



**HAL**  
open science

## Supervised learning of analysis-sparsity priors with automatic differentiation

Hashem Ghanem, Joseph Salmon, Nicolas Keriven, Samuel Vaiter

► **To cite this version:**

Hashem Ghanem, Joseph Salmon, Nicolas Keriven, Samuel Vaiter. Supervised learning of analysis-sparsity priors with automatic differentiation. IEEE Signal Processing Letters, 2023, 30, pp.339-343. 10.1109/LSP.2023.3244511 . hal-03518852

**HAL Id: hal-03518852**

**<https://hal.science/hal-03518852>**

Submitted on 10 Jan 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Supervised learning of analysis-sparsity priors with automatic differentiation

Hashem Ghanem, Joseph Salmon, Nicolas Keriven, and Samuel Vaiter

**Abstract**—Sparsity priors are commonly used in denoising and image reconstruction. For analysis-type priors, a dictionary defines a representation of signals that is likely to be sparse. In most situations, this dictionary is not known, and is to be recovered from pairs of ground-truth signals and measurements, by minimizing the reconstruction error. This defines a hierarchical optimization problem, which can be cast as a bi-level optimization. Yet, this problem is unsolvable, as reconstructions and their derivative *w.r.t.* the dictionary have no closed-form expression. However, reconstructions can be iteratively computed using the Forward-Backward splitting (FB) algorithm. In this paper, we approximate reconstructions by the output of the aforementioned FB algorithm. Then, we leverage automatic differentiation to evaluate the gradient of this output *w.r.t.* the dictionary, which we learn with projected gradient descent. Experiments show that our algorithm successfully learns the 1D Total Variation (TV) dictionary from piecewise constant signals. For the same case study, we propose to constrain our search to dictionaries of 0-centered columns, which removes undesired local minima and improves numerical stability.

**Index Terms**—Sparsity, dictionary learning, total variation, bi-level optimization, automatic differentiation.

## I. INTRODUCTION

**D**ENOISING is a widely-tackled problem that emerges in many fields, where the goal is to restore signals from noisy observations. This includes that the model of the imaging system is *a priori* known:  $\mathbf{y} = \mathbf{w} + \varepsilon$ , where  $\mathbf{w}, \mathbf{y} \in \mathbb{R}^p$  are the true and the measured signals, respectively, and  $\varepsilon \in \mathbb{R}^p$  is additive noise. In addition, a prior hypothesis on the nature of the signal, or an embedding of it, might be available, like sparsity [10]. Such extra knowledge can be incorporated in the optimization process to get better reconstructions (*e.g.*, a higher Signal to Noise Ratio (SNR)).

Sparsity priors exist in two forms [6]: *i*) synthesis (traditional) sparsity where  $\mathbf{w} = \mathbf{D}\mathbf{u}$ ,  $\mathbf{D}$  is a linear operator, and  $\mathbf{u}$  is sparse; *ii*) analysis sparsity where  $\mathbf{v} = \mathbf{D}^\top \mathbf{w} \in \mathbb{R}^m$  is sparse. This paper is interested in the latter. As a convex surrogate, this prior is enforced on  $\mathbf{w}$  by adding the term  $\|\mathbf{D}^\top \mathbf{w}\|_1$  to the (generally quadratic) loss function [9], where  $\|\cdot\|_1$  is the  $\ell_1$  norm. Putting all together, the problem consists in finding:  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{w}\|_2^2 + \lambda \|\mathbf{D}^\top \mathbf{w}\|_1$  for some regularization amplitude parameter  $\lambda \geq 0$ . The linear operator  $\mathbf{D}$  is either user-defined, or learnt directly from data.

