

---

# Inverse Problems

---

Matthias J. Ehrhardt and Lukas F. Lang

Last updated on: March 9, 2018

Lecture Notes  
Lent Term 2017/2018

This work is licensed under a Creative Commons “Attribution-NonCommercial-ShareAlike 3.0 Unported” license.





# Contents

<b>1</b>	<b>Introduction to inverse problems</b>	<b>7</b>
1.1	Examples . . . . .	8
1.1.1	Matrix inversion . . . . .	8
1.1.2	Differentiation . . . . .	9
1.1.3	Deconvolution . . . . .	9
1.1.4	Tomography . . . . .	10
<b>2</b>	<b>Linear inverse problems</b>	<b>13</b>
2.1	Generalised solutions . . . . .	15
2.2	Generalised inverse . . . . .	18
2.3	Compact operators . . . . .	21
2.4	Singular value decomposition of compact operators . . . . .	22
<b>3</b>	<b>Regularisation</b>	<b>27</b>
3.1	Parameter-choice strategies . . . . .	30
3.1.1	A-priori parameter choice rules . . . . .	31
3.1.2	A-posteriori parameter choice rules . . . . .	32
3.1.3	Heuristic parameter choice rules . . . . .	33
3.2	Spectral regularisation methods . . . . .	34
3.2.1	Convergence rates . . . . .	36
3.2.2	Truncated singular value decomposition . . . . .	37
3.2.3	Tikhonov regularisation . . . . .	37
3.2.4	Source-conditions . . . . .	37
3.2.5	Asymptotic regularisation . . . . .	39
3.2.6	Landweber iteration . . . . .	40
3.3	Tikhonov regularisation revisited . . . . .	44
<b>4</b>	<b>Variational regularisation</b>	<b>47</b>
4.1	Variational methods . . . . .	50
4.1.1	Background . . . . .	50
4.1.2	Minimisers . . . . .	55
4.1.3	Existence . . . . .	55
4.1.4	Uniqueness . . . . .	58
4.2	Well-posedness and regularisation properties . . . . .	59
4.2.1	Existence and uniqueness . . . . .	59
4.2.2	Continuity . . . . .	64
4.2.3	Convergent regularisation . . . . .	66
4.2.4	Convergence rates . . . . .	69

<b>5</b>	<b>Numerical Solutions</b>	<b>71</b>
5.1	More on derivatives in Banach spaces . . . . .	71
5.2	Minimization problems . . . . .	73
5.2.1	Gradient Descent . . . . .	73

These lecture notes are based on the course “Inverse Problems in Imaging”, which was held by Matthias J. Ehrhardt and Martin Benning in Michemas term 2016 at the University of Cambridge.<sup>1</sup> Complementary material can be found in the following books and lecture notes:

- (a) Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. *Regularization of Inverse Problems*. Vol. 375. Springer Science & Business Media, 1996.
- (b) Martin Burger. *Inverse Problems*. Lecture notes winter 2007/2008.  
[http://www.math.uni-muenster.de/num/Vorlesungen/IP\\_WS07/skript.pdf](http://www.math.uni-muenster.de/num/Vorlesungen/IP_WS07/skript.pdf)
- (c) Andreas Kirsch. *An Introduction to the Mathematical Theory of Inverse Problems*. Vol. 120. Springer Science & Business Media, 1996.
- (d) Kazufumi Ito and Bangti Jin. *Inverse Problems: Tikhonov Theory and Algorithms*. World Scientific, 2014.
- (e) Per Christian Hansen. *Discrete Inverse Problems: Insight and Algorithms*. Fundamentals of Algorithms, SIAM Philadelphia, 2010.
- (f) Otmar Scherzer, Markus Grasmair, Harald Grossauer, Markus Haltmeier and Frank Lenzen. *Variational Methods in Imaging*. Applied Mathematical Sciences, Springer New York, 2008.
- (g) Jennifer L. Mueller and Samuli Siltanen. *Linear and Nonlinear Inverse Problems with Practical Applications*. Vol. 10. SIAM, 2012.
- (h) Andreas Rieder. *Keine Probleme mit Inversen Problemen (in German)*. Vieweg+Teubner Verlag. 2003.
- (i) Christian Clason. *Inverse Probleme (in German)*, Lecture notes winter term 2016/2017  
<https://www.uni-due.de/~adf040p/skripte/InverseSkript16.pdf>

These lecture notes are under constant redevelopment and might contain typos or errors. We very much appreciate the finding and reporting of those (to ll542@cam.ac.uk or to m.j.ehrhardt@damtp.cam.ac.uk). Thanks!

---

<sup>1</sup><http://www.damtp.cam.ac.uk/research/cia/teaching/2016inverseproblems.html>



# Chapter 1

## Introduction to inverse problems

Solving an inverse problem is the task of computing an unknown quantity from observed (and potentially noisy) measurements. Typically, these two are related via a forward model. Inverse problems appear in a variety of fields such as physics, biology, medicine, engineering, and finance, and include—for instance—tomography (e.g. computed tomography (CT)), machine learning, computer vision, and image processing. In this lecture course we address mathematical aspects of linear inverse problems that are needed to find stable and meaningful solutions.

The main focus of this lecture is the solution of the operator equation

$$Ku = f \tag{1.1}$$

with given measurement data  $f$  for the unknown quantity  $u$ . Here,  $K : \mathcal{U} \rightarrow \mathcal{V}$  denotes a linear operator that maps from a space  $\mathcal{U}$  to a space  $\mathcal{V}$ . We will restrict ourselves to the study of bounded linear operators between Hilbert spaces.

Computing a solution to (1.1) is typically not straightforward in most relevant applications for three basic reasons:

- a solution might not exist,
- if it exists it might not be unique,
- small errors (such as noise) in the measurements get heavily amplified.

The latter has the potential to render solutions useless without proper treatment.

In the sense of Hadamard the problem (1.1) is called *well-posed* if

- for all input data there exists a solution to the problem, i.e. for all  $f \in \mathcal{V}$  there exists a  $u \in \mathcal{U}$  with  $Ku = f$ .
- for all input data this solution is unique, i.e.  $u \neq v$  implies  $Kv \neq f$ .
- the solution of the problem depends continuously on the input datum, i.e. for all  $\{u_k\}_{k \in \mathbb{N}}$  with  $Ku_k \rightarrow f$  implies  $u_k \rightarrow u$ .

If any of these conditions is violated, problem (1.1) is called *ill-posed*. In the following we will see that many relevant inverse problems are ill-posed.<sup>1</sup>

---

<sup>1</sup>In fact, the name ill-posed problems may be a more suitable name for this lecture, as the real challenge is to deal with the ill-posedness of these problems. However, the name inverse problems became more widely accepted.

## 1.1 Examples

In the following we are going to present various examples of inverse problems and highlight the challenges in solving them.

### 1.1.1 Matrix inversion

One of the most simple (class of) inverse problems that arises from (numerical) linear algebra is the solution of linear systems. These can be written in the form of (1.1) with  $u \in \mathbb{R}^n$  and  $f \in \mathbb{R}^n$  being  $n$ -dimensional vectors with real entries and  $K \in \mathbb{R}^{n \times n}$  being a matrix with real entries. We further assume  $K$  to be symmetric, positive definite.

We know from the spectral theory of symmetric matrices that there exist eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$  and corresponding (orthonormal) eigenvectors  $k_j \in \mathbb{R}^n$  for  $j \in \{1, \dots, n\}$  such that  $K$  can be written as

$$K = \sum_{j=1}^n \lambda_j k_j k_j^\top. \quad (1.2)$$

It is well known from numerical linear algebra that the condition number  $\kappa = \lambda_1/\lambda_n$  is a measure of how stable (1.1) can be solved, which we will illustrate in the following.

We assume that we measure  $f^\delta$  instead of  $f$ , with  $\|f - f^\delta\|_2 \leq \delta \|K\| = \delta \lambda_1$ , where  $\|\cdot\|_2$  denotes the Euclidean norm of  $\mathbb{R}^n$  and  $\|K\|$  the operator norm of  $K$  (which equals the largest eigenvalue of  $K$ ). Then, if we further denote with  $u^\delta$  the solution of  $Ku^\delta = f^\delta$ , the difference between  $u^\delta$  and the solution  $u$  to (1.1) is

$$u - u^\delta = \sum_{j=1}^n \lambda_j^{-1} k_j k_j^\top (f - f^\delta).$$

Therefore, we can estimate

$$\|u - u^\delta\|_2^2 = \sum_{j=1}^n \lambda_j^{-2} \underbrace{\|k_j\|_2^2}_{=1} |k_j^\top (f - f^\delta)|^2 \leq \lambda_n^{-2} \|f - f^\delta\|_2^2,$$

due to the orthonormality of eigenvectors, the Cauchy-Schwarz inequality, and  $\lambda_n \leq \lambda_j$ . Thus, taking square roots on both sides yields the estimate

$$\|u - u^\delta\|_2 \leq \lambda_n^{-1} \|f - f^\delta\|_2 \leq \kappa \delta.$$

Hence, we observe that in the worst case an error  $\delta$  in the data  $y$  is amplified by the condition number  $\kappa$  of the matrix  $K$ . A matrix with large  $\kappa$  is therefore called *ill-conditioned*. We want to demonstrate the effect of this error amplification with a small example.

**Example 1.1.** Let us consider the matrix

$$K = \begin{pmatrix} 1 & 1 \\ 1 & \frac{1001}{1000} \end{pmatrix},$$

which has eigenvalues  $\lambda_j = 1 + \frac{1}{2000} \pm \sqrt{1 + \frac{1}{2000^2}}$ , condition number  $\kappa \approx 4002 \gg 1$ , and operator norm  $\|K\| \approx 2$ . For given data  $f = (1, 1)^\top$  the solution to  $Ku = f$  is  $u = (1, 0)^\top$ .

Now let us instead consider perturbed data  $f^\delta = (99/100, 101/100)^\top$ . The solution  $u^\delta$  to  $Ku^\delta = f^\delta$  is then  $u^\delta = (-19.01, 20)^\top$ .

Let us reflect on the amplification of the measurement error. By our initial assumption we find that  $\delta = \|f - f^\delta\|/\|K\| \approx \|(0.01, -0.01)^\top\|/2 = \sqrt{2}/200$ . Moreover, the norm of the error in the reconstruction is then  $\|u - u^\delta\| = \|(20.01, 20)^\top\| \approx 20\sqrt{2}$ . As a result, the amplification due to the perturbation is  $\|u - u^\delta\|/\delta \approx 4000 \approx \kappa$ .



### 1.1.2 Differentiation

Another classic inverse problem is differentiation of data. Assume we are given a function  $f$  with  $f(0) = 0$  for which we want to compute  $u = f'$ . For  $f$  sufficiently smooth, these conditions are satisfied if and only if  $u$  and  $f$  satisfy the operator equation

$$f(y) = \int_0^y u(x) dx,$$

which can be written as the operator equation  $Ku = f$  with the linear operator  $(K\cdot)(y) := \int_0^y \cdot(x) dx$ .

As before, we assume that instead of  $f$  we measure a perturbed version  $f^\delta = f + n^\delta$  with  $f \in C^1([0, 1])$  and noise  $n^\delta \in L^\infty([0, 1])$ . It is obvious that the derivative  $u$  exists if the noise  $n^\delta$  is differentiable. However, even in the (unrealistic) case that  $n^\delta$  is differentiable, the error in the derivative can become arbitrarily large as we will see.

Consider a sequence of noise functions  $n^\delta \in C^1([0, 1]) \hookrightarrow L^\infty([0, 1])$  with

$$n^\delta(x) := \delta \sin\left(\frac{kx}{\delta}\right), \quad (1.3)$$

for a fixed but arbitrary  $k > 0$ . Then, the solution to  $Ku^\delta = f^\delta$  is

$$u^\delta(x) = f'(x) + k \cos\left(\frac{kx}{\delta}\right).$$

Observe that, for  $\|n^\delta\|_{L^\infty([0,1])} = \delta \rightarrow 0$ , we on the other hand have

$$\|u - u^\delta\|_{L^\infty([0,1])} = \|(n^\delta)'\|_{L^\infty([0,1])} = k.$$

Thus, despite the error in the data becoming arbitrarily small (in the  $L^\infty$  norm), the error in the derivative can become arbitrarily big (in dependence of  $k$ ). In any case, for  $k > 0$  we observe that the solution does not depend continuously on the data.

On the other hand, considering a decreasing error in the norm of the Banach space  $C^1([0, 1])$  yields a different result. If we have a sequence of noise functions (other than those defined in equation (1.3)) with  $\|n^\delta\|_{C^1([0,1])} \leq \delta \rightarrow 0$  instead, we can conclude

$$\|u - u^\delta\|_{L^\infty([0,1])} = \|(n^\delta)'\|_{L^\infty([0,1])} \leq \|n^\delta\|_{C^1([0,1])} \rightarrow 0.$$

In contrast to the previous example the sequence of functions  $n^\delta(x) := \delta \sin(kx)$  for instance satisfies

$$\|n^\delta\|_{C^1([0,1])} = \sup_{x \in [0,1]} |n^\delta(x)| + \sup_{x \in [0,1]} |(n^\delta)'(x)| = (1 + k)\delta \rightarrow 0.$$

However, for a fixed  $\delta$  the bound on  $\|u - u^\delta\|_{L^\infty([0,1])}$  can obviously still become fairly large compared to  $\delta$ , depending on how large  $k$  is.

### 1.1.3 Deconvolution

An interesting problem that occurs in many imaging, image- and signal processing applications is the *deblurring* or *deconvolution* of signals from a known, linear degradation. Deconvolution of a signal  $f$  can be modelled as solving the inverse problem of the convolution, which reads as

$$f(y) = (Ku)(y) := \int_{\mathbb{R}^n} u(x)g(y-x) dx, \quad (1.4)$$

Here,  $f$  denotes the blurry image,  $u$  is the (unknown) true image, and  $g$  is the function that models the degradation. Due to the Fourier convolution theorem we can rewrite (1.4) to

$$f = (2\pi)^{\frac{n}{2}} \mathcal{F}^{-1}(\mathcal{F}(u)\mathcal{F}(g)). \quad (1.5)$$

with  $\mathcal{F}$  denoting the Fourier transform

$$\mathcal{F}(u)(\xi) := (2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} u(x)e^{-ix \cdot \xi} dx \quad (1.6)$$

and  $\mathcal{F}^{-1}$  being the inverse Fourier transform

$$\mathcal{F}^{-1}(f)(x) := (2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} f(\xi)e^{ix \cdot \xi} d\xi \quad (1.7)$$

It is important to note that the inverse Fourier transform is indeed the unique, inverse operator of the Fourier transform in the Hilbert space  $L^2(\mathbb{R}^n)$  due to the theorem of Plancherel. If we rearrange (1.5) to solve for  $u$  we obtain

$$u = (2\pi)^{-\frac{n}{2}} \mathcal{F}^{-1} \left( \frac{\mathcal{F}(f)}{\mathcal{F}(g)} \right), \quad (1.8)$$

and hence, we allegedly can recover  $u$  by simple division in the Fourier domain. However, we will see that this inverse problem is ill-posed and the division will lead to heavy amplifications of small measurement errors.

Let  $u$  denote the image that satisfies (1.4). Further, we assume that instead of the blurry image  $f$  we observe  $f^\delta = f + n^\delta$  instead and that  $u^\delta$  is the solution of (1.8) with input datum  $f^\delta$ . Hence, by the linearity of (1.6) and (1.7), we observe

$$(2\pi)^{\frac{n}{2}} |u - u^\delta| = \left| \mathcal{F}^{-1} \left( \frac{\mathcal{F}(f - f^\delta)}{\mathcal{F}(g)} \right) \right| = \left| \mathcal{F}^{-1} \left( \frac{\mathcal{F}(n^\delta)}{\mathcal{F}(g)} \right) \right|. \quad (1.9)$$

As the convolution kernel  $g$  usually has compact support,  $\mathcal{F}(g)$  will tend to zero for high frequencies. Hence, the denominator of (1.9) becomes fairly small, whereas the numerator will be non-zero as the noise is of high frequency. Thus, in the limit the solution will not depend continuously on the data and the convolution problem therefore be ill-posed.

#### 1.1.4 Tomography

In almost any tomography application the underlying inverse problem is either the inversion of the Radon transform<sup>2</sup> or of the X-ray transform.

For  $u \in C_0^\infty(\mathbb{R}^n)$ ,  $s \in \mathbb{R}$ , and  $\theta \in S^{n-1}$  the *Radon transform*  $R : C_0^\infty(\mathbb{R}^n) \rightarrow C^\infty(S^{n-1} \times \mathbb{R})$  can be defined as the integral operator

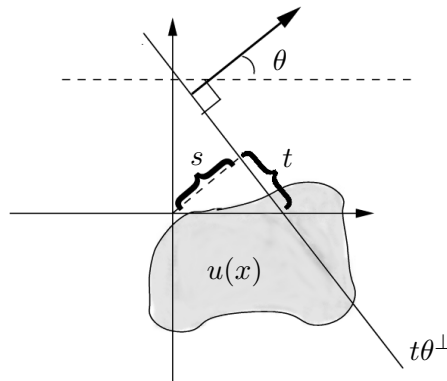
$$\begin{aligned} f(\theta, s) &= (\mathcal{R}u)(\theta, s) = \int_{x \cdot \theta = s} u(x) dx \\ &= \int_{\theta^\perp} u(s\theta + y) dy, \end{aligned} \quad (1.10)$$

which, for  $n = 2$ , coincides with the X-ray transform,

$$f(\theta, s) = (\mathcal{P}u)(\theta, s) = \int_{\mathbb{R}} u(s\theta + t\theta^\perp) dt,$$

---

<sup>2</sup>Named after the Austrian mathematician Johann Karl August Radon (16 December 1887 – 25 May 1956).



**Figure 1.1:** Visualization of the Radon transform in two dimensions (which coincides with the X-ray transform). The function  $u$  is integrated over the ray parametrized by  $\theta$  and  $s$ .<sup>3</sup>

for  $\theta \in S^{n-1}$  and  $\theta^\perp$  being the vector orthogonal to  $\theta$ . Hence, the X-ray transform (and therefore also the Radon transform in two dimensions) integrates the function  $u$  over lines in  $\mathbb{R}^n$ , see Fig. 1.1.

**Example 1.2.** Let  $n = 2$ . Then  $S^{n-1}$  is simply the unit sphere  $S^1 = \{\theta \in \mathbb{R}^2 \mid \|\theta\| = 1\}$ . We can choose for instance  $\theta = (\cos(\varphi), \sin(\varphi))^\top$ , for  $\varphi \in [0, 2\pi)$ , and parametrise the Radon transform in terms of  $\varphi$  and  $s$ , i.e.

$$f(\varphi, s) = (\mathcal{R}u)(\varphi, s) = \int_{\mathbb{R}} u(s \cos(\varphi) - t \sin(\varphi), s \sin(\varphi) + t \cos(\varphi)) dt. \quad (1.11)$$

Note that—with respect to the origin of the reference coordinate system— $\varphi$  determines the angle of the line along one wants to integrate, while  $s$  is the offset from that line from the centre of the coordinate system.

### X-ray Computed Tomography (CT)

In X-ray computed tomography (CT), the unknown quantity  $u$  represents a spatially varying density that is exposed to X-radiation from different angles, and that absorbs the radiation according to its material or biological properties.

The basic modelling assumption for the intensity decay of an X-ray beam is that within a small distance  $\Delta t$  it is proportional to the intensity itself, the density, and the distance, i.e.

$$\frac{I(x + (t + \Delta t)\theta) - I(x + t\theta)}{\Delta t} = -I(x + t\theta)u(x + t\theta),$$

for  $x \in \theta^\perp$ . By taking the limit  $\Delta t \rightarrow 0$  we end up with the ordinary differential equation

$$\frac{d}{dt}I(x + t\theta) = -I(x + t\theta)u(x + t\theta), \quad (1.12)$$

Let  $R > 0$  be the radius of the domain of interest centred at the origin. Then, we integrate (1.12) from  $t = -\sqrt{R^2 - \|x\|_2^2}$ , the position of the emitter, to  $t = \sqrt{R^2 - \|x\|_2^2}$ , the position

<sup>3</sup>Figure adapted from Wikipedia <https://commons.wikimedia.org/w/index.php?curid=3001440>, by Begemotv2718, CC BY-SA 3.0.

of the detector, and obtain

$$\int_{-\sqrt{R^2-\|x\|_2^2}}^{\sqrt{R^2-\|x\|_2^2}} \frac{\frac{d}{dt}I(x+t\theta)}{I(x+t\theta)} dt = - \int_{-\sqrt{R^2-\|x\|_2^2}}^{\sqrt{R^2-\|x\|_2^2}} u(x+t\theta) dt.$$

Note that, due to  $d/dx \log(f(x)) = f'(x)/f(x)$ , the left hand side in the above equation simplifies to

$$\int_{-\sqrt{R^2-\|x\|_2^2}}^{\sqrt{R^2-\|x\|_2^2}} \frac{\frac{d}{dt}I(x+t\theta)}{I(x+t\theta)} dt = \log \left( I \left( x + \sqrt{R^2 - \|x\|_2^2} \theta \right) \right) - \log \left( I \left( x - \sqrt{R^2 - \|x\|_2^2} \theta \right) \right).$$

As we know the radiation intensity at both the emitter and the detector, we therefore know  $f(x, \theta) := \log(I(x - \theta \sqrt{R^2 - \|x\|_2^2})) - \log(I(x + \theta \sqrt{R^2 - \|x\|_2^2}))$  and we can write the estimation of the unknown density  $u$  as the inverse problem of the X-ray transform (1.11) (if we further assume that  $u$  can be continuously extended to zero outside of the circle of radius  $R$ ).

### Positron Emission Tomography (PET)

In Positron Emission Tomography (PET) a so-called radioactive tracer (a positron emitting radionuclide on a biologically active molecule) is injected into a patient (or subject). The emitted positrons of the tracer will interact with the subjects' electrons after travelling a short distance (usually less than 1mm), causing the annihilation of both the positron and the electron, which results in a pair of gamma rays moving into (approximately) opposite directions. This pair of photons is detected by the scanner detectors, and an intensity  $f(\varphi, s)$  can be associated with the number of annihilations detected at the detector pair that forms the line with offset  $s$  and angle  $\varphi$  (with respect to the reference coordinate system). Thus, we can consider the problem of recovering the unknown tracer density  $u$  as a solution of the inverse problem (1.10) again. The line of integration is determined by the position of the detector pairs and the geometry of the scanner.

## Chapter 2

# Linear inverse problems

Throughout this lecture we deal with functional analytic operators. For the sake of brevity, we cannot recall all basic concepts of functional analysis but refer to popular textbooks that deal with this subject, like [4, 16]. Nevertheless, we want to recall a few important properties that will be important for this lecture.

In particular, we will focus mainly on inverse problems with *bounded, linear operators*  $K$  only, i.e.  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  with

$$\|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} := \sup_{u \in \mathcal{U} \setminus \{0\}} \frac{\|Ku\|_{\mathcal{V}}}{\|u\|_{\mathcal{U}}} = \sup_{\|u\|_{\mathcal{U}} \leq 1} \|Ku\|_{\mathcal{V}} < \infty.$$

For  $K: \mathcal{U} \rightarrow \mathcal{V}$  we further want to denote by

- (a)  $\mathcal{D}(K) := \mathcal{U}$  the domain
- (b)  $\mathcal{N}(K) := \{u \in \mathcal{U} \mid Ku = 0\}$  the kernel
- (c)  $\mathcal{R}(K) := \{f \in \mathcal{V} \mid f = Ku, u \in \mathcal{U}\}$  the range

of  $K$ , see Figure 2.1

We say that  $K$  is continuous in  $u \in \mathcal{U}$  if there exists a  $\delta > 0$  for all  $\varepsilon > 0$  with

$$\|Ku - Kv\|_{\mathcal{V}} \leq \varepsilon \text{ for all } v \in \mathcal{U} \text{ with } \|u - v\|_{\mathcal{U}} \leq \delta.$$

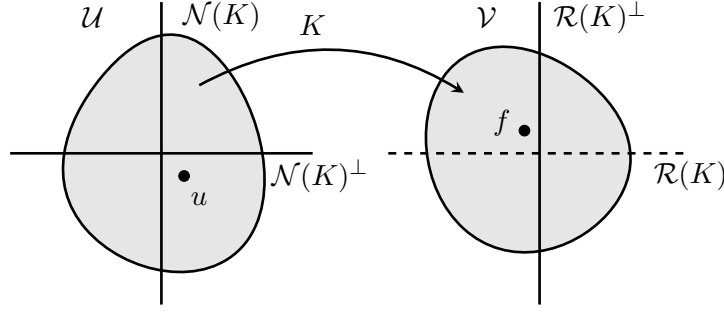
For linear  $K$  it can be shown that continuity is equivalent to the existence of a constant  $C > 0$  such that

$$\|Ku\|_{\mathcal{V}} \leq C\|u\|_{\mathcal{U}}$$

for all  $u \in \mathcal{U}$ . Note that this constant  $C$  actually equals the operator norm  $\|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})}$ .

For the first part of the lecture we only consider  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  with  $\mathcal{U}$  and  $\mathcal{V}$  being Hilbert spaces. From functional calculus we know that every Hilbert space is equipped with a *scalar product*, which we are going to denote by  $\langle \cdot, \cdot \rangle_{\mathcal{U}}$  (if  $\mathcal{U}$  denotes the corresponding Hilbert space). In analogy to the transpose of a matrix, this scalar product structure together with the theorem of Fréchet-Riesz [16, Section 2.10, Theorem 2.E] allows us to define the (unique) *adjoint operator* of  $K$ , denoted with  $K^*$ , as follows:

$$\langle Ku, v \rangle_{\mathcal{V}} = \langle u, K^*v \rangle_{\mathcal{U}}, \text{ for all } u \in \mathcal{U}, v \in \mathcal{V}.$$



**Figure 2.1:** Visualization of the setting for linear inverse problems where we want to solve the inverse problem (1.1). The operator  $K$  is a linear mapping between  $\mathcal{U}$  and  $\mathcal{V}$ . The kernel  $\mathcal{N}(K)$  and range  $\mathcal{R}(K)$  are used to analyse solutions to the inverse problem.

In addition to that, a scalar product allows to have a notion of orthogonality. Two elements  $u, v \in \mathcal{U}$  are said to be *orthogonal* if  $\langle u, v \rangle_{\mathcal{U}} = 0$ . For a subset  $\mathcal{X} \subset \mathcal{U}$  the *orthogonal complement* of  $\mathcal{X}$  in  $\mathcal{U}$  is defined as

$$\mathcal{X}^{\perp} := \{u \in \mathcal{U} \mid \langle u, v \rangle_{\mathcal{U}} = 0 \text{ for all } v \in \mathcal{X}\}.$$

One can show that  $\mathcal{X}^{\perp}$  is a closed subspace and that  $\mathcal{U}^{\perp} = \{0\}$ . Moreover, we have  $\mathcal{X} \subset (\mathcal{X}^{\perp})^{\perp}$ . If  $\mathcal{X}$  is a closed subspace we even have  $\mathcal{X} = (\mathcal{X}^{\perp})^{\perp}$ . In this case there exists the *orthogonal decomposition*

$$\mathcal{U} = \mathcal{X} \oplus \mathcal{X}^{\perp},$$

which means that every element  $u \in \mathcal{U}$  can uniquely be represented as

$$u = x + x^{\perp} \text{ with } x \in \mathcal{X} \text{ and } x^{\perp} \in \mathcal{X}^{\perp},$$

see for instance [16, Section 2.9, Corollary 1].

The mapping  $u \mapsto x$  defines a linear operator  $P_{\mathcal{X}} \in \mathcal{L}(\mathcal{U}, \mathcal{U})$  that is called *orthogonal projection* on  $\mathcal{X}$ .

**Lemma 2.1** (cf. [11, Section 5.16]). *Let  $\mathcal{X} \subset \mathcal{U}$  be a closed subspace. The orthogonal projection onto  $\mathcal{X}$  satisfies the following conditions:*

- (a)  $P_{\mathcal{X}}$  is self-adjoint, i.e.  $P_{\mathcal{X}}^* = P_{\mathcal{X}}$ ,
- (b)  $\|P_{\mathcal{X}}\|_{\mathcal{L}(\mathcal{U}, \mathcal{U})} = 1$  (if  $\mathcal{X} \neq \{0\}$ ),
- (c)  $I - P_{\mathcal{X}} = P_{\mathcal{X}^{\perp}}$ ,
- (d)  $\|u - P_{\mathcal{X}}u\|_{\mathcal{U}} \leq \|u - v\|_{\mathcal{U}}$  for all  $v \in \mathcal{X}$ ,
- (e)  $x = P_{\mathcal{X}}u$  if and only if  $x \in \mathcal{X}$  and  $u - x \in \mathcal{X}^{\perp}$ .

**Remark 2.1.** Note that for a non-closed subspace  $\mathcal{X}$  we only have  $(\mathcal{X}^{\perp})^{\perp} = \overline{\mathcal{X}}$ . For  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  we therefore have

- $\mathcal{R}(K)^{\perp} = \mathcal{N}(K^*)$  and thus  $\mathcal{N}(K^*)^{\perp} = \overline{\mathcal{R}(K)}$ ,
- $\mathcal{R}(K^*)^{\perp} = \mathcal{N}(K)$  and thus  $\mathcal{N}(K)^{\perp} = \overline{\mathcal{R}(K^*)}$ .

Hence, we can conclude the orthogonal decompositions

$$\mathcal{U} = \mathcal{N}(K) \oplus \overline{\mathcal{R}(K^*)} \text{ and } \mathcal{V} = \mathcal{N}(K^*) \oplus \overline{\mathcal{R}(K)}.$$

In the following we want to investigate the concept of generalised inverses of bounded, linear operators, before we will identify compactness of operators as the major source of ill-posedness. Subsequently, we are going to discuss this in more detail by analysing compact operators in terms of their singular value decomposition.

## 2.1 Generalised solutions

In order to overcome the issues of non-existence or non-uniqueness of (1.1) we want to generalise the concept of least squares solutions to linear operators in Hilbert spaces.

If we consider the generic inverse problem (1.1) again, we know that there does not exist a solution of the inverse problem if  $f \notin \mathcal{R}(K)$ . In that case it seems reasonable to find an element  $u \in \mathcal{U}$  for which  $\|Ku - f\|_{\mathcal{V}}$  gets minimal instead. If  $\mathcal{V} = L^2$  then  $u$  minimizes the squared error and thus motivates the name least squares solution.

However, for  $\mathcal{N}(K) \neq \{0\}$  there are infinitely many solutions that minimise  $\|Ku - f\|_{\mathcal{V}}$  of which we have to pick one. Picking the one with minimal norm  $\|u\|_{\mathcal{U}}$  brings us to the definition of the minimal norm solution.

**Definition 2.1.** We call  $u \in \mathcal{U}$  a least squares solution of the inverse problem (1.1), if

$$\|Ku - f\|_{\mathcal{V}} \leq \|Kv - f\|_{\mathcal{V}} \quad \text{for all } v \in \mathcal{U}. \quad (2.1)$$

Furthermore, we call  $u^\dagger \in \mathcal{U}$  a minimal norm solution of the inverse problem (1.1), if

$$\|u^\dagger\|_{\mathcal{U}} \leq \|v\|_{\mathcal{U}} \quad \text{for all least squares solutions } v. \quad (2.2)$$

**Remark 2.2.** Let  $u$  be a least squares solution to  $Ku = f$ . It is easy to see that each  $v \in \{u\} + \mathcal{N}(K)$  is a least squares solution as well.

Moreover, let  $u^\dagger$  be a minimal norm solution, then  $u^\dagger \in \mathcal{N}(K)^\perp$ . Assume to the contrary that this was not the case. Then, as  $\mathcal{N}(K)$  is closed for  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ , there exists elements  $x^\perp \in \mathcal{N}(K)^\perp$  and  $x \in \mathcal{N}(K)$  with  $\|x\|_{\mathcal{U}} > 0$  such that  $u^\dagger = x + x^\perp$ . Clearly,  $x^\perp$  is a least squares solution and by

$$\|u^\dagger\|_{\mathcal{U}}^2 = \|x + x^\perp\|_{\mathcal{U}}^2 = \|x^\perp\|_{\mathcal{U}}^2 + \underbrace{2\langle x^\perp, x \rangle_{\mathcal{U}}}_{=0} + \|x\|_{\mathcal{U}}^2 > \|x^\perp\|_{\mathcal{U}}^2$$

has smaller norm than  $u^\dagger$ , which contradicts that  $u^\dagger$  is of minimal norm, thus  $u^\dagger \in \mathcal{N}(K)^\perp$ .

