

Deep learning with stochastic neural networks

It has been recently shown that new insights on deep learning can be obtained by regarding the process of training a deep neural network as a discretization of an optimal control problem involving nonlinear differential equations [5, 4, 8]. One attractive feature of this formulation is that it allows to use tools from dynamical system theory to study deep learning from a rigorous mathematical perspective [12, 9, 14]. For instance, it has been recently shown that by idealizing deep residual networks (ResNet) as continuous-time dynamical systems it is possible to derive sufficient conditions for universal approximation in L^p , which can be understood as an approximation theory built on flow maps generated by dynamical systems [13].

In this note we present a formulation deep neural networks obtained by applying simple probabilistic tools to discrete dynamical systems. Specifically, we consider two types of neural network models:

- Neural networks perturbed by additive random noise;
- Neural networks with random weights and biases.

Modeling neural networks as discrete stochastic dynamical systems. Let us begin by modeling the input-output map of a neural network as a discrete dynamical system (see Figure 1)

$$\mathbf{X}_1 = \mathbf{F}_0(\mathbf{X}_0, \mathbf{w}_0) + \boldsymbol{\xi}_0 \quad \mathbf{X}_{n+1} = \mathbf{F}(\mathbf{X}_n, \mathbf{w}_n) + \boldsymbol{\xi}_n, \quad (1)$$

Here the index n labels a specific layer in the network, $\mathbf{X}_0 \in \mathbb{R}^d$ is the input, $\mathbf{X}_n \in \mathbb{R}^N$ ($n = 1, \dots, L$) represents output of the n -th layer, and $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ is set of statistically independent random vectors, or more generally a vector-valued Markov process. We allow the initial state \mathbf{X}_0 to be random as well, which can be directly connected to a data set in a training algorithm. A neural networks of the form (1) is called *recurrent*, to emphasize the fact that the mapping

$$\mathbf{F}(\mathbf{X}_n, \mathbf{w}_n) = \varphi(\mathbf{W}_n \mathbf{X}_n + \mathbf{b}_n) \quad \mathbf{w}_n = \{\mathbf{W}_n, \mathbf{b}_n\}, \quad (2)$$

between one layer and the next has the same functional form. In (2) $\varphi : \mathbb{R}^N \mapsto \mathbb{R}^N$ is the *activation function* of the network, \mathbf{W}_n is a $N \times N$ weight matrix and $\mathbf{b}_n \in \mathbb{R}^N$ is a bias vector.

In a supervised learning setting, the degrees of freedom

$$\mathbf{w} = \{\mathbf{w}_0, \dots, \mathbf{w}_{L-1}\}, \quad (3)$$

are determined by optimizing a suitable performance metric depending on the network output. For instance, if we are interested in using the network depicted in Figure 1 to approximate a multivariate function $g(\mathbf{x}) \in L^2([0, 1]^d)$ then we can identify the degrees of freedom (3) by minimizing, e.g., the non-convex functional

$$\{\boldsymbol{\alpha}, \mathbf{w}\} = \underset{\boldsymbol{\alpha}, \mathbf{w}}{\operatorname{argmin}} \|g(\mathbf{x}) - \boldsymbol{\alpha} \cdot \mathbb{E}[\mathbf{X}_L | \mathbf{X}_0 = \mathbf{x}]\|_{L^2([0, 1]^d)}^2, \quad (4)$$

where $\boldsymbol{\alpha}$ are the output weights, and $\mathbb{E}[\mathbf{X}_L | \mathbf{X}_0 = \mathbf{x}]$ conditional expectation of \mathbf{X}_L given $\mathbf{X}_0 = \mathbf{x}$. In this setting, it is clear that the process of training a neural network is basically an optimal control problem (the controls being the weights and biases) of a discrete stochastic differential equation.

A different stochastic neural network model can be defined by randomizing weights and biases [6, 21]. In this setting we have

$$\mathbf{X}_{n+1} = \varphi(\mathbf{W}_n(\omega) \mathbf{X}_n + \mathbf{b}_n(\omega)), \quad (5)$$

where $\mathbf{W}_n(\omega)$ are random weight matrices and $\mathbf{b}_n(\omega)$ are random bias vectors. We shall assume that \mathbf{W}_n and \mathbf{b}_n corresponding to different layers are statistically independent.

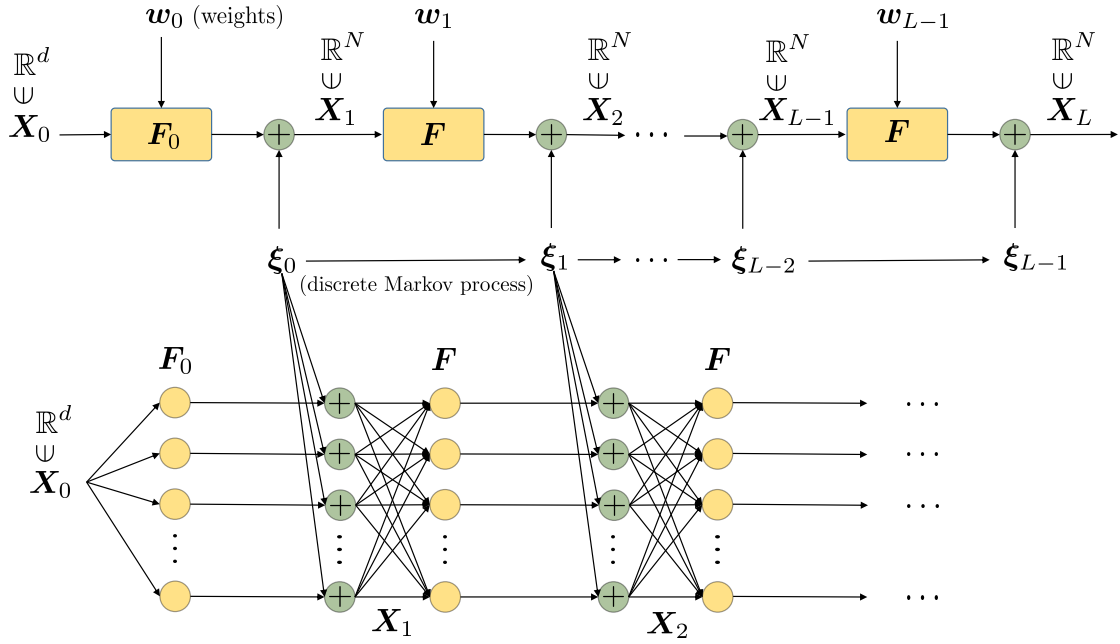


Figure 1: Sketch of the stochastic neural network model (1). Note that the transfer function F is the same in every layer (except the first one). This implies that the random vectors $\{\mathbf{X}_1, \dots, \mathbf{X}_L\}$ all have the same dimension.

By adding random noise to the output of each neural network layer, or by randomizing weights and biases, we are essentially adding an infinite number of degrees of freedom to our system. This allows us to rethink the process of training the neural network from a probabilistic perspective. For instance, random noise allows us to approximately encode secret messages in fully trained deterministic neural networks by selecting an appropriate transition probability for the noise process.

Composition and transfer operators

Let us now derive the composition and transfer operators associated with the neural network models (1) and (5), which map, respectively, the conditional expectation $\mathbb{E}\{\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_n = \mathbf{x}\}$ and $p_n(\mathbf{x})$ (the probability density of \mathbf{X}_n) forward and backward across the network. To this end, we assume that $\{\xi_0, \dots, \xi_{L-1}\}$ in (1) are independent random vectors. Similarly, we assume that the random matrices $\{\mathbf{W}_0, \dots, \mathbf{W}_{L-1}\}$ and the random vectors $\{\mathbf{b}_0, \dots, \mathbf{b}_{L-1}\}$ in (5) are statistically independent. These assumptions imply that the sequences of vectors $\{\mathbf{X}_0, \mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L\}$ generated by either (1) or (5) are *discrete Markov processes*¹. Therefore, the joint probability density function (PDF) of the random vectors $\{\mathbf{X}_0, \dots, \mathbf{X}_L\}$, i.e., joint PDF of the state of the entire neural network, can be factored² as

$$p(\mathbf{x}_0, \dots, \mathbf{x}_L) = p_{L|L-1}(\mathbf{x}_L|\mathbf{x}_{L-1})p_{L-1|L-2}(\mathbf{x}_{L-1}|\mathbf{x}_{L-2}) \cdots p_{1|0}(\mathbf{x}_1|\mathbf{x}_0)p_0(\mathbf{x}_0). \quad (6)$$

By using the identity

$$p(\mathbf{x}_{k+1}, \mathbf{x}_k) = p_{k+1|k}(\mathbf{x}_{k+1}|\mathbf{x}_k)p_k(\mathbf{x}_k) = p_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1})p_{k+1}(\mathbf{x}_{k+1}) \quad (7)$$

¹The independence assumption of the random noise vector $\{\xi_0, \dots, \xi_{L-1}\}$ or the random weights and biases $\{\mathbf{W}_1, \dots, \mathbf{W}_1\}$ and $\{\mathbf{b}_1, \dots, \mathbf{b}_m\}$ is not necessary for the process $\{\mathbf{X}_0, \mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L\}$ to be Markov. Such assumption just simplifies the expression of the transition density $p(\mathbf{x}_{i+1}|\mathbf{x}_i)$.

²In equation (6) we used the shorthand notation $p_{i|j}(\mathbf{x}|\mathbf{y})$ to denote the conditional probability density function of the random vector \mathbf{X}_i given $\mathbf{X}_j = \mathbf{y}$. With this notation we have that the conditional probability density of \mathbf{X}_i given $\mathbf{X}_i = \mathbf{y}$ is $p_{i|i}(\mathbf{x}|\mathbf{y}) = \delta(\mathbf{x} - \mathbf{y})$, where $\delta(\cdot)$ is the Dirac delta function.

we see that the chain of transition probabilities (6) can be reverted, yielding

$$p(\mathbf{x}_0, \dots, \mathbf{x}_L) = p_{0|1}(\mathbf{x}_0|\mathbf{x}_1)p_{1|2}(\mathbf{x}_1|\mathbf{x}_2) \cdots p_{L-1|L}(\mathbf{x}_{L-1}|\mathbf{x}_L)p_L(\mathbf{x}_L). \quad (8)$$

From these expression, it follows that

$$p_{n|q}(\mathbf{x}|\mathbf{y}) = \int p_{n|j}(\mathbf{x}|\mathbf{z})p_{j|q}(\mathbf{z}|\mathbf{y})d\mathbf{z}, \quad (9)$$

for all indices n, j and q in $\{0, \dots, L\}$, excluding $n = j = q$ (see footnote 2). The transition probability equation (9) is known as *discrete Chapman-Kolmogorov equation* and it allows us to define the transfer operator mapping the PDF $p_n(\mathbf{x}_n)$ into $p_{n+1}(\mathbf{x}_{n+1})$, together with the composition operator for the conditional expectation $\mathbb{E}\{\mathbf{u}(\mathbf{x}_L)|\mathbf{X}_n = \mathbf{x}_n\}$. As we shall see hereafter, the discrete composition and transfer operators are adjoint to one another.