The main problem we tackle in this work is to extract both  $\mathbf{D}$  and  $\lambda$  from data with supervised learning. Having that  $\lambda \|\mathbf{D}^\top \mathbf{w}\|_1 = \|(\lambda \mathbf{D})^\top \mathbf{w}\|_1$ , this problem is equivalent to extracting the product  $\lambda \mathbf{D}$  as one object, thus we stop notating

$\lambda$  explicitly from now on and keep  $\mathbf{D}$ . To avoid trivial solutions and undesired local minima, it is common to restrict the search space to an admissible set  $\mathcal{C} \subset \mathbb{R}^{p \times m}$ , where dictionaries have a specific property. In [12] for instance,  $\mathbf{D}$  is forced to be orthogonal, *i.e.*,  $\mathcal{C} = \{\mathbf{D}; \mathbf{D}\mathbf{D}^\top = \mathbf{I}_p\}$ . In [11],  $\mathbf{D}$  is constrained to be a convolution dictionary. However, it is still challenging to find an admissible set that performs well in all applications of this problem [17]. In this work, we don't commit to find such universal set. Formally, we suppose we have a dataset  $(\mathbf{y}_l, \mathbf{w}_l)_{l=1}^L$ , that includes  $L$  pairs of measurements and their associated ground truth signals. Knowing such pairs, our task is to find the operator  $\mathbf{D}$  that minimizes the mean squared error between reconstructions and ground truth signals:

$$\hat{\mathbf{D}} \in \arg \min_{\mathbf{D} \in \mathbb{R}^{p \times m}} \mathcal{E}(\mathbf{D}) = \sum_{l=1}^L \|\hat{\mathbf{w}}(\mathbf{D}, \mathbf{y}_l) - \mathbf{w}_l\|_2^2 + \iota_{\mathcal{C}}(\mathbf{D}) \quad (1a)$$

s.t.

$$\hat{\mathbf{w}}(\mathbf{D}, \mathbf{y}) = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{w}\|_2^2 + \|\mathbf{D}^\top \mathbf{w}\|_1, \quad (1b)$$

where  $\mathcal{C}$  is an admissible set of dictionaries, and  $\iota_{\mathcal{C}}$  is the indicator function:  $\iota_{\mathcal{C}}(\mathbf{D}) = 0$  if  $\mathbf{D} \in \mathcal{C}$  and  $+\infty$  otherwise.

One can recast Equation (1) as a *bilevel optimization problem*: in the outer problem, we learn the model  $\mathbf{D}$  as in Eq. (1a) while in the inner problem, we denoise measurements following Eq. (1b). We already know that the inner part can be solved applying the Forward-Backward splitting (FB) algorithm on the dual problem [5]. However, due to the  $\ell_1$  norm, neither the solution nor its gradient *w.r.t.*  $\mathbf{D}$  have closed-form expressions. Thus, the minimizer  $\hat{\mathbf{D}}$  cannot be derived analytically nor obtained with gradient-based methods.

To illustrate how our algorithm can recover a dictionary  $\hat{\mathbf{D}}$ , we consider the well-known problem of *1D piecewise constant signals reconstruction* as a case-study [5], where this prior indicates that  $(w_2 - w_1, \dots, w_1 - w_p)^\top$  is sparse. The estimator is often written as an instance of Eq. (1b), with  $\mathbf{D} = \mathbf{D}_{TV}$  is the dictionary associated to the 1D Total Variation (TV) regularization: for all  $i \in \{1, \dots, p\}$ ;  $\mathbf{D}_{i,i} = -1$ ,  $\mathbf{D}_{i+1,i} = 1$ , and 0 otherwise, up to rescaling.

**Contribution** We approximately recover the analysis-sparsity operator  $\hat{\mathbf{D}}$  by: *i*) replacing the true minimizer  $\hat{\mathbf{w}}(\mathbf{D}, \mathbf{y})$  by the output of the FB algorithm applied on the dual problem; *ii*) deploying *automatic differentiation*, a technique capable of evaluating the gradient of an algorithm *w.r.t.* input variables, to solve Eq. (1a) with projected gradient descent. We empirically show that our method recovers the TV dictionary  $\mathbf{D}_{TV}$  from piecewise constant signals. Finally, for the same case study, we reduce the admissible set  $\mathcal{C}$  to dictionaries with

The authors acknowledge the support of ANR Grava ANR-18-CE40-0005, ANR GRandMa ANR-21-CE23-0006 and ANER BFC RAGA.

HG is with CNRS, IMB, Univ. de Bourgogne; JS is with Univ. Montpellier; NK is with CNRS, GIPSA-lab; SV is with CNRS, LJAD, Univ. Côte d'Azur.

columns summing up to zero, and empirically prove that this increases stability and extracts the TV operator with higher quality than previous methods.

