Deep Neural Network Structures Solving Variational Inequalities

Patrick L. Combettes & Jean-Christophe Pesquet

Set-Valued and Variational Analysis

Theory and Applications

ISSN 1877-0533

Set-Valued Var. Anal DOI 10.1007/s11228-019-00526-z





Your article is protected by copyright and all rights are held exclusively by Springer Nature B.V.. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your article, please use the accepted manuscript version for posting on your own website. You may further deposit the accepted manuscript version in any repository, provided it is only made publicly available 12 months after official publication or later and provided acknowledgement is given to the original source of publication and a link is inserted to the published article on Springer's website. The link must be accompanied by the following text: "The final publication is available at link.springer.com".



Set-Valued and Variational Analysis https://doi.org/10.1007/s11228-019-00526-z

Deep Neural Network Structures Solving Variational Inequalities



Patrick L. Combettes¹ · Jean-Christophe Pesquet²

Received: 25 April 2019 / Accepted: 9 December 2019 / Published online: 13 February 2020 © Springer Nature B.V. 2020

Abstract

Motivated by structures that appear in deep neural networks, we investigate nonlinear composite models alternating proximity and affine operators defined on different spaces. We first show that a wide range of activation operators used in neural networks are actually proximity operators. We then establish conditions for the averagedness of the proposed composite constructs and investigate their asymptotic properties. It is shown that the limit of the resulting process solves a variational inequality which, in general, does not derive from a minimization problem. The analysis relies on tools from monotone operator theory and sheds some light on a class of neural networks structures with so far elusive asymptotic properties.

Keywords Averaged operator \cdot Deep neural network \cdot Monotone operator \cdot Nonexpansive operator \cdot Proximity operator \cdot Variational inequality

1 Introduction

A powerful tool from fixed point theory to analyze and solve optimization and inclusion problems in a real Hilbert space \mathcal{H} is the class of averaged nonexpansive operators, which was introduced in [3]. Let $T: \mathcal{H} \to \mathcal{H}$ be a nonexpansive operator, i.e., T is 1-Lipschitzian. Then $\alpha \in]0, 1]$ is an averagedness constant of T if $\mathrm{Id} + \alpha^{-1}(T - \mathrm{Id})$ remains nonexpansive, in which case we say that T is α -averaged; if $\alpha = 1/2$, T is firmly nonexpansive. The importance of firmly nonexpansive operators in convex optimization and variational methods has long been recognized [19, 27, 36, 41, 46]. The broader class of averaged operators

The work of P. L. Combettes was supported by the National Science Foundation under grant CCF-1715671. The work of J.-C. Pesquet was supported by Institut Universitaire de France.

Patrick L. Combettes plc@math.ncsu.edu

Jean-Christophe Pesquet jean-christophe@pesquet.eu

- Department of Mathematics, North Carolina State University, Raleigh, NC 27695-8205, USA
- CentraleSupélec, Center for Visual Computing, OPIS Inria Project Team, Université Paris-Saclay, 91190 Gif sur Yvette, France



was shown in [7] to play a prominent role in the analysis of convex feasibility problems. In this context, the underlying problem is to find a common fixed point of averaged operators. In [20], it was shown that many convex minimization and monotone inclusion problems reduce to the more general problem of finding a fixed point of compositions of averaged operators, which provided a unified analysis of various proximal splitting algorithms. Along these lines, several fixed point methods based on various combinations of averaged operators have since been devised, see [1, 2, 5, 9, 11, 13, 14, 17, 18, 24, 25, 38, 43, 47] for recent work. Motivated by deep neural network structures with thus far elusive asymptotic properties, we investigate in the present paper a novel averaged operator model involving a mix of nonlinear and linear operators.

Artificial neural networks have attracted considerable attention as a tool to better understand, model, and imitate the human brain [31, 37, 42]. In a Hilbertian setting [6], an (n + 1)-layer feed-forward neural network architecture acting on real Hilbert spaces $(\mathcal{H}_i)_{0 \le i \le n}$ is defined as the composition of operators $R_n \circ (W_n \cdot + b_n) \circ \cdots \circ R_1 \circ (W_1 \cdot + b_1)$ where, for every $i \in \{1, ..., n\}$, $R_i : \mathcal{H}_i \to \mathcal{H}_i$ is a nonlinear operator known as an activation operator, $W_i: \mathcal{H}_{i-1} \to \mathcal{H}_i$ is a linear operator, known as a weight operator, and $b_i \in \mathcal{H}_i$ is a so-called bias parameter. Deep neural networks feature a (possibly large) number n of layers. In recent years, they have been found to be quite successful in a wide array of classification, recognition, and prediction tasks; see [34] and its bibliography. Despite their success, the operational structure and properties of deep neural networks are not yet well understood from a mathematical viewpoint. In the present paper, we propose to analyze them within the following iterative model. We emphasize that our purpose is not to study the training of the network, which consists of optimally setting the weight operators and bias parameters from data samples, but to analyze mathematically such a structure once it is trained. Our model is also of general interest in constructive fixed point theory for monotone inclusion problems.

Model 1.1 Let $m \ge 1$ be an integer, let \mathcal{H} and $(\mathcal{H}_i)_{0 \le i \le m}$ be nonzero real Hilbert spaces, such that $\mathcal{H}_m = \mathcal{H}_0 = \mathcal{H}$. For every $i \in \{1, \ldots, m\}$ and every $n \in \mathbb{N}$, let $W_{i,n} : \mathcal{H}_{i-1} \to \mathcal{H}_i$ be a bounded linear operator, let $b_{i,n} \in \mathcal{H}_i$, and let $R_{i,n} : \mathcal{H}_i \to \mathcal{H}_i$. Let $x_0 \in \mathcal{H}$, let $(\lambda_n)_{n \in \mathbb{N}}$ be a sequence in $]0, +\infty[$, set

$$(\forall n \in \mathbb{N})(\forall i \in \{1, \dots, m\}) \quad T_{i,n} \colon \mathcal{H}_{i-1} \to \mathcal{H}_i \colon x \mapsto R_{i,n}(W_{i,n}x + b_{i,n}), \tag{1.1}$$

and iterate

for
$$n = 0, 1, ...$$

$$\begin{vmatrix} x_{1,n} = T_{1,n}x_n \\ x_{2,n} = T_{2,n}x_{1,n} \\ \vdots \\ x_{m,n} = T_{m,n}x_{m-1,n} \\ x_{n+1} = x_n + \lambda_n(x_{m,n} - x_n). \end{vmatrix}$$
(1.2)

In sharp contrast with existing algorithmic frameworks involving averaged operators (see cited works above), the operators involved in Model 1.1 are not necessarily all defined on the same Hilbert space and, in addition, they need not all be averaged. Let us also note that the relaxation parameters $(\lambda_n)_{n\in\mathbb{N}}$ in (1.2) allow us to model skip connections [44], in



the spirit of residual networks [33]. If $\lambda_n \equiv 1$, we obtain the standard feed-forward architecture [31].

Our contributions are articulated around the following findings.

- We show that a wide range of activation operators used in neural networks are actually
 proximity operators, which paves the way to the analysis of such networks via fixed
 point theory for monotone inclusion problems.
- We provide a new analysis of compositions of proximity and affine operators, establishing mild conditions that guarantee that the resulting operator is averaged.
- We show that, under suitable assumptions, the asymptotic output of the network converges to a point defined via a variational inequality. Furthermore, in general, this variational inequality does not derive from a minimization problem.

The remainder of the paper is organized as follows. In Section 2, we bring to light strong connections between the activation functions employed in neural networks and the theory of proximity operators in convex analysis. In Section 3, we derive new results on the averagedness properties of compositions of proximity and affine operators acting on different spaces. In Section 4, we investigate the asymptotic behavior of a class of deep neural networks structures and show that their fixed points solve a variational inequality. The main assumption on this subclass of Model 1.1 is that the structure of the network is periodic in the sense that a group of layers is repeated. Finally, in Section 5, the same properties are established for a broader class of networks.

Notation We follow standard notation from convex analysis and operator theory [8, 40]. Thus, the expressions $x_n \to x$ and $x_n \to x$ denote, respectively, weak and strong convergence of a sequence $(x_n)_{n\in\mathbb{N}}$ to x in \mathcal{H} , and $\Gamma_0(\mathcal{H})$ is the class of lower semicontinuous convex functions $\varphi \colon \mathcal{H} \to]-\infty, +\infty]$ such that dom $\varphi = \{x \in \mathcal{H} \mid \varphi(x) < +\infty\} \neq \varnothing$. Now let $\varphi \in \Gamma_0(\mathcal{H})$. The conjugate of φ is denoted by φ^* , its subdifferential by $\partial \varphi$, and its proximity operator is $\operatorname{prox}_{\varphi} \colon \mathcal{H} \to \mathcal{H} \colon x \mapsto \operatorname{argmin}_{y \in \mathcal{H}} (\varphi(y) + \|x - y\|^2/2)$. The symbols $\operatorname{ran} T$, dom T, Fix T, and zer T denote respectively the range, the domain, the fixed point set, and the set of zeros of an operator T. The space of bounded linear operators from a Banach space \mathcal{X} to a Banach space \mathcal{Y} is denoted by $\mathcal{B}(\mathcal{X}, \mathcal{Y})$. Finally, ℓ_+^1 denotes the set of summable sequences in $[0, +\infty[$.

2 Proximal Activation in Neural Networks

The following facts will be needed.

Lemma 2.1 Let $\varphi \in \Gamma_0(\mathcal{H})$. Then the following hold:

- (i) [8, Proposition 12.29] Fix $prox_{\varphi} = Argmin \varphi$.
- (ii) [8, Corollary 24.5] Let $g \in \Gamma_0(\mathcal{H})$ be such that $\varphi = g \|\cdot\|^2/2$. Then $\operatorname{prox}_{\varphi} = \nabla g^*$.

2.1 Activation Functions

An activation function is a function $\varrho \colon \mathbb{R} \to \mathbb{R}$ which models the firing activity of neurons. The simplest instance, that goes back to the perceptron machine [42], is that of a binary



firing model: the neuron is either firing or at rest. For instance, if the firing level is 1 and the rest state is 0, we obtain the binary step function

$$\varrho \colon \xi \mapsto \begin{cases} 1, & \text{if } \xi > 0; \\ 0, & \text{if } \xi \leqslant 0, \end{cases}$$
 (2.1)

which was initially proposed in [37]. As this discontinuous activation model may lead to unstable neural networks, various continuous approximations have been proposed. Our key observation is that a vast array of activation functions used in neural networks belong to the following class.

Definition 2.2 The set of functions from \mathbb{R} to \mathbb{R} which are increasing, 1-Lipschitzian, and take value 0 at 0 is denoted by $\mathcal{A}(\mathbb{R})$.

Remarkably, we can precisely characterize this class of activation functions as that of proximity operators.

Proposition 2.3 Let $\varrho \colon \mathbb{R} \to \mathbb{R}$. Then $\varrho \in \mathcal{A}(\mathbb{R})$ if and only if there exists a function $\phi \in \Gamma_0(\mathbb{R})$, which has 0 as a minimizer, such that $\varrho = \operatorname{prox}_{\phi}$.

Proof The fact that the class of increasing, 1-Lipschitzian functions from \mathbb{R} to \mathbb{R} coincides with that of proximity operators of functions in $\Gamma_0(\mathbb{R})$ is shown in [22, Proposition 2.4]. In view of Lemma 2.1(i) and Definition 2.2, the proof is complete.

To illustrate the above results, let us provide examples of common activation functions $\varrho \in \mathcal{A}(\mathbb{R})$, and identify the potential φ they derive from in Proposition 2.3 (see Fig. 1).

Example 2.4 The most basic activation function is $\varrho = \text{Id} = \text{prox}_0$. It is in particular useful in dictionary learning approaches, which correspond to the linear special case of Model 1.1 [45].

Example 2.5 The saturated linear activation function [31]

$$\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \begin{cases} 1, & \text{if } \xi > 1; \\ \xi, & \text{if } -1 \leqslant \xi \leqslant 1; \\ -1, & \text{if } \xi < -1 \end{cases}$$
 (2.2)

can be written as $\varrho = \text{prox}_{\phi}$, where ϕ is the indicator function of [-1, 1].

Example 2.6 The rectified linear unit (ReLU) activation function [39]

$$\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \begin{cases} \xi, & \text{if } \xi > 0; \\ 0, & \text{if } \xi \leqslant 0 \end{cases}$$
 (2.3)

can be written as $\varrho = \text{prox}_{\phi}$, where ϕ is the indicator function of $[0, +\infty[$.

Example 2.7 Let $\alpha \in [0, 1]$. The parametric rectified linear unit activation function [32] is

$$\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \begin{cases} \xi, & \text{if } \xi > 0; \\ \alpha \xi, & \text{if } \xi \leqslant 0. \end{cases}$$
 (2.4)

We have $\varrho = \operatorname{prox}_{\phi}$, where

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] \colon \xi \mapsto \begin{cases} 0, & \text{if } \xi > 0; \\ (1/\alpha - 1)\xi^2/2, & \text{if } \xi \leqslant 0. \end{cases}$$
 (2.5)



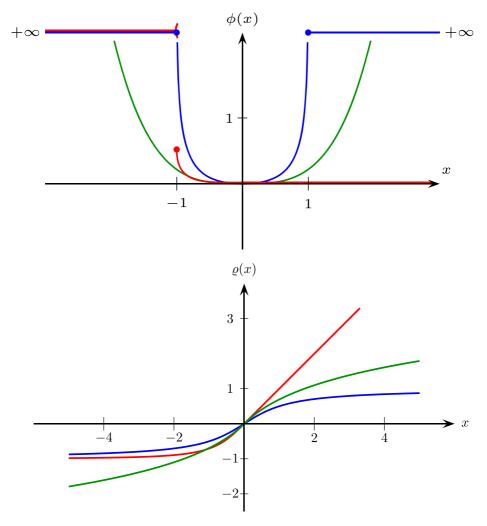


Fig. 1 The function ϕ (top) and the corresponding proximal activation function (bottom) ϱ in Proposition 2.3. Example 2.10 is in red, Example 2.11 is in blue, Example 2.17 is in green

Proof Let $\xi \in \mathbb{R}$. Then $\phi'(\xi) = 0$ if $\xi > 0$, and $\phi'(\xi) = (1/\alpha - 1)\xi$ if $\xi \leq 0$. In turn $(\mathrm{Id} + \phi')\xi = \xi$ if $\xi > 0$, and $(\mathrm{Id} + \phi')(\xi) = \xi/\alpha$ if $\xi \leq 0$. Hence, $\varrho = (\mathrm{Id} + \phi')^{-1}$ is given by (2.4).

