

Discretization of parameter identification in PDEs using Neural Networks

Barbara Kaltenbacher, University of Klagenfurt
and
Tram Thi Ngoc Nguyen, University of Graz

Abstract. We consider the ill-posed inverse problem of identifying a nonlinearity in a time-dependent PDE model. The nonlinearity is approximated by a neural network, and needs to be determined alongside other unknown physical parameters and the unknown state. Hence, it is not possible to construct input-output data pairs to perform a supervised training process. Proposing an all-at-once approach, we bypass the need for training data and recover all the unknowns simultaneously. In the general case, the approximation via a neural network can be realized as a discretization scheme, and the training with noisy data can be viewed as an ill-posed inverse problem. Therefore, we study discretization of regularization in terms of Tikhonov and projected Landweber methods for discretization of inverse problems, and prove convergence when the discretization error (network approximation error) and the noise level tend to zero.

Key words: Neural networks, unsupervised learning, discretization of regularization, parameter identification, nonlinear PDEs, Tikhonov regularization, Landweber iteration.

1 Introduction

Parameter identification in partial differential equations (PDEs) from indirect observation is a category of inverse problems that arises in numerous applications, such as medical imaging, geophysical prospection and nondestructive testing.

In this paper, we focus on transient models and the appearance of unknown nonlinearities. In order to find a finite dimensional representation of the latter, we make use of the powerful approximation properties and computational efficiency of *neural networks* (*NNs*). Due to the inherent ill-posedness of these inverse problems, however, regularization must be employed. We therefore study two such regularization methods: the variational Tikhonov method and the iterative projected Landweber method. These reconstruction methods are analyzed in the spirit of regularization theory, with a discretization by neural networks, as well as other general discretization schemes. We must therefore investigate the interplay between noise level, regularization parameter and discretization error (approximation error in case of discretization by NNs). In the language of machine learning, this

is the interplay between approximation error and optimization/estimation error, with the impact of ill-posedness and data noise additionally taken into consideration. The resulting convergence analysis hints at a dependence of the network size (discretization parameter) on the noise level. This constitutes one of the main contributions of the paper.

We point to the fact that a regularization theoretical viewpoint for the training problems has already been taken in [5], although there the focus is on linear problems. In [5], the authors solve $Au = y$ without making use of the linear forward map A , relying solely on the input-output training pairs $[u_m, y_m]_{m=1\dots M}$ satisfying $Au_m = y_m$. This data-driven approach is interpreted as regularization by projection, where the subspaces are spanned by the training data. Along this line, [13] investigates the supervised training problem of approximating a smooth function via one-layer feed-forward networks with noisy data as an ill-posed problem. This is shown to be equivalent to least-squares collocation for a linear integral equation. The core result is the derivation of an optimal choice of the network size depending upon on the data error δ . In the same spirit, our work focuses on the connection between machine learning and regularization. Our objective is to establish a convergence analysis of regularization methods under the influence of the network approximation error in a nonlinear PDE. As the considered problem already exhibits multi-faceted complexity, namely parameter identification, a nonlinear model and unsupervised training, the task of deriving convergence rates as in [13] is deferred to future research.

Our use of neural networks in parameter identification is inspired by [18]. There, the focus is on stationary problems, and the nonlinearity is represented by a neural network. In [18], the network is learned beforehand via supervised training, and then it is inserted into the PDE model underlying the parameter identification. The supervised learning thus requires exact and full measurements of the state u , as well as physical parameters of the PDE, in order to form the training pairs. In contrast to this, we here consider time-dependent models; in addition, by virtue of our *all-at-once formulation*, the supervised training is skipped, hence no access to the exact state and physical parameters is involved. The application of an all-at-once approach for unsupervised learning is another main contribution of our paper.

Another advantage of the all-at-once formulation lies in the fact that it avoids the evaluation of a parameter-to-state map, thus bypassing the need for a nonlinear PDE solver in a practical implementation. Additionally, this setting simplifies verification of the so-called tangential cone condition, a requirement for convergence guarantees of gradient-based methods, whose verification in many applications is neglected. This, may be considered another advantage of our approach.

This work is a continuation of [1], in which we also parametrize the unknown nonlinearities in time-dependent PDEs by NNs. There, a so-called *learning-informed parameter identification* was investigated by way of discretized inverse problems (i.e. when f is already approximated by some NN). The present study develops a theoretical framework for [1], in the sense that we show convergence of the regularized and discretized reconstructions towards a ground truth. The approximation/discretization is incorporated into the regularization, which places a stronger emphasis on regularization theory in the light of existing literature. The analysis applies not only to discretization by NNs, but also to more

general discretization schemes.

The field of deep learning for PDEs is well developed, with many novel results and techniques available in the literature. One such technique is *physics informed neural networks* (PINNs) [47], where one parametrizes the solution to the PDE, as opposed to using NNs to parametrize the unknown nonlinearities as in our case. A theoretical justification for using NNs to parametrize PDE solutions or parameter-to-solution maps can be found in [35].

Recovery of hidden physics laws from empirical observations is, in fact, an active field with a significant history. Recently, the rapid advances in computing power and data acquisition open the door for advanced techniques. For example [10, 49] are concerned with the recovery of the governing PDE from full measurements of the state u . These two papers suggest to first construct a rich library of possible basis elements and then optimize the corresponding coefficients using sparse regression. Adding deep learning techniques, [46] proposes to use two deep neural networks, one representing the solution u , the other representing the nonlinear dynamics $f := u_t - \mathcal{N}(t, x, u, \nabla u, \nabla u^2 \dots)$. Algorithmic differentiation is employed for computing the required derivatives. On the other hand, PDE-NET [38] represents another flexible framework, in which one approximates the model f by a feed-forward NN, while numerically approximating the differential operators $\nabla, \nabla^2 \dots$ by convolutional NNs. In all these mentioned studies, the problems are studied in a discrete setting and the collocation points (t_i, x_i) range over the entire time and space $[0, T] \times \Omega$. In this work, we put more emphasis on the aspect of the PDE model being derived from physical laws, and use data-driven methods solely to complement this approach. Further, our study combines a functional setting for the unknown physical parameters and states with a network parametrization for the unknown nonlinearity in an hybrid form.

While [5, 13, 18, 47, 10, 49, 46, 38, 1] are the recent publications that our work is most closely related to, there is clearly a vast amount of existing and emerging literature on the mathematics of machine learning. In the context of data-driven inverse problems, we refer to [4] for an excellent review. For a profound exposition on the theory of deep learning, we refer to the lecture series and associated upcoming publication [22].

1.1 The inverse problem

Quite often, the nonlinearity is not the only unknown quantity, but rather must be determined alongside other coefficients in the PDE, as exemplified in the following application.

Application. Consider the problem of recovering the unknown nonlinearity f , the potential c , the source φ , and the initial data u_0 in

$$\begin{aligned} \dot{u} - \Delta u + cu + h(u) &= f(u) + \varphi && \text{in } \Omega \times (0, T) \\ u(0) &= u_0 && \text{on } \Omega, \end{aligned} \tag{1}$$

from measurements y of the state u . The state u is a function on the finite time line $(0, T)$ and the bounded, smooth domain Ω with its time derivative being denoted by \dot{u} . In (1),

h is the known nonlinear part of the model; the unknown nonlinear part f , which needs to be determined, plays the role of a model correction, thus helping to refine the physical model. The additional data y available to identify the unknown quantities are observations of the state u expressed via some observation operator M (which could e.g., be the trace of u at the boundary over time or its values at some fixed times instance(s) in Ω)

$$Mu = y. \quad (2)$$

In more complex settings, for the purpose of identifying the unknown functions, several repeated or possibly also different observations $y^m = M^m u^m$, $m = 1 \dots K$ will be needed. These observations entail a variation in the data, and possibly also in some of the unknown coefficients, while the nonlinearity remains the same. Different measured data y^m correspond to the model (1) at different parameters, thus at different states, i.e. $(c, \varphi, u_0)^m, u^m, m = 1 \dots K$ may vary between observations, while the unknown function f describing the underlying physical law is fixed. There are several real life inverse problems obeying this setting. In medical imaging MRI, this situations appears when different patients are scanned, resulting in K sets of patient-dependent physical/body parameters. These patient-specific parameters, however, enter the same model governed by the same underlying physical law, e.g. the Bloch-Torrey equation model [8]. Thus, the unknown nonlinear response f can be considered as being fixed, while relaxation and diffusion parameters are allowed to vary between patients.

Remark 1 (uniqueness). *In the context of restricted measurements (2) such as boundary observations, as relevant in tomographic applications, the question arises whether f can be determined uniquely from these observations. Answers to this question can be found in the literature in the case of unknown f and known initial data u_0 as well as coefficients c, φ , see, e.g., [15, 19, 44], or in the case of known f and unknown initial data u_0 or coefficients c, φ , see, e.g., [26, 27], but more rarely on simultaneous identifiability of all these quantities. As the unknown f case is the most relevant for our study, we point to the fact that a range condition on the exact state u_{exact}*

$$u_{\text{exact}}(\Omega \times (0, T)) \subseteq u_{\text{exact}}(\omega \times (0, T))$$

is essential for establishing unique recovery of $f(u)$ from observations of u in some subset ω of Ω or its boundary.

We will study such inverse problems in a more general framework of the following form.

Inverse Problem. Before stating the inverse problem, we point out that in the model (5) below, the unknown nonlinearity $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ is identified with the corresponding *Nemitskii operator* (cf. Section 1.2)

$$f : \mathbb{R}^n \times \mathcal{V} \rightarrow \mathcal{W} \quad \text{via} \quad [f(\alpha, u)](x, t) := f(\alpha, u(x, t)), \quad (3)$$

and similarly, the known part $F : (0, T) \times X \times V \rightarrow W$ is identified with the Nemitskii operator

$$F : X \times \mathcal{V} \rightarrow \mathcal{W} \quad \text{via} \quad [F(\lambda, u)](t) := F(t, \lambda, u(t)), \quad (4)$$

with \mathcal{V} denoting the state space and \mathcal{W} denoting the image space of the model. Here, $\mathcal{V} \subseteq L^2(0, T; V)$ and $\mathcal{W} \subseteq L^2(0, T; W)$ are *Bochner spaces* (cf. Section 1.2) and with V being a space of x dependent functions, in (3) we make the identification $u(x, t) = (u(t))(x)$.

We now investigate the inverse problem of determining the physical parameters $\lambda^m \in X$ (parameter space), $u_0^m \in H$ (initial data space), $\alpha^m \in \mathbb{R}^n$, $m \in \{1, \dots, K\}$, and the nonlinearity $f \in \mathcal{C} \subseteq C(\mathbb{R}^{n+1}; \mathbb{R})$, in the evolution system

$$\begin{aligned} \dot{u} &= F(\lambda, u) + f(\alpha, u) \quad \text{in } (0, T) \\ u(0) &= u_0 \end{aligned} \quad (5)$$

from K noisy measurements $y^{m, \delta} \in \mathcal{Y}^m$ (data space) of the states $u^m \in \mathcal{V}$ (state space) under the measurement operator M^m according to

$$\begin{aligned} y^m &= M^m u^m, \quad M^m : \mathcal{V} \rightarrow \mathcal{Y}^m \\ 0 &\leq S^m(y^m, y^{m, \delta}) \leq \delta \quad m \in \{1, \dots, K\}, \end{aligned} \quad (6)$$

where u^m solves (5) with $(\lambda, u_0, \alpha) = (\lambda^m, u_0^m, \alpha^m)$. Here, the distance between the exact data y^m and the noisy data $y^{m, \delta}$ under the misfit measure S^m is assumed to be bounded by the noise level δ . Typical choices of S^m are norms (as in Section 3 below) or more general distance measures such as the Kullback-Leibler divergence. If (5) represents an evolutionary PDE such as (1), the spaces V, W, X, H are typically Banach spaces of space-dependent functions, and \mathcal{C} is a function space over \mathbb{R}^{n+1} that will be approximated by neural networks later on. The data spaces \mathcal{Y}^m are Banach spaces as well and, depending on what type of observations are made, may consist of space and/or time dependent functions.

This parameter identification problem can equivalently be written as the *all-at-once system* (cf. Section 1.2)

$$\left(\begin{array}{c} E(\zeta^m, u^m, f) \\ M^m u^m \end{array} \right)_{m=1}^K = \left(\begin{array}{c} 0 \\ y^m \end{array} \right)_{m=1}^K, \quad \zeta^m := (\lambda^m, u_0^m, \alpha^m) \in Z := X \times H \times \mathbb{R}^n \quad (7)$$

for the parameters ζ^m , the nonlinearity f , and the states u^m , with

$$E(\zeta, u, f) := \left(\begin{array}{c} \dot{u} - F(\lambda, u) - f(\alpha, u) \\ u(0) - u_0 \end{array} \right) \in \mathcal{W} \times H. \quad (8)$$

We denote by $(\zeta^\dagger, u^\dagger, f^\dagger)$ an exact solution to the inverse problem, that is

$$E(\zeta^\dagger, u^\dagger, f^\dagger) = 0, \quad M u^\dagger = y,$$

leading to a vanishing PDE residual and a perfect match of the measurements to the noise free data y .

In particular, the application (1) is a special case of (5) in the setting (7) with

$$K = 1, \quad n = 0, \quad \lambda = (c, \varphi), \quad F(\lambda, u) = \Delta u - cu - h(u) + \varphi. \quad (9)$$

1.2 Preliminaries

Before launching a detailed discussion of discretization for Tikhonov and Landweber regularization, we briefly elaborate on some concepts that have been mentioned in the preceding section.

Bochner spaces. Given a Banach space V , the *Bochner space* $L^p(0, T; V)$ [48, Section 1.5] consists of the Bochner integrable functions $u : [0, T] \rightarrow V$ satisfying $\int_0^T \|u(t)\|_V^p dt < +\infty$. It is a Banach space under the norm

$$\|u\|_{L^p(0, T; V)} := \left(\int_0^T \|u(t)\|_V^p dt \right)^{1/p} \quad 1 \leq p < \infty.$$

Likewise, the Bochner spaces $L^\infty(0, T; V)$ and $C(0, T; V)$ are Banach spaces under the respective norms

$$\|u\|_{L^\infty(0, T; V)} := \sup_{t \in [0, T]} \|u(t)\|_V, \quad \|u\|_{C(0, T; V)} := \max_{t \in [0, T]} \|u(t)\|_V.$$

Given a convex Banach space V_1 and a locally convex Banach space $V_2 \supset V_1$, we define the *Sobolev-Bochner space* $W^{1,p,q}(0, T; V_1, V_2)$ [48, Section 7.1], which itself is a Banach space, as

$$\begin{aligned} W^{1,p,q}(0, T; V_1, V_2) &:= \{u \in L^p(0, T; V_1) : \dot{u} \in L^q(0, T; V_2)\} \quad 1 \leq p, q \leq \infty, \\ \|u\|_{W^{1,p,q}(0, T; V_1, V_2)} &= \|u\|_{L^p(0, T; V_1)} + \|\dot{u}\|_{L^q(0, T; V_2)}. \end{aligned}$$

An example that is used in Section 3.3 is $W^{1,2,2}(0, T; V_1, V_2) = L^2(0, T; V_1) \cap H^1(0, T; V_2)$.

Nemitskii operators. A mapping $\mathfrak{f} : I \times A \rightarrow B$ with Banach spaces A, B and $I \subset \mathbb{R}^d$ is called a *Caratheodory mapping* if $\mathfrak{f}(\cdot, u)$ is measurable for all $u \in A$ and $\mathfrak{f}(z, \cdot)$ is continuous for a.e. $z \in I$. The so-called *Nemitskii operator* \mathfrak{F} assigns a function $v : I \rightarrow A$ to a function $w : I \rightarrow B$ by

$$[\mathfrak{F}(v)](z) := \mathfrak{f}(z, v(z)).$$

In (3), we have, for any fixed $\alpha \in \mathbb{R}^n$, $\mathfrak{f} = f(\alpha, \cdot)$, $v = u$, $I = \Omega \times (0, T)$, $A = B = \mathbb{R}$, while in (4), we have for any fixed $\lambda \in X$, $\mathfrak{f} = F(\lambda, \cdot)$, $v = u$, $I = (0, T)$, $A = V$, $B = W$.

For a detailed discussion on Nemitskii operators in Bochner spaces, we refer to [48, Sections 1.3, 1.4].

All-at-once formulation. The classical way to formulate the inverse problem (5)-(6) (for simplicity of exposition setting $K = 1$ and therefore skipping the superscripts m) is to construct the *reduced* forward operator

$$G : Z \times \mathcal{C} \rightarrow \mathcal{Y} \quad G(\zeta, f) := M \circ P(\zeta, f) = y,$$

which composes the observation operator M with the parameter-to-state map

$$P : Z \times \mathcal{C} \rightarrow \mathcal{V} \quad P(\zeta, f) = u, \text{ where } u \text{ solves (5).}$$

This formulation involves evaluating well-definedness of P via unique existence theory for the nonlinear PDE (5), and in practice requires solving this nonlinear equation.

Alternatively, the *all-at-once* approach formulates (5)-(6) into a system

$$\begin{aligned} E(\zeta, u, f) &= 0 \\ Mu &= y \end{aligned}$$

of model and observation equation as in (7)-(8). Hence, we can define the forward operator

$$\mathbf{G} : Z \times \mathcal{V} \times \mathcal{C} \rightarrow \mathcal{W} \times \mathcal{Y} \quad \mathbf{G}(\zeta, u, f) := (E(\zeta, u, f), Mu) = (0, y).$$

The all-at-once formulation bypasses the construction of the parameter-to-state map P , which is nonlinear and often requires restrictive assumptions on F, f . This formulation therefore allows more general classes of F, f , and is also advantageous in practical implementation, where a PDE solver is not needed. All-at-once approaches have been studied in PDE constrained optimization in [33, 34, 36, 55, 42, 51, 52] and more recently, for ill-posed inverse problems, in [11, 12, 23, 28, 29, 55]; a comparison between the reduced and all-at-once formulation for time dependent problems can be found in [31, 41].

