

**TRUNCATION AND CONVERGENCE ISSUES FOR BOUNDED
LINEAR INVERSE PROBLEMS IN HILBERT SPACE**

NOE CARUSO, ALESSANDRO MICHELANGELI, AND PAOLO NOVATI

ABSTRACT. We present a general discussion of the main features and issues that (bounded) inverse linear problems in Hilbert space exhibit when the dimension of the space is *infinite*. This includes the set-up of a consistent notation for inverse problems that are genuinely infinite-dimensional, the analysis of the finite-dimensional truncations, a discussion of the mechanisms why the error or the residual generically fail to vanish in norm, and the identification of practically plausible sufficient conditions for such indicators to be small in some weaker sense. The presentation is based on theoretical results together with a series of model examples and numerical tests.

1. INTRODUCTION AND OUTLOOK

In this note we discuss a number of features that are typical of bounded inverse linear problems set on *infinite-dimensional* Hilbert spaces, the infinite dimensionality being the source of phenomena that become most relevant in the numerical treatment, and are absent when the considered space instead has finite dimension.

More precisely, we shall focus on typical issues and behaviours of the sequence of truncated, finite-dimensional problems that arise from the discretisation of the original, infinite-dimensional one.

As we shall explain in a moment, for specific classes of infinite-dimensional inverse problems an already well-established insight is available in the literature concerning the solvability of the truncated problems and the convergence of the finite-dimensional solutions. However, for *generic* inverse problems the control of such issues is surely less developed and a systematic discussion is missing.

In this respect, we do not aim here at a comprehensive classification of infinite-dimensional inverse problems and we rather keep the point of view of presenting *generic* features and difficulties that look ‘unavoidable’ at the considered level of generality. In our intentions this should provide the setting for a future thorough analysis of classes of infinite-dimensional inverse problems.

For this reason, besides stating and proving our main results, the material will also be presented through several model examples (and counter-examples).

To fix the nomenclature and the notation, by an *inverse linear problem* in Hilbert space we shall mean the problem, given a Hilbert space \mathcal{H} , a linear operator A acting on \mathcal{H} , and a vector $g \in \mathcal{H}$, to determine the solution(s) $f \in \mathcal{H}$ to the linear equation

$$(1.1) \quad Af = g.$$

We shall say that: (1.1) is *solvable* if a solution f exists, namely if $g \in \text{ran}A$; (1.1) is *well-defined* if additionally the solution f is unique, i.e., if A is also injective (in which case one refers to f ‘*exact*’ solution); (1.1) is *well-posed* if there exists a

Date: November 20, 2018.

Key words and phrases. inverse linear problems, infinite-dimensional Hilbert space, ill-posed problems, orthonormal basis discretisation, bounded linear operators, Krylov subspaces, Krylov solution, GMRES, conjugate gradient, LSQR.

unique solution that depends continuously (i.e., in the norm of \mathcal{H}) on the datum g , equivalently, that $g \in \text{ran}A$ and A has bounded inverse on its range.

In applications, the linear law A that associates an input f to an output g is prescribed by some physical model, and hence within that model such a law is exactly known. Experimental measurements produce a possibly approximate knowledge of the output g , from which one wants to obtain information on the input f , which is the final object of interest.

Of course what is ‘exactly known’ of A is its domain and action as an operator acting on \mathcal{H} . Other relevant features of A might not be explicitly accessible, and only computable within some approximation: for example, if $A : \mathcal{H} \rightarrow \mathcal{H}$ is a (everywhere defined) Hilbert-Schmidt operator, one may know its integral kernel, based on the theoretical framework within which the problem is modelled, however it might not be possible to write explicitly (exactly) its singular value decomposition.

Although well-defined inverse linear problems are in a sense trivial theoretically, as the existence and uniqueness of the solution is not of concern, it is clear that there are at least two main issues arising when one aims at solving them numerically.

The first, which is typical already at the finite-dimensional level, namely when A is a matrix, is the fact that the measurement of g is in practice plagued by some noise, or error of sort: as a consequence, numerically one has to deal with the possibly ill-posed problem

$$(1.2) \quad Af = g + \nu,$$

where the ‘true’ physical output is some $g \in \text{ran}A$, however the actually measured output is $g + \nu$, with some small noise-like perturbation $\nu \in \mathcal{H}$ for which possibly $g + \nu \notin \text{ran}A$.

The second issue is actually typical of the *infinite-dimensional* setting, on which in fact we are going to focus most of our discussion, namely when $\dim \mathcal{H} = \infty$ and A is a genuine infinite-dimensional operator on \mathcal{H} . By this we mean, as customary [21, Sect. 1.4], that A is *not* reduced to $A = A_1 \oplus A_2$ by an orthogonal direct sum decomposition $\mathcal{H} = \mathcal{H}_1 \oplus \mathcal{H}_2$ with $\dim \mathcal{H}_1 < \infty$, $\dim \mathcal{H}_2 = \infty$, and $A_2 = \mathbb{O}$.

Clearly, non-trivial inverse linear problems in the sense just described are to be *truncated* to a finite-dimensional Hilbert space, in order to be treated numerically. This poses the questions on how close the solution(s) to the truncated problem are with respect to the exact solution, let alone on whether the truncated problem is solvable itself.

All this is very familiar and already under control for relevant classes of boundary value problems on $L^2(\Omega)$ for some domain $\Omega \subset \mathbb{R}^d$, the typical playground for Galerkin and Petrov-Galerkin finite element methods [8, 19]. In these cases A is an *unbounded* operator, say, of elliptic type [8, Chapter 3], [19, Chapter 4], of Friedrichs type [8, Sect. 5.2], [9, 1, 2], of parabolic type [8, Chapter 6], [19, Chapter 5], of ‘mixed’ (i.e., inducing saddle-point problems) type [8, Sect. 2.4 and Chapter 4], etc. Such A ’s are assumed to satisfy (and so they do in applications) some kind of coercivity, or more generally one among the various classical conditions that ensure the corresponding problem (1.1) to be well-posed, such as the Banach-Nečas-Babuška Theorem or the Lax-Milgram Lemma [8, Chapter 2].

For the above-mentioned classes of inverse linear problems, the finite-dimensional truncation and the infinite-dimensional error analysis are widely studied and well understood, as we shall comment further in due time. In that context, in order for the finite-dimensional solutions to converge strongly, one requires stringent yet often plausible conditions [8, Sect. 2.2-2.4], [19, Sect. 4.2] both on the truncation spaces, that need to approximate suitably well the ambient space \mathcal{H} (*‘approximability’*, thus the interpolation capability of finite elements), and on the behaviour of the reduced problems, that need admit solutions that are uniformly controlled by

the data (*'uniform stability'*), and that are suitably good approximate solutions of the original problem (*'asymptotic consistency'*), together with some suitable boundedness of the problem in appropriate topologies (*'uniform continuity'*).

As plausible as the above conditions are, they are *not* matched by several other types of inverse problem of applied interest. Mathematically this is the case whenever A does not have a 'good' inverse, for instance when A is a compact operator on \mathcal{H} with arbitrarily small singular values, or when the exact solution of the inverse problem does not belong to the corresponding Krylov space used for the finite-dimensional truncations.

For such an abstract level of generality, for compact and generic bounded inverse linear problems, in this work we set up the theoretical formalism and settle the analysis of the above-mentioned questions specifically when the dimension of the underlying Hilbert space is infinite.

As declared already, the purpose is to highlight non-trivial features typical of infinite dimensionality and discuss them through an amount of model examples that challenge the common intuition.

In particular, we carry on the point of view that error and residual may be controlled in a still informative way in some *weaker* sense than the expected norm topology of the Hilbert space. In this respect, we identify practically plausible sufficient conditions for the error or the residual to be small in such generalised senses and we discuss the *mechanisms* why the same indicators may actually fail to vanish in norm.

In the concluding part of the work, we investigate the main features discussed theoretically through a series of numerical tests, focusing on the truncation of infinite-dimensional inverse problems when the dimension of the truncation space increases.

General notation. Besides further notation that will be declared in due time, we shall keep the following convention. \mathcal{H} denotes a complex Hilbert space, that will be separable throughout this note, with norm $\|\cdot\|_{\mathcal{H}}$ and scalar product $\langle \cdot, \cdot \rangle$, anti-linear in the first entry and linear in the second. Bounded operators on \mathcal{H} are tacitly understood to be linear and everywhere defined. $\|\cdot\|_{\text{op}}$ denotes the corresponding operator norm. The space of bounded operators on \mathcal{H} is denoted with $\mathcal{B}(\mathcal{H})$. The spectrum of an operator A is denoted by $\sigma(A)$. $\mathbb{1}$ and $\mathbb{0}$ are, respectively, the identity and the zero operator, meant as finite matrices or infinite-dimensional operators depending on the context. An upper bar denotes the complex conjugate \bar{z} when $z \in \mathbb{C}$, and the norm closure $\bar{\mathcal{V}}$ of the span of the vectors in \mathcal{V} when \mathcal{V} is a subset of \mathcal{H} . For $\psi, \varphi \in \mathcal{H}$, by $|\psi\rangle\langle\psi|$ and $|\psi\rangle\langle\varphi|$ we shall denote the $\mathcal{H} \rightarrow \mathcal{H}$ rank-one maps acting respectively as $f \mapsto \langle\psi, f\rangle\psi$ and $f \mapsto \langle\varphi, f\rangle\psi$ on generic $f \in \mathcal{H}$. For identities such as $\psi(x) = \varphi(x)$ in L^2 -spaces we will tacitly understand the 'for almost every x ' specification in the equality.

2. FINITE-DIMENSIONAL TRUNCATION

2.1. Set up and notation.

Let us start with setting up a convenient formalism for the treatment of finite-dimensional truncations of linear inverse problems in infinite-dimensional Hilbert space. In the framework of Galerkin and Petrov-Galerkin methods this is customarily referred to as the *'approximation setting'* [8, Sect. 2.2.1].

Let $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$ be two orthonormal *systems* of the considered Hilbert space \mathcal{H} . They need not be orthonormal *bases*, although their completeness is crucial for the goodness of the approximation.

In practice these are two *explicitly known* sets of orthonormal vectors (unlike, for instance, the *possibly non-explicit* orthonormal bases expressing the singular value decomposition of a given compact operator) that are going to be used in a numerical algorithm. In the Petrov-Galerkin nomenclature [8, 19] the u_n 's and the v_n 's span respectively the so-called '*solution space*' (or '*trial space*') and the '*test space*' of the problem.

The choice of $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$ depends on the specific approach. In the framework of finite element methods they can be taken to be the global shape functions of the interpolation scheme [8, Chapter 1]. For Krylov subspace methods they are just the spanning vectors of the associated Krylov subspace [16, Chapter 2].

Correspondingly, for each $N \in \mathbb{N}$, the orthonormal projections in \mathcal{H} respectively onto $\text{span}\{u_1, \dots, u_N\}$ and $\text{span}\{v_1, \dots, v_N\}$ shall be

$$(2.1) \quad P_N := \sum_{n=1}^N |u_n\rangle\langle u_n|, \quad Q_N := \sum_{n=1}^N |v_n\rangle\langle v_n|.$$

Associated to a given well-defined linear inverse problem $Af = g$ in \mathcal{H} as (1.1), one considers the finite-dimensional truncations induced by P_N and Q_N , hence, for each N , the problem to find solutions $\widehat{f}^{(N)} \in P_N\mathcal{H}$ to the equation

$$(2.2) \quad (Q_N A P_N) \widehat{f}^{(N)} = Q_N g.$$

In (2.2) $Q_N g = \sum_{n=1}^N \langle v_n, g \rangle v_n$ is the datum and $\widehat{f}^{(N)} = \sum_{n=1}^N \langle u_n, \widehat{f}^{(N)} \rangle u_n$ is the unknown, and the compression $Q_N A P_N$ is only non-trivial as a map from $P_N\mathcal{H}$ to $Q_N\mathcal{H}$, its kernel containing at least the subspace $(\mathbb{1} - P_N)\mathcal{H}$.

Clearly, (2.2) (and more precisely (2.5) below) is nothing but the truncated problem arising from the oblique projection of the Petrov-Galerkin scheme. When the special choice $(u_n)_{n \in \mathbb{N}} = (v_n)_{n \in \mathbb{N}}$ is made, and hence $P_N = Q_N$ for all N 's, this is the orthogonal projection approach of the ordinary Galerkin scheme.

