

# DFG-SNF Research Group FOR916

Statistical Regularization and Qualitative Constraints

Klaus Frick

Philipp Marnitz

Axel Munk

## **Statistical Multiresolution Estimation in Imaging: Fundamental Concepts and Algorithmic Framework**

Preprint FOR916 11-2

Preprint-Series of the Research Group FOR916

# STATISTICAL MULTIREOLUTION ESTIMATION IN IMAGING: FUNDAMENTAL CONCEPTS AND ALGORITHMIC FRAMEWORK

KLAUS FRICK

*Institute for Mathematical Stochastics  
University of Göttingen  
Goldschmidtstraße 7, 37077 Göttingen*

PHILIPP MARNITZ

*Institute for Mathematical Stochastics  
University of Göttingen  
Goldschmidtstraße 7, 37077 Göttingen*

AXEL MUNK

*Institute for Mathematical Stochastics  
University of Göttingen  
Goldschmidtstraße 7, 37077 Göttingen*

*Max Planck Institute for Biophysical Chemistry  
Am Faßberg 11, 37077 Göttingen*

ABSTRACT. In this paper we introduce a general class of *statistical multiresolution estimators* and develop an algorithmic framework for computing those. By this we mean estimators that are defined as solutions of convex optimization problems with  $\ell_\infty$ -type constraints. We employ a combination of an alternating direction augmented Lagrangian technique with Dykstra's algorithm for computing orthogonal projections onto intersections of convex sets. The capability of the proposed method is illustrated by various examples from imaging.

---

*E-mail addresses:* `frick@math.uni-goettingen.de`, `marnitz@math.uni-goettingen.de`,  
`munk@math.uni-goettingen.de`.

*Key words and phrases.* Statistical Inverse Problems, Statistical Multiscale Analysis, Extreme-Value Statistics, Dykstra's Projection Algorithm, Total Variation, Statistical Imaging, Dantzig Selector, Alternating Direction Method.

Correspondence to `frick@math.uni-goettingen.de` .

## 1. INTRODUCTION

In numerous applications, the relation of observable data  $Y$  and the (unknown) signal of interest  $u^0$  can be modeled as an inverse linear regression problem. We shall assume that the data  $Y = \{Y_\nu\}$  is sampled on the equidistant grid  $X = \{1, \dots, m\}^d$ , with  $m, d \in \mathbb{N}$  and that  $u^0 \in U$  for some linear space  $U$ , such as the Euclidean space or a Sobolev class of functions. Hence the model can be formalized as

$$Y_\nu = (Ku^0)_\nu + \varepsilon_\nu, \quad \nu \in X. \quad (1)$$

Here we assume that  $\varepsilon = \{\varepsilon_\nu\}_{\nu \in X}$  are independent and identically distributed r.v. with  $\mathbf{E}(\varepsilon_\nu) = 0$  and  $\mathbf{E}(\varepsilon_\nu^2) = \sigma^2 > 0$  (white noise). Moreover,  $K : U \rightarrow (\mathbb{R}^m)^d$  denotes a linear operator that encodes the functional relation between the quantities that are accessible by experiment and the underlying signal. Often the operator  $K$  does not have a continuous inverse (or its inverse is ill-conditioned in a discrete setting), that is estimation of  $u^0$  given the data  $Y$  is an *ill-posed problem*. As a consequence, estimators for  $u^0$  can in general not be obtained by merely applying the inverse of  $K$  to an estimator of  $Ku^0$ . Instead, more sophisticated *statistical regularization* techniques have to be employed that, loosely speaking, are capable of simultaneously inverting  $K$  and solving the regression problem.

The application we primarily have in mind is the reconstruction of low-dimensional ( $d \leq 4$  say) signals  $u^0$  which are presumed to exhibit a strong neighborhood structure as it is characteristic of imaging or signal detection problems. These neighborhood relations are often modeled by prior smoothness or structural assumptions on  $u^0$  (e.g. on the texture of an image).

The aim of this paper is twofold. First, we will introduce the broad class of *statistical regularization estimators (SMRE)*. We claim that numerous regularization techniques, that were recently proposed for different problems in various branches of applied mathematics and statistics, can be considered as special cases of these. Among others, this includes the *Dantzig selector* [5] that was recently proposed in the context of high dimensional statistics. Our prior focus, however, will be put on imaging problems and it will turn out that the aforementioned neighborhood relations can be modeled within our SMRE framework in a straightforward manner. This will for example result in *locally adaptive* image reconstruction methods.

The high intrinsic structure of the signals that are typically under consideration in imaging is in drastic contrast to the situation in high-dimensional statistics. Here  $u^0$  is usually assumed to be unstructured and its consistent estimation is guaranteed by the prior assumption of sparsity (cf. [5, 6, 39]). For this reason (and as it will be stressed in Paragraph 1.2.2 below in more detail), algorithmic approaches for estimation of sparse signals have only limited use for imaging application and new methods are required. This will be the second core issue of this work.

**1.1. Statistical Multiresolution Estimation.** We will now introduce the announced class of estimators. To this end, let  $\mathcal{S}$  be some index set and  $\mathcal{W} = \{\omega^S\}_{S \in \mathcal{S}}$  be a set of given weight-functions on the grid  $X = \{1, \dots, m\}^d$ . A *statistical multiresolution estimator (SMRE)*, is defined as a solution of the constrained optimization problem

$$J(u) \rightarrow \inf! \quad \text{s.t.} \quad \max_{S \in \mathcal{S}} \left| \sum_{\nu \in X} \omega_\nu^S (\Lambda(Y - Ku))_\nu \right| \leq q. \quad (2)$$

Here,  $J : U \rightarrow \mathbb{R}$  denotes a regularization functional that incorporates a priori knowledge on the unknown signal  $u^0$  (such as smoothness or sparsity) and  $\Lambda : (\mathbb{R}^m)^d \rightarrow (\mathbb{R}^m)^d$  a possibly non-linear transformation. The constant  $q$  can be considered as a *regularization parameter* that governs the trade-off between regularity and data-fit of the reconstruction. In most practical situations  $q$  is chosen to be the  $\alpha$ -quantile  $q_\alpha$  of the *multiresolution (MR) statistic*  $T(\varepsilon)$ , where  $T : (\mathbb{R}^m)^d \rightarrow \mathbb{R}$  denotes the inequality constraint in (2), i.e.

$$T(v) = \max_{S \in \mathcal{S}} \left| \sum_{\nu \in X} \omega_\nu^S (\Lambda(v))_\nu \right|, \quad v \in (\mathbb{R}^m)^d. \quad (3)$$

This is particularly appealing since then the regularization parameter  $q$  admits a sound statistical interpretation: each solution  $\hat{u}_\alpha$  of (2) satisfies

$$\mathbb{P}(J(\hat{u}_\alpha) \leq J(u^0)) \geq \alpha$$

i.e. the estimator  $\hat{u}_\alpha$  is *more regular* (in terms of  $J$ ) than  $u^0$  with a probability of at least  $\alpha$ .

For a given estimator  $\hat{u}$  of  $u^0$ , the set  $\mathcal{W}$  is assumed to be rich enough in order to catch all relevant non-random signals that are visible in the residual  $Y - K\hat{u}$ . Then, the average function

$$\mu_S(v) = \left| \sum_{\nu \in X} \omega_\nu^S (\Lambda(v))_\nu \right| \quad (4)$$

evaluated at  $v = Y - K\hat{u}$  is supposed to be significantly larger than  $q$  for at least one  $\omega \in \mathcal{W}$ , whenever  $Y - K\hat{u}$  fails to resemble white noise. Put differently, the MR-statistic  $T(Y - K\hat{u})$  is bounded by  $q$ , whenever  $Y - K\hat{u}$  is accepted as white noise according to the *resolution* provided by  $\mathcal{W}$ . In fact, this is a key observations that reveals numerous potential application areas of the estimation method (2). The examples we have in mind are mainly from *statistical imaging*, where the index set  $\mathcal{S}$  is typically chosen to be an overlapping (redundant) system of subsets of the grid  $X$  and  $\omega^S$  is the normalized indicator function on  $S \in \mathcal{S}$ . Consequently the inequality constraint in (2) guarantees that the residual resembles white noise on all sets  $S \in \mathcal{S}$ . In other words, the SMRE approach in (2) applied to imaging yields a reconstruction method that *locally adapts the amount of regularization* according to the underlying image features. We illustrate this in Section 3 by various examples.

Summarizing, the optimization problem in (2) amounts to choose the most parsimonious among all estimators  $\hat{u}$  for which the residual  $Y - K\hat{u}$  resembles white noise according to the statistic  $T$ . If  $Y - K\hat{u}$  contains some non random signal, i.e.  $T(Y - K\hat{u})$  is likely to be larger than  $q$  and  $u$  happens to lie outside the admissible domain of (2). Thus, the multi-resolution constraint prevents too parsimonious reconstructions due to the minimization of  $J$ .

## 1.2. Aims and Related Work.

1.2.1. *Multiresolution Methods.* SMR estimation and related MR statistics have recently been studied in various contexts. We give a brief (but incomplete) overview.

Classical MR statistics are obtained from the general form in (3) by setting  $U = (\mathbb{R}^m)^d$  and  $\Lambda = \text{Id}$ . Moreover, one considers the system  $\mathcal{W}$  to contain indicator functions on cubes. To be more precise, define the index set  $\mathcal{S}$  to be the system of all  $d$ -dimensional cubes in  $X$  and set

$$\omega^S = \frac{1}{\sqrt{\#\mathcal{S}}} \chi_S. \quad (5)$$

Then, the MR-statistic in (3) reduces to

$$T(v) = \max_{S \in \mathcal{S}} \frac{1}{\sqrt{\#\mathcal{S}}} \left| \sum_{\nu \in S} v_\nu \right|.$$

This statistic was introduced in [38] (called scanning statistic there) in order to detect a signal against a noisy background. It was shown in [29] (see also [37]) that

$$\lim_{m \rightarrow \infty} \frac{T(\varepsilon)}{\sqrt{2d \log m}} = \sigma \quad \text{a.s.}$$

If the system  $\mathcal{S}$  is reduced to the set of all *dyadic* squares, then it was proved in [26] that (after suitable transformations)  $T$  also converges weakly to the Gumbel distribution. There, the authors also established a method for locally adaptive image denoising employing linear diffusion equations with spatially varying diffusivity. SMR estimators (2) have been studied recently for the case  $d = 1$  in [10] and [4], where total-variation penalty and the number of jumps in piecewise constant regression were considered as regularization functional  $J$ , respectively. In [19] consistency and convergence rates for SMR estimators have been studied in a general Hilbert space setting, including (2) as a special case.

SMR estimation with squared residuals, that is  $\Lambda(v)_\nu = v_\nu^2$  yields another class of estimators that have attracted much attention. Above all, the situation where  $\mathcal{S}$  consists of the single set  $X$  and  $\omega^X$  is chosen to be the constant 1 function is of special interest, since then (2) reduces to the *penalized least square estimation*. In particular (2) then can be rewritten into

$$J(u) + \lambda \sum_{\nu \in X} (Ku - Y)_\nu^2 \rightarrow \inf! \quad (6)$$

for a suitable multiplier  $\lambda > 0$ . Recently, also non-trivial choices of  $\mathcal{S}$  were considered. In [2]  $\mathcal{S}$  is chosen to consist of a partition of  $G$  which is obtained beforehand by a Mumford-Shah segmentation. In [13], a subset  $S \subset X$  is fixed and afterwards  $\mathcal{S}$  is defined as the collection of all translates of  $S$ .

In [15] MR-statistics are used for shape-constrained estimation based on testing qualitative hypothesis in nonparametric regression for  $d = 1$ . Here, the weight functions  $\omega^S$  incorporate qualitative features such as monotonicity or concavity. Similarly, MR-statistics are used in [16] in order to detect locations of local increase and decrease in density estimation. Much in the same spirit is the work in [14] where multiscale sign tests are employed for computing confidence bands for isotonic median curves.

As mentioned previously, the *Dantzig-selector* [5] is also covered by the general SMR estimation framework in (2). To see this, set  $U = \mathbb{R}^p$  (with typically  $p \gg m$ ),  $\Lambda = \text{Id}$  and define the weights

$$\omega^S = K \chi_S, \quad S \in \mathcal{S}.$$

Then, each solution of (2) can be considered as a generalized Dantzig selector. The matrix  $K \in \mathbb{R}^{p \times m}$  in this context is usually interpreted as *design matrix* of a high dimensional linear model. The classical Dantzig selector as introduced in [5] then results in the special case where  $\mathcal{S}$  only consists of single-elemented subsets of  $X$  and  $J$  is chosen to be the  $\ell_1$ -regularization functional

$$J(u) = |u|_1 = \sum_{i=1}^p |u_i|.$$

1.2.2. *Algorithmic Challenges.* From a computational point of view, SMR estimation amounts to solve the *constrained optimization problem* (2) which can be rewritten into

$$J(u) \rightarrow \inf! \quad \text{s.t.} \quad \mu_S(Y - Ku) \leq q, \quad \forall(S \in \mathcal{S}). \quad (2')$$

We note that in practical applications the number of constraints in (2'), that is the cardinality of the index set  $\mathcal{S}$ , can be quite large and that the inequalities (even for the simplest case where  $\Lambda = \text{Id}$ ) are mutually correlated. Both of these facts turn (2') into a numerically challenging problem, as a consequence of which standard approaches (such as interior point or conjugate gradient methods) perform far from satisfactorily.

