# Chapter 1
# Inverse Problems in Statistics

Laurent Cavalier

**Abstract** There exist many fields where inverse problems appear. Some examples are: astronomy (blurred images of the Hubble satellite), econometrics (instrumental variables), financial mathematics (model calibration of the volatility), medical image processing (X-ray tomography), and quantum physics (quantum homodyne tomography).

These are problems where we have indirect observations of an object (a function) that we want to reconstruct, through a linear operator $A$. Due to its indirect nature, solving an inverse problem is usually rather difficult.

For this reason, one needs regularization methods in order to get a stable and accurate reconstruction.

We present the framework of statistical inverse problems where the data are corrupted by some stochastic error. This white noise model may be discretized in the spectral domain using Singular Value Decomposition (SVD), when the operator $A$ is compact. Several examples of inverse problems where the SVD is known are presented (circular deconvolution, heat equation, tomography).

We explain some basic issues regarding nonparametric statistics applied to inverse problems. Standard regularization methods and their counterpart as estimation procedures by use of SVD are discussed (projection, Landweber, Tikhonov, ...). Several classical statistical approaches like minimax risk and optimal rates of convergence, are presented. This notion of optimality leads to some optimal choice of the tuning parameter.

However these optimal parameters are unachievable since they depend on the unknown smoothness of the function. This leads to more recent concepts like adaptive estimation and oracle inequalities. Several data-driven selection procedures of the regularization parameter are discussed in details, among these: model selection methods, Stein's unbiased risk estimation and the recent risk hull method.

Laurent Cavalier

Université Aix-Marseille 1, LATP, CMI, 39 rue Joliot-Curie, 13453 Marseille, France, e-mail: cavalier@cmi.univ-mrs.fr

# Preface

These notes are based on a mini-course which was given during the summer school *Stats in the Château* in August 2009. The first version of these notes was written for a course at Heidelberg University in 2007. Another course was given at *Ecole d' été en statistique* in Switzerland in September 2010. A longer version of the course is given to the graduate students at Université de Provence in Marseille.

I would like to thank the colleagues and students who attended these courses and asked questions, made comments and remarks.

Since these notes were written in several places, I would also like to thank Heidelberg University, Göttingen University, Sydney University, University College London and Princeton University.

Many thanks, for very helpful discussions, to Yuri Golubev, Markus Reiss and a special thank to Thorsten Hohage for giving to me his lecture notes.

The three referees also helped a lot, with their very interesting remarks and comments, in improving these notes.

I would like to dedicate these notes to two absentees:

To Marc Raimondo, I will not join you any more in Sydney to write a book on inverse problems and wavelets;

To my father, I know you would have loved...

Marseille, January 2011                                          Laurent Cavalier

## 1.1 Inverse Problems

### *1.1.1 Introduction*

There exist many fields of sciences where inverse problems appear. Some examples are: astronomy (blurred images of the Hubble satellite), econometrics (instrumental variables), financial mathematics (model calibration of the volatility), medical image processing (X-ray tomography), and quantum physics (quantum homodyne tomography)

These are problems where we have indirect observations of an object (a function) that we want to reconstruct. The common structure of all these problems, coming from very different fields, is that we only have access to indirect observations. Due to its indirect nature, solving an inverse problem is usually rather difficult. In fact, there is a need for accurate methods, called regularization methods, in order to solve such an inverse problem.

One example is the problem of X-ray tomography (see Section 1.1.6.5). In this framework, the goal is to reconstruct the internal structure of a human body, by use of external observations. Thus, the internal image cannot be observed directly, but only indirectly.

This notion of indirect observations of some function is usually modeled by use of an operator $A$. From a mathematical point of view, inverse problems usually correspond to the inversion of this operator.

Let $A$ be a bounded operator from $H$ into $G$, where $H$ and $G$ are two separable Hilbert spaces. The classical problem is the following.

$$\text{Given } g \in G, \text{ find } f \in H \text{ such that } Af = g. \qquad (1.1)$$

The terminology of inverse problem comes from the fact that one has to invert the operator $A$. A case of major interest is the case of ill-posed problems where the operator is not invertible. The issue is then to handle this inversion in order to obtain a precise reconstruction.

A classical definition is the following (see [65]).

**Definition 1.1.** A problem is called **well-posed** if

1. there exists a solution to the problem (existence);
2. there is at most one solution to the problem (uniqueness);
3. the solution depends continuously on the data (stability);

A problem which is not well-posed is called **ill-posed**.

One is usually not too much concerned with the existence. If the data space is defined as the set of solutions, existence is clear. Otherwise, the concept of solution may be slightly changed.

If uniqueness is not verified, this is more serious. If there exist several solutions then one has to decide which one is of interest or give additional information. However, the problem of uniqueness is usually relevant in inverse problems.

A standard way of solving the existence and uniqueness problems is by resorting to generalized inverses (see Section 1.1.4).

These two problems (existence and uniqueness) are similar to the standard problem of identifiability in statistics.

Nevertheless, the main issue is usually stability. Indeed, suppose $A^{-1}$ exists but is not bounded. Given a noisy version of $g$ called $g_\varepsilon$, the reconstruction $f_\varepsilon = A^{-1} g_\varepsilon$ may be far from the true $f$. Thus, one needs to invert the operator $A$ in a more stable way. Therefore, one has to develop regularization methods, in order to get fine reconstructions even in ill-posed problems.

A century ago it was generally believed that for natural problems the solution would always depend continuously on the data. Otherwise the mathematical model was believed to be inadequate. These problems therefore were called ill-posed. The idea was that the problem was genuiely not well-posed and that there was no chance to solve such a problem. Ill-posed problems were usually considered, more or less, as unsolvable problems.

Only sixty years ago, scientists realized that a large number of problems which appeared in sciences were ill-posed in any reasonable framework. The idea was developed that there was natural ill-posed problems, in the sense that these were ill-posed in any setting, but they could be however solved by use of regularization methods.

This initiated a lot of research in order to get accurate regularization methods, see for example [127, 123, 128, 108, 10, 49, 110, 117, 116, 126, 51, 72, 112].

### 1.1.2 Statistical Inverse Problems

Loosely speaking solving an inverse problem means recovering an object $f$ from indirect noisy observations $Y$. The object $f$ is usually modeled as a function (or a vector) that has been modified by an operator $A$; thus one observes a noisy version of $Af$. From a mathematical point of view, solving the inverse problem boils down to inverting the operator $A$. The problem is that $A$ may not be invertible or nearly so. This is the case of ill-posed problems and it is of great practical interest as it arises naturally in many fields such as geophysics, finance, astronomy, biology, …

Ill-posed problems are further compounded by the presence of errors (noise) in the data. Statistics enters inverse problems when at least one of the components of the inverse problem (usually the noise) is modeled as stochastic. The question is then to study statistical regularization methods that lead to a meaningful reconstruction despite the noise and ill-posedness.

In Section 1.1 we will present the standard framework of inverse problems focusing on linear operator and stochastic noise. Basic notions on operator theory will be recalled, especially the case of compact operators and singular value decomposition. However, the spectral theory and functional calculus will be defined even for non-compact operators. Several examples of standard inverse problems will be given.

In our opinion the inverse problem framework is better known among statisticians than its statistical approach is among the inverse problem community. For instance, the latter is well acquainted with the concepts of mean, variance and bias but is less familiar with classical concepts such as white noise model, risk estimation, minimax risk, model selection and optimal rates of convergence, which we will discuss in Section 1.2. In addition to these classical notions we will present in Section 1.3 some more recent concepts that have been developed since the 90s like adaptive estimation, oracle inequalities, model selection methods, Stein's unbiased risk estimation and the recent risk hull method. Section 1.4 is a conclusion. We will discuss on the topics which we think are important in the statistical study of inverse problems. Moreover, several open problems will be presented in order to go beyond the framework of these lectures.

All the statistical concepts will be defined and discussed in the framework of inverse problems. Although some of the techniques are specific to this field, some may also be used in more general situations. Other statistical methods not discussed in these notes may also have applications to inverse problems but one should be careful with their application given the intrinsic difficulty and instability of ill-posed problems.

In our mind this is one of the most appealing points of statistical inverse problems. Indeed, most of the standard problems in nonparametric statistics are present in this framework. One may study estimation methods, minimax estimation, rates

of convergence for different functional classes (Besov balls, Hölder balls, Sobolev balls), various risk assessments ($L^2$, $L^p$, pointwise risk). One may also study more recent notions like adaptive estimation, model selection, data-driven selection methods, oracle inequalities, and so on.

On the other hand, there exist also many problems which are specific to the framework of inverse problems. One can consider, noise in the operator, or the problem of choosing the best basis for a given operator. Moreover, due to the ill-posedness and the difficulty of inverse problems, building accurate estimators is usually much more involved here than in the direct problem.

The aim of these notes is to explain some standard theoretical issues regarding the statistical framework of inverse problems. These lectures provide a glimpse of modern nonparametric statistics in the context of inverse problems. Other topics and reviews may be found in [108, 110, 116, 126, 51, 81, 23].

### 1.1.3 Linear Inverse Problems with Random Noise

The classical framework for inverse problem is given by linear inverse problems between two Hilbert spaces.

Let $H$ and $G$ two separable Hilbert spaces. Let $A$ be a known linear bounded operator from the space $H$ to $G$.

Suppose that we have the following observation model

$$Y = Af + \varepsilon\xi, \tag{1.2}$$

where $Y$ is the observation, $f$ is an unknown element in $H$, $\xi$ is an error, $\varepsilon$ corresponds to the noise level. Our aim here is to estimate (or reconstruct) the unknown $f$ by use of the observation $Y$. The idea is that, at least when $\varepsilon$ is small, rather sharp reconstruction should be obtained.

The standard framework first considered by [127] and further studied by [128] corresponds to the case of inverse problems with deterministic noise. In this case, the noise $\xi$ is considered as some element in $G$, with $\|\xi\| \leqslant 1$. Since the noise is some unknown element of a ball in $G$, the results have to be obtained for any possible noise, i.e. for the worst noise. The study of deterministic noise is not the aim of these notes and may be found in Section 1.2.5.

Our framework is a statistical inverse problem, which was considered in [123]. Indeed we observe a noisy version (with random error) of $Af$ and we want to reconstruct $f$. Thus, three main difficulties appear:

- dealing with the noise in the observation (statistics);
- inverting the operator $A$ (inverse problems theory);
- deriving numerical implementations (computational mathematics);

Our aim is now to propose reasonable assumptions on the stochastic noise. The stochastic error is a Hilbert-space process, i.e. a bounded linear operator $\xi : G \to L^2(\Omega, \mathscr{A}, P)$ where $(\Omega, \mathscr{A}, P)$ is the underlying probability space and $L^2(\cdot)$ is the space of all square integrable measurable functions.

Thus, for all functions $g_1, g_2 \in G$, the random variables $\langle \xi, g_j \rangle$ $j = 1, 2$ are defined, by definition $\mathbf{E}\langle \xi, g_j \rangle = 0$ and define its covariance $\mathrm{Cov}_\xi$ as the bounded linear operator ($\|\mathrm{Cov}_\xi\| \leqslant 1$) from $G$ in $G$ such that $\langle \mathrm{Cov}_\xi g_1, g_2 \rangle = \mathrm{Cov}(\langle \xi, g_1 \rangle, \langle \xi, g_2 \rangle)$.

A Hilbert-space random variable $\varkappa$ is a measurable function: $\Omega \to G$. Any Hilbert-space random variable with a finite second moment may be identified with an Hilbert-space process by defining $\varphi \to \langle \varkappa, \varphi \rangle$. However, not all Hilbert-space processes are Hilbert-space random variables.

The action of an operator $A \in L(G, H)$ on some Hilbert-space process $\xi$ is given in Definition 1.3.

The standard hypothesis, which will be mainly considered in these notes, corresponds to the following assumption.

**Definition 1.2.** We say that $\xi$ is a **white noise** process in $G$, if $\mathrm{Cov}_\xi = I$ and the induced random variables are Gaussian:

for all functions $g_1, g_2 \in G$, the random variables $\langle \xi, g_j \rangle$ have distributions $\mathcal{N}(0, \|g_j\|^2)$ and $\mathrm{Cov}(\langle \xi, g_1 \rangle, \langle \xi, g_2 \rangle) = \langle g_1, g_2 \rangle$.

See for example [69].

The white noise is one of the more standard stochastic noise considered in statistics, see for example the Gaussian white noise model in Section 1.1.6.1.

One of the main property of a white noise process is the following.

**Lemma 1.1.** *Let $\xi$ be a white noise in $G$ and $\{\psi_k\}$ be an orthonormal basis in $G$. Define $\xi_k$ by $\xi_k = \langle \xi, \psi_k \rangle$. Then $\{\xi_k\}$ are i.i.d. standard Gaussian random variables.*

*Proof.* By definition $\xi_k \sim \mathcal{N}(0, \|\psi_k\|^2) = \mathcal{N}(0, 1)$. Moreover, we have $\mathbf{E}(\langle \xi, \psi_k \rangle, \langle \xi, \psi_\ell \rangle) = \langle \psi_k, \psi_\ell \rangle = \delta_{k\ell}$. Note also that $\{\xi_k\}$ is Gaussian.

*Remark 1.1.* This lemma is very important and almost characterizes a white noise. Indeed, by projection on some orthonormal basis $\{\psi_k\}$, one obtains a sequence of i.i.d. standard Gaussian random variables $\{\xi_k\}$. This is a way to understand the notion of white noise in applications. In a model with white noise, one obtains a standard Gaussian i.i.d. noise in each observed coefficient (see Section 1.1.5).

*Remark 1.2.* Another remark is that a white noise, as a Hilbert-space process, is not in general a Hilbert-space random variable; note also that $\|\xi\|_G = \infty$, thus $\xi$ is not an element of $G$. One main difference between the deterministic and the stochastic approaches of inverse problems is that the random noise is large compared to the deterministic one. This discussion is postponed to Section 1.2.5.

Note that when $\xi$ is a white noise, $Y$ does not belong to $G$, but acts on $G$, with the following definition, which follows from (1.2),

$$\forall \psi \in G, \ \langle Y, \psi \rangle = \langle Af, \psi \rangle + \varepsilon \langle \xi, \psi \rangle,$$

where $\langle \xi, \psi \rangle \sim \mathcal{N}(0, \|\psi\|^2)$.

*Remark 1.3.* White noise may also be identified with a generalized random variable. Indeed, it does not take its values in $G$ but acts on $G$, see [69].

### 1.1.4 Basic Notions on Operator Theory

Operator theory contains the basic mathematical tools that are needed in inverse problems. In this section we recall rather quickly some standard notions on operator theory which will be used through these lectures. We concentrate on linear bounded operators between Hilbert spaces.

Let $H$ and $G$ be two separable Hilbert spaces.

**Definition 1.3.** 1. $A$ is a **bounded (or continuous) linear operator** from $H$ to $G$ if it is a linear application from $D(A) = H$ to $G$ which is continuous on $H$.
2. Denote by $D(A)$ the definition domain of $A$, by $R(A) = A(H)$ its range, by $N(A) = \{\varphi \in H : A\varphi = 0\}$ its null-space, by $L(H, G)$ the space of linear bounded operators from $H$ to $G$ and by $\|A\|$ the operator norm $\|A\| = \sup\{\|A\varphi\| : \|\varphi\| = 1\}$.
3. The operator $A \in L(H, G)$ is said to be **invertible** if there exists $A^{-1}$ in $L(G, H)$ such that $AA^{-1} = I_G$ and $A^{-1}A = I_H$.
4. There exists $A^*$ such that

$$\langle A\varphi, \psi \rangle = \langle \varphi, A^*\psi \rangle, \ \forall \varphi \in H, \psi \in G.$$

The operator $A^*$ is called the **adjoint** of $A \in L(H, G)$.
5. An operator $A \in L(H, H) = L(H)$ is said to be **self-adjoint** if $A^* = A$. It is called **(strictly) positive** if

$$\langle A\varphi, \varphi \rangle \geqslant (>)0, \ \forall \varphi \in H.$$

6. An operator $U \in L(H, G)$ is said **unitary** if $U^*U = UU^* = I$.
7. One call **eigenvalues** $\lambda \in \mathbb{C}$ and **eigenfunctions** $\varphi \in H, \varphi \neq 0$, elements such that $A\varphi = \lambda\varphi$.
8. Define $A\xi$, the action of any operator $A \in L(G, H)$ on some Hilbert-space process $\xi : G \to L^2(\Omega, \mathscr{A}, P)$ by

$$\langle A\xi, \varphi \rangle = \langle \xi, A^*\varphi \rangle, \ \forall \varphi \in H.$$

Here are some standard results.

**Lemma 1.2.** 1. If $A \in L(H, G)$ and is bijective then $A$ is invertible (i.e. $A^{-1}$ is a linear bounded operator, $A^{-1} \in L(G, H)$).
2. If $A \in L(H, G)$ then $N(A) = R(A^*)^\perp$ and $\overline{R(A)} = N(A^*)^\perp$, where $\overline{(\cdot)}$ and $(\cdot)^\perp$ denote the closure and the orthogonal subspaces.
3. If $A \in L(H, G)$ then $A^* \in L(G, H)$.
4. If $A$ is injective so is $A^*A$.
5. If $A \in L(H, G)$ then $A^*A \in L(H)$ is self-adjoint and positive.
6. A self-adjoint operator is injective if and only if its range is dense in $H$.
7. A self-adjoint operator is invertible if and only if $R(A) = H$.
8. A self-adjoint operator then

$$\|A\| = \sup_{\|\varphi\|=1} |\langle A\varphi, \varphi \rangle|.$$

*9. If $U \in L(H, G)$ is unitary, then*

$$\langle U\varphi, U\psi \rangle = \langle \varphi, \psi \rangle, \ \forall \varphi, \psi \in H.$$

*Proof.* (1) A proof may be found in [73].

(2) We have $\langle \varphi, A^*\psi \rangle = \langle A\varphi, \psi \rangle = 0$ for all $\varphi \in N(A), \psi \in G$. Hence, $N(A) = R(A^*)^\perp$. Interchanging the roles of $A$ and $A^*$ gives $N(A^*) = R(A)^\perp$. Thus, $N(A^*)^\perp = (R(A)^\perp)^\perp = \overline{R(A)}$.

(3) Straightforward.

(4) We have $\langle A^*A\varphi, \varphi \rangle = \langle A\varphi, A\varphi \rangle = \|A\varphi\|^2$. If $\varphi_0 \in N(A^*A)$ then $\varphi_0 \in N(A)$.

(5) We have $\langle A^*A\varphi, \psi \rangle = \langle A\varphi, A\psi \rangle = \langle \varphi, A^*A\psi \rangle$. Note also that $\langle A^*A\varphi, \varphi \rangle = \|A\varphi\|^2 \geqslant 0$.

(6) $A$ injective if and only if $N(A) = \{0\}$ if and only if $N(A)^\perp = H$. We then use that $\overline{R(A)} = N(A^*)^\perp = N(A)^\perp$ by (2) and the fact that $A$ is self-adjoint.

(7) By (6), $A$ invertible is thus equivalent to $R(A) = H$.

(8) A proof may be found in [73].

(9) We have, since $U$ is unitary,

$$\langle U\varphi, U\psi \rangle = \langle U^*U\varphi, \psi \rangle = \langle \varphi, \psi \rangle, \ \forall \varphi, \psi \in H.$$

Some new definitions and properties concerning mostly compact operators are presented here. Compact operators are very important in inverse problems for several reasons.

First, a compact operator is not invertible, i.e. has no bounded inverse (see Lemma 1.3). Thus if $A$ is a compact operator the problem is naturally ill-posed in the sense of Definition 1.1. From a mathematical point of view, ill-posed problems are the more challenging.

Compact operators have simple spectra only composed of eigenvalues, see Theorem 1.1. This is a nice property of compact operators which gives rise to natural basis of functions to use, the singular value decomposition. By projection on this natural basis, we will obtain a sequence space model in Section 1.1.5. This model in the space of coefficients, is usually more easy to deal with from a statistical point of view.

**Definition 1.4.** 1. An operator $A$ from $H$ to $G$ is called **compact** if each bounded set in $H$ has an image by $A$ which is relatively compact in $G$, i.e. with a compact closure.

2. Denote by $K(H, G)$ the space of **compact linear bounded operator**.

3. The **strong convergence**, denoted $\rightarrow_s$, is the convergence with respect to the norm in $H$ or $G$.

4. The **weak convergence**, denoted $\rightarrow_w$, is the convergence with respect to $\langle \varphi, \cdot \rangle$ for all $\varphi \in H$ or $G$.

**Lemma 1.3.** *1. Let $A \in K(H,G)$, then there exist $A_n \in K(H,G)$, such that dim $R(A_n) < \infty$ and $\|A_n - A\| \to 0$, as $n \to \infty$.*
*2. $A \in K(H,G)$ is equivalent to $A^* \in K(G,H)$*
*3. $A \in K(H,G)$ is equivalent to $\forall \varphi_k \in H : \varphi_k \to_w \varphi$ implies $A\varphi_k \to_s A\varphi$.*
*4. If $A \in K(H,G)$ and $\dim(H) = \infty$ then $A^{-1}$ is not bounded.*
*5. $A \in K(H)$ is equivalent to the fact that for any orthonormal sequence $\{\varphi_k\}$, one has $\lim_{k \to \infty} \langle A\varphi_k, \varphi_k \rangle = 0$.*

*Proof.* A proof may be found in [73].

**Theorem 1.1.** *Let $A \in K(H)$ be self-adjoint. Then there exists a complete orthonormal system $E = \{\varphi_j : j \in I\}$ of $H$ consisting of eigenfunctions of $A$. Here $I$ is some index set and $A\varphi_j = \lambda_j \varphi_j$, for $j \in I$. The set $J = \{j \in I : \lambda_j \neq 0\}$ is countable and*

$$A\varphi = \sum_{j \in I} \lambda_j \langle \varphi, \varphi_j \rangle \varphi_j, \qquad (1.3)$$

*for all $\varphi \in H$. Moreover, for any $\delta > 0$ the set $J_\delta = \{j \in I : |\lambda_j| \geqslant \delta\}$ is finite.*

*Proof.* This proof may be found in [72]. First, we prove the existence of an eigenvalue for a self-adjoint compact operator (if $H \neq \{0\}$). Due to Lemma 1.2, there exists a sequence $\{\varphi_k\}$ with $\|\varphi_k\| = 1$ such that, for $\lambda = \pm\|A\|$, $\langle A\varphi_k, \varphi_k \rangle \to \lambda$ as $k \to \infty$. Remark that

$$0 \leqslant \|A\varphi_k - \lambda \varphi_k\|^2 = \|A\varphi_k\|^2 - 2\lambda \langle A\varphi_k, \varphi_k \rangle + \lambda^2 \|\varphi_k\|^2$$

$$\leqslant \|A\|^2 - 2\lambda \langle A\varphi_k, \varphi_k \rangle + \lambda^2 \to 0, \text{ as } k \to \infty.$$

Thus, $A\varphi_k \to \lambda \varphi_k$ as $k \to \infty$. Since $A$ is compact, there exists a subsequence such that $A\varphi_{k(n)} \to \psi$ as $n \to \infty$. It follows that $\lambda \varphi_{k(n)} \to \psi$ as $n \to \infty$. Denote $\varphi = \psi/\lambda$, therefore $\varphi_{k(n)} \to \psi$ as $n \to \infty$ and $A\varphi = \lambda \varphi$, since $A$ is bounded.

We then prove that the system is orthogonal. Suppose $\lambda_j \neq \lambda_k$. We have

$$\langle A\varphi_j, \varphi_k \rangle = \lambda_j \langle \varphi_j, \varphi_k \rangle.$$

Moreover, since $A$ is self-adjoint, we have

$$\langle A\varphi_j, \varphi_k \rangle = \langle \varphi_j, A\varphi_k \rangle = \lambda_k \langle \varphi_j, \varphi_k \rangle.$$

Thus $\varphi_j$ and $\varphi_k$ are orthogonal.

We now study the case where $\varphi_j$ and $\varphi_k$ are eigenfunctions with the same eigenvalue $\lambda$, suppose $\langle \varphi_j, \varphi_k \rangle = c \neq 0$. Thus, $\varphi_j - \varphi_k/c$ is still an eigenfunction related to $\lambda$ and orthogonal to $\varphi_j$. One may easily orthonormalize the system.

The last part consists in proving the completeness. By Zorn's Lemma, choose $E$ the maximal set of eigenfunctions of $A$. Let $S$ be the closed linear span of $E$. Obviously, $A(S) \subset S$. Moreover, $A(S^\perp) \subset S^\perp$, since $\langle As, \varphi \rangle = \langle s, A\varphi \rangle = 0$ for all $s \in S^\perp$ and all $\varphi \in S$. Remark that $A_{|S^\perp}$ is compact and self-adjoint. Hence, if $S^\perp \neq \{0\}$ there exists an eigenfunction $\psi \in S^\perp$ (by the first part of this proof). Since

this contradicts the maximality of $E$, we conclude that $S^\perp = \{0\}$. Therefore the orthonormal system is complete. To show (1.3) we apply $A$ to the representation

$$\varphi = \sum_{j \in I} \langle \varphi, \varphi_j \rangle \varphi_j. \tag{1.4}$$

Remark that only countable number of terms in (1.4) can be non-zero. Indeed, by Bessel's inequality we have

$$\sum_{\varphi_j \in E} |\langle \varphi, \varphi_j \rangle|^2 = \sup \left\{ \sum_{\varphi_j \in F} |\langle \varphi, \varphi_j \rangle|^2 : F \subset E, \operatorname{card}(F) < \infty \right\} \leqslant \|\varphi\|^2 < \infty.$$

Therefore, for any $k \in \mathbb{N}$, the set $S_k = \{\varphi_j \in E : |\langle \varphi, \varphi_j \rangle| \in [\|\varphi\|/(k+1), \|\varphi\|/k]\}$ is finite, and the union for all $k \in \mathbb{N}$ is then countable.

Assume that $J_\delta$ is infinite for some $\delta > 0$. Since $A$ is compact, there exists a subsequence $\{\varphi_{k(n)}\}$ of $\{\varphi_k\}$ such that $\{A\varphi_{k(n)}\} = \{\lambda_{k(n)}\varphi_{k(n)}\}$ is a Cauchy sequence. This is in contradiction since $\|\lambda_k \varphi_k - \lambda_j \varphi_j\|^2 = \lambda_k^2 + \lambda_j^2 \geqslant 2\delta^2$ for $j \neq k$ due to the orthonormality of $\{\varphi_k\}$.

*Remark 1.4.* A linear bounded self-adjoint compact operator between two Hilbert spaces may thus be seen as an infinite matrix. In applications, a large matrix could be modelized by a compact operator. However, due to Theorem 1.1, the eigenvalues $\lambda_j$ are going to 0. This is fundamental and characterizes the notion of ill-posed problems (see Definition 1.7). One observes a function through an operator $A$ which, in some sense, concentrates to 0. Thus, the inversion of such an operator has to be made carefully, otherwise, the reconstruction will explose.

In general inverse problems, we neither assume that $A$ is injective nor that $g \in R(A)$. Thus, we usually need some standard definitions of a generalized notion of inverse for the equation $Af = g$ (see [64]).

**Definition 1.5.** Let $A \in L(H, G)$.

1. We call $f$ a **least-squares solution** of the problem (1.1) if

$$\|Af - g\| = \inf\{\|A\varphi - g\| : \varphi \in H\}.$$

2. We call $f$ a **best approximate solution** of the problem (1.1) if it is a least-squares solution and if

$$\|f\| = \inf\{\|\varphi\| : \varphi \text{ is a least-squares solution}\}.$$

3. The **Moore-Penrose (generalized) inverse** $A^\dagger : D(A^\dagger) \to H$ of $A$ defined on $D(A^\dagger) = R(A) \oplus R(A)^\perp$ maps $g \in D(A^\dagger)$ to the best approximate solution of (1.1). The existence of a best approximate solution is guaranted by $g \in R(A) \oplus R(A)^\perp$.

We have

**Lemma 1.4.** *Let* $Q : G \to \overline{R(A)}$ *be the orthogonal projection onto* $\overline{R(A)}$. *Then the three statements are equivalent:*

1. $f \in H$ *is a least-squares solution of (1.1).*
2. $Af = Qg$.
3. *The* **normal equation** $A^*Af = A^*g$ *holds.*

   *We have in addition the following properties for* $g \in R(A) \oplus R(A)^\perp$.
4. *Any least-squares solution belongs to* $A^\dagger g + N(A)$.
5. *We also have that a best approximate solution exists, is unique and equals to* $A^\dagger g$.

*Proof.* Since $Q$ is an orthogonal projection on $\overline{R(A)}$, remark that $\langle Af - Qg, (I - Q)g \rangle = 0$. We then have

$$\|Af - g\|^2 = \|Af - Qg\|^2 + \|(I - Q)g\|^2.$$

This shows that (2) implies (1). Vice versa, if $f$ is a least-squares solution the last equation shows that $f$ is a minimum of $\|Af - Qg\|$. Again by property of the projection we obtain (2).

Moreover, $f$ is a least-squares solution if and only if $Af$ is the closest element in $R(A)$ to $g$, which is equivalent to $Af - g \in R(A)^\perp = N(A^*)$, i.e. $A^*(Af - g) = 0$.

(4) Suppose that $g \in R(A) \oplus R(A)^\perp$. Then $Qg \in R(A)$ and (2) is true and there exists at least one least-squares solution $f_0$. Moreover, due to (2), any element of $f_0 + N(A)$ is also a least-squares solution.

(5) Remark that for any $u \in N(A)$:

$$\|f_0 + u\|^2 = \|(I - P)(f_0 + u)\|^2 + \|P(f_0 + u)\|^2 = \|(I - P)f_0\|^2 + \|Pf_0 + u\|^2,$$

where $f_0$ is a least-squares solution, $P$ is the orthogonal projection on $N(A)$. This yields the uniqueness of the best approximate solution, which is equal to $(I - P)f_0$.

Obviously, if $A^{-1} \in L(G, H)$ exists then $A^{-1} = A^\dagger$.

Under assumptions of Lemma 1.4, the best approximate solution is in fact the least-squares solution with a null term in the null-space of $A$. Indeed, any $f \in N(A)$ is such that $Af = 0$, and cannot be observed through $A$. Thus, there is no real meaning in trying to reconstruct it.

The normal equation is a different way to express an inverse problem. Indeed one may multiply the first problem by $A^*$ and then get the equivalent normal equation.

*Remark 1.5.* In a statistical inverse problem, we observe a (random) noisy version of $Af$. Thus, if $A$ is injective then the unique best approximate solution is $f$ (by Lemma 1.4).

### *1.1.5 Singular Value Decomposition and Sequence Space Model*

Let $A \in L(H, G)$ be an injective and compact operator. We have, by applying Theorem 1.1 to $A^*A$, which is self-adjoint and strictly positive,

$$A^*Af = \sum_{k=1}^{\infty} \rho_k \langle f, \varphi_k \rangle \varphi_k,$$

where $\rho_k > 0$. Define the normalized image $\{\psi_k\} \in G$ of $\{\varphi_k\} \in H$ by

$$\psi_k = b_k^{-1} A \varphi_k,$$

where $b_k = \sqrt{\rho_k} > 0$. Remark that $\{\psi_k\}$ are orthogonal,

$$\langle \psi_k, \psi_\ell \rangle = b_k^{-1} b_\ell^{-1} \langle A\varphi_k, A\varphi_\ell \rangle = b_k^{-1} b_\ell^{-1} \langle A^*A\varphi_k, \varphi_\ell \rangle = b_k b_\ell^{-1} \langle \varphi_k, \varphi_\ell \rangle = \delta_{k\ell},$$

where $\delta_{k\ell}$ denotes the Kronecker symbol (0 if $k \neq \ell$, 1 if $k = \ell$). Note that this implies $\|\psi_k\|^2 = 1$. Thus, $\{\psi_k\}$ is an orthonormal system. Moreover

$$A^*\psi_k = b_k^{-1} A^*A\varphi_k = b_k^{-1} b_k^2 \varphi_k = b_k \varphi_k.$$

Thus, we have

$$A\varphi_k = b_k \psi_k, \ \ A^*\psi_k = b_k \varphi_k.$$

The $b_k > 0$ are called **singular values** of the operator $A$. Note also that, since $A^*A$ is compact and self-adjoint then $b_k \to 0$ as $k \to \infty$ by Theorem 1.1.