There is an obvious and non-relevant degeneracy (which is infinite when $\dim \mathcal{H} = \infty$) in (2.2) when it is regarded as a problem on the whole \mathcal{H} . The actual interest towards (2.2) is the problem resulting from the identification $P_N\mathcal{H} \cong \mathbb{C}^N \cong Q_N\mathcal{H}$, in terms of which $P_N f \in \mathcal{H}$ and $Q_N g \in \mathcal{H}$ are canonically identified with the vectors

$$(2.3) \quad f_N = \begin{pmatrix} \langle u_1, f \rangle \\ \vdots \\ \langle u_N, f \rangle \end{pmatrix} \in \mathbb{C}^N, \quad g_N = \begin{pmatrix} \langle v_1, g \rangle \\ \vdots \\ \langle v_N, g \rangle \end{pmatrix} \in \mathbb{C}^N,$$

and $Q_N A P_N$ with a $\mathbb{C}^N \rightarrow \mathbb{C}^N$ linear map represented by the $N \times N$ matrix $A_N = (A_{N;ij})_{i,j \in \{1, \dots, N\}}$

$$(2.4) \quad A_{N;ij} = \langle v_i, Q_N A P_N u_j \rangle.$$

The matrix A_N is what in the framework of finite element methods for partial differential equations is customarily referred to as the '*stiffness matrix*'.

We shall call the inverse linear problem

$$(2.5) \quad A_N f^{(N)} = g_N$$

with datum $g_N \in \mathbb{C}^N$ and unknown $f^{(N)} \in \mathbb{C}^N$, and matrix A_N defined by (2.4), the N -dimensional truncation of the original problem $Af = g$.

Let us stress the meaning of the present notation.

- $Q_N A P_N$, $P_N f$, and $Q_N g$ are objects (one operator and two vectors) referred to the whole Hilbert space \mathcal{H} , whereas A_N , $f^{(N)}$, f_N , and g_N are the analogues referred now to the space \mathbb{C}^N .

- Moreover, the subscript in A_N , f_N , and g_N indicates that the components of such objects are precisely the corresponding components, up to order N , respectively of A , f , and g , with respect to the tacitly declared bases $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$, through formulas (2.3)-(2.4).
- As opposite, the superscript in $f^{(N)}$ indicates that the components of the \mathbb{C}^N -vector $f^{(N)}$ are not necessarily to be understood as the first N components of the \mathcal{H} -vector f with respect to the basis $(u_n)_{n \in \mathbb{N}}$, and in particular for $N_1 < N_2$ the components of $f^{(N_1)}$ are not a priori equal to the first N_1 components of $f^{(N_2)}$. In fact, if $f \in \mathcal{H}$ is a solution to $Af = g$, it is evident from obvious counterexamples that in general the truncations A_N , f_N , g_N do *not* satisfy the identity $A_N f_N = g_N$, whence the notation $f^{(N)}$ for the unknown in (2.5).
- Last, for a \mathbb{C}^N -vector $f^{(N)}$ the notation $\widehat{f^{(N)}}$ indicates a vector in \mathcal{H} whose first N components, with respect to the basis $(u_n)_{n \in \mathbb{N}}$, are precisely those of $f^{(N)}$, all others being zero. Thus, as pedantic as it looks, $f^{(N)} = (\widehat{f^{(N)}})_N$ and $f_N = (\widehat{f_N})_N$, and of course in general $f \neq \widehat{f_N}$.

With A , g , $(u_n)_{n \in \mathbb{N}}$, and $(v_n)_{n \in \mathbb{N}}$ explicitly known, the truncated problem (2.5) is explicitly formulated and, being finite-dimensional, it is suited for numerical algorithms.

This poses the general question on *whether the truncated problem itself is solvable, and whether its exact or approximate solution $f^{(N)}$ is close to the exact solution f and in which (possibly quantitative) sense.*

Let us elaborate more on these two issues in the following two subsections.

2.2. Singularity of the truncated problem.

It is clear, first of all, that the question of the singularity of the truncated problem (2.5) makes sense here *eventually in N* , meaning for all N 's that are large enough. For a *fixed* value of N the truncation might drastically alter the problem so as to make it manifestly non-informative as compared to $Af = g$, such alteration then disappearing for larger values.

Yet, even when the solvability of $A_N f^{(N)} = g_N$ is inquired eventually in N , it is no surprise that the answer is generically negative.

Example 2.1. That the matrix A_N may remain singular for arbitrary N even when the operator A is injective can be seen, for example, with the truncation of the weighted (compact) right-shift operator $\mathcal{R} = \sum_{n=1}^{\infty} \sigma_n |e_{n+1}\rangle \langle e_n|$ on $\ell^2(\mathbb{N})$ (Sect. A.3) with respect to the basis $(e_n)_{n \in \mathbb{N}}$ itself: indeed,

$$(2.6) \quad \mathcal{R}_N = \begin{pmatrix} 0 & \cdots & \cdots & \cdots & 0 \\ \sigma_1 & 0 & \cdots & \cdots & 0 \\ 0 & \sigma_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & \cdots & \sigma_{N-1} & 0 \end{pmatrix}$$

is singular irrespectively of N , with $\ker \mathcal{R}_N = \text{span}\{e_N\}$. (See Lemma 2.3 below for a more general perspective on such an example.)

It is not difficult to cook up variations of the above example where the matrix A_N is alternatingly singular and non-singular as $N \rightarrow \infty$.

Example 2.2. Of course, on the other hand, it may also well happen that the truncated matrix is always non-singular: the truncation of the multiplication operator

(see Sect. A.1)

$$M = \sum_{n=1}^{\infty} \frac{1}{n} |e_n\rangle\langle e_n|$$

on $\ell^2(\mathbb{N})$ with respect to $(e_n)_{n \in \mathbb{N}}$ yields the matrix $M_N = \text{diag}(1, \frac{1}{2}, \dots, \frac{1}{N})$, which is a $\mathbb{C}^N \rightarrow \mathbb{C}^N$ bijection for every N .

In fact, ‘bad’ truncations are always possible, as the following mechanism shows.

Lemma 2.3. *Let \mathcal{H} be a separable Hilbert space with $\dim \mathcal{H} = \infty$, and let $A \in \mathcal{B}(\mathcal{H})$. There always exist two orthonormal bases $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$ of \mathcal{H} such that the corresponding truncated matrix A_N defined as in (2.4) is singular for every $N \in \mathbb{N}$.*

Proof. Let us pick an arbitrary orthonormal basis $(u_n)_{n \in \mathbb{N}}$ and construct the other basis $(v_n)_{n \in \mathbb{N}}$ inductively. When $N = 1$, it suffices to choose v_1 such that $v_1 \perp Au_1$ and $\|v_1\|_{\mathcal{H}} = 1$. Let now $(v_n)_{n \in \{1, \dots, N-1\}}$ be an orthonormal system in \mathcal{H} satisfying the thesis up to the order $N - 1$ and let us construct v_N so that $(v_n)_{n \in \{1, \dots, N\}}$ satisfies the thesis up to order N . To this aim, let us show that a choice of v_N is always possible so that the final row in the matrix A_N has all zero entries. In fact, $(A_N)_{ij} = (Q_N A P_N)_{ij} = \langle v_i, Au_j \rangle$ for $i \in \{1, \dots, N - 1\}$ and $j \in \{1, \dots, N\}$ and in order for $\langle v_N, Au_j \rangle = 0$ for $j \in \{1, \dots, N\}$ it suffices to take

$$v_N \perp \text{ran}(A P_N), \quad v_N \perp \text{ran} Q_{N-1}, \quad \|v_N\|_{\mathcal{H}} = 1,$$

where P_N and Q_{N-1} are the orthogonal projections defined in 2.1. Since $\text{ran}(A P_N)$ and $\text{ran} Q_{N-1}$ are finite-dimensional subspaces of \mathcal{H} , there is surely a vector $v_N \in \mathcal{H}$ with the above properties. \square

The occurrence described by Lemma 2.3 may happen both with an orthogonal and with an oblique projection scheme, namely both when $P_N = Q_N$ and when $P_N \neq Q_N$ eventually in N . In the standard framework of (Petrov-)Galerkin methods such an occurrence is prevented by suitable assumptions on A , a typical example being coercivity [8, Sect. 2.2], [19, Sect. 4.1].

As in our discussion we do not exclude a priori such an occurrence. We are compelled to regard $f^{(N)}$ as an approximate solution to the truncated problem, in the sense that

$$(2.7) \quad A_N f^{(N)} = g_N + \varepsilon^{(N)} \quad \text{for some } \varepsilon^{(N)} \in \mathbb{C}^N.$$

(We write $\varepsilon^{(N)}$ and not ε_N because there is no reason to claim that the residual $\varepsilon^{(N)}$ in the N -dimensional problem is the actual truncation for every N of the same infinite-dimensional vector $\varepsilon \in \mathcal{H}$.)

It would be desirable to assume that $\varepsilon^{(N)}$ is indeed small and asymptotically vanishing with N , or even that $\varepsilon^{(N)} = 0$ for N large enough, as is case in some applications. Morally (up to passing to the weak formulation of the inverse problem), this is the assumption of *asymptotic consistency* naturally made for approximations by Galerkin methods [8, Definition 2.15 and Theorem 2.24]. We shall make this assumption here too, observing that in the present abstract context it is motivated by the following property, whose proof is postponed to Section 4.

Lemma 2.4. *Let $A \in \mathcal{B}(\mathcal{H})$ and $g \in \text{ran} A$. Let A_N and g_N be defined as in (2.3)-(2.4) above. Then there always exists a sequence $(f^{(N)})_{N \in \mathbb{N}}$ such that*

$$f^{(N)} \in \mathbb{C}^N \quad \text{and} \quad \lim_{N \rightarrow \infty} \|A_N f^{(N)} - g_N\|_{\mathbb{C}^N} = 0.$$

In other words, there do exist approximate solutions $f^{(N)}$ to (2.5) actually satisfying (2.7) with $\|\varepsilon^{(N)}\|_{\mathbb{C}^N} \rightarrow 0$ as $N \rightarrow \infty$.

2.3. Convergence of the truncated problem: error and residual.

For an infinite-dimensional inverse problem the other major question is the vanishing, as $N \rightarrow \infty$, of the two natural indicators of the displacement between the infinite-dimensional inverse linear problem and its finite-dimensional truncation, namely the *infinite-dimensional error* \mathcal{E}_N and the *infinite-dimensional residual* \mathfrak{R}_N , defined respectively as

$$(2.8) \quad \begin{aligned} \mathcal{E}_N &:= f - \widehat{f^{(N)}} \\ \mathfrak{R}_N &:= g - A \widehat{f^{(N)}}. \end{aligned}$$

We qualify them as ‘infinite-dimensional’, although we shall drop this extra nomenclature when no confusion arises, in order to distinguish them from the error and residual at fixed N , which may be indexed by the number of steps in an iterative algorithm.

A first evident obstruction to the actual vanishing of \mathcal{E}_N when $\dim \mathcal{H} = \infty$ is the use of a *non-complete* orthonormal system $(u_n)_{n \in \mathbb{N}}$, that is, such that $\text{span}\{u_n \mid n \in \mathbb{N}\}$ is not dense in \mathcal{H} .

Example 2.5. If the weighted (compact) right-shift operator \mathcal{R} (Sect. A.3) is truncated with respect to

$$(u_n)_{n \in \mathbb{N}} = (e_n)_{n \in \mathbb{N}, n \geq 2}, \quad (v_n)_{n \in \mathbb{N}} = (e_n)_{n \in \mathbb{N}}$$

and the initial inverse problem is $\mathcal{R}f = g = e_2$, then the exact solution is $f = \frac{1}{\sigma_1} e_1$, yet the truncated problem can only produce approximate solutions

$$\widehat{f^{(N)}} \in \text{span}\{e_2, e_3, \dots\},$$

whence $\widehat{f^{(N)}} \perp f$ and $\|\widehat{f^{(N)}} - f\|_{\mathcal{H}} \geq \frac{1}{\sigma_1}$.

Truncations with respect to a potentially non-complete orthonormal system might appear unwise, but in certain contexts are natural. One is the vast framework of the Krylov subspace methods [16], where one searches for approximate solutions among the linear combinations of the vectors g, Ag, A^2g, \dots and hence to perform the truncation with respect to an orthonormal basis of the *Krylov subspace*

$$(2.9) \quad \mathcal{K}(A, g) := \text{span}\{A^k g \mid k \in \mathbb{N}_0\}$$

associated to $A \in \mathcal{B}(\mathcal{H})$ and $g \in \mathcal{H}$. Obviously, when $\dim \mathcal{K}(A, g) = \infty$ the subspace $\mathcal{K}(A, g)$ is open in \mathcal{H} . Its closure can be the whole \mathcal{H} , but also just a *proper* closed subspace of \mathcal{H} .

Example 2.6.

- (i) For the right-shift operator R on $\ell^2(\mathbb{N})$ (Sect. A.2) and the vector $g = e_{m+1}$ (one of the canonical basis vectors), $\overline{\mathcal{K}(R, e_{m+1})} = \text{span}\{e_1, \dots, e_m\}^\perp$, which is a proper subspace of $\ell^2(\mathbb{N})$ if $m \geq 1$, and instead is the whole $\ell^2(\mathbb{N})$ if $g = e_1$. Therefore the exact solution $f = e_m$ to $Rf = g$ does *not* belong to $\overline{\mathcal{K}(R, e_{m+1})}$.
- (ii) For the Volterra integral operator V on $L^2[0, 1]$ (Sect. A.5) and the function $g = \mathbf{1}$ (the constant function with value 1), it follows from (A.10) or (A.15) that the functions Vg, V^2g, V^3g, \dots are (multiples of) the polynomials x, x^2, x^3, \dots , therefore $\mathcal{K}(V, g)$ is the space of polynomials on $[0, 1]$, which is dense in $L^2[0, 1]$.

