

On projective stochastic-gradient type methods for solving large scale systems of nonlinear ill-posed equations: Applications to machine learning

J.C. Rabelo[†] Y. Saporito[‡] A. Leitão^{‡§} A. L. Madureira[¶]

January 30, 2025

Abstract

We propose and analyze a stochastic-gradient type method for solving systems of nonlinear ill-posed equations. The method considered here extends the SGD type iteration introduced in [2022, Rabelo et al., Inv. Probl. **38** 025003] for solving linear ill-posed systems.

A distinctive feature of our method resides in the [adaptive choice](#) of the stepsize, which promotes a relaxed orthogonal projection of the current iterate onto a conveniently chosen convex set. This characteristic distinguishes our method from other SGD type methods in the literature (where the stepsize is typically chosen *a priori*) and accounts for the faster convergence observed in the numerical experiments conducted in this manuscript.

The convergence analysis discussed here includes: monotonicity and mean square convergence of the iteration error (exact data case), stability and semi-convergence (noisy data case). In the later case, our method is coupled with an *a priori* stopping rule.

Numerical experiments are presented for two large scale nonlinear inverse problems in machine learning (both with real data): (i) we address, using neural networks, the big data problem of CO-concentration prediction considered in the above cited article; (ii) we tackle the classification problem for the MNIST database (<http://yann.lecun.com/exdb/mnist/>). Additionally, a parameter identification problem in a 3D elliptic PDE system is considered.

Keywords. Ill-posed problems; Nonlinear equations; SGD method; Landweber method; Projective method.

AMS Classification: 65J20, 47J06.

1 Introduction

In these notes we extend to nonlinear ill-posed problems the projective stochastic-gradient type method proposed in [29, 28] for solving large scale systems of linear equations.

Problems under consideration

The *inverse problem* under consideration consists of determining an unknown quantity $x \in X$ from a set of data $(y_0, \dots, y_{N-1}) = (y_i) \in Y^N$, where X and Y are (infinite dimensional) Hilbert spaces. The data (y_i) correspond to indirect observations of the parameter x , this process being

[†]Department of Mathematics, Federal Univ. of Piauí, 64049-550 Teresina, Brazil

[‡]EMAp, Getúlio Vargas Foundation, Praia de Botafogo 190, 22250-900 Rio de Janeiro, Brazil

[§]On leave from Department of Mathematics, Federal Univ. of St. Catarina, 88040-900 Florianópolis, Brazil

[¶]LNCC, Av. Getúlio Vargas 333, P.O. Box 95113, 25651-070 Petrópolis, Brazil

Emails: joelrabelo@ufpi.edu.br, yuri.saporito@gmail.com, acgleitao@gmail.com, alm@lncc.br.

described by the model $y_i = F_i(x)$, for $i = 0, \dots, N-1$; where $F_i : D(F_i) \subset X \rightarrow Y$ are ill-posed nonlinear operators [8, 10].

We are particularly interested in the situation where $N \gg 1$ is large, and the exact data (y_i) are not available; instead, only approximate data (y_i^δ) satisfying

$$\|y_i^\delta - y_i\|_Y \leq \delta_i, \quad i = 0, \dots, N-1, \quad (1)$$

are available. Here $\delta_i > 0$ are known noise levels; we write $\delta := (\delta_0, \dots, \delta_{N-1}) \in \mathbb{R}^N$.

The abstract formulation of the problems under consideration can be summarized as follows: Given inexact data (y_i^δ) and the levels of noise (δ_i) as in (1), find an approximate solution to the large scale system of nonlinear operator equations

$$F_i(x) = y_i^\delta, \quad i = 0, \dots, N-1. \quad (2)$$

A straightforward approach for solving the inverse problem (1), (2) consists in rewriting (2) as a single operator equation $\mathbf{F}(x) = \mathbf{y}^\delta$, with $\mathbf{F} := (F_0, \dots, F_{N-1}) : X \rightarrow Y^N$ and $\mathbf{y}^\delta := (y_0^\delta, \dots, y_{N-1}^\delta)$, and using standard regularization methods; e.g., *Iterative regularization* [1, 8, 14, 20, 22] or *Tikhonov regularization* [8, 27, 31, 32, 33, 30]. When using this functional analytical formulation, dealing with the numerical challenges of solving a high-dimensional ill-posed equation becomes inevitable. Specifically, when applied to $\mathbf{F}(x) = \mathbf{y}^\delta$, the above mentioned regularization methods often become numerically inefficient when $N \gg 1$.

In what follows we briefly discuss alternative approaches for solving the nonlinear inverse problem (1), (2) in a stable manner:

- **Kaczmarz type methods** (cyclic iterations): this technique is considered in [13, 11, 19], [7], [12], [2], [26] and [3] for the Landweber iteration, the Steepest-Descent iteration, the Expectation-Maximization iteration, the Levenberg-Marquardt iteration, the REGINN-Landweber iteration, and the Iteratively Regularized Gauss-Newton iteration respectively;
- **SGD type methods** (non-cyclic iterations): this stochastic technique is considered in [17] with *a priori* chosen stepsize and *a priori* stopping rule (see [16] for the linear case);¹ in [15] with *a priori* chosen stepsize and *a posteriori* stopping rule; in this manuscript **adaptively chosen** stepsize and *a priori* stopping rule (see [29] for the linear case).²

See also [18] for a SGD type method for solving linear ill-posed problems in Banach spaces.

The rationale behind our method

The stochastic-gradient (SGD) type method considered in this manuscript aims to compute, in a stable way, approximate solutions to (1), (2). Our method stands out due to the stepsize selection, which is inspired by the *projective Landweber (PLW) method* [24] and the *projective Landweber-Kaczmarz (PLWK) method* [23]. In the sequel, we briefly address these two methods:

- **The PLW method** was proposed in [24] for solving (1), (2) with $N = 1$, i.e. $F_0(x) = y_0^\delta$ with $\|y_0 - y_0^\delta\| \leq \delta$. A sequence (x_k^δ) is generated as follows: at each iteration k , a half space

$$H_{x_k^\delta} := \{z \in X, \langle z - x_k^\delta, F_0'(x_k^\delta)^* F_\delta(x_k^\delta) \rangle \leq -\|F_\delta(x_k^\delta)\|((1-\eta)\|F_\delta(x_k^\delta)\| - (1+\eta)\delta)\}$$

is defined, where $F_\delta(x) := F_0(x) - y_0^\delta$ (see (A2) in Section 2.1 for the definition of the constant η). Under appropriate assumptions, it is proven that $H_{x_k^\delta}$ contains all solutions of $F_0(x) = y_0$; moreover, if the norm of the residual $\|F_\delta(x_k^\delta)\|$ is above the trashold $(1+\eta)(1-\eta)^{-1}\delta$ then $H_{x_k^\delta}$ does not contain the iterate x_k^δ .³ The next iterate x_{k+1}^δ is defined as a (relaxed) orthogonal

¹That is, the stepsizes are determined prior to computing the iterates; the same applies to the computation of the stopping index.

²I.e. the stepsizes are determined while the iteration is being computed, whereas the stopping index is chosen before the computation of the iterates begins.

³In this situation we say that **the set $H_{x_k^\delta}$ separates the iterate x_k^δ from the solution set $F_0^{-1}(y_0)$** ; in other words $F_0^{-1}(y_0) \subset H_{x_k^\delta}$, while $x_k^\delta \notin H_{x_k^\delta}$.

projection of x_k^δ onto $H_{x_k^\delta}$ (see [24, Eq. (8)] for details). Summarizing, PLW corresponds to a Landweber type iteration [22, 8] with stepsize defined by (relaxed) orthogonal projections onto the separating sets $H_{x_k^\delta}$.

• **The PLWK method** was proposed in [23] for solving systems of nonlinear ill-posed equations as in (1), (2) with $N > 1$. It consists in coupling the PLW method with the Kaczmarz strategy and incorporating a bang-bang parameter. The corresponding iteration formula reads

$$x_{k+1}^\delta := x_k^\delta - \theta_k \lambda_k \omega_k F'_{[k]}(x_k^\delta)^* (F_{[k]}(x_k^\delta) - y_{[k]}^\delta),$$

where $[k] := (k \bmod N) \in \{0, \dots, N-1\}$, $\theta_k \in (0, 2)$ is a relaxation parameter and $\omega_k \in \{0, 1\}$ is a bang-bang parameter (see [23, Eq. (6)]). Moreover, $\lambda_k \geq 0$ (see [23, Eq. (12)]) gives the exact orthogonal projection of x_k^δ onto the half space $H_{[k], x_k^\delta}$, where

$$H_{i, x_k^\delta} := \left\{ z \in X, \langle z - x_k^\delta, F'_i(x_k^\delta)^* F_{i, \delta}(x_k^\delta) \rangle \leq -\|F_{i, \delta}(x_k^\delta)\| \left((1 - \eta) \|F_{i, \delta}(x_k^\delta)\| - (1 + \eta) \delta \right) \right\},$$

for $i = 0, \dots, N-1$; here $F_{i, \delta}(x) := F_i(x) - y_i^\delta$ (see [23, Eq. (11)]). Summarizing, PLWK corresponds to a Landweber-Kaczmarz (cyclic) type iteration [13] with stepsize defined by (relaxed) orthogonal projections onto the separating sets $H_{[k], x_k^\delta}$.

In [29, 28] the projective step of the PLW and PLWK iterations was used as starting point to derive a SGD type method for solving large scale systems of linear ill-posed equations.

The projective stochastic-gradient (pSGD) method

In this manuscript we build upon a well-established nonlinear assumption, namely the *weak tangential cone condition* (wTCC) [14, 8], to expand the method in [29, 28]. As a result, we create a new approach capable of efficiently solving large-scale systems of nonlinear equations of the form (1), (2). For obvious reasons, the method considered in these notes is named *projective stochastic-gradient* (pSGD) method.

Unlike the majority of SGD type methods found in the literature, our approach employs **adaptive stepsize** selection (see (6)). Additionally, in the noisy data case, our iterative method is combined with an *a priori* stopping rule (see (A6)), classifying pSGD as a regularization method as defined in [8].

Outline of the manuscript

In Section 2 we introduce the pSGD method; the main assumptions used in our analysis are presented. Section 3 is dedicated to convergence analysis of pSGD. In Section 3.1 the exact data case is considered: We estimate the *average gain* (Proposition 3.2), and prove monotonicity of the *average iteration error* (Corollary 3.3) as well as square summability of the *average residuals* (Corollary 3.4). Additionally, a convergence result is proven (Theorem 3.5). Section 3.2 is devoted to analyzing the noisy data case and regularization properties of pSGD. The key findings include a stability result (Theorem 3.9) and a semi-convergence result (Theorem 3.11).

In Section 4 numerical experiments are presented for solving two large scale systems of nonlinear equations with real data. Both inverse problems relate to parameter identification in neural network training, in detail:

- In Section 4.1 we tackle the big data problem of CO-concentration prediction in a gas sensor array [9, 29]. A special neural-network (NN) is used to model the related inverse problem (a variation of the saturated linear activation function is used). We prove in Lemma 4.2 that the nonlinear function modeling this NN satisfies the wTCC.

- In Section 4.2 we address the well-known classification problem for the MNIST database, consisting of images of handwritten digits (see https://en.wikipedia.org/wiki/MNIST_database).

- In Section 4.3 a parameter identification problem in a 3D elliptic PDE system is considered. Section 5 is devoted to final remarks and conclusions.

2 The method under investigation

This section presents the nonlinear pSGD method under consideration in this notes. We begin by addressing the main assumptions necessary for the analysis derived in the forthcoming sections.

2.1 Main assumptions

Throughout this work we assume that $\bigcap_i D(F_i)$ has nonempty interior, where $D(F_i) \subset X$ is the domain of definition of F_i . Additionally, the initial guess $x_0 \in X$ satisfies $B_\rho(x_0) \subset \bigcap_{i=0}^{N-1} D(F_i)$ for some $\rho > 0$. Moreover, the following assumptions are used:

(A1) Each operator F_i is Fréchet differentiable with continuous derivative F'_i . Moreover, there exists a constant $C > 0$ such that

$$\|F'_i(x)\| \leq C, \quad i = 0, \dots, N-1, \quad \forall x \in B_\rho(x_0); \quad (3)$$

(A2) The *weak Tangential Cone Condition* (wTCC) holds at $B_\rho(x_0)$, with $0 < \eta < 1$, i.e.

$$\|F_i(\bar{x}) - F_i(x) - F'_i(x)(\bar{x} - x)\|_Y \leq \eta \|F_i(\bar{x}) - F_i(x)\|_Y, \quad (4)$$

for $i = 0, \dots, N-1$ and $\forall x, \bar{x} \in B_\rho(x_0)$;

(A3) There exists $x^* \in B_{\rho/2}(x_0)$ such that $F_i(x^*) = y_i$ for $i = 0, \dots, N-1$, i.e. x^* is a (non necessarily unique) solution of (2) with exact data;

(A4) $(\theta_k) \in \mathbb{R}^+$ is a sequence satisfying $0 < \inf_k \theta_k$ and $\sup_k \theta_k < 2$;

(A5) $\gamma > 0$ is a constant satisfying $\gamma > \frac{1+\eta}{1-\eta}C$, with C as in (A1) and η as in (A2);

(A6) The stopping index $k_\delta^* = k^*(\delta) \in \mathbb{N}$, satisfies $\lim_{\delta \rightarrow 0} k_\delta^* = \infty$.

The following inequalities are immediate consequences of (A2):

$$(1 - \eta)\|F_i(\bar{x}) - F_i(x)\| \leq \|F'_i(x)(\bar{x} - x)\| \leq (1 + \eta)\|F_i(\bar{x}) - F_i(x)\|, \quad (5)$$

for $i = 0, \dots, N-1$ and $x, \bar{x} \in B_\rho(x_0)$ (see [8, Chapter 11] for further discussion).

2.2 Introducing the nonlinear pSGD iteration

Here we present the nonlinear pSGD method for solving (1), (2). In what follows we adopt the simplified notation: $F_i^\delta(x) := F_i(x) - y_i^\delta$, for $i = 0, \dots, N-1$; and define the polynomial function $p^\epsilon(t) := t((1 - \eta)t - (1 + \eta)\epsilon)$, for $\epsilon > 0$.