**Definition 1.6.** We say that $A$ admits a **singular value decomposition (SVD)** if, $\forall f \in H$,

$$A^*Af = \sum_{k=1}^{\infty} b_k^2 \theta_k \, \varphi_k,$$

where $\theta_k$ are the coefficients of $f$ in the orthonormal basis $\{\varphi_k\} \in H$, $\{b_k\}$ are the singular values.

The SVD is the natural basis for $A$ since it diagonalizes $A^*A$.

Now consider the projection of $Y$ on $\{\psi_k\}$

$$\langle Y, \psi_k \rangle = \langle Af, \psi_k \rangle + \varepsilon \langle \xi, \psi_k \rangle = \langle Af, b_k^{-1} A\varphi_k \rangle + \varepsilon \xi_k$$

$$= b_k^{-1} \langle A^*Af, \varphi_k \rangle + \varepsilon \xi_k = b_k \theta_k + \xi_k,$$

where $\xi_k = \langle \xi, \psi_k \rangle$.

Since $\xi$ is a white noise $\{\xi_k\}$ is a sequence of i.i.d. standard Gaussian random variables $\mathcal{N}(0,1)$ by Lemma 1.1.

Thus, under these assumptions, one has the equivalent discrete sequence observation model derived from (1.2):

$$y_k = b_k \theta_k + \varepsilon \xi_k, \ \ k = 1, 2, \ldots, \tag{1.5}$$

where $y_k$ stands for $\langle Y, \psi_k \rangle$. This model is called the **sequence space model**. The aim here is to estimate the sequence $\theta = \{\theta_k\}$ from the observations $y = \{y_k\}$.

One can see the influence of the ill-posedness of the inverse problem when $A$ is compact. Indeed, since $b_k$ are the singular values of a compact operator, then $b_k \to 0$

as $k \to \infty$. Thus, when $k$ increases the 'signal' $b_k \theta_k$ is weaker and it is clearly more difficult to estimate $\theta_k$.

Another comment concerns the fact that the aim is to estimate $\{\theta_k\}$ and not $\{b_k \theta_k\}$. Thus, one really has to consider the inverses of the $b_k$, i.e., to invert the operator $A$.

For this reason, the following equivalent model to (1.5) is more natural

$$X_k = \theta_k + \varepsilon \sigma_k \xi_k, \quad k = 1, 2, \ldots, \tag{1.6}$$

where $X_k = y_k / b_k$, and $\sigma_k = b_k^{-1} > 0$. Note that $\sigma_k \to \infty$. In this model the aim is to estimate $\{\theta_k\}$ from $\{X_k\}$. When $k$ is large the noise in $X_k$ may then be very large, making the estimation difficult.

The sequence space model (1.5) or (1.6) for statistical inverse problems was studied in many papers, see [39, 95, 78, 32], among others.

*Remark 1.6.* For ill-posed inverse problems we have $b_k \to 0$ and $\sigma_k \to \infty$, as $k \to \infty$. We can see that ill-posed problems are more difficult than the direct problem. Indeed, when $k$ is large, the noise $\varepsilon \sigma_k \xi_k$ will dominate. Thus, the estimation of $\{\theta_k\}$ from $\{X_k\}$ is more involved.

One can characterize linear inverse problems by the difficulty of the operator, i.e. with our notations, by the behaviour of the $\sigma_k$. If $\sigma_k \to \infty$, as $k \to \infty$, the problem is ill-posed.

**Definition 1.7.** An inverse problem is called **mildly ill-posed** if the sequence $\sigma_k$ has a polynomial behaviour when $k$ is large

$$\sigma_k \asymp k^\beta, \ k \to \infty,$$

and **severely ill-posed** if $\sigma_k$ tends to infinity at an exponential rate

$$\sigma_k \asymp \exp(\beta k), \ k \to \infty,$$

where $\beta > 0$ is called the **degree of ill-posedness** of the inverse problem.

A special case of inverse problems is the **direct problem** where

$$\sigma_k \asymp 1, \ k \to \infty,$$

which corresponds to $\beta = 0$.

Here and later, $a_n \asymp b_n$ means that there exist $0 < c_1 \leqslant c_2 < \infty$ such that, $c_1 \leqslant a_n / b_n \leqslant c_2$, as $n \to \infty$.

*Remark 1.7.* One may also consider inverse problems which are more difficult than severely ill-posed, in the case where $\sigma_k \asymp \exp(\beta k^r)$, where $\beta > 0$ and $r \geqslant 1$.

*Remark 1.8.* There exist more general definitions of the degree of ill-posedness related to the noise structure, smoothness assumptions on $f$, smoothing properties of $A$ (see [131, 103]). However, for the sake of simplicity, we prefer to deal with the simple notion defined above.

*Remark 1.9.* An important special case is the case where $A = I$. This corresponds to the **direct problem** where $f$ is directly observed (with noise) with no inverse problem, i.e. without the need of inverting some operator $A$. In this case $\sigma_k \equiv 1$ and the model in (1.6) corresponds to the classical sequence space model in statistics. The model is then related to the Gaussian white noise model and is very close to nonparametric regression with $\varepsilon = n^{-1/2}$ (see Section 1.1.6.1).

## 1.1.6 Examples

Here are some examples of ill-posed problems where the SVD may be applied. In each case, the SVD can be explicitly computed.

Moreover, from a practical point of view, methods based on SVD are usually rather expensive in term of computations. For these reasons, many populars methods nowadays do not use explicitely the SVD.

On the other hand, even for these methods, the spectral domain is often used in order to deal with the theoretical accuracy of the methods.

### 1.1.6.1 Standard Gaussian White Noise

One of most classical model in nonparametric statistics is the Gaussian white noise

$$dY(t) = f(t)dt + \varepsilon dW(t), \ t \in [0,1], \tag{1.7}$$

where one observes $\{Y(t), t \in [0,1]\}$, $f$ is an unknown function in $L^2[0,1]$, $W$ is a Wiener process, $\varepsilon > 0$ is the noise level. One may check easily that $dW$ corresponds to a white noise. Indeed, we obtain directly from the definition of integral against a Wiener process that for all $\varphi \in L^2[0,1]$,

$$\int_0^1 \varphi(t)dW(t) \sim \mathcal{N}\left(0, \int_0^1 |\varphi(t)|^2 dt\right).$$

We also obtain the property for the scalar product by the definition of the Wiener process

$$\text{Cov}\left(\int_0^1 \varphi_1(t)dW(t), \int_0^1 \varphi_2(t)dW(t)\right) = \int_0^1 \varphi_1(t)\varphi_2(t)dt = \langle \varphi_1, \varphi_2 \rangle,$$

for all $\varphi_1, \varphi_2 \in L^2[0,1]$.

This model is a very specific inverse problem since, in this case, the operator is $A = I$ and $H = G = L^2[0,1]$. However, most of the results on inverse problems will apply in this framework. This model is often called a **direct problem**, since from our definition we have at our disposal direct observations and not indirect ones.

In this case, the sequence space model may be obtained by projecting on any orthonormal basis $\{\psi_k\} \in L^2[0,1]$. Doing so, one obtains

$$\int_0^1 \psi_k(t)dY(t) = \int_0^1 \psi_k(t)f(t)dt + \varepsilon \int_0^1 \psi_k(t)dW(t),$$

which is equivalent to

$$y_k = \theta_k + \varepsilon \xi_k, \ k = 1,\ldots,$$

where the $\theta_k$ are the coefficients of $f$ in $\{\psi_k\}$ and

$$\xi_k = \int_0^1 \psi_k(t)dW(t) \sim \mathcal{N}(0,1)$$

with $\{\xi_k\}$ i.i.d. We then obtain a sequence space model where $b_k \equiv 1$.

It is well-known that the Gaussian white noise model defined in (1.7) is an idealized version of the more standard nonparametric regression

$$Y_i = f(X_i) + \xi_i, \ i = 1,\ldots,n, \tag{1.8}$$

where $(X_1,Y_1),..,(X_n,Y_n)$ are observed (we may assume $X_i \in [0,1]$), $f$ is an unknown function in $L^2[0,1]$, and $\{\xi_i\}$ are i.i.d. zero-mean Gaussian random variables with variance $\sigma^2$.

The Gaussian white noise model may be understood as a large sample limit of nonparametric regression in (1.8). Indeed, by projecting (1.7) on the intervals $I_i = [(i-1)/n, i/n]$, $i = 1,\ldots,n$, one obtains

$$n \int_{I_i} dY(t) = n \int_{(i-1)/n}^{i/n} f(t)dt + n\varepsilon \int_{(i-1)/n}^{i/n} dW(t).$$

Thus, if $f$ is smooth enough, and $\varepsilon^2 = \sigma^2/n$, one has an informal writting

$$Y_i \asymp f(i/n) + \xi_i, \ i = 1,\ldots,n,$$

where $\{\xi_i\}$ are i.i.d. zero-mean Gaussian random variables with variance $\sigma^2$.

This equivalence is proved in different frameworks and models (nonparametric regression, density, non-Gaussian noise ...) in [13, 107, 63, 113, 129].

Thus, under proper calibration, i.e. $\varepsilon^2 \asymp \sigma^2/n$, the asymptotics of model (1.8) as $n \to \infty$ and (1.7) as $\varepsilon \to 0$ are equivalent with the asymptotics of the latter being easier to derive.

In the inverse problem context, model (1.2) may be seen as an idealized version of the discrete sample model

$$Y_i = Af(X_i) + \xi_i, \ i = 1,\ldots,n, \tag{1.9}$$

where $(X_1,Y_1),..,(X_n,Y_n)$ are observed (we may assume $X_i \in [0,1]$), $f$ is an unknown function in $L^2[0,1]$, $A$ is an operator from $L^2[0,1]$ into $L^2[0,1]$, and $\xi_i$ are i.i.d. zero-mean Gaussian random variables with variance $\sigma^2$.

### 1.1.6.2 Derivation

Another related example, which does not exactly correspond to our framework, but is very important, is the estimation of a derivative. Suppose that we observe

$$Y = f + \varepsilon\xi, \tag{1.10}$$

where $H = L^2[0,1]$, $f$ is a $1-$periodic $C^\beta$ function in $L^2[0,1]$, i.e. $\beta$ continuously differentiable, $\beta \in \mathbb{N}$ and $\xi$ is a white noise. A standard problem in statistics is the estimation of the derivative $D^\beta f = f^{(\beta)}$ of $f$, or the function $f$ itself when $\beta = 0$ (which corresponds to the previous section). This problem is studied for example in [47].

One may use here the Fourier basis $\varphi_k(x) = e^{2\pi ikx}$, $k \in \mathbb{Z}$. Denote by $\theta_k$ the Fourier coefficients of $f$,

$$\theta_k = \int_0^1 f(x)e^{2\pi ikx}dx$$

and note that

$$D^\beta(e^{2\pi ik\cdot})(x) = (2\pi ik)^\beta e^{2\pi ikx}.$$

It is well-known that we then have

$$f^{(\beta)} = \sum_{k=-\infty}^{\infty} (2\pi ik)^\beta \theta_k \varphi_k.$$

We have the following equivalent model in the Fourier domain

$$y_k = \theta_k + \varepsilon\xi_k, \ k \in \mathbb{Z} \setminus \{0\},$$

and we want to estimate $v_k = \theta_k(2\pi ik)^\beta$. This is equivalent to, observing

$$y_k = (2\pi ik)^{-\beta}v_k + \varepsilon\xi_k, \ k \in \mathbb{Z} \setminus \{0\},$$

and estimating $\theta_k$.

Thus, derivation is a mildly ill-posed inverse problem of degree $\beta$.

### 1.1.6.3 Circular Deconvolution

The framework of (circular) deconvolution is perhaps one of the most well-known inverse problem. It is used in many applications as econometrics, physics, astronomy, medical image processing. For example, it corresponds to the problem of a blurred signal that one wants to recover from indirect data.

*Example 1.1.* One famous example of an inverse problem of deconvolution is the blurred images of the Hubble space telescope. In the early 1990, the Hubble satellite was launched into low-earth orbit outside of the disturbing atmosphere in order to provide images with a spatial resolution never achieved before. Unfortunately,

quickly after launch, a manufacturing error in the main mirror was detected, causing severe spherical aberrations in the images. Therefore, before the space shuttle Endeavour visited the telescope in 1993 to fix the error, astronomers employed inverse problem techniques to improve the blurred images (see [1]).

Consider the following convolution operator:

$$Af(t) = r * f(t) = \int_0^1 r(t-x)f(x)dx, \ \ t \in [0,1],$$

where $r$ is a known 1-periodic symmetric around 0 real-valued convolution kernel in $L^2[0,1]$. In this model, $A$ is a linear bounded operator from $L^2[0,1]$ to $L^2[0,1]$.

This operator is a Hilbert-Schmidt integral operator and it is an **Hilbert-Schmidt operator**, i.e., it is such that for some (and then any) orthonormal basis $\{e_k\}$ we have $\sum \|Ae_k\|^2 < \infty$. It is then a compact operator.

Remark that, if $\{\varphi_k\}$ and $\{\psi_k\}$ are the SVD bases defined in Section 1.1.5 then

$$\sum \|A\varphi_k\|^2 = \sum \langle A^*A\varphi_k, \varphi_k \rangle = \sum b_k^2 < \infty.$$

This shows that the singular values are decreasing rather fastly in this situation.

By simple computations one may see that the adjoint $A^*$ is also a Hilbert-Schmidt integral operator, with kernel $\overline{r(x-t)}$, where $\overline{(\cdot)}$ denotes the complex conjugate. Since $r$ is real-valued and symmetric around 0, the operator $A$ is also self-adjoint.

Define then the following model

$$Y(t) = r * f(t) + \varepsilon \, \xi(t), \ \ t \in [0,1], \tag{1.11}$$

where $Y$ is observed, $f$ is an unknown periodic function in $L^2[0,1]$ and $\xi(t)$ is a white noise.

This model is quite popular and has been studied in a large number of statistical papers, see [50, 39, 45, 32, 29].

Define here $\{\varphi_k(t)\}$ the real trigonometric basis on $[0,1]$:

$$\varphi_1(t) \equiv 1, \ \ \varphi_{2k}(t) = \sqrt{2}\cos(2\pi kt), \ \ \varphi_{2k+1}(t) = \sqrt{2}\sin(2\pi kt), \ \ k = 1,2,\ldots.$$

A function in $L^2[0,1]$ may be decomposed on $\{\varphi_k(t)\}$.

Remark now that by a simple change of variables

$$\int_0^1 r(t-x)e^{2\pi ikx}dx = e^{2\pi ikt}\int_{-t}^{1-t} r(-y)e^{2\pi iky}dy = e^{2\pi ikt}\int_0^1 r(x)e^{2\pi ikx}dx,$$

by periodicity.

The SVD basis is then clearly here the Fourier basis, i.e. $e^{2\pi ik\cdot}$.

We make the projection of (1.11) on $\{\varphi_k(t)\}$, in the Fourier domain, and obtain

$$y_k = b_k\theta_k + \varepsilon\xi_k,$$

where $b_k = \sqrt{2}\int_0^1 r(x)\cos(2\pi kx)dx$ for even $k$, $b_k = \sqrt{2}\int_0^1 r(x)\sin(2\pi kx)dx$ for odd $k$, $\theta_k$ are the Fourier coefficients of $f$, and $\xi_k$ are i.i.d. standard Gaussian random variables.

### 1.1.6.4 Heat Equation

Consider the following heat equation which describes the heat at time $t$ and position $x$ based on some initial conditions:

$$\frac{\partial}{\partial t}u(x,t) = \frac{\partial^2}{\partial x^2}u(x,t), \; u(x,0) = f(x), \; u(0,t) = u(1,t) = 0,$$

where $u(x,t)$ is defined for $x \in [0,1], t \in [0,T]$, and the initial condition $f$ is a 1-periodic function. The problem is the following: given the temperature $g(x) = u(x,T)$ at time $T$ find the initial temperature $f \in L^2[0,1]$ at time $t = 0$.

Due to the boundary conditions, one uses here the sine basis ($\{\sqrt{2}\sin(k\pi\cdot)\}$). Let $\theta_k(t) = \sqrt{2}\int_0^1 \sin(k\pi x)f(x)dx$ denote the Fourier coefficients of $f$ with respect to the complete orthonormal system $\{\varphi_k\}$ of $L^2[0,1]$.

In this case, one obtains an ordinary differential equation in the Fourier domain, which provides the following expression for $u$:

$$u(x,t) = \sqrt{2}\sum_{k=1}^{\infty}\theta_k e^{-\pi^2 k^2 t}\sin(k\pi x).$$

The problem is then: given the final temperature $u(x,T) = Af(x)$ to find the initial temperature $f$. Thus, we may write this problem as an inverse problem with the operator

$$Af(x) = \int_0^1 \sum_{k=1}^{\infty} e^{-\pi^2 k^2 T} 2\sin(k\pi x)\sin(k\pi y)f(y)dy.$$

Thus, $A$ is a linear bounded injective compact operator, whose SVD is given by the sine basis. The singular values $b_k$ are equal to $e^{-\pi^2 k^2 T/2}$ and the problem is therefore severely ill-posed.

The model is then the following:

$$Y(x) = u(x,T) + \varepsilon\,\xi(x), \; x \in [0,1],$$

where $\xi$ is a white noise in $L^2[0,1]$. We want here to recover $f \in L^2[0,1]$.

From a statistical point of view, the problem is the following: given a noisy version of the final temperature (at time $T$) find the unknown initial condition $f$ (at time 0).

This framework has been studied in [61] and [24].

By projection on the sine basis, one obtains the following sequence space model:

$$y_k = b_k\theta_k + \varepsilon\xi_k,$$

where $b_k = e^{-\pi^2 k^2 T/2}$, $\theta_k$ are the Fourier (sine) coefficients of $f$, and $\xi_k$ are i.i.d. standard Gaussian random variables.

Remark that here the problem is very difficult. Indeed, it is even worse than a severely ill-posed problem since the singular values are decreasing faster than exponentially.

From a practical point of view, one can see that an error of order $10^{-8}$ in the fifth Fourier coefficient of $u(x,T)$ may lead to an error of $1000C$ in the initial temperature $f(x) = u(x,0)$.

One has to be very careful when solving this kind of inverse problem.

### 1.1.6.5 Computerized Tomography

Computerized tomography is used in medical image processing and has been studied for a long time, see [104]. In medical X-ray tomography one tries to have an image of the internal structure of an object. This image is characterized by a function $f$. However, there is no direct observations of $f$. Suppose that one observes the attenuation of the X-rays. Denote by $I_0$ and $I_1$ the initial and final intensities, by $x$ the position on a given line $L$ and by $\Delta I(x)$ the attenuation for a small $\Delta x$. One then has

$$\Delta I(x) = -f(x)I(x)\Delta x,$$

which corresponds from a mathematical point of view to

$$\frac{I'(x)}{I(x)} = -f(x),$$

and then by integration

$$\log(I_1) - \log(I_0) = \log\left(\frac{I_1}{I_0}\right) = -\int_L f(x)dx.$$

Thus observing $I_1/I_0$ is equivalent to the observation of $\exp\left(-\int_L f(x)dx\right)$. By measuring attenuation of X-rays, one observes cross section of the body.

From a mathematical point of view this problem corresponds to the reconstruction of an unknown function $f$ in $\mathbb{R}^2$ (or in general $\mathbb{R}^d$) based on observations of its Radon transform $Rf$, i.e., of integrals over hyperplanes.

Let $B = \{x \in \mathbb{R} : \|x\| \leqslant 1\}$ be the unit ball in $\mathbb{R}^2$. Consider the integrals of a function $f : B \to \mathbb{R}$ over all the lines that intersect $B$. We parametrize the lines by the length $u \in [0,1]$ of the perpendicular from the origin to the line and by the orientation $s \in [0,2\pi)$ of this perpendicular with respect to the $x$-axis.

Suppose that the function $f$ belongs to $L^1(B) \cap L^2(B)$. Define the Radon transform $Rf$ of the function $f$ by

$$Rf(u,s) = \frac{\pi}{2(1-u^2)^{\frac{1}{2}}} \int_{-\sqrt{1-u^2}}^{\sqrt{1-u^2}} f(u\cos s - t\sin s, u\sin s + t\cos s)dt, \qquad (1.12)$$

where $(u,s) \in S = \{(u,s) : 0 \leqslant u \leqslant 1, \ 0 \leqslant s < 2\pi\}$. With this definition, the Radon transform $Rf(u,s)$ is $\pi$ times the average of $f$ over the line segment (parametrized by $(u,s)$) that intersects $B$. It is natural to consider $Rf$ as an element of $L^2(S,d\mu)$ where $\mu$ is the measure defined by $d\mu(u,s) = 2\pi^{-1}(1-u^2)^{\frac{1}{2}} du \, ds$. This measure $\mu$ is here in order to renormalize over lines.

In this case, the Radon operator $R$ is a linear, bounded and compact operator from $L^2(B)$ into $L^2(S,d\mu)$.

The SVD of the Radon transform was well-studied, e.g. by [37, 104]. To introduce it, define the set of double indices $\mathscr{L} = \{\ell = (j,k) : j \geqslant 0, k \geqslant 0\}$. An orthonormal complex-valued basis for $L^2(B)$ is given by

$$\tilde{\varphi}_\ell(r,t) = \pi^{-\frac{1}{2}}(j+k+1)^{\frac{1}{2}} Z_{j+k}^{|j-k|}(r) e^{i(j-k)t}, \ \ell = (j,k) \in \mathscr{L}, \ (r,t) \in B, \quad (1.13)$$

where $Z_a^b$ denotes the Zernike polynomial of degree $a$ and order $b$. The corresponding orthonormal functions in $L^2(S,d\mu)$ are

$$\tilde{\psi}_\ell(u,s) = \pi^{-\frac{1}{2}} U_{j+k}(u) e^{i(j-k)s}, \ \ell = (j,k) \in \mathscr{L}, \ (u,s) \in S, \quad (1.14)$$

where $U_m(\cos s) = \sin((m+1)s)/\sin s$ are the Chebyshev polynomials of the second kind. We have $R\tilde{\varphi}_\ell = b_\ell \tilde{\psi}_\ell$, with the singular values

$$b_\ell = \pi^{-1}(j+k+1)^{-\frac{1}{2}}, \ \ell = (j,k) \in \mathscr{L}. \quad (1.15)$$

Since we work with real functions, we identify the complex bases (1.13) and (1.14) with the equivalent real orthonormal bases $\{\varphi_\ell\}$, $\{\psi_\ell\}$ in a standard way,

$$\varphi_\ell = \begin{cases} \sqrt{2}\,\mathrm{Re}(\tilde{\varphi}_\ell) & \text{if } j > k, \\ \tilde{\varphi}_\ell & \text{if } j = k, \\ \sqrt{2}\,\mathrm{Im}(\tilde{\varphi}_\ell) & \text{if } j < k. \end{cases} \quad (1.16)$$

The problem of tomography in statistics is studied, for example, in [80, 83, 39, 21].

The model is the following

$$Y(u,s) = Rf(u,s) + \varepsilon\xi(u,s), \ (u,s) \in S,$$

where $\xi$ is a white noise in $G = L^2(S,d\mu)$.

The SVD basis is known for the Radon transform. However, this basis is very difficult to compute. By projection on $\{\psi_\ell\}$, one obtains the equivalent sequence space model,

$$y_\ell = b_\ell \theta_\ell + \varepsilon \, \xi_\ell, \ \ell = (j,k), \ j \geqslant 0, \ k \geqslant 0,$$

where $\theta_\ell = \langle f, \varphi_\ell \rangle$, and $\xi_\ell$ are i.i.d. standard Gaussian random variables.

*Remark 1.10.* In tomography, the problem is mildly ill-posed, since the singular values have a polynomial behaviour. The exact degree of ill-posedness is a bit different, since the problem is ill-posed, but is also a problem of estimation of a function in two

dimensions, which is known to be more difficult. One often considers that $\beta = 1/2$, due to (1.15).

There exist several models of tomography (X-rays tomography, positron emission tomography, discrete tomography, tomography in quantum physics and so on). The models each have their own specificities but are however all linked to the Radon operator.

## 1.1.7 Spectral Theory

In this section, we generalize the statistical study of inverse problems to the case of not only compact but also linear bounded operators. This extension is needed since there exist a lot of natural inverse problems where the operator is not compact, see for example the deconvolution on $\mathbb{R}$ in Section 1.1.7.2. In this situation, one needs other tools than the SVD. Moreover, the spectral Theorem and functional calculus may also be used in the case of compact operators.

### 1.1.7.1  The Spectral Theorem

The Halmos version of the spectral Theorem is convenient for the study of inverse problems (see [68]).

**Theorem 1.2.** *Let $A \in L(H)$ be a self-adjoint operator defined on a separable Hilbert space $H$. There exist a locally compact space $S$, a positive Borel measure $\Sigma$ on $S$, a unitary operator $U : H \rightarrow L^2(\Sigma)$, and a continuous function $\rho : S \rightarrow \mathbb{R}$ such that*

$$A = U^{-1} M_\rho U, \tag{1.17}$$

*where $M_\rho$ is the multiplication operator $M_\rho : L^2(\Sigma) \rightarrow L^2(\Sigma)$ defined $M_\rho \varphi = \rho \cdot \varphi$.*

*Proof.* A proof may be found in [125].

*Remark 1.11.* This fundamental result means that any self-adjoint linear bounded operator is similar to a multiplication in some $L^2$-space.

*Remark 1.12.* In the special case where $A$ is a compact operator, a well-known version of the spectral theorem (see Theorem 1.1) states that $A$ has a complete orthogonal system of eigenvectors $\{\varphi_k\}$ with corresponding eigenvalues $\rho_k$. This is a special case of (1.17) where $S = \mathbb{N}$, $\Sigma$ is the counting measure, $L^2(\Sigma) = \ell^2(\mathbb{N})$, and $\rho(k) = \rho_k$.

*Remark 1.13.* If $A$ is not self-adjoint then we use $A^*A$, where $A^*$ is the adjoint of $A$.

Using the spectral theorem one obtains the equivalent model to (1.2)

$$UY = U(Af + \varepsilon \xi) = UAf + \varepsilon U\xi = UU^{-1} M_\rho Uf + \varepsilon U\xi,$$

which gives in the spectral domain

$$Z = \rho \cdot \theta + \varepsilon \eta, \tag{1.18}$$

where $Z = UY$, $\theta = Uf$ and $\eta = U\xi$ is a white noise in $L^2(\Sigma)$ since $U$ is a unitary operator. Indeed, we have the following lemma.

**Lemma 1.5.** *Let $\xi$ be a white noise in $G$ and $\eta = U\xi$ where $U$ is a unitary operator. We have for all $\theta = Uf$ and $v = Uh$, where $f, h \in H$,*

$$\langle \eta, \theta \rangle \sim \mathcal{N}(0, \|\theta\|^2),$$

*and*

$$\mathbf{E}(\langle \eta, \theta \rangle \langle \eta, v \rangle) = \langle \theta, v \rangle.$$

*Thus $\eta$ is a white noise in $L^2(\Sigma)$.*

*Proof.* For all $\theta = Uf$ with $f \in H$, we have

$$\langle \eta, \theta \rangle = \langle U\xi, Uf \rangle = \langle \xi, U^*Uf \rangle = \langle \xi, f \rangle \sim \mathcal{N}(0, \|f\|^2) = \mathcal{N}(0, \|\theta\|^2). \tag{1.19}$$

In the same way if $\theta = Uf$ and $v = Uh$, where $f, h \in H$, we have

$$\mathbf{E}(\langle \eta, \theta \rangle \langle \eta, v \rangle) = \mathbf{E}(\langle \xi, f \rangle \langle \xi, h \rangle) = \langle f, h \rangle = \langle \theta, v \rangle, \tag{1.20}$$

since $U$ is unitary. Using (1.19) and (1.20), we obtain the lemma.

*Remark 1.14.* The model (1.18) really helps to understand the utility of spectral Theorem in inverse problems. Indeed, by use of the unitary transform $U$, one replaces the model (1.2), not always easy to handle with a general linear operator, by a multiplication by a function $\rho$. Moreover, since $U$ is unitary the noise is still a white noise.

### 1.1.7.2 Deconvolution on $\mathbb{R}$

In this section, we present an example of application, see for example [49, 116] when the operator is not compact. The operator considered here is

$$Af(t) = r * f(t) = \int_{-\infty}^{\infty} r(t - u)f(u)du,$$

where $r * f$ denotes the convolution through a known filter $r \in L^1(\mathbb{R})$. The aim is to reconstruct the unknown function $f$.

Deconvolution is one of the most standard inverse problems. The problem of circular convolution, i.e. with a periodic kernel $r$ on $[a, b]$, appears for example in [32] and in Section 1.1.6.3. The main difference is that for periodic convolution the operator is compact and the basis of eigenfunctions is the Fourier basis. It seems clear,

from a heuristic point of view, that the results could be extended to the case of convolution on $\mathbb{R}$ by using the Fourier transform on $L^2(\mathbb{R})$ instead of the Fourier series. This heuristic extension can be made formal by resorting to the spectral Theorem (Theorem 1.2).

Suppose that $r$ is a real-valued function symmetric around 0, then

$$\tilde{r}(\omega) = \int_{-\infty}^{\infty} e^{it\omega} r(t) dt = \int_{-\infty}^{\infty} \cos(t\omega) r(t) dt, \ \forall \omega \in \mathbb{R}.$$

Suppose also that $\tilde{r}(\omega) > 0$ for all $\omega \in \mathbb{R}$. It is straightforward to see that the operator $A$ is self-adjoint and strictly positive, since $r$ is real-valued and symmetric around 0 and $\tilde{r} > 0$.

Define the Fourier transform as a unitary operator from $L^2(\mathbb{R})$ into $L^2(\mathbb{R})$ by

$$(Ff)(\omega) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{it\omega} f(t) dt, \ \omega \in \mathbb{R}, \ f \in L^1(\mathbb{R}) \cap L^2(\mathbb{R}), \qquad (1.21)$$

and its continuous extension on $L^2(\mathbb{R})$.

We have that

$$F(r * f)(\omega) = \tilde{r}(\omega).(Ff)(\omega);$$

hence $A = F^{-1} M_{\tilde{r}} F$.

The model is then the following

$$Y(t) = r * f(t) + \varepsilon \xi(t), \ \forall t \in \mathbb{R}$$

where $f \in L^2(\mathbb{R})$ is unknown, $\xi$ is a white noise in $L^2(\mathbb{R})$.

By applying the Fourier transform we obtain

$$FY(\omega) = F(r * f)(\omega) + \varepsilon \ F\xi(\omega) = \tilde{r}(\omega).(Ff)(\omega) + \varepsilon \eta(\omega), \qquad (1.22)$$

where, by Lemma 1.5, $\eta$ is a white noise in $L^2(\mathbb{R})$.

### 1.1.7.3 Functional Calculus

In this section, the aim is to provide some important tools from operator theory linked to the spectral Theorem. Functional calculus is the main tool in order to modify the operator $A$, by applying functions to the operator. This result is crucial in the study of regularization methods. This section is based on [72].

**Definition 1.8.** Let $A \in L(H)$. The **resolvent** $\rho(A)$ of $A$ is the set of all $\lambda \in \mathbb{C}$ for which $(\lambda I - A)$ is invertible. The **spectrum** of $A$ is defined as $\sigma(A) = \mathbb{C} \setminus \rho(A)$.

Note that an eigenvalue is in the spectrum, but that not all points in the spectrum are eigenvalues. However, compact operators have spectra composed of eigenvalues.

It follows immediately that the spectrum is invariant by unitary transformations. Thus, $\sigma(A) = \sigma(M_\rho)$ with the notation of Theorem 1.2. One may also prove that $\sigma(A)$ is a closed and bounded set, hence is compact.

**Lemma 1.6.** *Let $A \in L(H)$ be self-adjoint.*

*1. For any $f \in C(S)$, i.e. a continuous function on S, we have*

$$\|M_f\| = \|f\|_{\infty,\mathrm{supp}\Sigma},$$

*where the norm is the restricted sup-norm on $\mathrm{supp}\Sigma$,*

$$\mathrm{supp}\Sigma = S \setminus \bigcup_{V\mathrm{open},\Sigma(V)=0} V.$$

*2. We have $\sigma(A) = \overline{\rho(\mathrm{supp}\Sigma)}$.*
*3. We have $\sigma(A) \subset \mathbb{R}$ and $\forall \lambda \in \mathbb{C}$*

$$\|(\lambda I - A)^{-1}\| \leqslant |\mathrm{Im}\lambda|^{-1}.$$

*4. Let*

$$m_- = \inf_{\|\varphi\|=1} \langle A\varphi, \varphi \rangle, \text{ and } m_+ = \sup_{\|\varphi\|=1} \langle A\varphi, \varphi \rangle,$$

*then $\sigma(A) \subset [m_-; m_+]$.*
*5. We have $\sigma(A^*A) \subset [0, \|A^*A\|]$.*

*Proof.* See [73].

