

Explorations on anisotropic regularisation of dynamic inverse problems by bilevel optimisation

Martin Benning¹, Carola-Bibiane Schönlieb¹, Tuomo Valkonen^{1,2}, Verner Vlačić¹

¹ Department of Applied Mathematics and Theoretical Physics, University of Cambridge, United Kingdom

² Department of Mathematical Sciences, University of Liverpool, United Kingdom

Abstract. We explore anisotropic regularisation methods in the spirit of [1]. Based on ground truth data, we propose a bilevel optimisation strategy to compute the optimal regularisation parameters of such a model for the application of video denoising. The optimisation poses a challenge in itself, as the dependency on one of the regularisation parameters is non-linear such that the standard existence and convergence theory does not apply. Moreover, we analyse numerical results of the proposed parameter learning strategy based on three exemplary video sequences and discuss the impact of these results on the actual modelling of dynamic inverse problems.

Submitted to: *Inverse Problems*

1. Introduction

In this paper, we employ bilevel optimisation to explore different choices of spatial- and temporal regularisation in a variational image reconstruction model.

Variational regularisation methods are extremely popular and versatile tools when it comes to computing approximate solutions to ill-posed inverse problems. Given the assumption of normal distributed noise, they usually have the form of a generalised Tikhonov-type regularisation, i.e.

$$u_{\alpha} \in R(\alpha) := \arg \min_{u \in L^2(\Omega) \cap \mathcal{U}} \left\{ \frac{1}{2} \|Ku - g\|_2^2 + G(u; \alpha) \right\}, \quad (1.1)$$

where $K \in \mathcal{L}(L^2(\Omega), L^2(\Sigma))$ is a bounded, linear operator mapping from the Hilbert space $L^2(\Omega)$ to the Hilbert space $L^2(\Sigma)$, with Ω and Σ being bounded and connected domains. The function $g \in L^2(\Sigma)$ represents the given measurement data, and the norm $\|\cdot\|_2$ simply denotes the $L^2(\Sigma)$ -norm. Given a signal $u \in \mathcal{U}$ and regularisation parameters $\alpha \in \mathcal{V}$, the functional $G : \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}$ represents the regulariser, for Banach spaces \mathcal{U} and \mathcal{V} .

Particularly in imaging and image processing applications, proper, lower semi-continuous (l.s.c), convex and non-smooth regularisers have attracted a lot of attention over the last two decades. Various types of total variation regularisation [2] and ℓ^1 regularisation of unitary transformed signals [3] have been proposed, in order to exploit sparsity of a signal with respect to a given representation.

Despite allowing for significant improvements in terms of reconstruction quality, non-smooth regularisation methods suffer from introducing systematic modelling artefacts like any other regularisation method; in case of total variation regularisation for instance, the regularisation method is well-known to introduce piecewise constant approximations of noisy, non-constant regions, which is known as the stair-casing effect (cf. [4, Section 4.2]).

In order to compensate for modelling artefacts, the concept of infimal convolutions can be used for combining the advantages of different regularisers into one. The infimal convolution of two proper, l.s.c. and convex functionals J_1 and J_2 is defined as

$$(J_1 \square J_2)(u) := \inf_{u=v+w} J_1(v) + J_2(w). \quad (\text{IC})$$

In [5, 6, 7, 8] it has been shown that an infimal convolution of the total variation and a higher-order total variation is beneficial for fighting the stair-casing phenomenon, and that infimal convolutions in general are useful in order to reconstruct functions that are additive compositions of functions which can individually be recovered by different regularisers.

Recently, infimal convolution has been considered as a suitable model for handling dynamic inverse problems [1]. Holler and Kunisch have proposed to use infimal convolution in order to combine regularisation functionals that are suitable for either spatial or temporal regularisation, in order to create an appropriate spatio-temporal regulariser. In particular, their model involves a regulariser G in (1.1) which constitutes an infimal-convolution of total variation functionals of weighted spatial and temporal derivatives. This will be specified in Section 2. The use of an infimal convolution between dominantly spatial and dominantly temporal regularisation terms not only allows the reconstruction of a regularised dynamic image sequence, but

also the decomposition of the latter into a sequence of images encoding dominantly temporal information and another sequence encoding dominantly spatial information.

In what follows we want to learn optimal decompositions in this regularisation in the context of video de-noising by an appropriate bilevel optimisation approach. In the context of (1.1) the associated bilevel learning problem we will discuss looks for an optimal parameter vector α that solves for some convex, proper, weak* lower semicontinuous cost functional $F : X \rightarrow \mathbb{R}$ the problem

$$\min_{\alpha \in \mathcal{P}} F(u_\alpha) \tag{P}$$

subject to u_α being a solution of (1.1).

In the context of computer vision and image processing bilevel optimisation is considered as a supervised learning method that optimally adapts itself to a given dataset of measurements and desirable solutions. In [9, 10, 11, 12], for instance the authors consider bilevel optimization for finite dimensional Markov random field models. In inverse problems the optimal inversion and experimental acquisition setup is discussed in the context of optimal model design in works by Haber, Horesh and Tenorio [13, 14], as well as Ghattas et al. [15, 16]. Recently parameter learning in the context of functional variational regularisation models (1.1) also entered the image processing community with works by the authors [17, 18, 19, 20, 21, 22], Kunisch, Pock and co-workers [23, 24, 25], Chung et al. [26], Hintermüller et al. [27] and others [28, 29, 30]. Interesting recent works also include bilevel learning approaches for image segmentation [31], learning and optimal setup of support vector machines [32] and learning discrete reaction-diffusion filters [33].

What we show is closest in flavour to recent applications of bilevel optimisation to supervised learning of optimal parameters in a total variation type regularisation model [34, 23, 35, 36]. The main difference in the theory and computational realisation to these works and the model discussed in this paper is due to the nonlinear dependence of the lower level problem on the parameter vector α that the model is optimised for. To our knowledge, this is the first publication that deals with learning of non-linear regularisation parameters in the context of regularisation of inverse problems.

The paper is organised as follows. First, we are going to recall the concept of spatio-temporal regularisation via infimal convolutions of regularisers. Then, we are going to present the bilevel optimisation framework for learning the regularisation parameters of the infimal convolution regularisers. Subsequently, we are going to address the numerical aspects of the parameter learning strategy. We then conclude with numerical examples and their discussion. We particularly want to address the question of how realistic the assumption of the coupling of the spatial and temporal regularisation is, given three different types of images sequences.

2. Regularisation of dynamic inverse problems

Let $g = (g_1, \dots, g_m) \in L^1(\Omega; \mathbb{R}^m)$. We denote by $\|\cdot\|_{2,1}$ the $L^1(\Omega)$ -norm of the two-norm of vector-valued functions, namely

$$\|g\|_{2,1} := \int_{\Omega} \|g(x)\|_2 dx,$$

where $\|g(x)\|_2 = \sqrt{g_1(x)^2 + \dots + g_m(x)^2}$. Based on [1], we introduce the anisotropic derivative ∇_κ and its negative adjoint div_κ defined for a scalar $\kappa \in (0, 1)$ as

$$\begin{aligned} \nabla_\kappa &= \left(\kappa \frac{\partial}{\partial x}, \kappa \frac{\partial}{\partial y}, (1 - \kappa) \frac{\partial}{\partial t} \right), \quad \text{and} \\ \text{div}_\kappa &= \kappa \frac{\partial}{\partial x} + \kappa \frac{\partial}{\partial y} + (1 - \kappa) \frac{\partial}{\partial t}, \end{aligned} \tag{2.1}$$