The authors in [2, 13, 26] approach the numerical solution of (2') by means of an analogon of (6) with spatially dependent multiplier  $\lambda \in (\mathbb{R}^m)^d$ , i.e.

$$J(u) + \sum_{\nu \in X} \lambda_\nu (Ku - Y)_\nu^2 \rightarrow \inf!$$

Starting from a (constant) initial parameter  $\lambda = \lambda_0$ , the parameter  $\lambda$  is iteratively adjusted by increasing it in regions which were poorly reconstructed before according to the MR-statistic  $T$ . This approach strongly depends on the special structure of  $\mathcal{S}$  that allows a straightforward identification of each set  $S \in \mathcal{S}$  with a unique point in the grid  $X$ . Put differently, it is not clear how to modify this paradigm in order to solve (2) for highly redundant systems  $\mathcal{S}$  as we have it in mind.

Recently a general algorithmic framework was introduced in [1] for the solutions of large-scale convex cone problems that formally cover (2'). The authors consider

$$J(u) \rightarrow \inf! \quad \text{s.t.} \quad Y - Ku \in \mathcal{K}$$

where  $\mathcal{K}$  is a convex cone in some Euclidean space. Hence, in order to recover (2') one has to consider the cone

$$\mathcal{K} = \left\{ (v, q) \in (\mathbb{R}^m)^d \times \mathbb{R} : \left| \sum_{\nu \in \mathcal{S}} \Lambda(v)_\nu \right| \leq q \quad \forall(S \in \mathcal{S}) \right\}$$

The approach in [1] employs the dual formulation of the problem which involves the computation of the dual cone  $\mathcal{K}^*$ . This approach is particularly appealing for the Dantzig selector since in this situation the cone  $\mathcal{K}$  coincides with the epi-graph of the  $\ell^\infty$ -norm and hence its dual cone is straightforward to compute. As it is argued in [1], this approach is capable of computing Dantzig selectors for large scale problems in contrast to previous approaches such as standard linear programming techniques [5] or homotopy methods such as DASSO [27] or [36]. However, for the applications we have in mind (such as locally adaptive imaging reconstruction), the approach in [1] is only of limited use since then the dual cone  $\mathcal{K}^*$  is not directly accessible.

The aim of this paper is to develop a general algorithmic framework that makes solutions of (2') numerically accessible for many applications. In order to do so we propose to introduce a slack variable in (2') and then use a alternating direction Uzawa type algorithm that decomposes the problem (2') into a  $J$ -penalized least squares problem for the primal variable and a orthogonal projection problem on the feasible set for the slack variable. Thus our work is much in the same spirit as [30], which considered an alternating direction method (ADM) for the computation of the Dantzig selector recently. In this case, however, the computation of the occurring orthogonal projections are available in closed form, whereas in our applications this is not the case due to the aforementioned dependencies.

In order to tackle the orthogonal projection problems we employ Dykstra’s projection method [3] which is capable of computing the projection onto the intersection of convex bodies by merely using the individual projections onto the latter. The efficiency of the proposed method hence increases considerably if the index set  $\mathcal{S}$  can be decomposed into “few” partitions that contain indices of mutually independent inequalities in (2’). In particular, by this approach we will be able to compute classical SMRE (as introduced in [10, 19]) in  $d = 2$  space dimensions which to our knowledge has never been done so far. This puts us into the position to study the performance of such estimators compared with other benchmark methods (such as *adaptive weights smoothing* cf. [35]). As it will turn out in Section 3 it will outperform these visually as well as quantitatively.

**1.3. Organization of the Paper.** The paper is organized as follows: In Section 2 we introduce a general algorithmic approach for computing SMREs. We will rewrite (2’) into a linearly constrained problem and compute a saddle point of the corresponding augmented Lagrangian by an alternating direction Uzawa-type algorithm in Paragraph 2.2. Under quite general assumption, we prove convergence of the algorithm in Theorem 2.2 and give some qualitative estimates for the iterates in Corollary 2.4. One of the occurring minimization steps amounts to the computation of an orthogonal projection onto a convex set in Euclidean space. In Paragraph 2.3, this problem will be tackled by means of Dykstra’s projection algorithm introduced in [3]. Finally, we illustrate the performance of some particular instances of SMR estimators in Section 3: we study problems in nonparametric regression, image denoising and deconvolution of fluorescence microscopy images and compare our results to other methods by means of simulations.

## 2. COMPUTATIONAL METHODOLOGY

In this section we will address the question on how to solve the linearly constrained optimization problem (2’). After discussing some notations and basic assumptions in Subsection 2.1, we will reformulate the problem in Paragraph 2.2 such that an Uzawa-type Algorithm can be employed as a solution method. As an effect, the task of computing a solution of (2’) is replaced by alternating

- i) solving an unconstrained penalized least squares problem that is *independent of the MR-statistic  $T$*  and
- ii) computing the orthogonal projection on a convex set in Euclidean space that is *independent of  $J$* .

This constitutes an appealing modular nature of our approach: the regularization functional  $J$  can easily be replaced once a method for the projection problem is settled. For the latter we will propose an iterative projection algorithm in Paragraph 2.3 that was introduced by Boyle and Dykstra in [3].

**2.1. Basic Assumptions and Notation.** From now on,  $X$  will stand for the  $d$ -dimensional grid  $\{1, \dots, m\}^d$  and agree upon  $H = \mathbb{R}^X \simeq (\mathbb{R}^m)^d$  being the space of all real valued functions  $v : X \rightarrow \mathbb{R}$ . Moreover, we assume that  $\mathcal{S}$  denotes some index set and that  $\mathcal{W} = \{\omega^S\}_{S \in \mathcal{S}}$  is a collection of elements in  $H$ . For two elements  $v, w \in H$  we will use the standard inner product and norm

$$\langle v, w \rangle = \sum_{\nu \in X} v_\nu w_\nu \quad \text{and} \quad \|v\| = \sqrt{\langle v, v \rangle}$$

respectively. Next, we assume that  $\Lambda : H \rightarrow H$  is continuous such that  $\Lambda(0) = 0$  and that for all  $S \in \mathcal{S}$  the mapping

$$v \mapsto \langle \omega^S, \Lambda(v) \rangle$$

is convex. With this notation, we can rewrite the average function in (4) in the compact form

$$\mu_S = |\langle \omega^S, \Lambda(v) \rangle|.$$

Furthermore, we define  $U$  to be a separable Hilbert-space with inner product  $\langle \cdot, \cdot \rangle_U$  and induced norm  $\|\cdot\|_U$ . The operator  $K : U \rightarrow H$  is assumed to be linear and bounded and the functional  $J : U \rightarrow \mathbb{R}$  is convex and lower semi-continuous, that is

$$\{u_n\}_{n \in \mathbb{N}} \subset U \text{ and } \lim_{n \rightarrow \infty} u_n =: u \in U \implies J(u) \leq \liminf_{n \rightarrow \infty} J(u_n).$$

Recall the definition of the MR-statistic in (3). Throughout this paper we will agree upon the following

**Assumption A.** *i) For all  $y \in H$  there exists  $u \in U$  such that  $T(Ku - y) < q$ .  
ii) For all  $y \in H$  and  $c \in \mathbb{R}$  the set*

$$\left\{ u \in U : \max_{S \in \mathcal{S}} \mu_S(Ku - y) + J(u) \leq c \right\}$$

*is bounded.*

Under Assumption A it follows from standard techniques in convex optimization, that a solution of (2') exists. As we will discuss in Section 2.2 it even follows that a saddle point of the corresponding Lagrangian exists (cf. Theorem 2.1 below). In this context Assumption A i) is often referred to as *Slater's constraint qualification* and is for instance satisfied if  $K(U)$  is dense in  $H$ . Moreover, Assumption A ii) will be needed in order to guarantee convergence of the algorithm for computing such a solution, as it is proposed in the upcoming section. This requirement is fulfilled if  $J$  is coercive i.e.

$$J(u) \rightarrow \infty \quad \text{if} \quad \|u\|_U \rightarrow \infty.$$

In many applications  $U$  is some function space and  $J$  a gradient based regularization method, such as the total variation semi-norm (cf. Section 3.2). Then a typical sufficient condition for Assumption A ii) is that  $K$  does not annihilate constant functions.

**2.2. Alternating Directions Uzawa Algorithm.** By introducing a slack variable  $v$  we rewrite (2') to the equivalent problem

$$J(u) + G(v) \rightarrow \inf! \quad \text{subject to} \quad Ku + v = Y. \quad (7)$$

Here,  $G$  denotes the characteristic function on the feasible region  $\mathcal{C}$  of (2'), that is,

$$\mathcal{C} = \{v \in H : \mu_S(v) \leq q \forall (S \in \mathcal{S})\} \quad \text{and} \quad G(v) = \begin{cases} 0 & \text{if } v \in \mathcal{C} \\ +\infty & \text{else.} \end{cases} \quad (8)$$

Note that due to the assumptions on  $\Lambda$ , the set  $\mathcal{S}$  is closed and convex. The technique of rewriting (2') into (7) is referred to as the *decomposition-coordination approach* (or *variable splitting*), see e.g. Fortin & Glowinski [18, Chap. III]. There, Lagrangian multiplier methods are used for solving (7). To this end, we recall the definition of the *augmented Lagrangian* of Problem (7), that is

$$L_\lambda(u, v; p) = \frac{1}{2\lambda} \|Ku + v - Y\|^2 + J(u) + G(v) - \langle p, Ku + v - Y \rangle, \quad \lambda > 0. \quad (9)$$

The name stems from the fact that the ordinary Lagrangian

$$L(u, v; p) = J(u) + G(v) - \langle p, Ku + v - Y \rangle$$

is augmented by the quadratic penalty term  $(2\lambda)^{-1} \|Ku + v - Y\|^2$  that fosters the fulfillment of the linear constraints in (7). The *augmented Lagrangian Method* consists in computing a saddle point  $(\hat{u}, \hat{v}, \hat{p})$  of  $L_\lambda$ , that is

$$L_\lambda(\hat{u}, \hat{v}; p) \leq L_\lambda(\hat{u}, \hat{v}; \hat{p}) \leq L_\lambda(u, v; \hat{p}), \quad \forall ((u, v, p) \in U \times H \times H)$$

We note that each saddle point  $(\hat{u}, \hat{v}, \hat{p})$  of the augmented Lagrangian  $L_\lambda$  is already a saddle point of  $L$  and vice versa and that in either case the pair  $(\hat{u}, \hat{v})$  is a solution of (7) (and thus  $\hat{u}$  is a desired solution of (2')). This follows e.g. from [18, Chap 3. Thm. 2.1]. Sufficient conditions for the existence of saddle points are usually harder to come up with. Assumption A summarizes a standard set of such conditions.

**Theorem 2.1.** *Assume that Assumption A holds. Then, there exists a saddle point  $(\hat{u}, \hat{v}, \hat{p})$  of  $L_\lambda$ .*

*Proof.* According to [17, Chap. III, Prop. 3.1 and Prop. 4.2] a saddle point of  $L$  exists, if there is an element  $u_0 \in U$  such that  $G$  is continuous at  $Ku_0 - Y$  and that

$$J(u) + G(Ku - Y) \rightarrow \infty \quad \text{as} \quad \|u\|_U \rightarrow \infty. \quad (10)$$

According to Assumption A i) and due to the continuity of  $\Lambda$  the first requirement is clearly satisfied. Further, the coercivity assumption (10) is a consequence of Assumption A ii).  $\square$

We will use an *Alternating Directions Uzawa Algorithm* (cf. Algorithm 1) as described in [18, Chap. III Sec. 3.2] for the computation of a saddle point of  $L_\lambda$  (and hence of a solution of (2')): successive minimization of the augmented Lagrangian  $L_\lambda$  w.r.t. the first and second variable followed by an explicit step for maximizing w.r.t. the third variable is performed. Convergence of this method is established in Theorem 2.2 which is a generalization of [18, Chap. III Thm. 4.1].