Thus, in Example 2.5 above the system $(u_n)_{n \in \mathbb{N}} = (e_n)_{n \in \mathbb{N}, n \geq 2}$ spans the Krylov subspace relative to \mathcal{R} and e_2 .

In standard (Petrov-)Galerkin methods an occurrence as in Example 2.5 or (2.6)(i) is ruled out by an ad hoc ‘*approximability*’ assumption [8, Definition 2.14

and Theorem 2.24] that can be rephrased as the request that $(u_n)_{n \in \mathbb{N}}$ is indeed an orthonormal *basis* of \mathcal{H} .

The approximability property is known to fail in situations of engineering interest, as is the case for the failure of the Lagrange finite elements in differential problems for electromagnetism [8, Sect. 2.3.3].

Even when (complete) orthonormal bases of \mathcal{H} are employed for the truncation, another feature of the infinite dimensionality must be taken into account, namely the possibility that error and residual are asymptotically small only in some weaker sense than the customary norm topology of \mathcal{H} .

There are indeed at least three meaningful senses in which the vanishing of \mathcal{E}_N or \mathfrak{R}_N , as $N \rightarrow \infty$, can be monitored in an informative way.

I. Strong (\mathcal{H} -norm) convergence. This is the vanishing $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$, resp., $\|\mathcal{E}_N\|_{\mathcal{H}} \rightarrow 0$ of the residual, resp., or the error. Obviously,

$$(2.10) \quad \|\mathfrak{R}_N\|_{\mathcal{H}} \leq \|A\|_{\text{op}} \|\mathcal{E}_N\|_{\mathcal{H}}.$$

II. Weak convergence. This is the vanishing $\mathfrak{R}_N \rightarrow 0$ or $\mathcal{E}_N \rightarrow 0$: recall that a sequence $(\xi_N)_{N \in \mathbb{N}}$ in \mathcal{H} converges weakly to $\xi \in \mathcal{H}$ as $N \rightarrow \infty$, $\xi_N \rightharpoonup \xi$, when $\langle \eta, \xi_N \rangle \rightarrow \langle \eta, \xi \rangle$ for any $\eta \in \mathcal{H}$.

III. Component-wise convergence. This is the vanishing of each component of the vector \mathfrak{R}_N or \mathcal{E}_N with respect to the considered basis. Recall that a sequence $(\xi_N)_{N \in \mathbb{N}}$ in \mathcal{H} converges component-wise to $\xi \in \mathcal{H}$ as $N \rightarrow \infty$ with respect to the orthonormal basis $(e_n)_{n \in \mathbb{N}}$ of \mathcal{H} , and we write $\xi_N \rightsquigarrow \xi$, when $\langle e_n, \eta_N \rangle \xrightarrow{N \rightarrow \infty} \langle e_n, \eta \rangle \forall n \in \mathbb{N}$. Thus, $\mathcal{E}_N \rightsquigarrow 0$ means that each n -th component $\langle u_n, f - \widehat{f^{(N)}} \rangle$ of \mathcal{E}_N vanishes as $N \rightarrow \infty$ and $\mathfrak{R}_N \rightsquigarrow 0$ means that each n -th component $\langle v_n, g - A\widehat{f^{(N)}} \rangle$ of \mathfrak{R}_N vanishes as $N \rightarrow \infty$, possibly with different vanishing rate depending on n .

Clearly,

$$(2.11) \quad \text{strong} \quad \Rightarrow \quad \text{weak} \quad \Rightarrow \quad \text{component-wise},$$

and these notions are all inequivalent when $\dim \mathcal{H} = \infty$ (whereas they are all equivalent when $\dim \mathcal{H} < \infty$). In fact, it is standard to check that

$$(2.12) \quad \eta_N \xrightarrow{\|\cdot\|_{\mathcal{H}}} \eta \quad \text{as } N \rightarrow \infty \quad \Leftrightarrow \quad \begin{cases} \eta_N \rightarrow \eta \\ \|\eta_N\|_{\mathcal{H}} \rightarrow \|\eta\|_{\mathcal{H}}, \end{cases}$$

and

$$(2.13) \quad \eta_N \rightharpoonup \eta \quad \text{as } N \rightarrow \infty \quad \Leftrightarrow \quad \begin{cases} \langle e_n, \eta_N \rangle \xrightarrow{N \rightarrow \infty} \langle e_n, \eta \rangle \quad \forall n \in \mathbb{N} \\ \sup_{N \in \mathbb{N}} \|\eta_N\|_{\mathcal{H}} < +\infty, \end{cases}$$

where $(e_n)_{n \in \mathbb{N}}$ is an orthonormal basis of \mathcal{H} .

Despite (2.11), a mere component-wise vanishing $\mathcal{E}_N \rightsquigarrow 0$ is in many respects already satisfactorily informative, for in this case each component of $\widehat{f^{(N)}}$ (with respect to the basis $(u_n)_{n \in \mathbb{N}}$) approximates the corresponding component of the exact solution f .

As a matter of fact, a strong control such as $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$ or $\|\mathcal{E}_N\|_{\mathcal{H}} \rightarrow 0$ is *not* generic and only holds under specific a priori conditions on the inverse linear problem.

Thus, as already recalled in the Introduction, for elliptic boundary value problems the standard Galerkin finite element method produces a *strong* vanishing of the error, provided that two crucial conditions are satisfied, namely a *careful* choice of the truncation space and the *coercivity* of the differential operator [19, Sect. 4.2.3]: when this is the case, the vanishing rate depends on the truncation basis and the regularity of the solution. More generally [8, Sect. 2.3.1], standard

Petrov-Galerkin methods give rise to a strong convergence of the approximate solution under the simultaneous validity of uniform stability, uniform boundedness and asymptotic consistency of the linear problem, and approximability by means of the chosen truncation spaces. When the differential operator is non-coercive, additional sufficient conditions have been studied for the stability of the truncated problem and for the quasi optimality of the discretization scheme [5, 4, 3].

On a related scenario, special classes of linear ill-conditioned problems (rank-deficient and discrete ill-posed problems) can be treated with regularisation methods in which the solution is stabilised [22, 13]. The most notable regularisation methods, namely the Tikhonov-Phillips method, the Landweber-Fridman iteration method, and the truncated singular value decomposition, produce indeed a strongly vanishing error [11, 17]. Yet, when the inverse linear problem $Af = g$ is governed by an infinite-rank compact operator A , it can be seen that the conjugate gradient method, as well as α -processes (in particular, the method of steepest descent) may have strongly divergent error and residual in the presence of noise [6] and one is forced to consider weaker forms of convergence. In fact, in [6] the presence of component-wise convergence is also alluded to.

3. THE COMPACT LINEAR INVERSE PROBLEM

Let us now examine, within the framework elaborated in the previous Section, the abstract truncation and convergence scheme for compact linear inverse problems.

When the operator A on the given Hilbert space \mathcal{H} is compact, it admits a ‘canonical’ decomposition, the ‘singular value decomposition’ [20, Theorem VI.17]

$$(3.1) \quad A = \sum_n \sigma_n |\psi_n\rangle\langle\varphi_n|,$$

where n runs in a finite or infinite subset of \mathbb{N} , $(\varphi_n)_n$ and $(\psi_n)_n$ are two orthonormal systems of \mathcal{H} , and $0 < \sigma_{n+1} < \sigma_n$ for all n , and the above series converges in operator norm. In the following we shall reserve the above notation for the singular value decomposition of the considered compact operator.

The injectivity of A is tantamount as $(\varphi_n)_{n \in \mathbb{N}}$ being an orthonormal basis. A is not necessarily surjective, but $\text{ran} A = \mathcal{H}$ if and only if $(\psi_n)_{n \in \mathbb{N}}$ is an orthonormal basis.

The inverse problem (1.1) for compact *and injective* A and $g \in \text{ran} A$ is well-defined: there exists a unique $f \in \mathcal{H}$ such that $Af = g$.

The compactness of A has two noticeable consequences here. First, since $\dim \mathcal{H} = \infty$, A is invertible on its range only, and cannot have an everywhere defined bounded inverse: $\text{ran} A$ can be dense in \mathcal{H} , as in the case of the Volterra operator on $L^2[0, 1]$ (Sect. A.5), or also dense in a closed proper subspace of \mathcal{H} , as for the weighted right-shift on $\ell^2(\mathbb{N})$ (Sect. A.3).

Furthermore, A and its compression (in the usual meaning of Sect. 2.1) are close in a robust sense, as the following standard Lemma shows.

Lemma 3.1. *With respect to an infinite-dimensional separable Hilbert space \mathcal{H} , let $A : \mathcal{H} \rightarrow \mathcal{H}$ be a compact operator and let $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$ be two orthonormal bases of \mathcal{H} . Then*

$$(3.2) \quad \|A - Q_N A P_N\|_{\text{op}} \xrightarrow{N \rightarrow \infty} 0,$$

P_N and Q_N being as usual the orthogonal projections (2.1).

Proof. Upon splitting

$$A - Q_N A P_N = (A - Q_N A) + Q_N (A - A P_N)$$

it suffices to prove that $\|A - AP_N\|_{\text{op}} \xrightarrow{N \rightarrow \infty} 0$ and $\|A - Q_N A\|_{\text{op}} \xrightarrow{N \rightarrow \infty} 0$. Let us prove the first limit (the second being completely analogous).

Clearly, it is enough to prove that $\|A - AP_N\|_{\text{op}}$ vanishes assuming further that A has finite rank. Indeed, the difference $(A - AP_N) - (\tilde{A} - \tilde{A}P_N)$, where \tilde{A} is a finite-rank approximant of the compact operator A , is controlled in operator norm by $2\|A - \tilde{A}\|_{\text{op}}$ and hence can be made arbitrarily small.

Thus, we consider non-restrictively $A = \sum_{k=1}^M \sigma_k |\psi_k\rangle\langle\varphi_k|$ for some integer M , where $(\varphi_k)_{k=1}^M$ and $(\psi_k)_{k=1}^M$ are two orthonormal systems, and $0 < \sigma_M < \dots < \sigma_1$. Now, for a generic $\xi = \sum_{n=1}^{\infty} \xi_n v_n \in \mathcal{H}$ one has

$$\begin{aligned} \|(A - AP_N)\xi\|_{\mathcal{H}}^2 &= \left\| \sum_{k=1}^M \sigma_k \left(\sum_{n=N+1}^{\infty} \xi_n \langle\varphi_k, v_n\rangle \right) \psi_k \right\|_{\mathcal{H}}^2 \\ &= \sum_{k=1}^M \sigma_k^2 \left| \sum_{n=N+1}^{\infty} \xi_n \langle\varphi_k, v_n\rangle \right|^2 \leq \|\xi\|_{\mathcal{H}}^2 \sum_{k=1}^M \sigma_k^2 \|(\mathbb{1} - P_N)\varphi_k\|_{\mathcal{H}}^2, \end{aligned}$$

therefore

$$\|A - AP_N\|_{\text{op}}^2 \leq M \sigma_1^2 \cdot \max_{k \in \{1, \dots, M\}} \|(\mathbb{1} - P_N)\varphi_k\|_{\mathcal{H}}^2 \xrightarrow{N \rightarrow \infty} 0,$$

since the above maximum is taken over M (hence, finitely many) quantities, each of which vanishes as $N \rightarrow \infty$. \square

In the following Theorem we describe the generic behaviour of well-defined compact inverse problem.

Theorem 3.2. *Consider*

- the linear inverse problem $Af = g$ in a separable Hilbert space \mathcal{H} for some compact and injective $A : \mathcal{H} \rightarrow \mathcal{H}$ and some $g \in \text{ran}A$;
- the finite-dimensional truncation A_N obtained by compression with respect to the orthonormal bases $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$ of \mathcal{H} .

Let $(f^{(N)})_{N \in \mathbb{N}}$ be a sequence of approximate solutions to the truncated problems in the quantitative sense

$$A_N f^{(N)} = g_N + \varepsilon^{(N)}, \quad f^{(N)}, \varepsilon^{(N)} \in \mathbb{C}^N, \quad \|\varepsilon^{(N)}\|_{\mathbb{C}^N} \xrightarrow{N \rightarrow \infty} 0$$

for every (sufficiently large) N . If $\widehat{f^{(N)}}$ is \mathcal{H} -norm bounded uniformly in N , then

$$\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0 \quad \text{and} \quad \mathcal{E}_N \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

Proof. We split

$$\begin{aligned} Af^{(N)} - g &= (A - Q_N AP_N) \widehat{f^{(N)}} \\ (*) \quad &+ Q_N AP_N \widehat{f^{(N)}} - Q_N g \\ &+ Q_N g - g. \end{aligned}$$

By assumption, $\|Q_N g - g\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0$ and

$$\begin{aligned} \|Q_N AP_N \widehat{f^{(N)}} - Q_N g\|_{\mathcal{H}} &= \|A_N f^{(N)} - g_N\|_{\mathbb{C}^N} \\ &= \|\varepsilon^{(N)}\|_{\mathbb{C}^N} \xrightarrow{N \rightarrow \infty} 0. \end{aligned}$$

Moreover, Lemma 3.1 and the uniform boundedness of $\widehat{f^{(N)}}$ imply

$$\|(A - Q_N AP_N) \widehat{f^{(N)}}\|_{\mathcal{H}} \leq \|A - Q_N AP_N\|_{\text{op}} \|\widehat{f^{(N)}}\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0$$

Plugging the three limits above into (*) proves $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$.