With these, and for $u \in W^{1,1}(\Omega)$ we define the following dynamic regularisation functionals: Here (IC-TVTV) is a direct adaptation of the ICTV functional proposed

$$G(u; \alpha_1, \alpha_2, \kappa) = \inf_{u=v+w} \alpha_1 \|\nabla_\kappa v\|_{2,1} + \alpha_2 \|\nabla_{1-\kappa} w\|_{2,1} \tag{IC-TVTV}$$

$$G(u; \alpha_1, \alpha_2, \kappa) = \inf_{u=v+w} \frac{\alpha_1}{2} \|\nabla_\kappa v\|_2^2 + \alpha_2 \|\nabla_{1-\kappa} w\|_{2,1} \tag{IC- L^2 TV}$$

$$G(u; \alpha_1, \alpha_2) = \alpha_1 \|\nabla u\|_{2,1} + \alpha_2 \left\| \frac{\partial}{\partial t} u \right\|_1 \tag{Rigid TVTV}$$

$$G(u; \alpha_1, \alpha_2) = \frac{\alpha_1}{2} \|\nabla u\|_2^2 + \alpha_2 \left\| \frac{\partial}{\partial t} u \right\|_1 \tag{Rigid L^2 TV}$$

Table 1: The different dynamic regularisers used throughout this paper.

in [1]. Regulariser (IC- L^2 TV) is a modification that allows for different regularisation models; in case we have picked two rather complementary regularisation functionals, with the L^2 norm of the gradient complementing the total variation regularisation. Regularisers (Rigid TVTV) and (Rigid L^2 TV) can almost be seen as limiting cases of (IC-TVTV) and (IC- L^2 TV), respectively. Choosing $\kappa \in \{0, 1\}$ and restricting solutions to $v = w$ converts (IC-TVTV) and (IC- L^2 TV) into (Rigid TVTV) and (Rigid L^2 TV).

The basic motivation for these models is a decomposition of a dynamic image sequence into a spatial and a temporal component, such that these components are penalised individually with suitable regularisation functionals. However, in order to allow for space-time correspondence in these sequences, an additional anisotropy parameter κ is introduced that ensures neither one of the components to be penalised by the spatial or the temporal regulariser alone. Speaking of spatial and temporal components, a spatial component is regularised in space only, whereas the temporal component is regularised in only in time. In (2.1) this corresponds to the extreme cases $\kappa = 0$ (only temporal regularisation) and $\kappa = 1$ (only spatial regularisation). The restriction to $\kappa \in (0, 1)$, which further ensures spatial and temporal regularity in both components, also guarantees well-posedness of (1.1).

A major challenge for successful regularisation of dynamic image sequences via any of the regularisers listed in Table 1 is the 'optimal' choice of parameters α_1 , α_2 and κ . On the one hand, this is due to the number of parameters itself, making it difficult to employ heuristic parameter choice rules that may succeed in case of single parameters. On the other hand, the dependency of the model on the regularisation parameter κ is non-linear, making it even harder to optimise for.

Following [34, 23, 35, 36], we propose a learning framework that allows to learn the regularisation parameters from training data. As mentioned before, the key difference is that the dependency of κ is non-linear.

3. Optimising the parameters—the theory

It is not immediately clear, which parameter choices for α_1 , α_2 and κ may be optimal for dynamic inverse problems regularised by one of the anisotropic regularisers given in Table 1. Parameter learning approaches, as discussed in the introduction, provide a means towards studying optimal parameter choices with respect to a known ground truth. We concentrate in particular on the bilevel optimisation framework, first presented in [34, 23] and further studied in the context of multi-parameter regularisers in [35, 36]. The general form therein is given as

$$\min_{\alpha \in \mathcal{P}} \frac{1}{2} \|R(\alpha) - g\|_2^2 \quad \text{s.t.} \quad R(\alpha) = \arg \min_{u \in L^1(\Omega)} \frac{1}{2} \|f - Ku\|_2^2 + G(u; \alpha). \quad (3.1)$$

Here α denotes the vector of parameters we are optimising for, and \mathcal{P} is our space of allowed parameters; generally $\alpha = (\alpha_1, \alpha_2, \kappa)$ with $\mathcal{P} = [0, \infty)^2 \times [0, 1]$ for the infimal convolution regularisers, and $\alpha = (\alpha_1, \alpha_2)$ with $\mathcal{P} = [0, \infty)^2$ for the rigid regularisers. Here, $\|R(\alpha) - g\|_2$ is the cost functional $F(R(\alpha))$ measuring the distance of the denoised solution $R(\alpha)$ from the ground truth original g , and G is a regulariser for the lower level problem of reconstructing u from noisy f .

3.1. Existence of solutions

A general existence and convergence theory for solutions α to (3.1) is presented in [35] when G is linear in α and defined over the space of bounded variation functions, i.e. $\|\nabla u\|_{2,1}$ in Table 1 is replaced by

$$\sup_{\substack{\varphi \in C_0^\infty(\Omega; \mathbb{R}^m) \\ \|\varphi\|_{2,\infty} \leq 1}} \int_{\Omega} u(x) (\operatorname{div}_{\kappa} \varphi)(x) dx,$$

where the notation $\|\cdot\|_{2,\infty}$ is analogous to $\|\cdot\|_{2,1}$, but with the $L^\infty(\Omega)$ -norm replacing the $L^1(\Omega)$ -norm. Our regularisers are not linear in κ however, so an extended theory would be needed. With an additional elliptic regularisation, however, we can still show existence of solutions easily.

Proposition 3.1. *Consider the problem*

$$\min_{\alpha \in \mathcal{P}} \frac{1}{2} \|R(\alpha) - g\|_2^2 \quad \text{s.t.} \quad R(\alpha) = \arg \min_{u \in L^1(\Omega)} \frac{1}{2} \|f - Ku\|_2^2 + G(u; \alpha) + \frac{\epsilon}{2} \|\nabla u\|^2, \quad (3.2)$$

where G is one of the regularisers from Table 1, \mathcal{P} the corresponding admissible set of parameters, and $\epsilon > 0$. Suppose constant functions are not in the null space of K . Then there exists an optimal solution $\alpha \in \mathcal{P}$ to (3.2).

Proof. Note that all regularisers presented in Table 1 are lower semi-continuous with respect to the mutual convergence of α in \mathbb{R}^3 and of u in L^1 . They are moreover continuous with respect to α for fixed u .

Let us first consider the inner problem

$$\arg \min_{u \in L^1(\Omega)} J(u; \alpha) := \frac{1}{2} \|f - Ku\|_2^2 + G(u; \alpha) + \frac{\epsilon}{2} \|\nabla u\|^2 \quad (3.3)$$

first. The term $\frac{\epsilon}{2}\|\nabla u\|^2$ and the fact that constant functions are not in the kernel of K guarantees weak convergence in $H^1(\Omega)$ of a (subsequence of a) minimising sequence $\{u^k\}$ regardless of the choice of α in the inner problem. This implies the convergence in L^1 , and consequently by the standard argument of calculus of variations, an existence of a solution $R(\alpha)$ to the inner problem.