Transfer operator. Let us denote by $p_q(\mathbf{x})$ the PDF of \mathbf{X}_q , i.e., the output of the q -th neural network layer. We first define the operator that maps $p_q(\mathbf{x})$ into $p_n(\mathbf{x})$. By integrating the joint probability density of \mathbf{X}_n and \mathbf{X}_q , i.e., $p_{n|q}(\mathbf{x}|\mathbf{y})p_q(\mathbf{y})$ with respect to \mathbf{y} we immediately obtain

$$p_n(\mathbf{x}) = \int p_{n|q}(\mathbf{x}|\mathbf{y})p_q(\mathbf{y})d\mathbf{y}. \quad (10)$$

At this point, it is convenient to define the operator

$$\mathcal{N}(n, q)f(\mathbf{x}) = \int p_{n|q}(\mathbf{x}|\mathbf{y})f(\mathbf{y})d\mathbf{y}. \quad (11)$$

$\mathcal{N}(n, q)$ is known as *transfer operator* [3]. From a mathematical viewpoint $\mathcal{N}(n, q)$ is an integral operator with kernel $p_{n|q}(\mathbf{x}, \mathbf{y})$, i.e., the transition density integrated “from the right”. It follows from the Chapman-Kolmogorov identity (9) that the set of integral operators $\{\mathcal{N}(n, q)\}$ forms a group. Namely,

$$\mathcal{N}(n, q) = \mathcal{N}(n, j)\mathcal{N}(j, q), \quad \mathcal{N}(j, j) = \mathcal{I}, \quad \forall n, j, q \in \{0, \dots, L\}. \quad (12)$$

The operator \mathcal{N} allows us to map the one-layer PDF, e.g., the PDF of \mathbf{X}_q , either forward or backward across the neural network (see Figure 2). As an example, consider a network with four layers with states \mathbf{X}_0 (input), \mathbf{X}_1 , \mathbf{X}_2 , \mathbf{X}_3 , and \mathbf{X}_4 (output). Then Eq. (11) implies that,

$$p_2(\mathbf{x}) = \underbrace{\mathcal{N}(2, 1)\mathcal{N}(1, 0)}_{\mathcal{N}(2, 0)}p_0(\mathbf{x}) = \underbrace{\mathcal{N}(2, 3)\mathcal{N}(3, 4)}_{\mathcal{N}(2, 4)}p_4(\mathbf{x}).$$

In summary, we have

$$p_n(\mathbf{x}) = \mathcal{N}(n, q)p_q(\mathbf{x}) \quad \forall n, q \in \{0, \dots, L\}, \quad (13)$$

where

$$\mathcal{N}(n, q)p_q(\mathbf{x}) = \int p_{n|q}(\mathbf{x}|\mathbf{y})p_q(\mathbf{y})d\mathbf{y}. \quad (14)$$

We emphasize that modeling PDF dynamics via neural networks has been studied extensively in machine learning, e.g., in the theory of normalizing flows for density estimation or variational inference [17, 10, 18].

Composition operator For any measurable deterministic function $\mathbf{u}(\mathbf{x})$, the conditional expectation of $\mathbf{u}(\mathbf{X}_j)$ given $\mathbf{X}_n = \mathbf{x}$ is defined as

$$\mathbb{E}\{\mathbf{u}(\mathbf{X}_j)|\mathbf{X}_n = \mathbf{x}\} = \int \mathbf{u}(\mathbf{y})p_{j|n}(\mathbf{y}|\mathbf{x})d\mathbf{y}. \quad (15)$$

A substitution of (9) into (15) yields

$$\mathbb{E}\{\mathbf{u}(\mathbf{X}_j)|\mathbf{X}_n = \mathbf{x}\} = \int \mathbb{E}\{\mathbf{u}(\mathbf{X}_j)|\mathbf{X}_q = \mathbf{y}\} p_{q|n}(\mathbf{y}|\mathbf{x}) d\mathbf{y}, \quad (16)$$

which holds for all $j, n, q \in \{0, \dots, L-1\}$. At this point we define the integral operator

$$\mathcal{M}(n, q)f(\mathbf{x}) = \int f(\mathbf{y})p_{q|n}(\mathbf{y}|\mathbf{x})d\mathbf{y}, \quad (17)$$

which is known as *composition* [3] or “stochastic Koopman” [19, 23] operator. Thanks to the Chapman-Kolmogorov identity (9), the set of operators $\{\mathcal{M}(q, j)\}$ forms a group, i.e.,

$$\mathcal{M}(n, q) = \mathcal{M}(n, j)\mathcal{M}(j, q), \quad \mathcal{M}(j, j) = \mathcal{I}, \quad \forall n, j, q \in \{0, \dots, L\}. \quad (18)$$

Equation (18) allows us to map the conditional expectation (15) of any measurable phase space function $\mathbf{u}(\mathbf{X}_j)$ forward or backward through the network. As an example, consider again a neural network with four layers and states $\{\mathbf{X}_0, \dots, \mathbf{X}_4\}$. We have

$$\begin{aligned} \mathbb{E}\{\mathbf{u}(\mathbf{X}_j)|\mathbf{X}_2 = \mathbf{x}\} &= \mathcal{M}(2, 3)\mathcal{M}(3, 4)\mathbb{E}\{\mathbf{u}(\mathbf{X}_j)|\mathbf{X}_4 = \mathbf{x}\} \\ &= \mathcal{M}(2, 1)\mathcal{M}(1, 0)\mathbb{E}\{\mathbf{u}(\mathbf{X}_j)|\mathbf{X}_0 = \mathbf{x}\}. \end{aligned} \quad (19)$$

Equation (19) holds for every $j \in \{0, \dots, 4\}$. Of particular interest in machine-learning context is the conditional expectation of $\mathbf{u}(\mathbf{X}_L)$ (network output) given $\mathbf{X}_0 = \mathbf{x}$, which can be computed as

$$\begin{aligned} \mathbb{E}\{\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_0 = \mathbf{x}\} &= \mathcal{M}(0, L)\mathbf{u}(\mathbf{x}), \\ &= \mathcal{M}(0, 1)\mathcal{M}(1, 2) \cdots \mathcal{M}(L-1, L)\mathbf{u}(\mathbf{x}), \end{aligned} \quad (20)$$

i.e., by propagating $\mathbf{u}(\mathbf{x}) = \mathbb{E}\{\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_L = \mathbf{x}\}$ *backward* through the neural network using single layer operators $\mathcal{M}(i-1, i)$. Similarly, we can compute, e.g., $\mathbb{E}\{\mathbf{u}(\mathbf{X}_0)|\mathbf{X}_L = \mathbf{x}\}$ as

$$\mathbb{E}\{\mathbf{u}(\mathbf{X}_0)|\mathbf{X}_L = \mathbf{x}\} = \mathcal{M}(L, 0)\mathbf{u}(\mathbf{x}). \quad (21)$$

For subsequent analysis, it is convenient to define

$$\mathbf{q}_n(\mathbf{x}) = \mathbb{E}\{\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_{L-n} = \mathbf{x}\}. \quad (22)$$

In this way, if $\mathbb{E}\{\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_n = \mathbf{x}\}$ is propagated *backward* through the network by $\mathcal{M}(n-1, n)$, then $\mathbf{q}_n(\mathbf{x})$ is propagated *forward* by the operator

$$\mathcal{G}(n, q) = \mathcal{M}(L-n, L-q). \quad (23)$$

In fact, equations (22)-(23) allow us to write (20) in the equivalent form

$$\begin{aligned} \mathbf{q}_L(\mathbf{x}) &= \mathcal{G}(L, L-1)\mathbf{q}_{L-1}(\mathbf{x}) \\ &= \mathcal{G}(L, L-1) \cdots \mathcal{G}(1, 0)\mathbf{q}_0(\mathbf{x}), \end{aligned} \quad (24)$$

i.e., as a forward propagation problem (see Figure 2). Note that we can write (24) (or (20)) explicitly in terms of iterated integrals involving single layer transition densities as

$$\begin{aligned} \mathbf{q}_L(\mathbf{x}) &= \int \mathbf{u}(\mathbf{y})p_{0|L}(\mathbf{y}|\mathbf{x})d\mathbf{y} \\ &= \int \mathbf{u}(\mathbf{y}) \left(\int \cdots \int p_{L|L-1}(\mathbf{y}|\mathbf{x}_{L-1}) \cdots p_{2|1}(\mathbf{x}_2|\mathbf{x}_1)p_{1|0}(\mathbf{x}_1|\mathbf{x})d\mathbf{x}_{L-1} \cdots d\mathbf{x}_1 \right) d\mathbf{y}. \end{aligned} \quad (25)$$

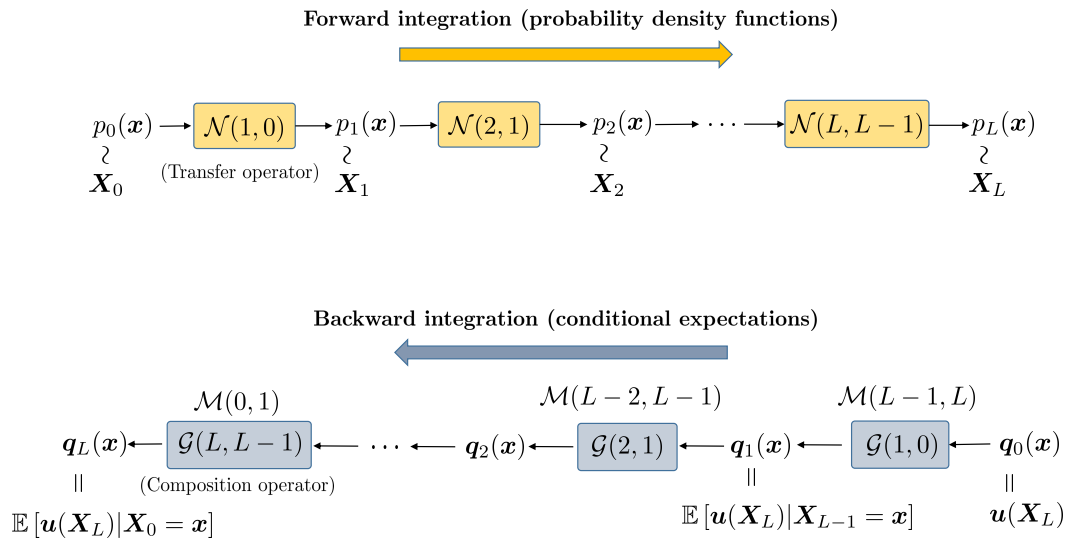


Figure 2: Sketch of the forward/backward integration process for probability density functions (PDFs) and conditional expectations. The transfer operator $\mathcal{N}(n+1, n)$ maps the PDF $p_n(\mathbf{x})$ of the state \mathbf{X}_n into $p_{n+1}(\mathbf{x})$ forward through the neural network. On the other hand, the composition operator \mathcal{M} maps the conditional expectation $\mathbb{E}[\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_{n+1} = \mathbf{x}]$ backwards to $\mathbb{E}[\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_n = \mathbf{x}]$. By defining the operator $\mathcal{G}(n, m) = \mathcal{M}(L-n, L-m)$ we can transform the backward propagation problem for $\mathbb{E}[\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_n = \mathbf{x}]$ into a forward propagation problem for $\mathbf{q}_n(\mathbf{x}) = \mathbb{E}[\mathbf{u}(\mathbf{X}_L)|\mathbf{X}_{L-n} = \mathbf{x}]$.