Given $x_0, \gamma > 0$ and (θ_k) as in Section 2.1, the iteration formula of the pSGD method reads

$$x_{k+1}^\delta = x_k^\delta - \theta_k \lambda_{I_k}^\delta F'_{I_k}(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta), \quad k = 0, 1, \dots, k_\delta^* - 1. \quad (6a)$$

Here the stepsize $\lambda_{I_k}^\delta \geq 0$ is a function of $(x_k^\delta, y_{I_k}^\delta, \delta_{I_k})$ and is defined by

$$\lambda_{I_k}^\delta := \begin{cases} p^{\delta_{I_k}}(\|F_{I_k}^\delta(x_k^\delta)\|) / \|F'_{I_k}(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta)\|^2 & , \text{ if } \|F'_{I_k}(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta)\| > \gamma \delta_{I_k} \\ 0 & , \text{ otherwise.} \end{cases} \quad (6b)$$

The (I_k) is an independent and identically distributed sequence of random indexes, taking values in $\{0, \dots, N-1\}$, in a fixed probability space $(\Omega, \mathcal{F}, \mathbb{P})$.⁴ For simplicity of the presentation we assume that $\mathbb{P}(I_k = i) = \frac{1}{N}$ for $i = 0, \dots, N-1$.

⁴Here Ω is the space of sequences taking values in $\{0, 1, \dots, N-1\}$ and \mathbb{P} is the probability given by the Kolmogorov extension theorem that extends the uniform distributed probability \mathbb{P}_k in $\{0, 1, \dots, N-1\}^k$ for any k . For a more detailed description, we refer the reader to [18, Section 3].

Remark 2.1. In the exact data case, we write $F_i^0(x) := F_i(x) - y_i$ and $p^0(t) := (1 - \eta)t^2$. The iteration formula of the pSGD method reads

$$x_{k+1} = x_k - \theta_k \lambda_{I_k} F'_{I_k}(x_k)^* F_{I_k}^0(x_k), \quad k = 0, 1, \dots \quad (7a)$$

where the stepsize $\lambda_{I_k} \geq 0$ is a function of (x_k, y_{I_k}) and is defined by

$$\lambda_{I_k} := \begin{cases} p^0(\|F_{I_k}^0(x_k)\|) / \|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\|^2 & , \text{ if } \|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| > 0 \\ 0 & , \text{ otherwise.} \end{cases} \quad (7b)$$

Remark 2.2 (Exact projections). If one takes $\theta_k \equiv 1$ in (6a), then x_{k+1}^δ corresponds to the orthogonal projection of x_k^δ onto $H_{I_k, x_k^\delta}^\delta$, where

$$H_{i, x}^\delta := \{z \in X \mid \langle z - x, F'_i(x)^* F_i^\delta(x) \rangle \geq \|F_i^\delta(x)\| ((1 - \eta)\|F_i^\delta(x)\| - (1 + \eta)\delta_i)\}.$$

Alternatively, if $\theta_k \in (0, 2)$ x_{k+1}^δ can be interpreted as a relaxed projection of x_k^δ onto H_{I_k, x_k^δ} .

Remark 2.3 (Lower bound for the stepsizes).

- In the exact data case, it holds $\lambda_{I_k} \geq (1 - \eta)C^{-2}$ whenever $\|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| > 0$. I.e. λ_{I_k} in (7b) is bounded by below, whenever x_k is not a solution of $F_{I_k}(x) = y_{I_k}$.
- In the noisy data case, Assumption (A5) implies $\lambda_{I_k}^\delta \geq C^{-2}(1 - \eta - C(1 + \eta)/\gamma) =: \lambda_{\min}$, whenever $\|F'_{I_k}(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta)\| > \gamma\delta_{I_k}$.

3 Convergence analysis

In this section, analytical properties of the nonlinear pSGD method in (6) are investigated. We start the discussion by considering the case of exact data, i.e. $\delta_i = 0$ for $i = 0, \dots, N - 1$.

3.1 The exact data case

In this case, the inverse problem (1), (2) can be written in the form

$$F_i(x) = y_i, \quad i = 0, \dots, N - 1, \quad (8)$$

or simply $\mathbf{F}(x) = \mathbf{y}$. Next we introduce relevant notation used in this manuscript.

Remark 3.1 (Notation).

- For $x^* \in X$ a solution of (8), the mean square iteration error $\mathbb{E}[\|x^* - x_k\|^2]$ is defined by the average error over all possible realizations of I_0, \dots, I_{k-1} that define x_k . E.g., for $k = 0$ and $k = 1$, it holds

$$\mathbb{E}[\|x^* - x_0\|^2] = \|x^* - x_0\|^2, \quad \mathbb{E}[\|x^* - x_1\|^2] = \frac{1}{N} \sum_{i=0}^{N-1} \|x^* - [x_0 - \theta_1 \lambda_i F'_i(x_0)^* F_i^0(x_0)]\|^2.$$

- Let $k \in \mathbb{N}$ be fixed, and denote by \mathcal{F}_k the σ -algebra generated by I_0, \dots, I_{k-1} . It holds

$$\mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2 | \mathcal{F}_k] = \frac{1}{N} \sum_{i=0}^{N-1} \lambda_i \|F_i^0(x_k)\|^2, \quad \mathbb{E}[\|x^* - x_k\|^2 | \mathcal{F}_k] = \|x^* - x_k\|^2,$$

$$\mathbb{E}[\|x^* - x_{k+1}\|^2 - \|x^* - x_k\|^2 | \mathcal{F}_k] = \frac{1}{N} \sum_{i=0}^{N-1} [\|x^* - [x_k - \theta_k \lambda_i F'_i(x_k)^* F_i^0(x_k)]\|^2 - \|x^* - x_k\|^2].$$

Moreover, by the law of iterated expectation, we have

$$\mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2] = \mathbb{E}[\mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2 | \mathcal{F}_k]] = \frac{1}{N} \sum_{i=0}^{N-1} \mathbb{E}[\lambda_i \|F_i^0(x_k)\|^2],$$

where the last expectation averages the residual of equation i times λ_i over all possible realizations of x_k . It is worth noticing that λ_i is a random variable (indeed, λ_i depends on the realization of x_k ; see (7b)).

In the next proposition we estimate the difference $\mathbb{E}[\|x^* - x_{k+1}\|^2] - \mathbb{E}[\|x^* - x_k\|^2]$, where $x^* \in X$ is a solution of (8). This is a quintessential result in the forthcoming analysis.

Proposition 3.2. *Let assumptions (A1), (A2) and (A3) hold and (x_k) be a sequence generated by the nonlinear pSGD method (7). If $x_k \in B_\rho(x_0)$, then for any x^* solution of (8) it holds*

$$\mathbb{E}[\|x^* - x_{k+1}\|^2] - \mathbb{E}[\|x^* - x_k\|^2] \leq \theta_k(\theta_k - 2)(1 - \eta) \mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2]. \quad (9)$$

Proof. If $F'_{I_k}(x_k)^* F_{I_k}^0(x_k) \neq 0$, we obtain from (7a) and (7b)

$$\begin{aligned} \|x^* - x_{k+1}\|^2 - \|x^* - x_k\|^2 &= 2\langle x_{k+1} - x_k, x_{k+1} - x^* \rangle - \|x_k - x_{k+1}\|^2 \\ &= -2\theta_k \lambda_{I_k} \langle F'_{I_k}(x_k)^* F_{I_k}^0(x_k), x_{k+1} - x^* \rangle - \theta_k^2 \lambda_{I_k}^2 \|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\|^2 \\ &= -2\theta_k \lambda_{I_k} \langle F_{I_k}^0(x_k), -F_{I_k}^0(x_k) - F'_{I_k}(x_k)(x^* - x_k) + F'_{I_k}(x_k)(x_{k+1} - x_k) + F_{I_k}^0(x_k) \rangle \\ &\quad - \theta_k^2 \lambda_{I_k}^2 \|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\|^2 \\ &= 2\theta_k \lambda_{I_k} [\langle -F_{I_k}^0(x_k), -F_{I_k}^0(x_k) - F'_{I_k}(x_k)(x^* - x_k) \rangle + \theta_k p^0(\|F_{I_k}^0(x_k)\|) - \|F_{I_k}^0(x_k)\|^2] \\ &\quad - \theta_k^2 \lambda_{I_k} p^0(\|F_{I_k}^0(x_k)\|). \end{aligned}$$

Since $-F_{I_k}^0(x_k) = F_{I_k}(x^*) - F_{I_k}(x_k)$ and $x_k \in B_\rho(x_0)$, we argue with (A2) to obtain

$$\begin{aligned} \|x^* - x_{k+1}\|^2 - \|x^* - x_k\|^2 &\leq 2\theta_k \lambda_{I_k} \left[(\eta - 1) \|F_{I_k}^0(x_k)\|^2 + \theta_k p^0(\|F_{I_k}^0(x_k)\|) \right] - \theta_k^2 \lambda_{I_k} p^0(\|F_{I_k}^0(x_k)\|) \\ &= 2\theta_k \lambda_{I_k} \left[(\eta - 1) \|F_{I_k}^0(x_k)\|^2 + \frac{1}{2} \theta_k (1 - \eta) \|F_{I_k}^0(x_k)\|^2 \right] \\ &= \theta_k (\theta_k - 2) (1 - \eta) \lambda_{I_k} \|F_{I_k}^0(x_k)\|^2. \end{aligned} \quad (10)$$

Otherwise, if $F'_{I_k}(x_k)^* F_{I_k}^0(x_k) = 0$ then $x_{k+1} = x_k$ and $F_{I_k}(x_k) = y_{I_k}$. Consequently, (10) holds also in this case.

Denoting by \mathcal{F}_k the σ -algebra generated by (I_0, \dots, I_{k-1}) , we conclude that x_k is measurable with respect to \mathcal{F}_k (while I_k is independent of it). Consequently, we derive from (10) the estimate

$$\mathbb{E}[\|x^* - x_{k+1}\|^2 - \|x^* - x_k\|^2 | \mathcal{F}_k] \leq \theta_k (\theta_k - 2) (1 - \eta) \mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2 | \mathcal{F}_k].$$

Averaging over all possible realizations of (I_0, \dots, I_{k-1}) in the last inequality yields (9). \square

If, additionally to (A1), \dots , (A3), assumption (A4) holds, then (10) implies that any sequence (x_k) generated by the pSGD method (7) satisfies $x_k \in B_{\rho/2}(x^*) \subset B_\rho(x_0)$, for $k = 0, 1, \dots$. Thus, under this additional assumption, Proposition 3.2 implies the monotonicity of the mean square iteration error $\mathbb{E}[\|x^* - x_k\|^2]$, where $x^* \in X$ is a solution of (8). This fact is summarized in

Corollary 3.3. *Let assumptions (A1), \dots , (A4) hold true and (x_k) be a sequence generated by the nonlinear pSGD method (7). For any solution x^* of (8) in $B_\rho(x_0)$ it holds*

$$\mathbb{E}[\|x^* - x_{k+1}\|^2] \leq \mathbb{E}[\|x^* - x_k\|^2], \quad k = 0, 1, \dots \quad (11)$$

An important consequence of Proposition 3.2 is discussed in the sequel (this result is used in the proof of the main convergence Theorem 3.5).

Corollary 3.4. *Let assumptions (A1), ..., (A4) hold true and (x_k) be a sequence generated by the nonlinear pSGD method (7). The series*

$$\sum_{k=0}^{\infty} \theta_k(2 - \theta_k)(1 - \eta)\mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2], \quad \sum_{k=0}^{\infty} \theta_k \mathbb{E}[\lambda_I \|F_I^0(x_k)\|^2] \quad \text{and} \quad \sum_{k=0}^{\infty} \mathbb{E}[\|F_I^0(x_k)\|^2]$$

are summable.

Proof. The summability of the first series follows from Proposition 3.2. The summability of the second and third series follows from (A4) together with the summability of the first series and Remark 2.3. \square

We are now ready to state and prove a convergence result for the pSGD method in (7).

Theorem 3.5 (Convergence for exact data). *Let assumptions (A1), ..., (A4) hold. Any sequence (x_k) generated by the nonlinear pSGD method (7) converges in mean square to some random element $x^* \in B_\rho(x_0)$, i.e. $\mathbb{E}[\|x_k - x^*\|^2] \rightarrow 0$ as $k \rightarrow \infty$, and x^* is a solution of (8) almost surely.*

Proof. We claim that (x_k) is a Cauchy sequence. It is enough to prove that $e_k := x^* - x_k$ is Cauchy, where x^* is defined as in (A3). From Corollary 3.3 follows

$$\lim_{k \rightarrow \infty} \mathbb{E}[\|e_k\|^2] = \varepsilon \geq 0, \quad (12)$$

Next we prove that

$$\mathbb{E}[\langle e_n - e_k, e_n \rangle] \rightarrow 0 \quad \text{and} \quad \mathbb{E}[\langle e_l - e_n, e_n \rangle] \rightarrow 0 \quad \text{as} \quad k, l \rightarrow \infty, \quad (13)$$

with $k \leq l$ for some $k \leq n \leq l$ (compare with [14, Theorem 2.3]). Notice that $\mathbb{E}[\langle \cdot, \cdot \rangle_X]$, $\mathbb{E}[\langle \cdot, \cdot \rangle_Y]$ define inner products in $L_2(\Omega; X)$ and $L_2(\Omega; Y)$ respectively.⁵

Notice that, for any fixed $k \leq l$, one can always choose an index n with $k \leq n \leq l$ such that

$$\mathbb{E}[\lambda_I \|F_I^0(x_n)\|^2] \leq \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2], \quad \forall k \leq j \leq l. \quad (14)$$

Next, arguing with (7a) and the Cauchy–Schwartz inequality (for random variables) we estimate

$$\begin{aligned} |\mathbb{E}[\langle e_n - e_k, e_n \rangle]| &= \left| \sum_{j=k}^{n-1} \mathbb{E}[\langle x_{j+1} - x_j, x^* - x_n \rangle] \right| = \left| \sum_{j=k}^{n-1} \mathbb{E}[\theta_j \lambda_I \langle F_I'(x_j)^* F_I^0(x_j), x_n - x^* \rangle] \right| \\ &= \left| \sum_{j=k}^{n-1} \theta_j \mathbb{E}[\lambda_I \langle F_I^0(x_j), F_I'(x_j)(x_n - x_j + x_j - x^*) \rangle] \right| \\ &= \left| \sum_{j=k}^{n-1} \theta_j \mathbb{E}[\langle \lambda_I^{\frac{1}{2}} F_I^0(x_j), \lambda_I^{\frac{1}{2}} F_I'(x_j)(x_n - x_j) \rangle] + \theta_j \mathbb{E}[\langle \lambda_I^{\frac{1}{2}} F_I^0(x_j), \lambda_I^{\frac{1}{2}} F_I'(x_j)(x_j - x^*) \rangle] \right| \\ &\leq \sum_{j=k}^{n-1} \left(\theta_j \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]^{\frac{1}{2}} \mathbb{E}[\lambda_I \|F_I'(x_j)(x_n - x_j)\|^2]^{\frac{1}{2}} \right. \\ &\quad \left. + \theta_j \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]^{\frac{1}{2}} \mathbb{E}[\lambda_I \|F_I'(x_j)(x_j - x^*)\|^2]^{\frac{1}{2}} \right). \end{aligned}$$

Thus, it follows from (5) that

$$\begin{aligned} |\mathbb{E}[\langle e_n - e_k, e_n \rangle]| &\leq (1 + \eta) \sum_{j=k}^{n-1} \theta_j \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]^{\frac{1}{2}} \mathbb{E}[\lambda_I \|F_I^0(x_n) - F_I^0(x_j)\|^2]^{\frac{1}{2}} \\ &\quad + (1 + \eta) \sum_{j=k}^{n-1} \theta_j \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]. \end{aligned}$$

⁵ $L_2(\Omega; X)$ is the space of square integrable random variables defined on Ω and taking values in X .