If $p(\rho)$ is a polynomial in $\rho$ then it is natural to define

$$p(A) = \sum_{j=0}^{k} p_j A^j. \tag{1.23}$$

The next theorem generalizes the idea of applying continuous functions to the operator $A$ by just applying continuous functions to the spectrum, i.e. in $C(\sigma(A))$.

**Theorem 1.3.** *With the notation of Theorem 1.2, define*

$$\Phi(A) = U^{-1} M_{\Phi \circ \rho} U, \tag{1.24}$$

*for a continuous real-valued $\Phi \in C(\sigma(A))$, where $\Phi \circ \rho(\omega) = \Phi(\rho(\omega))$. Then $\Phi(A) \in L(H)$ is self-adjoint and satisfies (1.23) if $\Phi$ is polynomial. The mapping $\Phi \to \Phi(A)$ is called* **functional calculus** *at A and is an isometric algebra homomorphism from $C(\sigma(A))$ to $L(H)$, i.e., for all $\Phi, \Psi \in C(\sigma(A))$ and $\alpha, \beta \in \mathbb{R}$ we have*

$$(\alpha\Phi + \beta\Psi)(A) = \alpha\Phi(A) + \beta\Psi(A), \tag{1.25}$$

$$(\Phi\Psi)(A) = \Phi(A)\Psi(A), \tag{1.26}$$

$$\|\Phi(A)\| = \|\Phi\|_{\infty}. \tag{1.27}$$

*The functional calculus is uniquely determined by (1.23) and (1.25)-(1.27).*

*Proof.* This proof may be found in [72]. By Lemma 1.6, $\Phi(A)$ is bounded. It is self-adjoint because $\Phi$ is real-valued and

$$(\Phi(A))^* = U^{-1}(M_{\Phi\circ\rho})^* U = \Phi(A).$$

For a polynomial $p$ we have $p(A) = U^{-1}p(M_\rho)U$ with (1.23). Since by definition of the multiplication operator $p(M_\rho) = M_{p\circ\rho}$, this corresponds to definition (1.24). The proof of (1.25) is clear. Moreover, we have

$$\Phi(A)\Psi(A) = U^{-1}M_{\Phi\circ\rho}UU^{-1}M_{\Psi\circ\rho}U = U^{-1}M_{\Phi\Psi\circ\rho}U = (\Phi\Psi)(A).$$

Finally, by Lemma 1.6 and continuity of $\Phi$ we obtain

$$\|\Phi(A)\| = \|M_{\Phi\circ\rho}\| = \|\Phi\circ\rho\|_\infty = \|\Phi\|_{\infty,\sigma(A)}.$$

Let $\Phi_A : C(\sigma(A)) \to L(H)$ be any isometric algebra homomorphism satisfying $\Phi_A(p) = p(A)$ for all polynomials. By the Weierstrass approximation Theorem, for any $\Phi \in C(\sigma(A))$ there exists a sequence of polynomials $p_k$ such that $\|\Phi - p_k\|_{\infty,\sigma(A)} \to 0$ as $k \to \infty$. Using the property for the norm we obtain the desired unicity

$$\Phi_A(\Phi) = \lim_{k\to\infty} \Phi_A(p_k) = \lim_{k\to\infty} p_k(A) = \Phi(A).$$

*Remark 1.15.* The functional calculus may be extended to the case of bounded functions on $\sigma(A)$. The isometry is then replaced by an upper bound in (1.27).

*Remark 1.16.* This theorem is a fundamental tool in analysis of inverse problems. It allows to apply functions to the operator and then to its spectrum. The aim is to study the behaviour of $A^{-1}$ or of more stable inverses (e.g. the regularized inverse).

*Example 1.2. Regularization methods.* Suppose that $A \in L(H)$ is self-adjoint and positive. In ill-posed problems, if $A$ is not invertible, then 0 is in the spectrum. Another way to understand this is by Theorem 1.3. Indeed, when 0 is in the spectrum then the function $\Phi(x) = 1/x$ is not even bounded on $\sigma(A) \subset [0, \|A\|]$, it explodes at point 0. Thus $\Phi(A) = A^{-1}$ is not a bounded operator by (1.27). There is a need to invert $A$ in a more stable way. This is exactly the role of regularization methods. One way to invert $A$, is by a small modification of the function $\Phi$. For example, one may use $\Phi_\gamma(x) = 1/(x+\gamma)$ where $\gamma > 0$, which is continuous and bounded on $\sigma(A)$. This is exactly the idea of the Tikhonov regularization method, see Section 1.2.2.

## 1.2 Nonparametric Estimation

The aim of nonparametric estimation is to estimate (reconstruct) a function $f$ (density or regression funtion) based on some observations. The main difference with parametric statistics is that the function $f$ is not in some parametric family of functions, for example, the family of Gaussian probability density functions $\{\mathcal{N}(\mu, 1),\ \mu \in \mathbb{R}\}$.

Instead of a general framework, the problem of nonparametric estimation will be described here in the setting of the sequence space model (1.5) which is related to the inverse problem with random noise (1.2).

### 1.2.1 Minimax Approach

Let $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \ldots)$ be an estimator of $\theta = (\theta_1, \theta_2, \ldots)$ based on the data $X = \{X_k\}$. An estimator of $\theta$ may be any measurable function of the observation $X = \{X_k\}$.

Then $f$ is estimated by $\hat{f} = \sum_k \hat{\theta}_k \varphi_k$, where $\{\varphi_k\}$ is a basis.

The first point is to define the accuracy of some given estimator $\hat{\theta}$. Since an estimator is by definition random, we will measure the squared difference between $\hat{\theta}$ and the true $\theta$, and then take the mathematical expectation.

Define the **mean integrated squared risk (MISE)** of $\hat{f}$ by

$$\mathscr{R}(\hat{f}, f) = \mathbf{E}_f \|\hat{f} - f\|^2 = \mathbf{E}_\theta \sum_{k=1}^{\infty} (\hat{\theta}_k - \theta_k)^2 = \mathbf{E}_\theta \|\hat{\theta} - \theta\|^2,$$

where the second equality follows from Parseval's Theorem (and relies on the fact that $\{\varphi_k\}$ is a basis), where the notation $\|\cdot\|$ stands for $\ell^2$-norm of $\theta$-vectors in the sequence space. Here and in the sequel $\mathbf{E}_f$ and $\mathbf{E}_\theta$ denote the expectations w.r.t. $Y$ or $X = (X_1, X_2, \ldots)$ for models (1.2) and (1.5) respectively. Analyzing the risk $\mathscr{R}(\hat{f}, f)$ of the estimator $\hat{f}$ is equivalent to analyze the corresponding sequence space risk

$$\mathscr{R}(\hat{\theta}, \theta) = \mathbf{E}_\theta \|\hat{\theta} - \theta\|^2.$$

The aim would be to find the estimator with the minimum risk. However, the risk of an estimator depends, by definition, on the unknown $f$ or $\theta$.

To that end, we assume that $f$ belongs to some class of function $\mathscr{F}$.

**Definition 1.9.** Define the **maximal risk** of the estimator $\hat{f}$ on $\mathscr{F}$ as

$$\sup_{f \in \mathscr{F}} \mathscr{R}(\hat{f}, f),$$

and the **minimax risk** as

$$r_\varepsilon(\mathscr{F}) = \inf_{\hat{f}} \sup_{f \in \mathscr{F}} \mathscr{R}(\hat{f}, f),$$

where the $\inf_{\hat{f}}$ is taken over all possible estimators of $f$.

It is usually not possible in nonparametric statistics to find estimators which attain the minimax risk. A more natural approach is to consider the asymptotic properties, i.e. when the noise level tends to 0 ($\varepsilon \to 0$).

**Definition 1.10.** Suppose that some estimator $\tilde{f}$ is such that there exist constants $0 < C_2 \leqslant C_1 < \infty$ with, as $\varepsilon \to 0$

$$\sup_{f \in \mathscr{F}} \mathscr{R}(\tilde{f}, f) \leqslant C_1 v_\varepsilon^2,$$

where the positive sequence $v_\varepsilon$ is such that $v_\varepsilon \to 0$ and

$$\inf_{\hat{f}} \sup_{f \in \mathscr{F}} \mathscr{R}(\hat{f}, f) \geqslant C_2 v_\varepsilon^2.$$

In this case the estimator $\tilde{f}$ is said to be **optimal** or to **attain the optimal rate of convergence** $v_\varepsilon^2$.

In the special case where $C_1 = C_2$ the estimator $\tilde{f}$ is said to be **minimax** or to **attain the exact constant**.

An optimal estimator is then an estimator whose risk is of the order of the best possible estimator.

Minimax estimation in nonparametric statistics is nowadays a classical approach. It goes back to [11, 122] and also [74]. Since the 80's, it has been obtained in many different models (nonparametric regression, Gaussian white noise, density estimation, spectral density estimation), with a varied form of estimators (kernels, projections, splines, wavelets), and for most classes of functions (Besov, Hölder, Sobolev, . . . ).

These kind of results are often considered as a first step in order to prove that a given method has good theoretical properties. Indeed, one has a criterion, which garantees that on some class of functions a given method is optimal.

Minimax estimation for statistical inverse problems (1.2) (or for its sequence space analogue (1.5)) was discussed in a number of papers. Optimal rates of convergence in this problem are obtained in [80, 84, 39, 83, 95, 47, 78, 9] and in related frameworks in [52, 19].

Exact asymptotics of the minimax $L^2$-risks are known in the deconvolution problem with somewhat different setups [50], in inverse problems for partial differential equations [61] and in tomography, for minimax $L^2$-risks among linear estimators [80]. Exact asymptotics for pointwise risks on the classes of analytic functions in tomography are due to [20].

A recent discussion of the different rates may be found in the review [93].

*Remark 1.17.* We only use the standard $L^2-$risk along these notes. However, many results may be obtained with different risks and loss functions, for example, the $L^p$, $L^\infty$ or the pointwise risk, see [43, 85, 20, 79].

## 1.2.2 Regularization Methods

### 1.2.2.1 Continuous Regularization Methods

The main part of ill-posed inverse problems is to find regularization methods which will help to get a fine reconstruction of $f$.

Recall that the normal equation, defined in Lemma 1.4, is

$$A^*Y = A^*Af + \varepsilon\, A^*\xi.$$

Formally, one has to estimate the solution $(A^*A)^{-1}A^*Y$. The problem in ill-posed situation is that the operator $A^*A$ is not (boundedly) invertible.

The idea is to get some continuous inversion by use of regularization methods. This allows to obtain much more stable reconstruction.

**Definition 1.11.** We call a **regularization method** an estimator defined by

$$\hat{f}_\gamma = \Phi_\gamma(A^*A)A^*Y,$$

where $\Phi_\gamma \in C(\sigma(A^*A))$, i.e a continuous function on $\sigma(A^*A)$ (or even bounded) depending on some **regularization parameter** $\gamma > 0$.

We are going to give some examples of regularization methods or estimators which are commonly used. All these methods are defined in the spectral domain even if some of them may be computed without using the whole spectrum.

### Spectral Cut-Off

This regularization method is very simple. The idea is to get rid of the high frequencies. In the spectral domain, by using the spectral cut-off, one just cut the frequencies over some threshold.

The definition of a **spectral cut-off** with parameter $\gamma > 0$ is the following

$$\Phi_\gamma(x) = \begin{cases} x^{-1}, & x \geqslant \gamma, \\ 0, & x < \gamma. \end{cases}$$

This notion may be well-defined by use of the functional calculus for bounded functions (instead of continuous ones).

The spectral cut-off is a very simple estimator. It is usually used as a benchmark since it attains the optimal rate of convergence. However, it is not a very precise estimator. Moreover, from a numerical point of view, it is usually time consuming since one has to compute the whole spectrum.

### Tikhonov Regularization

The Tikhonov method is one of the first and the most well-known regularization method in inverse problems.

The direct inversion of the operator $A^*A$ is not satisfying since it is not a (boundedly) invertible operator. The idea is to control the norm of the solution by using a penalty term.

Define now, the well-known **Tikhonov regularization method** (see [127]). In this method one wants to minimize the following functional $L_\gamma(\varphi)$:

$$\inf_{\varphi \in H} \left\{ \|A\varphi - Y\|^2 + \gamma \|\varphi\|^2 \right\}, \qquad (1.28)$$

where $\gamma > 0$ is some tuning parameter.

The Tikhonov method is very natural. Indeed, the idea is to choose an estimator which, due to the first term will fit the data, and which will be "stable", due to the second term, which is called the energy. As we will see in Section 1.3.3 the choice of $\gamma$ is very sensitive since it characterizes the balance between the fitting and the stability.

The functional $L_\gamma$ is strictly convex for any $\gamma > 0$. Its minimum is attained when its differential in $h \in H$

$$(L_\gamma)'_\varphi h = 2\langle A\varphi - Y, Ah \rangle + 2\gamma \langle \varphi, h \rangle, \qquad (1.29)$$

is zero, i.e.

$$\langle A^*(Y - A\varphi), h \rangle = \gamma \langle \varphi, h \rangle, \ \forall h \in H.$$

The minimum is then attained by

$$\hat{f}_\gamma = (A^*A + \gamma I)^{-1} A^* Y. \qquad (1.30)$$

In the spectral domain this method is defined by

$$\Phi_\gamma(x) = \frac{1}{x + \gamma}.$$

*Remark 1.18.* There exist some troubles with this simple Tikhonov regularization. For these reasons, several modifications of the Tikhonov have been defined.

**Variants of Tikhonov Regularization**

There exist several modified versions of the Tikhonov regularization.

The first variant is the **Tikhonov method with starting point** $\varphi_0$. It consists in giving a different starting point than 0. We have $\varphi_0 \in H$ and we penalize by $\|\varphi - \varphi_0\|$ instead of $\|\varphi\|$. By use of (1.29) one then obtains

$$\hat{f}_\gamma = (A^*A + \gamma I)^{-1} (A^* Y + \gamma \varphi_0). \qquad (1.31)$$

The second modified version, which already appears in [127], is the **Tikhonov method with a different prior**. It is is based on the idea that the function could be smoother. Thus, a penalty term of the form $\|Q^a \varphi\|$, where $Q^a, a > 0$, is some differential operator, would be more suited. A classical example is then $Q^a = (A^*A)^{-a}$. The estimator is then defined by

$$\hat{f}_\gamma = (A^*A + \gamma (Q^a)^* Q^a)^{-1} A^* Y.$$

With this method one is able to better estimate smoother functions. Indeed, the standard Tikhonov method, penalize only by $\|\varphi\|^2$. If the function is smoother, then it is natural to take this into account in the second term, by a smoothness constraint. This effect may be seen in the better qualification of the method (see Section 1.2.3.1).

A last variant is called **iterative Tikhonov method**. It consists in starting a first Tikhonov regularization with $\varphi_0 = 0$ and then obtain an estimator $\hat{f}_\gamma^1$. In the second iteration, one applies the Tikhonov method with a starting point $\hat{f}_\gamma^1$. We iterate this method several times. The estimate is then by (1.31)

$$\hat{f}_{m+1} = (A^*A + \gamma I)^{-1}(A^*Y + \gamma \hat{f}_m).$$

It may be shown by induction that

$$\hat{f}_m = (A^*A + \gamma I)^{-m}(A^*A)^{-1}((A^*A + \gamma I)^m - \gamma^m I)A^*Y.$$

For $m = 1$ this corresponds exactly to the standard Tikhonov regularization. With this method one is able to better estimate smoother functions, the qualification of the method is increased (see Section 1.2.3.1). From a numerical point of view, this method is not really much longer than the Tikhonov one, since the only operator to invert is $(A^*A + \gamma I)$.

In the spectral domain this method is defined by

$$\Phi_\gamma(x) = \frac{(x+\gamma)^m - \gamma^m}{x(x+\gamma)^m}.$$

**Landweber Iteration**

Another very standard method is based on the idea to minimize the functional $\|A\varphi - Y\|$ by the steepest descent method (i.e. Gradient descent algorithm). The idea then is to choose the direction $h$ equals to minus the gradient (in fact here the approximate gradient). Thus, we obtain $h = -A^*(A\varphi - Y)$ by (1.29). This leads to the recursion formula $\varphi_0 = \hat{f}_0 = 0$ and

$$\hat{f}_m = \hat{f}_{m-1} - \mu A^*(A\hat{f}_{m-1} - Y),$$

for some $\mu > 0$. This method is called **Landweber iteration**, see for example in [86].

It may be shown by induction that

$$\hat{f}_m = \sum_{j=0}^{m-1} (I - \mu A^*A)^j \mu A^*Y.$$

Indeed, it is clearly true for $m = 0$ and if true for $m$ then

$$\hat{f}_{m+1} = (I - \mu A^*A)\hat{f}_m + \mu A^*Y = \sum_{j=0}^{m}(I - \mu A^*A)^j \mu A^*Y.$$

The parameter $\mu$ has to be chosen such that $\mu \|A^*A\| \leqslant 1$ which has a strong influence on the speed of convergence. The regularization parameter is then linked to the number of iterations $m$. Formally, the number of iterations may be written as $\gamma^{-1}$.

In the spectral domain this method is defined, for $\mu = 1$ and $\|A^*A\| = 1$, by

$$\Phi_m(x) = \sum_{j=0}^{m-1}(1-x)^j.$$

There exist another version of this formula which will be used later. We have

$$\Phi_m(x) = \frac{1-(1-x)^m}{x}, \quad (\Phi_m(0) = m).$$

*Remark 1.19.* From a numerical point of view, this method is faster than Tikhonov method, since one does not need here the inversion of an operator (as in (1.30)).

However, Landweber iteration has some drawbacks as we will see in Section 1.2.3.1. Indeed, the number of iterations may be very large. For this reason, new methods, based on Landweber have been defined, as the semi-iterative procedures and $\nu$-methods.

## Semi-iterative Procedures and $\nu$-Methods

As we will see the Landweber iteration is not so efficient. One of the reason is, that this method uses only the previous iteration in order to compute the next one. A more general idea is to use all the previous iterations $\hat{f}_j$, $j = 1, \ldots, m-1$, to define $\hat{f}_m$.

This is the starting point of the so-called **semi-iterative procedures**. Let $\hat{f}_j$, $j = 1, \ldots, m-1$, and $\hat{f}_0 = 0$ then define

$$\hat{f}_m = \mu_{1,m}\hat{f}_{m-1} + \cdots + \mu_{m,m}\hat{f}_0 + \omega_m A^*(Y - A\hat{f}_{m-1}),$$

where $\sum_j \mu_{j,m} = 1$.

The semi-iterative methods are then defined by

$$\hat{f}_m = \Phi_m(A^*A)A^*Y,$$

where $\Phi_m$ is a polynomial of degree exactly $m-1$, which is called iteration polynomial.

Clearly, such a method is computationaly rather efficient but we use all the iterations and not only one.

A special case of such iterative method are the $\nu$-**methods** which only use two iterations. These methods were introduced in [10] and in the statistical literature

by [106]. It is defined as a semi-iterative procedure with a parameter $v > 0$,

$$\mu_1 = 1, \ \omega_1 = \frac{4v+2}{4v+1},$$

$$\mu_m = 1 + \frac{(m-1)(2m-3)(2m+2v-1)}{(m+2v-1)(2m+4v-1)(2m+2v-3)},$$

$$\omega_m = 4\frac{(2m+2v-1)(m+v-1)}{(m+2v-1)(2m+4v-1)},$$

and

$$\hat{f}_m = \mu_m \hat{f}_{m-1} + (1-\mu_m)\hat{f}_{m-2} + \omega_m A^*(Y - A\hat{f}_{m-1}).$$

*Remark 1.20.* We will see in Section 1.2.3.1 that $v$-methods, and many semi-iterative methods, are much faster than the Landweber method. We will explain, what is the idea behind these $v$-methods.

### Risk of Regularization Methods

A regularization method defined by $\Phi_\gamma$ may be decomposed as

$$\hat{f}_\gamma = \Phi_\gamma(A^*A)A^*Af + \varepsilon\Phi_\gamma(A^*A)A^*\xi, \tag{1.32}$$

since (1.2). Its risk may then be written as

$$\mathbf{E}_f\|\hat{f}_\gamma - f\|^2 = \|\mathbf{E}_f(\hat{f}_\gamma) - f\|^2 + \mathbf{E}_f\|\hat{f}_\gamma - \mathbf{E}_f(\hat{f}_\gamma)\|^2,$$

since $\mathbf{E}\Phi_\gamma(A^*A)A^*\xi = 0$. The first term is called the **approximation error** and the second is called **propagated noise error**.

Remark that by Theorem 1.3 we have

$$\Phi_\gamma(A^*A)A^*Af = U^{-1}M_{\Phi_\gamma(\rho)\rho}Uf.$$

Thus, $\Phi_\gamma(A^*A)$ should be an approximate inverse of $A^*A$. The study of the function $\Phi_\gamma(x)x$ in the spectral domain is then of major importance.

*Remark 1.21.* As for the estimation method, a key-point in regularization methods is to choose the parameter $\gamma$ in a proper way.

### 1.2.2.2  Estimation Procedures

### Equivalence in the Sequence Space Model

In order to get a framework more standard in statistics, suppose now that the operator $A$ is compact. Then, by using the SVD, one obtains the sequence space model (1.5). Usually statisticians prefer to work with the sequence space model (1.6).

In this context, many regularization methods may be expressed in a statistical framework, and usually correspond to some known estimation method in statistics. The notion of regularization is not really used in statistics. However, there exists a more standard definition which is related.

Let $\lambda = (\lambda_1, \lambda_2, \ldots)$ be a sequence of nonrandom weights. Every sequence $\lambda$ defines a **linear estimator** $\hat{\theta}(\lambda) = (\hat{\theta}_1, \hat{\theta}_2, \ldots)$ where

$$\hat{\theta}_k = \lambda_k X_k \text{ and } \hat{f}(\lambda) = \sum_{k=1}^{\infty} \hat{\theta}_k \, \varphi_k.$$

Remark also by use of the SVD in Theorem 1.1 or Theorem 1.3 one obtains for a general regularization method

$$\hat{f}(\lambda) = \Phi_\gamma(A^*A)A^*Y = \sum_{k=1}^{\infty} \Phi_\gamma(b_k^2)b_k y_k \varphi_k = \sum_{k=1}^{\infty} \Phi_\gamma(b_k^2)b_k^2 \, X_k \varphi_k,$$

which exactly corresponds to the special case of linear estimator with

$$\lambda_k = \Phi_\gamma(b_k^2)b_k^2. \tag{1.33}$$

### Truncated SVD

Examples of commonly used weights $\lambda_k$ are the projection weights $\lambda_k = I(k \leqslant N)$ where $I(\cdot)$ denotes the indicator function. These weights correspond to the **projection estimator** (also called **truncated SVD**).

$$\hat{\theta}(N) = \begin{cases} X_k, \ k \leqslant N, \\ 0, \ \ k > N. \end{cases}$$

The value $N$ is called the **bandwidth**.

The projection estimator is then defined by

$$\hat{f}_N = \sum_{k=1}^{N} X_k \varphi_k.$$

The truncated SVD is a very simple estimator. With this natural estimator, one estimates the first $N$ coefficients $\theta_k$ by their empirical counter-part $X_k$ and then estimate the remainder terms by 0 for $k > N$.

This is an estimator equivalent to the spectral cut-off, but expressed in a different way and in a different setting. From a numerical point of view, it is still usually time consuming since, one has to compute all the coefficients $X_k$.

One may easily check by using (1.33) that, for the case $\sigma_k = k^\beta$, the spectral cut-off is equivalent to the projection estimator with $N = [\gamma^{-1/2\beta}]$.

## Kernel Estimator

One of the most well-known method in statistics is the **kernel estimator** (see [115, 129]). In our context, kernel estimator could be defined in the special case of the direct problem; i.e. $A = I$. A kernel estimator is defined by its kernel function $K \in L^2$ (usually also $K \in L^1$) and

$$\hat{f}_\gamma = K_\gamma * Y,$$

where $*$ denotes the convolution product, $K_\gamma(\cdot) = \gamma^{-1} K(\cdot/\gamma)$ and $\gamma > 0$ is known as the bandwidth.

The idea of kernel estimators is to estimate the function $f$ by using a local (by the bandwidth) weighted mean of the data, i.e. a convolution.

Kernel estimators may also be defined in inverse problem framework in order to invert the operator, see for example the so-called deconvolution kernel in [52].

This method is also linked to the mollifier methods in inverse problems, see [94].

## The Tikhonov Estimator

The **Tikhonov estimator** is defined by the same minimization in (1.28) as for the Tikhonov regularization. In a more statistical framework, one may define the Tikhonov estimator by its equivalent form in the SVD domain:

$$\lambda_k = \frac{1}{1 + \gamma \sigma_k^2},$$

which is easy to verify by use of (1.33).

In the special case where $A = I$ this estimator is defined and computed as a modified version of the Tikhonov regularization and is called spline (see [131]).

In the parametric context of the standard linear regression, this method is called ridge regression, see [70]. It is known to improve on the standard least-squares estimator when the singular values of the design matrix are close to 0.

## The Landweber Method

The Landwber iteration is not really known under this name in statistics. However, there exists a well-known approach in the community of learning which is strongly related.

Boosting algorithms include a family of iterative procedures which improve the performance at each step. The $L^2$-boosting has been introduced in the context of regression and classification in [15].

The idea is to start from a weak learner, i.e. a rather rough estimator $\hat{f}_0$. The algorithm consists then in boosting this learner in a recursive iteration, which may be showed to correspond to Landweber iteration (see [9]).

**The Pinsker Estimator**

The Pinsker estimator has been defined in [109]. This special class of linear estimators is defined in the sequence space model by the following weights coefficients

$$\lambda_k = (1 - c_\varepsilon a_k)_+,$$

where $c_\varepsilon$ is the solution of the equation

$$\varepsilon^2 \sum_{k=1}^{\infty} \sigma_k^2 a_k (1 - c_\varepsilon a_k)_+ = c_\varepsilon L,$$

with $x_+ = \max(0, x)$ and $a_k > 0$.

   As we will see in Section 1.2.4, this class of estimators is defined in the context of estimation in ellipsoids where they attain the optimal rates of convergence, but also the minimax constants.

**Risk of a Linear Estimator**

Define now the $L^2$−risk of linear estimators :

$$\mathscr{R}(\hat{\theta}(\lambda), f) = R(\theta, \lambda) = \mathbf{E}_\theta \sum_{k=1}^{\infty} (\hat{\theta}_k(\lambda) - \theta_k)^2 = \sum_{k=1}^{\infty} (1 - \lambda_k)^2 \theta_k^2 + \varepsilon^2 \sum_{k=1}^{\infty} \sigma_k^2 \lambda_k^2.$$

(1.34)

The first term in the RHS is called **bias term** and the second term is called the **stochastic term** or **variance term**. The bias term is linked to the approximation error and measure if the chosen estimation procedure is a good approximation of the unknown $f$. On the other hand, the stochastic term measure the influence of the random noise and of the inverse problem in the accuracy of the method.

   In these lectures, we are going to study in details the projection estimators. This method is the most simple one and can be studied in a very easy way. The risk of a projection estimator with bandwidth $N$ is

$$R(\theta, N) = \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2.$$

In this case the decomposition is very simple. Indeed, we estimate the first $N$ coefficients by their empirical version $X_k$ and the other coefficients by 0. Thus, the bias term measure the influence of the remainder coefficients $\theta_k$, $k > N$, and the stochastic term is due to the random noise in the $N$ first coefficients. We can see now that one simple question is how to choose the bandwidth $N$?

*Remark 1.22.* Thus, we get to the key-point in nonparametric statistics. We have to choose $N$ (or $\gamma$ or $m$) in order to balance the bias term and the variance term. As we will see this choice will be difficult since the bias term depends on the unknown $f$.

## *1.2.3 Classes of Functions*

An important problem now is to define "natural" classes of functions on $\mathscr{F}$.

### 1.2.3.1 Source Conditions

A standard way to measure the smoothness of the function $f$ is relative to the smoothing properties of the operator $A$, more precisely in terms of $A^*A$. Let $\ell :$ $[0,\infty) \to [0,\infty)$ be a continuous, strictly increasing function with $\ell(0) = 0$ and assume that there exists a source $w \in H$ such that

$$f = \ell(A^*A)w, \; w \in H, \; \|w\|^2 \leqslant L. \tag{1.35}$$

This is called a **source condition**. The most standard choice for $\ell$ is the **Hölder type source condition** where $\ell(x) = x^\mu$, $\mu \geqslant 0$, i.e.

$$f = (A^*A)^\mu w, \; w \in H, \; \|w\|^2 \leqslant L. \tag{1.36}$$

Denote by $\mathscr{F}_\ell(L)$ the class of functions

$$\mathscr{F}_\ell(L) = \left\{ f = \ell(A^*A)w : w \in H, \; \|w\|^2 \leqslant L \right\}. \tag{1.37}$$

In order to take advantage of the source condition we assume that, for any regularization methods, there exists a constant $\nu_0$ called **qualification** and a constant $\overline{\nu}$ such that

$$\sup_{x \in \sigma(A^*A)} |x^\nu (1 - x\Phi_\gamma(x))| \leqslant \overline{\nu}\, \gamma^\nu, \; \forall\, \gamma > 0, \; \forall 0 \leqslant \nu \leqslant \nu_0. \tag{1.38}$$

We then get the following theorem in order to control the bias term, i.e. the approximation error.

**Theorem 1.4.** *Suppose that one has a regularization method $\hat{f}_\gamma$ checking condition (1.38). Define $\mathscr{F}_\ell(L)$ with $\ell(x) = x^\mu$. Then we have*

$$\sup_{f \in \mathscr{F}_\ell(L)} \|\mathbf{E}_f(\hat{f}_\gamma) - f\|^2 \leqslant \overline{\nu}^2 L\, \gamma^{2\mu},$$

*for all $0 \leqslant \mu \leqslant \nu_0$.*

*Proof.* We have

$$B(\hat{f}_\gamma) = \|\mathbf{E}_f(\hat{f}_\gamma) - f\|^2 = \|\Phi_\gamma(A^*A)A^*Af - f\|^2.$$

Using (1.35), (1.38) and the isometry of the functional calculus we obtain

$$B(\hat{f}_\gamma) = \|(\Phi_\gamma(A^*A)A^*A - I)(A^*A)^\mu w\|^2$$

$$\leqslant L \sup_{x\in\sigma(A^*A)} |x^\mu(1-x\Phi_\gamma(x))|^2 \leqslant \bar{v}^2 L\,\gamma^{2\mu}.$$

The qualification of a method is the largest source condition for which the bias of the method converges with the optimal rate.

For the Landweber iteration, suppose here that $\|A^*A\|=1$. Note that $\sigma(A^*A)\subset [0,1]$. The approximation error is

$$\sup_{x\in\sigma(A^*A)} |x^\mu(1-x\Phi_m(x))| \leqslant \sup_{x\in[0,1]} |x^\mu(1-x)^m|.$$

If we solve this problem, we then obtain that the supremum is attained at point $x_0 = \mu/(\mu+m) \in [0,1]$. Thus, the approximation error is bounded by, for any $\mu > 0$,

$$\sup_{x\in[0,1]} |x^\mu(1-x)^m| \leqslant \left(\frac{\mu}{\mu+m}\right)^\mu \leqslant Cm^{-\mu}.$$

The qualification of the Landweber method is then $\infty$, since this result is valid for any $\mu > 0$. Note that $\gamma = 1/m$ here.

The semi-iterative methods are defined via an iteration polynomial $\Phi_m$. The $v-$methods have, in fact, been defined such that they minimize

$$\sup_{x\in\sigma(A^*A)} |x^v(1-x\Phi_m(x))|,$$

for all polynomials of degree $m-1$.

One may prove then (see [49])

$$|1-x\Phi_m(x)| \leqslant c_v (1+m^2x)^{-v}.$$

Thus, we have

$$\sup_{x\in\sigma(A^*A)} |x^\mu(1-x\Phi_m(x))| \leqslant c_v \sup_{x\in[0,1]} |x^\mu(1+m^2x)^{-v}|.$$

The maximum is attained at point

$$x_0 = \begin{cases} \mu/(m^2(v-\mu)) & \text{if } \mu < v, \\ 1 & \text{if } \mu \geqslant v. \end{cases}$$

We finally obtain

$$\sup_{x\in\sigma(A^*A)} |x^\mu(1-x\Phi_m(x))| \leqslant \begin{cases} Cm^{-2\mu} & \text{if } \mu < v, \\ Cm^{-2v} & \text{if } \mu \geqslant v. \end{cases}$$

There is a saturation effect. The qualification of the $v-$method is then $v$.