Relation between composition and transfer operators. The integral operators \mathcal{M} and \mathcal{N} defined in (17) and (11) involve the same kernel function, i.e., the multi-layer transition probability density $p_{q|n}(\mathbf{x}, \mathbf{y})$. In particular, we noticed that $\mathcal{M}(n, q)$ integrates $p_{q|n}$ “from the left”, while $\mathcal{N}(q, n)$ integrates it “from the right”. It is easy to show that $\mathcal{M}(n, q)$ and $\mathcal{N}(q, n)$ are adjoint to each other relative to the standard inner product in L^2 (see [3] for the continuous-time case). In fact,

$$\begin{aligned}
 \mathbb{E}\{\mathbf{u}(\mathbf{X}_k)\} &= \int \mathbb{E}\{\mathbf{u}(\mathbf{X}_k)|\mathbf{X}_q = \mathbf{x}\} p_q(\mathbf{x}) d\mathbf{x} \\
 &= \int [\mathcal{M}(q, j) \mathbb{E}\{\mathbf{u}(\mathbf{X}_k)|\mathbf{X}_j = \mathbf{x}\}] p_q(\mathbf{x}) d\mathbf{x} \\
 &= \int \mathbb{E}\{\mathbf{u}(\mathbf{X}_k)|\mathbf{X}_j = \mathbf{x}\} \mathcal{N}(j, q) p_q(\mathbf{x}) d\mathbf{x}.
 \end{aligned} \tag{26}$$

Therefore

$$\mathcal{M}(q, j)^* = \mathcal{N}(j, q) \quad \forall q, j \in \{0, \dots, L\}, \tag{27}$$

where $\mathcal{M}(q, j)^*$ denotes the operator adjoint of $\mathcal{M}(q, j)$ with respect to the L^2 inner product. By invoking the definition (23), we can also write (27) as

$$\mathcal{G}(L-q, L-j)^* = \mathcal{N}(j, q), \quad \forall j, q \in \{0, \dots, L\}. \tag{28}$$

In Appendix A we show that if the cumulative distribution function of each random vector ξ_n in the noise process has partial derivatives that are Lipschitz in $\mathcal{R}(\xi_n)$ (range of ξ_n), then the composition and transfer operators defined in Eqs. (17) and 11 are *bounded* in L^2 (see Proposition 0.9 and Proposition 0.10). Moreover, is possible to choose the probability density of ξ_n such that the single layer composition and transfer operators become strict *contractions*.

Conditional transition density

We have seen that the composition and the transfer operators \mathcal{M} and \mathcal{N} defined in Eqs. (17) and (11), allow us to push forward and backward conditional expectations and probability densities across the entire neural network. Moreover such operators are adjoint to one another (see equation (27)) [3, 22, 2], and also have the same kernel, i.e., the transition density $p_{n|q}(\mathbf{x}_n|\mathbf{x}_q)$. In this section, we determine an explicit expression for such transition density. To this end, we first derive analytical formulas for the one-layer transition density $p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n)$ for various types of neural network models. The multi-layer transition density $p_{n|q}(\mathbf{x}_n|\mathbf{x}_q)$ is then obtained by composing one-layer transition densities as follows

$$p_{n|q}(\mathbf{x}_n|\mathbf{x}_q) = \int \cdots \int p_{n|n-1}(\mathbf{x}_n|\mathbf{x}_{n-1}) \cdots p_{q+1|q}(\mathbf{x}_{q+1}|\mathbf{x}_q) d\mathbf{x}_{n-1} \cdots d\mathbf{x}_{q+1}. \quad (29)$$

Neural network with additive noise. Let us first consider the neural network model

$$\mathbf{X}_{n+1} = \mathbf{F}(\mathbf{X}_n, \mathbf{w}_n) + \boldsymbol{\xi}_n, \quad (30)$$

where $\{\boldsymbol{\xi}_0, \boldsymbol{\xi}_1, \dots\}$ is a discrete (vector-valued) Markov process indexed by “ n ”. By (30), \mathbf{X}_{n+1} is the sum of two *independent* random vectors³, i.e., $\mathbf{F}(\mathbf{X}_n, \mathbf{w}_n)$ and $\boldsymbol{\xi}_n$. Given any measurable function $h(\mathbf{x})$ we clearly have

$$\begin{aligned} \mathbb{E}\{h(\mathbf{X}_{n+1})\} &= \int h(\mathbf{x}) p_{n+1}(\mathbf{x}) d\mathbf{x} \\ &= \int \int h(\mathbf{F}(\mathbf{x}, \mathbf{w}_n) + \boldsymbol{\xi}) p_n(\mathbf{x}) \rho_n(\boldsymbol{\xi}) d\mathbf{x} d\boldsymbol{\xi} \\ &= \int \int h(\mathbf{y}) \underbrace{\rho_n(\mathbf{x} - \mathbf{F}(\mathbf{y}, \mathbf{w}_n))}_{p_{n+1|n}(\mathbf{x}|\mathbf{y})} p_n(\mathbf{x}) d\mathbf{y} d\mathbf{x}, \end{aligned} \quad (31)$$

where $\rho_n(\mathbf{x})$ denotes the probability density of the random vector $\boldsymbol{\xi}_n$. Therefore, the one-layer transition density for the neural network model (30) is⁴

$$p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n) = \rho_n(\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)). \quad (34)$$

Note that such transition density depends on the PDF of the random noise ρ_n , the neural activation function \mathbf{F} , and the neural network weights \mathbf{w}_n .

Neural network with random weights and random biases. Next, consider the recurrent neural network model

$$\mathbf{X}_{n+1} = \boldsymbol{\varphi}(\mathbf{W}_n(\omega)\mathbf{X}_n + \mathbf{b}_n(\omega)), \quad (35)$$

³Recall that \mathbf{X}_n and $\boldsymbol{\xi}_n$ are statistically independent random vectors. Hence, $\mathbf{F}(\mathbf{X}_n, \mathbf{w}_n)$ and $\boldsymbol{\xi}_n$ are statistically independent random vectors.

⁴Equation (34) can be derived in a more general setting by recalling the conditional probability identity

$$p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n) = \int p_{\mathbf{X}_{n+1}|\mathbf{X}_n, \boldsymbol{\xi}_n}(\mathbf{x}_{n+1}|\mathbf{x}_n, \mathbf{z}) \rho_n(\mathbf{z}) d\mathbf{z}. \quad (32)$$

The conditional density of \mathbf{X}_{n+1} given $\mathbf{X}_n = \mathbf{x}_n$ and $\boldsymbol{\xi}_n = \mathbf{z}$, i.e., $p_{\mathbf{X}_{n+1}|\mathbf{X}_n, \boldsymbol{\xi}_n}(\mathbf{x}_{n+1}|\mathbf{x}_n, \mathbf{z})$, can be immediately computed by using (30) as

$$p_{\mathbf{X}_{n+1}|\mathbf{X}_n, \boldsymbol{\xi}_n}(\mathbf{x}_{n+1}|\mathbf{x}_n, \mathbf{z}) = \delta(\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n) - \mathbf{z}), \quad (33)$$

where $\delta(\mathbf{x})$ is the multivariate Dirac delta function. Substituting (33) into (32), and integrating over \mathbf{z} yields (34).

where $\mathbf{W}_n(\omega)$ are random weight matrices and $\mathbf{b}_n(\omega)$ are random bias vectors. The PDF of \mathbf{X}_{n+1} given $\mathbf{X}_n = \mathbf{x}_n$, i.e., the conditional density we are interested in, can be obtained first by computing the PDF of the random vector

$$\mathbf{Z}_n(\omega) = \mathbf{W}_n(\omega)\mathbf{X}_n + \mathbf{b}_n(\omega), \quad (36)$$

i.e., a linear mapping between independent random variables, and then computing the PDF of $\mathbf{X}_{n+1} = \varphi(\mathbf{Z}_n)$, where φ is the (invertible) activation function. By using the methods we have seen in chapter 1, it is rather straightforward to obtain an expression for the conditional density of \mathbf{X}_{n+1} given $\mathbf{X}_n = \mathbf{x}_n$ for specific probability distributions of $\mathbf{W}_n(\omega)$ (random matrix ensembles) and \mathbf{b}_n .

General neural networks. The transition density of a general neural network of the form

$$\mathbf{X}_{n+1} = \mathbf{H}(\mathbf{X}_n, \mathbf{w}_n, \boldsymbol{\xi}_n), \quad (37)$$

where $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ are statistically independent and do not depend on $\{\mathbf{X}_j\}$ can be written as (see, e.g., [7])

$$p(\mathbf{x}_{n+1}|\mathbf{x}_n) = \int \underbrace{\delta(\mathbf{x}_{n+1} - \mathbf{H}(\mathbf{x}_n, \mathbf{w}_n, \boldsymbol{\xi}_n))}_{p(\mathbf{x}_{n+1}|\mathbf{x}_n, \boldsymbol{\xi}_n)} p(\boldsymbol{\xi}_n) d\boldsymbol{\xi}_n. \quad (38)$$

Remark: The transition density (34) associated with the neural network model (30) can be computed explicitly once we choose a probability model for $\boldsymbol{\xi}_n \in \mathbb{R}^N$. For instance, if we assume that $\{\boldsymbol{\xi}_0, \boldsymbol{\xi}_1, \dots\}$ are i.i.d. Gaussian random vectors with PDF

$$\rho_n(\boldsymbol{\xi}) = \frac{1}{(2\pi)^{N/2}} e^{-\boldsymbol{\xi}^T \boldsymbol{\xi} / 2} \quad \text{for all } n = 0, \dots, L \quad (39)$$

then we can explicitly write the one-layer transition density (34) as

$$p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n) = \frac{1}{(2\pi)^{N/2}} \exp \left[-\frac{[\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)]^T [\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)]}{2} \right]. \quad (40)$$

In Appendix A we provide an analytical example of transition density for a neural network with two layers (one neuron per layer), $\tanh(\cdot)$ activation function, and uniformly distributed random noise.

