

BACHELOR THESIS

APPLICATION OF REGULARISATION  
METHODS FOR A DECONVOLUTION  
PROBLEM WITHIN THE KATRIN  
COLLABORATION

by Henrik Rose

WWU Münster

Matrikel 441476

Kettelerstraße 37

48147 Münster

[h\\_rose10@uni-muenster.de](mailto:h_rose10@uni-muenster.de)

supervised by Prof. Dr. Benedikt Wirth and Dr. Volker Hannen

submitted August 14, 2020



## Abstract

The *Karlsruhe Tritium Neutrino Experiment* (KATRIN) experiment attempts to reach a model independent neutrino mass estimate from kinematic observations of tritium  $\beta$ -decays. A major source of systematic uncertainty are energy losses of decay electrons by inelastically scattering off tritium molecules. It is therefore necessary to obtain information on the differential scattering cross section for this process. The method suggested by Hannen *et al.*[1] applies Truncated Singular Value Decomposition to experimental *scattering functions* obtained from an in-situ experiment. It is extended by means of Tikhonov regularisation, constrained optimisation and a parametrisation approach. Numerical experiments based upon simulations of the KATRIN experiment indicate that these methods allow for an improved reconstruction of the energy loss function that respects the assigned KATRIN error budget even at a higher noise level.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Physical Background</b>	<b>4</b>
2.1	Experimental Realisation . . . . .	6
2.2	Scattering Functions . . . . .	8
<b>3</b>	<b>Mathematical background</b>	<b>10</b>
3.1	Ill-Posedness for Compact Operators . . . . .	11
3.2	Generalised Inversion . . . . .	13
3.3	Discretisation . . . . .	15
<b>4</b>	<b>Regularisation Methods</b>	<b>16</b>
4.1	Truncated Singular Value Decomposition . . . . .	16
4.2	Tikhonov Regularisation . . . . .	19
4.3	L-curve criterion . . . . .	21
4.4	Model-adapted Minimisation and Parametrisation . . . . .	23
<b>5</b>	<b>Numerical Experiments</b>	<b>25</b>
<b>6</b>	<b>Conclusion and Outlook</b>	<b>32</b>
<b>7</b>	<b>Appendix</b>	<b>33</b>

---

# 1 Introduction

The standard model of particle physics summarises long established theories concerning the various elementary particles and their possible interactions. In recent years, experiments from various fields of physics provided strong evidence that it is incomplete, though. Proving physics beyond the standard model has therefore become one of the most intensely researched fields in experimental particle physics. The *Karlsruhe Tritium Neutrino Experiment* (KATRIN) attempts to contribute to these efforts by determining the mass of neutrinos. In the standard model, neutrinos are assumed to be massless, so neutrino mass determination is another import step towards physics beyond the standard model. The idea of the KATRIN experiment is to determine the shape of the energy spectrum of electrons created in tritium  $\beta$ -decays. In its endpoint, it follows a distribution that depends heavily on the energy (or equivalently mass) of neutrinos. Conversely, observing this distribution should in theory allow for a determination of the neutrino mass. However, multiple sources of uncertainty need to be considered in order to attain the aspired sensitivity of  $0.2 \text{ eV}/c^2$  at 90% confidence level. One of them is energy loss of electrons due to inelastic scattering within the experiment's tritium source. Hence, the differential scattering cross section needs to be known with sufficient precision in order to limit energy loss effects in KATRIN evaluation methods. By defining a corresponding energy loss function and scattering functions, Hannen *et al.* showed that this can be achieved by applying Truncated Singular Value Decomposition (TSVD) to data obtained from a KATRIN pre-experiment[1]. However, their method relies on a low noise level, i. e. very long measuring campaigns, and the resulting energy loss function still shows several nonphysical properties, including negative energies and strong oscillatory behaviour. It is therefore desirable to provide an augmented method that yields a more realistic representation of energy losses while respecting the associated KATRIN error budget.

The underlying deconvolution problem can be seen as a typical inverse problem where a well-defined forward operator  $O$ , applied to some input  $x$ , yields an output  $y$ . A measurement  $y^\delta$  of  $y$  is typically corrupted by some noise  $\delta$ , though. Since  $O$  frequently does not possess a well-defined, continuous inverse, it is in most physical applications highly unlikely to recover the true  $x$ . Instead, it is convenient to find some reasonable approximation  $x_\alpha$  to  $x$ . This can be achieved by determining it as

$$x_\alpha = \operatorname{argmin}_{x \in \Omega} \mathcal{D}(x) + \mathcal{R}_\alpha(x), \quad (1.1)$$

where  $\Omega$  is a suitable function space and the functional  $\mathcal{D}(x)$  denotes a solution's fidelity to the measured data.  $\mathcal{R}_\alpha(x)$  refers to the regularisation term that reflects on the solution's accordance with some desired properties.

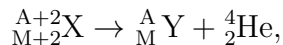
This thesis primarily investigates the choice of a suitable regularisation method for this deconvolution problem and the identification of a good regularisation parameter  $\alpha$ .

After giving a brief overview of the physical context, the meaning and some expected properties of the KATRIN energy loss function will be discussed in section 2. Section 3 provides

the necessary mathematical background to establish the regularisation methods discussed in section 4. The solution offered by TSVD will there be compared in terms of reconstruction quality and parameter choice to Tikhonov regularisation as its most "natural" refinement. Additionally, solutions through parametrisation and constrained minimisation are presented that make stronger model assumptions. Finally, numerical experiments on simulated KATRIN measurements are performed in section 5 that evaluate the expected influence of the different deconvolution methods on the KATRIN neutrino mass estimate.

## 2 Physical Background

The existence of neutrinos was first proposed in 1930 by Wolfgang Pauli in order to provide an explanation for the observation of continuous  $\beta$ -decay spectra that respects conservation of energy – a fundamental principle of physics that others like Bohr were willing to sacrifice at that time in the wake of quantum theory[2]. If  $\beta$ -decays were two particle decays like the  $\alpha$ -decay



conservation of momentum and energy would require them to provide a discrete energy spectrum (and spherically opposed detection locations with respect to the centre-of-mass system).  $\beta$ -decays by contrast, or more precisely the

$$\beta^- \text{-decay: } n \rightarrow p + e^- + \bar{\nu}_e \tag{2.1}$$

as the then only known form of a wider group of processes nowadays summarised by the former term, yield a continuous angular distribution and a rather characteristically shaped electron energy spectrum up to the region of the total decay energy. Although it took another 26 years until the famous poltergeist experiment by Cowan and Reines confirmed their existence by observing inverse  $\beta$ -decay events[3], neutrinos quickly gained acceptance within the scientific community. The corresponding theory has since been massively expanded, predicting and detecting neutrino flavours and corresponding antineutrinos, which were included into a wider theoretical framework commonly referred to as the standard model of particle physics. It associates the three different neutrino flavours with the respective particle generations in which they first occur.

Despite providing an abundance of predictions that could be experimentally validated to high precision, the standard model has proven to be an incomplete theory as it fails to account for some recent discoveries, most notably neutrino oscillations. In contrast to Pauli, who originally expected neutrinos to have masses similar to that of the concurrent  $\beta$ -particle, the standard model assumes neutrinos from all particle generations to be massless. Large-scale detectors like *Super-Kamiokande*, however, display strong evidence for a description of neutrinos or the respective neutrino flavours as superpositions of the different mass eigenstates. The fact that the rate changes, at which neutrinos of a specific flavour are detected, can then

---

be traced back to differences in the particular mass states. Thus detailed observations do not only show that neutrinos have non-vanishing masses but can provide first-order estimates for the mass difference between the three flavours. This is clear evidence for *physics beyond the standard model*; it does not provide absolute values for the neutrino mass, though[4].

Another approach for determining neutrino masses could be pursued by detecting neutrinoless double  $\beta$ -decays. If neutrinos were Majorana fermions, so neutrinos and antineutrinos were in fact identical particles (in contrast to the standard model), a simultaneous decay of two neutrons could occur that only emits two electrons by exchange of virtual neutrinos. In-depth considerations on the decay rates could allow for sensible, but highly model-sensitive mass estimates[5]. Claims of neutrinoless  $\beta$ -decay observations have been made, but met criticism on the statistical methodology and could not be reproduced[6].

KATRIN pursues a different strategy within this highly active field of research by attempting to determine the neutrino mass from purely kinematic, hence model-independent observations. The energy released in a  $\beta$ -decay is shared among the emitted  $\beta$ -particle and its associated neutrino in form of kinetic energy and the respective rest mass energy. A high-resolution measurement of a  $\beta$ -particle's energy spectrum near the total decay energy would provide valuable information because its shape will strongly depend on the neutrino's rest mass energy. While this could be achieved with a variety of decaying nuclei, the tritium  $\beta^-$ -decay



is especially well suited due to its low decay energy of 18.6 keV, as well as its simple atomic structure and tritium molecules' moderate half-time of about twelve years. This decay will from now on be referred to as *TBD* and the comprised electron antineutrino  $\bar{\nu}_e$  will be simply called neutrino ( $\nu$ ).

The shape of the electron's energy spectrum in a TBD can be derived from Fermi's golden rule as

$$\frac{d^2 N}{dE dt} = C \cdot F(E, Z) \cdot p_e E_e p_\nu E_\nu \Theta(E_0 - E - m_\nu c^2). \quad (2.3)$$

It links the electron transmission rate to both the electron's and neutrino's momentum and energy, respectively. The Fermi function  $F$  takes the 'charge value'  $Z = 2$  since a helium nucleus is created in this process. The Heaviside function  $\Theta$  guarantees conservation of the total decay energy  $E_0$ . By means of the relativistic energy-momentum relation, all energy and momentum quantities can be expressed in terms of the electron's kinetic energy  $E$ , the total decay energy  $E_0$  and the respective rest masses. The key idea is that around the spectrum's endpoint ( $E_0 \approx E$ ), the shape strongly depends on the neutrino mass  $m_\nu$  which can thus be recovered by KATRIN. Further modifications to this equation are necessary to account for superposition of the different neutrino flavours' mass states as a result of neutrino oscillations. The neutrino terms therefore need to be actually expressed by an incoherent sum over these states, weighted by the respective population probabilities. For a full derivation see for instance Otten and Weinheimer[7].

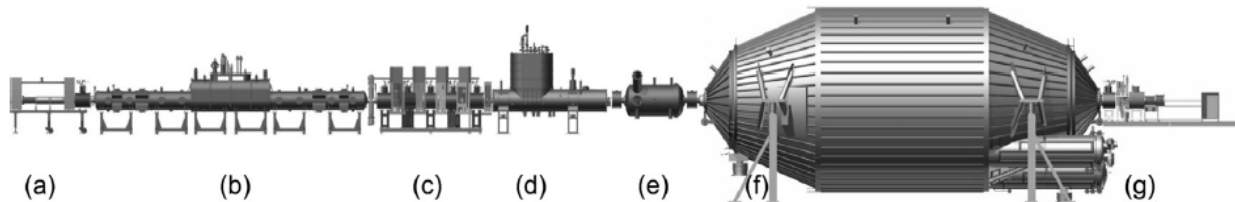


Figure 1: The KATRIN experiment comprises a calibration system (a), a Windowless Gaseous Tritium Source (b), a pumping section (c,d) to limit the amount of tritium in the MAC-E filter type pre- (e) and main spectrometer (f) and an electron detector section (g). From [1].

## 2.1 Experimental Realisation

The key components of KATRIN are depicted in fig. 1. Inside a Windowless Gaseous Tritium Source (WGTS) tritium beta decay occurs isotropically at a sufficiently high rate. Electrons emitted under suitable angles are guided through a pumping section, where tritium molecules are filtered in order to reduce background effects. They then reach the spectrometer section. By means of Magnetic Adiabatic Collimation combined with an Electrostatic filter (MAC-E filtering), the spectrometers can convert transversal energy into longitudinal energy and thus meet the requirements of high energy resolution as well as high luminosity, transmitting essentially all electrons to the detector that are emitted under acceptable angles  $\theta < \theta_{max}$  and with kinetic energies above an adjustable retarding potential. While a pre-spectrometer retains electrons below 18.3 keV, the main spectrometer should achieve the aspired KATRIN energy resolution.