**Related work** The bilevel problem was first posed in [11], where the authors smoothed the  $\ell_1$  norm so that the derivative of  $\hat{w}(\mathbf{D}, \mathbf{y})$  w.r.t.  $\mathbf{D}$  has a closed-form expression. Then, they applied gradient descent to find a local minimizer  $\mathbf{D}$ . However, the sought for sparsity is degraded with this methodology. In [14], a relaxation regime of  $\ell_1$  with the smooth  $\ell_2$  norm is adopted, with a relaxation parameter easy to assign. Recently in [10], a formula of  $\hat{w}(\mathbf{D}, \mathbf{y})$  is obtained under some conditions. This formula includes inverting a large matrix, which is computed iteratively while Automatic Differentiation (AD) is used to evaluate gradients. In [16], and without a proof of convergence, they restrict  $\mathbf{D}$  to have unit columns norm, that also satisfy  $\mathbf{D}\mathbf{D}^\top = \mathbf{I}$ , in order to avoid trivial solutions (e.g., repeated columns). The work in [4] solves the 2D piecewise signals reconstruction problem using AD, constrained by learning *convolution-type* dictionaries of kernels with small support, which improves the quality of standard 2D TV. Such strong constraints are not considered in our work, however we will see that simple column-centering suffices to learn a high-quality dictionary in the 1D version of this problem.

## II. AUTOMATIC DIFFERENTIATION (AD)

To compute gradients, one often manually writes down its analytical expression. Yet, this can only be performed for functions with closed-form expression, which is not the case in Eq. (1b). An alternative is symbolic differentiation, automatically performed by computer tools like Mathematica [7] when dealing with syntax tree expressions. A third approach consists in approximating gradients using numerical differentiation. It is easy to implement though exposed to round-off errors [8] and expensive in case of a high number of variables (here  $p \times m$ ).

Automatic Differentiation (AD), that is of interest in this work, mitigates the previous drawbacks. AD manipulates computation flow in a computer program, with all numerical computations reduced to compositions of elementary operations, for which derivative rules are known [15]. In such computer programs, a value must be assigned to each input variable, thus, any operation in the code will result in a variable with a numerical value. During the execution of the program, AD consists in [1]: 1) keeping a trace to all intermediate variables that are dependent on the ones we want to differentiate for; 2) once an operation is performed on a dependent variable, say  $v_i$ , to evaluate another  $v_j$ , directly computing the derivative value  $dv_j/dv_i$ ; 3) accumulating derivatives in step 2 through the chain rule. This gives the derivative value of the whole composition w.r.t. a chosen variable. AD can trace numerical computations in recursion functions and in control-flow (*if, while, for*) statements. Thus, AD efficiently differentiates not only closed-form formulas, but also implemented algorithms. This makes AD suitable for gradient-based optimization. AD can be implemented in two ways: *forward* mode and *reverse*

mode. In our work, we adopt the reverse mode<sup>1</sup>. Let us assume having  $\mathbf{y} = (y_1, \dots, y_m)$  as a function of the variable  $\mathbf{x} = (x_1, \dots, x_n)$ , with  $v$  as an intermediate variable. The backward scheme takes in input the vector  $\bar{\mathbf{y}} = (\bar{y}_1, \dots, \bar{y}_m)$ , and accumulates gradients as follows: first  $\bar{\mathbf{v}} = \mathbf{J}_y^\top(v)\bar{\mathbf{y}}$ , then  $\bar{\mathbf{x}} = \mathbf{J}_v^\top(x)\bar{\mathbf{v}}$ , where  $(\mathbf{J}_y(v))_{i,j} = dy_i/dv_j$ . The output  $\bar{\mathbf{x}}$  is nothing but the transpose Jacobian-vector product  $\mathbf{J}_y^\top(x)\bar{\mathbf{y}}$ .

## III. PROPOSED ALGORITHM

We solve the *dual* problem of Eq. (1b) with Forward-Backward splitting (FB). Using Automatic Differentiation, we obtain gradients of the previous FB algorithm w.r.t.  $\mathbf{D}$ . These gradients are used to learn a local minimum  $\hat{\mathbf{D}}$  using gradient descent, while projecting  $\mathbf{D}$  on the admissible set  $\mathcal{C}$  at every iteration.