The zero noise limit An important question is what happens to the neural network as we send the amplitude of the noise to zero. To answer this question consider the system (30) and introduce the parameter $\epsilon \geq 0$, i.e.,

$$\mathbf{X}_{n+1} = \mathbf{F}(\mathbf{X}_n, \mathbf{w}_n) + \epsilon \boldsymbol{\xi}_n, \quad (41)$$

We are interested in studying the orbits of this system as $\epsilon \rightarrow 0$. To this end, we assume the $\boldsymbol{\xi}_n$ to be independent random vectors each having the same density $\rho(\mathbf{x})$. This implies that for all $n = 0, \dots, L-1$, the PDF of $\epsilon \boldsymbol{\xi}_n$ is

$$\epsilon \boldsymbol{\xi}_n \sim \frac{1}{\epsilon} \rho_n \left(\frac{\mathbf{x}}{\epsilon} \right). \quad (42)$$

It is shown in [11, Proposition 10.6.1] that the operator $\mathcal{N}(n+1, n)$ defined in (11)

$$\begin{aligned} p_{n+1}(\mathbf{x}) &= \mathcal{N}(n+1, n) p_n(\mathbf{x}) \\ &= \int \frac{1}{\epsilon} \rho_n \left(\frac{\mathbf{x} - \mathbf{F}(\mathbf{z}, \mathbf{w}_n)}{\epsilon} \right) p_n(\mathbf{z}) d\mathbf{z} \end{aligned} \quad (43)$$

converges in norm to the Frobenius-Perron operator corresponding to $\mathbf{F}(\mathbf{X}_n, \mathbf{w}_n)$ as $\epsilon \rightarrow 0$. Indeed, in the limit $\epsilon \rightarrow 0$ we have, formally

$$\lim_{\epsilon \rightarrow 0} p_{n|n+1}(\mathbf{x}_{n+1}|\mathbf{x}_n) = \lim_{\epsilon \rightarrow 0} \int \frac{1}{\epsilon} \rho_n \left(\frac{\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)}{\epsilon} \right) = \delta(\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)). \quad (44)$$

Substituting this expression into (11), one gets,

$$p_{n+1}(\mathbf{x}) = \mathcal{N}(n+1, n) p_n(\mathbf{x}) = \int \delta(\mathbf{x} - \mathbf{F}(\mathbf{z}, \mathbf{w}_n)) p_n(\mathbf{z}) d\mathbf{z}. \quad (45)$$

Similarly, a substitution into equation (24) yields

$$\mathbf{q}_{n+1}(\mathbf{x}) = \mathcal{G}(n+1, n) \mathbf{q}_n(\mathbf{x}) = \mathbf{q}_n(\mathbf{F}(\mathbf{x}, \mathbf{w}_{L-n-1})), \quad (46)$$

i.e., the familiar function composition representation of neural network mappings

$$\mathbf{q}_{n+1} = \mathbf{q}_0 \left(\mathbf{F}(\mathbf{F}(\cdots \mathbf{F}(\mathbf{x}, \mathbf{w}_{L-n}) \cdots, \mathbf{w}_{L-1}), \mathbf{w}_L) \right). \quad (47)$$

Training over weights versus training over noise

By adding random noise to the output of each layer in a neural network we are essentially adding an infinite number of degrees of freedom to our system. This allows us to rethink the process of training the neural network from a probabilistic perspective. In particular, instead of optimizing a performance metric⁵ relative to the neural network weights $\mathbf{w} = \{\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_{L-1}\}$ for fixed noise, we can now optimize the transition density⁶ $p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n)$. Clearly, such transition density is connected to the neural network weights, e.g., by equation (34). Hence, if we prescribe the PDF of the random noise, i.e., $\rho_n(\cdot)$ in (34), then the transition density $p_{n+1|n}$ is uniquely determined by the functional form of the activation function \mathbf{F} , and by the weights \mathbf{w}_n . On the other hand, we can *optimize* ρ_n (probability density of the random noise ξ_n) while keeping the weights \mathbf{w}_n *fixed*. As we shall see hereafter, this process opens the possibility approximately encode and decode secret messages in a fully trained neural network using random noise.

Encoding secret messages in neural networks using random noise. An interesting question is whether random noise added to the output of each layer in the neural network can enhance features of the output, or allow us to encode/decode secret signals in the network. The interaction between random noise and the nonlinear dynamics modeled by the network can yield indeed many surprising results. For example, in stochastic resonance [16, 20] it is well known that random noise added to a properly tuned bi-stable system can induce a peak in the Fourier power spectrum of the output, hence effectively amplifying the signal. Similarly, random noise added to a neural network can have remarkable effects. In particular, it allows us to re-purpose (to some extent) a previously trained network by hiding a secret signal in it, which can be approximately encoded and decoded by using random noise. To this end, it is sufficient to optimize

⁵In a supervised learning setting the neural network weights are usually determined by minimizing a dissimilarity measure between the output of the network and a target function. Such measure may be an entropy measure, the Wasserstein distance, the Kullback–Leibler divergence, or other measures defined by classical L^p norms.

⁶In a deterministic setting, the transition density for a neural network model of the form $\mathbf{X}_{n+1} = \mathbf{F}(\mathbf{X}_n, \mathbf{w}_n)$ is simply

$$p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n) = \delta(\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)), \quad (48)$$

where $\delta(\cdot)$ is the Dirac delta function. Such density does not have any degree of freedom other than \mathbf{w}_n . On the other hand, in a stochastic setting we are free to *choose* the PDF of ξ_n . For a neural network model of the form $\mathbf{X}_{n+1} = \mathbf{F}(\mathbf{X}_n, \mathbf{w}_n) + \xi_n$ the transition density has the form

$$p_{n+1|n}(\mathbf{x}_{n+1}|\mathbf{x}_n) = \rho_n(\mathbf{x}_{n+1} - \mathbf{F}(\mathbf{x}_n, \mathbf{w}_n)), \quad (49)$$

where $\rho_n(\xi)$ is the PDF of ξ_n . We are clearly free to choose the functional form of ρ_n .

the PDF of the noise appropriately, and then train the conditional expectation of the output over such PDF.

To describe the method, suppose that we are given a fully trained deterministic neural network with only two layers, and weights chosen to represent an input-output map defined on some domain $\Omega \subseteq \mathbb{R}^d$. In the absence of noise we can write the output of the neural network as

$$q_2(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{F}(\mathbf{F}_0(\mathbf{x}, \mathbf{w}_0), \mathbf{w}_1) \quad (50)$$

where $\{\boldsymbol{\alpha}, \mathbf{w}_0, \mathbf{w}_1\}$ can be optimized to minimize the distance between $q_2(\mathbf{x})$ and a given target function $f(\mathbf{x})$ ($\mathbf{x} \in \Omega$). Injecting noise $\boldsymbol{\xi}_0$ in the output of the first layer yields the input-output map

$$h_2(\mathbf{x}) = \boldsymbol{\alpha}^T \int \mathbf{F}(\mathbf{y} + \mathbf{F}_0(\mathbf{x}, \mathbf{w}_0), \mathbf{w}_1) \rho_0(\mathbf{y}) d\mathbf{y}, \quad (51)$$

where $\{\boldsymbol{\alpha}, \mathbf{w}_1, \mathbf{w}_0\}$ here are fixed, and ρ_0 is the PDF of $\boldsymbol{\xi}_0$. Equation (51) resembles a Fredholm integral equation of the first kind. In fact, it can be written as

$$h_2(\mathbf{x}) = \int \kappa_2(\mathbf{x}, \mathbf{y}) \rho_0(\mathbf{y}) d\mathbf{y}, \quad (52)$$

where

$$\kappa_2(\mathbf{x}, \mathbf{y}) = \boldsymbol{\alpha}^T \mathbf{F}(\mathbf{y} + \mathbf{F}_0(\mathbf{x}, \mathbf{w}_0), \mathbf{w}_1). \quad (53)$$

However, differently from standard Fredholm equations of the first kind, here we have $\mathbf{x} \in \Omega \subseteq \mathbb{R}^d$ while $\mathbf{y} \in \mathbb{R}^N$, i.e., the integral operator with kernel κ_2 maps functions in N variables into functions in d variables. We are interested in finding a PDF $\rho_0(\mathbf{y})$ that solves (51) for a given function $h_2(\mathbf{x})$. In other words, we are re-purposing the neural network (50) with output $q_2(\mathbf{x}) \simeq f(\mathbf{x})$ to approximate now a different function $h_2(\mathbf{x}) \simeq v(\mathbf{x})$, without modifying the weights $\{\boldsymbol{\alpha}, \mathbf{w}_1, \mathbf{w}_0\}$ but rather simply adding noise $\boldsymbol{\xi}_0$ and averaging the output over the PDF ρ_0 of the noise (Eq. (52)). Equation (52) is unfortunately ill-posed in the space of probability distributions. In other words, for a given kernel κ_2 and a given target q_2 there is (in general) no PDF ρ_0 that satisfies (52) exactly. However, one can proceed by optimization. For instance, ρ_0 can be determined by solving the constrained least squares problem⁷

$$\rho_0 = \underset{\rho}{\operatorname{argmin}} \left\| h_2(\mathbf{x}) - \int \kappa_2(\mathbf{x}, \mathbf{y}) \rho(\mathbf{y}) d\mathbf{y} \right\|_{L^2(\Omega)} \quad \text{subject to} \quad \|\rho\|_{L^1(\mathbb{R}^N)} = 1 \quad \rho \geq 0. \quad (54)$$

⁷The optimization problem (54) is a quadratic program with linear constraints if we represent ρ_0 in the span of a basis made of positive functions, e.g., Gaussian kernels [1].

Appendix A: Functional setting

Let $(\mathcal{S}, \mathcal{F}, \mathcal{P})$ be a probability space. Consider the neural network model (see Figure 1)

$$\mathbf{X}_1 = \mathbf{F}_0(\mathbf{X}_0, \mathbf{w}_0) + \boldsymbol{\xi}_0 \quad \mathbf{X}_{n+1} = \mathbf{F}(\mathbf{X}_n, \mathbf{w}_n) + \boldsymbol{\xi}_n \quad n = 1, \dots, L-1, \quad (55)$$

where $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ is a discrete (vector-valued) Markov process. Suppose we are interested in using the model (55) to approximate a multivariate function $f(\mathbf{x})$. This is usually done by taking a linear combination of the network output, e.g., Eq. (4). In this setting, the neural network can be thought of as a process of constructing an adaptive basis by function composition. Here we consider the case where the function we are approximating is defined on a compact subset Ω_0 of \mathbb{R}^d . This means that the input vector of the neural network, i.e. \mathbf{X}_0 , is an element of Ω_0 . We assume that the following conditions are satisfied

1. $\mathbf{X}_0 \in \Omega_0 \subseteq \mathbb{R}^d$ (Ω compact), $\mathbf{X}_n \in \mathbb{R}^N$ for $n = 1, \dots, L-1$;
2. The image of \mathbf{F}_0 and \mathbf{F} is the hyper-cube $[-1, 1]^N$.