For an experiment of this scale and precision, an abundance of uncertainties and background effects naturally needs to be taken into account. One of these uncertainties is the electrons' energy loss from scattering processes within the WGTS. Because there is only an extremely low share of electrons that are emitted in TBD at the required high energies, the amount of  $T_2$  molecules inside the WGTS needs to be high enough to allow for sufficient luminosity and a feasible duration of measurements. Electrons transiting the WGTS then have a certain probability of interacting with tritium molecules on their path by elastic scattering and various inelastic processes, most importantly excitation of electronic states as well as dissociation and ionisation. These cause different levels of energy loss and systematically reduce the share of electrons passing the MAC-E filters, thus altering the shape of the acquired spectrum.

A suitable correction method therefore requires knowledge on the energy distribution of the respective scattering probabilities. It is useful to describe these in terms of the differential scattering cross section  $\frac{\partial\sigma}{\partial E}$ . Normalisation by the total inelastic scattering cross section  $\sigma_{tot}$  then yields an energy loss function  $f(E)$  that can be interpreted as a probability density distribution. It is not identified to sufficient precision for tritium, though: Aseev *et al.* used data from KATRIN's predecessor *Troitsk nu-mass* for a simple least-squares fit to a model that combines a Gaussian and a Lorentzian curve[8]. Rather sharp energy peaks from



## 2.1 Experimental Realisation

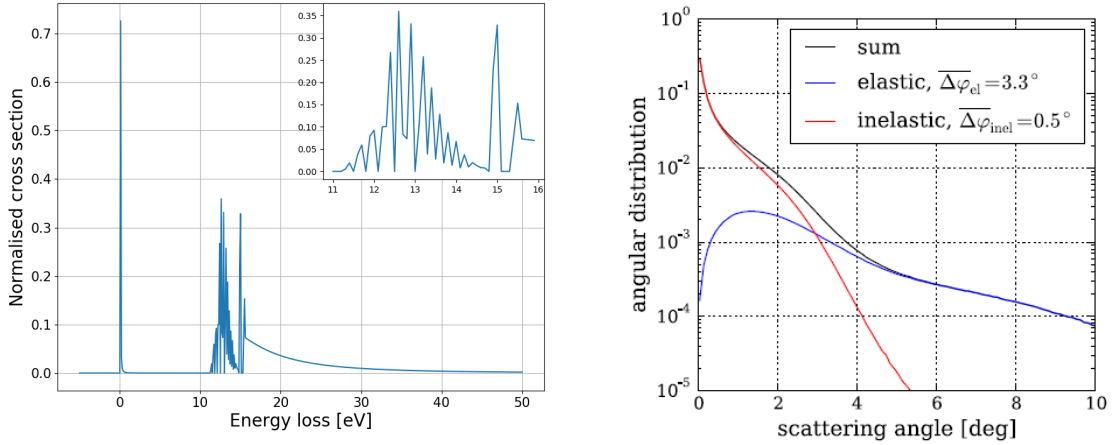


Figure 2: The discrete model energy loss function (left) for 18.6 keV electrons in tritium, based upon work by F. Glück. The widening of sharp peaks is due to the distribution of scattering angles (right, from [1]).

dissociation and electronic state excitation cannot be reconstructed by it due to low resolution and smearing effects. Alternatively, time-of-flight methods have recently been applied in the KATRIN context that can be used to ultimately operate the spectrometers as differential rather than integral filters. This comes at the expense of cutting away huge amounts of data, though, and thus reducing the statistical reliability of this approach. Therefore, the energy loss needs to be evaluated by a KATRIN pre-experiment.

A promising approach is the deconvolution of scattering functions as discussed in the following subsection. In order to determine to which level of precision this is feasible, deconvolution methods should be tested with a sufficiently accurate energy loss model. For this purpose detailed results for molecular hydrogen were assembled by F. Glück and combined to a software routine that constructs a more involved model. The Glück routine generates energy losses and scattering angles satisfying various stochastic expectations that are assembled to the spectral shape in fig. 2 with a binning of 0.1 eV[9]. Note that this shape is adapted to the static laboratory system: The discrete energy losses from elastic scattering as well as electronic excitation and dissociation would require a distributional description in a center-of-mass system. Due to the small recoil energy taken by tritium molecules in the laboratory system, however,  $f$  can be assumed to be continuous with the provided discrete representation. The imperfect angular resolution of KATRIN causes an additional smearing in the actual measurement. Several key features of the energy loss function can be identified:

- As it is a probability density,  $f(E)$  needs to be positive definite and normalised, so  $\int_{\mathbb{R}} f(E)dE = 1$ .
- A large fraction of interactions will be elastic processes, accounting for a sharp peak of

vanishing energy losses, weakly smeared due to recoiling tritium molecules.

- No further energy loss of less than  $\approx 10$  eV should be detected; any measurement that indicates otherwise needs to be attributed to background effects.
- Because electrons are indistinguishable particles, the continuous tail accounting for ionisation energies can take any value up to half the decay energy.
- Hence, the energy loss function  $f$  is a bounded function with compact support on the interval  $[0, E_{\text{tot}}/2]$ . It decays fast, though, as indicated by the BED-model[10], so high-energy loss ionisation events are negligible.

This model can be used to identify and evaluate suitable reconstruction methods.

## 2.2 Scattering Functions

The idea of the previous approach is to extract  $f(E)$  from a *scattering function*[1]. It is defined by convolution of  $f$  with the experimental KATRIN transmission function,

$$\epsilon_1(E) = (T^e * f)(E) = \int T^e(E - \hat{E})f(\hat{E})d\hat{E}. \quad (2.4)$$

It describes how probable it is for an electron to be detected given that it has been (elastically or inelastically) scattered exactly once.

As already suggested in the KATRIN design report[4],  $n$ -fold scattering functions  $\epsilon_n$  can be regained from measurement on the overall KATRIN response function: Consider the use of an electron gun in the calibration section, emitting electrons at 18.6 keV. This is only slightly above the total decay energy, so electrons should behave like electrons from TBD while any effects from actual TBD can be neglected. The observable KATRIN response function  $R(E)$  describes (after normalisation) the probability of an electron to be detected if the energy associated with the set retarding potential differs by  $E$  from the original electron energy. It is measured for multiple particle densities in the WGTS, conveniently expressed as column densities  $\rho d$ . In an empty WGTS, ideal responses merely depend on the filters' theoretical transmission function[8]

$$T(E) = \begin{cases} 0 & \text{if } E < 0 \\ \frac{1 - \sqrt{1 - \frac{E}{E+qU} \frac{B_e}{B_A}}}{1 - \sqrt{1 - \frac{E_{\perp max}}{E+qU} \frac{B_e}{B_A}}} & \text{if } 0 \leq E < E_{\perp max} \\ 1 & \text{if } E_{\perp max} \leq E. \end{cases}$$

Here  $B_e$  denotes the magnetic field at the electron source and  $B_A$  in the spectrometer plane. The quantity  $E_{\perp max}$  is the maximum remaining transversal energy of electrons emitted under  $\theta_{\text{max}}$ , limiting the spectrometer resolution. The measurable experimental transmission function  $T^e$  will mildly suffer from additional smearing effects.

## 2.2 Scattering Functions

---

Scattering off tritium molecules within a filled WGTS will necessarily reduce the value of the response function within the narrow region around the endpoint energy considered in the following. The experimental response function  $R$  can be split into contributions from  $n$ -foldly scattered electrons:

$$R(E) = P_0 \cdot T^e(E) + P_1 \cdot \underbrace{T^e * f(E)}_{\epsilon_1(E)} + \dots + P_n \cdot \underbrace{T^e(*f)^n(E)}_{\epsilon_n(E)} + \dots \quad (2.5)$$

Here  $P_n$  denotes the probability for an electron to be scattered exactly  $n$  times. Regarding the total scattering cross section  $\sigma_{\text{tot}}$  as the inverse of an electron's mean free column density  $\rho d_{\text{free}}$ , scattering events can be considered independent with a scattering rate

$$\lambda = \frac{\rho d}{\rho d_{\text{free}}} = \rho d \sigma_{\text{tot}}.$$

These probabilities are therefore assumed to largely obey a Poissonian distribution, though some adjustments need to be made in order to properly reflect scattering angles and the electrons' angular distribution at emission. Neglecting more than  $n$ -fold scattering,  $n$  additional measurements of  $R(E)$  at non-zero column densities provide a system of linear equations

$$\begin{aligned} R(E)^{(\rho d)_1} &= P_0^{(\rho d)_1} \cdot T^e(E) + P_1^{(\rho d)_1} \cdot \epsilon_1(E) + P_2^{(\rho d)_1} \epsilon_2(E) + \dots + P_n^{(\rho d)_1} \epsilon_n(E) \\ R(E)^{(\rho d)_2} &= P_0^{(\rho d)_2} \cdot T^e(E) + P_1^{(\rho d)_2} \cdot \epsilon_1(E) + P_2^{(\rho d)_2} \epsilon_2(E) + \dots + P_n^{(\rho d)_2} \epsilon_n(E) \\ &\vdots \\ R(E)^{(\rho d)_n} &= P_0^{(\rho d)_n} \cdot T^e(E) + P_1^{(\rho d)_n} \cdot \epsilon_1(E) + P_2^{(\rho d)_n} \epsilon_2(E) + \dots + P_n^{(\rho d)_n} \epsilon_n(E) \end{aligned}$$

that holds for all values of  $E$ . It can be expressed as a matrix equation

$$\vec{R}(E) - T^e(E) \vec{P}_0 = \mathbf{P} \vec{\epsilon}(E). \quad (2.6)$$

Since very small  $n$  are sufficient – Hannen *et al.* argue that  $n = 3$  provides satisfactory results at the aspired KATRIN column density  $\rho d = 5 \times 10^{17} \text{ cm}^{-2}$  – and its entries are of similar order,  $P$  can easily be inverted. This allows for a stable calculation of the single scattering function at reasonable precision. Monte Carlo simulations of the energy loss at various column densities support this, provided measurement campaigns are sufficiently long. The process simultaneously provides the first  $n$  scattering functions as shown in fig. 3. Although it would in theory be possible to use these for a cross validation on the reconstructed energy loss function, noise gets amplified in each convolution step by such amounts that this is not feasible. The method of simulation and the impact of increased noise levels on the observed response functions will be further investigated in section 5.

The extraction of the energy loss function from this scattering function proves to be a much more difficult problem. A trivial approach would be to apply the convolution theorem which

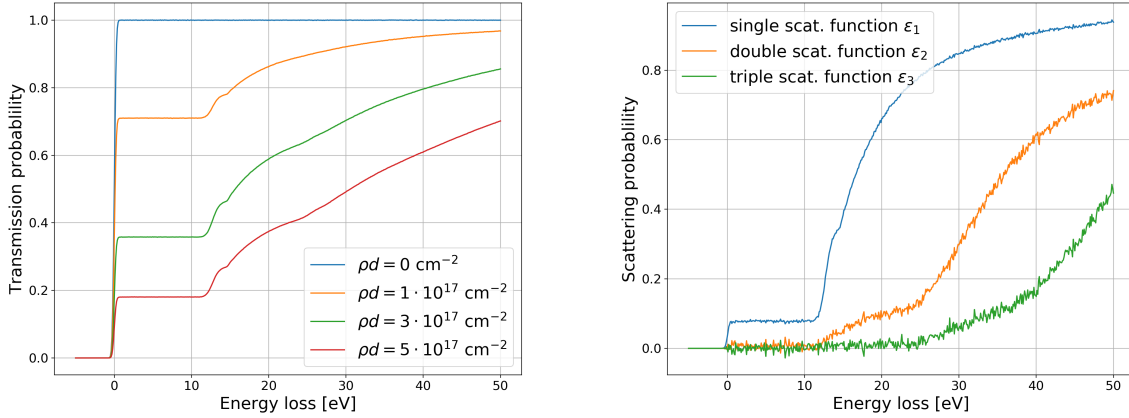


Figure 3: Simulated KATRIN response functions for four different column densities, affected by multiple background effects (left). These can be used to calculate transmission probabilities for up to triply scattered electrons by eq. (2.6) (right).

states that convolution in Euclidean space is equivalent to multiplication in Fourier space. So if  $\mathcal{F}$  denotes the Fourier transformation and  $\mathcal{F}^{-1}$  the inverse Fourier transformation,

$$f(E) = \mathcal{F}^{-1} \left( \frac{\mathcal{F}(\epsilon_1(E))}{\mathcal{F}(T_e(E))} \right)$$

holds. But  $\mathcal{F}(T_e(E))$  is oscillatory and varies strongly in magnitude, so the division operation is numerically unstable and an approach by Fast Fourier Transformation yields arbitrarily large results. The following section outlines why the described behaviour is in fact expected from a mathematical perspective.

