

Convex and Network Flow Optimization for Structured Sparsity

Julien Mairal^{* †}

JULIEN@STAT.BERKELEY.EDU

*Department of Statistics
University of California
Berkeley, CA 94720-1776, USA*

Rodolphe Jenatton^{* †}

RODOLPHE.JENATTON@INRIA.FR

Guillaume Obozinski[†]

GUILLAUME.OBOZINSKI@INRIA.FR

Francis Bach[†]

FRANCIS.BACH@INRIA.FR

INRIA - SIERRA Project-Team

Laboratoire d'Informatique de l'Ecole Normale Supérieure (INRIA/ENS/CNRS UMR 8548)

23, avenue d'Italie 75214 Paris CEDEX 13, France.

Abstract

We consider a class of learning problems regularized by a structured sparsity-inducing norm defined as the sum of ℓ_2 - or ℓ_∞ -norms over groups of variables. Whereas much effort has been put in developing fast optimization techniques when the groups are disjoint or embedded in a hierarchy, we address here the case of general overlapping groups. To this end, we present two different strategies: On the one hand, we show that the proximal operator associated with a sum of ℓ_∞ -norms can be computed exactly in polynomial time by solving a *quadratic min-cost flow problem*, allowing the use of accelerated proximal gradient methods. On the other hand, we use proximal splitting techniques, and address an equivalent formulation with non-overlapping groups, but in higher dimension and with additional constraints. We propose efficient and scalable algorithms exploiting these two strategies, which are significantly faster than alternative approaches. We illustrate these methods with several problems such as CUR matrix factorization, multi-task learning of tree-structured dictionaries, background subtraction in video sequences, image denoising with wavelets, and topographic dictionary learning of natural image patches.

Keywords: Convex optimization, proximal methods, sparse coding, structured sparsity, matrix factorization, network flow optimization, alternating direction method of multipliers.

1. Introduction

Sparse linear models have become a popular framework for dealing with various unsupervised and supervised tasks in machine learning and signal processing. In such models, linear combinations of small sets of variables are selected to describe the data. Regularization by the ℓ_1 -norm has emerged as a powerful tool for addressing this variable selection problem, relying on both a well-developed theory (see Tibshirani, 1996; Chen et al., 1999; Mallat, 1999; Bickel et al., 2009; Wainwright, 2009, and references therein) and efficient algorithms (Efron et al., 2004; Nesterov, 2007; Beck and Teboulle, 2009; Needell and Tropp, 2009; Combettes and Pesquet, 2010).

*. These authors contributed equally.

†. When most of this work was conducted, all authors were affiliated to INRIA, WILLOW Project-Team.

The ℓ_1 -norm primarily encourages sparse solutions, regardless of the potential structural relationships (e.g., spatial, temporal or hierarchical) existing between the variables. Much effort has recently been devoted to designing sparsity-inducing regularizations capable of encoding higher-order information about the patterns of non-zero coefficients (Cehver et al., 2008; Jenatton et al., 2009; Jacob et al., 2009; Zhao et al., 2009; He and Carin, 2009; Huang et al., 2009; Baraniuk et al., 2010; Micchelli et al., 2010), with successful applications in bioinformatics (Jacob et al., 2009; Kim and Xing, 2010), topic modeling (Jenatton et al., 2010a, 2011) and computer vision (Cehver et al., 2008; Huang et al., 2009; Jenatton et al., 2010b). By considering sums of norms of appropriate subsets, or *groups*, of variables, these regularizations control the sparsity patterns of the solutions. The underlying optimization is usually difficult, in part because it involves nonsmooth components.

Our first strategy uses proximal gradient methods, which have proven to be effective in this context, essentially because of their fast convergence rates and their ability to deal with large problems (Nesterov, 2007; Beck and Teboulle, 2009). They can handle differentiable loss functions with Lipschitz-continuous gradient, and we show in this paper how to use them with a regularization term composed of a sum of ℓ_∞ -norms. The second strategy we consider exploits proximal splitting methods (see Combettes and Pesquet, 2008, 2010; Goldfarb and Ma, 2009; Tomioka et al., 2011; Qin and Goldfarb, 2011; Boyd et al., 2011, and references therein), which builds upon an equivalent formulation with non-overlapping groups, but in a higher dimensional space and with additional constraints.¹ More precisely, we make four main contributions:

- We show that the *proximal operator* associated with the sum of ℓ_∞ -norms with overlapping groups can be computed efficiently and exactly by solving a *quadratic min-cost flow* problem, thereby establishing a connection with the network flow optimization literature.² This is the main contribution of the paper, which allows us to use proximal gradient methods in the context of structured sparsity.
- We prove that the dual norm of the sum of ℓ_∞ -norms can also be evaluated efficiently, which enables us to compute duality gaps for the corresponding optimization problems.
- We present proximal splitting methods for solving structured sparse regularized problems.
- We demonstrate that our methods are relevant for various applications whose practical success is made possible by our algorithmic tools and efficient implementations. First, we introduce a new CUR matrix factorization technique exploiting structured sparse regularization, built upon the links drawn by Bien et al. (2010) between CUR decomposition (Mahoney and Drineas, 2009) and sparse regularization. Then, we illustrate our algorithms with different tasks: video background subtraction, estimation of hierarchical structures for dictionary learning of natural image patches (Jenatton et al., 2010a, 2011), wavelet image de-

1. The idea of using this class of algorithms for solving structured sparse problems was first suggested to us by Jean-Christophe Pesquet and Patrick-Louis Combettes. It was also suggested to us later by Ryota Tomioka, who briefly mentioned this possibility in (Tomioka et al., 2011). It can also briefly be found in (Boyd et al., 2011), and in details in the work of Qin and Goldfarb (2011) which was conducted as the same time as ours. It was also used in a related context by Sprechmann et al. (2010) for solving optimization problems with hierarchical norms.

2. Interestingly, this is not the first time that network flow optimization tools have been used to solve sparse regularized problems with proximal methods. Such a connection was recently established by Chambolle and Darbon (2009) in the context of total variation regularization, and similarly by Hoeffling (2010) for the fused Lasso. One can also find the use of maximum flow problems for non-convex penalties in the work of Cehver et al. (2008) which combines Markov random fields and sparsity.

noising with a structured sparse prior, and topographic dictionary learning of natural image patches (Hyvärinen et al., 2001; Kavukcuoglu et al., 2009; Garrigues and Olshausen, 2010).

Note that this paper extends a shorter version published in *Advances in Neural Information Processing Systems* (Mairal et al., 2010b), by adding new experiments (CUR matrix factorization, wavelet image denoising and topographic dictionary learning), presenting the proximal splitting methods, providing the full proofs of the optimization results, and adding numerous discussions.

1.1 Notation

Vectors are denoted by bold lower case letters and matrices by upper case ones. We define for $q \geq 1$ the ℓ_q -norm of a vector \mathbf{x} in \mathbb{R}^m as $\|\mathbf{x}\|_q \triangleq (\sum_{i=1}^m |\mathbf{x}_i|^q)^{1/q}$, where \mathbf{x}_i denotes the i -th coordinate of \mathbf{x} , and $\|\mathbf{x}\|_\infty \triangleq \max_{i=1,\dots,m} |\mathbf{x}_i| = \lim_{q \rightarrow \infty} \|\mathbf{x}\|_q$. We also define the ℓ_0 -pseudo-norm as the number of nonzero elements in a vector:³ $\|\mathbf{x}\|_0 \triangleq \#\{i \text{ s.t. } \mathbf{x}_i \neq 0\} = \lim_{q \rightarrow 0^+} (\sum_{i=1}^m |\mathbf{x}_i|^q)$. We consider the Frobenius norm of a matrix \mathbf{X} in $\mathbb{R}^{m \times n}$: $\|\mathbf{X}\|_F \triangleq (\sum_{i=1}^m \sum_{j=1}^n \mathbf{X}_{ij}^2)^{1/2}$, where \mathbf{X}_{ij} denotes the entry of \mathbf{X} at row i and column j . Finally, for a scalar y , we denote $(y)_+ \triangleq \max(y, 0)$. For an integer $p > 0$, we denote by $2^{\{1,\dots,p\}}$ the powerset composed of the 2^p subsets of $\{1, \dots, p\}$.

The rest of this paper is organized as follows: Section 2 presents structured sparse models and related work. Section 3 is devoted to proximal gradient algorithms, and Section 4 to proximal splitting methods. Section 5 presents several experiments and applications demonstrating the effectiveness of our approach and Section 6 concludes the paper.

2. Structured Sparse Models

We are interested in machine learning problems where the solution is not only known beforehand to be sparse—that is, the solution has only a few non-zero coefficients, but also to form non-zero patterns with a specific structure. It is indeed possible to encode additional knowledge in the regularization other than just sparsity. For instance, one may want the non-zero patterns to be structured in the form of non-overlapping groups (Turlach et al., 2005; Yuan and Lin, 2006; Stojnic et al., 2009; Obozinski et al., 2010), in a tree (Zhao et al., 2009; Bach, 2009; Jenatton et al., 2010a, 2011), or in overlapping groups (Jenatton et al., 2009; Jacob et al., 2009; Huang et al., 2009; Baraniuk et al., 2010; Cehver et al., 2008; He and Carin, 2009), which is the setting we are interested in here.

As for classical non-structured sparse models, there are basically two lines of research, that either (A) deal with nonconvex and combinatorial formulations that are in general computationally intractable and addressed with greedy algorithms or (B) concentrate on convex relaxations solved with convex programming methods.

2.1 Nonconvex Approaches

A first approach introduced by Baraniuk et al. (2010) consists in imposing that the sparsity pattern of a solution (i.e., its set of non-zero coefficients) is in a predefined subset of groups of variables $\mathcal{G} \subseteq 2^{\{1,\dots,p\}}$. Given this a priori knowledge, a greedy algorithm (Needell and Tropp, 2009) is used

3. Note that it would be more proper to write $\|\mathbf{x}\|_0^0$ instead of $\|\mathbf{x}\|_0$ to be consistent with the traditional notation $\|\mathbf{x}\|_q$. However, for the sake of simplicity, we will keep this notation unchanged in the rest of the paper.

to address the following nonconvex structured sparse decomposition problem

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \text{ s.t. } \text{Supp}(\mathbf{w}) \in \mathcal{G} \text{ and } \|\mathbf{w}\|_0 \leq s,$$

where s is a specified sparsity level (number of nonzeros coefficients), \mathbf{y} in \mathbb{R}^m is an observed signal, \mathbf{X} is a design matrix in $\mathbb{R}^{m \times p}$ and $\text{Supp}(\mathbf{w})$ is the support of \mathbf{w} (set of non-zero entries).

In a different approach motivated by the minimum description length principle (see Barron et al., 1998), Huang et al. (2009) consider a collection of groups $\mathcal{G} \subseteq 2^{\{1, \dots, p\}}$, and define a ‘‘coding length’’ for every group in \mathcal{G} , which in turn is used to define a coding length for every pattern in $2^{\{1, \dots, p\}}$. Using this tool, they propose a regularization function $\text{cl} : \mathbb{R}^p \rightarrow \mathbb{R}$ such that for a vector \mathbf{w} in \mathbb{R}^p , $\text{cl}(\mathbf{w})$ represents the number of bits that are used for encoding \mathbf{w} . The corresponding optimization problem is also addressed with a greedy procedure:

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \text{ s.t. } \text{cl}(\mathbf{w}) \leq s,$$

Intuitively, this formulation encourages solutions \mathbf{w} whose sparsity patterns have a small coding length, meaning in practice that they can be represented by a union of a small number of groups. Even though they are related, this model is different from the one of Baraniuk et al. (2010).