For example, if \mathbf{F} in (55) is of the form

$$\mathbf{F}(\mathbf{x}, \mathbf{w}) = \tanh(\mathbf{W}\mathbf{x} + \mathbf{b}) \quad \mathbf{w} = \{\mathbf{W}, \mathbf{b}\}, \quad (56)$$

then conditions 1. and 2. imply that $\mathbf{W}_0 \in M_{N \times d}(\mathbb{R})$ and $\mathbf{W}_n \in M_{N \times N}(\mathbb{R})$ for $n = 1, \dots, L-1$, while the biases are $\mathbf{b}_n \in M_{N \times 1}(\mathbb{R})$ for $n = 0, \dots, L-1$. The random vectors $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ added to the output of each layer make $\{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ a discrete Markov process (each \mathbf{X}_i is a random vector). The range of \mathbf{X}_{n+1} depends essentially on the range of $\boldsymbol{\xi}_n$, as the image of \mathbf{F} is the hyper-cube $[-1, 1]^N$ (see condition 2. above). Let us define⁸

$$\begin{aligned} \Omega_{n+1} &= [-1, 1]^N + \mathcal{R}(\boldsymbol{\xi}_n) \\ &= \{\mathbf{c} \in \mathbb{R}^N : \mathbf{c} = \mathbf{a} + \mathbf{b} \quad \mathbf{a} \in [-1, 1]^N, \mathbf{b} \in \mathcal{R}(\boldsymbol{\xi}_n)\}, \end{aligned} \quad (58)$$

where $\mathcal{R}(\boldsymbol{\xi}_n)$ denotes the range of the random vector $\boldsymbol{\xi}_n$, i.e.,

$$\mathcal{R}(\boldsymbol{\xi}_n) = \{\boldsymbol{\xi}_n(\omega) \in \mathbb{R}^N : \omega \in \mathcal{S}\}. \quad (59)$$

Clearly, the range of the random vector \mathbf{X}_{n+1} is a subset⁹ of Ω_{n+1} , i.e., $\mathcal{R}(\mathbf{X}_{n+1}) \subseteq \Omega_{n+1}$. This implies the following lemma.

Lemma 0.1. *Let $\lambda(\Omega_{n+1})$ the Lebesgue measure of the set (58). The Lebesgue measure of the range of \mathbf{X}_{n+1} satisfies*

$$\lambda(\mathcal{R}(\mathbf{X}_{n+1})) \leq \lambda(\Omega_{n+1}). \quad (60)$$

Proof. The proof follows from the inclusion $\mathcal{R}(\mathbf{X}_{n+1}) \subseteq \Omega_{n+1}$. □

⁸The notation $[-1, 1]^N$ denotes a Cartesian product of N one-dimensional domains $[-1, 1]$, i.e.,

$$[-1, 1]^N = \prod_{k=1}^N [-1, 1] = \underbrace{[-1, 1] \times [-1, 1] \times \dots \times [-1, 1]}_{N \text{ times}}. \quad (57)$$

⁹We emphasize that if we are given a specific form of the activation function \mathbf{F} together with suitable bounds on the neural network weights and biases $\{\mathbf{W}, \mathbf{b}\}$ then we can easily identify a domain that is smaller than Ω_n , and that still contains $\mathcal{R}(\mathbf{X}_n)$. This allows us to construct a tighter bound for $\lambda(\mathcal{R}(\mathbf{X}_{n+1}))$ in Lemma 0.1, which depends on the activation function and on the bounds we set on neural network weights and biases.

The L^∞ norm of the random vector $\boldsymbol{\xi}$ is defined as the largest value of $r \geq 0$ that yields a nonzero probability on the event $\{\omega \in \mathcal{S} : \|\boldsymbol{\xi}(\omega)\|_\infty > r\} \in \mathcal{F}$, i.e.,

$$\|\boldsymbol{\xi}\|_\infty = \sup_{r \in \mathbb{R}} \{\mathcal{P}(\{\omega \in \mathcal{S} : \|\boldsymbol{\xi}(\omega)\|_\infty > r\}) > 0\}. \quad (61)$$

This definition allows us to bound the Lebesgue measure of Ω_{n+1} as follows.

Proposition 0.2. *The Lebesgue measure of the set Ω_{n+1} defined in (58) can be bounded as*

$$\lambda(\Omega_{n+1}) \leq \left(\sqrt{N} + \|\boldsymbol{\xi}_n\|_\infty\right)^N \frac{\pi^{N/2}}{\Gamma(1 + N/2)}, \quad (62)$$

where N is the number of neurons and $\Gamma(\cdot)$ is the Gamma function.

Proof. As is well known, the length of the diagonal of the hypercube $[-1, 1]^N$ is \sqrt{N} . Hence, $\sqrt{N} + \|\boldsymbol{\xi}_n\|_\infty$ is the radius of a ball that encloses all elements of Ω_{n+1} . The Lebesgue measure of such ball is obtained by multiplying the Lebesgue measure of the unit ball in \mathbb{R}^N , i.e., $\pi^{N/2}/\Gamma(1 + N/2)$ by the scaling factor $\left(\sqrt{N} + \|\boldsymbol{\xi}_n\|_\infty\right)^N$. □

Lemma 0.3. *If $\mathcal{R}(\boldsymbol{\xi}_n)$ is bounded then $\mathcal{R}(\mathbf{X}_{n+1})$ is bounded.*

Proof. The image of the activation function \mathbf{F} is a bounded set. If $\mathcal{R}(\boldsymbol{\xi}_n)$ is bounded then Ω_{n+1} in (58) is bounded. $\mathcal{R}(\mathbf{X}_{n+1}) \subseteq \Omega_{n+1}$ and therefore $\mathcal{R}(\mathbf{X}_{n+1})$ is bounded. □

Clearly, if $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ are i.i.d. random variables then there exists a domain $V = \Omega_1 = \dots = \Omega_L$ such that

$$\mathcal{R}(\boldsymbol{\xi}_n) \subseteq \mathcal{R}(\mathbf{X}_{n+1}) \subseteq V \quad \forall n = 0, \dots, L-1. \quad (63)$$

In fact, if $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ are i.i.d. random variables then we have

$$\mathcal{R}(\boldsymbol{\xi}_0) = \mathcal{R}(\boldsymbol{\xi}_1) = \dots = \mathcal{R}(\boldsymbol{\xi}_{L-1}), \quad (64)$$

which implies that all of Ω_i defined in (58) are the same. If the range of each random vector $\boldsymbol{\xi}_n$ is a tensor product of one-dimensional domain, e.g., if the components of $\boldsymbol{\xi}_n$ are statistically independent, then $V = \Omega_1 = \dots = \Omega_L$ becomes particularly simple, i.e., a hypercube.

Lemma 0.4. *Let $\{\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{L-1}\}$ be i.i.d. random variables with bounded range and suppose that each $\boldsymbol{\xi}_k$ has statistically independent components with range $[a, b]$. Then all domains $\{\Omega_1, \dots, \Omega_L\}$ defined in equation (58) are the same, and they are equivalent to*

$$V = \bigtimes_{k=1}^N [-1 + a, 1 + b]. \quad (65)$$

V includes the range of all random vectors \mathbf{X}_n ($n = 1, \dots, L$) and has Lebesgue measure

$$\lambda(V) = (2 + b - a)^N. \quad (66)$$

Proof. The proof is trivial and therefore omitted. □

Remark: It is worth noticing that if each ξ_k is a uniformly distributed random vector with statistically independent components in $[-1, 1]$, then for $N = 10$ (number of neurons) the upper bound in (62) is 3.98×10^6 while the exact result (66) gives 1.05×10^6 . Hence the estimate (62) is quite sharp in the case of uniform random vectors.

Boundedness of composition and transfer operators

Lemma 0.3 states that if we perturb the output of the n -th layer of a neural network by a random vector ξ_n with finite range then we obtain a random vector \mathbf{X}_{n+1} with finite range. In this hypothesis it is straightforward to show that the composition and transfer operators defined in (17) and (11) are bounded. We have seen that these operators can be written as

$$\mathcal{M}(n, n+1)v = \int_{\mathcal{R}(\mathbf{X}_{n+1})} v(\mathbf{y})p_{n+1|n}(\mathbf{y}|\mathbf{x})d\mathbf{y}, \quad \mathcal{N}(n+1, n)v = \int_{\mathcal{R}(\mathbf{X}_n)} p_{n+1|n}(\mathbf{x}|\mathbf{y})v(\mathbf{y})d\mathbf{y}, \quad (67)$$

where $p_{n+1|n}(\mathbf{y}|\mathbf{x}) = \rho_n(\mathbf{y} - \mathbf{F}(\mathbf{x}, \mathbf{w}_n))$ is the conditional transition density of \mathbf{X}_{n+1} given \mathbf{X}_n , and ρ_n is the joint PDF of the random vector ξ_n . The conditional transition density $p_{n+1|n}(\mathbf{y}|\mathbf{x})$ is always non-negative, i.e.,

$$p_{n+1|n}(\mathbf{y}|\mathbf{x}) \geq 0 \quad \forall \mathbf{y} \in \mathcal{R}(\mathbf{X}_{n+1}), \quad \forall \mathbf{x} \in \mathcal{R}(\mathbf{X}_n). \quad (68)$$

Moreover, the conditional density $p_{n+1|n}$ is defined on the set

$$\mathcal{B}_n = \{(\mathbf{x}, \mathbf{y}) \in \mathcal{R}(\mathbf{X}_n) \times \mathcal{R}(\mathbf{X}_{n+1}) : (\mathbf{y} - \mathbf{F}(\mathbf{x}, \mathbf{w}_n)) \in \mathcal{R}(\xi_n)\}. \quad (69)$$

It is also important to emphasize that $\mathbf{y} \in \mathcal{R}(\mathbf{X}_{n+1})$ and $\mathbf{x} \in \mathcal{R}(\mathbf{X}_n)$. Both $\mathcal{R}(\mathbf{X}_{n+1})$ and $\mathcal{R}(\mathbf{X}_n)$ depend on Ω_0 (domain of the input), the neural network weights, and the noise amplitude. Thanks to Lemma 0.1, we have that

$$\mathcal{B}_n \subseteq \Omega_n \times \Omega_{n+1}. \quad (70)$$

The Lebesgue measure of \mathcal{B}_n can be calculated as follows.

Lemma 0.5. *The Lebesgue measure of the set \mathcal{B}_n defined in (69) is equal to the product of the measure of $\lambda(\mathcal{R}(\mathbf{X}_n))$ and the measure of $\mathcal{R}(\xi_n)$, i.e.,*

$$\lambda(\mathcal{B}_n) = \lambda(\mathcal{R}(\mathbf{X}_n))\lambda(\mathcal{R}(\xi_n)). \quad (71)$$

Moreover, $\lambda(\mathcal{B}_n)$ is bounded by $\lambda(\mathcal{R}(\Omega_n))\lambda(\mathcal{R}(\xi_n))$, which is independent of the neural network weights.