The term $\mathbb{E}[\lambda_I \|F_I^0(x_n) - F_I^0(x_j)\|^2]^{\frac{1}{2}}$ on the right hand side of the last estimate can be estimated using (14). Indeed, for each $j = k, \dots, n-1$ it holds

$$\begin{aligned} \mathbb{E}[\lambda_I \|F_I^0(x_n) - F_I^0(x_j)\|^2]^{\frac{1}{2}} &\leq \left(2\mathbb{E}[\lambda_I \|F_I^0(x_n)\|^2] + 2\mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]\right)^{\frac{1}{2}} \\ &\leq \sqrt{2} \left(\mathbb{E}[\lambda_I \|F_I^0(x_n)\|^2] + \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]\right)^{\frac{1}{2}} \leq 2 \left[\mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2]\right]^{\frac{1}{2}}. \end{aligned}$$

Consequently,

$$|\mathbb{E}[\langle e_n - e_k, e_n \rangle]| \leq 3(1 + \eta) \sum_{j=k}^{n-1} \theta_j \mathbb{E}[\lambda_I \|F_I^0(x_j)\|^2].$$

Now, Corollary 3.4 allow us to conclude that $\mathbb{E}[\langle e_n - e_k, e_n \rangle] \rightarrow 0$ as $k, l \rightarrow \infty$. Analogously one proves that $\mathbb{E}[\langle e_l - e_n, e_n \rangle] \rightarrow 0$ as $k, l \rightarrow \infty$, establishing (13).

Finally, one argues with (13), (12), inequality $\mathbb{E}[\|e_l - e_k\|^2]^{\frac{1}{2}} \leq \mathbb{E}[\|e_l - e_n\|^2]^{\frac{1}{2}} + \mathbb{E}[\|e_n - e_k\|^2]^{\frac{1}{2}}$ and identities

$$\begin{aligned} \mathbb{E}[\|e_l - e_n\|^2] &= 2\mathbb{E}[\langle e_n - e_l, e_n \rangle] + \mathbb{E}[\|e_l\|^2] - \mathbb{E}[\|e_n\|^2], \\ \mathbb{E}[\|e_n - e_k\|^2] &= 2\mathbb{E}[\langle e_n - e_k, e_n \rangle] + \mathbb{E}[\|e_k\|^2] - \mathbb{E}[\|e_n\|^2], \end{aligned}$$

to conclude that $\mathbb{E}[\|e_l - e_k\|^2] \rightarrow 0$, as $k, l \rightarrow \infty$; i.e. (e_k) is a Cauchy sequence in $L_2(\Omega; X)$.

Since, (x_k) is Cauchy in $L_2(\Omega; X)$ it has an accumulation point $x^* \in L_2(\Omega; X)$. It remains to verify that this x^* is a solution of (8). It follows from Corollary 3.4 that the mean square residuals $\mathbb{E}[\|F_I^0(x_k)\|^2]$ converge to zero as $k \rightarrow \infty$. Arguing with the Assumptions (A1) and (A2), one concludes that $\mathbb{E}[\|F_I^0(x^*)\|^2] = 0$. Thus, $\sum_{i=0}^{N-1} \mathbb{E}[\|F_i(x^*) - y_i\|^2] = 0$ and, consequently, $\mathbb{E}[\|F_i(x^*) - y_i\|^2] = 0$ for $i = 0, \dots, N-1$. Since x^* is a random element, this implies that $\|F_i(x^*) - y_i\|^2 = 0$ for $i = 0, \dots, N-1$ almost surely, from where we conclude that x^* is a solution of (8) almost surely. \square

3.2 The noisy data case

In this section we investigate regularization properties of the nonlinear pSGD method in (6) coupled with the (*a priori*) stopping criterion in (A6).

In the next result, the nonlinear residual norm $\|F_i^\delta(x)\|$ is compared with the norm of the linearization $\|F_i^\delta(x) + F_i'(x)(x^* - x)\|$ for $x \in B_\rho(x_0)$ and $x^* \in B_\rho(x_0)$ a solution of (8).

Lemma 3.6. *Let Assumptions (A1), (A2) and (A3) hold. For all $x, \bar{x} \in B_\rho(x_0)$ we have*

$$\| -F_i^\delta(x) - F_i'(x)(\bar{x} - x) \| \leq \eta \|F_i^\delta(x)\| + (1 + \eta) \|F_i^\delta(\bar{x})\|, \quad i = 0, \dots, N-1.$$

In particular, if $\bar{x} = x^ \in B_\rho(x_0)$ is a solution of (2) it holds*

$$\| -F_i^\delta(x) - F_i'(x)(x^* - x) \| \leq \eta \|F_i^\delta(x)\| + (1 + \eta) \delta_i, \quad \forall x \in B_\rho(x_0).$$

Proof. Given $x, \bar{x} \in B_\rho(x_0)$ we conclude from (A2) that

$$\begin{aligned} \| -F_i^\delta(x) - F_i'(x)(\bar{x} - x) \| &= \| y_i^\delta - F_i(x) - F_i'(x)(\bar{x} - x) \| \\ &= \| y_i^\delta - F_i(\bar{x}) + F_i(\bar{x}) - F_i(x) - F_i'(x)(\bar{x} - x) \| \\ &\leq \eta \|F_i(\bar{x}) - F_i(x)\| + \|y_i^\delta - F_i(\bar{x})\| \\ &\leq \eta \|y_i^\delta - F_i(x)\| + (1 + \eta) \|y_i^\delta - F_i(\bar{x})\|, \end{aligned}$$

proving the first assertion. The second assertion is a direct consequence of the first one. \square

In the sequel we estimate the difference $\mathbb{E}[\|x^* - x_{k+1}^\delta\|^2] - \mathbb{E}[\|x^* - x_k^\delta\|^2]$, extending the estimate (9) in Proposition 3.2 to the noisy data case.

Proposition 3.7. *Let Assumptions (A1), (A2) and (A3) hold, (x_k^δ) be a sequence generated by the nonlinear pSGD method (6), and $x^* \in B_\rho(x_0)$ be a solution of (8). If $x_k^\delta \in B_\rho(x_0)$ for some $0 \leq k \leq k_\delta^*$, then*

$$\mathbb{E}[\|x^* - x_{k+1}^\delta\|^2] - \mathbb{E}[\|x^* - x_k^\delta\|^2] \leq \theta_k(\theta_k - 2) \left[(1 - \eta) \mathbb{E}[\lambda_I^\delta \|F_I^\delta(x_k^\delta)\|^2] - (1 + \eta) \mathbb{E}[\delta_I \lambda_I^\delta \|F_I^\delta(x_k^\delta)\|] \right]. \quad (15)$$

Proof. If $\|F_{I_k}'(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta)\| > \gamma \delta_{I_k}$, it follows from (6a) and (6b)

$$\begin{aligned} \|x^* - x_{k+1}^\delta\|^2 - \|x^* - x_k^\delta\|^2 &= 2 \langle x_{k+1}^\delta - x_k^\delta, x_{k+1}^\delta - x^* \rangle - \|x_{k+1}^\delta - x_k^\delta\|^2 \\ &= 2\theta_k \lambda_{I_k}^\delta \left\langle -F_{I_k}^\delta(x_k^\delta), -F_{I_k}^\delta(x_k^\delta) - F_{I_k}'(x_k^\delta)(x^* - x_k^\delta) + F_{I_k}'(x_k^\delta)(x_{k+1}^\delta - x_k^\delta) + F_{I_k}^\delta(x_k^\delta) \right\rangle \\ &\quad - \|x_{k+1}^\delta - x_k^\delta\|^2 \\ &= 2\theta_k \lambda_{I_k}^\delta \left[\langle -F_{I_k}^\delta(x_k^\delta), -F_{I_k}^\delta(x_k^\delta) - F_{I_k}'(x_k^\delta)(x^* - x_k^\delta) \rangle + \langle -F_{I_k}'(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta), x_{k+1}^\delta - x_k^\delta \rangle \right. \\ &\quad \left. - \|F_{I_k}^\delta(x_k^\delta)\|^2 \right] - \|x_{k+1}^\delta - x_k^\delta\|^2 \\ &= 2\theta_k \lambda_{I_k}^\delta \left[\langle -F_{I_k}^\delta(x_k^\delta), -F_{I_k}^\delta(x_k^\delta) - F_{I_k}'(x_k^\delta)(x^* - x_k^\delta) \rangle + \theta_k p^{\delta_{I_k}} (\|F_{I_k}^\delta(x_k^\delta)\|) - \|F_{I_k}^\delta(x_k^\delta)\|^2 \right] \\ &\quad - \theta_k^2 \lambda_{I_k}^\delta p^{\delta_{I_k}} (\|F_{I_k}^\delta(x_k^\delta)\|). \end{aligned}$$

Thus, arguing with the Cauchy-Schwartz inequality, Lemma 3.6 and the definition of $p^\delta(\cdot)$ follows

$$\begin{aligned} \|x^* - x_{k+1}^\delta\|^2 - \|x^* - x_k^\delta\|^2 &\leq 2\theta_k \lambda_{I_k}^\delta \left[\|F_{I_k}^\delta(x_k^\delta)\| \| -F_{I_k}^\delta(x_k^\delta) - F_{I_k}'(x_k^\delta)(x^* - x_k^\delta) \| + \frac{1}{2} \theta_k p^{\delta_{I_k}} (\|F_{I_k}^\delta(x_k^\delta)\|) - \|F_{I_k}^\delta(x_k^\delta)\|^2 \right] \\ &\leq 2\theta_k \lambda_{I_k}^\delta \left[\eta \|F_{I_k}^\delta(x_k^\delta)\|^2 + (1 + \eta) \delta_{I_k} \|F_{I_k}^\delta(x_k^\delta)\| + \frac{1}{2} \theta_k p^{\delta_{I_k}} (\|F_{I_k}^\delta(x_k^\delta)\|) - \|F_{I_k}^\delta(x_k^\delta)\|^2 \right] \\ &= 2\theta_k \lambda_{I_k}^\delta \left[(\eta - 1) \|F_{I_k}^\delta(x_k^\delta)\|^2 + (1 + \eta) \delta_{I_k} \|F_{I_k}^\delta(x_k^\delta)\| + \frac{1}{2} \theta_k ((1 - \eta) \|F_{I_k}^\delta(x_k^\delta)\|^2 \right. \\ &\quad \left. - (1 + \eta) \delta_{I_k} \|F_{I_k}^\delta(x_k^\delta)\|) \right] \\ &= 2\theta_k \lambda_{I_k}^\delta \left[(1 - \eta) \left(\frac{\theta}{2} - 1 \right) \|F_{I_k}^\delta(x_k^\delta)\|^2 + (1 + \eta) \left(1 - \frac{\theta}{2} \right) \delta_{I_k} \|F_{I_k}^\delta(x_k^\delta)\| \right] \\ &= \theta_k (\theta_k - 2) \lambda_{I_k}^\delta \left[(1 - \eta) \|F_{I_k}^\delta(x_k^\delta)\|^2 - (1 + \eta) \delta_{I_k} \|F_{I_k}^\delta(x_k^\delta)\| \right]. \quad (16) \end{aligned}$$

Otherwise, if $\|F_{I_k}'(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta)\| \leq \gamma \delta_{I_k}$, then $\lambda_{I_k}^\delta = 0$ and $x_{k+1}^\delta = x_k^\delta$ (see (6a) and (6b)). Therefore, (16) holds also in this case. By employing a similar reasoning as in the final part of the proof of Proposition 3.2, we establish the validity of (15). \square

It is worth noticing that the right hand side of (16) can be rewritten in the form

$$\begin{aligned} \theta_k (\theta_k - 2) \lambda_{I_k}^\delta \left[(1 - \eta) \|F_{I_k}^\delta(x_k^\delta)\|^2 - (1 + \eta) \delta_{I_k} \|F_{I_k}^\delta(x_k^\delta)\| \right] &= \theta_k (\theta_k - 2) \lambda_{I_k}^\delta p^{\delta_{I_k}} (\|F_{I_k}^\delta(x_k^\delta)\|) \\ &= \theta_k (\theta_k - 2) (\lambda_{I_k}^\delta)^2 \|F_{I_k}'(x_k^\delta)^* F_{I_k}^\delta(x_k^\delta)\|^2 \geq 0 \end{aligned}$$

(see Assumption (A4)), from where we conclude that $\|x^* - x_{k+1}^\delta\| \leq \|x^* - x_k^\delta\|$. This inequality and (A3) allow us to conclude that (x_k^δ) satisfies $\|x^* - x_{k+1}^\delta\| \leq \|x^* - x_k^\delta\|$ and $x_k^\delta \in B_\rho(x_0)$, for $k = 0, \dots, k_\delta^*$. Summarizing, we have

Corollary 3.8. *Under the assumptions of Proposition 3.7 it holds $\|x^* - x_{k+1}^\delta\| \leq \|x^* - x_k^\delta\|$, for $k = 0, \dots, k_\delta^*$. Consequently, $(x_k^\delta) \subset B_\rho(x_0)$. Additionally, for any x^* solution of (8) in $B_\rho(x_0)$ it holds*

$$\mathbb{E}[\|x^* - x_{k+1}^\delta\|^2] \leq \mathbb{E}[\|x^* - x_k^\delta\|^2], \quad k = 0, \dots, k_\delta^*.$$

We are now ready to state and prove a stability result (Theorem 3.9) and a semi-convergence result (Theorem 3.11) for the pSGD method in (6).