Example 2.8 The bent identity activation function $\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto (\xi + \sqrt{\xi^2 + 1} - 1)/2$ satisfies $\varrho = \operatorname{prox}_{\phi}$, where

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] \colon \xi \mapsto \begin{cases} \xi/2 - \left(\ln(\xi + 1/2)\right)/4, & \text{if } \xi > -1/2; \\ +\infty, & \text{if } \xi \leqslant -1/2. \end{cases}$$
 (2.6)

Proof This follows from [23, Lemma 2.6 and Example 2.18].



Example 2.9 The inverse square root unit activation function [16] is $\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \xi/\sqrt{1+\xi^2}$. We have $\varrho = \operatorname{prox}_{\phi}$, where

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] \colon \xi \mapsto \begin{cases} -\xi^2/2 - \sqrt{1-\xi^2}, & \text{if } |\xi| \leqslant 1; \\ +\infty, & \text{if } |\xi| > 1. \end{cases}$$
 (2.7)

Proof Let $\xi \in]-1$, $1[=\dim \nabla \phi = \dim \partial \phi = \operatorname{ran}\operatorname{prox}_{\phi}$. Then $\xi + \phi'(\xi) = \xi/\sqrt{1-\xi^2}$ and therefore $\operatorname{prox}_{\phi} = (\operatorname{Id} + \phi')^{-1} \colon \mu \mapsto \mu/\sqrt{1+\mu^2}$.

Example 2.10 The inverse square root linear unit activation function [16]

$$\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \begin{cases} \xi, & \text{if } \xi \geqslant 0; \\ \frac{\xi}{\sqrt{1+\xi^2}}, & \text{if } \xi < 0 \end{cases}$$
 (2.8)

can be written as $\varrho = \operatorname{prox}_{\phi}$, where

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] : \xi \mapsto \begin{cases} 0, & \text{if } \xi \geqslant 0; \\ 1 - \xi^2/2 - \sqrt{1 - \xi^2}, & \text{if } -1 \leqslant \xi < 0; \\ +\infty, & \text{if } \xi < -1. \end{cases}$$
 (2.9)

Proof Let
$$\xi \in]-1, +\infty[=\operatorname{dom} \nabla \phi = \operatorname{ran}\operatorname{prox}_{\phi}$$
. Then $\xi + \phi'(\xi) = \xi$ if $\xi \geqslant 0$, and $\xi + \phi'(\xi) = \xi/\sqrt{1 - \xi^2}$ if $\xi < 0$. Hence, $\varrho = (\operatorname{Id} + \phi')^{-1}$ is given by (2.8).

Example 2.11 The arctangent activation function $(2/\pi)$ arctan is the proximity operator of

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] \colon \xi \mapsto \begin{cases} -\frac{2}{\pi} \ln\left(\cos\left(\frac{\pi\xi}{2}\right)\right) - \frac{1}{2}\xi^2, & \text{if } |\xi| < 1; \\ +\infty, & \text{if } |\xi| \geqslant 1. \end{cases}$$
 (2.10)

Proof Let
$$\xi \in]-1$$
, $1[=\operatorname{dom} \nabla \phi = \operatorname{ran} \operatorname{prox}_{\phi}$. Then $\xi + \phi'(\xi) = \tan(\pi \xi/2)$ and therefore $\varrho = (\operatorname{Id} + \phi')^{-1} = (2/\pi) \operatorname{arctan}$.

Example 2.12 The hyperbolic tangent activation function tanh [35] is the proximity operator of

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] \colon \xi \mapsto \begin{cases} \frac{(1+\xi)\ln(1+\xi) + (1-\xi)\ln(1-\xi) - \xi^{2}}{2} & \text{if } |\xi| < 1; \\ \ln(2) - 1/2 & \text{if } |\xi| = 1; \\ +\infty, & \text{if } |\xi| > 1. \end{cases}$$
(2.11)

Proof Let $\xi \in]-1$, $1[=\operatorname{dom} \nabla \phi = \operatorname{ran} \operatorname{prox}_{\phi}$. Then $\xi + \phi'(\xi) = \operatorname{arctanh}(\xi)$ and therefore $\varrho = (\operatorname{Id} + \phi')^{-1} = \operatorname{tanh}$.



Example 2.13 The unimodal sigmoid activation function [30]

$$\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \frac{1}{1 + e^{-\xi}} - \frac{1}{2} \tag{2.12}$$

is the proximity operator of

 $\phi \colon \mathbb{R} \to]-\infty, +\infty]$

$$\xi \mapsto \begin{cases} (\xi+1/2) \ln(\xi+1/2) + (1/2-\xi) \ln(1/2-\xi) - \frac{1}{2}(\xi^2+1/4) & \text{if } |\xi| < 1/2; \\ -1/4, & \text{if } |\xi| = 1/2; \\ +\infty, & \text{if } |\xi| > 1/2. \end{cases}$$

$$(2.13)$$

 $\begin{array}{l} \textit{Proof} \ \ \text{Let} \ \xi \in \]-1/2, \ 1/2[\ = \ \text{dom} \ \nabla \phi = \ \text{ran prox}_{\phi}. \ \text{Then} \ \xi + \phi'(\xi) = \ln((1+2\xi)/(1-2\xi)) \\ \text{and therefore} \ \text{prox}_{\phi} = (\text{Id} + \phi')^{-1} \colon \mu \mapsto (1/2)(e^{\mu} - 1)/(e^{\mu} + 1) = 1/(1+e^{-\mu}) - 1/2. \end{array} \ \ \Box$

Remark 2.14 Examples 2.12 and 2.13 are closely related in the sense that the function of (2.12) can be written as $\varrho = (1/2) \tanh(\cdot/2)$.

Example 2.15 The Elliot activation function is [28] $\varrho: \mathbb{R} \to \mathbb{R}: \xi \mapsto \xi/(1+|\xi|)$ can be written as $\varrho = \operatorname{prox}_{\phi}$, where

$$\phi \colon \mathbb{R} \to]-\infty, +\infty]
\xi \mapsto \begin{cases}
-|\xi| - \ln(1 - |\xi|) - \frac{\xi^2}{2}, & \text{if } |\xi| < 1; \\
+\infty, & \text{if } |\xi| \geqslant 1.
\end{cases}$$
(2.14)

Proof Let $\xi \in]-1$, $1[=\dim \nabla \phi = \operatorname{ran}\operatorname{prox}_{\phi}$. Then $\xi + \phi'(\xi) = \xi/(1-|\xi|)$ and therefore $\operatorname{prox}_{\phi} = (\operatorname{Id} + \phi')^{-1} : \mu \mapsto \mu/(1+|\mu|)$.

Example 2.16 The inverse hyperbolic sine activation function arcsinh is the proximity operator of $\phi = \cosh - |\cdot|^2/2$.

Proof Let $\xi \in \mathbb{R}$. Then $\xi + \phi'(\xi) = \sinh \xi$ and therefore $\operatorname{prox}_{\phi} = (\operatorname{Id} + \phi')^{-1} = \operatorname{arcsinh}$.

Example 2.17 The logarithmic activation function [10]

$$\rho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \operatorname{sign}(\xi) \ln \left(1 + |\xi| \right) \tag{2.15}$$

is the proximity operator of

$$\phi \colon \mathbb{R} \to]-\infty, +\infty] : \xi \mapsto e^{|\xi|} - |\xi| - 1 - \frac{\xi^2}{2}.$$
 (2.16)

Proof We have $\phi' : \xi \mapsto \operatorname{sign}(\xi)(e^{|\xi|} - 1) - \xi$. Hence $(\operatorname{Id} + \phi') : \xi \mapsto \operatorname{sign}(\xi)(e^{|\xi|} - 1)$ and, in turn, $\operatorname{prox}_{\phi} = (\operatorname{Id} + \phi')^{-1} : \xi \mapsto \operatorname{sign}(\xi) \ln(1 + |\xi|)$.



The class of activation functions $\mathcal{A}(\mathbb{R})$ has interesting stability properties.

Proposition 2.18 *The following hold:*

- (i) Let $\alpha \in]0, +\infty[$ and $\beta \in]0, +\infty[$ be such that $\alpha\beta \leqslant 1$, and let $\varrho \in \mathcal{A}(\mathbb{R})$. Then $\alpha\varrho(\beta \cdot) \in \mathcal{A}(\mathbb{R})$.
- (ii) Let $(\varrho_i)_{i\in I}$ be a finite family in $\mathcal{A}(\mathbb{R})$ and let $(\omega_i)_{i\in I}$ be real numbers in]0,1] such that $\sum_{i\in I} \omega_i = 1$. Then $\sum_{i\in I} \omega_i \varrho_i \in \mathcal{A}(\mathbb{R})$.
- (iii) Let $\varrho_1 \in \mathcal{A}(\mathbb{R})$ and $\varrho_2 \in \mathcal{A}(\mathbb{R})$. Then $\varrho_1 \circ \varrho_2 \in \mathcal{A}(\mathbb{R})$.
- (iv) Let $\varrho \in \mathcal{A}(\mathbb{R})$. Then $\mathrm{Id} \varrho \in \mathcal{A}(\mathbb{R})$.
- (v) Let $\varrho_1 \in \mathcal{A}(\mathbb{R})$ and $\varrho_2 \in \mathcal{A}(\mathbb{R})$. Then $(\varrho_1 \varrho_2 + \mathrm{Id})/2 \in \mathcal{A}(\mathbb{R})$.
- (vi) Let $\varrho_1 \in \mathcal{A}(\mathbb{R})$ and $\varrho_2 \in \mathcal{A}(\mathbb{R})$. Then $\varrho_1 \circ (2\varrho_2 \mathrm{Id}) + \mathrm{Id} \varrho_2 \in \mathcal{A}(\mathbb{R})$.

Proof (i)–(iii): This follows at once from Definition 2.2.

(iv)–(v): The fact that the resulting operators are proximity operators is established in [21, Section 3.3]. The fact that they are proximity operators of a function $\phi \in \Gamma_0(\mathcal{H})$ that is minimal at 0 is equivalent to the fact that $\operatorname{prox}_{\phi} 0 = 0$ Lemma 2.1(i). This identity is easily seen to hold in each instance.

(vi): Set $\varrho = \varrho_1 \circ (2\varrho_2 - \mathrm{Id}) + \mathrm{Id} - \varrho_2$. Then ϱ is firmly nonexpansive [8, Proposition 4.31(ii)]. It is therefore increasing and nonexpansive. Finally, $\varrho(0) = 0$.

Remark 2.19 Using Proposition 2.18, the above examples can be combined to obtain additional activation functions. For instance, it follows from Example 2.5 and Proposition 2.18(iv) that the soft thresholder

$$\varrho \colon \mathbb{R} \to \mathbb{R} \colon \xi \mapsto \begin{cases} \xi - 1, & \text{if } \xi > 1; \\ 0, & \text{if } -1 \leqslant \xi \leqslant 1; \\ \xi + 1, & \text{if } \xi < -1 \end{cases}$$
 (2.17)

belongs to $\mathcal{A}(\mathbb{R})$. It was proposed as an activation function in [48].

2.2 Activation Operators

In Section 2.1, we have described activation functions which model neuronal activity in terms of a scalar function. In this section, we extend this notion to more general activation operators.

Definition 2.20 Let \mathcal{H} be a real Hilbert space and let $R: \mathcal{H} \to \mathcal{H}$. Then R belongs to the class $\mathcal{A}(\mathcal{H})$ if there exists a function $\varphi \in \Gamma_0(\mathcal{H})$ which is minimal at the zero vector and such that $R = \operatorname{prox}_{\varphi}$.

Property (ii) below shows that activation operators in $\mathcal{A}(\mathcal{H})$ have strong stability properties. On the other hand, the boundedness property (iv) is important in neural network-based functional approximation [26, 29].

Proposition 2.21 *Let* \mathcal{H} *be a real Hilbert space and let* $R \in \mathcal{A}(\mathcal{H})$ *. Then the following hold:*

- (i) R0 = 0.
- (ii) Let x and y be in H. Then $||Rx Ry||^2 \le ||x y||^2 ||x y Rx + Ry||^2$.
- (iii) Let $x \in \mathcal{H}$. Then $||Rx|| \leq ||x||$.



(iv) Let $\varphi \in \Gamma_0(\mathcal{H})$ be such that $R = \operatorname{prox}_{\varphi}$. Then ran R is bounded if and only if dom φ is bounded.

Proof (i): This follows from Lemma 2.1(i).

- (ii): This follows from the firm nonexpansiveness of proximity operators [8, Proposition 12.28].
 - (iii): Set y = 0 in (ii) and use (i).
- (iv): We have ran $R = \text{ran} (\text{Id} + \partial \varphi)^{-1} = \text{dom} (\text{Id} + \partial \varphi) = \text{dom} \partial \varphi$. On the other hand, dom $\partial \varphi$ is a dense subset of dom φ [8, Corollary 16.39].

Proposition 2.22 *Let* \mathcal{H} *and* \mathcal{G} *be real Hilbert spaces. Then the following hold:*

- (i) Let $L \in \mathcal{B}(\mathcal{H}, \mathcal{G})$ be such that $||L|| \le 1$ and let $R \in \mathcal{A}(\mathcal{H})$. Then $L^* \circ R \circ L \in \mathcal{A}(\mathcal{H})$.
- (ii) Let $(R_i)_{i\in I}$ be a finite family in $\mathcal{A}(\mathcal{H})$ and let $(\omega_i)_{i\in I}$ be real numbers in]0,1] such that $\sum_{i\in I} \omega_i = 1$. Then $\sum_{i\in I} \omega_i R_i \in \mathcal{A}(\mathcal{H})$.
- (iii) Let $R \in \mathcal{A}(\mathcal{H})$. Then $\mathrm{Id} R \in \mathcal{A}(\mathcal{H})$.
- (iv) Let $R_1 \in \mathcal{A}(\mathcal{H})$ and $R_2 \in \mathcal{A}(\mathcal{H})$. Then $(R_1 R_2 + \mathrm{Id})/2 \in \mathcal{A}(\mathcal{H})$.

Proof The fact that the resulting operators are proximity operators is established in [21, Section 3.3]. In addition, 0 is clearly a fixed point of the resulting operators. In view of Lemma 2.1(i), the proof is complete.

Example 2.23 The softmax activation operator [15] is

$$R: \mathbb{R}^N \to \mathbb{R}^N: (\xi_k)_{1 \leqslant k \leqslant N} \mapsto \left(\exp(\xi_k) / \sum_{j=1}^N \exp(\xi_j) \right)_{1 < k < N} - u, \qquad (2.18)$$

where $u = (1, ..., 1)/N \in \mathbb{R}^N$. We have $R = \text{prox}_{\varphi}$, where $\varphi = \psi(\cdot + u) + \langle \cdot \mid u \rangle$ and

$$\psi \colon \mathbb{R}^{N} \to]-\infty, +\infty]$$

$$(\xi_{k})_{1 \leqslant k \leqslant N} \mapsto \begin{cases} \sum_{k=1}^{N} \left(\xi_{k} \ln \xi_{k} - \frac{\xi_{k}^{2}}{2} \right), & \text{if } (\xi_{k})_{1 \leqslant i \leqslant N} \in [0, 1]^{N} \text{ and } \sum_{k=1}^{N} \xi_{k} = 1; \\ +\infty, & \text{otherwise,} \end{cases}$$

with the convention $0 \ln 0 = 0$.