In numerical linear algebra it is a well known fact that the normal equations can be considered to compute least squares solutions. The same holds true in the infinite-dimensional case.

**Theorem 2.1.** Let  $f \in \mathcal{V}$  and  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ . Then, the following three assertions are equivalent.

(a)  $u \in \mathcal{U}$  satisfies  $Ku = P_{\overline{\mathcal{R}(K)}}f$ .

(b)  $u$  is a least squares solution of the inverse problem (1.1).

(c)  $u$  solves the normal equation

$$K^*Ku = K^*f. \quad (2.3)$$

**Remark 2.3.** The name normal equation is derived from the fact that for any solution  $u$  its residual  $Ku - f$  is orthogonal (normal) to  $\mathcal{R}(K)$ . This can be readily seen, as we have for any  $v \in \mathcal{U}$  that

$$0 = \langle v, K^*(Ku - f) \rangle_{\mathcal{U}} = \langle Kv, Ku - f \rangle_{\mathcal{V}}$$

which shows  $Ku - f \in \mathcal{R}(K)^\perp$ .

*Proof of Theorem 2.1.* For (a)  $\Rightarrow$  (b): Let  $u \in \mathcal{U}$  such that  $Ku = P_{\overline{\mathcal{R}(K)}}f$  and let  $v \in \mathcal{U}$  be arbitrary. With the basic properties of the orthogonal projection, Lemma 2.1 (d), we have

$$\|Ku - f\|_{\mathcal{V}}^2 = \|(I - P_{\overline{\mathcal{R}(K)}})f\|_{\mathcal{V}}^2 \leq \inf_{g \in \overline{\mathcal{R}(K)}} \|g - f\|_{\mathcal{V}}^2 \leq \inf_{v \in \mathcal{U}} \|Kv - f\|_{\mathcal{V}}^2,$$

which shows that  $u$  is a least squares solution. Here, the last inequality follows from  $\mathcal{R}(K) \subset \overline{\mathcal{R}(K)}$ .

For (b)  $\Rightarrow$  (c): Let  $u \in \mathcal{U}$  be a least squares solution and let  $v \in \mathcal{U}$  an arbitrary element. We define the quadratic polynomial  $F: \mathbb{R} \rightarrow \mathbb{R}$ ,

$$F(\lambda) := \|K(u + \lambda v) - f\|_{\mathcal{V}}^2 = \lambda^2 \|Kv\|_{\mathcal{V}}^2 - 2\lambda \langle Kv, f - Ku \rangle_{\mathcal{V}} + \|f - Ku\|_{\mathcal{V}}^2.$$

A necessary condition for  $u \in \mathcal{U}$  to be a least squares solution is  $F'(0) = 0$ , which leads to  $\langle v, K^*(f - Ku) \rangle_{\mathcal{U}} = 0$ . As  $v$  was arbitrary, it follows that the normal equation (2.3) must hold.

For (c)  $\Rightarrow$  (a): From the normal equation it follows that  $K^*(f - Ku) = 0$ , which is equivalent to  $f - Ku \in \mathcal{R}(K)^\perp$ , see Remark 2.3. Since  $\mathcal{R}(K)^\perp = \left(\overline{\mathcal{R}(K)}\right)^\perp$  and  $Ku \in \mathcal{R}(K) \subset \overline{\mathcal{R}(K)}$ , the assertion follows from Lemma 2.1 (e):

$$Ku = P_{\overline{\mathcal{R}(K)}}f \Leftrightarrow Ku \in \overline{\mathcal{R}(K)} \text{ and } f - Ku \in \left(\overline{\mathcal{R}(K)}\right)^\perp.$$

□

**Lemma 2.2.** Let  $f \in \mathcal{V}$  and let  $\mathbb{L}$  be the set of least squares solutions to the inverse problem (1.1). Then,  $\mathbb{L}$  is non-empty if and only if  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ .

*Proof.* Let  $u \in \mathbb{L}$ . It is easy to see that  $f = Ku + (f - Ku) \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$  as the normal equations are equivalent to  $f - Ku \in \mathcal{R}(K)^\perp$ .

Consider now  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ . Then there exists  $u \in \mathcal{U}$  and  $g \in \mathcal{R}(K)^\perp = \left(\overline{\mathcal{R}(K)}\right)^\perp$  such that  $f = Ku + g$  and thus  $P_{\overline{\mathcal{R}(K)}}f = P_{\overline{\mathcal{R}(K)}}Ku + P_{\overline{\mathcal{R}(K)}}g = Ku$  and the assertion follows from Theorem 2.1 (a). □

**Remark 2.4.** If the dimensions of  $\mathcal{U}$  and  $\mathcal{R}(K)$  are finite, then  $\mathcal{R}(K)$  is closed, i.e.  $\overline{\mathcal{R}(K)} = \mathcal{R}(K)$ . Thus, in a finite dimensional setting, there always exists a least squares solution.

It is natural to ask whether there are always least squares solutions. From the above remark it is clear that we have to look for an example in infinite dimensional spaces. The answer is negative as we see from the following counter example.

**Example 2.1.** Let  $\mathcal{U} = \ell^2$ ,  $\mathcal{V} = \ell^2$ , where the space  $\ell^2$  is the space of all square summable sequences, i.e.

$$\ell^2 := \left\{ \{x_j\}_{j \in \mathbb{N}} \mid x_j \in \mathbb{R}, \sum_{j=1}^{\infty} x_j^2 < \infty \right\}.$$



It is a Hilbert space with inner product and norm given by

$$\langle x, y \rangle_{\ell^2} := \sum_{j=1}^{\infty} x_j y_j \quad \text{and} \quad \|x\|_{\ell^2} := \left( \sum_{j=1}^{\infty} x_j^2 \right)^{1/2}, \quad \text{respectively.}$$

For more information see, for instance, [4].

Consider the inverse problem  $Kx = f$ , where the linear operator  $K: \ell^2 \rightarrow \ell^2$  is defined by

$$(Kx)_j := \frac{x_j}{j}.$$

and the data by  $f_j := j^{-1}$ . It is easy to see that  $K$  is linear and bounded, i.e.  $K \in \mathcal{L}(\ell^2, \ell^2)$  and  $f \in \ell^2$ .

We will show that  $f \in \overline{\mathcal{R}(K)} \setminus \mathcal{R}(K)$  and thus  $f \notin \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ . With Lemma 2.2 it follows then that there are no least squares solutions.

First we show that  $f \notin \mathcal{R}(K)$  by contradiction. Assume that  $f \in \mathcal{R}(K)$ , then there exists  $x \in \ell^2$  such that  $Kx = f$  and thus  $j^{-1}x_j = j^{-1}$  for all  $j \in \mathbb{N}$ . Therefore, we have  $x_j = 1$  and  $x \notin \ell^2$ .

Next we show that  $f \in \overline{\mathcal{R}(K)}$ . Let  $\{x^k\}_{k \in \mathbb{N}} \subset \ell^2$  be a sequence in  $\ell^2$  (each element is a sequence as well), with

$$(x^k)_j := \begin{cases} 1, & j \leq k \\ 0, & j > k \end{cases}.$$

It is easy to see that  $x^k \in \ell^2$  as it has only finitely many non-negative components. In addition, we have

$$f^k := Kx^k, \quad (f^k)_j = \begin{cases} \frac{1}{j}, & j \leq k \\ 0, & j > k \end{cases}$$

and therefore

$$\|f - f^k\|_{\ell^2}^2 = \sum_{j=k+1}^{\infty} f_j^2 = \sum_{j=1}^{\infty} f_j^2 - \sum_{j=1}^k f_j^2 \rightarrow 0 \quad \text{as } k \rightarrow \infty$$

by definition of a convergent series. Therefore,  $f^k \rightarrow f$  in  $\ell^2$  and thus  $f \in \overline{\mathcal{R}(K)}$ .

**Theorem 2.2.** *Let  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ . Then there exists a unique minimal norm solution  $u^\dagger$  to the inverse problem (1.1) and all least squares solutions are given by  $\{u^\dagger\} + \mathcal{N}(K)$ .*

*Proof.* From Lemma 2.2 we know that there exist least squares solutions and denote any arbitrary two of them (not necessarily different) by  $u, v \in \mathcal{U}$ . Then there exist  $\varphi, \psi \in \mathcal{N}(K)^\perp$  and  $x, y \in \mathcal{N}(K)$  such that  $u = \varphi + x$  and  $v = \psi + y$ . As we noted in Remark 2.2  $\varphi$  and  $\psi$  are least squares solutions as well. With Theorem 2.1 we conclude

$$K(\varphi - \psi) = K\varphi - K\psi = P_{\overline{\mathcal{R}(K)}}f - P_{\overline{\mathcal{R}(K)}}f = 0, \quad (2.4)$$

which shows that  $\varphi - \psi \in \mathcal{N}(K)$ . But as  $\varphi - \psi \in \mathcal{N}(K)^\perp$  and  $\mathcal{N}(K) \cap \mathcal{N}(K)^\perp = \{0\}$  we see that  $\varphi = \psi$ . Therefore all least squares solutions are of the form  $\{\varphi\} + \mathcal{N}(K)$ .

Moreover, we know that  $u^\dagger$  is a least squares solution and that  $u^\dagger \in \mathcal{N}(K)^\perp$ , see Remark 2.2. Thus we have that  $u^\dagger = \varphi$ , which completes the proof.  $\square$

**Corollary 2.1.** *The minimal norm solution is the unique solution of the normal equation in  $\mathcal{N}(K)^\perp$ .*

## 2.2 Generalised inverse

We have seen that, for arbitrary  $f \in \mathcal{V}$ , a least squares solution does not need to exist if  $\mathcal{R}(K)$  is not closed. If, however, a least squares solution exists, then we have shown that the minimum norm solution is unique. We will see in the following that the minimum norm solution can be computed via the *Moore-Penrose* generalised inverse.

**Definition 2.2.** Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  and let

$$\tilde{K} := K|_{\mathcal{N}(K)^\perp} : \mathcal{N}(K)^\perp \rightarrow \mathcal{R}(K)$$

denote the restriction of  $K$  to  $\mathcal{N}(K)^\perp$ . The Moore-Penrose inverse  $K^\dagger$  is defined as the unique linear extension of  $\tilde{K}^{-1}$  to

$$\mathcal{D}(K^\dagger) = \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$$

with

$$\mathcal{N}(K^\dagger) = \mathcal{R}(K)^\perp.$$

**Remark 2.5.** Due to the restriction to  $\mathcal{N}(K)^\perp$  and  $\mathcal{R}(K)$  we have that  $\tilde{K}$  is injective and surjective. Hence,  $\tilde{K}^{-1}$  is linear and exists and—as a consequence— $K^\dagger$  is well-defined on  $\mathcal{R}(K)$  and linear.

Moreover, due to the orthogonal decomposition  $\mathcal{D}(K^\dagger) = \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ , there exists for arbitrary  $f \in \mathcal{D}(K^\dagger)$  elements  $f_1 \in \mathcal{R}(K)$  and  $f_2 \in \mathcal{R}(K)^\perp$  with  $f = f_1 + f_2$ . Therefore, we have

$$K^\dagger f = K^\dagger f_1 + K^\dagger f_2 = K^\dagger f_1 = \tilde{K}^{-1} f_1 = \tilde{K}^{-1} P_{\overline{\mathcal{R}(K)}} f, \quad (2.5)$$

where we used that  $f_2 \in \mathcal{R}(K)^\perp = \mathcal{N}(K^\dagger)$ . Thus,  $K^\dagger$  is well-defined on the entire domain  $\mathcal{D}(K^\dagger)$ .

Note that, if  $K$  is bijective we have that  $K^\dagger = K^{-1}$ . Moreover, we highlight that the extension  $K^\dagger$  is not necessarily continuous.

**Example 2.2.** To illustrate the definition of the Moore-Penrose inverse we consider a simple example in finite dimensions. Let the linear operator  $K: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  be given by

$$Kx = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 2x_1 \\ 0 \end{pmatrix}.$$

It is easy to see that  $\mathcal{R}(K) = \{f \in \mathbb{R}^2 \mid f_2 = 0\}$  and  $\mathcal{N}(K) = \{x \in \mathbb{R}^3 \mid x_1 = 0\}$ . Thus,  $\mathcal{N}(K)^\perp = \{x \in \mathbb{R}^3 \mid x_2, x_3 = 0\}$ . Therefore,  $\tilde{K}: \mathcal{N}(K)^\perp \rightarrow \mathcal{R}(K)$ , given by  $x \mapsto (2x_1, 0)^\top$ , is bijective and its inverse  $\tilde{K}^{-1}: \mathcal{R}(K) \rightarrow \mathcal{N}(K)^\perp$  is given by  $f \mapsto (f_1/2, 0, 0)^\top$ .

As the orthogonal projection onto  $\mathcal{R}(K)$  is given by  $f = (f_1, f_2) \mapsto (f_1, 0)$ , the Moore-Penrose inverse of  $K$  is  $K^\dagger: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ ,

$$K^\dagger f = \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} f_1 \\ f_2 \end{pmatrix} = \begin{pmatrix} f_1/2 \\ 0 \\ 0 \end{pmatrix}.$$

Let us consider data  $\tilde{f} = (8, 1)^\top \notin \mathcal{R}(K)$ . Then,  $K^\dagger \tilde{f} = K^\dagger(8, 1)^\top = (4, 0, 0)^\top$ .

It can be shown that  $K^\dagger$  can be characterized by the Moore-Penrose equations.

**Lemma 2.3.** *The Moore-Penrose inverse  $K^\dagger$  satisfies  $\mathcal{R}(K^\dagger) = \mathcal{N}(K)^\perp$  and the Moore-Penrose equations*

- (a)  $KK^\dagger K = K$ ,
- (b)  $K^\dagger KK^\dagger = K^\dagger$ ,
- (c)  $K^\dagger K = I - P_{\mathcal{N}(K)}$ ,
- (d)  $KK^\dagger = P_{\overline{\mathcal{R}(K)}}|_{\mathcal{D}(K^\dagger)}$ ,

where  $P_{\mathcal{N}(K)}$  and  $P_{\overline{\mathcal{R}(K)}}$  denote the orthogonal projections on  $\mathcal{N}(K)$  and  $\overline{\mathcal{R}(K)}$ , respectively.

*Proof.* First we prove that  $\mathcal{R}(K^\dagger) = \mathcal{N}(K)^\perp$ . Let  $u \in \mathcal{R}(K^\dagger)$ . Then, there exists a  $f \in \mathcal{D}(K^\dagger)$  with  $u = K^\dagger f$  and according to (2.5) we observe that  $u = K^\dagger f = \tilde{K}^{-1} P_{\overline{\mathcal{R}(K)}} f$ . Hence,  $u \in \mathcal{R}(\tilde{K}^{-1}) = \mathcal{N}(K)^\perp$  and therefore  $\mathcal{R}(K^\dagger) \subseteq \mathcal{N}(K)^\perp$ . To prove  $\mathcal{N}(K)^\perp \subseteq \mathcal{R}(K^\dagger)$ , let  $u \in \mathcal{N}(K)^\perp$  and it holds  $u = \tilde{K}^{-1} \tilde{K} u = K^\dagger K u$ . Thus,  $u \in \mathcal{R}(K^\dagger)$  showing set equality.

It remains to prove the Moore-Penrose equations:

(d): For  $f \in \mathcal{D}(K^\dagger)$  it follows from (2.5) and  $K = \tilde{K}$  on  $\mathcal{N}(K)^\perp$  that

$$KK^\dagger f = K \tilde{K}^{-1} P_{\overline{\mathcal{R}(K)}} f = \tilde{K} \tilde{K}^{-1} P_{\overline{\mathcal{R}(K)}} f = P_{\overline{\mathcal{R}(K)}} f.$$

(c): According to the definition of  $K^\dagger$  we have  $K^\dagger K u = \tilde{K}^{-1} K u$  for all  $u \in \mathcal{U}$  and thus

$$K^\dagger K u = \underbrace{\tilde{K}^{-1} K P_{\mathcal{N}(K)} u}_{=0} + \underbrace{\tilde{K}^{-1} K (I - P_{\mathcal{N}(K)}) u}_{=P_{\mathcal{N}(K)^\perp}} = (I - P_{\mathcal{N}(K)}) u,$$

where we have used Lemma 2.1 (c) and the fact that  $\mathcal{N}(K)$  is closed.

(b): Inserting (d) into (2.5) yields

$$K^\dagger f = K^\dagger P_{\overline{\mathcal{R}(K)}} f = K^\dagger K K^\dagger f.$$

(a): With (c) we have

$$KK^\dagger K = K(I - P_{\mathcal{N}(K)}) = K - K P_{\mathcal{N}(K)} = K.$$

□

The following theorem states that minimum norm solutions can be computed via the generalised inverse.

**Theorem 2.3.** *For each  $f \in \mathcal{D}(K^\dagger)$ , the minimal norm solution  $u^\dagger$  to the inverse problem (1.1) is given via*

$$u^\dagger = K^\dagger f.$$

*Proof.* As  $f \in \mathcal{D}(K^\dagger)$ , we know from Theorem 2.2 that the minimal norm solution  $u^\dagger$  exists and is unique. With  $u^\dagger \in \mathcal{N}(K)^\perp$ , Lemma 2.3, and Theorem 2.1 we conclude that

$$u^\dagger = (I - P_{\mathcal{N}(K)}) u^\dagger = K^\dagger K u^\dagger = K^\dagger P_{\overline{\mathcal{R}(K)}} f = K^\dagger K K^\dagger f = K^\dagger f.$$

□

As a consequence of Theorem 2.3 and Theorem 2.1, we find that the minimum norm solution  $u^\dagger$  of  $Ku = f$  is a minimum norm solution of the normal equation (2.3), i.e.

$$u^\dagger = (K^*K)^\dagger K^* f.$$

Thus, in order to compute  $u^\dagger$  we can equivalently consider finding the minimum norm solution of the normal equation.

At the end of this section we further want to analyse the domain of the generalised inverse in more detail. Due to the construction of the Moore-Penrose inverse we have  $\mathcal{D}(K^\dagger) = \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ . As orthogonal complements are always closed we can conclude

$$\overline{\mathcal{D}(K^\dagger)} = \overline{\mathcal{R}(K)} \oplus \mathcal{R}(K)^\perp = \mathcal{V},$$

and hence,  $\mathcal{D}(K^\dagger)$  is dense in  $\mathcal{V}$ . Thus, if  $\mathcal{R}(K)$  is closed it follows that  $\mathcal{D}(K^\dagger) = \mathcal{V}$  and on the other hand,  $\mathcal{D}(K^\dagger) = \mathcal{V}$  implies  $\mathcal{R}(K)$  is closed.

Moreover, for  $f \in \mathcal{R}(K)^\perp = \mathcal{N}(K^\dagger)$  the minimum norm solution is  $u^\dagger = 0$ . Therefore, for given  $f \in \overline{\mathcal{R}(K)}$ , the important question to address is when  $f$  also satisfies  $f \in \mathcal{R}(K)$ . If this is the case,  $K^\dagger$  has to be continuous. However, the existence of a single element  $f \in \overline{\mathcal{R}(K)} \setminus \mathcal{R}(K)$  is enough already to prove that  $K^\dagger$  is discontinuous.

**Definition 2.3.** Let  $\mathcal{V}$  and  $\mathcal{U}$  be Hilbert spaces and consider  $A: \mathcal{V} \rightarrow \mathcal{U}$ . We call the graph of  $A$ ,

$$\text{gr}(A) := \{(f, u) \in \mathcal{V} \times \mathcal{U} \mid Af = u\},$$

closed if for any sequence  $\{(f_j, u_j)\}_{j \in \mathbb{N}}$  with  $u_j = Af_j$ ,  $f_j \rightarrow f \in \mathcal{V}$ , and  $u_j \rightarrow u \in \mathcal{U}$  we have that  $Af = u$ .

**Theorem 2.4** (Closed graph theorem [14, Proposition 2.14 and Theorem 2.15]). Let  $\mathcal{V}$  and  $\mathcal{U}$  be Hilbert spaces and let  $A: \mathcal{V} \rightarrow \mathcal{U}$  be a linear mapping with a closed graph. Then  $A \in \mathcal{L}(\mathcal{V}, \mathcal{U})$ .

**Theorem 2.5.** Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ . Then  $K^\dagger$  is continuous, i.e.  $K^\dagger \in \mathcal{L}(\mathcal{D}(K^\dagger), \mathcal{U})$ , if and only if  $\mathcal{R}(K)$  is closed.

*Proof.* We will show first that the graph of the Moore-Penrose inverse is closed. To this end, let  $\{(f_j, u_j)\}_{j \in \mathbb{N}} \subset \text{gr}(K^\dagger)$  be a sequence in the graph of the Moore-Penrose inverse, i.e.  $u_j = K^\dagger f_j$ , and  $f_j \rightarrow f$  and  $u_j \rightarrow u$ . Then, because of the continuity of  $K$ , Lemma 2.3 (d), and the continuity of the orthogonal projection, we have

$$Ku = \lim_{j \rightarrow \infty} Ku_j = \lim_{j \rightarrow \infty} KK^\dagger f_j = \lim_{j \rightarrow \infty} P_{\overline{\mathcal{R}(K)}} f_j = P_{\overline{\mathcal{R}(K)}} f.$$

Thus, by Theorem 2.1,  $u$  is a least squares solution. As  $K^\dagger f_j \in \mathcal{N}(K)^\perp$  and  $\mathcal{N}(K)^\perp$  is closed, we have  $u \in \mathcal{N}(K)^\perp$  and it follows from the uniqueness of the minimal norm solution that  $u = K^\dagger f$ . This shows that the graph of  $K^\dagger$  is closed.

For the proof of the theorem, assume first that  $\mathcal{R}(K)$  is closed so that  $\mathcal{D}(K^\dagger) = \mathcal{V}$ . Then, by the closed graph theorem (Theorem 2.4),  $K^\dagger$  is bounded and therefore continuous.

Conversely, let  $K^\dagger$  be continuous. As  $\mathcal{D}(K^\dagger)$  is dense in  $\mathcal{V}$ , there is a unique continuous extension  $A$  of  $K^\dagger$  to  $\mathcal{V}$ ,

$$Af := \lim_{j \rightarrow \infty} K^\dagger f_j \quad \text{for } \{f_j\}_{j \in \mathbb{N}} \subset \mathcal{D}(K^\dagger) \text{ with } f_j \rightarrow f \in \mathcal{V}.$$

Now let  $f \in \overline{\mathcal{R}(K)}$  and let  $\{f_j\}_{j \in \mathbb{N}} \subset \mathcal{D}(K^\dagger)$  with  $f_j \rightarrow f$ . Then, from Lemma 2.3 (d) we find that

$$f = P_{\overline{\mathcal{R}(K)}} f = \lim_{j \rightarrow \infty} P_{\overline{\mathcal{R}(K)}} f_j = \lim_{j \rightarrow \infty} KK^\dagger f_j = KAf \in \mathcal{R}(K)$$

and thus  $\overline{\mathcal{R}(K)} = \mathcal{R}(K)$  showing that  $\mathcal{R}(K)$  is closed.  $\square$

In the next section we are going to discover that the class of compact operators is a class for which the Moore-Penrose inverses are discontinuous.

## 2.3 Compact operators

Compact operators are very common in inverse problems. In fact, almost all (linear) inverse problems involve the inversion of compact operators. Compact operators are defined as follows.

**Definition 2.4.** *Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ . Then  $K$  is said to be compact if the image of a bounded sequence  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  contains a convergent subsequence  $\{Ku_{j_k}\}_{k \in \mathbb{N}} \subset \mathcal{V}$ . We denote the space of compact operators by  $\mathcal{K}(\mathcal{U}, \mathcal{V})$ .*

**Remark 2.6.** We can equivalently define an operator  $K$  to be compact if and only if for any bounded set  $B$ , the closure of its image  $\overline{K(B)}$  is compact.

**Example 2.3** (Follows from e.g. [17, p. 49]). Let  $I: \mathcal{U} \rightarrow \mathcal{U}$  be the identity operator on  $\mathcal{U}$ , i.e.  $u \mapsto u$ . Then  $I$  is compact if and only if the dimension of  $\mathcal{U}$  is finite.

**Example 2.4** (e.g. [17, p. 286, Proposition 5] or [4, p. 186]). Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ . If the range of  $K$  is finite dimensional, then  $K$  is compact.

**Example 2.5** ([10, p. 230]). The operator  $K: \ell^2 \rightarrow \ell^2$ ,  $(Kx)_j = j^{-1}x_j$  from Example 2.1 is compact.

**Example 2.6** ([10, p. 231]). Let  $\emptyset \neq \Omega \subset \mathbb{R}^n$  be compact. Let  $k \in L^2(\Omega \times \Omega)$  and define the integral operator  $K: L^2(\Omega) \rightarrow L^2(\Omega)$  with

$$(Ku)(x) = \int_{\Omega} k(x, y)u(y) dy.$$

Then,  $K$  is compact.

**Example 2.7** ([12, p. 38]). Let  $B := \{x \in \mathbb{R}^2 \mid \|x\| \leq 1\}$  denote the unit ball in  $\mathbb{R}^2$  and  $Z := [-1, 1] \times [0, \pi)$ . Moreover, let  $\theta(\varphi) := (\cos(\varphi), \sin(\varphi))^\top$ ,  $\theta^\perp(\varphi) := (\sin(\varphi), -\cos(\varphi))^\top$  be the unit vectors pointing in the direction described by  $\varphi$  and orthogonal to it. Then, the Radon transform/X-ray transform is defined as the operator  $R: L^2(B) \rightarrow L^2(Z)$  with

$$(Ru)(s, \varphi) := \int_{-\sqrt{1-s^2}}^{\sqrt{1-s^2}} u(s\theta(\varphi) + t\theta^\perp(\varphi)) dt.$$

It can be shown that the Radon transform is linear and continuous, i.e.  $R \in \mathcal{L}(L^2(B), L^2(Z))$ , and even compact, i.e.  $R \in \mathcal{K}(L^2(B), L^2(Z))$ .

Compact operators can be seen as the infinite dimensional analogue to ill-conditioned matrices. Indeed it can be seen that compactness is a main source of ill-posedness in infinite dimensions, confirmed by the following result.

**Theorem 2.6.** *Let  $K \in \mathcal{K}(\mathcal{U}, \mathcal{V})$  with an infinite dimensional range. Then, the Moore-Penrose inverse of  $K$  is discontinuous.*

*Proof.* As the range  $\mathcal{R}(K)$  is of infinite dimension, we can conclude that  $\mathcal{U}$  and  $\mathcal{N}(K)^\perp$  are also infinite dimensional. We can therefore find a sequence  $\{u_j\}_{j \in \mathbb{N}}$  with  $u_j \in \mathcal{N}(K)^\perp$ ,  $\|u_j\|_{\mathcal{U}} = 1$  and  $\langle u_j, u_k \rangle_{\mathcal{U}} = 0$  for  $j \neq k$ . Since  $K$  is a compact operator the sequence  $f_j = Ku_j$  has a convergent subsequence, hence, for all  $\delta > 0$  we can find  $j, k$  such that  $\|f_j - f_k\|_{\mathcal{V}} < \delta$ . However, we also obtain

$$\begin{aligned} \|K^\dagger f_j - K^\dagger f_k\|_{\mathcal{U}}^2 &= \|K^\dagger Ku_j - K^\dagger Ku_k\|_{\mathcal{U}}^2 \\ &= \|u_j - u_k\|_{\mathcal{U}}^2 = \|u_j\|_{\mathcal{U}}^2 - 2\langle u_j, u_k \rangle_{\mathcal{U}} + \|u_k\|_{\mathcal{U}}^2 = 2, \end{aligned}$$

which shows that  $K^\dagger$  is discontinuous. Here, the second identity follows from Lemma 2.3 (c) and the fact that  $u_j, u_k \in \mathcal{N}(K)^\perp$ .  $\square$

To have a better understanding of when we have  $f \in \overline{\mathcal{R}(K)} \setminus \mathcal{R}(K)$  for compact operators  $K$ , we want to consider the singular value decomposition of compact operators.

## 2.4 Singular value decomposition of compact operators

We want to characterise the Moore-Penrose inverse of compact operators in terms of a spectral decomposition. Like in the finite dimensional case of matrices, we can only expect a spectral decomposition to exist for self-adjoint operators.

**Theorem 2.7** ([10, p. 225, Theorem 9.16]). *Let  $\mathcal{U}$  be a Hilbert space and  $K \in \mathcal{K}(\mathcal{U}, \mathcal{U})$  be self-adjoint. Then there exists an orthonormal basis  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  of  $\overline{\mathcal{R}(K)}$  and a sequence of eigenvalues  $\{\lambda_j\}_{j \in \mathbb{N}} \subset \mathbb{R}$  with  $|\lambda_1| \geq |\lambda_2| \geq \dots > 0$  such that for all  $u \in \mathcal{U}$  we have*

$$Ku = \sum_{j=1}^{\infty} \lambda_j \langle u, u_j \rangle_{\mathcal{U}} u_j.$$

The sequence  $\{\lambda_j\}_{j \in \mathbb{N}}$  is either finite or we have  $\lambda_j \rightarrow 0$ .

**Remark 2.7.** The notation in the theorem above only makes sense if the sequence  $\{\lambda_j\}_{j \in \mathbb{N}}$  is infinite. For the case that there are only finitely many  $\lambda_j$  the sum has to be interpreted as a finite sum.

Moreover, as the eigenvalues are sorted by absolute value  $|\lambda_j|$ , we have  $\|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{U})} = |\lambda_1|$ .

Due to Theorem 2.1 we can consider  $K^*K$  instead of  $K$ , which brings us to the singular value decomposition of linear, compact operators.

**Theorem 2.8.** *Let  $K \in \mathcal{K}(\mathcal{U}, \mathcal{V})$ . Then there exists*

- (a) *a not-necessarily infinite null sequence  $\{\sigma_j\}_{j \in \mathbb{N}}$  with  $\sigma_1 \geq \sigma_2 \geq \dots > 0$ ,*
- (b) *an orthonormal basis  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  of  $\mathcal{N}(K)^\perp$ ,*
- (c) *an orthonormal basis  $\{v_j\}_{j \in \mathbb{N}} \subset \mathcal{V}$  of  $\overline{\mathcal{R}(K)}$  with*

$$Ku_j = \sigma_j v_j, \quad K^* v_j = \sigma_j u_j, \quad \text{for all } j \in \mathbb{N}. \quad (2.6)$$

Moreover, for all  $w \in \mathcal{U}$  we have the representation

$$Kw = \sum_{j=1}^{\infty} \sigma_j \langle w, u_j \rangle_{\mathcal{U}} v_j. \quad (2.7)$$

The sequence  $\{(\sigma_j, u_j, v_j)\}$  is called singular system or singular value decomposition (SVD) of  $K$ .

*Proof.* As  $K$  is compact we have that  $K^*K: \mathcal{U} \rightarrow \mathcal{U}$  is compact and self-adjoint. By Theorem 2.7 there exists a decreasing (in terms of absolute values) null sequence  $\{\lambda_j\}_{j \in \mathbb{N}} \subset \mathbb{R} \setminus \{0\}$  and an orthonormal basis  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  of  $\overline{\mathcal{R}(K^*K)}$  with  $K^*Ku = \sum_{j=1}^{\infty} \lambda_j \langle u, u_j \rangle_{\mathcal{U}} u_j$  for all  $u \in \mathcal{U}$ .

Due to

$$\lambda_j = \lambda_j \|u_j\|_{\mathcal{U}}^2 = \langle \lambda_j u_j, u_j \rangle_{\mathcal{U}} = \langle K^*K u_j, u_j \rangle_{\mathcal{U}} = \langle K u_j, K u_j \rangle_{\mathcal{V}} = \|K u_j\|_{\mathcal{V}}^2 > 0$$

we can define

$$\sigma_j := \sqrt{\lambda_j} \quad \text{and} \quad v_j := \sigma_j^{-1} K u_j \in \mathcal{V} \quad \text{for all } j \in \mathbb{N}.$$

Further, we observe

$$K^*v_j = \sigma_j^{-1} K^*K u_j = \sigma_j^{-1} \lambda_j u_j = \sigma_j u_j,$$

which proves Equation (2.6).