Next, in terms of the singular value decomposition (3.1) of A , where now $(\varphi_n)_{n \in \mathbb{N}}$ is an orthonormal basis of \mathcal{H} , $(\psi_n)_{n \in \mathbb{N}}$ is an orthonormal system, and $0 < \sigma_{n+1} < \sigma_n \forall n \in \mathbb{N}$, we write

$$\widehat{f^{(N)}} = \sum_{n \in \mathbb{N}} f_n^{(N)} \varphi_n, \quad \widehat{f} = \sum_{n \in \mathbb{N}} f_n \varphi_n,$$

whence

$$0 = \lim_{N \rightarrow \infty} \|A\widehat{f^{(N)}} - g\|_{\mathcal{H}}^2 = \lim_{N \rightarrow \infty} \sum_{n \in \mathbb{N}} \sigma_n^2 |f_n^{(N)} - f_n|^2.$$

Then necessarily $\widehat{f^{(N)}}$ converges to f component-wise ($\mathcal{E}_N \rightsquigarrow 0$).

On the other hand, $\widehat{f^{(N)}}$ is uniformly bounded in \mathcal{H} , thus, owing to (2.13), $\widehat{f^{(N)}}$ converges to f weakly ($\mathcal{E}_N \rightarrow 0$). \square

Theorem 3.2 provides sufficient conditions for some form of vanishing of the error and the residual. The key assumptions are:

- *injectivity* of A ,
- *asymptotic solvability of the truncated problems*, i.e., asymptotic smallness of the finite-dimensional residual $A_N f^{(N)} - g_N$,
- *uniform boundedness of the approximate solutions* $f^{(N)}$.

In fact, injectivity was only used in the analysis of the error in order to conclude $\mathcal{E}_N \rightarrow 0$; instead, the conclusion $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$ follows irrespectively of injectivity.

To further understand the impact of such assumptions, a few remarks are in order.

Remark 3.3 (Genericity). Under the conditions of Theorem 3.2, the occurrence of the *strong* vanishing of the residual ($\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$) and the *weak* vanishing of the error ($\mathcal{E}_N \rightarrow 0$) as $N \rightarrow \infty$ is a *generic behaviour*. For example, the compact inverse problem $\mathcal{R}f = 0$ in $\ell^2(\mathbb{N})$ associated with the weighted right-shift \mathcal{R} (Sect. A.3) has exact solution $f = 0$. The truncated problem $\mathcal{R}_N f^{(N)} = 0$ with respect to the same basis $(e_n)_{n \in \mathbb{N}}$, \mathcal{R}_N being the matrix (2.6), is solved by the \mathbb{C}^N -vectors whose first $N-1$ components are zero, i.e., $\widehat{f^{(N)}} = e_N$. The sequence $(\widehat{f^{(N)}})_{N \in \mathbb{N}} \equiv (e_N)_{N \in \mathbb{N}}$ converges weakly to zero in $\ell^2(\mathbb{N})$, whence indeed $\mathcal{E}_N \rightarrow 0$, and also, by compactness, $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$. However, $\|\mathcal{E}_N\|_{\mathcal{H}} = 1$ for every N , thus the error cannot vanish in the \mathcal{H} -norm.

Remark 3.4 ('Bad' approximate solutions). The example considered in Remark 3.3 is also instructive to understand that generically one may happen to select 'bad' approximate solutions $\widehat{f^{(N)}}$ such that, despite the 'good' property $\|A_N \widehat{f^{(N)}} - g_N\|_{\mathbb{C}^N} \rightarrow 0$, have the unsatisfactory feature $\|f^{(N)}\|_{\mathbb{C}^N} = \|\widehat{f^{(N)}}\|_{\mathcal{H}} \rightarrow +\infty$: this is the case if one chooses, for instance, $\widehat{f^{(N)}} = N e_N$. Thus, the uniform boundedness of $\widehat{f^{(N)}}$ in \mathcal{H} required in Theorem 3.2 is *not* redundant. (This also shows, in view of the proof of Theorem 3.2, that whereas by compactness $\widehat{f^{(N)}} \rightharpoonup f$ implies $\|A\widehat{f^{(N)}} - Af\| \rightarrow 0$, the opposite implication is not true in general.)

Remark 3.5 (The density of $\text{ran} A$ does not help). Even if the genericity discussed in Remarks 3.3 and 3.4 is referred to compact injective operators with non-dense range, requiring $\overline{\text{ran} A} = \mathcal{H}$ does not improve the convergence in general. For instance, the compact inverse problem associated with the weighted right-shift \mathcal{R} in $\ell^2(\mathbb{Z})$ (Sect. A.4) involves an operator that is compact, injective, and with dense range, but its compression with $Q_N := P_N := \sum_{n=-N}^N |e_N\rangle\langle e_N|$ produces for every N a $(2N+1) \times (2N+1)$ square matrix that is singular and for which, therefore, all the considerations of Remarks 3.3 and 3.4 can be repeated verbatim.

Remark 3.6 (‘Bad’ truncations and ‘good’ truncations). We saw in Lemma 2.3 that ‘bad’ truncations (i.e., leading to matrices A_N that are, eventually in N , all singular) are always possible. On the other hand, there always exists a “good” choice for the truncation – although such a choice might not be identifiable explicitly – which makes the infinite-dimensional residual and error vanish in a stronger sense than what stated in Theorem 3.2, and without the extra assumption of uniform boundedness on the approximate solutions. For instance, in terms of the singular value decomposition (3.1) of A , it is enough to choose

$$(u_n)_{n \in \mathbb{N}} = (\varphi_n)_{n \in \mathbb{N}}, \quad (v_n)_{n \in \mathbb{N}} = (\psi_n)_{n \in \mathbb{N}},$$

in which case $Q_N A P_N = \sum_{n=1}^N \sigma_n |\psi_n\rangle\langle\varphi_n|$ and $A_N = \text{diag}(\sigma_1, \dots, \sigma_N)$, and for given $g = \sum_{n \in \mathbb{N}} g_n \psi_n$ one has $\widehat{f^{(N)}} = \sum_{n=1}^N \frac{g_n}{\sigma_n} \varphi_n$, where the sequence $(\frac{g_n}{\sigma_n})_{n \in \mathbb{N}}$ belongs to $\ell^2(\mathbb{N})$ owing to the assumption $g \in \text{ran} A$, whence

$$\|f - \widehat{f^{(N)}}\|_{\mathcal{H}}^2 = \sum_{n=N+1}^{\infty} \left| \frac{g_n}{\sigma_n} \right|^2 \xrightarrow{N \rightarrow \infty} 0.$$

4. THE BOUNDED LINEAR INVERSE PROBLEM

It is instructive to compare the findings of the previous Section with the more general case of a bounded (not necessarily compact) inverse linear problem.

When $\dim \mathcal{H} = \infty$ and a generic bounded linear operator $A : \mathcal{H} \rightarrow \mathcal{H}$ is compressed (in the usual sense of Sect. 2) between the spans of the first N vectors of the orthonormal bases $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$, then surely $Q_N A P_N \rightarrow A$ as $N \rightarrow \infty$ in the strong operator topology, that is, $\|Q_N A P_N \psi - A \psi\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0 \forall \psi \in \mathcal{H}$, yet the convergence *may fail to occur in the operator norm*.

The first statement is an obvious consequence of the inequality

$$\|(A - Q_N A P_N) \psi\|_{\mathcal{H}} \leq \|(\mathbb{1} - Q_N) A \psi\|_{\mathcal{H}} + \|A\|_{\text{op}} \|\psi - P_N \psi\|_{\mathcal{H}}$$

valid for any $\psi \in \mathcal{H}$. The lack of operator norm convergence is clear, for instance, when one compresses the identity operator (or any bounded, non-compact operator): the operator norm limit of finite-rank operators can only be compact.

For this reason, the control of the infinite-dimensional inverse problem in terms of its finite-dimensional truncated versions is in general less strong.

As a counterpart of Theorem 3.2 above, let us discuss the following generic behaviour of *well-posed* bounded linear inverse problems.

Theorem 4.1. *Consider*

- *the linear inverse problem $Af = g$ in a Hilbert space \mathcal{H} for some bounded and injective $A : \mathcal{H} \rightarrow \mathcal{H}$ and some $g \in \mathcal{H}$;*
- *the finite-dimensional truncation A_N obtained by compression with respect to the orthonormal bases $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$ of \mathcal{H} .*

Let $(f^{(N)})_{N \in \mathbb{N}}$ be a sequence of approximate solutions to the truncated problems in the quantitative sense

$$A_N f^{(N)} = g_N + \varepsilon^{(N)}, \quad f^{(N)}, \varepsilon^{(N)} \in \mathbb{C}^N, \quad \|\varepsilon^{(N)}\|_{\mathbb{C}^N} \xrightarrow{N \rightarrow \infty} 0$$

for every (sufficiently large) N . Assume further that $\widehat{f^{(N)}}$ converges strongly in \mathcal{H} , equivalently, that $\|f^{(N)} - f^{(M)}\|_{\mathbb{C}^{\max\{N, M\}}} \xrightarrow{N, M \rightarrow \infty} 0$. Then

$$\|\mathcal{E}_N\|_{\mathcal{H}} \rightarrow 0 \quad \text{and} \quad \|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

Proof. Since

$$\begin{aligned}
(*) \quad Af^{(N)} - g &= (A - Q_N AP_N) \widehat{f^{(N)}} \\
&+ Q_N AP_N \widehat{f^{(N)}} - Q_N g \\
&+ Q_N g - g,
\end{aligned}$$

and since by assumption $\|Q_N g - g\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0$ and

$$\begin{aligned}
\|Q_N AP_N \widehat{f^{(N)}} - Q_N g\|_{\mathcal{H}} &= \|A_N f^{(N)} - g_N\|_{\mathbb{C}^N} \\
&= \|\varepsilon^{(N)}\|_{\mathbb{C}^N} \xrightarrow{N \rightarrow \infty} 0,
\end{aligned}$$

then the strong vanishing of $Af^{(N)} - g$ is tantamount as the strong vanishing of $(A - Q_N AP_N) \widehat{f^{(N)}}$.

Since in addition $\|f^{(N)} - \tilde{f}\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0$ for some $\tilde{f} \in \mathcal{H}$, then

$$\begin{aligned}
\|(A - Q_N AP_N) \widehat{f^{(N)}}\|_{\mathcal{H}} &\leq \|(A - Q_N AP_N) \tilde{f}\|_{\mathcal{H}} + 2\|A\|_{\text{op}} \|\tilde{f} - f^{(N)}\|_{\mathcal{H}} \\
&\xrightarrow{N \rightarrow \infty} 0
\end{aligned}$$

(the first summand in the r.h.s. above vanishing due to the operator strong convergence $Q_N AP_N \rightarrow A$), and (*) thus implies $\|\mathfrak{R}_N\|_{\mathcal{H}} = \|Af^{(N)} - g\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0$.

Moreover, $Af^{(N)} \rightarrow g$ (as proved right now) and $Af^{(N)} \rightarrow A\tilde{f}$ (by continuity), whence $A\tilde{f} = g = Af$ and also (by injectivity) $f = \tilde{f}$. This shows that $\|\mathcal{E}_N\|_{\mathcal{H}} = \|f - f^{(N)}\|_{\mathcal{H}} = \|\tilde{f} - f^{(N)}\|_{\mathcal{H}} \rightarrow 0$. \square

We observe that also here injectivity was only used in the analysis of the error, whereas it is not needed to conclude that $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$.

As compared to Theorem 3.2, Theorem 4.1 now relies on the following hypotheses:

- *injectivity* of A ,
- *asymptotic solvability of the truncated problems*,
- *convergence of the approximate solutions* $f^{(N)}$.

The first two assumptions are the same as in the compact case: the first guarantees the existence of a unique solution and the second is a natural working hypothesis, by virtue of Lemma 2.4. Under such assumptions, we thus see that, in passing from a (well-defined) compact to a generic (well-defined) bounded inverse problem, one has to strengthen the hypothesis of uniform boundedness of the $f^{(N)}$'s to their actual strong convergence, in order for the residual \mathfrak{R}_N to vanish strongly (in which case, as a by-product, also the error \mathcal{E}_N vanishes strongly).