For a minimising sequence $\{\alpha^k\}$ of the whole problem (3.2) we therefore get weak convergence in L^2 of $u^k := R(\alpha^k)$ to some \hat{u} . But since $\|u^k - g\|_2^2 \rightarrow \|u - g\|_2^2$ for such a minimising sequence, we have so-called strict convergence of $u^k - g$ to $u - g$. In L^2 this implies strong convergence of first $u^k - g$ to $u - g$, and then of u^k to \hat{u} . Suppose also $\alpha^k \rightarrow \hat{\alpha}$. We then compute

$$\begin{aligned} J(u; \hat{\alpha}) &= \liminf_{k \rightarrow \infty} (J(u; \alpha^k) + G(u; \hat{\alpha}) - G(u; \alpha^k)) \\ &\geq \liminf_{k \rightarrow \infty} J(u; \alpha^k) + \liminf_{k \rightarrow \infty} (G(u; \hat{\alpha}) - G(u; \alpha^k)). \end{aligned}$$

Using the continuity of G with respect to α , we therefore obtain

$$J(u; \hat{\alpha}) = \liminf_{k \rightarrow \infty} J(u; \alpha^k) \geq \liminf_{k \rightarrow \infty} J(u^k; \alpha^k) \geq J(\hat{u}; \hat{\alpha}).$$

This shows that $\hat{u} = R(\hat{\alpha})$, and hence that $\hat{\alpha}$ solves (3.2). \square

Remark 3.1. Note that we only used ϵ to show existence of solutions to the inner problem. In particular, we do not require $\epsilon > 0$ to be constant, but we can allow $\epsilon = \epsilon(\alpha)$ to vary continuously with respect to α .

3.2. Derivative of the solution map

In order to use standard optimisation methods, such as BFGS, to minimise (3.1) we need to calculate a gradient of the solution map R ; equivalently, following the PDE approach of [34], we need to solve a so-called adjoint equation. A solution is given by the classical implicit function theorem [37, Theorem 4.E], and in particular the version in [38, Corollary 4.34] with relaxed assumptions. Let us suppose G is differentiable, such that $R(\alpha)$ is given as the solution $u = u_\alpha$ to

$$0 = S(u, \alpha) := K^*(Ku - f) + \nabla_u G(u; \alpha).$$

Then, if S is strictly differentiable, and $\nabla_u S(u, \alpha)$ is invertible, we have

$$\nabla R(\alpha) = [\nabla_u S(u, \alpha)]^{-1} \nabla_\alpha S(u, \alpha). \quad (3.4)$$

The strict differentiability of S may be achieved by a ‘‘second-degree’’ Huber regularisation [39]. The invertibility of $\nabla_u S(u, \alpha)$ can be guaranteed by an additional elliptic regularisation term $\frac{\epsilon}{2}\|\nabla u\|_2^2$. Both extra regularisations are the same as already employed in [34], where an alternative adjoint equation route was taken to avoid direct construction of ∇R . We now take a closer look at the derivatives of the solution maps for the regularisers (IC-TVTV) and (IC- L^2 TV). Furthermore, as we are going to restrict ourselves to the regularisation of dynamic video sequences, we consider K to be the identity operator mapping from $L^2(\Omega)$ to $L^2(\Omega)$ throughout the remainder of this paper.

3.3. Derivative of the IC TVTV solution map

To use the formula (3.4), we need the derivative of the map

$$R(\alpha_1, \alpha_2, \kappa) = P_1 \hat{R}(\alpha_1, \alpha_2, \kappa)$$

for $P_1(u, w) := u$, and

$$\hat{R}(\alpha_1, \alpha_2, \kappa) := \arg \min_{(u, w)} \frac{1}{2} \|u - g\|_2^2 + \alpha_1 \|\nabla_\kappa(u - w)\|_{2,1} + \alpha_2 \|\nabla_{1-\kappa}(w)\|_{2,1}.$$

Clearly

$$\nabla R(\alpha_1, \alpha_2, \kappa) = P_1 \nabla \hat{R}(\alpha_1, \alpha_2, \kappa) \quad (3.5)$$

if the latter derivative exists. For $\hat{R}(\alpha_1, \alpha_2, \kappa)$ to have a unique minimiser and to be differentiable, we further replace $\hat{R}(\alpha_1, \alpha_2, \kappa)$ with

$$\hat{R}_{\gamma, \epsilon}(\alpha_1, \alpha_2, \kappa) = \arg \min_{(u, w)} J_{\gamma, \epsilon}(u, w, \alpha_1, \alpha_2, \kappa) \quad (3.6)$$

for

$$J_{\gamma, \epsilon}(u, w, \alpha_1, \alpha_2, \kappa) := \frac{1}{2} \|u - g\|_2^2 + \alpha_1 \Theta(\nabla_\kappa(u - w)) + \alpha_2 \Theta(\nabla_{1-\kappa}(w)) + \frac{1}{2} \left(\int_\Omega w \, dx \right)^2.$$

Here we use the notation

$$\Theta(v) := \|v\|_\gamma + \frac{\epsilon}{2} \|v\|_2^2 \quad (3.7)$$

for the sum of the Huber-regularised 1-norm $\|v\|_\gamma = \int_\Omega H_\gamma(\|v(x)\|_2) \, dx$ with

$$H_\gamma(r) = \begin{cases} |r| - \frac{\gamma}{2} & |r| \geq \gamma \\ \frac{1}{2\gamma} |r|^2 & |r| < \gamma \end{cases}$$

and parameter $\gamma > 0$, and additional Hilbert space regularisation with parameter $\epsilon > 0$. The last term of (3.6) eliminates the translational invariance in \mathbb{R} with respect to w , thus forcing unique solutions as needed for the application of (3.4).

Note that existence of a solution (u, w) to (3.6) is guaranteed by simple application of Proposition 3.1 and Remark 3.1 provided $\kappa \notin \{0, 1\}$. In that case $\Theta(\nabla_\kappa u) \geq \epsilon' \|\nabla u\|_2^2$ for some $\epsilon' = \epsilon'(\kappa, \epsilon)$.

For the following propositions we define the shorthand notations

$$\Psi_\kappa^1(u) := [\nabla \Theta](\nabla_\kappa(u - w)), \quad \Psi_\kappa^2(u) := [\nabla^2 \Theta](\nabla_\kappa(u - w))$$

and

$$Q := \begin{pmatrix} I & -I \\ 0 & I \end{pmatrix},$$

where we denote by I the identity operator. We also model by $c := \chi_\Omega$ the term $\int_\Omega w \, dx = \langle c, w \rangle$ in (3.6).

Now we can derive the derivatives of the solution map.

Proposition 3.2 (Derivative of the IC TVTV solution map). *Let $\hat{R}_{\gamma,\epsilon}$ be defined by (3.6) for some $\gamma, \epsilon > 0$, and*

$$R(\alpha_1, \alpha_2, \kappa) := P_1 \hat{R}_{\gamma,\epsilon}(\alpha_1, \alpha_2, \kappa).$$

Then

$$\nabla R(\alpha_1, \alpha_2, \kappa) = -P_1 \left(\frac{\partial S}{\partial(u, w)} \right)^{-1} \frac{\partial S}{\partial(\alpha_1, \alpha_2, \kappa)}, \quad (3.8a)$$

where $S = \nabla_{(u, w)} J_{\gamma, \epsilon}$ satisfies

$$\frac{\partial S}{\partial(u, w)} := Q^* \begin{pmatrix} I - \alpha_1 \operatorname{div}_\kappa \Psi_\kappa^2(u - w) \nabla_\kappa & I \\ I & I - \alpha_2 \operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^2(w) \nabla_{1-\kappa} + c \otimes c \end{pmatrix} Q, \quad (3.8b)$$

$$\frac{\partial S}{\partial(\alpha_1, \alpha_2, \kappa)} := Q^* \begin{pmatrix} -\operatorname{div}_\kappa \Psi_\kappa^1(u - w) & 0 & T_1 \\ 0 & -\operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^1(w) & T_2 \end{pmatrix}, \quad (3.8c)$$

where

$$T_1 := -\alpha_1 \left(\operatorname{div} - \frac{\partial}{\partial t^*} \right) \Psi_\kappa^1(u - w) - \alpha_1 \operatorname{div}_\kappa \Psi_\kappa^2(u - w) \left(\nabla - \frac{\partial}{\partial t} \right) (u - w), \quad \text{and} \quad (3.8d)$$

$$T_2 := -\alpha_2 \left(\frac{\partial}{\partial t^*} - \operatorname{div} \right) \Psi_{1-\kappa}^1(w) - \alpha_2 \operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^2(w) \left(\frac{\partial}{\partial t} - \nabla \right) w. \quad (3.8e)$$

The problem of finding the derivative of the solution map thus reduces to the problem of solving three linear systems with the operator (3.8b) for the right-hand side vectors the columns of (3.8c).