Proof. Let χ_n be the indicator function of the set $\mathcal{R}(\xi_n)$, $\mathbf{y} \in \mathcal{R}(\mathbf{X}_{n+1})$ and $\mathbf{x} \in \mathcal{R}(\mathbf{X}_n)$. Then

$$\begin{aligned} \lambda(\mathcal{B}_n) &= \int_{\mathcal{R}(\mathbf{X}_{n+1})} \int_{\mathcal{R}(\mathbf{X}_n)} \chi_n(\mathbf{y} - \mathbf{F}(\mathbf{x}, \mathbf{w}_n))d\mathbf{x}d\mathbf{y} \\ &= \lambda(\mathcal{R}(\xi_n)) \int_{\mathcal{R}(\mathbf{X}_n)} d\mathbf{x} \\ &= \lambda(\mathcal{R}(\mathbf{X}_n))\lambda(\mathcal{R}(\xi_n)). \end{aligned} \quad (72)$$

By using Lemma 0.1 we conclude that $\lambda(\mathcal{B}_n)$ is bounded from above by $\lambda(\mathcal{R}(\Omega_{L-m}))\lambda(\mathcal{R}(\xi_n))$, which is independent of the neural network weights. □

Remark: The result (71) has a straightforward geometrical interpretation in two dimensions. Pick a ruler of length $r = \lambda(\mathcal{R}(\xi_n))$ with endpoints that can leave markings if we slide it on a rectangular table with side lengths $s_b = \lambda(\mathcal{R}(X_{n+1}))$ (horizontal side) $s_h = \lambda(\mathcal{R}(X_n))$ (vertical side). Then slide the ruler from the top to the bottom of the table, while keeping it horizontal, i.e., parallel to the horizontal sides of the table (see Figure 3). The area of the domain defined by the two curves drawn by the endpoints of the ruler is always $r \times s_h$ independently of the way we slide the ruler laterally – provided the ruler never gets out of the table.

Lemma 0.6. *If the range of ξ_{n-1} is a bounded subset of \mathbb{R}^N then the transition density $p_{n+1|n}(\mathbf{y}|\mathbf{x})$ is an element of $L^1(\mathcal{R}(\mathbf{X}_{n+1}) \times \mathcal{R}(\mathbf{X}_n))$.*

Proof. Note that

$$\int_{\mathcal{R}(\mathbf{X}_{n+1})} \int_{\mathcal{R}(\mathbf{X}_n)} p_{n+1|n}(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x} = \lambda(\mathcal{R}(\mathbf{X}_n)) \leq \lambda(\Omega_n). \quad (73)$$

The Lebesgue measure $\lambda(\Omega_n)$ can be bounded as (see Proposition 0.2)

$$\lambda(\Omega_n) \leq \left(\sqrt{N} + \|\xi_{n-1}\|_\infty \right)^N \frac{\pi^{N/2}}{\Gamma(1 + N/2)}. \quad (74)$$

Since the range of ξ_{n-1} is bounded by hypothesis we have that there exists a finite real number $M > 0$ such that $\|\xi_{n-1}\|_\infty \leq M$. This implies that the integral in (73) is finite, i.e., that the transition kernel $p_{n+1|n}(\mathbf{y}|\mathbf{x})$ is in $L^1(\mathcal{R}(\mathbf{X}_{n+1}) \times \mathcal{R}(\mathbf{X}_n))$. □

Theorem 0.7. *Let $C_{\xi_n}(\mathbf{x})$ be the cumulative distribution function ξ_n . If $C_{\xi_n}(\mathbf{x})$ is Lipschitz continuous on $\mathcal{R}(\xi_n)$ and the partial derivatives $\partial C_{\xi_n}/\partial x_k$ ($k = 1, \dots, N$) are Lipschitz continuous in x_1, x_2, \dots, x_N , respectively, then the joint probability density function of ξ_n is bounded on $\mathcal{R}(\xi_n)$.*

Proof. By using Rademacher's theorem we have that if $C_{\xi_n}(\mathbf{x})$ is Lipschitz on $\mathcal{R}(\xi_n)$ then it is differentiable almost everywhere on $\mathcal{R}(\xi_n)$ (except on a set with zero Lebesgue measure). Therefore the partial derivatives $\partial C_{\xi_n}/\partial x_k$ exist almost everywhere on $\mathcal{R}(\xi_n)$. If, in addition, we assume that $\partial C_{\xi_n}/\partial x_k$ are Lipschitz continuous with respect to x_k (for all $k = 1, \dots, N$) then by applying [15, Theorem 9] recursively we conclude that the joint probability density function of ξ_n is bounded. □

Lemma 0.8. *Under the same assumptions of Theorem 0.7 we have that the conditional PDF $p_{n+1|n}(\mathbf{y}|\mathbf{x}) = \rho_n(\mathbf{y} - \mathbf{F}(\mathbf{x}, \mathbf{w}))$ is bounded on $\mathcal{R}(\mathbf{X}_{n+1}) \times \mathcal{R}(\mathbf{X}_n)$.*

Proof. Theorem 0.7 states that ρ_n is a bounded function. This implies that the conditional density $p_{n+1|n}(\mathbf{y}|\mathbf{x}) = \rho_n(\mathbf{y} - \mathbf{F}(\mathbf{x}, \mathbf{w}))$ is bounded on $\mathcal{R}(\mathbf{X}_{n+1}) \times \mathcal{R}(\mathbf{X}_n)$. □

Proposition 0.9. *Let $\mathcal{R}(\xi_n)$ and $\mathcal{R}(\xi_{n-1})$ be bounded subsets of \mathbb{R}^N . Then, under the same assumptions of Theorem 0.7, we have that the composition and the transfer operators defined in (67) are bounded in L^2 .*

Proof. Let us first prove that $\mathcal{M}(n, n+1)$ is a bounded linear operator from $L^2(\mathcal{R}(\mathbf{X}_{n+1}))$ into $L^2(\mathcal{R}(\mathbf{X}_n))$. To this end, note that

$$\begin{aligned} \|\mathcal{M}(n, n+1)v\|_{L^2(\mathcal{R}(\mathbf{X}_n))}^2 &= \int_{\mathcal{R}(\mathbf{X}_n)} \left| \int_{\mathcal{R}(\mathbf{X}_{n+1})} v(\mathbf{y}) p_{n+1|n}(\mathbf{y}|\mathbf{x}) d\mathbf{y} \right|^2 d\mathbf{x} \\ &\leq \|v\|_{L^2(\mathcal{R}(\mathbf{X}_{n+1}))}^2 \underbrace{\int_{\mathcal{R}(\mathbf{X}_n)} \int_{\mathcal{R}(\mathbf{X}_{n+1})} p_{n+1|n}(\mathbf{y}|\mathbf{x})^2 d\mathbf{y} d\mathbf{x}}_{K_n} \\ &= K_n \|v\|_{L^2(\mathcal{R}(\mathbf{X}_{n+1}))}^2. \end{aligned} \quad (75)$$

Clearly, $K_n < \infty$. In fact, if $\mathcal{R}(\boldsymbol{\xi}_n)$ and $\mathcal{R}(\boldsymbol{\xi}_{n-1})$ are bounded then $\mathcal{R}(\mathbf{X}_{n+1})$ and $\mathcal{R}(\mathbf{X}_n)$ are bounded. Moreover, thanks to Lemma 0.8 we have that $p_{n+1|n}(\mathbf{y}|\mathbf{x})$ is bounded on $\mathcal{R}(\mathbf{X}_{n+1}) \times \mathcal{R}(\mathbf{X}_n)$. Hence, K_n is the integral of the square of a bounded function defined on a bounded domain, and therefore it is finite. By following the same steps it is straightforward to show that the transfer operator \mathcal{N} is a bounded linear operator. Alternatively, simply recall that \mathcal{N} is the adjoint of \mathcal{M} , and the adjoint of a bounded linear operator is bounded. Specifically we have,

$$\|\mathcal{N}(n+1, n)p\|_{L^2(\mathcal{R}(\mathbf{X}_{n+1}))}^2 \leq K_n \|p\|_{L^2(\mathcal{R}(\mathbf{X}_n))}^2. \quad (76)$$

□

Remark: The integrals

$$K_n = \int_{\mathcal{R}(\mathbf{X}_n)} \int_{\mathcal{R}(\mathbf{X}_{n+1})} p_{n+1|n}(\mathbf{y}|\mathbf{x})^2 d\mathbf{y} d\mathbf{x} \quad (77)$$

can be computed by noting that

$$p_{n+1|n}(\mathbf{y}|\mathbf{x}) = \rho_n(\mathbf{y} - \mathbf{F}(\mathbf{x}, \mathbf{w})) \quad (78)$$

is essentially a *shift* of the PDF ρ_n by a quantity $\mathbf{F}(\mathbf{x}, \mathbf{w})$ that depends on \mathbf{x} and \mathbf{w} (see, e.g., Figure 3). Such a shift does not influence the integral with respect to \mathbf{y} , meaning that the integral of $p_{n+1|n}(\mathbf{y}|\mathbf{x})$ or $p_{n+1|n}(\mathbf{y}|\mathbf{x})^2$ with respect to \mathbf{y} is the same for all \mathbf{x} . Hence, by changing variables we have that the integral (77) is equivalent to

$$K_n = \lambda(\mathcal{R}(\mathbf{X}_n)) \int_{\mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x})^2 d\mathbf{x}, \quad (79)$$

where $\lambda(\mathcal{R}(\mathbf{X}_n))$ is the Lebesgue measure of $\mathcal{R}(\mathbf{X}_n)$, and $\mathcal{R}(\boldsymbol{\xi}_n)$ is the range of $\boldsymbol{\xi}_n$. Note that K_n depends on the neural net weights only through the Lebesgue measure of $\mathcal{R}(\mathbf{X}_n)$. Clearly, since the set Ω_n includes $\mathcal{R}(\mathbf{X}_n)$ we have by Lemma 0.1 that $\lambda(\mathcal{R}(\mathbf{X}_n)) \leq \lambda(\Omega_n)$. This implies that

$$K_n \leq \lambda(\Omega_n) \int_{\mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x})^2 d\mathbf{x}. \quad (80)$$

The upper bound here does not depend on the neural network weights. The following lemma summarizes all these remarks.

Proposition 0.10. *Under the same assumptions of Theorem 0.7, we have that the composition and the transfer operators defined in (67) can be bounded as*

$$\|\mathcal{M}(n, n+1)\|^2 \leq K_n, \quad \|\mathcal{N}(n+1, n)\|^2 \leq K_n, \quad (81)$$

where

$$K_n = \lambda(\mathcal{R}(\mathbf{X}_n)) \int_{\mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x})^2 d\mathbf{x}. \quad (82)$$

Moreover, K_n can be bounded as

$$K_n \leq \lambda(\Omega_n) \int_{\mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x})^2 d\mathbf{x}, \quad (83)$$

where Ω_n is defined in (58) and ρ_n is the PDF of $\boldsymbol{\xi}_n$. The upper bound in (83) does not depend on the neural network weights and biases.

Under additional assumptions on the PDF $\rho_n(\mathbf{x})$ it is also possible to bound the integrals at the right hand side of (82) and (83). Specifically we have the following sharp bound.