### A. Deriving the dual problem of Eq. (1b)

The term  $\|\mathbf{D}^\top \mathbf{w}\|_1$  in Eq. (1b) is neither differentiable w.r.t.  $\mathbf{w}$ , nor has a simple proximal operator that can be efficiently performed. Hence, we cannot apply the gradient descent or the FB algorithm directly to evaluate  $\hat{w}(\mathbf{D}, \mathbf{y})$ . However, the latter is possible if we tackle this optimization from the dual perspective [13]. In fact, one can prove that Eq. (1b) is equivalent to its dual problem [5]:

$$\hat{z}(\mathbf{D}, \mathbf{y}) = \arg \min_{\mathbf{z} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{D}\mathbf{z} - \mathbf{y}\|_2^2 + \iota_{B_\infty}(\mathbf{z}), \quad (2)$$

where  $B_\infty = \{\mathbf{z} \in \mathbb{R}^m; |z_i| \leq 1, \forall i \in [m]\}$  is the unit ball of the  $\ell_\infty$  norm. The target recovery  $\hat{\mathbf{w}}$  is then given by:

$$\hat{\mathbf{w}}(\mathbf{D}, \mathbf{y}) = \mathbf{y} - \mathbf{D} \hat{z}(\mathbf{D}, \mathbf{y}). \quad (3)$$

### B. Solving the dual problem with the FB algorithm

Fortunately, the term  $\iota_{B_\infty}(\mathbf{z})$  has a simple proximal function given by  $\Pi_{B_\infty}$ : the orthogonal projection on the ball  $B_\infty$ . So indeed, we solve Eq. (2) with an accelerated FB algorithm, namely the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) [2]. At each iteration we update  $\mathbf{z}$  as follows:

$$\begin{aligned} \mathbf{z}_{i+1} &= \text{prox}_{\iota_{B_\infty}} \left( \mathbf{q}_i - \eta_1 \nabla \left( \frac{1}{2} \|\mathbf{D}\mathbf{q}_i - \mathbf{y}\|_2^2 \right) \right) \\ &= \Pi_{B_\infty} \left( \mathbf{q}_i - \eta_1 \mathbf{D}^\top (\mathbf{D}\mathbf{q}_i - \mathbf{y}) \right) \end{aligned} \quad (4)$$

s.t.  $\mathbf{q}_1 = \mathbf{z}_0$  is the initialization of  $\mathbf{z}$ ,  $\eta_1$  is the step size, and:

$$\begin{aligned} \mathbf{q}_{i+1} &= \mathbf{z}_i + \frac{t_i - 1}{t_{i+1}} (\mathbf{z}_i - \mathbf{z}_{i-1}) \\ t_{i+1} &= \frac{1}{2} (1 + \sqrt{1 + 4t_i^2}); t_1 = 1. \end{aligned} \quad (5)$$

To conclude with the inner part: having the matrix  $\mathbf{D}$  and the vector  $\mathbf{y}$  as *input*, one can compute  $\hat{z}(\mathbf{D}, \mathbf{y})$  with a sufficient number of updates as in Eq. (4), then get  $\hat{\mathbf{w}}(\mathbf{D}, \mathbf{y})$  in *output* using Eq. (3).