Proof. The optimality conditions for the solution $u = R_{\gamma,\epsilon}(\alpha_1, \alpha_2, \kappa)$ are obtained from its Euler-Lagrange equation $S(u, w, \alpha_1, \alpha_2, \kappa) = 0$, which we split into

$$S_1 := u - f - \alpha_1 \operatorname{div}_\kappa \Psi_\kappa^1(u - w) = 0, \quad \text{and} \quad (3.9a)$$

$$S_2 := \alpha_1 \operatorname{div}_\kappa \Psi_\kappa^1(u - w) - \alpha_2 \operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^1(w) + \int_\Omega w \, dx = 0 \quad (3.9b)$$

Now $(S_1, S_2) = 0$ defines (u, w) as an implicit function of $(\alpha_1, \alpha_2, \kappa)$. By (3.4), the Jacobian of the function $\hat{R}_{\gamma,\epsilon} : (\alpha_1, \alpha_2, \kappa) \mapsto (u, w)$ is given by

$$\frac{\partial(u, w)}{\partial(\alpha_1, \alpha_2, \kappa)} = - \left(\frac{\partial S}{\partial(u, w)} \right)^{-1} \frac{\partial S}{\partial(\alpha_1, \alpha_2, \kappa)} \quad (3.10)$$

Calculating all the partial derivatives and using (3.5), we obtain (3.8). \square

3.4. Derivative of the IC L^2 TV solution map

Again we have a problem of the form

$$R(\alpha_1, \alpha_2, \kappa) = P_1 \hat{R}(\alpha_1, \alpha_2, \kappa)$$

for $P_1(u, w) := u$, and

$$\hat{R}(\alpha_1, \alpha_2, \kappa) = \arg \min_{(u, w)} \frac{1}{2} \|u - g\|_2^2 + \frac{\alpha_1}{2} \|\nabla_\kappa(u - w)\|_{2,1}^2 + \alpha_2 \|\nabla_{1-\kappa}(w)\|_{2,1}. \quad (3.11)$$

To employ the formula (3.4), we need to regularise $\hat{R}(\alpha_1, \alpha_2, \kappa)$ by replacing it with

$$\hat{R}_{\gamma,\epsilon}(\alpha_1, \alpha_2, \kappa) = \arg \min_{u, w} J_{\gamma,\epsilon}(u, w, \alpha_1, \alpha_2, \kappa) \quad (3.12)$$

for

$$J_{\gamma,\epsilon}(u, w, \alpha_1, \alpha_2, \kappa) := \frac{1}{2}\|u - g\|_2^2 + \frac{\alpha_1}{2}\|\nabla_\kappa(u - w)\|_2^2 + \alpha_2\Theta(\nabla_{1-\kappa}(w)) + \frac{1}{2}\left(\int_\Omega w \, dx\right)^2$$

Similar to Section 3.3, we obtain the following result.

Proposition 3.3 (Derivative of the IC L^2 TV solution map). *Let $\hat{R}_{\gamma,\epsilon}$ be defined by (3.12) for some $\gamma, \epsilon > 0$, and*

$$R(\alpha_1, \alpha_2, \kappa) := P_1 \hat{R}_{\gamma,\epsilon}(\alpha_1, \alpha_2, \kappa).$$

Then

$$\nabla R(\alpha_1, \alpha_2, \kappa) = -P_1 \left(\frac{\partial S}{\partial(u, w)} \right)^{-1} \frac{\partial S}{\partial(\alpha_1, \alpha_2, \kappa)}, \quad (3.13a)$$

where $S = \nabla_{(u, w)} J_{\gamma,\epsilon}$ satisfies

$$\frac{\partial S}{\partial(u, w)} := Q^* \begin{pmatrix} I - \alpha_1 \operatorname{div}_\kappa \nabla_\kappa & I \\ I & I - \alpha_2 \operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^2(w) \nabla_{1-\kappa} + c \times c \end{pmatrix} Q \quad (3.13b)$$

$$\frac{\partial S}{\partial(\alpha_1, \alpha_2, \kappa)} := Q^* \begin{pmatrix} -\operatorname{div}_\kappa \nabla_\kappa(u - w) & 0 & T_1 \\ 0 & -\operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^1(w) & T_2 \end{pmatrix} \quad (3.13c)$$

where

$$T_2 := -\alpha_2 \left(\frac{\partial}{\partial t^*} - \operatorname{div} \right) \Psi_{1-\kappa}^1(w) - \alpha_2 \operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^2(w) \left(\frac{\partial}{\partial t} - \nabla \right) w, \quad (3.13d)$$

$$T_2 := -\alpha_1 \left(\operatorname{div} - \frac{\partial}{\partial t^*} \right) \nabla_\kappa(u - w) - \alpha_1 \operatorname{div}_\kappa \left(\nabla - \frac{\partial}{\partial t} \right) (u - w). \quad (3.13e)$$

Here Q and c are as in Section 3.3.

Proof. As in Section 3.3, we find the optimal conditions for the minimization problem in (3.12) are

$$S_1 := u - f - \alpha_1 \operatorname{div}_\kappa \nabla_\kappa(u - w) = 0 \quad (3.14a)$$

$$S_2 := \alpha_1 \operatorname{div}_\kappa \nabla_\kappa(u - w) - \alpha_2 \operatorname{div}_{1-\kappa} \Psi_{1-\kappa}^1(w) + \int_\Omega w \, dV = 0 \quad (3.14b)$$

The Jacobian of the function $R_{\gamma,\epsilon} : (\alpha_1, \alpha_2, \kappa) \rightarrow u$ is again given by

$$\frac{\partial u}{\partial(\alpha_1, \alpha_2, \kappa)} = -P_1 \left(\frac{\partial S}{\partial(u, w)} \right)^{-1} \frac{\partial S}{\partial(\alpha_1, \alpha_2, \kappa)}. \quad (3.15)$$

We thus calculate the partial derivatives with respect to all the variables to obtain (3.13). \square

4. Computational realisation

The upper level problem of the bilevel optimisation (3.1) is solved numerically via the BFGS algorithm with backtracking line search and curvature verification, where

we update the quasi-Hessian only if it remains positive semi-definite. The following Armijo condition

$$F(\boldsymbol{\alpha} + \sigma \mathbf{d}) - F(\boldsymbol{\alpha}) \leq \sigma c \nabla F(\boldsymbol{\alpha}) \cdot \mathbf{d} \quad (4.1)$$

is used, where $F(\boldsymbol{\alpha}) = \frac{1}{2} \|R(\boldsymbol{\alpha}) - g\|_2^2$, \mathbf{d} is the search direction, σ the step-length and c a positive constant. The relative residual is used as a stopping criterion, so that the algorithm terminates if

$$\|\boldsymbol{\alpha}_i - \boldsymbol{\alpha}_{i-1}\|_2 < \rho \|\boldsymbol{\alpha}_i\|_2 \quad (4.2)$$

is satisfied for a fixed, positive parameter ρ . As all models have to be differentiable in order to solve the upper level problem, we used the L^2 norm $\frac{\epsilon}{2} \|\cdot\|_2^2$ with $\epsilon = 10^{-8}$ and the Huberised L^1 norm $\|\cdot\|_\gamma$ with $\gamma = 0.01$ as our cost functionals in order to optimise the parameters $(\alpha_1, \alpha_2, \kappa)$ for the regularisers in Table 1. We have further used MATLAB's inbuilt `gmres` function with a diagonally compensated incomplete Cholesky preconditioner in order to solve (3.10) and (3.15), respectively.