### 3 Mathematical background

**Definition 3.1** (Well- and ill-posed problems). A problem is called well-posed, if

- a unique solution to the problem exists, and
- the solution depends continuously on given boundary conditions.

Otherwise the problem is called ill-posed.

This definition is commonly attributed to Hadamard who – building upon previous work – introduced it somewhat implicitly in a lecture series on Cauchy problems[11] where he tried to demarcate problems he deemed not worth a thorough consideration. This perception has

shifted thanks to the huge practical impact of ill-posed problems that very frequently occur in the context of inverse problems. Although no canonical mathematical definition exists, it is convenient to provide the application driven

**Definition 3.2** (Inverse Problem). The problem

Find  $x$  so that for a well-defined (forward) operator  $O$   
and true data  $y$  the equation  $Ox = y$  holds

is called an *inverse problem* if it is ill-posed and has no trivial solution.

While the uniqueness of  $x$  can often be shown (at least in some appropriate sense), the second condition in the above definition is frequently violated, such that exact solutions become useless when only noisy data  $y^\delta$  is available.

By means of functional analysis it is possible to rigorously verify the ill-posedness of eq. (2.4) as well as to establish a notion of singular systems that can be used in a first attempt to overcome it. In order to do so, notations and descriptions are adopted from the classic textbook by Werner[12].

### 3.1 Ill-Posedness for Compact Operators

Equation (2.4) may be identified as a general linear integral equation of the form

$$\int_I k(x, y)u(y)dy = v(x) \quad \forall x. \tag{3.1}$$

It is also known as a Fredholm integral equation of the first kind and can be expressed in terms of an integral operator  $K$  as  $Ku = v$ .

This equation is well-defined if  $k \cdot u$  is integrable for all values of  $x$ . This can be guaranteed by limiting  $k$  and  $u$  to appropriate function spaces.

**Definition 3.3** ( $L^p$ -spaces). For any  $p$  in the interval  $[1, \infty)$ ,

$$\|f\|_{L^p(I)} := \left( \int_I |f|^p \right)^{\frac{1}{p}}$$

defines a semi-norm for any real-valued, measurable function on  $I$ . By the equivalence relation  $f \sim g \Leftrightarrow \|f\|_p = \|g\|_p$  one can define equivalence classes  $[f]$  with representatives  $f$ . Restricted to these representatives, the Lebesgue space  $L^p(I)$  is then the normed vector space defined by

$$L^p(I) := \{f : I \mapsto \mathbb{R} \text{ is measurable with } \|f\|_{L^p} < \infty\}.$$

$L^p$ -spaces are complete and hence Banach spaces; in the case  $p = 2$  it is even a Hilbert space. For real-valued functions the norm is induced by the scalar product

$$\langle f, g \rangle = \int_I f \cdot g$$

which will mostly be the case in the upcoming considerations. If  $k$  and  $u$  are required to be in  $L^2$ , the Hölder inequality directly shows that  $Ku$  is well-defined.

Let for the moment  $X, Y$  denote general Banach spaces.

**Definition 3.4** (Compact Operator). A linear operator  $\tilde{K} : X \rightarrow Y$  is called compact, if for any bounded subset  $S \subset X$  the closure of its image is compact in  $Y$ .

This statement immediately allows for an equivalent definition by

**Lemma 3.5.** *An operator  $\tilde{K} : X \rightarrow Y$  is compact if, and only if, the image  $(\tilde{K}x_n)_{n \in \mathbb{N}}$  of any bounded sequence  $(x_n)_{n \in \mathbb{N}}$  possesses a convergent subsequence.*

Hence it can be shown that  $K$  in eq. (3.1) is a compact operator:

**Lemma 3.6.** *Let  $\tilde{K}_n : X \rightarrow Y$  denote a sequence of bounded linear operators with finite-dimensional image. If there is an operator  $\tilde{K} : X \rightarrow Y$  with  $\tilde{K}_n \rightarrow \tilde{K}$  in the operator norm, then  $\tilde{K}$  is compact.*

This statement follows from the observation that the space  $K(X, Y)$  of compact operators from  $X$  to  $Y$  is a Banach space, too, and that bounded linear operators with finite dimensional image spaces are always compact. Because a finite-dimensional approximation by "step functions"  $k_n \rightarrow k$  can be constructed for  $k \in L^2$  (by calculating weighted averages of  $k$  on increasingly narrow partitions of its domain),  $\tilde{K}$  is a compact operator.

This observation, however, implies that  $K$  cannot be continuously inverted. To see this, consider

**Lemma 3.7** (Riesz). *Let  $U \subsetneq X$  be a closed subspace. For each  $\delta \in (0, 1)$  there is always an  $x_\delta$  satisfying  $\|x_\delta\| = 1$  and  $\|x_\delta - u\| \geq 1 - \delta \quad \forall u \in U$ .*

For  $X$  and  $Y$  with infinite dimension assume now a compact  $\tilde{K} : X \rightarrow Y$  had a continuous inverse. For every bounded sequence  $(x_i)_{i \in \mathbb{N}}$  there would be a subsequence  $(x_{i_k})_{k \in \mathbb{N}}$  with converging image  $(y_{i_k})_{k \in \mathbb{N}}$  in  $Y$  by compactness of  $\tilde{K}$ . Due to continuity of the inverse this would imply that  $(x_i)_{i \in \mathbb{N}}$  already had a converging subsequence. But with  $U_k = \text{span}\{x_1, x_2, \dots, x_{k-1}\}$  Riesz's lemma allows for the construction of a bounded sequence  $(x_k)$  in  $X$  without any converging subsequences. Thus there cannot be a continuous inverse to  $\tilde{K}$ .

Any solution  $x^\delta$  to  $Kx = y^\delta$  will thus most likely be meaningless in the sense that it does not convey any information about the true solution in the noise-free case. The problem should therefore be re-framed to finding an approximate solution

$$x = \operatorname{argmin}_{\tilde{x} \in X} \frac{1}{2} \left\| K\tilde{x} - y^\delta \right\|^2 + \mathcal{R}(\tilde{x}), \tag{3.2}$$

where the latter *regularisation term* penalises undesired properties of  $\tilde{x}$  occurring in  $x^\delta$  and transforms it into a well-posed problem. While there are no general restrictions on the norm, it is often convenient to choose it as induced by the  $L^2$  scalar product. In particular, this choice comes naturally when assuming  $x, y \in L^2$ .

In the context of the KATRIN energy loss function this shows that eq. (2.4) is certainly ill-posed:  $T^e$  and  $f$  are bounded functions on closed intervals and thus  $L^2$ -integrable; the convolution in eq. (2.4) is therefore a well-defined compact operator. Since  $L^p$  spaces have infinite dimension, it has no continuous inverse and the second condition in definition 3.1 is violated.

### 3.2 Generalised Inversion

Apart from the fact that inversion can become discontinuous for a mapping on infinite-dimensional spaces, it often is found that a mapping is not injective in the first place, so an inverse is ambiguous. However, a notion of inversion can be derived by means of a singular value decomposition.

**Definition 3.8** (Spectrum). For a linear operator  $L : X \mapsto Y$  on vector spaces, the spectrum is defined as the set

$$\sigma(L) = \{\lambda \in \mathbb{R} : \lambda \text{Id} - L \text{ has no continuous inverse}\}. \quad (3.3)$$

This definition generalises the notion of the spectrum for finite-dimensional operators;  $\lambda \in \sigma$  is likewise called an eigenvalue.

Applying its adjoint to a linear operator  $L : X \rightarrow Y$  clearly yields a normal operator, in which case one can always find an eigensystem comprising an orthonormal set  $\{x_1, x_2 \dots\} \in X$  that can be expanded to an orthonormal basis and eigenvalues  $\{\lambda_i\} \subset \mathbb{R}$  (or another field over which  $X$  is defined). Based upon Riesz's theorems, it can be shown for any compact operator  $K$  that these eigenvalues are at most countable with 0 as the only possible accumulation point. Because the product of a compact and a linear operator is always compact,  $K^*K$  is a normal, compact operator. Because it is self-adjoint, its eigenvalues are nonnegative. Therefore they form a null sequence that can be represented as  $\sigma_1^2 \geq \sigma_2^2 \geq \dots \rightarrow 0$ . A set of linearly independent, orthonormal vectors is then defined by  $y_i := \frac{Kx_i}{\sigma_i}, \sigma_i \neq 0$  which can be canonically expanded to an orthonormal basis of  $Y$ . From these definitions the following identities arise directly:

$$Kx_i = \sigma_i y_i \quad (3.4)$$

$$K^*y_i = \sigma_i x_i \quad (3.5)$$

$$Kx = \sum_i \sigma_i \langle x, x_i \rangle y_i \quad (3.6)$$

**Definition 3.9** (Singular Value Decomposition). The set  $\{x_i, y_i, \sigma_i\}_{i \in \mathbb{N}}$  is called a *singular system* (to the operator  $K$ ) with the *singular values*  $\sigma_i$  if it satisfies eqs. (3.4) to (3.6).

Equation (3.6) is called the *singular value decomposition (SVD)* of operator  $K$ ; it is a series in case of infinite dimensions.

Assume for the moment that the operator  $K$  is bijective, so all singular values are positive. Since  $x = \sum \langle x, x_i \rangle x_i$ , SVD allows for an inversion of  $Kx = y$  as

$$x = \sum_i \frac{\langle x, \sigma_i x_i \rangle}{\sigma_i} x_i = \sum_i \frac{\langle x, K^* y_i \rangle}{\sigma_i} x_i = \sum_i \frac{\langle y, y_i \rangle}{\sigma_i} x_i. \quad (3.7)$$

However, this solution is only well-defined if the series is convergent (it is then said to satisfy the Picard criterion). This is not generally the case for a compact operator as demonstrated in the previous section. If the amount of singular values is finite, though, or the series is truncated for some  $n$ , a well-defined result is obtained: If  $\sigma_i > 0 \Leftrightarrow i \leq n$ , let  $K_n$  denote the operator defined by

$$K_n x = \sum_{i=1}^n \sigma_i \langle x, x_i \rangle y_i \quad (3.8)$$

on the finite dimensional subspaces  $\text{span}\{x_1, \dots, x_n\} = X_n \subset X$ ,  $\text{span}\{y_1, \dots, y_n\} = Y_n \subset Y$ . Then there are canonical projections for any elements  $x \in X$ ,  $y \in Y$ :

$$\bar{x} = \sum_{i=1}^n \langle x, x_i \rangle x_i \quad \bar{y} = \sum_{i=1}^n \langle y, y_i \rangle y_i. \quad (3.9)$$

These allow for an approximate solution to the original problem by  $K_n \bar{x} = \bar{y}$ . Note that this projection is well-defined for operators between vector spaces of dimensions  $k \neq l$ , too. For  $k, l < \infty$  a generalised inverse is then conveniently defined as the unique solution (by eq. (3.7)) to the problem  $K_m \bar{x} = \bar{y}$  where  $m$  denotes the smallest positive singular value.

If the  $K_n$  are finite-dimensional operators, they can be identified with matrices. In order to simplify the distinction between matrix and operator "language", let the SVD then be given in matrix notation as  $A = U \Sigma V^T$  with orthonormal matrices  $U$  and  $V^T$ , where the columns of  $U$  and rows of  $V$  shall be identified in the finite-dimensional case with basis vectors  $\{y_i\}$  and  $\{x_i\}$ , respectively. The matrix  $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots)$  contains the singular values. Finally, let  $\Sigma^{-1}$  be the matrix obtained by inverting the non-zero entries and subsequently transposing  $\Sigma$ . Then the generalised solution to the matrix equation  $Ax = y$  is given by

$$\bar{x} = V \Sigma^{-1} U^T y =: \tilde{A}^{-1} y. \quad (3.10)$$

$\tilde{A}^{-1}$  is called the generalised inverse of  $A$ .