<sup>1</sup>this is the mode implemented in the PyTorch framework

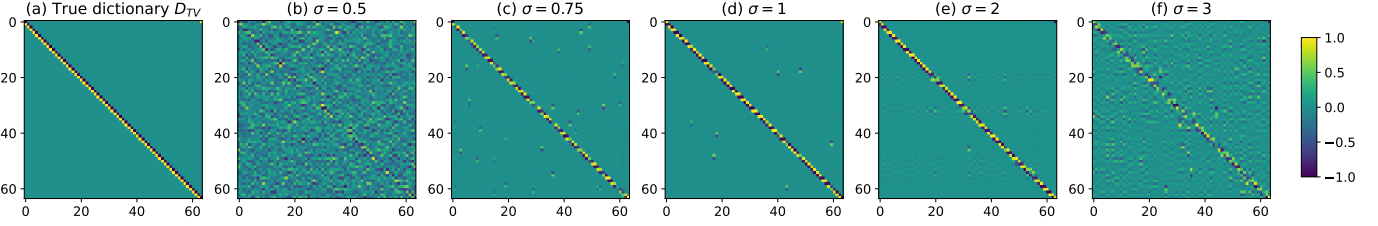


Fig. 1: Performance of our projected gradient descent algorithm 1, *w.r.t.* to noise level. We plot a sorted view of the dictionary  $\hat{D}$ . We rescale all dictionaries to  $[-1, 1]$  for a better visualization of the structure recovered in  $\hat{D}$ . (a)  $D_{TV}$  our algorithm is expected to learn. From (b) to (f):  $\hat{D}$  for different values of  $\sigma$ , the standard deviation of noise. Jumps in the dataset have the same amplitude:  $0 \rightarrow 10$  or  $10 \rightarrow 0$ .

### C. Outer loop design to learn the dictionary

Let us first notice that when the admissible set  $\mathcal{C} = \mathbb{R}^{p \times m}$ , the solution  $\hat{D}$  is *not unique* as  $\|\cdot\|_1$  is invariant to the coefficients order in a vector, *e.g.*,  $\|(u_1, u_2)^\top\|_1 = \|(u_2, u_1)^\top\|_1$ . Thus, permuting columns in  $D$  will lead to the same  $\hat{w}(D, y)$  in Eq. (1b), which means the same cost  $\mathcal{E}(D)$ . In this paper, we consider any solution that is a minimizer of  $\mathcal{E}(D)$ . We also assume that the second dimension of  $D$  is given, while optimizing for it might be the subject of a future work.

Towards our goal, we use AD to get the gradient of the mean squared error (MSE) term in  $\mathcal{E}(D)$ , by tracing the FB algorithm that solves Eq. (1b). In fact, to reduce computational cost, we do not compute the full MSE but sample a batch of training signals at each iteration (see Alg. 1). We denote this term by  $\mathcal{E}_{MSE}(D)$ . Since proving the convergence of AD's Jacobian to the variational one is complicated in such non-smooth setting, we replace the true reconstruction  $\hat{w}(D, y)$  with the output produced by the FB algorithm, and empirically show that it is a good proxy. For simplicity, we keep the same notation for this output. We randomly initialize  $D_0$ , then we start each iteration  $t$  by setting *PyTorch AD framework* to record operations on  $D_t$ . We can then compute the output  $\hat{w}(D_t, y)$ , and  $\mathcal{E}_{MSE}(D_t)$ . Now, we use the *PyTorch AD* to get the gradient  $\nabla \mathcal{E}_{MSE}(D_t)$ , and we update the estimated dictionary as

$$D_{t+1} = D_t - \eta_2 \nabla \mathcal{E}_{MSE}(D_t), \quad (6)$$

where  $\eta_2 > 0$  is a step size. Lastly, we project  $D_t$  on the admissible set  $\mathcal{C}$  by computing  $\Pi_{\mathcal{C}}(D_t)$ . We obtain Algorithm 1.

## IV. EXPERIMENTS

We consider the 1D piecewise constant signals reconstruction problem, with  $p = m = 64$  in all experiments. Ground-truth examples  $w$  are generated s.t. they have 4 discontinuities, thus 4 constant pieces. The coefficients where discontinuities take place are randomly chosen in each  $w$ . Their amplitude is either fixed ( $0 \rightarrow 10$  or  $10 \rightarrow 0$ ), or randomly sampled s.t. they happen between two values in  $[0, 10]$ , to be mentioned when necessary. Observations  $y$  are constructed by adding a noise vector to each ground-truth signal, such vector is sampled from  $\mathcal{N}(0, \sigma^2 I_p)$ , where  $\sigma$  varies through experiments. The “true” underlying dictionary  $D_{TV}$  is shown in Fig. 1 (a), up to permuting its columns, and to rescaling with  $\lambda$ , which we compute in each experiment with a grid search solving Eq. (1a)

### Algorithm 1: AD-based projected gradient descent

**Input:**  $\{(w_l, y_l)\}_{l \in \{1, \dots, L\}}$ : dataset.

**Output:**  $\hat{D}$ : minimizer to  $\mathcal{E}(D)$  in Eq. (1a).

**Params:**  $m, \eta_1, \eta_2, max\_itr1, max\_itr2, batch\_sz$

**Algorithm:**

Initialize  $D$  iid from  $\mathcal{N}(0, 10^{-4})$ .