To compute numerical solutions u (and w) of the lower-level problem for fixed $\boldsymbol{\alpha}$, the lower level problem of the bilevel optimisation (3.1) is solved via the primal-dual hybrid gradient method (PDHGM) as presented in [40]. In order to apply the PDHGM to the lower level problem, we need to recast it into a saddle-point formulation. In case of (IC- L^2 TV) for instance, this saddle-point formulation reads as

$$\min_{(u,w)} \max_{(p,q)} \frac{1}{2} \|u - g\|_2^2 + \langle A \begin{pmatrix} u \\ w \end{pmatrix}, \begin{pmatrix} p \\ q \end{pmatrix} \rangle - \frac{1}{2\alpha_1} \|p\|_2^2 - \delta_{\alpha_2 P}(q), \quad (4.3)$$

where $P = \{p \mid \|p\|_{2,\infty} \leq 1\}$ is the unit ball with respect to the supremum norm and A is a linear operator given by

$$A = \begin{pmatrix} \nabla_\kappa & -\nabla_\kappa \\ 0 & \nabla_{1-\kappa} \end{pmatrix}. \quad (4.4)$$

To determine the step sizes in the PDHGM, we need to find a bound on $\|A\|$. Assuming $\kappa \in [0, 1]$, it is easy to show that $\|A\|^2 < 24$. The other formulations can be cast to saddle-point formulations in the same fashion.

In order to discretise ∇_κ and $\nabla_{1-\kappa}$, we simply use forward finite differences.

Note that the parameters ϵ, γ in (3.7) are numerical regularisation parameters only. As we run a relatively small number of PDHGM iterations to solve each inner problem, the solutions will be numerically inaccurate. This therefore allows us to ignore γ and ϵ in the inner denoising problem, and to solve the original non-smooth problem instead. In the outer problem, we further do not restrict $\kappa \in [0, 1]$, as this is not strictly necessary. For reporting the results in a uniform fashion, we use the identity

$$\|\nabla_\kappa v\|_1 = |2\kappa - 1| \|\nabla_{\frac{\kappa}{2\kappa-1}} v\|_1, \quad (4.5)$$

which holds for the unsmoothed inner problems. Therefore every triple $(\alpha_1, \alpha_2, \kappa)$ with $\kappa \notin [0, 1]$ corresponds to a triple with $\kappa \in (0, 1)$.

5. Numerical results

For the implementation of the BFGS scheme we use the following numerical setup. The constant c in (4.1) is set to $c = 10^{-4}$ throughout all experiments, whereas we use $\rho = 10^{-8}$ in (4.2). In all cases the parameters α_1 and α_2 were constrained to

lie in $(10^{-5}, 100)$, and κ was constrained to lie in $(-50, 50)$. We ran BFGS with 100 different initialisations of the parameters α drawn from a 3rd order χ^2 -squared distribution with mean rescaled to the values $(0.15, 0.15)$ for TVTV and $(3.9, 0.15)$ for L^2 TV, obtained by previous experimentation. The χ^2 -squared distribution is used in order to force the sampled parameters to be positive, but to not impose an upper bound. We report the afterwards optimised parameter α for which we obtained the best PSNR values for visual and quantitative comparison in the following figures and tables.

In order to solve the lower level problem, the PDHGM is run with a fixed number of iterations - 50 iterations in case the accelerated variant is applicable, otherwise 200 iterations, respectively. These were applied to the non-relaxed variants of the regularisers, as this allows us to leave κ unconstrained (see (4.5)). Each iteration uses warm initialisation with the optimal solution from the previous iteration.

We use three different $2 + 1D$ video sequences for our denoising experiments: the sequences 'hand', 'flight' and 'harlem shake'. The sequence 'hand', showing a hand falling onto a table, is of grid-size 54×96 and consists of 65 frames. It contains steady objects (like the table) and a moving object (the hand) that, at first not seen, moves onto the table where it becomes a steady object for the rest of the scene. The camera is fixed to a specific position. The sequence 'flight' is filmed from a flying glider. Here the background scenery passes by, while also the camera is at movement. The movie has grid-size 96×54 and consists of 90 frames. The third sequence 'harlem shake' is taken from https://www.youtube.com/watch?v=-_ZG2xgNAr4. The original RGB video has grid-size 540×720 and consists of 711 frames. It has been converted to gray-scale double precision, and down-sampled to grid-size 70×93 and 73 frames in order to make processing in a reasonable amount of time possible. The video shows a room inside a lodge in which the majority of people are in a rather steady position, whereas one person is dancing (and therefore moving). Approximately half-way through the video the scene changes, and everyone is at movement. As in case of the video 'hand', the camera is in a fixed position.

For all videos the underlying mesh-sizes are considered to be $h = 1$ in each dimension. Further, noisy versions of the video sequences have been created by perturbing the original sequences with Gaußian noise with mean zero and variance $\sigma^2 = 0.02$ respectively.

5.1. Discussion of results

We want to start discussing the results starting with the video 'hands'. Figure 1a shows six frames that are selected from the original sequence at the times displayed, and its noisy counterparts for variance $\sigma^2 = 0.02$. We clearly see the features described in the previous section, starting with a relatively steady scene and a moving hand emerging half-way through. Figure 1a shows the spatial and temporal components of the TV-TV reconstruction ((3.2) with (IC-TVTV) as regulariser) with optimal parameters that are given in Table 2, as well as the sum of both components. Note that we define the temporal component as u if $\kappa \leq 0.5$, and $u - w$ otherwise. We observe that in particular for the steady parts of the scene a lot of the noise has been removed. Most of the remaining artefacts seem to be present in the moving object, the hand. The temporal component seems to mostly contain the moving parts of the hand, whereas the spatial component contains a rather piecewise constant transition from desk without to desk with hand. Figure 1b on the other hand shows the result

Table 2: “Hand” test video optimal results. The star * in the optimal parameter means that the original κ for the optimal result was outside the range $[0, 1]$, and the conversion (4.5) has been used to derive the presented values.

| Model | $(\alpha_1, \alpha_2, \kappa)$ | Opt. value | PSNR | SSIM |
|-------|--------------------------------|------------|-------|--------|
| TVTV | $(0.162, 0.0844, 0.0466)^*$ | 104.5 | 32.08 | 0.9271 |
| L2TV | $(9.85, 0.0674, 0.0529)^*$ | 127.5 | 31.22 | 0.9125 |

Table 3: “Flight” test video optimal results. The star * in the optimal parameter means that the original κ for the optimal result was outside the range $[0, 1]$, and the conversion (4.5) has been used to derive the presented values.

| Model | $(\alpha_1, \alpha_2, \kappa)$ | Opt. value | PSNR | SSIM |
|-------|--------------------------------|------------|-------|--------|
| TVTV | $(0.0141, 0.017, 0.337)^*$ | 46.41 | 36.91 | 0.9569 |
| L2TV | $(4.81, 0.0163, 0.542)$ | 46.87 | 32.84 | 0.942 |