Theorem 3.9 (Stability). *Let Assumptions (A1), (A2) and (A4) hold, $(\delta^j) = (\delta_0^j, \dots, \delta_{N-1}^j) \in (\mathbb{R}^+)^N$ be a sequence with $\|\delta^j\| \rightarrow 0$ as $j \rightarrow \infty$, and $(y^{\delta^j}) = (y_0^{\delta^j}, \dots, y_{N-1}^{\delta^j}) \in Y^N$ be a corresponding sequence of noisy data satisfying (1). Moreover, let $(x_l)_{l \in \mathbb{N}}$ and $(x_l^{\delta^j})_{l=0}^{k_*}$ be sequences generated by the nonlinear pSGD method in the case of exact and noisy data respectively; all sequences are generated using the same (I_0, \dots, I_k, \dots) . For each $k \in \mathbb{N}$ it holds*

$$\lim_{j \rightarrow \infty} \mathbb{E}[\|x_k^{\delta^j} - x_k\|^2] = 0. \quad (17)$$

Prior to establishing this theorem, we examine an auxiliary result:

Lemma 3.10. *Under assumptions of Theorem 3.9, if $\|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| > 0$ and $\lim_j x_k^{\delta^j} = x_k$ for some $k \in \mathbb{N}$, then $\lim_j |\lambda_{I_k}^{\delta^j} - \lambda_{I_k}| = 0$.*

Proof. From (A1), (1), $\lim_j \|\delta^j\| = 0$ and $\lim_j x_k^{\delta^j} = x_k$ follow

$$\lim_{j \rightarrow \infty} p^{\delta_{I_k}}(\|F_{I_k}^{\delta^j}(x_k^{\delta^j})\|) = p^0(\|F_{I_k}^0(x_k)\|). \quad (18)$$

On the other hand, from (A1), (1), $\lim_j \|\delta^j\| = 0$ and $\lim_j x_k^{\delta^j} = x_k$ we conclude that

$$\lim_j \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| = \|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\|.$$

Therefore, the hypothesis $\|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| > 0$ (together with the fact $\lim_j \delta_{I_k}^j = 0$) allow us to conclude that $\|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| > \frac{1}{2} \|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| > \gamma \delta_{I_k}^j$ for sufficiently large j . Consequently, the lemma follows from (18) together with the definitions of $\lambda_{I_k}^{\delta^j}$ and λ_{I_k} in (6b) and (7b) respectively. \square

Proof. (of Theorem 3.9)

We give an inductive proof in k . Since $x_0 = x_0^{\delta^j}$ for all $j \in \mathbb{N}$, (17) trivially holds for $k = 0$. Assume that $\lim_j \mathbb{E}[\|x_l^{\delta^j} - x_l\|^2] = 0$ for $l = 0, \dots, k$. We aim to prove $\lim_j \mathbb{E}[\|x_{k+1}^{\delta^j} - x_{k+1}\|^2] = 0$. The argumentation is divided in three steps:

Step 1: We claim that $\lim_j \|x_l^{\delta^j} - x_l\|^2 = 0$ for $l = 1, \dots, k$.

Arguing as in the proof of [29, Theorem 4.3] one proves that, for each realization (I_0, \dots, I_{l-1}) the inequality $\|x_l^{\delta^j} - x_l\|^2 \leq (\frac{1}{N})^l E[\|x_l^{\delta^j} - x_l\|^2]$ holds for $l = 1, \dots, k$ and $j \in \mathbb{N}$. Thus, for each fixed $l \in \{1, \dots, k\}$ the inductive hypothesis guarantees that $\lim_j \|x_l^{\delta^j} - x_l\|^2 = 0$.

Step 2: We claim that $\lim_j \|x_{k+1}^{\delta^j} - x_{k+1}\|^2 = 0$. Two distinct cases are considered:

(I) $\|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| > 0$. From the iteration formulas (6a) and (7a) follow

$$\begin{aligned} x_{k+1}^{\delta^j} - x_{k+1} &= \\ &= x_k^{\delta^j} - x_k - \theta_k [\lambda_{I_k}^{\delta^j} F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j}) - \lambda_{I_k} F'_{I_k}(x_k)^* F_{I_k}^0(x_k)] \\ &= x_k^{\delta^j} - x_k - \theta_k [(\lambda_{I_k}^{\delta^j} - \lambda_{I_k}) F'_{I_k}(x_k)^* F_{I_k}^0(x_k) + \lambda_{I_k}^{\delta^j} (F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j}) - F'_{I_k}(x_k)^* F_{I_k}^0(x_k))]. \end{aligned}$$

Therefore, arguing with (A1) we estimate

$$\begin{aligned} \|x_{k+1}^{\delta^j} - x_{k+1}\| &\leq \\ &\|x_k^{\delta^j} - x_k\| + 2C |\lambda_{I_k}^{\delta^j} - \lambda_{I_k}| \|F_{I_k}^0(x_k)\| + 2|\lambda_{I_k}^{\delta^j}| \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j}) - F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\|. \end{aligned}$$

Taking the limit $j \rightarrow \infty$ in the above inequality and arguing with Step 1 (for $l = k$), Lemma 3.10 and (A1), we conclude that the three terms on the right hand side do converge to zero;⁶ proving our claim in Step 2, in case (I).

(II) $\|F'_{I_k}(x_k)^* F_{I_k}^0(x_k)\| = 0$. In this case $\lambda_{I_k} = 0$ and $F_{I_k}^0(x_k) = 0$. From (6a) and (7a) follow

$$\|x_{k+1}^{\delta^j} - x_{k+1}\| \leq \|x_k^{\delta^j} - x_k\| + 2\lambda_{I_k}^{\delta^j} \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\|. \quad (19)$$

Given $j \in \mathbb{N}$, if $\|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \leq \gamma \delta_{I_k}^j$ then $\lambda_{I_k}^{\delta^j} = 0$ and we derive from (19)

$$\|x_{k+1}^{\delta^j} - x_{k+1}\| \leq \|x_k^{\delta^j} - x_k\|. \quad (\dagger)$$

Otherwise, if $\|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| > \gamma \delta_{I_k}^j$ we argue with Lemma 3.6 to estimate

$$\begin{aligned} \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\|^2 &= \langle F_{I_k}(x_k^{\delta^j}) - y_{I_k}^{\delta^j}, -y_{I_k}^{\delta^j} + F_{I_k}(x_k^{\delta^j}) - F'_{I_k}(x_k^{\delta^j})(x_k - x_k^{\delta^j}) + F'_{I_k}(x_k^{\delta^j})(x_k - x_k^{\delta^j}) \rangle \\ &\leq \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \| -F_{I_k}^{\delta^j}(x_k^{\delta^j}) - F'_{I_k}(x_k^{\delta^j})(x_k - x_k^{\delta^j}) \| + \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \|x_k - x_k^{\delta^j}\| \\ &\leq \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \left[\eta \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| + (1 + \eta) \|F_{I_k}^{\delta^j}(x_k)\| \right] + \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \|x_k - x_k^{\delta^j}\| \\ &= \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \left[\eta \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| + (1 + \eta) \|y_{I_k}^{\delta^j} \pm y_{I_k} - F_{I_k}(x_k)\| \right] + \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \|x_k - x_k^{\delta^j}\| \\ &\leq \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \left[\eta \|F_{I_k}^{\delta^j}(x_k^{\delta^j})\| + (1 + \eta) (\delta_{I_k}^j + \|F_{I_k}^0(x_k)\|) \right] + \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \|x_k - x_k^{\delta^j}\|. \end{aligned}$$

Since $F_{I_k}^0(x_k) = 0$ holds in case (II), the last inequality and the definition of $p^\delta(\cdot)$ allow us to conclude that $p^{\delta_{I_k}^j}(\|F_{I_k}^{\delta^j}(x_k^{\delta^j})\|) \leq \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| \|x_k - x_k^{\delta^j}\|$. Thus, it follows from (6b)

$$\lambda_{I_k}^{\delta^j} \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\| = p^{\delta_{I_k}^j}(\|F_{I_k}^{\delta^j}(x_k^{\delta^j})\|) \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^{\delta^j}(x_k^{\delta^j})\|^{-1} \leq \|x_k^{\delta^j} - x_k\|.$$

From this inequality and (19) follow

$$\|x_{k+1}^{\delta^j} - x_{k+1}\| \leq 3\|x_k^{\delta^j} - x_k\|. \quad (\ddagger)$$

From inequalities (\dagger) and (\ddagger) together with Step 1 (for $l = k$) follow $\lim_j \|x_{k+1}^{\delta^j} - x_{k+1}\| \leq 3 \lim_j \|x_k^{\delta^j} - x_k\| = 0$, proving our claim in Step 2, in case (II).

Step 3: We claim that $\lim_j \mathbb{E}[\|x_{k+1}^{\delta^j} - x_{k+1}\|^2] = 0$, concluding the inductive proof.

Indeed, notice that

$$\mathbb{E}[\|x_{k+1}^{\delta} - x_{k+1}\|^2] = \left(\frac{1}{N}\right)^{k+1} \sum_{\substack{i_0=0 \\ \vdots \\ i_k=0}}^N \|x_{k;i_0,\dots,i_k}^{\delta} - x_{k;i_0,\dots,i_k}\|^2,$$

where $x_{k;i_0,\dots,i_k}^{\delta}$ is defined by (6a) taking $(I_0, \dots, I_{k-1}) = (i_0, \dots, i_{k-1})$, and $x_{k;i_0,\dots,i_{k-1}}$ is defined by (7a) analogously. Taking the average in Step 2 over all possible realizations (I_0, \dots, I_k) , one concludes that $\lim_j \mathbb{E}[\|x_{k+1}^{\delta^j} - x_{k+1}\|^2] = 0$. \square

Theorem 3.11 (Semi-convergence). *Let Assumptions (A1), ..., (A6) hold true; $(\delta^j) = (\delta_0^j, \dots, \delta_{N-1}^j) \in \mathbb{R}^N$ be a sequence with $\lim_j \|\delta_j\| = 0$ and $(\mathbf{y}^{\delta^j}) = (y_0^{\delta^j}, \dots, y_{N-1}^{\delta^j}) \in Y^N$ be a corresponding sequence of noisy data satisfying (1). Additionally, for each $j \in \mathbb{N}$, let $(x_k^{\delta^j})_{k=0}^{k^*(\delta^j)}$ be the corresponding sequence generated by the nonlinear pSGD method (6). There exists a random element $x^* \in B_\rho(x_0)$, which is a solution of (8) almost surely, such that*

$$\lim_{j \rightarrow \infty} \mathbb{E}[\|x_{k^*(\delta^j)}^{\delta^j} - x^*\|^2] = 0. \quad (20)$$

⁶Notice that Lemma 3.10 guarantees the boundedness of the sequence $(\lambda_{I_k}^{\delta^j})_j$.

Proof. Let $\hat{x} \in B_\rho(x_0)$ be a solution of (8) (the existence of \hat{x} follows from (A3)). We claim that, for every fixed $k \in \mathbb{N}$ it holds

$$\|\hat{x} - x_k^{\delta^j}\|^2 - \|\hat{x} - x_{k+1}^{\delta^j}\|^2 \geq 0, \quad (21)$$

for sufficiently large j .

From (A6) we conclude that $k_{\delta_j}^* \geq k$ for sufficiently large j . Consequently, $x_k^{\delta^j}$ and $x_{k+1}^{\delta^j}$ are defined for sufficiently large j .

If $\|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^\delta(x_k^{\delta^j})\| \leq \gamma \delta_{I_k}^j$, then $\lambda_{I_k}^{\delta^j} = 0$, $x_{k+1}^{\delta^j} = x_k^{\delta^j}$ and (21) holds trivially. Otherwise, it follows from Remark 2.3, (16) and (6b)

$$\begin{aligned} \|\hat{x} - x_k^{\delta^j}\|^2 - \|\hat{x} - x_{k+1}^{\delta^j}\|^2 &\geq \theta_k(2 - \theta_k) \lambda_{I_k} \mathcal{P}^{\delta^j} (\|F_{I_k}^\delta(x_k^{\delta^j})\|) \\ &= \theta_k(2 - \theta_k) (\lambda_{I_k})^2 \|F'_{I_k}(x_k^{\delta^j})^* F_{I_k}^\delta(x_k^{\delta^j})\|^2 \\ &> \theta_k(2 - \theta_k) \lambda_{\min}^2 (\gamma \delta_{I_k}^j)^2 \\ &\geq \theta_k(2 - \theta_k) \lambda_{\min}^2 \gamma^2 (\delta_{\min}^j)^2 \end{aligned}$$

(notice that Assumption (A4) guarantees $\theta_k(2 - \theta_k) > 0$). Thus, $\|\hat{x} - x_k^{\delta^j}\|^2 - \|\hat{x} - x_{k+1}^{\delta^j}\|^2 \geq 0$, establishing our claim (21).

Now, taking the average in (21) over all possible realizations (I_0, \dots, I_k) and arguing with Remark 3.1 we obtain

$$\mathbb{E}[\|\hat{x} - x_k^{\delta^j}\|^2] - \mathbb{E}[\|\hat{x} - x_{k+1}^{\delta^j}\|^2] \geq 0.$$

From (A6) we may assume that $k_{\delta_j}^*$ increases strict monotonically with j . Given $m < n$, we derive from the last inequality, with $j = n$ and $k = k_{\delta_m}^*, \dots, k_{\delta_n}^* - 1$, the estimate

$$\mathbb{E}[\|\hat{x} - x_{k_n}^{\delta^n}\|^2] \leq \mathbb{E}[\|\hat{x} - x_{k_m}^{\delta^n}\|^2] \leq 2\mathbb{E}[\|\hat{x} - x_{k_m}^*\|^2] + 2\mathbb{E}[\|x_{k_m}^* - x_{k_m}^{\delta^n}\|^2] \quad (22)$$

(we adopted the simplified notation $k_j^* = k_{\delta_j}^*$). Here (x_k) is the sequence generated by pSGD (7) using exact data and the same (I_0, \dots, I_k, \dots) as the sequences $(x_k^{\delta^j})$.

Let x^* be a solution of (8) satisfying $\lim_k \mathbb{E}[\|x^* - x_k\|^2] = 0$ (the existence of x^* is ensured by Theorem 3.5). From (21) follows $\|x^* - x_k^{\delta^j}\|^2 - \|x^* - x_{k+1}^{\delta^j}\|^2 \geq 0$ almost surely. Thus, by taking the full expectation gives $\mathbb{E}[\|x^* - x_k^{\delta^j}\|^2] - \mathbb{E}[\|x^* - x_{k+1}^{\delta^j}\|^2] \geq 0$. Arguing as above we conclude that (22) holds with \hat{x} replaced by x^* . Therefore, there is a large enough m , s.t. the first term on the rhs of (22) is smaller than $\varepsilon/2$. Additionally, from Theorem 3.9 with $k = k_m^*$ we conclude that the second term on the rhs of (22) becomes smaller than $\varepsilon/2$ for large enough n . This concludes the proof. \square

4 Numerical experiments

In Sections 4.1 and 4.2 the pSGD method from (6) is applied to solve two large-scale systems of nonlinear operator equations using real-world data. The corresponding inverse problems relate to parameter identification in neural network training. In Section 4.3 a parameter identification problem in a 3D elliptic PDE system is considered. Computations are performed using MATLAB[®] R2017a, running in a Intel[®] Core[™] i9-10900 CPU (10 cores, 20 threads).