Neural networks (NNs). In the setting of this paper, we make use of the *feedforward neural network* of *depth* L on $(\alpha, u(x, t)) \in \mathbb{R}^{n+1}$, expressed as a function of the form

$$\mathcal{N} : \mathbb{R}^{n+1} \rightarrow \mathbb{R} \quad \mathcal{N}(\alpha, u(x, t)) := A_L \circ \dots \circ A_1(\alpha, u(x, t)), \quad A_\ell(\mathbf{z}) := \sigma_\ell(w_\ell \mathbf{z} + b_\ell),$$

where the matrix $w_\ell \in \mathcal{L}(\mathbb{R}^{p_{\ell-1}}, \mathbb{R}^{p_\ell})$ and the vector $b_\ell \in \mathbb{R}^{p_\ell}$ are the so-called *hyperparameters* at *layer* $\ell = 1$ (input) $\dots L$ (output). The *activations* $\sigma_\ell : \mathbb{R} \rightarrow \mathbb{R}$ are nonlinear point-wise functions allowed to differ between layers, and $\sigma_L = \text{Id}$. In summary, at layer $\ell - 1$ the affine operator A_ℓ transforms an input vector in $\mathbb{R}^{p_{\ell-1}}$ into one in \mathbb{R}^{p_ℓ} , applies the activation σ_ℓ pointwise, and returns the input to the next layer ℓ . Some standard activation functions include the RELU function $\sigma(z) = \max\{z, 0\}$, tansig function $\sigma(z) = \tanh(z)$, softsign function $\sigma(z) = \frac{z}{1+|z|}$ and softplus function $\sigma(z) = \ln(1 + e^z)$.

Based on the universal approximation theorem for smooth functions [25], we use standard feedforward neural networks to approximate the nonlinearity $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ in the finite dimensional set \mathcal{C}_N (cf. (11)), whose number of hyperparameters is $N = \sum_{\ell=1}^L (p_{\ell-1} + 1)p_\ell$ with $p_0 = n + 1, p_{L+1} = 1$. Fitting this into the formulation of the inverse problem, we identify f as a Nemitskii operator between Bochner spaces, as introduced in (3).

Notation

We will use shortcut notations $L^2(L^2) = L^2(0, T; L^2(\Omega))$, $C(L^2) = C(0, T; L^2(\Omega))$, $L^2 = L^2(\Omega)$, and analogously for some further Sobolev spaces, when they appear as subscripts in some norms or constants.

We will make use of boundedness of some Sobolev embeddings according to, e.g., [3, Chapter 4], [37, Chapter 11] and more generally denote embedding constants between spaces X and Y by $C_{X \rightarrow Y}$. Generic constants will be denoted by $C > 0$, and continuity or compactness of embeddings is indicated by $X \hookrightarrow Y$ or $X \hookrightarrow\hookrightarrow Y$, respectively.

Partial derivatives are denoted by subscripts, e.g., f_α , f_u , while ordinary or total derivatives by a prime, e.g. f' .

The remainder of this paper is organized as follows. In Sections 2, we prove convergence of Tikhonov regularization with an appropriate choice of the regularization parameter. Section 3 presents convergence results for Landweber regularization with an appropriate stopping index. In both approaches, the discretization level N needs to be chosen too, in order to achieve convergence as the noise level δ tends to zero. The required conditions are thoroughly discussed and interpreted for the particular Application 1.

2 Tikhonov regularization

With positive definite model and data misfit as well as regularization functionals

$$\begin{aligned} Q : \mathcal{W} &\rightarrow [0, \infty] & s.t. & \quad Q(w) = 0 \Leftrightarrow w = 0, \\ S : Y^2 &\rightarrow [0, \infty] & s.t. & \quad S(y_1, y_2) = 0 \Leftrightarrow y_1 = y_2, \\ R_1 : Z \times \mathcal{V} &\rightarrow [0, \infty], & R_2 : \mathcal{C} &\rightarrow [0, \infty], \end{aligned}$$

consider the objective functional T_γ^δ given by

$$T_\gamma^\delta(\vec{\zeta}, \vec{u}, f) := \sum_{m=1}^K \left(Q(E(\zeta^m, u^m, f)) + S(M^m u^m, y^{m,\delta}) + \gamma R_1(\zeta^m, u^m) \right) + \gamma R_2(f).$$

Here, K is the number of parameters and states corresponding to K different observations of the data $y^{m,\delta}$, while f is the common nonlinearity across all experiments. The objective functional T_γ^δ depends on the noise level δ , measured data $y^{m,\delta}$ and the regularization parameter $\gamma > 0$. We then define *regularized approximations* as minimizers of T_γ^δ , that is,

$$(\vec{\zeta}^{\gamma,\delta}, \vec{u}^{\gamma,\delta}, f^{\gamma,\delta}) \in \operatorname{argmin}_{(\zeta^1, u^1), \dots, (\zeta^K, u^K) \in (Z \times \mathcal{V})^K, f \in \mathcal{C}} T_\gamma^\delta(\vec{\zeta}, \vec{u}, f). \quad (10)$$

The unknown nonlinearity f is approximated by NNs, that is, within the finite dimensional set

$$\mathcal{C}_N := \{\text{neural networks on } \mathbb{R}^{n+1} \text{ with } N \text{ parameters}\} \subseteq \mathcal{C}. \quad (11)$$

Denoting by N the discretization parameter, we define *partially discretized regularized approximations* as

$$(\vec{\zeta}^{\gamma,\delta,N}, \vec{u}^{\gamma,\delta,N}, f^{\gamma,\delta,N}) \in \operatorname{argmin}_{(\zeta^1, u^1), \dots, (\zeta^K, u^K) \in (Z \times \mathcal{V})^K, f \in \mathcal{C}_N} T_\gamma^\delta(\vec{\zeta}, \vec{u}, f). \quad (12)$$

In comparison to (10), the discretization parameter N in (12) enters the minimization as another parameter, which needs to be properly controlled. The focus of this section is on deriving a rule for the regularization parameter γ and the discretization parameter N with respect to the noise level δ , such that convergence of the Tikhonov regularization method is guaranteed. For simplicity of exposition, we set $K = 1$, and mention in passing that an alternative way to take into account multiple observations, as opposed to summing over them in the Tikhonov functional, is the use of Kaczmarz methods. That is, implementing a cyclic iteration over the individual observations, see, e.g., [41] for the all-at-once setting relevant here, as well as the references therein.

2.1 Convergence

We now study convergence of the Tikhonov regularized approximations in the sense of regularization, so as $\delta \rightarrow 0$ with an appropriate choice of regularizer parameter $\gamma(\delta)$ and discretization parameter $N(\delta)$.

Assumption 1. *There exist topologies $\tau_1 := \tau_{Z \times \mathcal{V}}$ on $Z \times \mathcal{V}$ and $\tau_2 := \tau_{\mathcal{C}}$ on \mathcal{C} such that the following holds:*

(T1) *sublevel sets of R_1 are $\tau_{Z \times \mathcal{V}}$ compact, and sublevel sets of R_2 are $\tau_{\mathcal{C}}$ compact;*

(T2) *$(Q \circ E, S \circ M)$ is $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ sequentially closed:*

$$\forall (\zeta^j, u^j, f^j, y^j)_{j \in \mathbb{N}} \subseteq Z \times \mathcal{V} \times \mathcal{C} \times \mathcal{Y} :$$

$$\begin{aligned} & \left((\zeta^j, u^j, f^j) \xrightarrow{\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}} (\bar{\zeta}, \bar{u}, \bar{f}) \text{ and } Q(E(\zeta^j, u^j, f^j)) \rightarrow 0 \right. \\ & \quad \left. \text{and } S(Mu^j, y^j) \rightarrow 0 \text{ and } S(y, y^j) \rightarrow 0 \right) \\ & \implies \left(Q(E(\bar{\zeta}, \bar{u}, \bar{f})) = 0 \text{ and } S(M\bar{u}, y) = 0 \right); \end{aligned}$$

(T3) *$(Q \circ E, S \circ M)$, R_1, R_2 are $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ lower semicontinuous:*

$$\forall (\zeta^j, u^j, f^j)_{j \in \mathbb{N}} \subseteq Z \times \mathcal{V} \times \mathcal{C} :$$

$$\begin{aligned} & (\zeta^j, u^j, f^j) \xrightarrow{\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}} (\bar{\zeta}, \bar{u}, \bar{f}) \\ & \implies \left(Q(E(\bar{\zeta}, \bar{u}, \bar{f})) \leq \liminf_{j \rightarrow \infty} Q(E(\zeta^j, u^j, f^j)) \text{ and } S(M\bar{u}, y^\delta) \leq \liminf_{j \rightarrow \infty} S(Mu^j, y^\delta) \right. \\ & \quad \left. \text{and } R_1(\bar{\zeta}, \bar{u}) \leq \liminf_{j \rightarrow \infty} R_1(\zeta^j, u^j) \text{ and } R_2(\bar{f}) \leq \liminf_{j \rightarrow \infty} R_2(f^j) \right). \end{aligned}$$

The choice of R_1 and R_2 is dictated by the continuity requirements (T2) and (T3). Accordingly, the topology $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ needs to be sufficiently strong. This topology is then linked to R_1, R_2 via the constraint on compactness of the sublevel sets expressed in (T1). Overall, these assumptions are thus criteria to choose the regularizers R_1 and R_2 . The latter two requirements in (T3) are automatically satisfied if R_1, R_2 are defined by norms on $Z \times \mathcal{V}$ and \mathcal{C} , provided the spaces are reflexive or duals of separable spaces and τ_i is defined by the corresponding weak(*) topology.

Proposition 1. *Under the assumptions (T1)-(T3), the discrete minimization problems (12) admit minimizers.*

Proof. The proof follows from standard results [20, 54] that essentially assume compactness of sublevel sets of $R_i, i = 1, 2$, $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ closedness of $(Q \circ E, S \circ M)$, and lower semicontinuity of T_γ^δ .

◇

Thus, under these assumptions, the method is well defined by (12). In order to prove that it actually defines a convergent regularization method, we need further assumptions on the approximation quality and on the choice of the regularization and discretization parameters.

Moreover, we allow for *inexact minimization* by introducing the tolerance $\eta \geq 0$ in the relaxed definition

$$T_\gamma^\delta(\vec{\zeta}^{\gamma, \delta, N}, \vec{u}^{\gamma, \delta, N}, f^{\gamma, \delta, N}) \leq T_\gamma^\delta(\vec{\zeta}, \vec{u}, f) + \eta \quad \forall (\zeta^1, u^1), \dots, (\zeta^K, u^K) \in (Z \times \mathcal{V})^K, f \in \mathcal{C}_N. \quad (13)$$

This definition actually does not even require existence of a minimizer.

Assumption 2. .

(T4) *approximation by NNs:*

$$\begin{aligned} q_N &:= \inf_{f_N \in \mathcal{C}_N} Q(E(\zeta^\dagger, u^\dagger, f_N)) + \gamma_N(R_2(f_N) - R_2(f^\dagger)) \\ &= \inf_{f_N \in \mathcal{C}_N} (Q(E(\zeta^\dagger, u^\dagger, f_N)) - Q(E(\zeta^\dagger, u^\dagger, f^\dagger))) + \gamma_N(R_2(f_N) - R_2(f^\dagger)) \rightarrow 0, \\ &\text{as } N \rightarrow \infty; \end{aligned}$$

(T5) *asymptotics of the parameters as $\delta \rightarrow 0$: There exists $C > 0$ such that*

$$\gamma(\delta) \rightarrow 0, \quad \frac{\delta}{\gamma(\delta)} \leq C, \quad \frac{q_N(\delta)}{\gamma(\delta)} \leq C, \quad \frac{\eta}{\gamma(\delta)} \leq C, \quad N(\delta) \rightarrow \infty, \quad \eta(\delta) \rightarrow 0 \quad \text{as } \delta \rightarrow 0$$

Proposition 2. *Under Assumptions 1, 2 we have $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ subsequential convergence of $(\zeta^{\gamma(\delta), \delta, N(\delta)}, u^{\gamma(\delta), \delta, N(\delta)}, f^{\gamma(\delta), \delta, N(\delta)})$ to a solution of the inverse problem (7), (8) as $\delta \rightarrow 0$, i.e., every sequence $(\zeta^{\gamma(\delta_j), \delta, N(\delta_j)}, u^{\gamma(\delta_j), \delta, N(\delta_j)}, f^{\gamma(\delta_j), \delta, N(\delta_j)})$ with $\delta_j \rightarrow 0$ as $j \rightarrow \infty$ has a $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ convergent subsequence, and the limit of every $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ convergent subsequence solves the inverse problem.*

Note that unlike [45] we assume convergence $q_N \rightarrow 0$ of the “discretization error” to zero only at the exact solution, not uniformly over all elements of $Z \times \mathcal{V} \times \mathcal{C}$. Also, it is not necessary to assume any vector space structure on \mathcal{C}_N and S does not need to satisfy a triangle inequality.

Proof. By minimality, that is, $T_\gamma^\delta(\zeta^{\gamma,\delta,N}, u^{\gamma,\delta,N}, f^{\gamma,\delta,N}) \leq T_\gamma^\delta(\zeta, u, f_N) + \eta$ for all $\zeta \in Z$, $u \in \mathcal{V}$, $f_N \in \mathcal{C}_N$, setting $\zeta = \zeta^\dagger$, $u = u^\dagger$, and thus $Q(E(\zeta^\dagger, u^\dagger, f^\dagger)) = 0$, $Mu^\dagger = y$, shows that

$$\begin{aligned}
& T_\gamma^\delta(\zeta^{\gamma,\delta,N}, u^{\gamma,\delta,N}, f^{\gamma,\delta,N}) \\
& \leq \inf_{f_N \in \mathcal{C}_N} \left(Q(E(\zeta^\dagger, u^\dagger, f_N)) + S(Mu^\dagger, y^\delta) + \gamma R_1(\zeta^\dagger, u^\dagger) + \gamma R_2(f_N) + \eta \right) \\
& \leq Q(E(\zeta^\dagger, u^\dagger, f^\dagger)) + S(y, y^\delta) + \gamma R_1(\zeta^\dagger, u^\dagger) + \gamma R_2(f^\dagger) + \eta \\
& \quad + \inf_{f_N \in \mathcal{C}_N} \left((Q(E(\zeta^\dagger, u^\dagger, f_N)) - Q(E(\zeta^\dagger, u^\dagger, f^\dagger))) + \gamma (R_2(f_N) - R_2(f^\dagger)) \right) \\
& = \delta + \gamma R_1(\zeta^\dagger, u^\dagger) + \gamma R_2(f^\dagger) + q_N + \eta.
\end{aligned} \tag{14}$$

By employing $S \geq 0$, $Q \geq 0$, and dividing by $\gamma > 0$, we then obtain

$$R_1(\zeta^{\gamma,\delta,N}, u^{\gamma,\delta,N}) + R_2(f^{\gamma,\delta,N}) \leq R_1(\zeta^\dagger, u^\dagger) + R_2(f^\dagger) + \frac{q_N}{\gamma} + \frac{\delta}{\gamma} + \frac{\eta}{\gamma}, \tag{15}$$

which by (T1)–(T3) implies existence of a $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ convergent subsequence $(\zeta^j, u^j, f^j)_{j \in \mathbb{N}}$ of $(\zeta^{\gamma(\delta), \delta, N(\delta)}, u^{\gamma(\delta), \delta, N(\delta)}, f^{\gamma(\delta), \delta, N(\delta)})_{\delta > 0}$ with limit (ζ^*, u^*, f^*) . Since from the same minimality estimate, by $R_1 \geq 0$, $R_2 \geq 0$ we also get

$$\begin{aligned}
& Q(E((\zeta^j, u^j, f^j))) + S(Mu^j, y^{\delta^j}) \\
& \leq \gamma(\delta^j) \left(R_1(\zeta^\dagger, u^\dagger) + R_2(f^\dagger) \right) + q_N(\delta^j) + \delta^j + \eta(\delta^j) \rightarrow 0 \text{ as } j \rightarrow \infty
\end{aligned}$$

for any such $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ convergent subsequence, from (T4) we conclude that (ζ^*, u^*, f^*) solves the inverse problem $E(\zeta^*, u^*, f^*) = 0$, $Mu^* = y$.

◇

Note that due to estimate (15) and τ_i -lower semiconituity of R_i in (T3), if in addition to (T5)

$$\frac{\delta}{\gamma(\delta)} \rightarrow 0, \quad \frac{q_N(\delta)}{\gamma(\delta)} \rightarrow 0, \quad \frac{\eta(\delta)}{\gamma(\delta)} \rightarrow 0, \quad \text{as } \delta \rightarrow 0,$$

then the limit according to Proposition 2 is even an R_1 , R_2 minimizing solution of the inverse problem, that is, $(\xi^*, u^*, f^*) = \min(R_1(\xi, u) + R_2(f))$, where the minimum is taken over all (ξ, u, f) solving the inverse problem.

Discussion 1 (Quantitative approximation error). *Assumption (T4) follows by application of the universal approximation theorem [25] in our setting. This theorem states that for any continuous function f on a compact domain, there exists a NN with a sufficiently large*

number of neurons approximating f with arbitrary prescribed accuracy. Recently, advanced studies on quantifying the size of NNs have been carried out, even in terms of width and depth, to obtain approximation rates. Seminal results in this direction include [6, 7], which show an asymptotic approximation rate $\mathcal{O}(1/\sqrt{N})$ in the L^2 -norm of NNs with N neurons and sigmoidal activation to any target function f with finite Fourier moments. The study on approximation rates has greatly evolved in recent years [40, 56, 43, 16]. For a full survey on approximation theory, we refer to [17], as well as [35, Section 1.4.2], [39, Table 1] for brief summaries.

Incorporating these approximation rates into q_N in (T4) enables an analysis for the convergence rate of Tikhonov regularization, under so-called source conditions, see, e.g., [20, 50]. The asymptotics of the parameters (T5) shows that when $\delta \rightarrow 0$, the NNs size should increase accordingly, that is, $N(\delta) \rightarrow \infty$. The relation reveals a choice of the network size dependent on the noise level δ . This potentially reduces the overfitting problem caused by noisy training data. By virtue of (14)-(15), the approximation errors both w.r.t the model $Q \circ E$ and w.r.t the regularizer R_2 contribute to the total convergence rate, hinting at a possible mutual effect of these two factors in the overall rate. If a convergence rate analysis for Tikhonov regularization integrating the quantitative approximation error can be carried out, a choice of N with respect to δ (c.f. [13]), as well as $Q \circ E$ and R_2 can be made explicit. We leave this interesting task for future research.