We also observe that  $\{v_j\}_{j \in \mathbb{N}}$  form an orthonormal basis due to

$$\langle v_i, v_j \rangle_{\mathcal{V}} = \frac{1}{\sigma_i \sigma_j} \langle K u_i, K u_j \rangle_{\mathcal{V}} = \frac{1}{\sigma_i \sigma_j} \langle K^*K u_i, u_j \rangle_{\mathcal{U}} = \frac{\lambda_i}{\sigma_i \sigma_j} \langle u_i, u_j \rangle_{\mathcal{U}} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{else.} \end{cases}$$

We know that  $\{u_j\}_{j \in \mathbb{N}}$  is an orthonormal basis of  $\overline{\mathcal{R}(K^*K)}$  and we want to show that it is also an orthonormal basis of  $\mathcal{N}(K)^\perp$ . As made apparent in Remark 2.1, we have  $\overline{\mathcal{R}(K^*)} = \mathcal{N}(K)^\perp$  and thus it is sufficient to show that  $\overline{\mathcal{R}(K^*K)} = \overline{\mathcal{R}(K^*)}$ .

It is clear that  $\overline{\mathcal{R}(K^*K)} = \overline{\mathcal{R}(K^*|_{\mathcal{R}(K)})} \subseteq \overline{\mathcal{R}(K^*)}$ , such that we are left to prove that  $\overline{\mathcal{R}(K^*)} \subseteq \overline{\mathcal{R}(K^*K)}$ .

Let  $u \in \overline{\mathcal{R}(K^*)}$  and let  $\varepsilon > 0$ . Then, there exists  $f \in \mathcal{N}(K^*)^\perp$  with  $\|K^*f - u\|_{\mathcal{U}} < \varepsilon/2$ . As  $\mathcal{N}(K^*)^\perp = \overline{\mathcal{R}(K)}$  (again see Remark 2.1), there exists  $x \in \mathcal{U}$  such that  $\|Kx - f\|_{\mathcal{V}} < \varepsilon/(2\|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})})$ . Putting these together we have

$$\begin{aligned} \|K^*Kx - u\|_{\mathcal{U}} &\leq \|K^*Kx - K^*f\|_{\mathcal{U}} + \|K^*f - u\|_{\mathcal{U}} \\ &\leq \underbrace{\|K^*\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} \|Kx - f\|_{\mathcal{V}}}_{< \varepsilon/2} + \underbrace{\|K^*f - u\|_{\mathcal{U}}}_{< \varepsilon/2} < \varepsilon \end{aligned}$$

which shows that  $u \in \overline{\mathcal{R}(K^*K)}$  and thus also  $\overline{\mathcal{R}(K^*)} \subseteq \overline{\mathcal{R}(K^*K)}$ .

To show (2.7), observe that we have an orthonormal basis  $\{u_j\}_{j \in \mathbb{N}}$  of  $\overline{\mathcal{R}(K^*)}$ , which we can extend to an orthonormal basis  $V$  of  $\mathcal{U}$ . Since  $\mathcal{U} = \mathcal{N}(K) \oplus \mathcal{N}(K)^\perp$  and  $\mathcal{N}(K)^\perp = \overline{\mathcal{R}(K^*)}$  we need to consider elements from  $\mathcal{N}(K)$  for the extension.

Then,

$$\begin{aligned} Ku &= \sum_{v \in V} \langle u, v \rangle_{\mathcal{U}} K v = \sum_{j \in \mathbb{N}} \langle u, u_j \rangle_{\mathcal{U}} K u_j = \sum_{j \in \mathbb{N}} \langle u, u_j \rangle_{\mathcal{U}} \sigma_j v_j \\ &= \sum_{j \in \mathbb{N}} \langle u, K^*v_j \rangle_{\mathcal{U}} v_j = \sum_{j \in \mathbb{N}} \langle K u, v_j \rangle_{\mathcal{V}} v_j \end{aligned}$$

The first line shows (2.7) and the second line shows that  $\{v_j\}_{j \in \mathbb{N}}$  is an orthonormal basis of  $\overline{\mathcal{R}(K)}$ . □

**Remark 2.8.** Since Eigenvalues of  $K^*K$  with Eigenvectors  $u_j$  are also Eigenvalues of  $KK^*$  with Eigenvectors  $v_j$ , we further obtain a singular value decomposition of  $K^*$ , i.e.

$$K^*z = \sum_{j=1}^{\infty} \sigma_j \langle z, v_j \rangle_{\mathcal{V}} u_j.$$

A singular system allows us to characterize elements in the range of the operator.

**Theorem 2.9.** *Let  $K \in \mathcal{K}(\mathcal{U}, \mathcal{V})$  with singular system  $\{(\sigma_j, u_j, v_j)\}_{j \in \mathbb{N}}$ , and  $f \in \overline{\mathcal{R}(K)}$ . Then  $f \in \mathcal{R}(K)$  if and only if the Picard criterion*

$$\sum_{j=1}^{\infty} \frac{|\langle f, v_j \rangle_{\mathcal{V}}|^2}{\sigma_j^2} < \infty \quad (2.8)$$

is met.

*Proof.* Let  $f \in \mathcal{R}(K)$ , thus there is a  $u \in \mathcal{U}$  such that  $Ku = f$ . It is easy to see that we have

$$\langle f, v_j \rangle_{\mathcal{V}} = \langle Ku, v_j \rangle_{\mathcal{V}} = \langle u, K^*v_j \rangle_{\mathcal{U}} = \sigma_j \langle u, u_j \rangle_{\mathcal{U}}$$

and therefore

$$\sum_{j=1}^{\infty} \sigma_j^{-2} |\langle f, v_j \rangle_{\mathcal{V}}|^2 = \sum_{j=1}^{\infty} |\langle u, u_j \rangle_{\mathcal{U}}|^2 \leq \|u\|_{\mathcal{U}}^2 < \infty.$$

Now let the Picard criterion (2.8) hold and define  $u := \sum_{j=1}^{\infty} \sigma_j^{-1} \langle f, v_j \rangle_{\mathcal{V}} u_j \in \mathcal{U}$ . It is well-defined by the Picard criterion (2.8) and we conclude

$$Ku = \sum_{j=1}^{\infty} \sigma_j^{-1} \langle f, v_j \rangle_{\mathcal{V}} K u_j = \sum_{j=1}^{\infty} \langle f, v_j \rangle_{\mathcal{V}} v_j = P_{\overline{\mathcal{R}(K)}} f = f,$$

which shows  $f \in \mathcal{R}(K)$ . □

**Remark 2.9.** The Picard criterion is a condition on the decay of the coefficients  $\langle f, v_j \rangle_{\mathcal{V}}$ . As the singular values  $\sigma_j$  decay to zero as  $j \rightarrow \infty$ , the Picard criterion is only met if the coefficients  $\langle f, v_j \rangle_{\mathcal{V}}$  decay sufficiently fast.

In case the singular system is given by the Fourier basis, then the coefficients  $\langle f, v_j \rangle_{\mathcal{V}}$  are just the Fourier coefficients of  $f$ . Therefore, the Picard criterion is a condition on the decay of the Fourier coefficients which is equivalent to the smoothness of  $f$ .

We can now derive a representation of the Moore-Penrose inverse in terms of the singular value decomposition.

**Theorem 2.10.** *Let  $K \in \mathcal{K}(\mathcal{U}, \mathcal{V})$  with singular system  $\{(\sigma_j, v_j, u_j)\}_{j \in \mathbb{N}}$  and  $f \in \mathcal{D}(K^\dagger)$ . Then the Moore-Penrose inverse of  $K$  can be written as*

$$K^\dagger f = \sum_{j=1}^{\infty} \sigma_j^{-1} \langle f, v_j \rangle_{\mathcal{V}} u_j. \quad (2.9)$$



*Proof.* As  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$  there exist  $u \in \mathcal{N}(K)^\perp$  and  $g \in \mathcal{R}(K)^\perp$  such that  $f = Ku + g$ . As  $\{u_j\}_{j \in \mathbb{N}}$  is an orthonormal basis of  $\mathcal{N}(K)^\perp$  we have that

$$\begin{aligned} u &= \sum_{j=1}^{\infty} \langle u, u_j \rangle_{\mathcal{U}} u_j = \sum_{j=1}^{\infty} \sigma_j^{-1} \langle u, \sigma_j u_j \rangle_{\mathcal{U}} u_j = \sum_{j=1}^{\infty} \sigma_j^{-1} \langle u, K^* v_j \rangle_{\mathcal{U}} u_j \\ &= \sum_{j=1}^{\infty} \sigma_j^{-1} \langle Ku, v_j \rangle_{\mathcal{V}} u_j = \sum_{j=1}^{\infty} \sigma_j^{-1} \langle f - g, v_j \rangle_{\mathcal{V}} u_j = \sum_{j=1}^{\infty} \sigma_j^{-1} \langle f, v_j \rangle_{\mathcal{V}} u_j \end{aligned}$$

where we used for the last equality that  $g \in \mathcal{R}(K)^\perp$  and  $v_j \in \overline{\mathcal{R}(K)}$ .

Moreover, in addition to  $u \in \mathcal{N}(K)^\perp$  we have that  $u$  satisfies the normal equation

$$K^*Ku = \sum_{j=1}^{\infty} \sigma_j^2 \sigma_j^{-1} \langle f, v_j \rangle_{\mathcal{V}} u_j = \sum_{j=1}^{\infty} \sigma_j \langle f, v_j \rangle_{\mathcal{V}} u_j = K^*f$$

and is therefore the minimal norm solution to the inverse problem  $Ku = f$  (1.1). With Theorem 2.3 we conclude that  $u = K^\dagger f$ .  $\square$

From representation (2.9) we can see what happens in case of noisy measurements. Assume we are given  $f^\delta = f + \delta v_j$  and denote by  $u^\dagger$  and  $u_\delta^\dagger$  the minimal norm solutions of  $Ku = f$  and  $Ku = f^\delta$ . Then we observe

$$\|u^\dagger - u_\delta^\dagger\|_{\mathcal{U}} = \|K^\dagger f - K^\dagger f^\delta\|_{\mathcal{U}} = \delta \|K^\dagger v_j\|_{\mathcal{U}} = \frac{\delta}{\sigma_j} \rightarrow \infty \text{ for } j \rightarrow \infty.$$

For fixed  $j$  we see that the amplification of the error  $\delta$  depends on how small  $\sigma_j$  is. Hence, the faster the singular values decay, the stronger the amplification of errors. For that reason, one distinguishes between two classes of ill-posed problems:

**Definition 2.5.** *We say that an ill-posed inverse problem (1.1) is severely ill-posed if the singular values decay as  $\sigma_j = \mathcal{O}(\exp(-j))$ , where the “Big-O-notation” means that there exists  $j_0$  and  $c > 0$  such that for all  $j \geq j_0$  there is  $\sigma_j \leq c \exp(-j)$ . We call the ill-posed inverse problem mildly ill-posed if it is not severely ill-posed.*

**Example 2.8.** Let us consider the example of differentiation again, as introduced in Section 1.1.2. The operator  $K: L^2([0, 1]) \rightarrow L^2([0, 1])$  of the inverse problem (1.1) of differentiation is given as

$$(Ku)(y) = \int_0^y u(x) dx = \int_0^1 k(x, y) u(x) dx,$$

with  $k: [0, 1] \times [0, 1] \rightarrow \mathbb{R}$  defined as

$$k(x, y) := \begin{cases} 1 & x \leq y \\ 0 & \text{else} \end{cases}.$$

This is a special case of the integral operators as introduced in Example 2.6 due to its kernel  $k$  being square integrable and thus  $K$  is compact.

In order to compute the singular value decomposition of  $K$  we compute its adjoint  $K^*$  first, which is characterised via

$$\langle Ku, v \rangle_{L^2([0,1])} = \langle u, K^*v \rangle_{L^2([0,1])}.$$

Hence, we obtain

$$\langle Ku, v \rangle_{L^2([0,1])} = \int_0^1 \int_0^1 k(x, y) u(x) dx v(y) dy = \int_0^1 u(x) \int_0^1 k(x, y) v(y) dy dx.$$

Hence, the adjoint operator  $K^*$  is given via

$$(K^*v)(x) = \int_0^1 k(x, y) v(y) dy = \int_x^1 v(y) dy. \quad (2.10)$$

Now we want to compute the Eigenvalues and Eigenvectors of  $K^*K$ , i.e. we look for  $\lambda > 0$  and  $u \in L^2([0, 1])$  with

$$\lambda u(x) = (K^*Ku)(x) = \int_x^1 \int_0^y u(z) dz dy.$$

We immediately observe  $u(1) = 0$  and further

$$\lambda u'(x) = \frac{d}{dx} \int_x^1 \int_0^y u(z) dz dy = - \int_0^x u(z) dz,$$

from which we conclude  $u'(0) = 0$ . Taking the derivative another time thus yields the ordinary differential equation

$$\lambda u''(x) + u(x) = 0,$$

for which solutions are of the form

$$u(x) = c_1 \sin(\sigma^{-1}x) + c_2 \cos(\sigma^{-1}x),$$

with  $\sigma := \sqrt{\lambda}$  and constants  $c_1, c_2$ . In order to satisfy the boundary conditions  $u(1) = c_1 \sin(\sigma^{-1}) + c_2 \cos(\sigma^{-1}) = 0$  and  $u'(0) = c_1 = 0$ , we chose  $c_1 = 0$  and  $\sigma$  such that  $\cos(\sigma^{-1}) = 0$ . Hence, we have

$$\sigma_j = \frac{2}{(2j-1)\pi} \text{ for } j \in \mathbb{N},$$

and by choosing  $c_2 = \sqrt{2}$  we obtain the following normalised representation of  $u_j$ :

$$u_j(x) = \sqrt{2} \cos\left(\left(j - \frac{1}{2}\right) \pi x\right).$$

According to (2.6) we further obtain

$$v_j(x) = \sigma_j^{-1} (Ku_j)(x) = \left(j - \frac{1}{2}\right) \pi \int_0^x \sqrt{2} \cos\left(\left(j - \frac{1}{2}\right) \pi y\right) dy = \sqrt{2} \sin\left(\left(j - \frac{1}{2}\right) \pi x\right),$$

and hence, for  $f \in L^2([0, 1])$  the Picard criterion becomes

$$2 \sum_{j=1}^{\infty} \sigma_j^{-2} \left( \int_0^1 f(x) \sin(\sigma_j^{-1}x) dx \right)^2 < \infty.$$

Thus, the Picard criterion holds if  $f$  is differentiable and  $f' \in L^2([0, 1])$ .

From the decay of the singular values we see that this inverse problem is mildly ill-posed.

## Chapter 3

# Regularisation

We have seen in the previous section that the major source of ill-posedness of inverse problems of the type (1.1) is a fast decay of the singular values of  $K$ . An idea to overcome this issue is to define approximations of  $K^\dagger$  in the following fashion. Consider the family of operators

$$R_\alpha f := \sum_{j=1}^{\infty} g_\alpha(\sigma_j) \langle f, v_j \rangle_{\mathcal{V}} u_j, \quad (3.1)$$

with functions  $g_\alpha : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$  that converge to  $1/\sigma_j$  as  $\alpha$  converges to zero. We are going to see that such an operator  $R_\alpha$  is what is called a *regularisation* (of  $K^\dagger$ ), if  $g_\alpha$  is bounded, i.e.

$$g_\alpha(\sigma) \leq C_\alpha \text{ for all } \sigma \in \mathbb{R}_{>0}. \quad (3.2)$$

In case (3.2) holds true, we immediately observe

$$\|R_\alpha f\|_{\mathcal{U}}^2 = \sum_{j=1}^{\infty} g_\alpha(\sigma_j)^2 |\langle f, v_j \rangle_{\mathcal{V}}|^2 \leq C_\alpha^2 \sum_{j=1}^{\infty} |\langle f, v_j \rangle_{\mathcal{V}}|^2 \leq C_\alpha^2 \|f\|_{\mathcal{V}}^2,$$

which means that  $C_\alpha$  is a bound for the norm of  $R_\alpha$  and thus  $R_\alpha \in \mathcal{L}(\mathcal{V}, \mathcal{U})$ .

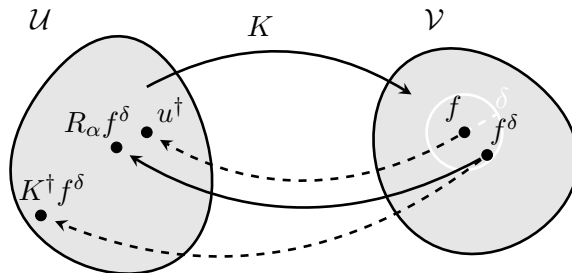
**Example 3.1** (Truncated singular value decomposition). As a first example for a spectral regularisation of the form (3.1) we want to consider the so-called *truncated singular value decomposition*. As the name suggests, the idea is to discard all singular values below a certain threshold  $\alpha$

$$g_\alpha(\sigma) = \begin{cases} \frac{1}{\sigma} & \sigma \geq \alpha \\ 0 & \sigma < \alpha \end{cases}. \quad (3.3)$$

Note that for all  $\sigma > 0$  we naturally obtain  $\lim_{\alpha \rightarrow 0} g_\alpha(\sigma) = 1/\sigma$ . Equation (3.1) then reads as

$$R_\alpha f = \sum_{\sigma_j \geq \alpha} \frac{1}{\sigma_j} \langle f, v_j \rangle_{\mathcal{V}} u_j, \quad (3.4)$$

for all  $f \in \mathcal{V}$ . Note that (3.4) is always well-defined (i.e. finite) for  $\alpha > 0$  as zero is the only accumulation point of singular vectors of compact operators. From (3.3) we immediately observe  $g_\alpha(\sigma) \leq 1/\alpha$  so that  $\|R_\alpha\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} \leq 1/\alpha$ .



**Figure 3.1:** Visualization of reconstruction from noisy data. While the Moore–Penrose inverse reconstructs optimally from noiseless data, its noise amplification renders it useless when small errors are present in the data. A regularisation operator gives a robust solution while still approximating the Moore–Penrose inverse.

**Example 3.2** (Tikhonov regularisation). The main idea behind Tikhonov regularisation<sup>1</sup> is to shift the singular values of  $K^*K$  by a constant factor, which will be associated with the regularisation parameter  $\alpha$ . This shift can be realised via the function

$$g_\alpha(\sigma) = \frac{\sigma}{\sigma^2 + \alpha}. \quad (3.5)$$

Again, we immediately observe that for all  $\sigma > 0$  we have  $\lim_{\alpha \rightarrow 0} g_\alpha(\sigma) = 1/\sigma$ . Further, we can estimate  $g_\alpha(\sigma) \leq 1/(2\sqrt{\alpha})$  due to  $\sigma^2 + \alpha \geq 2\sqrt{\alpha}\sigma$ . The corresponding Tikhonov regularisation (3.1) reads as

$$R_\alpha f = \sum_{j=1}^{\infty} \frac{\sigma_j}{\sigma_j^2 + \alpha} \langle f, v_j \rangle_{\mathcal{V}} u_j. \quad (3.6)$$

After getting an intuition about regularisation of the form (3.1) via examples, we want to define what a regularisation actually is, and what properties come along with it.

**Definition 3.1.** Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  be a bounded operator. A family  $\{R_\alpha\}_{\alpha>0}$  of continuous operators is called regularisation (or regularisation operator) of  $K^\dagger$  if

$$R_\alpha f \rightarrow K^\dagger f = u^\dagger$$

for all  $f \in \mathcal{D}(K^\dagger)$  as  $\alpha \rightarrow 0$ .

**Definition 3.2.** We further call  $\{R_\alpha\}_{\alpha>0}$  a linear regularisation, if Definition 3.1 is satisfied together with the additional assumption

$$R_\alpha \in \mathcal{L}(\mathcal{V}, \mathcal{U}),$$

for all  $\alpha \in \mathbb{R}_{>0}$ .

Hence, a regularisation is a pointwise approximation of the Moore–Penrose inverse with continuous operators, see Figure 3.1 for an illustration. As in the interesting cases the Moore–Penrose inverse may not be continuous we cannot expect that the norms of a regularisation stay bounded as  $\alpha \rightarrow 0$ . This is confirmed by the following results.

<sup>1</sup>Named after the Russian mathematician Andrey Nikolayevich Tikhonov (30 October 1906 - 7 October 1993)

**Theorem 3.1** (Banach–Steinhaus e.g. [4, p. 78], [17, p. 173]). *Let  $\mathcal{U}, \mathcal{V}$  be Hilbert spaces and  $\{K_j\}_{j \in \mathbb{N}} \subset \mathcal{L}(\mathcal{U}, \mathcal{V})$  a family of point-wise bounded operators, i.e. for all  $u \in \mathcal{U}$  there exists a constant  $C(u) > 0$  with  $\sup_{j \in \mathbb{N}} \|K_j u\|_{\mathcal{V}} \leq C(u)$ . Then*

$$\sup_{j \in \mathbb{N}} \|K_j\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} < \infty.$$

**Corollary 3.1** ([17, p. 174]). *Let  $\mathcal{U}, \mathcal{V}$  be Hilbert spaces and  $\{K_j\}_{j \in \mathbb{N}} \subset \mathcal{L}(\mathcal{U}, \mathcal{V})$ . Then the following two conditions are equivalent:*

(a) *There exists  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  such that*

$$Ku = \lim_{j \rightarrow \infty} K_j u \quad \text{for all } u \in \mathcal{U}.$$

(b) *There is a dense subset  $\mathcal{X} \subset \mathcal{U}$  such that  $\lim_{j \rightarrow \infty} K_j u$  exists for all  $u \in \mathcal{X}$  and*

$$\sup_{j \in \mathbb{N}} \|K_j\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} < \infty.$$

**Theorem 3.2.** *Let  $\mathcal{U}, \mathcal{V}$  be Hilbert spaces,  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  and  $\{R_\alpha\}_{\alpha > 0}$  a linear regularisation as defined in Definition 3.2. If  $K^\dagger$  is not continuous,  $\{R_\alpha\}_{\alpha > 0}$  cannot be uniformly bounded. In particular this implies the existence of an element  $f \in \mathcal{V}$  with  $\|R_\alpha f\|_{\mathcal{U}} \rightarrow \infty$  for  $\alpha \rightarrow 0$ .*

*Proof.* We prove the theorem by contradiction and assume that  $\{R_\alpha\}_{\alpha > 0}$  is uniformly bounded. Hence, there exists a constant  $C$  with  $\|R_\alpha\|_{\mathcal{L}(\mathcal{V}, \mathcal{U})} \leq C$  for all  $\alpha > 0$ . Due to Definition 3.1, we have  $R_\alpha \rightarrow K^\dagger$  on  $\mathcal{D}(K^\dagger)$ . Corollary 3.1 then already implies  $K^\dagger \in \mathcal{L}(\mathcal{V}, \mathcal{U})$ , which is a contradiction to the assumption that  $K^\dagger$  is not continuous.

It remains to show the existence of the element  $f \in \mathcal{V}$  with  $\|R_\alpha f\|_{\mathcal{U}} \rightarrow \infty$  for  $\alpha \rightarrow 0$ . If such an element would not exist, we could conclude  $\{R_\alpha\}_{\alpha > 0} \subset \mathcal{L}(\mathcal{V}, \mathcal{U})$ . However, Theorem 3.1 then implies that  $\{R_\alpha\}_{\alpha > 0}$  has to be uniformly bounded, which contradicts the first part of the proof.  $\square$

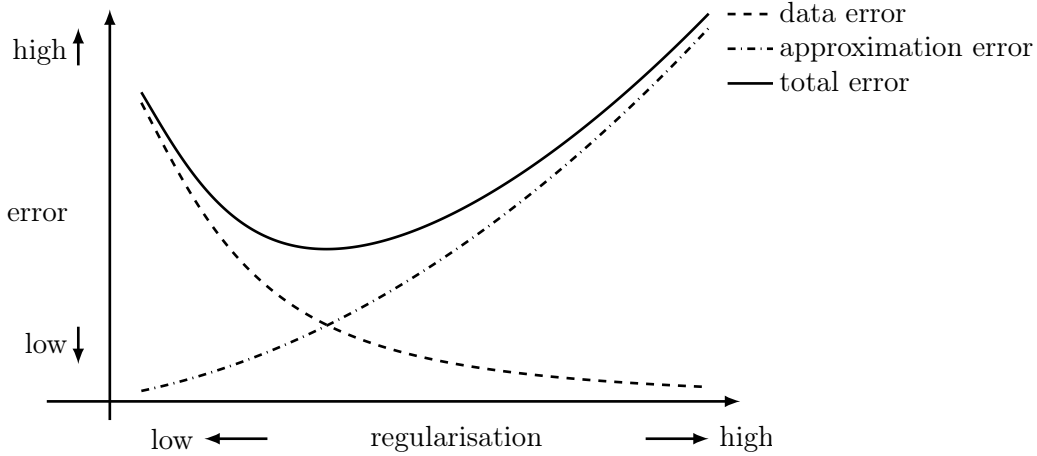
With the additional assumption that  $\|KR_\alpha\|_{\mathcal{L}(\mathcal{V}, \mathcal{V})}$  is bounded, we can even show that  $R_\alpha f$  diverges for all  $f \notin \mathcal{D}(K^\dagger)$ .

**Theorem 3.3.** *Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  and  $\{R_\alpha\}_{\alpha > 0}$  be a linear regularisation of  $K^\dagger$ , and define  $u_\alpha := R_\alpha f$ . If*

$$\sup_{\alpha > 0} \|KR_\alpha\|_{\mathcal{L}(\mathcal{V}, \mathcal{V})} < \infty,$$

*then  $\|u_\alpha\|_{\mathcal{U}} \rightarrow \infty$  for  $f \notin \mathcal{D}(K^\dagger)$ .*

*Proof.* The convergence in case of  $f \in \mathcal{D}(K^\dagger)$  simply follows from Definition 3.1. We therefore only need to consider the case  $f \notin \mathcal{D}(K^\dagger)$ . We assume that there exists a sequence  $\alpha_k \rightarrow 0$  such that  $\|u_{\alpha_k}\|_{\mathcal{U}}$  is uniformly bounded. Then there exists a weakly convergent subsequence  $u_{\alpha_{k_l}}$  with some limit  $u \in \mathcal{U}$ , cf. [9, Section 2.2, Theorem 2.1]. As continuous linear operators are also weakly continuous, we further have  $Ku_{\alpha_{k_l}} \rightharpoonup Ku$ . However, as  $KR_\alpha$  are uniformly bounded operators, we also conclude  $Ku_{\alpha_{k_l}} = KR_{\alpha_{k_l}} f \rightharpoonup P_{\overline{\mathcal{R}(K)}} f$  for all  $f \in \mathcal{V}$  (and not just  $f \in \mathcal{D}(K^\dagger)$ ), because of Corollary 3.1. Hence, we further conclude  $f \in \mathcal{R}(K)$  and therefore  $f \in \mathcal{D}(K^\dagger)$  in contradiction to the assumption  $f \notin \mathcal{D}(K^\dagger)$ .  $\square$



**Figure 3.2:** The *total error* between a regularised solution and the minimal norm solution decomposes into the *data error* and the *approximation error*. These two errors have opposing trends: For a small regularisation parameter  $\alpha$  the error in the data gets amplified through the ill-posedness of the problem and for large  $\alpha$  the operator  $R_\alpha$  is a poor approximation of the Moore–Penrose inverse.

Usually we cannot expect  $f \in \mathcal{D}(K^\dagger)$  for most applications, due to measurement and modelling errors. However, we assume that there exists  $f \in \mathcal{D}(K^\dagger)$  such that we have

$$\|f - f^\delta\|_{\mathcal{V}} \leq \delta$$

for measured data  $f^\delta \in \mathcal{V}$ . For linear regularisations we can split the *total error* between the regularised solution of the noisy problem  $R_\alpha f^\delta$  and the minimal norm solution of the noise-free problem  $u^\dagger = K^\dagger f$  as

$$\begin{aligned} \|R_\alpha f^\delta - u^\dagger\|_{\mathcal{U}} &\leq \|R_\alpha f^\delta - R_\alpha f\|_{\mathcal{U}} + \|R_\alpha f - u^\dagger\|_{\mathcal{U}} \\ &\leq \underbrace{\delta \|R_\alpha\|_{\mathcal{L}(\mathcal{V}, \mathcal{U})}}_{\text{data error}} + \underbrace{\|R_\alpha f - K^\dagger f\|_{\mathcal{U}}}_{\text{approximation error}}. \end{aligned} \quad (3.7)$$

The first term of (3.7) is the *data error*; this term unfortunately does not stay bounded for  $\alpha \rightarrow 0$ , which we can conclude from Theorem 3.2. The second term, known as the *approximation error*, however vanishes for  $\alpha \rightarrow 0$ , due to the pointwise convergence of  $R_\alpha$  to  $K^\dagger$ . Hence it becomes evident from (3.7) that a good choice of  $\alpha$  depends on  $\delta$ , and needs to be chosen such that the approximation error becomes as small as possible, whilst the data error is being kept at bay. See Figure 3.2 for a visualisation of this situation. In the following we are going to discuss typical strategies for choosing  $\alpha$  appropriately.

### 3.1 Parameter-choice strategies

In this section we want to discuss three standard rules for the choice of the regularisation parameter  $\alpha$  and whether they lead to (convergent) regularisation methods.

**Definition 3.3.** A function  $\alpha : \mathbb{R}_{>0} \times \mathcal{V} \rightarrow \mathbb{R}_{>0}$ ,  $(\delta, f^\delta) \mapsto \alpha(\delta, f^\delta)$  is called parameter choice rule. We distinguish between

- (a) a-priori parameter choice rules, if they depend on  $\delta$  only;
- (b) a-posteriori parameter choice rules, if they depend on  $\delta$  and  $f^\delta$ ;
- (c) heuristic parameter choice rules, if they depend on  $f^\delta$  only.

In case of (a) or (c) we would simply write  $\alpha(\delta)$ , respectively  $\alpha(f^\delta)$ , instead of  $\alpha(\delta, f^\delta)$ .

**Definition 3.4.** If  $\{R_\alpha\}_{\alpha>0}$  is a regularisation of  $K^\dagger$  and  $\alpha$  is a parameter choice rule, then the pair  $(R_\alpha, \alpha)$  is called convergent regularisation, if for all  $f \in \mathcal{D}(K^\dagger)$  there exists a parameter choice rule  $\alpha : \mathbb{R}_{>0} \times \mathcal{V} \rightarrow \mathbb{R}_{>0}$  such that

$$\limsup_{\delta \rightarrow 0} \left\{ \left\| R_\alpha f^\delta - K^\dagger f \right\|_{\mathcal{U}} \mid f^\delta \in \mathcal{V}, \left\| f - f^\delta \right\|_{\mathcal{V}} \leq \delta \right\} = 0 \quad (3.8)$$

and

$$\limsup_{\delta \rightarrow 0} \left\{ \alpha(\delta, f^\delta) \mid f^\delta \in \mathcal{V}, \left\| f - f^\delta \right\|_{\mathcal{V}} \leq \delta \right\} = 0 \quad (3.9)$$

are guaranteed.

### 3.1.1 A-priori parameter choice rules

First of all we want to discuss a-priori parameter choice rules in more detail. In fact, it can be shown that for every regularisation an a-priori parameter choice rule, and thus, a convergent regularisation, exists.

**Theorem 3.4.** Let  $\{R_\alpha\}_{\alpha>0}$  be a regularisation of  $K^\dagger$ , for  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ . Then there exists an a-priori parameter choice rule, such that  $(R_\alpha, \alpha)$  is a convergent regularisation.

*Proof.* Let  $f \in \mathcal{D}(K^\dagger)$  be arbitrary but fixed. We can find a monotone increasing function  $\gamma : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$  with  $\lim_{\varepsilon \rightarrow 0} \gamma(\varepsilon) = 0$  such that for every  $\varepsilon > 0$  we have

$$\left\| R_{\gamma(\varepsilon)} f - K^\dagger f \right\|_{\mathcal{U}} \leq \frac{\varepsilon}{2},$$

due to the pointwise convergence  $R_\alpha \rightarrow K^\dagger$ .