These two approaches are encoding a priori knowledge on the shape of non-zero patterns that the solution of a regularized problem should have. A different point of view consists of modelling the zero patterns of the solution—that is, define groups of variables that should be encouraged to be set to zero together. After defining a set $\mathcal{G} \subseteq 2^{\{1, \dots, p\}}$ of such groups of variables, the following penalty can naturally be used as a regularization to induce the desired property

$$\psi(\mathbf{w}) \triangleq \sum_{g \in \mathcal{G}} \eta_g \delta^g(\mathbf{w}), \text{ with } \delta^g(\mathbf{w}) \triangleq \begin{cases} 1 & \text{if there exists } j \in g \text{ such that } \mathbf{w}_j \neq 0, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where the η_g 's are positive weights. This penalty was considered by Bach (2010), who showed that the convex envelope of such nonconvex functions (more precisely strictly positive, non-increasing submodular functions of $\text{Supp}(\mathbf{w})$, see Fujishige, 2005) when restricted on the unit ℓ_∞ -ball, are in fact types of structured sparsity-inducing norms which are the topic of the next section.

2.2 Convex Approaches with Sparsity-Inducing Norms

In this paper, we are interested in convex regularizations which induce structured sparsity. Generally, we consider the following optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w}) + \lambda \Omega(\mathbf{w}), \quad (2)$$

where $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is a convex function (usually an empirical risk in machine learning and a data-fitting term in signal processing), and $\Omega : \mathbb{R}^p \rightarrow \mathbb{R}$ is a structured sparsity-inducing norm, defined as

$$\Omega(\mathbf{w}) \triangleq \sum_{g \in \mathcal{G}} \eta_g \|\mathbf{w}_g\|, \quad (3)$$

where $\mathcal{G} \subseteq 2^{\{1, \dots, p\}}$ is a set of groups of variables, the vector \mathbf{w}_g in $\mathbb{R}^{|g|}$ represents the coefficients of \mathbf{w} indexed by g in \mathcal{G} , the scalars η_g are positive weights, and $\|\cdot\|$ denotes the ℓ_2 - or ℓ_∞ -norm. We now consider different cases:

- When \mathcal{G} is the set of singletons—that is $\mathcal{G} \triangleq \{\{1\}, \{2\}, \dots, \{p\}\}$, and all the η_g are equal to one, Ω is the ℓ_1 -norm, which is well known to induce sparsity. This leads for instance to the Lasso (Tibshirani, 1996) or equivalently to basis pursuit (Chen et al., 1999).
- If \mathcal{G} is a partition of $\{1, \dots, p\}$, i.e. the groups do not overlap, variables are selected in groups rather than individually. When the coefficients of the solution are known to be organized in such a way, explicitly encoding the a priori group structure in the regularization can improve the prediction performance and/or interpretability of the learned models (Turlach et al., 2005; Yuan and Lin, 2006; Roth and Fischer, 2008; Stojnic et al., 2009; Huang and Zhang, 2010; Obozinski et al., 2010). Such a penalty is commonly called group-Lasso penalty.
- When the groups overlap, Ω is still a norm and sets groups of variables to zero together (Jenatton et al., 2009). The latter setting has first been considered for hierarchies (Zhao et al., 2009; Kim and Xing, 2010; Bach, 2009; Jenatton et al., 2010a, 2011), and then extended to general group structures (Jenatton et al., 2009). Solving Eq. (2) in this context is a challenging problem which is the topic of this paper.

Note that other types of structured-sparsity inducing norms have also been introduced, notably the approach of Jacob et al. (2009), which penalizes the following quantity

$$\Omega'(\mathbf{w}) \triangleq \min_{\xi=(\xi^g)_{g \in \mathcal{G}} \in \mathbb{R}^{p \times |\mathcal{G}|}} \sum_{g \in \mathcal{G}} \eta_g \|\xi^g\| \quad \text{s.t. } \mathbf{w} = \sum_{g \in \mathcal{G}} \xi^g \quad \text{and } \forall g, \text{Supp}(\xi^g) \subseteq g.$$

This penalty, which is also a norm, can be seen as a convex relaxation of the regularization introduced by Huang et al. (2009), and encourages the sparsity pattern of the solution to be a union of a small number of groups. Even though both Ω and Ω' appear under the terminology of “structured sparsity with overlapping groups”, they have in fact significantly different purposes and algorithmic treatments. For example, Jacob et al. (2009) consider the problem of selecting genes in a gene network which can be represented as the union of a few predefined pathways in the graph (groups of genes), which overlap. In this case, it is natural to use the norm Ω' instead of Ω . On the other hand, we present a matrix factorization task in Section 5.3, where the set of zero-patterns should be a union of groups, naturally leading to the use of Ω . Dealing with Ω' is therefore relevant, but out of the scope of this paper.

2.3 Convex Optimization Methods Proposed in the Literature

Generic approaches to solve Eq. (2) mostly rely on subgradient descent schemes (see Bertsekas, 1999), and interior-point methods (Boyd and Vandenberghe, 2004). These generic tools do not scale well to large problems and/or do not naturally handle sparsity (the solutions they return may have small values but no “true” zeros). These two points prompt the need for dedicated methods.

To the best of our knowledge, only a few recent papers have addressed problem Eq. (2) with dedicated optimization procedures, and in fact, only when Ω is a linear combination of ℓ_2 -norms. In this setting, a first line of work deals with the non-smoothness of Ω by expressing the norm as the minimum over a set of smooth functions. At the cost of adding new variables (to describe the set of smooth functions), the problem becomes more amenable to optimization. In particular, reweighted- ℓ_2 schemes consist of approximating the norm Ω by successive quadratic upper bounds (Argyriou et al., 2008; Rakotomamonjy et al., 2008; Jenatton et al., 2010b; Michelli et al., 2010). It is possible

to show for instance that

$$\Omega(\mathbf{w}) = \min_{(z_g)_{g \in \mathcal{G}} \in \mathbb{R}_+^{|\mathcal{G}|}} \frac{1}{2} \left\{ \sum_{g \in \mathcal{G}} \frac{\eta_g^2 \|\mathbf{w}_g\|_2^2}{z_g} + z_g \right\}.$$

Plugging the previous relationship into Eq. (2), the optimization can then be performed by alternating between the updates of \mathbf{w} and the additional variables $(z_g)_{g \in \mathcal{G}}$.⁴ When the norm Ω is defined as a linear combination of ℓ_∞ -norms, we are not aware of the existence of such variational formulations.

Problem (2) has also been addressed with working-set algorithms (Bach, 2009; Jenatton et al., 2009; Schmidt and Murphy, 2010). The main idea of these methods is to solve a sequence of increasingly larger subproblems of (2). Each subproblem consists of an instance of Eq. (2) reduced to a specific subset of variables known as the *working set*. As long as some predefined optimality conditions are not satisfied, the working set is augmented with selected inactive variables (for more details, see Bach et al., 2011).

The last approach we would like to mention is that of Chen et al. (2010), who used a smoothing technique introduced by Nesterov (2005). A smooth approximation Ω_μ of Ω is used, when Ω is a sum of ℓ_2 -norms, and μ is a parameter controlling the trade-off between smoothness of Ω_μ and quality of the approximation. Then, Eq. (2) is solved with accelerated gradient techniques (Beck and Teboulle, 2009; Nesterov, 2007) but Ω_μ is substituted to the regularization Ω . Depending on the required precision for solving the original problem, this method provides a natural choice for the parameter μ , with a known convergence rate. A drawback is that it requires to choose the precision of the optimization beforehand. Moreover, since a ℓ_1 -norm is added to the smoothed Ω_μ , the solutions returned by the algorithm might be sparse but possibly without respecting the structure encoded by Ω . This should be contrasted with other smoothing techniques, e.g., the reweighted- ℓ_2 scheme we mentioned above, where the solutions are only approximately sparse.

3. Optimization with Proximal Gradient Methods

We address in this section the problem of solving Eq. (2) under the following assumptions:

- *f is differentiable with Lipschitz-continuous gradient.* For machine learning problems, this hypothesis holds when f is for example the square, logistic or multi-class logistic loss (see Shawe-Taylor and Cristianini, 2004).
- *Ω is a sum of ℓ_∞ -norms.* Even though the ℓ_2 -norm is sometimes used in the literature (Jenatton et al., 2009), and is in fact used later in Section 4, the ℓ_∞ -norm is piecewise linear, and we take advantage of this property in this work.

To the best of our knowledge, no dedicated optimization method has been developed for this setting. Following Jenatton et al. (2010a, 2011) who tackled the particular case of hierarchical norms, we propose to use proximal gradient methods, which we now introduce.

4. Note that such a scheme is interesting only if the optimization with respect to \mathbf{w} is simple, which is typically the case with the square loss function (Bach et al., 2011). Moreover, for this alternating scheme to be provably convergent, the variables $(z_g)_{g \in \mathcal{G}}$ have to be bounded away from zero, resulting in solutions whose entries may have small values, but not “true” zeros.

3.1 Proximal Gradient Methods

Proximal methods have drawn increasing attention in the signal processing (e.g., Wright et al., 2009b; Combettes and Pesquet, 2010, and numerous references therein) and the machine learning communities (e.g., Bach et al., 2011, and references therein), especially because of their convergence rates (optimal for the class of first-order techniques) and their ability to deal with large nonsmooth convex problems (e.g., Nesterov, 2007; Beck and Teboulle, 2009).

These methods are iterative procedures that can be seen as an extension of gradient-based techniques when the objective function to minimize has a nonsmooth part. The simplest version of this class of methods linearizes at each iteration the function f around the current estimate $\tilde{\mathbf{w}}$, and this estimate is updated as the (unique by strong convexity) solution of the *proximal* problem, defined as:

$$\min_{\mathbf{w} \in \mathbb{R}^p} f(\tilde{\mathbf{w}}) + (\mathbf{w} - \tilde{\mathbf{w}})^\top \nabla f(\tilde{\mathbf{w}}) + \lambda \Omega(\mathbf{w}) + \frac{L}{2} \|\mathbf{w} - \tilde{\mathbf{w}}\|_2^2.$$

The quadratic term keeps the update in a neighborhood where f is close to its linear approximation, and $L > 0$ is a parameter which is an upper bound on the Lipschitz constant of ∇f . This problem can be equivalently rewritten as:

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \left\| \tilde{\mathbf{w}} - \frac{1}{L} \nabla f(\tilde{\mathbf{w}}) - \mathbf{w} \right\|_2^2 + \frac{\lambda}{L} \Omega(\mathbf{w}),$$

Solving *efficiently* and exactly this problem allows to attain the fast convergence rates of proximal methods, i.e., reaching a precision of $O(\frac{L}{k^2})$ in k iterations.⁵ In addition, when the nonsmooth term Ω is not present, the previous proximal problem exactly leads to the standard gradient update rule. More generally, we define the *proximal operator*:

Definition 1 (Proximal Operator)

The proximal operator associated with our regularization term $\lambda \Omega$, which we denote by $\text{Prox}_{\lambda \Omega}$, is the function that maps a vector $\mathbf{u} \in \mathbb{R}^p$ to the unique solution of

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{u} - \mathbf{w}\|_2^2 + \lambda \Omega(\mathbf{w}). \quad (4)$$

This operator was initially introduced by Moreau (1962) to generalize the projection operator onto a convex set. What makes proximal methods appealing to solve sparse decomposition problems is that this operator can often be computed in closed form. For instance,

- When Ω is the ℓ_1 -norm—that is $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1$ —the proximal operator is the well-known elementwise soft-thresholding operator,

$$\forall j \in \{1, \dots, p\}, \quad \mathbf{u}_j \mapsto \text{sign}(\mathbf{u}_j) (|\mathbf{u}_j| - \lambda)_+ = \begin{cases} 0 & \text{if } |\mathbf{u}_j| \leq \lambda \\ \text{sign}(\mathbf{u}_j) (|\mathbf{u}_j| - \lambda) & \text{otherwise.} \end{cases}$$

- When Ω is a group-Lasso penalty with ℓ_2 -norms—that is, $\Omega(\mathbf{u}) = \sum_{g \in \mathcal{G}} \|\mathbf{u}_g\|_2$, with \mathcal{G} being a partition of $\{1, \dots, p\}$, the proximal problem is *separable* in every group, and the solution is a generalization of the soft-thresholding operator to groups of variables:

$$\forall g \in \mathcal{G}, \quad \mathbf{u}_g \mapsto \mathbf{u}_g - \Pi_{\|\cdot\|_2 \leq \lambda}[\mathbf{u}_g] = \begin{cases} 0 & \text{if } \|\mathbf{u}_g\|_2 \leq \lambda \\ \frac{\|\mathbf{u}_g\|_2 - \lambda}{\|\mathbf{u}_g\|_2} \mathbf{u}_g & \text{otherwise,} \end{cases}$$

5. Note, however, that fast convergence rates can also be achieved while solving approximately the proximal problem (see Schmidt et al., 2011, for more details).

where $\Pi_{\|\cdot\|_2 \leq \lambda}$ denotes the orthogonal projection onto the ball of the ℓ_2 -norm of radius λ .