### 3.3 Discretisation

A matrix description is especially useful if discretisation is necessary. Let  $K$  specifically denote convolution with the kernel  $k$  on an interval  $[a, b]$ , so  $Ku = v$  becomes

$$Ku(x) = \int_a^b k(x-y)u(y)dy = v(x) \quad \forall x. \quad (3.11)$$

If the value of  $k$  is only accessible for a discrete set of arguments  $Z = \{z_1, z_2, \dots, z_n\} \subset \mathbb{R}$ , e.g. experimental measurement points, integration is not well-defined. Instead, a discrete operator  $K^h$  needs to be considered that replaces integration by summation:

$$K^h u^h(x_i) := \sum_{j=1}^{m-1} k^h(z = x_i - y_j) u^h(y_j) \cdot (y_{j+1} - y_j) = v^h(x_i) \quad \text{for } i = 1, \dots, l. \quad (3.12)$$

The superscript  $h$  indicates the discrete setting which can only use/provide knowledge on  $v$  or  $u$  in discrete arguments, too.

If the  $y$  have equal step width  $\Delta y$  (and consequently the  $x$  and  $z$ , too), eq. (3.12) can be expressed – using physical vector notation for clarity – as a matrix equation  $\vec{v}^h = \Delta y \cdot \mathbf{K}^h \vec{u}^h$  with

$$\mathbf{K}^h = \begin{pmatrix} k^h(z_l = x_l - y_{m-1}) & k^h(z_{l-1}) & \dots & k^h(z_2 = x_2 - y_{m-1}) & k^h(z_1) \\ k^h(z_{l+1} = x_l - y_{m-2}) & k^h(z_l) & \dots & k^h(z_3 = x_2 - y_{m-3}) & k^h(z_2) \\ & \ddots & & \ddots & \\ k^h(z_{l+m-2} = x_l - y_2) & k^h(z_{l+m-3}) & \dots & k^h(z_{m-1} = x_2 - y_2) & k^h(z_{m-2}) \\ k^h(z_{l+m-1} = x_l - y_1) & k^h(z_{l+m-2}) & \dots & k^h(z_m = x_2 - y_1) & k^h(z_{m-1}) \end{pmatrix}. \quad (3.13)$$

Note that since these arguments are subject to the condition

$$x_i - y_j \in Z \quad \forall i, j \quad (3.14)$$

(otherwise  $k^h$  had no defined value),  $l+m-1 \leq n$  must hold. Thus, the goal of high precision of the solution  $u^h$ , which is based on knowledge on the boundary condition  $v^h$  and therefore relies on a large  $m$ , conflicts with the goal of high generality of the solution, i.e. a large  $l$ . In the context of the KATRIN energy loss this would mean that data for the single scattering function needed to be discarded – by eq. (2.6) it has the same domain as the experimental transmission function  $T^e$ . But this is not necessary when instead  $T^h(E) \equiv T(E)$  is assumed as a reliable model where no measurement data is available. Since measurement will be performed at  $N + M + 1$  retarding potentials  $U_i$ ,  $i = -N, \dots, M$ , eq. (2.4) is expressed by the matrix equation

$$\vec{\epsilon}_1 = \mathbf{T}^e \vec{f} e \Delta U \quad \text{with} \quad (3.15)$$

$$\mathbf{T}^e = \begin{pmatrix} T_0^e & T_{-1}^e & \cdots & T_{-N+1}^e & T_{-N}^e & 0 & 0 & \cdots & 0 \\ T_1^e & T_0^e & T_{-1}^e & \cdots & T_{-N+1}^e & T_{-N}^e & 0 & \cdots & 0 \\ \vdots & & \ddots & & \ddots & & \ddots & & \vdots \\ T_{M-1}^e & \cdots & T_1^e & T_0^e & T_{-1}^e & \cdots & T_{-N+1}^e & T_{-N}^e & 0 \\ T_M^e & T_{M-1}^e & \cdots & T_1^e & T_0^e & T_{-1}^e & \cdots & T_{-N+1}^e & T_{-N}^e \\ 1 & T_M^e & T_{M-1}^e & \cdots & T_1^e & T_0^e & T_{-1}^e & \cdots & T_{-N+1}^e \\ \vdots & & \ddots & & \ddots & & \ddots & & \vdots \\ 1 & \cdots & 1 & T_M^e & T_{M-1}^e & \cdots & T_1^e & T_0^e & T_{-1}^e \\ 1 & \cdots & 1 & 1 & T_M^e & T_{M-1}^e & \cdots & T_1^e & T_0^e \end{pmatrix}. \quad (3.16)$$

Equation (3.13) represents a finite-dimensional approximation to the convolution operation; it is associated with a bounded operator of finite rank. Increasing the amount of measurement points would expand it and thus define a family of linear operators in the sense of lemma 3.6. They are continuous, but the difference in magnitude between the smallest and largest singular value will quickly increase. Therefore the ill-posedness of eq. (3.1) makes a more precise matrix approximation to it ill-conditioned, inflating rounding and data errors in numeric computation. Consequently, the finite-dimensional discretisation requires some regularisation method, too, and eq. (3.12) is restated to search for

$$u^h = \operatorname{argmin} \frac{1}{2} \left\| \mathbf{K}^h \tilde{u}^h \Delta y^h - v^h \right\|^2 + \mathcal{R}^h(\tilde{u}^h) \quad (3.17)$$

where  $\mathcal{R}^h$  denotes a discrete regularisation term. The discrete norm is understood as the euclidean 2-norm which is the canonical choice for a discrete analogue to the  $L^2$ -norm.

## 4 Regularisation Methods

Let again  $A$  denote any real matrix and let  $Ax = b$  be a corresponding matrix equation while the letter  $K$  shall be reserved for operators.

### 4.1 Truncated Singular Value Decomposition

Just like the divergence of eq. (3.7), the instability of a solution for the matrix equation  $Ax = U\Sigma V^T x = b$  gained by SVD,

$$\tilde{x} = \sum_{i=1}^n \frac{\langle u_i, b \rangle}{\sigma_i} v_i, \quad (4.1)$$

mostly originates from or is amplified by small singular values. It is therefore reasonable to assume that a better approximation is achieved by truncating the sum once singular values

## 4.1 Truncated Singular Value Decomposition

---

fall below a specific threshold, i.e. the regularisation parameter  $\alpha$ . This approach is therefore called *Truncated Singular Value Decomposition* (TSVD).

Although compelling in its simplicity, it is unclear how the parameter  $\alpha$  needs to be chosen and to what extent results gained by this approach can be considered reliable. An answer to the latter question is given by

**Theorem 4.1** (TSVD as Minimisation). *The Truncated Singular Value Decomposition solution can be identified with a least norm solution to the minimisation problem*

$$x_{\alpha,TSVD} = \underset{x \in X}{\operatorname{argmin}} \quad \|A_k x - b\|^2 \quad \text{subject to} \quad (4.2)$$

$$A_k = \underset{\tilde{A} \in \mathbb{R}^{n \times n} : \operatorname{rk}(A_k) \leq k}{\operatorname{argmin}} \quad \left\| \tilde{A} - A \right\|^2. \quad (4.3)$$

*Proof.* Let  $A = U\Sigma V^T$  denote the standard singular value decomposition of a matrix and let  $\Sigma_k$  denote the diagonal matrix containing the first  $k$  singular values. Firstly,

$$A_k := U\Sigma_k V^T = \underset{\operatorname{rk}(B) \leq k}{\operatorname{argmin}} \|A - B\| \quad (4.4)$$

is the best approximation to  $A$  with at most rank  $k$ . This was first shown in the Frobenius norm[13], but the statement holds likewise for any norm invariant under unitary transformation as proven by Mirsky[14]. For the spectral norm as the induced matrix norm by the standard scalar product, this is especially easy to see: While  $\|A - A_k\|^2 = \sigma_{k+1}^2$ , for any  $\tilde{A}$  of same rank there is a normalised vector  $x \in \ker \tilde{A} \cap \operatorname{span}\{v_i\}_{i=1, \dots, k+1}$  so that

$$\begin{aligned} \left\| A - \tilde{A} \right\|^2 &\geq \left\| (A - \tilde{A})x \right\|^2 = \|Ax\|^2 \\ &= \sum_{i=1}^{k+1} \sigma_i^2 \|x_i\|^2 \geq \sigma_{k+1}^2 \sum_{i=1}^{k+1} \|x_i\|^2 = \sigma_{k+1}^2. \end{aligned} \quad (4.5)$$

Secondly,  $x_k := x_{\alpha,TSVD}$  is a best approximation because for any  $\tilde{x}$

$$\begin{aligned} \|A_k \tilde{x} - b\|^2 &= \|A_k(\tilde{x} - x_k) + A_k x_k - b\|^2 \\ &= \|A_k(\tilde{x} - x_k)\|^2 + \|A_k x_k - b\|^2 \\ &\quad + 2 \underbrace{\langle A_k(\tilde{x} - x_k), A_k x_k - b \rangle}_{= \langle \tilde{x} - x_k, A_k^T A_k x_k - A_k^T b \rangle} \geq \|A_k x_k - b\|^2. \end{aligned} \quad (4.6)$$

This is true since the last summand vanishes by definition of  $x_k$ . Finally, knowing  $A_k$  has no full rank, this statement also holds for any  $\tilde{x}_k = x_k + \hat{x}_k$  where  $\hat{x}_k \in \ker(A_k)$ . But  $x_k$  is the least norm solution to eq. (4.2) due to the orthogonality of  $\{v_{k+1}, \dots, v_n\}$ , a generating set of  $\ker(A_k)$ .  $\square$

The proof works analogously in the more general case of compact operators, too. It needs to be emphasised that in eqs. (4.2) and (4.3) two subsequent minimisations are performed, so

$$\|A_k x_k - Ax\| \neq \min_{\text{rk}(\tilde{A})=k, \tilde{x} \in X} \|\tilde{A}\tilde{x} - Ax\| \quad (4.7)$$

if  $k \neq n$ . The theorem underlines that solutions presented by this method are, although numerically stable, sensible to noise on the the input.

An improved solution can therefore be obtained by smoothing the input data. For the KATRIN energy loss, for instance, the single scattering function's interpretation as a probability distribution can be used. This requires the true scattering function to be monotonously increasing, so if noise on the transmission function measurements described in section 2 is additive with mean 0,

$$\tilde{\epsilon}_1 = \underset{\epsilon' \geq 0}{\text{argmin}} \|\epsilon - \epsilon_1\| \quad (4.8)$$

provides an improved approximation to the true scattering function.

Smoothing data is also useful for a preliminary answer to the open question of appropriate parameter choices. A reasonable a posteriori estimate for the regularisation parameter can be derived from the discrepancy principle, popularised by Morozov[15] and therefore frequently attributed to him:

**Definition 4.2** (Morozov's discrepancy principle). Let  $O$  denote a forward operator and  $y^\delta$  the data resulting from a measurement of some true data  $y$ , corrupted by noise  $\delta$ . The regularisation parameter  $\alpha(y^\delta, \delta)$  shall be chosen such that for a reconstruction of input data  $x_\alpha$

$$\|O(x_\alpha) - y^\delta\| \approx \delta. \quad (4.9)$$

The idea behind this definition is that no reconstruction  $x_\alpha$  can be expected to provide a lower noise level in the result than what is obtainable from an actual observation. This requires some knowledge about the noise level present in the measurement, which is well estimated by a comparison of the original and the smoothed data. For the discrete KATRIN setting, eq. (4.9) suggests to use an  $\alpha$  that satisfies

$$\|T^e f_\alpha - \epsilon_1\|_2 \approx \|\tilde{\epsilon}_1 - \epsilon_1\|_2. \quad (4.10)$$

This leaves plenty of room for interpretation, though. Since the above error estimate assumes noise to be additive with 0 mean, it necessarily results in an underestimation of the actual noise level. While this makes it plausible to chose a rather generous  $\alpha$ , no strict relation can be derived. Eventually, the discrepancy merely provides a reasonable order of magnitude for the regularisation parameter, subsequent adjustment is often necessary and therefore susceptible to expectation biases.

## 4.2 Tikhonov Regularisation

Another approach is given by the so called Tikhonov regularisation. By discarding singular values below the regularisation threshold, the TSVD solution to the issue of high-frequency noise in the equation  $Kx = y$  implies that any high frequency data *is* noise. The method comes hence at the expense of discarding any information on ground truth associated with these small singular values.