As the operator  $R_{\gamma(\varepsilon)}$  is continuous for fixed  $\varepsilon$ , there exists  $\rho(\varepsilon) > 0$  with

$$\left\| R_{\gamma(\varepsilon)} g - R_{\gamma(\varepsilon)} f \right\|_{\mathcal{U}} \leq \frac{\varepsilon}{2} \text{ for all } g \in \mathcal{V} \text{ with } \|g - f\|_{\mathcal{V}} \leq \rho(\varepsilon).$$

Without loss of generality we can assume  $\rho$  to be a continuous, strictly monotone increasing function with  $\lim_{\varepsilon \rightarrow 0} \rho(\varepsilon) = 0$ . Then, due to the inverse function theorem there exists a strictly monotone and continuous function  $\rho^{-1}$  on  $\mathcal{R}(\rho)$  with  $\lim_{\delta \rightarrow 0} \rho^{-1}(\delta) = 0$ . We continuously extend  $\rho^{-1}$  on  $\mathbb{R}_{>0}$  and define our a-priori strategy as

$$\alpha : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}, \quad \delta \mapsto \gamma(\rho^{-1}(\delta)).$$

Then  $\lim_{\delta \rightarrow 0} \alpha(\delta) = 0$  follows. Furthermore, there exists  $\delta := \rho(\varepsilon)$  for all  $\varepsilon > 0$ , such that with  $\alpha(\delta) = \gamma(\varepsilon)$

$$\left\| R_{\alpha(\delta)} f^\delta - K^\dagger f \right\|_{\mathcal{U}} \leq \left\| R_{\gamma(\varepsilon)} f^\delta - R_{\gamma(\varepsilon)} f \right\|_{\mathcal{U}} + \left\| R_{\gamma(\varepsilon)} f - K^\dagger f \right\|_{\mathcal{U}} \leq \varepsilon$$

follows for all  $f^\delta \in \mathcal{V}$  with  $\|f - f^\delta\|_{\mathcal{V}} \leq \delta$ . Thus,  $(R_\alpha, \alpha)$  is a convergent regularisation method.  $\square$

For linear regularisations we can characterise a-priori parameter choice strategies that lead to convergent regularisation methods via the following theorem.

**Theorem 3.5.** *Let  $\{R_\alpha\}_{\alpha>0}$  be a linear regularisation, and  $\alpha : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$  an a-priori parameter choice rule. Then  $(R_\alpha, \alpha)$  is a convergent regularisation method if and only if*

$$(a) \lim_{\delta \rightarrow 0} \alpha(\delta) = 0$$

$$(b) \lim_{\delta \rightarrow 0} \delta \|R_{\alpha(\delta)}\|_{\mathcal{L}(\mathcal{V}, \mathcal{U})} = 0$$

*Proof.*  $\Leftarrow$ : Let condition a) and b) be fulfilled. From (3.7) we then observe

$$\left\| R_{\alpha(\delta)} f^\delta - K^\dagger f \right\|_{\mathcal{U}} \rightarrow 0 \text{ for } \delta \rightarrow 0.$$

Hence,  $(R_\alpha, \alpha)$  is a convergent regularisation method.

$\Rightarrow$ : Now let  $(R_\alpha, \alpha)$  be a convergent regularisation method. We prove that conditions 1 and 2 have to follow from this by showing that violation of either one of them leads to a contradiction to  $(R_\alpha, \alpha)$  being a convergent regularisation method. If condition a) is violated, (3.9) is violated and hence,  $(R_\alpha, \alpha)$  is not a convergent regularisation method. If condition a) is fulfilled but condition b) is violated, there exists a null sequence  $\{\delta_k\}_{k \in \mathbb{N}}$  with  $\delta_k \|R_{\alpha(\delta_k)}\|_{\mathcal{L}(\mathcal{V}, \mathcal{U})} \geq C > 0$ , and hence, we can find a sequence  $\{g_k\}_{k \in \mathbb{N}} \subset \mathcal{V}$  with  $\|g_k\|_{\mathcal{V}} = 1$  and  $\delta_k \|R_{\alpha(\delta_k)} g_k\|_{\mathcal{U}} \geq \tilde{C}$  for some  $\tilde{C}$ . Let  $f \in \mathcal{D}(K^\dagger)$  be arbitrary and define  $f_k := f + \delta_k g_k$ . Then we have on the one hand  $\|f - f_k\|_{\mathcal{V}} \leq \delta_k$ , but on the other hand the norm of

$$R_{\alpha(\delta_k)} f_k - K^\dagger f = R_{\alpha(\delta_k)} f - K^\dagger f + \delta_k R_{\alpha(\delta_k)} g_k$$

cannot converge to zero, as the second term  $\delta_k R_{\alpha(\delta_k)} g_k$  is bounded from below by construction. Hence, (3.8) is violated for  $f^\delta = g_k$  and thus,  $(R_\alpha, \alpha)$  is not a convergent regularisation method.  $\square$

### 3.1.2 A-posteriori parameter choice rules

In the following sections we are going to see that Theorem 3.5 basically means that  $\alpha(\delta)$  cannot converge too quickly to zero in relation to  $\delta$ ; typical parameter choice strategies will be of the form  $\alpha(\delta) = \delta^p$ . However, finding an optimal choice of  $p$  often requires additional information about  $u^\dagger$ , for instance in terms of source conditions that we are going to discuss in Section 3.2.4. A-posteriori parameter choice rules have the advantage that they do not require this additional information. The basic idea is as follows. We again have  $f \in \mathcal{D}(K^\dagger)$  and  $f^\delta$  with  $\|f - f^\delta\|_{\mathcal{V}} \leq \delta$ , and now consider the *residual* between  $f^\delta$  and  $u_\alpha := R_\alpha f^\delta$ , i.e.

$$\|K u_\alpha - f^\delta\|_{\mathcal{V}}.$$

If we assume that  $u^\dagger$  is the minimal norm solution and  $f$  is given via  $f = K u^\dagger$ , we immediately observe that  $u^\dagger$  satisfies

$$\|K u^\dagger - f^\delta\|_{\mathcal{V}} = \|f - f^\delta\|_{\mathcal{V}} = \delta.$$

Hence, it appears not to be useful to choose  $\alpha(\delta, f^\delta)$  with  $\|K u_{\alpha(\delta, f^\delta)} - f^\delta\|_{\mathcal{V}} < \delta$ , which motivates Morozov's discrepancy principle.



**Definition 3.5** (Morozov's discrepancy principle). *Let  $\alpha(\delta, f^\delta)$  be chosen such that*

$$\|Ku_{\alpha(\delta, f^\delta)} - f^\delta\|_{\mathcal{V}} \leq \eta\delta \quad (3.10)$$

*is satisfied, for given  $\delta, f^\delta$ , and a fixed constant  $\eta > 1$ . Then  $u_{\alpha(\delta, f^\delta)} = R_{\alpha(\delta, f^\delta)}f^\delta$  is said to satisfy Morozov's discrepancy principle.*

**Remark 3.1.** It is important to point out that (3.10) may never be fulfilled, as is the case for  $f \in \mathcal{R}(K)^\perp$ . Following Lemma 2.3 (d), even for exact data  $f^\delta = f$  we observe

$$\|Ku^\dagger - f\|_{\mathcal{V}} = \|KK^\dagger f - f\|_{\mathcal{V}} = \|P_{\overline{\mathcal{R}(K)}}f - f\|_{\mathcal{V}} = \|f\|_{\mathcal{V}} > \delta$$

in this case, for  $\delta$  being small enough. In order to avoid this scenario, we ideally ensure that  $\mathcal{R}(K)$  is dense in  $\mathcal{V}$ , as this already implies  $\mathcal{R}(K)^\perp = \{0\}$  due to Remark 2.1.

Practical a-posteriori regularisation strategies are usually designed as follows. We pick a null sequence  $\{\alpha_j\}_{j \in \mathbb{N}}$  and iteratively compute  $u_{\alpha_j} = R_{\alpha_j}f^\delta$  for  $j \in \{1, \dots, j^*\}$ ,  $j^* \in \mathbb{N}$ , until  $u_{\alpha_{j^*}}$  satisfies (3.10). This procedure is justified by the following theorem.

**Theorem 3.6.** *Let  $\{R_\alpha\}_{\alpha > 0}$  be a regularisation of  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ , and let  $\mathcal{R}(K)$  be dense in  $\mathcal{V}$ . Further, let  $\{\alpha_j\}_{j \in \mathbb{N}}$  be a strictly monotonically decreasing null sequence, and let  $\eta > 1$ . If the family of operators  $\{KR_\alpha\}_{\alpha > 0}$  is uniformly bounded, there exists a finite index  $j^* \in \mathbb{N}$  such that for all  $f \in \mathcal{D}(K^\dagger)$  and  $f^\delta$  with  $\|f - f^\delta\|_{\mathcal{V}} \leq \delta$  the inequalities*

$$\|Ku_{\alpha_{j^*}} - f^\delta\|_{\mathcal{V}} \leq \eta\delta < \|Ku_{\alpha_j} - f^\delta\|_{\mathcal{V}}$$

*are satisfied for all  $j < j^*$ .*

*Proof.* We know that  $KR_\alpha$  converges pointwise to  $KK^\dagger = P_{\overline{\mathcal{R}(K)}}$  in  $\mathcal{D}(K^\dagger)$ , which together with the uniform boundedness assumption already implies pointwise convergence in  $\mathcal{V}$ , as we have already shown in the proof of Theorem 3.2. Hence, for all  $f \in \mathcal{D}(K^\dagger)$  and  $f^\delta \in \mathcal{V}$  with  $\|f - f^\delta\|_{\mathcal{V}} \leq \delta$  we can conclude

$$\begin{aligned} \lim_{j \rightarrow \infty} \|Ku_{\alpha_j} - f^\delta\|_{\mathcal{V}} &= \lim_{j \rightarrow \infty} \|KR_{\alpha_j}f^\delta - f^\delta\|_{\mathcal{V}} = \left\| P_{\overline{\mathcal{R}(K)}}f^\delta - f^\delta \right\|_{\mathcal{V}} \\ &= \inf_{g \in \overline{\mathcal{R}(K)}} \|g - f^\delta\|_{\mathcal{V}} \leq \|f - f^\delta\|_{\mathcal{V}} \leq \delta. \end{aligned}$$

□

We are going to demonstrate later that (3.10) in combination with specific regularisations is indeed a regularisation method. Before we do so, we want to conclude the discussion of parameter choice strategies by investigating heuristic regularisation methods.

### 3.1.3 Heuristic parameter choice rules

Heuristic parameter choice rules do not require knowledge of the noise level  $\delta$ , which makes them popular strategies in practice. In the following we give three examples of popular heuristic parameter choice rules.

**Quasi-optimality principle** For the first  $n$  elements of a null sequence, i.e.  $\{\alpha_j\}_{j \in \{1, \dots, n\}}$ , we choose  $\alpha(f^\delta) = \alpha_{j^*}$  with

$$j^* = \arg \min_{1 \leq j < n} \|u_{\alpha_{j+1}} - u_{\alpha_j}\|_{\mathcal{U}}.$$

**Hanke-Raus rule** The parameter  $\alpha(f^\delta)$  is chosen via

$$\alpha(f^\delta) = \arg \min_{\alpha > 0} \frac{1}{\sqrt{\alpha}} \|Ku_\alpha - f^\delta\|_{\mathcal{V}}.$$

**L-curve method** The parameter  $\alpha(f^\delta)$  is chosen via

$$\alpha(f^\delta) = \arg \min_{\alpha > 0} \|u_\alpha\|_{\mathcal{U}} \|Ku_\alpha - f^\delta\|_{\mathcal{V}}.$$

Despite their popularity and the fact that they do not require any knowledge about  $\delta$ , heuristic parameter choice rules have one significant theoretical disadvantage. While any regularisation can be equipped with an a-priori parameter choice rule to form a convergent regularisation as seen in Theorem 3.4, heuristic parameter choice rules cannot lead to convergent regularisations, a result that has become famous as the so-called Bakushinskiĭ veto [2].

**Theorem 3.7.** *Let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  with  $\mathcal{R}(K) \neq \overline{\mathcal{R}(K)}$ . Then for any regularisation  $\{R_\alpha\}_{\alpha > 0}$  and any heuristic parameter choice rule  $\alpha(f^\delta)$  the pair  $(\{R_\alpha\}, \alpha)$  is not a convergent regularisation.*

*Proof.* Assume that  $(\{R_\alpha\}, \alpha)$  is a convergent regularisation method and that the parameter choice rule is heuristic, i.e.  $\alpha = \alpha(f^\delta)$ . Then it follows from (3.8) that

$$\limsup_{\delta \rightarrow 0} \left\{ \left\| R_{\alpha(f^\delta)} f^\delta - K^\dagger f \right\|_{\mathcal{U}} \mid f^\delta \in \mathcal{V}, \left\| f - f^\delta \right\|_{\mathcal{V}} \leq \delta \right\} = 0$$

and in particular  $R_{\alpha(f)} f = K^\dagger f$  for all  $f \in \mathcal{D}(K^\dagger)$ . Thus, for any sequence  $\{f_j\}_{j \in \mathbb{N}} \subset \mathcal{D}(K^\dagger)$  which converges to  $f \in \mathcal{D}(K^\dagger)$  we have that

$$\lim_{j \rightarrow \infty} K^\dagger f_j = \lim_{j \rightarrow \infty} R_{\alpha(f_j)} f_j = K^\dagger f$$

which shows that  $K^\dagger$  is continuous. It follows from Theorem 2.5 that the range of  $K$  is closed, which contradicts the assumption.  $\square$

**Remark 3.2.** We want to point out that Theorem 3.7 does not automatically make any heuristic parameter choice rule useless, for two reasons. Firstly, because Theorem 3.7 applies to infinite dimensional problems. Hence, discretised, ill-conditioned problems can still benefit from heuristic parameter choice rules. Secondly, the proof of Theorem 3.7 explicitly uses perturbed data  $f_j \in \mathcal{D}(K^\dagger)$  to show the contradiction. For actual perturbed data  $f^\delta$  however, it is quite unusual that they will satisfy  $f^\delta \in \mathcal{D}(K^\dagger)$ . It can indeed be shown that, under the additional assumption  $f^\delta \notin \mathcal{D}(K^\dagger)$ , a lot of regularisation strategies together with a whole class of heuristic parameter choice strategies can be turned into convergent regularisations.

## 3.2 Spectral regularisation methods

Now we revisit (3.1) and finally prove that these methods are regularisation methods for piecewise continuous functions  $g_\alpha$  satisfying (3.2).

**Theorem 3.8.** Let  $g_\alpha : \mathbb{R}_{>0} \rightarrow \mathbb{R}$  be a piecewise continuous function satisfying (3.2),  $\lim_{\alpha \rightarrow 0} g_\alpha(\sigma) = \frac{1}{\sigma}$  and

$$\sup_{\alpha, \sigma} \sigma g_\alpha(\sigma) \leq \gamma, \quad (3.11)$$

for some constant  $\gamma > 0$ . If  $R_\alpha$  is defined as in (3.1), we have

$$R_\alpha f \rightarrow K^\dagger f \text{ as } \alpha \rightarrow 0,$$

for all  $f \in \mathcal{D}(K^\dagger)$ .

*Proof.* From the singular value decomposition of  $K^\dagger$  and the definition of  $R_\alpha$  we obtain

$$R_\alpha f - K^\dagger f = \sum_{j=1}^{\infty} \left( g_\alpha(\sigma_j) - \frac{1}{\sigma_j} \right) \langle f, v_j \rangle_{\mathcal{V}} u_j = \sum_{j=1}^{\infty} (\sigma_j g_\alpha(\sigma_j) - 1) \langle u^\dagger, u_j \rangle_{\mathcal{U}} u_j.$$

From (3.11) we can conclude

$$\left| (\sigma_j g_\alpha(\sigma_j) - 1) \langle u^\dagger, u_j \rangle_{\mathcal{U}} \right| \leq (1 + \gamma) \|u^\dagger\|_{\mathcal{U}},$$

and hence, each element of the sum stays bounded. Thus, we can also estimate

$$\begin{aligned} \|R_\alpha f - K^\dagger f\|_{\mathcal{U}}^2 &= \sum_{j=1}^{\infty} |\sigma_j g_\alpha(\sigma_j) - 1|^2 \left| \langle u^\dagger, u_j \rangle_{\mathcal{U}} \right|^2 \leq (1 + \gamma)^2 \sum_{j=1}^{\infty} \left| \langle u^\dagger, u_j \rangle_{\mathcal{U}} \right|^2 \\ &= (1 + \gamma)^2 \|u^\dagger\|_{\mathcal{U}}^2 < \infty \end{aligned}$$

and conclude that  $\|R_\alpha f - K^\dagger f\|_{\mathcal{U}}$  is bounded from above. This allows the application of the reverse Fatou lemma, which yields the estimate

$$\begin{aligned} \limsup_{\alpha \rightarrow 0} \|R_\alpha f - K^\dagger f\|_{\mathcal{U}}^2 &\leq \limsup_{\alpha \rightarrow 0} \sum_{j=1}^{\infty} |\sigma_j g_\alpha(\sigma_j) - 1|^2 \left| \langle u^\dagger, u_j \rangle_{\mathcal{U}} \right|^2 \\ &\leq \sum_{j=1}^{\infty} \left| \lim_{\alpha \rightarrow 0} \sigma_j g_\alpha(\sigma_j) - 1 \right|^2 \left| \langle u^\dagger, u_j \rangle_{\mathcal{U}} \right|^2. \end{aligned}$$

Due to the pointwise convergence of  $g_\alpha(\sigma_j)$  to  $1/\sigma_j$  we obtain  $\lim_{\alpha \rightarrow 0} \sigma_j g_\alpha(\sigma_j) - 1 = 0$ . Hence, we have  $\|R_\alpha f - K^\dagger f\|_{\mathcal{U}} \rightarrow 0$  for  $\alpha \rightarrow 0$  for all  $f \in \mathcal{D}(K^\dagger)$ .  $\square$

**Proposition 3.1.** Let the same assumptions hold as in Theorem 3.8. Further, let  $\alpha$  be an a-priori parameter choice rule. Then  $(R_{\alpha(\delta)}, \alpha(\delta))$  is a convergent regularisation method if

$$\lim_{\delta \rightarrow 0} \delta C_{\alpha(\delta)} = 0$$

is guaranteed.

*Proof.* The result follows immediately from  $\|R_{\alpha(\delta)}\|_{\mathcal{L}(\mathcal{V}, \mathcal{U})} \leq C_{\alpha(\delta)}$  and Theorem 3.5.  $\square$

### 3.2.1 Convergence rates

Knowing that spectral regularisation methods of the form (3.1) together with (3.2) represent convergent regularisation methods, we now want to understand how the error in the data propagates to the error in the reconstruction.

**Theorem 3.9.** *Let the same assumptions hold for  $g_\alpha$  as in Theorem 3.8. If we define  $u_\alpha := R_\alpha f$  and  $u_\alpha^\delta := R_\alpha f^\delta$ , with  $f \in \mathcal{D}(K^\dagger)$ ,  $f^\delta \in \mathcal{V}$  and  $\|f - f^\delta\|_{\mathcal{V}} \leq \delta$ , then*

$$\|Ku_\alpha - Ku_\alpha^\delta\|_{\mathcal{V}} \leq \gamma\delta, \quad (3.12)$$

and

$$\|u_\alpha - u_\alpha^\delta\|_{\mathcal{U}} \leq C_\alpha\delta \quad (3.13)$$

hold true.

*Proof.* From the singular value decomposition we can estimate

$$\begin{aligned} \|Ku_\alpha - Ku_\alpha^\delta\|_{\mathcal{V}}^2 &\leq \sum_{j=1}^{\infty} \sigma_j^2 g_\alpha(\sigma_j)^2 |\langle f - f^\delta, v_j \rangle_{\mathcal{V}}|^2 \\ &\leq \gamma^2 \sum_{j=1}^{\infty} |\langle f - f^\delta, v_j \rangle_{\mathcal{V}}|^2 = \gamma^2 \|f - f^\delta\|_{\mathcal{V}}^2 \leq \gamma^2 \delta^2, \end{aligned}$$

which yields (3.12). In the same fashion we can estimate

$$\begin{aligned} \|u_\alpha - u_\alpha^\delta\|_{\mathcal{U}}^2 &\leq \sum_{j=1}^{\infty} g_\alpha(\sigma_j)^2 |\langle f - f^\delta, v_j \rangle_{\mathcal{V}}|^2 \\ &\leq C_\alpha^2 \sum_{j=1}^{\infty} |\langle f - f^\delta, v_j \rangle_{\mathcal{V}}|^2 = C_\alpha^2 \|f - f^\delta\|_{\mathcal{V}}^2 \leq C_\alpha^2 \delta^2, \end{aligned}$$

to obtain (3.13).  $\square$

**Remark 3.3.** At first glance (3.13) gives the impression as if the error in the reconstruction is also of order  $\delta$ . This, however, is not the case, as  $C_\alpha$  also depends on  $\delta$ , as we have seen in Proposition 3.1. The condition  $\lim_{\delta \rightarrow 0} \delta C_\alpha = 0$  will in particular force  $C_\alpha$  to decay more quickly than  $\delta$ . Hence,  $C_\alpha\delta$  will be of order  $\delta^\nu$ , with  $0 < \nu < 1$ .

Combining the assertions of Theorem 3.8, Proposition 3.1 and Theorem 3.9, we obtain the following convergence results of the regularised solutions.

**Proposition 3.2.** *Let the assumptions of Theorem 3.8, Proposition 3.1 and Theorem 3.9 hold true. Then,*

$$u_{\alpha(\delta)} \rightarrow u^\dagger$$

is guaranteed as  $\delta \rightarrow 0$ .

### 3.2.2 Truncated singular value decomposition

As a first example for a spectral regularisation of the form (3.1) we have considered the so-called truncated singular value decomposition in Example 3.1. From (3.3) we immediately observe  $g_\alpha(\sigma) \leq C_\alpha = 1/\alpha$ . Thus, according to Proposition 3.1 the truncated singular value decomposition, together with an a-priori parameter choice strategy satisfying  $\lim_{\delta \rightarrow 0} \alpha(\delta) = 0$ , is a convergent regularisation method if  $\lim_{\delta \rightarrow 0} \delta/\alpha(\delta) = 0$ .

Moreover, we observe  $\sup_{\sigma, \alpha} \sigma g_\alpha(\sigma) = \gamma = 1$  and hence, we obtain the error estimates  $\|Ku_\alpha - Ku_\alpha^\delta\|_{\mathcal{V}} \leq \delta$  and  $\|u_\alpha - u_\alpha^\delta\|_{\mathcal{U}} \leq \delta/\alpha(\delta)$  as a consequence of Theorem 3.9.

Let  $K \in \mathcal{K}(\mathcal{U}, \mathcal{V})$  with singular system  $\{\sigma_j, u_j, v_j\}_{j \in \mathbb{N}}$ , and choose for  $\delta > 0$  an index function  $j^* : \mathbb{R}_{>0} \rightarrow \mathbb{N}$  with  $j^*(\delta) \rightarrow \infty$  for  $\delta \rightarrow 0$  and  $\lim_{\delta \rightarrow 0} \delta/\sigma_{j^*(\delta)} = 0$ . We can then choose  $\alpha(\delta) = \sigma_{j^*(\delta)}$  as our a-priori parameter choice rule to obtain a convergent regularisation.

Note that in practice a larger  $\delta$  implies that more and more singular values have to be cut off in order to guarantee a stable recovery that successfully suppresses the data error.

### 3.2.3 Tikhonov regularisation

The second example we were considering was Tikhonov regularisation in Example 3.2, where we have shifted the singular values of  $K^*K$  by a constant factor, which will be associated with the regularisation parameter  $\alpha$ .

In case of  $g_\alpha$  as defined in (3.5) we observe  $\lim_{\alpha \rightarrow 0} g_\alpha(\sigma) = 1/\sigma$  for  $\sigma > 0$ . Further, we can estimate  $g_\alpha(\sigma) \leq 1/(2\sqrt{\alpha}) = C_\alpha$  due to  $\sigma^2 + \alpha \geq 2\sqrt{\alpha}\sigma$ . Moreover, we discover  $\sigma g_\alpha(\sigma) = \sigma^2/(\sigma^2 + \alpha) < 1 =: \gamma$  for  $\alpha > 0$ . Consequently, we have to ensure  $\delta/(2\sqrt{\alpha(\delta)}) \rightarrow 0$  for  $\delta \rightarrow 0$  to obtain a convergent regularisation, and in that case get the estimates  $\|Ku_\alpha - Ku_\alpha^\delta\|_{\mathcal{V}} \leq \delta$  and  $\|u_\alpha - u_\alpha^\delta\|_{\mathcal{U}} \leq \delta/(2\sqrt{\alpha(\delta)})$ . Thus, equipping  $R_{\alpha(\delta)}$  for instance with the a-priori parameter choice rule  $\alpha(\delta) = \delta/4$  will lead to a convergent regularisation for which we have  $\|u_\alpha - u_\alpha^\delta\|_{\mathcal{U}} = \mathcal{O}(\sqrt{\delta})$ .

Note that Tikhonov regularisation can be computed without knowledge of the singular system. Considering the equation  $(K^*K + \alpha I)u_\alpha$  in terms of the singular value decomposition, we observe

$$\begin{aligned} & \sum_{j=1}^{\infty} \frac{\sigma_j}{\sigma_j^2 + \alpha} \langle f, v_j \rangle_{\mathcal{V}} \underbrace{K^* K u_j}_{=\sigma_j^2 u_j} + \sum_{j=1}^{\infty} \frac{\alpha \sigma_j}{\sigma_j^2 + \alpha} \langle f, v_j \rangle_{\mathcal{V}} u_j \\ &= \sum_{j=1}^{\infty} \frac{\sigma_j(\sigma_j^2 + \alpha)}{\sigma_j^2 + \alpha} \langle f, v_j \rangle_{\mathcal{V}} u_j = \sum_{j=1}^{\infty} \sigma_j \langle f, v_j \rangle_{\mathcal{V}} u_j = K^* f. \end{aligned}$$

Hence, the Tikhonov-regularised solution  $u_\alpha$  can be obtained by solving

$$(K^*K + \alpha I)u_\alpha = K^* f \tag{3.14}$$

for  $u_\alpha$ . The advantage in computing  $u_\alpha$  via (3.14) is that its computation does not require the singular value decomposition of  $K$ , but only involves the inversion of a linear, well-posed operator equation with a symmetric, positive definite operator.

### 3.2.4 Source-conditions

Before we continue to investigate other examples of regularisations we want to briefly address the question of the convergence speed of a regularisation method. From Theorem

3.9 we have already obtained a convergence rate result; however, with additional regularity assumptions on the (unknown) minimal norm solution we are able to improve those. The regularity assumptions that we want to consider are known as *source conditions*, and are of the form

$$\exists w \in \mathcal{U} : u^\dagger = (K^*K)^\mu w. \quad (3.15)$$

The power  $\mu > 0$  of the operator is understood in the sense of the consider the  $\mu$ -th power of the singular values of the operator  $K^*K$ , i.e.

$$(K^*K)^\mu w = \sum_{j=1}^{\infty} \sigma_j^{2\mu} \langle w, u_j \rangle_{\mathcal{U}} u_j.$$

**Example 3.3** (Differentiation). We want to take a look at what (3.15) actually means in the case of a specific example. We therefore again consider the inverse problem of differentiation, i.e.

$$(Ku)(y) = \int_0^y u(x) dx.$$

In case of  $\mu = 1$  (3.15) reads as

$$u^\dagger(x) = \int_x^1 \int_0^y w(z) dz dy.$$

due to (2.10). Hence, (3.15) does simply imply that  $u^\dagger$  has to be twice weakly differentiable. It becomes even more obvious if we look at twice differentiable  $u^\dagger$ . In that case applying the Leibniz differentiation rule for parameter integrals leaves us with

$$(u^\dagger)''(x) = -w(x).$$

Hence, any twice differentiable  $u^\dagger$  automatically satisfies the source condition (3.15) for  $\mu = 1$ .

Similar results follow for different choices of  $\mu \in \mathbb{N}$ .

The rate of convergence of a regularisation scheme to the minimal norm solution now depends on the specific choice of  $g_\alpha$ . We assume that  $g_\alpha$  satisfies

$$\sigma^{2\mu} |\sigma g_\alpha(\sigma) - 1| \leq \omega_\mu(\alpha),$$

for all  $\sigma > 0$ . In case of the truncated singular value decomposition we would for instance have  $\omega_\mu(\alpha) = \alpha^{2\mu}$ . With this additional assumption, we can improve the estimate in Theorem 3.8 as follows:

$$\begin{aligned} \|R_\alpha f - K^\dagger f\|_{\mathcal{V}}^2 &\leq \sum_{j=1}^{\infty} |\sigma_j g_\alpha(\sigma_j) - 1|^2 |\langle u^\dagger, u_j \rangle_{\mathcal{U}}|^2 \\ &= \sum_{j=1}^{\infty} |\sigma_j g_\alpha(\sigma_j) - 1|^2 \sigma_j^{4\mu} |\langle w, u_j \rangle_{\mathcal{U}}|^2 \\ &\leq \omega_\mu(\alpha)^2 \|w\|_{\mathcal{U}}^2 \end{aligned}$$

Hence, we have obtained the estimate

$$\|u_\alpha - u^\dagger\|_{\mathcal{U}} \leq \omega_\mu(\alpha) \|w\|_{\mathcal{U}}.$$

Together with (3.7) we can further estimate

$$\|u_{\alpha(\delta)} - u^\dagger\|_{\mathcal{U}} \leq \omega_\mu(\alpha) \|w\|_{\mathcal{U}} + C_\alpha \delta. \quad (3.16)$$

**Example 3.4.** In case of the truncated singular value decomposition we know from Section 3.2.2 that  $C_\alpha = 1/\alpha$ , and we can further conclude  $\omega_\mu(\alpha) = \alpha^{2\mu}$ . Hence, (3.16) simplifies to

$$\|u_{\alpha(\delta)} - u^\dagger\|_{\mathcal{U}} \leq \alpha^{2\mu} \|w\|_{\mathcal{U}} + \delta \alpha^{-1} \quad (3.17)$$

in this case. In order to make the right-hand-side of (3.17) as small as possible, we have to choose  $\alpha$  such that

$$\alpha = \left( \frac{\delta}{2\mu \|w\|_{\mathcal{U}}} \right)^{\frac{1}{2\mu+1}}.$$

With this choice of  $\alpha$  we estimate

$$\begin{aligned} \|u_{\alpha(\delta)} - u^\dagger\|_{\mathcal{U}} &\leq \underbrace{2^{\frac{1-2\mu}{1+2\mu}}}_{\leq 2} \underbrace{\mu^{\frac{1-2\mu}{1+2\mu}}}_{\leq 1} \delta^{\frac{2\mu}{2\mu+1}} \|w\|_{\mathcal{U}}^{\frac{1}{2\mu+1}} \\ &\leq 2\delta^{\frac{2\mu}{2\mu+1}} \|w\|_{\mathcal{U}}^{\frac{1}{2\mu+1}}. \end{aligned}$$

It is important to note that no matter how large  $\mu$  is, the rate of convergence  $\delta^{\frac{2\mu}{2\mu+1}}$  will always be slower than  $\delta$ , due to the ill-posedness of the inversion of  $K$ .

### 3.2.5 Asymptotic regularisation

Another form of regularisation is asymptotic regularisation of the form

$$\begin{aligned} \partial_t u(t) &= K^*(f - Ku(t)) \\ u(0) &= 0 \end{aligned} \quad (3.18)$$

As the linear operator  $K$  does not change with respect to the time  $t$ , we can make the Ansatz of writing  $u(t)$  in terms of the singular value decomposition of  $K$  as

$$u(t) = \sum_{j=1}^{\infty} \gamma_j(t) u_j, \quad (3.19)$$

for some function  $\gamma : \mathbb{R} \rightarrow \mathbb{R}$ . From the initial conditions we immediately observe  $\gamma(0) = 0$ . From the singular value decomposition and (3.18) we further see

$$\sum_{j=1}^{\infty} \gamma_j'(t) u_j = \sum_{j=1}^{\infty} \sigma_j \left( \langle f, v_j \rangle_{\mathcal{V}} - \underbrace{\sigma_j \gamma(t) \langle u_j, u_j \rangle_{\mathcal{U}}}_{=\|u_j\|_{\mathcal{U}}^2} \right) u_j.$$

Hence, by equating the coefficients we get

$$\gamma_j'(t) = \sigma_j \langle f, v_j \rangle_{\mathcal{V}} - \sigma_j^2 \gamma_j(t),$$

and together with  $\gamma_j(0)$  we obtain

$$\gamma_j(t) = \left(1 - e^{-\sigma_j^2 t}\right) \frac{1}{\sigma_j} \langle f, v_j \rangle_{\mathcal{V}}$$

as a solution for all  $j$  and hence, (3.19) reads as

$$u(t) = \sum_{j=1}^{\infty} \left(1 - e^{-\sigma_j^2 t}\right) \frac{1}{\sigma_j} \langle f, v_j \rangle_{\mathcal{V}} u_j.$$

If we substitute  $t = 1/\alpha$ , we obtain the regularisation

$$u_{\alpha} = \sum_{j=1}^{\infty} \left(1 - e^{-\frac{\sigma_j^2}{\alpha}}\right) \frac{1}{\sigma_j} \langle f, v_j \rangle_{\mathcal{V}} u_j$$

with  $g_{\alpha}(\sigma) = \left(1 - e^{-\frac{\sigma^2}{\alpha}}\right) \frac{1}{\sigma}$ . We immediately see that  $g_{\alpha}(\sigma)\sigma \leq 1 =: \gamma$ , and due to  $e^x \geq 1 + x$  we further observe  $1 - e^{-\frac{\sigma^2}{\alpha}} \leq \sigma^2/\alpha$  and therefore  $(1 - e^{-\frac{\sigma^2}{\alpha}})/\sigma \leq \max_j \sigma_j/\alpha = \sigma_1/\alpha = \|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})}/\alpha =: C_{\alpha}$ .