- When Ω is a group-Lasso penalty with ℓ_∞ -norms—that is, $\Omega(\mathbf{u}) = \sum_{g \in \mathcal{G}} \|\mathbf{u}_g\|_\infty$, with \mathcal{G} being a partition of $\{1, \dots, p\}$, the solution is a different group-thresholding operator:

$$\forall g \in \mathcal{G}, \quad \mathbf{u}_g \mapsto \mathbf{u}_g - \Pi_{\|\cdot\|_1 \leq \lambda}[\mathbf{u}_g],$$

where $\Pi_{\|\cdot\|_1 \leq \lambda}$ denotes the orthogonal projection onto the ℓ_1 -ball of radius λ , which can be solved in $O(p)$ operations (Brucker, 1984; Maculan and de Paula, 1989). Note that when $\|\mathbf{u}_g\|_1 \leq \lambda$, we have a group-thresholding effect, with $\mathbf{u}_g - \Pi_{\|\cdot\|_1 \leq \lambda}[\mathbf{u}_g] = 0$.

- When Ω is a tree-structured sum of ℓ_2 - or ℓ_∞ -norms as introduced by Zhao et al. (2009)—meaning that two groups are either disjoint or one is included in the other, the solution admits a closed form. Let \preceq be a total order on \mathcal{G} such that for g_1, g_2 in \mathcal{G} , $g_1 \preceq g_2$ if and only if either $g_1 \subset g_2$ or $g_1 \cap g_2 = \emptyset$.⁶ Then, if $g_1 \preceq \dots \preceq g_{|\mathcal{G}|}$, and if we define Prox^g as (a) the proximal operator $\mathbf{u}_g \mapsto \text{Prox}_{\lambda \eta_g \|\cdot\|}(\mathbf{u}_g)$ on the subspace corresponding to group g and (b) the identity on the orthogonal, Jenatton et al. (2010a, 2011) showed that:

$$\text{Prox}_{\lambda \Omega} = \text{Prox}^{g_m} \circ \dots \circ \text{Prox}^{g_1}, \quad (5)$$

which can be computed in $O(p)$ operations. It also includes the sparse group Lasso (sum of group-Lasso penalty and ℓ_1 -norm) of Friedman et al. (2010) and Sprechmann et al. (2010).

The first contribution of our paper is to address the case of general overlapping groups with ℓ_∞ -norm.

3.2 Dual of the Proximal Operator

We now show that, for a set \mathcal{G} of general overlapping groups, a convex dual of the proximal problem (4) can be reformulated as a *quadratic min-cost flow problem*. We then propose an efficient algorithm to solve it exactly, as well as a related algorithm to compute the dual norm of Ω . We start by considering the dual formulation to problem (4) introduced by Jenatton et al. (2010a, 2011):

Lemma 1 (Dual of the proximal problem, Jenatton et al., 2010a, 2011)

Given \mathbf{u} in \mathbb{R}^p , consider the problem

$$\min_{\xi \in \mathbb{R}^{p \times |\mathcal{G}|}} \frac{1}{2} \|\mathbf{u} - \sum_{g \in \mathcal{G}} \xi^g\|_2^2 \quad \text{s.t.} \quad \forall g \in \mathcal{G}, \|\xi^g\|_1 \leq \lambda \eta_g \quad \text{and} \quad \xi_j^g = 0 \text{ if } j \notin g, \quad (6)$$

where $\xi = (\xi^g)_{g \in \mathcal{G}}$ is in $\mathbb{R}^{p \times |\mathcal{G}|}$, and ξ_j^g denotes the j -th coordinate of the vector ξ^g . Then, every solution $\xi^* = (\xi^{*g})_{g \in \mathcal{G}}$ of Eq. (6) satisfies $\mathbf{w}^* = \mathbf{u} - \sum_{g \in \mathcal{G}} \xi^{*g}$, where \mathbf{w}^* is the solution of Eq. (4) when Ω is a weighted sum of ℓ_∞ -norms.

Without loss of generality,⁷ we assume from now on that the scalars \mathbf{u}_j are all non-negative, and we constrain the entries of ξ to be so. Such a formulation introduces $p|\mathcal{G}|$ dual variables which can be much greater than p , the number of primal variables, but it removes the issue of overlapping regularization. We now associate a graph with problem (6), on which the variables ξ_j^g , for g in \mathcal{G} and j in g , can be interpreted as measuring the components of a flow.

6. For a tree-structured set \mathcal{G} , such an order exists.

7. Let ξ^* denote a solution of Eq. (6). Optimality conditions of Eq. (6) derived in Jenatton et al. (2010a, 2011) show that for all j in $\{1, \dots, p\}$, the signs of the non-zero coefficients ξ_j^{*g} for g in \mathcal{G} are the same as the signs of the entries \mathbf{u}_j . To solve Eq. (6), one can therefore flip the signs of the negative variables \mathbf{u}_j , then solve the modified dual formulation (with non-negative variables), which gives the magnitude of the entries ξ_j^{*g} (the signs of these being known).

3.3 Graph Model

Let G be a directed graph $G = (V, E, s, t)$, where V is a set of vertices, $E \subseteq V \times V$ a set of arcs, s a source, and t a sink. For all arcs in E , we define a non-negative capacity constant, and as done classically in the network flow literature (Ahuja et al., 1993; Bertsekas, 1998), we define a *flow* as a non-negative function on arcs that satisfies capacity constraints on all arcs (the value of the flow on an arc is less than or equal to the arc capacity) and conservation constraints on all vertices (the sum of incoming flows at a vertex is equal to the sum of outgoing flows) except for the source and the sink. For every arc e in E , we also define a real-valued cost function, which depends on the value of the flow on e . We now introduce the *canonical* graph G associated with our optimization problem:

Definition 2 (Canonical Graph)

Let $\mathcal{G} \subseteq \{1, \dots, p\}$ be a set of groups, and $(\eta_g)_{g \in \mathcal{G}}$ be positive weights. The canonical graph $G = (V, E, s, t)$ is the unique graph defined as follows:

1. $V = V_u \cup V_{gr}$, where V_u is a vertex set of size p , one vertex being associated to each index j in $\{1, \dots, p\}$, and V_{gr} is a vertex set of size $|\mathcal{G}|$, one vertex per group g in \mathcal{G} . We thus have $|V| = |\mathcal{G}| + p$. For simplicity, we identify groups g in \mathcal{G} and indices j in $\{1, \dots, p\}$ with vertices of the graph, such that one can from now on refer to “vertex j ” or “vertex g ”.
2. For every group g in \mathcal{G} , E contains an arc (s, g) . These arcs have capacity $\lambda\eta_g$ and zero cost.
3. For every group g in \mathcal{G} , and every index j in g , E contains an arc (g, j) with zero cost and infinite capacity. We denote by ξ_j^g the flow on this arc.
4. For every index j in $\{1, \dots, p\}$, E contains an arc (j, t) with infinite capacity and a cost $\frac{1}{2}(\mathbf{u}_j - \bar{\xi}_j)^2$, where $\bar{\xi}_j$ is the flow on (j, t) .

Examples of canonical graphs are given in Figures 1a-c for three simple group structures. The flows ξ_j^g associated with G can now be identified with the variables of problem (6). Since we have assumed the entries of \mathbf{u} to be non-negative, we can now reformulate Eq. (6) as

$$\min_{\xi \in \mathbb{R}_+^{p \times |\mathcal{G}|}, \bar{\xi} \in \mathbb{R}^p} \sum_{j=1}^p \frac{1}{2} (\mathbf{u}_j - \bar{\xi}_j)^2 \quad \text{s.t.} \quad \bar{\xi} = \sum_{g \in \mathcal{G}} \xi^g \quad \text{and} \quad \forall g \in \mathcal{G}, \left\{ \sum_{j \in g} \xi_j^g \leq \lambda\eta_g \quad \text{and} \quad \text{Supp}(\xi^g) \subseteq g \right\}. \quad (7)$$

Indeed,

- the only arcs with a cost are those leading to the sink, which have the form (j, t) , where j is the index of a variable in $\{1, \dots, p\}$. The sum of these costs is $\sum_{j=1}^p \frac{1}{2} (\mathbf{u}_j - \bar{\xi}_j)^2$, which is the objective function minimized in Eq. (7);
- by flow conservation, we necessarily have $\bar{\xi}_j = \sum_{g \in \mathcal{G}} \xi_j^g$ in the canonical graph;
- the only arcs with a capacity constraints are those coming out of the source, which have the form (s, g) , where g is a group in \mathcal{G} . By flow conservation, the flow on an arc (s, g) is $\sum_{j \in g} \xi_j^g$ which should be less than $\lambda\eta_g$ by capacity constraints;
- all other arcs have the form (g, j) , where g is in \mathcal{G} and j is in g . Thus, $\text{Supp}(\xi^g) \subseteq g$.

Therefore we have shown that finding a flow *minimizing the sum of the costs* on such a graph is equivalent to solving problem (6). When some groups are included in others, the canonical graph can be simplified to yield a graph with a smaller number of edges. Specifically, if h and g are groups with $h \subset g$, the edges (g, j) for $j \in h$ carrying a flow ξ_j^g can be removed and replaced by a single edge (g, h) of infinite capacity and zero cost, carrying the flow $\sum_{j \in h} \xi_j^g$. This simplification is illustrated in Figure 1d, with a graph equivalent to the one of Figure 1c. This does not change the optimal value of $\bar{\xi}^*$, which is the quantity of interest for computing the optimal primal variable \mathbf{w}^* . We present in Appendix A a formal definition of equivalent graphs. These simplifications are useful in practice, since they reduce the number of edges in the graph and improve the speed of our algorithms.

3.4 Computation of the Proximal Operator

Quadratic min-cost flow problems have been well studied in the operations research literature (Hochbaum and Hong, 1995). One of the simplest cases, where \mathcal{G} contains a single group as in Figure 1a, is solved by an orthogonal projection on the ℓ_1 -ball of radius $\lambda \eta_g$. It has been shown, both in machine learning (Duchi et al., 2008) and operations research (Hochbaum and Hong, 1995; Brucker, 1984), that such a projection can be computed in $O(p)$ operations. When the group structure is a tree as in Figure 1d, strategies developed in the two communities are also similar (Jenatton et al., 2010a; Hochbaum and Hong, 1995),⁸ and solve the problem in $O(pd)$ operations, where d is the depth of the tree.