*Remark 1.23.* One very important point here, is that, for the same number of iterations, the approximation error is much better for the $v-$method than for the Landweber one. In other terms, one needs $m^2$ iterations with Landweber and $m$ iterations

with $v-$method for the same accuracy. The Landweber method attains the optimal rates of convergence but with much more iterations than $v-$method. This is one drawback of the Landweber method in applications.

Some direct computations show that the qualification of the different methods are (see following table).

**Table 1.1** Qualification of regularization methods

| Method | Qualification |
|---|---|
| Spectral cut-off | $\infty$ |
| Tikhonov | 1 |
| Tikhonov with prior $a$ | $1 + 2a$ |
| $m$-iterated Tikhonov | $m$ |
| Landweber | $\infty$ |
| $v$-method | $v$ |

The main aim now is to understand the precise meaning of source condition on some well-known examples.

### 1.2.3.2 Ellipsoid of Coefficients

Suppose here that the operator $A$ is compact.

Assuming Hölder type source condition $f = (A^*A)^\mu w$ is then equivalent in the SVD domain to, by functional calculus,

$$f = (A^*A)^\mu w = \sum_{k=1}^{\infty} b_k^{2\mu} w_k \varphi_k,$$

since $w \in H$, where $w_k = \langle w, \varphi_k \rangle$ is in $\ell^2$. Denote by $\langle f, \varphi_k \rangle = \theta_k = b_k^{2\mu} w_k$ and since $\|w\|^2 \leqslant L$, we then obtain

$$\|w\|^2 = \sum_{k=1}^{\infty} w_k^2 = \sum_{k=1}^{\infty} b_k^{-4\mu} \theta_k^2 \leqslant L. \tag{1.39}$$

Thus, in the inverse problem framework with compact operator, the source conditions will correspond to the assumption that the coefficients of $f$ belong to some ellipsoid in $\ell^2$.

Assume that $f$ belongs to the functional class corresponding to ellipsoids $\Theta$ in the space of coefficients $\{\theta_k\}$:

$$\Theta = \Theta(a, L) = \left\{ \theta : \sum_{k=1}^{\infty} a_k^2 \theta_k^2 \leqslant L \right\}, \tag{1.40}$$

where $a = \{a_k\}$ is a non-negative sequence that tends to infinity with $k$, and $L > 0$. This means that for large values of $k$ the coefficients $\theta_k$ will have (a negative) polynomial behaviour in $k$ and will be small.

*Remark 1.24.* Assumptions on the coefficients $\theta_k$ will be usually related to some properties (smoothness) on $f$. One difficulty in using SVD in inverse problems is that the basis $\{\varphi_k\}$ is defined by the operator $A$. One then has to hope good properties for the coefficients $\theta_k$ of $f$ in this specific basis.

### 1.2.3.3 Classes of Functions

Suppose that we are exactly in the setting of the periodic convolution of Section 1.1.6.3. Then the operator is compact and the SVD basis is exactly the Fourier basis.

In the special cases where the SVD basis is the Fourier basis, hypothesis on $\{\theta_k\}$ may be precisely written in terms of smoothness for $f$.

Such classes arise naturally in various inverse problems, they include as special cases the Sobolev classes and classes of analytic functions. In fact, we consider balls of size $L > 0$ in functions spaces.

Let $\{\varphi_k(t)\}$ be the real trigonometric basis on $[0,1]$:

$$\varphi_1(t) \equiv 1, \quad \varphi_{2k}(t) = \sqrt{2}\cos(2\pi kt), \quad \varphi_{2k+1}(t) = \sqrt{2}\sin(2\pi kt), \quad k = 1,2,\ldots.$$

Introduce the **Sobolev classes of functions** (see [12])

$$\mathscr{W}(\alpha,L) = \left\{ f = \sum_{k=1}^{\infty} \theta_k \varphi_k : \theta \in \Theta(\alpha,L) \right\}$$

where $\Theta(\alpha,L)$ with the sequence $a = \{a_k\}$ such that $a_1 = 0$ and

$$a_k = \begin{cases} (k-1)^{\alpha} & \text{for } k \text{ odd,} \\ k^{\alpha} & \text{for } k \text{ even,} \end{cases} \quad k = 2,3,\ldots,$$

where $\alpha > 0, L > 0$.

If $\alpha$ is an integer, this corresponds to the equivalent definition, see Proposition 1.14 in [129],

$$\mathscr{W}(\alpha,L) = \left\{ f \in L^2[0,1] : \int_0^1 (f^{(\alpha)}(t))^2 dt \leqslant L, f^{(j)}(0) = f^{(j)}(1) = 0, j = 0,\ldots,\alpha-1 \right\}$$

where $f$ is 1-periodic and $f^{(\alpha)}$ denotes the weak derivative of $f$ of order $\alpha$.

In the case where the problem is mildly ill-posed with $b_k = k^{-\beta}$, by (1.39), Hölder type source conditions correspond to

$$\sum_{k=1}^{\infty} b_k^{-4\mu} \theta_k^2 = \sum_{k=1}^{\infty} k^{4\mu\beta} \theta_k^2 \leqslant L,$$

and then $\theta \in \Theta(\alpha, L)$ with $\alpha = 2\mu\beta$.

One may also consider more restrictive conditions and the classes of functions

$$\mathscr{A}(\alpha, L) = \left\{ f = \sum_{k=1}^{\infty} \theta_k \varphi_k : \theta \in \Theta_{\mathscr{A}}(\alpha, L) \right\}$$

where $a_k = \exp(\alpha k)$, $\alpha > 0$, and $L > 0$. This corresponds to the usual classes of **analytical functions**. These functions admit an analytical continuation into a band of the complex plane, see [75]. These functions are thus very smooth ($C^{\infty}$).

In the case where the problem is severely ill-posed with $b_k = \exp(-\beta k)$, by (1.39), Hölder type source conditions correspond to

$$\sum_{k=1}^{\infty} b_k^{-4\mu} \theta_k^2 = \sum_{k=1}^{\infty} e^{4\mu\beta k} \theta_k^2 \leqslant L,$$

and then $\theta \in \Theta_{\mathscr{A}}(\alpha, L)$ with $\alpha = 2\mu\beta$.

### *1.2.4 Rates of Convergence*

#### 1.2.4.1 SVD Setting

In this setting of ill-posed inverse problems with compact operator and functions with coefficients in some ellipsoid, several results have been obtained.

As in Definition 1.9, denote by

$$r_{\varepsilon}(\Theta) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathscr{R}(\hat{\theta}, \theta), \tag{1.41}$$

where the $\inf_{\hat{\theta}}$ is taken for all estimators of $f$, the *minimax risk* on the class of coefficients $\Theta$ and the *linear minimax risk*

$$r_{\varepsilon}^{\ell}(\Theta) = \inf_{\hat{\theta}^{\ell}} \sup_{\theta \in \Theta} \mathscr{R}(\hat{\theta}, \theta),$$

where the $\inf_{\hat{\theta}^{\ell}}$ is among all linear estimators.

There exists a famous result by [109] which exhibits an estimator which is even minimax, i.e. which attains not only the optimal rate, but also the exact constant. This estimator is called the Pinsker estimator.

The following theorem is due to [109].

**Theorem 1.5.** *Let $\{a_k\}$ be a sequence of non-negative numbers and let $\sigma_k > 0$, $k = 1, 2, \ldots$ Then the linear minimax estimator $\lambda = \{\lambda_k\}$ on $\Theta(a, L)$ is given by*

$$\lambda_k = (1 - c_{\varepsilon} a_k)_+, \tag{1.42}$$

*where $c_{\varepsilon}$ is the solution of the equation*

$$\varepsilon^2 \sum_{k=1}^{\infty} \sigma_k^2 a_k (1 - c_\varepsilon a_k)_+ = c_\varepsilon L$$

*and the linear minimax risk is*

$$r_\varepsilon^\ell(\Theta) = \varepsilon^2 \sum_{k=1}^{\infty} \sigma_k^2 (1 - c_\varepsilon a_k)_+. \tag{1.43}$$

*Furthermore, if*

$$\frac{\max_{k:a_k<T} \sigma_k^2}{\sum_{k:a_k<T} \sigma_k^2} = o(1), \quad T \to \infty, \tag{1.44}$$

*then*

$$r_\varepsilon(\Theta) = r_\varepsilon^\ell(\Theta)(1 + o(1)), \tag{1.45}$$

*as $\varepsilon \to 0$.*

*Proof.* A proof may be found in [109, 7].

Thus, under the condition (1.44), the linear minimax estimator given by (1.42) is asymptotically minimax among all estimators.

This result has been also generalized to the very specific case of severely ill-posed problems with analytic functions, i.e. when (1.44) is not verified, in [61, 62].

The optimal rates of convergence may also be found, for example in [7] and [32]. The function $f$ is supposed to have Fourier coefficients in some ellipsoid, and the problem is mildly, severely ill-posed or even direct. The rates appear in Table 2.

*Example 1.3.* All the rates are given for the estimation of a function in one dimension ($d = 1$). Otherwise, in a multidimensional framework, it is well-known that the minimax rates depend on the dimension $d$.

There exist also many optimal rates results in inverse problems see for example the deconvolution problem in [50, 52], the tomography problem studied in the papers [80, 85, 83, 20], general inverse problems [84, 39, 95, 47, 78, 9] and in related frameworks [19].

A recent review of the different rates in rather general inverse problems may be found in [93].

**Table 1.2** Optimal rates of convergence

| Inverse Problem/Functions | Sobolev | Analytic |
|---|---|---|
| Direct problem | $\varepsilon^{\frac{4\alpha}{2\alpha+1}}$ | $\varepsilon^2 \log \frac{1}{\varepsilon}$ |
| Mildly ill-posed | $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta+1}}$ | $\varepsilon^2 \left(\log \frac{1}{\varepsilon}\right)^{2\beta+1}$ |
| Severely ill-posed | $\left(\log \frac{1}{\varepsilon}\right)^{-2\alpha}$ | $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta}}$ |

## Comments

*Ill-posedness.* We may remark that the rates usually depend strongly on the smoothness $\alpha$ of the function $f$ and on the degree of ill-posedness $\beta$. When $\beta$ increases the rates are slower. This is a very important point, which characterizes the influence of the ill-posedness in the results. In ill-posed problems the rates are slower, making estimation in these models more difficult.

*Direct model/Sobolev.* We get the standard rates for nonparametric estimation. Indeed, with the relation $\varepsilon^2 \asymp 1/n$, one really obtains the usual $n^{-\frac{2\alpha}{2\alpha+1}}$ rate for estimating a $\alpha$ smooth function in a nonparametric context, see [74, 122].

The more standard cases for inverse problems are, mildly ill-posed/Sobolev, or severely ill-posed/Analytic. Indeed, they correspond to the natural setting of Hölder source conditions (see Section 1.2.3.2).

*Mildly ill-posed/Sobolev.* This is, in a way, the more standard framework. One has a not so difficult inverse problem with smooth functions. The rate is then $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta+1}}$. One may see the loss in the rate due to ill-posedness $\beta$, compared to the rate in the direct problem $\varepsilon^{\frac{4\alpha}{2\alpha+1}}$. The rate is polynomial in $\varepsilon$ and slower than $\varepsilon^2$, as usual in nonparametric statistics.

*Severely ill-posed/Analytic.* In this context, the problem is very difficult, but the functions are then very smooth. The rate is then still polynomial $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta}}$. This rate is slightly different from the previous case, but related.

The three other cases are very specific problems. The rates are then not polynomial.

*Direct model/Analytic.* This framework is rather easy. Indeed, the problem is direct, and the functions are very smooth. The rate is then almost parametric, i.e. $\varepsilon^2$. One just looses a logarithmic term compared to the parametric context. From a statistical point of view, the situation is very specific. Indeed, there is no trade-off between bias and variance, the variance term is dominating.

*Severely ill-posed/Sobolev.* This case corresponds to a very difficult inverse problem with not smooth enough functions. The rate is logarithmic, and thus very slow. From a theoretical point of view, this context might be considered as too difficult. Here, the bias is dominating.

*Mildly ill-posed/Analytic.* In this case, a mildly ill-posed problem with very smooth functions, the rate is almost the parametric rate $\varepsilon^2$. The variance term is dominating. The functions are so smooth that the inverse problem has almost no influence. Indeed, the degree of ill-posedness appears only in the logarithmic term.

*Remark 1.25.* One may also consider inverse problems where $\sigma_k \asymp \exp(\beta k^r)$, where $\beta > 0$ and $r \geqslant 1$, for example Heat equation or convolution by a Gaussian kernel. Here the rates will be worse. For example, in the case of Sobolev functions, the rate will be $(\log \frac{1}{\varepsilon})^{-2\alpha/r}$.

*Remark 1.26.* In the problem of tomography presented in Section 1.1.6.5, the situation is slightly different. Indeed, this is a two dimensional problem. The optimal rate of convergence is given, in [80], and corresponds to $2\alpha/(2\alpha+3)$. This rate has to be compared to the optimal rate of estimating a function in $d$ dimensions, which is $2\alpha/(2\alpha+d)$. Thus, in dimension $d=2$, one really sees the ill-posedness $\beta=1/2$, in obtaining the rate $2\alpha/(2\alpha+3)$.

In the sequel, we will use quite often the two very standard results,

$$\sum_{k=1}^{n} k^p \approx \frac{n^{p+1}}{p+1}, \; p > -1, \text{ as } n \to \infty \tag{1.46}$$

and

$$\sum_{k=1}^{n} e^{pk} \approx \frac{e^{p(n+1)}}{e^p - 1}, \; p > 0, \text{ as } n \to \infty, \tag{1.47}$$

where $a_n \approx b_n$ means that $a_n/b_n \to 1$ as $n \to \infty$.
    In this framework we obtain the following theorem.

**Theorem 1.6.** *Consider now the case where $\sigma_k \asymp k^\beta, \beta \geqslant 0$ and $\theta$ belongs to the ellipsoid $\Theta(\alpha,L)$, where $a_k = k^\alpha, \alpha > 0$. Then the projection estimator with $N \asymp \varepsilon^{-2/(2\alpha+2\beta+1)}$ verifies as $\varepsilon \to 0$*

$$\sup_{\theta \in \Theta(\alpha,L)} R(\theta,N) \leqslant C\varepsilon^{4\alpha/(2\alpha+2\beta+1)}.$$

*This rate is optimal (see Theorem 1.5).*

*Proof.* We have,

$$\sup_{\theta \in \Theta(\alpha,L)} R(\theta,N) = \sup_{\theta \in \Theta(\alpha,L)} \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2.$$

We bound the first term as follows,

$$\sup_{\theta \in \Theta(\alpha,L)} \sum_{k=N+1}^{\infty} \theta_k^2 \leqslant \sup_{\theta \in \Theta(\alpha,L)} \sum_{k=N+1}^{\infty} k^{2\alpha} \theta_k^2 k^{-2\alpha}$$

$$\leqslant N^{-2\alpha} \sup_{\theta \in \Theta(\alpha,L)} \sum_{k=1}^{\infty} k^{2\alpha} \theta_k^2 \leqslant L N^{-2\alpha}.$$

The variance term is controlled by

$$\varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 \asymp \varepsilon^2 \sum_{k=1}^{N} k^{2\beta} \asymp \frac{\varepsilon^2 N^{2\beta+1}}{2\beta+1},$$

when $N$ is large, by use of (1.46). Thus,

$$\sup_{\theta \in \Theta(\alpha,L)} R(\theta,N) \leqslant L\, N^{-2\alpha} + \frac{\varepsilon^2 N^{2\beta+1}}{2\beta+1}.$$

If we want to attain the optimal rate of convergence we have to choose $N$ of order $\varepsilon^{-2/(2\alpha+2\beta+1)}$ as $\varepsilon \to 0$. This choice corresponds to the trade-off between the bias term and the variance term.

*Remark 1.27.* This proof is very simple and only concerns the rate of convergence for a given estimator, the so-called upper bound. The proof of an upper bound for some estimator is usually rather easy. There is no proof here of the lower bound, i.e. showing that no estimator has a risk converging faster. The lower bound is proved by Theorem 1.5. Nevertheless, lower bounds are very important in nonparametric statistics. Indeed, it is the lower bound which proves that the estimator is optimal, i.e. one of the best estimator in a given model. For a discussion in details of the standard methods, see [129].

*Remark 1.28.* Considering the minimax point of view, we may remark that there exists an optimal choice for $N$ which corresponds to the balance between the bias and the variance. However, this choice depends very precisely on the smoothness $\alpha$ and on the degree of ill-posedness of the inverse problem $\beta$.

Even in the case where the operator $A$ (and then its degree $\beta$) is known, it has no real meaning to consider that we know the smoothness of the unknown function $f$.

These remarks lead to the notion of adaptation and also oracle inequalities, i.e. how to choose the bandwidth $N$ without strong a priori assumptions on $f$ (see Section 1.3).

### 1.2.4.2  Deconvolution on $\mathbb{R}$

Assume that we are in the special inverse problem of deconvolution on $\mathbb{R}$ (see Section 1.1.7.2). Consider only the case of spectral cut-off regularization. We estimate in the Fourier domain $Ff$ by

$$\frac{FY(\omega)}{\tilde{r}(\omega)} I(\omega : \tilde{r}^2(\omega) > \gamma),$$

and then the spectral cut-off regularization is

$$\hat{f}_\gamma^{SC} = F^{-1}\left( \frac{FY(\omega)}{\tilde{r}(\omega)} I(\omega : \tilde{r}^2(\omega) > \gamma) \right).$$

As in the SVD case, the bias term (approximation error) is usually controlled by the source conditions. In this framework, the Hölder source condition (1.36) is equivalent to, in the Fourier domain,

$$Ff = \tilde{r}^{2\mu} Fw, \; w \in L^2(\mathbb{R}), \; \|w\|^2 \leqslant L.$$

$$\int_{\mathbb{R}} |Fw(\omega)|^2 d\omega = \int_{\mathbb{R}} |Ff(\omega)|^2 \tilde{r}^{-4\mu}(\omega) d\omega \leqslant L. \qquad (1.48)$$

Similarly to the previous section, if $\tilde{r}(\omega) = |\omega|^{-\beta}$, the problem is then mildly ill-posed. In this case, the source conditions correspond to some Sobolev class of functions on $\mathbb{R}$ (see [10])

$$\mathscr{W}(\alpha, L) = \left\{ f \in L^2(\mathbb{R}) : \int_{\mathbb{R}} |\omega|^2 |Ff(\omega)|^2 \leqslant L \right\},$$

which is equivalent to, for $\alpha \in \mathbb{N}$,

$$\mathscr{W}(\alpha, L) = \left\{ f \in L^2(\mathbb{R}) : \int_{\mathbb{R}} (f^{(\alpha)}(t))^2 dt \leqslant L \right\}.$$

We then obtain

$$(\mathbf{E}_f \hat{f}_\gamma^{SC}(x) - f(x)) = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-i\omega x} \left( Ff(\omega) I(\omega : \tilde{r}^2(\omega) > \gamma)) - Ff(\omega) \right) d\omega,$$

and then for the bias

$$\int_{\mathbb{R}} (\mathbf{E}_f \hat{f}_\gamma^{SC}(x) - f(x))^2 dx \leqslant \frac{1}{2\pi} \int_{|\omega| > \gamma^{-1/2\beta}} |Ff(\omega)|^2 d\omega$$

$$\leqslant \gamma^{4\mu\beta/2\beta} \int_{\mathbb{R}} |Ff(\omega)|^2 |\omega|^{4\mu\beta} d\omega \leqslant L\gamma^{2\mu}.$$

We need now to bound the stochastic term. Using (1.32), we have

$$\mathbf{E}_f \|\hat{f}_\gamma^{SC} - \mathbf{E}_f(\hat{f}_\gamma^{SC})\|^2 \leqslant \mathbf{E}\|\varepsilon \Phi_\gamma(A^*A)A^*\xi\|^2.$$

Using (1.22) and Lemma 1.5, we may bound the variance term

$$\mathbf{E}_f \|\hat{f}_\gamma^{SC} - \mathbf{E}_f(\hat{f}_\gamma^{SC})\|^2 \leqslant \frac{\varepsilon^2}{2\pi} \mathbf{E} \int_{|\omega| < \gamma^{-1/2\beta}} \left| \frac{\eta(\omega)}{\tilde{r}(\omega)} \right|^2 d\omega$$

$$\leqslant \frac{\varepsilon^2}{2\pi} \int_{|\omega| < \gamma^{-1/2\beta}} |\omega|^{2\beta} d\omega \asymp \varepsilon^2 (\gamma^{-1/2\beta})^{2\beta+1} = \varepsilon^2 \gamma^{-\frac{2\beta+1}{2\beta}}.$$

The risk of the spectral cut-off is then bounded by

$$\mathbf{E}_f \|\hat{f}_\gamma^{SC} - f\|^2 \leqslant L\gamma^{2\mu} + C\gamma^{-\frac{2\beta+1}{2\beta}},$$

the optimal choice is then $\gamma^* \asymp (\varepsilon^2)^{2\beta/(4\mu\beta+1)}$ which corresponds to the rate

$$\mathbf{E}_f \|\hat{f}_{\gamma^*}^{SC} - f\|^2 \leqslant C(\varepsilon^2)^{\frac{4\mu\beta}{4\mu\beta+2\beta+1}}.$$

This rate may be shown to be optimal.

*Remark 1.29.* Recall that here $N = [\gamma^{-1/2\beta}]$. The rates are in fact the same than in the compact case of Section 1.2.4.1.

In the case of severely ill-posed problems, i.e. $\tilde{r}(\omega) = \exp(-\beta|\omega|)$, the class coming from Hölder source condition is then different. By using (1.48), we have

$$\int_{\mathbb{R}} |Fw(\omega)|^2 d\omega = \int_{\mathbb{R}} |Ff(\omega)|^2 \tilde{r}^{-4\mu}(\omega) d\omega \leqslant L.$$

Thus,

$$\int_{\mathbb{R}} |Ff(\omega)|^2 \exp(4\mu\beta|\omega|) d\omega \leqslant L.$$

which corresponds to the the the class of analytic functions, i.e. which admits an analytic continuation into a band of the complex plane, see for example [75].

*Remark 1.30.* In the context of general inverse problems, with general regularization methods, it is also possible to obtain results concerning rates of convergence (see [9]).

## *1.2.5 Comparison Between Deterministic and Stochastic Noise*

In this section, consider the model of inverse problems with deterministic noise. This model is, in some sense, the historical model of inverse problems. It appears for example in [127] and [128]. The analog of the stochastic model (1.2), in the deterministic framework is the following. We have

$$Y = Af + \varepsilon h, \tag{1.49}$$

where the noise $h$ is considered as some deterministic element $h \in G$, with $\|h\| \leqslant 1$. Since the noise is some unknown element of a ball in $G$, then the results have to be obtained for any possible noise, i.e. for the worst noise.

Compare the deterministic model in (1.49) and the stochastic model in (1.2) where $\xi$ is a white noise. At first glance, it may seem, that the main difference between the two models concerns the nature of the noise, deterministic against stochastic.

In fact, it is more the level of the two noises which are not the same.

The first main difference, since $\xi$ is a white noise, is that $Y$ in (1.2) is not really "observed". Indeed, $\xi$ does not take its values in $G$. We only observe its projection on some basis. Indeed, $\xi$ as a white noise, is not a Hilbert-space random variable in $G$ but a Hilbert-space process acting on $G$. Formally, we have $\|\xi\|_G = \infty$, thus $\xi$ is not an element of $G$. On the other hand, the deterministic noise $h$ belongs to $G$, and $\|h\| \leqslant 1$. The deterministic noise is then "small" compared to the stochastic one.

This fact, has been already noted in [38].

In order to have a more comprehensive study, consider the class of linear injective and compact operators which admit a singular value decomposition (SVD) (see Section 1.1.5).

The analog of the sequence space model in (1.6) may be written as

$$X_k = \theta_k + \varepsilon \sigma_k \, h_k, \quad k = 1, 2, \ldots, \tag{1.50}$$

where $\{h_k\}$ are the coefficients of $h$ in the basis $\{\psi_k\}$.

A natural way of studying the two frameworks is to compare the accuracy of estimation (reconstruction). Define two standard criteria, in order to measure the error or risk, for any estimator $\hat{f}$ (or regularization method). For the stochastic noise model, use the maximal risk defined in Definition 1.9. For the deterministic noise model, define the worst noise risk

$$\sup_{f \in \mathscr{F}} \sup_{\|h\| \leqslant 1} \|\hat{f} - f\|^2,$$

where $f$ belongs to some class of functions $\mathscr{F}$.

The goal is to compare the optimal rates of convergence in each model, i.e. the order of the risk of the best possible estimator as $\varepsilon \to 0$. Indeed, this rate defines a notion of difficulty of estimation in a given model. Two models with the same optimal rates of convergence are usually thought to be close, at least from the estimation point of view.

One difference between deterministic and stochastic cases, is that since $\|h\| \leqslant 1$ (i.e. $\sum h_k^2 \leqslant 1$), the noise $h_k$ decreases in (1.50) as $k$ increases. In the stochastic case, the level of the noise $\xi_k$ is the same in each coefficient $X_k$. Thus, the stochastic noise seems to be larger.

It is well-known, that the rates of convergence depend on difficulty of the inverse problem and smoothness conditions on the function $f$ (see Section 1.2.4). For the inverse problems, the two standard cases are $\sigma_k \asymp k^\beta$ or $\sigma_k \asymp e^{\beta k}$, $\beta > 0$ which correspond to mildly or severely ill-posed respectively. The parameter $\beta$ denotes the degree of ill-posedness.

Concerning smoothness properties of $f$, associated with the behaviour of its coefficients $\theta_k$, consider the ellipsoid of coefficients in $\ell^2$ as in Section 1.2.3.2. Consider the two standard cases, Sobolev ($a_k = k^\alpha$) and Analytic ($a_k = e^{\alpha k}$), where $\alpha > 0$ is the smoothness of $f$.

For the stochastic noise, the optimal rates of convergence may be found in Table 2. Concerning the deterministic noise, rates of convergence may be obtained, for example in [49].

Consider here the two more natural cases, *polynomial* ($\sigma_k \asymp k^\beta$ and $a_k = k^\alpha$) and *exponential* ($\sigma_k \asymp e^{\beta k}$ and $a_k = e^{\alpha k}$).

Consider also a third case: the direct problem, where $\sigma_k \equiv 1$ (i.e. $A = I$) and $f$ belongs to a Sobolev ball ($a_k = k^\alpha$).

All these rates are given in the following table:

Remark that in the exponential case, rates of convergence are the same for the two models.

**Table 1.3** Rates for deterministic and stochastic model

|             | Deterministic                    | Stochastic                         |
| ----------- | -------------------------------- | ---------------------------------- |
| Direct      | $\varepsilon^2$                  | $\varepsilon^{\frac{4\alpha}{2\alpha+1}}$ |
| Polynomial  | $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta}}$ | $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta+1}}$ |
| Exponential | $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta}}$ | $\varepsilon^{\frac{4\alpha}{2\alpha+2\beta}}$ |

On the other hand, rates are different in the polynomial case, which is more standard. There is a small difference between $2\alpha/(2\alpha + 2\beta)$ and $2\alpha/(2\alpha + 2\beta + 1)$ which could be thought as not very important. However, this is fundamental.

In order to understand well this phenomenon, consider what happens when $\beta \to 0$. The problem is less and less ill-posed and becomes close to the case $\beta = 0$, i.e. to the direct case where $\sigma_k \equiv 1$ and $A = I$ is the identity. In the deterministic problem, the rate will attain $\varepsilon^2$ ($a = 1$) in the direct case. In the stochastic framework, the rate will be $\varepsilon^{4\alpha/(2\alpha+1)}$.

The fundamental difference now appears. In the stochastic direct problem, the rate depends on the smoothness $\alpha$ of the estimated function $f$. This is not true for the deterministic framework.

In the stochastic case, in order to estimate the function $f$, one needs to balance the approximation error and the stochastic error. This is the usual trade-off in nonparametric statistics between the bias and the variance.

Everything is different in the deterministic case. The function $f$ will be directly estimated by $Y$, which attains the rate $\varepsilon^2$. There is no trade-off, the whole series $\{X_k\}$ is used to estimate $\{\theta_k\}$. The rate $\varepsilon^2$ is usually obtained in statistics in the parametric case, i.e. when estimating a vector $\theta$ of finite dimension. In the stochastic case, one cannot use directly $Y$ which has infinite risk.

In the direct case, the two models are thus totally different. Indeed, the deterministic noise is smaller than the stochastic one, because it is bounded. In (1.50) the errors $h_k$ become small with $k$, whereas the stochastic errors $\xi_k$ are of the same order in (1.6).

From a statistical point of view such a small error would not really make sense. Indeed, statistics study the effect of stochastic errors and these errors should be important enough. However, from a numerical point of view, it could make sense to neglect the noise, or at least to consider it as small. Thus, the difference is more in the level of the noise than its nature (deterministic or stochastic).

In order to explain more clearly the influence of noise, consider the simple projection estimator,

$$\hat{\theta}_k = I(k \leqslant N) \frac{y_k}{b_k},$$

where $I(\cdot)$ is the indicator function and $N$ is some integer. It is known that this family of estimators attains, for a correct choice of $N$, the optimal rate of convergence on $\Theta$ (see Theorem 1.6).

The $\ell^2-$risk of this estimator is, in the stochastic model,

$$\mathbf{E}_\theta \|\hat{\theta} - \theta\|^2 = \mathbf{E}_\theta \sum_{k=1}^\infty (\hat{\theta}_k - \theta_k)^2 = \sum_{k=N+1}^\infty \theta_k^2 + \varepsilon^2 \sum_{k=1}^N \sigma_k^2, \qquad (1.51)$$

and the $\ell^2-$error of the reconstruction method, in the deterministic model, is

$$\sup_{\|h\|\leqslant 1} \|\hat{\theta} - \theta\|^2 = \sup_{\|h\|\leqslant 1} \left( \sum_{k=N+1}^\infty \theta_k^2 + \varepsilon^2 \sum_{k=1}^N h_k^2 \sigma_k^2 \right) = \sum_{k=N+1}^\infty \theta_k^2 + \varepsilon^2 \sigma_N^2, \quad (1.52)$$

in the case of increasing $\sigma_k$.

The influence of the inverse problem is only on the variance term, i.e. the second term in the right-hand side of (1.51) and (1.52). The approximation error $\sum_{k>N} \theta_k^2$ is not modified by the ill-posedness of the inverse problem.

The following table gives the order, as $\varepsilon \to 0$, of the variance term $\varepsilon^2 \sum_{k=1}^N \sigma_k^2$ or $\varepsilon^2 \sigma_N^2$, in the various settings.

**Table 1.4** Variances for deterministic and stochastic model

| | Deterministic | Stochastic |
|---|---|---|
| Direct | $\varepsilon^2$ | $\varepsilon^2 N$ |
| Mildly | $\varepsilon^2 N^{2\beta}$ | $\varepsilon^2 N^{2\beta+1}$ |
| Severely | $\varepsilon^2 e^{\beta N}$ | $\varepsilon^2 e^{\beta N}$ |

The direct case corresponds to $b_k \equiv 1$. In the deterministic model, since $h \in \ell^2$, the variance term is $\varepsilon^2$, and does not depend on $N$. Thus, there is no trade-off, $N$ can be chosen as $\infty$, or any choice such that $\sum_{k>N} \theta_k^2 = O(\varepsilon^2)$.

In the stochastic case, the variance term is $\varepsilon^2 N$. Thus, we have to balance the bias and the variance, as usually in nonparametric statistics, and find the optimal choice of $N$.

The variance terms stay different in the mildly ill-posed (polynomial) case. The ratio between the two variance terms is again $N$. However, this difference is less important as $\beta$ increases. Indeed, when $\beta$ is large $N^{2\beta+1}$ is close to $N^{2\beta}$.

The main point is that the variance term is larger in the case of ill-posed problems. The degree of ill-posedness $\beta$ appears directly in the variance term. The variance increases with $\beta$.

Thus, in the case of ill-posed inverse problems, the deterministic error has more influence than for the direct case. The presence of $\beta$ increases the variance term.