In Section 4.1 we revisit, using neural networks, an inverse problem discussed in [29], namely the big data problem of CO-concentration prediction in a gas sensor array.⁷ In Section 4.2 we consider a well-known benchmark problem in machine learning, namely the classification problem for the MNIST database (see https://en.wikipedia.org/wiki/MNIST_database).

⁷Differently from the approach used here, in [29, Section 5.2] this inverse problem was modeled as a multiple linear regression problem.

4.1 Prediction of CO-concentration in a gas sensor array

For the application considered in this section, we use a data set collected in a gas delivery platform facility at the ChemoSignals Laboratory in the BioCircuits Institute at University of California, San Diego. The data set contains the readings of 16 distinct chemical sensors, which were exposed to the mixture of Ethylene and CO at varying concentrations in air. For this gas mixture, the measurement was constructed by the continuous acquisition of the 16-sensor array signals for a duration of approximately 12 hours without interruption (we refer the reader to [9, 29] for a detailed description of the experiment). The real data used in this section is available at the UC Irvine Machine Learning Repository <https://archive.ics.uci.edu/ml/index.php>, dataset *Gas sensor array under dynamic gas mixtures*.

In [29] the following experimental setting was considered: readings from the last sensor (sensor #16) are used as the response variable, and readings from sensors $\{\#1, \#3, \#4, \dots, \#15\}$ are used as covariates (readings from sensor #2 are disregarded due to strong lack of accuracy); each sensor data consists of 4,188,262 scalar measurements.⁸ The inverse problem considered in [29] is a *multiple linear regression* problem. It consists of finding an approximate solution to the linear system $Ax = y^\delta$ (with unknown noise level $\delta > 0$), where $A = (A_i)_{i=0}^{N-1} \in \mathbb{R}^{N \times M}$, with $N = 4\,188\,262$ and $M = 15$. Here $x \in \mathbb{R}^M$ is the unknown vector of regression coefficients, $y^\delta \in \mathbb{R}^N$ contains the readings from sensor #16, and A_i (the i^{th} -row of A) is such that: the first 14 columns of A_i contain the i^{th} -readings from sensors $\{\#1, \#3, \#4, \dots, \#15\}$ and the last column is equal to one.

An inverse problem in machine learning

In this section we consider the problem of predicting the reading from sensor #16 based on the readings from the previous 14 sensors. However, differently from the multiple linear regression approach in [29], we use here a neural network (NN) that inputs the readings from the first sensors and outputs a scalar value, which predicts the reading of the last sensor.

The structure of proposed NN reads:

- **Input:** $z \in \mathbb{R}^{14}$, readings of the first 14 sensors;
- **Output:** $NN(z; W, b) = \sigma(Wz + (1 - \varepsilon)b) \in \mathbb{R}$, where $W \in \mathbb{R}^{1 \times 14}$ is a matrix of weights, $\varepsilon > 0$ is a small constant, $b \in \mathbb{R}$ is a scalar bias and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an activation function.

Notice that this is a very simple NN with only one layer (the output layer); the dimension of the corresponding parameter space is 15, the size of (W, b) . If the activation function σ is linear, this approach becomes equivalent to the multiple linear regression used in [29].⁹

The inverse problem under consideration is a NN training problem, i.e. one aims to find an approximate solution (a pair of parameters (W, b)) to the nonlinear system

$$F_i(W, b) = y_i^\delta, \quad i = 0, \dots, N_t - 1, \quad (23)$$

where $F_i(W, b) := NN(z_i; W, b) = \sigma(Wz_i + (1 - \varepsilon)b)$. Here $N_t < N$ is the size of the training set and $z_i = ((A_i)_j)_{j=1}^{14}$, where A_i is the i^{th} -row of A .

Once the parameters (W, b) are chosen, the performance \mathcal{P} of the corresponding neural network $NN(\cdot; W, b)$ is defined by

$$\mathcal{P}(NN(\cdot; W, b)) := 1 - \frac{1}{N_T} \sum_{i=N_t}^{N_t+N_T-1} \frac{\|NN(z_i; W, b) - y_i^\delta\|}{\|y_i^\delta\|}. \quad (24)$$

The sum in the above definition gives the average (relative) misfit between the predicted value $NN(z_i; W, b)$ and y_i^δ , evaluated over the test set $\{z_i, N_t \leq i < N_t + N_T - 1\}$. Clearly it holds $0 \leq \mathcal{P}(NN(\cdot; W, b)) \leq 1$ for all (W, b) , and $\mathcal{P}(NN(\cdot; W, b)) = 1$ is the best possible performance.

⁸See [29, Figure 3] for scatter plots of sensor #i readings against sensor #16, for $i = 1, \dots, 15$.

⁹We choose the nonlinear activation function σ s.t. its range contains all possible readings of sensor #16.

Remark 4.1 (Choosing the training set and test set). *At the beginning of the experiment, the rows A_i are arranged in a random order. Consequently, the 'training set' and 'test set' are comprised of random samples with sizes of N_t and N_T respectively. In our numerical experiments we use $N_t = 4,000,000$ and $N_T = 100,000$ (notice that $N_t + N_T < N$ is satisfied).*

The choice of the activation function

The activation function σ used in the definition of our NN is a variation of the *saturated linear activation function* [6]. Here $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is defined by

$$\sigma(t) = \begin{cases} 2\sqrt{t+1} - 2 & , t \geq 0 \\ 2 - 2\sqrt{1-t} & , t < 0 \end{cases} . \quad (25)$$

The choice of this particular activation function is motivated by the next lemma (a proof is postponed to Appendix A).

Lemma 4.2. *For each $a > 1$ the real function σ in (25) satisfies wTCC (4) in the interval $(1 - a^2, a^2 - 1)$ for $\eta = \frac{1}{2}(a - 1)$, i.e.*

$$\|\sigma(\bar{t}) - \sigma(t) - \sigma'(t)(\bar{t} - t)\| \leq \frac{a-1}{2} \|\sigma(\bar{t}) - \sigma(t)\|, \quad \forall t, \bar{t} \in (1 - a^2, a^2 - 1).$$

A direct consequence of Lemma 4.2 (with $a = 3$) is that the function σ in (25) satisfies the wTCC (4) in the interval $(-8, 8)$ with $\eta = 1$. In the sequel we prove that, for every fixed z_i , the above defined neural network $NN(z_i; \cdot, \cdot)$ does satisfy wTCC (4).

Lemma 4.3. *The function $F_i : (W, b) \mapsto NN(z_i; W, b) = \sigma(Wz_i + (1 - \varepsilon)b)$, with σ as in (25), satisfies the wTCC (4) for $\eta = \max\{\|z_i\|, |1 - \varepsilon|\}$ in a suitable neighborhood V_i of $(0, 0) \in \mathbb{R}^{14} \times \mathbb{R}$.*

Proof. If f and g are functions (with $D(f) \supset Rg(g)$) satisfying wTCC (4) for constants η_f and η_g respectively, then

$$\begin{aligned} \|f(g(\bar{t})) - f(g(t)) - f'(g(t))g'(t)(\bar{t} - t)\| &\leq \\ &\leq \|f(g(\bar{t})) - f(g(t)) - f'(g(t))[g(\bar{t}) - g(t)] + f'(g(t))[g(\bar{t}) - g(t) - g'(t)(\bar{t} - t)]\| \\ &\leq \eta_f \|g(\bar{t}) - g(t)\| + \eta_g \|f'(g(t))\| \|\bar{t} - t\|. \end{aligned} \quad (26)$$

Notice that $F_i = f \circ g_i$ with $f(t) = \sigma(t)$ and $g_i(W, b) = Wz_i + (1 - \varepsilon)b$. Since g_i is linear, it satisfies wTCC (4) for $\eta_g = 0$. Consequently, it follows from (26) and Lemma 4.2

$$\begin{aligned} \|F_i(\bar{W}, \bar{b}) - F_i(W, b) - F_i'(W, b)[(\bar{W}, \bar{b}) - (W, b)]\| &\leq \|(\bar{W}z_i - (1 - \varepsilon)\bar{b}) - (Wz_i - (1 - \varepsilon)b)\| \\ &= \|(\bar{W} - W)z_i - (1 - \varepsilon)(\bar{b} - b)\| \leq [\|z_i\| \|\bar{W} - W\| + |1 - \varepsilon| \|\bar{b} - b\|], \end{aligned}$$

for $(\bar{W}, \bar{b}), (W, b) \in V_i := \{(W, b) \in \mathbb{R}^{14} \times \mathbb{R}; \|Wz_i + (1 - \varepsilon)b\| < 8\}$. The assertion follows from the last inequality. \square

Numerical implementations

In what follows the pSGD method in (6) is implemented for solving problem (23). The sensor readings $(z_i, y_i^\delta) \in \mathbb{R}^{14} \times \mathbb{R}$ on the training set are scaled by the factor $(1 - \varepsilon)^{-1} \max_{i \leq N_t} \|z_i\|$. An analogous procedure is performed on the test set. Consequently, after scaling, it holds $\|z_i\| \leq (1 - \varepsilon)$, for $i = 0, \dots, N_t + N_T$. From Lemma 4.3 it follows that each operator $F_i(W, b)$ satisfy Assumption (A2) for the same constant $\eta = 1 - \varepsilon$ at the corresponding neighborhood V_i .

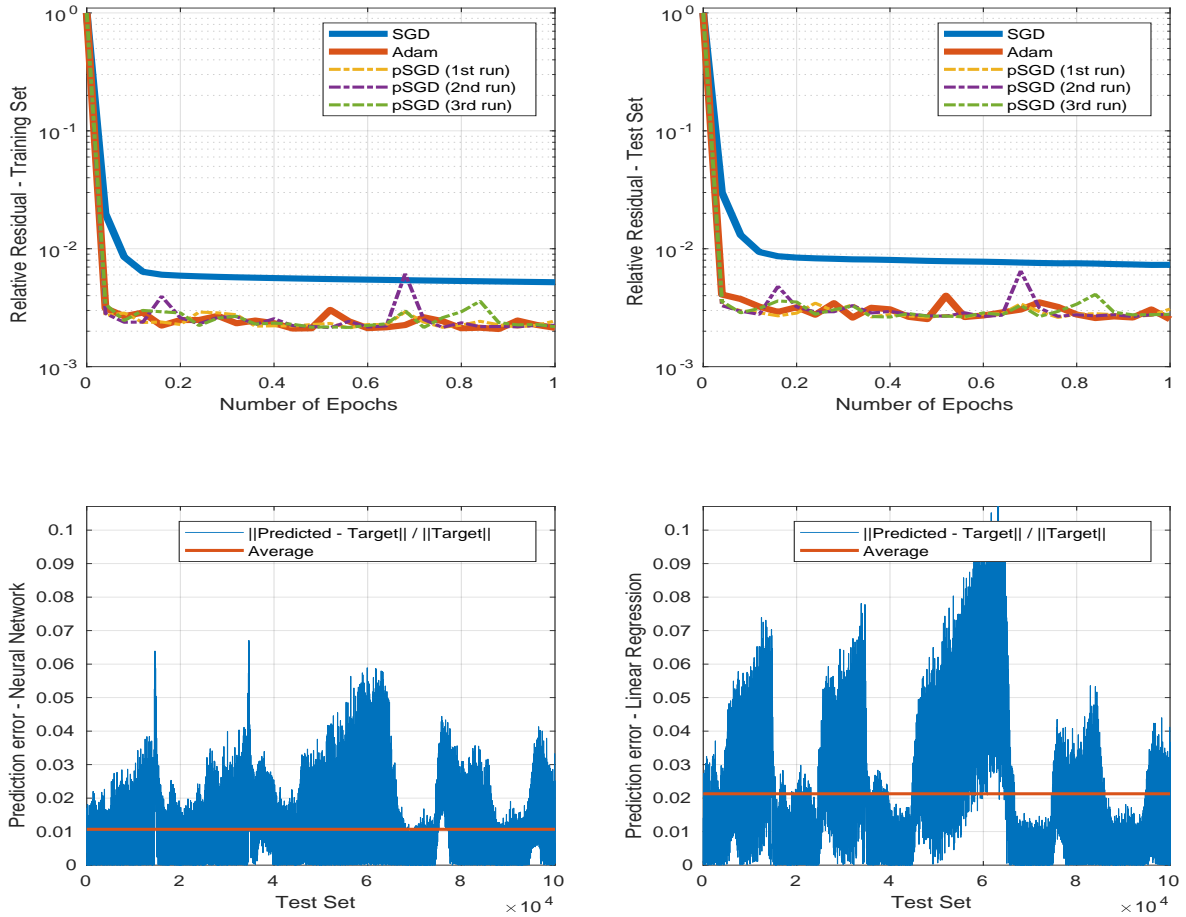


Figure 1: Prediction of CO-concentration in a gas sensor array. (TOP) Evolution of the relative residual: Training set (Left) and Test set (Right); (BOTTOM) Accuracy of the prediction: Neural network (Left) and Linear regression (Right).

In our experiments the initial guess (W_0, b_0) is a random vector with coordinate values ranging in $(-1, 1)$ and $\varepsilon = 0.1$. Moreover, we choose the sequence $\theta_k \equiv 1$ in (A4). The iteration (W_k, b_k) is computed for one epoch, i.e. for $k = 1, \dots, N_t$.

Three different runs of the pSGD method were computed. In each one of them we evaluate $W_k z_{I_k} + (1 - \varepsilon)b_k$ and observe that $\|W_k z_{I_k} + (1 - \varepsilon)b_k\| < 2$ for $k = 1, \dots, N_t$. From Lemma 4.3 it follows that $(W_k, b_k) \in V_{I_k}$ for $k = 1, \dots, N_t$.

Since the noise level δ is not known, we set $p^\delta(t) = (1 - \eta)t^2$ in (6b). The computed numerical results are summarized in Figure 1:

(TOP-LEFT) Evolution of relative residual on the training set: $\sum_{i=0}^{N_t-1} \frac{\|NN(z_i; W_k, b_k) - y_i^\delta\|}{\|NN(z_i; W_0, b_0) - y_i^\delta\|}$,

(TOP-RIGHT) Relative residual evolution on the test set: $\sum_{i=N_t}^{N_t+N_T-1} \frac{\|NN(z_i; W_k, b_k) - y_i^\delta\|}{\|NN(z_i; W_0, b_0) - y_i^\delta\|}$,

For comparison purposes, two relevant concurrent methods were implemented: (i) the stochastic gradient descent (SGD) method, which corresponds to the choice $\theta_k = 1$ and $\lambda_k = C^{-2}$ in (6a);¹⁰ (ii) the advanced stochastic algorithm (ADAM) [21], which is a well established extension of SGD. For the implementation of ADAM we use the choice of parameters described in [21], namely $\alpha = 0.001$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

At regular intervals of $\frac{1}{100}N_t$ steps, $\mathcal{P}(NN(\cdot; W_k, b_k))$ is computed. The index $0 \leq k^* \leq N_t$ is chosen such that (W_{k^*}, b_{k^*}) exhibits the highest performance among the evaluated ones. The additional task of evaluating $\mathcal{P}(NN(\cdot; W_k, b_k))$ a hundred times (per epoch) leads to a 10% increase in the overall computation time of the pSGD method.¹¹

The prediction accuracy of $NN(\cdot; W_{k^*}, b_{k^*})$ is investigated in Figure 1 (BOTTOM-LEFT):

¹⁰Here $C > 0$ is the constant in Assumption (A1).