**Theorem 2.2.** *Every sequence  $\{(u_k, v_k)\}_{k \geq 1}$  that is generated by Algorithm 1 is bounded in  $U \times H$  and every weak cluster point is a solution of (7). Moreover,*

$$\|Ku_k + v_k - Y\| = o(k^{-1/2}) \quad \text{and} \quad \|K(u_k - u_{k-1})\| = o(k^{-1/2}).$$

**Remark 2.3.** i) Theorem 2.2 implies, that each weak cluster point of  $\{u_k\}_{k \geq 1}$  is a solution of (2'). In particular, if the solution  $u^\dagger$  of (2') is unique (e.g. if  $J$  is strictly convex), then  $u_k \rightharpoonup u^\dagger$ .

ii) Note in particular that (11) is independent of the choice of  $J$ , while (12) is independent of the multiresolution statistic being used. This decomposition gives the proposed method a neat modular appeal: once an efficient solution method for the projection problem (11) is established (see e.g. Section 2.3), the regularization functional  $J$  in (2) can easily be replaced by simply providing an algorithm for the penalized least squares problem (12) (which is often at hand).

For a given tolerance  $\tau > 0$ , Theorem 2.2 implies that Algorithm 1 terminates and outputs approximate solution  $u[\tau]$  and  $v[\tau]$  of (7). However, the breaking condition in Algorithm 1 merely guarantees that the linear constraint in (7) is approximated sufficiently well. Moreover, we know from construction that  $v[\tau] \in \mathcal{C}$ , which implies  $G(v[\tau]) = 0$ . So, it remains to evaluate the validity of  $u[\tau]$  which is done in Corollary 2.4.

---

**Algorithm 1** Alternating Directions Uzawa Algorithm
 

---

**Require:**  $Y \in H$  (data),  $\lambda > 0$  (step size),  $\tau \geq 0$  (tolerance).

**Ensure:**  $(u[\tau], v[\tau])$  is an approximate solution of (7) computed in  $k[\tau]$  iteration steps.

$u_0 \leftarrow \mathbf{0}_U$  and  $v_0 = p_0 \leftarrow \mathbf{0}_H$

$r \leftarrow \|Ku_0 + v_0 - Y\|$  and  $k \leftarrow 0$ .

**while**  $r > \tau$  **do**

$k \leftarrow k + 1$ .

$v_k \leftarrow \tilde{v}$  where  $\tilde{v} \in \mathcal{C}$  satisfies

$$\|\tilde{v} - (Y + \lambda p_{k-1} - Ku_{k-1})\|^2 \leq \|v - (Y + \lambda p_{k-1} - Ku_{k-1})\|^2 \quad \forall (v \in \mathcal{C}). \quad (11)$$

$u_k \leftarrow \tilde{u}$  where  $\tilde{u}$  satisfies

$$\frac{1}{2} \|K\tilde{u} - (Y + \lambda p_{k-1} - v_k)\|^2 + \lambda J(\tilde{u}) \leq \frac{1}{2} \|Ku - (Y + \lambda p_{k-1} - v_k)\|^2 + \lambda J(u) \quad \forall (u \in U). \quad (12)$$

$p_k \leftarrow p_{k-1} - (Ku_k + v_k - Y)/\lambda$ .

$r \leftarrow \max(\|Ku_k + v_k - Y\|, \|K(u_k - u_{k-1})\|)$ .

**end while**

$u[\tau] \leftarrow u_k$  and  $v[\tau] \leftarrow v_k$  and  $k[\tau] \leftarrow k$ .

---

**Corollary 2.4.** Let  $(\hat{u}, \hat{v}, \hat{p}) \in U \times H \times H$  be any saddle point of  $L_\lambda$ . Then,

$$0 \leq J(u[\tau]) - J(\hat{u}) - \langle K^* \hat{p}, u[\tau] - \hat{u} \rangle_U \leq \left( \frac{\tau + \|K\hat{u}\|}{\lambda} + 3 \|\hat{p}\| \right) \tau \quad \forall (\tau > 0).$$

In particular if  $J(u) = \frac{1}{2} \|Lu\|_V^2$ , where  $V$  is a further Hilbert-space and  $L : U \supset D \rightarrow V$  is a linear, densely-defined and closed operator, then  $\|L(u[\tau] - \hat{u})\| = \mathcal{O}(\sqrt{\tau})$ .

The results in Theorem 2.2 and Corollary 2.4 show how the accuracy of the approximate solutions  $u[\tau]$  and  $v[\tau]$  depends on  $\tau$ . Moreover, it follows from the definition of  $L_\lambda$  in (9) and Corollary 2.4 that a small value for  $\lambda$  fosters the linear constraint in (7) but may result in slow decay of the objective functional  $J$ .

**Example 2.5** (Dantzig-selector). As already mentioned in the introduction, SMR estimation (i.e. finding solutions of (2)) reduces to the computation of *Dantzig selectors* for the particular setting  $d = 1$ ,  $U = \mathbb{R}^p$  (with usually  $p \gg m$ ) and

$$J(u) = |u|_1.$$

When applying Algorithm 1 the subproblem (12) amounts to compute

$$u_k \in \operatorname{argmin}_{u \in \mathbb{R}^p} \frac{1}{2} \|Ku - (Y + \lambda p_{k-1} - v_k)\|^2 + \lambda |u|_1.$$

This is the well known *least absolute shrinkage and selection operator (LASSO)* estimator [39]. For the classical Dantzig selector, one chooses  $\mathcal{S} = \{1, \dots, p\}$  and defines for  $S \in \mathcal{S}$  the weight  $\omega^S = K\chi_{\{S\}}$ . Hence, the subproblem (11) in this case consists in the orthonormal projection of  $Y_k = Y + \lambda p_{k-1} - Ku_{k-1}$  onto the set

$$\mathcal{C} = \left\{ v \in \mathbb{R}^m : \left| \sum_{1 \leq j \leq m} \omega_j^S v_j \right| \leq q \text{ for } 1 \leq S \leq p \right\}.$$

The implications of Corollary 2.4 in the present case are in general rather weak. If the saddle point  $\hat{u}$  is known to be  $S$ -sparse and when  $K$  restricted to the support of  $\hat{u}$  is injective, then it can be shown that  $\|u[\tau] - \hat{u}\|_1 = \mathcal{O}(\tau)$ .

We finally note that for this particular situation a slightly different decomposition than proposed in (2.2) is favorable. To be more precise, define  $\tilde{K} = K^T K$  and  $\tilde{Y} = K^T Y$  and consider

$$J(u) + \tilde{G}(v) \quad \text{subject to} \quad \tilde{K}u - v = \tilde{Y}.$$

where  $\tilde{G}$  is the characteristic function on the set  $\{v \in H : \|v\|_\infty \leq q\}$ . Algorithm 1 applied to this modified decomposition then results in the ADM method as introduced in [30]. In this case the projection in step (11) has a closed form.

**2.3. The Projection Problem.** Algorithm 1 resolves the constrained convex optimization problem (2') into a quadratic program (11) and an unconstrained optimization problem (12). The quadratic program (11) in the  $k$ -th step of Algorithm 1 can be written as a projection problem:

$$\|v - Y_k\|^2 \rightarrow \inf! \quad \text{subject to} \quad \mu_S(v) \leq q \quad \forall (s \in \mathcal{S}) \quad (13)$$

where  $Y_k = Y + \lambda p_{k-1} - K u_{k-1}$ . We reformulate the side conditions to

$$v \in \mathcal{C} = \bigcap_{S \in \mathcal{S}} C_S \quad \text{where} \quad C_S = \{v \in H : \mu_S(v) \leq q\}. \quad (14)$$

The sets  $C_S$  are closed and convex and problem (13) thus amounts to compute the projection  $P_{\mathcal{C}}(Y_k)$  of  $Y_k$  onto the intersection  $\mathcal{C}$  of closed and convex sets. According to this interpretation, we use Dykstra's projection algorithm as introduced in [3] to solve (13). This algorithm takes an element  $v \in H$  and convex sets  $D_1, \dots, D_M \subset H$  as arguments. It then creates a sequence converging to the projection of  $v$  onto the intersection of the  $D_j$  by successively performing projections onto individual  $D_j$ 's. To this end, let  $P_D(\cdot)$  denote the projection onto  $D \subset H$  and  $S_D = P_D - \text{Id}$  be the corresponding projection step.

---

#### Algorithm 2 Dykstra's Algorithm

---

**Require:**  $h \in H$ ,  $D_1, \dots, D_M \subset H$  are closed and convex sets,

**Ensure:** A sequence  $\{h_k\}_{k \in \mathbb{N}}$  that converges strongly to  $P_{\mathcal{D}}(h)$  where  $\mathcal{D} = \bigcap_{j=1, \dots, M} D_j$

$h_{0,0} \leftarrow h$

**for**  $j = 1$  to  $M$  **do**

$h_{0,j} \leftarrow P_{D_j}(h_{0,j-1})$  and  $Q_{0,j} \leftarrow S_{D_j}(h_{0,j-1})$

**end for**

$h_1 \leftarrow h_{0,M}$  and  $k \leftarrow 1$

**for**  $k \geq 1$  **do**

$h_{k,0} \leftarrow h_k$

**for**  $j = 1$  to  $M$  **do**

$h_{k,j} \leftarrow P_{D_j}(h_{k,j-1} - Q_{k-1,j})$  and  $Q_{k,j} \leftarrow S_{D_j}(h_{k,j-1} - Q_{k-1,j})$

**end for**

$h_{k+1} \leftarrow h_{k,M}$  and  $k \leftarrow k + 1$

**end for**

---

A natural explanation of the algorithm in a primal-dual framework as well as a proof that the sequence  $\{h(M, k)\}_{k \in \mathbb{N}}$  converges to  $P_{\mathcal{D}}(h)$  in norm can be found in [11, 21]. For the case when  $\mathcal{D}$  constitutes a polyhedron even explicit error estimates are at hand (cf. [42]):

**Theorem 2.6.** Let  $\{h_k\}_{k \in \mathbb{N}}$  be the sequence generated by Algorithm 2 and  $P_{\mathcal{D}}(h)$  be the projection of the input  $h$  onto  $\mathcal{D}$ . Then there exist constants  $\rho > 0$  and  $0 \leq c < 1$  such that for all  $k \in \mathbb{N}$

$$\|h_{M,k} - P_{\mathcal{D}}(h)\| \leq \rho c^k.$$

**Remark 2.7.** The constant  $c$  increases with the number  $M$  of convex sets which intersection form the set  $\mathcal{D}$  that  $h$  is to be projected on. The convergence rate therefore improves with decreasing  $M$ . For further details and estimates for the constants  $\rho$  and  $c$ , we refer to [42].

Note that application of Dykstra's algorithm is particularly appealing if the projections  $P_{D_j}$  can be easily computed or even stated explicitly, as it is the case in the following examples.

**Example 2.8.** Assume that  $\Lambda = \text{Id}$ . Then the sets  $C_S$  are the rectangular cylinders

$$C_S = \{v \in H : |\langle \omega^S, v \rangle| \leq q\}.$$

The projection can therefore be explicitly computed as

$$P_{C_S}(v) = \begin{cases} v - \text{sign}(\langle \omega^S, v \rangle) \frac{\omega^S}{\|\omega^S\|} \left( \frac{|\langle \omega^S, v \rangle|}{\|\omega^S\|} - q \right) & \text{if } \mu_S(v) > q \\ v & \text{else} \end{cases}. \quad (15)$$

The left image in Figure 1 depicts an example for  $\mathcal{C}$  for  $H = \mathbb{R}^2$ . For a detailed geometric interpretation of the MR-statistic we also refer to [32].

**Example 2.9.** Assume that  $\Lambda(v)_{\nu} = v_{\nu}^2$ . Then, it follows that  $v \mapsto \langle \omega^S, \Lambda(v) \rangle$  is convex if and only if  $\omega_{\nu}^S \geq 0$  for all  $\nu \in X$ . In this case, the sets  $C_S$  are elliptic cylinders

$$C_S = \left\{ v \in H : \sum_{\nu \in X} \omega_{\nu}^S v_{\nu}^2 \leq q \right\}.$$

Moreover, if  $\omega_{\nu}^S \in \{0, 1\}$  for all  $\nu \in X$ , then the projection  $P_{C_S}$  can be explicitly computed as

$$P_{C_S}(v) = \begin{cases} \frac{q}{\langle \omega^S, \Lambda(v) \rangle} v \chi_{\{\omega^S=1\}} + v \chi_{\{\omega^S=0\}} & \text{if } \mu_S(v) > q \\ v & \text{else} \end{cases}. \quad (16)$$

The right image in Figure 1 depicts an example of  $\mathcal{C}$  for  $H = \mathbb{R}^2$ .