Moreover, the proof of Theorem 4.1 shows that, under injectivity of A and asymptotic solvability of the truncated problems, the residual \mathfrak{R}_N vanishes strongly, or weakly or component-wise, if and only if so does $(A - Q_N AP_N) \widehat{f^{(N)}}$. In the compact case, $A - Q_N AP_N \rightarrow \mathbb{O}$ in operator norm (Lemma 3.1), and it suffices that the $f^{(N)}$'s are uniformly bounded (or, in principle, have increasing norm $\|f^{(N)}\|_{\mathcal{H}}$ compensated by the vanishing of $\|A - Q_N AP_N\|_{\text{op}}$), in order for $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$. In the general bounded case we controlled the vanishing of $\|(A - Q_N AP_N) \widehat{f^{(N)}}\|_{\mathcal{H}}$ by requiring additionally that the $f^{(N)}$'s converge strongly.

If instead the sequence of the $f^{(N)}$'s *does not* converge strongly, Theorem 4.1 is not applicable, and in general one has to expect only weak vanishing of the residual, $\mathfrak{R}_N \rightharpoonup 0$, which in turn prevents the error to vanish strongly – for otherwise

$\|\mathcal{E}_N\|_{\mathcal{H}} \rightarrow 0$ would imply $\|\mathfrak{R}_N\|_{\mathcal{H}} \rightarrow 0$, owing to (2.10). The following example shows such a possibility.

Example 4.2. For the right-shift R on $\ell^2(\mathbb{N})$ (Sect. A.2), an actual injective operator, the inverse problem $Rf = g$ with $g = 0$ admits the unique solution $f = 0$. The truncated finite-dimensional problems induced by the bases $(u_n)_{n \in \mathbb{N}} = (v_n)_{n \in \mathbb{N}} = (e_n)_{n \in \mathbb{N}}$, where $(e_n)_{n \in \mathbb{N}}$ is the canonical basis of $\ell^2(\mathbb{N})$, is governed by the sub-diagonal matrix

$$R_N = \begin{pmatrix} 0 & \cdots & \cdots & \cdots & 0 \\ 1 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & \cdots & 1 & 0 \end{pmatrix}.$$

Let us consider the sequence $(\widehat{f^{(N)}})_{N \in \mathbb{N}}$ with $\widehat{f^{(N)}} := e_N$ for each N . Then:

- $R_N f^{(N)} = 0 = g_N$ (the truncated problems are solved exactly),
- $\widehat{f^{(N)}} \rightharpoonup 0$ (only weakly, not strongly),
- $\mathfrak{R}_N = g - R\widehat{f^{(N)}} = -e_{N+1} \rightharpoonup 0$ (only weakly, not strongly).

Of course, what discussed so far emphasizes features of generic bounded inverse problems (as compared to compact ones). Ad hoc analyses for special classes of bounded inverse problems are available and complement the picture of Theorem 4.1. This is the case, to mention one example, when A is an *algebraic operator*, namely $p(A) = \mathbb{D}$ for some polynomial p (which includes finite-rank A 's) and one treats the inverse problem with the generalised minimal residual method (GMRES) [10].

In retrospect, the arguments developed in this Section allow us to prove Lemma 2.4.

Proof of Lemma 2.4. Let f solve $Af = g$. The sequence $(f^{(N)})_{N \in \mathbb{N}}$ defined by

$$f^{(N)} := (P_N f)_N = f_N \quad (\text{that is, } \widehat{f^{(N)}} = P_N f)$$

does the job, and that is a straightforward consequence of the fact that, as argued already at the beginning of this Section, $Q_N A P_N \rightarrow A$ strongly in the operator topology. Indeed, one has by adding and subtracting Af

$$\begin{aligned} \|A_N f^{(N)} - g_N\|_{\mathbb{C}^N} &= \|Q_N A P_N \widehat{f^{(N)}} - Q_N g\|_{\mathcal{H}} \\ &\leq \|(Q_N A P_N - A)f\|_{\mathcal{H}} + \|(1 - Q_N)Af\|_{\mathcal{H}}. \end{aligned}$$

The strong limit yields the conclusion. \square

5. COMPARISON TO CONJUGATE GRADIENT SCHEMES

In this Section we further discuss the scope of Theorems 3.2 (compact case) and 4.1 (bounded case) in application to conjugate gradient schemes for bounded, self-adjoint, positive semi-definite inverse linear problems [7, Chapt. 7], [8, Sect. 9.3.2], [19, Sect. 7.2.2]. Thus, throughout this Section $A = A^* \in \mathcal{B}(\mathcal{H})$ with $\langle h, Ah \rangle \geq 0 \forall h \in \mathcal{H}$.

In particular, we provide an additional insight on the key role played by the assumption of *uniform boundedness* (or even *strong convergence*) of the finite-dimensional approximants.

Under the above assumption on A , and with $g \in \text{ran} A$, the problem $Af = g$ admits solution(s) in \mathcal{H} , which form the (non-empty) manifold

$$(5.1) \quad \mathcal{S}(A, g) := \{f \in \mathcal{H} \mid Af = g\}.$$

Clearly, if A is injective, which in this case amounts to A being positive definite, then $\mathcal{S}(A, g)$ only consists of the unique solution to the inverse problem. Moreover, any f in the solution manifold $\mathcal{S}(A, g)$ can be variationally characterised as

$$(5.2) \quad \Phi[f] = \min_{h \in \mathcal{H}} \Phi[h], \quad \Phi[h] := \langle h, Ah \rangle - 2\langle h, g \rangle,$$

that is, f is the minimiser of the functional $\Phi[h]$ (which, in specific contexts, is referred to as the ‘energy functional’ of the problem).

Based on such properties, in the framework of conjugate gradient schemes one builds a sequence $(f^{[N]})_{N \in \mathbb{N}_0}$, the so-called ‘conjugate gradient iterates’, by taking $f^{[0]}$ to be an arbitrary vector in \mathcal{H} , and $f^{[N]}$, for $N \geq 1$, to be the minimiser of the problem

$$(5.3) \quad \min_{h \in \mathcal{Q}_N} \Phi[h], \quad \begin{aligned} \mathcal{Q}_N &:= \{f^{[0]}\} + \text{span}\{r_0, Ar_0, \dots, A^{N-1}r_0\} \\ r_0 &:= Af^{[0]} - g. \end{aligned}$$

Here ‘iterates’ refers to the fact that the $f^{[N]}$ ’s can be equivalently obtained by means of certain iterative procedures [14, 18].

The notation for the superscript in $f^{[N]}$ is chosen to avoid confusion with the special meaning already reserved to $f^{(N)}$ and $\widehat{f^{(N)}}$ in the general setting of Sect. 2.1, although it is clear that the $f^{[N]}$ ’s here are to be considered on the same conceptual footing as the $\widehat{f^{(N)}}$ ’s, that is, they can be naturally regarded as *approximate* solutions, expected to satisfy $Af^{[N]} \approx g$ in a suitable sense. This is suggested by the very construction (5.3) and the variational characterisation (5.2) of the solution(s) f .

That the above expectation is correct is expressed in rigorous terms by Theorem 5.1 below, a classical result by Nemirovskiy and Polyak [18] (with a precursor version by Kammerer and Nashed [15]), and discussed in more recent terms in [7, Sect. 7.2] and [12, Sect. 3.2]. In order to state it, let us introduce the map $\mathcal{P}_{\mathcal{S}} : \mathcal{H} \rightarrow \mathcal{S}(A, g)$ that associates to a point $h \in \mathcal{H}$ the nearest point $\mathcal{P}_{\mathcal{S}}h$ of the solution manifold. Then one has the following.

Theorem 5.1. (Nemirovskiy and Polyak [18, Theorem 7].)

Let $A = A^* \in \mathcal{B}(\mathcal{H})$ with $\langle h, Ah \rangle \geq 0 \ \forall h \in \mathcal{H}$, and let the sequence $(f^{[N]})_{N \in \mathbb{N}_0}$ in \mathcal{H} be defined by (5.3) above. Then

$$(5.4) \quad \lim_{N \rightarrow \infty} \|f^{[N]} - \mathcal{P}_{\mathcal{S}}f^{[N]}\|_{\mathcal{H}} = 0,$$

and moreover, for every $\gamma > 0$,

$$(5.5) \quad \|f^{[N]} - \mathcal{P}_{\mathcal{S}}f^{[N]}\|_{\mathcal{H}} \leq \left(\frac{C_{f^{[0]}, \gamma}}{2N + 1} \right)^{\gamma}$$

for some constant $C_{f^{[0]}, \gamma} > 0$ depending on $f^{[0]}$ and γ , provided that the problem $A^{\gamma/2}u = f^{[0]} - \mathcal{P}_{\mathcal{S}}f^{[0]}$ admits a solution $u \in \mathcal{H}$.

When A is injective and hence $\mathcal{S}(A, g)$ only consists of the unique solution f to $Af = g$, (5.4) reads $\|f^{[N]} - f\|_{\mathcal{H}} \rightarrow 0$ as $N \rightarrow \infty$. In the analogy with the analysis of Theorem 4.1, the sequence of approximate solutions is convergent and the error \mathcal{E}_N indeed vanishes strongly, and so does, necessarily, the residual \mathfrak{R}_N . We can thus understand Theorem 5.1 in view of our Theorem 4.1.

Equally instructive is the case when A is not injective and hence the solution manifold $\mathcal{S}(A, g)$ contains infinitely many vectors. Again, (5.4) indicates that the approximate solutions $f^{[N]}$ ’s are asymptotically close, in the \mathcal{H} -norm topology, to solutions of the considered inverse problem. However, now this does not necessarily imply the actual convergence to a *fixed solution*: both the $f^{[N]}$ ’s and the corresponding $\mathcal{P}_{\mathcal{S}}f^{[N]}$ ’s might in principle have arbitrarily large norm – in complete

analogy to what one would have in Example 4.2 if one considered approximate solutions $\widehat{f^{(N)}} = Ne_N$, instead of just $\widehat{f^{(N)}} = e_N$. In order to deduce from (5.4) that $f^{[N]} \rightarrow f$ for some solution f , an additional information is needed, for example the property that the $f^{[N]}$'s are uniformly norm bounded. This sheds further light on the requirement of strong convergence of the approximate solutions made in Theorem 4.1 needed to deduce the strong vanishing of the error.

6. COUNTERPART REMARKS ON LINEAR INVERSE PROBLEMS WITH NOISE

Let us reconsider the typical occurrence, mentioned in the Introduction, when

- within the modelling of the phenomenon under investigation, the linear inverse problem $Af = g$ is well-defined (or even well-posed), and thus, there is a unique ‘input’ f for given ‘output’ g and with an explicitly known law $f \xrightarrow{A} g$;
- however, the knowledge of g obtained from measurements is disturbed by various forms of uncertainty.

In view of the general discussion developed so far, we can make here a few remarks on such an occurrence.

Now the problem $Af = g$ cannot be studied directly, and instead one deals with the inverse problem

$$(6.1) \quad Af \approx \tilde{g}$$

in the new unknown \tilde{f} for some given (measured) $\tilde{g} := g + \nu \in \mathcal{H}$, where the ‘noise’ vector ν is present albeit not known explicitly, but is typically small – for instance a small bound on $\|\nu\|_{\mathcal{H}}$ may be known a priori.

If ν (and g) belongs to $\text{ran}A$, so does \tilde{g} , and there exist an actual (possibly non-unique) solution \tilde{f} to (6.1). Theorems 3.2 and 4.1 are then applicable, replacing g with $g + \nu$, and with analogous notation we may speak of an approximate solution $f^{(N)} \in \mathbb{C}^N$ such that

$$(6.2) \quad A_N f^{(N)} = g_N + \nu_N + \varepsilon^{(N)}, \quad \|\varepsilon^{(N)}\|_{\mathbb{C}^N} \xrightarrow{N \rightarrow \infty} 0.$$

This way, Theorems 3.2 and 4.1 produce a control on the ‘residual with noise’ $(g + \nu) - Af^{(N)}$ and on the ‘error with noise’ $\tilde{f} - f^{(N)}$. This only determines the ‘solution with noise’, namely \tilde{f} , and not the exact solution f , but that can be still informative if ν is sufficiently small. For example, if A is bounded and with everywhere defined bounded inverse, then $\tilde{f} = A^{-1}(g + \nu)$, whence $\|\tilde{f} - f\|_{\mathcal{H}} \leq \|A^{-1}\|_{\text{op}} \|\nu\|_{\mathcal{H}}$, and the smallness of $\|\nu\|_{\mathcal{H}}$, in terms of $\|A^{-1}\|_{\text{op}}$, provides an estimate on how close f and \tilde{f} are.

If, on the other hand, $\nu \notin \text{ran}A$, then the problem with noise loses solvability: there is no exact solution to (6.1) and one can only think of an approximate solution \tilde{f} satisfying $A\tilde{f} \approx \tilde{g}$ in some sense (whence also $A\tilde{f} \approx g$, since ν is conveniently small).