2.2 On Assumptions 1, 2: Discussion, Examples, and Application

Discussion 2 (on (T1)-(T2)). Assume the following:

- (i) $\exists \tau_{\mathcal{W} \times H}$ such that E is $\tau_{Z \times V} \times \tau_{\mathcal{C}}$ -to- $\tau_{\mathcal{W} \times H}$ continuous at the exact solutions of the PDE.
- (ii) $\exists \tau_{\mathcal{Y}}$ such that M is $\tau_{\mathcal{V}}$ -to- $\tau_{\mathcal{Y}}$ continuous at the exact states. In addition, $S(y, y^j) \rightarrow 0$ implies $y^j \rightarrow y$ in $\tau_{\mathcal{Y}}$.
- (iii) (Q, S) is $(\tau_{\mathcal{W} \times H}, (\tau_{\mathcal{Y}} \times \tau_{\mathcal{Y}}))$ lower semicontinuous.
- (iv) (R_1, R_2) is $(\tau_{Z \times V}, \tau_{\mathcal{C}})$ lower semicontinuous and its sublevel sets are $(\tau_{Z \times V}, \tau_{\mathcal{C}})$ compact.

Then (T1)-(T2) hold.

First, $Q \circ E$ is lower semicontinuous as it is a composition of a lower semicontinuous function and a continuous function assumed in (i)-(iii). Next, $Q \circ E$ is closed since by positivity of Q , lower semicontinuity of $Q \circ E$ and the premise of (T2), one has

$$0 \leq Q(E(\bar{\zeta}, \bar{u}, \bar{f})) \leq \liminf_{j \rightarrow \infty} Q(E(\zeta^j, u^j, f^j)) \leq \lim_{j \rightarrow \infty} Q(E(\zeta^j, u^j, f^j)) = 0$$

thus $Q(E(\bar{\zeta}, \bar{u}, \bar{f})) = 0$.

Note that closedness in the sense of (T2) is weaker than in the standard definition, (see, e.g., [45] and the references therein), as we require the closedness property only at the exact solutions of the PDE, i.e. at $Q(E(\bar{\zeta}, \bar{u}, \bar{f})) = 0$.

Furthermore, if $S(y, y^j) \rightarrow 0$, that is the premise of (T2), induces $y^j \rightarrow y$ in $\tau_{\mathcal{Y}}$ (cf. (ii)), then $S \circ M$ is closed due to

$$0 \leq S(M\bar{u}, y) \leq \lim_{j \rightarrow \infty} S(Mu^j, y^j) = 0 \implies \text{thus } S(M\bar{u}, y) = 0,$$

provided that S is lower semicontinuous in its two arguments (see (iii)).

Remark 2 (on (T3)). In Discussion 2, if E is $\tau_{Z \times \mathcal{V}} \times \tau_{\mathcal{C}}$ - to - $\tau_{\mathcal{W} \times H}$ continuous on the whole space $Z \times \mathcal{V}$ and M is $\tau_{\mathcal{V}}$ - to - $\tau_{\mathcal{Y}}$ continuous on \mathcal{V} , then lower semicontinuity of T_{γ}^{δ} on $Z \times \mathcal{V} \times \mathcal{C}_N$ ((T3)) holds. In some particular examples where convexity of $R_i, i = 1, 2$ is given, e.g. $R_i = \|\cdot\|^p, i = 1, 2$, for some $p \in [1, \infty]$ weaker conditions on continuity of E, M might be sufficient.

Remark 3. In case of full measurement, the term $S(Mu, y) = S(u, y)$ can play the role of a regularizer on u with $\tau_{\mathcal{V}} = \tau_{\mathcal{Y}}$.

Discussion 3 (on (T4)). The topology $\tau_{\mathcal{C}}$ induced by R_2 could be chosen as the weak* topology induced by the L^{∞} -norm to make use of available approximation rates of deep neural networks to smooth functions. In particular, these rates are with respect to arbitrary depths (number of layers) and widths (number of neurons per layer) [39, Table 1] to which N in (12) generally refers.

This and the discretization error assumption (T4)

$$q_N := \inf_{f_N \in \mathcal{C}_N} Q(E(\zeta^{\dagger}, u^{\dagger}, f_N)) = \inf_{f_N \in \mathcal{C}_N} Q(E(\zeta^{\dagger}, u^{\dagger}, f_N(u^{\dagger}))) \rightarrow 0$$

require uniform boundedness only on the exact state u^{\dagger} . Therefore, a candidate for R_2 is $R_2(f_N) = \|f_N\|_{L^{\infty}(\Omega_{y^{\dagger}})}$ with $\Omega_{y^{\dagger}} = u^{\dagger}((0, T) \times \Omega)$.

In the following examples of settings satisfying Assumptions 1, 2, we consider reflexive spaces or duals of separable spaces.

Example 1. Let

$$\begin{aligned} Q(E(\zeta, u, f_N)) &= \|E(\zeta, u, f_N)\|_{\mathcal{W} \times H}^2, & S(Mu, y) &= \|Mu - y\|_{\mathcal{Y}}^2, \\ R_1(\zeta, u) &= \|(\zeta, u)\|_{Z \times \mathcal{V}}^2, & R_2(f) &= \|f\|_{W^{1, \infty}(\Omega_y)}^2, \end{aligned}$$

with a bounded interval $\Omega_y \subset \mathbb{R}$ containing $\Omega_{y^{\dagger}} = u^{\dagger}((0, T) \times \Omega)$ as detailed below. Since the required compactness and continuity properties are straightforward on the finite dimensional space \mathbb{R}^n , for simplicity of exposition we skip α as an argument of f .

Then:

- Let $\tau_X \times \tau_{\mathcal{V}}$ be the weak topology on $X \times \mathcal{V}$ and assume that

$$\tilde{E} := \left(\frac{d}{dt} - F \right) \text{ is } (X \text{ weak}) \times (\mathcal{V} \text{ weak}) \text{ - to - } (\mathcal{W} \text{ weak}) \quad (16)$$

continuous at the exact solution.

This weak continuity, thus closedness (T2), depends on the PDE models and the choice of function spaces.

- Now, with $\tau_{\mathcal{C}}$ being the weak topology on $W^{1,\infty}(\Omega_y)$, we show continuity of the rest of E . In particular, we prove that $(u_N, f_N) \rightarrow (u, f)$ in the topology $(\mathcal{V} \text{ weak}) \times (W^{1,\infty}(\Omega_y) \text{ weak})$ implies $f_N(u_N) \rightarrow f(u)$ weakly in \mathcal{W} under appropriate conditions on \mathcal{V} , \mathcal{W} to be derived here. First, we observe

$$\begin{aligned} f_N(u_N(x, t)) - f(u(x, t)) &= (f_N(u_N(x, t)) - f_N(u(x, t))) + (f_N(u(x, t)) - f(u(x, t))) \\ &= \int_0^1 (f_N)'(u(x, t) + \theta(u_N(x, t) - u(x, t))) d\theta (u_N(x, t) - u(x, t)) \\ &\quad + (f_N(u(x, t)) - f(u(x, t))) \end{aligned} \quad (17)$$

$$\begin{aligned} &\langle f_N(u_N(x, t)) - f(u(x, t)), \psi \rangle_{\mathcal{W}, \mathcal{W}^*} \\ &\leq \|(f_N)'\|_{L^\infty(\Omega_y)} \|u_N - u\|_{L^p((0, T) \times \Omega)} \|\psi\|_{L^{p^*}((0, T) \times \Omega)} + \|f_N - f\|_{L^\infty(\Omega_y)} \|\psi\|_{L^1((0, T) \times \Omega)}, \end{aligned} \quad (18)$$

for any $\psi \in \mathcal{W}^*$, with $p \in [1, \infty]$ and p^* being the conjugate index of p . If

$$\mathcal{V} \subset L^\infty((0, T) \times \Omega), \quad (19)$$

then for $u_N \xrightarrow{\mathcal{V}} u$, one has $\|u_N\|_{L^\infty(0, T) \times \Omega}, \|u\|_{L^\infty(0, T) \times \Omega} \leq C, \forall N \in \mathbb{N}$, and may set $\Omega_y := [-C - 1, C + 1]$. Note that the inclusion $\mathcal{V} \subset L^\infty((0, T) \times \Omega)$ allows us to apply the fact that neural networks are dense in the space of smooth functions on compact sets. Next, $f_N \rightarrow f$ in $W^{1,\infty}(\Omega_y)$ shows that $\|(f_N)'\|_{L^\infty(\Omega_y)}$ is bounded for all N , and due to $W^{1,\infty}(\Omega_y) \hookrightarrow L^\infty(\Omega_y)$ we have $f_N \rightarrow f$ in $L^\infty(\Omega_y)$. If

$$\mathcal{V} \hookrightarrow L^{p_W}((0, T) \times \Omega) \subset \mathcal{W} \quad (20)$$

for some $p_W \in [1, \infty]$, then (17) shows $f_N(u_N) \rightarrow f(u)$ in \mathcal{W} , meaning $\tau_{\mathcal{V}} \times \tau_{\mathcal{C}} \text{-to-} \tau_{\mathcal{W}}$ continuity of $(u, f) \mapsto f(u)$ on $\mathcal{V} \times \mathcal{C}$.

Recall that for closedness of $(u, f) \mapsto f(u)$, we require only its continuity at exact solutions (u^\dagger, f^\dagger) of the PDE. Therefore, by

$$\begin{aligned} &f_N(u_N(x, t)) - f^\dagger(u^\dagger(x, t)) \\ &= (f^\dagger(u_N(x, t)) - f^\dagger(u^\dagger(x, t))) + (f_N(u_N(x, t)) - f^\dagger(u_N(x, t))) \end{aligned}$$

we only need to assume boundedness of $\|(f^\dagger)'\|_{L^\infty(\Omega_y)}$ thus can choose the weaker $R_2(f) = \|f\|_{\mathcal{C}}^2, \mathcal{C} \hookrightarrow L^\infty(\Omega_y)$. Note that the inclusions (19), (20) are still needed.

- The part on the initial condition $E_0(u_0, u) = u(0) - u_0$ is linear, thus requiring just the embedding $\mathcal{V} \hookrightarrow C(0, T; H)$ and the regularizer $R_1(u_0) = \|u_0\|_H$ induces τ_H , the weak topology on H .
- Regarding the observation M , if M is linear and bounded, then it is \mathcal{V} weak $-$ to $-$ \mathcal{Y} weak continuous.

Example 2 (norm of the hyperparameter as R_2). In the previous example, we consider the Sobolev $W^{1,\infty}$ -norm for $f, f_N \in \mathcal{C}$. As the discretized regularization is carried out for $f_N \in \mathcal{C}^N$, the space of neural networks of N hyperparameters, a natural question is whether one can replace the Sobolev norm by some equivalent norm on the hyperparameters. The answer in the general case is no. Consider e.g. the function $f(x) = x$ expressed via a 2 layers neural network of identity activation $f(x) = x = 1 \cdot \text{Id}(x + b) - b$. The hyperparameters are $\theta := (w_2, b_2, w_1, b_1) = (1, b, 1, -b)$. So, when b tends to infinity $\|\theta\| \rightarrow \infty$ while $\|f\|_{L^\infty(\Omega)} < \infty$ for any bounded domain Ω .

Let us study a standard case of a neural network with fixed depth two

$$f_N(y) = W_2^N \cdot \sigma(W_1^N y + b_1^N) + b_2^N,$$

where $y \in \mathbb{R}$, $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, $W_1^N \in \mathbb{R}^{N \times 1}$, $b \in \mathbb{R}^N$, $W_2^N \in \mathbb{R}^{1 \times N}$, $b_2^N \in \mathbb{R}$, and \cdot denotes matrix multiplication. This means when $N \rightarrow \infty$, the width of the neural network sequence tends to infinity. Assuming that σ is Lipschitz continuous with Lipschitz constant L_σ , we have

$$|f_N(u_N) - f_N(u)| \leq L_\sigma |W_2^N| \cdot |W_1^N| |u_N - u|, \quad (21)$$

where $|\cdot|$ represents element-wise absolute value. The class of Lipschitz activations used in practice is large; some examples include ReLU $\sigma(x) := \max\{0, x\}$ (with approximation rates), tansig $\sigma(x) := \tanh(x)$, softplus $\sigma(x) := \ln(1 + e^x)$, sigmoid or soft step $\sigma(x) := \frac{1}{1+e^{-x}}$, softsign $\sigma := \frac{x}{1+|x|}$ etc.

Furthermore, we assume that σ is coercive in the sense that $\exists C_\sigma > 0 : |y| \leq C_\sigma |\sigma(y)|, \forall y \in \mathbb{R}$, $\sigma(y) \geq 0$ for $y \geq 0$, all hyperparameters are nonnegative. We can then estimate

$$\begin{aligned} |W_2^N| \cdot |W_1^N| &\leq |W_2^N| \cdot (|W_1^N| + b_1^N) \leq C_\sigma W_2^N \cdot \sigma(W_1^N + b_1^N) \\ &\leq C_\sigma (W_2^N \cdot \sigma(W_1^N 1 + b_1^N) + b_2^N) = C_\sigma f_N(1) \leq \sup_{x \in [0,1]} C_\sigma |f_N(x)| = C_\sigma \|f_N\|_{L^\infty([0,1])}. \end{aligned}$$

Combining this with (21) and (17), we can replace boundedness of $\|(f_N)'\|_{L^\infty(\Omega)}$ by boundedness of $\|f_N\|_{L^\infty(\Omega_y)}$, and use the weaker regularizer $R_2 = \|\cdot\|_{L^\infty(\Omega_y)}^2$ instead of $\|\cdot\|_{W^{1,\infty}(\Omega_y)}^2$. When considering \mathbb{R}^+ , some examples for Lipschitz continuous and coercive activation functions are: ReLU, Leaky ReLU (coercive on \mathbb{R}), softplus etc. Assume further that the exact f^\dagger can be expressed exactly via a neural network, possibly with infinitely many hyperparameters, say $f^\dagger \in \mathcal{C}_\infty$ with $f^\dagger(0) = 0$, similar to (21) we have

$$|f^\dagger(y)| = |f^\dagger(y) - f^\dagger(0)| \leq L_\sigma |W_2| \cdot |W_1| |y| \leq |\Omega_y| L_\sigma |W_2| \cdot |W_1| \leq C \|W_2\|_{\ell^2} \|W_1\|_{\ell^2} \leq C \|\theta\|_{\ell^2}^2.$$

Then one can also use the stronger norm $\|\cdot\|_{\ell^2}$ in the regularizer R_2 , alternatively $\|\cdot\|_{\ell^1}$, due to norm equivalence in finite dimensional hyperparameter spaces. The application of sparsity-promoting techniques, such as incorporating ℓ^1 regularizers, has been proven as one of the remedies for overfitting in machine learning in practice. Indeed, the sparse optimization performs feature selection, yielding more interpretable trained models [53].

Example 3. Consider $R_1(u) = TV(u)$, the total variation of u on $\Omega_T := (0, T) \times \Omega, \Omega \subset \mathbb{R}^d$.

In order for R_1 to be lower semicontinuous and have $\tau_{\mathcal{V}}$ compact sublevel sets, we have some options:

1. $\tau_{\mathcal{V}}$ is the weak* topology on $BV(\Omega_T)$, the space of functions of bounded variation on $(0, T) \times \Omega$. Recall that $u^j \xrightarrow{*} u$ in $BV(\Omega_T)$ is defined as $u^j \xrightarrow{L^1} u$, $TV(u^j) \rightarrow TV(u)$. TV is weak* lower semicontinuous on $BV(\Omega_T)$, and its sublevel sets are weak* compact in $BV(\Omega_T)$ [14].
2. $\tau_{\mathcal{V}}$ is the strong topology on $L^{\frac{d+1}{d}-\epsilon}(\Omega_T)$ for arbitrary small $\epsilon > 0$. Application of the compact embedding $BV(\Omega_T) \hookrightarrow L^{\frac{d+1}{d}-\epsilon}(\Omega_T)$ yields that TV is lower semicontinuous on $L^{\frac{d+1}{d}-\epsilon}(\Omega_T)$, and its sublevel sets are compact in $L^{\frac{d+1}{d}-\epsilon}(\Omega_T)$ [2, Theorem 2.5].
3. $\tau_{\mathcal{V}}$ is the weak topology on $L^{\frac{d+1}{d}}(\Omega_T)$. Weak compactness of the sublevel sets is clear from the compact embedding mentioned above. Weak lower semicontinuity of TV was shown, e.g. in [2, Theorem 2.3].

Let us consider, for instance, the second case where $\tau_{\mathcal{V}}$ is the strong topology on $\mathcal{V} = L^{\frac{d+1}{d}-\epsilon}(\Omega_T)$ with $\tilde{E} = \dot{u} - \Delta u$, and assume $u^j \xrightarrow{\mathcal{V}} u$, $u^j(T) \xrightarrow{L^1(\Omega)} u(T)$, $u^j(0) = u(0) = 0$, $u^j(\partial\Omega) = u(\partial\Omega) = 0$. Let $\tilde{\epsilon} = \epsilon d^2 / (1 - \epsilon d)$, then due to the estimate

$$\begin{aligned} \langle \dot{u} - \dot{u}^j - \Delta(u - u^j), \psi \rangle_{\mathcal{W}, \mathcal{W}^*} &= \int_{\Omega_T} (u - u^j)(-\dot{\psi} - \Delta\psi) dx dt + \int_{\Omega} (u - u^j)(T)\psi(T) dx \\ &\leq \|u - u^j\|_{L^{\frac{d+1}{d}-\epsilon}(\Omega_T)} \|\dot{\psi} + \Delta\psi\|_{L^{d+1+\tilde{\epsilon}}(\Omega_T)} \\ &\quad + C_{W^{1,d+1+\tilde{\epsilon}}(\Omega) \rightarrow L^\infty(\Omega)} \|u(T) - u^j(T)\|_{L^1(\Omega)} \|\psi(T)\|_{W^{1,d+1+\tilde{\epsilon}}(\Omega)}, \end{aligned}$$

for any $\psi \in \mathcal{W}^*$, one can chose $\tau_{\mathcal{W}}$ as the strong topology on

$$\mathcal{W} := \left(L^{d+1+\tilde{\epsilon}}(0, T; W^{2,d+1+\tilde{\epsilon}}) \cap W^{1,d+1+\tilde{\epsilon}}(0, T; L^{d+1+\tilde{\epsilon}}(\Omega)) \right)^*.$$

Note that continuity of the embedding $\mathcal{W}^* \hookrightarrow C(0, T; W^{1,d+1+\tilde{\epsilon}}(\Omega))$ [48, Lemma 7.3] implies finiteness of $\|\psi(T)\|_{W^{1,d+1+\tilde{\epsilon}}(\Omega)}$.