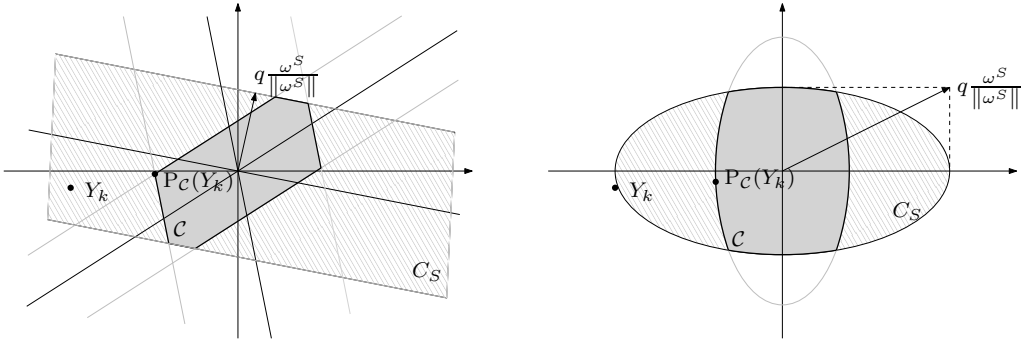


FIGURE 1. Admissible set  $\mathcal{C}$  for the projection problem (13) as in Example 2.8 (left) and Example 2.9 (right).

A first approach to use Dykstra's algorithm to solve (13) is to set  $M = \#\mathcal{S}$  and identify  $D_j$  with  $C_S$  for all  $j = 1, \dots, M$ . In view of Remark 2.7, however, it is clearly desirable to decrease the number  $M$  of convex sets that enter Dykstra's algorithm. In order to do so, we take a slightly more sophisticated approach than the one just presented. We partition the set  $\mathcal{S}$  into  $\mathcal{S}_1, \dots, \mathcal{S}_M$  such that for all  $S \neq \tilde{S} \in \mathcal{S}_j$

$$\omega^S \perp \omega^{\tilde{S}} \quad \text{and} \quad \frac{\partial}{\partial v_\nu} \Lambda(\cdot)_{\tilde{\nu}} \equiv 0, \forall (\nu \in S, \tilde{\nu} \in \tilde{S}) \quad (17)$$

and regroup  $\{C_S\}_{S \in \mathcal{S}}$  into  $\{D_1, \dots, D_M\}$  with

$$D_j = \bigcap_{S \in \mathcal{S}_j} C_S. \quad (18)$$

Given the projections  $P_{C_S}$ , the projection onto  $D_j$  can be easily computed: For  $v \in H$  identify the set

$$V_j = \{S \in \mathcal{S}_j : \mu_S(v) > q\}$$

of indices for which  $v$  violates the side condition (14) and set

$$P_{D_j}(v) = v - \sum_{S \in V_j} (P_{C_S} - \text{Id})v \quad (19)$$

To keep  $M$  small, we choose  $\mathcal{S}_1 \subset \{1, \dots, N\}$  as the biggest set such that (17) holds for all  $S, \tilde{S} \in \mathcal{S}_1$ . We then choose  $\mathcal{S}_2 \subset \mathcal{S} \setminus \mathcal{S}_1$  with the same property and continue in this way until all indices are utilized. While this procedure does not necessarily result into  $M$  being minimal with the desired property, it still yields a distinct reduction of  $N$  in many practical situations. We will illustrate this approach for SMR estimation in imaging in Section 3.

### 3. APPLICATIONS

In this section we will illustrate the capability of Algorithm 1 for computing SMR estimators in some practical situations: in Section 3.1 we will study a simply one-dimensional regression problem as it was also studied in [10], yet with a different penalty function  $J$ . In Section 3.2 we illustrate how SMR estimation performs in image denoising. In both cases we compare our results to other methods. Finally, we will apply the SMRE technique to the problem of image deblurring in confocal fluorescence microscopy in Section 3.3.

Before we study the aforementioned examples, we clarify some common notation. We will henceforth assume that  $U = H = (\mathbb{R}^m)^d$  with  $d = 1$  (Section 3.1) and  $d = 2$  (Sections 3.2 and 3.3), respectively. Moreover, we will employ gradient based regularization functionals of the form

$$J(u) = TV_p(u) := \frac{1}{p} \sum_{\nu \in X} |Du_\nu|_2^p \quad \text{with } p \in \{1, 2\} \quad (20)$$

where  $|\cdot|_2$  is the Euclidean norm in  $\mathbb{R}^d$  and  $D$  denotes the forward difference operator defined by

$$(Du_\nu)_i = \begin{cases} u_{\nu+e_i} - u_\nu & \text{if } 1 \leq \nu_i \leq n-1 \\ 0 & \text{else.} \end{cases}$$

For the case  $p = 2$  the minimization problem (12) amounts to solve an implicit time step of the  $d$ -dimensional diffusion equation with initial value  $(Y + \lambda p_{k-1} - v_k)$  and time step size  $\lambda$ . This can be solved by a simple (sparse) matrix inversion.

For the case  $p = 1$ ,  $TV_1$  is better known as *total-variation semi-norm*. It was shown in [31] (see also [22] for similar results in the continuous setting) that the *taut-string algorithm* (as introduced in [9]) constitutes an efficient solution method for (12) in the case  $d = 1$ . In the general case  $d \geq 1$ , we employ the fixed point approach for solving the Euler-Lagrange equations for (12) described in [12]. We finally note that the functional  $TV_1$  fails to be differentiable; a fact that leads to serious numerical problems when trying to compute the Euler-Lagrange conditions for (12). Hence, we will use in our simulations a regularized version of  $TV_1$  defined by

$$TV_1^\beta(u) = \sum_{\nu \in X} \sqrt{(Du_\nu)_i^2 + \beta^2} \quad (21)$$

for a small constant  $\beta > 0$ .

*Evaluation.* In order to evaluate the performance of SMR estimation, we will employ various distance measures between an estimator  $\hat{u}$  and the true signal  $u^0$ . On the one hand, we will use standard measures such as *mean integrated squared error (MISE)* and the *mean integrated absolute error (MIAE)* which are given by

$$\text{MISE} = \mathbf{E} \left( \frac{1}{m^d} \sum_{\nu \in X} (\hat{u}_\nu - u_\nu^0)^2 \right) \quad \text{and} \quad \text{MIAE} = \mathbf{E} \left( \frac{1}{m^d} \sum_{\nu \in X} |\hat{u}_\nu - u_\nu^0| \right),$$

respectively. On the other hand, we also intend to measure how well an estimator  $\hat{u}$  matches the “smoothness” of the true signal  $u^0$ , where smoothness is characterized by the regularization functional  $J$ . To this end, we introduce the *symmetric Bregman divergence*

$$D_J^{\text{sym}}(\hat{u}, u^0) = \frac{1}{m^d} \sum_{\nu \in X} (\nabla J(\hat{u})_\nu - \nabla J(u^0)_\nu) (\hat{u}_\nu - u_\nu^0),$$

where  $\nabla J$  denotes the gradient of the regularization functional  $J$ . Clearly,  $D_J^{\text{sym}}(\hat{u}, u^0)$  is symmetric and since  $J$  is assumed to be convex, also non-negative. However, the symmetric Bregman divergence usually does not satisfy the triangle inequality and hence in general does not define a (semi-) metric on  $U$  [7]. The following examples shed some light on how the Bregman divergence incorporates the functional  $J$  in order to measure the distance of  $\hat{u}$  and  $u^0$ .

**Example 3.1.** Let  $J(u) = TV_p$  as in (20).

i) If  $p = 2$ , then

$$D_J^{\text{sym}}(\hat{u}, u^0) = \sum_{\nu \in X} |D\hat{u}_\nu - Du_\nu^0|_2^2.$$

In other words, the symmetric Bregman distance w.r.t. to  $TV_2$  is the mean squared distance of the *derivatives* of  $\hat{u}$  and  $u^0$ .

ii) If  $p = 1$ , then

$$\begin{aligned} D_J^{\text{sym}}(\hat{u}, u^0) &= \frac{1}{m^2} \sum_{\nu \in X} \left( \frac{D\hat{u}_\nu}{|D\hat{u}_\nu|} - \frac{Du_\nu^0}{|Du_\nu^0|} \right) \cdot (D\hat{u}_\nu - Du_\nu^0) \\ &= \frac{1}{m^2} \sum_{\nu \in X} (|D\hat{u}_\nu| + |Du_\nu^0|) \left( 1 - \frac{D\hat{u}_\nu}{|D\hat{u}_\nu|} \cdot \frac{Du_\nu^0}{|Du_\nu^0|} \right) \\ &= \frac{1}{m^2} \sum_{\nu \in X} (|D\hat{u}_\nu| + |Du_\nu^0|) (1 - \cos \gamma_\nu), \end{aligned}$$

where  $\gamma_\nu$  denotes the angle between the level lines of  $\hat{u}$  and  $u^0$  at the point  $x_\nu$ . Put differently, the symmetric Bregman divergence w.r.t the total variation semi-norm  $TV_1$  is small if for sufficiently many points  $x_\nu$  either both  $\hat{u}$  and  $u^0$  are constant in a neighborhood of  $x_\nu$  or the level lines of  $\hat{u}$  and  $u^0$  through  $x_\nu$  are parallel.

In practice rather  $TV_1^\beta$  in (21) (for a small  $\beta > 0$ ) instead of  $TV_1$  is used in order to avoid singularities. Then, the above formulas are slightly more complicated.

We will use the mean symmetric Bregman divergence (MSB) given by

$$\text{MSB} = \mathbf{E} \left( D_J^{\text{sym}}(\hat{u}, u^0) \right)$$

as an additional evaluation method. In all our simulations we approximate the expectations above by the empirical means of 500 trials.

*Comparison with other methods.* We will compare the SMR estimators to other regression methods. Firstly, we will consider estimators obtained by the *global* penalized least squares method:

$$\hat{u}(\lambda) := \underset{u \in H}{\operatorname{argmin}} \frac{1}{2} \sum_{\nu \in X} (u_\nu - Y)^2 + \lambda J(u), \quad \lambda > 0. \quad (22)$$

In particular, we focus on estimators  $\hat{u}(\lambda)$  that are closest (in some sense) to the true function  $u^0$ . We call such estimators *oracles*. We define the  $L^2$ - and Bregman-oracle by  $\hat{u}_{L^2} = \hat{u}(\lambda_2)$  and  $\hat{u}_B = \hat{u}(\lambda_B)$ , where

$$\lambda_2 := \mathbf{E} \left( \underset{\lambda > 0}{\operatorname{argmin}} \|u^0 - \hat{u}(\lambda)\| \right) \quad \text{and} \quad \lambda_B := \mathbf{E} \left( \underset{\lambda > 0}{\operatorname{argmin}} D_J^{\text{sym}}(u^0, \hat{u}(\lambda)) \right)$$

respectively. Of course, oracles are not available in practice, since the true signal  $u^0$  is unknown. However, they represent ideal instances within the class of estimators given by (22) that usually perform better than any data-driven parameter choices (such as cross-validation) and hence may serve as a reference.

Secondly, we also compare our approach to *adaptive weights smoothing (AWS)* [35] which constitutes a benchmark technique for data-driven, spatially adaptive regression. We compute these estimators by means of the official R-package<sup>1</sup> and denote them by  $\hat{u}_{\text{aws}}^{\ker}$ , where  $\ker \in \{\text{Gaussian}, \text{Triangle}\}$  decodes the shape of the underlying regression kernel.

**3.1. Non-parametric Regression.** In this section we apply the SMR estimation technique to a nonparametric regression problem in  $d = 1$  dimensions, i.e. the noise model (1) becomes

$$Y_\nu = u_\nu^0 + \varepsilon_\nu \quad \nu = 1, \dots, m, \quad (23)$$

where we assume that  $\varepsilon_\nu$  are independently and normally distributed r.v. with  $\mathbf{E}(\varepsilon_\nu) = 0$  and  $\mathbf{E}(\varepsilon_\nu^2) = \sigma^2$ . The upper left image in Figure 2 depicts the true signal  $u^0$  (solid line) and the data  $Y$ , with  $m = 1024$  and  $\sigma = 0.5$ . The application we have in mind with this example arises in NMR spectroscopy, where the NMR spectra provide structural information on the number and type of chemical entities in a molecule.

In this context the concept of SMR-estimation was considered in [10] (see also [8]) for the case where  $J = TV_1$ , which is well suited for the reconstruction of piecewise constant signals. In the present situation, however, we suggest to choose  $J = TV_2$ , since the true signal  $u^0$  is rather smooth.