Proof Set

$$g: \mathbb{R}^N \to]-\infty, +\infty]$$

$$(\xi_k)_{1 \leqslant k \leqslant N} \mapsto \begin{cases} \sum_{k=1}^N \xi_k \ln \xi_k, & \text{if } (\xi_k)_{1 \leqslant k \leqslant N} \in [0, 1]^N \text{ and } \sum_{k=1}^N \xi_k = 1; \\ +\infty, & \text{otherwise.} \end{cases}$$

$$(2.20)$$

Then $\psi = g - \|\cdot\|^2/2$ and [40, Section 16] asserts that

$$g^* \colon \mathbb{R}^N \to \mathbb{R} \colon (\xi_k)_{1 \leqslant k \leqslant N} \mapsto \ln \left(\sum_{k=1}^N \exp(\xi_k) \right).$$
 (2.21)

Since $\nabla g^* = R + u$, according to Lemma 2.1(ii), $R = \text{prox}_{\psi} - u$. We complete the proof by invoking the shift properties of proximity operators [8, Proposition 24.8(iii)].



Separable activation operators supply another important instance of activation operators.

Proposition 2.24 *Let* \mathcal{H} *be a separable real Hilbert space, let* $(e_k)_{k \in \mathbb{K} \subset \mathbb{N}}$ *be an orthonormal basis of* \mathcal{H} *, and let* $(\phi_k)_{k \in \mathbb{K}}$ *be a family of functions in* $\Gamma_0(\mathbb{R})$ *such that* $(\forall k \in \mathbb{K})$ $\phi_k \geqslant \phi_k(0) = 0$. *Define*

$$R: \mathcal{H} \to \mathcal{H}: x \mapsto \sum_{k \in \mathbb{K}} (\operatorname{prox}_{\phi_k} \langle x \mid e_k \rangle) e_k.$$
 (2.22)

Then $R \in \mathcal{A}(\mathcal{H})$.

Proof The fact that R is the proximity operator of the $\Gamma_0(\mathcal{H})$ function $\varphi \colon x \mapsto \sum_{k \in \mathbb{K}} \phi_k(\langle x \mid e_k \rangle)$ is established in [23, Example 2.19]. In addition, it is clear that φ is minimal at 0.

3 Compositions of Firmly Nonexpansive and Affine Operators

Our analysis will revolve around the following property for a family of linear operators $(W_i)_{1 \le i \le m+1}$.

Condition 3.1 Let $m \ge 0$ be an integer, let $(\mathcal{H}_i)_{0 \le i \le m}$ be real Hilbert spaces, set $\mathcal{H}_{m+1} = \mathcal{H}_0$, and let $\alpha \in [1/2, 1]$. For every $i \in \{1, ..., m+1\}$, let $W_i \in \mathcal{B}(\mathcal{H}_{i-1}, \mathcal{H}_i)$ and set

$$L_i: \mathcal{H}_0 \times \dots \times \mathcal{H}_{i-1} \to \mathcal{H}_i: (x_k)_{0 \leqslant k \leqslant i-1} \mapsto \sum_{k=0}^{i-1} (W_i \circ \dots \circ W_{k+1}) x_k. \tag{3.1}$$

It is required that, for every $\mathbf{x} = (x_i)_{0 \le i \le m} \in \mathcal{H}_0 \times \cdots \times \mathcal{H}_m$ such that

$$(\forall i \in \{0, \dots, m\}) \quad \|x_i\| \leqslant \begin{cases} 1, & \text{if } i = 0; \\ \|L_i(x_0, \dots, x_{i-1})\|, & \text{if } i \geqslant 1, \end{cases}$$
 (3.2)

there holds

$$||L_{m+1}\mathbf{x} - 2^{m+1}(1-\alpha)x_0|| + ||L_{m+1}\mathbf{x}|| \le 2^{m+1}\alpha||x_0||.$$
(3.3)

Remark 3.2 In Condition 3.1, we take $\alpha \ge 1/2$ because, if $\mathbf{x} = (x_i)_{0 \le i \le m} \in (\mathcal{H}_0 \setminus \{0\}) \times \mathcal{H}_1 \times \cdots \times \mathcal{H}_m$ satisfies (3.3), then $2^{m+1}(1-\alpha)\|x_0\| \le \|L_{m+1}\mathbf{x} - 2^{m+1}(1-\alpha)x_0\| + \|L_{m+1}\mathbf{x}\| \le 2^{m+1}\alpha\|x_0\|$.

We establish some preliminary results before providing properties that imply Condition 3.1.

Lemma 3.3 Let $m \ge 1$ be an integer, let $(\mathcal{H}_i)_{0 \le i \le m}$ be real Hilbert spaces, and set $\theta_0 = 1$. For every $i \in \{1, ..., m\}$, let $W_i \in \mathcal{B}(\mathcal{H}_{i-1}, \mathcal{H}_i)$ and set

$$\theta_i = ||W_i \circ \cdots \circ W_1||$$

$$+\sum_{k=1}^{i-1}\sum_{1\leqslant j_{1}<...< j_{k}\leqslant i-1}\|W_{i}\circ\cdots\circ W_{j_{k}+1}\|\|W_{j_{k}}\circ\cdots\circ W_{j_{k-1}+1}\|\cdots\|W_{j_{1}}\circ\cdots\circ W_{1}\|.$$
(3.4)



Let $(x_i)_{0 \le i \le m} \in \mathcal{H}_0 \times \cdots \times \mathcal{H}_m$ be such that (3.2) is satisfied. Then the following hold:

(i)
$$(\forall i \in \{1, ..., m\}) \theta_i = \sum_{k=0}^{i-1} \theta_k || W_i \circ \cdots \circ W_{k+1} ||$$
.

(ii) $(\forall i \in \{1, ..., m\}) \|x_i\| \leq \theta_i \|x_0\|.$

Proof (i): This follows recursively from (3.4).

(ii): For every $i \in \{1, \ldots, m\}$, let L_i be as in (3.1). We proceed by induction on m. We first observe that the inequality is satisfied if m = 1 since $||x_1|| \le ||L_1x_0|| = ||W_1x_0|| \le ||W_1|| ||x_0|| = \theta_1 ||x_0||$. Now assume that $m \ge 2$ and that the inequalities hold for (x_1, \ldots, x_{m-1}) . Then, since (i) yields

$$\theta_m = \|W_m \circ \dots \circ W_1\| + \sum_{k=1}^{m-1} \theta_k \|W_m \circ \dots \circ W_{k+1}\|, \tag{3.5}$$

we obtain

$$||x_{m}|| \leq ||L_{m}(x_{0}, \dots, x_{m-1})|| = \left\| \sum_{k=0}^{m-1} (W_{m} \circ \dots \circ W_{k+1}) x_{k} \right\|$$

$$\leq \sum_{k=0}^{m-1} ||W_{m} \circ \dots \circ W_{k+1}|| \, ||x_{k}||$$

$$\leq \left(||W_{m} \circ \dots \circ W_{1}|| + \sum_{k=1}^{m-1} \theta_{k} ||W_{m} \circ \dots \circ W_{k+1}|| \right) ||x_{0}||$$

$$= \theta_{m} ||x_{0}||, \tag{3.6}$$

which concludes the proof.

Lemma 3.4 Let \mathcal{H} be a real Hilbert space, and let x and y be in \mathcal{H} . Then

$$||x|| ||y|| - \langle x | y \rangle \le (||x|| + ||y|| - ||x + y||)(||x|| + ||y||). \tag{3.7}$$

Proof Since $||x + y||^2 - 2||x + y||(||x|| + ||y||) + (||x|| + ||y||)^2 \ge 0$, we have $||x||^2 + ||y||^2 + \langle x | y \rangle + ||x|| ||y||$ $= ||x||^2 + ||y||^2 + \frac{||x + y||^2 - ||x||^2 - ||y||^2}{2} + \frac{(||x|| + ||y||)^2 - ||x||^2 - ||y||^2}{2}$ $= \frac{||x + y||^2 + (||x|| + ||y||)^2}{2}$ $\ge ||x + y||(||x|| + ||y||), \tag{3.8}$

as claimed.

Notation 3.5 Let $m \ge 0$ be an integer, and let $(\mathcal{H}_i)_{0 \le i \le m}$ be real Hilbert spaces. Let \mathcal{X} be the standard vector space $\mathcal{H}_0 \times \cdots \times \mathcal{H}_m$ equipped with the norm $\|\cdot\|_{\mathcal{X}} : \mathbf{x} = (x_i)_{0 \le i \le m} \mapsto \max_{0 \le i \le m} \|x_i\|$ and let \mathcal{Y} be the standard vector space $\mathcal{H}_0 \times \mathcal{H}_0$ equipped with the norm $\|\cdot\|_{\mathcal{Y}} : \mathbf{y} = (y_1, y_2) \mapsto \|y_1\| + \|y_2\|$. Henceforth, the norm of $\mathbf{M} \in \mathcal{B}(\mathcal{X}, \mathcal{Y})$ is denoted by $\|\mathbf{M}\|_{\mathcal{X}, \mathcal{Y}}$.

Proposition 3.6 Let $m \ge 0$ be an integer, let $(\mathcal{H}_i)_{0 \le i \le m}$ be nonzero real Hilbert spaces, set $\mathcal{H}_{m+1} = \mathcal{H}_0$, and use Notation 3.5. For every $i \in \{1, ..., m+1\}$, let $W_i \in \mathcal{B}(\mathcal{H}_{i-1}, \mathcal{H}_i)$.



Further, let $\alpha \in [1/2, 1]$, let $\theta_0 = 1$, let $(\theta_i)_{1 \le i \le m+1}$ be as in (3.4), and set

$$W = W_{m+1} \circ \cdots \circ W_1 \tag{3.9a}$$

$$\mu = \inf_{x \in \mathcal{H}_0, \|x\| = 1} \langle Wx \mid x \rangle \tag{3.9b}$$

$$M: \mathcal{X} \to \mathcal{Y}: \mathbf{x} \mapsto \frac{1}{2^{m+1}\alpha} (M\mathbf{x} - 2^{m+1}(1-\alpha)x_0, M\mathbf{x}).$$
 (3.9d)

Suppose that one of the following holds:

- There exists $i \in \{1, ..., m+1\}$ such that $W_i = 0$. (i)
- (ii) $\|\boldsymbol{M}\|_{\boldsymbol{\mathcal{X}},\boldsymbol{\mathcal{Y}}} \leqslant 1.$
- (ii) $\|W\|_{\mathcal{X}, \mathcal{Y}} \le 1$. (iii) $\|W 2^{m+1}(1-\alpha)\operatorname{Id}\| \|W\| + 2\theta_{m+1} \le 2^{m+1}\alpha$.
- $\alpha \neq 1$, for every $i \in \{1, \dots, m+1\}$ $W_i \neq 0$, and there exists $\eta \in [0, \alpha/((1-\alpha)\theta_{m+1})]$

$$\begin{cases} \theta_{m+1} \leqslant 2^{m+1} \alpha \\ \alpha \theta_{m+1} + (1-\alpha)(\|\operatorname{Id} - \eta W\| - \eta \|W\|)(\theta_{m+1} - \|W\|) \leqslant 2^m (2\alpha - 1) + (1-\alpha)\mu. \end{cases}$$
(3.10)

Then $(W_i)_{1 \le i \le m+1}$ satisfies Condition 3.1.

Proof We use the operators $(L_i)_{1 \le i \le m+1}$ introduced in Condition 3.1. Per Notation 3.5 and (3.9d),

$$\sup_{\substack{\mathbf{y} \in \mathcal{X} \\ \max \|y_i\| \leqslant 1}} \frac{\|M\mathbf{y} - 2^{m+1}(1 - \alpha)y_0\| + \|M\mathbf{y}\|}{2^{m+1}\alpha} = \sup_{\substack{\mathbf{y} \in \mathcal{X} \\ \|\mathbf{y}\|_{\mathcal{X}} \leqslant 1}} \|M\mathbf{y}\|_{\mathcal{Y}} = \|M\|_{\mathcal{X}, \mathcal{Y}} \quad (3.11)$$

and therefore

$$(\forall \mathbf{y} \in \mathcal{X}) \quad \max_{0 \leqslant i \leqslant m} \|y_i\| \leqslant 1 \quad \Rightarrow \quad \|M\mathbf{y} - 2^{m+1}(1-\alpha)y_0\| + \|M\mathbf{y}\| \leqslant 2^{m+1}\alpha \|\mathbf{M}\|_{\mathcal{X}, \mathcal{Y}}.$$
(3.12)

Now let $x \in \mathcal{X}$ be such that

$$(\forall i \in \{0, \dots, m\}) \quad \|x_i\| \leqslant \begin{cases} 1, & \text{if } i = 0; \\ \|L_i(x_0, \dots, x_{i-1})\|, & \text{if } i \geqslant 1. \end{cases}$$
 (3.13)

(i): We assume that $m \ge 1$. For every $k \in \{i, \ldots, m\}$, it follows from (3.4) that $\theta_k = 0$ and in turn from Lemma 3.3(ii) and (3.13) that $x_k = 0$. Therefore,

$$L_{m+1}x = \sum_{k=0}^{m} (W_{m+1} \circ \dots \circ W_{k+1})x_k = \sum_{k=0}^{i-1} (W_{m+1} \circ \dots \circ W_{k+1})x_k = 0,$$
 (3.14)

and (3.3) clearly holds.

(ii): In view of (i), we assume that, if $m \ge 1$, $(\forall i \in \{1, ..., m\})$ $W_i \ne 0$. We then derive from (3.4) that $(\forall i \in \{1, ..., m\})$ $\theta_i \ge \prod_{k=1}^i ||W_k|| > 0$. If $x_0 = 0$, (3.3) trivially follows from Lemma 3.3(ii), we therefore assume otherwise. Now set

$$(\forall i \in \{0, \dots, m\}) \quad y_i = \frac{x_i}{\theta_i \|x_0\|}.$$
 (3.15)

According to Lemma 3.3(ii), $(\forall i \in \{0, ..., m\}) \|y_i\| \le 1$. On the other hand, it follows from (3.9c), (3.15), and (3.1) that $My = L_{m+1}x/\|x_0\|$. Altogether, we deduce from (3.12) that (3.3) holds.