Lavrentiev's method resolves this by addressing an alternative equation  $(K + \alpha I)x = y$ , shifting all singular values by a constant parameter  $\alpha$ , so

$$x \approx x_{\alpha, \text{Lav}} = \sum_{i=1}^n \frac{1}{\alpha + \sigma_i} \langle y_i, y \rangle x_i. \quad (4.11)$$

This would allow for a preservation of data while guaranteeing small rounding errors in machine division. The influence of this shift on medium singular values, however, is fairly substantial. It can be reduced when following a method suggested by Tikhonov[16]:

$$x \approx x_{\alpha, \text{Tikh}} = \sum_{i=1}^n \frac{\sigma_i}{\alpha + \sigma_i^2} \langle y_i, y \rangle x_i. \quad (4.12)$$

Equation (4.12) can be rewritten in terms of a minimisation problem:

**Theorem 4.3.** *The vector  $x_\alpha$  defined by eq. (4.12) is a minimiser of the functional*

$$E_\alpha(x) := \frac{1}{2} \|Kx - y\|^2 + \frac{\alpha}{2} \|x\|^2. \quad (4.13)$$

*Proof.*  $E_\alpha$  is minimised by a solution of

$$K^*Kx + \alpha x = K^*y, \quad (4.14)$$

because for any feasible  $\tilde{x}$

$$\begin{aligned} E_\alpha(\tilde{x}) &= \frac{1}{2} \|K\tilde{x} - y\|^2 + \frac{\alpha}{2} \|\tilde{x}\|^2 \\ &= \frac{1}{2} \|K(\tilde{x} - x) + Kx - y\|^2 + \frac{\alpha}{2} \|(\tilde{x} - x) + x\|^2 \\ &= E_\alpha(x) + \frac{1}{2} \|K(\tilde{x} - x)\|^2 + \frac{\alpha}{2} \|\tilde{x} - x\|^2 \\ &\quad + \alpha \langle \tilde{x} - x, x \rangle + \underbrace{\langle \tilde{x} - x, K^*Kx - K^*y \rangle}_{\stackrel{4.14}{=} \langle \tilde{x} - x, -\alpha x \rangle} \geq E_\alpha(x). \end{aligned}$$

Using eqs. (3.4) to (3.6) on the other hand yields

$$\begin{aligned} K^* K x_\alpha + \alpha x_\alpha &\stackrel{4.12}{=} \sum_{i=1}^{\infty} \frac{\sigma_i}{\alpha + \sigma_i^2} \langle y_i, y \rangle \underbrace{K^* K x_i}_{=\sigma_i^2 x_i} + \sum_{i=1}^{\infty} \frac{\sigma_i}{\alpha + \sigma_i^2} \langle y_i, y \rangle x_i \\ &= \sum_{i=1}^{\infty} \sigma_i \langle y_i, y \rangle x_i = \sum_{i=1}^{\infty} \langle y_i, y \rangle K^* y_i = K^* y. \end{aligned}$$

Thus  $x_\alpha$  is in fact a solution to eq. (4.14) and therefore minimises  $E_\alpha$ . □

This shows that the Tikhonov regularisation is a minimisation problem in the style of eq. (3.2). The solution gained through singular value decomposition balances data fidelity and the solution's norm, thus disallowing strong oscillations. It is hence a reasonable choice if smooth solutions are expected. The identification of the Tikhonov solution as

$$x_{\alpha, \text{Tikh}} = \operatorname{argmin} E_\alpha(x) \tag{4.15}$$

allows for a more intuitive understanding of its deficiencies, too. As Tikhonov regularisation imposes a penalty on the solution's norm, it is systematically biased towards "smaller" solutions.

In comparison to a TSVD solution it can be expected that Tikhonov solutions yield less oscillations, but fail to recover sharp peaks, e.g. for elastic scattering in the TBD energy loss model. In more general terms, any regularisation is called a Tikhonov regularisation, if it uses a regularisation functional

$$\mathcal{R}_\alpha(x) = \frac{\alpha}{2} \|Dx\|^2 \tag{4.16}$$

with an operator  $D$  - commonly a differential operator. In case a solution gained by direct Tikhonov regularisation as above is oscillatory, this will naturally only provide an improvement with respect to a constant offset (or a low degree polynomial for higher derivatives).

Another motivation for Tikhonov regularisation is based on the analysis of noise in a measurement. By the Bayesian interpretation of inverse problems, outlined for instance by Tarantola[17], it is seen that Tikhonov regularisation is associated with Gaussian noise:

**Theorem 4.4.** *For a measurement  $y^\delta$  (via operator  $K$ ) let the true cause  $x$  be a realisation of a  $\mathcal{N}(0, \sigma_x^2)$ -distributed random variable  $X$  and let  $y^\delta$  be described by a  $\mathcal{N}(y, \sigma_\delta^2)$ -distributed random variable  $Y$ . Then the maximum a posteriori (MAP) estimator  $x_{MAP} := \operatorname{argmax} P(X = x | Y = y^\delta)$  can be expressed in the style*

$$x_{MAP} = \operatorname{argmin} \left\| Kx - y^\delta \right\|_2^2 + \alpha \|x\|_2^2. \tag{4.17}$$

The idea behind this theorem is that from a stochastic point of view the acceptance of a specific solution  $x$  to cause the observation  $y$  implies a judgement on how probable/reliable

the observation itself is. A suitable solution should therefore not only reflect on how well it reproduces an observation, but also on how likely that solution is to occur in itself. This is expressed by the MAP estimator that applies *a posteriori* knowledge concerning the actual observation.

*Proof.* By Bayes' theorem, the conditional probability for  $x$  to be the correct solution, given the observed data  $y^\delta$ , can be expressed in terms of the unconditional probabilities  $P(X = x)$ ,  $P(Y = y^\delta)$  and the data probability  $P(Y = y^\delta|X = x)$ . Monotony of the logarithm implies

$$x_{MAP} = \operatorname{argmax}_x P(X = x|Y = y^\delta) = \operatorname{argmax}_x \log P(X = x|Y = y^\delta) \quad (4.18)$$

$$\stackrel{\text{Bayes}}{=} \operatorname{argmin}_x -\log P(Y = y^\delta|X = x) - \log P(X = x) + \log P(Y = y^\delta) \quad (4.19)$$

$$= \operatorname{argmin}_x c_1 \left\| Kx - y^\delta \right\|^2 + c_2 \|x\|^2 \quad (4.20)$$

$$= \operatorname{argmin}_x \left\| Kx - y^\delta \right\|^2 + \frac{c_2}{c_1} \|x\|^2 \quad (4.21)$$

where the penultimate equality follows from the normal distribution of  $X$  and  $Y$ .  $\square$

The variances are often not sufficiently well known, so just like in the case of the Truncated Singular Value Decomposition it is necessary to choose the regularisation parameter consistently.

### 4.3 L-curve criterion

In theorem 4.3 the balance between residual norm  $R(\alpha) := \|Kx_\alpha - y^\delta\|^2$  and solution norm  $N(\alpha) := \|x_\alpha\|^2$  is mediated by the regularisation parameter. Choosing  $\alpha$  very small will severely underregularise the solution; the damping effect of  $\alpha$  in eq. (4.12) causes almost constant solutions if it is selected too large. A suitable regularisation parameter can then be identified by how well it balances the conflicting requirements imposed by the regularisation method. As had been observed earlier, but established beyond heuristics by Hansen[18], a logarithmic  $N(\alpha)$  vs.  $R(\alpha)$  plot takes in many cases a characteristic L-shape. The kink in the plot shows where a good balance between data fidelity and smoothness in the solution is obtained; the *L-curve criterion* therefore simply suggests to select the  $\alpha$  associated with it. For a sufficiently well-behaved L-curve – notorious counterexamples can be constructed – this corner can be specified as the instance where the curvature is maximised. Denoting  $\hat{R} = \log R(\alpha)$ ,  $\hat{N} = \log N(\alpha)$ , it is generally defined as

$$\kappa(\alpha) = \frac{\hat{R}'\hat{N}'' - \hat{R}''\hat{N}'}{(\hat{R}'^2 + \hat{N}'^2)^{\frac{3}{2}}} \quad (4.22)$$

A numerical calculation of this expression is very expensive, though. Following an observation by Vogel[19] (with slightly different notations and hence different formula), however, it can be greatly simplified in the case of Tikhonov regularisation:

**Theorem 4.5.** *For a Tikhonov-regularised problem, the L-curve's curvature is given by*

$$\kappa = RN \frac{RN/|N'| + \alpha(R + \alpha N)}{(R^2 + \alpha^2 N^2)^{\frac{3}{2}}} \quad (4.23)$$

*Proof.* One may first note that by means of the Singular Value Decomposition and eq. (4.12)

$$N(\alpha) = \|x_\alpha\|^2 \stackrel{4.12}{=} \left\| \sum_{i=1}^{\infty} \frac{\sigma_i}{\sigma_i^2 + \alpha} \langle y_i, y \rangle x_i \right\|^2 \stackrel{\langle x_i, x_j \rangle = \delta_{i,j}}{=} \sum_{i=1}^{\infty} \frac{\sigma_i^2}{(\sigma_i^2 + \alpha)^2} \langle y_i, y \rangle^2 \quad (4.24)$$

$$\begin{aligned} R(\alpha) &= \|Kx_\alpha - y\|^2 = \left\| \sum_{i=1}^{\infty} \frac{\sigma_i}{(\sigma_i^2 + \alpha)} \langle y_i, y \rangle Kx_i - \langle y_i, y \rangle y_i \right\|^2 \\ &\stackrel{3.4}{=} \left\| \sum_{i=1}^{\infty} \left(1 - \frac{\sigma_i^2}{(\sigma_i^2 + \alpha)}\right) \langle y_i, y \rangle y_i \right\|^2 = \sum_{i=1}^{\infty} \frac{\alpha^2}{(\sigma_i^2 + \alpha)^2} \langle y_i, y \rangle^2. \end{aligned} \quad (4.25)$$

From this it becomes obvious that

$$\frac{dR}{d\alpha} = -\alpha \frac{dN}{d\alpha} \text{ and consequently} \quad (4.26)$$

$$\frac{d^2 R}{d\alpha^2} = -\frac{dN}{d\alpha} - \alpha \frac{d^2 N}{d\alpha^2}. \quad (4.27)$$

Furthermore making use of  $\hat{R}' = R'/R$  and  $\hat{N}' = N'/N$ , this can be expressed as

$$\kappa = \frac{\frac{R'}{R} \frac{N'' N - N'^2}{N^2} - \frac{R'' R - R'^2}{R^2} \frac{N'}{N}}{\left(\frac{R'}{R} + \frac{N'}{N}\right)^{\frac{3}{2}}} \quad (4.28)$$

$$= \frac{-\frac{\alpha N' R}{R^2} \frac{N'' N - N'^2}{N^2} + \frac{\alpha^2 N'^2 + R N' + \alpha R N''}{R^2} \frac{N' N}{N^2}}{(\alpha^2 N^2 + R^2)^{\frac{3}{2}} \frac{|N'|^3}{R N^3}} \quad (4.29)$$

$$= RN \frac{RN/|N'| + \alpha(R + \alpha N)}{\alpha^2 N^2 + R^2}. \quad (4.30)$$

This shows the desired result.  $\square$

The reduction in computational effort by this approach is obvious. Unfortunately, explicit Singular Value Decomposition might become prohibitively expensive in large scale implementations. As mentioned by Hansen[20], the result is still useful where Tikhonov solutions are



obtained as minimisers of eq. (4.14), though. Since

$$N' = \frac{2}{\alpha} \langle z_\alpha, x_\alpha \rangle, \text{ where } z_\alpha = \sum_i^\infty \frac{\alpha \sigma_i}{(\sigma_i^2 + \alpha)^2},$$

it follows analogously to theorem 4.3 that

$$z_\alpha = \operatorname{argmin} \frac{1}{2} \|Kz - (Kx_\alpha - y)\|^2 + \frac{\alpha}{2} \|z\|^2, \quad (4.31)$$

so the same minimisation algorithm as for the original Tikhonov solution can be applied, using the residual from the obtained solution as the new data term.

The method again has certain disadvantages, most notably that it hardly profits from (or even breaks at) a reduced noise level: While a reasonable regularisation method should allow for

$$\|Kx_\alpha - y\| \xrightarrow{\alpha \rightarrow 0} 0, \quad (4.32)$$

that is to say it should reproduce ground truth in the limit of no regularisation, an ideal method of parameter identification for  $\alpha = \alpha(y^\delta, \delta)$  should yield

$$\|Kx_\alpha - y\| \xrightarrow{y^\delta \rightarrow y} 0. \quad (4.33)$$

As shown by Bakushinskii[21], however, this cannot be achieved when  $\alpha$  is chosen irrespective of the noise level  $\delta$ . Hanke[22] also constructed examples that illustrate well how the quality of Tikhonov solutions identified by the L-curve criterion is reduced with respect to ground truth for low noise on the measurement data. In case of moderately high noise, however, a fair quality of solution is obtained.