### 3.2.6 Landweber iteration

If we approximate (3.18) via a forward finite-difference discretisation, we end up with the iterative procedure

$$\begin{aligned} \frac{u^{k+1} - u^k}{\tau} &= K^* \left( f - K u^k \right), & (3.20) \\ \Leftrightarrow u^{k+1} &= u^k + \tau K^* \left( f - K u^k \right), \\ \Leftrightarrow u^{k+1} &= (I - \tau K^* K) u^k + \tau K^* f, \end{aligned}$$

for some  $\tau > 0$  and  $u^0 \equiv 0$ . Iteration (3.20) is known as the so-called Landweber iteration. We assume  $f \in \mathcal{D}(K^{\dagger})$  first, and with the singular value decomposition of  $K$  and  $K^*$  we obtain

$$\sum_{j=1}^{\infty} \langle u^{k+1}, u_j \rangle_{\mathcal{U}} u_j = \sum_{j=1}^{\infty} \left( (1 - \tau \sigma_j^2) \langle u^k, u_j \rangle_{\mathcal{U}} + \tau \sigma_j \langle f, v_j \rangle_{\mathcal{V}} \right) u_j, \quad (3.21)$$

and hence, by equating the individual summands

$$\langle u^{k+1}, u_j \rangle_{\mathcal{U}} = (1 - \tau \sigma_j^2) \langle u^k, u_j \rangle_{\mathcal{U}} + \tau \sigma_j \langle f, v_j \rangle_{\mathcal{V}}. \quad (3.22)$$

Assuming  $u^0 \equiv 0$ , summing up equation (3.22) yields

$$\langle u^k, u_j \rangle_{\mathcal{U}} = \tau \sigma_j \langle f, v_j \rangle_{\mathcal{V}} \sum_{i=1}^k (1 - \tau \sigma_j^2)^{k-i}. \quad (3.23)$$

The following Lemma will help us simplifying (3.23).

**Lemma 3.1.** *For  $k \in \mathbb{N} \setminus \{1\}$  we have*

$$\sum_{i=1}^k (1 - \tau \sigma^2)^{k-i} = \frac{1 - (1 - \tau \sigma^2)^k}{\tau \sigma^2}. \quad (3.24)$$



*Proof.* Equation (3.24) can simply be verified via induction. We immediately see that

$$\sum_{i=1}^2 (1 - \tau\sigma^2)^{2-i} = 1 + (1 - \tau\sigma^2) = \frac{1 - (1 - 2\tau\sigma^2 + \tau^2\sigma^4)}{\tau\sigma^2} = \frac{1 - (1 - \tau\sigma^2)^2}{\tau\sigma^2}$$

serves as our induction base. Considering  $k \rightarrow k + 1$ , we observe

$$\begin{aligned} \sum_{i=1}^{k+1} (1 - \tau\sigma^2)^{k+1-i} &= 1 + \sum_{i=1}^k (1 - \tau\sigma^2)^{k+1-i} \\ &= 1 + (1 - \tau\sigma^2) \sum_{i=1}^k (1 - \tau\sigma^2)^{k-i} \\ &= 1 + (1 - \tau\sigma^2) \frac{1 - (1 - \tau\sigma^2)^k}{\tau\sigma^2} \\ &= \frac{1 - (1 - \tau\sigma^2)^{k+1}}{\tau\sigma^2}, \end{aligned}$$

and we are done.  $\square$

If we now insert (3.24) into (3.23) we therefore obtain

$$\langle u^k, u_j \rangle_{\mathcal{U}} = \left(1 - (1 - \tau\sigma_j^2)^k\right) \frac{1}{\sigma_j} \langle f, v_j \rangle_{\mathcal{V}}. \quad (3.25)$$

The important consequence of Equation (3.25) is that we now immediately see that  $\langle u^k, u_j \rangle_{\mathcal{U}} \rightarrow \langle u^\dagger, u_j \rangle_{\mathcal{U}}$  if we ensure  $(1 - \tau\sigma_j^2)^k \rightarrow 0$ . In other words, we need to choose  $\tau$  such that  $|1 - \tau\sigma_j^2| < 1$  (respectively  $0 < \tau\sigma_j < 2$ ) for all  $j$ . As in the case of asymptotic regularisation we exploit that  $\sigma_1 = \|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} > \sigma_j$  for all  $j$  and select  $\tau$  such that

$$0 < \tau < \frac{2}{\|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})}^2} \quad (3.26)$$

is satisfied. If we interpret the iteration number as the regularisation parameter  $\alpha := 1/k$ , we obtain the regularisation method

$$u_\alpha = R_\alpha f = \sum_{j=1}^{\infty} \left(1 - (1 - \tau\sigma_j^2)^\alpha\right) \frac{1}{\sigma_j} \langle f, v_j \rangle_{\mathcal{V}}$$

with  $g_\alpha(\sigma) = \left(1 - (1 - \tau\sigma^2)^\alpha\right) / \sigma$ .

### Landweber Iteration & the discrepancy principle

To conclude this section on the Landweber iteration we want to prove convergence rates given  $u^\dagger$  satisfies a source condition. We further want to demonstrate that Landweber iteration in combination with the a-posteriori parameter choice rule defined in Definition 3.5 is a sensible strategy that ensures  $u^k \rightarrow u^\dagger$  as long as the discrepancy principle is violated. Following the introduction of the source condition in Section 3.2.4, we want to assume a source condition similar (3.15) for  $\mu = 1/2$ , i.e. there exists a  $w \in \mathcal{V}$  such that

$$u^\dagger = K^* w \quad (3.27)$$

is satisfied. Under that additional assumption we can conclude the following convergence rate in the case of noise-free data  $f^\delta = f$ .

**Lemma 3.2.** *Let (3.27) be satisfied. Then the Landweber iterates (3.20) satisfy*

$$\|u^k - u^\dagger\|_{\mathcal{U}} = \mathcal{O}\left(\frac{1}{\sqrt{k}}\right) = \mathcal{O}(\sqrt{\alpha}),$$

for  $f = Ku^\dagger$ .

*Proof.* We start proving this statement by showing that the inner product of  $u^k - u^\dagger$  with a singular vector  $u_j$  simplifies to

$$\begin{aligned} \langle u^k - u^\dagger, u_j \rangle_{\mathcal{U}} &= \langle u^k, u_j \rangle_{\mathcal{U}} - \langle u^\dagger, u_j \rangle_{\mathcal{U}} \\ &= \left(1 - (1 - \tau\sigma_j^2)^k\right) \langle u^\dagger, u_j \rangle_{\mathcal{U}} - \langle u^\dagger, u_j \rangle_{\mathcal{U}} \\ &= (1 - \tau\sigma_j^2)^k \langle u^\dagger, u_j \rangle_{\mathcal{U}} \\ &= \underbrace{\sigma_j (1 - \tau\sigma_j^2)^k}_{=: r(\sigma_j)} \langle w, u_j \rangle_{\mathcal{U}}, \end{aligned}$$

with the second equality following from Equation (3.25). As our next step, we want to find an upper bound for  $r(\sigma_j)$ . We therefore analyse the concave function  $r(\sigma) = \sigma(1 - \tau\sigma^2)^k$  by computing its first derivative, setting it to zero and inserting the resulting argument that maximises  $r$ . This yields

$$\max_{\sigma} r(\sigma) = \frac{\left(\frac{2k}{2k+1}\right)^k}{\sqrt{\tau(2k+1)}} \leq \frac{1}{\sqrt{\tau(2k+1)}}$$

for  $k \in \mathbb{N}$ . Hence, we obtain the estimate

$$\left| \langle u^k - u^\dagger, u_j \rangle_{\mathcal{U}} \right| \leq \frac{|\langle w, u_j \rangle_{\mathcal{U}}|}{\sqrt{\tau(2k+1)}},$$

and consequently

$$\|u^k - u^\dagger\|_{\mathcal{U}} = \sqrt{\sum_{j=1}^{\infty} |\langle u^k - u^\dagger, u_j \rangle_{\mathcal{U}}|^2} \leq \frac{1}{\sqrt{\tau(2k+1)}} \sqrt{\sum_{j=1}^{\infty} |\langle w, u_j \rangle_{\mathcal{U}}|^2} = \frac{\|w\|_{\mathcal{U}}}{\sqrt{\tau(2k+1)}}.$$

□

Together with the stepsize-constraint (3.26) we can further conclude convergence of the iterates to a least squares solution.

**Lemma 3.3.** *Let (3.26) be satisfied. Then the iterates (3.20) satisfy*

$$\|Ku^{k+1} - f\|_{\mathcal{V}} \leq \|Ku^k - f\|_{\mathcal{V}},$$

for  $f = Ku^\dagger$  and all  $k \in \mathbb{N}$ , where equality only holds if  $u^k$  already satisfies the normal equation (2.3).

*Proof.* We easily estimate

$$\begin{aligned}
\|Ku^{k+1} - f\|_{\mathcal{V}}^2 &= \|K(I - \tau K^* K)u^k - (I - \tau K^*)f\|_{\mathcal{V}}^2 \\
&= \|Ku^k - f - \tau K K^*(Ku^k - f)\|_{\mathcal{V}}^2 \\
&= \|Ku^k - f\|_{\mathcal{V}}^2 - 2\tau \langle K^*(Ku^k - f), K^*(Ku^k - f) \rangle_{\mathcal{U}} + \tau^2 \|K K^*(Ku^k - f)\|_{\mathcal{V}}^2 \\
&= \|Ku^k - f\|_{\mathcal{V}}^2 + \tau \left( \tau \|K K^*(Ku^k - f)\|_{\mathcal{V}}^2 - 2 \|K^*(Ku^k - f)\|_{\mathcal{U}}^2 \right) \\
&\leq \|Ku^k - f\|_{\mathcal{V}}^2 + \tau \|K^*(Ku^k - f)\|_{\mathcal{U}}^2 \underbrace{\left( \tau \|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})}^2 - 2 \right)}_{<0} \\
&\leq \|Ku^k - f\|_{\mathcal{V}}^2,
\end{aligned}$$

which proves the statement.  $\square$

Lemma 3.2 and Lemma 3.3 allow us to conclude the following proposition.

**Proposition 3.3.** *The Landweber iteration is a linear regularisation in the sense of Definition 3.2.*

In order to show that the Landweber iteration (3.20) in combination with the discrepancy principle (3.10) is also a convergent regularisation, we obviously have to look at the case of noisy data  $f^\delta$  with  $\|f^\delta - f\|_{\mathcal{V}} \leq \delta$  for  $f = Ku^\dagger$ . We denote the solution of (3.20) in case of noisy data  $f^\delta$  as  $u_\delta^k$  for all  $k \in \mathbb{N}$  and obtain the following estimate for the norm between  $u_\delta^k$  and  $u^\dagger$ .

**Lemma 3.4.** *Let (3.27) be satisfied. Then the Landweber iterates (3.20) satisfy*

$$\|u_\delta^k - u^\dagger\|_{\mathcal{U}} \leq \tau k \delta \|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} + \frac{\|w\|_{\mathcal{U}}}{\sqrt{\tau(2k-1)}} \quad (3.28)$$

for  $k \in \mathbb{N} \setminus \{1\}$ ,  $f = Ku^\dagger$ ,  $f^\delta \in \mathcal{V}$  and  $\|f^\delta - f\|_{\mathcal{V}} \leq \delta$ .

*Proof.* Similar to the proof of Lemma 3.2 we consider the inner product between  $u_\delta^k - u^\dagger$  and a singular vector  $u_j$ , which yields

$$\begin{aligned}
\langle u_\delta^k - u^\dagger, u_j \rangle_{\mathcal{U}} &= \frac{1}{\sigma_j} \left( \left( 1 - (1 - \tau \sigma_j^2)^k \right) \langle f^\delta, v_j \rangle_{\mathcal{V}} - \langle f, v_j \rangle_{\mathcal{V}} \right) \\
&= \frac{1}{\sigma_j} \left( 1 - (1 - \tau \sigma_j^2)^k \right) \langle f^\delta - f, v_j \rangle_{\mathcal{V}} - \sigma_j (1 - \tau \sigma_j^2)^k \langle w, v_j \rangle_{\mathcal{V}}.
\end{aligned}$$

Hence, for  $k > 1$  we can use (3.24) to estimate

$$\begin{aligned}
\frac{1}{\sigma_j} \left( 1 - (1 - \tau \sigma_j^2)^k \right) \left| \langle f^\delta - f, v_j \rangle_{\mathcal{V}} \right| &= \tau \sigma_j \sum_{j=1}^k (1 - \tau \sigma_j^2)^{k-j} \left| \langle f^\delta - f, v_j \rangle_{\mathcal{V}} \right| \\
&\leq \tau k \sigma_j \left| \langle f^\delta - f, v_j \rangle_{\mathcal{V}} \right| \leq \tau k \sigma_1 \left| \langle f^\delta - f, v_j \rangle_{\mathcal{V}} \right|.
\end{aligned}$$

Together with the result from Lemma 3.2 we conclude

$$\|u_\delta^k - u^\dagger\|_{\mathcal{U}} \leq \tau k \delta \|K\|_{\mathcal{L}(\mathcal{U}, \mathcal{V})} + \frac{\|w\|_{\mathcal{U}}}{\sqrt{\tau(2k-1)}}.$$

$\square$

Note that the decrease of the residual in Lemma 3.3 holds true for all  $f \in \mathcal{V}$ . As we obviously do not want to iterate until infinity – this would blow up the data error in (3.28) – this decrease together with the stepsize-constraint (3.26) motivates the use of (3.10) as a stopping criterion. The following lemma shows that with (3.20) we indeed minimise the difference between  $u_\delta^k$  and  $u^\dagger$  (in terms of the  $\mathcal{U}$  norm) as long as (3.10) is violated.

**Lemma 3.5.** *Let (3.26) be satisfied. Then the iterates of (3.20) satisfy*

$$\|u_\delta^{k+1} - u^\dagger\|_{\mathcal{U}} \leq \|u_\delta^k - u^\dagger\|_{\mathcal{U}}$$

for  $k \leq k^*$ ,  $f = Ku^\dagger$  and  $f^\delta \in \mathcal{V}$  with  $\|f^\delta - f\|_{\mathcal{V}} \leq \delta$ . Here,  $k^*$  satisfies the discrepancy principle (3.10) for  $\eta = 2/(2 - \tau\|K\|_{\mathcal{L}(\mathcal{U},\mathcal{V})}^2) > 1$ . Moreover, equality can only be attained for  $\delta = 0$  and  $u_\delta^k$  satisfying the normal equation (2.3).

*Proof.* We prove the statement by showing that  $\|u_\delta^{k+1} - u^\dagger\|_{\mathcal{U}}^2 - \|u_\delta^k - u^\dagger\|_{\mathcal{U}}^2$  is negative whilst the discrepancy principle is not violated. We estimate

$$\begin{aligned} \|u_\delta^{k+1} - u^\dagger\|_{\mathcal{U}}^2 - \|u_\delta^k - u^\dagger\|_{\mathcal{U}}^2 &= \|u_\delta^k - \tau K^*(Ku_\delta^k - f^\delta) - u^\dagger\|_{\mathcal{U}}^2 - \|u_\delta^k - u^\dagger\|_{\mathcal{U}}^2 \\ &= \tau^2 \|K^*(Ku_\delta^k - f^\delta)\|_{\mathcal{U}}^2 - 2\tau \langle Ku_\delta^k - f^\delta, Ku_\delta^k - f \rangle_{\mathcal{V}} \\ &\leq \tau^2 \|K\|_{\mathcal{L}(\mathcal{U},\mathcal{V})}^2 \|Ku_\delta^k - f^\delta\|_{\mathcal{V}}^2 - 2\tau \underbrace{\langle Ku_\delta^k - f^\delta, Ku_\delta^k - f + f^\delta - f^\delta \rangle_{\mathcal{V}}}_{= \|Ku_\delta^k - f^\delta\|_{\mathcal{V}}^2 + \langle Ku_\delta^k - f^\delta, f^\delta - f \rangle_{\mathcal{V}}} \\ &= \tau \left( \tau \|K\|_{\mathcal{L}(\mathcal{U},\mathcal{V})}^2 - 2 \right) \|Ku_\delta^k - f^\delta\|_{\mathcal{V}}^2 + 2\tau \langle f - f^\delta, Ku_\delta^k - f^\delta \rangle_{\mathcal{V}} \\ &\leq \tau \left( \tau \|K\|_{\mathcal{L}(\mathcal{U},\mathcal{V})}^2 - 2 \right) \|Ku_\delta^k - f^\delta\|_{\mathcal{V}}^2 + 2\tau \delta \|Ku_\delta^k - f^\delta\|_{\mathcal{V}} \\ &= -\tau \|Ku_\delta^k - f^\delta\|_{\mathcal{V}} \left( \left( 2 - \tau \|K\|_{\mathcal{L}(\mathcal{U},\mathcal{V})}^2 \right) \|Ku_\delta^k - f^\delta\|_{\mathcal{V}} - 2\delta \right) \\ &= -\frac{2\tau}{\eta} \|Ku_\delta^k - f^\delta\|_{\mathcal{V}} \left( \|Ku_\delta^k - f^\delta\|_{\mathcal{V}} - \eta\delta \right). \end{aligned}$$

Hence, for  $k \leq k_*$  we conclude  $\|u_\delta^{k+1} - u^\dagger\|_{\mathcal{U}} \leq \|u_\delta^k - u^\dagger\|_{\mathcal{U}}$ .  $\square$

### 3.3 Tikhonov regularisation revisited

We conclude this chapter by showing that Tikhonov regularisation can not just be interpreted as the spectral regularisation (3.6) and the solution of the well-posed operator equation (3.14), but also as the minimiser of a functional.

**Theorem 3.10.** *For  $f \in \mathcal{V}$  the Tikhonov-regularised solution  $u_\alpha = R_\alpha f$  with  $R_\alpha$  as defined in (3.6) is uniquely determined as the global minimiser of the Tikhonov-functional*

$$T_\alpha(u) := \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \frac{\alpha}{2} \|u\|_{\mathcal{U}}^2. \quad (3.29)$$

*Proof.*  $\Rightarrow$ : Let  $u_\alpha$  be the Tikhonov-regularised solution and we show that it is also a global minimiser. A global minimiser  $\hat{u} \in \mathcal{U}$  of  $T_\alpha(\hat{u})$  is characterised via  $T_\alpha(\hat{u}) \leq T_\alpha(u)$  for all

$u \in \mathcal{U}$ . Hence, it follows from

$$\begin{aligned}
T_\alpha(u) - T_\alpha(u_\alpha) &= \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \frac{\alpha}{2} \|u\|_{\mathcal{U}}^2 - \frac{1}{2} \|Ku_\alpha - f\|_{\mathcal{V}}^2 - \frac{\alpha}{2} \|u_\alpha\|_{\mathcal{U}}^2 \\
&= \frac{1}{2} \|Ku\|_{\mathcal{V}}^2 - \langle Ku, f \rangle + \frac{\alpha}{2} \|u\|_{\mathcal{U}}^2 - \frac{1}{2} \|Ku_\alpha\|_{\mathcal{V}}^2 + \langle Ku_\alpha, f \rangle - \frac{\alpha}{2} \|u_\alpha\|_{\mathcal{U}}^2 \\
&\quad + \underbrace{\langle (K^*K + \alpha I)u_\alpha - K^*f, u_\alpha - u \rangle}_{=0} \\
&= \frac{1}{2} \|Ku - Ku_\alpha\|_{\mathcal{V}}^2 + \frac{\alpha}{2} \|u - u_\alpha\|_{\mathcal{U}}^2 \\
&\geq 0
\end{aligned}$$

that  $u_\alpha$  is a global minimiser of  $T_\alpha$ .

$\Leftarrow$ : Let now  $\hat{u}$  be a global minimiser. If we have  $T_\alpha(\hat{u}) \leq T_\alpha(u)$  (for all  $u \in \mathcal{U}$ ), it follows with  $u = \hat{u} + \tau v$  for arbitrary  $\tau > 0$  and fixed  $v \in \mathcal{U}$  that

$$0 \leq T_\alpha(u) - T_\alpha(\hat{u}) = \frac{\tau^2}{2} \|Kv\|_{\mathcal{V}}^2 + \frac{\tau^2\alpha}{2} \|v\|_{\mathcal{U}}^2 + \tau \langle (K^*K + \alpha I)\hat{u} - K^*f, v \rangle_{\mathcal{U}}$$

holds true. Dividing by  $\tau$  and subsequent consideration of the limit  $\tau \downarrow 0$  thus yields

$$\langle (K^*K + \alpha I)\hat{u} - K^*f, v \rangle_{\mathcal{U}} \geq 0, \text{ for all } v \in \mathcal{U}.$$

Thus  $(K^*K + \alpha I)\hat{u} - K^*f = 0$  and we conclude  $\hat{u} = u_\alpha$ , i.e. a global minimiser is the Tikhonov-regularised solution. This also shows that the global minimiser of the Tikhonov functional (3.29) is unique.  $\square$

This result paves the way for a generalisation of Tikhonov regularisation to a much broader class of regularisation methods that we want to discuss in the following chapter.



## Chapter 4

# Variational regularisation

At the end of the last chapter we have seen that Tikhonov regularisation<sup>1</sup>  $R_\alpha f$  can be characterised as the solution of the minimisation problem

$$R_\alpha f = \arg \min_{u \in \mathcal{U}} \left\{ \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \frac{\alpha}{2} \|u\|_{\mathcal{U}}^2 \right\}.$$

It is well known that the solution to an unconstrained minimisation problem has a vanishing derivative. In function spaces, the (Gâteaux-) derivative is also called the “first variation” such that minimisation problems are also called *variational problems* and methods that rely on minimising a functional *variational methods*. In this section we want to investigate variational methods for regularisation of linear inverse problems. To do so we will generalise Tikhonov regularisation by choosing different regularisation functionals  $J: \mathcal{U} \rightarrow \mathbb{R}$  and compute regularised solutions by minimising the functional

$$\Phi_{\alpha, f}(u) := \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \alpha J(u).$$

Regularisation of this form is sometimes called *Tikhonov-type regularisation* but we will refer to this as *variational regularisation*. Before we have a look at the theory behind variational regularisation such as the existence and uniqueness of minimisers we will discuss several examples of *regularisation functionals*  $J$ .

**Example 4.1** (Tikhonov-Philipps regularisation). The easiest way to extend classical Tikhonov regularisation to a more general regularisation method is to replace  $\frac{1}{2} \|u\|_{\mathcal{U}}^2$  by  $\frac{1}{2} \|Du\|_{\mathcal{Z}}^2$  where  $D: \mathcal{U} \rightarrow \mathcal{Z}$  is a linear (not necessarily bounded) operator and we thus minimise

$$\frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \frac{\alpha}{2} \|Du\|_{\mathcal{Z}}^2,$$

which became known as *Tikhonov-Philipps regularisation*. While Tikhonov regularisation penalises the norm of  $u$ , in Tikhonov-Philipps regularisation only certain features of  $u$  (depending on the choice of  $D$ ) are penalised. The most frequent used operator  $D$  in imaging applications is the gradient operator  $\nabla$  such that the regulariser  $J$  corresponds to the semi-norm on  $H^1(\Omega)$  which is the Sobolev space of functions  $u \in L^2(\Omega)$  such that the weak derivative  $\nabla u$  exists and  $\nabla u \in L^2(\Omega, \mathbb{R}^n)$ . By using this regulariser, only the variations in  $u$  but not the actual intensities are penalised which helps to control noise without a bias of the intensities towards zero.

---

<sup>1</sup>This regularisation is called *ridge regression* in the statistical literature.

If the operator  $D$  is given by  $Du = (u, \nabla u)$  and  $\mathcal{Z} = L^2(\Omega) \times L^2(\Omega, \mathbb{R}^n)$  is equipped with the natural inner product for product spaces, then

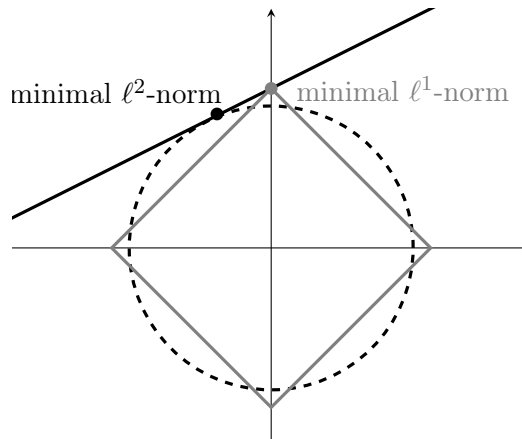
$$J(u) = \frac{1}{2} \|Du\|_{\mathcal{Z}}^2 = \frac{1}{2} \|u\|_{L^2}^2 + \frac{1}{2} \|\nabla u\|_{L^2}^2$$

is the norm on  $H^1(\Omega)$  and it corresponds to classical Tikhonov regularisation on  $H^1(\Omega)$ .

**Example 4.2** ( $\ell^1$ -regularisation). When it comes to non-injective operators  $K \in \mathcal{L}(\ell^1, \ell^2)$  between sequence spaces, the  $\ell^1$ -norm, i.e.  $\|u\|_{\ell^1} := \sum_{j=1}^{\infty} |u_j|$  is often used as a regulariser, in order to enforce sparse solutions. The corresponding minimisation problem reads as

$$\min_{u \in \ell^1} \left\{ \frac{1}{2} \|Ku - f\|_{\ell^2}^2 + \alpha \|u\|_1 \right\}. \quad (4.1)$$

This problem is also called *lasso* in the statistical literature. One can show that minimisers of (4.1) are always *sparse* in the sense that they have finite support, i.e.  $|\text{supp}(u)| < \infty$  with  $\text{supp}(u) = \{i \in \mathbb{N} \mid u_i \neq 0\}$ . This is in contrast to solutions of the Tikhonov regularised problem which may not be sparse. For a finite dimensional example see Figure 4.1.



**Figure 4.1:** Non-injective operators have a non-trivial kernel such that the inverse problem has more than one solution and the solutions form an affine subspace visualised by the solid line. Different regularisation functionals favour different solutions. The circle and the diamond indicate all points with constant  $\ell^2$ -norm, respectively  $\ell^1$ -norm, and the minimal  $\ell^2$ -norm and  $\ell^1$ -norm solutions are the intersections of the line with the circle, respectively the diamond. As it can be seen, the minimal  $\ell^2$ -norm solution has two non-zero components while the minimal  $\ell^1$ -norm solution has only one non-zero component and thus is *sparser*.

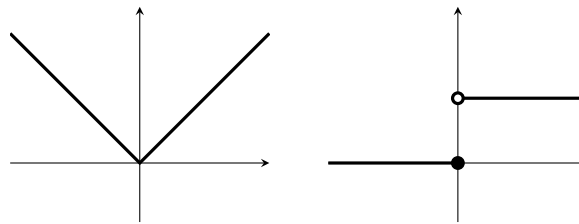
**Example 4.3** (Elastic net). Another regularisation method from statistics is the *elastic net*, where the regulariser is the weighted sum of the  $\ell^1$ -norm and the squared  $\ell^2$ -norm:

$$J(u) = \|u\|_{\ell^1} + \frac{\beta}{2} \|u\|_{\ell^2}^2.$$

Here the idea is to combine two favorable models in order to get sparse solutions with more stability. As  $\ell^1 \subset \ell^2$  we could either consider the elastic net on the Banach space  $\ell^1$  or on the Hilbert space  $\ell^2$ . In case we decide to do the latter, we can extend the elastic net such that

$$J(u) = \begin{cases} \|u\|_{\ell^1} + \frac{\beta}{2} \|u\|_{\ell^2}^2 & \text{if } u \in \ell^1 \\ \infty & \text{if } u \in \ell^2 \setminus \ell^1 \end{cases}.$$





**Figure 4.2:** The absolute value function on the left is in  $H^{1,1}(\Omega)$ ,  $\Omega = [-1, 1]$ , while the Heaviside function on the right is not. The solid dot at a jump indicates the value that the function takes. However, the Heaviside function is in  $BV(\Omega)$  which shows that  $BV(\Omega)$  is larger than  $H^{1,1}(\Omega)$ . Moreover, it shows that  $BV(\Omega)$  includes function with discontinuities which is a favourable model for images with sharp edges.

Intuitively, the value  $\infty$  makes sure that a minimiser will never be in  $\ell^2 \setminus \ell^1$  but we will discuss this aspect in more detail later.

**Example 4.4** (Total variation). Total variation as a regulariser has originally been introduced for image-denoising and -restoration applications with the goal to preserve edges in images, respectively discontinuities in signals [13]. For smooth signals  $u \in H^{1,1}(\Omega)$ , i.e.  $u \in L^1(\Omega)$  and has a weak derivative  $\nabla u \in L^1(\Omega, \mathbb{R}^n)$ , the total variation is simply defined as the semi-norm on the Sobolev space  $H^{1,1}(\Omega)$

$$J(u) = \text{TV}(u) := \int_{\Omega} \|\nabla u(x)\|_2 dx.$$

However, functions in  $H^{1,1}(\Omega)$  may not allow discontinuities which are useful in imaging applications to model images with sharp edges.

To allow discontinuities while still preserving some regularity (otherwise we could model images in  $L^1(\Omega)$  for instance) we generalise the definition of the total variation. It is well-known (e.g. Cauchy–Schwarz inequality) that for  $x, v \in \mathbb{R}^n$  with  $\|v\|_2 \leq 1$  we have that  $\langle v, x \rangle \leq \|x\|_2$ . Thus, for any *test function*  $\varphi \in \mathcal{D}(\Omega, \mathbb{R}^n)$  with

$$\mathcal{D}(\Omega, \mathbb{R}^n) := \left\{ \varphi \in C_0^\infty(\Omega; \mathbb{R}^n) \mid \sup_{x \in \Omega} \|\varphi(x)\|_2 \leq 1 \right\}$$

we have that

$$\text{TV}(u) = \int_{\Omega} \|\nabla u(x)\|_2 dx \geq \int_{\Omega} \langle \nabla u(x), \varphi(x) \rangle dx = - \int_{\Omega} u(x) \text{div } \varphi(x) dx$$

where the last equality is due to partial integration (Gauss’ divergence theorem). In fact one can show that

$$\text{TV}(u) = \sup_{\varphi \in \mathcal{D}(\Omega, \mathbb{R}^n)} \int_{\Omega} u(x) \text{div } \varphi(x) dx,$$

which gives rise to the definition of *functions of bounded variation*.

$$BV(\Omega) := \left\{ u \in L^1(\Omega) \mid \|u\|_{BV} := \|u\|_{L^1} + \text{TV}(u) < \infty \right\}$$

It can be shown that  $BV(\Omega)$  is much larger than  $H^{1,1}(\Omega)$  and contains functions with discontinuities, see examples in Figure 4.2.