(iii) \Rightarrow (ii): Take $y \in \mathcal{X}$ such that $||y||_{\mathcal{X}} \leq 1$. Then it follows from (3.9c) and Lemma 3.3(i) that

$$\|My - 2^{m+1}(1 - \alpha)y_0\| + \|My\|$$

$$\leq \|W - 2^{m+1}(1 - \alpha)\operatorname{Id}\| \|y_0\| + \|W\| \|y_0\| + 2\sum_{i=1}^{m} \theta_i \|W_{m+1} \circ \cdots \circ W_{i+1}\| \|y_i\|$$

$$\leq \|W - 2^{m+1}(1 - \alpha)\operatorname{Id}\| - \|W\| + 2\theta_{m+1}$$

$$\leq 2^{m+1}\alpha. \tag{3.16}$$

In turn, (3.11) yields $||M||_{\mathcal{X},\mathcal{V}} \leq 1$.

(iv) \Rightarrow (ii): Let $\mathbf{y} = (y_0, \dots, y_m) \in \mathcal{X}$ be such that $||y_0|| = \dots = ||y_m|| = 1$, and set

$$u = \begin{cases} \sum_{i=1}^{m} \theta_i(W_{m+1} \circ \dots \circ W_{i+1}) y_i, & \text{if } m \neq 0; \\ 0, & \text{if } m = 0. \end{cases}$$
 (3.17)

The assumptions and (3.9b) imply that

$$\begin{cases} \eta \theta_{m+1} \leq \alpha/(1-\alpha) \\ \theta_{m+1} \leq 2^{m+1}\alpha \\ \alpha \theta_{m+1} + (1-\alpha)(\|\operatorname{Id} - \eta W\| - \eta \|W\|)(\theta_{m+1} - \|W\|) \\ \leq 2^{m}(2\alpha - 1) + (1-\alpha)\langle W y_{0} | y_{0} \rangle. \end{cases}$$
(3.18)

On the other hand.

$$\alpha \| W y_{0} + u \| - (1 - \alpha) \langle y_{0} | u \rangle$$

$$= \alpha \| W y_{0} + u \| - (1 - \alpha) \langle \eta W y_{0} + (\operatorname{Id} - \eta W) y_{0} | u \rangle$$

$$\leq \alpha \| W y_{0} + u \| - \eta (1 - \alpha) \langle W y_{0} | u \rangle + (1 - \alpha) \| (\operatorname{Id} - \eta W) y_{0} \| \| u \|.$$
 (3.19)

Since, by Lemma 3.3(i) and (3.18),

$$\eta \sum_{i=0}^{m} \theta_{i} \| W_{m+1} \circ \dots \circ W_{i+1} \| = \eta \theta_{m+1} \leqslant \frac{\alpha}{1-\alpha}, \tag{3.20}$$

we deduce from (3.17) that

$$\eta(1-\alpha)(\|Wy_0\| + \|u\|) \le \alpha. \tag{3.21}$$

However, by Lemma 3.4,

$$||Wy_0|| ||u|| - \langle Wy_0 | u \rangle \le (||Wy_0|| + ||u|| - ||Wy_0 + u||)(||Wy_0|| + ||u||). \tag{3.22}$$

In view of (3.21), this yields

$$\eta(1-\alpha)\big(\|Wy_0\| \|u\| - \langle Wy_0 \mid u \rangle\big) \leqslant \alpha(\|Wy_0\| + \|u\| - \|Wy_0 + u\|), \tag{3.23}$$

that is.

$$\alpha \|Wy_0 + u\| - \eta(1 - \alpha)\langle Wy_0 \mid u \rangle \le \alpha (\|Wy_0\| + \|u\|) - \eta(1 - \alpha)\|Wy_0\| \|u\|.$$
 (3.24)

Therefore, since (3.21) implies that $\alpha - \eta(1-\alpha)\|u\| \ge 0$, it results from (3.19) that

$$\alpha \|Wy_{0} + u\| - (1 - \alpha)\langle y_{0} | u\rangle$$

$$\leq \alpha (\|Wy_{0}\| + \|u\|) - \eta(1 - \alpha)\|Wy_{0}\| \|u\| + (1 - \alpha)\|(\operatorname{Id} - \eta W)y_{0}\| \|u\|$$

$$= \alpha \|u\| + (\alpha - \eta(1 - \alpha)\|u\|)\|Wy_{0}\| + (1 - \alpha)\|(\operatorname{Id} - \eta W)y_{0}\| \|u\|$$

$$\leq \alpha \|u\| + (\alpha - \eta(1 - \alpha)\|u\|)\|W\| + (1 - \alpha)\|(\operatorname{Id} - \eta W)y_{0}\| \|u\|$$

$$= \alpha \|W\| + (\alpha - \eta(1 - \alpha)\|W\|)\|u\| + (1 - \alpha)\|\operatorname{Id} - \eta W\| \|u\|. \tag{3.25}$$

However, since (3.20) implies that $\alpha - \eta(1-\alpha)||W|| \ge 0$, while (3.17) implies that $||u|| \le \theta_{m+1} - ||W||$, we derive from (3.25) that

$$\alpha \| W y_0 + u \| - (1 - \alpha) \langle y_0 | u \rangle$$

$$\leq \alpha \| W \| + (\alpha - \eta (1 - \alpha) \| W \|) (\theta_{m+1} - \| W \|) + (1 - \alpha) \| \operatorname{Id} - \eta W \| (\theta_{m+1} - \| W \|).$$
(3.26)

We also have

$$||Wy_0 + u|| \le ||W|| + ||u|| \le \theta_{m+1}. \tag{3.27}$$

Hence, using (3.26), (3.27), (3.9c), (3.9a), and (3.9d) we obtain

$$(3.18) \Rightarrow \begin{cases} \|Wy_{0} + u\| \leq 2^{m+1}\alpha \\ \alpha \|Wy_{0} + u\| - (1 - \alpha)\langle y_{0} | Wy_{0} + u \rangle \leq 2^{m}(2\alpha - 1) \end{cases}$$

$$\Leftrightarrow \begin{cases} \|My\| \leq 2^{m+1}\alpha \\ \alpha \|My\| - (1 - \alpha)\langle y_{0} | My \rangle \leq 2^{m}(\alpha^{2} - (1 - \alpha)^{2}) \end{cases}$$

$$\Leftrightarrow \begin{cases} \|My\| \leq 2^{m+1}\alpha \\ \|My - 2^{m+1}(1 - \alpha)y_{0}\|^{2} \leq (2^{m+1}\alpha - \|My\|)^{2} \end{cases}$$

$$\Leftrightarrow \|My - 2^{m+1}(1 - \alpha)y_{0}\| + \|My\| \leq 2^{m+1}\alpha$$

$$\Leftrightarrow \|My\|_{\mathbf{Y}} \leq 1. \tag{3.28}$$

Now set $C = \{ y \in \mathcal{X} \mid \|y_0\| = \dots = \|y_m\| = 1 \}$. Then, in view of (3.11), (3.28), and [8, Proposition 11.1(ii)], we conclude that $\|M\|_{\mathcal{X},\mathcal{Y}} = \sup_{y \in \text{conv } C} \|My\|_{\mathcal{Y}} = \sup_{y \in C} \|My\|_{\mathcal{Y}} \leq 1$.

The next result establishes a link between deep neural network structures and the operators introduced in (3.1).

Lemma 3.7 Let $m \ge 1$ be an integer and let $(\mathcal{H}_i)_{0 \le i \le m+1}$ be nonzero real Hilbert spaces. For every $i \in \{1, ..., m+1\}$, let $W_i \in \mathcal{B}(\mathcal{H}_{i-1}, \mathcal{H}_i)$ and let L_i be as in (3.1). Further, for every $i \in \{1, ..., m\}$, let $P_i : \mathcal{H}_i \to \mathcal{H}_i$ be firmly nonexpansive. Set

$$T_m = W_{m+1} \circ P_m \circ W_m \circ \dots \circ P_1 \circ W_1, \tag{3.29}$$

let x and y be distinct points in \mathcal{H}_0 , and set $v_0 = (x - y)/\|x - y\|$. Then there exists $(v_1, \ldots, v_m) \in \mathcal{H}_1 \times \cdots \times \mathcal{H}_m$ such that

$$\begin{cases}
(\forall i \in \{1, \dots, m\}) & \|v_i\| \leq \|L_i(v_0, \dots, v_{i-1})\| \\
\frac{2^m(T_m x - T_m y)}{\|x - y\|} = L_{m+1}(v_0, \dots, v_m).
\end{cases}$$
(3.30)

Proof For every $i \in \{1, ..., m\}$, since P_i is firmly nonexpansive, there exists a nonexpansive operator $Q_i : \mathcal{H}_i \to \mathcal{H}_i$ such that

$$P_i = \frac{\operatorname{Id} + Q_i}{2}. (3.31)$$

We proceed by induction on m. Suppose that m = 1 and set

$$v_1 = \frac{Q_1(W_1 x) - Q_1(W_1 y)}{\|x - y\|},$$
(3.32)



which implies that $||v_1|| \le ||W_1(x-y)||/||x-y|| = ||L_1v_0||$. Then

$$2(T_1x - T_1y) = (W_2 \circ W_1)(x - y) + (W_2 \circ Q_1 \circ W_1)x - (W_2 \circ Q_1 \circ W_1)y$$

= $||x - y||((W_2 \circ W_1)v_0 + W_2v_1)).$ (3.33)

Thus, (3.30) holds for m=1. Next, we assume that m>1 and that there exists $(v_1,\ldots,v_{m-1})\in\mathcal{H}_1\times\cdots\times\mathcal{H}_{m-1}$ such that

$$\begin{cases}
(\forall i \in \{1, \dots, m-1\}) & \|v_i\| \leq \|L_i(v_0, \dots, v_{i-1})\| \\
\frac{2^{m-1}(T_{m-1}x - T_{m-1}y)}{\|x - y\|} = L_m(v_0, \dots, v_{m-1}),
\end{cases} (3.34)$$

and we set

$$v_m = \frac{2^{m-1} \left((Q_m \circ T_{m-1}) x - (Q_m \circ T_{m-1}) y \right)}{\|x - y\|}.$$
 (3.35)

Then (3.29), (3.31), and (3.34) yield

$$T_{m}x - T_{m}y = \frac{(W_{m+1} \circ T_{m-1})x - (W_{m+1} \circ T_{m-1})y}{2} + \frac{(W_{m+1} \circ Q_{m} \circ T_{m-1})x - (W_{m+1} \circ Q_{m} \circ T_{m-1})y}{2}$$

$$= \frac{\|x - y\|}{2^{m}} ((W_{m+1} \circ L_{m})(v_{0}, \dots, v_{m-1}) + W_{m+1}v_{m})$$

$$= \frac{\|x - y\|}{2^{m}} L_{m+1}(v_{0}, \dots, v_{m}). \tag{3.36}$$

In addition, it follows from (3.34) and (3.35) that

$$||v_m|| \leqslant \frac{2^{m-1}||T_{m-1}x - T_{m-1}y||}{||x - y||} = ||L_m(v_0, \dots, v_{m-1})||,$$
(3.37)

which completes the proof.

We now establish connections between Condition 3.1 for linear operators and the concept of averagedness for composite nonlinear operators.

Theorem 3.8 Let $m \ge 1$ be an integer, let $(\mathcal{H}_i)_{0 \le i \le m-1}$ be nonzero real Hilbert spaces, set $\mathcal{H}_m = \mathcal{H}_0$, and let $\alpha \in [1/2, 1]$. For every $i \in \{1, \ldots, m\}$, let $W_i \in \mathcal{B}$ $(\mathcal{H}_{i-1}, \mathcal{H}_i)$ and let $P_i : \mathcal{H}_i \to \mathcal{H}_i$ be firmly nonexpansive. Suppose that $(W_i)_{1 \le i \le m}$ satisfies Condition 3.1. Then $P_m \circ W_m \circ \cdots \circ P_1 \circ W_1$ is α -averaged.

Proof Set $T = P_m \circ W_m \circ \cdots \circ P_1 \circ W_1$. We must show that

$$Q = \left(1 - \frac{1}{\alpha}\right)\operatorname{Id} + \frac{1}{\alpha}T\tag{3.38}$$

is nonexpansive. By assumption, for every $i \in \{1, ..., m\}$, there exists a nonexpansive operator $Q_i : \mathcal{H}_i \to \mathcal{H}_i$ such that (3.31) holds. Let $(L_i)_{1 \le i \le m}$ be as in (3.1) and let x and



y be distinct points in \mathcal{H}_0 . According to Lemma 3.7, there exists $\mathbf{v} = (v_0, \dots, v_{m-1}) \in \mathcal{H}_0 \times \dots \times \mathcal{H}_{m-1}$ such that

$$\begin{cases}
v_0 = \frac{x - y}{\|x - y\|} \\
(\forall i \in \{1, \dots, m - 1\}) \quad \|v_i\| \leq \|L_i(v_0, \dots, v_{i-1})\| \\
\frac{2^{m-1} \left((W_m \circ P_{m-1} \circ \dots \circ P_1 \circ W_1) x - (W_m \circ P_{m-1} \dots \circ P_1 \circ W_1) y \right)}{\|x - y\|} = L_m v.
\end{cases}$$
(3.39)

Condition 3.1 imposes that

$$||L_m \mathbf{v} - 2^m (1 - \alpha) v_0|| + ||L_m \mathbf{v}|| \le 2^m \alpha ||v_0|| = 2^m \alpha, \tag{3.40}$$

which is equivalent to

$$||(W_{m} \circ P_{m-1} \circ \cdots \circ P_{1} \circ W_{1})x - (W_{m} \circ P_{m-1} \cdots \circ P_{1} \circ W_{1})y - 2(1-\alpha)(x-y)|| + ||(W_{m} \circ P_{m-1} \circ \cdots \circ P_{1} \circ W_{1})x - (W_{m} \circ P_{m-1} \cdots \circ P_{1} \circ W_{1})y|| \leq 2\alpha ||x-y||.$$
(3.41)

In turn, we derive from (3.38) and (3.31) that

$$\|Qx - Qy\|$$

$$\leq \frac{1}{\alpha} \| \left(\frac{\text{Id} + Q_m}{2} \circ W_m \circ \cdots \circ P_1 \circ W_1 \right) x - \left(\frac{\text{Id} + Q_m}{2} \circ W_m \circ \cdots \circ P_1 \circ W_1 \right) y$$

$$- (1 - \alpha)(x - y) \|$$

$$\leq \frac{1}{2\alpha} \left(\| (W_m \circ P_{m-1} \circ \cdots \circ P_1 \circ W_1) x - (W_m \circ P_{m-1} \cdots \circ P_1 \circ W_1) y \right)$$

$$- 2(1 - \alpha)(x - y) \| + \| (Q_m \circ W_m \circ P_{m-1} \circ \cdots \circ P_1 \circ W_1) x$$

$$- (Q_m \circ W_m \circ P_{m-1} \cdots \circ P_1 \circ W_1) y \| \right)$$

$$\leq \frac{1}{2\alpha} \left(\| (W_m \circ P_{m-1} \circ \cdots \circ P_1 \circ W_1) x - (W_m \circ P_{m-1} \cdots \circ P_1 \circ W_1) y \right)$$

$$- 2(1 - \alpha)(x - y) \| + \| (W_m \circ P_{m-1} \circ \cdots \circ P_1 \circ W_1) x$$

$$- (W_m \circ P_{m-1} \cdots \circ P_1 \circ W_1) y \| \right)$$

$$\leq \|x - y\|,$$

$$(3.42)$$

which establishes the nonexpansiveness of Q.