¹¹Since $N_T \ll N_t$, the computational cost associated with calculating $\mathcal{P}(NN(\cdot; W_k, b_k))$ is significantly lower compared to the cost of calculating the average relative residual on the training set.

the relative prediction error $\|NN(z_i; W_{k^*}, b_{k^*}) - y_i^\delta\|/\|y_i^\delta\|$ is plotted for the test set $\{z_i, N_t \leq i < N_t + N_T - 1\}$ (BLUE), the average value (RED) is 0.010. **The performance $\mathcal{P}(NN(\cdot; W_{k^*}, b_{k^*}))$ amounts to 99%.**

For comparison purpose we plot in Figure 1 (BOTTOM-RIGHT) the prediction accuracy of the linear regression approach [29] for the same test set (BLUE), the average value is 0.021 (RED). **The performance of this approach amounts to 97.7%.**

4.2 Classification problem for the MNIST database

The *Modified National Institute of Standards and Technology* (MNIST) database consists of images of handwritten digits (Figure 2). Each image is accompanied by a corresponding label indicating the digit it represents. This dataset is commonly used in the field of machine learning for developing neural network architectures, and for testing training algorithms for neural networks.

The MNIST database contains 60,000 training images (along with 10,000 testing images) of the ten digits. Each image consists of a 28×28 pixel array of grayscale levels.



Figure 2: Sample images from the MNIST database (source Wikipedia).

The corresponding data-files are accessible from many different sources. For the experiments conducted here, the files were downloaded from <http://yann.lecun.com/exdb/mnist/>.

In this section we consider the well-known classification problem for the MNIST database. In order to model this problem, we use here a NN that inputs a 28×28 pixel array of grayscale levels (i.e. a vector in \mathbb{R}^{784} with coordinates ranging from zero (black) to 255 (white)) and outputs a vector in \mathbb{R}^{10} . The classification of the handwritten digit depicted in the 28×28 image is given by the coordinate of this output vector with maximal absolute value (for alternative NN architectures with outputs in \mathbb{R}^{10} we refer the reader to [5] for a Deep NN, or [4] for a Convolutional NN).

The architecture of the NN used in our experiments is as follows:

- **Input:** $z \in \mathbb{R}^{784}$, pixel array from the MNIST database;
- **Hidden layer:** $\tilde{z} := \sigma_1(W_1 z + b_1) \in \mathbb{R}^{64}$, where $W_1 \in \mathbb{R}^{64,784}$ and $b_1 \in \mathbb{R}^{64}$;
- **Output:** $NN(z; W_1, b_1, W_2, b_2) := \sigma_2(W_2 \tilde{z} + b_2) \in \mathbb{R}^{10}$, where $W_2 \in \mathbb{R}^{10,64}$ and $b_2 \in \mathbb{R}^{10}$.

Here W_1, W_2 are weight matrices and b_1, b_2 are biases vectors. Moreover, $\sigma_1 : \mathbb{R}^{64} \rightarrow \mathbb{R}^{64}$ and $\sigma_2 : \mathbb{R}^{10} \rightarrow \mathbb{R}^{10}$ are nonlinear activation functions.

The classification of the input image z is given by the scalar value $j \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ defined by $j := \arg \max_{0 \leq i \leq 9} |NN_i(z; W_1, b_1, W_2, b_2)|$.¹²

This simple NN has only 2-layers 784–64–10 (one hidden layer and the output layer); the dimension of the corresponding parameter space is $50,890 = 64(784 + 1) + 10(64 + 1)$, i.e. the dimension of the set of parameters (W_1, b_1, W_2, b_2) .

Typically, much larger NN's are used for solving the MNIST classification problem (e.g., the Deep NN in [5] has 6-layers 784-2500-2000-1500-1000-500-10 and achieves an accuracy rate of 99.65%). Our goal with this experiment is not to investigate state-of-the-art NN-architectures. Instead, we aim to test the efficiency of the pSGD method in (6) as a training algorithm. For this purpose, the above described NN-architecture is rich enough to define a challenging inverse problem as we shall see next.

¹²We adopt here the notation $NN(\cdot) = [NN_i(\cdot)]_{i=0}^9 \in \mathbb{R}^{10}$.

The inverse problem under consideration is a NN training problem, i.e. one aims to find an approximate solution (W_1, b_1, W_2, b_2) to the nonlinear system

$$F_i(W_1, b_1, W_2, b_2) = y_i, \quad i = 0, \dots, N_t - 1. \quad (27)$$

Here $N_t = 60,000$ is the size of the training set and $F_i(W_1, b_1, W_2, b_2) := NN(z_i; W_1, b_1, W_2, b_2) = \sigma_2(W_2 \sigma_1(W_1 z_i + b_1) + b_2)$, where z_i is the i^{th} -image of the MNIST database for $i = 0, \dots, N_t - 1$. The right hand side in (27) is a vector of the type $y_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^{10}$, where the index of the coordinate with value “1” indicates the digit depicted in the image z_i (e.g., if z_i is an image of the digit 2, then $y_i = (0, 1, 0, \dots, 0)$). Note that the data in (27) is exact, i.e. the noise level is $\delta = 0$.

The activation functions $\sigma_1, \sigma_2 : \mathbb{R} \rightarrow \mathbb{R}$ used in the above NN are variations of the *sigmoid function* [6], namely: $\sigma_1(t) = \frac{1}{2} \tanh(t/10)$ and $\sigma_2(t) = 2 \tanh(t/10)$.

Numerical implementations

The pSGD method in (6) is implemented for solving the NN training problem (27). It is worth mentioning that the operators F_i in (27) do not satisfy the wTCC (4).

In our numerical experiments, the initial guess $(W_1^0, b_1^0, W_2^0, b_2^0)$ consists of random matrices/vectors with coordinate values ranging in $(-1, 1)$. Since the noise level is $\delta = 0$, we set $p^\delta(t) = t^2$ in (6b).

Three different runs of the pSGD method are presented in Figure 3. In each one of them the iteration $(W_1^k, b_1^k, W_2^k, b_2^k)$ is computed for 20 epochs, i.e. for $k = 1, \dots, 20N_t$. In the first 2 runs we choose the sequence $\theta_k \equiv 1$ in (A4), while in the last run a random sequence $\theta_k \in (0, 2)$ is chosen. The numerical results plotted in Figure 3 show:

$$\text{(TOP) Evolution of relative residual on the training set: } \sum_{i=0}^{N_t-1} \frac{\|NN(z_i; W_1^k, b_1^k, W_2^k, b_2^k) - y_i\|}{\|NN(z_i; W_1^0, b_1^0, W_2^0, b_2^0) - y_i\|},$$

$$\text{(BOTTOM) Evolution of relative residual on the test set: } \sum_{i=N_t}^{N_t+N_T-1} \frac{\|NN(z_i; W_1^k, b_1^k, W_2^k, b_2^k) - y_i\|}{\|NN(z_i; W_1^0, b_1^0, W_2^0, b_2^0) - y_i\|}.$$

Here $N_T = 10,000$ is the number of images in the MNIST database test set.

For comparison purposes, the SGD and the ADAM methods were implemented for solving (27) (the ADAM method is implemented using the choice of parameters described in Section 4.1). The evolution of the corresponding residuals are plotted in Figure 3.

Following every $\frac{1}{10}N_t$ steps, the average relative residual is computed on the test set; the index $0 \leq k^* \leq 20N_t$ is chosen such that $(W_1^{k^*}, b_1^{k^*}, W_2^{k^*}, b_2^{k^*})$ exhibits the smallest relative residual. After selecting the set of parameters $(W_1^{k^*}, b_1^{k^*}, W_2^{k^*}, b_2^{k^*})$, the accuracy rate of the corresponding neural network $NN(\cdot; W_1^{k^*}, b_1^{k^*}, W_2^{k^*}, b_2^{k^*})$ is calculated using the test set. In the experiment above, $k^* = 19.4N_t$ is obtained from the first run of the pSGD method (run 1 in Figure 3). The accuracy rate of $NN(\cdot; W_1^{k^*}, b_1^{k^*}, W_2^{k^*}, b_2^{k^*})$ is 95.96%.

Some remarks regarding the numerical experiments:

- **Computation of k^* :** Since $N_T \ll N_t$, the computational cost associated with calculating the (average) residual on the test set is insignificant compared to the cost of computing the (average) residual on the training set. The computation required in order to determine k^* is performed 10 times per epoch; this additional task does not impose a significant numerical burden on the implementation of the pSGD method.
- **Choice of (θ_k) :** The experiments presented above suggest that the choice of the relaxation parameters (θ_k) does not significantly impact the decay rate of the residual. Different runs of the pSGD method using sequences $\theta_k \in (0, 1)$ (under relaxation), or $\theta_k \in (1, 2)$ (over relaxation), or random $\theta_k \in (0, 2)$ produce similar numerical results.
- **Accuracy rate:** The accuracy rate of $NN(\cdot; W_1^{k^*}, b_1^{k^*}, W_2^{k^*}, b_2^{k^*})$ is given by the trace of the confusion matrix divided by N_T . The confusion matrix is a table that is used to evaluate the

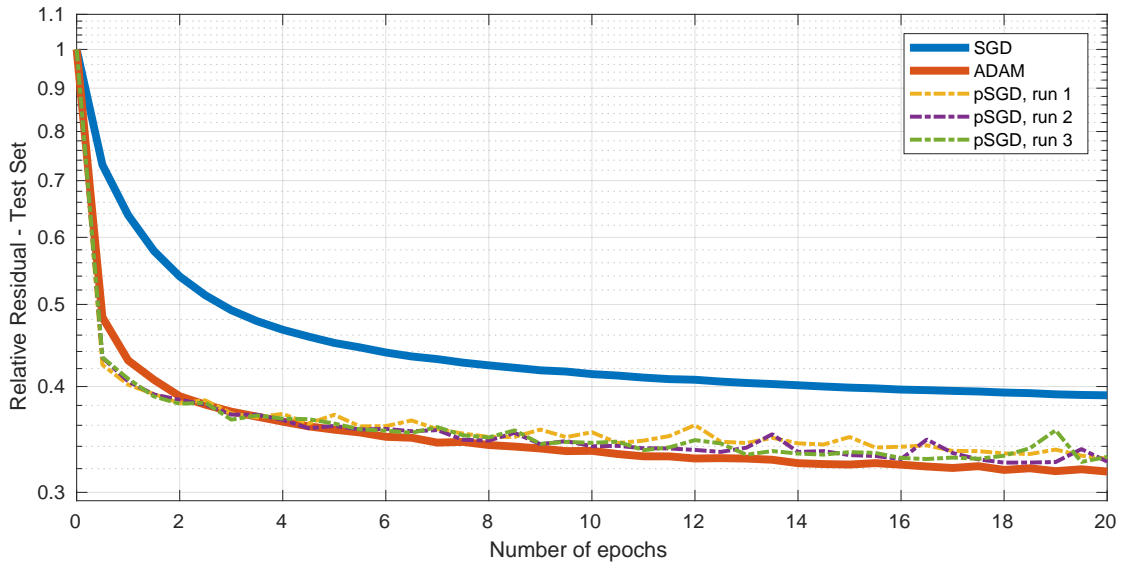
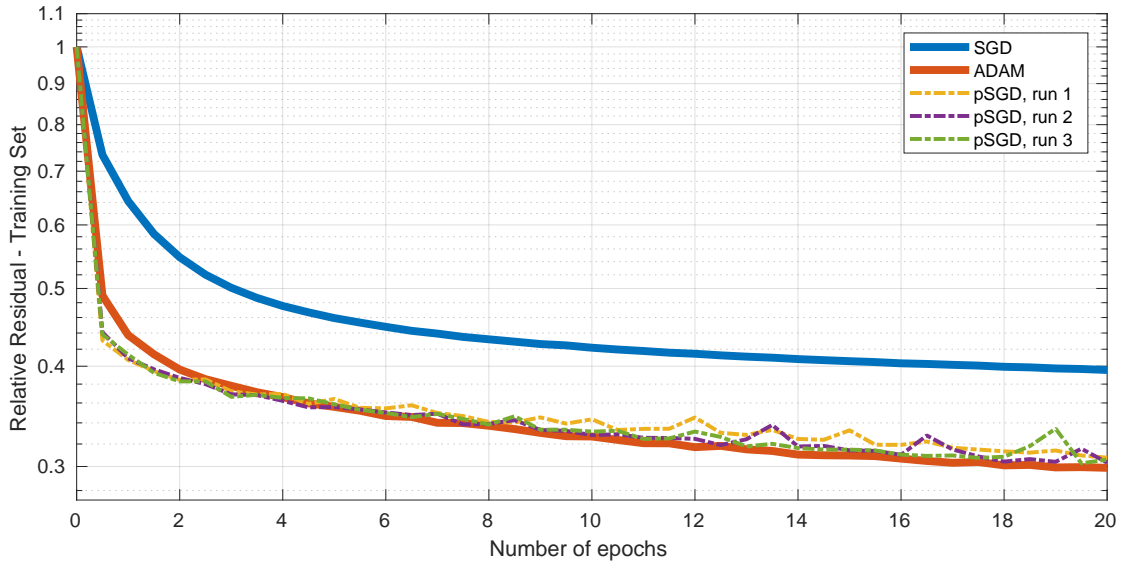


Figure 3: MNIST classification problem. (TOP) Evolution of the relative residual for training set; (BOTTOM) Evolution of the relative residual for test set.

performance of a classification model. It provides a summary of how well the model has classified the different classes in a dataset. It is typically used for problems like the MNIST classification, where the output of the model can belong to multiple classes. It displays the actual class labels of the data against the predicted class labels generated by the model. The main diagonal in Figure 4 represents the correctly classified instances, while the off-diagonal elements represent misclassifications. The final entry in a row/column represents the cumulative sum of all preceding elements in that particular row/column.

4.3 Parameter identification in a 3D elliptic PDE system

The previous applications in Sections 4.1 and 4.2 are related to the training of neural networks using real data. Despite the practical importance of these applications, the exact solution to the inverse problem and the level of noise is not known in both cases, making it only possible to analyze the evolution of the residual and not the iteration error.