Let us comment on the typical behaviour of the residual \mathfrak{R}_N and the error \mathcal{E}_N associated with f , $\widehat{f^{(N)}}$, g , for simplicity in the case where A is compact and injective, with $g \in \text{ran}A$ (thus, with f unique solution to $Af = g$).

6.1. Typical behaviour of \mathfrak{R}_N with noise.

When the truncated problem with noise is solved in the approximate sense (6.2), and the $\widehat{f^{(N)}}$'s are uniformly bounded in \mathcal{H} , then necessarily

$$(6.3) \quad \|\mathfrak{R}_N\|_{\mathcal{H}} = \|A\widehat{f^{(N)}} - g\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} \|\nu\|_{\mathcal{H}}.$$

This is seen by splitting as usual

$$\mathfrak{R}_N = (Q_N A P_N - A) \widehat{f^{(N)}} + (Q_N g - Q_N A P_N \widehat{f^{(N)}}) + (g - Q_N g),$$

and observing that $\|(Q_N A P_N - A) \widehat{f^{(N)}}\|_{\mathcal{H}} \leq \|Q_N A P_N - A\|_{\text{op}} \|\widehat{f^{(N)}}\|_{\mathcal{H}} \rightarrow 0$ (owing to Lemma 3.1), $\|g - Q_N g\|_{\mathcal{H}} \rightarrow 0$, and

$$\|Q_N g - Q_N A P_N \widehat{f^{(N)}}\|_{\mathcal{H}} = \|A_N f^{(N)} - g_N\|_{\mathbb{C}^N} = \|\nu_N + \varepsilon^{(N)}\|_{\mathbb{C}^N} \rightarrow \|\nu\|_{\mathcal{H}}.$$

Clearly, based on the above argument, one actually has

$$(6.4) \quad \|\mathfrak{R}_N - \nu\|_{\mathcal{H}} \xrightarrow{N \rightarrow \infty} 0,$$

which is in fact stronger than (6.3). Thus, ‘*the residual vanishes up to the noise threshold*’.

6.2. Typical behaviour of \mathcal{E}_N with noise.

In the presence of noise one cannot expect that $\widehat{f^{(N)}}$, even just component-wise, converges to f ; in particular, the possibility that $\|\mathcal{E}_N\|_{\mathcal{H}} \rightarrow 0$ or $\mathcal{E}_N \rightarrow 0$ would violate (6.4).

Thus, $\|\mathcal{E}_N\|_{\mathcal{H}}$ stays strictly above zero, uniformly in N , in fact with a typical behaviour that $\|\mathcal{E}_N\|_{\mathcal{H}}$ *initially decreases for not to large N , reaches a minimum, then for larger N eventually increases, possibly blowing up*. (This differs from the behaviour of $\|\mathfrak{R}_N\|_{\mathcal{H}}$, which typically decreases monotonically to $\|\nu\|_{\mathcal{H}}$.) The minimum for $\|\mathcal{E}_N\|_{\mathcal{H}}$, say, when $N = N_0$, provides the best approximant of f in \mathcal{H} , namely $\widehat{f^{(N_0)}}$.

For concreteness, let us consider the case in which the Petrov-Galerkin projection to (6.2) is performed with the same bases $(\varphi_n)_{n \in \mathbb{N}}$ and $(\psi_n)_{n \in \mathbb{N}}$ of the canonical singular value decomposition (3.1) of A . Let us also assume that $\nu \in \text{ran} A$ (the generalisation of what follows to the case $\nu \notin \text{ran} A$ is straightforward). These simplifications guarantee that for all N the matrix $A_N = \text{diag}(\sigma_1, \dots, \sigma_N)$ is non-singular on \mathbb{C}^N , because now $Q_N A P_N = \sum_{n=1}^N \sigma_n |\psi_n\rangle \langle \varphi_n|$, and that (6.2) is exactly solved by

$$\widehat{f^{(N)}} = \sum_{n=1}^N \frac{g_n + \nu_n}{\sigma_n} \varphi_n,$$

having decomposed

$$\nu = \sum_{n=1}^{\infty} \nu_n \psi_n, \quad g = \sum_{n=1}^{\infty} g_n \psi_n, \quad f = \sum_{n=1}^{\infty} f_n \varphi_n, \quad g_n = \sigma_n f_n.$$

Thus, $A_N \widehat{f^{(N)}} = g_N + \nu_N$ ($\varepsilon^{(N)} = 0$). Then

$$\begin{aligned} \|\mathfrak{R}_N\|_{\mathcal{H}}^2 &= \|g - A \widehat{f^{(N)}}\|_{\mathcal{H}}^2 = \sum_{n=1}^N |\nu_n|^2 + \sum_{n=N+1}^{\infty} |g_n|^2 \xrightarrow{N \rightarrow \infty} \|\nu\|_{\mathcal{H}}^2, \\ \|\mathcal{E}_N\|_{\mathcal{H}}^2 &= \|f - \widehat{f^{(N)}}\|_{\mathcal{H}}^2 = \sum_{n=1}^N \frac{|\nu_n|^2}{\sigma_n^2} + \sum_{n=N+1}^{\infty} |f_n|^2 =: \alpha(N) + \beta(N). \end{aligned}$$

It is clear that $\beta(N)$ decreases monotonically to zero as $N \rightarrow \infty$, whereas $\alpha(N)$ is monotone increasing with N . This can produce the typical initial decrease of $\|\mathcal{E}_N\|_{\mathcal{H}}$, driven by a substantial decrease of $\beta(N)$ as opposite to a mild increase of $\alpha(N)$, which is the case when f is mainly supported on low modes φ_n 's and ν instead has a substantial tail on high modes ψ_n 's. For N sufficiently large, $\alpha(N)$ then becomes leading, which would produce the typical inversion of the curve of $\|\mathcal{E}_N\|_{\mathcal{H}}$ versus N . Having assumed $\nu \in \text{ran} A$, necessarily $\alpha(N) \rightarrow \|A^{-1} \nu\|_{\mathcal{H}}^2$, thus

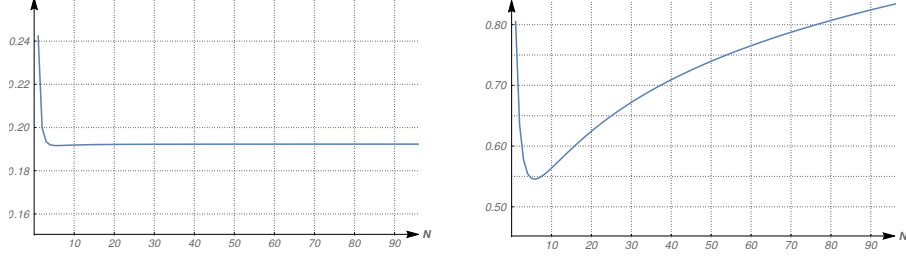


FIGURE 1. Typical behaviour of the residual $\|\mathfrak{R}_N\|_{\mathcal{H}}^2$ (left) and of the error $\|\mathcal{E}_N\|_{\mathcal{H}}^2$ (right) for increasing size of the finite-dimensional truncation, relative to the problem $Af = g$ considered in Example 6.1, with the choice $\sigma_n = \frac{1}{n}$, $g_n = \frac{1}{n^2}$, $\nu_n = \frac{0.4}{n^{3/2}}$.

with no blow-up of $\|\mathcal{E}_N\|_{\mathcal{H}}$. Reasoning as above with $\nu \notin A$ one would conclude instead that the series defining $\alpha(N)$ diverges.

Example 6.1. Take, $\forall n \in \mathbb{N}$,

$$\sigma_n = n^{-1}, \quad g_n = n^{-2}, \quad \nu_n = n^{-\frac{3}{2}}.$$

Thus, A is an injective Hilbert-Schmidt operator, $\|\nu\|_{\mathcal{H}}^2 = \zeta(3) \simeq 1.20$ (where $\zeta(x)$ denotes the Riemann zeta function), and $\nu \notin \text{ran}A$. Then $f_n = n^{-1}$, $\|f\|_{\mathcal{H}}^2 = \beta(0) = \frac{\pi^2}{6}$, and

$$\beta(N) \leq (N+1)^{-2} \rightarrow 0, \quad \alpha(N) \sim \ln N \rightarrow +\infty.$$

Figure 1 displays the behaviour of residual and error in this case.

7. NUMERICAL TESTS: EFFECTS OF CHANGING THE TRUNCATION BASIS

In this final Section we examine some of the features discussed theoretically so far through a few numerical tests concerning different choices of the truncation bases. We employed a Legendre, complex Fourier, and a Krylov basis to truncate the problems.

The two model operators that we considered are the Volterra operator V in $L^2[0, 1]$ (Sect. A.5) and the self-adjoint multiplication operator $M : L^2[1, 2] \rightarrow L^2[1, 2]$, $\psi \mapsto x\psi$. We examined the following two inverse problems.

First problem: $Vf_1 = g_1$, with $g_1(x) = \frac{1}{2}x^2$.

The problem has unique solution

$$(7.1) \quad f_1(x) = x, \quad \|f_1\|_{L^2[0,1]} = \frac{1}{\sqrt{3}} \simeq 0.5774$$

and f_1 is a Krylov solution, i.e., $f_1 \in \overline{\mathcal{K}(V, g)}$, although $f_1 \notin \mathcal{K}(V, g)$. To prove the first fact, let us observe that $\mathcal{K}(V, g)$ is spanned by the monomials x^2, x^3, x^4, \dots , i.e., $\mathcal{K}(V, g) = \{x^2 p \mid p \text{ is a polynomial on } [0, 1]\}$; therefore, if $h \in \mathcal{K}(V, g)^\perp$, then $0 = \int_0^1 \overline{h(x)} x^2 p(x) dx$ for any polynomial p ; the L^2 -density of polynomials on $[0, 1]$ implies necessarily that $x^2 h = 0$, whence also $h = 0$; this proves that $\mathcal{K}(V, g)^\perp = \{0\}$ and hence $\overline{\mathcal{K}(V, g)} = L^2[0, 1]$. The fact that $f_1 \notin \mathcal{K}(V, g)$ follows from $f(x) = x^2 \cdot \frac{1}{x}$ and $\frac{1}{x} \notin L^2[0, 1]$.

Second problem: $Mf_2 = g_2$, with $g_2(x) = x^2$.

The problem has unique solution

$$(7.2) \quad f_2(x) = x, \quad \|f_2\|_{L^2[1,2]} = \sqrt{\frac{7}{3}} \simeq 1.5275$$

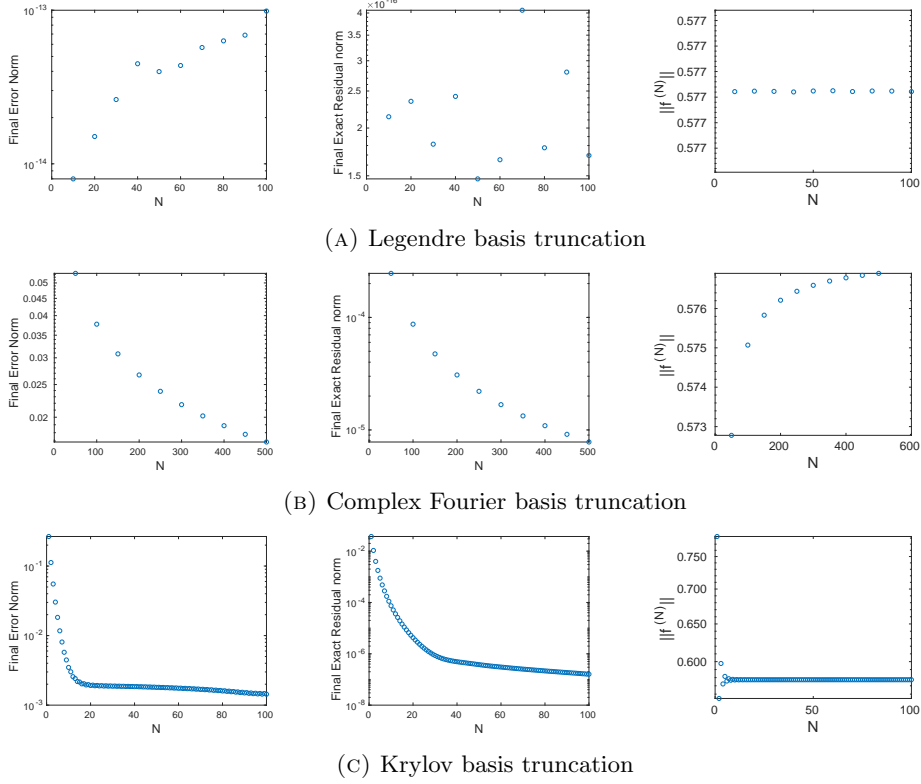


FIGURE 2. Norm of the infinite-dimensional error and residual, and of the approximated solution for the Volterra inverse problem truncated with the Legendre, complex Fourier, and Krylov bases.

and f_2 is a Krylov solution. Indeed, $\mathcal{K}(M, g) = \{x^2 p \mid p \text{ is a polynomial on } [1, 2]\}$ and $\overline{\mathcal{K}(M, g)} = \{x^2 h(x) \mid h \in L^2[1, 2]\} = L^2[1, 2]$, whence $f_2 \in \overline{\mathcal{K}(M, g)}$ and $f_2 \notin \mathcal{K}(M, g)$.