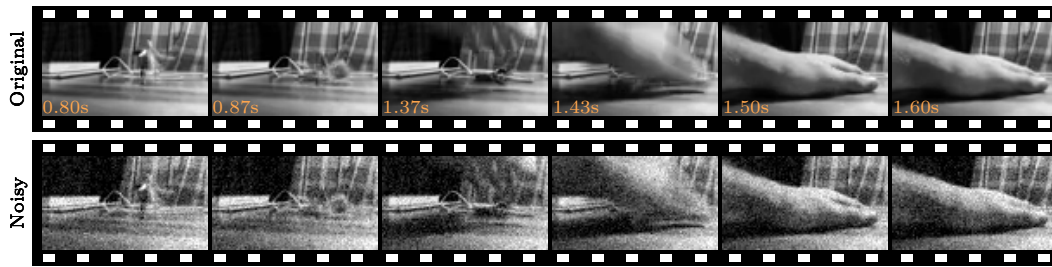
of the numerical reconstruction of (3.2) with (IC- L^2 TV) as a regulariser, for optimal parameters that are also given in Table 2. The results are similar to the TV-TV case; however, the spatial component shows a much smoother transition and therefore picks up the hand much earlier than TV-TV does. This also explains the quality-measure results in Table 2, as those tell us that TV-TV outperforms L_2 -TV for this sequence both in terms of PSNR and SSIM.

For the next sequence, ‘flight’, the situation looks quite similar in terms of PSNR and SSIM, however, the components are very different in this case. Figure 2a shows six consecutive frames of the original sequence, and the same frames from the sequence contaminated with noise. Figure 2b shows the reconstructions with (IC-TVTV) as a regulariser. Clearly, the temporal component in this case fails to pick up any information about the sequence, whereas all the information is stored in the spatial component. In case of (IC- L^2 TV), the optimal value for κ is approximately 0.5 according to Table 3. Hence, we can hardly speak of a temporal and a spatial component in this case, but rather of an L^2 - and a TV-penalised component that are shown in Figure 2c. The first component is rather blurry, and more or less approximating the homogeneous sky and a blurred version of the cockpit. The second component on the other hand approximates the sharp features and structures, but not the homogeneous background parts.

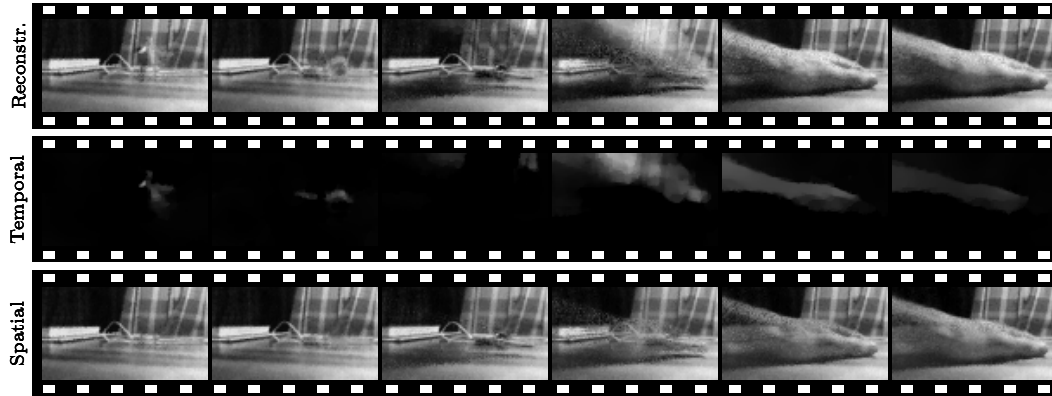
For the last sequence shown in Figure 3, Harlem shake, we figure out that both (IC-TVTV) and (IC- L^2 TV) seem to perform almost equally well, indicating a small value of κ which is confirmed by Table 4. Given the nature of the scene, the temporal part captures most of the scene, as there are only very few pixels that remain (almost) unchanged over time. Some of those are captured in the spatial component, like parts of the table in the background on the left hand side for instance.

6. Conclusions & Outlook

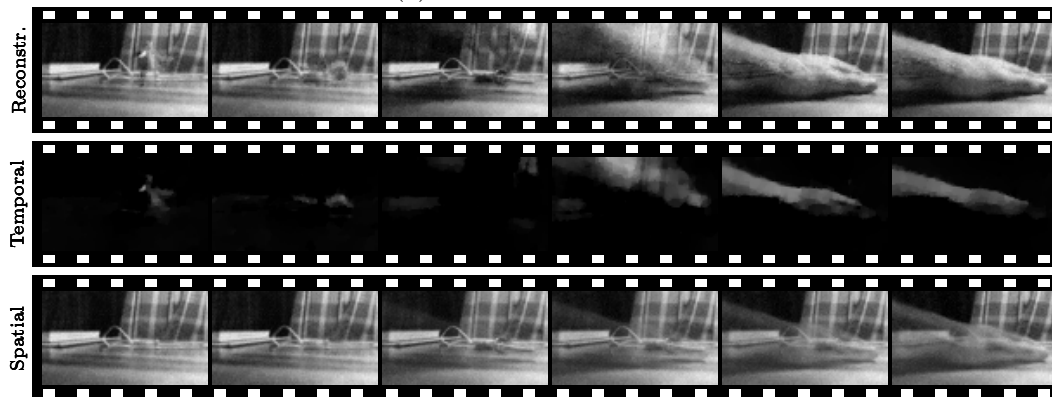
We have presented a bilevel optimisation strategy for optimising the regularisation parameters of infimal convolution-type regularisation methods for dynamic image regularisation. We have shown existence of solutions under additional regularity assumptions, and demonstrated the numerical performance of the proposed method



(a) Original and noisy videos

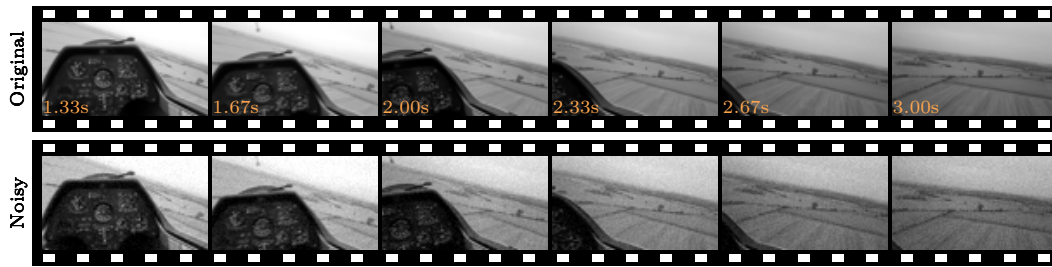


(b) TVTV reconstruction

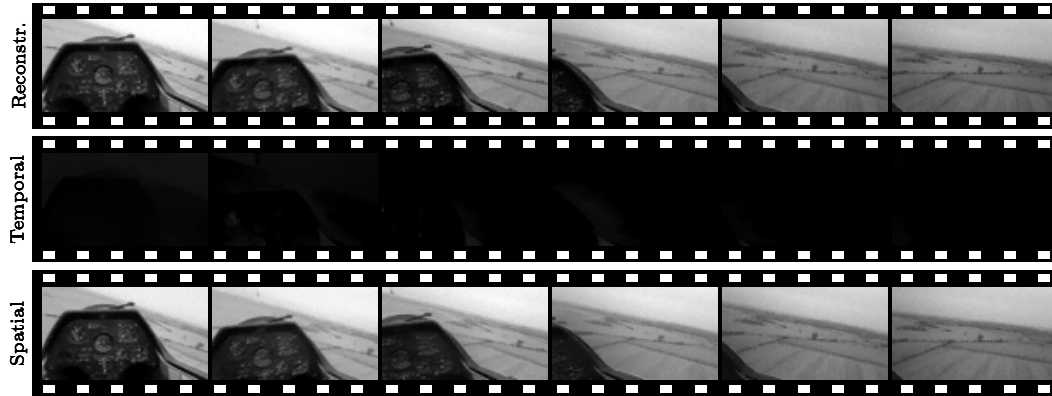


(c) L^2TV reconstruction

Figure 1: “Hand” test video reconstructions with optimally learned parameters using (3.2) for the ICTVTV and ICL^2TV regularisation models. The reconstructions are depicted together with their temporal and spatial components. Note how L^2TV picks up the hand in the spatial component much earlier than TVTV.