The general case of overlapping groups is more difficult. Hochbaum and Hong (1995) have shown that *quadratic min-cost flow problems* can be reduced to a specific *parametric max-flow* problem, for which an efficient algorithm exists (Gallo et al., 1989).⁹ While this generic approach could be used to solve Eq. (6), we propose to use Algorithm 1 that also exploits the fact that our graphs have non-zero costs only on edges leading to the sink. As shown in Appendix D, it has a significantly better performance in practice. This algorithm clearly shares some similarities with existing approaches in network flow optimization such as the simplified version of Gallo et al. (1989) presented by Babenko and Goldberg (2006) that uses a divide and conquer strategy. Moreover, an equivalent algorithm exists for minimizing convex functions over polymatroid sets (Groenevelt, 1991). This equivalence, a priori non trivial, is uncovered through a representation of structured sparsity-inducing norms via submodular functions, which was recently proposed by Bach (2010).

The intuition behind our algorithm, `computeFlow` (see Algorithm 1), is the following: since $\bar{\xi} = \sum_{g \in \mathcal{G}} \xi^g$ is the only value of interest to compute the solution of the proximal operator $\mathbf{w} = \mathbf{u} - \bar{\xi}$, the first step looks for a candidate value γ for $\bar{\xi}$ by solving the following relaxed version of problem (7):

$$\arg \min_{\gamma \in \mathbb{R}^p} \sum_{j \in V_u} \frac{1}{2} (\mathbf{u}_j - \gamma_j)^2 \quad \text{s.t.} \quad \sum_{j \in V_u} \gamma_j \leq \lambda \sum_{g \in V_{gr}} \eta_g. \quad (8)$$

The cost function here is the same as in problem (7), but the constraints are weaker: Any feasible point of problem (7) is also feasible for problem (8). This problem can be solved in linear time (Brucker, 1984). Its solution, which we denote γ for simplicity, provides the lower bound $\|\mathbf{u} - \gamma\|_2^2/2$ for the optimal cost of problem (7).

8. Note however that, while Hochbaum and Hong (1995) only consider a tree-structured sum of ℓ_∞ -norms, the results from Jenatton et al. (2010a) also apply for a sum of ℓ_2 -norms.

9. By definition, a parametric max-flow problem consists in solving, for every value of a parameter, a max-flow problem on a graph whose arc capacities depend on this parameter.

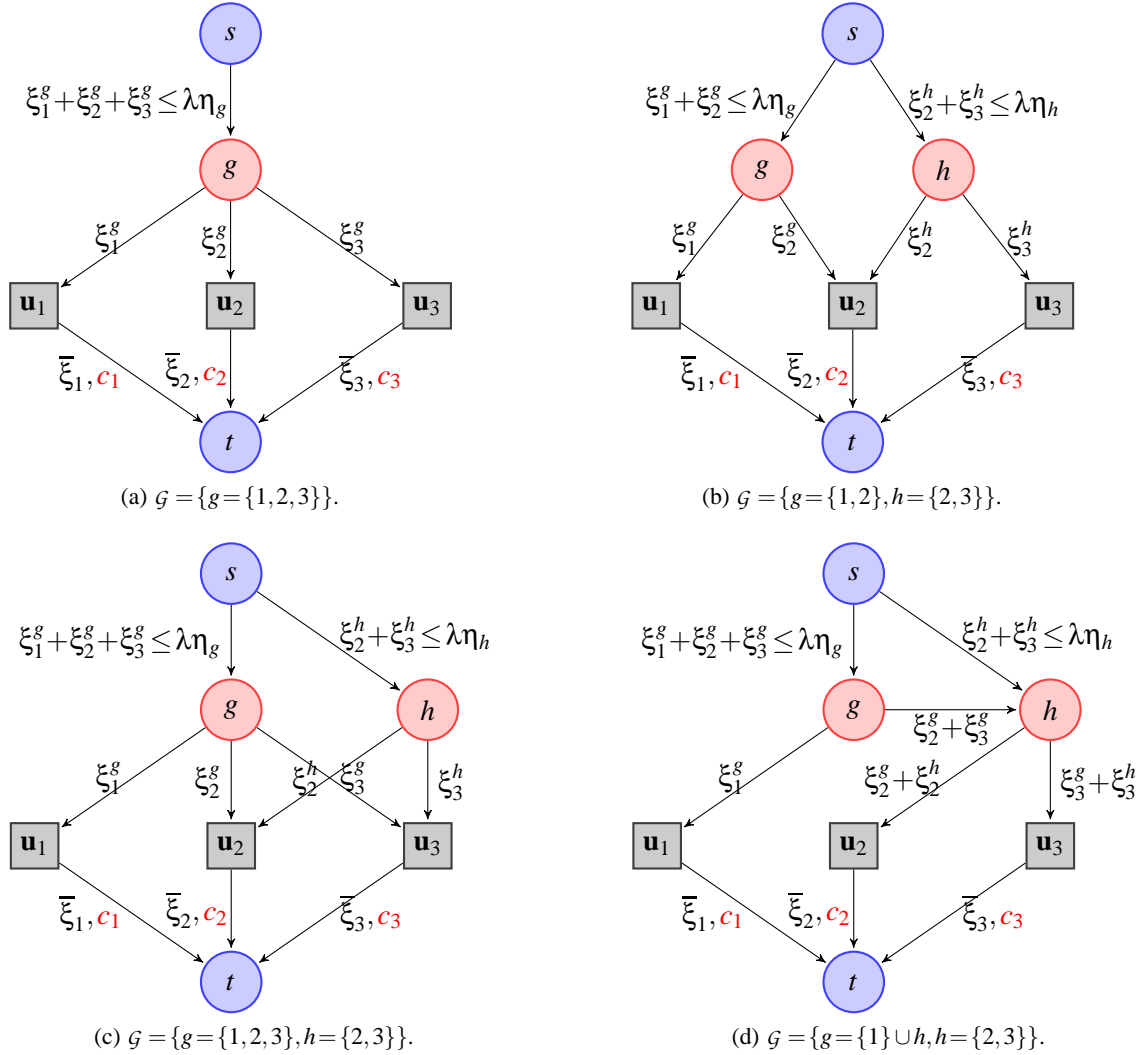


Figure 1: Graph representation of simple proximal problems with different group structures \mathcal{G} . The three indices 1, 2, 3 are represented as grey squares, and the groups g, h in \mathcal{G} as red discs. The source is linked to every group g, h with respective maximum capacity $\lambda\eta_g, \lambda\eta_h$ and zero cost. Each variable \mathbf{u}_j is linked to the sink t , with an infinite capacity, and with a cost $c_j \triangleq \frac{1}{2}(\mathbf{u}_j - \bar{\xi}_j)^2$. All other arcs in the graph have zero cost and infinite capacity. They represent inclusion relations in-between groups, and between groups and variables. The graphs (c) and (d) correspond to a special case of tree-structured hierarchy in the sense of Jenatton et al. (2010a). Their min-cost flow problems are equivalent.

Algorithm 1 Computation of the proximal operator for overlapping groups.

input $\mathbf{u} \in \mathbb{R}^p$, a set of groups \mathcal{G} , positive weights $(\eta_g)_{g \in \mathcal{G}}$, and λ (regularization parameter).

- 1: Build the initial graph $G_0 = (V_0, E_0, s, t)$ as explained in Section 3.4.
- 2: Compute the optimal flow: $\bar{\xi} \leftarrow \text{computeFlow}(V_0, E_0)$.
- 3: **Return:** $\mathbf{w} = \mathbf{u} - \bar{\xi}$ (optimal solution of the proximal problem).

Function $\text{computeFlow}(V = V_u \cup V_{gr}, E)$

- 1: Projection step: $\gamma \leftarrow \arg \min_{\gamma} \sum_{j \in V_u} \frac{1}{2} (\mathbf{u}_j - \gamma_j)^2$ s.t. $\sum_{j \in V_u} \gamma_j \leq \lambda \sum_{g \in V_{gr}} \eta_g$.
 - 2: For all nodes j in V_u , set γ_j to be the capacity of the arc (j, t) .
 - 3: Max-flow step: Update $(\bar{\xi}_j)_{j \in V_u}$ by computing a max-flow on the graph (V, E, s, t) .
 - 4: **if** $\exists j \in V_u$ s.t. $\bar{\xi}_j \neq \gamma_j$ **then**
 - 5: Denote by (s, V^+) and (V^-, t) the two disjoint subsets of (V, s, t) separated by the minimum (s, t) -cut of the graph, and remove the arcs between V^+ and V^- . Call E^+ and E^- the two remaining disjoint subsets of E corresponding to V^+ and V^- .
 - 6: $(\bar{\xi}_j)_{j \in V_u^+} \leftarrow \text{computeFlow}(V^+, E^+)$.
 - 7: $(\bar{\xi}_j)_{j \in V_u^-} \leftarrow \text{computeFlow}(V^-, E^-)$.
 - 8: **end if**
 - 9: **Return:** $(\bar{\xi}_j)_{j \in V_u}$.
-

The second step tries to construct a feasible flow $(\xi, \bar{\xi})$, satisfying additional capacity constraints equal to γ_j on arc (j, t) , and whose cost matches this lower bound; this latter problem can be cast as a max-flow problem (Goldberg and Tarjan, 1986). If such a flow exists, the algorithm returns $\bar{\xi} = \gamma$, the cost of the flow reaches the lower bound, and is therefore optimal. If such a flow does not exist, we have $\bar{\xi} \neq \gamma$, the lower bound is not achievable, and we build a minimum (s, t) -cut of the graph (Ford and Fulkerson, 1956) defining two disjoint sets of nodes V^+ and V^- ; V^+ is the part of the graph which is reachable from the source (for every node j in V^+ , there exists a non-saturated path from s to j), whereas all paths going from s to nodes in V^- are saturated. More details about these properties can be found at the beginning of Appendix B. At this point, it is possible to show that the value of the optimal min-cost flow on all arcs between V^+ and V^- is necessary zero. Thus, removing them yields an equivalent optimization problem, which can be decomposed into two independent problems of smaller sizes and solved recursively by the calls to $\text{computeFlow}(V^+, E^+)$ and $\text{computeFlow}(V^-, E^-)$. A formal proof of correctness of Algorithm 1 and further details are relegated to Appendix B.