Example 3.9 Consider Theorem 3.8 with m=2. In view of Proposition 3.6(iii), $P_2 \circ W_2 \circ P_1 \circ W_1$ is α -averaged if $\|W_2 \circ W_1 - 4(1-\alpha)\operatorname{Id}\| + \|W_2 \circ W_1\| + 2\|W_2\| \|W_1\| \leqslant 4\alpha$. In particular, if $\alpha=1$, this condition is obviously less restrictive than requiring that W_1 and W_2 be nonexpansive.

4 A Variational Inequality Model

In this section, we first investigate an autonomous version of Model 1.1.



Model 4.1 This is the special case of Model 1.1 in which, for every $i \in \{1, ..., m\}$, there exist $R_i \in \mathcal{A}(\mathcal{H}_i)$, say $R_i = \text{prox}_{\varphi_i}$ for some $\varphi_i \in \Gamma_0(\mathcal{H}_i)$ with $\varphi_i(0) = \inf \varphi_i(\mathcal{H}_i)$, $W_i \in \mathcal{B}(\mathcal{H}_{i-1}, \mathcal{H}_i)$, and $b_i \in \mathcal{H}_i$ such that $(\forall n \in \mathbb{N})$ $R_{i,n} = R_i$, $W_{i,n} = W_i$, $b_{i,n} = b_i$. We set

$$(\forall i \in \{1, \dots, m\}) \quad T_i \colon \mathcal{H}_{i-1} \to \mathcal{H}_i \colon x \mapsto R_i(W_i x + b_i) \tag{4.1}$$

and

$$\begin{cases}
F = \operatorname{Fix} (T_{m} \circ \cdots \circ T_{1}) \\
\mathcal{H} = \mathcal{H}_{1} \oplus \cdots \oplus \mathcal{H}_{m-1} \oplus \mathcal{H}_{m} \\
\overrightarrow{\mathcal{H}} = \mathcal{H}_{m} \oplus \mathcal{H}_{1} \oplus \cdots \oplus \mathcal{H}_{m-1} \\
S : \mathcal{H} \to \overrightarrow{\mathcal{H}} : (x_{1}, \dots, x_{m-1}, x_{m}) \mapsto (x_{m}, x_{1}, \dots, x_{m-1}) \\
W : \overrightarrow{\mathcal{H}} \to \mathcal{H} : (x_{m}, x_{1}, \dots, x_{m-1}) \mapsto (W_{1}x_{m}, W_{2}x_{1}, \dots, W_{m}x_{m-1}) \\
\varphi : \mathcal{H} \to]-\infty, +\infty] : \mathbf{x} \mapsto \sum_{i=1}^{m} \varphi_{i}(x_{i}) \\
\psi : \mathcal{H} \to]-\infty, +\infty] : \mathbf{x} \mapsto \sum_{i=1}^{m} (\varphi_{i}(x_{i}) - \langle x_{i} \mid b_{i} \rangle) \\
F = \left\{ \mathbf{x} \in \mathcal{H} \mid x_{1} = T_{1}x_{m}, x_{2} = T_{2}x_{1}, \dots, x_{m} = T_{m}x_{m-1} \right\}, \\
\mathbf{x} = (x_{1}, \dots, x_{m}) \text{ denotes a generic element in } \mathcal{H}.
\end{cases} \tag{4.2}$$

where $\mathbf{x} = (x_1, \dots, x_m)$ denotes a generic element in \mathcal{H} .

4.1 Static Analysis

We start with a property of the compositions of the operators $(T_i)_{1 \le i \le m}$ of (4.1).

Proposition 4.2 Consider the setting of Model 4.1, let i and j be integers such that $1 \le 1$ $j \leq i \leq m$, and let $x \in \mathcal{H}_{i-1}$. Then

$$\|(T_i \circ \dots \circ T_j)x\| \leqslant \|x\| \prod_{k=i}^i \|W_k\| + \sum_{q=i}^i \left(\|b_q\| \prod_{k=q+1}^i \|W_k\| \right). \tag{4.3}$$

Proof In view of (4.1), the property is satisfied when i = j. We now assume that i > j. Since $R_i \in \mathcal{A}(\mathcal{H}_i)$, Proposition 2.21(i) yields

$$\|(T_{i} \circ \cdots \circ T_{j})x\| = \|R_{i}(W_{i}(T_{i-1} \circ \cdots \circ T_{j})x + b_{i})\|$$

$$= \|R_{i}(W_{i}(T_{i-1} \circ \cdots \circ T_{j})x + b_{i}) - R_{i}0\|$$

$$\leq \|W_{i}(T_{i-1} \circ \cdots \circ T_{j})x + b_{i}\|$$

$$\leq \|W_{i}\| \|(T_{i-1} \circ \cdots \circ T_{i})x\| + \|b_{i}\|. \tag{4.4}$$

We thus obtain (4.3) recursively.

Next, we establish a connection between Model 4.1 and a variational inequality.

Proposition 4.3 In the setting of Model 4.1, consider the variational inequality problem

find
$$\overline{x}_1 \in \mathcal{H}_1, \dots, \overline{x}_m \in \mathcal{H}_m$$
 such that
$$\begin{cases}
b_1 \in \overline{x}_1 - W_1 \overline{x}_m + \partial \varphi_1(\overline{x}_1) \\
b_2 \in \overline{x}_2 - W_2 \overline{x}_1 + \partial \varphi_2(\overline{x}_2) \\
\vdots \\
b_m \in \overline{x}_m - W_m \overline{x}_{m-1} + \partial \varphi_m(\overline{x}_m).
\end{cases}$$
(4.5)

Then the following hold:

(i) The set of solutions to (4.5) is F.



- (ii) $F = \operatorname{zer}(\operatorname{Id} W \circ S + \partial \psi) = \operatorname{Fix}(\operatorname{prox}_{\psi} \circ W \circ S).$
- (iii) $\mathbf{F} = \{ (T_1 \overline{x}_m, (T_2 \circ T_1) \overline{x}_m, \dots, (T_{m-1} \circ \dots \circ T_1) \overline{x}_m, \overline{x}_m) \mid \overline{x}_m \in F \}.$
- (iv) Suppose that $(W_i)_{1 \le i \le m}$ satisfies Condition 3.1 for some $\alpha \in [1/2, 1]$. Then F is closed and convex.
- (v) Suppose that $(W_i)_{1 \le i \le m}$ satisfies Condition 3.1 for some $\alpha \in [1/2, 1]$ and that one of the following holds:
 - (a) $ran(T_m \circ \cdots \circ T_1)$ is bounded.
 - (b) There exists $j \in \{1, ..., m\}$ such that dom φ_j is bounded.

Then F and F are nonempty.

- (vi) Suppose that $\mathbf{Id} \mathbf{W} \circ \mathbf{S}$ is monotone. Then \mathbf{F} is closed and convex. In addition, \mathbf{F} and \mathbf{F} are nonempty if any of the following holds:
 - (a) **Id** $-\mathbf{W} \circ \mathbf{S} + \partial \boldsymbol{\varphi}$ is surjective.
 - (b) $\partial \varphi W \circ S$ is maximally monotone.
 - (c) $\max_{1 \le i \le m} ||W_i|| \le 1$, $S^* W$ has closed range, and $\ker(S W^*) = \{0\}$.
 - (d) $\max_{1 \leq i \leq m} ||W_i|| \leq 1$ and, for every $i \in \{1, ..., m\}$, dom $\varphi_i^* = \mathcal{H}_i$.
 - (e) For every $i \in \{1, ..., m\}$, dom $\varphi_i = \mathcal{H}$ and dom $\varphi_i^* = \mathcal{H}_i$.
 - (f) $S^* W$ has closed range, $\ker(S W^*) = \{0\}$, and, for every $i \in \{1, ..., m\}$, dom $\varphi_i = \mathcal{H}_i$.
 - (g) For every $i \in \{1, ..., m\}$, dom φ_i is bounded.

Proof We first observe that $S \in \mathcal{B}(\mathcal{H}, \overrightarrow{\mathcal{H}}), W \in \mathcal{B}(\overrightarrow{\mathcal{H}}, \mathcal{H}), \varphi \in \Gamma_0(\mathcal{H}), \text{ and } \psi \in \Gamma_0(\mathcal{H}).$ (i): Let $x \in \mathcal{H}$. Then

$$x \text{ solves (4.5)} \Leftrightarrow \begin{cases} W_{1}x_{m} + b_{1} \in x_{1} + \partial \varphi_{1}(x_{1}) \\ W_{2}x_{1} + b_{2} \in x_{2} + \partial \varphi_{2}(x_{2}) \\ \vdots \\ W_{m}x_{m-1} + b_{m} \in x_{m} + \partial \varphi_{m}(x_{m}). \end{cases}$$

$$\Leftrightarrow \begin{cases} x_{1} = \operatorname{prox}_{\varphi_{1}}(W_{1}x_{m} + b_{1}) = T_{1}x_{m} \\ x_{2} = \operatorname{prox}_{\varphi_{2}}(W_{2}x_{1} + b_{2}) = T_{2}x_{1} \\ \vdots \\ x_{m} = \operatorname{prox}_{\varphi_{m}}(W_{m}x_{m-1} + b_{m}) = T_{m}x_{m-1}. \end{cases}$$

$$(4.6)$$

(ii): Let $x \in \mathcal{H}$. Using (4.2), we obtain

$$x \text{ solves } (4.5) \Leftrightarrow \mathbf{0} \in x - W(Sx) + \partial \psi(x) \Leftrightarrow x = \operatorname{prox}_{\psi}(W(Sx)).$$
 (4.8)

- (iii): Clear from the definitions of F and F.
- (iv): Define m firmly nonexpansive operators by $(\forall i \in \{1, ..., m\})$ $P_i : \mathcal{H}_i \to \mathcal{H}_i : y \mapsto R_i(y+b_i)$. Then it follows from (4.1) and Theorem 3.8 applied to $(P_i)_{1 \leqslant i \leqslant m}$ that $T_m \circ \cdots \circ T_1$ is nonexpansive. In turn, we derive from [8, Corollary 4.24] that its fixed point set F is closed and convex.
- (v): Thanks to (iii), it is enough to show that $F \neq \emptyset$. Set $T = T_m \circ \cdots \circ T_1$ and recall that it is nonexpansive by virtue of Theorem 3.8.
- (a): Let C be a closed ball such that ran $T \subset C$ and set $S = T|_C$. Then $S: C \to C$ is nonexpansive and therefore [8, Proposition 4.29] asserts that Fix $T = \text{Fix } S \neq \emptyset$.



(b) \Rightarrow (a): We have ran $T_j \subset \operatorname{ran} R_j = \operatorname{ran} \operatorname{prox}_{\varphi_j} = \operatorname{dom} (\operatorname{Id} + \partial \varphi_j) = \operatorname{dom} \partial \varphi_j \subset \operatorname{dom} \varphi_j$. Hence ran T_j is bounded and Proposition 4.2 (with i = m) implies that

$$\operatorname{ran} T \subset \begin{cases} \operatorname{ran} T_m, & \text{if } j = m; \\ (T_m \circ \dots \circ T_{j+1})(\operatorname{ran} T_j), & \text{if } 1 \leqslant j \leqslant m-1 \end{cases}$$
(4.9)

is likewise.

(vi): Set $A = \mathbf{Id} - \mathbf{W} \circ \mathbf{S} + \partial \psi$. Since $\mathbf{Id} - \mathbf{W} \circ \mathbf{S}$ is monotone and continuous, it is maximally monotone [8, Corollary 20.28], with \mathcal{H} as its domain. Since $\partial \psi$ is also maximally monotone [8, Theorem 20.25], \mathbf{A} is likewise [8, Corollary 25.5(i)] and hence $\mathbf{F} = \operatorname{zer} \mathbf{A}$ is closed and convex [8, Proposition 23.39]. Next, we note that, in view of (iii), $F \neq \emptyset \Leftrightarrow \mathbf{F} \neq \emptyset$.

(a): The hypothesis implies that $(b_i)_{1 \le i \le m} \in \text{ran} (\mathbf{Id} - \mathbf{W} \circ \mathbf{S} + \partial \boldsymbol{\varphi})$ and therefore that (4.5) has a solution, i.e., $\mathbf{F} \ne \emptyset$.

(b) \Rightarrow (a): The claim follows from Minty's theorem [8, Theorem 21.1].

(c) \Rightarrow (a): We have $\|W \circ S\| = \|W\| = \max_{1 \le i \le m} \|W_i\| \le 1$. Therefore, $-W \circ S$ is nonexpansive, which implies that $(\mathbf{Id} - W \circ S)/2$ is firmly nonexpansive [8, Corollary 4.5], that is $(\forall x \in \mathcal{H}) \langle x - W(Sx) | x \rangle \ge \|x - W(Sx)\|^2/2$. Consequently, $\mathbf{Id} - W \circ S$ is 3^* monotone [8, Proposition 25.16], while $\partial \varphi$ is also 3^* monotone [8, Example 25.13]. Finally, since S is unitary,

$$\operatorname{ran}\left(\operatorname{Id}-W\circ S\right)=\operatorname{ran}\left(S^{*}-W\right)=\overline{\operatorname{ran}}\left(S-W^{*}\right)^{*}=\left(\operatorname{ker}\left(S-W^{*}\right)\right)^{\perp}=\mathcal{H},\ (4.10)$$

which shows that $\mathbf{Id} - W \circ S$ is surjective. Altogether, since [8, Corollary 25.5(i)] implies that $\mathbf{Id} - W \circ S + \partial \varphi$ is maximally monotone, it follows from [8, Corollary 25.27(i)] that $\mathbf{Id} - W \circ S + \partial \varphi$ is surjective.

(d) \Rightarrow (a): We have dom $\varphi^* = \mathcal{H}$. Therefore, since int dom $\varphi^* \subset \text{dom } \partial \varphi^*$ [8, Proposition 16.27], we have ran $\partial \varphi = \text{dom } (\partial \varphi)^{-1} = \text{dom } \partial \varphi^* = \mathcal{H}$. Hence, $\partial \varphi$ is surjective. We conclude using the same arguments as in (c): $\partial \varphi$ and $\text{Id} - W \circ S$ are both 3* monotone and their sum is maximally monotone, which allows us to invoke [8, Corollary 25.27(i)].