<sup>1</sup>available at <http://cran.r-project.org/web/packages/aws/index.html>

Finally, we discuss the MR-statistic  $T$  in (3). We choose  $\Lambda = \text{Id}$  and the index set  $\mathcal{S}$  to consist of all discrete intervals with side lengths ranging from 1 to 100. For an interval  $S \in \mathcal{S}$  we set  $\omega^S = (\#S)^{-1/2}\chi_S$ . Thus, each SMR estimator solves the constrained optimization problem

$$TV_2(u) \rightarrow \inf! \quad \text{s.t.} \quad \frac{1}{\sqrt{\#S}} \left| \sum_{\nu \in S} (Y - u)_\nu \right| \leq q \quad \forall (S \in \mathcal{S}). \quad (24)$$

We choose  $q$  to be the  $\alpha$ -quantile of the MR-statistic  $T$ , that is

$$q_\alpha = \inf \left\{ q \in \mathbb{R} : \mathbb{P} \left( \max_{S \in \mathcal{S}} \frac{1}{\sqrt{\#S}} \left| \sum_{\nu \in S} \varepsilon_\nu \right| \leq q \right) \geq \alpha \right\} \quad \alpha \in (0, 1). \quad (25)$$

We note that except for few special cases (cf. [26, 28]) closed form expressions for the distribution of the MT-statistic  $T$  are usually not at hand. In practice one rather considers the empirical distribution of  $T$  where the variance  $\sigma^2$  can be estimated at a rate  $\sqrt{md}$  (cf. [33]).

We will henceforth denote by  $\hat{u}_\alpha$  a solution of (24) with  $q = q_\alpha$ . As argued in Section 1.1,  $\hat{u}_\alpha$  is smoother (i.e. has smaller value  $TV_2$ ) than the true signal  $u^0$  with a probability of at least  $\alpha$  while it satisfies the constraint that the multiresolution statistic  $T$  does not exceed  $q_\alpha$ . This is a sound statistical interpretation of the regularization parameter  $\alpha$ .

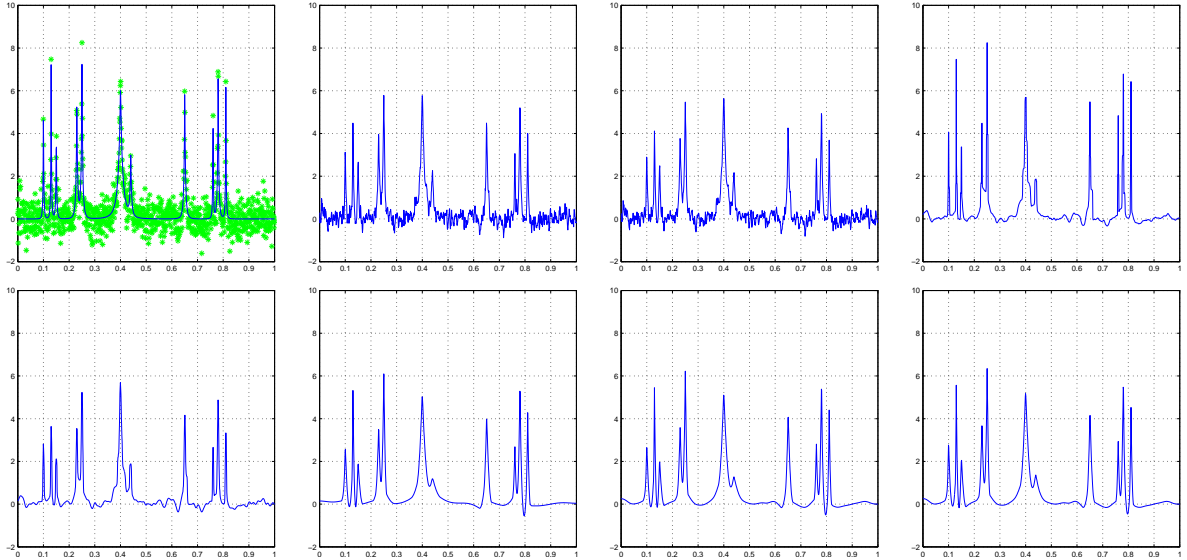


FIGURE 2. Row-wise from top left: true signal  $u^0$  (solid) with data  $Y$ , oracles  $\hat{u}_{L^2}$  and  $\hat{u}_B$ , AWS estimators  $\hat{u}_{\text{aws}}^{\text{Triangle}}$  and  $\hat{u}_{\text{aws}}^{\text{Gaussian}}$ , and the SMREs  $\hat{u}_{0.9}$ ,  $\hat{u}_{0.75}$  and  $\hat{u}_{0.5}$ .

*Numerical results and simulations.* In Figure 2 the oracles  $\hat{u}_{L^2}$  and  $\hat{u}_B$ , the AWS-estimators  $\hat{u}_{\text{aws}}^{\text{Triangle}}$  and  $\hat{u}_{\text{aws}}^{\text{Gaussian}}$  as well as the SMRE  $\hat{u}_{0.9}$ ,  $\hat{u}_{0.75}$  and  $\hat{u}_{0.5}$  are depicted. It is evident that the SMRE matches the smoothness of the true object much better than the other estimators while the essential features of the signal (such as peak location and peak height) are preserved. In particular, almost no additional local extrema are generated by our approach which stays

in obvious contrast to the other methods. Moreover, we point out that the SMRE are quite robust w.r.t. the choice of the confidence level  $\alpha$ .

We verify this behavior by a simulation study in Table 1. For different noise levels ( $\sigma = 0.1, 0.3$  and  $0.5$ ) we compare the MISE, MIAE and MSB. Additionally, we compute the *mean number of local maxima (MLM)* of  $\hat{u}$  relative to the number of local maxima in  $u^0$  (which is 11). Here  $\hat{u}$  is any of the above estimators. Note that the latter measure (similar to the MSB) takes into account the smoothness of the estimators where a value  $\text{MLM} \gg 1$  indicates too many local maxima and hence a lack of regularity whereas  $\text{MLM} < 1$  implies severe oversmoothing.

	$\sigma = 0.1$				$\sigma = 0.3$				$\sigma = 0.5$			
	MISE	MSB	MIAE	MLM	MISE	MSB	MIAE	MLM	MISE	MSB	MIAE	MLM
$\hat{u}_{L^2}$	0.009	0.008	0.071	11.881	0.046	0.027	0.156	10.915	0.091	0.037	0.213	9.860
$\hat{u}_{\mathbb{B}}$	0.009	0.007	0.070	11.700	0.048	0.026	0.149	10.359	0.094	0.036	0.206	9.135
$\hat{u}_{\text{aws}}^{\text{Triangle}}$	0.007	0.007	0.048	2.551	0.040	0.035	0.112	3.053	0.078	0.058	0.162	3.141
$\hat{u}_{\text{aws}}^{\text{Gauss}}$	0.054	0.040	0.068	1.971	0.062	0.041	0.107	2.230	0.079	0.043	0.149	2.330
$\hat{u}_{0.9}$	0.008	0.004	0.047	1.336	0.056	0.019	0.127	1.273	0.134	0.034	0.207	1.194
$\hat{u}_{0.75}$	0.007	0.004	0.044	1.342	0.050	0.018	0.121	1.290	0.120	0.032	0.196	1.241
$\hat{u}_{0.5}$	0.007	0.004	0.043	1.366	0.046	0.017	0.116	1.290	0.109	0.030	0.186	1.238

TABLE 1. Simulation studies for one dimensional peak data set.

As it becomes apparent from Table 1, the SMR estimators are performing similarly well when compared to the reference estimators as far as the standard measures MISE and MIAE are concerned. For small noise levels ( $\sigma = 0.1$ ) SMR estimation even proves to be superior. The distance measures MSB and MLM, however, are significantly smaller for SMR estimators which indicates that these meet the smoothness of the true object  $u^0$  much better than the reference estimators (cf. Example 3.1 i)). All in all, the simulation results confirm our visual impressions above.

*Implementation Details.* The current index set  $\mathcal{S}$  results in an overall number of constraints in (24) of

$$\#\mathcal{S} = \sum_{i=1}^{100} (1024 - i + 1) = 97.450.$$

As pointed out in Section 2.3, the efficiency of Dykstra's Algorithm can be increased by grouping independent side-conditions, that is side-conditions corresponding to intervals in  $\mathcal{S}$  with empty intersection. For example, the system  $\mathcal{S}$  can be grouped such that the intersection of the corresponding sets  $D_1, \dots, D_M$  in (18) form  $\mathcal{C}$  with

$$M = \sum_{i=1}^{100} i = 5.050.$$

In all our simulations we set  $\tau = 10^{-4}$  and  $\lambda = 1.0$  in Algorithm 1 which results in  $k[\tau] \approx 100$  iterations and an overall computation time of approximately 20 minutes for each SMRE. We note, however, that more than 95% of the computation time is needed for the projection step (11) and that a considerable speed up for the latter could be achieved by parallelization.

**3.2. Image denoising.** In this section we apply SMR estimation to the problem of image denoising, that is non-parametric regression in  $d = 2$  dimensions. In other words, we consider the noise model (23) as in Section 3.1, where the index  $\nu$  ranges over the discrete square  $\{1, \dots, m\}^2$ . In Figure 3 two typical examples for images  $u^0$  and noisy observations  $Y$  are depicted ( $m = 512$  and  $\sigma = 0.1$ , where  $u^0$  is scaled between 0 (black) and 1 (white)).

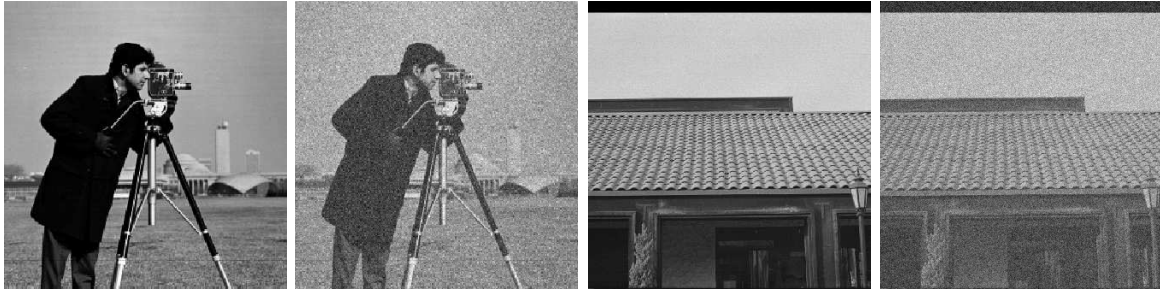


FIGURE 3. Standard test images “cameraman” and “roof” and their noisy counterparts.

We will use the total-variation semi-norm  $J = TV_1^\beta$  as regularization functional ( $\beta = 10^{-8}$ ). Moreover, we choose  $\Lambda$  to be defined as

$$\Lambda(v)_\nu = v_\nu^2, \quad \forall(\nu \in 1, \dots, m^2). \tag{26}$$

The index set  $\mathcal{S}$  is defined to be the collection of all discrete squares with side lengths ranging from 1 to 25 and we set  $\omega^S = c_S \chi_S$  with yet to be defined constants  $c_S$ . Thus, each SMR estimator solves the constrained optimization problem

$$TV_1^\beta(u) \rightarrow \inf! \quad \text{s.t.} \quad \sum_{\nu \in S} c_S (Y - u)_\nu^2 \leq q \quad \forall(S \in \mathcal{S}). \tag{27}$$

We agree upon  $q = 1$  and specify the constants  $c_S$ . To this end, compute for  $s = 1, \dots, 25$  the quantile values

$$q_{\alpha,s} = \inf \left\{ q \in \mathbb{R} : \mathbb{P} \left( \max_{\substack{S \in \mathcal{S} \\ \#S=s}} \sum_{\nu \in S} \varepsilon_\nu^2 \leq q \right) \geq 1 - \alpha \right\} \quad \alpha \in (0, 1)$$

and set  $c_S = q_{\alpha, \#S}^{-1}$ . In other words, the definition of  $c_S$  implies that the true signal  $u^0$  satisfies the constraints in (27) for squares of a fixed side length  $s$  with probability at least  $\alpha$ . We will henceforth denote by  $\hat{u}_\alpha$  a solution of (27). We remark on this particular choice of the parameters  $\omega_S$  below.