Regarding f , since $\mathcal{V} \not\subset L^\infty(\Omega_T)$, in order to obtain uniform boundedness of u^j, u , we invoke full measurement data in a sufficiently strong observation space, e.g. $M = Id, \mathcal{Y} = L^\infty(\Omega_T)$. Then observe that the inclusions $\mathcal{V} \hookrightarrow L^{\frac{d+1}{d}-\epsilon}(\Omega_T) \subset \mathcal{W}$ hold, so convergence of the neural network sequence $f^j(u^j) \xrightarrow{j \rightarrow \infty} f^\dagger(u^\dagger)$, as discussed in Example 1, is guaranteed.

As such, we have two types of convergence for the sequence u^j : the strong convergence in $\mathcal{V} = L^{\frac{d+1}{d}-\epsilon}(\Omega_T)$, and the weak* convergence in $\mathcal{Y} = L^\infty(\Omega_T)$. These types of convergence are in general not equivalent. An example for this is the sequence of Rademacher functions $f_n : [0, 1] \rightarrow \{0, 1\}$ [45, Example 4.13]

$$f_n(x) = (-1)^{i+1} \quad \text{for } x \in [(i-1)/2^n, i/2^n], 1 \leq i \leq 2^n,$$

which weak* converges to zero in $L^\infty([0, 1])$, but not in the L^1 -norm, thus not in the $L^{\frac{d+1}{d}-\epsilon}$ -norm.

Example 4. Consider $S = KL$, the Kullback–Leibler divergence defined by

$$\text{for } y \in L^1(\Omega_T), KL(g, y) := \begin{cases} \int_{\Omega_T} y \left(\frac{g}{y} - \log \left(\frac{g}{y} \right) - 1 \right) dx dt & g, y \geq 0 \text{ a.e.} \\ \infty & \text{else} \end{cases}$$

It is clear that $S = KL$ does not satisfy a triangle inequality, a situation that is taken into account in this work. Positivity of KL is obvious as $(g/y - 1) \geq \log(g/y)$ and, $KL(g, y) = 0$ iff $g = y = 0$.

[9, Lemma A.2] states that $KL(y, y^j) \rightarrow 0$ implies $\|y^j - y\|_{L^1(\Omega')} \rightarrow 0$ for some Ω' , and for $\{g^j\} \in L^1(\Omega')$ with $g^j \rightarrow g$ in $L^1(\Omega')$ as $j \rightarrow \infty$, then $KL(g, y^j) \leq \liminf_{j \rightarrow \infty} KL(g^j, y^j)$. Fitting into our framework, in particular for existence of $\tau_{\mathcal{Y}}$ in Discussion 2, from $S(y, y^j) = KL(y, y^j) \rightarrow 0$ inducing $y^j \rightarrow y$ in $L^1(\Omega_T)$, one can choose $\tau_{\mathcal{Y}}$ as the strong topology on $L^1(\Omega_T)$. Also by this lemma, S is $\tau_{\mathcal{Y}} \times \tau_{\mathcal{Y}}$ lower semicontinuous. Therefore, we need M to be $\tau_{\mathcal{Y}}$ -to- $L^1(\Omega_T)$ continuous; this condition is very much obtainable in practice. In case $M = \text{Id}$, an estimate similar to the one in Example 3 could be carried out for $u^j \rightarrow u$ in $L^1(\Omega_T)$. Still, convergence of the neural network part requires $R_1(u) = \|u\|_{\mathcal{V}}$ with $\mathcal{V} \subset L^\infty((0, T) \times \Omega)$ and $L^1((0, T) \times \Omega) \subset \mathcal{W}$.

Application. We now return to Application (1), (2) and from Propositions 1 and 2 conclude a result for Tikhonov regularization in the setting of Example 1

$$(c^{\gamma, \delta, N}, \varphi^{\gamma, \delta, N}, u_0^{\gamma, \delta, N}, u^{\gamma, \delta, N}, f^{\gamma, \delta, N}) \in \operatorname{argmin}_{(c, \varphi, u_0, u, f_N) \in X_c \times X_\varphi \times H \times \mathcal{V} \times \mathcal{C}_N} T_\gamma^\delta(c, \varphi, u_0, u, f_N) \quad (22)$$

for

$$T_\gamma^\delta(c, \varphi, u_0, u, f_N) = \| \dot{u} - \Delta u + cu + h(u) - \varphi - f_N(u) \|_{W \times H}^2 + \| Mu - y \|_{\mathcal{Y}}^2 + \gamma \| (c, \varphi, u_0, u) \|_{X_c \times X_\varphi \times H \times \mathcal{V}}^2 + \gamma \| f_N \|_{W^{1, \infty}(\Omega_y)}^2, \quad (23)$$

where we can replace the $W^{1, \infty}(\Omega_y)$ norm of f_N by the hyperparameter norm according to Example 2.

For this purpose, recall the following requirements on the underlying spaces: (16), (19), (20), as well as boundedness of $M : \mathcal{V} \rightarrow \mathcal{Y}$ and of $\operatorname{tr}_{t=0} : \mathcal{V} \rightarrow H$. Here, we have the

operator $\tilde{E}u = \dot{u} - \Delta u + cu + h(u) - \varphi$, and the X space is decomposed as $X = X_c \times X_\varphi$; recall that h is known. We use the spaces

$$\begin{aligned} H &= W^{t_V, q_V}(\Omega), \quad X_c = L^r(\Omega), \quad X_\varphi = \mathcal{W} = W^{-s, p}(0, T; W^{-t, q}(\Omega)), \\ \mathcal{V} &= W^{1-s, p}(0, T; W^{-t, q}(\Omega)) \cap W^{-s, p}(0, T; W^{2-t, q}(\Omega)) \cap W^{s_V, p_V}(0, T; W^{t_V, q_V}(\Omega)) \end{aligned} \quad (24)$$

with

$$s_V > \frac{1}{p_V}, \quad t_V > \frac{d}{q_V}, \quad r \leq p_W := \min\{p, q\} \quad (25)$$

to satisfy (19), (20) and part of (16). To see the latter for the c part of the operator \tilde{E} , observe that for any sequence (c_n, u_n) converging weakly to (c, u) in $X_c \times \mathcal{V}$, by our choice of s_V, p_V, t_V, q_V there exists a subsequence (c_{n_k}, u_{n_k}) such that c_{n_k} converges weakly in $L^r(\Omega)$ and u_{n_k} converges strongly in $L^\infty(0, T; L^\infty(\Omega))$, so that for any $\psi \in \mathcal{W}^* \subseteq L^{p_W^*}((0, T) \times \Omega)$ we have $\int_0^T \psi u \in L^{r^*}(\Omega)$ and thus

$$\int_0^T \int_\Omega (c_{n_k} u_{n_k} - cu) \psi \, dx \, dt = \int_0^T \int_\Omega c_{n_k} (u_{n_k} - u) \psi \, dx \, dt + \int_\Omega (c_{n_k} - c) \int_0^T u \psi \, dt \, dx \rightarrow 0.$$

Likewise, it is straightforward to see that on the strength of the embeddings available for \mathcal{V} , it suffices to assume continuity of the real function h to achieve $(\mathcal{V} \text{ weak}) - \text{to} - (\mathcal{W} \text{ weak})$ continuity of the mapping $u \mapsto h(u)$ contained in \tilde{E} . Note that continuity of $\text{tr}_{t=0} : \mathcal{V} \rightarrow H$ also holds for any subspace H in which $W^{t_V, q_V}(\Omega)$ is continuously embedded.

Corollary 1. *With the spaces according to (24), (25), $M : \mathcal{V} \rightarrow \mathcal{Y}$ linear and bounded as well as $h : \mathbb{R} \rightarrow \mathbb{R}$ continuous, Tikhonov regularization is well-defined by (22), (23), and with a choice of $\gamma(\delta)$, $N(\delta)$ according to Assumption 2 (T5) we have weak^(*) subsequential convergence of $(c^{\gamma(\delta), \delta, N(\delta)}, \varphi^{\gamma(\delta), \delta, N(\delta)}, u_0^{\gamma(\delta), \delta, N(\delta)}, u^{\gamma(\delta), \delta, N(\delta)}, f^{\gamma(\delta), \delta, N(\delta)})$ to a solution of the inverse problem (1), (2) as $\delta \rightarrow 0$.*

3 Landweber iteration

In this section, for simplicity of exposition, collecting all unknowns $\lambda^m, u_0^m, \alpha^m, u^m, f$ in a single variable \mathbf{x} and setting $\mathbf{y} = \begin{pmatrix} 0 \\ y^m \end{pmatrix}_{m=1}^K$ we rewrite (7) as an operator equation

$$\mathbb{F}(\mathbf{x}) = \mathbf{y}. \quad (26)$$

Moreover, we restrict the setting to Hilbert spaces \mathbb{X}, \mathbb{Y} with Hilbert space adjoints denoted by a superscript $*$.

Landweber iteration defines regularized approximations as gradient descent steps for the least squares cost functional $\|\mathbb{F}(\mathbf{x}) - \mathbf{y}^\delta\|^2$, explicitly,

$$\mathbf{x}_{k+1}^\delta = \mathbf{x}_k^\delta - \mathbb{F}'(\mathbf{x}_k^\delta)^*(\mathbb{F}(\mathbf{x}_k^\delta) - \mathbf{y}^\delta).$$

Here, the stopping index $k = k_*(\delta, \mathbf{y}^\delta)$, which depends on the noise level δ and data \mathbf{y}^δ , acts as a regularization parameter. In order to accommodate additional constraints, e.g., on the magnitude or sign of \mathbf{x} , we consider a subset \mathbb{M} of \mathbb{X} . Constraints formulated by membership in the subset \mathbb{M} can be incorporated by projection via

$$\mathbf{x}_{k+1}^\delta = P_{\mathbb{M}}\left(\mathbf{x}_k^\delta - \mathbb{F}'(\mathbf{x}_k^\delta)^*(\mathbb{F}(\mathbf{x}_k^\delta) - \mathbf{y}^\delta)\right), \quad (27)$$

where the metric projection operator $P_{\mathbb{M}}$ onto a closed convex set \mathbb{M} is characterized by the variational inequality

$$\mathbf{x} = P_{\mathbb{M}}(\tilde{\mathbf{x}}) \Leftrightarrow (\mathbf{x} \in \mathbb{M} \text{ and } \forall \mathbf{z} \in \mathbb{M} : \langle \tilde{\mathbf{x}} - \mathbf{x}, \mathbf{z} - \mathbf{x} \rangle \leq 0). \quad (28)$$

$P_{\mathbb{M}}$ is nonexpansive and monotone, that is, for all $\mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{X}$,

$$\|P_{\mathbb{M}}(\mathbf{x}) - P_{\mathbb{M}}(\tilde{\mathbf{x}})\| \leq \|\mathbf{x} - \tilde{\mathbf{x}}\| \quad (29)$$

and

$$\langle P_{\mathbb{M}}(\mathbf{x}) - P_{\mathbb{M}}(\tilde{\mathbf{x}}), \mathbf{x} - \tilde{\mathbf{x}} \rangle \geq \|P_{\mathbb{M}}(\mathbf{x}) - P_{\mathbb{M}}(\tilde{\mathbf{x}})\|^2 \quad (30)$$

as well as continuous and, in general, nonlinear.

Discretization by restriction to a linear subspace $\mathbb{X}_N \subseteq \mathbb{X}$ can be easily done by replacing $\mathbb{F} : \mathbb{X} \rightarrow \mathbb{Y}$ by its restriction

$$\mathbb{F}_N := \mathbb{F}|_{\mathbb{X}_N} : \mathbb{X}_N \rightarrow \mathbb{Y}.$$

In our case, $\mathbb{X}_N = (Z \times \mathcal{V})^K \times \mathcal{C}_N$, so \mathbb{X}_N is a linear space in case of a linear activation function σ albeit not necessarily finite dimensional (for approximation on manifolds, see e.g. [21]). It yields the k -th iterate

$$\mathbf{x}_{N,k+1}^\delta = P_{\mathbb{M}}\left(\mathbf{x}_{N,k} - \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^*(\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta)\right) \quad (31)$$

in \mathbb{X}_N . Doing so, we use the Hilbert space adjoint of $\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta) : \mathbb{X}_N \rightarrow \mathbb{Y}$ that is uniquely determined by the identity

$$\langle \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* \mathbf{y}, \mathbf{x}_N \rangle = \langle \mathbf{y}, \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta) \mathbf{x}_N \rangle \quad \text{for all } \mathbf{y} \in \mathbb{Y}, \mathbf{x}_N \in \mathbb{X}_N.$$

Therefore, the adjoint in the discretized and projected Landweber (31) equals to

$$\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* : \mathbb{Y} \rightarrow \mathbb{X}_N \quad \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* = P_{\mathbb{X}_N} \mathbb{F}'(\mathbf{x}_{N,k}^\delta)^*, \quad (32)$$

the concatenation of $\mathbb{F}'(\mathbf{x}_{N,k}^\delta)^*$ with the orthogonal projection $P_{\mathbb{X}_N} : \mathbb{X} \rightarrow \mathbb{X}_N$ onto \mathbb{X}_N in the Hilbert space \mathbb{X} .

3.1 Convergence

Also for the discretized and projected Landweber iteration, we will show that with an appropriate choice of the stopping index $k_*(\delta)$ and discretization parameter $N(\delta)$ it is a regularization method.

We denote by $\mathbf{x}^\dagger \in \mathbb{X}$ a solution of the inverse problem with exact data, that is, $\mathbb{F}(\mathbf{x}^\dagger) = \mathbf{y}$, by $\mathbf{x}_{\infty,k}$ the iterates in \mathbb{X} according to (27), and make the following assumptions

Assumption 3. (L1) *Approximation by \mathbb{X}_N : There exists a sequence $(\mathbf{x}_N^\dagger)_{N \in \mathbb{N}}$, $\mathbf{x}_N^\dagger \in \mathbb{X}_N \cap \mathbb{M}$ such that for some \bar{d}*

$$d_N := \|\mathbf{x}_N^\dagger - \mathbf{x}^\dagger\| \leq \bar{d}, \quad m_N := \|\mathbb{F}_N(\mathbf{x}_N^\dagger) - \mathbf{y}\| = \|\mathbb{F}(\mathbf{x}_N^\dagger) - \mathbb{F}(\mathbf{x}^\dagger)\| \rightarrow 0, \\ \text{and } s_N := \sup_{\mathbf{x} \in B_R(\mathbf{x}^\dagger)} \|(I - P_{\mathbb{X}_N})\mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{y})\| \rightarrow 0 \text{ as } N \rightarrow \infty;$$

(L2) *Convergence and boundedness of the starting values*

$$\|\mathbf{x}_{N,0}^\delta - \mathbf{x}_{\infty,0}\| \leq \rho_N \rightarrow 0 \text{ as } N \rightarrow \infty, \\ \|\mathbf{x}_{N,0}^\delta - \mathbf{x}^\dagger\| \leq \rho, \quad \|\mathbb{F}(\mathbf{x}_{N,0}^\delta) - \mathbf{y}\| \leq \tilde{\rho}, \quad \text{for all } N \in \mathbb{N}$$

(e.g., by setting $\mathbf{x}_{N,0}^\delta := P_{\mathbb{X}_N} \mathbf{x}_0$);

(L3) *Local boundedness and tangential cone condition on \mathbb{F} as well as Lipschitz continuity of \mathbb{F}' : There exists $R > \rho + 2\bar{d}$, $\mu_R > 0$, $M_R > 0$, $K_R > 0$, $L_R > 0$, such that for all $x \in B_R(\mathbf{x}^\dagger)$ and for all $N \in \mathbb{N}$*

$$\|\mathbb{F}'(\mathbf{x})\| \leq M_R \leq \sqrt{2} \tag{33}$$

as well as

$$2\langle \mathbb{F}(\mathbf{x}) - \mathbb{F}(\mathbf{x}_N^\dagger), \mathbb{F}'(\mathbf{x})(\mathbf{x} - \mathbf{x}_N^\dagger) \rangle \geq (M_R^2 + \mu_R) \|\mathbb{F}(\mathbf{x}) - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 \\ \|\mathbb{F}'(\mathbf{x})(\mathbf{x} - \mathbf{x}_N^\dagger)\| \leq K_R \|\mathbb{F}(\mathbf{x}) - \mathbb{F}(\mathbf{x}_N^\dagger)\|; \tag{34}$$

$$\|\mathbb{F}'(\mathbf{x}) - \mathbb{F}'(\tilde{\mathbf{x}})\| \leq L_R \|\mathbf{x} - \tilde{\mathbf{x}}\|; \tag{35}$$

(L4) *Asymptotics of the parameters as $\delta \rightarrow 0$:*

$$k_*(\delta) \rightarrow \infty, \quad N(\delta) \rightarrow \infty \\ \left(\frac{4}{\mu_R} K_R + \left(1 + \frac{4}{\mu_R} M_R^2\right) M_R^2\right) k_*(\delta) (m_{N(\delta)} + \delta)^2 \leq (R - \bar{d})^2 - (\rho + \bar{d})^2;$$

(L5) *The mapping $\mathbf{x} \mapsto P_{\mathbb{M}}\left(\mathbf{x} - \mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{y})\right) - \mathbf{x}$ is weakly sequentially closed;*

(L6) *For all $y \in \mathbf{B}_{\bar{d}}(\mathbf{y})$ the mapping $x \mapsto \mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{z})$ is Lipschitz continuous with constant $L < 2$.*

Remark 4 (On (L1), approximation by NNs). *Note that the first part of (L1) (boundedness of d_N and convergence of m_N) only requires approximation of the single element \mathbf{x}^\dagger . By smoothness assumptions on \mathbf{x}^\dagger , this assumption therefore can be achieved, even with rates for discretization by NNs, as mentioned in Discussion 1.*

The second part of (L1) which is supposed to hold for all $\mathbf{x} \in B_R(\mathbf{x}^\dagger)$ can be obtained by using the fact that $\mathbb{F}'(\mathbf{x})^$ is a smoothing operator and therefore even norm convergence $\|(I - P_{\mathbb{X}_N})\mathbb{F}'(\mathbf{x})^*\| \rightarrow 0$ follows from error estimates of $I - P_{\mathbb{X}_N}$ under a priori regularity conditions.*

Remark 5 (on (L3), tangential cone condition). *A sufficient condition for (34) is the classical tangential cone condition (cf. [30])*

$$\|\mathbb{F}(\mathbf{x}) - \mathbb{F}(\mathbf{x}_N^\dagger) - \mathbb{F}'(\mathbf{x})(\mathbf{x} - \mathbf{x}_N^\dagger)\| \leq c_{tc}\|\mathbb{F}(\mathbf{x}) - \mathbb{F}(\mathbf{x}_N^\dagger)\|. \quad (36)$$

for some $c_{tc} < 1$ independent of $x \in B_R(\mathbf{x}^\dagger)$ with and all $N \in \mathbb{N}$, since by the inverse triangle inequality it is readily checked that we can then set $K_R = 1 + c_{tc}$ and $M_R + \mu_r = 1 - c_{tc}^2 + (1 - c_{tc})^2$.