The approach of Hochbaum and Hong (1995); Gallo et al. (1989) which recasts the quadratic min-cost flow problem as a parametric max-flow is guaranteed to have the same worst-case complexity as a single max-flow algorithm. However, we have experimentally observed a significant discrepancy between the worst case and empirical complexities for these flow problems, essentially because the empirical cost of each max-flow is significantly smaller than its theoretical cost. Despite the fact that the worst-case guarantees for our algorithm is weaker than theirs (up to a factor $|V|$), it is more adapted to the structure of our graphs and has proven to be much faster in our experiments (see Appendix D).¹⁰ Some implementation details are also crucial to the efficiency of the algorithm:

10. The best theoretical worst-case complexity of a max-flow is achieved by Goldberg and Tarjan (1986) and is $O(|V||E| \log(|V|^2/|E|))$. Our algorithm achieves the same worst-case complexity when the cuts are well balanced—

- **Exploiting maximal connected components:** When there exists no arc between two subsets of V , the solution can be obtained by solving two smaller optimization problems corresponding to the two disjoint subgraphs. It is indeed possible to process them independently to solve the global min-cost flow problem. To that effect, before calling the function `computeFlow(V, E)`, we look for maximal connected components $(V_1, E_1), \dots, (V_N, E_N)$ and call sequentially the procedure `computeFlow(V_i, E_i)` for i in $\{1, \dots, N\}$.
- **Efficient max-flow algorithm:** We have implemented the “push-relabel” algorithm of Goldberg and Tarjan (1986) to solve our max-flow problems, using classical heuristics that significantly speed it up in practice; see Goldberg and Tarjan (1986) and Cherkassky and Goldberg (1997). We use the so-called “highest-active vertex selection rule, global and gap heuristics” (Goldberg and Tarjan, 1986; Cherkassky and Goldberg, 1997), which has a worst-case complexity of $O(|V|^2|E|^{1/2})$ for a graph (V, E, s, t) . This algorithm leverages the concept of *pre-flow* that relaxes the definition of flow and allows vertices to have a positive excess.
- **Using flow warm-restarts:** The max-flow steps in our algorithm can be initialized with any valid pre-flow, enabling warm-restarts. This is also a key concept in the parametric max-flow algorithm of Gallo et al. (1989).
- **Improved projection step:** The first line of the procedure `computeFlow` can be replaced by $\gamma \leftarrow \arg \min_{\gamma} \sum_{j \in V_u} \frac{1}{2} (\mathbf{u}_j - \gamma_j)^2$ s.t. $\sum_{j \in V_u} \gamma_j \leq \lambda \sum_{g \in V_{gr}} \eta_g$ and $|\gamma_j| \leq \lambda \sum_{g \ni j} \eta_g$. The idea is to build a relaxation of Eq. (7) which is closer to the original problem than the one of Eq. (8), but that still can be solved in linear time. The structure of the graph will indeed not allow $\bar{\xi}_j$ to be greater than $\lambda \sum_{g \ni j} \eta_g$ after the max-flow step. This modified projection step can still be computed in linear time (Brucker, 1984), and leads to better performance.

3.5 Computation of the Dual Norm

The dual norm Ω^* of Ω , defined for any vector κ in \mathbb{R}^p by

$$\Omega^*(\kappa) \triangleq \max_{\Omega(\mathbf{z}) \leq 1} \mathbf{z}^\top \kappa,$$

is a key quantity to study sparsity-inducing regularizations in many respects. For instance, dual norms are central in working-set algorithms (Jenatton et al., 2009; Bach et al., 2011), and arise as well when proving theoretical estimation or prediction guarantees (Negahban et al., 2009).

In our context, we use it to monitor the convergence of the proximal method through a duality gap, hence defining a proper optimality criterion for problem (2). As a brief reminder, the duality gap of a minimization problem is defined as the difference between the primal and dual objective functions, evaluated for a feasible pair of primal/dual variables (see Section 5.5, Boyd and Vandenberghe, 2004). This gap serves as a certificate of (sub)optimality: if it is equal to zero, then the optimum is reached, and provided that strong duality holds, the converse is true as well (see Section 5.5, Boyd and Vandenberghe, 2004). A description of the algorithm we use in the experiments (Beck and Teboulle, 2009) along with the integration of the computation of the duality gap is given in Appendix C.

that is $|V^+| \approx |V^-| \approx |V|/2$, but we lose a factor $|V|$ when it is not the case. The practical speed of such algorithms is however significantly different than their theoretical worst-case complexities (see Boykov and Kolmogorov, 2004).

We now denote by f^* the Fenchel conjugate of f (Borwein and Lewis, 2006), defined by $f^*(\kappa) \triangleq \sup_{\mathbf{z}} [\mathbf{z}^\top \kappa - f(\mathbf{z})]$. The duality gap for problem (2) can be derived from standard Fenchel duality arguments (Borwein and Lewis, 2006) and it is equal to

$$f(\mathbf{w}) + \lambda \Omega(\mathbf{w}) + f^*(-\kappa) \text{ for } \mathbf{w}, \kappa \text{ in } \mathbb{R}^p \text{ with } \Omega^*(\kappa) \leq \lambda.$$

Therefore, evaluating the duality gap requires to compute efficiently Ω^* in order to find a feasible dual variable κ (the gap is otherwise equal to $+\infty$ and becomes non-informative). This is equivalent to solving another network flow problem, based on the following variational formulation:

$$\Omega^*(\kappa) = \min_{\xi \in \mathbb{R}^{p \times |\mathcal{G}|}} \tau \quad \text{s.t.} \quad \sum_{g \in \mathcal{G}} \xi^g = \kappa, \text{ and } \forall g \in \mathcal{G}, \|\xi^g\|_1 \leq \tau \eta_g \text{ with } \xi_j^g = 0 \text{ if } j \notin g. \quad (9)$$

In the network problem associated with (9), the capacities on the arcs (s, g) , $g \in \mathcal{G}$, are set to $\tau \eta_g$, and the capacities on the arcs (j, t) , $j \in \{1, \dots, p\}$, are fixed to κ_j . Solving problem (9) amounts to finding the smallest value of τ , such that there exists a flow saturating all the capacities κ_j on the arcs leading to the sink t . Equation (9) and Algorithm 2 are proven to be correct in Appendix B.

Algorithm 2 Computation of the dual norm.

input $\kappa \in \mathbb{R}^p$, a set of groups \mathcal{G} , positive weights $(\eta_g)_{g \in \mathcal{G}}$.

- 1: Build the initial graph $G_0 = (V_0, E_0, s, t)$ as explained in Section 3.5.
- 2: $\tau \leftarrow \text{dualNorm}(V_0, E_0)$.
- 3: **Return:** τ (value of the dual norm).

Function $\text{dualNorm}(V = V_u \cup V_{gr}, E)$

- 1: $\tau \leftarrow (\sum_{j \in V_u} \kappa_j) / (\sum_{g \in V_{gr}} \eta_g)$ and set the capacities of arcs (s, g) to $\tau \eta_g$ for all g in V_{gr} .
 - 2: Max-flow step: Update $(\bar{\xi}_j)_{j \in V_u}$ by computing a max-flow on the graph (V, E, s, t) .
 - 3: **if** $\exists j \in V_u$ s.t. $\bar{\xi}_j \neq \kappa_j$ **then**
 - 4: Define (V^+, E^+) and (V^-, E^-) as in Algorithm 1, and set $\tau \leftarrow \text{dualNorm}(V^-, E^-)$.
 - 5: **end if**
 - 6: **Return:** τ .
-

4. Optimization with Proximal Splitting Methods

We now present proximal splitting algorithms (see Combettes and Pesquet, 2008, 2010; Tomioka et al., 2011; Boyd et al., 2011, and references therein) for solving Eq. (2). Differentiability of f is not required here and the regularization function can either be a sum of ℓ_2 - or ℓ_∞ -norms. However, we assume that:

- (A) either f can be written $f(\mathbf{w}) = \sum_{i=1}^n \tilde{f}_i(\mathbf{w})$, where the functions \tilde{f}_i are such that $\text{prox}_{\gamma \tilde{f}_i}$ can be obtained in closed form for all $\gamma > 0$ and all i —that is, for all \mathbf{u} in \mathbb{R}^m , the following problems admit closed form solutions: $\min_{\mathbf{v} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_2^2 + \gamma \tilde{f}_i(\mathbf{v})$.
- (B) or f can be written $f(\mathbf{w}) = \tilde{f}(\mathbf{X}\mathbf{w})$ for all \mathbf{w} in \mathbb{R}^p , where \mathbf{X} in $\mathbb{R}^{n \times p}$ is a design matrix, and one knows how to efficiently compute $\text{prox}_{\gamma \tilde{f}}$ for all $\gamma > 0$.

It is easy to show that this condition is satisfied for the square and hinge loss functions, making it possible to build linear SVMs with a structured sparse regularization. These assumptions are not the same as the ones of Section 3, and the scope of the problems addressed is therefore slightly different. Proximal splitting methods seem indeed to offer more flexibility regarding the regularization function, since they can deal with sums of ℓ_2 -norms.¹¹ However, proximal gradient methods, as presented in Section 3, enjoy a few advantages over proximal splitting methods, namely: automatic parameter tuning with line-search schemes (Nesterov, 2007), known convergence rates (Nesterov, 2007; Beck and Teboulle, 2009), and ability to provide sparse solutions (approximate solutions obtained with proximal splitting methods often have small values, but not “true” zeros).

4.1 Algorithms

We consider a class of algorithms which leverage the concept of variable splitting (see Combettes and Pesquet, 2010; Bertsekas and Tsitsiklis, 1989; Tomioka et al., 2011). The key is to introduce additional variables \mathbf{z}^g in $\mathbb{R}^{|g|}$, one for every group g in \mathcal{G} , and equivalently reformulate Eq. (2) as

$$\min_{\substack{\mathbf{w} \in \mathbb{R}^p \\ \mathbf{z}^g \in \mathbb{R}^{|g|} \text{ for } g \in \mathcal{G}}} f(\mathbf{w}) + \lambda \sum_{g \in \mathcal{G}} \eta_g \|\mathbf{z}^g\| \quad \text{s.t. } \forall g \in \mathcal{G}, \mathbf{z}^g = \mathbf{w}_g, \quad (10)$$

The issue of overlapping groups is removed, but new constraints are added, and as in Section 3, the method introduces additional variables which induce a memory cost of $O(\sum_{g \in \mathcal{G}} |g|)$.

To solve this problem, it is possible to use the so-called alternating direction method of multipliers (ADMM) (see Combettes and Pesquet, 2010; Bertsekas and Tsitsiklis, 1989; Tomioka et al., 2011; Boyd et al., 2011).¹² It introduces dual variables \mathbf{v}^g in $\mathbb{R}^{|g|}$ for all g in \mathcal{G} , and defines the augmented Lagrangian:

$$\mathcal{L}(\mathbf{w}, (\mathbf{z}^g)_{g \in \mathcal{G}}, (\mathbf{v}^g)_{g \in \mathcal{G}}) \triangleq f(\mathbf{w}) + \sum_{g \in \mathcal{G}} [\lambda \eta_g \|\mathbf{z}^g\| + \mathbf{v}^{g\top} (\mathbf{z}^g - \mathbf{w}_g) + \frac{\gamma}{2} \|\mathbf{z}^g - \mathbf{w}_g\|_2^2],$$

where $\gamma > 0$ is a parameter. It is easy to show that solving Eq. (10) amounts to finding a saddle-point of the augmented Lagrangian.¹³ The ADMM algorithm finds such a saddle-point by iterating between the minimization of \mathcal{L} with respect to each primal variable, keeping the other ones fixed, and gradient ascent steps with respect to the dual variables. More precisely, it can be summarized as:

1. Minimize \mathcal{L} with respect to \mathbf{w} , keeping the other variables fixed.

-
11. We are not aware of any efficient algorithm providing the exact solution of the proximal operator associated to a sum of ℓ_2 -norms, which would be necessary for using (accelerated) proximal gradient methods. An iterative algorithm could possibly be used to compute it approximately (e.g., see Jenatton et al., 2010a, 2011), but such a procedure would be computationally expensive and would require to be able to deal with approximate computations of the proximal operators (e.g., see Combettes and Pesquet, 2010; Schmidt et al., 2011, and discussions therein). We have chosen not to consider this possibility in this paper.
 12. This method is used by Sprechmann et al. (2010) for computing the proximal operator associated to hierarchical norms, and independently in the same context as ours by Boyd et al. (2011) and Qin and Goldfarb (2011).
 13. The augmented Lagrangian is in fact the classical Lagrangian (see Boyd and Vandenberghe, 2004) of the following optimization problem which is equivalent to Eq. (10):

$$\min_{\mathbf{w} \in \mathbb{R}^p, (\mathbf{z}^g \in \mathbb{R}^{|g|})_{g \in \mathcal{G}}} f(\mathbf{w}) + \lambda \sum_{g \in \mathcal{G}} \eta_g \|\mathbf{z}^g\| + \frac{\gamma}{2} \|\mathbf{z}^g - \mathbf{w}_g\|_2^2 \quad \text{s.t. } \forall g \in \mathcal{G}, \mathbf{z}^g = \mathbf{w}_g.$$

2. Minimize \mathcal{L} with respect to the \mathbf{z}^g 's, keeping the other variables fixed. The solution can be obtained in closed form: for all g in G , $\mathbf{z}^g \leftarrow \text{prox}_{\frac{\lambda \eta_g}{\gamma} \|\cdot\|} [\mathbf{w}_g - \frac{1}{\gamma} \mathbf{v}^g]$.
3. Take a gradient ascent step on \mathcal{L} with respect to the \mathbf{v}^g 's: $\mathbf{v}^g \leftarrow \mathbf{v}^g + \gamma(\mathbf{z}^g - \mathbf{w}_g)$.
4. Go back to step 1.