*Numerical results and simulations.* In Figures 4 and 5 the oracles  $\hat{u}_{L^2}$  and  $\hat{u}_B$ , the AWS-estimators  $\hat{u}_{\text{aws}}^{\text{Triangle}}$  and  $\hat{u}_{\text{aws}}^{\text{Gauss}}$  as well as the SMRE  $\hat{u}_{0.9}$  are depicted (for the “cameraman” and “roof” test image respectively). It is rather obvious that the  $L^2$ -oracles are not favorable: although relevant details in the image are preserved, smooth parts (as e.g. the sky) still contain random structures. In contrast, the estimator  $\hat{u}_{\text{aws}}^{\text{Gauss}}$  preserves smooth areas but loses essential details. The aws-estimator with triangular kernel performs much better, however, it gives piecewise constant reconstructions of smoothly varying portions of the image, which is clearly undesirable. The SMR estimator and the Bregman-oracle visually perform superior to the other methods. The good performance of the Bregman-oracle indicates that

the symmetric Bregman distance is a good measure for comparing images. In contrast to the Bregman-oracle, the SMRE adapts the amount of smoothing to the underlying image structure: constant image areas are smoothed nicely (e.g. sky portions), while oscillating patterns (e.g. the grass part in the “cameraman” image or the roof tiles in the “roof” image) are recovered.



FIGURE 4. Reconstructions “cameraman” (from left to right):  $L^2$ -oracle  $\hat{u}_{L^2}$ , Bregman-oracle  $\hat{u}_B$ , AWS estimators  $\hat{u}_{\text{aws}}^{\text{Triangle}}$  and  $\hat{u}_{\text{aws}}^{\text{Gaussian}}$ , and SMR estimator  $\hat{u}_{0.9}$ .

We evaluate the performance of the SMR estimators by means of a simulation study. To this end, we compute the MISE, MIAE and MSB and compare these values with the reference estimators. We note, however, that in particular the MISE and MIAE are not well suited in order to measure the distance of images for they are inconsistent with human eye perception. In [41] the *structural similarity index (SSIM)* was introduced for image quality assessment that takes into account luminance, contrast and structure of the images at the same time. We use the author’s implementation<sup>2</sup> which is normalized such that the SSIM lies in the interval  $[-1, 1]$  and is 1 in case of a perfect match. We denote by MSSIM the empirical mean of the SSIM in our simulations.

In Table 2 the simulation results are listed. A first striking fact is the good performance of the  $L^2$ -oracle w.r.t. the MISE and MIAE which is supposed to imply reconstruction properties superior to the other methods. Keeping in mind the visual comparison in Figures 4 and 5, however, this is rather questionable. On the other hand, it becomes evident that the  $L^2$ -oracle has a rather poor performance w.r.t. the MSB which is more suited for measuring image distances. It is therefore remarkable that the SMRE performs equally good as the Bregman-oracle which, in contrast to the SMRE, is not accessible (since  $u^0$  is usually unknown). As far as the structural similarity measure MSSIM is concerned our approach proves to be superior

<sup>2</sup>available at <https://www.ece.uwaterloo.ca/~z70wang/research/ssim/>

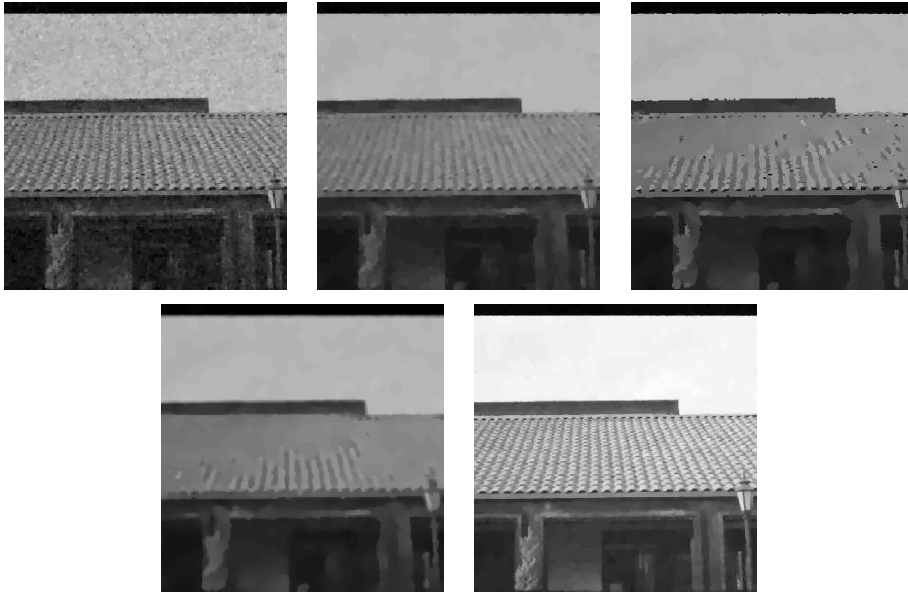


FIGURE 5. Reconstructions “roof” (from left to right):  $L^2$ -oracle  $\hat{u}_{L^2}$ , Bregman-oracle  $\hat{u}_B$ , AWS estimators  $\hat{u}_{aws}^{\text{Triangle}}$  and  $\hat{u}_{aws}^{\text{Gaussian}}$ , and SMR estimator  $\hat{u}_{0.9}$ .

	“cameraman”				“roof”			
	MISE	MSB	MIAE	MSSIM	MISE	MSB	MIAE	MSSIM
$\hat{u}_{L^2}$	0.0017	0.0314	0.0276	0.7739	0.0029	0.0499	0.0383	0.6700
$\hat{u}_B$	0.0023	0.0256	0.0275	0.7995	0.0038	0.0405	0.0391	0.6607
$\hat{u}_{aws}^{\text{Triangle}}$	0.0032	0.0482	0.0308	0.7657	0.0046	0.0702	0.0416	0.6205
$\hat{u}_{aws}^{\text{Gauss}}$	0.0046	0.0470	0.0360	0.7284	0.0053	0.0686	0.0457	0.5668
$\hat{u}_{0.9}$	0.0021	0.0252	0.0297	0.8024	0.0033	0.0374	0.0407	0.7003

TABLE 2. Simulation studies for the test images “cameraman” and “roof”.

to all others. Finally, the simulation results indicate that aws estimation is not favourable for denoising of natural images.

*Notes on the choice of  $\Lambda$  and  $\omega^S$ .* In general, a proper choice of the transformation  $\Lambda$  and of the weight-functions  $\omega^S$  can be achieved by including prior structural information on the true image to be estimated. Substantial parts of natural images, such as photographs, consists of oscillating patterns (as e.g. fabric, wood, hair, grass etc.). This becomes obvious in the standard test images depicted in Figure 3. We claim that for signals that exhibit oscillating patterns, a quadratic transformation  $\Lambda$  as in (26) is favorable, since it yields (compared to the linear statistic studied in Section 3.1) a larger power of the local test statistic on small scales.

In order to illustrate this, we simulate noisy observations  $Y$  of the test images  $u$  in Figure 3 as in (23) with  $\sigma = 0.1$  and compute a *global* estimator  $\hat{u}$  by computing a minimizer of the ROF-functional (22) (with  $\lambda = 0.1$ ). We intend to examine how well over-smoothed regions in  $\hat{u}$  are detected by the MR-statistic  $T(Y - \hat{u})$  as in (3) with two different average functions

(cf. (4))

$$\mu_{1,S}(v) = \left| \sum_{\nu \in S} v_\nu \right| \quad \text{and} \quad \mu_{2,S}(v) = \sum_{\nu \in S} v_\nu^2$$

respectively. For the sake of simplicity we restrict for the moment our considerations on the index set  $\mathcal{S}$  of all  $5 \times 5$  sub-squares in  $\{1, \dots, m\}^2$ . In Figure 6 the local means  $\mu_{i,S}$  of the residuals  $v = Y - \hat{u}$  for the “roof”-image are depicted. To be more precise, the center coordinate of each square  $S \in \mathcal{S}$  is colored according to  $\mu_{i,S}$ . Hence, large values indicate locations where the estimator  $\hat{u}$  is considered over-smoothed according to the statistic. It becomes visually clear that the localization of oversmoothed regions is better for  $\mu_{2,S}$ . This is a good motivation for incorporating the local means of the squared residuals in the SMRE model (2’).

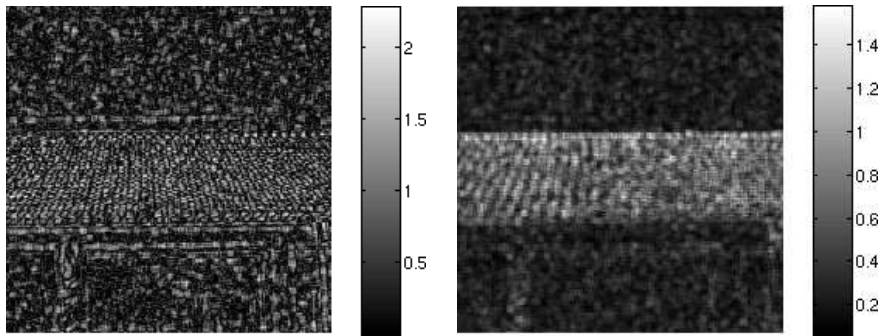


FIGURE 6. Local means  $\mu_{1,S}$  (left) and  $\mu_{2,S}$  (right) of the residuals for “roof” image.

We finally comment on the choice of  $c_S$ . Since  $\varepsilon_\nu$  are independent and normally distributed random variables, the (scaled) average function

$$\sigma^{-2} \mu_S(\varepsilon) = \sum_{\nu \in S} \left( \frac{\varepsilon_\nu}{\sigma} \right)^2$$

is  $\chi^2$  distributed with  $\#S$  degrees of freedom. Note that the distribution of  $\sigma^{-2} \mu_S(\varepsilon)$  is identical only for sets  $S$  of the same scale  $\#S$ . As a consequence of this, it is likely that certain scales dominate the supremum in the MR-statistic  $T$  which spoils the multiscale properties of our approach. As a way out, we compute normalizing constants *for each scale separately*.

An alternative approach would be to search for transformations that turn  $\mu_S(\varepsilon)$  into almost identically distributed random variables. Logarithmic and  $p$ -root transformations are often employed for this purpose (see e.g. [23]), our approach, however, has proved to be superior to these methods.

*Implementation Details.* The current index set  $\mathcal{S}$  results in an overall number of constraints in (27) of

$$\#\mathcal{S} = \sum_{i=1}^{25} (512 - i + 1)^2 = 6.251.300.$$

Again by grouping independent side-conditions, the system  $\mathcal{S}$  can be grouped such that the intersection of the corresponding sets  $D_1, \dots, D_M$  in (18) form  $\mathcal{C}$  with

$$M = \sum_{i=1}^{25} i^2 = 5.525.$$

In all our simulations we set  $\tau = 10^{-4}$  and  $\lambda = 0.25$  in Algorithm 1 which results in  $k[\tau] \approx 30$  iterations and a overall computation time of approximately 2 hours for each SMRE. Hence, parallelization is clearly desirable in this case.

**3.3. Deconvolution.** Another interesting class of problems which can be approached by means of SMR-estimation are deconvolution problems. To be more precise, we assume that  $K$  is a convolution operator, that is

$$(Ku)_\nu = (k * u)_\nu = \sum_{m \in \mathbb{R}^d} k_{\nu-m} u_m$$

where  $k$  is a square-summable kernel on the lattice  $\mathbb{Z}^d$  and  $u \in H$  is extended by zero-padding. We will focus on the situation where  $k$  is a circular Gaussian kernel with standard deviation  $\sigma$  given by

$$k_\nu = \frac{1}{(\sqrt{2\pi}\sigma)^d} e^{-\frac{\sum_{i=1}^d \nu_i^2}{2\sigma^2}}. \quad (28)$$

With  $Z = Y + \lambda p_{k-1} + v_k$ , the primal step (12) in Algorithm 1 amounts to solve

$$u_k \leftarrow \operatorname{argmin}_{u \in H} \frac{1}{2} \sum_{\nu \in X} ((Ku)_\nu - Z_\nu)^2 + \lambda J(u),$$

where we choose  $J$  to be as in (20) and apply the techniques described in [40] for the numerical solution.

In order to illustrate the performance of our approach in practical applications, we give an example from confocal microscopy, nowadays a standard technique in fluorescence microscopy (cf. [34]). When recording images with this kind of microscope, the original object gets blurred by a Gaussian kernel (in first order). The observations (photon counts) can be modeled as a Poisson process, i.e.

$$Y_\nu = \operatorname{Pois}((Ku^0)_\nu), \quad \nu \in X. \quad (29)$$

The image depicted in Figure 7(a) shows a recording of a PtK2 cell taken from the kidney of *potorous tridactylus*. Before the recording, the protein  $\beta$ -tubulin was tagged with a fluorescent marker such that it can be traced by the microscope. The image in 7(a) shows an area of  $18 \times 18 \mu\text{m}^2$  at a resolution of  $798 \times 798$  pixel. The point spread function of the optical system can be modeled as a Gaussian kernel with full width at half maximum of 230nm, which corresponds to  $\sigma = 4.3422$  in (28).