## 4.4 Model-adapted Minimisation and Parametrisation

While the above shows that Tikhonov regularisation can be expected to yield more stable results than TSVD, this direct method cannot meet the expected features of the KATRIN energy loss function: Favouring small solutions, the sharp peaks in the model from elastic scattering and electronic excitation/dissociation events are fully lost, as will be visualised in section 5. Additionally, compensatory negative values are introduced that lack a valid physical interpretation. This section therefore focuses on strategies to overcome the described issues. Knowledge about the expected result can be used to define a model-adapted constrained minimisation:

$$\begin{aligned} x &= \operatorname{argmin} \frac{1}{2} \|K\tilde{x} - y\|^2 \\ &\text{subject to } \operatorname{supp}(\tilde{x}) \subset [a, b] \\ &\quad p < \tilde{x} < r \\ &\quad s < \tilde{x}' < t \end{aligned} \quad (4.34)$$

where  $p, r, s, t$  are suitable bounds and constraints to  $\tilde{x}$  and its derivative. This could be expanded by other constraints. In case of the energy loss function, a valid formulation is

$$\begin{aligned}
 f = \operatorname{argmin} & & \frac{1}{2} \left\| T^e \tilde{f} - \epsilon_1 \right\|^2 \\
 \text{subject to} & & \operatorname{supp}(\tilde{f}) \subset [0, E_{\text{tot}}/2] \\
 & & \tilde{f} \geq 0 \\
 & & \frac{\tilde{f}}{dE} < 0 \quad \forall E \in [0, E_{\text{tot}}/2] \setminus (11 \text{ eV}, 16 \text{ eV}).
 \end{aligned} \tag{4.35}$$

Sufficiently efficient algorithms to tackle large-scale problems of this sort exist.

Alternatively, it can be worthwhile to enhance the previous parametrisation approach by Aseev *et al.*[8]: The energy loss function can be modelled as a combination of a (truncated) Lorentzian curve for the ionisation tail and a Gaussian curve each for elastic scattering, excitation of electronic states and dissociation, respectively. Further enforcing truncation on these provides a parametrisation that will satisfy all of the requirements defined in section 2.1:

$$\begin{aligned}
 f_p(x) = & \Theta(x) \cdot a_{el} \exp\left(-\frac{(x - x_{el})^2}{w_{el}}\right) + a_{exc} \exp\left(-\frac{(x - x_{exc})^2}{w_{exc}}\right) \\
 & + a_{diss} \exp\left(-\frac{(x - x_{diss})^2}{w_{diss}}\right) + \Theta(x - x_{ion}) \cdot \frac{a_{ion}}{(x - b_{ion})^2 + w_{ion}^2}.
 \end{aligned} \tag{4.36}$$

Even higher precision is possible if superpositions and combinations of multiple Gaussians or Lorentzians are allowed. This construction based upon model assumptions suffers from inherent bias towards the specific model, though. In addition to that, raising the amount of free parameters makes the parametrisation process itself increasingly ill-posed. So while it can be refined arbitrarily towards a high-precision solution for the hydrogen-based model problem, the same might fail when applied to tritium. One should therefore refrain from pursuing deceptively fine approximations.

A parametrisation can also be useful in order to counteract the deficiencies of Tikhonov regularisation. Recall that it favours solutions with small norm. When strong peaks are expected, these cannot be accurately recovered. By contrast, consider a decomposition of an ill-posed problem in the form  $Kx = y$  into a

$$\text{well-posed problem } K\hat{x} = \hat{y} \tag{4.37}$$

$$\text{and an ill-posed residual problem } K\partial x = \partial y, \quad \partial y = y - \hat{y}, \tag{4.38}$$

yielding the solution  $x = \hat{x} + \partial x$ . If necessary, the operator  $K$  could also be replaced in one of the above equations by a suitable approximation. Finding a preliminary parametrisation with far fewer (and non-redundant) free parameters than data points will usually be well-posed. A very rough parametrisation can then yield a residual problem with a smaller solution. As a matter of fact, it is crucial to use only a mediocre parametrisation: Since it hardly

---

reduces noise, the noise level in the residual problem increases. For an ideal parametrisation that reproduces ground truth, the residual problem would simply be noise and the Tikhonov solution would become completely meaningless. In consequence, this approach mainly makes sense when the noise level is well under control, but a strong parametrisation is difficult to obtain.

## 5 Numerical Experiments

The overall aim for an appropriate deconvolution method within the context of the KATRIN energy loss function is to provide an approximation to the true energy loss function  $f$  such that systematic errors in the ultimate neutrino mass estimate do not exceed  $0.0075 \text{ eV}^2$ , the foreseen share of the KATRIN error budget[4]. Therefore any solution derived through one of the methods discussed in section 4 should not only possess the expected characteristics noted in section 2.1, but also reliably satisfy this error limit.

In order to evaluate the expected level of systematics, Monte Carlo experiments are performed. The routine was already applied by Hannen *et al.* and is based on work from the KATRIN predecessor experiment in Mainz. It is contained in the KEloss software package[9]. 1000 KATRIN measurements are simulated that provide a spectral shape according to eq. (2.3), modified by several experimental influences. These include physical effects like the final state distribution of TBD, as well as statistical fluctuations and up to four-fold electron scattering according to the original Glück model. For sake of simplicity,  $m(\bar{\nu}_e)^2=0.2 \text{ eV}^2$  is assumed in all simulations. Using the squared neutrino mass, the endpoint energy, the background rate and the spectrum's amplitude as free parameters, the obtained data is then fitted to eq. (2.3). In this routine the same physical effects are applied as before with the exception of the energy loss model. It is now replaced by a deconvoluted energy loss according to the aforementioned methods. The mean of the thereby obtained values for  $m(\bar{\nu}_e)^2$  differs by some  $\mu$  from the input; this is ultimately assumed to provide the desired estimate for the uncertainty due to energy losses by scattering events. The uncertainty of the mean is for all simulations  $\pm 5 \times 10^{-4} \text{ eV}^2$ ; it will be omitted for simplicity.

For the purpose of this thesis, three variants of the discrete single scattering function (SSF) are considered:

- The high-fidelity estimate  $\epsilon_{hf}^h$  was already used by Hannen *et al* and has been displayed in fig. 3. It is gained from a KATRIN-adapted simulation routine of the transmission process, assuming observation of 10 million electrons per energy step and column density[9, 1].
- A noisy estimate  $\epsilon_n^h$  is constructed by adding  $\mathcal{N}(0, 0.001)$ -distributed noise to  $\epsilon_{hf}^h$ . It could result from shorter calibration measurements or unexpectedly strong background effects.

## 5 NUMERICAL EXPERIMENTS

- The (almost) ideal  $\epsilon_{id}^h$  is obtained by convolution of the Glück model  $f_{mod}^h$  with the simulated transmission function from the high-fidelity estimate.

Simulations of the transmission function are performed for 551 evenly spaced retarding potentials that span the range of 18550-18605 keV. Any deconvolution has been performed with respect to that domain, using the very robust algorithms provided by the *scipy.optimize* library[23] in Python. The uncertainty estimates from the above described evaluation methods are displayed in table 1.

Hannen *et al.* applied TSVD to  $\epsilon_{hf}^h$ ; their results are shown in fig. 4. They found an optimal regularisation parameter for TSVD as  $\alpha=0.3\%$  of the maximum singular value both by minimising  $\|f_{\alpha}^h - f_{mod}^h\|^2$  and by the above evaluation method, yielding a satisfying  $\mu = 0.0053 \text{ eV}^2$ . (Note that an optimisation of the regularisation parameter with respect to the evaluation method itself would be desirable in order to investigate what smallest offset is theoretically obtainable for the respective scattering functions; but it is fully impractical due to the long run-time of the evaluation routine.) These solutions agree with the general shape

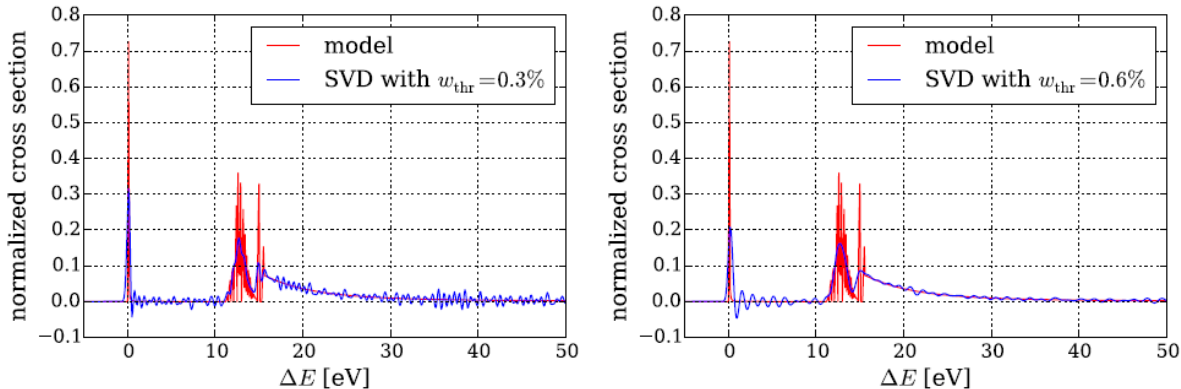


Figure 4: Energy loss functions obtained by by Truncated Singular Value Decomposition. Note that the regularisation parameter is expressed as a percentage of the largest singular value. From [1]

of the energy loss function  $f$ . However, they are strongly oscillatory and neither meet the expected non-negativity of domain nor codomain. A more exotic approach by non-negative matrix factorisation [24] cannot resolve this due to the density of the convolution matrix. More importantly, this cannot hold as an *a priori* parameter choice where measurements yield poor data quality. If noise is further reduced, on the other hand, the potential for increased reconstruction quality is not exhausted. Applying TSVD with  $\alpha = 0.003\sigma_1$  to  $\epsilon_n^h$ , for instance, provides a prohibitively large  $\mu = 0.0145 \text{ eV}^2$ , while the optimal solution for a high precision measurement is not attained. Reconsidering this method with a modified SSF satisfying the natural monotony constraint as defined by eq. (4.8) can eventually reduce oscillations; yet, the result fails at limiting the uncertainty resulting from energy losses with

SSF	Truncated Singular Value Decomposition		Tikhonov Regularisation		Minimisation		Parametrisation with Tikhonov Regularisation	
	Method	$\overline{m(\bar{\nu}_e)^2}$	Method	$\overline{m(\bar{\nu}_e)^2}$	Method	$\overline{m(\bar{\nu}_e)^2}$	Method	$\overline{m(\bar{\nu}_e)^2}$
$\epsilon_{hf}^h$	$\alpha = 0.003\sigma_1$	0.0053[1]	L-curve	0.0048	SLSQP	0.0045	pure	0.038
$\tilde{\epsilon}_{hf}^h$	Model comp.	0.0053[1]	Model comp.	0.0055	trust-constr	0.0029	L-curve	0.0034
	$\alpha = 0.003\sigma_1$	0.0085	L-curve	0.0044	SLSQP	0.0062		
$\epsilon_n^h$	Model comp.	0.0045	Model comp.	0.0063	trust-constr	0.0029		
	Morozov	0.0145	L-curve	0.0283	SLSQP	0.0093	pure	0.154
$\epsilon_{id}^h$			Model comp.	0.0258	trust-constr	0.0071		
	Morozov	0.0000	L-curve	<i>fails</i>	SLSQP	0.0068		
	Model comp.	0.0000	Model comp.	0.0000	trust-constr	0.0032	pure	0.0137

Table 1: Expected systematic uncertainties  $\mu = \overline{m(\bar{\nu}_e)^2}$  in  $\text{eV}^2$ , each with an uncertainty of  $\pm 5 \times 10^{-4} \text{eV}^2$ , for the different regularisation methods and respective scattering models. 'L-curve' indicates a solution obtained by the L-curve criterion, 'Morozov' by application of the discrepancy principle. 'Model comp.' solutions are optimal approximations to the energy loss model with respect to the  $L^2$  norm. 'SLSQP' and 'trust-constr' denote the applied algorithm for constrained optimisation. The method 'pure' refers to the mere identification of optimal parameters in eq. (4.36).

## 5 NUMERICAL EXPERIMENTS

$\mu = 0.0085 \text{ eV}^2$ . Equally, an *a posteriori* estimate by application of Morozov’s discrepancy principle does not provide an improvement as it merely can identify an order of magnitude for the regularisation parameter.