For large $\beta$, the two models give almost the same rates. Finally, for severely ill-posed problems (exponential case), these rates are the same.

The ill-posedness of the inverse problem hides, in some sense, the difference between the two kinds of noise, by increasing the small deterministic noise. When $\beta$ is large, the main part of the variance term $\varepsilon^2 N^{2\beta(+1)}$ is due to the inverse problem and not to the nature of the noise. The inverse problem makes these two models more close, when for the direct problem they are completely different.

The difference between deterministic and stochastic noise is in its level and not really in its nature. Thus, a stochastic model with a small noise could be considered. The model is the following,

$$X_k = \theta_k + \varepsilon \; \sigma_k e_k \xi_k, \;\; k = 1, 2, \dots, \tag{1.53}$$

where $\{\xi_k\}$ are independent standard Gaussian random variables, and $\{e_k\} \in \ell_2$, $\|e\| = 1$. The risk of a projection estimator is then

$$\mathbf{E}_\theta \|\hat{\theta} - \theta\|^2 = \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} e_k^2 \sigma_k^2. \tag{1.54}$$

In the direct case, the variance term is then $\varepsilon^2 \sum_{k=1}^{N} e_k^2$, bounded by $\varepsilon^2$. The optimal rate is so $\varepsilon^2$, as for the deterministic case.

In the case of ill-posed problem, hypothesis should be more precise in order to obtain explicit rate of convergence. Indeed, in the deterministic case we study the worst noise, i.e. $\sup_{\|h\| \leqslant 1}$. Thus, we have to consider a noise in $\ell^2$ but rather large, almost on the "edge". Some example is $e_k = (\sqrt{k} \log(k+1))^{-1}$, which is in $\ell^2$. It is clear that dividing $X_k$ by $e_k$, one obtains a model equivalent to (1.53), with a new $\sigma_k' = \sigma_k e_k$.

With this choice of $\{e_k\}$, in the mildly ill-posed case, the variance term is then $\varepsilon^2 \sum_{k=1}^{N} k^{2\beta-1} \log^{-1}(k+1)$, which is equivalent (up to a log term) to $\varepsilon^2 N^{2\beta}$. Thus, the risk in (1.54) is of the same order than in the deterministic case. Looking at the rate for the stochastic case with $\beta - 1/2$, we obtain the rate with $\beta$ for deterministic case.

In the exponential case, $\{e_k\}$ has no real influence.

Thus, using a model of inverse problem with stochastic noise with a "small" noise, we obtain the same rate of convergence than for the deterministic case (up sometimes to some log term). A "small" stochastic noise is in fact a Hilbert-space random variable, and not only a Hilbert-space process. It is random, but really takes its values in $G$.

In conclusion, the main difference between the two approaches comes more from the level of the noise and not so much from its nature.

However, this short study is not at all exhaustive. A more precise approach, based not only on the comparison between the optimal rates, but also the exact constants in the risk, would highlight more differences. For all that, such a technical comparison would not really make sense, since at this precision level, any models are different.

A more sensible framework, in order to compare deterministic and stochastic noise, concerns the construction of adaptive estimators, i.e. which do not depend on the smoothness $\alpha$ of the function to reconstruct (see the following Section 1.3).

In this case the nature of the noise would have more influence. Indeed, the methods could then be very different, for example the discrepancy principle [49] for deterministic noise, or cross-validation, unbiased risk estimator (see Section 1.3.3.1) or the Lepski method [87, 101] for stochastic noise. In the deterministic case, one crucial point is that the error in the data $(\varepsilon)$ is precisely known, and then, one can

reject reconstruction $\hat{f}$ such that $\|A\hat{f} - Y\| > \varepsilon$. In the stochastic case, the main idea of adaptation is to use large deviations for the noise. Usually, one find values such that the noise will have a very small probability to fall beyond, as in Lemma 1.7 (see also for some examples of adaptivity results [78] and [25]).

In conclusion, this study is not claiming that the two approaches present no difference. The two frameworks are similar in some ways. The differences coming more from the level of the noise than from its nature.

## 1.3 Adaptation and Oracle Inequalities

One of the most important point in nonparametric statistics is then typically linked to the problem of calibrating by the data the tuning parameter $(N, \gamma$ or $m)$ in any class of estimators. For example, we have seen that this choice is very sensitive if we want to attain the optimal rate of convergence.

This problem leads to the notion of adaptation and oracle inequalities, i.e. how to construct truly data-driven estimators which have good theoretical properties.

This framework is very important, in theory, but also in applications. Indeed, the notion of, rates of convergence, smoothness $\alpha$ of the function to reconstruct, degree of ill-posedness $\beta$ of the inverse problem, are very interesting. They help to understand, the difficulty of an ill-posed problem, the influence of smoothness on the rates and so on... The (minimax) optimality of an estimator is also very important. Indeed it shows that no estimator may do better in a given class of functions.

However, they are just mathematical and asymptotical tools. The degree $\beta$ of a given inverse problem is usually not known. It is even worse concerning the smoothness $\alpha$ of the target function $f$. One has no chance to have any idea of it.

Definitely, one cannot rely on some unknown smoothness, and asymptotic relationship, in order to make the choice of the tuning parameter $(N$ or $\gamma)$. One has to really construct data-driven methods in order to calibrate the tuning parameter. Then, the main goal is to prove that this data-driven method has a good behaviour, from a mathematical point of view.

This problem of adaptation is presented in the framework of the sequence space model defined in (1.6) and directly linked, by use of the SVD, to some inverse problem with a compact operator $A$.

### 1.3.1 Minimax Adaptive Procedures

The starting point of the approach of **minimax adaptation** is a collection $\mathscr{G} = \{\Theta_\alpha\}$ of classes $\Theta_\alpha \subset \ell^2$. The statistician knows that $\theta$ belongs to some member $\Theta_\alpha$ of the collection $\mathscr{G}$, but he does not know exactly which one. If $\Theta_\alpha$ is a smoothness class, this assumption can be interpreted as follows: the statistician knows that the underlying function has some smoothness, but he does not know the degree of smoothness.

**Definition 1.12.** An estimator $\theta^\star$ is called **minimax adaptive** on the scale of classes $\mathscr{G}$ if for every $\Theta_\alpha \in \mathscr{G}$ the estimator $\theta^\star$ attains the optimal rate of convergence.

An estimator $\theta^\star$ is called **sharp minimax adaptive** on the scale of classes $\mathscr{G}$ if it also attains the exact minimax constant.

The idea of choosing the tuning parameter (bandwidth) of an estimator in a data-driven way is a very standard idea in nonparametric statistics. However, the main difficulty then concerns the mathematical behaviour of such an estimator. Only quite recently, this idea has been formalized in a rigorous way by [87].

Lepski, in [88], has developed a method in order to construct adaptive estimators, i.e. an estimator which attains the optimal rate for any class $\Theta_\alpha$.

In some cases, no estimator attains (exactly) the optimal rate on the whole scale. One has often to pay a price for adaptation [89]. This cost in the accuracy for construction of adaptive estimator is usually the loss of a logarithmic term in the rate of convergence.

Since the beginning of the 90's, adaptive estimation is really one of the leading topics in nonparametric statistics. Many adaptive (or almost adaptive) estimators have been constructed, in very different frameworks, and various classes of functions.

One may use very different procedures in order to construct adaptive estimators, for example, Lepski's algorithm in [88], model selection in [4], unbiased risk estimation in [82], or wavelets thresholding in [40].

Adaptive minimax estimation in statistical inverse problems as (1.2) has been studied quite recently. This has been done for many inverse problems (deconvolution, heat equation, tomography...).

There exist also a very vast literature on adaptation in inverse problems by Wavelet-Vaguelette Decomposition (WVD) on the Besov scale of classes, see [39, 83, 78, 28, 34, 79, 29, 71].

Lepski's procedure has been also used in inverse problems in several papers [55, 56, 21, 5, 101].

The unbiased risk estimation is also quite popular in inverse problems, see [25, 97].

The model selection is considered in inverse problems [35, 91].

Other adaptive results may be found in [46, 47, 48, 58, 24, 57].

*Remark 1.31.* Minimax adaptive estimators are really important in statistics from a theoretical and from a practical point of view. Indeed, it implies that these estimators are optimal for any possible parameter in the collection $\mathscr{G}$. From a more practical point of view it garantees a good accuracy of the estimator for a very large choice of functions.

Thus, we have an estimator which automatically adapts to the unknown smoothness of the underlying function. The estimator is then completely data-driven and automatic. However, it behaves as if it knew the true smoothness. This notion is very important since this smoothness is almost never known.

### *1.3.2 Oracle Inequalities*

Consider now a linked, but different point of view. Assume that a class of estimators is fixed, i.e. that the class of possible weights $\Lambda$ is given. Define the **oracle** $\lambda^0$ as

$$R(\theta, \lambda^0) = \inf_{\lambda \in \Lambda} R(\theta, \lambda). \qquad (1.55)$$

The oracle corresponds to the best possible choice in $\Lambda$, i.e. the one which minimizes the risk. However, this is not an estimator since the risk depends on $\theta$, the oracle will also depend on this unknown $\theta$. For this reason, it is called oracle since it is the best one in the family, but it knows the true $\theta$. Another important point is to note that the oracle $\lambda^0$ usually depends really on the family $\Lambda$. As an infimum, the oracle is not necessarily unique or may not be exactly attained. However, this has no influence on the results. Indeed, one only considers the risk of the oracle $\inf_{\lambda \in \Lambda} R(\theta, \lambda)$.

The goal is then to find a data-driven sequence of weights $\lambda^\star$ with values in $\Lambda$ such that the estimator $\theta^\star = \hat{\theta}(\lambda^\star)$ satisfies an **oracle inequality**, for any $\varepsilon > 0$ and any $\theta \in \ell^2$, there exits $\tau_\varepsilon > 0$,

$$\mathbf{E}_\theta \|\theta^\star - \theta\|^2 \leqslant (1 + \tau_\varepsilon) \inf_{\lambda \in \Lambda} R(\theta, \lambda) + \Omega_\varepsilon, \qquad (1.56)$$

where $\Omega_\varepsilon$ is some positive remainder term and $\tau_\varepsilon > 0$ (close to 0 if possible). If the remainder term is small, i.e. smaller than the main term $R(\theta, \lambda^0)$ then an oracle inequality proves that the estimator has a risk of the order of the oracle.

A standard remainder term is $\Omega_\varepsilon = c\varepsilon^2$, where $c$ is uniform positive constant. In this case, the remainder term is really considered as "small". Indeed, in most of the nonparametric frameworks, the rates of convergence are worse than $\varepsilon^2$, which is the parametric rate (see Table 2). In an asymptotic point of view, the risk of the oracle, will then be larger than the remainder term. Thus, the leading term of the inequality will be the risk of the oracle.

A more precise result is the following. The estimator $\theta^\star = \hat{\theta}(\lambda^\star)$ satisfies an **exact oracle inequality**, for any $\varepsilon > 0$, any $\theta \in \ell^2$, and for all $\tau_\varepsilon > 0$,

$$\mathbf{E}_\theta \|\theta^\star - \theta\|^2 \leqslant (1 + \tau_\varepsilon) \inf_{\lambda \in \Lambda} R(\theta, \lambda) + \Omega_\varepsilon, \qquad (1.57)$$

where $\Omega_\varepsilon \geqslant 0$ and usually $\Omega_\varepsilon$ depends on $\tau_\varepsilon$.

*Remark 1.32.* We are interested in data-driven methods, and thus automatic, which more or less mimic the oracle.

One may obtain some asymptotic results when $\varepsilon \to 0$. We call an **asymptotic exact oracle inequality** on the class $\Lambda$, as $\varepsilon \to 0$,

$$\mathbf{E}_\theta \|\theta^\star - \theta\|^2 \leqslant (1 + o(1)) \inf_{\lambda \in \Lambda} R(\theta, \lambda), \qquad (1.58)$$

for every $\theta$ within some large subset $\Theta_0 \subseteq \ell^2$.

In other words, the estimator $\theta^\star$ asymptotically precisely mimics the oracle on $\Lambda$ for any sequence $\theta \in \Theta_0$.

An important question is how large is the class $\Theta_0$ for which (1.58) can be guaranteed. Ideally, we would like to have (1.58) for all $\theta \in \ell^2$ and with $o(1)$ that is uniform over $\theta \in \ell^2$ (i.e. $\Theta_0 = \ell^2$). This property can be obtained for some classes $\Lambda$ (see, for example, [25]), but with restrictions on $\Lambda$ that do not allow correct rates of the oracle risk $R(\theta, \lambda^0)$ for "very smooth" $\theta$, i.e. analytic functions. If we choose $\Lambda$ large enough to allow all the spectrum of rates for the oracle risk, up to the parametric rate $\varepsilon^2$, we cannot have (1.58) for all $\theta \in \ell^2$ and with $o(1)$ that is uniform over $\theta \in \ell^2$. Although, slightly restricted versions of (1.58) are possible. In particular, $\Theta_0$ can be either the set of all $\theta \neq 0$, or the set $\ell^2_-$, i.e. the subspace of $\ell^2$ containing all the sequences with infinitely many non-zero coefficients (i.e. "nonparametric" sequences), or the set $\{\theta : \|\theta\| \geqslant r_0\}$ for some small $r_0 > 0$. Also, the uniformity of $o(1)$ in $\theta$ is not always granted if both classes $\Lambda$ and $\Theta_0$ are large.

One of first to really see the importance of oracle inequalities are Donoho and Johnstone in [40] where they introduced also the name *oracle*.

During the end of 90's, the oracle was still mainly seen as just a tool in order to prove adaptation. However, nowadays, this point of view has really changed. Oracle inequalities are often considered as the main results for a given estimator. The oracle approach has also modified the statisticians behaviour. For example, non-asymptotic point of view is much more common now.

To our knowledge, one of the first exact oracle inequalities were obtained for the classes of "ordered linear smoothers" in [82]. In particular, Kneip's result applies to projection estimators and to spline smoothing.

The work of Birgé and Massart on model selection is also strongly related to the notion of oracle inequalities, usually in a slightly different form with a penalized version of an oracle inequality, see [4, 8, 100].

Oracle inequalities are nowadays popular, in the nonparametric statistics literature, see [41, 17, 105, 31, 114].

The earlier papers of [119, 90, 59, 111] also contain, although implicitly, oracle inequalities for some classes $\Lambda$. All these papers use the unbiased risk estimation method (see Section 1.3.3.1).

A very interesting review on the topic is [18].

The oracle approach is quite recent in inverse problems. However, the oracle point of view, was growing at the same times than the statistical study of inverse problems. Thus, there is now a rather large interest on oracle inequalities in the statistical inverse problem community, see [78, 25, 60, 35, 91, 98].

## Comments

*Oracle/minimax.* The oracle approach is in some sense the opposite of the minimax approach. Here, we fix a family of estimators and choose the best one among them. In the minimax approach, on the other hand, one tries to get the best accuracy for functions which belong to some function class. The oracle approach is really based on classes of estimators, when the minimax approach is built on classes of functions.

*Non-asymptotic oracle.* The oracle inequalities, are true for any $\theta$, and are non asymptotic. This fact has really changed the point of view concerning nonparametric statistics. Nowadays, non-asymptotic results are really popular.

*Oracle: tool for adaptation.* The oracle approach is often used as a tool in order to obtain adaptive estimators. Indeed, the oracle in a given class often attains the optimal rate of convergence. Moreover, the estimator does not depend on any smoothness assumptions on $f$. Thus, by proving an oracle inequality, one often obtains, a minimax adaptive estimator, see for example, Theorems 1.8 and 1.10. During a quite long time, oracle inequalities were mainly considered as just a tool in order to get minimax adaptive results. Already, [40] pointed out that minimax adaptation can be proved as a consequence of oracle inequalities. They also showed that the method of Stein's unbiased risk estimator is minimax sharp adaptive (or almost minimax sharp adaptive) on some Besov classes.

*Minimax: justification for oracle.* On the other hand, nowadays, the minimax theory may be viewed as a justification for oracle inequality. Indeed, one may ask if the given family of estimators is satisfying. One possible mathematical answer comes from minimax results, which prove that a given family gives optimal estimators. However, in applications, scientists are usually convinced that their favourite method (Tikhonov, projection, $\nu-$method,...) is satisfying.

### 1.3.3 Model Selection

Usually, one key assumption in this approach of oracle inequality, is that $\lambda^\star$ is restricted to take its values in the same class $\Lambda$ that appears in the RHS of (1.56). A **model selection** interpretation of (1.56) is the following: in a given class of models $\Lambda$ we pick the model $\lambda^\star$ that is the closest to the true parameter $\theta$ in terms of the risk $R(\theta, \lambda)$.

The framework of model selection is very popular in statistics, and may have several meaning depending on the topics. We consider the model selection approach to the problem of choosing, among a given family of models $\Lambda$ (estimators), the best possible one. This choice should be made based on the data and not due to some a priori information on the unknown function $f$.

#### 1.3.3.1 Unbiased Risk Estimation

The definition of the oracle in (1.55) is that it minimizes the risk. Since $\theta$ is unknown, the risk is also, and so is the oracle.

A very natural idea in statistics is to estimate this unknown risk by a function of the observations, and then to minimize this estimator of the risk. A classical approach to this minimization problem is based on the principle of **Unbiased Risk Estimation** (URE).

The idea of unbiased risk estimation was developed in [2] and also in [96, 121]. This problem was originally studied in the framework of parametric estimation where the dimension of the model had to choosen.

Mallows, in [96], introduced the $C_p$ in the specific context of regression and the problem of selecting the number of variables that one wants to use in the model.

Akaike, in [2], proposed the Akaike Information Criteria (AIC) in a rather general setting. The idea is to choose the number of parameters $N$ in order to minimize $-2\mathscr{L}_N + 2N$ where $\mathscr{L}_N$ is the maximal value of the log-likelihood, see [3]. In our framework of Gaussian white noise and sequence space model, i.e. a Gaussian noise with a known variance, then AIC and $C_p$ are equivalent. There exist now a very large number of criteria many of them related to AIC, see [119, 90, 111] or the Bayesian Information Criteria (BIC) in [118].

Stein, in [121], proposed his well-known version of URE as the *Stein Unbiased Risk Estimation* (SURE). The results are specific to the Gaussian framework.

Nowadays, $C_p$, AIC, BIC are all used as basic data-driven choices for many statistical models and in several standard softwares.

This idea appears also in all the cross-validation techniques, see the Generalized Cross-Validation (GCV) in [36].

In inverse problems, the URE method is studied in [25], where exact oracle inequalities for the mean square risk were obtained.

In this setting, the functional

$$\mathscr{U}(X,\lambda) = \sum_{k=1}^{\infty}(1-\lambda_k)^2(X_k^2 - \varepsilon^2\sigma_k^2) + \varepsilon^2\sum_{k=1}^{\infty}\sigma_k^2\lambda_k^2 \qquad (1.59)$$

is an unbiased estimator of $R(\theta,\lambda)$ defined in (1.34).

$$R(\theta,\lambda) = \mathbf{E}_\theta\,\mathscr{U}(X,\lambda), \ \forall\lambda. \qquad (1.60)$$

The principle of unbiased risk estimation suggests to minimize over $\lambda \in \Lambda$ the functional $\mathscr{U}(X,\lambda)$ in place of $R(\theta,\lambda)$. This leads to the following data-driven choice of $\lambda$:

$$\lambda_{ure}^\star = \arg\min_{\lambda\in\Lambda}\mathscr{U}(X,\lambda) \qquad (1.61)$$

and the estimator $\theta_{ure}^\star$ defined by

$$\theta_k^\star = \lambda_k^\star X_k. \qquad (1.62)$$

Let the following assumptions hold. For any $\lambda \in \Lambda$

$$(\mathbf{A1}) \quad 0 < \sum_{k=1}^{\infty}\sigma_k^2\lambda_k^2 < \infty, \quad \max_{\lambda\in\Lambda}\sup_k|\lambda_k| \leqslant 1,$$

and, there exists a constant $C > 0$ such that,

$$(\mathbf{A2}) \qquad \sum_{k=1}^{\infty}\sigma_k^4\lambda_k^2 \leqslant C\sum_{k=1}^{\infty}\sigma_k^4\lambda_k^4.$$

Assumptions (A1) and (A2) are rather mild, and they are satisfied in most of the interesting examples. For example, they are trivialy true for projection estimators. Since $|\lambda_k| \leqslant 1$, we also have

$$\sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^4 \leqslant \sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^2,$$

and Assumption (A2) means that both sums are of the same order. The sums $\varepsilon^4 \sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^4$ and $\varepsilon^4 \sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^2$ are the main terms of the variance of $\mathscr{U}(X, \lambda)$.

The Assumption (A1) is quite natural. The first part of (A1) is just to claim that any estimator in $\Lambda$ has a finite variance. The second point follows from (1.34) the remark that the estimator $\hat{\theta}(\lambda)$ with at least one $\lambda_k \notin [0, 1]$ is inadmissible. However, we included the case of negative bounded $\lambda_k$ since it corresponds to a number of well-known estimators, such as some kernel ones.

Denote

$$\rho(\lambda) = \sup_k \sigma_k^2 |\lambda_k| \left\{ \sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^4 \right\}^{-1/2}$$

and

$$\rho = \max_{\lambda \in \Lambda} \rho(\lambda).$$

Although the main results of this section hold for general $\rho$, usually think of $\rho$ as being small (for small $\varepsilon$).

Denote also

$$S = \left( \frac{\max_{\lambda \in \Lambda} \sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^2}{\min_{\lambda \in \Lambda} \sum_{k=1}^{\infty} \sigma_k^4 \lambda_k^2} \right)^{1/2},$$

$$M = \sum_{\lambda \in \Lambda} \exp\{-1/\rho(\lambda)\},$$

and

$$L_\Lambda = \log(DS) + \rho^2 \log^2(MS).$$

Note that $L_\Lambda$ is a term that measure the complexity of the family $\Lambda$ and not only its cardinality $D$.

We obtain the following oracle inequality.

**Theorem 1.7.** *Suppose $\sigma_k \asymp k^\beta$, $\beta \geqslant 0$. Assume that $\Lambda$ is finite with cardinality $D$ and checking Assumptions (A1)-(A2). There exist constants $\gamma_1, \gamma_2 > 0$ such that for every $\theta \in \ell^2$ and for the estimator $\theta_{ure}^\star$ defined in (1.62), we have for B large enough,*

$$\mathbf{E}_\theta \|\theta_{ure}^\star - \theta\|^2 \leqslant (1 + \gamma_1 B^{-1}) \min_{\lambda \in \Lambda} R(\theta, \lambda) + \gamma_2 B \varepsilon^2 L_\Lambda \ \omega(B^2 L_\Lambda), \qquad (1.63)$$

*where*

$$\omega(x) = \max_{\lambda \in \Lambda} \sup_k \left( \sigma_k^2 \lambda_k^2 I\left\{ \sum_{i=1}^{\infty} \sigma_i^2 \lambda_i^2 \leqslant x \sup_k \sigma_k^2 \lambda_k^2 \right\} \right), \quad x > 0.$$

*Proof.* The proof of this theorem may be found in [25].

This result has been extended to the non-compact case in [22].

Function $\omega(x)$ may appear a bit unclear. It depends on the degree of ill-posedness $\beta$ of the inverse problem and the family of estimators. However, in many examples, it is bounded (up to a constant) by $x^{2\beta}$ (see Examples in [25]). Thus the remainder term in the oracle inequality is usually of order $\varepsilon^2 L_\Lambda^{2\beta+1}$.

By assuming hypothesis on the behaviour of $D$ and $S$ when $\varepsilon$ is large, one may obtain an asymptotic exact oracle inequality.

Consider the following family of projection estimators.

*Example 1.4. Projection estimators.* Let $1 \leqslant N_1 < \ldots < N_D$ be integers. Consider the projection filters $\lambda^s = (\lambda_1^s, \lambda_2^s, \ldots)$ defined by

$$\lambda_k^1 = I(k \leqslant N_1), \ \ \lambda_k^2 = I(k \leqslant N_2), \ \ldots \ , \ \lambda_k^D = I(k \leqslant N_D), \ \ k = 1, 2, \ldots \quad (1.64)$$

Suppose also, a polynomial behaviour for $S = O(\varepsilon^{-t})$, for some $t > 0$ and $D = O(\varepsilon^{-v})$, for some $v > 0$. We have $\log(DS) = O(\log(1/\varepsilon))$. As noted Assumptions (A1) and (A2) are always true for projection estimators. Note also that here $\omega(x) \leqslant Cx^{2\beta}$ and

$$L_\Lambda \leqslant C \left( \log(DN_D/N_1) + N_1^{-1} \log^2(N_D/N_1) \right).$$

We have the following corollary.

**Corollary 1.1.** *Assume that $\Lambda = (\lambda^1, \ldots, \lambda^D)$ is the set of projection weights defined in (1.64). If $D = D(\varepsilon)$ and $N_1 = N_1(\varepsilon)$, $N_D = N_D(\varepsilon)$ are such that*

$$\lim_{\varepsilon \to 0} \frac{\log(DN_D/N_1)}{N_1} = 0 \quad (1.65)$$

*then for every $\theta \in \ell^2$ and for the estimator $\theta_{ure}^\star$ defined in (1.62), we have*

$$\mathbf{E}_\theta \| \theta_{ure}^\star - \theta \|^2 \leqslant (1 + o(1)) \inf_{\lambda \in \Lambda} R(\theta, \lambda),$$

*where $o(1) \to 0$ uniformly in $\theta \in \ell^2$.*

*Proof.* The proof of this theorem may be found in [25].

In other words, Corollary 1.1 states that the data-driven selection method $\lambda_{ure}^\star$ behaves itself asymptotically at least as good as the best projection estimator in $\Lambda$.

As noted, a major contribution of oracle inequalities is that they usually allow to construct rather easily minimax adaptive estimators. One just has to construct carefully a family of projection estimators which allows to attain the optimal rate of convergence.

**Theorem 1.8.** *Suppose $\sigma_k \asymp k^\beta$, $\beta \geqslant 0$. Assume that $\Lambda = (\lambda^1, \ldots, \lambda^D)$ is the set of projection weights defined in (1.64). Choose $N_j = j$, $j = 1, \ldots, \varepsilon^{-2}$. Assume that $\theta$ belongs to the ellipsoid $\Theta(\alpha, L)$, where $a_k = k^\alpha$, $\alpha > 0$, $L > 0$, defined in (1.40).*

*Then the URE estimator $\theta_{ure}^{\star}$ defined in (1.62) verifies, for any $\alpha > 0$ and $L > 0$, as $\varepsilon \to 0$,*

$$\sup_{\theta \in \Theta(\alpha, L)} \mathbf{E}_{\theta} \|\theta_{ure}^{\star} - \theta\|^2 \leqslant C\varepsilon^{4\alpha/(2\alpha + 2\beta + 1)}.$$

*This rate is optimal (see Theorem 1.5).*

   *Thus, the URE estimator is then minimax adaptive on the class of ellipsoid.*

*Proof.* The first part of the proof is based on Theorem 1.7. As noted Assumptions (A1) and (A2) are always true for projection estimators. Moreover, here $S = O(\varepsilon^{-2\beta-1})$, $D = O(\varepsilon^{-2})$ and $\omega(x) \leqslant Cx^{2\beta}$. The remainder term is then of order $\varepsilon^2 \log^{2\beta+1}(1/\varepsilon)$.

   The second part is just checking that the best projection estimator in $\Lambda$ attains the optimal rate of convergence. This is true by Theorem 1.6 which gives the optimal choice $N \asymp \varepsilon^{-2/(2\alpha+2\beta+1)}$. Remark also that the remainder term is then much smaller than the optimal rate of convergence.

*Remark 1.33.* Theorem 1.8 may very easily be modified in order a sharp adaptive estimator, i.e. minimax adaptive which also the exact constant. One just has to replace the projection family by the Pinsker family, which is minimax on ellipsoids (see Theorem 1.5).

*Remark 1.34.* One may note that even if we have obtained a very precise oracle inequality in Theorem 1.7, the URE method is in fact not so satisfying in simulations. In the case where the problem is really ill-posed, the URE method is in fact not stable enough (see Section 1.3.3.3).

   This behaviour, may also be understood, by looking at the results and the proof of Theorem 1.7. These remarks lead to the idea of choosing the bandwidth $N$ by a more stable approach (see Section 1.3.3.2).

## Comments

   *Data-driven choices.* One of main difficulties in adaptation or oracle results is that we deal with data-driven choices of $N$. Thus, the risk of the estimator is very difficult to control since it depends on the observations through $X_k$ and also through the data-driven choice of $\lambda^{\star}(X)$. This really changes the structure of the estimator. For example, a linear estimator $\hat{\theta}(\lambda)$ with a data-driven choice of $\lambda^{\star}$ is no more linear. The same remark is true for the unbiased risk estimator, which is no more unbiased for a data-driven choice $\lambda^{\star}$.

   *More difficult proofs.* This remark is clearly one of the main difficulty when dealing with data-driven choices of $N$. Thus, adaptive estimator or oracle inequality are usually more difficult to obtain than rates of convergence results for a given estimator.

   *Proof of an oracle.* We will see this influence in the proof of Theorem 1.9. Indeed, one has to deal carefully with remainder terms depending on a data-driven choice $\lambda^{\star}$.

The very important following lemma is used in the proofs of Theorems 1.7 and 1.9. It may be found in [25]. It allows to control the deviation of the centered stochastic term. This version is not sharp enough to obtain very precise results (see proof of Theorem 1.9). However, it allows to understand the behaviour of the main stochastic term.

This kind of lemma linked to large deviations and exponential inequalities is usually very important in adaptation or oracle inequality results. One needs to study more carefully the behaviour of the stochastic term, and not only control its variance, which is usually enough in rates of convergences results. These inequalities are also linked to the concentration inequalities, see [124].

Let

$$\bar{\eta}_v = (\sqrt{2}\|v\|)^{-1} \sum_{i=1}^{\infty} v_k(\xi_i^2 - 1)$$

where the sums $\|v\|^2$ and $\sum_{k=1}^{\infty} v_i(\xi_i^2 - 1)$ are understood in the sense of mean squared convergence. Define

$$m(v) = \sup |v_i| / \|v\|.$$

**Lemma 1.7.** *We have, for $\kappa > 0$*

$$\mathbf{P}(\bar{\eta}_v > x) \leqslant \begin{cases} \exp\left(-\frac{x^2}{2(1+\kappa)}\right) & \text{for } 0 \leqslant \sqrt{2}m(v)x \leqslant \kappa, \\ \exp\left(-\frac{x}{2\sqrt{2}(1+\kappa^{-1})m(v)}\right) & \text{for } \sqrt{2}m(v)x > \kappa. \end{cases} \tag{1.66}$$

*Proof.* Using the Markov inequality and the formula

$$-\log(1-x) = \sum_{k=1}^{\infty} \frac{x^k}{k}$$

one obtains, for any $0 < t < [\sqrt{2}m(v)]^{-1}$, since $\{\xi_i\}$ are i.i.d. standard Gaussian,

$$\begin{aligned}
\mathbf{P}\{\bar{\eta}_v > x\} &\leqslant \exp(-tx)\mathbf{E}\exp(t\bar{\eta}_v) \\
&= \exp(-tx)\prod_{i=1}^{\infty}\exp\left\{-\frac{tv_i}{\sqrt{2}\|v\|} - \frac{1}{2}\log\left(1 - \frac{\sqrt{2}tv_i}{\|v\|}\right)\right\} \\
&= \exp(-tx)\exp\left\{\sum_{k=2}^{\infty}\sum_{i=1}^{\infty}\frac{1}{2k}\left(\frac{\sqrt{2}tv_i}{\|v\|}\right)^k\right\} \\
&= \exp(-tx)\exp\left\{\sum_{k=2}^{\infty}\frac{(\sqrt{2}t)^k}{2k}\sum_{i=1}^{\infty}\left(\frac{v_i}{\|v\|}\right)^2\left(\frac{v_i}{\|v\|}\right)^{k-2}\right\} \\
&\leqslant \exp(-tx)\exp\left\{\frac{1}{m^2(v)}\sum_{k=2}^{\infty}\frac{1}{2k}[\sqrt{2}tm(v)]^k\right\} \\
&\leqslant \exp(-tx)\exp\left\{-\frac{1}{2m^2(v)}\log[1 - \sqrt{2}tm(v)] - \frac{t}{\sqrt{2}m(v)}\right\}.
\end{aligned}$$

Minimization of the last expression with respect to $t$ yields

$$\mathbf{P}\{\bar{\eta}_v > x\} \leqslant \exp[\varphi_v(x)], \qquad \varphi_v(x) = \frac{1}{2m^2(v)}\log[1 + \sqrt{2}xm(v)] - \frac{x}{\sqrt{2}m(v)}.$$

Note that for $u \geqslant 0$ we have

$$\log(1 + u) - u = u\int_0^1 \left(-\frac{\tau u}{1 + \tau u}\right)d\tau \leqslant -\int_0^1 \frac{\tau u^2}{1 + u}d\tau = -\frac{u^2}{2(1 + u)}.$$

Thus

$$\varphi_v(x) \leqslant -\frac{x^2}{2(1 + \sqrt{2}xm(v))},$$

and we obtain

$$\mathbf{P}\{\bar{\eta}_v > x\} \leqslant \exp\left\{-\frac{x^2}{2(1 + \sqrt{2}xm(v))}\right\}, \qquad \forall x > 0. \tag{1.67}$$

It is easy to see that

$$-\frac{x^2}{2(1 + \sqrt{2}xm(v))} \leqslant \begin{cases} -x^2/2(1 + \kappa), & \sqrt{2}m(v)x \leqslant \kappa, \\ -x/[2\sqrt{2}(1 + \kappa^{-1})m(v)], & \sqrt{2}m(v)x > \kappa. \end{cases}$$

*Remark 1.35.* There exist two different behaviours for $\eta_v$.