**Example 4.5** (Maximum-entropy regularisation). Maximum-entropy regularisation is of particular interest if solutions of the inverse problem are assumed to be probability density functions (pdf), i.e. functions in the set

$$\text{PDF}(\Omega) := \left\{ u \in L^1(\Omega) \mid \int_{\Omega} u(x) dx = 1, u \geq 0 \right\}.$$

The set  $\text{PDF}(\Omega)$  is a convex subset but it is not a subspace as differences of pdfs are not necessarily pdfs. The (differential) entropy used in physics and information theory is defined as the functional  $\text{PDF}(\Omega) \rightarrow \mathbb{R}$  with

$$u \mapsto - \int_{\Omega} u(x) \log(u(x)) dx$$

and the convention  $0 \log(0) := 0$ . As  $\text{PDF}(\Omega)$  is not a vector space, we extend it to the whole  $L^1(\Omega)$  and define the *negative* entropy regularisation as

$$J(u) = \begin{cases} \int_{\Omega} u(x) \log(u(x)) dx & \text{if } u \in \text{PDF}(\Omega) \\ \infty & \text{else} \end{cases}.$$

To summarise the introduction, variational regularisation aims at finding approximations to the solution of the inverse problem (1.1) by minimising appropriate functionals of the form

$$\Phi_{\alpha, f}(u) := \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \alpha J(u), \quad (4.2)$$

where  $J : \mathcal{U} \rightarrow \mathbb{R} \cup \{\infty\}$  represents a functional over the Banach space  $\mathcal{U}$ ,  $\mathcal{V}$  is a Hilbert space and  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  a linear and continuous operator, and  $\alpha > 0$  is a real, positive constant. The term  $D(u) := \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2$  is usually named *fidelity* or *data term*, as it measures the deviation between the measured data  $f$  and the forward model  $Ku$ . The functional  $J$  is the *regularisation term* or *regulariser* as it will impose certain regularity conditions on the unknown  $u$ . The *regularisation parameter* will balance between both terms. Next, we will study some general theory on variational methods that will tell us under which conditions we can expect existence and uniqueness of solutions to those minimisation problems.

## 4.1 Variational methods

### 4.1.1 Background

#### Banach spaces and weak convergence

To cover all the examples of the beginning of this chapter we have to extend our setting to include *Banach spaces*. These are complete, normed vector spaces (as Hilbert spaces) but they may not have an inner product. For every Banach space  $\mathcal{U}$ , we can define the space of linear and continuous functionals which is called the *dual space*  $\mathcal{U}^*$  of  $\mathcal{U}$ , i.e.  $\mathcal{U}^* := \mathcal{L}(\mathcal{U}, \mathbb{R})$ . Let  $u \in \mathcal{U}$  and  $p \in \mathcal{U}^*$ , then we usually write the *dual product*  $\langle p, u \rangle$  instead of  $p(u)$ . Obviously, the dual product is not symmetric (in contrast to the inner product of Hilbert spaces). Moreover, for any  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$  there exists a unique operator  $K^* : \mathcal{V}^* \rightarrow \mathcal{U}^*$ , called the *adjoint* of  $K$  such that for all  $u \in \mathcal{U}$  and  $p \in \mathcal{V}^*$  we have

$$\langle K^* p, u \rangle = \langle p, Ku \rangle.$$

It is easy to see that either side of the equation are well-defined, e.g.  $K^*p \in \mathcal{U}^*$  and  $u \in \mathcal{U}$ .

As the dual space is a Banach space, it has a dual space as well which we will call the bi-dual space of  $\mathcal{U}$  and denote it with  $\mathcal{U}^{**} := (\mathcal{U}^*)^*$ . As every  $u \in \mathcal{U}$  defines a continuous and linear mapping on the dual space  $\mathcal{U}^*$  by

$$\langle E(u), p \rangle := \langle p, u \rangle,$$

the mapping  $E: \mathcal{U} \rightarrow \mathcal{U}^{**}$  is well-defined. It can be shown that  $E$  is a linear and continuous isometry (and thus injective). In the special case when  $E$  is surjective, we call  $\mathcal{U}$  *reflexiv*. Examples of reflexive Banach spaces include Hilbert spaces and  $L^q, \ell^q$  spaces with  $1 < q < \infty$ . We call the space  $\mathcal{U}$  *separable* if there exists a set  $\mathcal{X} \subset \mathcal{U}$  of at most countable cardinality such that  $\overline{\mathcal{X}} = \mathcal{U}$ .

A problem in infinite dimensional spaces is that bounded sequences may fail to have convergent subsequences. An example is for instance in  $\ell^2$  the sequence  $\{u^k\}_{k \in \mathbb{N}} \subset \ell^2$ ,  $u_j^k = 1$  if  $k = j$  and 0 otherwise. It is easy to see that  $\|u^k\|_{\ell^2} = 1$  and that there is no  $u \in \ell^2$  such that  $u^k \rightarrow u$ . To circumvent this problem, we define a weaker topology on  $\mathcal{U}$ . We say that  $\{u^k\}_{k \in \mathbb{N}} \subset \mathcal{U}$  *converges weakly* to  $u \in \mathcal{U}$  if and only if for all  $p \in \mathcal{U}^*$  the sequence of real numbers  $\{\langle p, u^k \rangle\}_{k \in \mathbb{N}}$  converges and

$$\langle p, u_j \rangle \rightarrow \langle p, u \rangle.$$

We will denote weak convergence by  $u^k \rightharpoonup u$ . On a dual space  $\mathcal{U}^*$  we could define another topology (in addition to the strong topology induced by the norm and the weak topology as the dual space is a Banach space as well). We say a sequence  $\{p^k\}_{k \in \mathbb{N}} \subset \mathcal{U}^*$  *converges in weak-\** to  $p \in \mathcal{U}^*$  if and only if

$$\langle p^k, u \rangle \rightarrow \langle p, u \rangle \quad \text{for all } u \in \mathcal{U}$$

and we denote weak-\* convergence by  $p^k \xrightarrow{*} p$ . Similarly, for any topology  $\tau$  on  $\mathcal{U}$  we denote the convergence in that topology by  $u^k \xrightarrow{\tau} u$ .

With these two new notions of convergence, we can solve the problem of bounded sequences:

**Theorem 4.1** (Sequential Banach-Alaoglu Theorem, e.g. [14, p. 70] or [15, p. 141]). *Let  $\mathcal{U}$  be a separable normed vector space. Then every bounded sequence  $\{u^k\}_{k \in \mathbb{N}} \subset \mathcal{U}^*$  has a weak-\* convergent subsequence.*

**Theorem 4.2** ([17, p. 64]). *Each bounded sequence  $\{u^k\}_{k \in \mathbb{N}}$  in a reflexive Banach space  $\mathcal{U}$  has a weakly convergent subsequence.*

### Infinity calculus

We will look at functionals  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  whose range is modelled to be the *extended real line*  $\mathbb{R}_\infty := \mathbb{R} \cup \{+\infty\}$  where the symbol  $\infty$  denotes an element that is not part of the real line that is by definition larger than any other element of the reals, i.e.

$$x < \infty$$

for all  $x \in \mathbb{R}$ . This is useful to model constraints: For instance, if we were trying to minimise  $E: [-1, \infty) \rightarrow \mathbb{R}, x \mapsto x^2$  we could remodel this minimisation problem by  $\tilde{E}: \mathbb{R} \rightarrow \mathbb{R}_\infty$

$$\tilde{E}(x) = \begin{cases} x^2 & \text{if } x \geq -1 \\ \infty & \text{else} \end{cases}.$$

Obviously both functionals have the same minimiser but  $\tilde{E}$  is defined on a vector space and not only on a subset. This has two important consequences: On the one hand, it makes many theoretical arguments easier as we do not need to worry whether  $E(x+y)$  is defined or not. On the other hand, it makes practical implementations easier as we are dealing with unconstrained optimisation instead of constrained optimisation. This comes at a cost that some algorithms are not applicable any more, e.g. the function  $\tilde{E}$  is not differentiable everywhere whereas  $E$  is (in the interior of its domain).

It is useful to note that one can calculate on the extended real line  $\mathbb{R}_\infty$  as we are used to on the real line  $\mathbb{R}$  but the operations with  $\infty$  need yet to be defined. As  $\infty$  is larger than any other element it makes sense that it dominates any other calculation, i.e. for all  $x \in \mathbb{R}$  and  $\lambda > 0$ , we have

$$\begin{aligned} x + \infty &:= \infty + x := \infty & \lambda \cdot \infty &:= \infty \cdot \lambda := \infty \\ x/\infty &:= 0 & \infty + \infty &:= \infty. \end{aligned}$$

However, care needs to be taken as some calculations are *not defined*, e.g.

$$\infty - \infty, \quad 0 \cdot \infty \quad \text{and} \quad \infty \cdot \infty.$$

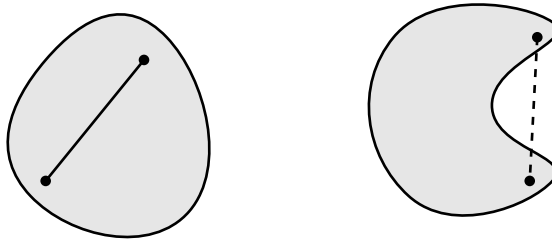
**Definition 4.1.** Let  $\mathcal{U}$  be a vector space and  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  a functional. Then the effective domain of  $E$  is

$$\text{dom}(E) := \{u \in \mathcal{U} \mid E(u) < \infty\}.$$

### Convex calculus

A property of fundamental importance of sets and functions is convexity.

**Definition 4.2.** Let  $\mathcal{U}$  be a vector space. A subset  $\mathcal{C} \subset \mathcal{U}$  is called convex, if  $\lambda u + (1-\lambda)v \in \mathcal{C}$  for all  $\lambda \in (0, 1)$  and all  $u, v \in \mathcal{C}$ .

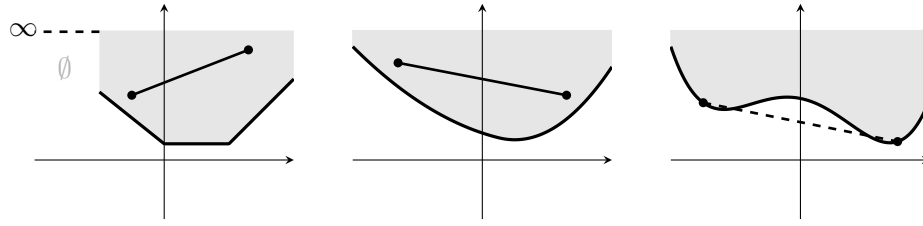


**Figure 4.3:** Example of a convex set (left) and non-convex set (right).

In analogy we can define convex functionals with the help of their epigraph which are all points that lie “above” its graph.

**Definition 4.3.** The epigraph of a functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is defined as the set

$$\text{epi}(E) := \left\{ (u, \lambda) \in \mathcal{U} \times \mathbb{R} \mid E(u) \leq \lambda \right\}.$$



**Figure 4.4:** Example of a convex function (left), a strictly convex function (middle) and a non-convex function (right). Their epigraph are shaded in grey.

**Definition 4.4.** A functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is called convex if its epigraph is convex in  $\mathcal{U} \times \mathbb{R}$ .

It can be shown that this definition is equivalent to the following more common definition.

**Definition 4.5.** A functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is called convex, if

$$E(\lambda u + (1 - \lambda)v) \leq \lambda E(u) + (1 - \lambda)E(v)$$

for all  $\lambda \in (0, 1)$  and all  $u, v \in \text{dom}(E)$  with  $u \neq v$ . It is called strictly convex if the inequality is strict.

**Example 4.6.** The absolute value function  $\mathbb{R} \rightarrow \mathbb{R}, x \mapsto |x|$  is convex but not strictly convex while the quadratic function  $x \mapsto x^2$  is strictly convex. For other examples, see Figure 4.4.

**Example 4.7.** Let  $\mathcal{C} \subset \mathcal{U}$  be a set. Then the characteristic functional  $\chi_{\mathcal{C}}: \mathcal{U} \rightarrow \mathbb{R}_\infty$  with

$$\chi_{\mathcal{C}}(u) := \begin{cases} 0 & u \in \mathcal{C} \\ \infty & u \in \mathcal{U} \setminus \mathcal{C} \end{cases} \quad (4.3)$$

is convex if and only if  $\mathcal{C}$  is a convex set. To see the convexity, let  $u, v \in \text{dom}(\chi_{\mathcal{C}}) = \mathcal{C}$ . Then by the convexity of  $\mathcal{C}$  the convex combination  $\lambda u + (1 - \lambda)v$  is as well in  $\mathcal{C}$  and both the left and the right hand side of the desired inequality are zero.

**Lemma 4.1.** Let  $\alpha \geq 0$  and  $E, F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be two convex functionals. Then  $E + \alpha F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is convex. Furthermore, if  $\alpha > 0$  and  $F$  strictly convex, then  $E + \alpha F$  is strictly convex.

*Proof.* The proof shall be done as an exercise.  $\square$

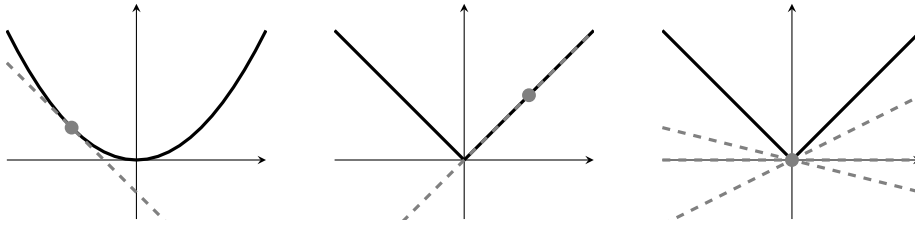
**Definition 4.6.** A functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is called subdifferentiable at  $u \in \mathcal{U}$ , if there exists an element  $p \in \mathcal{U}^*$  such that

$$E(v) \geq E(u) + \langle p, v - u \rangle$$

holds, for all  $v \in \mathcal{U}$ . Furthermore, we call  $p$  a subgradient at position  $u$ . The collection of all subgradients at position  $u$ , i.e.

$$\partial E(u) := \{p \in \mathcal{U}^* \mid E(v) \geq E(u) + \langle p, v - u \rangle, \forall v \in \mathcal{U}\},$$

is called subdifferential of  $E$  at  $u$ .



**Figure 4.5:** Visualisation of the subdifferential. Linear approximations of the functional have to lie completely underneath the function. For points where the function is not differentiable there may be more than one such approximation.

**Remark 4.1.** Let  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be a convex functional. Then the subdifferential is non-empty at all  $u \in \text{dom}(E)$ . If  $\text{dom}(E) \neq \emptyset$ , then for all  $u \notin \text{dom}(E)$  the subdifferential is empty, i.e.  $\partial E(u) = \emptyset$ .

For non-differentiable functionals the subdifferential is multivalued; we want to consider the subdifferential of the absolute value function as an illustrative example.

**Example 4.8.** Let  $E: \mathbb{R} \rightarrow \mathbb{R}$  be the absolute value function  $E(u) = |u|$ . Then, the subdifferential of  $E$  at  $u$  is given by

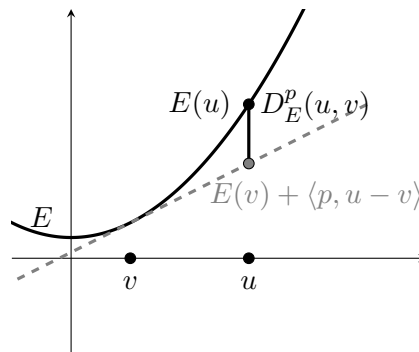
$$\partial E(u) = \begin{cases} \{1\} & \text{for } u > 0 \\ [-1, 1] & \text{for } u = 0 \\ \{-1\} & \text{for } u < 0 \end{cases},$$

which you will prove as an exercise. A visual explanation is given in Figure 4.5.

It turns out that convex functions naturally define some distance measure that became known as the Bregman distance.

**Definition 4.7.** Let  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be a functional. Moreover, let  $u, v \in \mathcal{U}$ ,  $E(v) < \infty$  and  $p \in \partial E(v)$ . Then the (generalised) Bregman distance of  $E$  between  $u$  and  $v$  is defined as

$$D_E^p(u, v) := E(u) - E(v) - \langle p, u - v \rangle. \quad (4.4)$$



**Figure 4.6:** Visualization of the Bregman distance.

**Remark 4.2.** It is easy to check that a Bregman distance somewhat resembles a metric as for all  $u, v \in \mathcal{U}$ ,  $p \in \partial E(v)$  there is  $D_E^p(u, v) \geq 0$  and  $D_E^p(v, v) = 0$ . There are functionals where the Bregman distance (up to a square root) is actually a metric; e.g.  $E(u) := \frac{1}{2}\|u\|_{\mathcal{U}}^2$  for Hilbert space  $\mathcal{U}$ , then  $D_E^p(u, v) = \frac{1}{2}\|u - v\|_{\mathcal{U}}^2$ . However, there are functionals  $E$  where  $D_E^p(u, v) = 0$  does not imply  $u = v$ , as you will see on the example sheets.

### 4.1.2 Minimisers

**Definition 4.8.** Let  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be a functional. We say that  $u^* \in \mathcal{U}$  solves the minimisation problem

$$\min_{u \in \mathcal{U}} E(u)$$

if and only if  $E(u^*) < \infty$  and  $E(u^*) \leq E(u)$ , for all  $u \in \mathcal{U}$ . We call  $u^*$  a minimiser of  $E$ .

We will now review two properties that are necessary for the well-definedness of a minimisation problem.

**Definition 4.9.** A functional  $E$  is called proper if the effective domain  $\text{dom}(E)$  is not empty.

**Definition 4.10.** A functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is called bounded from below if there exists a constant  $C > -\infty$  such that for all  $u \in \mathcal{U}$  we have  $E(u) \geq C$ .

This condition is obviously necessary for the existence of the infimum  $\inf_{u \in \mathcal{U}} E(u)$ .

Finally we characterise minimisers of functionals.

**Theorem 4.3.** An element  $u \in \mathcal{U}$  is a minimiser of the functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  if and only if  $0 \in \partial E(u)$ .

*Proof.* By definition,  $0 \in \partial E(u)$  if and only if for all  $v \in \mathcal{U}$  it holds

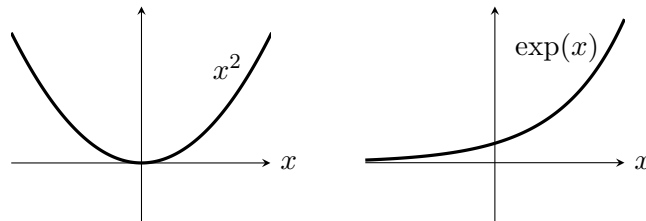
$$E(v) \geq E(u) + \langle 0, v - u \rangle = E(u),$$

which is by definition the case if and only if  $u$  is a minimiser of  $E$ .  $\square$

### 4.1.3 Existence

If all minimising sequences (that converge to the infimum assuming it exists) are unbounded, then there cannot exist a minimiser. A sufficient condition to avoid such a scenario is *coercivity*.

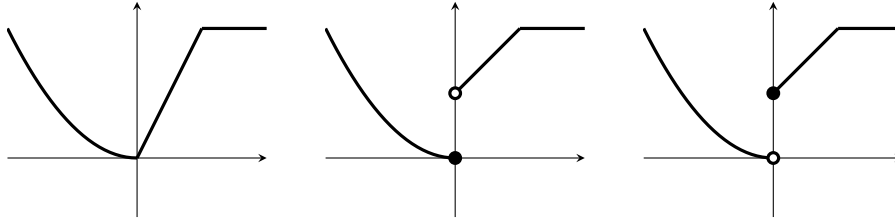
**Definition 4.11.** A functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is called coercive, if for all  $\{u_j\}_{j \in \mathbb{N}}$  with  $\|u_j\|_{\mathcal{U}} \rightarrow \infty$  we have  $E(u_j) \rightarrow \infty$ .



**Figure 4.7:** While the coercive function on the left has a minimiser, it is easy to see that the non-coercive function on the right does not have a minimiser.

**Remark 4.3.** Coercivity is equivalent to its negated statement which is “if the function values  $\{E(u_j)\}_{j \in \mathbb{N}} \subset \mathbb{R}$  are bounded, so is the sequence  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$ ”.

Although coercivity is not strictly speaking necessary, it is sufficient that all minimising sequences are bounded.



**Figure 4.8:** Visualisation of lower semi-continuity. The solid dot at a jump indicates the value that the function takes. The function on the left is continuous and thus lower semi-continuous. The functions in the middle and on the right are discontinuous. While the function in the middle is lower semi-continuous, the function on the right is not (due to the limit from the left at the discontinuity).

**Lemma 4.2.** *Let  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be a proper, coercive functional and bounded from below. Then the infimum  $\inf_{u \in \mathcal{U}} E(u)$  exists in  $\mathbb{R}$ , there are minimising sequences, i.e.  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  with  $E(u_j) \rightarrow \inf_{u \in \mathcal{U}} E(u)$ , and all minimising sequences are bounded.*

*Proof.* As  $E$  is proper and bounded from below, there exists a  $C_1 > 0$  such that we have  $-\infty < -C_1 < \inf_u E(u) < \infty$  which also guarantees the existence of a minimising sequence. Let  $\{u_j\}_{j \in \mathbb{N}}$  be any minimising sequence, i.e.  $E(u_j) \rightarrow \inf_u E(u)$ . Then there exists a  $j_0 \in \mathbb{N}$  such that for all  $j > j_0$  we have

$$E(u_j) \leq \underbrace{\inf_u E(u)}_{=: C_2} + 1 < \infty.$$

With  $C := \max\{C_1, C_2\}$  we have that  $|E(u_j)| < C$  for all  $j > j_0$  and thus from the coercivity it follows that  $\{u_j\}_{j > j_0}$  is bounded, see Remark 4.3. Including a finite number of elements does not change its boundedness which proves the assertion.  $\square$

More importantly we are going to need that functionals are sequentially lower semi-continuous. Roughly speaking this means that the functional values for arguments near an argument  $u$  are either close to  $E(u)$  or greater than  $E(u)$ .

**Definition 4.12.** *Let  $\mathcal{U}$  be a Banach space with topology  $\tau_{\mathcal{U}}$ . The functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is said to be sequentially lower semi-continuous with respect to  $\tau_{\mathcal{U}}$  ( $\tau_{\mathcal{U}}$ -l.s.c.) at  $u \in \mathcal{U}$  if*

$$E(u) \leq \liminf_{j \rightarrow \infty} E(u_j)$$

for all sequences  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  with  $u_j \rightarrow u$  in the topology  $\tau_{\mathcal{U}}$  of  $\mathcal{U}$ .

**Remark 4.4.** For topologies that are not induced by a metric we have to differ between a topological property and its sequential version, e.g. continuous and sequentially continuous. If the topology is induced by a metric, then these two are the same. However, for instance the weak and weak-\* topology are generally not induced by a metric.

**Example 4.9.** The functional  $\|\cdot\|_1: \ell^2 \rightarrow \mathbb{R}_\infty$  with

$$\|u\|_1 = \begin{cases} \sum_{j=1}^{\infty} |u_j| & \text{if } u \in \ell^1 \\ \infty & \text{else} \end{cases}$$

is lower semi-continuous with respect to  $\ell^2$ .



*Proof.* Let  $\{u^j\}_{j \in \mathbb{N}} \subset \ell^2$  be a convergent sequence with  $u^j \rightarrow u \in \ell^2$ . As strong convergence implies weak convergence, we have with  $\delta_k : \ell^2 \rightarrow \mathbb{R}$ ,  $\langle \delta_k, v \rangle = v_k$  that for all  $k \in \mathbb{N}$  that

$$u_k^j = \langle \delta_k, u^j \rangle \rightarrow \langle \delta_k, u \rangle = u_k.$$

The assertion follows then with Fatou's lemma

$$\|u\|_1 = \sum_{k=1}^{\infty} |u_k| = \sum_{k=1}^{\infty} \lim_{j \rightarrow \infty} |u_k^j| \leq \liminf_{j \rightarrow \infty} \sum_{k=1}^{\infty} |u_k^j| = \liminf_{j \rightarrow \infty} \|u^j\|_1.$$

Note that it is not clear whether both the left and the right hand side are finite. □

**Example 4.10.** Let  $\Omega \subset \mathbb{R}^n$  be open and bounded. Then, the total variation is lower semi-continuous with respect to  $L^1$ .

*Proof.* Recall that the total variation was defined by means of the test functions

$$\mathcal{D}(\Omega, \mathbb{R}^n) := \left\{ \varphi \in C_0^\infty(\Omega; \mathbb{R}^n) \mid \sup_{x \in \Omega} \|\varphi(x)\|_2 \leq 1 \right\}$$

as

$$\text{TV}(u) = \sup_{\varphi \in \mathcal{D}(\Omega, \mathbb{R}^n)} \int_{\Omega} u(x) \operatorname{div} \varphi(x) dx.$$

Let  $\{u_j\}_{j \in \mathbb{N}} \subset \text{BV}(\Omega)$  be a sequence converging in  $L^1(\Omega)$  with  $u_j \rightarrow u$  in  $L^1(\Omega)$ . Then for any test function  $\varphi \in \mathcal{D}(\Omega, \mathbb{R}^n)$  there is

$$\int_{\Omega} [u(x) - u_j(x)] \operatorname{div} \varphi(x) dx \leq \underbrace{\int_{\Omega} |u(x) - u_j(x)| dx}_{=\|u-u_j\|_{L^1} \rightarrow 0} \underbrace{\sup_{x \in \Omega} |\operatorname{div} \varphi(x)|}_{< \infty} \rightarrow 0$$

and thus

$$\int_{\Omega} u(x) \operatorname{div} \varphi(x) dx = \liminf_{j \rightarrow \infty} \int_{\Omega} u_j(x) \operatorname{div} \varphi(x) dx \leq \liminf_{j \rightarrow \infty} \text{TV}(u_j).$$

Taking the supremum over all test functions shows the assertion. Note that again the left and right hand side may not be finite. □

We have now all ingredients in place for a positive answer about the existence of minimisers also known as the “direct method” or “fundamental theorem of optimisation”.

**Theorem 4.4** (“Direct method”, David Hilbert, around 1900). *Let  $\mathcal{U}$  be a Banach space and  $\tau_{\mathcal{U}}$  a topology (not necessarily the one induced by the norm) on  $\mathcal{U}$  such that bounded sequences have  $\tau_{\mathcal{U}}$ -convergent subsequences. Let  $E : \mathcal{U} \rightarrow \mathbb{R}_{\infty}$  be proper, bounded from below, coercive and  $\tau_{\mathcal{U}}$ -l.s.c. Then  $E$  has a minimiser.*

*Proof.* From Lemma 4.2 we know that  $\inf_{u \in \mathcal{U}} E(u)$  is finite, minimising sequences exist and that they are bounded. Let  $\{u_j\}_{j \in \mathbb{N}} \in \mathcal{U}$  be a minimising sequence. Thus, from the assumption on the topology  $\tau_{\mathcal{U}}$  there exists a subsequence  $\{u_{j_k}\}_{k \in \mathbb{N}}$  and  $u^* \in \mathcal{U}$  with  $u_{j_k} \xrightarrow{\tau_{\mathcal{U}}} u^*$  for  $k \rightarrow \infty$ . From the sequential lower semi-continuity of  $E$  we obtain

$$E(u^*) \leq \liminf_{k \rightarrow \infty} E(u_{j_k}) = \lim_{j \rightarrow \infty} E(u_j) = \inf_{u \in \mathcal{U}} E(u) < \infty,$$

which shows that  $E(u^*) < \infty$  and  $E(u^*) \leq E(u)$  for all  $u \in \mathcal{U}$ ; thus  $u^*$  minimises  $E$ . □

The above theorem is very general but its conditions are hard to verify but the situation is easier in *reflexive* Banach spaces (thus also in Hilbert spaces).

**Corollary 4.1.** *Let  $\mathcal{U}$  be a reflexive Banach space and  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be a functional which is proper, bounded from below, coercive and l.s.c. with respect to the weak topology. Then there exists a minimiser of  $E$ .*

*Proof.* The statement follows from the direct method, Theorem 4.4, as in reflexive Banach spaces bounded sequences have weakly convergent subsequences, see Theorem 4.2.  $\square$

**Remark 4.5.** For convex functionals on reflexive Banach spaces, the situation is even easier. It can be shown that a convex function is l.s.c. with respect to the weak topology if and only if it is l.s.c. with respect to the strong topology (see e.g. [7, Corollary 2.2., p. 11] or [3, p. 149] for Hilbert spaces).

**Remark 4.6.** It is easy to see that the key ingredient for the existence of minimisers is that bounded sequences have a convergent subsequence which is difficult to prove in practical situations. Another option is to change the space and consider a space in which  $\mathcal{U}$  is compactly embedded in, i.e. the mapping  $\mathcal{U} \rightarrow \mathcal{V}, u \mapsto u$  is compact. Then (by definition) every bounded sequence in  $\mathcal{U}$  has a convergent subsequence in  $\mathcal{V}$ .

#### 4.1.4 Uniqueness

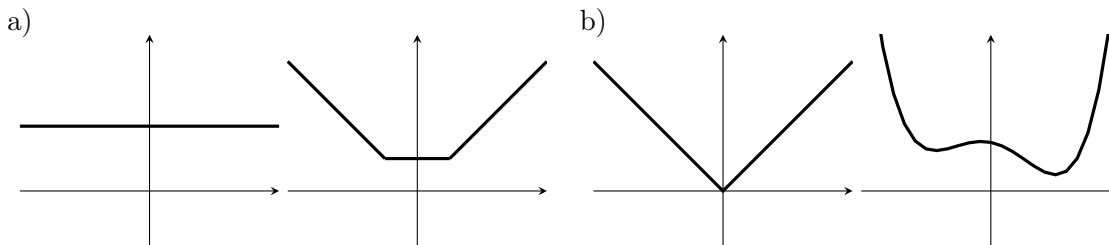
**Theorem 4.5.** *Assume that the functional  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  has at least one minimiser and is strictly convex. Then the minimiser is unique.*

*Proof.* Let  $u, v$  be two minimisers of  $E$  and assume that they are different, i.e.  $u \neq v$ . Then it follows from the minimising properties of  $u$  and  $v$  as well as the strict convexity of  $E$  that

$$E(u) \leq E\left(\frac{1}{2}u + \frac{1}{2}v\right) < \frac{1}{2}E(u) + \frac{1}{2}\underbrace{E(v)}_{\leq E(u)} \leq E(u)$$

which is a contradiction. Thus,  $u = v$  and the assertion is proven.  $\square$

**Example 4.11.** Convex (but not strictly convex) functions may have more than one minimiser, examples include constant and trapezoidal functions, see Figure 4.9. On the other hand, convex (and even non-convex) functions may have a unique minimiser, see Figure 4.9.



**Figure 4.9:** a) Convex functions may not have a unique minimiser. b) Neither strict convexity nor convexity is necessary for the uniqueness of a minimiser.

## 4.2 Well-posedness and regularisation properties

The aim of this section is to have a detailed look at the model  $R_\alpha: \mathcal{V} \rightarrow \mathcal{U}$  with

$$R_\alpha f := \arg \min_{u \in \mathcal{U}} \left\{ \Phi_{\alpha, f}(u) := \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \alpha J(u) \right\}. \quad (4.5)$$

We will establish conditions on the spaces  $\mathcal{U}, \mathcal{V}$ , the functional  $J$  and the operator  $K$  under which a minimiser exists and is unique and therefore the mapping  $R_\alpha$  is well-defined. We will analyse the continuity of the mapping  $R_\alpha$  which means that the solution depends continuously on the data and thus can handle small variations due to noise. We also show that there are parameter choice rules that make  $R_\alpha$  a convergent regularisation in a modified sense (that we will define later) and prove convergence rates under a source condition.

### 4.2.1 Existence and uniqueness

#### Existence

**Lemma 4.3.** *Let  $\mathcal{U}$  be a Banach space and  $\tau_{\mathcal{U}}$  a topology on it. Let  $E: \mathcal{U} \rightarrow \mathbb{R}$  and  $F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be proper functionals that are both  $\tau_{\mathcal{U}}$ -l.s.c. and bounded from below. Then  $E + F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  is proper,  $\tau_{\mathcal{U}}$ -l.s.c. and bounded from below.*

*Proof.* First of all, as  $F$  is proper, there exists  $u \in \mathcal{U}$  such that  $F(u) < \infty$  and as  $E(u) < \infty$  it is clear that  $(E + F)(u) < \infty$  which shows that  $E + F$  is proper.