Set *PyTorch AD* to track computations on  $D$ .

**for**  $t$  in  $\{0, \dots, max\_itr2 - 1\}$ :

**do**

$\mathcal{E}_{MSE}(D) \leftarrow 0$

**for**  $l \in \{t * batch\_sz, \dots, (t + 1) * batch\_sz - 1\}$

**do**

Initialize  $q$  iid from  $\mathcal{N}(0, 1)$ .

**for**  $i$  in  $\{1, \dots, max\_itr1\}$ :

**do**

$z \leftarrow \Pi_{B_\infty}(q - \eta_1 D^\top (Dq - y_l))$ .

Update  $q$  as in Eq. (5).

$\hat{w}(D, y_l) \leftarrow y_l - Dz$ .

$\mathcal{E}_{MSE}(D) \leftarrow \mathcal{E}_{MSE}(D) + \|\hat{w}(D, y_l) - w_l\|_2^2$ .

$\nabla \mathcal{E}_{MSE}(D) \leftarrow$  *PyTorch AD* gradient.

$D \leftarrow \Pi_{\mathcal{C}}(D - \frac{\eta_2}{batch\_sz} \nabla \mathcal{E}_{MSE}(D))$ .

**Return**  $\hat{D} = D$

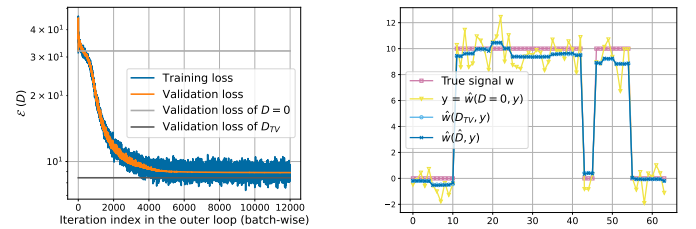


Fig. 2: Performance curves of the learnt  $\hat{D}$  in the case  $\sigma = 1$  in Fig. 1. Left: training and validation losses as in Eq. (1a) as a function of the iteration index. We plot the loss incurred by  $D_{TV}$  and  $D = 0$  on the validation set. Right: Recovered signals with  $D_{TV}$ ,  $D = 0$  and  $\hat{D}$  as in Eq. (1b) from a random signal  $y$  from the validation set.

*w.r.t.*  $\lambda$ . Matrices  $\hat{D}$  shown in this section have their columns sorted by magnitudes, to ease the comparison with  $D_{TV}$ . We consider *stochastic gradient descent* updates with batch size 64. Training set size: 640000, validation set size: 256. We adopt random white noise initialization with varying variance.

**Denoisier setup:** The value of  $\eta_1$  in Eq. (4) is assigned auto-

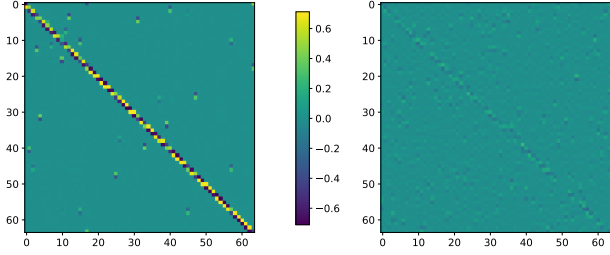


Fig. 3: Showing the difference made by our proposed admissible set  $\mathcal{C}$ . We plot  $\hat{D}$  produced by our algorithm. Left: with our proposed  $\mathcal{C}$ . Right: with  $\mathcal{C} = \mathbb{R}^{p \times p}$ . Dataset: Jumps in the dataset can randomly occur between any two values in  $[0, 10]$ , thus they have random amplitudes.  $\sigma = 0.5$ .