(a) Original and noisy videos



(b) TVTV reconstruction

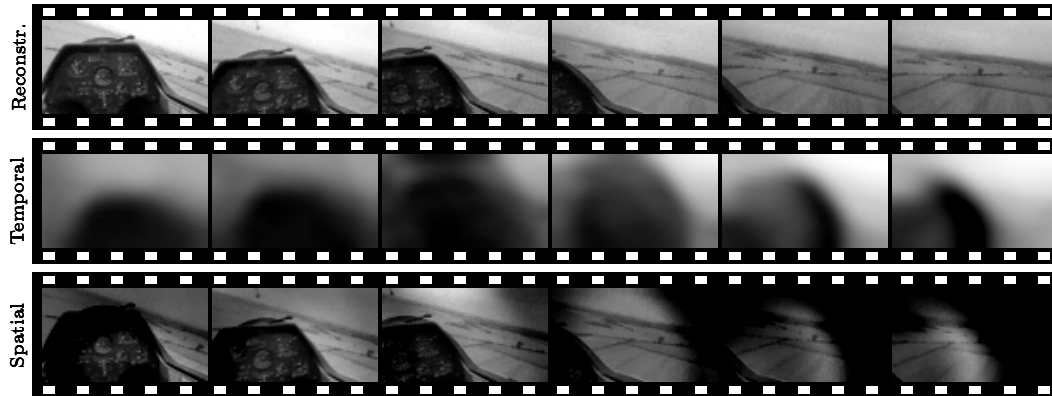
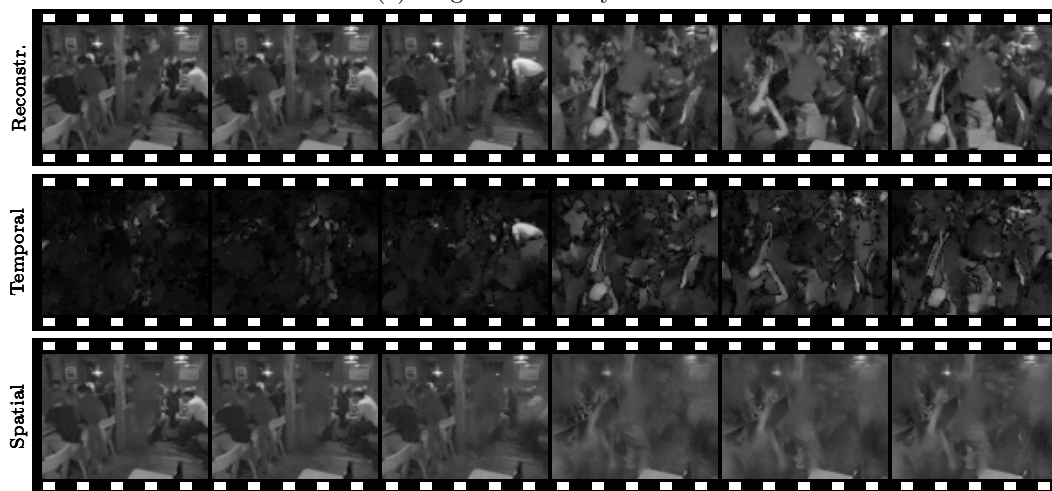
(c) L^2TV reconstruction

Figure 2: “Flight” test video reconstructions with optimally learned parameters using (3.2) for the ICTVTV and ICL^2TV regularisation models. The reconstructions are depicted together with their temporal and spatial components. Note how TVTV manages to extract no temporal component, while L^2TV manages to extract the spatially stable rough features in the temporal component. Recall that we define the temporal component as u if $\kappa \leq 0.5$, and $u - w$ otherwise. Since κ is approximately 0.5, see Table 3, it can be argued that the nomenclature “spatial” and “temporal” components should be swapped for L^2TV , and the spatial component should be the blurry one, corresponding to the very approximate spatially constant information.



(a) Original and noisy videos



(b) TVTV reconstruction

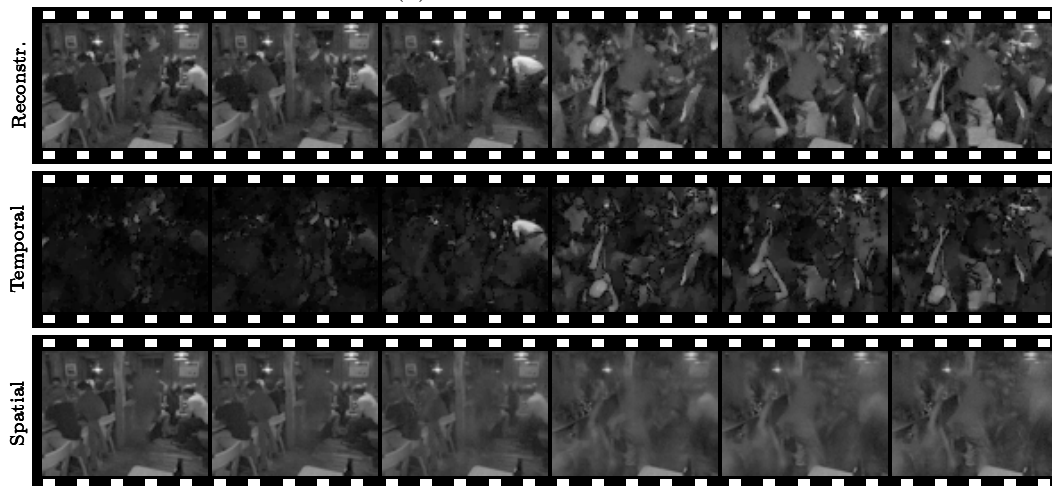
(c) L^2 TV reconstruction

Figure 3: “Harlem shake” video reconstructions with optimally learned parameters using (3.2) for the ICTVTV and ICL^2 TV regularisation models. The reconstructions are depicted together with their temporal and spatial components. The displayed images have been gamma-corrected with factor $\gamma = 0.6$ to improve legibility on paper.

Table 4: “Harlem shake” test video optimal results. The star * in the optimal parameter means that the original κ for the optimal result was outside the range $[0, 1]$, and the conversion (4.5) has been used to derive the presented values.

| Model | $(\alpha_1, \alpha_2, \kappa)$ | Opt. value | PSNR | SSIM |
|-------|--------------------------------|------------|-------|--------|
| T VTV | (0.0649, 0.0122, 0.0319) | 39.52 | 37.67 | 0.9621 |
| L2TV | (23, 0.0114, 0.0265) | 41.02 | 37.53 | 0.9604 |

for three distinctive video sequences. The sequences considered allow the assumption, that the use of the infimal convolution models for video denoising highly depends on the corresponding video sequence. As for sequences with stationary backgrounds, the decomposition into spatial and temporal components seems to work well, with the optimal κ being close to 0 or 1 allowing for a clear distinction of a temporal and a spatial component. For sequences like the ‘flight’-sequence that lack stationary parts, the distinction is not clear at all, which is also underpinned by the optimal κ value being close to 0.5. In this setup the modelling assumption of multiple infimal convolutions of the same functional also breaks down, as they will be the same. A choice of κ close to 0.5 only makes sense for infimal convolutions of functionals that promote very different information (like the L^2 TV model in our case).