The first one is a Gaussian behaviour $\bar{\eta}_v \sim \mathcal{N}(0, 1)$, when $x$ is small, i.e. for moderate deviations.

If $\bar{\eta}_v$ was really $\mathcal{N}(0, 1)$, the exponential term should be with a constant $1/2$ and not $1/2(1 + \kappa)$.

The second behaviour, for large $x$, i.e. for large deviations, is a Chi-square, centered and dilated by influence of $v_i$ (exponential).

### 1.3.3.2 Risk Hull Method

In order to present the risk hull minimization, which is an improvement of the URE method, we restrict ourselves to the class of projection estimators. In this case, the URE criterion may be written

$$\mathscr{U}(X, N) = \sum_{k=N+1}^{\infty} (X_k^2 - \varepsilon^2\sigma_k^2) + \varepsilon^2\sum_{k=1}^{N}\sigma_k^2.$$

This corresponds in fact to the minimization in $N$ of

$$\bar{R}(X, N) = \sum_{k=N+1}^{\infty} (X_k^2 - \varepsilon^2\sigma_k^2) + \varepsilon^2\sum_{k=1}^{N}\sigma_k^2 - \sum_{k=1}^{\infty}(X_k^2 - \varepsilon^2\sigma_k^2)$$

and then

$$\bar{R}(X,N) = -\sum_{k=1}^{N} X_k^2 + 2\varepsilon^2 \sum_{k=1}^{N} \sigma_k^2.$$

There exists a more general approach which is very close to the URE. This method is called **method of penalized empirical risk**, and in the context of our problem it provides us with the following bandwidth choice

$$N = \arg\min_{N \geqslant 1} \bar{R}_{pen}(X,N), \quad \bar{R}_{pen}(X,N) = \left\{ -\sum_{k=1}^{N} X_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 + \text{pen}(N) \right\}, \quad (1.68)$$

where pen($N$) is a penalty function. The modern literature on this method is very vast and we refer interested reader to [8]. The main idea at the heart of this approach is that severe penalties permit to improve substantially the performance of URE. However, it should be mentioned that the principal difficulty of this method is related to the choice of the penalty function pen($N$). In this context, the URE criterion corresponds to a specific penalty called the URE penalty

$$\text{pen}_{ure}(N) = \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2.$$

The idea is usually to choose a heavier penalty, but the choice of such a penalty is a very sensitive problem, and as we will see later, especially in the inverse problems context.

In [26], a more general approach is proposed, called **Risk Hull Minimization** (RHM) which gives a relatively good strategy for the choice of the penalty. The goal is to present heuristic and mathematical justifications of this method.

The heuristic motivation of the RHM approach is based on the oracle approach.

Consider here only the family of projection estimators $\hat{\theta}(N), N \geqslant 1$. Suppose there is an oracle which provides us with $\theta_k$. In this case the oracle bandwidth is evidently given by

$$N_{or} = \arg\min_{N} r(X,N), \text{ where } r(X,N) = \|\hat{\theta}(N) - \theta\|^2.$$

This oracle mimimizes the loss and is even better than the oracle of the risk. Let us try to mimic this bandwidth choice. At the first glance this problem seems hopeless since in the decomposition

$$r(X,N) = \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 \xi_k^2, \quad (1.69)$$

neither $\theta_k^2$ nor $\xi_k^2$ are really known. However, suppose for a moment, that we know all $\theta_k^2$, and try to minimize $r(X,N)$. Since $\xi_k^2$ are assumed to be unknown, we want to find an upper bound. It means that we minimize the following non-random functional

$$l(\theta,N) = \sum_{k=N+1}^{\infty} \theta_k^2 + V(N), \tag{1.70}$$

where $V(N)$ bounds from above the stochastic term $\varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 \xi_k^2$. It seems natural to choose this function such that

$$\mathbf{E}\sup_{N}\left[\varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 \xi_k^2 - V(N)\right] \leqslant 0, \tag{1.71}$$

since then we can easily control the risk of any projection estimator with any data-driven bandwidth $N^\star$

$$\mathbf{E}_\theta \|\hat{\theta}(N^\star) - \theta\|^2 \leqslant \mathbf{E}_\theta l(\theta,N^\star). \tag{1.72}$$

This motivation leads to the following definition:

**Definition 1.13.** A non random function $\ell(\theta,N)$ is called **risk hull** if

$$\mathbf{E}_\theta \sup_{N}[r(X,N) - \ell(\theta,N)] \leqslant 0.$$

Thus, we can say that $l(\theta,N)$ defined by (1.70) and (1.71) is a risk hull. Evidently, we want to have the upper bound (1.72) as small as possible. So, we are looking for a rather small hull. Note that this hull strongly depends on $\sigma_k^2$.

Once $V(N)$ satisfying (1.71) has been chosen, the minimization of $l(\theta,N)$ can be completed by the standard way using the unbiased estimation. Note that our problem is reduced to minimization of $-\sum_{k=1}^{N} \theta_k^2 + V(N)$. Replacing the unknown $\theta_k^2$ by their unbiased estimates $X_k^2 - \varepsilon^2 \sigma_k^2$, we arrive at the following method of adaptive bandwidth choice

$$\bar{N} = \arg\min_{N}\left[-\sum_{k=1}^{N} X_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 + V(N)\right].$$

In the framework of the empirical risk minimization in inverse problems, the RHM can be defined as follows. Let the penalty in (1.68) be for any $\alpha > 0$

$$V(N) = \text{pen}_{rhm}(N) = \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 + (1+\alpha)U_0(N), \tag{1.73}$$

where

$$U_0(N) = \inf\left\{t > 0: \ \mathbf{E}\left(\eta_N I(\eta_N \geqslant t)\right) \leqslant \varepsilon^2 \sigma_1^2\right\}, \tag{1.74}$$

with

$$\eta_N = \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 (\xi_k^2 - 1). \tag{1.75}$$

This RHM penalty corresponds in fact to the URE penalty plus some term $(1+\alpha)U_0(N)$. One may prove that (see [26]) when $N \to \infty$

$$U_0(N) \approx \left( 2\varepsilon^4 \sum_{k=1}^{N} \sigma_k^4 \log \left( \frac{\sum_{k=1}^{N} \sigma_k^4}{2\pi \sigma_1^4} \right) \right)^{1/2}. \tag{1.76}$$

The RHM chooses the bandwidth $N_{rhm}$ according to (1.68) with the penalty function defined by (1.73) and (1.74). The estimator $\theta_{rhm}^{\star}$ is then defined by

$$\theta_k^{\star} = I(k \leqslant N_{rhm}) X_k. \tag{1.77}$$

The following oracle inequality provides an upper bound for the mean square risk of this approach.

**Theorem 1.9.** *Suppose that $\sigma_k \asymp k^\beta$. Let RHM bandwidth choice $N_{rhm}$ according to (1.68) with the penalty function defined by (1.73) and $\theta_{rhm}^{\star}$ the associated projection estimator defined in (1.77).*

*There exist constants $C_* > 0$ and $\delta_0 > 0$ such that for all $\delta \in (0, \delta_0]$ and $\alpha > 1$*

$$\mathbf{E}_\theta \|\theta_{rhm}^{\star} - \theta\|^2 \leqslant (1+\delta) \inf_{N \geqslant 1} R_\alpha(\theta, N) + C_* \varepsilon^2 \left( \frac{1}{\delta^{4\beta+1}} + \frac{1}{\alpha - 1} \right), \tag{1.78}$$

*where*

$$R_\alpha(\theta, N) = \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 + (1+\alpha) U_0(N).$$

*Proof.* Many of the details are deleted, in order to keep only the idea behind the risk hull. The proof in its full length can be found in [26].

The proof is now in two parts:

The first part is to prove the following lemma.

**Lemma 1.8.** *We have, for any $\alpha > 0$,*

$$l_\alpha(\theta, N) = \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 + (1+\alpha) U_0(N) + \frac{C\varepsilon^2}{\alpha}.$$

*is a risk hull, where $C > 0$ is a positive constant.*

*Proof.* Using (1.69) and (1.70), remark that

$$\mathbf{E} \sup_{N} (\eta_N - (1+\alpha) U_0(N))_+ \leqslant \frac{C\varepsilon^2}{\alpha}$$

implies

$$\mathbf{E}_\theta \sup_{N} (r(X, N) - l_\alpha(\theta, N))_+ \leqslant 0.$$

We have

$$\mathbf{E} \sup_{N} (\eta_N - (1+\alpha) U_0(N))_+ \leqslant \sum_{N=1}^{\infty} \mathbf{E} (\eta_N - (1+\alpha) U_0(N))_+. \tag{1.79}$$

The definition of $U_0(N)$ in (1.74) implies

$$\mathbf{E}(\eta_N - U_0(N))_+ \leqslant \mathbf{E}(\eta_N I(\eta_N \geqslant U_0(N))) \leqslant \varepsilon^2 \sigma_1^2.$$

Moreover, by integrating by parts we obtain

$$\mathbf{E}(\eta_N - (1+\alpha)U_0(N))_+ = \int_{(1+\alpha)U_0(N)}^{\infty} \mathbf{P}(\eta_N > x)dx. \qquad (1.80)$$

Denote by

$$M_N = \varepsilon^2 \max_{k=1,\dots,N} \sigma_k^2$$

and

$$\Sigma_N = \varepsilon^4 \sum_{k=1}^{N} \sigma_k^4.$$

Since the inverse problem is mildly ill-posed, one obtains, by use of (1.46), as $N \to \infty$,

$$M_N \asymp \varepsilon^2 N^{2\beta}, \qquad (1.81)$$

$$\Sigma_N \asymp \varepsilon^2 \sum_{k=1}^{N} k^{4\beta} \asymp \varepsilon^4 N^{4\beta+1} \qquad (1.82)$$

and, using (1.76),

$$U_0(N) \asymp \varepsilon^2 N^{2\beta+1/2} \sqrt{\log N}. \qquad (1.83)$$

Considering only the family of projection estimators, we get another version of Lemma 1.7 with $\kappa = 1/4$.

**Lemma 1.9.** *We have*

$$\mathbf{P}(\eta_N > x) \leqslant \begin{cases} \exp\left(-\frac{x^2}{5\Sigma_N}\right) & 0 \leqslant x \leqslant \frac{\Sigma_N}{4M_N}, \\ \exp\left(-\frac{x}{20M_N}\right) & x > \frac{\Sigma_N}{4M_N}. \end{cases} \qquad (1.84)$$

Remark that, due to (1.81)-(1.83), $U_0(N) \leqslant \Sigma_N/4M_N$, when $N$ is large.

We can then divide in two parts the integral in (1.80),

$$\int_{(1+\alpha)U_0(N)}^{\infty} \mathbf{P}(\eta_N > x)dx =$$

$$= \int_{(1+\alpha)U_0(N)}^{\frac{\Sigma_N}{4M_N}} \mathbf{P}(\eta_N > x)dx + \int_{\frac{\Sigma_N}{4M_N}}^{\infty} \mathbf{P}(\eta_N > x)dx. \qquad (1.85)$$

When $x > \Sigma_N/4M_N$, we have, when $N \to \infty$,

$$\int_{\Sigma_N/4M_N}^{\infty} \exp\left(-\frac{x}{20M_N}\right)dx \leqslant CM_N \exp\left(-C\frac{\Sigma_N}{M_N^2}\right) \asymp C\varepsilon^2 N^{2\beta} \exp(-CN). \qquad (1.86)$$

Moreover

$$\int_{(1+\alpha)U_0(N)}^{\infty} \exp\left(-\frac{x^2}{5\Sigma_N}\right) dx \leqslant \int_{(1+\alpha)U_0(N)}^{\infty} \frac{x}{(1+\alpha)U_0(N)} \exp\left(-\frac{x^2}{5\Sigma_N}\right) dx$$

$$\leqslant \frac{5\Sigma_N}{2(1+\alpha)U_0(N)} \exp\left(-\frac{(1+\alpha)^2 U_0(N)^2}{5\Sigma_N}\right).$$

Thus, using (1.76), we obtain

$$\int_{(1+\alpha)U_0(N)}^{\infty} \exp\left(-\frac{x^2}{5\Sigma_N}\right) dx$$

$$\leqslant C\sqrt{\Sigma_N} \exp\left(-\frac{2}{5}(1+\alpha)^2 \log\left(\frac{\sum_{k=1}^{N}\sigma_k^4}{2\pi\sigma_1^4}\right)\right). \tag{1.87}$$

Using (1.82), remark that the term in (1.86) is smaller than the one in (1.87), as $N \to \infty$. Using (1.79), (1.80) and (1.85), we then obtain

$$\mathbf{E}\sup_N \left(\eta_N - (1+\alpha)U_0(N)\right)_+ \leqslant \sum_{N=1}^{\infty} C\varepsilon^2 \exp\left(-\left(\frac{2}{5}(1+\alpha)^2 - \frac{1}{2}\right)\log(N)\right).$$

Thus, for $\alpha$ large enough ($\alpha > 2$), the term is then summable in $N$ and we obtain

$$\mathbf{E}\sup_N \left(\eta_N - (1+\alpha)U_0(N)\right)_+ \leqslant \frac{C\varepsilon^2}{\alpha}.$$

The proof for $\alpha > 0$ small is much more technical and based on chaining arguments (see [26]).

In the second part of the proof of Theorem 1.9, we need to prove that we are able to minimize this risk hull based on the data. Since $l_\mu(\theta, N)$ is a risk hull for any $\mu > 0$ we have

$$l_\mu(\theta, N) = \sum_{k=N+1}^{\infty} \theta_k^2 + \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 + (1+\mu)U_0(N) + \frac{C\varepsilon^2}{\mu}, \tag{1.88}$$

and therefore

$$\mathbf{E}_\theta \|\hat{\theta}(N_{rhm}) - \theta\|^2 \leqslant \mathbf{E}_\theta l_\mu(\theta, N_{rhm}). \tag{1.89}$$

On the other hand, since $N_{rhm}$ minimizes $\bar{R}_{pen}(X, N)$, we have for any integer $N$

$$\mathbf{E}_\theta \bar{R}_{pen}(X, N_{rhm}) \leqslant \mathbf{E}_\theta \bar{R}_{pen}(X, N) = R_\alpha(\theta, N) - \|\theta\|^2. \tag{1.90}$$

In order to combine the inequalities (1.89) and (1.90), we rewrite $l_\mu(\theta, N_{rhm})$ in terms of $\bar{R}_{pen}(X, N_{rhm})$,

$$l_\mu(\theta, N_{rhm}) = \bar{R}_{pen}(X, N_{rhm}) + \|\theta\|^2 + \frac{C\varepsilon^2}{\mu}$$

$$+2\varepsilon \sum_{k=1}^{N_{rhm}} \sigma_k \theta_k \xi_k + \varepsilon^2 \sum_{k=1}^{N_{rhm}} \sigma_k^2(\xi_k^2 - 1) - (\alpha - \mu)U_0(N_{rhm}).$$

Therefore, using this equation, (1.89) and (1.90), we obtain that for any integer $N$

$$\mathbf{E}_\theta \|\hat{\theta}(N_{rhm}) - \theta\|^2 \leqslant R_\alpha(\theta, N) + \frac{C\varepsilon^2}{\mu} + \mathbf{E}_\theta 2\varepsilon \sum_{k=1}^{N_{rhm}} \sigma_k \theta_k \xi_k$$

$$+ \mathbf{E}_\theta \left[ \varepsilon^2 \sum_{k=1}^{N_{rhm}} \sigma_k^2(\xi_k^2 - 1) - (\alpha - \mu)U_0(N_{rhm}) \right].$$

The next step is to control the last two terms in the above equation. This part of proof is not done here (see [26]).

This control should be done for any data-driven choice $N^\star$ (or $N_{rhm}$), this is why these terms are difficult to control. Moreover, to get a sharp oracle inequality, one has to be rather precise.

The first term, of the last two terms, may be included in the left term (the risk of the RHM estimator) and in the remainder term. However, this part of the proof is one of the more delicate. One really has to control this stochastic term for any data-driven $N^\star$ (see [26]).

The second term, of the last two terms, is very close to Lemma 1.8 and its proof. Thus, we may use again the risk hull in order to control it.

As noted, a major contribution of oracle inequalities is that they usually allow to construct rather easily minimax adaptive estimators. Here the proof is very simple because the family of estimators corresponds to all possible choices of $N$.

**Theorem 1.10.** *Suppose $\sigma_k \asymp k^\beta$, $\beta \geqslant 0$. Let RHM bandwidth choice $N_{rhm}$ according to (1.68) with the penalty function defined by (1.73) and $\theta_{rhm}^\star$ the associated projection estimator defined in (1.77).*

*Assume that $\theta$ belongs to the ellipsoid $\Theta(\alpha, L)$, where $a_k \asymp k^\alpha$, $\alpha > 0$, $L > 0$, defined in (1.40). Then the RHM estimator $\theta_{rhm}^\star$ verifies, for any $\alpha > 0$ and $L > 0$, as $\varepsilon \to 0$,*

$$\sup_{\theta \in \Theta(\alpha, L)} \mathbf{E}_\theta \|\theta_{rhm}^\star - \theta\|^2 \leqslant C\varepsilon^{4\alpha/(2\alpha+2\beta+1)}.$$

*This rate is optimal (see Theorem 1.5).*

*Thus, the RHM estimator is then minimax adaptive on the class of ellipsoid.*

*Proof.* The proof is a direct consequence of Theorem 1.6 and 1.9. One has to note that, due to (1.83), $U_0(N) = o(\varepsilon^2 \sum_{k=1}^N \sigma_k^2)$ as $N \to \infty$. Asymptotically, the RHM penalty is negligible as compared to the URE penalty. Thus, the penalized oracle $R_\alpha(\theta, N)$ on the right hand side of Theorem 1.9 still attains the optimal rate of convergence.

*Remark 1.36.* In order to construct a sharp adaptive estimator on ellipsoids, one has to obtain results for the Pinsker family. The RHM method has been extended to the Pinsker family in [99].

## Comments

*Penalized oracle.* We have an oracle inequality but with a penalty term on the RHS. This is usually called a (penalized) oracle inequality. This is standard in the penalized empirical risk approach. At the first sight, the result may look weaker than in Theorem 1.7. Indeed, the main term is a penalized oracle here when it was the true oracle in Theorem 1.7. However, here the remainder term is better. In Theorem 1.7, the remainder term depends on the cardinality and on the complexity of the family of estimators. In Theorem 1.9, there is no such price, and moreover the family may be infinite. However, as will be explained by simulations in Section 1.3.3.3, the even more important point is that the constant is much more under control than in Theorem 1.7.

*Natural penalty.* By (1.76) the penalty $U_0(N)$ is almost of the order of the standard deviation of the empirical risk. This seems rather natural, since it really controls the behaviour of the empirical risk, i.e. not only its expectation but also its standard deviation.

*Second order penalty.* We have $U_0(N) = o(\varepsilon^2 \sum_{k=1}^{N} \sigma_k^2)$ as $N \to \infty$, since (1.83). We add a penalty (see (1.73)) which is small compared to the URE penalty. In fact, the RHM penalty may be thought as the URE penalty plus a second order penalty. From an asymptotical point of view, there is no real difference between the URE and the RHM. Thus, the two methods should be very close. A consequence of the previous remark, is that the (penalized) oracle inequality is then (asymptotically) as sharp as the one in Theorem 1.7. Asymptotically, one may obtain exactly the same results, since the penalty is smaller. Thus, minimax adaptive estimators may be constructed directly (see Theorem 1.10).

*Direct problem.* In the direct problem ($A = I$), i.e. in Gaussian white noise model, due to (1.76), the penalty is then:

$$Pen_{rhm}(N) = \varepsilon^2 N + (1 + \alpha)U_0(N),$$

where

$$U_0(N) \approx \left( 2\varepsilon^4 N \log \frac{N}{2\pi} \right)^{1/2}.$$

One may see that we really add a second order penalty.

*Difference between RHM and URE.* On the one hand, the previous remarks show that the RHM penalty is equal to the URE penalty plus a small term (compared to the URE penalty). On the other hand, there exist main differences between the two estimators, especially in the case of inverse problems. RHM is much more stable than URE (see Section 1.3.3.3). Moreover, in the simulations, it is always more accurate, even in the direct problem. However, the difference is much more important in ill-posed framework.

*Asymptotics in inverse problems.* One of the reason for its instability is that URE is based on some asymptotical ideas. In inverse problems, usually $N$ is not very large, due to the increasing noise. Indeed, in the ill-posed context, the

term $\sigma_k \to \infty$. It means that the noise is really increasing with $k$. One has to be very careful with high frequencies. More or less, it is very difficult to choose a large number of coefficients $N$. On the one hand, the minimax theory, claims that the optimal choice of $N$ is going to infinity in nonparametric statistics (see for example Theorem 1.6). On the other hand, the choice of $N$ cannot be too large, otherwise, in real inverse problems the noise will explose.

Thus, one has to be very careful with asymptotics in inverse problems.

*Penalty computed by Monte Carlo.* The penalty $U_0(N)$ may be computed by Monte Carlo simulations. Indeed, the definition of $U_0(N)$ in (1.74) has no explicit solution. There exists an approximation of $U_0(N)$ in (1.76), but it is true for $N$ large enough. As noted, $N$ is not so large in inverse problems. Thus, a more careful and accurate way to compute $U_0(N)$ is by use of Monte Carlo. It is a bit time consuming, but it is done only once for one given inverse problem.

*Explicit penalty.* By use of RHM we obtain, an explicit penalty which comes from the proof of Theorem 1.9. It is really by looking inside the proof of Lemma 1.8 that one may understand the penalty form. The constraint, that one wants to have a risk hull really help in choosing such a penalty.

Another very important point, is that after, this penalty may be used directly in simulations. The method, really gives, an explicit penalty. There is no gap between the penalty needed in Theorem 1.9 and the one used in the simulation study.

### 1.3.3.3 Simulations

In order to illustrate the difference between direct and inverse estimation, we will carry out a very simple numerical experiment. Obviously, we cannot compute it for all $\theta \in \ell^2$. Therefore, let us take $\theta_k \equiv 0$ and compute the ratio between the risk and the risk of the oracle for two cases $\sigma_k \equiv 1$ and $\sigma_k = k$. The first case corresponds to classical function estimation (direct estimation), whereas the second is related to the estimation of the first order derivative of a function (inverse estimation). Notice that in both cases the risk of the oracle is clearly $\inf_N R(0,N) = \varepsilon^2$ since $\operatorname{argmin}_N R(0,N) = 1$.

In order to study the performance of the URE, we generate 2000 independent random vectors of $y^j$, $j = 1, ..., 2000$ with the components defined by (1.5). For each vector we compute $N_{ure}(y^j)$ and the normalized error $\| \hat{\theta}[N_{ure}(y^j)] - \theta \|^2 / \varepsilon^2$ and plot these values as a stem diagram. We also compute the mean empirical bandwidth $N_{emp}$ and the normalized mean empirical risk $R_{emp}$ by

$$N_{emp} = \frac{1}{2000} \sum_{j=1}^{2000} N_{ure}(y^j), \quad R_{emp} = \frac{1}{2000\varepsilon^2} \sum_{j=1}^{2000} \| \hat{\theta}[N_{ure}(y^j)] - \theta \|^2.$$

Let us discuss briefly the numerical results of this experiment shown on Figure 1.1. The first display (direct estimation) shows that the URE method works reasonably well. Almost all bandwidths $N_{ure}(y^j)$ are relatively small (their mean is 1.98)
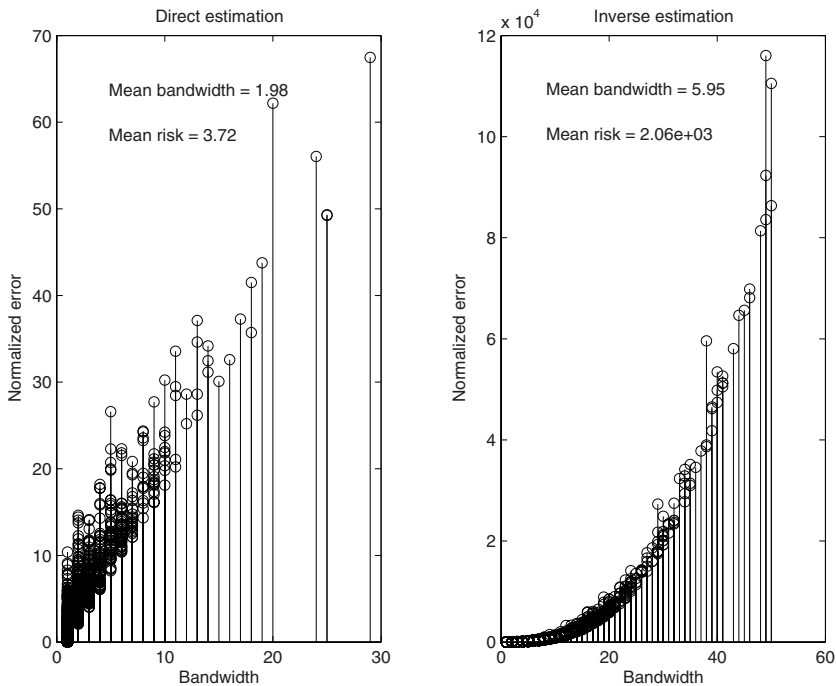
**Fig. 1.1** The method of unbiased risk estimation

and the normalized error is 3.72. The second display shows that the distribution of $N_{ure}(y^j)$ changed essentially. Now the mean bandwidth is 5.95 and there are sufficiently many bandwidths $N_{ure}(y^j)$ greater than 20. This results in a catastrophic normalized error around 2000.

In this section, we present some numerical properties of the RHM approach. We will study in a more general context than the previous no-signal one, i.e. $\theta_k \equiv 0$. Numerical testing of nonparametric statistical methods is a very difficult and delicate problem. The goal of this section is to illustrate graphically Theorem 1.7 and Theorem 1.9. To do that, we propose to measure statistical performance of a method $N^\star$ by its *oracle efficiency* defined by

$$e_{or}(\theta, N^\star) = \frac{\inf_N \mathbf{E}_\theta \|\hat{\theta}(N) - \theta\|^2}{\mathbf{E}_\theta \|\hat{\theta}(N^\star) - \theta\|^2}.$$

If the oracle efficiency of a method is close to 1 then the risk is very close to the risk of the oracle.

It should be mentioned that we use the inverse of the previous ratio since we want to get a good graphical representation of the performance. We have just seen in the previous part that the ratio can vary from 1 to 2000 for the URE method. This results

in a degenerate plot. Therefore, in order to avoid this effect, we use this definition of the oracle efficiency $e_{or}(\theta, N^\star)$.

Since it is evidently impossible to compute the oracle efficiency for all $\theta \in \ell^2$, we choose a sufficiently representative family of vectors $\theta$. In what follows we will use the following family, with polynomial decreasing,

$$\theta_k^a = \frac{a\varepsilon}{1 + (k/W)^m},$$

where $\varepsilon$ is the noise level, $a$ is called amplitude, $W$ bandwidth, and $m$ smoothness.

We shall vary $a$ in a large range and plot $e_{or}(\theta^a, N^\star)$ as a function of $a$ which is directly related with the signal-to-noise ratio in the model considered. In a statistical framework $a^2$ would be $n$ the number of observations. The parameters $m = 6$ and $W = 6$ are fixed. Many other examples of $(W, m)$ were looked at, simulations showed that the oracle efficiency exhibits similar behaviour.

Two methods of data-driven bandwidth choice will be compared: the URE and the RHM with $\alpha = 1.1$. One may note that for these methods $e_{or}(\theta^a, N^\star)$ does not depend on $\varepsilon$. This function was computed by the Monte Carlo method with 40000 replications.

We start with the direct estimation where $\sigma_k \equiv 1$. Figure 1.2 shows the oracle efficiency of the URE (left panel) and the oracle efficiency of the RHM (right panel). Comparing these plots, one can say that both methods work reasonably well. Both efficiencies are very close to 1. The risk of URE method is around $1/0.75 = 1.33$ times the risk of oracle, when RHM method is around $1/0.82 = 1.22$ times the oracle. Thus RHM is always better than URE, but the ratio is something like 5% to 10%.

However, if we deal with an inverse problem ($\sigma_k = k$), we can already see a significant difference between these methods. The corresponding oracle efficiencies are plotted on the left and on the right panels of Figure 1.3. For small values $a$ the performance of the URE is very poor, whereas the RHM demonstrates a very stable behaviour. For very large $a = 500$ the oracle efficiency of the URE is of order 0.16, which means that its risk is around 6 times the one of the oracle. For smaller $a = 100$, it is around 10 times the oracle. In the meantime, the RHM has always an efficiency greater than 0.4 and usually around 0.5, i.e. 2 times the risk of the oracle.

The last Figure 1.4 deals with the case when the inverse problem is more ill-posed ($\sigma_k = k^2$) . In this situation the URE fails completely. Its maximal oracle efficiency is of order $3 * 10^{-4}$, i.e. 10000 times the oracle. Nevertheless, the RHM has a good efficiency (greater than 0.3). Its risk is then around 3 times those of the oracle.

Another remark is that the RHM is really stable compared to the increasing degree of ill-posedness $\beta$ of the problem. The efficiency is worse when the inverse problem is more difficult, but it is always reasonable. The behaviour of URE is completely different, it really exploses with $\beta$.

One may also see that URE is really based on asymptotic ideas. Indeed, its oracle efficiency is highly increasing with the amplitude $a$. On the other hand, RHM is stable, and does not rely on large values of $a$.
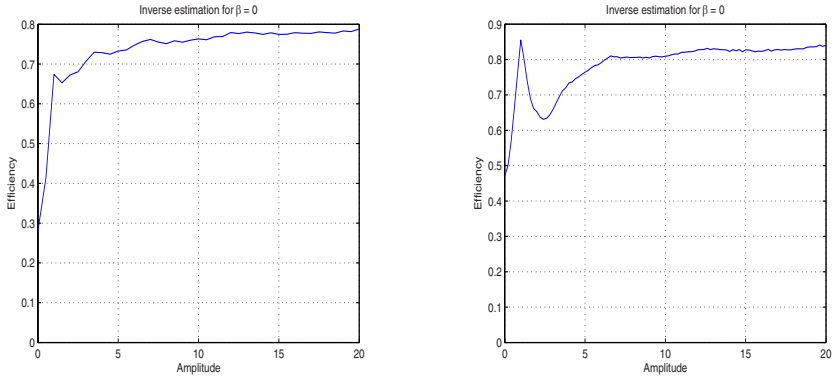
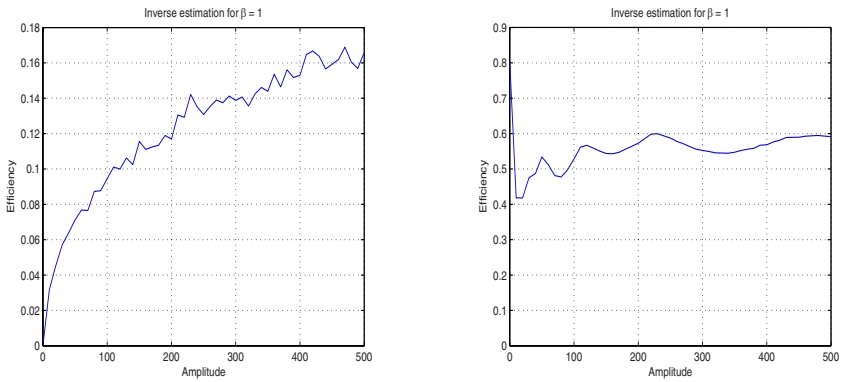**Fig. 1.2** Oracle efficiency of URE (left) and RHM (right) for direct estimation.