(e) \Rightarrow (a): As seen in (d), $\partial \varphi$ is surjective. We have $\mathcal{H} = \operatorname{int} \operatorname{dom} \varphi \subset \operatorname{dom} \partial \varphi$ [8, Proposition 16.27]. Consequently, $\mathcal{H} = \operatorname{dom} (\operatorname{Id} - W \circ S) \subset \operatorname{dom} \partial \varphi$. Altogether, since $\partial \varphi$ is 3* monotone, it follows from [8, Corollary 25.27(ii)] that $\operatorname{Id} - W \circ S + \partial \varphi$ is surjective.

(f) \Rightarrow (a): As seen in (c), $\mathbf{Id} - \mathbf{W} \circ \mathbf{S}$ is surjective and $\partial \varphi$ is 3^* monotone. In addition, dom ($\mathbf{Id} - \mathbf{W} \circ \mathbf{S}$) \subset dom $\partial \varphi$ since $\mathcal{H} = \operatorname{int} \operatorname{dom} \varphi \subset \operatorname{dom} \partial \varphi$ [8, Proposition 16.27]. Altogether, it follows from [8, Corollary 25.27(ii)] that $\mathbf{Id} - \mathbf{W} \circ \mathbf{S} + \partial \varphi$ is surjective.

(g): Here dom $A = \text{dom } \partial \varphi \subset \text{dom } \varphi = \times_{i=1}^m \text{dom } \varphi_i$ is bounded. Hence, $F = \text{zer } A \neq \emptyset$ [8, Proposition 23.36(iii)].

Remark 4.4 In Proposition 4.3(vi), it is required that $\mathbf{Id} - W \circ S$ be monotone, or equivalently, that its self-adjoint part $\mathbf{Id} - (W \circ S + S^* \circ W^*)/2$ be positive. In a finite-dimensional setting, this just means that the eigenvalues of the matrix $WS + S^*W^*$ are in $]-\infty, 2]$.

Remark 4.5 Let $\overline{x} \in \mathcal{H}$ be a solution to the variational inequality (4.5). A natural question is whether \overline{x} solves a minimization problem. In general the answer is negative. For instance, for $m \geqslant 3$ layers, even if the Hilbert spaces $(\mathcal{H}_i)_{1 \leqslant i \leqslant m}$ are identical, $W = \mathbf{Id}$, the vectors $(b_i)_{1 \leqslant i \leqslant m}$ are zero, and the functions $(\varphi_i)_{1 \leqslant i \leqslant m}$ are indicator functions of closed convex sets $(C_i)_{1 \leqslant i \leqslant m}$, the solutions to (4.5) do not minimize any function $\Phi \colon \mathcal{H} \to \mathbb{R}$ [4]. A



rather restrictive scenario in which the answer is positive is when $\mathbf{Id} - W \circ S$ is monotone and $W \circ S$ is self-adjoint. Then \overline{x} is a minimizer of $\Phi \colon x \mapsto (1/2)\langle x - W(Sx) \mid x \rangle + \psi(x)$.

Example 4.6 In Model 4.1, suppose that, for every $i \in \{1, \ldots, m\}$, $\mathcal{H}_i = \mathbb{R}^{N_i}$ for some strictly positive integer N_i . In addition, assume that, for every $i \in \{1, \ldots, m\}$, R_i is a separable activation operator with respect to the canonical basis of \mathbb{R}^{N_i} (see Proposition 2.24), and that it employs the ReLU activation functions of Example 2.6. For every $i \in \{1, \ldots, m\}$, let $x_i = (\xi_{i,k})_{1 \leqslant k \leqslant N_i} \in \mathbb{R}^{N_i}$ and set $b_i = (\beta_{i,k})_{1 \leqslant k \leqslant N_i}$. Then it follows from Proposition 4.3(i) that $(x_1, \ldots, x_m) \in \mathbf{F}$ if and only if, for every $i \in \{1, \ldots, m\}$, $x_i \in [0, +\infty[^{N_i}]$ and

$$\begin{cases} (\forall k \in \{1, \dots, N_1\}) & [W_1 x_m]_k + \beta_{1,k} - \xi_{1,k} \in \mathcal{I}(\xi_{1,k}) \\ (\forall k \in \{1, \dots, N_2\}) & [W_2 x_1]_k + \beta_{2,k} - \xi_{2,k} \in \mathcal{I}(\xi_{2,k}) \\ & \vdots \\ (\forall k \in \{1, \dots, N_{m-1}\}) & [W_{m-1} x_{m-2}]_k + \beta_{m-1,k} - \xi_{m-1,k} \in \mathcal{I}(\xi_{m-1,k}) \\ (\forall k \in \{1, \dots, N_m\}) & [W_m x_{m-1}]_k + \beta_{m,k} - \xi_{m,k} \in \mathcal{I}(\xi_{m,k}) \end{cases}$$

$$(4.11)$$

where, given $x \in \mathcal{H}_{i-1}$, $[W_i x]_k$ is the kth component of $W_i x$ and

$$(\forall \xi \in [0, +\infty[) \quad \mathcal{I}(\xi) = \begin{cases} \{0\}, & \text{if } \xi \in]0, +\infty[; \\]-\infty, 0], & \text{if } \xi = 0. \end{cases}$$
(4.12)

Altogether, we conclude that F is a closed convex polyhedron.

4.2 Asymptotic Analysis

We investigate the asymptotic behavior of (1.2) in the context of Model 4.1.

Theorem 4.7 In the setting of Model 4.1, set $T = T_m \circ \cdots \circ T_1$, let $\alpha \in [1/2, 1]$, and suppose that the following hold:

- (a) $F \neq \emptyset$.
- (b) $(W_i)_{1 \le i \le m}$ satisfies Condition 3.1 with parameter α .
- (c) One of the following is satisfied:
 - (i) $\lambda_n \equiv 1/\alpha = 1$ and $Tx_n x_n \to 0$.
 - (ii) $(\lambda_n)_{n\in\mathbb{N}}$ lies in $]0, 1/\alpha[$ and $\sum_{n\in\mathbb{N}} \lambda_n(1-\alpha\lambda_n) = +\infty.$

Then $(x_n)_{n\in\mathbb{N}}$ converges weakly to a point $\overline{x}_m\in F$ and $(T_1\overline{x}_m, (T_2\circ T_1)\overline{x}_m, \ldots, (T_{m-1}\circ \cdots \circ T_1)\overline{x}_m, \overline{x}_m)$ solves (4.5). Now suppose that, in addition, any of the following holds:

- (iii) For every $i \in \{1, ..., m-1\}$, R_i is weakly sequentially continuous.
- (iv) For every $i \in \{1, ..., m-1\}$, R_i is a separable activation operator in the sense of Proposition 2.24.
- (v) For every $i \in \{1, ..., m-1\}$, \mathcal{H}_i is finite-dimensional.
- (vi) For some $\varepsilon \in]0, 1/2[$, $(\lambda_n)_{n \in \mathbb{N}}$ lies in $[\varepsilon, (1-\varepsilon)(\varepsilon+1/\alpha)]$ and, for every $i \in \{1, \ldots, m\}$, $\mathcal{H}_i = \mathcal{H}$ and there exists $\beta_i \in]0, 1[$ such that $||W_i 2(1-\beta_i)| + ||W_i|| \leq 2\beta_i$.

Then, for every $i \in \{1, ..., m-1\}$, $(x_{i,n})_{n \in \mathbb{N}}$ converges weakly to $\overline{x}_i = (T_i \circ \cdots \circ T_1)\overline{x}_m$ and $(\overline{x}_1, ..., \overline{x}_m)$ solves (4.5).



Proof We first derive from (1.2) and Model 4.1 that

$$(\forall n \in \mathbb{N}) \quad x_{n+1} = x_n + \lambda_n (Tx_n - x_n). \tag{4.13}$$

Now set $(\forall i \in \{1, ..., m\})$ $P_i : \mathcal{H}_i \to \mathcal{H}_i : y \mapsto R_i(y + b_i)$. Then (4.1) yields $T = P_m \circ W_m \circ \cdots \circ P_1 \circ W_1$ and, since the operators $(R_i)_{1 \leqslant i \leqslant m}$ are firmly nonexpansive, the operators $(P_i)_{1 \leqslant i \leqslant m}$ are likewise. Hence, it follows from (b), Theorem 3.8, and (4.2) that

T is
$$\alpha$$
-averaged and Fix $T = F$. (4.14)

- (i): In view of (4.14), T is nonexpansive and hence we derive from [8, Theorem 5.14(i)] that $(x_n)_{n\in\mathbb{N}}$ converges weakly to a point in F. The second assertion follows from Proposition 4.3(iii).
- (ii): In view of (4.13) and (4.14), [8, Theorem 5.15(iii) and Proposition 5.16(iii)] imply that $(x_n)_{n\in\mathbb{N}}$ converges weakly to a point in F, and we conclude by invoking Proposition 4.3(iii).

We now prove the convergence of the individual sequences under each assumption.

- (iii): We have already established that $x_n
 ightharpoonup \overline{x}_m$. Since W_1 is weakly continuous as a bounded linear operator, so is T_1 in (4.1). Hence, (1.2) implies that $x_{1,n} = T_1 x_n
 ightharpoonup T_1 \overline{x}_m = \overline{x}_1$. Likewise, we obtain successively $x_{2,n} = T_2 x_{1,n}
 ightharpoonup T_2 \overline{x}_1 = \overline{x}_2$, $x_{3,n} = T_3 x_{2,n}
 ightharpoonup T_3 \overline{x}_2 = \overline{x}_3, \dots, x_{m,n} = T_m x_{m-1,n}
 ightharpoonup T_m \overline{x}_{m-1} = \overline{x}_m$.
 - (iv) \Rightarrow (iii): See [8, Proposition 24.12(iii)].
- (v)⇒(iii): A proximity operator is nonexpansive and therefore continuous, hence weakly continuous in a finite-dimensional setting.
- (vi): As shown above, $x_n \rightarrow \overline{x}_m \in F$. It follows from Proposition 3.6(iii) and Theorem 3.8 (applied with m = 1) that, for every $i \in \{1, ..., m\}$, T_i is β_i -averaged. Hence, upon applying [24, Theorem 3.5(ii)] with α as an averaging constant of T, we infer that

$$\begin{cases}
(\operatorname{Id} - T_{1})x_{n} - (\operatorname{Id} - T_{1})\overline{x}_{m} \to 0 \\
(\operatorname{Id} - T_{2})(T_{1}x_{n}) - (\operatorname{Id} - T_{2})(T_{1}\overline{x}_{m}) \to 0 \\
\vdots \\
(\operatorname{Id} - T_{m})((T_{m-1} \circ \cdots \circ T_{1})x_{n}) - (\operatorname{Id} - T_{m})((T_{m-1} \circ \cdots \circ T_{1})\overline{x}_{m}) \to 0.
\end{cases} (4.15)$$

Thus, $x_{1,n} - x_n = T_1 x_n - x_n \to T_1 \overline{x}_m - \overline{x}_m$, which implies that $x_{1,n} = (x_{1,n} - x_n) + x_n \to (T_1 \overline{x}_m - \overline{x}_m) + \overline{x}_m = T_1 \overline{x}_m$. However, since $x_{2,n} - x_{1,n} = (T_2 \circ T_1) x_n - T_1 x_n \to (T_2 \circ T_1) \overline{x}_m - T_1 \overline{x}_m$, we obtain $x_{2,n} \to (T_2 \circ T_1) \overline{x}_m$. Continuing this telescoping process yields the claim.

The next result covers the case when the variational inequality problem (4.5) has no solution.

Proposition 4.8 In the setting of Model 4.1, suppose that $(W_i)_{1 \le i \le m}$ satisfies Condition 3.1 with $\alpha \in [1/2, 1]$, that $(\lambda_n)_{n \in \mathbb{N}}$ lies in $[\varepsilon, (1/\alpha) - \varepsilon]$, for some $\varepsilon \in]0, 1/2[$, and that $F = \emptyset$. Then $||x_n|| \to +\infty$.

Proof We derive from (4.13) and (4.14) that, for every $n \in \mathbb{N}$, $x_{n+1} = x_n + \mu_n(Qx_n - x_n)$, where $Q = (1 - 1/\alpha) \operatorname{Id} + (1/\alpha)T$ is nonexpansive and such that Fix Q = F, and $\mu_n = \alpha \lambda_n \in]0$, 1[. Hence the claims follows from [8, Proposition 4.29] and [12, Corollary 9(b)].

Remark 4.9 When assumptions (a)–(c) in Theorem 4.7 are satisfied, the neural network described in Model 1.1 is robust to perturbations of its input. Indeed, since T is α -averaged in (4.13), we can write the updating rule as $x_{n+1} = Q_n x_n$, where Q_n is nonexpansive. In



turn, if x_0 and \widetilde{x}_0 are two inputs in \mathcal{H}_0 , for a given $n \in \mathbb{N}$, the resulting outputs x_n and \widetilde{x}_n are such that $||x_n - \widetilde{x}_n|| \le ||x_0 - \widetilde{x}_0||$.

Remark 4.10 In connection with Theorem 4.7 and Remark 4.5, let us underline that in general the weak limit \overline{x}_m of $(x_n)_{n\in\mathbb{N}}$ does not solve a minimization problem. A very special case in which it does is the following. Suppose that m=2, $\mathcal{H}_1=\mathcal{H}$, $\|W_1\|\leqslant 1$, and $W_2=W_1^*$. Set $\psi_1=\varphi_1-\langle\cdot\mid b_1\rangle$ and $\psi_2=\varphi_2-\langle\cdot\mid b_2\rangle$, and let $\overline{x}_2\in F$, i.e., $\overline{x}_2=(\operatorname{prox}_{\psi_2}\circ W_1^*\circ\operatorname{prox}_{\psi_1}\circ W_1)\overline{x}_2$. It follows from [21, Remark 3.10(iv)] that there exists a function $\vartheta\in\Gamma_0(\mathcal{H})$ such that $W_1^*\circ\operatorname{prox}_{\psi_1}\circ W_1=\operatorname{prox}_{\vartheta}$. Thus, \overline{x}_2 is a fixed point of the backward-backward operator $\operatorname{prox}_{\psi_2}\circ\operatorname{prox}_{\vartheta}$. It then follows from [20, Remark 6.13] that \overline{x}_2 is a minimizer of $\mathbb{I}^{\vartheta}+\psi_2$, where $\mathbb{I}^{\vartheta}:x\mapsto\inf_{y\in\mathcal{H}}(\vartheta(y)+\|x-y\|^2/2)$ is the Moreau envelope of ϑ .