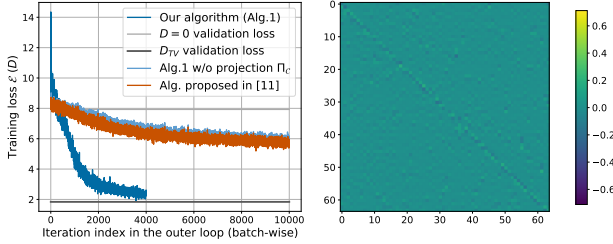


Fig. 4: Benchmark against the unconstrained optimization with  $\ell_1$ -smoothing proposed in [11], with identical problem setting to Fig. 3. Left: training loss curves of different algorithms. Right:  $\hat{D}$  produced by the algorithm in [11]. Dataset: same as in Fig. 3.

matically while guaranteeing convergence of  $z_i$ . This value is  $\eta_1 = 0.95 / \|\mathbf{D}^\top \mathbf{D}\|_2$ , where  $\|\mathbf{D}^\top \mathbf{D}\|_2$  is the Lipschitz constant of  $\nabla_{\frac{1}{2}} \|\mathbf{D}\mathbf{z} - \mathbf{y}\|_2^2$ . In all experiments, a tolerance threshold is used to determine the number of iterations, specifically, we assume convergence if  $\|z_{i+1} - z_i\|_\infty / \|z_i\|_\infty < 10^{-4}$ .

#### A. Projection proposed for this case study

We reduce the admissible set  $\mathcal{C}$  in our proposed algorithm 1 to dictionaries whose columns sum up to zero, *i.e.*,  $\mathcal{C} = \{\mathbf{D} | \mathbf{1}_p^\top \mathbf{D} = \mathbf{0}_m\}$ . This property is seen in the prior operator  $\mathbf{D}_{TV}$  our algorithm is expected to learn, and in the more general family of problems known as graph total variation [3]. This very simple prior greatly improves the results by filtering out many local minima. Different priors for other cases will be the goal of future investigations. Referring to the  $c$ -th column of  $\mathbf{D}$  by  $\mathbf{D}[:, c]$ , and by  $\text{mean}(\mathbf{D}[:, c])$  to the mean value of this column, we project  $\mathbf{D}_{t+1}$  as follows:

$$\mathbf{D}_{t+1}[:, c] = \mathbf{D}_{t+1}[:, c] - \text{mean}(\mathbf{D}_{t+1}[:, c]), \quad \forall c \in [m]. \quad (7)$$

In Fig. 3 and Fig. 4(left), we show the essential role of the centering projection in our proposed algorithm. We run our algorithm twice: *i*) with the projection; *ii*) ignoring the projection step; on the same dataset, and we plot the learned dictionary  $\hat{D}$  in both cases. Discontinuities in the used dataset are of random magnitudes (between any two values in  $[0, 10]$ ), and the noise has a standard deviation  $\sigma = 0.5$ . Our algorithm coupled with the projection successfully captures  $\mathbf{D}_{TV}$ -like structure from the dataset, unlike when the projection is not considered, which shows its capability in this problem setting.

#### B. Sensitivity w.r.t. to the noise level

When our observations are noise-free, *i.e.*,  $\sigma = 0$  and  $\mathbf{y} = \mathbf{w}$ , it is clear  $\mathbf{D} = \mathbf{0}$  is the optimal dictionary, as no regularization is required to retrieve the true signal. Therefore, when  $\sigma$  is small, we don't expect our algorithm to learn the structure in  $\mathbf{D}_{TV}$ , since in such case, its magnitude is of the same level as numerical errors, see Fig. 1 (b).

On the other hand, when the noise level is high, the piecewise constant prior is degraded and is poorly seen in observations. As a result, the learnt dictionary is a distorted version of  $\mathbf{D}_{TV}$ , as in Fig. 1 (f).

In between, it is important to verify that our algorithm is stable, and can extract  $\mathbf{D}_{TV}$  out from data. As shown in Fig. 1 (c-e), this is indeed true when  $\sigma \in [0.75, 2]$ , which spans a non-trivial range in SNR scale  $[12, 89]$ , given that our signals  $\mathbf{w}$  take values in  $\{0, 10\}$ .

For  $\sigma = 1$ , we show in Fig. 2 the evolution of training/validation losses as a function of the batch index through training, while benchmarking it against the same loss incurred when  $\mathbf{D} = \mathbf{D}_{TV}$  with optimal  $\lambda$  and  $\mathbf{D} = \mathbf{0}$  as a reference. We also present the same benchmark but *w.r.t.* the denoising performance, as in Eq. (1b), applied on a signal from the validation set.