Note that (29) does not fall immediately into the range of models covered by (1). We will adapt the present situation to the SMRE methodology described in Section 1 by standardization and consider instead of (2) the modified problem

$$J(u) \rightarrow \inf! \quad \text{s.t.} \quad T \left( \frac{Y - Ku}{\sqrt{Ku}} \right) \leq 1 \quad (30)$$

where the division is understood pointwise. Clearly, the problem of finding a solution of (30) is much more involved than solving (2) for the constraints being *nonconvex*: firstly, the functional  $G$  as defined in (8) is nonconvex as a consequence of which the convergence result

in Theorem 2.2 does not apply and secondly Dykstra’s projection algorithm as described in Section 2 cannot be employed.

We propose the following ansatz in order to circumvent this problem: instead of projecting onto the intersection  $\mathcal{C}$  of sets  $C_S$  as described in (14), we now project in the  $k$ -th step of Algorithm 1 onto

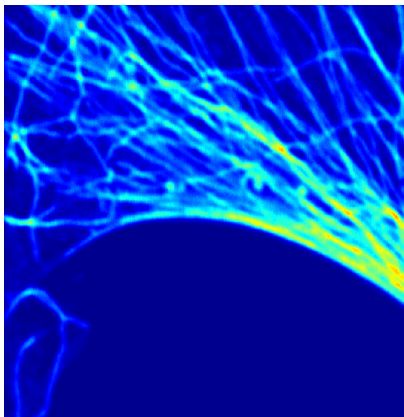
$$\mathcal{C}_P[k] = \bigcap_{n=1}^N C_{P,S}[k] \quad \text{where} \quad C_{P,S}[k] = \left\{ v \in H : \mu_S \left( v / \sqrt{K u_k} \right) \leq q \right\}. \quad (31)$$

with a pointwise division by the square root of  $K u_k$ . Put differently, in the  $k$ -th step of Algorithm 1 we use the previous estimate  $u_k$  of  $u^0$  as a *lagged standardization* in order to approximate the constraints in (30).

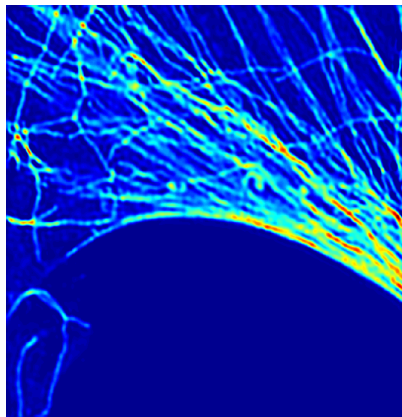
We note that while with this modification Dykstra’s algorithm becomes applicable again, the projection problem (11) now changes in each iteration step of Algorithm 1. As a consequence, Theorem 2.2 does not hold anymore after this modification, either. So far, we have not come up with a similar convergence analysis.

We compute the SMRE  $\hat{u}$  by employing Algorithm 1 with the modifications described above. As in the denoising examples in Section 3.2 the index set  $\mathcal{S}$  consists of all squares with the side-lengths  $\{1, \dots, 25\}$  and we choose  $\omega^S = \chi_S$  and  $\Lambda = \text{Id}$ . We note, that this results in an overall number of  $\#\mathcal{S} = 95\,436\,200$  inequality constraints. The constant  $q$  are chosen as in (25), where we assume that  $\varepsilon_{\nu}$  are independent and standard normally distributed r.v.

In Algorithm 1 we set  $\lambda = 0.05$  and compute 100 steps. We observe that after a few iterations ( $\sim 15$ ) the error  $\tau$  falls below  $10^{-3}$  and almost stagnates thereafter, which is due to the fact that we do not increase the accuracy in the subroutines for (12) and (11). Each iteration step in Algorithm 1 approximately takes 10 minutes, where 90% of the computation time is needed for (12). The result is depicted in Figure 7(b).



(a) Fluorescence microscopy data of a PtK2 cell in *potorous tridactylus* kidney. The bright filaments indicate the location of the protein  $\beta$ -tubulin.



(b) SMRE  $\hat{u}$ : fully automated and locally adaptive deconvolution of microscopy data.

FIGURE 7. Reconstruction of confocal microscopy data.

The benefits of our method are twofold:

- i) The amount of regularization is chosen in a *completely data driven way*. The only parameter to be selected is the level  $\alpha$  in (25). Note that the parameter  $\lambda$  in Algorithm 1 has no effect on the output (though it has an effect on the number of iterations needed and the numerical stability).
- ii) The reconstruction has an appealing locally adaptive behavior which in the present example mainly concerns the gaps between the protein filaments: whereas the marked  $\beta$ -tubulin is concentrated in regions of basically one scale, the gaps in between actually make up the multiscale nature of the image.

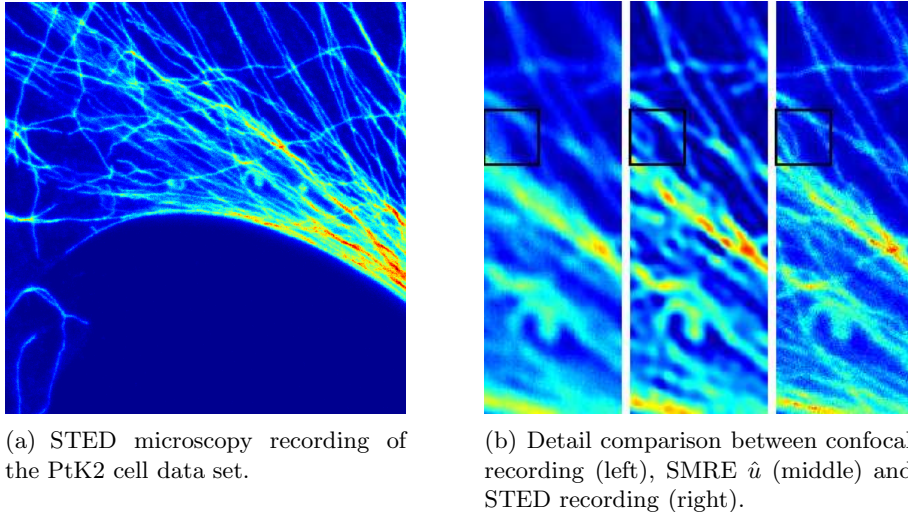


FIGURE 8. Comparison with high-resolution STED microscopy data.

In the present situation we are in the comfortable position to have a reference image at hand by means of which we can evaluate the result of our method: STED (STimulated Emission Depletion) microscopy constitutes a relatively new method, that is capable of recording images at a physically 5-10 times higher resolution as confocal microscopy (see [24, 25]). Hence a STED image of this object may serve as “gold standard” reference image.

Figure 8(a) depicts a STED recording of the PtK2 cell data set in Figure 7(a). The comparison of the SMRE  $\hat{u}(0.90)$  with the STED recording in Figure 8(b) shows that our SMRE technique chooses a reasonable amount of regularization: no artifacts due to under-regularization are generated and on the other hand almost all relevant geometrical features that are present in the high-resolution STED recording become visible in the reconstruction. In particular, we note that filament bifurcations (one such bifurcation is marked by a black box in Figure 8(b)) become apparent in our reconstruction that are not visible in the recorded data.

#### ACKNOWLEDGEMENT

K.F. is supported by the DFG-SNF Research Group FOR916 *Statistical Regularization* (Z-Project). P.M is supported by the BMBF project 03MUPAH6 *INVERS*. A.M and P.M. are supported by the SFB755 *Photonic Imaging on the Nanoscale* and the SFB803 *Functionality Controlled by Organization in and between Membranes*. The authors would like to thank

S. Hell, A. Egner and A. Schoenle (Department of NanoBiophotonics, Max Planck Institute for Biophysical Chemistry, Göttingen) for providing the microscopy data and L. Dümbgen (University of Bern) for stimulating discussions.

#### APPENDIX A. PROOFS

In this section we shall give the proofs of Theorem 2.2 and Corollary 2.4. We note, that convergence of Algorithm 1 is a classical subject in optimization theory and a proof can e.g. be found in [18, Chap. III Thm. 4.1]. However, in order to apply these results, it is necessary that certain regularity conditions for  $J$  hold, that are not realistic for our purposes (as e.g. in the case of total-variation regularization). The assertion of Theorem 2.2 is a modification of the standard results.

To this end we remind the reader of the definition of the *subdifferential* (or generalized derivative)  $\partial F$  of a convex function  $F : V \rightarrow \mathbb{R}$  on a real Hilbert-space  $V$ :

$$\xi \in \partial F(v) \iff F(w) \geq F(v) + \langle \xi, w - v \rangle_V \quad \forall (w \in V).$$

If  $\xi \in \partial F(v)$ , then  $\xi$  is called *subgradient* of  $F$  at  $v$ . It follows from [17, Chap III, Prop. 3.1 and Prop. 4.1] that the Lagrangian  $L$  (and hence also the augmented Lagrangian  $L_\lambda$ ) has a saddle-point  $(\hat{u}, \hat{v}, \hat{p}) \in U \times H \times H$  if and only if

$$K\hat{u} + \hat{v} = Y, \quad K^*\hat{p} \in \partial J(\hat{u}) \quad \text{and} \quad \hat{p} \in \partial G(\hat{v}). \quad (32)$$

*Proof of Theorem 2.2.* Let us assume that  $(\hat{u}, \hat{v}, \hat{p})$  is a saddle point of the augmented Lagrangian  $L_\lambda(u, v, p)$  as defined in (9) and that  $\{(u_k, v_k, p_k)\}_{k \in \mathbb{N}}$  is a sequence generated by Algorithm 1. Further, we introduce the notation

$$\bar{u}_k := u_k - \hat{u}, \quad \bar{v}_k := v_k - \hat{v} \quad \text{and} \quad \bar{p}_k := p_k - \hat{p}.$$

From now on, we assume that  $k \geq 1$ . By repeating the steps (5.6)-(5.25) in the proof of [18, Chap. III Thm. 4.1], it follows that

$$\begin{aligned} & \left( \|\bar{p}_{k-1}\|^2 + \lambda^{-2} \|K\bar{u}_{k-1}\|^2 \right) - \left( \|\bar{p}_k\|^2 + \lambda^{-2} \|K\bar{u}_k\|^2 \right) \\ & \geq \lambda^{-2} \left( \|K\bar{u}_k + \bar{v}_k\|^2 + \|K\bar{u}_{k-1} - K\bar{u}_k\|^2 \right) \end{aligned} \quad (33)$$

Summing over  $k$  and keeping in mind that  $K\bar{u}_k + \bar{v}_k = Ku_k + v_k - Y$  and  $K\bar{u}_{k-1} - K\bar{u}_k = Ku_{k-1} - Ku_k$  shows

$$\sum_{k=1}^{\infty} \|Ku_k + v_k - Y\|^2 + \|Ku_{k-1} - Ku_k\|^2 \leq \lambda^2 \|\hat{p}\|^2 + \|K\hat{u}\|^2 < \infty \quad (34)$$

where we have used that  $\bar{u}_0 = \hat{u}$  and  $\bar{p}_0 = \hat{p}$ . In other words, this shows that

$$\|Ku_k + v_k - Y\| = o(k^{-1/2}) \quad \text{and} \quad \|Ku_{k-1} - Ku_k\| = o(k^{-1/2}).$$

Furthermore, it follows from (33) that  $\|\bar{p}_k\|^2 + \lambda^{-2} \|K\bar{u}_k\|^2$  is nonincreasing and thus bounded. This together with  $\|Ku_k + v_k - Y\| = o(k^{-1/2})$  implies that

$$\max(\|Ku_k\|, \|v_k\|, \|p_k\|) = \mathcal{O}(1).$$

Together with the optimality condition for (12) this in turn implies that for an arbitrary  $u \in H$

$$J(u_k) \leq J(u) + \lambda^{-1} \langle Ku_k + v_k - Y - \lambda p_{k-1}, Ku - Ku_k \rangle = \mathcal{O}(1). \quad (35)$$

Summarizing, we find that

$$\max_{S \in \mathcal{S}} \mu_S(Ku_k - Y) + J(u_k) \leq \max_{S \in \mathcal{S}} \|\omega^S\| \|\Lambda(Ku_k - Y)\| + J(u_k) \leq c < \infty$$

for a suitably chosen constant  $c \in \mathbb{R}$ , since  $\Lambda$  is supposed to be continuous. Thus, it follows from Assumption A that  $\{u_k\}_{k \in \mathbb{N}}$  is bounded and hence sequentially weakly compact. Now, let  $(\tilde{u}, \tilde{v}, \tilde{p})$  be a weak cluster point of  $\{(u_k, v_k, p_k)\}_{k \in \mathbb{N}}$  and recall that  $(\hat{u}, \hat{v}, \hat{p})$  was assumed to be a saddle point of the augmented Lagrangian  $L_\lambda$ . Setting  $u = \hat{u}$  in (35) thus results in

$$\begin{aligned} J(u_k) &\leq J(\hat{u}) + \lambda^{-1} \langle Ku_k + v_k - Y, K\hat{u} - Ku_k \rangle + \langle p_{k-1}, Ku_k - K\hat{u} \rangle \\ &= J(\hat{u}) + \langle p_{k-1}, Ku_k - K\hat{u} \rangle + o(k^{-1/2}) \end{aligned} \quad (36)$$