Remark 4.11 To model closely existing deep neural networks, we have chosen the activation operators in Definition 2.20 and Model 4.1 to be proximity operators. However, as is clear from the results of Section 3 and in particular the central Theorem 3.8, an activation operator $R_i: \mathcal{H}_i \to \mathcal{H}_i$ could more generally be a firmly nonexpansive operator that admits 0 as a fixed point. By [8, Corollary 23.9], this means that R_i is the resolvent of some maximally monotone operator $A_i: \mathcal{H}_i \to 2^{\mathcal{H}_i}$ (i.e., $R_i = (\mathrm{Id} + A_i)^{-1}$) such that $0 \in A_i$ 0. In this context, the variational inequality (4.5) assumes the more general form of a system of monotone inclusions, namely,

find
$$\overline{x}_1 \in \mathcal{H}_1, \dots, \overline{x}_m \in \mathcal{H}_m$$
 such that
$$\begin{cases}
b_1 \in \overline{x}_1 - W_1 \overline{x}_m + A_1 \overline{x}_1 \\
b_2 \in \overline{x}_2 - W_2 \overline{x}_1 + A_2 \overline{x}_2 \\
\vdots \\
b_m \in \overline{x}_m - W_m \overline{x}_{m-1} + A_m \overline{x}_m.
\end{cases}$$
(4.16)

5 Analysis of Nonperiodic Networks

We analyze the deep neural network described in Model 1.1 in the following scenario.

Assumption 5.1 In the setting of Model 1.1, there exist sequences $(\omega_n)_{n\in\mathbb{N}}\in\ell^1_+$, $(\rho_n)_{n\in\mathbb{N}}\in\ell^1_+$, $(\eta_n)_{n\in\mathbb{N}}\in\ell^1_+$, and $(\nu_n)_{n\in\mathbb{N}}\in\ell^1_+$ for which the following hold for every $i\in\{1,\ldots,m\}$:

- (i) There exists $W_i \in \mathcal{B}(\mathcal{H}_{i-1}, \mathcal{H}_i)$ such that $(\forall n \in \mathbb{N}) \|W_{i,n} W_i\| \leq \omega_n$.
- (ii) There exists $R_i \in \mathcal{A}(\mathcal{H}_i)$ such that $(\forall n \in \mathbb{N})(\forall x \in \mathcal{H}_i) \|R_{i,n}x R_ix\| \le \rho_n \|x\| + \eta_n$.
- (iii) There exists $b_i \in \mathcal{H}_i$ such that $(\forall n \in \mathbb{N}) ||b_{i,n} b_i|| \leq \nu_n$.

In addition, we set

$$(\forall i \in \{1, \dots, m\}) \quad T_i : \mathcal{H}_{i-1} \to \mathcal{H}_i : x \mapsto R_i(W_i x + b_i). \tag{5.1}$$

Proposition 5.2 In the setting of Model 1.1, suppose that Assumption 5.1 is satisfied, let $i \in \{1, ..., m\}$, and set

$$(\forall n \in \mathbb{N}) \quad \chi_{i,n} = \rho_n \|W_{i,n}\| + \omega_n \quad and \quad \zeta_{i,n} = \rho_n \|b_{i,n}\| + \eta_n + \nu_n. \tag{5.2}$$

Then $(\chi_{i,n})_{n\in\mathbb{N}} \in \ell^1_+$, $(\zeta_{i,n})_{n\in\mathbb{N}} \in \ell^1_+$, and $(\forall n \in \mathbb{N})(\forall x \in \mathcal{H}_{i-1}) \|T_{i,n}x - T_ix\| \leq \chi_{i,n}\|x\| + \zeta_{i,n}$.



as claimed.

Proof According to Assumptions 5.1(i) and 5.1(ii), $\sup_{n\in\mathbb{N}} \|W_{i,n}\| < +\infty$ and $\sup_{n\in\mathbb{N}} \|b_{i,n}\| < +\infty$. It then follows from (5.2) that $(\chi_{i,n})_{n\in\mathbb{N}} \in \ell^1_+$ and $(\zeta_{i,n})_{n\in\mathbb{N}} \in \ell^1_+$. Hence, we deduce from (1.1), (5.1), the nonexpansiveness of R_i , and Assumption 5.1 that

$$(\forall n \in \mathbb{N})(\forall x \in \mathcal{H}_{i-1}) \| T_{i,n}x - T_{i}x \|$$

$$\leq \| R_{i,n}(W_{i,n}x + b_{i,n}) - R_{i}(W_{i,n}x + b_{i,n}) \| + \| R_{i}(W_{i,n}x + b_{i,n}) - R_{i}(W_{i}x + b_{i}) \|$$

$$\leq \rho_{n} \| W_{i,n}x + b_{i,n} \| + \eta_{n} + \| W_{i,n}x + b_{i,n} - W_{i}x - b_{i} \|$$

$$\leq \rho_{n} (\| W_{i,n} \| \| x \| + \| b_{i,n} \|) + \eta_{n} + \| W_{i,n} - W_{i} \| \| x \| + \| b_{i,n} - b_{i} \|$$

$$\leq \rho_{n} (\| W_{i,n} \| \| x \| + \| b_{i,n} \|) + \eta_{n} + \omega_{n} \| x \| + \nu_{n}$$

$$= \chi_{i,n} \| x \| + \zeta_{i,n},$$

$$(5.3)$$

Proposition 5.3 *In the setting of Model 1.1, suppose that Assumption 5.1 is satisfied. Then, for every* $i \in \{1, ..., m\}$ *, there exist* $(\tau_{i,n})_{n \in \mathbb{N}} \in \ell^1_+$ *and* $(\theta_{i,n})_{n \in \mathbb{N}} \in \ell^1_+$ *such that*

$$(\forall n \in \mathbb{N})(\forall x \in \mathcal{H}) \quad \|(T_{i,n} \circ \cdots \circ T_{1,n})x - (T_i \circ \cdots \circ T_1)x\| \leqslant \tau_{i,n}\|x\| + \theta_{i,n}. \tag{5.4}$$

Proof For every $i \in \{1, ..., m\}$, define $(\chi_{i,n})_{n \in \mathbb{N}}$ and $(\zeta_{i,n})_{n \in \mathbb{N}}$ as in (5.2), According to Proposition 5.2, (5.4) is satisfied for i = 1 by setting $(\forall n \in \mathbb{N})$ $\tau_{1,n} = \chi_{1,n}$ and $\theta_{1,n} = \zeta_{1,n}$. Next, let us assume that (5.4) holds for $i \in \{1, ..., m-1\}$ and set

$$(\forall n \in \mathbb{N}) \begin{cases} \tau_{i+1,n} = (\|W_{i+1}\| + \chi_{i+1,n})\tau_{i,n} + \chi_{i+1,n} \prod_{k=1}^{i} \|W_{k}\| \\ \theta_{i+1,n} = (\|W_{i+1}\| + \chi_{i+1,n})\theta_{i,n} + \chi_{i+1,n} \sum_{j=1}^{i} \left(\|b_{j}\| \prod_{k=j+1}^{i} \|W_{k}\|\right) + \zeta_{i+1,n}. \end{cases}$$

$$(5.5)$$

Then the sequences $(\tau_{i+1,n})_{n\in\mathbb{N}}$ and $(\theta_{i+1,n})_{n\in\mathbb{N}}$ belong to ℓ_+^1 . Now let $n\in\mathbb{N}$ and $x\in\mathcal{H}$. Upon invoking Proposition 5.2, the nonexpansiveness of R_{i+1} , and Proposition 4.2, we obtain

$$\|(T_{i+1,n} \circ \cdots \circ T_{1,n})x - (T_{i+1} \circ \cdots \circ T_{1})x\|$$

$$\leq \|(T_{i+1,n} \circ T_{i,n} \circ \cdots \circ T_{1,n})x - (T_{i+1} \circ T_{i,n} \circ \cdots \circ T_{1,n})x\|$$

$$+ \|(T_{i+1} \circ T_{i,n} \circ \cdots \circ T_{1,n})x - (T_{i+1} \circ T_{i} \circ \cdots \circ T_{1})x\|$$

$$\leq \chi_{i+1,n} \|(T_{i,n} \circ \cdots \circ T_{1,n})x\| + \zeta_{i+1,n}$$

$$+ \|(T_{i+1} \circ T_{i,n} \circ \cdots \circ T_{1,n})x - (T_{i+1} \circ T_{i} \circ \cdots \circ T_{1})x\|$$

$$\leq \chi_{i+1,n} (\|(T_{i,n} \circ \cdots \circ T_{1,n})x - (T_{i} \circ \cdots \circ T_{1})x\| + \|(T_{i} \circ \cdots \circ T_{1})x\|) + \zeta_{i+1,n}$$

$$+ \|R_{i+1} ((W_{i+1} \circ T_{i,n} \circ \cdots \circ T_{1,n})x + b_{i+1}) - R_{i+1} ((W_{i+1} \circ T_{i} \circ \cdots \circ T_{1})x + b_{i+1})\|$$

$$\leq (\|W_{i+1}\| + \chi_{i+1,n}) \|(T_{i,n} \circ \cdots \circ T_{1,n})x - (T_{i} \circ \cdots \circ T_{1})x\|$$

$$+ \chi_{i+1,n} \|(T_{i} \circ \cdots \circ T_{1})x\| + \zeta_{i+1,n}$$

$$\leq (\|W_{i+1}\| + \chi_{i+1,n}) (\tau_{i,n}\|x\| + \theta_{i,n})$$

$$+ \chi_{i+1,n} \left(\|x\| \prod_{k=1}^{i} \|W_{k}\| + \sum_{j=1}^{i} \left(\|b_{j}\| \prod_{k=j+1}^{i} \|W_{k}\|\right)\right) + \zeta_{i+1,n}$$

$$= \tau_{i+1,n} \|x\| + \theta_{i+1,n}, \tag{5.6}$$

which proves the result by induction.



We can now present the main result of this section on the asymptotic behavior of Model 1.1. The proof of this result relies on Theorem 4.7, which it extends.

Theorem 5.4 Consider the setting of Model 1.1 and let $\alpha \in [1/2, 1]$. Suppose that Assumption 5.1 is satisfied as well as the following:

- $F = \text{Fix} T \neq \emptyset$, where $T = T_m \circ \cdots \circ T_1$.
- (b) $(W_i)_{1 \leq i \leq m}$ satisfies Condition 3.1 with parameter α .
- One of the following is satisfied:

 - (i) $\lambda_n \equiv \alpha = 1 \text{ and } Tx_n x_n \to 0.$ (ii) $(\lambda_n)_{n \in \mathbb{N}} \text{ lies in }]0, 1/\alpha[\text{ and } \sum_{n \in \mathbb{N}} \lambda_n (1 \alpha \lambda_n) = +\infty.$

Then $(x_n)_{n\in\mathbb{N}}$ converges weakly to a point $\overline{x}_m\in F$ and $(T_1\overline{x}_m,(T_2\circ T_1)\overline{x}_m,\ldots,(T_{m-1}\circ T_m)\overline{x}_m)$ $\cdots \circ T_1)\overline{x}_m, \overline{x}_m$) solves (4.5). Now suppose that, in addition, any of the following holds:

- For every $i \in \{1, ..., m-1\}$, R_i is weakly sequentially continuous.
- (iv) For every $i \in \{1, ..., m-1\}$, R_i is a separable activation function in the sense of Proposition 2.24.
- For every $i \in \{1, ..., m-1\}$, \mathcal{H}_i is finite-dimensional.
- For some $\varepsilon \in (0, 1/2)$, $(\lambda_n)_{n \in \mathbb{N}}$ lies in $[\varepsilon, (1-\varepsilon)(\varepsilon+1/\alpha)]$ and, for every $i \in (0, 1/2)$ $\{1,\ldots,m\},\ \mathcal{H}_i=\mathcal{H}\ and\ there\ exists\ \beta_i\in]0,1[\ such\ that\ \|W_i-2(1-\beta_i)\ \mathrm{Id}\ \|+1\}$ $||W_i|| \leq 2\beta_i$.

Then, for every $i \in \{1, ..., m-1\}$, $(x_{i,n})_{n \in \mathbb{N}}$ converges weakly to $\overline{x}_i = (T_i \circ \cdots \circ T_1)\overline{x}_m$ and $(\overline{x}_1, \ldots, \overline{x}_m)$ solves (4.5).

Proof Let $(y_n)_{n\in\mathbb{N}}$ be the sequence defined by $y_0=x_0$ and

for
$$n = 0, 1, ...$$

$$\begin{vmatrix} y_{1,n} = T_1 y_n \\ y_{2,n} = T_2 y_{1,n} \\ \vdots \\ y_{m,n} = T_m y_{m-1,n} \\ y_{n+1} = y_n + \lambda_n (y_{m,n} - y_n). \end{vmatrix}$$
(5.7)

For every $n \in \mathbb{N}$, set $S_n = T_{m,n} \circ \cdots \circ T_{1,n}$. We derive from (1.2) and (5.7) that

$$(\forall n \in \mathbb{N}) \quad \|x_{n+1} - y_{n+1}\| = \|x_n + \lambda_n (S_n x_n - x_n) - y_n - \lambda_n (T y_n - y_n)\|$$

$$\leq \lambda_n \|S_n x_n - T x_n\| + \|x_n - y_n + \lambda_n (T x_n - T y_n - x_n + y_n)\|.$$
(5.8)

At the same time, by Proposition 5.3, there exist $(\tau_{m,n})_{n\in\mathbb{N}}\in\ell_+^1$ and $(\theta_{m,n})_{n\in\mathbb{N}}\in\ell_+^1$ such that

$$(\forall n \in \mathbb{N}) \quad ||S_n x_n - T x_n|| \leqslant \tau_{m,n} ||x_n|| + \theta_{m,n}$$

$$\leqslant \tau_{m,n} (||x_n - y_n|| + ||y_n||) + \theta_{m,n}.$$
(5.9)



On the other hand, by Theorem 3.8, Assumption 5.1(ii), and (b), T is α -averaged. Hence, there exists a nonexpansive operator $Q: \mathcal{H} \to \mathcal{H}$ such that $T = (1 - \alpha) \operatorname{Id} + \alpha Q$. Since (c) implies that $(\lambda_n)_{n \in \mathbb{N}}$ lies in $]0, 1/\alpha]$, we deduce that

$$(\forall n \in \mathbb{N}) \quad \|x_{n} - y_{n} + \lambda_{n} (Tx_{n} - Ty_{n} - x_{n} + y_{n})\|$$

$$= \|(1 - \alpha \lambda_{n})(x_{n} - y_{n}) + \alpha \lambda_{n} (Qx_{n} - Qy_{n})\|$$

$$\leq (1 - \alpha \lambda_{n}) \|x_{n} - y_{n}\| + \alpha \lambda_{n} \|Qx_{n} - Qy_{n}\|$$

$$\leq \|x_{n} - y_{n}\|.$$
(5.11)