Future research should address different error measures for the comparison of the denoised video sequences to the ground truth, in order to see, if different error measures will lead to similar or different conclusions. Further can the proposed research be easily extended to general, bounded linear operators K , which has been omitted here for the sake of brevity. Research on infimal convolutions of different, complementary regularisation functionals seems to be another promising direction that future research can head for.

Acknowledgements

M. Benning, C.-B. Schönlieb and T. Valkonen acknowledge support from the EPSRC grant Nr. EP/M00483X/1 and from the Leverhulme grant ‘Breaking the non-convexity barrier’. V. Vlačić was supported by a Bridgewater summer internship.

A data statement for the EPSRC

The code and data will be put into a repository as the final version is submitted.

References

- [1] Holler M and Kunisch K 2014 *SIAM Journal on Imaging Sciences* **7** 2258–2300
- [2] Rudin L I, Osher S and Fatemi E 1992 *Physica D: Nonlinear Phenomena* **60** 259–268
- [3] Donoho D L 2006 *Information Theory, IEEE Transactions on* **52** 1289–1306
- [4] Burger M and Osher S 2013 A guide to the tv zoo *Level Set and PDE Based Reconstruction Methods in Imaging* (Springer) pp 1–70
- [5] Chambolle A and Lions P L 1997 *Numerische Mathematik* **76** 167–188
- [6] Benning M, Brune C, Burger M and Müller J 2013 *Journal of Scientific Computing* **54** 269–310
- [7] Benning M 2011 *Singular regularization of inverse problems* Ph.D. thesis
- [8] Müller J 2013 *Advanced image reconstruction and denoising: Bregmanized (higher order) total variation and application in PET* Ph.D. thesis

- [9] Roth and Black M J 2005 Fields of experts: A framework for learning image priors *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on* vol 2 (IEEE) pp 860–867
- [10] Tappen M F 2007 Utilizing variational optimization to learn Markov random fields *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on* (IEEE) pp 1–8
- [11] Domke J 2012 Generic methods for optimization-based modeling *International Conference on Artificial Intelligence and Statistics* pp 318–326
- [12] Chen Y, Ranftl R and Pock T 2014 *Image Processing, IEEE Transactions on* To appear
- [13] Haber E and Tenorio L 2003 *Inverse Problems* **19** 611
- [14] Haber E, Horesh L and Tenorio L 2010 *Inverse Problems* **26** 025002
- [15] Bui-Thanh T, Willcox K and Ghattas O 2008 *SIAM Journal on Scientific Computation* **30** 3270–3288
- [16] Biegler L, Biros G, Ghattas O, Heinkenschloss M, Keyes D, Mallick B, Tenorio L, van Bloemen Waanders B, Willcox K and Marzouk Y 2011 *Large-scale inverse problems and quantification of uncertainty* vol 712 (John Wiley & Sons)
- [17] De los Reyes J C and Schönlieb C B 2013 *Inverse Problems and Imaging* **7**
- [18] Calatroni L, De los Reyes J C and Schönlieb C B 2014 Dynamic sampling schemes for optimal noise learning under multiple nonsmooth constraints *System Modeling and Optimization* (Springer Verlag) pp 85–95
- [19] Reyes J C D L, Schönlieb C B and Valkonen T 2015 *Journal of Mathematical Analysis and Applications* Accepted (Preprint arXiv:1505.01953)
- [20] Reyes J C D L, Schönlieb C B and Valkonen T 2015 Bilevel parameter learning for higher-order total variation regularisation models submitted (Preprint arXiv:1508.07243)
- [21] Calatroni L, De los Reyes J C and Schönlieb C B A variational model for mixed noise distribution in preparation
- [22] Chung C V and De los Reyes J C Learning optimal spatially-dependent regularization parameters in total variation image restoration in preparation
- [23] Kunisch K and Pock T 2013 *SIAM Journal on Imaging Sciences* **6** 938–983
- [24] Chen Y, Pock T and Bischof H 2012 Learning ℓ_1 -based analysis and synthesis sparsity priors using bi-level optimization *Workshop on Analysis Operator Learning vs. Dictionary Learning, NIPS 2012*
- [25] Ochs P, Ranftl R, Brox T and Pock T 2015 Bilevel Optimization with Nonsmooth Lower Level Problems *International Conference on Scale Space and Variational Methods in Computer Vision (SSVM)* to appear
- [26] Chung J, Español M I and Nguyen T 2014 *arXiv preprint arXiv:1407.1911*
- [27] Hintermüller M and Wu T 2014 Bilevel optimization for calibrating point spread functions in blind deconvolution preprint
- [28] Baus F, Nikolova M and Steidl G 2014 *Journal of Mathematical Imaging and Vision* **48** 295–307
- [29] Schmidt U and Roth S 2014 Shrinkage fields for effective image restoration *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on* (IEEE) pp 2774–2781
- [30] Fehrenbach J, Nikolova M, Steidl G and Weiss P 2015 Bilevel image denoising using gaussianity tests *Scale Space and Variational Methods in Computer Vision (Lecture Notes in Computer Science* vol 9087) ed Aujol J F, Nikolova M and Papadakis N (Springer International Publishing) pp 117–128 URL http://dx.doi.org/10.1007/978-3-319-18461-6_10
- [31] Ranftl R and Pock T 2014 A deep variational model for image segmentation *36th German Conference on Pattern Recognition (GCPR)*
- [32] Klatzer T and Pock T 2015 Continuous Hyper-parameter Learning for Support Vector Machines *Computer Vision Winter Workshop (CVWW)*
- [33] Chen Y, Yu W and Pock T 2015 On learning optimized reaction diffusion processes for effective image restoration *IEEE Conference on Computer Vision and Pattern Recognition* to appear
- [34] Reyes J C D L and Schönlieb C B 2013 *Inverse Problems and Imaging* **7** 1183–1214 (Preprint arXiv:1207.3425)
- [35] Reyes J C D L, Schönlieb C B and Valkonen T 2015 *Journal of Mathematical Analysis and Applications* In press (Preprint arXiv:1505.01953) URL <http://iki.fi/tuomov/mathematics/interior.pdf>
- [36] Reyes J C D L, Schönlieb C B and Valkonen T 2015 Bilevel parameter learning for higher-order total variation regularisation models submitted (Preprint arXiv:1508.07243) URL http://iki.fi/tuomov/mathematics/tgv_learn.pdf
- [37] Zeidler E 2012 *Applied Functional Analysis: Applications to Mathematical Physics* Applied Mathematical Sciences (Springer New York) ISBN 9781461208150
- [38] Mordukhovich B S 2006 *Variational Analysis and Generalized Differentiation I: Basic Theory*

(*Grundlehren der mathematischen Wissenschaften* vol 330) (Springer-Verlag)

[39] De los Reyes J C 2011 *SIAM Journal on Control and Optimization* **49** 1629–1658

[40] Chambolle A and Pock T 2011 *Journal of Mathematical Imaging and Vision* **40** 120–145