Using the relation  $K\hat{u} + \hat{v} = Y$  we further find

$$\begin{aligned} \langle p_{k-1}, Ku_k - K\hat{u} \rangle &= \langle p_{k-1}, Ku_k - Y + \hat{v} \rangle \\ &= \langle p_{k-1}, Ku_k + v_k - Y \rangle - \langle p_{k-1}, v_k - \hat{v} \rangle = o(k^{-1/2}) - \langle p_{k-1}, v_k - \hat{v} \rangle \end{aligned} \quad (37)$$

From the definition of  $v_k$  in (11) it follows that

$$\langle Y + \lambda p_{k-1} - (Ku_{k-1} + v_k), \hat{v} - v_k \rangle \leq 0$$

which in turn implies that

$$\begin{aligned} -\langle p_{k-1}, v_k - \hat{v} \rangle &\leq \lambda^{-1} \langle Y - (Ku_{k-1} + v_k), v_k - \hat{v} \rangle \\ &= \lambda^{-1} \langle Y - (Ku_k + v_k), v_k - \hat{v} \rangle + \lambda^{-1} \langle Ku_k - Ku_{k-1}, v_k - \hat{v} \rangle = o(k^{-1/2}) \end{aligned} \quad (38)$$

Combining (36), (37) and (38) gives

$$\limsup_{k \rightarrow \infty} J(u_k) \leq J(\hat{u}).$$

Now, choose a subsequence  $\{u_{\rho(k)}\}_{k \in \mathbb{N}}$  such that  $u_{\rho(k)} \rightharpoonup \tilde{u}$ . Since  $J$  is convex and lower semi-continuous it is also weakly lower semi-continuous and hence the previous estimate yields

$$J(\tilde{u}) \leq \liminf_{k \rightarrow \infty} J(u_{\rho(k)}) \leq J(\hat{u}).$$

Moreover, we have that  $v_{\rho(k)} \in \mathcal{C}$  for all  $k \in \mathbb{N}$ . Since  $\mathcal{C}$  is closed we conclude that  $\hat{v} \in \mathcal{C}$ . Since  $K\tilde{u} + \tilde{v} = Y$  this shows that  $(\tilde{u}, \tilde{v})$  solves (7) and thus  $J(\tilde{u}) = J(\hat{u})$ .  $\square$

*Proof of Corollary 2.4.* First, observe that the estimate in (33) implies that the sequence  $\|\bar{p}_k\|^2 + \lambda^{-2} \|K\bar{u}_k\|^2$  is nonincreasing. Since  $u_0 = p_0 = 0$ , we have that

$$\|p_k\| \leq 2 \|\hat{p}\| + \lambda^{-1} \|K\hat{u}\|,$$

where  $(\hat{u}, \hat{v}, \hat{p})$  is an arbitrary saddle point of  $L_\lambda(u, v, p)$ . Assume that  $\tau > 0$  and that  $k = k[\tau]$  is such that

$$\max(\|Ku_k + v_k - Y\|, \|Ku_{k-1} - Ku_k\|) \leq \tau.$$

Then, it follows from (36) that

$$\begin{aligned} J(u_k) &\leq J(\hat{u}) + \lambda^{-1} \langle Ku_k + v_k - Y, K\hat{u} - Ku_k \rangle + \langle p_{k-1}, Ku_k - K\hat{u} \rangle \\ &\leq J(\hat{u}) + \lambda^{-1} \tau^2 + \|p_{k-1}\| \tau \\ &\leq J(\hat{u}) + \left( \frac{\tau + \|K\hat{u}\|}{\lambda} + 2 \|\hat{p}\| \right) \tau. \end{aligned} \quad (39)$$

Now, observe that from the definition of the subgradient and (32), it follows that  $J(u_k) \geq J(\hat{u}) + \langle K^* \hat{p}, u_k - \hat{u} \rangle$  and that  $\langle \hat{p}, v_k - \hat{v} \rangle \leq 0$ . This and the fact that  $K\hat{u} + \hat{v} = Y$  implies that

$$\begin{aligned} 0 &\leq J(u_k) - J(\hat{u}) - \langle K^* \hat{p}, u_k - \hat{u} \rangle \\ &= J(u_k) - J(\hat{u}) - \langle \hat{p}, Ku_k + v_k - Y \rangle + \langle \hat{p}, K\hat{u} + \hat{v} - Y \rangle + \langle \hat{p}, v_k - \hat{v} \rangle \\ &\leq J(u_k) - J(\hat{u}) + \|\hat{p}\| \tau. \end{aligned}$$

This together with (39) finally proves the first part of the assertion.

Now, assume that  $J(u) = \frac{1}{2} \|Lu\|_V^2$ . Then it follows (see e.g. [20, Lem. 2.4]) that the subdifferential  $\partial J(\hat{u})$  consists of the single element  $L^*L\hat{u}$ . Hence the extremality relations (32) imply that  $K^*\hat{p} = L^*L\hat{u}$ . Now it is easy to observe that

$$J(u_k) - J(\hat{u}) - \langle K^* \hat{p}, u_k - \hat{u} \rangle = \frac{1}{2} \|L(u_k - \hat{u})\|_V^2.$$

□

#### REFERENCES

- [1] S. R. Becker, E. J. Candès, and M. Grant. Templates for convex cone problems with applications to sparse signal recovery, 2010. Available at <http://arxiv.org/abs/1009.2065>.
- [2] M. Bertalmio, V. Caselles, B. Rougé, and A. Solé. TV based image restoration with local constraints. *J. Sci. Comput.*, 19(1-3):95–122, 2003. Special issue in honor of the sixtieth birthday of Stanley Osher.
- [3] J. P. Boyle and R. L. Dykstra. A method for finding projections onto the intersection of convex sets in Hilbert spaces. In *Advances in order restricted statistical inference (Iowa City, Iowa, 1985)*, volume 37 of *Lecture Notes in Statist.*, pages 28–47. Springer, Berlin, 1986.
- [4] L. Boysen, A. Kempe, V. Liebscher, A. Munk, and O. Wittich. Consistencies and rates of convergence of jump-penalized least squares estimators. *Ann. Statist.*, 37(1):157–183, 2009.
- [5] E. Candès and T. Tao. The Dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *Ann. Statist.*, 35(6):2313–2351, 2007.
- [6] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159, 2001.
- [7] I. Csiszár. Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems. *Ann. Statist.*, 19(4):2032–2066, 1991.
- [8] P. L. Davies, U. Gather, M. Meise, D. Mergel, and T. Mildenerger. Residual-based localization and quantification of peaks in x-ray diffractograms. *Ann. Appl. Stat.*, 2(3):861–886, 2008.
- [9] P. L. Davies and A. Kovac. Local extremes, runs, strings and multiresolution. *Ann. Statist.*, 29(1):1–65, 2001. With discussion and rejoinder by the authors.
- [10] P. L. Davies, A. Kovac, and M. Meise. Nonparametric regression, confidence regions and regularization. *Ann. Statist.*, 37(5B):2597–2625, 2009.
- [11] F. Deutsch and H. Hundal. The rate of convergence of Dykstra’s cyclic projections algorithm: the polyhedral case. *Numer. Funct. Anal. Optim.*, 15(5-6):537–565, 1994.
- [12] D. C. Dobson and C. R. Vogel. Convergence of an iterative method for total variation denoising. *SIAM J. Numer. Anal.*, 34(5):1779–1791, 1997.

- [13] Y. Dong, M. Hintermüller, and M. Rincon-Camacho. Automated regularization parameter selection in a multi-scale total variation model for image restoration. Technical report, Institute of Mathematics and Scientific Computing, 2008. IFB Report 22.
- [14] L. Dümbgen and R. B. Johns. Confidence bands for isotonic median curves using sign tests. *Journal of Computational and Graphical Statistics*, 13(2):519–533, 2004.
- [15] L. Dümbgen and V. G. Spokoiny. Multiscale testing of qualitative hypotheses. *Ann. Statist.*, 29(1):124–152, 2001.
- [16] L. Dümbgen and G. Walther. Multiscale inference about a density. *Ann. Statist.*, 36(4):1758–1785, 2008.
- [17] I. Ekeland and R. Temam. *Convex analysis and variational problems*, volume 1 of *Studies in Mathematics and its Applications*. North-Holland Publishing Co., Amsterdam-Oxford, 1976.
- [18] M. Fortin and R. Glowinski. *Augmented Lagrangian methods*, volume 15 of *Studies in Mathematics and its Applications*. North-Holland Publishing Co., Amsterdam, 1983. Applications to the numerical solution of boundary value problems, Translated from the French by B. Hunt and D. C. Spicer.
- [19] K. Frick, P. Marnitz, and A. Munk. Shape constrained regularisation by statistical multiresolution for inverse problems, 2010. Available at <http://arxiv.org/abs/1003.3323>.
- [20] K. Frick and O. Scherzer. Regularization of ill-posed linear equations by the non-stationary Augmented Lagrangian Method. *J. Integral Equations Appl.*, 22(2):217–257, 2010.
- [21] N. Gaffke and R. Mathar. A cyclic projection algorithm via duality. *Metrika*, 36:29–54, 1989.
- [22] M. Grasmair. The equivalence of the taut string algorithm and BV-regularization. *J. Math. Imaging Vision*, 27(1):59–66, 2007.
- [23] D. M. Hawkins and R. Wixley. A note on the transformation of chi-squared variables to normality. *Amer. Statist.*, 40(4):296–298, 1986.
- [24] S. W. Hell. Far-Field Optical Nanoscopy. *Science*, 316(5828):1153–1158, 2007.
- [25] S. W. Hell and J. Wichmann. Breaking the diffraction resolution limit by stimulated emission: stimulated-emission-depletion fluorescence microscopy. *Opt. Lett.*, 19(11):780–782, 1994.
- [26] T. Hotz, P. Marnitz, R. Stichtenoth, L. Davies, Z. Kabluchko, and A. Munk. Locally adaptive image denoising by a statistical multiresolution criterion. available at <http://arxiv.org/abs/1001.5447>, 2009.
- [27] G. M. James, P. Radchenko, and J. Lv. Dasso: connections between the dantzig selector and lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 71:127–142.
- [28] Z. Kabluchko. Title: Extremes of the standardized gaussian noise. available at <http://arxiv.org/abs/1007.0312>, 2010.
- [29] Z. Kabluchko and A. Munk. Shao’s theorem on the maximum of standardized random walk increments for multidimensional arrays. *ESAIM Probab. Stat.*, 13:409–416, 2009.
- [30] Z. Lu, T. K. Pong, and Y. Zhang. An alternating direction method for finding Dantzig selectors, 2010. Available at <http://arxiv.org/abs/1011.4604v1>.
- [31] E. Mammen and S. van de Geer. Locally adaptive regression splines. *Ann. Statist.*, 25(1):387–413, 1997.
- [32] T. Mildenerger. A geometric interpretation of the multiresolution criterion in nonparametric regression. *J. Nonparametr. Stat.*, 20(7):1048–5252, 2008.

- [33] A. Munk, N. Bissantz, T. Wagner, and G. Freitag. On difference-based variance estimation in nonparametric regression when the covariate is high dimensional. *Journal Of The Royal Statistical Society Series B*, 67(1):19–41, 2005.
- [34] J. B. Pawley. *Handbook of Biological Confocal Microscopy*. Springer, 2006.
- [35] J. Polzehl and V. G. Spokoiny. Adaptive weights smoothing with applications to image restoration. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 62(2):335–354, 2000.
- [36] J. K. Romberg. The dantzig selector and generalized thresholding. In *CISS*, pages 22–25, 2008.
- [37] Q. M. Shao. On a conjecture of Révész. *Proc. Amer. Math. Soc.*, 123(2):575–582, 1995.
- [38] D. Siegmund and B. Yakir. Tail probabilities for the null distribution of scanning statistics. *Bernoulli*, 6(2):191–213, 2000.
- [39] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1994.
- [40] C. R. Vogel. *Computational methods for inverse problems*, volume 23 of *Frontiers in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. With a foreword by H. T. Banks.
- [41] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- [42] S. Xu. Estimation of the convergence rate of Dykstra’s cyclic projections algorithm in polyhedral case. *Acta Math. Appl. Sinica (English Ser.)*, 16(2):217–220, 2000.