We start with an estimate on the propagated noise and discretization error.

Lemma 1. *Under conditions (L1), (33), (35), for any $k \in \mathbb{N}$ the estimates*

$$\|\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}\| \leq (1 + \frac{5}{2}M_R L_R R)^k \rho_N + \frac{2}{5M_R L_R R} (1 + \frac{5}{2}M_R L_R R)^k (s_N + M_R \delta)$$

and

$$\begin{aligned} & \|\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^*(\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta) - \mathbb{F}'(\mathbf{x}_{\infty,k})^*(\mathbb{F}(\mathbf{x}_{\infty,k}) - \mathbf{y})\| \\ & \leq s_N + M_R(M_R + \frac{3}{2}L_R R)\|\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}\| + M_R \delta \end{aligned}$$

hold, provided that for all $\ell \leq k$, $\mathbf{x}_{N,\ell}^\delta, \mathbf{x}_{\infty,\ell} \in B_R(\mathbf{x}^\dagger)$.

Proof. We make use of the recursions

$$\begin{aligned} \mathbf{x}_{N,k+1}^\delta - \mathbf{x}^\dagger &= \left(I - P_{\mathbb{X}_N} K_{N,k}^\delta * \bar{K}_{N,k}^\delta \right) (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) + P_{\mathbb{X}_N} K_{N,k}^\delta * (\mathbf{y}^\delta - \mathbf{y}) \\ \mathbf{x}_{\infty,k+1} - \mathbf{x}^\dagger &= \left(I - K_{\infty,k} * \bar{K}_{\infty,k} \right) (\mathbf{x}_{\infty,k} - \mathbf{x}^\dagger) \end{aligned}$$

where $K_{N,k}^\delta = \mathbb{F}'(\mathbf{x}_{N,k}^\delta)$, $\bar{K}_{N,k}^\delta = \int_0^1 \mathbb{F}'(\mathbf{x}^\dagger + \theta(\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger)) d\theta$, $K_{\infty,k} = \mathbb{F}'(\mathbf{x}_{\infty,k})$, $\bar{K}_{\infty,k} = \int_0^1 \mathbb{F}'(\mathbf{x}^\dagger + \theta(\mathbf{x}_{\infty,k} - \mathbf{x}^\dagger)) d\theta$. This yields

$$\begin{aligned} & \mathbf{x}_{N,k+1}^\delta - \mathbf{x}_{\infty,k+1} \\ &= \left(I - K_{\infty,k} * K_{\infty,k} \right) (\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}) + K_{\infty,k} * (K_{\infty,k} - \bar{K}_{\infty,k}) (\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}) \\ & \quad + \left(K_{\infty,k} * (\bar{K}_{\infty,k} - \bar{K}_{N,k}^\delta) + (K_{\infty,k} - K_{N,k}^\delta) * \bar{K}_{N,k}^\delta \right) (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) \\ & \quad + (I - P_{\mathbb{X}_N}) K_{N,k}^\delta * \bar{K}_{N,k}^\delta (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) + P_{\mathbb{X}_N} K_{N,k}^\delta * (\mathbf{y}^\delta - \mathbf{y}) \end{aligned}$$

thus by $M_R \leq \sqrt{2}$, which implies $\|I - K_{\infty,k}^* K_{\infty,k}\| \leq 1$

$$\begin{aligned} \|\mathbf{x}_{N,k+1}^\delta - \mathbf{x}_{\infty,k+1}\| &\leq (1 + \frac{5}{2}M_R L_R R) \|\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}\| + s_N + M_R \delta \\ &\leq (1 + \frac{5}{2}M_R L_R R)^{k+1} \|\mathbf{x}_{N,0}^\delta - \mathbf{x}_{\infty,0}\| + \sum_{j=0}^k (1 + \frac{5}{2}M_R L_R R)^j (s_N + M_R \delta) \end{aligned}$$

Moreover,

$$\begin{aligned} &\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^*(\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta) - \mathbb{F}'(\mathbf{x}_{\infty,k})^*(\mathbb{F}(\mathbf{x}_{\infty,k}) - \mathbf{y}) \\ &= P_{\mathbb{X}_N} K_{N,k}^\delta \bar{K}_{N,k}^\delta (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) - K_{\infty,k}^* \bar{K}_{\infty,k} (\mathbf{x}_{\infty,k} - \mathbf{x}^\dagger) + P_{\mathbb{X}_N} K_{N,k}^\delta (\mathbf{y} - \mathbf{y}^\delta) \\ &= -(I - P_{\mathbb{X}_N}) K_{N,k}^\delta \bar{K}_{N,k}^\delta (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) + K_{N,k}^\delta \bar{K}_{N,k}^\delta (\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}) + P_{\mathbb{X}_N} K_{N,k}^\delta (\mathbf{y} - \mathbf{y}^\delta) \\ &\quad + \left(K_{N,k}^\delta (\bar{K}_{N,k}^\delta) - \bar{K}_{\infty,k} \right) + (K_{N,k}^\delta - K_{\infty,k})^* \bar{K}_{\infty,k} (\mathbf{x}_{\infty,k} - \mathbf{x}^\dagger) \end{aligned}$$

◇

Remark 6. In the linear case $\mathbb{F}(\mathbf{x}) = K\mathbf{x}$ with $\|K\| \leq 1$ the much better estimates

$$\|\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}\| \leq \|\mathbf{x}_{N,0}^\delta - \mathbf{x}_0\| + k(\|(I - P_{\mathbb{X}_N})K^*K\|R + \delta)$$

and

$$\begin{aligned} &\|\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^*(\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta) - \mathbb{F}'(\mathbf{x}_{\infty,k})^*(\mathbb{F}(\mathbf{x}_{\infty,k}) - \mathbf{y})\| \\ &\leq \frac{1}{k+1} \|\mathbf{x}_{N,0}^\delta - \mathbf{x}_0\| + \left(1 + \sum_{j=0}^{k-1} \frac{1}{j+1}\right) (\|(I - P_{\mathbb{X}_N})K^*K\|R + \delta) \end{aligned}$$

can be easily verified by means of spectral theoretic methods. More precisely, we use the fact that $\|K^*K(I - K^*K)^j\| \leq \frac{1}{j+1}$, and the identities

$$\begin{aligned} \mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k} &= (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) - (\mathbf{x}_{\infty,k} - \mathbf{x}^\dagger) \\ &= (I - K^*K)^k (\mathbf{x}_{N,0}^\delta - \mathbf{x}^\dagger) - (I - K^*K)^k (\mathbf{x}_{\infty,0} - \mathbf{x}^\dagger) \\ &\quad + \sum_{j=0}^{k-1} (I - K^*K)^j \left((I - P_{\mathbb{X}_N})K^*K (\mathbf{x}_{N,k-j-1}^\delta - \mathbf{x}^\dagger) + P_{\mathbb{X}_N} K^* (\mathbf{y}^\delta - \mathbf{y}) \right) \end{aligned}$$

$$\begin{aligned} &P_{\mathbb{X}_N} K^* (K\mathbf{x}_{N,k}^\delta - \mathbf{y}^\delta) - K^* (\mathbf{x}_{\infty,k} - \mathbf{y}) \\ &= K^* K (\mathbf{x}_{N,k}^\delta - \mathbf{x}_{\infty,k}) - (I - P_{\mathbb{X}_N}) K^* K (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) - P_{\mathbb{X}_N} K^* (\mathbf{y}^\delta - \mathbf{y}) \\ &= K^* K (I - K^*K)^k (\mathbf{x}_{N,0}^\delta - \mathbf{x}_{\infty,0}) \\ &\quad + \sum_{j=0}^{k-1} K^* K (I - K^*K)^j \left((I - P_{\mathbb{X}_N}) K^* K (\mathbf{x}_{N,k-j-1}^\delta - \mathbf{x}^\dagger) + P_{\mathbb{X}_N} K^* (\mathbf{y}^\delta - \mathbf{y}) \right) \\ &\quad - (I - P_{\mathbb{X}_N}) K^* K (\mathbf{x}_{N,k}^\delta - \mathbf{x}^\dagger) - P_{\mathbb{X}_N} K^* (\mathbf{y}^\delta - \mathbf{y}) \end{aligned}$$

They can actually be transferred to the nonlinear setting under an adjoint range invariance condition on \mathbb{F} , which is a stronger assumption than the tangential cone condition,

similarly to the convergence rates estimates in [24]. However, in our example, this assumption does not seem to be verifiable, whereas the tangential cone condition can be established, see below.

While uniform boundedness of the iterates can be shown under the assumptions (L1)-(L4), in order to control the propagated noise in the iterates, we will therefore have to additionally impose

Assumption 4.

$$(1 + \frac{5}{2}M_R L_R R)^k \rho_{N(k)} \rightarrow 0, \quad (1 + \frac{5}{2}M_R L_R R)^k s_{N(k)} \rightarrow 0 \text{ as } k \rightarrow \infty \quad (37)$$

in case of exact data $\delta = 0$ and

$$(1 + \frac{5}{2}M_R L_R R)^{k_*(\delta)} \rho_{N(\delta)} \rightarrow 0, \quad (1 + \frac{5}{2}M_R L_R R)^{k_*(\delta)} (s_{N(\delta)} + M_R \delta) \rightarrow 0 \text{ as } \delta \rightarrow 0. \quad (38)$$

Proposition 3. *Under the above assumptions (L1)-(L6) with \mathbb{M} closed and convex, the iterates are well-defined by (31) and remain in $B_R(\mathbf{x}^\dagger)$.*

Under the additional condition (38), we also have weak subsequential convergence of $\mathbf{x}_{N(\delta), k_(\delta)}^\delta$ to a solution of (26) as $\delta \rightarrow 0$.*

With exact data $\delta = 0$ and $N = N(k)$ chosen according to (37), we have weak subsequential convergence of $\mathbf{x}_{N(k), k}$ to a solution of (26) as $k \rightarrow \infty$.

Proof. We follow the classical monotonicity proof from [24], see also [32], but do so with F_N instead of F so that we can exploit the identity $\langle \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* (\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta), (\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger) \rangle = \langle \mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta, \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)(\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger) \rangle$ in the first equality below. It is also for this reason that we had to introduce the auxiliary variable x_N^\dagger as a substitute for x^\dagger in \mathbb{X}_N . Therewith we obtain, for arbitrary $N \in \mathbb{N}$, using the fact that we can skip the subscript N when applying \mathbb{F}_N to an element of \mathbb{X}_N and nonexpansivity (29) together with the fact that $P_{\mathbb{M}}(\mathbf{x}_N^\dagger) = \mathbf{x}_N^\dagger$

$$\begin{aligned} & \|\mathbf{x}_{N,k+1}^\delta - \mathbf{x}_N^\dagger\|^2 - \|\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger\|^2 \\ &= \|\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* (\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta)\|^2 - 2\langle \mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta, \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)(\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger) \rangle \\ &= \|\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* (\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbb{F}_N(\mathbf{x}_N^\dagger))\|^2 - 2\langle \mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbb{F}_N(\mathbf{x}_N^\dagger), \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)(\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger) \rangle \\ & \quad + 2\langle \mathbf{y}^\delta - \mathbb{F}_N(\mathbf{x}_N^\dagger), \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta) \left((\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger) - \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* (\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbb{F}_N(\mathbf{x}_N^\dagger)) \right) \rangle \\ & \quad + \|\mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^* (\mathbf{y}^\delta - \mathbb{F}_N(\mathbf{x}_N^\dagger))\|^2 \\ &\leq -\mu_R \|\mathbb{F}(\mathbf{x}_{N,k}^\delta) - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 + \left(\frac{2}{\epsilon} + M_R^2\right) \|\mathbf{y}^\delta - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 \\ & \quad + \epsilon \|\mathbb{F}'(\mathbf{x}_{N,k}^\delta)(\mathbf{x}_{N,k}^\delta - \mathbf{x}_N^\dagger)\|^2 + \epsilon M_R^4 \|\mathbb{F}(\mathbf{x}_{N,k}^\delta) - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 \\ &\leq -(\mu_R - \epsilon K_R - \epsilon M_R^4) \|\mathbb{F}(\mathbf{x}_{N,k}^\delta) - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 + \left(\frac{2}{\epsilon} + M_R^2\right) (\delta + m_N)^2 \\ &\leq -\frac{\mu_R}{2} \|\mathbb{F}(\mathbf{x}_{N,k}^\delta) - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 + \left(\frac{4}{\mu_R} K_R + \left(1 + \frac{4}{\mu_R} M_R^2\right) M_R^2\right) (\delta + m_N)^2 \end{aligned}$$

provided $\mathbf{x}_{N,k}^\delta \in B_R(\mathbf{x}^\dagger)$. Here we have employed Young's inequality in the form $2a(b+c) \leq \frac{2}{\epsilon}a^2 + \frac{\epsilon}{2}(b+c)^2 \leq \frac{2}{\epsilon}a^2 + \epsilon b^2 + \epsilon c^2$ with $\epsilon = \frac{\mu_R}{2(K_R + M_R^4)}$.

Summing up for k from zero to $\tilde{k} - 1$ we obtain

$$\begin{aligned} & \frac{\mu_R}{2} \sum_{k=0}^{\tilde{k}-1} \|\mathbb{F}(\mathbf{x}_{N,k}^\delta) - \mathbb{F}(\mathbf{x}_N^\dagger)\|^2 + \|\mathbf{x}_{N,\tilde{k}}^\delta - \mathbf{x}_N^\dagger\|^2 \\ & \leq \|\mathbf{x}_{N,0}^\delta - \mathbf{x}_N^\dagger\|^2 + \tilde{k} \left(\frac{4}{\mu_R} K_R + \left(1 + \frac{4}{\mu_R} M_R^2\right) M_R^2 \right) (\delta + m_N)^2, \end{aligned} \quad (39)$$

which by (L1), (L2) and (L4) inductively implies that the iterates $\mathbf{x}_{N(\delta),\tilde{k}}^\delta$ remain in $B_R(\mathbf{x}^\dagger)$ for all $\tilde{k} \leq k_*(\delta)$. Thus $(\mathbf{x}_{N(\delta),k_*(\delta)}^\delta)_{\delta>0}$ has a weakly convergent subsequence

$$\mathbf{x}_{N(\delta^j),k_*(\delta^j)}^{\delta^j} \rightharpoonup \bar{\mathbf{x}}. \quad (40)$$

and in case $\delta = 0$, with $N = N(k)$ in place of $N = N(\delta)$

$$\mathbf{x}_{N(k_j),k_j} \rightharpoonup \bar{\mathbf{x}}. \quad (41)$$

Since \mathbb{M} is closed and convex, hence weakly closed, $\bar{\mathbf{x}}$ is contained in \mathbb{M} .

To prove that this limit solves the inverse problem, like in [32, Lemma 3.1] with

$$J_N^\delta(\mathbf{x}) = \frac{1}{2} \|\mathbb{F}_N(\mathbf{x}) - \mathbf{y}^\delta\|^2, \quad \Delta_{N,k}^\delta := \mathbf{x}_{N,k+1}^\delta - \mathbf{x}_{N,k}^\delta = P_{\mathbb{M}} \left(\mathbf{x}_{N,k}^\delta - J_N^{\delta'}(\mathbf{x}_{N,k}) \right) - \mathbf{x}_{N,k}^\delta$$

and (L6), which implies that for all $\mathbf{x}_N, \tilde{\mathbf{x}}_N \in B_R(\mathbf{x}^\dagger)$

$$\begin{aligned} & \|J_N^{\delta'}(\mathbf{x}_N) - J_N^{\delta'}(\tilde{\mathbf{x}}_N)\| = \|\mathbb{F}'_N(\mathbf{x}_N)^*(\mathbb{F}_N(\mathbf{x}_N) - \mathbf{y}^\delta) - \mathbb{F}'_N(\tilde{\mathbf{x}}_N)^*(\mathbb{F}_N(\tilde{\mathbf{x}}_N) - \mathbf{y}^\delta)\| \\ & = \|P_{\mathbb{X}_N} \left(\mathbb{F}'_N(\mathbf{x}_N)^*(\mathbb{F}_N(\mathbf{x}_N) - \mathbf{y}^\delta) - \mathbb{F}'_N(\tilde{\mathbf{x}}_N)^*(\mathbb{F}_N(\tilde{\mathbf{x}}_N) - \mathbf{y}^\delta) \right)\| \leq L \|\mathbf{x}_N - \tilde{\mathbf{x}}_N\| \end{aligned}$$

We then obtain, for any k and for both the discretized problem in \mathbb{X}_N as well as non-discretized problem in \mathbb{X} (i.e. $N = \infty$),

$$\begin{aligned} J_N^\delta(\mathbf{x}_{N,k+1}) - J_N^\delta(\mathbf{x}_{N,k}) &= \int_0^1 J_N^{\delta'}(\mathbf{x}_{N,k}^\delta + \theta \Delta_{N,k}^\delta) \Delta_{N,k}^\delta d\theta \\ &\leq -\|\Delta_{N,k}^\delta\|^2 + \int_0^1 (J_N^{\delta'}(\mathbf{x}_{N,k}^\delta + \theta \Delta_{N,k}^\delta) - J_N^{\delta'}(\mathbf{x}_{N,k}^\delta)) \Delta_{N,k}^\delta d\theta \leq -(1 - \frac{L}{2}) \|\Delta_{N,k}^\delta\|^2, \end{aligned}$$

where we have used the fact that monotonicity (30) with $\mathbf{x} = \mathbf{x}_{N,k}^\delta - J_N^{\delta'}(\mathbf{x}_{N,k}^\delta)$, $\tilde{\mathbf{x}} = \mathbf{x}_{N,k}^\delta = P_{\mathbb{M}}(\mathbf{x}_{N,k}^\delta)$ since $\mathbf{x}_{N,k}^\delta \in \mathbb{M}$ implies

$$-J_N^{\delta'}(\mathbf{x}_{N,k}^\delta) \Delta_{N,k}^\delta \geq \|\Delta_{N,k}^\delta\|^2.$$

After summation and by (L2) and $J_N^\delta \geq 0$ this implies that

$$\sup_{\delta \in (0, \bar{\delta}], N \in \mathbb{N}} \sum_{k=0}^{\infty} \|\Delta_{N,k}^\delta\|^2 \leq \frac{1}{2-L} \sup_{N \in \mathbb{N}} \|\mathbb{F}(\mathbf{x}_{N,0}^\delta) - \mathbf{y}^\delta\|^2 \leq \frac{1}{2-L} (\tilde{\rho} + \bar{\delta})^2, \quad (42)$$

where

$$\begin{aligned}\Delta_{N,k}^\delta &= P_{\mathbb{M}}\left(\mathbf{x}_{N,k} - \mathbb{F}'_N(\mathbf{x}_{N,k}^\delta)^*(\mathbb{F}_N(\mathbf{x}_{N,k}^\delta) - \mathbf{y}^\delta)\right) - \mathbf{x}_{N,k} \\ &= P_{\mathbb{M}}\left(\mathbf{x}_{N,k} - P_{\mathbb{X}_N}\mathbb{F}'(\mathbf{x}_{N,k})^*(\mathbb{F}(\mathbf{x}_{N,k}) - \mathbf{y}^\delta)\right) - \mathbf{x}_{N,k}\end{aligned}$$

In particular, in case $\delta = 0$ (thus skipping the superscript δ and setting $N = N(k)$) with nonexpansivity of $P_{\mathbb{M}}$

$$\begin{aligned}\|\Delta_{N,k}^\delta\| &= \|P_{\mathbb{M}}\left(\mathbf{x}_{N(k),k} - \mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y})\right) - \mathbf{x}_{N(k),k}\| \\ &= \|P_{\mathbb{M}}\left(\mathbf{x}_{N(k),k} - P_{\mathbb{X}_N}\mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y})\right) - \mathbf{x}_{N(k),k} \\ &\quad + P_{\mathbb{M}}\left(\mathbf{x}_{N(k),k} - \mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y})\right) \\ &\quad - P_{\mathbb{M}}\left(\mathbf{x}_{N(k),k} - P_{\mathbb{X}_N}\mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y})\right)\| \\ &\leq \|\Delta_{N(k),k}\| + \|(I - P_{\mathbb{X}_N(k)})\mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y})\| \\ &\leq \|\Delta_{\infty,k}\| + 2\|(I - P_{\mathbb{X}_N(k)})\mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y})\| + 2\|\mathbf{x}_{N(k),k} - \mathbf{x}_{\infty,k}\| \\ &\quad + \|\mathbb{F}'(\mathbf{x}_{N(k),k})^*(\mathbb{F}(\mathbf{x}_{N(k),k}) - \mathbf{y}) - \mathbb{F}'(\mathbf{x}_{\infty,k})^*(\mathbb{F}(\mathbf{x}_{\infty,k}) - \mathbf{y})\| \\ &\rightarrow 0 \text{ as } k \rightarrow \infty,\end{aligned}$$

due to (42), (L1) and (37), according to Lemma 1.