We treated both problems with three different orthonormal bases: the Legendre polynomials and the complex Fourier modes (on the intervals $[0, 1]$ or $[1, 2]$, depending on the problem) solved using the QR factorisation algorithm, and the Krylov basis generated using the GMRES algorithm.

Computationally speaking, generating accurate representations of the Legendre polynomials is quite demanding and accuracy can be lost rather soon due to their highly oscillatory nature, particularly at the end points. For this reason we limited our investigation up to $N = 100$ when considering the Legendre basis, but $N = 500$ when considering the complex Fourier basis. It is expected that there is no significant numerical error from the computation of the Legendre basis, as the $L^2[0, 1]$ and $L^2[1, 2]$ norms of the basis polynomials have less than 1% error compared to their exact unit value.

For each problem and each choice of the basis, we monitored the norm of the infinite-dimensional error $\|\mathcal{E}_N\|_{L^2} = \|f - \widehat{f^{(N)}}\|_{L^2}$ ($f = f_1$ or f_2), of the infinite-dimensional residual $\|\mathfrak{R}_N\|_{L^2} = \|g - A \widehat{f^{(N)}}\|_{L^2}$ ($g = g_1$ or g_2 ; $A = V$ or M), and of the approximated solution $\|\widehat{f^{(N)}}\|_{L^2} = \|f^{(N)}\|_{\mathbb{C}^N}$.

Figures 2 and 4 highlight the difference between the computation in the three bases for the Volterra operator.

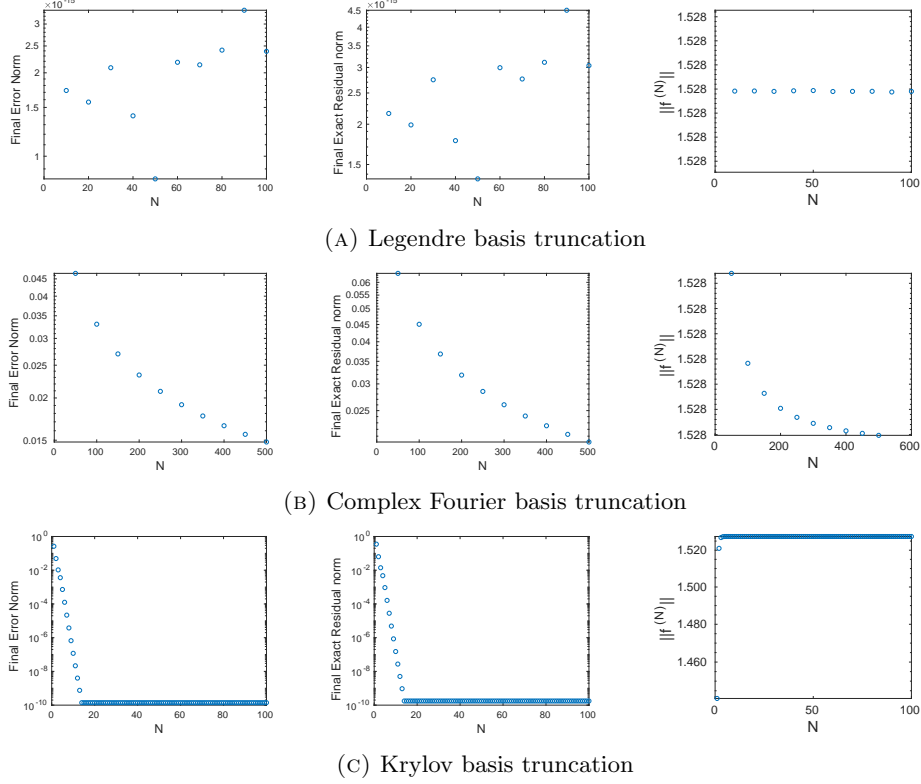


FIGURE 3. Norm of the infinite-dimensional error, residual, and approximated solution for the M -multiplication inverse problem truncated with the Legendre, complex Fourier, and Krylov bases.

- In the Legendre basis, $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ are almost zero. $\|\widehat{f^{(N)}}\|_{L^2}$ stays bounded and constant with N and matches the expected value (7.1). The approximated solutions reconstruct the exact solution f_1 at any truncation number.
- In the complex Fourier basis, both $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ are some orders of magnitude *larger* than in the Legendre basis and decrease monotonically with N ; in fact, $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ display an evident convergence to zero, however attaining values that are more than ten orders of magnitude larger than the corresponding error and residual norms for the same N in the Legendre case. $\|\widehat{f^{(N)}}\|_{L^2}$, on the other hand, increases monotonically and appears to approach the theoretical value (7.1). These quite stringent differences in the error and residual may be attributable to the Gibbs phenomenon. In fact, reconstructing f_1 using the Krylov approximated solutions produces a vector that shows a highly oscillatory behaviour near the end points, confirming the presence of the Gibbs phenomenon.
- In the Krylov basis $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ decrease monotonically, relatively fast for small N 's, then rather slowly with N . Such quantities are smaller than in the Fourier basis. $\|\widehat{f^{(N)}}\|_{L^2}$ displays some initial highly oscillatory behaviour, but quickly approaches the theoretical value (7.1). On the other hand, the reconstruction appears to be quite good with some noticeable oscillations at the end points.

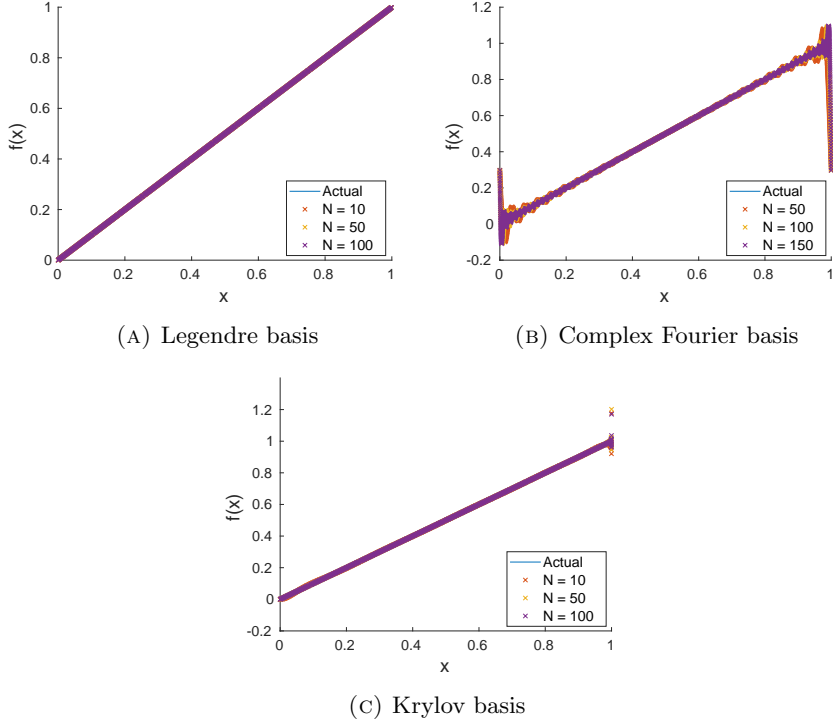


FIGURE 4. Reconstruction of the exact solution $f_1(x) = x$ from the solutions for the problem $V f_1 = g_1$. The Fourier basis produces an inaccurate reconstruction due to high oscillations, resulting in higher errors.

Thus, among the considered truncations the Legendre basis yields the most accurate reconstruction and the complex Fourier basis yields the least accurate reconstruction of the exact solution.

In contrast, Figures 3 and 5 highlight the difference between the computation in the three bases for the M -multiplication operator.

- In the Legendre basis, $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ are again almost zero. $\|\widehat{f^{(N)}}\|_{L^2}$ is constant with N at the expected value (7.2). The approximated solutions reconstruct the exact solution f_2 at any truncation number.
- In the Fourier basis the behaviour of the above indicators is again qualitatively the same, and again with a much milder convergence rate in N to the asymptotic values as compared with the Legendre case. $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ still display an evident convergence to zero. Again the higher error compared to the Legendre case is likely due to the nature of the approximation of the exact solution f_2 by oscillatory functions and the Gibbs phenomenon.
- The Krylov basis displays a fast initial decrease of both $\|\mathcal{E}_N\|_{L^2}$ and $\|\mathfrak{R}_N\|_{L^2}$ to the tolerance level of 10^{-10} that was set for the residual. $\|\widehat{f^{(N)}}\|_{L^2}$ also increases rapidly and remains constant at the expected value (7.2). The reconstruction of the solution is excellent, but still not quite as good as the Legendre case.

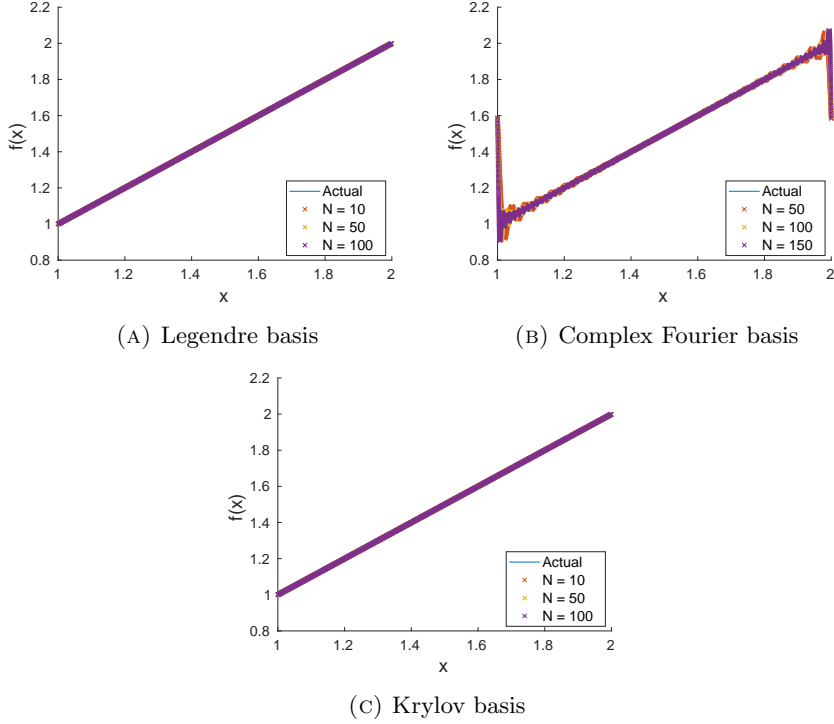


FIGURE 5. Reconstruction of the exact solution $f_2(x) = x$ from the solutions for the problem $Mf_2 = g_2$.

All this gives numerical evidence that the choice of the truncation basis *does* affect the sequence of solutions. The Legendre basis is best suited to these problems as f_1 , f_2 , g_1 and g_2 are perfectly representable by the first few basis vectors.

APPENDIX A. SOME PROTOTYPICAL EXAMPLE OPERATORS

Let us describe in this Appendix a few operators in Hilbert space that were useful in the course of our discussion, both as a source of examples or counter-examples, and as a playground to understand certain mechanisms typical of the infinite dimensionality.

A.1. The multiplication operator on $\ell^2(\mathbb{N})$.

Let us denote with $(e_n)_{n \in \mathbb{N}}$ the canonical orthonormal basis of $\ell^2(\mathbb{N})$. For a given bounded sequence $a \equiv (a_n)_{n \in \mathbb{N}}$ in \mathbb{C} , the multiplication by a is the operator $M^{(a)} : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$ defined by $M^{(a)}e_n = a_n e_n \forall n \in \mathbb{N}$ and then extended by linearity and density, in other words the operator given by the series

$$(A.1) \quad M^{(a)} = \sum_{n=1}^{\infty} a_n |e_n\rangle \langle e_n|$$

(that converges strongly in the operator sense).

$M^{(a)}$ is bounded with norm $\|M^{(a)}\|_{\text{op}} = \sup_n |a_n|$ and spectrum $\sigma(M^{(a)})$ given by the closure in \mathbb{C} of the set $\{a_1, a_2, a_3 \dots\}$. Its adjoint is the multiplication by a^* . Thus, $M^{(a)}$ is normal. $M^{(a)}$ is self-adjoint whenever a is real and it is compact if $\lim_{n \rightarrow \infty} a_n = 0$.

A.2. The right-shift operator on $\ell^2(\mathbb{N})$.

The operator $R : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$ defined by $Re_n = e_{n+1} \forall n \in \mathbb{N}$ and then extended by linearity and density, in other words the operator given by the series

$$(A.2) \quad R = \sum_{n=1}^{\infty} |e_{n+1}\rangle\langle e_n|$$

(that converges strongly in the operator sense), is called the right-shift operator.

