approximations of conservation laws viewed here as Markov chain network training procedures. Our main results have been exemplified with several illustrative examples for which we have explicitly derived stability and consistency conditions.

Acknowledgements

I wish to thank my colleagues from Mathematics Departments at the University of South Australia for fruitful discussions on the topics of this paper. Support of the Mads Clausen Foundation is also gratefully acknowledged.

References

- [1] A. Agogino, K. Stanley, R. Mikkulainen, Online interactive neuro-evolution, *Neural Process. Lett.* 11 (2000) 29–37.
- [2] J. Aguilar, Learning algorithm retrieval process for the multiple classes random neural network model, *Neural Process. Lett.* 13 (2001) 81–91.
- [3] M.S. Ahmed, S.H. Riyaz, Dynamic observers – a neural net approach, *J. Intell. Fuzzy Syst.* 9 (2000) 113–127.
- [4] S.N. Balakrishnan, V. Biega, Adaptive-critic-based neural networks for aircraft optimal control, *J. Guid. Control Dynam.* 19 (1996) 893–899.
- [5] A.R. Barron, Universal approximation bounds for superpositions of a sigmoidal function, *IEEE Trans. Inform. Theory* 39 (1993) 930–945.
- [6] C. Bishop, *Neural Networks for Pattern Recognition*, Clarendon Press, Oxford, 1995.
- [7] J.B.D. Cabrera, K.S. Narendra, Issues in the application of neural networks for tracking based on inverse control, *IEEE Trans. Autom. Contr.* 44 (1999) 2007–2027.
- [8] A. Catala, C. Angulo, A comparison between the Tikhonov and the Bayesian approaches to calculate regularisation matrices, *Neural Process. Lett.* 11 (2000) 185–195.
- [9] G.J. Chaitin, Godel’s theorem and information, *Int. J. Theor. Phys.* 22 (1982) 941–954.
- [10] F. Corinto, M. Biey, M. Gilli, Non-linear coupled CNN models for multiscale image analysis, *Int. J. Circuit Theory Appl.* 34 (1) (2006) 77–88.
- [11] P. Costa, P. Larzabal, Initialization of supervised training for parametric estimation, *Neural Process. Lett.* 9 (1999) 53–61.
- [12] H.T. Croft, K.J. Falconer, R.K. Guy, *Unsolved Problems in Geometry*, Springer-Verlag, New York, 1991, p. 3.
- [13] G. Cybenko, Approximation by superposition of sigmoidal functions, *MCSS* 2 (1989) 303–314.
- [14] B. DasGupta, B. Hammer, On approximate learning by multi-layered feedforward circuits, *Theor. Comput. Sci.* 348 (1) (2005) 95–127.
- [15] M.O. Efe, O. Kaynak, Stabilizing and robustifying the learning mechanisms of ANNs in control engineering applications, *Int. J. Intell. Syst.* 15 (2000) 365–388.
- [16] D. Gawin, M. Lefik, B.A. Schrefler, ANN approach to sorption hysteresis within a coupled hydro-thermo-mechanical FE analysis, *Int. J. Numer. Method Eng.* 50 (2001) 299–323.
- [17] R.B. Guenther, J.W. Lee, *Partial Differential Equations of Mathematical Physics and Integral Equations*, Prentice Hall, NJ, 1988.
- [18] K. Hornik, M. Stinchcombe, H. White, Multi-layer feedforward networks are universal approximators, *Neural Network* 2 (1989) 359–366.
- [19] L. Lapidus, G.F. Pinder, *Numerical Solution of PDEs in Science and Engineering*, John Wiley & Sons, 1999.
- [20] A.U. Levin, K.S. Narendra, Control of nonlinear dynamical systems using neural networks. Part I, *IEEE Trans. Neural Network* 4 (1993) 192–206;
A.U. Levin, K.S. Narendra, Control of nonlinear dynamical systems using neural networks. Part II, *IEEE Trans. Neural Network* 7 (1996) 30–42.
- [21] L.-Z. Liao, A recurrent neural network for N-stage optimal control problems, *Neural Process. Lett.* 10 (1999) 195–200.
- [22] R.V.N. Melnik, A hierarchy of hyperbolic macrodynamic equations as a model for network training, in: *Proceedings of the IEEE International Symposium on Information Theory*, Ulm, Germany, June 29–July 4, 1997, p. 322.