Thus from (41), $\bar{x} \in \mathbb{M}$ and (L5) we get $P_{\mathbb{M}}\left(\bar{\mathbf{x}} - \mathbb{F}'(\bar{\mathbf{x}})^*(\mathbb{F}(\bar{\mathbf{x}}) - \mathbf{y})\right) - \bar{\mathbf{x}} = 0$, hence due to (28) with $\tilde{\mathbf{x}} = \bar{\mathbf{x}} - \mathbb{F}'(\bar{\mathbf{x}})^*(\mathbb{F}(\bar{\mathbf{x}}) - \mathbf{y})$, $\mathbf{x} = \bar{\mathbf{x}}$, $\mathbf{z} = \mathbf{x}^\dagger$ and (L3)

$$0 \geq \langle \mathbb{F}'(\bar{\mathbf{x}})(\bar{\mathbf{x}} - \mathbf{x}^\dagger), \mathbb{F}(\bar{\mathbf{x}}) - \mathbb{F}(\mathbf{x}^\dagger) \rangle \geq \frac{M_R^2 + \mu_R}{2} \|\mathbb{F}(\bar{\mathbf{x}}) - \mathbb{F}(\mathbf{x}^\dagger)\|^2.$$

This gives subsequential convergence of $\mathbf{x}_{N(\delta),k(\delta)}$ to a solution x^\dagger of (26) as $k \rightarrow \infty$ with exact data, for both the discretized and the nondiscretized problem.

Convergence with noisy data can be concluded from Lemma 1 under the more restrictive assumption (38). Indeed, in the decomposition

$$\mathbf{x}_{N(\delta),k(\delta)}^\delta - \mathbf{x}^\dagger = (\mathbf{x}_{N(\delta),k(\delta)}^\delta - \mathbf{x}_{\infty,k(\delta)}) + (\mathbf{x}_{\infty,k(\delta)} - \mathbf{x}^\dagger),$$

convergence of the first term follows from Lemma 1 and the rule (38), while weak convergence of the second term is a consequence of the result with exact data in the nondiscretized setting $N = \infty$, that we have just proven above.

◇

3.2 Adjoint

We will now write out $\mathbb{F}'(\mathbf{x})^*$ and the expression $\mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{y})$ that plays a role both in the definition of the Landweber iteration and in the verification of the conditions (L5), (L6).

To do so, we recall the setting

$$\begin{aligned} \mathbb{F}(\mathbf{x}) &= \begin{pmatrix} \dot{u} - F(\lambda, u) - f(\alpha, u) \\ u(0) - u_0 \\ Mu \end{pmatrix} \in \mathcal{V} \times H \times \mathcal{Y}, \quad \mathbf{y} = \begin{pmatrix} 0 \\ 0 \\ y \end{pmatrix}, \quad \tilde{\mathbf{y}} = \begin{pmatrix} \tilde{w} \\ \tilde{h} \\ \tilde{y} \end{pmatrix}, \\ \mathbb{F}'(\mathbf{x})\tilde{\mathbf{x}} &= \begin{pmatrix} -F_\lambda(\lambda, u)\tilde{\lambda} - f_\alpha(\alpha, u)\tilde{\alpha} + \dot{\tilde{u}} - F_u(\lambda, u)\tilde{u} - f_u(\alpha, u)\tilde{u} - \tilde{f}(\alpha, u) \\ \tilde{u}(0) - \tilde{u}_0 \\ M\tilde{u} \end{pmatrix}, \\ \mathbf{x} &= \begin{pmatrix} \lambda \\ u_0 \\ \alpha \\ u \\ f \end{pmatrix}, \quad \tilde{\mathbf{x}} = \begin{pmatrix} \tilde{\lambda} \\ \tilde{u}_0 \\ \tilde{\alpha} \\ \tilde{u} \\ \tilde{f} \end{pmatrix}, \quad \mathbb{F}'(\mathbf{x})^*\tilde{\mathbf{y}} = \begin{pmatrix} \tilde{\mu} \\ \tilde{v}_0 \\ \tilde{\beta} \\ \tilde{v} \\ \tilde{g} \end{pmatrix} \in X \times H \times \mathbb{R}^n \times \mathcal{V} \times \mathcal{C}, \\ \mathcal{V} &= H^1(0, T; V^*) \cap L^2(0, T; V), \quad \mathcal{W} = L^2(0, T; V^*), \quad \mathcal{Y} = L^2(0, T; Y), \end{aligned} \quad (43)$$

$$\begin{aligned} \langle \tilde{u}, \tilde{v} \rangle_{\mathcal{V}} &= \int_0^T \left(\langle \dot{\tilde{u}}(t), \dot{\tilde{v}}(t) \rangle_{V^*} + \langle \tilde{u}(t), \tilde{v}(t) \rangle_V \right) dt \\ &= \int_0^T \left(\langle \dot{\tilde{u}}(t), I_V \dot{\tilde{v}}(t) \rangle_{V^*, V} + \langle D_V \tilde{v}(t), \tilde{u}(t) \rangle_{V^*, V} \right) dt, \\ \langle w, \tilde{w} \rangle_{\mathcal{W}} &= \int_0^T \langle w(t), \tilde{w}(t) \rangle_{V^*} dt = \int_0^T \langle w(t), I_V \tilde{w}(t) \rangle_{V^*, V} dt \end{aligned}$$

with the Riesz isomorphisms $I_V : V^* \rightarrow V$, $D_V : V \rightarrow V^*$ and $V \hookrightarrow H \hookrightarrow V^*$ forming a Gelfand triple, and a Hilbert parameter space X . We use the integration by parts identity

$$\int_0^T \left(\langle \dot{\tilde{u}}(t), z(t) \rangle_{V^*, V} + \langle \dot{z}(t), \tilde{u}(t) \rangle_{V^*, V} \right) dt = \langle \tilde{u}(T), z(T) \rangle_H - \langle \tilde{u}(0), z(0) \rangle_H.$$

Moreover, in order to work in a Hilbert space setting, we will use the Bochner Sobolev space

$$\mathcal{C} = H^\ell(\mathbb{R}^n, H^r(\mathbb{R})), \quad \langle \tilde{f}, \tilde{g} \rangle_{\mathcal{C}} := \int_{\mathbb{R}^n} (1 + |\kappa|^2)^\ell \int_{\mathbb{R}} (1 + |\omega|^2)^r (\mathcal{F}\tilde{f})(\kappa, \omega) \overline{(\mathcal{F}\tilde{g})(\kappa, \omega)} d\omega d\kappa \quad (44)$$

with ℓ, r large enough to allow for $\mathcal{C} \subseteq C(\mathbb{R}^n, \mathbb{R}) \cap C(\mathbb{R}^n; W^{1, \infty}(\mathbb{R}))$, see (50) below, where \mathcal{F} denotes the Fourier transform. Therewith, the defining identity for the Hilbert space adjoint $\mathbb{F}'(\mathbf{x})^*\tilde{\mathbf{y}}$, that is,

$$0 = \langle \mathbb{F}'(\mathbf{x})^*\tilde{\mathbf{y}}, \tilde{\mathbf{x}} \rangle - \langle \tilde{\mathbf{y}}, \mathbb{F}'(\mathbf{x})\tilde{\mathbf{x}} \rangle \quad \text{for all } \mathbf{x} \in \mathbb{X}$$

reads as follows:

$$0 = \langle \tilde{\lambda}, \tilde{\mu} \rangle_X + \langle \tilde{u}_0, \tilde{v}_0 \rangle_H + \langle \tilde{\alpha}, \tilde{\beta} \rangle_{\mathbb{R}^n} + \int_0^T \left(\langle \dot{\tilde{u}}(t), I_V \dot{\tilde{v}}(t) \rangle_{V^*, V} + \langle D_V \tilde{v}(t), \tilde{u}(t) \rangle_{V^*, V} \right) dt + \langle \tilde{f}, \tilde{g} \rangle_{\mathcal{C}}$$

$$\begin{aligned}
& + \int_0^T \langle F_\lambda(\lambda, u)\tilde{\lambda} + f_\alpha(\alpha, u)\tilde{\alpha} - \dot{\tilde{u}} + F_u(\lambda, u)\tilde{u} + f_u(\alpha, u)\tilde{u} + \tilde{f}(\alpha, u), I_V\tilde{w} \rangle_{V^*, V} dt \\
& - \langle \tilde{u}(0) - \tilde{u}_0, \tilde{h} \rangle_H - \int_0^T \langle M\tilde{u}, \tilde{y} \rangle_Y dt \\
& = \int_0^T \left(\langle -I_V\ddot{\tilde{v}} + I_V\dot{\tilde{w}} + F_u(\lambda, u)^*I_V\tilde{w} + f_u(\alpha, u)^*I_V\tilde{w} - M^*\tilde{y} + D_V\tilde{v}, \tilde{u} \rangle_{V^*, V} dt \right. \\
& + \langle I_V\dot{\tilde{v}}(T) - I_V\tilde{w}(T), \tilde{u}(T) \rangle_H - \langle I_V\dot{\tilde{v}}(0) - I_V\tilde{w}(0) - \tilde{h}, \tilde{u}(0) \rangle_H + \langle \tilde{h} + \tilde{v}_0, \tilde{u}_0 \rangle_H \\
& + \langle \int_0^T F_\lambda(\lambda, u)^*I_V\tilde{w} dt + \tilde{\mu}, \tilde{\lambda} \rangle_X + \langle \int_0^T f_\alpha(\alpha, u)^*I_V\tilde{w} dt + \tilde{\beta}, \tilde{\alpha} \rangle_{\mathbb{R}^n} \\
& + \int_{\mathbb{R}^n} \int_{\mathbb{R}} \left((1 + |\kappa|^2)^\ell (1 + |\omega|^2)^r \overline{\mathcal{F}\tilde{g}(\kappa, \omega)} \right. \\
& \quad \left. + \frac{1}{2\pi^{(n+1)/2}} \left[\int_0^T \int_\Omega e^{i\kappa \cdot \alpha} e^{i\omega u(x, t)} (I_V\tilde{w})(x, t) dx dt \right] \right) \mathcal{F}\tilde{f}(\beta, \omega) d\omega d\kappa
\end{aligned}$$

where we have rewritten

$$\tilde{f}(\alpha, u(x, t)) = \frac{1}{2\pi^{(n+1)/2}} \int_{\mathbb{R}^n} \int_{\mathbb{R}} e^{i\kappa \cdot \alpha} e^{i\omega u(x, t)} \mathcal{F}\tilde{f}(\kappa, \omega) d\omega d\kappa$$

by the definition of the Fourier transform. This leads us to defining

$$\begin{aligned}
\mathbb{F}'(\mathbf{x})^* \tilde{\mathbf{y}} & = \left(\tilde{\mu} \quad \tilde{v}_0 \quad \tilde{\beta} \quad \tilde{v} \quad \tilde{g} \right)^T \in X \times H \times \mathbb{R}^n \times \mathcal{V} \times \mathcal{C}, \\
\tilde{\mu} & = - \int_0^T F_\lambda(\lambda, u)^* I_V \tilde{w} dt, \\
\tilde{v}_0 & = -\tilde{h} \\
\tilde{\beta} & = - \int_0^T f_\alpha(\alpha, u)^* I_V \tilde{w} dt, \\
\tilde{v} & = I_V^{-1} \tilde{z}, \\
\tilde{g} & = - \frac{1}{2\pi^{(n+1)/2}} \mathcal{F}^{-1} \left[(1 + |\kappa|^2)^{-\ell} (1 + |\omega|^2)^{-r} \left(\int_0^T \int_\Omega e^{-i\kappa \cdot \alpha} e^{-i\omega u(x, t)} (I_V \tilde{w})(x, t) dx dt \right) \right],
\end{aligned} \tag{45}$$

where \tilde{z} solves the two point boundary value problem

$$\begin{aligned}
\ddot{\tilde{z}} - D_V I_V^{-1} \tilde{z} & = I_V \dot{\tilde{w}} + F_u(\lambda, u)^* I_V \tilde{w} + f_u(\alpha, u)^* I_V \tilde{w} - M^* \tilde{y} \\
\dot{\tilde{z}}(0) & = I_V \tilde{w}(0) + \tilde{h}, \quad \dot{\tilde{z}}(T) = I_V \tilde{w}(T)
\end{aligned} \tag{46}$$

and $F_u(\lambda, u)^*, f_u(\alpha, u)^* : V \rightarrow V^*, F_\lambda(\lambda, u)^* : V \rightarrow X, f_\alpha^*(\alpha, u) : V \rightarrow \mathbb{R}^n$ and $M^* : Y \rightarrow V^*$ are Banach space adjoints. (Note that in case $V = H_0^1(\Omega)$, we have $D_V = I_V^{-1} = -\Delta$ and so the above is a wave equation with the bi-Laplace operator.)

3.3 Discussion of the Assumptions for Application 1

We focus on the special case from the Application (1), (2), that is,

$$\mathbb{F}(\mathbf{x}) = \mathbb{F}(c, \varphi, u_0, u, f) = \begin{pmatrix} \dot{u} - \Delta u + cu + h(u) - f(u) - \varphi \\ u(0) - u_0 \\ Mu \end{pmatrix}, \quad (47)$$

with

$$H = L^2(\Omega), \quad V = H_0^1(\Omega), \quad X = X_c \times X_\varphi, \quad X_c = L^2(\Omega), \quad X_\varphi = V^*, \quad (48)$$

(cf. (43), (44) for the resulting spaces \mathcal{V} , \mathcal{W} , \mathcal{Y} , \mathcal{C}) and the known nonlinearity $h \in W^{2,\infty}(B)$.

At the end of this section, we will conclude convergence of Landweber iteration and also of Tikhonov regularization for this application from the analysis of the requirements in the following Sections 3.3.1–3.3.3.