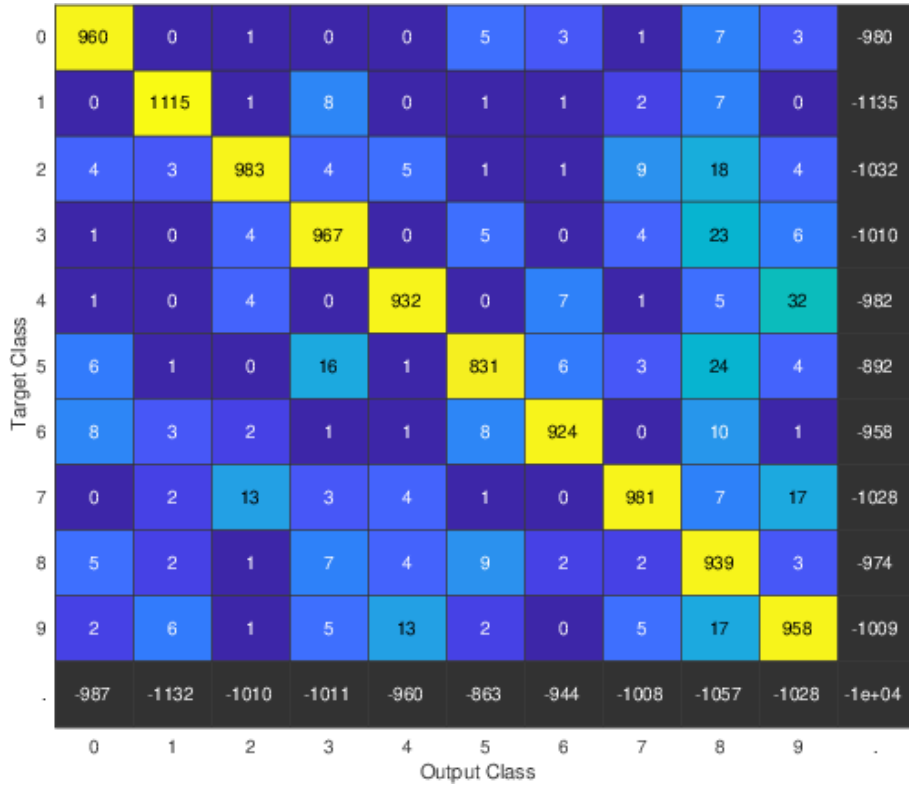


Figure 4: MNIST classification problem. Confusion matrix for $NN(\cdot; W_1^{k*}, b_1^{k*}, W_2^{k*}, b_2^{k*})$.

In this section we address a problem of parameter identification in a system of three-dimensional elliptic PDEs. Although this is a benchmark problem with synthetic data, it has certain features that are convenient for testing the efficiency of our method: (i) the exact solution is known; (ii) the level of noise is known; (iii) each operator in the system of elliptic PDEs is known to satisfy (A2).

The underlying inverse problem

We address the problem of determining the non-negative coefficient $c \in L^\infty(\Omega)$ in the system of elliptic PDEs defined on the three-dimensional domain $\Omega = (0, 1)^3$ with homogeneous Dirichlet boundary conditions

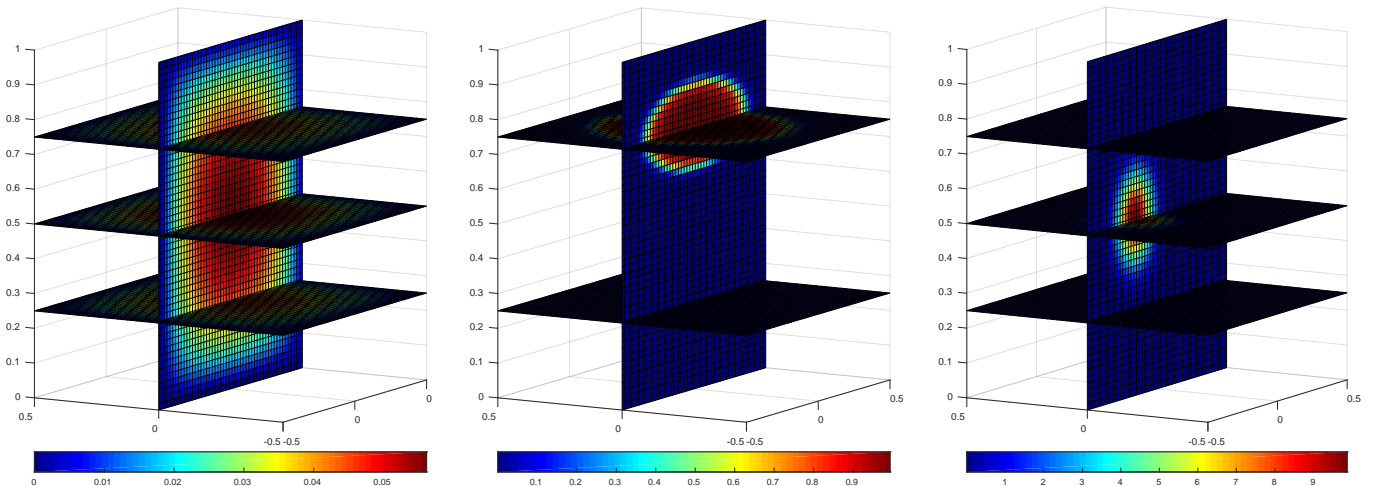


Figure 5: Section 4.3. Problem setup: (LEFT) Initial guess c_0 ; (CENTER) Ground truth c^* ; (RIGHT) One of the 27 functions f_i on the right-hand-side of system (28).

$$-\Delta u_i + cu_i = f_i, \text{ in } \Omega \quad u_i = 0, \text{ on } \partial\Omega \quad i = 0, \dots, N-1. \quad (28)$$

Here the function $c \geq 0$ almost everywhere is to be identified from the knowledge of $u_i \in H^1(\Omega) \subset L^2(\Omega)$, $i = 0, \dots, N-1$, on the entire domain Ω ,¹³ the right-hand sides $f_i \in L^2(\Omega)$ are known. Notice that for each $f_i \in L^2(\Omega)$ and $c \in L^\infty(\Omega)$ the boundary value problem in (28) admits a unique solution $u \in H_0^1(\Omega)$ [25].

This inverse problem can be written in the form of system (2), where $F_i : c \mapsto u_i$, u_i solves (28), for $i = 0, \dots, N-1$ (since $\Omega \subset \mathbb{R}^3$ is bounded, it holds $L^\infty(\Omega) \subset L^2(\Omega)$). In the case $N = 1$ this ill-posed benchmark problem is considered in [14, Example 4.2]. In particular the authors prove that each operator F_i does satisfy (A2).¹⁴

Numerical implementations

The pSGD method in (6) is implemented for solving the system (28). In our numerical experiments $N = 27$ equations, the initial guess c_0 is the solution of the boundary value problem $-\Delta u = 1$ in Ω and $u = 0$ on $\partial\Omega$, and the ground truth $c^* \in L^2(\Omega)$ is the non-negative function $c^*(x, y, z) = \frac{2}{\pi} \arctan[100 \exp(-[40x^2 + 40y^2 + 320(z - 0.75)^2])]$ (see Figure 5). The functions f_i are defined by $f_i(x, y, z) = 10 \exp[-8(x - x_i)^2 - 8(y - y_i)^2 - 8(z - z_i)^2]$, for $i = 0, \dots, 26$, where $(x_i, y_i, z_i) \in \Omega$ form a $3 \times 3 \times 3$ -array of equally spaced points within Ω . In Figure 5 (RIGHT) one of the 27 functions f_i on the right-hand-side of system (28) is depicted.

Regarding the generation of noisy data used in our experiments, the exact u_i (solutions of (28) for $c = c^*$) are perturbed by adding uniformly distributed random noise. The level of noise is $\delta = 0.5\%$. In order to avoid inverse crimes, finite element adaptively refined meshes (with average 6,300 elements) are used in the implementation of the pSGD method. These meshes are coarser than the meshes used to generate the data for the inverse problem (with average 11,000 elements).

For the computation of $\lambda_{I_k}^\delta$ in (6b), the step of pSGD, we use a conservative estimate for η in (A2), namely $\eta = 0.99$. Moreover, we set $\theta_k \equiv 1$ in (6a). The evolution of the relative iteration error $\|c_k^\delta - c^*\|/\|c^*\|$ for the pSGD method is presented in Figure 6.

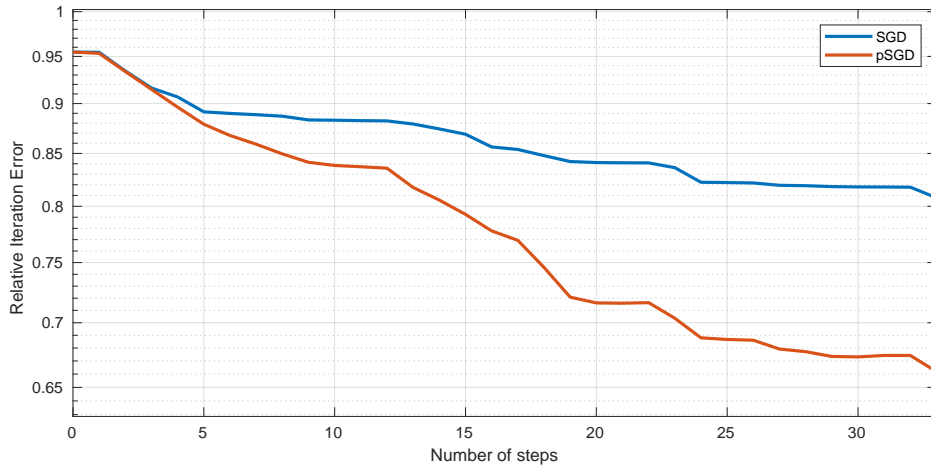


Figure 6: Section 4.3. Evolution of relative iteration error for pSGD and SGD methods.

¹³Since $H^1(\Omega) \subset L^2(\Omega)$, this does not conflict with the standard assumption in inverse problems that the data should belong to $L^2(\Omega)$.

¹⁴Note that the coefficient c in this inverse problem can be identified from a single measurement ($N = 1$) [14]. On the other hand, the quality of the numerical reconstruction of c can be significantly enhanced when additional data is available ($N > 1$), especially if the source functions f_i are strategically selected. These facts fully characterize the interest in the inverse problem described in system (28), both in the case $N = 1$ and $N > 1$.

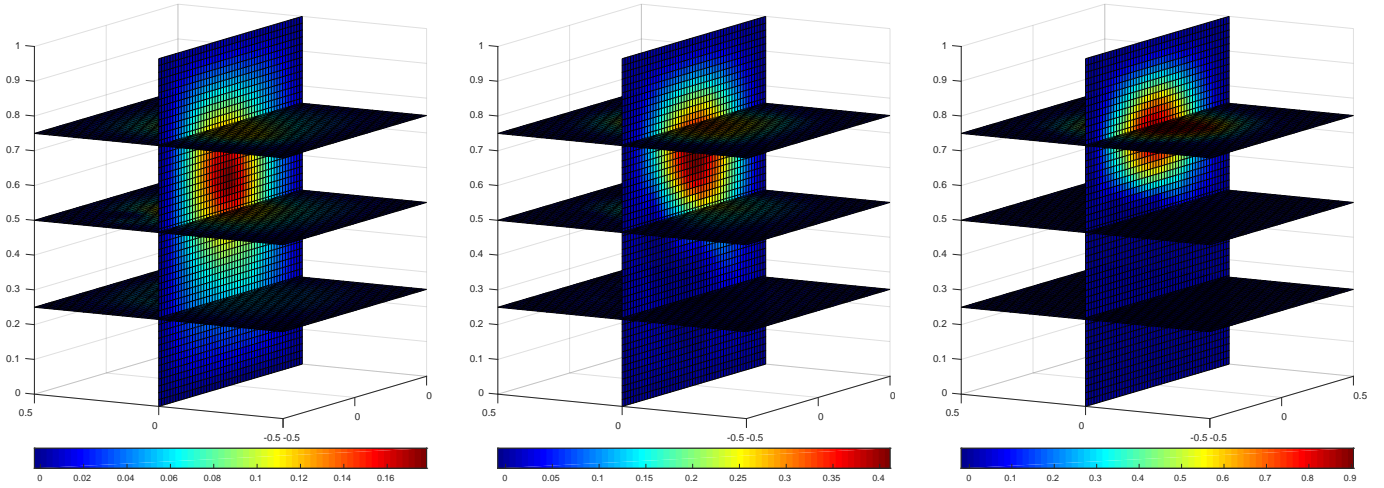


Figure 7: Section 4.3. Iterates c_4^δ (LEFT), c_{16}^δ (CENTER) and c_{64}^δ (RIGHT) of the pSGD method.

For comparison purposes, the SGD method was implemented for solving (28). In this example pSGD clearly outperforms the classical SGD method. The iterates c_4^δ , c_{16}^δ and c_{64}^δ of the pSGD method are presented in Figure 7.

A second experiment with $N = 1$ and $\delta = 5\%$ is conducted. Thus, system (28) reduces to the single equation $F_0(c) = u_0^\delta$, and the evolution of both relative iteration error $\|c_k^\delta - c^*\|/\|c^*\|$ and relative residual $\|F_0(c_k^\delta) - u_0^\delta\|/\|F_0(c_0) - u_0^\delta\|$ can be monitored (see Figure 8). The initial guess c_0 and the ground truth c^* are as before. An analogous procedure is used to generate the noisy data and to avoid inverse crimes. Again we use $\eta = 0.99$ and $\theta_k \equiv 1$ to compute the step of the pSGD method. In Figure 8 one observes the typical semi-convergence phenomenon for both pSGD and SGD methods.