Such a procedure is guaranteed to converge to the desired solution for all value of $\gamma > 0$ (however, tuning γ can greatly influence the convergence speed), but solving efficiently step 1 can be difficult. To cope with this issue, we propose two variations exploiting assumptions **(A)** and **(B)**.

4.1.1 SPLITTING THE LOSS FUNCTION f

We assume condition **(A)**—that is, we have $f(\mathbf{w}) = \sum_{i=1}^n \tilde{f}_i(\mathbf{w})$. For example, when f is the square loss function $f(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$, where \mathbf{X} in $\mathbb{R}^{n \times p}$ is a design matrix and \mathbf{y} is in \mathbb{R}^n , we would define for all i in $\{1, \dots, n\}$ the functions $\tilde{f}_i: \mathbb{R}^p \rightarrow \mathbb{R}$ such that $\tilde{f}_i(\mathbf{w}) \triangleq \frac{1}{2} (\mathbf{y}_i - \mathbf{x}_i^\top \mathbf{w})^2$, where \mathbf{x}_i is the i -th row of \mathbf{X} .

We now introduce new variables \mathbf{v}^i in \mathbb{R}^p for $i = 1, \dots, n$, and replace $f(\mathbf{w})$ in Eq. (10) by $\sum_{i=1}^n \tilde{f}_i(\mathbf{v}^i)$, with the additional constraints that $\mathbf{v}^i = \mathbf{w}$. The resulting equivalent optimization problem can now be tackled using the ADMM algorithm, following the same methodology presented above. It is easy to show that every step can be obtained efficiently, as long as one knows how to compute the proximal operator associated to the functions \tilde{f}_i in closed form. This is in fact the case for the square and hinge loss functions, where n is the number of training points. The main problem of this strategy is the possible high memory usage it requires when n is large.

4.1.2 DEALING WITH THE DESIGN MATRIX

If we assume condition **(B)**, another possibility consists of introducing a new variable \mathbf{v} in \mathbb{R}^n , such that one can replace the function $f(\mathbf{w}) = \tilde{f}(\mathbf{X}\mathbf{w})$ by $\tilde{f}(\mathbf{v})$ in Eq. (10) with the additional constraint $\mathbf{v} = \mathbf{X}\mathbf{w}$. Using directly the ADMM algorithm to solve the corresponding problem implies adding a term $\kappa^\top (\mathbf{v} - \mathbf{X}\mathbf{w}) + \frac{\gamma}{2} \|\mathbf{v} - \mathbf{X}\mathbf{w}\|_2^2$ to the augmented Lagrangian \mathcal{L} , where κ is a new dual variable. The minimization of \mathcal{L} with respect to \mathbf{v} is now obtained by $\mathbf{v} \leftarrow \text{prox}_{\frac{1}{\gamma} \tilde{f}} [\mathbf{X}\mathbf{w} - \kappa]$, which is easy to compute according to **(B)**. However, the design matrix \mathbf{X} in the quadratic term makes the minimization of \mathcal{L} with respect to \mathbf{w} more difficult. To overcome this issue, we adopt a strategy presented by Zhang et al. (2011), which replaces at iteration k the quadratic term $\frac{\gamma}{2} \|\mathbf{v} - \mathbf{X}\mathbf{w}\|_2^2$ in the augmented Lagrangian by an additional proximity term: $\frac{\gamma}{2} \|\mathbf{v} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{w}^k\|_{\mathbf{Q}}^2$, where \mathbf{w}^k is the current estimate of \mathbf{w} , and $\|\mathbf{w} - \mathbf{w}^k\|_{\mathbf{Q}}^2 = (\mathbf{w} - \mathbf{w}^k)^\top \mathbf{Q} (\mathbf{w} - \mathbf{w}^k)$, where \mathbf{Q} is a symmetric positive definite matrix. By choosing $\mathbf{Q} \triangleq \delta \mathbf{I} - \mathbf{X}^\top \mathbf{X}$, with δ large enough, minimizing \mathcal{L} with respect to \mathbf{w} becomes simple, while convergence to the solution is still ensured. More details can be found in Zhang et al. (2011).

5. Applications and Experiments

In this section, we present various experiments demonstrating the applicability and the benefits of our methods for solving large-scale sparse and structured regularized problems.

5.1 Speed Benchmark

We consider a structured sparse decomposition problem with overlapping groups of ℓ_∞ -norms, and compare the proximal gradient algorithm FISTA (Beck and Teboulle, 2009) with our proximal operator presented in Section 3 (referred to as ProxFlow), two variants of proximal splitting methods, (ADMM) and (Lin-ADMM) respectively presented in Section 4.1.1 and 4.1.2, and two generic optimization techniques, namely a subgradient descent (SG) and an interior point method,¹⁴ on a regularized linear regression problem. SG, ProxFlow, ADMM and Lin-ADMM are implemented in C++.¹⁵ Experiments are run on a single-core 2.8 GHz CPU. We consider a design matrix \mathbf{X} in $\mathbb{R}^{n \times p}$ built from overcomplete dictionaries of discrete cosine transforms (DCT), which are naturally organized on one- or two-dimensional grids and display local correlations. The following families of groups \mathcal{G} using this spatial information are thus considered: (1) every contiguous sequence of length 3 for the one-dimensional case, and (2) every 3×3 -square in the two-dimensional setting. We generate vectors \mathbf{y} in \mathbb{R}^n according to the linear model $\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, 0.01\|\mathbf{X}\mathbf{w}_0\|_2^2)$. The vector \mathbf{w}_0 has about 20% percent nonzero components, randomly selected, while respecting the structure of \mathcal{G} , and uniformly generated in $[-1, 1]$.

In our experiments, the regularization parameter λ is chosen to achieve the same level of sparsity (20%). For SG, ADMM and Lin-ADMM, some parameters are optimized to provide the lowest value of the objective function after 1000 iterations of the respective algorithms. For SG, we take the step size to be equal to $a/(k+b)$, where k is the iteration number, and (a, b) are the pair of parameters selected in $\{10^{-3}, \dots, 10\} \times \{10^2, 10^3, 10^4\}$. Note that a step size of the form $a/(\sqrt{t}+b)$ is also commonly used in subgradient descent algorithms. In the context of hierarchical norms, both choices have led to similar results (Jenatton et al., 2011). The parameter γ for ADMM is selected in $\{10^{-2}, \dots, 10^2\}$. The parameters (γ, δ) for Lin-ADMM are selected in $\{10^{-2}, \dots, 10^2\} \times \{10^{-1}, \dots, 10^8\}$. For interior point methods, since problem (2) can be cast either as a quadratic (QP) or as a conic program (CP), we show in Figure 2 the results for both formulations. On three problems of different sizes, with $(n, p) \in \{(100, 10^3), (1024, 10^4), (1024, 10^5)\}$, our algorithms ProxFlow, ADMM and Lin-ADMM compare favorably with the other methods, (see Figure 2), except for ADMM in the large-scale setting which yields an objective function value similar to that of SG after 10^4 seconds. Among ProxFlow, ADMM and Lin-ADMM, ProxFlow is consistently better than Lin-ADMM, which is itself better than ADMM. Note that for the small scale problem, the performance of ProxFlow and Lin-ADMM is similar. In addition, note that QP, CP, SG, ADMM and Lin-ADMM do not obtain sparse solutions, whereas ProxFlow does.¹⁶

5.2 Wavelet Denoising with Structured Sparsity

We now illustrate the results of Section 3, where a single large-scale proximal operator ($p \approx 250\,000$) associated to a sum of ℓ_∞ -norms has to be computed. We choose an image denoising task with an orthonormal wavelet basis, following an experiment similar to one proposed in Jenatton et al. (2011). Specifically, we consider the following formulation

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \Omega(\mathbf{w}), \quad (11)$$

14. In our simulations, we use the commercial software Mosek, <http://www.mosek.com/>

15. Our implementation of ProxFlow is available at <http://www.di.ens.fr/willow/SPAMS/>.

16. To reduce the computational cost of this experiment, the curves reported are the results of one single run. Similar types of experiments with several runs have shown very small variability (Bach et al., 2011).

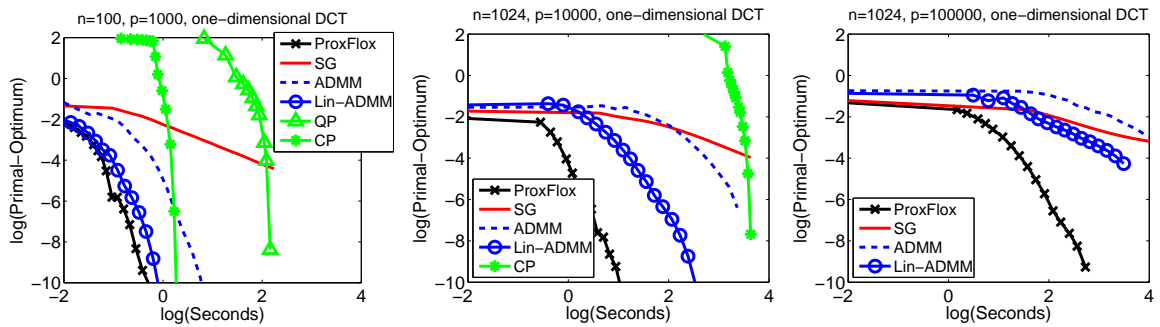


Figure 2: Speed comparisons: distance to the optimal primal value versus CPU time (log-log scale). Due to the computational burden, QP and CP could not be run on every problem.

where \mathbf{y} in \mathbb{R}^p is a noisy input image, \mathbf{w} represents wavelets coefficients, \mathbf{X} in $\mathbb{R}^{p \times p}$ is an orthonormal wavelet basis, $\mathbf{X}\mathbf{w}$ is the estimate of the denoised image, and Ω is a sparsity-inducing norm. Since here the basis is orthonormal, solving the decomposition problem boils down to computing $\mathbf{w}^* = \text{prox}_{\lambda\Omega}[\mathbf{X}^\top \mathbf{y}]$. This makes of Algorithm 1 a good candidate to solve it when Ω is a sum of ℓ_∞ -norms. We compare the following candidates for the sparsity-inducing norms Ω :

- the ℓ_1 -norm, leading to the wavelet soft-thresholding of Donoho and Johnstone (1995).
- a sum of ℓ_∞ -norms with a hierarchical group structure adapted to the wavelet coefficients, as proposed in Jenatton et al. (2011). Considering a natural quad-tree for wavelet coefficients (see Mallat, 1999), this norm takes the form of Eq. (3) with one group per wavelet coefficient that contains the coefficient and all its descendants in the tree. We call this norm Ω_{tree} .
- a sum of ℓ_∞ -norms with overlapping groups representing 2×2 spatial neighborhoods in the wavelet domain. This regularization encourages neighboring wavelet coefficients to be set to zero together, which was also exploited in the past in block-thresholding approaches for wavelet denoising (Cai, 1999). We call this norm Ω_{grid} .