3.3.1 Tangential cone condition

$$\begin{aligned} & \|\mathbb{F}(u, f) - \mathbb{F}(\tilde{u}, \tilde{f}) - \mathbb{F}'(u, f)(u - \tilde{u}, f - \tilde{f})\|_{\mathcal{W} \times H \times \mathcal{Y}} \\ & \leq \|(c - \tilde{c})(u - \tilde{u})\|_{\mathcal{W}} + \|h(u) - h(\tilde{u}) - h'(u)(u - \tilde{u})\|_{\mathcal{W}} \\ & \quad + \|f(u) - \tilde{f}(\tilde{u}) - f'(u)(u - \tilde{u}) - (f - \tilde{f})(u)\|_{\mathcal{W}} \\ & = I + II + III, \end{aligned}$$

where

$$\begin{aligned} III & = \|\tilde{f}(u) - \tilde{f}(\tilde{u}) - f'(u)(u - \tilde{u})\|_{\mathcal{W}} \\ & = \left\| \int_0^1 \left(\tilde{f}'(u + \theta(\tilde{u} - u)) - f'(u) \right) d\theta (u - \tilde{u}) \right\|_{\mathcal{W}} \\ & = \left\| \int_0^1 \left(\tilde{f}'(u + \theta(\tilde{u} - u)) - \tilde{f}'(u) \right) d\theta (u - \tilde{u}) + (f - \tilde{f})'(u)(u - \tilde{u}) \right\|_{\mathcal{W}} \\ & = \left\| \int_0^1 \int_0^1 \tilde{f}''(u + s\theta(\tilde{u} - u)) ds \theta d\theta (u - \tilde{u})^2 + (f - \tilde{f})'(u)(u - \tilde{u}) \right\|_{\mathcal{W}}. \end{aligned}$$

So with full observations $Mu = u$ and a choice of spaces

$$\mathcal{Y} = L^2(0, T; L^2(\Omega)), \quad L^p(0, T; L^p(\Omega)) \subseteq \mathcal{W} \text{ for } p \in \{1, 2\}, \quad \mathcal{C} \subseteq W^{1,\infty}(B) \quad (49)$$

with the embedding constant C_p , $p \in \{1, 2\}$, and $\text{supp}(u) \cup \text{supp}(\tilde{u}) \subset B$, we obtain

$$\begin{aligned} & \|\mathbb{F}(u, f) - \mathbb{F}(\tilde{u}, \tilde{f}) - \mathbb{F}'(u, f)(u - \tilde{u}, f - \tilde{f})\|_{\mathcal{W} \times H \times \mathcal{Y}} \\ & \leq C_1 \|c - \tilde{c}\|_{L^2} \|u - \tilde{u}\|_{L^1(L^2)} + \frac{1}{2} C_1 \|h''\|_{L^\infty(B)} \|u - \tilde{u}\|_{L^2(L^2)}^2 \\ & \quad + \frac{1}{2} C_1 \|\tilde{f}''\|_{L^\infty(B)} \|u - \tilde{u}\|_{L^2(L^2)}^2 + C_2 \|(f - \tilde{f})'\|_{L^\infty(B)} \|u - \tilde{u}\|_{L^2(L^2)} \\ & \leq \left(C_1 \sqrt{T} \|c - \tilde{c}\|_{L^2} + \frac{1}{2} C_1 (\|h''\|_{L^\infty(B)} + \|\tilde{f}''\|_{L^\infty(B)}) C_{\mathcal{V} \rightarrow \mathcal{Y}} \|u - \tilde{u}\|_{\mathcal{V}} \right. \\ & \quad \left. + C_2 \|f - \tilde{f}\|_{\mathcal{C}} \right) \|Mu - M\tilde{u}\|_{\mathcal{Y}} \\ & \leq c_{tc} \|\mathbb{F}(u, f) - \mathbb{F}(\tilde{u}, \tilde{f})\|_{\mathcal{W} \times H \times \mathcal{Y}} \end{aligned} \quad (50)$$

for all $(\tilde{u}, \tilde{f}) = (u_N^\dagger, f_N^\dagger)$, $N \in \mathbb{N}$, $(u, f) \in B_R^{\mathcal{V} \times \mathcal{C}}(u^\dagger, f^\dagger)$, provided R , and hence ρ , are small enough so that

$$\sup_{N \in \mathbb{N}} \left(C_1 \sqrt{T} + \frac{1}{2} C_1 (\|h''\|_{L^\infty(B)} + \|f_N^{\dagger\prime\prime}\|_{L^\infty(B)}) C_{\mathcal{V} \rightarrow \mathcal{Y}} + C_2 \right) (R + \bar{d}) \leq c_{tc}.$$

3.3.2 Weak sequential closedness of $\mathbf{x} \mapsto -\mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{y})$

In this section, we study the more general case, namely the application (47) with $f = f(\alpha, u)$. We will derive weak closedness via weak continuity.

In the following, we frequently employ the embeddings [48, Theorems 1.20, 1.21, Lemmas 7.3, 7.7], [3, Chapter 4], [37, Chapter 11] cf. (43), (48)

$$\begin{aligned} \mathcal{V} &\hookrightarrow C(0, T; L^2(\Omega)), & \mathcal{V} &\hookrightarrow L^2(0, T; L^{6-\epsilon'}(\Omega)), & 0 < \epsilon' \leq 5, & \mathcal{W}^* &\hookrightarrow L^2(0, T; L^6(\Omega)), \\ \mathcal{C} &= H^s(\mathbb{R}^{n+1}) \hookrightarrow C_b(\mathbb{R}^{n+1}), & & & s > (n+1)/2 \\ \mathcal{C} &= H^s(\mathbb{R}^{n+1}) \hookrightarrow W^{2,\infty}(\mathbb{R}^{n+1}), & & & s > (n+1)/2 + 2 \end{aligned}$$

as well as the Hölder inequalities

$$\int_{\Omega} abc \, dx \leq \|a\|_{L^{3/2}} \|b\|_{L^6} \|c\|_{L^6}, \quad \int_{\Omega} abc \, dx \leq \|a\|_{L^2} \|b\|_{L^3} \|c\|_{L^6}.$$

Let $u_n \xrightarrow{\mathcal{V}} u$, $f_n \xrightarrow{\mathcal{C}} f$, $(c, \varphi)_n \xrightarrow{\mathcal{X}} (c, \varphi)$, $(u_0)_n \xrightarrow{H} u_0$, $\alpha_n \xrightarrow{\mathbb{R}^n} \alpha$. We first show weak continuity of the model operator $\mathbf{x} \mapsto (\mathbb{F}(\mathbf{x}) - y)$.

Proposition 4. *The operator \mathbb{F} defined by (47) is weakly continuous on the spaces (48).*

Proof. Assuming

$$|h(x) - h(y)| \leq C|x - y|^{1-\epsilon}(1 + |x|^{4/3} + |y|^{4/3}), \forall x, y \in \mathbb{R} \quad (51)$$

for some $C > 0$, $0 < \epsilon < 1$, we have

$$|\langle h(u_n) - h(u), v \rangle_{\mathcal{W}, \mathcal{W}^*}| \leq \underbrace{C(\|u\|_{C(L^2)}^{\frac{4}{3}}, \|u_n\|_{C(L^2)}^{\frac{4}{3}})}_{< \infty} \underbrace{\|u_n - u\|_{L^2(L^{6-6\epsilon})}^{1-\epsilon}}_{\rightarrow 0} \|v\|_{L^2(L^6)} \rightarrow 0.$$

Next,

$$\begin{aligned} &|\langle f_n(\alpha_n, u_n) - f(\alpha, u), v \rangle_{\mathcal{W}, \mathcal{W}^*}| \\ &= \left| \int_0^T \int_{\Omega} (f_n - f)(\alpha, u) v \, dx \, dt + \int_0^T \int_{\Omega} (f_n(\alpha_n, u_n) - f_n(\alpha, u)) v \, dx \, dt \right| \\ &\leq \underbrace{\left| \int_0^T \int_{\Omega} (f_n - f)(\alpha, u) v \, dx \, dt \right|}_{=: A_n} + \underbrace{\|(f_n)'_{\alpha, u}\|_{L^\infty(\mathbb{R}^{n+1})}}_{< \infty} \underbrace{(|\alpha_n - \alpha| + \|u_n - u\|_{L^2(L^2)})}_{\rightarrow 0} \|v\|_{L^2(L^2)}. \end{aligned}$$

In A_n , for fixed $u \in \mathcal{V}$ and each $v \in \mathcal{W} \subset L^1((0, T) \times \Omega)$, we observe that $\mu_v \in (L^\infty(\mathbb{R}^{n+1}))^*$ with $\|\mu_v\| = \|v\|_{L^1((0, T) \times \Omega)}$ by defining $\mu_v := \int_0^T \int_\Omega (\cdot)(\alpha, u)v \, dx \, dt$. Since $f_n \xrightarrow{L^\infty(\mathbb{R}^{n+1})} f$, it yields $A_n = \mu_v(f_n - f) \rightarrow 0$. Now, the bilinear term is estimated as

$$\begin{aligned} |\langle c_n u_n - cu, v \rangle_{\mathcal{W}, \mathcal{W}^*}| &= \left| \int_0^T \int_\Omega c_n (u_n - u)v \, dx \, dt + \int_0^T \int_\Omega (c_n - c)uv \, dx \, dt \right| \\ &\leq \underbrace{\|c_n\|_{L^2}}_{< \infty} \underbrace{\|u - u_n\|_{L^2(L^3)}}_{\rightarrow 0} \|v\|_{L^2(L^6)} + \left| \int_\Omega \underbrace{(c_n - c)}_{\rightarrow 0 \text{ in } L^2(\Omega)} \underbrace{\int_0^T uv \, dt}_{\in L^2(\Omega)} \, dx \right| \rightarrow 0. \end{aligned} \quad (52)$$

Weak continuity of the remaining part $(u, \varphi, u_0) \mapsto (\dot{u} - \Delta u - \varphi, u(0) - u_0)$ is straightforward, as it is a linear, bounded operator from $\mathcal{V} \times X_\varphi \times H$ to $\mathcal{W} \times H$. Altogether, we claim weak continuity of $\mathbf{x} \mapsto (\mathbb{F}(\mathbf{x}) - \mathbf{y}) := \tilde{w} \in \mathcal{W}$, thus of $\mathbf{x} \mapsto I_V \tilde{w} \in L^2(0, T; V)$.

◇

With this result, we now study the weak sequential continuity of $\mathbf{x} \mapsto -\mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{y})$.

Weak continuity of $\mathbf{x} \mapsto \tilde{\mu}$ in (45). For $\lambda = (c, \varphi)$, $\tilde{\mu} = (\tilde{\mu}_c, \tilde{\mu}_\varphi)$, $\tilde{\mu}_c = -\int_0^T F_c(\lambda, u)^* I_V \tilde{w} \, dt$ where $\tilde{w} = \mathbb{F}(\mathbf{x}) - \mathbf{y}$ as above, we write

$$\begin{aligned} |\langle \tilde{\mu}_c^n - \tilde{\mu}_c, \tilde{c} \rangle_{X_c}| &:= \left| \int_0^T \int_\Omega (u_n I_V \tilde{w}^n - u I_V \tilde{w}) \tilde{c} \, dx \, dt \right| \\ &\leq \underbrace{\|u_n - u\|_{L^2(L^{6-\epsilon})}}_{\rightarrow 0} \underbrace{\|(I_V \tilde{w}^n) \tilde{c}\|_{L^2(L^{(6-\epsilon)/(5-\epsilon)})}}_{< \infty} + \left| \int_0^T \int_\Omega \underbrace{u \tilde{c}}_{\in L^2(V^*)} \underbrace{I_V (\tilde{w}^n - \tilde{w})}_{\rightarrow 0 \text{ in } L^2(V)} \, dx \, dt \right| \rightarrow 0 \end{aligned}$$

for any $\tilde{c} \in X_c$, thus showing weak continuity of $\mathbf{x} \mapsto \tilde{\mu}_c(\mathbf{x})$. Weak continuity of $\mathbf{x} \mapsto \tilde{\mu}_\varphi(\mathbf{x})$ could be obtained in a similar way, replacing $F_\varphi(\lambda, u)^* = \text{Id}$.

Weak continuity of $\mathbf{x} \mapsto \tilde{v}_0$ in (45). As $(u, u_0) \mapsto \tilde{v}_0 := u(0) - u_0$ is linear and bounded, its weak continuity is clear.

Weak continuity of $\mathbf{x} \mapsto \tilde{\beta}$ in (45). We consider

$$\begin{aligned} |\langle \tilde{\beta}^n - \tilde{\beta}, \zeta \rangle_{R^n}| &:= \left| -\int_0^T \int_\Omega (f_n)'_\alpha(\alpha_n, u_n) \zeta I_V \tilde{w}^n \, dx \, dt + \int_0^T \int_\Omega f'_\alpha(\alpha, u) \zeta I_V \tilde{w} \, dx \, dt \right| \\ &= \left| -\int_0^T \int_\Omega [(f_n)'_\alpha(\alpha_n, u_n) - (f_n)'_\alpha(\alpha, u)] \zeta I_V \tilde{w}^n \, dx \, dt - \int_0^T \int_\Omega f'_\alpha(\alpha, u) \zeta I_V (\tilde{w}^n - \tilde{w}) \, dx \, dt \right. \\ &\quad \left. - \int_0^T \int_\Omega [(f_n)'_\alpha(\alpha, u) - f'_\alpha(\alpha, u)] \zeta I_V \tilde{w}^n \, dx \, dt \right| \end{aligned}$$

$$\leq \underbrace{\|(f_n)'_\alpha\|_{W^{1,\infty}(\mathbb{R}^{n+1})}}_{< \infty} \underbrace{(|\alpha - \alpha_n| + \|u - u_n\|_{C(L^2)})}_{\rightarrow 0} \underbrace{\|\zeta I_V \tilde{w}^n\|_{L^1(L^2)}}_{< \infty} \quad (53)$$

$$+ \left| \int_0^T \int_\Omega \underbrace{f'_\alpha(\alpha, u)}_{\in L^2(V^*)} \underbrace{\zeta I_V(\tilde{w}^n - \tilde{w})}_{\rightarrow 0 \text{ in } L^2(V)} dx dt \right| \quad (54)$$

$$+ \underbrace{\|\zeta I_V \tilde{w}^n\|_{L^2(L^2)}}_{< \infty} \underbrace{\sqrt{\int_0^T \int_\Omega |(f_n)'_\alpha(\alpha, u) - f'_\alpha(\alpha, u)|^2 dx dt}}_{:= A'_n}$$

Regarding A'_n , let us fix (α, u) , then set $\Omega_T := \{(t, x) \in (0, T) \times \Omega : |u(t, x)| < \infty\}$. Now Ω_T has nonzero measure, as $u \in \mathcal{V} \subset L^1((0, T) \times \Omega)$. Moreover, $|((0, T) \times \Omega) \setminus \Omega_T| = 0$. Next, for each $(t, x) \in \Omega_T$, the functional defined by $\mu_{t,x} := (\cdot)(\alpha, u(t, x))$ belongs to $C_b(\mathbb{R}^{n+1})^*$ with $\|\mu_{t,x}\| = 1$. From this, we ascertain

$$(f_n - f)'_\alpha(\alpha, u)(x, t) = (f_n - f)'_\alpha(\alpha, u(x, t)) = \mu_{t,x}((f_n - f)'_\alpha) \rightarrow 0 \quad \forall^{\text{a.e.}} (t, x) \in (0, T) \times \Omega$$

for $(f_n)'_\alpha \rightharpoonup f'_\alpha$ in $C_b(\mathbb{R}^{n+1})$. In addition,

$$\sup_n \|(f_n)'_\alpha(\alpha, u)\|_{L^\infty((0,T) \times \Omega)} + \|f'_\alpha(\alpha, u)\|_{L^\infty((0,T) \times \Omega)} \leq \sup_n \|(f_n)'_\alpha\|_{L^\infty(\mathbb{R}^{n+1})} + \|f'_\alpha\|_{L^\infty(\mathbb{R}^{n+1})} < \infty.$$

Applying the Dominated Convergence Theorem yields $A'_n = \|(f_n)'_\alpha(\alpha, u) - f'_\alpha(\alpha, u)\|_{L^2((0,T) \times \Omega)} \rightarrow 0$. Note that this argument remains valid even for $\|\cdot\|_{L^p((0,T) \times \Omega)}$, $1 \leq p < \infty$. This demonstrates the weak convergence of $\mathbf{x} \mapsto \tilde{\beta}(\mathbf{x})$.

Weak continuity of $\mathbf{x} \mapsto \tilde{v} = I_V^{-1} \tilde{z}$ in (45). In the first step, testing (46) with $\Delta \tilde{z}$ where $\tilde{z} \in \mathcal{V}' := L^2(0, T; H^3(\Omega) \cap H_0^2(\Omega)) \cap H^1(0, T; H^1(\Omega))$, yields

$$\begin{aligned} & \|\nabla \dot{\tilde{z}}\|_{L^2(L^2)}^2 + \|\nabla \Delta \tilde{z}\|_{L^2(L^2)}^2 \\ &= \langle \tilde{\varphi}, \Delta z \rangle_{\mathcal{V}^*, \mathcal{V}} + \int_\Omega (I_V \tilde{w}(t) - \dot{\tilde{w}}(t)) \Delta z(t) \Big|_{t=0}^{t=T} dx - \int_0^T \int_\Omega I_V \tilde{w} \Delta \dot{\tilde{z}} dx dt \\ &\leq \|\tilde{\varphi}\|_{\mathcal{V}^*} \|\Delta \tilde{z}\|_{\mathcal{V}} + \|\tilde{h}\|_{L^2} \|\Delta \tilde{z}(0)\|_{L^2} + \|\nabla I_V \tilde{w}\|_{L^2(L^2)} \|\nabla \dot{\tilde{z}}\|_{L^2(L^2)} \\ &\leq \left(\|\tilde{\varphi}\|_{\mathcal{V}^*} + C_{\mathcal{V}' \rightarrow C(H^2)} \|\tilde{h}\|_{L^2} + \|I_V \tilde{w}\|_{L^2(V)} \right) \|\tilde{z}\|'_{\mathcal{V}}, \end{aligned} \quad (55)$$

where $\tilde{\varphi} := F'_u(\lambda, u)^* I_V \tilde{w} + f_u(\alpha, u)^* I_V \tilde{w} - M^* \tilde{y}$ is the right hand side of the wave equation (46) without the first term $I_V \dot{\tilde{w}}$. Here $F'_u(\lambda, u) = -\Delta + c + h'(u)$ under the assumption $|h'(x) - h'(y)| \leq C|x - y|^{1-\epsilon}(1 + |x|^{1/3} + |y|^{1/3})$, $\forall x, y \in \mathbb{R}$.

As previously, when $\mathbf{x}_n \rightharpoonup \mathbf{x}$, one has $I_V \tilde{w}^n \rightharpoonup I_V \tilde{w}$ in $L^2(0, T; V)$ and $\tilde{h}_n \rightharpoonup \tilde{h}$ in $L^2(\Omega)$. We now show $\tilde{\varphi}_n \rightharpoonup \tilde{\varphi}$ in \mathcal{V}^* . Indeed,

$$\langle -\Delta^*(I_V \tilde{w}^n - I_V \tilde{w}), v \rangle_{\mathcal{V}^*, \mathcal{V}} = - \int_0^T \int_\Omega \underbrace{I_V(\tilde{w}^n - \tilde{w})}_{\rightarrow 0 \text{ in } L^2(0,T;V)} \underbrace{\Delta v}_{\in \mathcal{W}} dx dt \rightarrow 0,$$

$\langle c_n I_V \tilde{w}^n - c I_V \tilde{w}, v \rangle_{\mathcal{V}^*, \mathcal{V}} \rightarrow 0$ similarly to (52) with $I_V \tilde{w}$ in place of u ,

$|\langle h'(u_n)^* I_V \tilde{w}^n - h'(u)^* I_V \tilde{w}, v \rangle_{\mathcal{V}^*, \mathcal{V}}|$

$$\leq \underbrace{\|h'(u_n) - h'(u)\|_{L^2(L^3)}}_{\leq C(\|u\|_{C(L^2)}^{\frac{1}{3}})\|u_n - u\|_{L^2(L^{6-6\epsilon})}^{1-\epsilon} \rightarrow 0} \underbrace{\|I_V \tilde{w}^n\|_{L^2(L^6)} \|v\|_{C(L^2)}}_{< \infty} + \left| \int_0^T \int_{\Omega} \underbrace{I_V(\tilde{w}^n - \tilde{w})}_{\rightarrow 0 \text{ in } L^2(0,T;V)} \underbrace{h'(u)v}_{\in \mathcal{W}} dx dx \right|,$$

$\langle f'_u(\alpha_n, u_n)^* I_V \tilde{w}^n - f'_u(\alpha, u)^* I_V \tilde{w}, v \rangle_{\mathcal{V}^*, \mathcal{V}} \rightarrow 0$

similarly to (53) with f'_u in place of f'_α , v in place of ζ .