Lemma 0.11. *Let*

$$s_n = \inf_{\mathbf{x} \in \mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x}), \quad S_n = \sup_{\mathbf{x} \in \mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x}). \quad (84)$$

If $s_n > 0$ then under the same assumptions of Theorem 0.7 we have that

$$\int_{\mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x})^2 d\mathbf{x} \leq \frac{1}{\lambda(\mathcal{R}(\boldsymbol{\xi}_n))} \frac{(S_n + s_n)^2}{4S_n s_n}. \quad (85)$$

Proof. First we notice that if the random vector $\boldsymbol{\xi}_n$ satisfies the assumptions of Theorem 0.7 then the upper bound S_n is finite. By using the definition (84) we have

$$(\rho_n(\mathbf{x}) - s_n)(S_n - \rho_n(\mathbf{x})) \geq 0 \quad \text{for all } \mathbf{x} \in \mathcal{R}(\boldsymbol{\xi}_n). \quad (86)$$

This implies

$$\int_{\mathcal{R}(\boldsymbol{\xi}_n)} \rho_n(\mathbf{x})^2 d\mathbf{x} \leq (S_n + s_n) - S_n s_n \lambda(\mathcal{R}(\boldsymbol{\xi}_n)), \quad (87)$$

where we used the fact that the PDF ρ_n integrates to one over $\mathcal{R}(\boldsymbol{\xi}_n)$. Next, define

$$R_n = \frac{1}{\lambda(\mathcal{R}(\boldsymbol{\xi}_n))} \frac{(S_n + s_n)^2}{4S_n s_n}. \quad (88)$$

Clearly,

$$R_n \left(1 - \frac{2S_n s_n}{s_n + S_n} \lambda(\mathcal{R}(\boldsymbol{\xi}_n)) \right)^2 = R_n - (S_n + s_n) + S_n s_n \lambda(\mathcal{R}(\boldsymbol{\xi}_n)) \geq 0 \quad (89)$$

which implies that

$$(S_n + s_n) - S_n s_n \lambda(\mathcal{R}(\boldsymbol{\xi}_n)) \leq R_n. \quad (90)$$

A substitution of (90) into (87) yields (85). □

An example. Let $X_0 \in \Omega_0 = [-1, 1]$ and consider

$$X_1 = \tanh(X_0 + 3) + \xi_0, \quad X_2 = \tanh(2X_1 - 1) + \xi_1, \quad (91)$$

where ξ_0 and ξ_1 are uniform random variables with range $\mathcal{R}(\xi_0) = \mathcal{R}(\xi_1) = [-2, 2]$. In this setting,

$$\begin{aligned} \mathcal{R}(X_1) &= [\tanh(2) - 2, \tanh(4) + 2], \\ \mathcal{R}(X_2) &= [\tanh(2 \tanh(2) - 5) - 2, \tanh(2 \tanh(4) + 3) + 2]. \end{aligned}$$

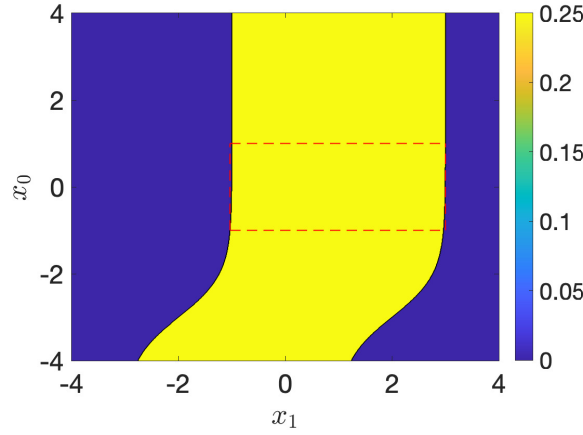


Figure 3: Conditional probability density function $p_{1|0}(x_1|x_0)$ defined in equation (92). The domain $\mathcal{R}(X_1) \times \mathcal{R}(X_0)$ is the interior of the rectangle delimited by dashed red lines.

The conditional density of X_1 given X_0 is given by

$$p_{1|0}(x_1|x_0) = \begin{cases} \frac{1}{4} & \text{if } |x_1 - \tanh(x_0 + 3)| \leq 2 \\ 0 & \text{otherwise} \end{cases} \quad (92)$$

This function is plotted in Figure 3 together with the domain $\mathcal{R}(X_1) \times \mathcal{R}(X_0)$ (interior of the rectangle delimited by dashed red lines). Clearly, the integral of the conditional PDF (92) is

$$\int_{\mathcal{R}(X_0)} \int_{\mathcal{R}(X_1)} p_{1|0}(x_1|x_0) dx_1 dx_0 = \lambda(\mathcal{R}(X_0)) = 2, \quad (93)$$

where $\lambda(\mathcal{R}(X_0))$ is the Lebesgue measure of $\mathcal{R}(X_0) = [-1, 1]$. The L^2 norm of the operators \mathcal{N} and \mathcal{M} is bounded by¹⁰

$$K_0 = \int_{\mathcal{R}(X_0)} \int_{\mathcal{R}(X_1)} p_{1|0}(x_1|x_0)^2 dx_1 dx_0 = \frac{\lambda(\mathcal{R}(X_0))}{\lambda(\mathcal{R}(\xi_0))} = \frac{1}{2}. \quad (96)$$

Hence, both operators $\mathcal{N}(1, 0)$ and $\mathcal{M}(0, 1)$ are contractions (Proposition 0.10). On the other hand,

$$K_1 = \frac{\lambda(\mathcal{R}(X_1))}{\lambda(\mathcal{R}(\xi_1))} = 1 + \frac{\tan(4) - \tan(2)}{4} > 1. \quad (97)$$

Next, define V as in Lemma 0.4, i.e., $V = [-3, 3]$. Clearly, both $\mathcal{R}(X_0)$ and $\mathcal{R}(X_1)$ are subsets of V . If we integrate the conditional PDF shown in Figure 3 in $V \times V$ we obtain

$$\int_V \int_V p_{1|0}(x_1|x_0)^2 dx_1 dx_0 = \frac{\lambda(V)}{\lambda(\mathcal{R}(\xi_0))} = \frac{3}{2}. \quad (98)$$

¹⁰For uniformly distributed random variables we have that

$$\int_{\mathcal{R}(\xi_n)} \rho_n(\mathbf{x})^2 d\mathbf{x} = \frac{1}{\lambda(\mathcal{R}(\xi_n))}. \quad (94)$$

Therefore equation (82) yields

$$K_n = \frac{\lambda(\mathcal{R}(\mathbf{X}_n))}{\lambda(\mathcal{R}(\xi_n))} \leq \frac{\lambda(\Omega_n)}{\lambda(\mathcal{R}(\xi_n))}. \quad (95)$$

Depending on the ratio between the Lebesgue measure of $\mathcal{R}(\mathbf{X}_n)$ and $\mathcal{R}(\xi_n)$ one can have K_n smaller or larger than 1.

Random noise can induce contractions

In this section we prove a result on neural networks perturbed by random noise of increasing amplitude which states that it is possible to make both operators \mathcal{N} and \mathcal{M} in (67) contractions¹¹ if the noise is properly chosen. To this end, we begin with the following lemma.

Lemma 0.12. *Let*

$$\|\rho_n\|_{L^2(\mathcal{R}(\xi_n))}^2 = \int_{\mathcal{R}(\xi_n)} \rho_n(\mathbf{x})^2 d\mathbf{x}. \quad (99)$$

If

$$\|\rho_n\|_{L^2(\mathcal{R}(\xi_n))}^2 \leq \frac{\kappa}{\lambda(\Omega_n)} \quad 0 \leq \kappa < 1 \quad (100)$$

then $\mathcal{M}(n, n+1)$ and $\mathcal{N}(n+1, n)$ are contractions. The condition (100) is independent of the neural network weights.

Proof. The proof follows from equation (83). □

Hereafter we specialize Lemma 0.12 to neural network perturbed by uniformly distributed random noise.

Proposition 0.13. *Let $\{\xi_0, \dots, \xi_{L-1}\}$ be independent random vectors. Suppose that the components of each ξ_n are zero-mean i.i.d. uniform random variables with range $[-b_n, b_n]$ ($b_n > 0$). If*

$$b_0 \geq \frac{1}{2} \left(\frac{\lambda(\Omega_0)}{\kappa} \right)^{1/N} \quad \text{and} \quad b_n \geq \frac{b_{n-1} + 1}{\kappa^{1/N}} \quad n = 1, \dots, L-1, \quad (101)$$

where Ω_0 is the domain of the neural network input, $0 \leq \kappa < 1$, and N is the number of neurons in each layer, then both operators $\mathcal{M}(n, n+1)$ and $\mathcal{N}(n+1, n)$ defined in (67) are contractions for all $n = 0, \dots, L-1$, i.e., their norm can be bounded by a constant $K_n \leq \kappa$, independently of the weights of the neural network.

Proof. If ξ_n is uniformly distributed then from (82) we have that

$$K_n = \frac{\lambda(\mathcal{R}(\mathbf{X}_n))}{\lambda(\mathcal{R}(\xi_n))}. \quad (102)$$

By using Lemma 0.4 we can bound K_n as

$$K_n \leq \left(\frac{1 + b_{n-1}}{b_n} \right)^N, \quad (103)$$

where N is the number of neurons in each layer of the neural network. Therefore, if $b_n \geq (b_{n-1} + 1)/\kappa^{1/N}$ ($n = 1, \dots, L-1$) we have that K_n is bounded by a quantity κ smaller than one. Regarding b_0 , we notice that

$$K_0 = \frac{\lambda(\mathcal{R}(\mathbf{X}_0))}{\lambda(\mathcal{R}(\xi_0))} = \frac{\lambda(\Omega_0)}{(2b_0)^N}, \quad (104)$$

where Ω_0 is the domain of the neural network input. Hence, if b_0 satisfies (102) then $K_0 \leq \kappa$. □

¹¹An linear operator is called a contraction if its operator norm is smaller than one.