A more stable result is obtained through Tikhonov regularisation. Since the gradient and Jacobian of the minimisation functional in theorem 4.3 can be easily implemented, computation with the Newton conjugate gradient trust-region algorithm[25] is faster than the minimisation of eq. (4.23) by means of singular value decomposition. The fact that this algorithm fails

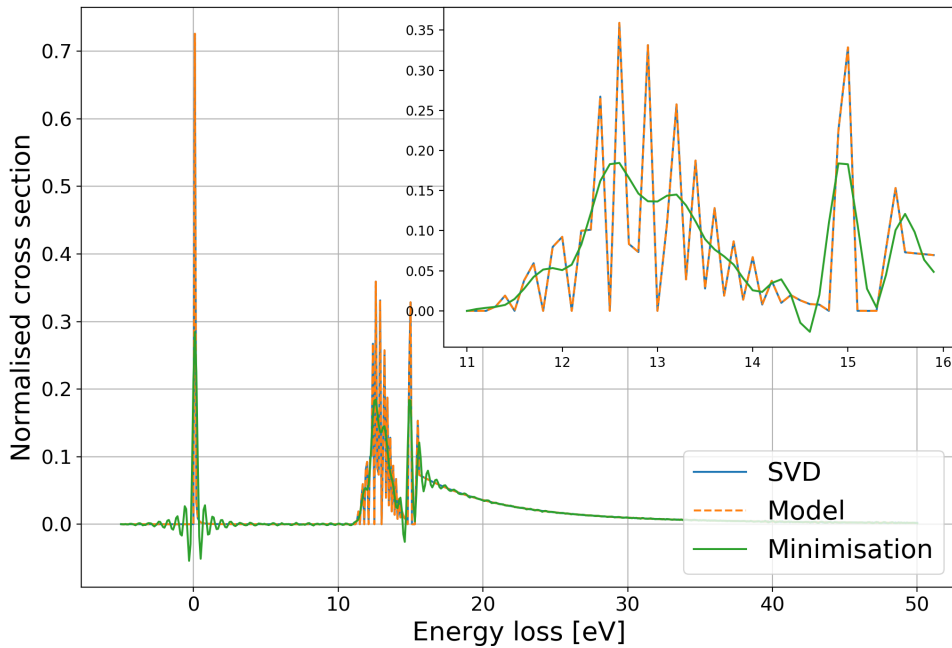


Figure 5: The calculation of Tikhonov solutions by Singular Value Decomposition (SVD) is slower, but more stable in the limit  $\alpha \rightarrow 0$  than a solution obtained by the Newton-CG trust region algorithm. Displayed is the case  $\alpha = 10^{-13}$  for the input  $\epsilon_{id}^h$ , where the SVD regularised solution is not discernible from the input model. The box shows the enlarged excitation/dissociation segment.

at handling very small regularisation parameters can be neglected for the real experiment. Only when applied to  $\epsilon_{id}^h$  it introduces undesired, artificial features for  $\alpha \lesssim 10^{-9}$ . Figure 5 shows solutions for the regularisation parameter  $\alpha = 10^{-13}$ . The optimal theoretical parameter, chosen again from minimising  $\|f_\alpha^h - f_{mod}^h\|^2$ , yields an almost perfect reconstruction of  $\epsilon_{id}^h$  and it is sufficiently strong in the case of the experimental inputs ( $\mu = 0.0055 \text{ eV}^2$ ). More importantly, a parameter chosen by means of the L-curve criterion can provide a comparably good deconvolution without the artificial knowledge of the original model. This indicates that good parameter identification is possible for the real experiment when no model is available.

However, comparing the solution displayed in fig. 6 to the already presented TSVD solution does not show a huge improvement in visual terms. The optimal Tikhonov parameters with respect to the input model and the L-curve both offer no substantial improvements on  $\mu$ . Additionally, they still exhibit the same unphysical features as the TSVD solution, although to a smaller extend.

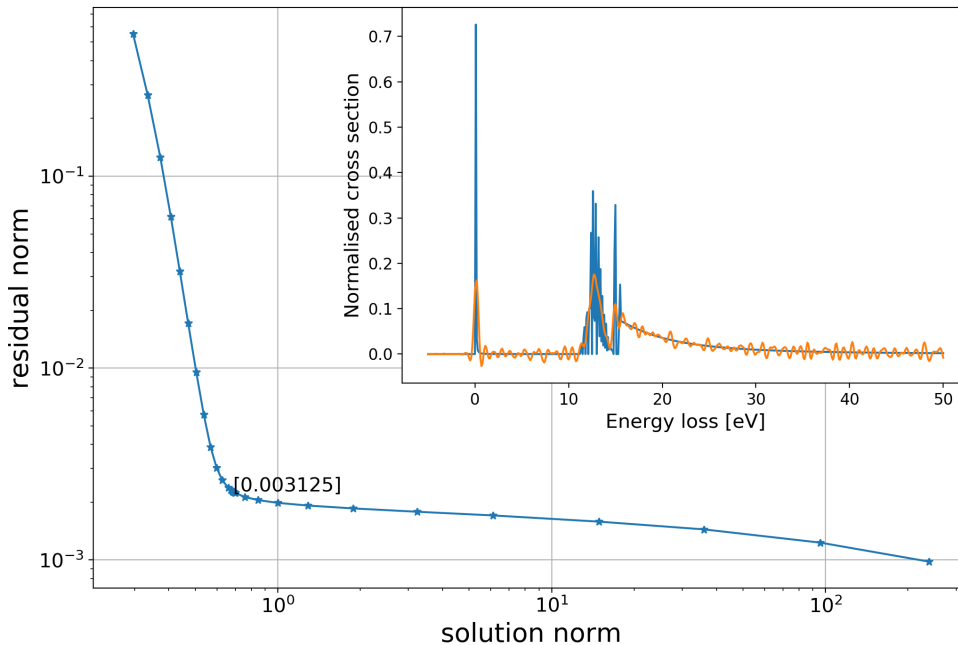


Figure 6: Tikhonov regularisation provides a good solution (box) that is reliably obtained by means of the L-curve criterion.

These shortcomings do not occur when a constrained minimisation is performed that enforces non-negativity of the function and fidelity to the requirement of compact support on  $[0, E_{tot}/2]$ . As an additional constrained, non-positivity of the derivative is imposed, except for the region of electronic excitation and dissociation between 11-16eV. These constraints are much more difficult to handle, though, and are computationally very expensive. Suitable algorithms are the highly versatile constrained trust region (trust-constr) algorithm and the slightly faster SLSQP method (Sequential Least Squares Programming)[23]. The former provides the best available offset for the neutrino mass at  $0.0029 \text{ eV}^2$  in the simulations; the deconvoluted energy loss functions from  $\epsilon_{hf}^h$  are displayed in fig. 7. Moreover, trust-constr is the only applied method for which the evaluation algorithm yields an offset within the error budget bounds even for  $\epsilon_n^h$ . The output obtained from constrained optimisation, however, is highly sensitive to the provided initial guess and introduces new undesired features, most notably a staircasing effect in the ionisation tail. Additionally, the trust-constr algorithm

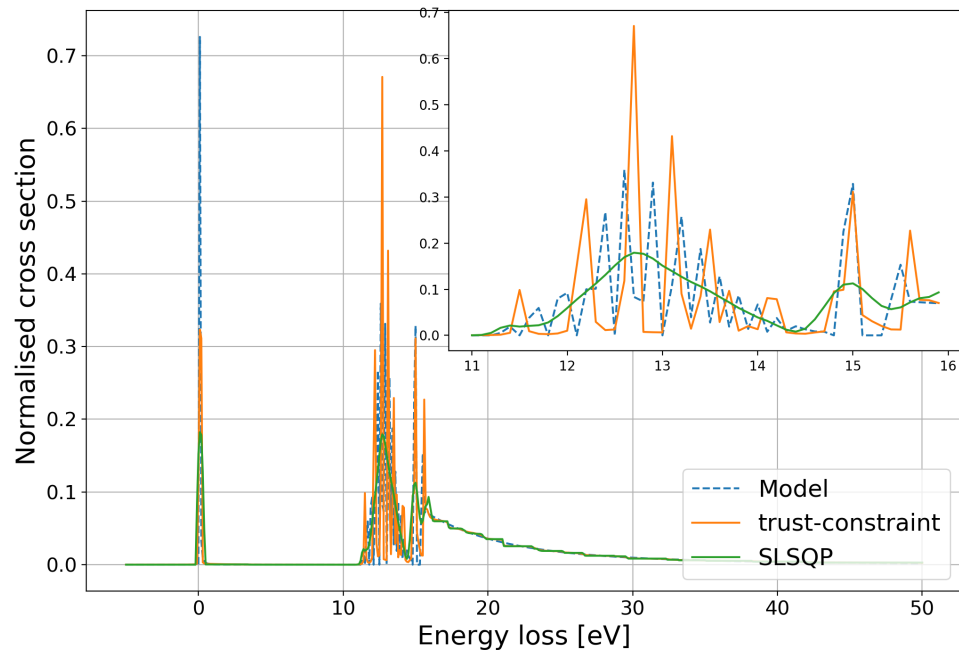


Figure 7: The solutions obtained by constrained optimisation provide the best available values for  $\mu$ , depending on the applied algorithm. The box shows the enlarged excitation/dissociation segment.



provides an apparent reconstruction of the electronic excitation/dissociation section that is essentially arbitrary, even when true input is provided. Thus the used algorithms represent no valid regularisation method in the sense of eq. (4.32).

For the purpose of parametrisations, a suitable method must only rely on the integral data of the SSF, so the Levenberg-Marquardt algorithm is arguably *the* canonical choice for this non-linear problem. In order to maintain physical plausibility, some bounds have to be imposed (e.g. non-negativity of the amplitudes), which can be effectively handled by the very stable implementation *leastsq* in the *lmfit*-package[26]. As can be seen for the input  $\epsilon_{hf}^h$  in fig. 8, it provides visually compelling parametrisations that obey in principal all of the required properties. Unfortunately, in the context of the suggested evaluation method, it

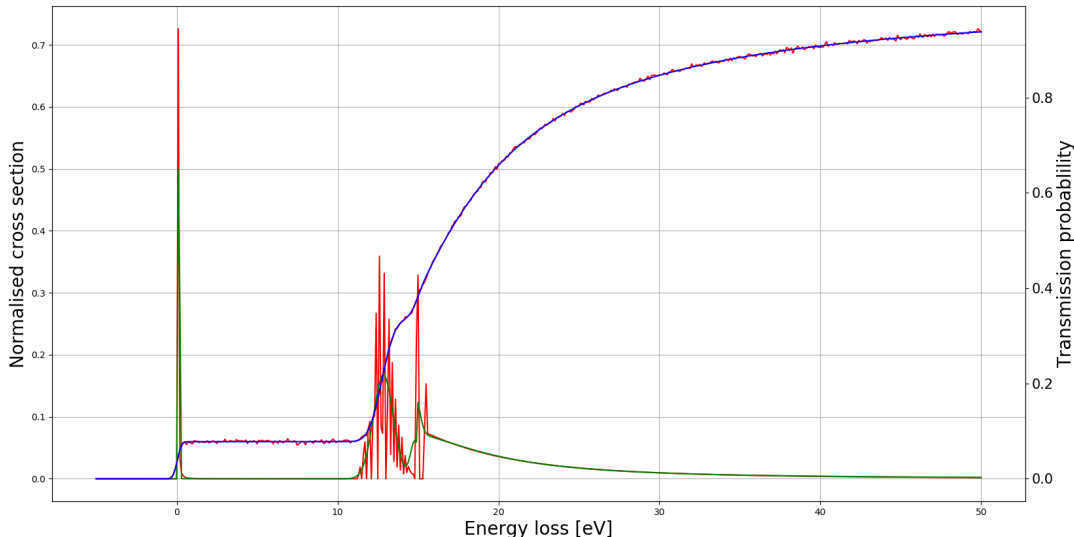


Figure 8: An optimal parametrisation (green) according to eq. (4.36), compared to the energy loss model (red). It is obtained from the input SSF  $\epsilon_{hf}^h$  (red) that is extremely well met when the parametrisation is convoluted with the corresponding transmission function again (blue).

generates extremely poor results. This can be mitigated by applying Tikhonov regularisation to a residual problem. Heuristically, a surprisingly good solution in terms of the evaluation method can be found for

$$\begin{aligned}
 f^h &= 0.3 * f_p^h + \partial f^h \quad \text{subject to} \\
 \partial f^h &= \operatorname{argmin}_{\partial \tilde{f}^h} \left\| T^e \partial \tilde{f}^h - (\epsilon_{hf}^h - 0.3 T^e f_p^h) \right\|^2 + \alpha \left\| \partial \tilde{f}^h \right\|^2
 \end{aligned} \tag{5.1}$$

where  $f_p^h$  denotes an optimal parametrisation of eq. (4.36) with respect to the input  $\epsilon_{hf}^h$  and

the regularisation parameter  $\alpha$  is determined by the L-curve criterion. However, this solution is a poorer approximation to the input energy model in any other respect and it gives reason to assume some instabilities in the evaluation software which is extremely sensible to the numeric integral of the deconvoluted energy loss functions.