The last estimate is analogous to (53), but modifies the upper bound for the term involving A'_n to $\|v\|_{C(L^2)} \|I_V \tilde{w}^n\|_{L^2(L^6)} \|(f_n)'_u(\alpha, u) - f'_u(\alpha, u)\|_{L^2(L^3)}$. This yields $\tilde{\varphi}_n \rightarrow \tilde{\varphi}$ in \mathcal{V}^* when $\mathbf{x} \rightarrow 0$, as claimed.

In order to form the full \mathcal{V}' -norm on the left hand side of (55), we test (46) by \tilde{z} . By then applying Young's inequality with $\epsilon > 0$, we eventually obtain

$$(1 - 3\epsilon) \|\tilde{z}_n - \tilde{z}\|_{\mathcal{V}'}^2 \leq \frac{1}{\epsilon} \left(\|\tilde{\varphi}_n - \tilde{\varphi}\|_{\mathcal{V}^*}^2 + (C_{\mathcal{V}' \rightarrow C(H^2)})^2 \|\tilde{h}_n - \tilde{h}\|_{L^2}^2 + \|I_V(\tilde{w}^n - \tilde{w})\|_{L^2(V)}^2 \right). \quad (56)$$

Using Galerkin approximation, one can show that for each $(\tilde{\varphi}, \tilde{h}, \tilde{w}) \in \mathcal{V}^* \times L^2(\Omega) \times \mathcal{W}$, there exists a unique $\tilde{z} \in \mathcal{V}'$ solving (46). Moreover, \tilde{z} depends continuously on the data $(\tilde{\varphi}, \tilde{h}, \tilde{w})$ through the expression (56). Since $\mathcal{V}^* \times L^2(\Omega) \times \mathcal{W} \ni (\tilde{\varphi}, \tilde{h}, \tilde{w}) \mapsto \tilde{z} \in \mathcal{V}'$ is linear and bounded, it is weakly continuous. In conclusion, when $\mathbf{x}_n \rightarrow \mathbf{x}$, we have $\tilde{z}_n \rightarrow \tilde{z}$ in \mathcal{V}' , equivalently $I_V^{-1} \tilde{z}_n \rightarrow I_V^{-1} \tilde{z}$ in \mathcal{V} , proving weak continuity of $\mathbf{x} \mapsto \tilde{v} := I_V^{-1} \tilde{z}$ in (45).

Weak continuity of $\mathbf{x} \mapsto \tilde{g}$ in (45). For \tilde{g} as in (45) and setting $C_\pi := \frac{1}{2\pi^{(n+1)/2}}$, we evaluate, for any $\psi \in \mathcal{C}$

$$\begin{aligned} \langle \tilde{g}^n - \tilde{g}, \psi \rangle_{\mathcal{C}} &= \int_{\mathbb{R}^{n+1}} (1 + |\kappa|^2)^\ell (1 + |\omega|^2)^r \overline{\mathcal{F}\psi(\kappa, \omega)} \mathcal{F}(\tilde{g}^n - \tilde{g})(\kappa, \omega) d\kappa d\omega \\ &= C_\pi \int_{\mathbb{R}^{n+1}} \overline{\mathcal{F}\psi(\kappa, \omega)} \left[\int_0^T \int_{\Omega} e^{-i\kappa \cdot \alpha_n - i\omega u_n(x,t)} I_V \tilde{w}^n - e^{-i\kappa \cdot \alpha - i\omega u(x,t)} I_V \tilde{w} dx dt \right] d\kappa d\omega \\ &= C_\pi \int_{\mathbb{R}^{n+1}} \overline{\mathcal{F}\psi(\kappa, \omega)} \left[- \int_0^T \int_{\Omega} \underbrace{e^{-i\kappa \cdot \alpha - i\omega u}}_{\in L^2(V^*)} \underbrace{I_V(\tilde{w}^n - \tilde{w})}_{\rightarrow 0 \text{ in } L^2(V)} dx dt \right. \\ &\quad \left. + \int_0^T \int_{\Omega} \underbrace{\int_0^1 e^{-i\kappa \cdot (\alpha - \theta(\alpha_n - \alpha)) - i\omega(u + \theta(u_n - u))} d\theta}_{\in [-1, 1]} \underbrace{(-i\kappa \cdot (\alpha_n - \alpha) - i\omega(u_n - u)) I_V \tilde{w}^n}_{\rightarrow 0 \text{ in } L^1(0,T;L^1(\Omega))} dx dt \right] d\kappa d\omega \\ &=: C_\pi \int_{\mathbb{R}^{n+1}} \overline{\mathcal{F}\psi(\kappa, \omega)} B_n(\kappa, \omega) d\kappa d\omega. \end{aligned}$$

and deduce pointwise convergence in (κ, ω) of B_n . Together with uniform boundedness via

$$\begin{aligned} & \int_{\mathbb{R}^{n+1}} \overline{\mathcal{F}\psi(\kappa, \omega)} B_n(\kappa, \omega) d\kappa d\omega \\ & \leq \int_{\mathbb{R}^{n+1}} |(1 + |\kappa|^2)^{1/2} (1 + |\omega|^2)^{1/2} \overline{\mathcal{F}\psi(\kappa, \omega)}| \underbrace{\left(\frac{C}{(1 + |\kappa|^2)^{1/2} (1 + |\omega|^2)^{1/2}} \right)}_{\in L^2(\mathbb{R}^{n+1})} d\kappa d\omega \\ & \leq C \|\psi\|_{H^1(\mathbb{R}^{n+1})} \leq C \|\psi\|_c, \end{aligned}$$

and application of the Dominated Convergence Theorem, we conclude $\langle \tilde{g}^n - \tilde{g}, \psi \rangle_c \rightarrow 0$, yielding weak continuity of $\mathbf{x} \mapsto \tilde{g}(\mathbf{x})$.

3.3.3 Lipschitz continuity of $\mathbf{x} \mapsto \mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{z})$

Above, we have verified weak continuity of $\mathbf{x} \mapsto -\mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{y})$, where all the estimates were written in the form of $\|\mathbf{x}_n - \mathbf{x}\|$. Therefore, Lipschitz continuity of $\mathbb{F}'(\mathbf{x})^*(\mathbb{F}(\mathbf{x}) - \mathbf{z})$ could be established in the same manner, the Lipschitz constant $L < 2$ being obtained in the ball $B_R(\mathbf{x}^\dagger)$ with sufficiently small R .

As a consequence, we can conclude from Propositions 4, 3 the following convergence results on Tikhonov regularization and Landweber iteration.

Corollary 2 (Tikhonov). *For the operator \mathbb{F} defined by (47) on the spaces (43), (48), (49), $\mathcal{C} = H^s(\mathbb{R})$, $s > \frac{5}{2}$ with $M : \mathcal{V} \rightarrow \mathcal{Y}$ linear and bounded and h satisfying (51) we have subsequential convergence of $(c^{\gamma(\delta), \delta, N(\delta)}, \varphi^{\gamma(\delta), \delta, N(\delta)}, u_0^{\gamma(\delta), \delta, N(\delta)}, u^{\gamma(\delta), \delta, N(\delta)}, f^{\gamma(\delta), \delta, N(\delta)})$ to a solution of the inverse problem (1), (2) as $\delta \rightarrow 0$.*

Corollary 3 (Landweber). *Let the assumptions of Corollary 2 hold and additionally assume full measurements $M = id_{\mathcal{Y} \rightarrow \mathcal{Y}}$ and (38). Then we have weak subsequential convergence of $(c_{N(\delta), k_*(\delta)}^\delta, \varphi_{N(\delta), k_*(\delta)}^\delta, u_{0, N(\delta), k_*(\delta)}^\delta, u_{N(\delta), k_*(\delta)}^\delta, f_{N(\delta), k_*(\delta)}^\delta)$ to a solution of the inverse problem (1), (2) as $\delta \rightarrow 0$.*

4 Outlook

In this study, we have carried out a convergence analysis for discretizations of Tikhonov and projected Landweber regularization using Neural Networks. The convergence analysis is based on a priori choice of the regularization and discretization parameters, where the latter relates to the network approximation error. Our analysis is applicable not only for discretization by NNs, but also for general discretization schemes. As an application, we have presented a parameter identification problem for a time-dependent PDE, whose unknown nonlinearity is approximated by a neural network. Our all-at-once approach does not require a training process for learning the nonlinearities beforehand, instead simultaneously determining it alongside the unknown coefficients and the solution of the PDE.

This paper focuses on the theoretical aspects. Numerical results for the regularization with neural networks can be found in [1]. Also in [1], further details on the discretized problem are discussed, such as differentiability of the forward mapping, unique existence for the *learning-informed PDEs* (NN as a reaction term in the PDE), the tangential cone condition for networks and so forth. A potential extension to our study is the inclusion of further components in the unknown nonlinear response, e.g. $f(\alpha, u, \nabla u, \nabla^2 u \dots)$, for more flexible models, as was done in [10, 38, 46].

On the analytical side, an open problem is determining convergence rates for Tikhonov regularization, based on quantified approximation results for neural networks as in [13] (see also Discussion 1 in the introduction). These rates require enhanced regularity of the exact solution in terms of so-called source conditions, whose interpretation for the problem setting considered here is another interesting task.

Acknowledgments. The work of the first author was supported by the Austrian Science Fund FWF under the grants P30054 and DOC 78. Moreover, we wish to thank both reviewers for fruitful comments leading to an improved version of the manuscript.

References

- [1] C. Aarset, M. Holler, and T. T. N. Nguyen. Learning-informed parameter identification in nonlinear time-dependent PDEs. arXiv:2202.10915 [math.OC].
- [2] R. Acar and C. R. Vogel. Analysis of bounded variation penalty methods for ill-posed problems. *Inverse Problems*, 10(6):1217–1229, 1994.
- [3] R.A. Adams and J.F. Fourier. *Sobolev Spaces*. Elsevier, Oxford, 2003.
- [4] Simon R. Arridge, Peter Maass, Ozan Öktem, and Carola-Bibiane Schönlieb. Solving inverse problems using data-driven models. *Acta Numer.*, 28:1–174, 2019.
- [5] Andrea Aspri, Yury Korolev, and Otmar Scherzer. Data driven regularization by projection. *Inverse Problems*, 36(12):125009, dec 2020.
- [6] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, 1993.
- [7] A. R. Barron. Approximation and estimation bounds for artificial neural networks. *Machine learning*, 14(1):115–133, 1994.
- [8] M. Benning and M. J. Ehrhardt. *Lecture notes on Inverse Problems in Imaging*. Online; accessed 2016.
- [9] Kristian Bredies and Martin Holler. Higher-order total variation approaches and generalisations. *Inverse Problems*, 36(12):123001, dec 2020.

- [10] Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *PNAS*, 113(15):3932–3937, 2016.
- [11] M. Burger and W. Mühlhuber. Iterative regularization of parameter identification problems by sequential quadratic programming methods. *Inverse Problems*, 18:943–969, 2002.
- [12] M. Burger and W. Mühlhuber. Numerical approximation of an SQP-type method for parameter identification. *SIAM J. Numer. Anal.*, 40:1775–1797, 2002.
- [13] Martin Burger and Heinz W. Engl. Training neural networks with noisy data as an ill-posed problem. *Advances in Computational Mathematics volume, pages 33*, 13:335–354, 2020.
- [14] Martin Burger and Stanley Osher. *A Guide to the TV Zoo*, pages 1–70. Springer International Publishing, Cham, 2013.
- [15] J. R. Cannon and Paul DuChateau. Structural identification of an unknown source term in a heat equation. *Inverse Problems*, 14(3):535–551, 1998.
- [16] F. Cao, T. Xie, and Z. Xu. The estimate for approximation error of neural networks: a constructive approach. *Neurocomputing*, 71(4–5):626–630, 2008.
- [17] R. DeVore, B. Hanin, and G. Petrova. Neural network approximation. *Acta Numerica*, 30:327–444, 2021.
- [18] Guozhi Dong, Michael Hintermüller, and Kostas Papafitsoros. Optimization with learning-informed differential equation constraints and its applications. arXiv:2008.10893 [math.OA].
- [19] Paul DuChateau and William Rundell. Unicity in an inverse problem for an unknown reaction term in a reaction-diffusion equation. *J. Differential Equations*, 59(2):155–164, 1985.
- [20] H.W. Engl, M. Hanke, and A. Neubauer. *Regularization of Inverse Problems*. Kluwer, Dordrecht, 1996.
- [21] Á. Arroyo G. S. Alberti and M. Santacesaria. Inverse problems on low-dimensional manifolds. arXiv:2009.00574v1 [math.FA].
- [22] P. Grohs and G. Kutyniok. *Mathematical Aspects of Deep Learning*. Cambridge University Press, upcoming.
- [23] Eldad Haber and Uri M Ascher. Preconditioned all-at-once methods for large, sparse parameter estimation problems. *Inverse Problems*, 17(6):1847, 2001.

- [24] M. Hanke, A. Neubauer, and O. Scherzer. A convergence analysis of the Landweber iteration for nonlinear ill-posed problems. *Numer. Math.*, 72:21–37, 1995.
- [25] Kurt Hornik, Tinchcombe Maxwell, and White Halbert. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2:1989, 359–366.
- [26] Victor Isakov. *Inverse source problems*, volume 34 of *Mathematical Surveys and Monographs*. American Mathematical Society, Providence, RI, 1990.
- [27] Victor Isakov. *Inverse Problems for Partial Differential Equations*. Springer, New York, second edition, 2006.
- [28] B. Kaltenbacher. Regularization based on all-at-once formulations for inverse problems. *SIAM Journal of Numerical Analysis*, 54:2594–2618, 2016.
- [29] B. Kaltenbacher, A. Kirchner, and B. Vexler. Goal oriented adaptivity in the IRGNM for parameter identification in PDEs II: all-at once formulations. *Inverse Problems*, 30, 2014. 045002.
- [30] B. Kaltenbacher, A. Neubauer, and O. Scherzer. *Iterative Regularization Methods for Nonlinear Problems*. de Gruyter, Berlin, New York, 2008. Radon Series on Computational and Applied Mathematics.
- [31] Barbara Kaltenbacher. All-at-once versus reduced iterative methods for time dependent inverse problems. *Inverse Problems*, 33, 2017.
- [32] Stefan Kindermann. Convergence of the gradient method for ill-posed problems. *Inverse Problems & Imaging*, 11(4):703–720, 2017.
- [33] K. Kunisch and E. W. Sachs. Reduced SQP methods for parameter identification problems. *SIAM Journal on Numerical Analysis*, 29(6):1793–1820, 1992.
- [34] FS Kupfer and EW Sachs. Numerical solution of a nonlinear parabolic control problem by a reduced SQP method. *Computational Optimization and Applications*, 1(1):113–135, 1992.
- [35] Gitta Kutyniok, Philipp Petersen, Mones Raslan, and Reinhold Schneider. A theoretical analysis of deep neural networks and parametric pdes. *Constructive approximation*, 55:73–125, 2022.
- [36] F. Leibfritz and E. W. Sachs. Inexact SQP interior point methods and large scale optimal control problems. *SIAM Journal on Control and Optimization*, 38(1):272–293, 1999.
- [37] Giovanni Leoni. *A first course in Sobolev spaces*, volume 105 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2009.

- [38] Z. Long, Y. Lu, X. Ma, and B. Dong. Pde-net: Learning pdes from data. *In International Conference on Machine Learning*, pages 3208–3216, PMLR, 2018.
- [39] Jianfeng Lu, Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation for smooth functions, 2021.
- [40] H. Mhaskar. Approximation properties of a multilayered feedforward artificial neural network. *Adv. Comput. Math.*, 1(1):61–80, 1993.
- [41] T. T. N. Nguyen. Landweber-Kaczmarz for parameter identification in time-dependent inverse problems: All-at-once versus reduced version. *Inverse Problems*, 35, 2019. Art. ID. 035009.
- [42] Carlos E Orozco and Omar N Ghattas. A reduced SAND method for optimal design of non-linear structures. *International Journal for Numerical Methods in Engineering*, 40(15):2759–2774, 1997.
- [43] P. C. Petersen and F. Voigtlaender. Optimal approximation of piecewise smooth functions using deep relu neural networks. *Neural Netw.*, 180:296–330, 2018.
- [44] Michael S. Pilant and William Rundell. An inverse problem for a nonlinear parabolic equation. *Comm. Partial Differential Equations*, 11(4):445–457, 1986.
- [45] C. Pöschl, E. Resmerita, and O. Scherzer. Discretization of variational regularization in Banach spaces. *Inverse Probl.*, 26(10):105017, 2010.
- [46] M. Raissi. Deep hidden physics models: Deep learning of nonlinear partial differential equations. *Journal of Machine Learning Research*, 19:24 pp, 2018.
- [47] M. Raissi, P. Perdikaris, and G. E. arniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [48] T Roubíček. *Nonlinear Partial Differential Equations with Applications*. Springer Basel, 2013.
- [49] Hayden Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proc. R. Soc. A*, 473, 2016.
- [50] Thomas Schuster, Barbara Kaltenbacher, Bernd Hofmann, and Kamil S. Kazimierski. *Regularization Methods in Banach Spaces*. Walter de Gruyter, Berlin, 2012.
- [51] A. R. Shenoy, M. Heinkenschloss, and E. M. Cliff. Airfoil design by an all-at-once method. *International Journal for Computational Fluid Mechanics*, 11:3–25, 1998.
- [52] Shlomo Ta’asan. ”one shot” methods for optimal control of distributed parameter systems: I finite dimensional control. Technical report, Institute for Computer Applications in Science and Engineering : NASA Langley Research Center, 1991.

- [53] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc.*, 58:267–288, 1996.
- [54] F. Tröltzsch. *Optimal Control of Partial Differential Equations: Theory, Methods and Applications*. American Mathematical Society, 2010.
- [55] T van Leeuwen and F J Herrmann. A penalty method for PDE-constrained optimization in inverse problems. *Inverse Problems*, 32(1):015007, 2016.
- [56] D. Yarotsky. Error bounds for approximations with deep relu networks. *Neural Netw.*, 94:103–114, 2017.