Second, for all  $u \in \mathcal{U}$  we have from the boundedness from below of  $E$  and  $F$  that  $E(u) \geq C_1$  and  $F(u) \geq C_2$  and thus,

$$(E + F)(u) = E(u) + F(u) \geq C_1 + C_2 > -\infty.$$

Finally, let  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  be a sequence and  $u \in \mathcal{U}$  with  $u_j \rightarrow u$  in  $\tau_{\mathcal{U}}$ . Then by  $\tau_{\mathcal{U}}$ -l.s.c. we have that

$$\begin{aligned} (E + F)(u) &\leq \liminf_{j \rightarrow \infty} E(u_j) + \liminf_{j \rightarrow \infty} F(u_j) \\ &\leq \liminf_{j \rightarrow \infty} (E(u_j) + F(u_j)) = \liminf_{j \rightarrow \infty} (E + F)(u_j) \end{aligned}$$

which shows that  $E + F$  is  $\tau_{\mathcal{U}}$ -l.s.c. and all assertions are proven.  $\square$

**Lemma 4.4.** *Let  $\mathcal{U}$  be a Banach space and  $E, F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be functionals. Let  $E$  be coercive and  $F$  be bounded from below, then  $E + F$  is coercive.*

*Proof.* From the boundedness from below of  $F$ , there exists a constant  $C \in \mathbb{R}$  such that  $F(u) > C$  for all  $u \in \mathcal{U}$ . Thus we see that

$$(E + F)(u) = E(u) + F(u) \geq E(u) + C \rightarrow \infty$$

as  $\|u\|_{\mathcal{U}} \rightarrow \infty$  which proves that  $E + F$  is coercive.  $\square$

In many situations of interest, the lemma above does not apply because the coercivity comes jointly from the data term and the prior as we will see in the following example.

**Example 4.12.** Let  $\Omega \subset \mathbb{R}^n$  be bounded and consider total variation regularisation, i.e.  $J = \text{TV}$  and  $\mathcal{U} = \text{BV}(\Omega)$ . One can easily see that constant functions have zero total variation, i.e.  $\text{TV}(u) = 0$ , for all  $u \equiv c, c \in \mathbb{R}$ . Notice that this implies that  $J$  is not coercive on the whole space  $\mathcal{U}$  as  $u_j(x) = j/|\Omega|, |\Omega| := \int_{\Omega} 1 \, dx$  defines a sequence such that  $\|u_j\|_{L^1} = j$  and  $\text{TV}(u_j) = 0$ . However, we can make use of a Poincaré–Wirtinger inequality for BV.

**Proposition 4.1** ([8, p. 24]). *Let  $\Omega \subset \mathbb{R}^n$  be a Lipschitz domain (non-empty, open, connected and bounded with Lipschitz boundary). There exists a constant  $C > 0$  such that for all  $u \in \text{BV}(\Omega)$  the Poincaré–Wirtinger type inequality is satisfied*

$$\|u - u_\Omega\|_{L^1} \leq C \text{TV}(u)$$

where  $u_\Omega := \frac{1}{|\Omega|} \int_\Omega u(x) dx$  is the mean-value of  $u$  over  $\Omega$ .

**Continuation of Example 4.12.** Let  $\Omega$  now fulfil the conditions of Proposition 4.1. Furthermore, let  $p_0 \in \mathcal{U}^*$  with

$$\langle p_0, u \rangle := u_\Omega = \frac{1}{|\Omega|} \int_\Omega u(x) dx$$

and denote the space of zero-mean functions by  $\mathcal{U}_0 := \{u \in \mathcal{U} \mid u \in \mathcal{N}(p_0)\}$ . By the Poincaré–Wirtinger inequality it is clear that the total variation is coercive on  $\mathcal{U}_0$  and the data term has to make sure that the functional  $\Phi_{\alpha, f}$  is coercive on the whole space  $\mathcal{U}$ . As we will see in the next, Lemma 4.5, the condition  $1 \notin \mathcal{N}(K)$  is sufficient to guarantee coercivity in this scenario. This will follow from a more general result.

**Lemma 4.5.** *Let  $\mathcal{U}, \mathcal{V}$  be Banach spaces,  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ ,  $J: \mathcal{U} \rightarrow [0, \infty]$  and  $f \in \mathcal{V}$ . Let  $p_0 \in \mathcal{U}^*$ ,  $u_0 \in \mathcal{U}$ ,  $\langle p_0, u_0 \rangle = 1$ ,*

$$\mathcal{U}_0 := \{u \in \mathcal{U} \mid u \in \mathcal{N}(p_0)\}$$

so that  $u_0 \notin \mathcal{N}(K)$  and  $J$  is coercive on  $\mathcal{U}_0$  in the sense that

$$\|u - \langle p_0, u \rangle u_0\|_{\mathcal{U}} \rightarrow \infty \quad \text{implies} \quad J(u) \rightarrow \infty.$$

Then the variational regularisation functional  $\Phi_{\alpha, f}$  defined by (4.5) is coercive.

*Proof.* Any  $u \in \mathcal{U}$  can be decomposed into

$$u = v + w$$

where  $v := u - \langle p_0, u \rangle u_0 \in \mathcal{U}_0$  and  $w := \langle p_0, u \rangle u_0 \in \text{span}(u_0)$ .

Now, let  $\{u^j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  be a sequence with  $\|u^j\|_{\mathcal{U}} \rightarrow \infty$ . On the one hand, if  $\|v^j\|_{\mathcal{U}} \rightarrow \infty$ , then by the coercivity of  $J$  on  $\mathcal{U}_0$  and the boundedness from below of the data term, we have that  $\Phi_{\alpha, f}(u^j) \rightarrow \infty$ .

On the other hand, if  $\|v^j\|_{\mathcal{U}} < C$  for some  $C > 0$ , then from

$$\|u^j\|_{\mathcal{U}} \leq \|v^j\|_{\mathcal{U}} + \|w^j\|_{\mathcal{U}} < C + |\langle p_0, u^j \rangle| \|u_0\|_{\mathcal{U}}$$

it follows that  $|\langle p_0, u^j \rangle| \rightarrow \infty$ . Therefore,

$$\begin{aligned} \|K u^j - f\|_{\mathcal{V}} &= \|K(v^j + w^j) - f\|_{\mathcal{V}} \\ &\geq \|K w^j\|_{\mathcal{V}} - \|K v^j - f\|_{\mathcal{V}} \\ &\geq \underbrace{\|K u_0\|_{\mathcal{V}}}_{>0} \underbrace{|\langle p_0, u^j \rangle|}_{\rightarrow \infty} \underbrace{- \|K\| C - \|f\|_{\mathcal{V}}}_{\text{constant}} \rightarrow \infty \end{aligned}$$

and thus  $\Phi_{\alpha, f}(u^j) \rightarrow \infty$  as the regularisation functional  $J$  is bounded from below.  $\square$

**Remark 4.7.** A natural question here is whether the coercivity can also come completely from the data term  $D(u) = \frac{1}{2}\|Ku - f\|_{\mathcal{V}}^2$ . On the one hand if  $K$  is not injective, then the kernel is non-trivial, thus we cannot expect the data term to be coercive.

On the other hand, even if  $K$  was injective we cannot expect coercivity in general. Assume that the data term  $D$  was coercive,  $\mathcal{U}$  a Hilbert space, the topologies  $\tau_{\mathcal{U}}, \tau_{\mathcal{V}}$  the weak topologies on  $\mathcal{U}, \mathcal{V}$ ,  $f \in \overline{\mathcal{R}(K)} \setminus \mathcal{R}(K)$  and  $D$   $\tau_{\mathcal{U}}$ -l.s.c. (e.g. see Lemma 4.6). Then we can apply the direct method on the data term  $D$  and would prove the existence of a minimiser. This is by definition a least squares solution; a contradiction to Lemma 2.2.

The remark will be illustrated by the following example.

**Example 4.13.** Let us consider the Example 2.1 again where the operator was  $K: \ell^2 \rightarrow \ell^2$ ,  $(Ku)_j := u_j/j$  and the data  $f \in \ell^2$  with  $f_j := 1/j$ . Then the  $\{u^k\}_{k \in \mathbb{N}} \subset \ell^2$  with

$$u_j^k := \begin{cases} 1 & j \leq k \\ 0 & \text{else} \end{cases}$$

defines a sequence  $\{Ku^k\}_{k \in \mathbb{N}}$  with  $Ku^k \rightarrow f$  in  $\ell^2$ .

Note that  $K$  is injective and  $\{u^k\}_{k \in \mathbb{N}}$  is a minimising sequence of the data term, i.e.  $\|Ku^k - f\|_{\ell^2} \rightarrow 0$ . However, as  $\|u^k\|_{\ell^2} = k$  the sequence  $\{u^k\}_{k \in \mathbb{N}}$  is unbounded and thus  $u \mapsto \|Ku - f\|_{\ell^2}$  cannot be coercive.

**Lemma 4.6.** *Let  $\mathcal{U}$  and  $\mathcal{V}$  be Banach spaces with topologies  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$ . Moreover, let the norm on  $\mathcal{V}$  be  $\tau_{\mathcal{V}}$ -l.s.c., the operator  $K: \mathcal{U} \rightarrow \mathcal{V}$  be sequentially continuous with respect to the topologies  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$  and let  $\{f_j\}_{j \in \mathbb{N}} \subset \mathcal{V}$  be convergent in  $\tau_{\mathcal{V}}$  with  $f_j \rightarrow f \in \mathcal{V}$ . Then for any  $\tau_{\mathcal{U}}$ -convergent sequence  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  with  $u_j \rightarrow u \in \mathcal{U}$  with respect to  $\tau_{\mathcal{U}}$ , we have*

$$\frac{1}{2}\|Ku - f\|_{\mathcal{V}}^2 \leq \liminf_{j \rightarrow \infty} \frac{1}{2}\|Ku_j - f_j\|_{\mathcal{V}}^2.$$

*In particular, if  $f_j \equiv f$ , then  $D: \mathcal{U} \rightarrow \mathbb{R}, u \mapsto \frac{1}{2}\|Ku - f\|_{\mathcal{V}}^2$  is  $\tau_{\mathcal{U}}$ -l.s.c.*

*Proof.* Let  $\{u_j\}_{j \in \mathbb{N}}$  be a  $\tau_{\mathcal{U}}$ -convergent sequence and denote its limit by  $u \in \mathcal{U}$ , i.e.  $u_j \rightarrow u$  in  $\tau_{\mathcal{U}}$ . Because  $K$  is continuous with respect to  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$  we have that  $Ku_j \rightarrow Ku$  in  $\tau_{\mathcal{V}}$  and thus  $Ku_j - f_j \rightarrow Ku - f$  in  $\tau_{\mathcal{V}}$ . Thus, the assertion follows from the  $\tau_{\mathcal{V}}$ -l.s.c. of the norm.  $\square$

**Remark 4.8.** If the topologies  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$  are the weak topologies, then the situation is much simpler as continuity in the strong topologies implies continuity in the weak topologies [5, chapter IV.3]. Thus the assumptions of Lemma 4.6 are met if  $K$  is continuous.

Note that the norm is convex (will be proven in what follows) such that weak-l.s.c. is equivalent to l.s.c.

Now we are in a position to state sufficient assumptions for the existence of minimisers.

**Assumption 4.1.** *Sufficient assumptions for the existence of minimisers of  $\Phi_{\alpha, f}$  are:*

- (a) *The Banach space  $\mathcal{U}$  and Hilbert space  $\mathcal{V}$  are associated with the topologies  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$ . The pair  $(\mathcal{U}, \tau_{\mathcal{U}})$  has the property that bounded sequences have  $\tau_{\mathcal{U}}$ -convergent subsequences. Moreover, the norm on  $\mathcal{V}$  is  $\tau_{\mathcal{V}}$ -l.s.c.*
- (b) *The operator  $K: \mathcal{U} \rightarrow \mathcal{V}$  is linear and sequentially continuous with respect to the topologies  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$ .*

(c) The functional  $J: \mathcal{U} \rightarrow [0, \infty]$  is proper and  $\tau_{\mathcal{U}}$ -l.s.c.

(d) Either  $J$  is coercive or the pair  $(K, J)$  fulfil the assumptions of Lemma 4.5.

**Theorem 4.6.** *Let the Assumptions 4.1 hold and let  $f \in \mathcal{V}, \alpha > 0$ . Then the variational regularisation functional  $\Phi_{\alpha, f}$  defined by (4.5) has a minimiser.*

*Proof.* It follows from the assumptions by Lemmata 4.3 and 4.6 that  $\Phi_{\alpha, f}$  is proper,  $\tau_{\mathcal{U}}$ -l.s.c. and bounded from below. Moreover, from Lemmata 4.4 or 4.5 (depending on assumption 4.1 (d))  $\Phi_{\alpha, f}$  is coercive. Then from the direct method, Theorem 4.4, it follows that there exists a minimiser.  $\square$

### Uniqueness

**Lemma 4.7.** *Let  $\mathcal{U}$  be a Banach space and  $\mathcal{V}$  a Hilbert space. Furthermore, let  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V}), f \in \mathcal{V}$  and  $D: \mathcal{U} \rightarrow \mathbb{R}_{\infty}$  be defined as  $D(u) := \frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2$ . Then  $D$  is convex. Furthermore,  $D$  is strictly convex if and only if  $K$  is injective.*

*Proof.* Let  $\lambda \in (0, 1)$  and  $u, v \in \mathcal{U}$  with  $u \neq v$ .

Long and straightforward calculations (good exercise) yield

$$D(\lambda u + (1 - \lambda)v) = \lambda D(u) + (1 - \lambda)D(v) - \underbrace{\frac{\lambda(1 - \lambda)}{2} \|K(u - v)\|_{\mathcal{V}}^2}_{\geq 0},$$

which shows that  $D$  is convex.

Note that  $\|K(u - v)\|_{\mathcal{V}} > 0$  if  $K$  is injective and  $u \neq v$ . On the other hand, if  $K$  is not injective, then we can find  $u, v \in \mathcal{U}$  with  $u \neq v$  so that  $u - v \in \mathcal{N}(K)$  and thus  $\|K(u - v)\|_{\mathcal{V}} = 0$  which shows that  $D$  is not strictly convex.  $\square$

**Remark 4.9.** The lemma is in general not true if  $\mathcal{V}$  is a Banach space. As an example, consider  $K = I, f = 0, \mathcal{V} = \mathbb{R}^2$  with  $\|\cdot\|_1$  as a norm. Then  $D$  is not strictly convex.

**Example 4.14.** Let  $\mathcal{U}$  be continuously embedded into the Hilbert space  $\mathcal{Z}$  (in symbols  $\mathcal{U} \hookrightarrow \mathcal{Z}$ ), i.e. there exists a constant  $C > 0$  such that for all  $u \in \mathcal{U}$  there is  $\|u\|_{\mathcal{Z}} \leq C\|u\|_{\mathcal{U}}$ . Furthermore, let  $\beta > 0$  and  $J$  be convex. Then the functional  $\Phi_{\alpha, f} + \frac{\beta}{2} \|\cdot\|_{\mathcal{Z}}^2$  is always strictly convex independent of  $K$ .

Consider the product space  $\mathcal{V} \times \mathcal{Z}$  which is a Hilbert space with the inner product

$$\langle (v_1, z_1), (v_2, z_2) \rangle_{\mathcal{V} \times \mathcal{Z}} := \langle v_1, v_2 \rangle_{\mathcal{V}} + \langle z_1, z_2 \rangle_{\mathcal{Z}}.$$

Then we can rewrite  $\frac{1}{2} \|Ku - f\|_{\mathcal{V}}^2 + \frac{\beta}{2} \|u\|_{\mathcal{Z}}^2$  as

$$\frac{1}{2} \left\| \begin{pmatrix} K \\ \sqrt{\beta}I \end{pmatrix} u - \begin{pmatrix} f \\ 0 \end{pmatrix} \right\|_{\mathcal{V} \times \mathcal{Z}}^2 = \frac{1}{2} \|\tilde{K}u - \tilde{f}\|_{\mathcal{V}}^2$$

where the modified operator  $\tilde{K}$  is injective. Therefore, adding the term  $\frac{\beta}{2} \|u\|_{\mathcal{Z}}^2$  can be seen as a regularisation of the linear operator  $K$  directly.

**Theorem 4.7.** *Let the Assumptions 4.1 hold and let  $J$  be convex. Moreover, let either  $K$  be injective or  $J$  be strictly convex. Then for any  $f \in \mathcal{V}$  and  $\alpha > 0$  the variational regularisation model  $R_{\alpha}$  is well-defined (there exists a unique minimiser of the functional  $\Phi_{\alpha, f}$  defined by (4.5)).*

*Proof.* Existence follows immediately from Theorem 4.6. For the uniqueness, notice that both  $\frac{1}{2}\|Ku - f\|_{\mathcal{V}}^2$  and  $\alpha J$  are convex and either of them is strictly convex by assumption, see Lemma 4.7. Thus by Lemma 4.1 the whole functional  $\Phi_{\alpha,f}$  is strictly convex and therefore the minimiser is unique, see Theorem 4.5.  $\square$

**Example 4.15.** Let  $\alpha, \beta > 0, K \in \mathcal{L}(\ell^2, \ell^2)$  and consider the elastic net variational regularisation model  $\frac{1}{2}\|Ku - f\|_{\ell^2}^2 + \alpha J(u)$  with

$$J(u) = \begin{cases} \|u\|_1 + \frac{\beta}{2}\|u\|_2^2 & \text{if } u \in \ell^1 \\ \infty & \text{else} \end{cases}.$$

We check the assumptions of Theorem 4.7 in this setting to conclude that  $R_\alpha$  is well-defined.

Let  $\mathcal{U}, \mathcal{V} = \ell^2$  and as these are Hilbert spaces we choose the topologies  $\tau_{\mathcal{U}}, \tau_{\mathcal{V}}$  to be both the weak topology on  $\ell^2$ . By Remark 4.8 the Assumptions 4.1 (a) and (b) are fulfilled and we can show lower semi-continuity in the strong topology rather than the weak one. It is easy to see that the prior  $J$  is strictly convex, proper and coercive. It remains to show that  $J$  is lower semi-continuous with respect to  $\ell^2$ . The squared  $\ell^2$ -norm is continuous, thus lower semi-continuous and the lower semi-continuity of the  $\ell^1$ -norm has been proven in Example 4.9 such that the elastic net is lower semi-continuous by Lemma 4.3. As it is convex, l.s.c. is equivalent to weak-l.s.c. and all assumptions of Theorem 4.7 hold.

For the example of the total variation we need to have some knowledge about compact embeddings of  $\text{BV}(\Omega)$ .

**Theorem 4.8** (Rellich-Kandrachov, [1, p. 168]). *Let  $\Omega \subset \mathbb{R}^n$  be a Lipschitz domain and either*

$$\begin{aligned} & n > mp \quad \text{and} \quad p^* := np/(n - mp) \\ \text{or} \quad & n \leq mp \quad \text{and} \quad p^* := \infty. \end{aligned}$$

*Then the embedding  $H^{m,p}(\Omega) \rightarrow L^q(\Omega)$  is continuous if  $1 \leq q \leq p^*$  and compact if in addition  $q < p^*$ .*

Due to approximations of  $u \in \text{BV}(\Omega)$  by smooth functions the Rellich-Kandrachov Theorem (for  $m = 1, p = 1$ ) gives us compactness for  $\text{BV}(\Omega)$ .

**Corollary 4.2** ([8, p. 17]). *For any Lipschitz domain  $\Omega \subset \mathbb{R}^n$  the embedding*

$$\text{BV}(\Omega) \rightarrow L^1(\Omega)$$

*is compact.*

**Example 4.16.** Let  $\Omega \subset \mathbb{R}^n$  be a Lipschitz domain,  $\alpha > 0, K \in \mathcal{L}(L^1(\Omega), L^2(\Omega))$  be injective and consider the total variation regularised model  $R_\alpha: L^2(\Omega) \rightarrow \text{BV}(\Omega), R_\alpha f := \arg \min_{u \in \text{BV}(\Omega)} \Phi_{\alpha,f}(u)$  with

$$\Phi_{\alpha,f}(u) = \frac{1}{2}\|Ku - f\|_{L^2}^2 + \alpha \text{TV}(u).$$

This time we are neither in a Hilbert nor in a reflexive Banach space setting but from Corollary 4.2 we see that  $\text{BV}(\Omega)$  is compactly embedded in  $L^1(\Omega)$ . Thus every sequence bounded in  $\text{BV}(\Omega)$  has a convergent subsequence in  $L^1(\Omega)$ .

Let  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$  be the topologies induced by the  $L^1$ -norm, respectively  $L^2$ -norm. It is clear that the assumptions on the spaces and topologies are met. The lower semi-continuity of TV with respect to  $L^1$  was shown in Example 4.10. Moreover, it can be shown that TV is proper and convex. The injectivity of  $K$  implies that  $1 \notin \mathcal{N}(K)$ . From Example 4.12 and  $1 \notin \mathcal{N}(K)$  it follows that  $\Phi_{\alpha,f}$  is coercive. Thus, a minimiser exists. The injectivity of the operator  $K$  guarantees strict convexity and therefore the uniqueness of the minimiser.

### 4.2.2 Continuity

We have seen that under some assumptions the variational regularisation  $R_\alpha$  is well-defined (solutions exist and are unique). In this section we show that variational regularisation is continuous with respect to the data, i.e. small variations in the data do not lead to arbitrary large distortions in the solution. To establish the main result we have to prove auxiliary lemmata.

**Lemma 4.8.** *Let  $\mathcal{V}$  be a normed space. For all  $f, g \in \mathcal{V}$  there is*

$$\|f + g\|_{\mathcal{V}}^2 \leq 2\|f\|_{\mathcal{V}}^2 + 2\|g\|_{\mathcal{V}}^2.$$

*Proof.* For any  $f, g \in \mathcal{V}$  we have with the monotonicity of  $[0, \infty) \rightarrow [0, \infty), x \mapsto x^2$  that

$$\begin{aligned} \|f + g\|_{\mathcal{V}}^2 &\leq \left(\|f\|_{\mathcal{V}} + \|g\|_{\mathcal{V}}\right)^2 \\ &= \|f\|_{\mathcal{V}}^2 + 2\|f\|_{\mathcal{V}}\|g\|_{\mathcal{V}} + \|g\|_{\mathcal{V}}^2. \end{aligned}$$

We complete the proof with the observation that  $2ab \leq a^2 + b^2$  for all  $a, b \in \mathbb{R}$ .  $\square$

**Lemma 4.9.** *Let  $\mathcal{U}, \mathcal{V}$  be normed spaces. For all  $u \in \mathcal{U}$  and  $f, g \in \mathcal{V}$  there is*

$$\Phi_{\alpha, f}(u) \leq 2\Phi_{\alpha, g}(u) + \|f - g\|_{\mathcal{V}}^2.$$

*Proof.* Using Lemma 4.8 and  $J(u) \geq 0$ , we have

$$\begin{aligned} \Phi_{\alpha, f}(u) &= \frac{1}{2}\|Ku - f\|_{\mathcal{V}}^2 + \alpha J(u) = \frac{1}{2}\|Ku - g + (g - f)\|_{\mathcal{V}}^2 + \alpha J(u) \\ &\leq \|Ku - g\|_{\mathcal{V}}^2 + \|g - f\|_{\mathcal{V}}^2 + 2\alpha J(u) \\ &= 2\left(\frac{1}{2}\|Ku - g\|_{\mathcal{U}}^2 + \alpha J(u)\right) + \|f - g\|_{\mathcal{V}}^2 \\ &= 2\Phi_{\alpha, g}(u) + \|f - g\|_{\mathcal{V}}^2. \end{aligned}$$

$\square$

**Theorem 4.9** (Continuity). *Assume the setting of Theorem 4.7 that guarantees the existence and uniqueness of minimisers of  $\Phi_{\alpha, f}(u) := \frac{1}{2}\|Ku - f\|_{\mathcal{V}}^2 + \alpha J(u)$  for any  $f \in \mathcal{V}$  and  $\alpha > 0$ . Moreover, let the topology  $\tau_{\mathcal{V}}$  on  $\mathcal{V}$  be weaker than the norm topology in the sense that convergence in norm implies convergence in  $\tau_{\mathcal{V}}$ . Then, the mapping  $R_\alpha: \mathcal{V} \rightarrow \mathcal{U}, R_\alpha f := \arg \min_{u \in \mathcal{U}} \Phi_{\alpha, f}(u)$  is sequentially strong- $\tau_{\mathcal{U}}$  continuous, i.e. for all sequences  $\{f_j\}_{j \in \mathbb{N}} \subset \mathcal{V}$  with  $f_j \rightarrow f$  we have*

$$R_\alpha f_j \rightarrow R_\alpha f \quad \text{in } \tau_{\mathcal{U}}.$$

*Moreover, the function values of the regulariser converge, i.e.  $J(R_\alpha f_j) \rightarrow J(R_\alpha f)$ .*

*Proof.* Let  $\{f_j\}_{j \in \mathbb{N}} \subset \mathcal{V}$  be a convergent sequence with  $f_j \rightarrow f$  and let  $u_j := R_\alpha f_j$  be the minimiser of  $\Phi_{\alpha, f_j}$ .

We first show that  $\{\Phi_{\alpha, f}(u_j)\}_{j \in \mathbb{N}} \subset \mathbb{R}$  is bounded. To see this, as  $J$  is proper, there exists  $\tilde{u} \in \mathcal{U}$  such that  $J(\tilde{u}) < \infty$  and we denote  $C := 2\|K\tilde{u}\|_{\mathcal{V}}^2 + 2\alpha J(\tilde{u})$ . With Lemma 4.9 and the minimising property of  $u_j$  we have that

$$\begin{aligned} \Phi_{\alpha, f}(u_j) &\leq 2\Phi_{\alpha, f_j}(u_j) + \|f - f_j\|_{\mathcal{V}}^2 \\ &\leq 2\Phi_{\alpha, f_j}(\tilde{u}) + \|f - f_j\|_{\mathcal{V}}^2 = \|K\tilde{u} - f_j\|_{\mathcal{V}}^2 + 2\alpha J(\tilde{u}) + \|f - f_j\|_{\mathcal{V}}^2 \end{aligned}$$



With Lemma 4.8 we can further estimate

$$\begin{aligned} \|K\tilde{u} - f_j\|_{\mathcal{V}}^2 + 2\alpha J(\tilde{u}) + \|f - f_j\|_{\mathcal{V}}^2 &\leq 2\|K\tilde{u}\|_{\mathcal{V}}^2 + 2\|f_j\|_{\mathcal{V}}^2 + 2\alpha J(\tilde{u}) + \|f - f_j\|_{\mathcal{V}}^2 \\ &\leq 2\|f_j\|_{\mathcal{V}}^2 + \|f - f_j\|_{\mathcal{V}}^2 + C. \end{aligned}$$

As  $f_j$  converges to  $f$ , there exists a  $j_0 \in \mathbb{N}$  such that for all  $j > j_0$  there is

$$\Phi_{\alpha,f}(u_j) \leq 2 \underbrace{\|f_j\|_{\mathcal{V}}^2}_{\|f\|_{\mathcal{V}}^2+1} + \underbrace{\|f - f_j\|_{\mathcal{V}}^2}_{\leq 1} + C \leq 2\|f\|_{\mathcal{V}}^2 + C + 3 < \infty.$$

By the coercivity of  $\Phi_{\alpha,f}$  we know that the sequence  $\{u_j\}_{j \in \mathbb{N}} \subset \mathcal{U}$  is bounded, see Remark 4.3. Thus there exist  $\tau_{\mathcal{U}}$ -convergent subsequences and let  $\{u_{j_k}\}_{k \in \mathbb{N}} \subset \mathcal{U}$  be any one of those. We denote its limit by  $u^* \in \mathcal{U}$ , i.e.  $u_{j_k} \rightarrow u^*$  in  $\tau_{\mathcal{U}}$ .

From Lemma 4.6 and as  $J$  is  $\tau_{\mathcal{U}}$ -l.s.c. we have that

$$\frac{1}{2}\|Ku^* - f\|_{\mathcal{V}}^2 \leq \liminf_{k \rightarrow \infty} \frac{1}{2}\|Ku_{j_k} - f_{j_k}\|_{\mathcal{V}}^2 \quad \text{and} \quad J(u^*) \leq \liminf_{k \rightarrow \infty} J(u_{j_k}) \quad (4.6)$$

and therefore

$$\begin{aligned} \Phi_{\alpha,f}(u^*) &= \frac{1}{2}\|Ku^* - f\|_{\mathcal{V}}^2 + \alpha J(u^*) \\ &\leq \liminf_{k \rightarrow \infty} \frac{1}{2}\|Ku_{j_k} - f_{j_k}\|_{\mathcal{V}}^2 + \alpha \liminf_{k \rightarrow \infty} J(u_{j_k}) \\ &\leq \liminf_{k \rightarrow \infty} \left( \frac{1}{2}\|Ku_{j_k} - f_{j_k}\|_{\mathcal{V}}^2 + \alpha J(u_{j_k}) \right) = \liminf_{k \rightarrow \infty} \Phi_{\alpha,f_{j_k}}(u_{j_k}). \end{aligned}$$

With the minimising property of  $u_{j_k}$  we arrive at

$$\begin{aligned} \Phi_{\alpha,f}(u^*) &\leq \liminf_{k \rightarrow \infty} \Phi_{\alpha,f_{j_k}}(u_{j_k}) \\ &\leq \liminf_{k \rightarrow \infty} \Phi_{\alpha,f_{j_k}}(u) = \lim_{k \rightarrow \infty} \Phi_{\alpha,f_{j_k}}(u) = \Phi_{\alpha,f}(u) \end{aligned} \quad (4.7)$$

for all  $u \in \mathcal{U}$ . In particular, this holds for  $u = R_{\alpha}f$ . Thus, as the minimiser of  $\Phi_{\alpha,f}$  is unique, we have that  $u^* = R_{\alpha}f$ . Repeating the same arguments as above for any subsequence of  $\{u_j\}_{j \in \mathbb{N}}$  instead of  $\{u_j\}_{j \in \mathbb{N}}$ , we see that every subsequence has a convergent subsequence that converges to  $R_{\alpha}f$  in  $\tau_{\mathcal{U}}$ . Thus,  $\{u_j\}_{j \in \mathbb{N}}$  is convergent in  $\tau_{\mathcal{U}}$  and we have  $u_j \rightarrow R_{\alpha}f$  in  $\tau_{\mathcal{U}}$  and the first assertion is proven.

Moreover, from (4.7) it follows with  $u_j \rightarrow R_{\alpha}f$  in  $\tau_{\mathcal{U}}$  that

$$\lim_{j \rightarrow \infty} \Phi_{\alpha,f_j}(u_j) = \Phi_{\alpha,f}(R_{\alpha}f).$$

and therefore

$$\begin{aligned} \limsup_{j \rightarrow \infty} \alpha J(u_j) &= \limsup_{j \rightarrow \infty} \left( \frac{1}{2}\|Ku_j - f_j\|_{\mathcal{V}}^2 + \alpha J(u_j) - \frac{1}{2}\|Ku_j - f_j\|_{\mathcal{V}}^2 \right) \\ &\leq \limsup_{j \rightarrow \infty} \underbrace{\left( \frac{1}{2}\|Ku_j - f_j\|_{\mathcal{V}}^2 + \alpha J(u_j) \right)}_{=\lim_{j \rightarrow \infty} \Phi_{\alpha,f_j}(u_j) = \Phi_{\alpha,f}(R_{\alpha}f)} + \limsup_{j \rightarrow \infty} \left( -\frac{1}{2}\|Ku_j - f_j\|_{\mathcal{V}}^2 \right) \\ &= \frac{1}{2}\|KR_{\alpha}f - f\|_{\mathcal{V}}^2 + \alpha J(R_{\alpha}f) - \underbrace{\liminf_{j \rightarrow \infty} \frac{1}{2}\|Ku_j - f_j\|_{\mathcal{V}}^2}_{\leq -\frac{1}{2}\|KR_{\alpha}f - f\|_{\mathcal{V}}^2 \text{ by (4.6)}} \\ &\leq \alpha J(R_{\alpha}f). \end{aligned}$$

The  $\tau_{\mathcal{U}}$ -l.s.c. of  $J$  proves the second assertion

$$\limsup_{j \rightarrow \infty} \alpha J(u_j) \leq \alpha J(R_\alpha f) \leq \liminf_{j \rightarrow \infty} \alpha J(u_j).$$

□

**Remark 4.10.** In the theorem above we could only prove convergence in  $\tau_{\mathcal{U}}$ . If  $J$  satisfies the *Radon-Riesz property* with respect to the topology  $\tau_{\mathcal{U}}$ , i.e.  $u_j \rightarrow u$  in  $\tau_{\mathcal{U}}$  and  $J(u_j) \rightarrow J(u)$  imply  $\|u_j - u\|_{\mathcal{U}} \rightarrow 0$ , then the convergence is in the norm topology. An example of a functional satisfying the Radon-Riesz property is  $\|\cdot\|_{L^p}^p / \|\cdot\|_{\ell^p}^p$  with  $1 < p < \infty$  if the underlying space is  $L^p / \ell^p$  and  $\tau_{\mathcal{U}}$  being the weak topology.