**Fig. 1.3** Oracle efficiency of URE and of RHM for inverse problem ($\beta = 1$).
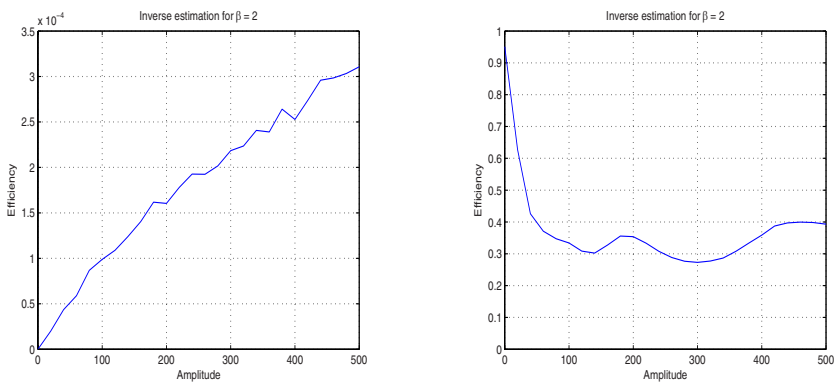


**Fig. 1.4** Oracle efficiency of URE and of RHM for inverse problem ($\beta = 2$).

This simulation study, shows that there is a huge difference between the two methods, at least in inverse problems. This may be surprising, since the RHM penalty, was supposed to be of second order. Then, the two methods should be closely related. However, this point of view, mainly relies on asymptotic ideas. As noted before, in inverse problems, one has to be really careful with asymptotics. This may really be seen here, where these two methods have a very different behaviour.

In the context of Theorem 1.7 and Theorem 1.9, this example shows also that the constants which appear in the remainder terms are quite different. The one in Theorem 1.9, seems to be small and really under control. While the one in Theorem 1.7, $C_*$ is in fact really large. Unfortunately, it means that the terms which are asymptotically small in Theorem 1.7 may dominate the risk of oracle.

## 1.3.4 Universal Optimality

### 1.3.4.1 Blockwise Estimators

In this section, we present a more general approach to optimality. Namely, we construct a sequence of weights $\lambda_{pbs}^{\star}$ such that the penalized blockwise Stein estimator $\theta_{pbs}^{\star} = \hat{\theta}(\lambda_{pbs}^{\star})$ satisfies both some exact oracle inequalities (for typical examples of classes $\Lambda$) and the (sharp) minimax adaptivity property (Definition 1.12) (for a large scale of classes $\Theta_\alpha$).

An important fact is that the estimator does not belong to either of the typical classes $\Lambda$ but it outperforms the oracles $\lambda_0$ corresponding to these classes. This property can be called **universal optimality** over a large scale of classes $\Lambda$. This point of view is different from the model selection ideas in Section 1.3.3, where the data-driven choice take its values in the family $\Lambda$.

An important point here, is to find a large family $\Lambda$ in order to obtain oracle inequalities valid for many different estimators. The first step is close to the approach of unbiased risk estimation. Indeed, one would like to minimize the criteria $\mathscr{U}(X, \lambda)$ on such a family.

What is the reasonable set of $\lambda$ where the minimization of $\mathscr{U}(X, \lambda)$ should be done? Minimizing $\mathscr{U}(X, \lambda)$ with respect to all possible $\lambda$ yields $\lambda_k = (1 - \varepsilon^2/X_k^2)$ or $\lambda_k = (1 - \varepsilon^2/X_k^2)_+$ if we restrict the minimization to $\lambda_k \in [0, 1]$. It is easy to see that the risk of the estimator $\{\lambda_k X_k\}$ is diverging if the sum is taken over all $k$ and is at least as great as $\varepsilon^2 N$ if one considers the sum over $k \leqslant N$ for some integer $N$ in the definition of $\mathscr{U}(X, \lambda)$. Since $N$ should be chosen in advance, such an estimator has poor adaptation properties, and minimizing over all $\lambda$ makes no sense.

A more fruitful idea is to minimize $\mathscr{U}(X, \lambda)$ in a restricted class of sequences, for example over one of the classes $\Lambda$ discussed in Section 1.2.2.

Choose $\lambda^{\star}$ as a minimizer of $\mathscr{U}(X, \lambda)$ over $\lambda \in \Lambda$ in order to mimic the linear oracle on $\Lambda$. However, this principle is difficult to apply for huge classes, such as $\Lambda_{mon}$, the class of monotone weights. [59] suggests the minimization of $\mathscr{U}(X, \lambda)$ on a truncated version of the class $\Lambda_{mon}$.

The search for more economic but yet huge enough subclasses of weight sequences $\lambda$ leads in particular to the family of blockwise constant weights which can interpreted as sieves over various sets of $\lambda$. Blockwise constant weights have been discussed in statistical literature starting from [44], and more recently by [47, 105]; for wavelets, see [40, 78, 67].

The key feature of our estimator is that it "mimics" the *monotone oracle* $\lambda_0^{mon}$ defined as a solution of

$$R(\theta, \lambda_0^{mon}) = \min_{\lambda \in \Lambda_{mon}} R(\theta, \lambda), \qquad (1.91)$$

where $\Lambda_{mon}$ is the class of monotone sequences. Consider the class of monotone weights sequences

$$\Lambda_{mon} = \{\lambda = \{\lambda_k\} \in \ell^2 : 1 \geqslant \lambda_1 \geqslant \ldots \geqslant \lambda_k \ldots \geqslant 0\},$$

and the class of *monotone estimators*

$$\hat{\theta}_k = \lambda_k \, X_k,$$

where $\{\lambda_k\} \in \Lambda_{mon}$ and $X_k$ is defined in (1.6).

If the coefficients $\theta_k$ are monotone non-increasing, remark that the monotone oracle is equal to the linear oracle.

Restrict the attention to the class $\Lambda_{mon}$ since it contains the most interesting examples of weight sequences $\{\lambda_k\}$. The projection weights and the Tikhonov weights belong to $\Lambda_{mon}$ (see Section 1.2.2). Next, typically $\sigma_k$ are monotone non-decreasing and $a_k$ in the definition of the ellipsoid in (1.40) are monotone non-decreasing. The Pinsker weights also belong to $\Lambda_{mon}$. It can be shown that some minimax solutions on other bodies in $\ell^2$ than ellipsoids (e.g. parallelepipeds) are also in $\Lambda_{mon}$, see [31].

We are looking for an adaptive estimator $\theta^\star = (\theta_1^\star, \theta_2^\star, \ldots)$ of the form

$$\theta_k^\star = \lambda_k^\star \, X_k,$$

where $\lambda_k^\star$ are some data-driven weights.

A well-known idea of choosing $\lambda^\star$ is based on the unbiased estimation of the risk by minimizing criteria $\mathscr{U}(X, \lambda)$ defined in (1.59) among the family $\Lambda$ (see Section 1.3.3.1). The difference here is that the class $\Lambda$ is not some given class of estimators (projection, Tikhonov,...) but the very large class $\Lambda_{mon}$ of monotone estimators.

However, as noted before, this class $\Lambda_{mon}$ is maybe too large. Consider instead, the class $\Lambda_b$ of coefficients with piecewise constant $\lambda_k$ over suitably chosen blocks.

Define the class of **blockwise estimators**

$$\hat{\theta}_k = \lambda_k \, X_k,$$

where $\lambda \in \Lambda_b$ is the set of piecewise constant sequences,

$$\Lambda_b = \{\lambda \in \ell^2 : 0 \leqslant \lambda_k \leqslant 1, \lambda_k = \lambda_{\kappa_j}, \forall k \in I_j, \lambda_k = 0, k > N_{max}\},$$

where $I_j$ denote the block $I_j = \{k \in [\kappa_{j-1}, \kappa_j - 1]\}, j = 0, \ldots, J-1$ and $J, N_{max},$ $\kappa_j, \ j = 0, \ldots, J,$ are integers such that $\kappa_0 = 1, \kappa_J = N_{max} + 1, \kappa_j > \kappa_{j-1}.$

Denote also by $T_j = \kappa_j - \kappa_{j-1}$ the size of the blocks $I_j$, for $j = 1, \ldots, J.$

### 1.3.4.2 Stein's Estimator

In this section, we change to a slightly different model in order to present and discuss the so-called Stein phenomenon. This problem goes back to the work of [120], and has been extended since. However, it is still one of the most surprising result in statistics. This section is based on [18, 129].

Consider the following model, which is a finite version of the sequence space model in the direct case (i.e. $b_k \equiv 1$),

$$y_k = \theta_k + \varepsilon \xi_k, \ k = 1, \ldots, d, \tag{1.92}$$

where $d$ is some integer, $\{\xi_k\}$ are i.i.d. $\mathcal{N}(0,1)$. The statistical problem is to estimate $\theta$ based on the data $y$.

In this simple situation, the Maximum Likelihood Estimator is then $\hat{\theta}_{mle} = y$. This estimator was believed, to be the best possible estimator in this context. Its risk is

$$R(\theta, \hat{\theta}_{mle}) = \varepsilon^2 d, \ \forall \theta \in \mathbb{R}^d.$$

However, [120] discovered a very strange phenomenon. Indeed, he constructed an estimator, the **Stein estimator**

$$\hat{\theta}_S = \left(1 - \frac{\varepsilon^2 d}{\|y\|^2}\right) y, \tag{1.93}$$

for which we have, see proof of Lemma 3.10 in [129],

$$\mathbf{E}_\theta \|\hat{\theta}_S - \theta\|^2 = \varepsilon^2 d - \varepsilon^4 d(d-4)\mathbf{E}_\theta \left(\frac{1}{\|y\|^2}\right),$$

and

$$\mathbf{E}_\theta \|\hat{\theta}_S - \theta\|^2 \leqslant \varepsilon^2 d - \frac{\varepsilon^4 d(d-4)}{\|\theta\|^2 + \varepsilon^2 d}. \tag{1.94}$$

Thus, the main result in [120] is that if $d \geqslant 5$,

$$\mathbf{E}_\theta \|\hat{\theta}_S - \theta\|^2 < \mathbf{E}_\theta \|y - \theta\|^2, \ \forall \theta \in \mathbb{R}^d.$$

This very surprising result proves that the MLE estimator $y$ is not even admissible (for $d \geqslant 5$).

Written in a slightly different framework, [120] discovered that the Stein estimator is better at each point $\theta \in \mathbb{R}^d$ than the mean $X$ (for $d \geqslant 5$).

Looking carefully at (1.94), note that the improvement on $y$ is by a constant at point $\theta_k \equiv 0$, but also if $\|\theta\| \asymp \varepsilon$. However, when $\theta$ is larger, the improvement is

of second order. Nevertheless, this is an asymptotical point of view, and the gain is valid for any $\theta \in \mathbb{R}^d$.

Several versions of the Stein estimator have been defined since then, many of them which improved on the basic estimator. One example is the **positive Stein estimator**

$$\hat{\theta}_s = \left(1 - \frac{\varepsilon^2 d}{\|y\|^2}\right)_+ y. \tag{1.95}$$

The following result may be found in Lemma 3.9 in [129], for all $d \geqslant 1$,

$$\mathbf{E}_\theta \|\hat{\theta}_s - \theta\|^2 < \mathbf{E}_\theta \|\hat{\theta}_S - \theta\|^2, \quad \forall \theta \in \mathbb{R}^d.$$

Another famous version is the **James-Stein estimator** (and its positive version),

$$\hat{\theta}_{JS} = \left(1 - \frac{\varepsilon^2(d-2)}{\|y\|^2}\right) y,$$

see [76], which is better than the MLE estimator $y$ even for $d \geqslant 3$.

*Remark 1.37.* The (positive) Stein estimator has an effect, even if still very surprising, which may be understood. On the one hand, when the whole signal $\|y\|^2$ is large (compared to $\varepsilon^2 d$), then one may rely on the data, and estimate $\theta$ by something very close to $y$. On the other hand, if $\|y\|^2$ is small, then one estimates $\theta$ by something close to 0, or even equal to 0 if $\|y\|^2 \leqslant \varepsilon^2 d$. The information on the whole sequence $\{y_k\}$ helps in estimating a single coefficient $\theta_k$ in a better way than just by using $y_k$.

Moreover, the Stein estimator has a role of moving the data $y$ to 0 by some factor. This effect is known nowadays as the *Stein shrinkage*. The idea is that one shrinks the observations, more or less, towards 0 in order to improve on $y$.

The ideas of Stein have been very successful and popular among statisticians. Moreover, since Stein's result is valid in large dimensions $d$, his ideas are still the topic of a vast literature in nonparametric statistics, where $d$ is very large, even infinite, see [41, 14, 77, 17, 31, 18, 97, 114, 129]. The main common point among these papers, is to try to estimate the infinite sequence $\{\theta_k\}$ by using block estimators, as sieves, see Section 1.3.4.1. Then on each of these blocks, the idea is to estimate the coefficients $\theta_k$ by use of the Stein estimator. As already noted, the nonparametric context is well suited, since then the blocks will be large, and Stein's estimator successful. One of the main difficulties is then related to the choice of the size of the blocks.

### 1.3.4.3 Blockwise Stein's Rules

The construction of the estimator $\theta^\star_{pbs}$ is the following.

Divide the set of coefficients $\theta_k$ into blocks in a proper way, and apply a penalized version of Stein's estimator on each block. The penalty should be rather small but non-zero. The same construction with non-penalized Stein's estimators can be also implemented, but leads to more limited results (see [31]).

Note that the solution $\lambda_{bs}^{\star}$ of the minimization problem

$$\mathscr{U}(X,\lambda_{bs}^{\star}) = \min_{\lambda \in \Lambda_b} \mathscr{U}(X,\lambda)$$

is given by $\lambda_{bs}^{\star} = (\lambda_1^{\star}, \lambda_2^{\star}, \dots)$, where

$$\lambda_k^{\star} = \begin{cases} \left(1 - \dfrac{\Gamma_{(j)}^2}{\|X\|_{(j)}^2}\right)_+ , & k \in I_j, \ j = 1,\dots,J, \\ 0 , & k > N_{max}, \end{cases} \qquad (1.96)$$

with $x_+ = \max(0,x)$,

$$\Gamma_{(j)}^2 = \varepsilon^2 \sum_{k \in I_j} \sigma_k^2, \qquad \|X\|_{(j)}^2 = \sum_{k \in I_j} X_k^2,$$

and

$$\Delta_{(j)} = \frac{\max_{k \in I_j} \sigma_k^2}{\sum_{k \in I_j} \sigma_k^2}.$$

The weights (1.96) define a *blockwise Stein rule*. The **blockwise Stein estimator** is

$$\theta_k^{\star} = \lambda_k^{\star} X_k,$$

where $\lambda_{bs}^{\star} = \{\lambda_k^{\star}\}$ is defined in (1.96).

However, for mildly ill-posed inverse problems, the estimator $\theta_{bs}^{\star}$ can be modified to have better properties.

We now modify the weights $\lambda_{bs}^{\star}$ and define $\lambda_{pbs}^{\star}$ by

$$\lambda_k^{\star} = \begin{cases} \left(1 - \dfrac{\Gamma_{(j)}^2(1+p_j)}{\|X\|_{(j)}^2}\right)_+ , & k \in I_j, \ j = 1,\dots,J, \\ 0 , & k > N_{max}, \end{cases}$$

where $0 \leqslant p_j \leqslant 1$ is some penalty term.

Finally, the estimator has the form $\theta_{pbs}^{\star} = (\theta_1^{\star}, \theta_2^{\star}, \dots)$ where

$$\theta_k^{\star} = \begin{cases} \left(1 - \dfrac{\Gamma_{(j)}^2(1+p_j)}{\|X\|_{(j)}^2}\right)_+ X_k , & k \in I_j, \ j = 1,\dots,J, \\ 0 , & k > N_{max}. \end{cases} \qquad (1.97)$$

This estimator is called the **penalized blockwise Stein estimator**.

*Remark 1.38.* The penalizing factor $(1+p_j)$ forces the estimator to contain fewer nonzero coefficients $\theta_k^{\star}$ than for the usual blockwise Stein's rule (1.96): our estimator is more "sparse". The general choice of the penalty $p_j$ will be $p_j = \Delta_j^a$, where $0 < a < 1/2$. The assumption $a < 1/2$ is important. Intuitively, this effect is easy to explain. If $b_k$ decreases as a power of $k$ we have:

$$\text{standard deviation}(Z_j)/\text{expectation}(Z_j) \asymp \Delta_{(j)}^{1/2}$$

where $Z_j$ is the stochastic error term corresponding to $j$th block. Hence, to control the variability of stochastic terms, one needs a penalty that is slightly larger than $\Delta_{(j)}^{1/2}$.

More general penalties are presented in [31].

### 1.3.4.4 Construction of Blocks

Introduce now a *special construction of blocks* $I_j$ which may be called *weakly geometrically increasing blocks*. In Theorem 1.11 we will show that with this construction the penalized blockwise Stein estimator verifies an oracle inequality. This construction (or some versions of it) is used by [105] but also in [58, 32, 114].

Let $\nu_\varepsilon$ be an integer valued function of $\varepsilon$ such that $\nu_\varepsilon \geqslant 5$ and $\nu_\varepsilon \to \infty$ as $\varepsilon \to 0$. A typical choice would be $\nu_\varepsilon \asymp \log(1/\varepsilon)$ or $\nu_\varepsilon \asymp \log\log(1/\varepsilon)$. Let

$$\rho_\varepsilon = \frac{1}{\log \nu_\varepsilon}.$$

Clearly, $\rho_\varepsilon \to 0$ as $\varepsilon \to 0$. Define the sequence $\{\kappa_j\}$ by

$$\kappa_j = \begin{cases} 1 & j = 0, \\ \nu_\varepsilon & j = 1, \\ \kappa_{j-1} + \lfloor \nu_\varepsilon \rho_\varepsilon (1+\rho_\varepsilon)^{j-1} \rfloor & j = 2,\ldots, \end{cases} \tag{1.98}$$

where $\lfloor x \rfloor$ is the maximal integer that is strictly less than $x$. Let $\bar{N}$ be any integer satisfying

$$\bar{N} \geqslant \max\{N : \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 \leqslant \rho_\varepsilon^{-3}\}. \tag{1.99}$$

Then, for $\varepsilon$ small enough, $\bar{N} \geqslant \max\{N : \varepsilon^2 \sum_{k=1}^{N} \sigma_k^2 \leqslant r^2 \rho_\varepsilon^{-2}\}$, $\forall r > 0$.

*Remark 1.39.* The term $\bar{N}$ is the final value of $k$. After that, the estimator is always fixed at 0. This $\bar{N}$ is fixed with the idea that the variance of a projection estimator $\hat{\theta}(N)$, i.e. $\varepsilon^2 \sum_{k=1}^{N} \sigma_k^2$, cannot be too large. Otherwise it is not even useful to consider larger values of $N$. Indeed, a good projection estimator should have a variance going to zero.

In this special construction assume the following:

**(B1)** *The blocks are $I_j = [\kappa_{j-1}, \kappa_j - 1]$ such that the values $\kappa_j$ satisfy (1.98), and $J = \min\{j : \kappa_j > \bar{N}\}$ where $\bar{N}$ satisfies (1.99).*
   Clearly, $N_{max} = k_J - 1 \geqslant \bar{N}$ if (B1) holds.
**(B2)** *The penalty is $p_j = \Delta_{(j)}^a$, where $0 < a < 1/2$.*

We also assume that the singular values $b_k$ decrease precisely as a power of $k$:

**(B3)** *The coefficients $b_k$ are positive and there exist $\beta \geqslant 0, b_* > 0$ such that*

$$b_k = b_* k^{-\beta}(1 + o(1)), \ k \to \infty.$$

**Theorem 1.11.** *Let $\theta^\star_{pbs}$ be the penalized blockwise Stein estimator defined in (1.97). Assume (B1),(B2) and (B3), and let $r > 0$ be fixed. Then:*

*(i) For any $\theta \in \ell^2$ such that $\|\theta\| \leqslant r$ and any $0 < \varepsilon < 1$ such that $\Delta_{(j)} \leqslant (1 - p_j)/4$ for all $j$, we have*

$$\mathbf{E}_\theta \|\theta^\star_{pbs} - \theta\|^2 \leqslant (1 + \tau_\varepsilon) \inf_{h \in \Lambda_{mon}} R(\theta, \lambda) + c\varepsilon^2 v_\varepsilon^{2\beta+1},$$

*where $c > 0$ does not depend on $\theta, \varepsilon$, and $\tau_\varepsilon = o(1)$, $\varepsilon \to 0$, $\tau_\varepsilon$ does not depend on $\theta$.*

*(ii) For any $\lambda \in \Lambda_{mon}$ and $\theta \in \ell^2$ such that $R(\theta, \lambda) \leqslant r^2$ and any $0 < \varepsilon < 1$ such that $\Delta_{(j)} \leqslant (1 - p_j)/4$ for all $j$, we have*

$$\mathbf{E}_\theta \|\theta^\star_{pbs} - \theta\|^2 \leqslant (1 + \tau_\varepsilon) R(\theta, \lambda) + c\varepsilon^2 v_\varepsilon^{2\beta+1}.$$

*Proof.* A proof may be found in [32].

[44] consider their own block estimator, and show its sharp minimax adaptivity on the classes of ellipsoids.

   A very long discussion, concerning, size of blocks, different penalties, several classes of functions where the estimator is minimax adaptive, may be found in [31]. The family of weakly geometrically increasing blocks is not in fact, the more precise choice in order to get very sharp results.

   Other interesting results about the penalized Stein rule may be found in [14, 17] in the wavelet case and with heavy penalties $p_j$ that do not tend to 0 as $T_j \to \infty$. In particular, [14] propose to take $p_j = 1/2 - 3/T_j$ and $T_j = 2^j$, while [17] considers small blocks with constant length $T_j \sim \log(1/\varepsilon)$ and $p_j > 4$. These penalties are too large to get exact oracle inequalities or sharp minimax adaptation, but they are sufficient for oracle inequality and then minimax adaptivity.

*Remark 1.40.* Since $\tau_\varepsilon = o(1)$, the oracle inequality of Theorem 1.11 may lead to some asymptotic exact oracle inequality. One needs to prove then that $\varepsilon^2 v_\varepsilon^{2\beta+1}$ is small.

In this part, we apply Theorem 1.11 to show that the penalized blockwise Stein estimator with the given special construction of blocks $I_j$ is sharp minimax adaptive on the classes of ellipsoids.

**Theorem 1.12.** *Let $\Theta = \Theta(a, L)$ be an ellipsoid defined in (1.40) with monotone non-decreasing $a = \{a_k\}$, $a_k \to \infty$ and $L > 0$. Let the blocks $I_j$ satisfy (B1), the penalties $p_j$ satisfy (B2), and the singular values $b_k$ satisfy (B3). Assume also that $v_\varepsilon$ is chosen so that*

$$\frac{\varepsilon^2 v_\varepsilon^{2\beta+1}}{r_\varepsilon(\Theta)} = o(1), \ \varepsilon \to 0. \tag{1.100}$$

*Then the penalized blockwise Stein estimator $\theta_{pbs}^\star = \{\theta_k^\star\}$ defined in (1.97) is asymptotically minimax on $\Theta$ among all estimators, i.e.*

$$\sup_{\theta \in \Theta} \mathbf{E}_\theta \|\theta_{pbs}^\star - \theta\|^2 = r_\varepsilon(\Theta)(1 + o(1)), \tag{1.101}$$

*as $\varepsilon \to 0$.*

*Proof.* This is a simple consequence of Theorem 1.5 and Theorem 1.11. Note that under the assumptions of Theorem 1.12, the minimax sequence of Pinsker weights $\lambda$ defined in (1.42) belongs to $\Lambda_{mon}$. Next, since $a_k$ is monotone non-decreasing, $a_k \to \infty$, and $b_k$ satisfies (B3), we have $r_\varepsilon(\Theta) \to 0$, as $\varepsilon \to 0$, by Theorem 1.5. Hence,

$$\sup_{\theta \in \Theta} R(\theta, \lambda) = r_\varepsilon^\ell(\Theta) = r_\varepsilon(\Theta)(1 + o(1)) = o(1),$$

as $\varepsilon \to 0$ where we used (1.45). Thus, the assumptions of Theorem 1.11 (ii) are satisfied for $\lambda = \lambda^p$ the Pinsker weights, $\theta \in \Theta$ and $r = 1$ if $\varepsilon$ is small enough, and we may write

$$\sup_{\theta \in \Theta} \mathbf{E}_\theta \|\theta_{pbs}^\star - \theta\|^2 \leqslant (1 + o(1)) \sup_{\theta \in \Theta} R(\theta, \lambda^p) + c\varepsilon^2 v_\varepsilon^{2\beta+1}. \tag{1.102}$$

This, together with (1.100), yields

$$\sup_{\theta \in \Theta} \mathbf{E}_\theta \|\theta_{pbs}^\star - \theta\|^2 \leqslant r_\varepsilon^\ell(\Theta)(1 + o(1)),$$

which is equivalent to (1.101), in view of (1.45) and of the definition of $r_\varepsilon(\Theta)$.

Remark that Theorem 1.12 states the sharp adaptivity property of $\theta_{pbs}^\star$: this estimator is sharp asymptotically minimax on every ellipsoid $\Theta = \Theta(a, L)$ satisfying (1.100), while no prior knowledge about $a$ and $L$ is required to define $\theta_{pbs}^\star$.

Note also that the condition (1.100) is quite weak. It suffices to choose $v_\varepsilon$ smaller than some iterated logarithm of $1/\varepsilon$, in order to satisfy these conditions for most of usual examples of ellipsoids $\Theta$.

**Corollary 1.2.** *Let $\Theta = \Theta(a, L)$ be any ellipsoid with monotone non-decreasing $a = \{a_k\}$ such that $k^{\alpha_1} \leqslant a_k \leqslant \exp(\alpha_2 k)$, $\forall k$, for some $\alpha_1 > 0, \alpha_2 > 0, L > 0$. Assume (B1), (B2) and (B3) with $v_\varepsilon = \max(\lfloor \log\log 1/\varepsilon \rfloor, 5)$. Then the estimator $\theta_{pbs}^\star$ defined in (1.97) satisfies (1.101).*

*Remark 1.41.* The penalized blockwise Stein estimator is thus minimax adaptive on a very large scale of ellipsoids.

### 1.3.4.5  Model Selection Versus Universal Optimality

**Comments**

The approach of universal optimality, and the penalized blockwise Stein estimator, presented in Section 1.3.4 has very general and sharp properties.

*Universal optimality.* Theorem 1.11 shows that penalized blockwise Stein's estimator defined in (1.97) satisfies an oracle inequality on the class of all monotone sequences $\Lambda_{mon}$. In other words, it mimics the monotone oracle in (asymptotically) exact way. This immediately entails oracle inequalities on all the subclasses $\Lambda' \subset \Lambda_{mon}$. In particular, the estimator $\theta^{\star}_{pbs}$ is asymptotically at least as good as the optimal projection estimator, the optimal Tikhonov estimator or the optimal Pinsker estimator (see Section 1.2.2).

In a sense, this is a stronger property than oracle inequalities for the "model selection" estimators in Section 1.3.3 or [111, 82, 8, 25]. In those papers it was possible to treat in each occasion only one class $\Lambda'$.

This point is really crucial for the model selection approach. Among a family of estimators, one select the best possible one, by a data-driven selection method $\lambda^{\star}$ which takes its values in $\Lambda$.

The penalized blockwise Stein estimator at least mimics (and in fact outperforms) simultaneously the oracles on all these classes $\Lambda'$. This behaviour may be called *universal optimality*. One has then a universal estimator which is as good as most of the standard families of linear estimators.

*Universal adaptivity.* Another point is that, no "ellipsoidal" structure appears in the definition of $\theta^{\star}_{pbs}$. In fact, minimax results similar to Theorem 1.12 can be formulated for other classes than ellipsoids (for example, for parallelepipeds), provided the minimax solution $\lambda$ is a monotone non-increasing sequence, see [31]. The penalized blockwise Stein estimator is thus minimax adaptive on a very large scale of classes of functions.

In a way, it is *universally adaptive*.

*Non-linearity property.* A last remark, is that the penalized blockwise Stein estimator is in fact, a non-linear estimator. Moreover, $\theta^{\star}_{pbs}$ does not even belong to the class $\Lambda_{mon}$.

It is well-known that on some classes of Besov classes with rather unsmooth functions, one needs non-linear estimators, for example wavelet thresholding, since linear ones are suboptimal, see [41]. [31] showed that $\theta^{\star}_{pbs}$ is (almost) optimal on these classes of unsmooth functions. The penalized blockwise Stein estimator, due to the shrinkage and the blocks, has, in some sense, the behaviour of a non-linear estimator.

Nevertheless, the penalized blockwise Stein estimator has some drawbacks.

*Instability in inverse problems.* The first one is almost the same than the URE estimator of Section 1.3.3.1. Due to the increase in the penalty the penalized version of blockwise Stein's estimator is less unstable than the URE estimator. However, one really needs a condition as a fixed $N_{max}$ defined in (1.99) in order

to avoid too large choices of blocks. Without this condition, the method is rather unstable in simulations.

*Universal method: a constraint?* A second drawback, of this universal approach is in fact its nature. Indeed, quite often in applications, scientists want to use their favourite method (Tikhonov, projection, $\nu$−method,...). They know, or believe, that this method works well in their field. In a way, the model selection approach answers to their problem. It allows to calibrate in a data-driven way (by choosing $\gamma$, $N$ or $m$) their favourite method.

On the other hand, the universal approach, by its universal definition, does not really answer to their question. The universal optimality just proves that one very specific universal method, the penalized blockwise Stein estimator, is as good as their favourite method. However, this could be disappointing since they cannot use their own method.

*Penalization in inverse problems.* In a way, the penalized blockwise Stein estimator already contained the idea that in inverse problems, penalizing slightly more than the URE penalty was needed. Such a choice improves the accuracy of the method. The paper [31] was in fact written after [32]. In the inverse problems framework presented in [32] already appeared the need of stronger penalties than URE. This point is true in theory, but also in simulations where one has to be very careful with too large choices of number of coefficients $N$.

Nevertheless, after some times, the idea of penalizing slightly more than URE was found to be successful even in the context of the direct problem, i.e. Gaussian white noise. Thus, non-penalized blockwise Stein's rule leads to an oracle inequality which is similar, but less accurate than that of Theorem 1.11, see [31]. The study in [31] was also, in a way, deeper than in the inverse problems context. Indeed, there is a rather long discussion concerning, the different penalties, block sizes, and functional classes that may be studied.

The main idea was that one needed to penalize more than the URE penalty, especially in inverse problems. However, in order to get sharp theoretical results, but also a method accurate in simulations, this penalty did not need to be too large. Thus, this was, in a way, the first step from unbiased risk estimation to risk hull method.

## 1.4 Conclusion

### 1.4.1 Summary

A very promising approach to inverse problems is the statistical framework. It is based on a model where observations contain a random noise. This does not correspond to the historical framework of [127] where the error is deterministic.

The optimal rates of convergence are different in the statistical and deterministic frameworks (see Section 1.2.5).

We have studied, in Section 1.1, the white noise model discretized in the spectral domain by use of the SVD, when the operator $A$ is compact. This allows to define a measure of ill-posedness of an inverse problem, with influence on the rates of convergence.

Several examples of inverse problems where the SVD is known were presented (circular deconvolution, heat equation, tomography,...).

The spectral theory for non-compact operators was also developped with the example of deconvolution on $\mathbb{R}$.

In Section 1.2, the nonparametric approach and minimax point of view were presented. This notion corresponds to the asymptotic minimax optimality as the noise level goes to zero.

Several examples of standard regularization methods, and their counterpart as estimation procedures by use of SVD, were discussed (projection, Landweber, Tikhonov,...).

The notion of source condition was introduced, with its link with ellipsoid in $\ell^2$ and standard classes of functions (Sobolev and analytic functions). The optimal rates of convergence were given. These rates depend on the smoothness of the function to reconstruct and on the degree of ill-posedness of the inverse problem.

In ill-posed inverse problems the rates are slower than in the direct problem, corresponding to the standard nonparametric statistics framework.