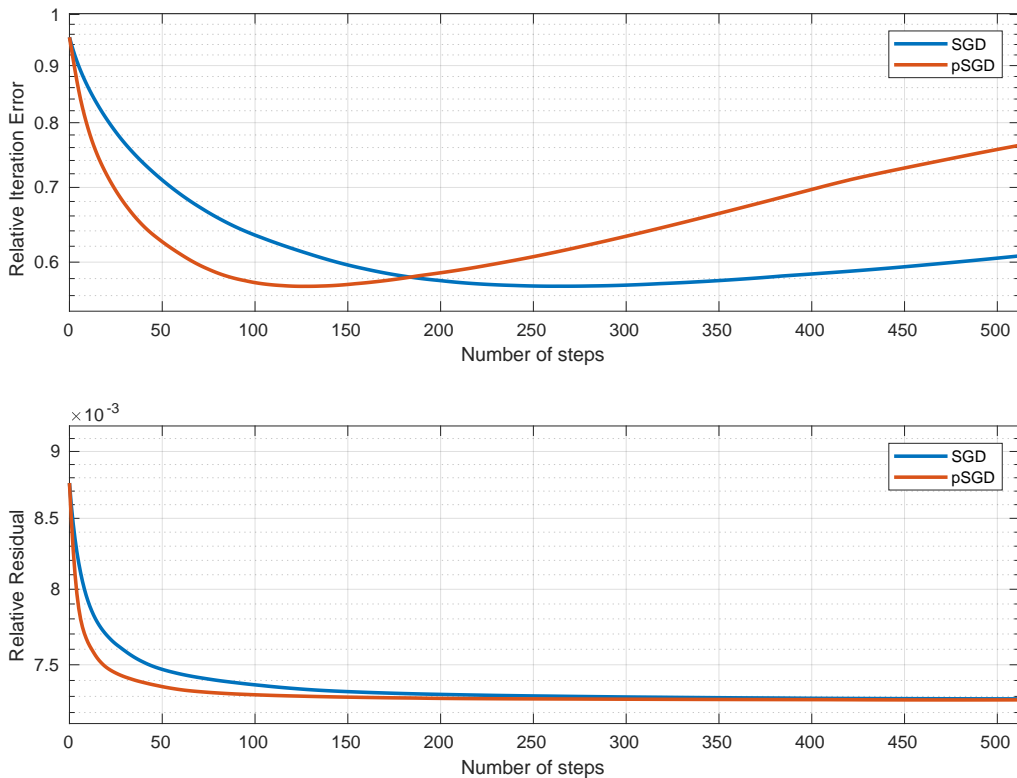


Figure 8: Section 4.3. Noise level $\delta = 5\%$. (TOP) Evolution of relative iteration error; (BOTTOM) Evolution of relative residual.

5 Conclusions

In this manuscript we investigate a nonlinear *projective stochastic-gradient* (pSGD) method for computing stable approximate solutions to large scale systems of nonlinear ill-posed equations.

We build upon a well-established nonlinear assumption, namely the weak tangential cone condition (wTCC), see (A2), to expand the method in [29, 28]. As a result, we create a new approach capable of efficiently solving large-scale systems of nonlinear equations.

Our method stands out due to the stepsize selection, which is inspired by the *projective Landweber (PLW) method* [24] and the *projective Landweber-Kaczmarz (PLWK) method* [23].

Highlighted among the key findings established in this manuscript are: (i) Estimates for the *average gain* and monotonicity results for the *average iteration error*; (ii) A convergence result for the pSGD method in the exact data case (Theorem 3.5); (iii) Regularization properties of the pSGD method: a stability result (Theorem 3.9) and a semi-convergence result (Theorem 3.11); (iv) In Lemma 4.3 we prove that the neural-network used to model the inverse problem in Section 4.1 satisfies the wTCC.

Numerical experiments are presented for two large scale nonlinear problems in machine learning: (i) the big data problem of CO-concentration prediction in a gas sensor array considered in [9, 29]; (ii) the classification problem for the MNIST database. In these experiments, pSGD performs on par with ADAM [21], one of the most efficient methods for solving large scale systems; this fact alone underscores the relevance of pSGD.

A third numerical experiment is conducted for a benchmark ill-posed problem of parameter identification in a system of three-dimensional elliptic PDEs, which satisfies the theoretical assumptions in this manuscript. The exact solution and level of noise are known; this enables us to analyze the convergence of the iterations toward the exact solution.

A Appendix: Proof of Lemma 4.2

In what follows we prove that, given a constant $a > 1$, the activation function σ in (25) (see Figure A) satisfies wTCC (4) in the interval $(1-a^2, a^2-1)$ for the constant $\eta = \frac{1}{2}(a-1)$.

The first step is to verify that the real function $h : x \mapsto \sqrt{x}$ does satisfy wTCC in the interval $[1, a^2]$ for the constant $\eta = \frac{1}{2}(a-1)$.

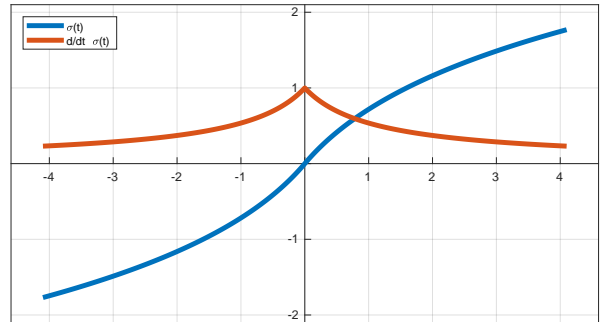


Figure A: Function $\sigma(t)$ in (25) and its derivative.

Given $x, y > 0$ it holds $h(x) - h(y) - h'(y)(x - y) = \sqrt{x} - \sqrt{y} - \frac{1}{2}(x - y)/\sqrt{y} = \sqrt{x} - \sqrt{y}/2 - x/(2\sqrt{y}) \leq 0$ (the inequality follows from the inequality of arithmetic and geometric means). Thus,

$$\begin{aligned} \|h(x) - h(y) - h'(y)(x - y)\| &= \frac{1}{2}\sqrt{y} - \sqrt{x} + \frac{x}{2\sqrt{y}} = \frac{1}{2}(\sqrt{y} - \sqrt{x}) + \frac{1}{2}\left(\frac{x}{\sqrt{y}} - \sqrt{x}\right) \\ &= -\frac{1}{2}(\sqrt{x} - \sqrt{y}) + \frac{\sqrt{x}}{2\sqrt{y}}(\sqrt{x} - \sqrt{y}) = \frac{1}{2}\left(\frac{\sqrt{x}}{\sqrt{y}} - 1\right)(\sqrt{x} - \sqrt{y}) \leq \frac{1}{2}\left|\frac{\sqrt{x}}{\sqrt{y}} - 1\right| \|h(x) - h(y)\|. \end{aligned}$$

Since $\left|\frac{\sqrt{x}}{\sqrt{y}} - 1\right| \leq a - 1$ for all $x, y \in [1, a^2]$, we conclude that $h(x) = \sqrt{x}$ indeed satisfies wTCC in the interval $[1, a^2]$ for the constant $\eta = \frac{1}{2}(a - 1)$.

Proof. (of Lemma 4.2)

(i) From the definition of σ in (25) follows immediately $\sigma(x) = 2h(x+1) - 2$, for $x \geq 0$. Consequently, σ satisfies the wTCC in the intervall $[0, a^2 - 1]$ with the same constant $\eta = \frac{1}{2}(a-1)$ as the function $h(x) = \sqrt{x}$.

(ii) On the other hand, since σ is an odd function it follows by symmetry that σ satisfies the wTCC in the intervall $[1 - a^2, 0]$ with the constant $\eta = \frac{1}{2}(a-1)$.

(iii) Define the real functions $f(x, y) := |\sigma(x) - \sigma(y) - \sigma'(y)(x-y)|/|\sigma(x) - \sigma(y)|$ and $g(x, y) := |\sigma(y) - \sigma(x) - \sigma'(x)(y-x)|/|\sigma(y) - \sigma(x)|$. Note that, for $x < 0$ and $y > 0$ it holds

$$f(x, y) = |4 - 2\sqrt{1-x} - 2\sqrt{y+1} - (x-y)/\sqrt{y+1}| / |4 - 2\sqrt{1-x} - 2\sqrt{y+1}|,$$

$$g(x, y) = |2\sqrt{y+1} + 2\sqrt{1-x} - 4 - (y-x)/\sqrt{1-x}| / |2\sqrt{y+1} + 2\sqrt{1-x} - 4|.$$

Defining the set $S := [1 - a^2, 0] \times [0, a^2 - 1]$, a direct calculation shows that

$$\operatorname{argmax}_{(x,y) \in S} f(x, y) = (1 - a^2, 0) \quad \text{and} \quad \operatorname{argmax}_{(x,y) \in S} g(x, y) = (0, a^2 - 1).$$

Since $f(1 - a^2, 0) = g(0, a^2 - 1) = \frac{1}{2}(a-1)$, this reasoning allow us to conclude that

$$|\sigma(x) - \sigma(y) - \sigma'(y)(x-y)| \leq \frac{1}{2}(a-1)|\sigma(x) - \sigma(y)|,$$

$$|\sigma(y) - \sigma(x) - \sigma'(x)(y-x)| \leq \frac{1}{2}(a-1)|\sigma(y) - \sigma(x)|,$$

for all $(x, y) \in [1 - a^2, 0] \times [0, a^2 - 1]$.

Adding up items (i), (ii) and (iii) we come to the conclusion that σ satisfies wTCC (4) in $[1 - a^2, a^2 - 1]$ for $\eta = \frac{1}{2}(a-1)$, concluding the proof. \square

References

- [1] A.B. Bakushinsky and M.Y. Kokurin. *Iterative Methods for Approximate Solution of Inverse Problems*, volume 577 of *Mathematics and Its Applications*. Springer, Dordrecht, 2004.
- [2] J. Baumeister, B. Kaltenbacher, and A. Leitão. On levenberg-marquardt-kaczmarz iterative methods for solving systems of nonlinear ill-posed equations. *Inverse Probl. Imaging*, 4(3):335–350, 2010.
- [3] M. Burger and B. Kaltenbacher. Regularizing Newton-Kaczmarz methods for nonlinear ill-posed problems. *SIAM J. Numer. Anal.*, 44:153–182, 2006.
- [4] D. Cireşan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3642–3649, 2012.
- [5] D.C. Cireşan, U. Meier, L.M. Gambardella, and J. Schmidhuber. Deep, big, simple neural nets for handwritten digit recognition. *Neural Computation*, 22(12):3207–3220, 2010.
- [6] P.L. Combettes and J.-C. Pesquet. Deep neural network structures solving variational inequalities. *Set-Valued Var. Anal.*, 28(3):491–518, 2020.
- [7] A. De Cezaro, M. Haltmeier, A. Leitão, and O. Scherzer. On steepest-descent-Kaczmarz methods for regularizing systems of nonlinear ill-posed equations. *Appl. Math. Comput.*, 202(2):596–607, 2008.
- [8] H.W. Engl, M. Hanke, and A. Neubauer. *Regularization of Inverse Problems*. Kluwer Academic Publishers, Dordrecht, 1996.

- [9] J. Fonollosa, S. Sheik, R. Huerta, and S. Marco. Reservoir computing compensates slow response of chemosensor arrays exposed to fast varying gas concentrations in continuous monitoring. *Sensors and Actuators B: Chemical*, 215:618–629, 2015.
- [10] C. W. Groetsch. *Stable Approximate Evaluation of Unbounded Operators*, volume 1894 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 2007.
- [11] M. Haltmeier, R. Kowar, A. Leitão, and O. Scherzer. Kaczmarz methods for regularizing nonlinear ill-posed equations. II. Applications. *Inverse Probl. Imaging*, 1(3):507–523, 2007.
- [12] M. Haltmeier, A. Leitão, and E. Resmerita. On regularization methods of EM-Kaczmarz type. *Inverse Problems*, 25:075008, 2009.
- [13] M. Haltmeier, A. Leitão, and O. Scherzer. Kaczmarz methods for regularizing nonlinear ill-posed equations. I. convergence analysis. *Inverse Probl. Imaging*, 1(2):289–298, 2007.
- [14] M. Hanke, A. Neubauer, and O. Scherzer. A convergence analysis of Landweber iteration for nonlinear ill-posed problems. *Numerische Mathematik*, 72:21–37, 1995.
- [15] Tim Jahn and Bangti Jin. On the discrepancy principle for stochastic gradient descent. *Inverse Problems*, 36(9):095009, 2020.
- [16] Bangti Jin and Xiliang Lu. On the regularizing property of stochastic gradient descent. *Inverse Problems*, 35(1):015004, 2018.
- [17] Bangti Jin, Zehui Zhou, and Jun Zou. On the convergence of stochastic gradient descent for nonlinear ill-posed problems. *SIAM Journal on Optimization*, 30(2):1421–1450, 2020.
- [18] Qinian Jin, Xiliang Lu, and Liuying Zhang. Stochastic mirror descent method for linear ill-posed problems in banach spaces. *Inverse Problems*, 39(6):065010, 2023.
- [19] Qinian Jin and Wei Wang. Landweber iteration of kaczmarz type with general non-smooth convex penalty functionals. *Inverse Problems*, 29(8):085011, 2013.
- [20] B. Kaltenbacher, A. Neubauer, and O. Scherzer. *Iterative regularization methods for nonlinear ill-posed problems*, volume 6 of *Radon Series on Computational and Applied Mathematics*. Walter de Gruyter GmbH & Co. KG, Berlin, 2008.
- [21] D.P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2017.
- [22] L. Landweber. An iteration formula for Fredholm integral equations of the first kind. *Amer. J. Math.*, 73:615–624, 1951.
- [23] A. Leitão and B.F. Svaiter. On projective Landweber-Kaczmarz methods for solving systems of nonlinear ill-posed equations. *Inverse Problems*, 32(1):025004, 2016.
- [24] A. Leitão and B.F. Svaiter. On a family of gradient type projection methods for nonlinear ill-posed problems. *Numerical Functional Analysis and Optimization*, 39(2-3):1153–1180, 2018.
- [25] J.L. Lions and E. Magenes. *Non-Homogeneous Boundary Value Problems and Applications*, volume 1. Springer, New York, 1972.
- [26] F. Margotti, A. Rieder, and A. Leitão. A Kaczmarz version of the reginn-Landweber iteration for ill-posed problems in Banach spaces. *SIAM J. Numer. Anal.*, 52(3):1439–1465, 2014.

- [27] V.A. Morozov. *Regularization Methods for Ill-Posed Problems*. CRC Press, Boca Raton, 1993.
- [28] J.C. Rabelo and A. Leitão. Addendum: On stochastic Kaczmarz type methods for solving large scale systems of ill-posed equations (2022 Inverse Problems **38** 025003). *Inverse Problems*, 38(5):Paper No. 059401, 3, 2022.
- [29] J.C. Rabelo, Y. Saporito, and A. Leitão. On stochastic kaczmarz type methods for solving large scale systems of ill-posed equations. *Inverse Problems*, 38(2):025003, 2022.
- [30] O. Scherzer. Convergence rates of iterated Tikhonov regularized solutions of nonlinear ill-posed problems. *Numerische Mathematik*, 66(2):259–279, 1993.
- [31] T.I. Seidman and C.R. Vogel. Well posedness and convergence of some regularisation methods for non-linear ill posed problems. *Inverse Probl.*, 5:227–238, 1989.
- [32] A.N. Tikhonov. Regularization of incorrectly posed problems. *Soviet Math. Dokl.*, 4:1624–1627, 1963.
- [33] A.N. Tikhonov and V.Y. Arsenin. *Solutions of Ill-Posed Problems*. John Wiley & Sons, Washington, D.C., 1977. Translation editor: Fritz John.