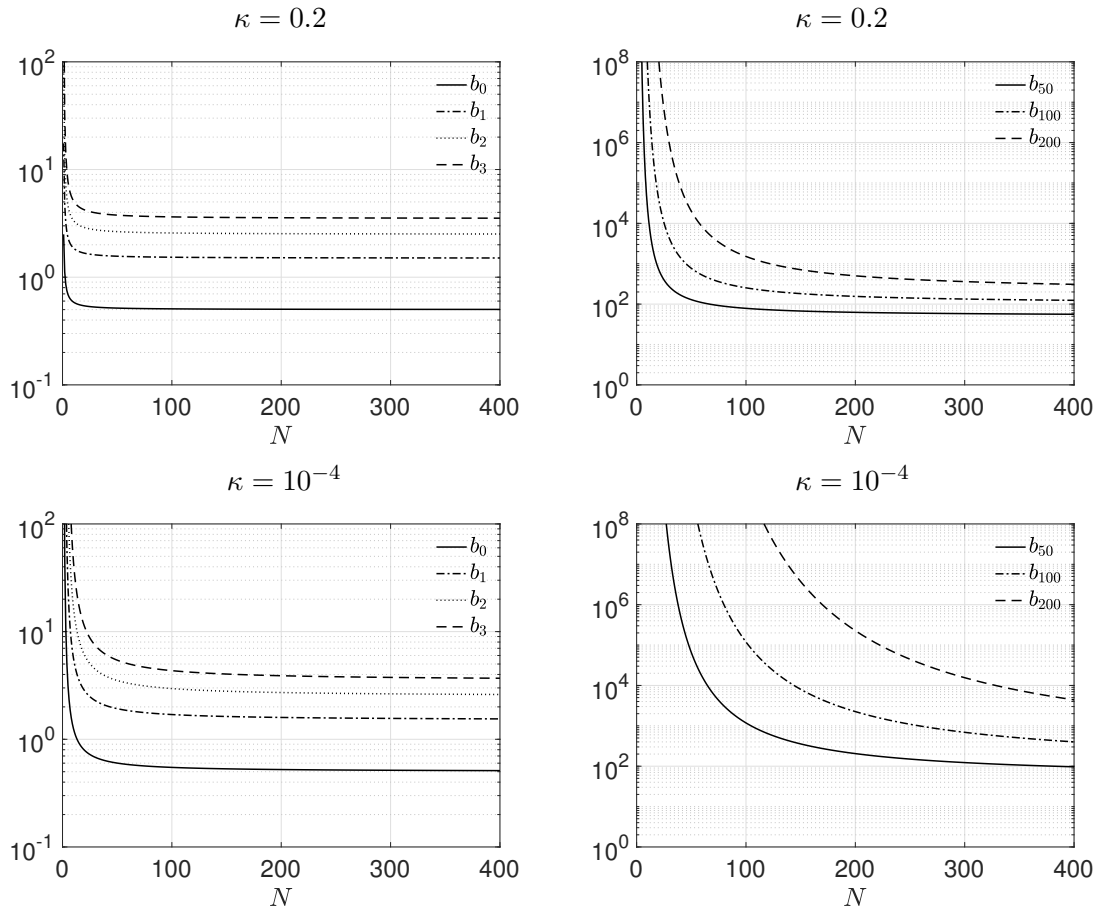


Figure 4: Lower bound on the coefficients b_n defined in (106) for $\lambda(\Omega_0) = 1$ as a function of the number of neurons N and number of layers of the neural network. With such values of b_n the operator $\mathcal{G}(L-n+1, L-n)$ is a contraction satisfying $\|\mathcal{G}(L-n+1, L-n)\|^2 \leq \kappa$. Shown are results for $\kappa = 0.2$ and $\kappa = 10^{-4}$ (contraction index).

One consequence of Proposition 0.13 is that the L^2 norm of the neural network output decays with both the number of layers and the number of neurons if the noise amplitude from one layer to the next increases as in (101). For example, if we represent the input-output map as a sequence of conditional expectations (see (20)), and set $u(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{x}$ (linear output) then we have

$$q_0(\mathbf{x}) = \mathcal{M}(0, 1)\mathcal{M}(1, 2) \cdots \mathcal{M}(L-1, L)(\boldsymbol{\alpha}^T \mathbf{x}). \quad (105)$$

By iterating the inequalities (101) in Proposition 0.13 we find that

$$b_n \geq \frac{1}{2\kappa^{n/N}} \left(\frac{\lambda(\Omega_0)}{\kappa} \right)^{1/N} + \sum_{k=1}^n \frac{1}{\kappa^{k/N}} \quad n = 0, \dots, L-1, \quad (106)$$

In Figure 4 we plot the lower bound at the right hand side of (106) for $\kappa = 0.2$ and $\kappa = 10^{-4}$ as a function of the number of neurons (N). With b_n given in (106) we have that the operator norms of $\mathcal{M}(n, n+1)$ and $\mathcal{N}(n+1, n)$ ($n = 0, \dots, L-1$) are bounded exactly by κ (see Lemma 0.12). Hence, by taking the L^2 norm of (105), and recalling that $\|\mathcal{M}(n, n+1)\|^2 \leq \kappa$ we obtain

$$\|q_0\|_{L^2(\Omega_0)}^2 \leq Z^2 \|\boldsymbol{\alpha}\|_2^2 \kappa^L, \quad (107)$$

where¹²

$$Z^2 = \sum_{k=1}^N \int_{\mathcal{R}(\mathbf{x}_L)} x_k^2 d\mathbf{x} \quad \text{and} \quad \|\boldsymbol{\alpha}\|_2^2 = \sum_{k=1}^N \alpha_k^2. \quad (109)$$

The inequality (107) shows that the 2-norm of the vector of weights $\boldsymbol{\alpha}$ must increase exponentially fast with the number of layers L if we chose the noise amplitude as in (106). As shown in the following Lemma, the growth rate of b_n that guarantees that both \mathcal{M} and \mathcal{N} are contractions is linear (asymptotically with the number of neurons).

Lemma 0.14. *Consider a neural network satisfying the hypotheses of Proposition 0.13. Then, in the limit of an infinite number of neurons ($N \rightarrow \infty$), the noise amplitude (106) satisfies*

$$\lim_{N \rightarrow \infty} b_n = \frac{1}{2} + n, \quad (110)$$

independently of the contraction factor κ and the domain Ω_0 . This means that for a finite number of neurons the noise amplitude b_n that guarantees that $\|\mathcal{M}(n, n+1)\| \leq \kappa$ is bounded from below ($\kappa < 1$) or from above ($\kappa > 1$) by a function that increases linearly with the number of layers.

Proof. The proof follows by taking the limit of (106) for $N \rightarrow \infty$. □

An example: Set $k = 10^{-4}$, $L = 4$ (four layers) and $N = 10$ neurons per layer. The factor κ^L in (107) is 10^{-16} . If we are interested in representing a d -dimensional function in the unit cube $\Omega_0 = [0, 1]^d$ then we have¹³

$$\lambda(\Omega_0) = 1 \quad Z^2 \leq N \frac{2^N (1 + b_3)^{N+2}}{3}, \quad (112)$$

where $b_3 = 3.684$ (see Figure 4 for $N = 10$). Hence, the norm of the output (107) is bounded by

$$\|q_L\|_{L^2(\Omega_0)} \leq C \|\boldsymbol{\alpha}\|_2 \quad C = 10^{-16} \sqrt{10 \frac{2^{10} (1 + b_3)^{12}}{3}} \simeq 6.17 \times 10^{-11}. \quad (113)$$

This means that the 2-norm of the output coefficients $\boldsymbol{\alpha}$ has to be of the order of 10^{11} to represent, e.g., a two-dimensional function of norm about one on the square $\Omega_0 = [0, 1]^2$.

¹²In equation (107) we used the Cauchy-Schwarz inequality

$$\|\boldsymbol{\alpha}^T \mathbf{x}\|_{L^2(\mathcal{R}(\mathbf{x}_L))}^2 \leq Z^2 \|\boldsymbol{\alpha}\|_2^2. \quad (108)$$

¹³In fact,

$$\sum_{i=1}^N \int_{\mathcal{R}(\mathbf{x}_4)} x_i^2 d\mathbf{x} \leq \sum_{i=1}^N \int_{\Omega_4} x_i^2 d\mathbf{x} = \sum_{i=1}^N \underbrace{\int_{-1-b_3}^{1+b_3} \cdots \int_{-1-b_3}^{1+b_3}}_{N \text{ times}} x_i^2 d\mathbf{x} = \frac{2^N N (1 + b_3)^{N+2}}{3}. \quad (111)$$

References

- [1] Z. I. Botev, J. F. Grotowski, and D. P. Kroese. Kernel density estimation via diffusion. *Ann. Stat.*, 38(5):2916–2957, 2010.
- [2] C. Brennan and D. Venturi. Data-driven closures for stochastic dynamical systems. *J. Comp. Phys.*, 372:281–298, 2018.
- [3] J. M. Dominy and D. Venturi. Duality and conditional expectations in the Nakajima-Mori-Zwanzig formulation. *J. Math. Phys.*, 58(8):082701, 2017.
- [4] W. E. A proposal on machine learning via dynamical systems. *Commun. Math. Stat.*, 5:1–10, 2017.
- [5] W. E, J. Han, and Q. Li. A mean-field optimal control formulation of deep learning. *Res. Math. Sci.*, 6:1–41, 2019.
- [6] L. Gonon, L. Grigoryeva, and J.-P. Ortega. Risk bounds for reservoir computing. *JMLR*, 21(240):1–61, 2020.
- [7] N. J. Gordon, D. J. Salmond, and A. F. M. Smith. Novel approach to nonlinear/non-Gaussian Bayesian state estimation. *IEE Proceedings F. Radar and Signal Processing*, 140(2):107–113, 1993.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, pages 770–778, 2016.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, pages 630–645. Springer, 2016.
- [10] I. Kobyzev, S. J. D. Prince, and M. A. Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.
- [11] A. Lasota and M. C. Mackey. *Chaos, fractals and noise: stochastic aspects of dynamics*. Springer–Verlag, second edition, 1994.
- [12] Q. Li, L. Chen, C. Tai, and W. E. Maximum principle based algorithms for deep learning. *JMLR*, 18:1–29, 2018.
- [13] Q. Li, T. Lin, and Z. Shen. Deep learning via dynamical systems: An approximation perspective. *J. Eur. Math. Soc.*, (published online first), 2022.
- [14] Y. Lu, Zhong A, Q. Li, and B. Dong. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations. *arXiv:1710.10121*, 2017.
- [15] E. Minguzzi. The equality of mixed partial derivatives under weak differentiability conditions. *eal Anal. Exch.*, 40(1):81–98, 2014/2015.
- [16] D. Nozaki, D. J. Mar, P. Grigg, and J. J. Collins. Effects of colored noise on stochastic resonance in sensory neurons. *Phys. Rev. Lett.*, 82(11):2402–2405, 1999.
- [17] D. Rezende and S. Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- [18] E. G. Tabak and E. Vanden-Eijnden. Density estimation by dual ascent of the log-likelihood. *Comm. Math. Sci.*, 8(1):217–233, 2010.
- [19] N. Črnjarić-Žic, S. Maćešić, and I. Mezić. Koopman operator spectrum for random dynamical systems. *J. Nonlinear Sci.*, 30:2007–2056, 2020.

- [20] D. Venturi and G. E. Karniadakis. Convolutionless Nakajima-Zwanzig equations for stochastic analysis in nonlinear dynamical systems. *Proc. R. Soc. A*, 470(2166):1–20, 2014.
- [21] T. Yu, Y. Yang, D. Li, T. Hospedales, and T. Xiang. Simple and effective stochastic neural networks. In *Proc. Innov. Appl. Artif. Intell. Conf.*, volume 35, pages 3252–3260, 2021.
- [22] Y. Zhu, J. M. Dominy, and D. Venturi. On the estimation of the Mori-Zwanzig memory integral. *J. Math. Phys.*, 59(10):103501, 2018.
- [23] Y. Zhu and D. Venturi. Hypocoellipticity and the Mori-Zwanzig formulation of stochastic differential equations. *J. Math. Phys.*, 62:1035051, 2021.