Altogether (5.8), (5.10), and (5.11) yield

$$(\forall n \in \mathbb{N}) \quad \|x_{n+1} - y_{n+1}\| \le \left(1 + \frac{\tau_{m,n}}{\alpha}\right) \|x_n - y_n\| + \frac{1}{\alpha} \left(\tau_{m,n} \|y_n\| + \theta_{m,n}\right). \tag{5.12}$$

However, Theorem 4.7 guarantees that $\delta = \sup_{n \in \mathbb{N}} \|y_n\| < +\infty$ and therefore that

$$(\forall n \in \mathbb{N}) \quad \|x_{n+1} - y_{n+1}\| \le \left(1 + \frac{\tau_{m,n}}{\alpha}\right) \|x_n - y_n\| + \frac{1}{\alpha} \left(\tau_{m,n} \delta + \theta_{m,n}\right). \tag{5.13}$$

Since $(\tau_{m,n})_{n\in\mathbb{N}}$ and $(\tau_{m,n}\delta + \theta_{m,n})_{n\in\mathbb{N}}$ are in ℓ_+^1 , there exists $\nu \in [0, +\infty[$ such that $\|x_n - y_n\| \to \nu$ [8, Lemma 5.31]. Consequently, $\delta' = \sup_{n\in\mathbb{N}} \|x_n\| \leqslant \delta + \sup_{n\in\mathbb{N}} \|x_n - y_n\| < +\infty$. Now, set

$$(\forall n \in \mathbb{N}) \quad e_n = \frac{1}{\alpha} (S_n x_n - T x_n). \tag{5.14}$$

Then it follows from (5.9) that

$$\sum_{n\in\mathbb{N}} \|e_n\| \leqslant \frac{1}{\alpha} \sum_{n\in\mathbb{N}} \left(\tau_{m,n} \|x_n\| + \theta_{m,n} \right) \leqslant \frac{\delta'}{\alpha} \sum_{n\in\mathbb{N}} \tau_{m,n} + \frac{1}{\alpha} \sum_{n\in\mathbb{N}} \theta_{m,n} < +\infty. \tag{5.15}$$

In view of (1.2), we have

$$(\forall n \in \mathbb{N})$$
 $x_{n+1} = x_n + \mu_n(Qx_n + e_n - x_n)$, where $\mu_n = \alpha \lambda_n \in]0, 1[$. (5.16)

(i): The weak convergence of $(x_n)_{n\in\mathbb{N}}$ to a point $\overline{x}_m \in \text{Fix } Q = F$ follows from (5.16) and [8, Theorem 5.33(iv)] by arguing as in the proof of [8, Theorem 5.14(i)].

(ii): It follows from (5.16) that $\sum_{n\in\mathbb{N}} \mu_n(1-\mu_n) = +\infty$. Hence [8, Proposition 5.34(iii)] implies that $(x_n)_{n\in\mathbb{N}}$ converges weakly to a point $\overline{x}_m \in \text{Fix } Q = F$.

In (i)–(ii) above, Proposition 4.3(iii) ensures that $(T_1\overline{x}_m, (T_2 \circ T_1)\overline{x}_m, \dots, (T_{m-1} \circ \cdots \circ T_1)\overline{x}_m, \overline{x}_m)$ solves (4.5).

(iii)–(v): If one of these assumptions holds, by proceeding as in the proof of Theorem 4.7(iii)–(v), we obtain that, for every $i \in \{1, \ldots, m-1\}$, $(T_i \circ \cdots \circ T_1)x_n \to \overline{x}_i = (T_i \circ \cdots \circ T_1)\overline{x}_m$ and that, furthermore, $(\overline{x}_1, \ldots, \overline{x}_m)$ solves (4.5). However, Proposition 5.3 asserts that, for every $i \in \{1, \ldots, m-1\}$, there exist $(\tau_{i,n})_{n \in \mathbb{N}} \in \ell^1_+$ and $(\theta_{i,n})_{n \in \mathbb{N}} \in \ell^1_+$ such that, for every $n \in \mathbb{N}$,

$$||x_{i,n} - (T_i \circ \cdots \circ T_1)x_n|| = ||(T_{i,n} \circ \cdots \circ T_{1,n})x_n - (T_i \circ \cdots \circ T_1)x_n|| \leqslant \tau_{i,n}||x_n|| + \theta_{i,n}.$$
 (5.17)

Since $(x_n)_{n\in\mathbb{N}}$ is bounded, $x_{i,n}-(T_i\circ\cdots\circ T_1)x_n\to 0$ and therefore $x_{i,n}\to \overline{x}_i$. (vi): For every $i\in\{1,\ldots,m\}$, set

$$(\forall n \in \mathbb{N}) \quad e_{i,n} = (T_{i,n} \circ T_{i-1,n} \circ \cdots \circ T_{1,n}) x_n - (T_i \circ T_{i-1,n} \circ \cdots \circ T_{1,n}) x_n, \quad (5.18)$$

and let $(\chi_{i,n})_{n\in\mathbb{N}}$ and $(\zeta_{i,n})_{n\in\mathbb{N}}$ be defined as in (5.2). By Propositions 4.2, 5.2, and 5.3, we have

$$(\forall n \in \mathbb{N}) \quad \|e_{1,n}\| \leqslant \chi_{1,n} \|x_n\| + \zeta_{1,n} \tag{5.19}$$

and

$$(\forall i \in \{2, \dots, m\})(\exists (\tau_{i-1,n})_{n \in \mathbb{N}} \in \ell_{+}^{1})(\exists (\theta_{i-1,n})_{n \in \mathbb{N}} \in \ell_{+}^{1})(\forall n \in \mathbb{N})$$

$$\|e_{i,n}\| \leq \chi_{i,n}\|(T_{i-1,n} \circ \cdots \circ T_{1,n})x_{n}\| + \zeta_{i,n}$$

$$\leq \chi_{i,n}(\|(T_{i-1,n} \circ \cdots \circ T_{1,n})x_{n} - (T_{i-1} \circ \cdots \circ T_{1})x_{n}\| + \|(T_{i-1} \circ \cdots \circ T_{1})x_{n}\|) + \zeta_{i,n}$$

$$\leq \chi_{i,n}\left(\tau_{i-1,n}\|x_{n}\| + \theta_{i-1,n} + \|x_{n}\| \prod_{k=1}^{i-1} \|W_{k}\| + \sum_{j=1}^{i-1} \|b_{j}\| \left(\prod_{k=j+1}^{i-1} \|W_{k}\|\right)\right) + \zeta_{i,n}.$$

$$(5.20)$$

Thus, since $(x_n)_{n\in\mathbb{N}}$ is bounded,

$$(\forall i \in \{1, \dots, m\}) \quad (\|e_{i,n}\|)_{n \in \mathbb{N}} \in \ell_+^1.$$
 (5.21)

In addition, by (5.18) and (1.2),

$$(\forall n \in \mathbb{N}) \ x_{n+1} = x_n + \lambda_n \left(T_m (T_{m-1} (\cdots T_2 (T_1 x_n + e_{1,n}) + e_{2,n} \cdots) + e_{m-1,n}) + e_{m,n} - x_n \right).$$

$$(5.22)$$

Thus, since Proposition 3.6(iii) and Theorem 3.8 imply that the operators $(T_i)_{1 \le i \le m}$ are averaged, the proof can be completed as that of Theorem 4.7(vi) since [24, Theorem 3.5(ii)] asserts that (4.15) remains valid under (5.21).

References

- Aragón Artacho, F.J., Campoy, R.: A new projection method for finding the closest point in the intersection of convex sets. Comput. Optim. Appl. 69, 99–132 (2018)
- Attouch, H., Peypouquet, J., Redont, P.: Backward-forward algorithms for structured monotone inclusions in Hilbert spaces. J. Math. Anal. Appl. 457, 1095–1117 (2018)
- 3. Baillon, J.-B., Bruck, R.E., Reich, S.: On the asymptotic behavior of nonexpansive mappings and semigroups in Banach spaces. Houston J. Math. 4, 1–9 (1978)
- Baillon, J.-B., Combettes, P.L., Cominetti, R.: There is no variational characterization of the cycles in the method of periodic projections. J. Funct. Anal. 262, 400–408 (2012)
- Bargetz, C., Reich, S., Zalas, R.: Convergence properties of dynamic string-averaging projection methods in the presence of perturbations. Numer. Algor. 77, 185–209 (2018)
- Barron, A.R.: Universal approximation bounds for superpositions of a sigmoidal function. IEEE Trans. Inform. Theory 39, 930–941 (1993)
- Bauschke, H.H., Borwein, J.M.: On projection algorithms for solving convex feasibility problems. SIAM Rev. 38, 367–426 (1996)
- 8. Bauschke, H.H., Combettes, P.L.: Convex Analysis and Monotone Operator Theory in Hilbert Spaces, 2nd edn. Springer, New York (2017)
- Bauschke, H.H., Noll, D., Phan, H.M.: Linear and strong convergence of algorithms involving averaged nonexpansive operators. J. Math. Anal. Appl. 421, 1–20 (2015)
- Bilski, J.: The backpropagation learning with logarithmic transfer function. In: Proc. 5th Conf. Neural Netw. Soft Comput., pp. 71–76 (2000)
- 11. Borwein, J.M., Li, G., Tam, M.K.: Convergence rate analysis for averaged fixed point iterations in common fixed point problems. SIAM J. Optim. 27, 1–33 (2017)
- Borwein, J., Reich, S., Shafrir, I.: Krasnoselski-Mann iterations in normed spaces. Canad. Math. Bull. 35, 21–28 (1992)
- Boţ, R.I., Csetnek, E.R.: A dynamical system associated with the fixed points set of a nonexpansive operator. J. Dynam. Diff. Equ. 29, 155–168 (2017)
- Bravo, M., Cominetti, R.: Sharp convergence rates for averaged nonexpansive maps. Israel J. Math. 227, 163–188 (2018)
- Bridle, J.S.: Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition. In: Neurocomputing, NATO ASI Series, Series F, vol. 68, pp. 227–236. Springer, Berlin (1990)



- Carlile, B., Delamarter, G., Kinney, P., Marti, A., Whitney, B.: Improving deep learning by inverse square root linear units (ISRLUs). https://arxiv.org/abs/1710.09967 (2017)
- Cegielski, A.: Iterative Methods for Fixed Point Problems in Hilbert Spaces. Lecture Notes in Mathematics, vol. 2057. Springer, Heidelberg (2012)
- Censor, Y., Mansour, R.: New Douglas–Rachford algorithmic structures and their convergence analyses. SIAM J. Optim. 26, 474–487 (2016)
- Combettes, P.L.: Construction d'un point fixe commun à une famille de contractions fermes. C. R. Acad. Sci. Paris Sér. I Math., 320, 1385–1390 (1995)
- Combettes, P.L.: Solving monotone inclusions via compositions of nonexpansive averaged operators. Optimization 53, 475–504 (2004)
- Combettes, P.L.: Monotone operator theory in convex optimization. Math. Programming B170, 177–206 (2018)
- 22. Combettes, P.L., Pesquet, J.-C.: Proximal thresholding algorithm for minimization over orthonormal bases. SIAM J. Optim. 18, 1351–1376 (2007)
- Combettes, P.L., Wajs, V.R.: Signal recovery by proximal forward-backward splitting. Multiscale Model. Simul. 4, 1168–1200 (2005)
- Combettes, P.L., Yamada, I.: Compositions and convex combinations of averaged nonexpansive operators. J. Math. Anal. Appl. 425, 55–70 (2015)
- Condat, L.: A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. J. Optim. Theory Appl. 158, 460–479 (2013)
- Cybenko, G.: Approximation by superposition of sigmoidal functions. Math. Control Signals Syst. 2, 303–314 (1989)
- Eckstein, J., Bertsekas, D.P.: On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. Math. Program. 55, 293–318 (1992)
- Elliot, D.L.: A better activation function for artificial neural networks, Institute for Systems Research, University of Maryland, Tech. Rep., pp. 93–8 (1993)
- Funahashi, K.-I.: On the approximate realization of continuous mappings by neural networks. Neural Netw. 2, 183–192 (1989)
- Glorot, X., Bordes, A., Bengio, Y.: Deep sparse rectifier neural networks. In: Proc. 14th Int. Conf. Artificial Intell. Stat., pp. 315–323 (2011)
- 31. Haykin, S. Neural Networks: A Comprehensive Foundation, 2nd edn. Pearson Education, Singapore (1998)
- 32. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proc. Int. Conf. Comput. Vision, pp. 1026–1034 (2015)
- 33. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. IEEE Conf. Comput. Vision Pattern Recogn., pp. 770–778 (2016)
- 34. LeCun, Y.A., Bengio, Y., Hinton, G.: Deep learning. Nature **521**, 436–444 (2015)
- LeCun, Y.A., Bottou, L., Orr, G.B., Müller, K.-R.: Efficient backprop. Lect. Notes Comput. Sci. 1524, 9–50 (1998)
- Martinet, B.: Détermination approchée d'un point fixe d'une application pseudo-contractante. Cas de l'application prox. C. R. Acad. Sci. Paris A274, 163–165 (1972)
- McCulloch, W.S., Pitts, W.H.: A logical calculus of the ideas immanent in nervous activity. Bull. Math. Biophys. 5, 115–133 (1943)
- Moursi, W.M.: The forward-backward algorithm and the normal problem. J. Optim. Theory Appl. 176, 605–624 (2018)
- Nair, V., Hinton, G.E.: Rectified linear units improve restricted Boltzmann machines. In: Proc. 27st Int. Conf. Machine Learn., pp. 807–814 (2010)
- 40. Rockafellar, R.T.: Convex Analysis. Princeton University Press, Princeton (1970)
- Rockafellar, R.T.: Monotone operators and the proximal point algorithm. SIAM J. Control Optim. 14, 877–898 (1976)
- 42. Rosenblatt, F.: The perceptron: A probabilistic model for information storage and organization in the brain. Psychological Rev. **65**, 386–408 (1958)
- Ryu, E.K., Hannah, R., Yin, W.: Scaled relative graph: Nonexpansive operators via 2D Euclidean geometry. https://arxiv.org/abs/1902.09788
- Srivastava, R.K., Greff, K., Schmidhuber, J.: Training very deep networks. Proc. Neural Inform. Process. Syst. Conf. 28, 2377–2385 (2015)
- Tariyal, S., Majumdar, A., Singh, R., Vatsa, M.: Deep dictionary learning. IEEE Access 4, 10096–10109 (2016)
- Tseng, P.: On the convergence of products of firmly nonexpansive mappings. SIAM J. Optim. 2, 425–434 (1992)



- 47. Yamagishi, M., Yamada, I.: Nonexpansiveness of a linearized augmented Lagrangian operator for hierarchical convex optimization. Inverse Problems, vol. 33, art. 044003, 35 pp. (2017)
- 48. Zhang, X.-P.: Thresholding neural network for adaptive noise reduction. IEEE Trans. Neural Netw. 12, 567–584 (2001)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