We consider Daubechies3 wavelets (see Mallat, 1999) for the matrix \mathbf{X} , use 12 classical standard test images,¹⁷ and generate noisy versions of them corrupted by a white Gaussian noise of variance σ^2 . For each image, we test several values of $\lambda = 2^i \sigma \sqrt{\log p}$, with i taken in the range $\{-15, -14, \dots, 15\}$. We then keep the parameter λ giving the best reconstruction error on average on the 12 images. The factor $\sigma \sqrt{\log p}$ is a classical heuristic for choosing a reasonable regularization parameter (see Mallat, 1999). We provide reconstruction results in terms of PSNR in Table 1.¹⁸ Unlike Jenatton et al. (2011), who set all the weights η_g in Ω equal to one, we tried exponential weights of the form $\eta_g = \rho^k$, with k being the depth of the group in the wavelet tree, and ρ is taken in $\{0.25, 0.5, 1, 2, 4\}$. As for λ , the value providing the best reconstruction is kept. The wavelet transforms in our experiments are computed with the matlabPyrTools software.¹⁹ Interestingly, we observe in Table 1 that the results obtained with Ω_{grid} are significantly better than those obtained

17. These images are used in classical image denoising benchmarks. See Mairal et al. (2009).

18. Denoting by MSE the mean-squared-error for images whose intensities are between 0 and 255, the PSNR is defined as $\text{PSNR} = 10 \log_{10}(255^2 / \text{MSE})$ and is measured in dB. A gain of 1dB reduces the MSE by approximately 20%.

19. <http://www.cns.nyu.edu/~eero/steerpyr/>.

σ	PSNR			IPSNR vs. ℓ_1		
	ℓ_1	Ω_{tree}	Ω_{grid}	ℓ_1	Ω_{tree}	Ω_{grid}
5	35.67	35.98	36.15	$0.00 \pm .0$	$0.31 \pm .18$	$0.48 \pm .25$
10	31.00	31.60	31.88	$0.00 \pm .0$	$0.61 \pm .28$	$0.88 \pm .28$
25	25.68	26.77	27.07	$0.00 \pm .0$	$1.09 \pm .32$	$1.38 \pm .26$
50	22.37	23.84	24.06	$0.00 \pm .0$	$1.47 \pm .34$	$1.68 \pm .41$
100	19.64	21.49	21.56	$0.00 \pm .0$	$1.85 \pm .28$	$1.92 \pm .29$

Table 1: PSNR measured for the denoising of 12 standard images when the regularization function is the ℓ_1 -norm, the tree-structured norm Ω_{tree} , and the structured norm Ω_{grid} , and improvement in PSNR compared to the ℓ_1 -norm (IPSNR). Best results for each level of noise and each wavelet type are in bold. The reported values are averaged over 5 runs with different noise realizations.

with Ω_{tree} , meaning that encouraging spatial consistency in wavelet coefficients is more effective than using a hierarchical coding. We also note that our approach is relatively fast, despite the high dimension of the problem. Solving exactly the proximal problem with Ω_{grid} for an image with $p = 512 \times 512 = 262144$ pixels (and therefore approximately the same number of groups) takes approximately $\approx 4 - 6$ seconds on a single core of a 3.07GHz CPU.

5.3 CUR-like Matrix Factorization

In this experiment, we show how our tools can be used to perform the so-called CUR matrix decomposition (Mahoney and Drineas, 2009). It consists of a low-rank approximation of a data matrix \mathbf{X} in $\mathbb{R}^{n \times p}$ in the form of a product of three matrices—that is, $\mathbf{X} \approx \mathbf{C}\mathbf{U}\mathbf{R}$. The particularity of the CUR decomposition lies in the fact that the matrices $\mathbf{C} \in \mathbb{R}^{n \times c}$ and $\mathbf{R} \in \mathbb{R}^{r \times p}$ are constrained to be respectively a subset of c columns and r rows of the original matrix \mathbf{X} . The third matrix $\mathbf{U} \in \mathbb{R}^{c \times r}$ is then given by $\mathbf{C}^+ \mathbf{X} \mathbf{R}^+$, where \mathbf{A}^+ denotes a Moore-Penrose generalized inverse of the matrix \mathbf{A} (Horn and Johnson, 1990). Such a matrix factorization is particularly appealing when the interpretability of the results matters (Mahoney and Drineas, 2009). For instance, when studying gene-expression datasets, it is easier to gain insight from the selection of actual patients and genes, rather than from linear combinations of them.

In Mahoney and Drineas (2009), CUR decompositions are computed by a sampling procedure based on the singular value decomposition of \mathbf{X} . In a recent work, Bien et al. (2010) have shown that *partial* CUR decompositions, i.e., the selection of either rows or columns of \mathbf{X} , can be obtained by solving a convex program with a group-Lasso penalty. We propose to extend this approach to the simultaneous selection of both rows and columns of \mathbf{X} , with the following convex problem:

$$\min_{\mathbf{W} \in \mathbb{R}^{p \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{X}\|_F^2 + \lambda_{\text{row}} \sum_{i=1}^n \|\mathbf{W}^i\|_\infty + \lambda_{\text{col}} \sum_{j=1}^p \|\mathbf{W}_j\|_\infty. \quad (12)$$

In this formulation, the two sparsity-inducing penalties controlled by the parameters λ_{row} and λ_{col} set to zero some entire rows and columns of the solutions of problem (12). Now, let us denote by $\mathbf{W}_{\mathbf{I}\mathbf{J}}$ in $\mathbb{R}^{|\mathbf{I}| \times |\mathbf{J}|}$ the submatrix of \mathbf{W} reduced to its nonzero rows and columns, respectively indexed by $\mathbf{I} \subseteq \{1, \dots, p\}$ and $\mathbf{J} \subseteq \{1, \dots, n\}$. We can then readily identify the three components of the CUR decomposition of \mathbf{X} , namely

$$\mathbf{X}\mathbf{W}\mathbf{X} = \mathbf{C}\mathbf{W}_{\mathbf{I}\mathbf{J}}\mathbf{R} \approx \mathbf{X}.$$

Problem (12) has a smooth convex data-fitting term and brings into play a sparsity-inducing norm with overlapping groups of variables (the rows and the columns of \mathbf{W}). As a result, it is a particular instance of problem (2) that can therefore be handled with the optimization tools introduced in this paper. We now compare the performance of the sampling procedure from Mahoney and Drineas (2009) with our proposed sparsity-based approach. To this end, we consider the four gene-expression datasets `9_Tumors`, `Brain_Tumors1`, `Leukemia1` and `SRBCT`, with respective dimensions $(n, p) \in \{(60, 5727), (90, 5921), (72, 5328), (83, 2309)\}$.²⁰ In the sequel, the matrix \mathbf{X} is normalized to have unit Frobenius-norm while each of its columns is centered. To begin with, we run our approach²¹ over a grid of values for λ_{row} and λ_{col} in order to obtain solutions with different sparsity levels, i.e., ranging from $|\mathbf{I}| = p$ and $|\mathbf{J}| = n$ down to $|\mathbf{I}| = |\mathbf{J}| = 0$. For each pair of values $[|\mathbf{I}|, |\mathbf{J}|]$, we then apply the sampling procedure from Mahoney and Drineas (2009). Finally, the variance explained by the CUR decompositions is reported in Figure 3 for both methods. Since the sampling approach involves some randomness, we show the average and standard deviation of the results based on five initializations. The conclusions we can draw from the experiments match the ones already reported in Bien et al. (2010) for the partial CUR decomposition. We can indeed see that both schemes perform similarly. However, our approach has the advantage not to be randomized, which can be less disconcerting in the practical perspective of analyzing a single run of the algorithm. It is finally worth being mentioned that the convex approach we develop here is flexible and can be extended in different ways. For instance, we can imagine to add further low-rank/sparsity constraints on \mathbf{W} thanks to sparsity-promoting convex regularizations.

5.4 Background Subtraction

Following Cehver et al. (2008); Huang et al. (2009), we consider a background subtraction task. Given a sequence of frames from a fixed camera, we try to segment out foreground objects in a new image. If we denote by $\mathbf{y} \in \mathbb{R}^n$ this image composed of n pixels, we model \mathbf{y} as a sparse linear combination of p other images $\mathbf{X} \in \mathbb{R}^{n \times p}$, plus an error term \mathbf{e} in \mathbb{R}^n , i.e., $\mathbf{y} \approx \mathbf{X}\mathbf{w} + \mathbf{e}$ for some sparse vector \mathbf{w} in \mathbb{R}^p . This approach is reminiscent of Wright et al. (2009a) in the context of face recognition, where \mathbf{e} is further made sparse to deal with small occlusions. The term $\mathbf{X}\mathbf{w}$ accounts for *background* parts present in both \mathbf{y} and \mathbf{X} , while \mathbf{e} contains specific, or *foreground*, objects in \mathbf{y} . The resulting optimization problem is given by

$$\min_{\mathbf{w} \in \mathbb{R}^p, \mathbf{e} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{e}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \{\|\mathbf{e}\|_1 + \Omega(\mathbf{e})\}, \text{ with } \lambda_1, \lambda_2 \geq 0. \quad (13)$$

In this formulation, the only ℓ_1 -norm penalty does not take into account the fact that neighboring pixels in \mathbf{y} are likely to share the same label (background or foreground), which may lead to scattered pieces of foreground and background regions (Figure 4). We therefore put an additional structured regularization term Ω on \mathbf{e} , where the groups in \mathcal{G} are all the overlapping 3×3 -squares on the image. For the sake of comparison, we also consider the regularization $\tilde{\Omega}$ where the groups are *non-overlapping* 3×3 -squares.

This optimization problem can be viewed as an instance of problem (2), with the particular design matrix $[\mathbf{X}, \mathbf{I}]$ in $\mathbb{R}^{n \times (p+n)}$, defined as the columnwise concatenation of \mathbf{X} and the identity

20. The datasets are freely available at <http://www.gems-system.org/>.

21. More precisely, since the penalties in problem (12) shrink the coefficients of \mathbf{W} , we follow a two-step procedure: We first run our approach to determine the sets of nonzero rows and columns, and then compute $\mathbf{W}_{\mathbf{IJ}} = \mathbf{C}^+ \mathbf{X}\mathbf{R}^+$.