R is an isometry (i.e., it is norm-preserving) with closed range $\text{ran}R = \{e_1\}^\perp$. In particular, it is bounded with $\|R\|_{\text{op}} = 1$, yet not compact, it is injective, and invertible on its range, with bounded inverse

$$(A.3) \quad R^{-1} : \text{ran}R \rightarrow \mathcal{H}, \quad R^{-1} = \sum_{n=1}^{\infty} |e_n\rangle\langle e_{n+1}|.$$

The adjoint of R on \mathcal{H} is the so-called left-shift operator, namely the everywhere defined and bounded operator $L : \mathcal{H} \rightarrow \mathcal{H}$ defined by the (strongly convergent, in the operator sense) series

$$(A.4) \quad L = \sum_{n=1}^{\infty} |e_n\rangle\langle e_{n+1}|, \quad L = R^*.$$

Thus, L inverts R on $\text{ran}R$, i.e., $LR = \mathbb{1}$, yet $RL = \mathbb{1} - |e_1\rangle\langle e_1|$. One has $\ker R^* = \text{span}\{e_1\}$.

R and L have the same spectrum $\sigma(R) = \sigma(L) = \{z \in \mathbb{C} \mid |z| \leq 1\}$, but R has no eigenvalue, whereas the eigenvalue of L form the open unit ball $\{z \in \mathbb{C} \mid |z| < 1\}$.

A.3. The compact (weighted) right-shift operator on $\ell^2(\mathbb{N})$.

This is the operator $\mathcal{R} : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$ defined by the operator-norm convergent series

$$(A.5) \quad \mathcal{R} = \sum_{n=1}^{\infty} \sigma_n |e_{n+1}\rangle\langle e_n|,$$

where $\sigma \equiv (\sigma_n)_{n \in \mathbb{N}}$ is a given bounded sequence with $0 < \sigma_{n+1} < \sigma_n \forall n \in \mathbb{N}$ and $\lim_{n \rightarrow \infty} \sigma_n = 0$. Thus, $\mathcal{R}e_n = \sigma_n e_{n+1}$.

\mathcal{R} is injective and compact, and (A.5) is its singular value decomposition, with norm $\|\mathcal{R}\|_{\text{op}} = \sigma_1$, $\overline{\text{ran}}\mathcal{R} = \{e_1\}^\perp$, and adjoint

$$(A.6) \quad \mathcal{R}^* = \mathcal{L} = \sum_{n=1}^{\infty} \sigma_n |e_n\rangle\langle e_{n+1}|.$$

Thus, $\mathcal{L}\mathcal{R} = M^{(\sigma^2)}$, the operator of multiplication by $(\sigma_n^2)_{n \in \mathbb{N}}$, whereas $\mathcal{R}\mathcal{L} = M^{(\sigma^2)} - \sigma_1^2 |e_1\rangle\langle e_1|$.

A.4. The compact (weighted) right-shift operator on $\ell^2(\mathbb{Z})$.

This is the operator $\mathcal{R} : \ell^2(\mathbb{Z}) \rightarrow \ell^2(\mathbb{Z})$ defined by the operator-norm convergent series

$$(A.7) \quad \mathcal{R} = \sum_{n \in \mathbb{Z}} \sigma_{|n|} |e_{n+1}\rangle\langle e_n|,$$

where $\sigma \equiv (\sigma_n)_{n \in \mathbb{N}_0}$ is a given bounded sequence with $0 < \sigma_{n+1} < \sigma_n \forall n \in \mathbb{N}_0$ and $\lim_{n \rightarrow \infty} \sigma_n = 0$. Thus, $\mathcal{R}e_n = \sigma_{|n|} e_{n+1}$.

\mathcal{R} is injective and compact, with $\text{ran}\mathcal{R}$ dense in \mathcal{H} and norm $\|\mathcal{R}\|_{\text{op}} = \sigma_0$. (A.7) gives the singular value decomposition. The adjoint of \mathcal{R} is

$$(A.8) \quad \mathcal{R}^* = \mathcal{L} = \sum_{n \in \mathbb{Z}} \sigma_{|n|} |e_n\rangle\langle e_{n+1}|.$$

Thus, $\mathcal{LR} = M^{(\sigma^2)} = \mathcal{RL}$.

The ‘inverse of \mathcal{R} on its range’ is the densely defined, surjective, unbounded operator $\mathcal{R}^{-1} : \text{ran } \mathcal{R} \rightarrow \mathcal{H}$ acting as

$$(A.9) \quad \mathcal{R}^{-1} = \sum_{n \in \mathbb{Z}} \frac{1}{\sigma_{|n|}} |e_n\rangle \langle e_{n+1}|$$

as a series that converges on $\text{ran } \mathcal{R}$ in the strong operator sense.

A.5. The Volterra operator on $L^2[0, 1]$.

This is the operator $V : L^2[0, 1] \rightarrow L^2[0, 1]$ defined by

$$(A.10) \quad (Vf)(x) = \int_0^x f(y) dy, \quad x \in [0, 1].$$

V is compact and injective with spectrum $\sigma(V) = \{0\}$ (thus, the spectral point 0 is not an eigenvalue) and norm $\|V\|_{\text{op}} = \frac{2}{\pi}$. Its adjoint V^* acts as

$$(A.11) \quad (V^*f)(x) = \int_x^1 f(y) dy, \quad x \in [0, 1],$$

therefore $V + V^*$ is the rank-one orthogonal projection

$$(A.12) \quad V + V^* = |\mathbf{1}\rangle \langle \mathbf{1}|$$

onto the function $\mathbf{1}(x) = 1$.

The singular value decomposition of V is

$$(A.13) \quad V = \sum_{n=0}^{\infty} \sigma_n |\psi_n\rangle \langle \varphi_n|, \quad \begin{aligned} \sigma_n &= \frac{2}{(2n+1)\pi} \\ \varphi_n(x) &= \sqrt{2} \cos \frac{(2n+1)\pi}{2} x \\ \psi_n(x) &= \sqrt{2} \sin \frac{(2n+1)\pi}{2} x, \end{aligned}$$

where both $(\varphi_n)_{n \in \mathbb{N}_0}$ and $(\psi_n)_{n \in \mathbb{N}_0}$ are orthonormal bases of $L^2[0, 1]$.

Thus, $\text{ran } V$ is dense, but strictly contained in \mathcal{H} : for example, $\mathbf{1} \notin \text{ran } V$.

In fact, V is invertible on its range, but does not have (everywhere defined) bounded inverse; yet $V - z\mathbb{1}$ does, for any $z \in \mathbb{C} \setminus \{0\}$ (recall that $\sigma(V) = \{0\}$), and

$$(A.14) \quad (z\mathbb{1} - V)^{-1}\psi = z^{-1}\psi + z^{-2} \int_0^x e^{\frac{x-y}{z}} \psi(y) dy \quad \forall \psi \in \mathcal{H}, z \in \mathbb{C} \setminus \{0\}.$$

The explicit action of the powers of V is

$$(A.15) \quad (V^n f)(x) = \frac{1}{(n-1)!} \int_0^x (x-y)^{n-1} f(y) dy, \quad n \in \mathbb{N}.$$

A.6. The multiplication operator on an annulus in $L^2(\Omega)$.

This is the operator $M_z : L^2(\Omega_r) \rightarrow L^2(\Omega_r)$, $f \mapsto zf$, where

$$(A.16) \quad \Omega_r := \{z \in \mathbb{C} \mid r < |z| < 1\}, \quad r \in (0, 1).$$

M_z is a normal bounded bijection with norm $\|M_z\|_{\text{op}} = 1$, spectrum $\sigma(M_z) = \overline{\Omega}_r$, and adjoint given by $M_z^* f = \bar{z}f$.

REFERENCES

- [1] N. ANTONIĆ, K. S. BURAZIN, I. CRNJAC, AND M. ERCEG, *Complex Friedrichs systems and applications*, J. Math. Phys., 58 (2017), pp. 101508, 22.
- [2] N. ANTONIĆ, M. ERCEG, AND A. MICHELANGELI, *Friedrichs systems in a Hilbert space framework: Solvability and multiplicity*, J. Differential Equations, 263 (2017), pp. 8264–8294.
- [3] A. BUFFA, *Remarks on the discretization of some noncoercive operator with applications to heterogeneous Maxwell equations*, SIAM J. Numer. Anal., 43 (2005), pp. 1–18.
- [4] A. BUFFA AND S. H. CHRISTIANSEN, *The electric field integral equation on Lipschitz screens: definitions and numerical approximation*, Numer. Math., 94 (2003), pp. 229–267.

- [5] S. CAORSI, P. FERNANDES, AND M. RAFFETTO, *On the convergence of Galerkin finite element approximations of electromagnetic eigenproblems*, SIAM J. Numer. Anal., 38 (2000), pp. 580–607.
- [6] B. EICKE, A. K. LOUIS, AND R. PLATO, *The instability of some gradient methods for ill-posed problems*, Numer. Math., 58 (1990), pp. 129–134.
- [7] H. W. ENGL, M. HANKE, AND A. NEUBAUER, *Regularization of inverse problems*, vol. 375 of Mathematics and its Applications, Kluwer Academic Publishers Group, Dordrecht, 1996.
- [8] A. ERN AND J.-L. GUERMOND, *Theory and practice of finite elements*, vol. 159 of Applied Mathematical Sciences, Springer-Verlag, New York, 2004.
- [9] A. ERN, J.-L. GUERMOND, AND G. CAPLAIN, *An intrinsic criterion for the bijectivity of Hilbert operators related to Friedrichs' systems*, Comm. Partial Differential Equations, 32 (2007), pp. 317–341.
- [10] M. G. GASPARO, A. PAPINI, AND A. PASQUALI, *Some properties of GMRES in Hilbert spaces*, Numer. Funct. Anal. Optim., 29 (2008), pp. 1276–1285.
- [11] C. W. GROETSCH, *The theory of Tikhonov regularization for Fredholm equations of the first kind*, vol. 105 of Research Notes in Mathematics, Pitman (Advanced Publishing Program), Boston, MA, 1984.
- [12] M. HANKE, *Conjugate gradient type methods for ill-posed problems*, vol. 327 of Pitman Research Notes in Mathematics Series, Longman Scientific & Technical, Harlow, 1995.
- [13] P. C. HANSEN, *Rank-deficient and discrete ill-posed problems*, SIAM Monographs on Mathematical Modeling and Computation, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1998. Numerical aspects of linear inversion.
- [14] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*, J. Research Nat. Bur. Standards, 49 (1952), pp. 409–436 (1953).
- [15] W. J. KAMMERER AND M. Z. NASHED, *On the convergence of the conjugate gradient method for singular linear operator equations*, SIAM J. Numer. Anal., 9 (1972), pp. 165–181.
- [16] J. LIESEN AND Z. E. STRAKOŠ, *Krylov subspace methods*, Numerical Mathematics and Scientific Computation, Oxford University Press, Oxford, 2013. Principles and analysis.
- [17] A. K. LOUIS, *Inverse und schlecht gestellte Probleme*, Teubner Studienbücher Mathematik. [Teubner Mathematical Textbooks], B. G. Teubner, Stuttgart, 1989.
- [18] A. S. NEMIROVSKIY AND B. T. POLYAK, *Iterative methods for solving linear ill-posed problems under precise information. I*, Izv. Akad. Nauk SSSR Tekhn. Kibernet., (1984), pp. 13–25, 203.
- [19] A. QUARTERONI, *Numerical models for differential problems*, vol. 16 of MS&A. Modeling, Simulation and Applications, Springer, Cham, 2017. Third edition.
- [20] M. REED AND B. SIMON, *Methods of Modern Mathematical Physics*, vol. 1, New York Academic Press, 1972.
- [21] K. SCHMÜDGEN, *Unbounded self-adjoint operators on Hilbert space*, vol. 265 of Graduate Texts in Mathematics, Springer, Dordrecht, 2012.
- [22] A. N. TIKHONOV AND V. Y. ARSENIN, *Solutions of ill-posed problems*, V. H. Winston & Sons, Washington, D.C.: John Wiley & Sons, New York-Toronto, Ont.-London, 1977. Translated from the Russian, Preface by translation editor Fritz John, Scripta Series in Mathematics.

(N. Caruso) INTERNATIONAL SCHOOL FOR ADVANCED STUDIES – SISSA, VIA BONOMEA 265, 34136 TRIESTE (ITALY).

Email address: ncaruso@sissa.it

(A. Michelangeli) INTERNATIONAL SCHOOL FOR ADVANCED STUDIES – SISSA, VIA BONOMEA 265, 34136 TRIESTE (ITALY).

Email address: alemiche@sissa.it

(P. Novati) UNIVERSITÀ DEGLI STUDI DI TRIESTE, PIAZZALE EUROPA 1, 34127 TRIESTE (ITALY).

Email address: novati@units.it