This notion of optimality leads to some optimal choice of the tuning parameter ($N$, $\gamma$, or $m$).

However these optimal parameters are unachievable since they depend on the unknown smoothness of the function.

This remark leads to the main point of Section 1.3. The goal is to find data-driven choices of the tuning parameter (adaptive methods). In applications, this choice is just done by simulations in a very empirical way. For example, one uses known phantom images in order to calibrate, the estimator. Usually, there is no theoretical results in order to validate this approach. Moreover, this could be very unstable when the observed functions are different from the phantom.

The minimax adaptive approach is concerned with the construction of estimators which attain the optimal rates of convergence for any smoothness $\alpha$ of the function $f$.

The oracle approach is a second step in the problem of data-driven selection method. The oracle is the best possible choice, in a given family of estimators, provided we knew the unknown function. However, such a procedure cannot be constructed, since it is not an estimator. The aim of oracle inequalities is to prove that the estimator accuracy is close from the oracle behaviour.

There exist many different methods in order to construct data-driven choice of the tuning parameter. One of the more natural is the idea of minimizing an estimate of the risk (URE). The theoretical results concerning this method are satisfying.

Nevertheless, in simulations, the URE method is usually not stable enough in inverse problems. The approach of penalized empirical risk, may be better than URE provided the penalty function is chosen appropriately. The risk hull method (RHM) provides one way to find a good and explicit penalty function.

Another, adaptive method is considered, based on the blockwise Stein estimator. Again, with a slightly stronger penalty, this method is rather satisfying.


## 1.4.2 Discussion

The statistical approach to inverse problems is nowadays quite popular and successful.

There exist some differences between the two frameworks, stochastic and deterministic. For example, the optimal rates of convergence are not the same (see Section 1.2.5). Nevertheless, this difference in the optimal rates is not so important. In a way, the two frameworks are rather related.

However, one of the major advantages of the statistical approach is that it allows to obtain oracle inequalities and to construct adaptive estimators.

The oracle approach is thus very interesting in inverse problems. Indeed, one can construct procedures in order to choose the best estimator among a given family of regularization methods. This really gives some answer to a very natural problem, the data-driven choice of the tuning parameter ($N$ or $\gamma$). From a practical point of view, this choice is usually just done by simulations in a very empirical way. Usually, by calibrating the method on some known phantom image. This approach may give a rather unstable procedure.

From a mathematical point of view, the oracle approach is very interesting. Indeed, the statistical theory is here able to give some answer to the very sensitive problem of data-driven choice of the tuning parameter.

Another important remark is that inverse problems are difficult problems. Indeed, we have to invert an operator in order to get the reconstruction. A main issue is then to get very precise oracle inequalities, i.e. with a good control on the constants of the main term, but also of the remainder term. The degree of ill-posedness of the problem appears in the rates of convergence, but at some point, in the oracle inequalities as well, which are thus sensitive to the difficulty of the problem.

Thus, in statistical inverse problems one has to define very precise model selection methods, or choice of the regularization parameter, otherwise, due to the difficulty of the problem, the estimator will not be accurate.

This remark is rather satisfying concerning the interest of the inverse problem framework in statistics. Indeed, due to the natural difficulty of the ill-posed problems, the statistical study is thus very challenging. In some sense, many estimators, or adaptive procedures, may be satisfying in the direct problem. Nevertheless, in the ill-posed context, one has to be much more careful, and the statistical study could really be more difficult.

## *1.4.3 Open Problems*

In these lectures, the results have been obtained for a very specific and restrictive model. The model is a white noise model, with an additive and Gaussian noise. Moreover, a strong assumption is related to the use of SVD.

There exist many different approaches in order to extend the results or to deal with other kind of problems.

The goal of the present section is to discuss problems which are not presented in these notes. Several of these topics have been already well-studied in the literature, others remain more open.

### Noisy Operators

One very restrictive assumption is that the operator $A$ is perfectly known. Indeed, in many applications, the operator is not known, or at least not completely known. For example, in astronomical observations, point spread function may be changing due to unknown physical conditions. This problem is also related to the well-known problem of blind deconvolution, where one has to estimate also the convolution kernel.

From a theoretical point of view this problem is also quite challenging. Indeed, the operator, by its spectral behaviour, characterizes the optimal rates of convergence. Thus, it is not clear, if any modification on the operator would change the rates or not.

The case of fully unknown operator $A$ is usually difficult. Indeed, one would need to estimate both the operator $A$ and the function $f$ by using the same data.

A more natural framework is the case of noisy operator, where the operator is not completely known and estimated using other data. This very important topic has been the subject of several recent statistical works, see for example in [48, 97, 66, 71]. In this framework, there exist two noises, one on the operator $A$ and one on the inverse problem data. The main conclusion here is that, usually, the rate is the worse possible between these two noises.

A more specific model may also be considered, where the SVD basis is known, but the singular values are noisy. This setting appears for example in circular convolution model where the SVD is always the Fourier basis, but where the convolution kernel has to be estimated. In this situation, sharp oracle inequalities may be obtained, see [27, 29, 98].

### Nondiagonal Case

One of the main drawbacks is that all these methods are linked to the spectral approach. We have intensively used the SVD to diagonalize the operators. The different regularization methods were presented for the spectral domain, even if, many of them can be computed without the explicit use of their SVD.

However, there is a more general situation where the operator $A$ cannot be represented by a diagonal matrix. For example, one uses a basis, but which does not diagonalize the operator.

In this case, several results have been obtained, such that, optimal rates of convergence, adaptive estimation, oracle inequalities, see for example in [102, 97, 91].

## Wavelets and Sparsity

As noted above, most of the methods are linked to the spectral approach. In many problems, this leads to the Fourier domain. Thus, due for example to the source conditions, the function to be reconstructed should have good properties in the Fourier domain.

Another very popular approach is based on wavelets, see for example in [42]. By using wavelets, one may usually deal with functions which are not very smooth, by replacing Sobolev classes by Besov classes. Indeed, there exist Besov classes which contain functions which are really unsmooth. Moreover, wavelets bases have the nice property that rather few coefficients are large, i.e. they give sparse representations. Thus, the standard estimator is constructed by using a threshold estimator of wavelets coefficients. This method allows to obtain adaptive estimators.

In inverse problems, wavelets have usually very good properties related to the operator $A$. Wavelets bases are not the exact SVD of a given operator. However, wavelets bases almost diagonalize many operators. Moreover by using thresholding they have good adaptability properties, see the Wavelet Vaguelette Decomposition (WVD) approach in [39]. This framework is thus strongly related to the previous nondiagonal case.

There exist a very large literature in inverse problems with wavelets, see for example in [39, 83, 78, 19, 28, 34, 79] and, with the framework of noisy operator, [71, 29].

## RHM for Other Methods

The RHM is presented here for the family of projection estimators. There exist many other regularization methods (Landweber, Tikhonov, $\nu-$methods). These methods usually attain the optimal rates of convergence, see for example [9]. The RHM approach has been very recently extended to these families of estimators (see [99]).

The RHM is also valid in the framework of noisy singular values, see [98].

## Nonlinear Operators

All the results given here are valid for linear inverse problems. In the case of nonlinear operator, the problem is much more difficult. This framework has been intensively studied in the deterministic context, see [49].

However, this problem is not yet well understood in statistics. Due to the stochastic nature of the noise the nonlinear operator is more difficult to handle. Moreover, adaptive estimation and oracle inequalities are even more involved in this framework.

Some recent papers concerning the statistical study of nonlinear inverse problems may be found in [6, 92].

### Error in Variables

There exist a rather popular topic in statistics which is very closely related to our framework, the error in variables problem. In this context, one observes

$$Y_i = X_i + \xi_i, \;\; i = 1, \ldots, n,$$

where $\{X_i\}$ and $\{\xi_i\}$ are i.i.d. random variables, independent, and usually defined on $\mathbb{R}$. The goal is to estimate the probability density $f$ of the random variable $X$. Since $X$ and $\xi$ are independent, the probability density of $Y$ is well-known to be a convolution of the two densities of $X$ and $\xi$.

The exist two main differences here. The first one is that the model is a density type model, and not any more a white noise model. The second point is that the operator is usually not compact. Indeed, the convolution is on the whole $\mathbb{R}$ due to the random variables which take their values on $\mathbb{R}$. However, it is well-known that a white noise model may be considered as an idealized version of a density model, there even exists a formal equivalence, see [107]. Thus, usually the rates of convergence are the same in these two models, even if the mathematical proofs could be quite different.

Formally, this model could be considered as a special case of the model (1.2). However, the noise $\xi$ should have a very specific behaviour, which is not true in the standard white noise case, see [9].

This model of error in variables is then really an inverse problem, and is often called the deconvolution problem in the statistical literature, see for example [52, 35, 16, 30].

### Econometrics

Nowadays, the topic of inverse problems has also a growing interest in econometrics. The problem of intrumental variables is closely related to the framework of inverse problems.

An economic relationship between a response variable $Y$ and a vector $X$ of explanatory variables is often represented by the following equation

$$Y_i = f(X_i) + U_i, \; i = 1, \ldots, n,$$

where the function $f$ has to be estimated and $U_i$ are the errors. This model does not characterize the function $f$ if $U$ is not constrained. The problem is solved if $\mathbf{E}(U|X) = 0$.

However, in many structural econometrics models, the parameter of interest is a relation between $Y$ and $X$, where some components of $X$ are endogeneous. This situation arises frequently in economics. For example, suppose that $Y$ denotes the wages and the $X$ includes, level of education, among other variables. The error $U$ includes, ability, which is not observed, but influences the wages. If a high ability tends to choose high level of education, then education and ability are correlated, and thus $X$ and $U$ also.

Nevertheless, suppose that we observe another set of data, $W_i$ where $W$ is called an instrumental variable for which

$$\mathbf{E}(U|W) = \mathbf{E}(Y - f(X)|W) = 0.$$

This equation characterizes $f$ by a Fredholm equation of the first kind. Estimation of the function $f$ is in fact an ill-posed inverse problems.

Since the years 2000, the framework of inverse problems has been the topic of many articles in the econometrics literature, see [54, 66, 33, 53], see also Chapter 2 by Jean-Pierre Florens in the present book.


## Inverse Problems in Applications

These lectures mainly consider inverse problems from a theoretical point of view. This is satisfying from a mathematical perspective. Indeed, one can define this kind of problems with a rather general description. The difficulty of an inverse problem, i.e. its degree of ill-posedness $\beta$, is characterized by the spectral behaviour of the operator $A$ when $k \to \infty$. Moreover, the smoothness $\alpha$ of the function $f$ is also important. These parameters give the optimal rates of convergence.

This general description is rather important. Indeed, it allows to understand the difference between inverse problems, and the influence of the smoothness on the accuracy of the reconstruction.

However, all these concepts are mainly just mathematical tools. They are based on asymptotics, when $k$ is large.

In a more applied point of view, there is mainly no difference between an inverse problem of degree $\beta = 2$ and a severely ill-posed problem. Moreover, many problems which are almost unsolvable are, in applications, rather easy to deal with. For example, a deconvolution by a Gaussian kernel ($\mathcal{N}(0, \sigma^2)$) is even worse that severely ill-posed. However, if this convolution kernel, has a small variance $\sigma^2$, then the problem is very easy to solve.

On the other hand, many of the problems which appear in applications are much more difficult than our framework with an idealized model. Even in the simple circular deconvolution, but with boundary effect, the SVD basis is not the Fourier basis anymore. The number of data in applications is finite $n$ and does not go to infinity.

There is a lot of frameworks where identifiability of the model, i.e. existence and unicity, is the main problem, before any stability results. In more realistic models, most of the operators are not completely known and not even observed with some additive noise.

Tomography, even based on the same operator, the Radon transform, is a world by itself. There exist conferences and articles, on computerized tomography, positron emission tomography, discrete tomography, quantum homodyne tomography, see [104].

It is rather common to say, that each inverse problem is in fact a specific case, see [117].

## Numerical Aspects

The numerical aspect of the different regularization methods was not so much discussed. However, this point is of importance, especially in inverse problems. As noted before, many of the regularization methods are expressed in the spectral domain (SVD) but many of them are in fact computed in a different way, without using the whole spectrum. For example, the Tikhonov regularization is computed by minimizing the functional (1.28). In deconvolution problems, the SVD will be the Fourier basis, which may really be computed quite fast by use of Fast Fourier Transform (FFT). In more difficult problems, for example Radon transform in tomography, the SVD could be much slower to compute, see [104, 49, 81, 130].

Moreover, iterative methods are rather popular because they also avoid the inversion of a large matrix as in (1.30). This is one the reason of the interest in all these iterative procedures, see [10].

## Well-Posed Questions

One of the drawbacks of the study of inverse problems is its intrinsic difficulty. Indeed, the optimal rates of convergence may be rather slow (see Section 1.2.4). Even in the case of mildly ill-posed problems when the degree of ill-posedness $\beta$ is large (even $\beta = 2$), the optimal rates will be quite slow. In the severely ill-posed context it is even worse and the rates could be logarithmic. Moreover, these rates are optimal, they cannot be improved on a given class of functions.

In a way some inverse problems are really too difficult (for example the heat equation). One may think that in a given model there is no hope to get better results.

A rather natural idea when a model is too difficult, is to change the goal of the problem. One tries to answer to problems that could be solved in a more satisfying way. The main point is thus to solve more easy problems than estimating the whole function $f$. For example, estimating linear functionals, level sets or change points, or solving testing or classification problems. It is well-known, that all these problems are more easy, i.e. have a better rate of convergence, than estimating the whole function $f$, see [56].

This point of view makes sense in problems where estimation of the whole function $f$ seems almost beyond the scope. Thus, the idea is to find more simple tasks to deal with. In fact, one is looking to well-posed questions in ill-posed problems, see [117]. These are questions that may be answered in a satisfying way.

# References

1. Adorf, H.M.: Hubble space telescope image restoration in its fourth year. Inverse Problems **11**, 639–653 (1995)
2. Akaike, H.: Information theory and an extension of the maximum likelihood principle. In: B. Petrov, F. Czáki (eds.) Proceedings of the Second International Symposium on Information Theory, pp. 267–281. Akademiai Kiadó, Budapest (1973)
3. Akaike, H.: A new look at the statistical model identification. IEEE Trans. Automat. Control **19**, 716–723 (1974)
4. Barron, A., Birgé, L., Massart, P.: Risk bounds for model selection via penalization. Probab. Theory Related Fields **113**, 301–413 (1999)
5. Bauer, F., Hohage, T.: A Lepski-type stopping rule for regularized Newton methods. Inverse Problems **21**, 1975–1991 (2005)
6. Bauer, F., Hohage, T., Munk, A.: Iteratively regularized Gauss-Newton method for nonlinear inverse problems with random noise. SIAM J. Numer. Anal. **47**, 1827–1846 (2009)
7. Belitser, E., Levit, B.: On minimax filtering on ellipsoids. Math. Methods Statist. **4**, 259–273 (1995)
8. Birgé, L., Massart, P.: Gaussian model selection. J. Eur. Math. Soc. **3**, 203–268 (2001)
9. Bissantz, N., Hohage, T., Munk, A., Ruymgaart, F.: Convergence rates of general regularizations methods for statistical inverse problems and application. SIAM J. Numer. Anal. **45**, 2610–2636 (2007)
10. Brakage, H.: On ill-posed problems and the method of conjugate gradients. In: Inverse and ill-posed problems. Academic Press, Orlando (1987)
11. Bretagnolle, J., Huber, C.: Estimation des densités : risque minimax. Z. Wahrsch. Verw. Gebiete **47**, 199–237 (1976)
12. Brezis, H.: Analyse fonctionnelle, Théorie et applications. Dunod, Paris (1999)
13. Brown, L., Low, M.: Asymptotic equivalence of nonparametric regression and white noise. Ann. Statist. **24**, 2384–2398 (1996)
14. Brown, L., Low, M., Zhao, L.: Superefficiency in nonparametric function estimation. Ann. Statist. **25**, 898–924 (1997)
15. Bühlmann, P., Yu, B.: Boosting with $\ell^2$-loss: regression and classification. J. Amer. Statist. Assoc. **98**, 324–339 (2003)
16. Butucea, C., Tsybakov, A.: Sharp optimality in density deconvolution with dominating bias. Theory Probab. Appl. **52**, 24–39 (2008)
17. Cai, T.: Adaptive wavelet estimation: a block thresholding and oracle inequality approach. Ann. Statist. **27**, 2607–2625 (1999)
18. Candès, E.: Modern statistical estimation via oracle inequalities. Acta Numer. **15**, 257–325 (2006)
19. Candès, E., Donoho, D.: Recovering edges in ill-posed inverse problems: Optimality of curvelet frames. Ann. Statist. **30**, 784–842 (2002)
20. Cavalier, L.: Efficient estimation of a density in a problem of tomography. Ann. Statist. **28**, 330–347 (2000)
21. Cavalier, L.: On the problem of local adaptive estimation in tomography. Bernoulli **7**, 63–78 (2001)
22. Cavalier, L.: Inverse problems with non-compact operator. J. of Statist. Plann. Inference **136**, 390–400 (2006)

23. Cavalier, L.: Nonparametric statistical inverse problems. Inverse Problems **24**, 1–19 (2008)
24. Cavalier, L., Golubev, G., Lepski, O., Tsybakov, A.: Block thresholding and sharp adaptive estimation in severely ill-posed inverse problems. Theory Probab. Appl. **48**, 426–446 (2003)
25. Cavalier, L., Golubev, G., Picard, D., Tsybakov, A.: Oracle inequalities in inverse problems. Ann. Statist. **30**, 843–874 (2002)
26. Cavalier, L., Golubev, Y.: Risk hull method and regularization by projections of ill-posed inverse problems. Ann. Statist. **34**, 1653–1677 (2006)
27. Cavalier, L., Hengartner, N.: Adaptive estimation for inverse problems with noisy operators. Inverse Problems **21**, 1345–1361 (2005)
28. Cavalier, L., Koo, J.Y.: Poisson intensity estimation for tomographic data using a wavelet shrinkage approach. IEEE Trans. Inform. Theory **48**, 2794–2802 (2002)
29. Cavalier, L., Raimondo, M.: Wavelet deconvolution with noisy eigenvalues. IEEE Trans. Signal Process. **55**, 2414–2424 (2007)
30. Cavalier, L., Raimondo, M.: Multiscale density estimation with errors in variables. J. Korean Statist. Soc. (2010)
31. Cavalier, L., Tsybakov, A.: Penalized blockwise Stein's method, monotone oracles and sharp adaptive estimation. Math. Methods Statist. **10**, 247–282 (2001)
32. Cavalier, L., Tsybakov, A.: Sharp adaptation for inverse problems with random noise. Probab. Theory Related Fields **123**, 323–354 (2002)
33. Chen, X., Reiss, M.: On rate optimality for ill-posed inverse problems in econometrics (2010). In press
34. Cohen, A., Hoffmann, M., Reiss, M.: Adaptive wavelet Galerkin method for linear inverse problems. SIAM J. Numer. Anal. **42**, 1479–1501 (2004)
35. Comte, F., Rozenholc, Y., Taupin, M.L.: Penalized contrast estimator for adaptive density deconvolution. Canad. J. Statist. **34**, 431–452 (2006)
36. Craven, P., Wahba, G.: Smoothing noisy data with spline functions: estimating the correct degree of smoothing by the method of generalized cross-validation. Numer. Math. **31**, 377–403 (1979)
37. Deans, S.: The Radon Transform and some of its Applications. Wiley, New York (1983)
38. Donoho, D.: Statistical estimation and optimal recovery. Ann. Statist. **22**, 238–270 (1994)
39. Donoho, D.: Nonlinear solution of linear inverse problems by wavelet-vaguelette decomposition. Appl. Comput. Harmon. Anal. **2**, 101–126 (1995)
40. Donoho, D., Johnstone, I.: Ideal spatial adaptation via wavelet shrinkage. Biometrika **81**, 425–445 (1994)
41. Donoho, D., Johnstone, I.: Adapting to unknown smoothness via wavelet shrinkage. J. Amer. Statist. Assoc. **90**, 1200–1224 (1995)
42. Donoho, D., Johnstone, I.: Minimax estimation via wavelet shrinkage. Ann. Statist. **26**, 879–921 (1998)
43. Donoho, D., Low, M.: Renormalization exponents and optimal pointwise rates of convergence. Ann. Statist. **20**, 944–970 (1992)
44. Efroimovich, S., Pinsker, M.: Learning algorithm for nonparametric filtering. Autom. Remote Control **11**, 1434–1440 (1984)
45. Efromovich, S.: Robust and efficient recovery of a signal passed through a filter and then contaminated by non-Gaussian noise. IEEE Trans. Inform. Theory **43**, 1184–1191 (1997)
46. Efromovich, S.: Nonparametric Curve Estimation. Springer, New York (1998)
47. Efromovich, S.: Simultaneous sharp adaptive estimation of functions and their derivatives. Ann. Statist. **26**, 273–278 (1998)
48. Efromovich, S., Koltchinskii, V.: On inverse problems with unknown operators. IEEE Trans. Inform. Theory **47**, 2876–2893 (2001)
49. Engl, H., Hanke, M., Neubauer, A.: Regularization of Inverse Problems. Kluwer Academic Publishers (1996)
50. Ermakov, M.: Minimax estimation of the solution of an ill-posed convolution type problem. Problems Inform. Transmission **25**, 191–200 (1989)
51. Evans, S., Stark, P.: Inverse problems as statistics. Inverse Problems **18**, 55–97 (2002)

52. Fan, J.: On the optimal rates of convergence for nonparametric deconvolution problems. Ann. Statist. **19**, 1257–1272 (1991)
53. Florens, J., Johannes, J., Van Bellegem, S.: Identification and estimation by penalization in nonparametric instrumental regression. Econ. Theory (2010). In press
54. Florens, J.P.: Inverse problems and structural econometrics: the example of instrumental variables. In: Advances in Economics and Econometrics: Theory and Applications, vol. 2, pp. 284–311 (2003)
55. Goldenshluger, A.: On pointwise adaptive nonparametric deconvolution. Bernoulli **5**, 907–925 (1999)
56. Goldenshluger, A., Pereverzev, S.: Adaptive estimation of linear functionals in Hilbert scales from indirect white noise observations. Probab. Theory Related Fields **118**, 169–186 (2000)
57. Goldenshluger, A., Spokoiny, V.: On the shape-from-moments problem and recovering edges from noisy Radon data. Probab. Theory Related Fields **128**, 123–140 (2004)
58. Goldenshluger, A., Tsybakov, A.: Adaptive prediction and estimation in linear regression with infinitely many parameters. Ann. Statist. **29**, 1601–1619 (2001)
59. Golubev, G.: Quasi-linear estimates of signals in $L^2$. Problems Inform. Transmission **26**, 15–20 (1990)
60. Golubev, G.: The principle of penalized empirical risk in severely ill-posed problems. Probab. Theory Related Fields **130**, 18–38 (2004)
61. Golubev, G., Khasminskii, R.: A statistical approach to some inverse problems for partial differential equations. Problems Inform. Transmission **35**, 51–66 (1999)
62. Golubev, G., Khasminskii, R.: A statistical approach to the Cauchy problem for the Laplace equation. Lecture Notes Monograph Series **36**, 419–433 (2001)
63. Grama, I., Nussbaum, M.: Asymptotic equivalence for nonparametric regression. Math. Methods Statist. **11**, 1–36 (2002)
64. Groetsch, C.: Generalized Inverses of Linear Operators: Representation and Approximation. Dekker, New York (1977)
65. Hadamard, J.: Le problème de Cauchy et les équations aux dérivées partielles hyperboliques. Hermann, Paris (1932)
66. Hall, P., Horowitz, J.: Nonparametric methods for inference in the presence of instrumental variables. Ann. Statist. **33**, 2904–2929 (2005)
67. Hall, P., Kerkyacharian, G., Picard, D.: Block threshold rules for curve estimation using kernel and wavelet methods. Ann. Statist. **26**, 922–942 (1998)
68. Halmos, P.: What does the spectral theorem say? Amer. Math. Monthly **70**, 241–247 (1963)
69. Hida, T.: Brownian Motion. Springer-Verlarg, New York-Berlin (1980)
70. Hoerl, A.: Application of ridge analysis to regression problems. Chem. Eng. Progress **58**, 54–59 (1962)
71. Hoffmann, M., Reiss, M.: Nonlinear estimation for linear inverse problems with error in the operator. Ann. Statist. **36**, 310–336 (2008)
72. Hohage, T.: Lecture notes on inverse problems (2002). Lectures given at the University of Göttingen
73. Hutson, V., Pym, J.: Applications of Functional Analysis and Operator Theory. Academic Press, London (1980)
74. Ibragimov, I., Khasminskii, R.: Statistical Estimation: Asymptotic Theory. Springer, New York (1981)
75. Ibragimov, I., Khasminskii, R.: On nonparametric estimation of the value of a linear functional in Gaussian white noise. Theory Probab. Appl. **29**, 19–32 (1984)
76. James, W., Stein, C.: Estimation with quadratic loss. In: Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability, pp. 361–380. University of California Press (1961)
77. Johnstone, I.: Function estimation in Gaussian noise: sequence models (1998). Draft of a monograph
78. Johnstone, I.: Wavelet shrinkage for correlated data and inverse problems: adaptivity results. Statist. Sinica **9**, 51–83 (1999)

79. Johnstone, I., Kerkyacharian, G., Picard, D., Raimondo, M.: Wavelet deconvolution in a periodic setting. J. R. Stat. Soc. Ser. B Stat. Methodol. **66**, 547–573 (2004)
80. Johnstone, I., Silverman, B.: Speed of estimation in positron emission tomography and related inverse problems. Ann. Statist. **18**, 251–280 (1990)
81. Kaipio, J., Somersalo, E.: Statistical and Computational Inverse Problems. Springer (2004)
82. Kneip, A.: Ordered linear smoothers. Ann. Statist. **22**, 835–866 (1994)
83. Kolaczyk, E.: A wavelet shrinkage approach to tomographic image reconstruction. J. Amer. Statist. Assoc. **91**, 1079–1090 (1996)
84. Koo, J.Y.: Optimal rates of convergence for nonparametric statistical inverse problems. Ann. Statist. **21**, 590–599 (1993)
85. Korostelev, A., Tsybakov, A.: Optimal rates of convergence of estimators in a probabilistic setup of tomography problem. Probl. Inf. Transm. **27**, 73–81 (1991)
86. Landweber, L.: An iteration formula for Fredholm equations of the first kind. Amer. J. Math. **73**, 615–624 (1951)
87. Lepskii, O.: One problem of adaptive estimation in Gaussian white noise. Theory Probab. Appl. **35**, 459–470 (1990)
88. Lepskii, O.: Asymptotic minimax adaptive estimation. 1. Upper bounds. Theory Probab. Appl. **36**, 654–659 (1991)
89. Lepskii, O.: Asymptotic minimax adaptive estimation. 2. Statistical model without optimal adaptation. Adaptive estimators. Theory Probab. Appl. **37**, 468–481 (1992)
90. Li, K.C.: Asymptotic optimality of $C_P, C_L$, cross-validation and generalized cross-validation: Discrete index set. Ann. Statist. **15**, 958–976 (1987)
91. Loubes, J.M., Ludena, C.: Adaptive complexity regularization for linear inverse problems. Electron. J. Stat. **2**, 661–677 (2008)
92. Loubes, J.M., Ludena, C.: Penalized estimators for nonlinear inverse problems. ESAIM Probab. Stat. (2010)
93. Loubes, J.M., Rivoirard, V.: Review of rates of convergence and regularity conditions for inverse problems. Int. J. Tomogr. and Stat. (2009)
94. Louis, A., Maass, P.: A mollifier method for linear operator equations of the first kind. Inverse Problems **6**, 427–440 (1990)
95. Mair, B., Ruymgaart, F.: Statistical estimation in Hilbert scale. SIAM J. Appl. Math. **56**, 1424–1444 (1996)
96. Mallows, C.: Some comments on $C_p$. Technometrics **15**, 661–675 (1973)
97. Marteau, C.: Regularization of inverse problems with unknown operator. Math. Methods Statist. **15**, 415–443 (2006)
98. Marteau, C.: On the stability of the risk hull method for projection estimators. J. Statist. Plann. Inference **139**, 1821–1835 (2009)
99. Marteau, C.: Risk hull method for general family of estimators. ESAIM Probab. Stat. (2010)
100. Massart, P.: Concentration Inequalities and Model Selection. Lecture Notes in Mathematics, Springer, Berlin (2007)
101. Mathé, P.: The Lepskii principle revisited. Inverse Problems **22**, 11–15 (2006)
102. Mathé, P., Pereverzev, S.: Optimal discretization of inverse problems in Hilbert scales. Regularization and self-regularization of projection methods. SIAM J. Numer. Anal. **38**, 1999–2021 (2001)
103. Mathé, P., Pereverzev, S.: Regularization of some linear ill-posed problems with discretized random noisy data. Math. Comp. **75**, 1913–1929 (2006)
104. Natterer, F.: The Mathematics of Computerized Tomography. J. Wiley, Chichester (1986)
105. Nemirovski, A.: Topics in Non-Parametric Statistics. Lecture Notes in Mathematics, Springer (2000)
106. Nemirovskii, A., Polyak, B.: Iterative methods for solving linear ill-posed problems under precise information I. Engrg. Cybernetics **22**, 1–11 (1984)
107. Nussbaum, M.: Asymptotic equivalence of density estimation and Gaussian white noise. Ann. Statist. **24**, 2399–2430 (1996)
108. O'Sullivan, F.: A statistical perspective on ill-posed problems. Statist. Sci. **1**, 502–527 (1986)

109. Pinsker, M.: Optimal filtering of square integrable signals in Gaussian white noise. Problems Inform. Transmission **16**, 120–133 (1980)
110. Plaskota, L.: Noisy Information and Computational Complexity. Cambridge University Press (1996)
111. Polyak, B., Tsybakov, A.: Asymptotic optimality of the $C_p$-test for the orthogonal series estimation of regression. Theory Probab. Appl. **35**, 293–306 (1990)
112. Raus, T., Hamarik, U., Palm, R.: Use of extrapolation in regularization methods. J. Inverse Ill-Posed Probl. **15**, 277 (2007)
113. Reiss, M.: Asymptotic equivalence for nonparametric regression with multivariate and random design. Ann. Statist. **36**, 1957–1982 (2008)
114. Rigollet, P.: Adaptive density estimation using the blockwise Stein method. Bernoulli **12**, 351–370 (2006)
115. Rosenblatt, M.: Remarks on some nonparametric estimates of a density function. Ann. Math. Statist. **27**, 832–837 (1956)
116. Ruymgaart, F.: A short introduction to inverse statistical inference (2001). Lecture given at Institut Henri Poincaré, Paris
117. Sabatier, P.: Past and future of inverse problems. J. Math. Phys. **41**, 4082–4124 (2000)
118. Schwarz, G.: Estimating the dimension of a model. Ann. Statist. **6**, 461–464 (1978)
119. Shibata, R.: An optimal selection of regression variables. Biometrika **68**, 45–54 (1981)
120. Stein, C.: Inadmissibility of the usual estimator of the mean of a multivariate distribution. In: Proceedings of the 3rd Berkeley Symposium on Mathematical Statistics and Probability, pp. 197–206. University of California Press (1956)
121. Stein, C.: Estimation of the mean of a multivariate normal distribution. Ann. Statist. **9**, 1135–1151 (1981)
122. Stone, C.: Optimal rates of convergence for nonparametric estimators. Ann. Statist. **8**, 1348–1360 (1980)
123. Sudakov, V., Khalfin, L.: Statistical approach to ill-posed problems in mathematical physics. Soviet Math. Dokl. **157**, 1094–1096 (1964)
124. Talagrand, M.: Concentration of measure and isoperimetric inequalities in product spaces. Publ. Math. IHES **81**, 73–205 (1995)
125. Taylor, M.: Partial differential equations, vol. 2. Springer, New York (1996)
126. Tenorio, L.: Statistical regularization of inverse problems. SIAM Rev. **43**, 347–366 (2001)
127. Tikhonov, A.: Regularization of incorrectly posed problems. Soviet Math. Dokl. **4**, 1624–1627 (1963)
128. Tikhonov, A., Arsenin, V.: Solution of Ill-posed Problems. Winston & Sons (1977)
129. Tsybakov, A.: Introduction to Nonparametric Estimation. Springer series in statistics (2009)
130. Vogel, C.: Computational Methods for Inverse Problems. SIAM, Philadelphia (2002)
131. Wahba, G.: Spline Models for Observational Data. SIAM, Philadelphia (1990)