- [23] R.V.N. Melnik, On consistent regularities of control and value functions, *Numer. Funct. Anal. Optim.* 18 (1997) 401–426.
- [24] R.V.N. Melnik, Dynamic system evolution and Markov chain approximation, *Discr. Dynam. Nat. Soc.* 2 (1998) 7–39.
- [25] R.V.N. Melnik, Deterministic and stochastic dynamics with hyperbolic HJB-type equations, *Dynam. Contin. Discr. Impul. Syst. Ser. A: Math. Anal.* 10 (3) (2003) 317–330.
- [26] I. Müller, T. Ruggeri, *Extended Thermodynamics*, Springer-Verlag, 1993.
- [27] K.S. Narendra, Neural networks for control: theory and practice, *Proc. IEEE* 84 (1996) 1385–1406.
- [28] K.S. Narendra, K. Parthasarathy, Identification and control of dynamical systems using neural networks, *IEEE Trans. Neural Network* 1 (1990) 4–27.
- [29] K.S. Narendra, C. Xiang, Adaptive control of discrete-time systems using multiple models, *IEEE Trans. Autom. Contr.* 45 (2000) 1669–1686.
- [30] A. Quarteroni, R. Sacco, F. Saleri, *Numerical Mathematics*, Springer, 2000.
- [31] S. Santini, A. Del Bimbo, Properties of block feedback neural networks, *Neural Network* 8 (4) (1995) 579–596.
- [32] B. Sendhoff, M. Kreutz, A model for the dynamic interaction between evolution and learning, *Neural Process. Lett.* 10 (1999) 181–193.
- [33] C.W. Ulmer II et al., Computational neural network and the rational design of polymeric materials, *Comput. Theor. Polym. Sci.* 8 (1998) 311–321.
- [34] G. Wahba, Generalization and regularization of nonlinear learning systems, *Handbook of Brain Theory and Neural Networks*, second ed., M. Arib, 2000.
- [35] J. Zhang et al., Prediction of polymer quality in batch polymerisation reactions using robust neural networks, *Chem. Eng. J.* 69 (1998) 135–143.
- [36] T. Zhang, S.S. Ge, C.C. Hang, Neural-based direct adaptive control for a class of general nonlinear systems, *Int. J. Syst. Sci.* 28 (1997) 1011–1020.

Further reading

- [1] O. Adetona, E. Garcia, L.H. Keel, A new method for the control of discrete nonlinear dynamic systems using neural networks, *IEEE Trans. Neural Network* 11 (2000) 102–112.
- [2] K.S. Narendra, J. Balakrishnan, Adaptive control using multiple models, *IEEE Trans. Automat. Contr.* 42 (1997) 171–187;
K.S. Narendra, S. Mukhopadhyaya, Adaptive control using neural networks and approximate models, *IEEE Trans. Neural Network* 8 (1997) 475–485.
- [3] G. Puscasu, C. Milos, A new approach to multilayer neural networks, in: *Proceedings of the International Conference on Technical Informatics: CONTI'96*, Timisoara, 1996, pp. 41–48.
- [4] I.W. Sandberg, G.J.J. van Zyl, The spectral coefficients of the response of nonlinear systems to asymptotically almost periodic inputs, *IEEE Trans. Circ. Syst. I-Fundam. Theor. Appl.* 48 (2001) 170–176.
- [5] E. Santos Jr., J.D. Young, Probabilistic temporal networks: a unified framework for reasoning with time and uncertainty, *Int. J. Approx. Reason.* 20 (1999) 263–291.