## 6 Conclusion and Outlook

This thesis has shown that the retrieval of the tritium energy loss function for 18.6 keV electrons by deconvolution of scattering functions is an ill-posed problem which can be solved with sufficient reliability by various methods, as long as the level of noise in the energy loss function is moderate. The approaches perform acceptably for the simulated data in use, the observed sensibility of the evaluation routine suggests that the principal path to low systematic errors in the KATRIN neutrino mass estimate due to inelastic scattering remains very smooth data for the scattering functions. Canonical methods, most notably Truncated Singular Value Decomposition and Tikhonov regularisation, provide reasonably good estimates, although they do not present satisfactory models of physical reality. The L-curve criterion can be trusted to provide good regularisation parameters for a larger range of noise levels and is especially efficient to apply for Tikhonov regularisation.

The more involved solutions obtained through constrained minimisation are stable and respect the KATRIN error budget best. Nevertheless, they produce certain misleading properties, too. Energy loss functions obtained from optimal parametrisations arguably offer the visually most compelling deconvolution results but fail with respect to the the estimated uncertainty limit for electron scattering. This can be overcome by combination with Tikhonov regularisation; a method for consistent parameter identification that simultaneously balances the parametrisation problem and residual problem is beyond the scope of this thesis, though. Meanwhile, recalling that the applied Tikhonov regularisation is associated with Gaussian noise, it could be useful to pursue regularisation terms optimised to other noise types instead. Some more recent removal methods, e.g. for Cauchy noise[27], are promising. For the purpose of the KATRIN experiment, the present methods should yet be sufficient to make sure that systematic uncertainties from inelastic scattering in the tritium source will contribute by no more than  $0.0075 \text{ eV}^2$  to the experiment's error budget.

---

## 7 Appendix

The code that was used for this thesis is provided as a 'Kelopy' package. For convenience, parameters are handled by txt-files:

The folder 'parameters' contains the parameter files that specify which input and output is handled. For each step, there is already a respective parameter file that can easily be adjusted to the desired use.

In order to obtain the scattering functions, note the relative path to the data files for normed transmission probabilities and scattering probabilities (including zero-column density) as well as for the output files in 'parameters/scatinversion.txt' and execute 'scatinversion.py'.

In order to perform deconvolution by truncated singular value decomposition or Tikhonov regularisation, note the relative path of the single scattering function, the transmission function measurement, the desired output file and, if available, the model to cross check, in 'parameters/deconvolution.txt'. Also note the desired method ('tikh' for Tikhonov regularisation, 'tsvd' for TSVD) and the regularisation parameter. If a parameter should be automatically chosen, the parameter can be replaced by 'morozov' for parameter choice according to the discrepancy principle or 'lcurve' according to the L-curve criterion. The method 'modelcompare' can be used to choose it with respect to a trusted model. Execute 'deconv.py' to obtain the result.

For constrained optimisation, choose the desired method ('trust-constr' or 'SLSQP') in 'optimisation.txt' with the input and output filenames. A monotony constraint can be relieved in a chosen interval. Execute 'constrop.py' to obtain the result.

'parametrisation.py' performs a parametrisation of a convolved function. The file 'parametrisation.txt' further requires the specification of some free parameters.

## List of Figures

1	KATRIN experimental setup . . . . .	6
2	Energy loss function according to Glück . . . . .	7
3	Simulated Transmission functions and single scattering functions . . . . .	10
4	Energy loss by Truncated Singular Value Decomposition . . . . .	26
5	Optimal Tikhonov regularisation . . . . .	28
6	L-curve criterion . . . . .	29
7	Constrained optimisation . . . . .	30
8	Parametrisation . . . . .	31

## References

- [1] V. Hannen et al. “Deconvolution of the energy loss function of the KATRIN experiment”. In: *Astroparticle Physics* 89 (2017), pp. 30–38. ISSN: 0927-6505. DOI: <https://doi.org/10.1016/j.astropartphys.2017.01.010>.
- [2] L. M. Brown. “The idea of the neutrino”. In: *Phys. Today; (United States)* 31:9 (Sept. 1978). DOI: [10.1063/1.2995181](https://doi.org/10.1063/1.2995181).
- [3] C. L. Cowan et al. “Detection of the Free Neutrino: a Confirmation”. In: *Science* 124.3212 (1956), pp. 103–104. ISSN: 0036-8075. DOI: [10.1126/science.124.3212.103](https://doi.org/10.1126/science.124.3212.103). URL: <https://science.sciencemag.org/content/124/3212/103>.
- [4] KATRIN Collaboration. *KATRIN design report 2004*. Tech. rep. 51.54.01; LK 01. Forschungszentrum, Karlsruhe, 2005. DOI: [10.5445/IR/270060419](https://doi.org/10.5445/IR/270060419).
- [5] W. Rodejohann. “Neutrino-less Double Beta Decay and Particle Physics”. In: *Int. J. Mod. Phys. E* 20 (2011), pp. 1833–1930. DOI: [10.1142/S0218301311020186](https://doi.org/10.1142/S0218301311020186). arXiv: [1106.1334](https://arxiv.org/abs/1106.1334) [hep-ph].
- [6] M. Agostini et al. “Results on Neutrinoless Double- $\beta$ Decay of Ge76 from Phase I of the GERDA Experiment”. In: *Physical Review Letters* 111.12 (Sept. 2013). ISSN: 1079-7114. DOI: [10.1103/physrevlett.111.122503](https://doi.org/10.1103/physrevlett.111.122503).
- [7] E. W. Otten and C. Weinheimer. “Neutrino mass limit from tritium  $\beta$  decay”. In: *Reports on Progress in Physics* 71.8 (July 2008), p. 086201. DOI: [10.1088/0034-4885/71/8/086201](https://doi.org/10.1088/0034-4885/71/8/086201).
- [8] V. Aseev et al. “Energy loss of 18 keV electrons in gaseous T and quench condensed D films”. In: *The European Physical Journal D* 10 (Mar. 2000), pp. 39–52. DOI: [10.1007/s100530050525](https://doi.org/10.1007/s100530050525).
- [9] V. Hannen et al. *KEloss code package*. Private communication. 2017.
- [10] Y.-K. Kim and M. E. Rudd. “Binary-encounter-dipole model for electron-impact ionization”. In: *Physical Review A* 50.5 (Nov. 1994), pp. 3954–3967. DOI: [10.1103/PhysRevA.50.3954](https://doi.org/10.1103/PhysRevA.50.3954).
- [11] J. Hadamard. *Lectures on Cauchy’s Problem in Linear Partial Differential Equations*. Isha Books, 2013. ISBN: 978-9333140355.
- [12] D. Werner. *Funktionalanalysis*. Springer-Lehrbuch. Springer Berlin Heidelberg, 2018. ISBN: 978-3-662-55407-4.
- [13] C. Eckart and G. Young. “The approximation of one matrix by another of lower rank”. In: *Psychometrika* 1.3 (1936), pp. 211–218. DOI: <https://doi.org/10.1007/BF02288367>.
- [14] L. Mirsky. “Symmetric Gauge Functions and Unitarily Invariant Norms”. In: *The Quarterly Journal of Mathematics* 11.1 (Jan. 1960), pp. 50–59. URL: <https://doi.org/10.1093/qmath/11.1.50>.

## REFERENCES

---

- [15] V.A. Morozov. *Methods for Solving Incorrectly Posed Problems*. Physics and astronomy online library. Springer-Verlag New York, 1984. ISBN: 978-1-4612-5280-1.
- [16] A. N. Tikhonov et al. “Regularization methods”. In: *Numerical Methods for the Solution of Ill-Posed Problems*. Dordrecht: Springer Netherlands, 1995, pp. 7–63. DOI: [10.1007/978-94-015-8480-7\\_2](https://doi.org/10.1007/978-94-015-8480-7_2).
- [17] A. Tarantola. *Inverse Problem Theory and Methods for Model Parameter Estimation*. Other Titles in Applied Mathematics. Society for Industrial and Applied Mathematics Philadelphia, 2005. ISBN: 978-0-89871-792-1.
- [18] P. C. Hansen. “Analysis of Discrete Ill-Posed Problems by Means of the L-Curve”. In: *SIAM Review* 34.4 (1992), pp. 561–580. URL: <https://doi.org/10.1137/1034115>.
- [19] C. R. Vogel. “Non-convergence of the L-curve regularization parameter selection method”. In: *Inverse Problems* 12.4 (Aug. 1996), pp. 535–547. DOI: [10.1088/0266-5611/12/4/013](https://doi.org/10.1088/0266-5611/12/4/013).
- [20] P. C. Hansen. “The L-Curve and its Use in the Numerical Treatment of Inverse Problems”. In: *Computational Inverse Problems in Electrocardiology*, ed. P. Johnston, *Advances in Computational Bioengineering*. WIT Press, 2000, pp. 119–142.
- [21] A.B. Bakushinskii. “Remarks on choosing a regularization parameter using the quasi-optimality and ratio criterion”. In: *USSR Computational Mathematics and Mathematical Physics* 24.4 (1984), pp. 181–182. ISSN: 0041-5553. DOI: [https://doi.org/10.1016/0041-5553\(84\)90253-2](https://doi.org/10.1016/0041-5553(84)90253-2).
- [22] M. Hanke. “Limitations of the L-curve method in ill-posed problems”. In: *Bit Numer Math* 36 (1996), pp. 287–301. URL: <https://doi.org/10.1007/BF01731984>.
- [23] P. Virtanen et al. “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python”. In: *Nature Methods* 17 (2020), pp. 261–272. DOI: <https://doi.org/10.1038/s41592-019-0686-2>.
- [24] W. Liu et al. “Nonnegative singular value decomposition for microarray data analysis of spermatogenesis”. In: *2008 International Conference on Information Technology and Applications in Biomedicine*. 2008, pp. 225–228.
- [25] J. Nocedal and S. Wright. “Large-Scale Unconstrained Optimization”. In: *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. New York, NY: Springer New York, 2006, pp. 164–192. ISBN: 978-0-387-40065-5. DOI: [10.1007/978-0-387-40065-5\\_7](https://doi.org/10.1007/978-0-387-40065-5_7).
- [26] M. Newville et al. *LMFIT: Non-Linear Least-Square Minimization and Curve-Fitting for Python*. Version 0.8.0. Sept. 2014. DOI: [10.5281/zenodo.11813](https://doi.org/10.5281/zenodo.11813).
- [27] G. Kim, J. Cho, and M. Kang. “Cauchy Noise Removal by Weighted Nuclear Norm Minimization”. In: *Journal of Scientific Computing* 83 (Apr. 2020). DOI: [10.1007/s10915-020-01203-2](https://doi.org/10.1007/s10915-020-01203-2).

## Plagiatserklärung der / des Studierenden

Hiermit versichere ich, dass die vorliegende Arbeit über „Application of Regularisation Methods for a Deconvolution Problem within the KATRIN Collaboration“ selbstständig verfasst worden ist, dass keine anderen Quellen und Hilfsmittel als die angegebenen benutzt worden sind und dass die Stellen der Arbeit, die anderen Werken – auch elektronischen Medien – dem Wortlaut oder Sinn nach entnommen wurden, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht worden sind.

Potsdam, 12. August 2020 \_\_\_\_\_  
(Datum, Unterschrift)

Ich erkläre mich mit einem Abgleich der Arbeit mit anderen Texten zwecks Auffindung von Übereinstimmungen sowie mit einer zu diesem Zweck vorzunehmenden Speicherung der Arbeit in einer Datenbank einverstanden.

Potsdam, 12. August 2020 \_\_\_\_\_  
(Datum, Unterschrift)