#### C. Benchmark against the algorithm proposed in [11]

In Fig. 4, we compare our output to the one produced by the unconstrained learning algorithm based on smoothing  $\ell_1$  in [11]. We fix the  $\ell_1$ -smoothing parameter  $\epsilon = 10^{-3}$ . The dataset has discontinuities of random magnitudes in  $[0, 10]$ , and the noise standard deviation  $\sigma = 0.5$ , *i.e.*, same setup as in Fig. 3. Unlike our projected gradient descent, the opponent algorithm fails to inspect the  $\mathbf{D}_{TV}$  structure from data.

## V. CONCLUSION AND FUTURE WORK

To alleviate the difficulty of deriving gradients to learn analysis-sparsity dictionaries, that are optimized through a bilevel problem, we proposed to make use of automatic differentiation, a technique shown to have high proficiency in machine learning and deep learning. In the absence of a theoretical link between the automatic differentiation and the analytical one in our setting, our experiments on the piecewise constant signals reconstruction problem showed a proof-of-concept, and a promising methodology that can be applied in other problems setting. For the same case study, we also incorporated a simple column-wise centering projection, which significantly increased the stability of the algorithm and the quality of the learned dictionary.

Further investigations can be done to generalize our algorithm to 2D signals, particularly the 2D piecewise constant signals reconstruction, while trying to combine other constraints. Moreover, it would be interesting to study the problem again when the dimensions of the dictionary are not given, for instance to learn graph incidence matrices.

## REFERENCES

- [1] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind. Automatic differentiation in machine learning: a survey. *J. Mach. Learn. Res.*, 18, 2018.
- [2] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.
- [3] P. Berger, G. Hannak, and G. Matz. Graph signal recovery via primal-dual algorithms for total variation minimization. *IEEE Journal of Selected Topics in Signal Processing*, 11(6):842–855, 2017.
- [4] A. Chambolle and T. Pock. Learning consistent discretizations of the total variation. 2020.
- [5] A. Chambolle, V. Caselles, D. Cremers, M. Novaga, and T. Pock. An introduction to total variation for image analysis. *Theoretical foundations and numerical methods for sparse recovery*, 9(263-340):227, 2010.
- [6] M. Elad, P. Milanfar, and R. Rubinstein. Analysis versus synthesis in signal priors. *Inverse problems*, 23(3):947, 2007.
- [7] J. Grabmeier and E. Kaltofen. *Computer Algebra Handbook: Foundations, Applications, Systems*. Springer Science & Business Media, 2003.
- [8] M. E. Jerrell. Automatic differentiation and interval arithmetic for estimation of disequilibrium models. *Computational Economics*, 10(3):295–316, 1997.
- [9] L. Mancera and J. Portilla. L0-norm-based sparse representation through alternate projections. In *2006 International Conference on Image Processing*, pages 2089–2092, 2006.
- [10] M. T. McCann and S. Ravishankar. Supervised learning of sparsity-promoting regularizers for denoising. *arXiv preprint arXiv:2006.05521*, 2020.
- [11] G. Peyré and J. M. Fadili. Learning analysis sparsity priors. In *Sampta'11*, pages 4–pp, 2011.
- [12] S. Ravishankar and Y. Bresler. Sparsifying transform learning with efficient optimal updates and convergence guarantees. *IEEE Transactions on Signal Processing*, 63(9):2389–2404, 2015.
- [13] R. T. Rockafellar. *Conjugate duality and optimization*. SIAM, 1974.
- [14] P. Sprechmann, R. Litman, T. Ben Yakar, A. M. Bronstein, and G. Sapiro. Supervised sparse analysis and synthesis operators. *Advances in Neural Information Processing Systems*, 26:908–916, 2013.
- [15] A. Verma. An introduction to automatic differentiation. *Current Science*, pages 804–807, 2000.
- [16] M. Yaghoobi, S. Nam, R. Gribonval, and M. E. Davies. Analysis operator learning for overcomplete cospase representations. In *2011 19th European Signal Processing Conference*, pages 1470–1474. IEEE, 2011.
- [17] M. Yaghoobi, S. Nam, R. Gribonval, and M. E. Davies. Noise aware analysis operator learning for approximately cospase signals. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5409–5412. IEEE, 2012.