### 4.2.3 Convergent regularisation

Note that variational regularisation for general  $J$  is not necessarily a regularisation in the sense of Definition 3.1, as we cannot expect  $R_\alpha f = \arg \min_{u \in \mathcal{U}} \Phi_{\alpha, f}(u) \rightarrow u^\dagger$  for  $\alpha \rightarrow 0$  where  $u^\dagger$  is the minimal norm solution. However, we can generalise Definition 2.1 of a minimal norm solution (and a least squares solution) to justify calling  $R_\alpha$  a regularisation.

**Definition 4.13.** Let  $\mathcal{U}$  and  $\mathcal{V}$  be normed spaces and  $f \in \mathcal{V}$ . We call  $u \in \mathcal{U}$  a least squares solution of the inverse problem (1.1), if

$$u \in \arg \min_{v \in \mathcal{U}} \|Kv - f\|_{\mathcal{V}} \quad (4.8)$$

As in the case of Hilbert spaces, we denote by  $\mathbb{L}$  the set of all least squares solutions (it might be empty). Furthermore, we call  $u^\dagger \in \mathcal{U}$  a  $J$ -minimising solution of the inverse problem (1.1), if

$$u^\dagger \in \arg \min_{u \in \mathbb{L}} J(u). \quad (4.9)$$

**Remark 4.11.** If  $\mathcal{V}$  is a Hilbert space (as in our setting in Assumption 4.1), then most of the statements from Chapter 2 about least squares solutions still hold. In particular Lemma 2.2, which states that  $\mathbb{L} \neq \emptyset$  if and only if  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ . However, the minimal norm solution ( $J$ -minimising solution for  $J = \|\cdot\|_{\mathcal{U}}$ ) may not be unique any more.

**Lemma 4.10.** Let  $\mathcal{U}$  be a vector space and  $E: \mathcal{U} \rightarrow \mathbb{R}_\infty$  a convex functional. Then the set of all minimisers  $\mathbb{S}$ , i.e.

$$\mathbb{S} := \arg \min_{u \in \mathcal{U}} E(u),$$

is convex.

*Proof.* For any  $u, v \in \mathbb{S}, u \neq v$  and  $\lambda \in (0, 1)$ , it is easy to see that

$$\begin{aligned} E(\lambda u + (1 - \lambda)v) &\leq \lambda E(u) + (1 - \lambda)E(v) \\ &\leq \lambda \inf_{w \in \mathcal{U}} E(w) + (1 - \lambda) \inf_{w \in \mathcal{U}} E(w) = \inf_{w \in \mathcal{U}} E(w) \end{aligned}$$

which shows that  $\lambda u + (1 - \lambda)v \in \mathbb{S}$  and thus  $\mathbb{S}$  is convex. □

**Corollary 4.3.** Let  $\mathcal{U}$  and  $\mathcal{V}$  be normed spaces,  $f \in \mathcal{V}$  and  $K: \mathcal{U} \rightarrow \mathcal{V}$  be linear. Then the set of least squares solutions  $\mathbb{L}$  is convex.

*Proof.* One can show that  $u \mapsto \|Ku - f\|_{\mathcal{V}}$  is convex so that the lemma applies. □

**Lemma 4.11.** *Let  $\mathcal{U}$  be a Banach space,  $\mathcal{V}$  be a Hilbert space,  $f \in \mathcal{V}$  and  $K: \mathcal{U} \rightarrow \mathcal{V}$  be linear and injective. Then the set of least squares solutions  $\mathbb{L}$  is at most a singleton.*

*Proof.* Assume that least squares solutions exist which are equivalently characterised by

$$u \in \arg \min_{v \in \mathcal{U}} \left\{ D(v) := \frac{1}{2} \|Kv - f\|_{\mathcal{V}}^2 \right\}. \quad (4.10)$$

From Lemma 4.7 we know that  $D$  is strictly convex and thus by Theorem 4.5 the minimiser is unique.  $\square$

**Proposition 4.2.** *Let the assumptions of Theorem 4.7 hold and  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$ . Then a  $J$ -minimising solution exists and is unique.*

*Proof.* By Remark 4.11, the condition  $f \in \mathcal{R}(K) \oplus \mathcal{R}(K)^\perp$  guarantees the existence of least squares solutions, i.e.  $\mathbb{L} \neq \emptyset$ . For the existence of  $J$ -minimising solutions via the direct method, Theorem 4.4, we see that only the coercivity on  $\mathbb{L}$  may not be guaranteed by the assumptions.

If  $J$  is coercive, then it is obviously also coercive on  $\mathbb{L}$ . If  $J$  is only coercive on  $\mathcal{U}_0$  (see Lemma 4.5 for a definition), then a similar calculation as in the proof of Lemma 4.5 shows that for any sequence  $\{u^j\}_{j \in \mathbb{N}} \subset \mathbb{L}$  we have with

$$u^j = \underbrace{u^j - \langle p_0, u^j \rangle u_0}_{=: v^j} + \underbrace{\langle p_0, u^j \rangle u_0}_{=: w^j}$$

that

$$\begin{aligned} \|f\|_{\mathcal{V}} &= \|K0 - f\|_{\mathcal{V}} \geq \|Kw^j - f\|_{\mathcal{V}} \\ &= \|K(v^j + w^j) - f\|_{\mathcal{V}} \\ &\geq \|Kw^j\|_{\mathcal{V}} - \|f - Kv^j\|_{\mathcal{V}} > \|Ku_0\|_{\mathcal{V}} |\langle p_0, u^j \rangle| - \|f\|_{\mathcal{V}} - \|K\| \|v^j\|_{\mathcal{V}}. \end{aligned}$$

Let  $\{u^j\}_{j \in \mathbb{N}} \subset \mathbb{L}$  with  $\|u^j\|_{\mathcal{U}} \rightarrow \infty$ . If  $\|v^j\|_{\mathcal{V}}$  was bounded, then so would be  $|\langle p_0, u^j \rangle|$  which contradicts the unboundedness of  $\{u^j\}_{j \in \mathbb{N}}$ . Thus

$$\|v^j\|_{\mathcal{U}} = \|u^j - \langle p_0, u^j \rangle u_0\|_{\mathcal{U}} \rightarrow \infty$$

and  $J(u^j) \rightarrow \infty$  by the coercivity of  $J$  on  $\mathcal{U}_0$ .

For the uniqueness, either  $J$  is strictly convex (and thus a minimiser is unique) or  $K$  is injective and only one least squares solution exists.  $\square$

**Definition 4.14** (Regularisation). *Let  $\mathcal{U}, \mathcal{V}$  be Banach spaces,  $\tau_{\mathcal{U}}$  a topology on  $\mathcal{U}$ ,  $f \in \mathcal{V}$  and  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ . Moreover, let  $u^\dagger$  be the  $J$ -minimising solution (assuming it exists and is unique). We call the family of operators  $\{R_\alpha\}_{\alpha > 0}, R_\alpha: \mathcal{V} \rightarrow \mathcal{U}$  a regularisation (with respect to  $\tau_{\mathcal{U}}$ ) of the inverse problem (1.1), if  $R_\alpha$  is sequentially strong- $\tau_{\mathcal{U}}$  continuous for all  $\alpha > 0$  and*

$$R_\alpha f \rightarrow u^\dagger \quad \text{in } \tau_{\mathcal{U}} \text{ as } \alpha \rightarrow 0.$$

**Theorem 4.10.** *Let the assumptions of Theorem 4.7 hold and assume (for simplicity) that the clean data is in the range, i.e.  $f \in \mathcal{R}(K)$ , thus the  $J$ -minimising solution  $u^\dagger$  exists and is unique. Moreover, assume that the topology  $\tau_{\mathcal{V}}$  is weaker than the norm topology on  $\mathcal{V}$ . Let  $\alpha: (0, \infty) \rightarrow (0, \infty)$  be a parameter choice rule with*

$$\alpha(\delta) \rightarrow 0, \quad \text{and} \quad \frac{\delta^2}{\alpha(\delta)} \rightarrow 0 \quad \text{as } \delta \rightarrow 0.$$

Let  $\{\delta_j\}_{j \in \mathbb{N}} \subset [0, \infty)$  be a sequence of noise levels with  $\delta_j \rightarrow 0$  and  $\{f_j\}_{j \in \mathbb{N}} \subset \mathcal{V}$  be a sequence of noisy observations with  $\|f - f_j\|_{\mathcal{V}} \leq \delta_j$ . Set  $\alpha_j := \alpha(\delta_j)$  and let  $\{u_j\}_{j \in \mathbb{N}}$  be the sequence of minimisers of  $\Phi_{\alpha_j, f_j}$ , i.e.  $u_j := R_{\alpha_j} f_j$ .

Then  $u_j \rightarrow u^\dagger$  in  $\tau_{\mathcal{U}}$  and  $J(u_j) \rightarrow J(u^\dagger)$ . In particular,  $R_\alpha$  is a regularisation.

*Proof.* From minimising property of  $u_j$  and  $Ku^\dagger = f$ , it follows that

$$\begin{aligned} 0 &\leq \frac{1}{2} \|Ku_j - f_j\|_{\mathcal{V}}^2 + \alpha_j J(u_j) \\ &\leq \frac{1}{2} \|Ku^\dagger - f_j\|_{\mathcal{V}}^2 + \alpha_j J(u^\dagger) \\ &= \frac{1}{2} \|f - f_j\|_{\mathcal{V}}^2 + \alpha_j J(u^\dagger) \leq \frac{\delta_j^2}{2} + \alpha_j J(u^\dagger) \rightarrow 0 \end{aligned} \quad (4.11)$$

as  $j \rightarrow \infty$  ( $\delta_j, \alpha_j \rightarrow 0$ ). In particular  $\lim_{j \rightarrow \infty} \|Ku_j - f_j\|_{\mathcal{V}} = 0$ , and then

$$\|Ku_j - f\|_{\mathcal{V}} \leq \|Ku_j - f_j\|_{\mathcal{V}} + \|f_j - f\|_{\mathcal{V}} \leq \|Ku_j - f_j\|_{\mathcal{V}} + \delta_j \rightarrow 0. \quad (4.12)$$

Similarly, we see from (4.11) that

$$\limsup_{j \rightarrow \infty} J(u_j) \leq \limsup_{j \rightarrow \infty} \frac{\delta_j^2}{2\alpha_j} + J(u^\dagger) = J(u^\dagger). \quad (4.13)$$

Let  $\alpha^+ := \max_{j \in \mathbb{N}} \alpha_j$  be the largest regularisation parameter (which exists as  $\alpha_j \rightarrow 0$ ), then

$$\begin{aligned} \limsup_{j \rightarrow \infty} \Phi_{\alpha^+, f}(u_j) &= \limsup_{j \rightarrow \infty} \left( \frac{1}{2} \|Ku_j - f\|_{\mathcal{V}}^2 + \alpha^+ J(u_j) \right) \\ &\leq \underbrace{\limsup_{j \rightarrow \infty} \frac{1}{2} \|Ku_j - f\|_{\mathcal{V}}^2}_{=0} + \underbrace{\limsup_{j \rightarrow \infty} \alpha^+ J(u_j)}_{\leq \alpha^+ J(u^\dagger)} \leq \alpha^+ J(u^\dagger) =: C < \infty \end{aligned}$$

This shows that there exists a  $j_0 \in \mathbb{N}$  such that for all  $j \geq j_0$  we have that  $\Phi_{\alpha^+, f}(u_j) \leq C+1$ .

From the coercivity of  $\Phi_{\alpha^+, f}$  it follows that  $\{u_j\}_{j \in \mathbb{N}}$  is bounded and therefore has a  $\tau_{\mathcal{U}}$ -convergent subsequence  $\{u_{j_k}\}_{k \in \mathbb{N}}$  with  $u_{j_k} \rightarrow \hat{u}$  in  $\tau_{\mathcal{U}}$ . By the continuity of  $K$  with respect to  $\tau_{\mathcal{U}}$  and  $\tau_{\mathcal{V}}$  we have that  $Ku_{j_k} \rightarrow K\hat{u}$  in  $\tau_{\mathcal{V}}$ . With the  $\tau_{\mathcal{V}}$ -l.s.c. of the norm of  $\mathcal{V}$  we conclude that  $K\hat{u} = f$  as

$$\|K\hat{u} - f\|_{\mathcal{V}} \leq \liminf_{k \rightarrow \infty} \|Ku_{j_k} - f\|_{\mathcal{V}} = 0.$$

Thus  $\hat{u}$  is a least squares solution.

From (4.13) and the  $\tau_{\mathcal{U}}$ -l.s.c. of  $J$  we have that

$$J(\hat{u}) \leq \liminf_{k \rightarrow \infty} J(u_{j_k}) \leq \limsup_{k \rightarrow \infty} J(u_{j_k}) \leq J(u^\dagger) \leq J(u) \quad (4.14)$$

for all  $u \in \mathcal{U}$ . Thus,  $\hat{u}$  is a  $J$ -minimising solution, which implies by its uniqueness that  $\hat{u} = u^\dagger$ . Moreover, from (4.14) we deduce  $J(u_{j_k}) \rightarrow J(u^\dagger)$ .

As in the proof of Theorem 4.9, all arguments can be applied to any subsequence of  $\{u_j\}_{j \in \mathbb{N}}$ , which shows that  $u_j \rightarrow u^\dagger$  in  $\tau_{\mathcal{U}}$  and  $J(u_j) \rightarrow J(u^\dagger)$ .  $\square$

**Remark 4.12.** Similar to the stability we can get strong convergence if  $J$  satisfies the Radon-Riesz property.

#### 4.2.4 Convergence rates

In the last section we have proven convergence of the regularisation method in the topology  $\tau_{\mathcal{U}}$  and not in the norm as in Chapter 3. Thus, we cannot expect to prove convergence rates in the norm. However, it turns out we can prove convergence rates in the Bregman distance.

**Theorem 4.11.** *Assume the setting of Theorem 4.7 that guarantees that the mapping  $R_\alpha$  is well-defined. Let  $f \in \mathcal{R}(K)$  be clean data and  $u^\dagger$  be a solution of the inverse problem, i.e.  $f = Ku^\dagger$ , and consider noisy data  $f^\delta \in \mathcal{V}$  with  $\|f - f^\delta\|_{\mathcal{V}} \leq \delta$ . Moreover, let  $u^\dagger$  satisfy the source condition*

$$p = K^*w \in \partial J(u^\dagger)$$

and denote  $u_\alpha^\delta := R_\alpha f^\delta$ . Then,

$$\begin{aligned} (a) \quad & D_J^p(u_\alpha^\delta, u^\dagger) \leq \frac{\delta^2}{2\alpha} + \|w\|_{\mathcal{V}^*}\delta + \frac{\alpha\|w\|_{\mathcal{V}^*}^2}{2}, \\ (b) \quad & \frac{1}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 \leq \delta^2 + 2\alpha\|w\|_{\mathcal{V}^*}\delta + 2\alpha^2\|w\|_{\mathcal{V}^*}^2, \text{ and} \\ (c) \quad & J(u_\alpha^\delta) \leq \frac{\delta^2}{2\alpha} + J(u^\dagger). \end{aligned}$$

Moreover, for the a-priori parameter choice rule  $\alpha(\delta) = \delta$  we have

$$D_J^p(u_\alpha^\delta, u^\dagger) = \mathcal{O}(\delta), \quad \|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}} = \mathcal{O}(\delta), \quad \text{and} \quad J(u_\alpha^\delta) \leq J(u^\dagger) + \mathcal{O}(\delta).$$

*Proof.* From the minimising property of  $u_\alpha^\delta$  and  $Ku^\dagger = f$  it follows that

$$\frac{1}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 + \alpha J(u_\alpha^\delta) \leq \frac{1}{2}\|Ku^\dagger - f^\delta\|_{\mathcal{V}}^2 + \alpha J(u^\dagger) \leq \frac{\delta^2}{2} + \alpha J(u^\dagger). \quad (4.15)$$

From the non-negativity of the data term and (4.15) we derive assertion (c) as

$$J(u_\alpha^\delta) \leq \frac{1}{\alpha} \left( \frac{1}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 + \alpha J(u_\alpha^\delta) \right) \leq \frac{\delta^2}{2\alpha} + J(u^\dagger).$$

Moreover, by reordering the terms of (4.15) and completing Bregman distance, we get

$$\frac{1}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 + \alpha D_J^p(u_\alpha^\delta, u^\dagger) \leq \frac{\delta^2}{2} - \alpha \langle p, u_\alpha^\delta - u^\dagger \rangle$$

where we can further estimate

$$\begin{aligned} -\langle p, u_\alpha^\delta - u^\dagger \rangle &= -\langle w, K(u_\alpha^\delta - u^\dagger) \rangle = -\langle w, Ku_\alpha^\delta - f \rangle \\ &\leq \|w\|_{\mathcal{V}^*}\|Ku_\alpha^\delta - f\|_{\mathcal{V}} \leq \|w\|_{\mathcal{V}^*}(\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}} + \delta) \end{aligned}$$

Combining the two yields

$$\begin{aligned} \frac{1}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 + \alpha D_J^p(u_\alpha^\delta, u^\dagger) &\leq \frac{\delta^2}{2} + \alpha\|w\|_{\mathcal{V}^*}\delta + \alpha\|w\|_{\mathcal{V}^*}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}} \\ &\leq \frac{\delta^2}{2} + \alpha\|w\|_{\mathcal{V}^*}\delta + \frac{\alpha^2\|w\|_{\mathcal{V}^*}^2}{2\gamma} + \frac{\gamma}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 \end{aligned}$$

where we used  $ab \leq \frac{1}{2}a^2 + \frac{1}{2}b^2$  for the second inequality. Thus, we derive

$$(1 - \gamma)\frac{1}{2}\|Ku_\alpha^\delta - f^\delta\|_{\mathcal{V}}^2 + \alpha D_J^p(u_\alpha^\delta, u^\dagger) \leq \frac{\delta^2}{2} + \alpha\|w\|_{\mathcal{V}^*}\delta + \frac{\alpha^2\|w\|_{\mathcal{V}^*}^2}{2\gamma}.$$

Choosing  $\gamma = 1$  and  $\gamma = 1/2$  yields the assertions (a) and (b).  $\square$

**Remark 4.13.** Note that we did not use the source condition for the assertion (c), thus it is true for all solutions  $u^\dagger$  of the inverse problem, i.e.  $Ku^\dagger = f$ .

**Remark 4.14.** We did not explicitly assume that  $u^\dagger$  is a  $J$ -minimising solution. However, let  $\mathcal{U}$  be a Hilbert space and  $J(u) = \frac{1}{2}\|u\|_{\mathcal{U}}^2$ . Then the source condition is equivalent to  $K^*w = u^\dagger$  which is in turn equivalent to  $u^\dagger \in \mathcal{R}(K^*) = \mathcal{N}(K)^\perp$ . Thus, by Corollary 2.1,  $u^\dagger$  is the minimal norm solution.

# Chapter 5

## Numerical Solutions

### 5.1 More on derivatives in Banach spaces

In order to make better sense out of the optimality conditions, we have to discuss some more properties of the subdifferential.

**Definition 5.1.** Let  $E: \mathcal{U} \rightarrow \mathbb{R}$  be a mapping from the Banach space  $\mathcal{U}$  and  $u \in \mathcal{U}$ . If there exists an operator  $A \in \mathcal{L}(\mathcal{U}, \mathbb{R}) = \mathcal{U}^*$  that

$$\lim_{h \rightarrow 0} \frac{|E(u+h) - E(u) - Ah|}{\|h\|_{\mathcal{U}}} = 0,$$

holds true, then  $E$  is called Fréchet differentiable in  $x$  and  $E'(u) := A$  the Fréchet derivative in  $u$ . If the Fréchet derivative exists for all  $u \in \mathcal{U}$ , the operator  $E' : \mathcal{U} \rightarrow \mathcal{U}^*$  is called Fréchet derivative.

**Example 5.1.** Let  $\mathcal{U}$  be a Banach space and  $p \in \mathcal{U}^*$ . Then the Fréchet derivative of  $p$  is given by  $p' = p$ .

**Example 5.2.** Let  $\mathcal{U}$  be a Hilbert space and  $M \in \mathcal{L}(\mathcal{U}, \mathcal{U})$ . Then the Fréchet derivative of  $E: \mathcal{U} \rightarrow \mathbb{R}$ ,

$$E(u) = \|u\|_M^2 := \langle Mu, u \rangle$$

at any  $u \in \mathcal{U}$  is given by

$$E'(u) = \langle (M + M^*)u, \cdot \rangle,$$

and thus by the Riesz representation theorem can be identified with  $(M + M^*)u$ .

Moreover, if  $M$  is self-adjoint then  $E'(u) = 2Mu$ .

*Proof.* Simple calculations show that

$$E(u+h) - E(u) = \langle (M + M^*)u, h \rangle + \langle Mh, h \rangle$$

which shows that

$$\frac{|E(u+h) - E(u) - \langle (M + M^*)u, h \rangle|}{\|h\|_{\mathcal{U}}} = \frac{|\langle Mh, h \rangle|}{2\|h\|_{\mathcal{U}}} \leq \frac{\|M\| \|h\|_{\mathcal{U}}}{2} \rightarrow 0$$

for  $\|h\| \rightarrow 0$ . □

**Example 5.3.** Let  $\mathcal{U}, \mathcal{V}$  be Hilbert spaces,  $K \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ ,  $f \in \mathcal{V}$  and  $E: \mathcal{U} \rightarrow \mathbb{R}$  be defined as  $E(u) := \frac{1}{2}\|Ku - f\|_{\mathcal{U}}^2$ . Then the Fréchet derivative of  $E$  can be identified with

$$E'(u) = K^*(Ku - f).$$

*Proof.* For any  $u \in \mathcal{U}$ , an easy calculation shows that

$$\frac{1}{2}\|K(u+h) - f\|_{\mathcal{U}}^2 - \frac{1}{2}\|Ku - f\|_{\mathcal{U}}^2 = \langle Ku - f, Kh \rangle_{\mathcal{U}} + \frac{1}{2}\|Kh\|_{\mathcal{U}}^2$$

and thus with  $Ah := \langle K^*(Ku - f), h \rangle_{\mathcal{U}}$  we have that

$$\begin{aligned} \frac{|E(u+h) - E(u) - Ah|}{\|h\|_{\mathcal{U}}} &= \frac{|\frac{1}{2}\|K(u+h) - f\|_{\mathcal{U}}^2 - \frac{1}{2}\|Ku - f\|_{\mathcal{U}}^2 - \langle K^*(Ku - f), h \rangle_{\mathcal{U}}|}{\|h\|_{\mathcal{U}}} \\ &= \frac{|\langle Ku - f, Kh \rangle_{\mathcal{U}} + \frac{1}{2}\|Kh\|_{\mathcal{U}}^2 - \langle Ku - f, Kh \rangle_{\mathcal{U}}|}{\|h\|_{\mathcal{U}}} \\ &= \frac{1}{2}\|Kh\|_{\mathcal{U}} \leq \frac{1}{2}\|K\|\|h\|_{\mathcal{U}} \rightarrow 0 \end{aligned}$$

as  $\|h\|_{\mathcal{U}} \rightarrow 0$ . □

**Proposition 5.1.** Let  $\mathcal{U}$  be a Banach space and  $E: \mathcal{U} \rightarrow \mathbb{R}$  be a convex functional that is Fréchet differentiable in  $u \in \mathcal{U}$ . Then

$$\partial E(u) = \{E'(u)\}.$$

*Proof.* Let  $p \in \partial E(u)$ . Then for every  $v \in \mathcal{U}, h > 0$  there is

$$\langle p, v \rangle = \frac{1}{h}\langle p, u + hv \rangle \leq \frac{1}{h}[E(u + hv) - E(u)]$$

and similar for  $h < 0$  there is

$$\langle p, v \rangle \geq \frac{1}{h}[E(u + hv) - E(u)].$$

Thus,

$$\begin{aligned} \langle E'(u), v \rangle &= \lim_{h \uparrow 0} \frac{1}{h}[E(u + hv) - E(u)] \\ &\leq \langle p, v \rangle \leq \lim_{h \downarrow 0} \frac{1}{h}[E(u + hv) - E(u)] = \langle E'(u), v \rangle \end{aligned}$$

which shows  $p = E'(u)$ .

On the other hand, let  $v \in \mathcal{U}$  and  $h \in (0, 1]$  there is

$$\begin{aligned} \frac{1}{h}[E(u + h(v - u)) - E(u)] &= \frac{1}{h}[E((1-h)u + hv) - E(u)] \\ &\leq \frac{1}{h}[(1-h)E(u) + hE(v) - E(u)] = -E(u) + E(v) \end{aligned}$$

and thus

$$\begin{aligned} E(u) + \langle E'(u), u - v \rangle &= E(u) + \lim_{h \downarrow 0} \frac{1}{h}[E(u + h(v - u)) - E(u)] \\ &\leq E(u) - E(u) + E(v) = E(v) \end{aligned}$$

thus  $E'(u) \in \partial E(u)$ . □



**Theorem 5.1** (e.g. [6, p. 279]). *Let  $\mathcal{U}, \mathcal{V}$  be normed spaces,  $E, F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be proper and convex. Then the following rules for the subdifferential hold.*

(a)  $\partial(\lambda E) = \lambda \partial E$  for  $\lambda > 0$ .

(b)  $\partial(E \circ T_v) = \partial E(u+v)$  for any  $v \in \mathcal{U}$  and the translation operator  $T_v: \mathcal{U} \rightarrow \mathcal{U}, T_v(u) = u + v$ .

(c)  $\partial E + \partial F \subset \partial(E + F)$  and equality if there exists  $u \in \text{dom}(E) \cap \text{dom}(F)$  such that  $E$  is continuous in  $u$ .

**Corollary 5.1.** *Let  $\mathcal{U}$  be a normed space,  $E: \mathcal{U} \rightarrow \mathbb{R}$  be convex and Fréchet differentiable and  $F: \mathcal{U} \rightarrow \mathbb{R}_\infty$  be convex. Then for all  $u \in \text{dom}(E + F) = \text{dom}(F)$  it holds*

$$\partial(E + F)(u) = E'(u) + \partial F(u).$$

## 5.2 Minimization problems

### 5.2.1 Gradient Descent

In this section we will analyse the iteration

$$u^{k+1} = u^k - \tau E'(u^k) \tag{5.1}$$

called *gradient descent* which is one of the most popular iterations to solve smooth minimisation problems.

**Lemma 5.1** (Descent Lemma). *Let  $E: \mathcal{U} \rightarrow \mathbb{R}$  be Fréchet differentiable and  $E'$  Lipschitz continuous with constant  $L \in \mathbb{R}$  (which we will call  $L$ -smooth in what is to follow). Then for all  $x, y \in \mathcal{U}$  there is*

$$E(x) \leq E(y) + \langle E'(y), x - y \rangle + \frac{L}{2} \|x - y\|^2.$$

*Proof.* For any  $t \in [0, 1]$  define  $g(t) := E(y + t(x - y))$  for which we obviously have  $g(1) = E(x)$  and  $g(0) = E(y)$ . Then we have that

$$\begin{aligned} \int_0^1 \langle E'(y + t(x - y)) - E'(y), x - y \rangle dt &\leq \int_0^1 \|E'(y + t(x - y)) - E'(y)\| \|x - y\| dt \\ &\leq \int_0^1 Lt \|x - y\|^2 dt \\ &= \frac{L}{2} \|x - y\|^2 \end{aligned}$$

and can further estimate

$$\begin{aligned} E(x) - E(y) = g(1) - g(0) &= \int_0^1 g'(t) dt \\ &= \int_0^1 \langle E'(y + t(x - y)), x - y \rangle dt \\ &= \int_0^1 \langle E'(y), x - y \rangle dt + \int_0^1 \langle E'(y + t(x - y)) - E'(y), x - y \rangle dt \\ &\leq \langle E'(y), x - y \rangle + \frac{L}{2} \|x - y\|^2. \end{aligned}$$

□

**Remark 5.1.** If  $E$  is convex, then the inequality of the lemma can also be written in terms of the Bregman distance as  $D_E^{E'(y)}(x, y) \leq \frac{L}{2}\|x - y\|^2$ .

**Theorem 5.2** (Convergence of gradient descent). *Let  $E$  be  $L$ -smooth and the step size of gradient descent be chosen as*

$$\tau < \frac{2}{L}.$$

*Then gradient descent monotonically decreases the function value, i.e.*

$$E(u^{k+1}) \leq E(u^k).$$

*Moreover, if  $E$  is bounded from below, then the gradients convergence to zero, i.e.*

$$\|E'(u^k)\| \rightarrow 0,$$

*with rate (for some  $C > 0$ )*

$$\min_{k=0, \dots, K-1} \|E'(u^k)\| \leq \frac{C}{K^{1/2}}.$$

*Proof.* Choosing  $x = u^{k+1}$  and  $y = u^k$  in the Descent Lemma yields

$$\begin{aligned} E(u^{k+1}) - E(u^k) &\leq \langle E'(u^k), -\tau E'(u^k) \rangle + \frac{L}{2} \|\tau E'(u^k)\|^2 \\ &= -\tau \|E'(u^k)\|^2 + \frac{\tau^2 L}{2} \|E'(u^k)\|^2 = -\frac{c}{2} \|E'(u^k)\|^2 \end{aligned} \tag{5.2}$$

with  $c := \tau L \left(\frac{2}{L} - \tau\right) > 0$  which shows the monotonic descent.

Moreover, summing (5.2) over  $k = 0, \dots, K - 1$  yields

$$E(u^K) - E(u^0) \leq -c \sum_{k=0}^{K-1} \|E'(u^k)\|^2$$

and after rearranging

$$\sum_{k=0}^{K-1} \|E'(u^k)\|^2 \leq \frac{E(u^0) - E(u^K)}{c} \leq \frac{E(u^0) - \inf_{u \in \mathcal{U}} E(u)}{c} \leq C^2.$$

Thus, letting  $K \rightarrow \infty$  we have that

$$\|E'(u^k)\| \rightarrow 0$$

and the convergence is with rate

$$\min_{k=0, \dots, K-1} \|E'(u^k)\|^2 \leq \frac{1}{K} \sum_{k=0}^{K-1} \|E'(u^k)\|^2 \leq \frac{C^2}{K}.$$

Taking the square root completes the proof.  $\square$

**Remark 5.2.** It follows from the theorem that if  $\{u^k\}_k$  converges, then it converges to a stationary point  $u^* \in \mathcal{U}$  with  $E'(u^*) = 0$ .

# Bibliography

- [1] R. A. Adams and J. J. F. Fournier. *Sobolev Spaces*. Elsevier Science, Singapore, 2003.
- [2] A. B. Bakushinskii. Remarks on the choice of regularization parameter from quasioptimality and relation tests. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 24(8):1258–1259, 1984.
- [3] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. 2011.
- [4] B. Bollobás. *Linear Analysis: An Introductory Course*. Cambridge University Press, Cambridge, second edition, 1999.
- [5] N. Bourbaki. *Topological Vector Spaces*. Éléments de mathématique. Springer-Verlag, 1987.
- [6] K. Bredies and D. A. Lorenz. *Mathematische Bildverarbeitung: Einführung in Grundlagen und moderne Theorie (German)*. Vieweg+Teubner Verlag, 2011.
- [7] I. Ekeland and R. Témam. *Convex Analysis and Variational Problems*. 1976.
- [8] E. Giusti. *Minimal Surfaces and Functions of Bounded Variation*. Birkhaeuser, Basel, Boston, Stuttgart, 1984.
- [9] C. W. Groetsch. *Stable approximate evaluation of unbounded operators*. Springer, 2006.
- [10] J. Hunter and B. Nachtergaele. *Applied Analysis*. World Scientific Publishing Company Incorporated, 2001.
- [11] A. W. Naylor and G. R. Sell. *Linear Operator Theory in Engineering and Science*. Springer Science & Business Media, 2000.
- [12] A. Rieder. *Keine Probleme mit Inversen Problemen: Eine Einführung in ihre stabile Lösung*. Vieweg+Teubner Verlag, 2003.
- [13] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, 1992.
- [14] W. Rudin. *Functional Analysis*. International series in pure and applied mathematics. McGraw-Hill, 1991.
- [15] T. Tao. *Epsilon of Room, One*, volume 1. American Mathematical Soc., 2010.
- [16] E. Zeidler. *Applied Functional Analysis: Applications to Mathematical Physics*, volume 108 of *Applied Mathematical Sciences Series*. Springer, 1995.

- [17] E. Zeidler. *Applied Functional Analysis: Main Principles and Their Applications*, volume 109 of *Applied Mathematical Sciences Series*. Springer, 1995.