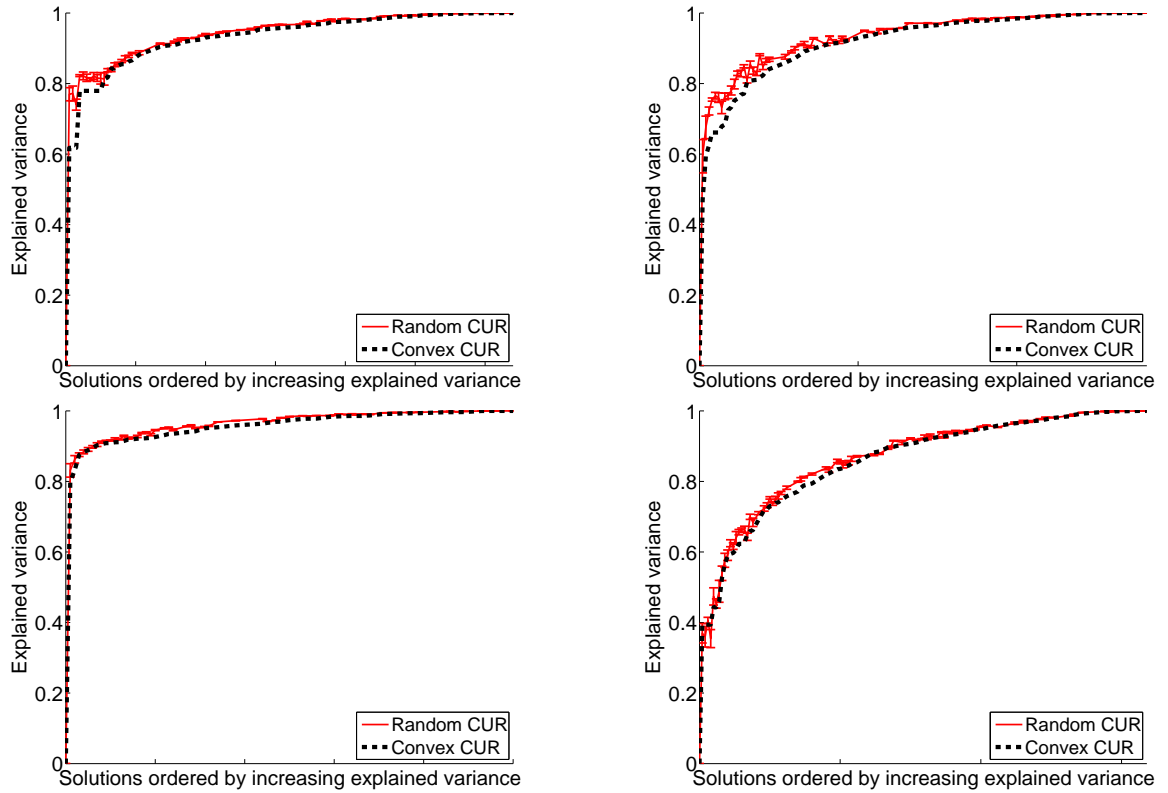


Figure 3: Explained variance of the CUR decompositions obtained for our sparsity-based approach and the sampling scheme from Mahoney and Drineas (2009). For the latter, we report the average and standard deviation of the results based on five initializations. From left to right and top to bottom, the curves correspond to the datasets 9_Tumors, Brain_Tumors1, Leukemia1 and SRBCT.

matrix. As a result, we could directly apply the same procedure as the one used in the other experiments. Instead, we further exploit the specific structure of problem (13): Notice that for a fixed vector \mathbf{e} , the optimization with respect to \mathbf{w} is a standard Lasso problem (with the vector of observations $\mathbf{y} - \mathbf{e}$),²² while for \mathbf{w} fixed, we simply have a proximal problem associated to the sum of Ω and the ℓ_1 -norm. Alternating between these two simple and computationally inexpensive steps, i.e., optimizing with respect to one variable while keeping the other one fixed, is guaranteed to converge to a solution of (13).²³ In our simulations, this alternating scheme has led to a significant speed-up compared to the general procedure.

A dataset with hand-segmented images is used to illustrate the effect of Ω .²⁴ For simplicity, we use a single regularization parameter, i.e., $\lambda_1 = \lambda_2$, chosen to maximize the number of pixels matching the ground truth. We consider $p = 200$ images with $n = 57600$ pixels (i.e., a resolution

22. Since successive frames might not change much, the columns of \mathbf{X} exhibit strong correlations. Consequently, we use the LARS algorithm (Efron et al., 2004) whose complexity is independent of the level of correlation in \mathbf{X} .

23. More precisely, the convergence is guaranteed since the non-smooth part in (13) is *separable* with respect to \mathbf{w} and \mathbf{e} (Tseng, 2001). The result from Bertsekas (1999) may also be applied here, after reformulating (13) as a smooth convex problem under separable conic constraints.

24. <http://research.microsoft.com/en-us/um/people/jckrumm/wallflower/testimages.htm>

of 120×160 , times 3 for the RGB channels). As shown in Figure 4, adding Ω improves the background subtraction results for the two tested images, by removing the scattered artifacts due to the lack of structural constraints of the ℓ_1 -norm, which encodes neither spatial nor color consistency. The group sparsity regularization $\tilde{\Omega}$ also improves upon the ℓ_1 -norm but introduces block-artefacts corresponding to the non-overlapping group structure.

5.5 Topographic Dictionary Learning

Let us consider a set $\mathbf{Y} = [\mathbf{y}^1, \dots, \mathbf{y}^n]$ in $\mathbb{R}^{m \times n}$ of n signals of dimension m . The problem of dictionary learning, originally introduced by Olshausen and Field (1996), is a matrix factorization problem which aims at representing these signals as linear combinations of *dictionary elements* that are the columns of a matrix $\mathbf{X} = [\mathbf{x}^1, \dots, \mathbf{x}^p]$ in $\mathbb{R}^{m \times p}$. More precisely, the dictionary \mathbf{X} is *learned* along with a matrix of *decomposition coefficients* $\mathbf{W} = [\mathbf{w}^1, \dots, \mathbf{w}^n]$ in $\mathbb{R}^{p \times n}$, so that $\mathbf{y}^i \approx \mathbf{X}\mathbf{w}^i$ for every signal \mathbf{y}^i . Typically, n is large compared to m and p . In this experiment, we consider for instance a database of $n = 100\,000$ natural image patches of size $m = 12 \times 12$ pixels, for dictionaries of size $p = 400$. Adapting the dictionary to specific data has proven to be useful in many applications, including image restoration (Elad and Aharon, 2006; Mairal et al., 2009), learning image features in computer vision (Kavukcuoglu et al., 2009). The resulting optimization problem we are interested in can be written

$$\min_{\mathbf{X} \in \mathcal{C}, \mathbf{W} \in \mathbb{R}^{p \times n}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{y}^i - \mathbf{X}\mathbf{w}^i\|_2^2 + \lambda \Omega(\mathbf{w}^i), \quad (14)$$

where \mathcal{C} is a convex set of matrices in $\mathbb{R}^{m \times p}$ whose columns have ℓ_2 -norms less than or equal to one,²⁵ λ is a regularization parameter and Ω is a sparsity-inducing norm. When Ω is the ℓ_1 -norm, we obtain a classical formulation, which is known to produce dictionary elements that are reminiscent of Gabor-like functions, when the columns of \mathbf{Y} are whitened natural image patches (Olshausen and Field, 1996).

Another line of research tries to put a structure on decomposition coefficients instead of considering them as independent. Jenatton et al. (2010a, 2011) have for instance embedded dictionary elements into a tree, by using a hierarchical norm (Zhao et al., 2009) for Ω . This model encodes a rule saying that a dictionary element can be used in the decomposition of a signal only if its ancestors in the tree are used as well. In the related context of independent component analysis (ICA), Hyvärinen et al. (2001) have arranged independent components (corresponding to dictionary elements) on a two-dimensional grid, and have modelled spatial dependencies between them. When learned on whitened natural image patches, this model exhibits ‘‘Gabor-like’’ functions which are smoothly organized on the grid, which the authors call a topographic map. As shown by Kavukcuoglu et al. (2009), such a result can be reproduced with a dictionary learning formulation, using a structured norm for Ω . Following their formulation, we organize the p dictionary elements on a $\sqrt{p} \times \sqrt{p}$ grid, and consider p overlapping groups that are 3×3 or 4×4 spatial neighborhoods on the grid (to avoid boundary effects, we assume the grid to be cyclic). We define Ω as a sum of ℓ_2 -norms over these groups, since the ℓ_∞ -norm has proven to be less adapted for this task. Another formulation achieving a similar effect was also proposed by Garrigues and Olshausen (2010) in the context of sparse coding with a probabilistic model.

25. Since the quadratic term in Eq. (14) is invariant by multiplying \mathbf{X} by a scalar and \mathbf{W} by its inverse, constraining the norm of \mathbf{X} has proven to be necessary in practice to prevent it from being arbitrarily large.

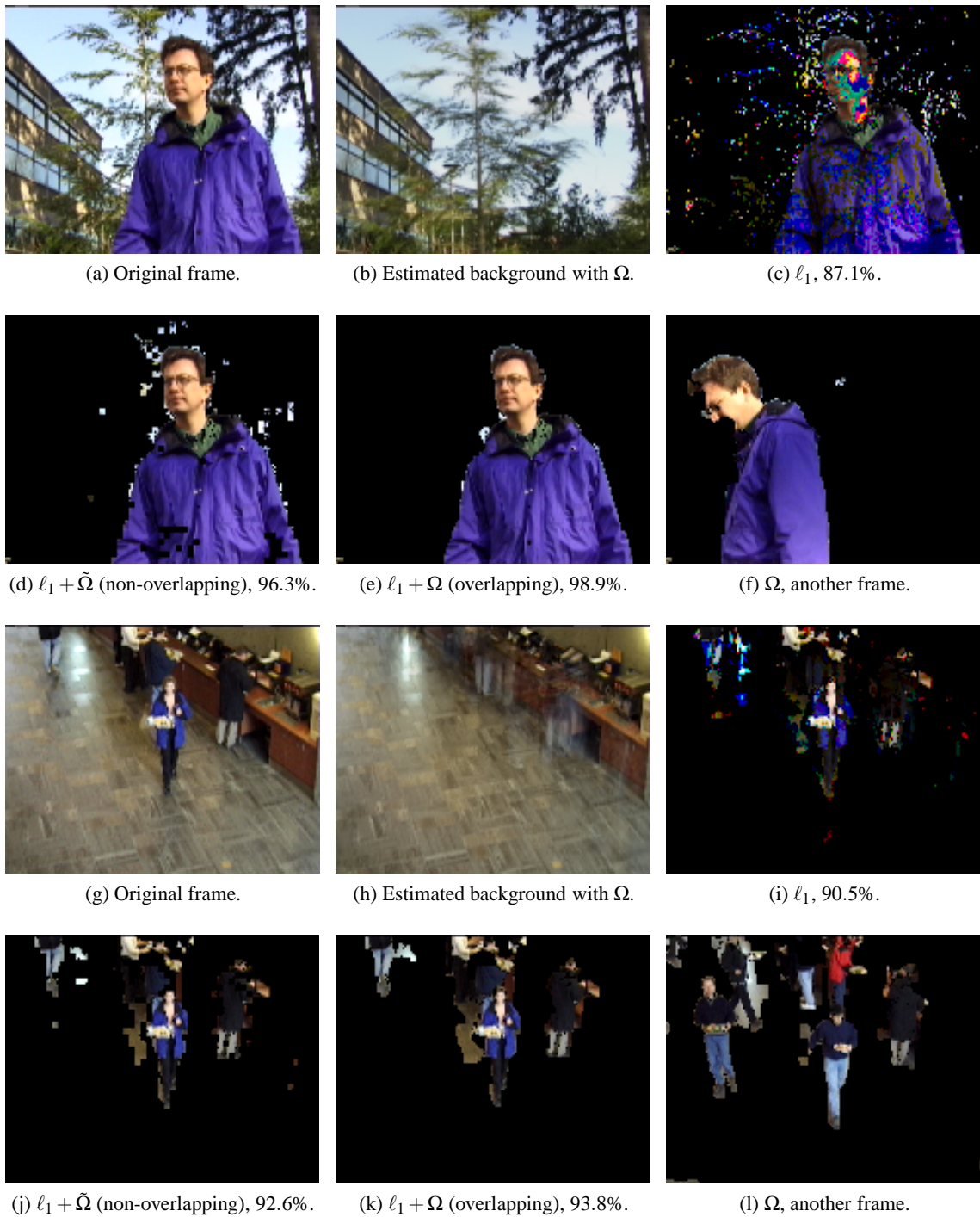


Figure 4: Background subtraction results. For two videos, we present the original image \mathbf{y} , the estimated background (i.e., $\mathbf{X}\mathbf{w}$) reconstructed by our method, and the foreground (i.e., the sparsity pattern of \mathbf{e} as a mask on the original image) detected with ℓ_1 , $\ell_1 + \tilde{\Omega}$ (non-overlapping groups) and with $\ell_1 + \Omega$. Figures (f) and (l) present another foreground found with Ω , on a different image, with the same values of λ_1, λ_2 as for the previous image. Best seen in color.

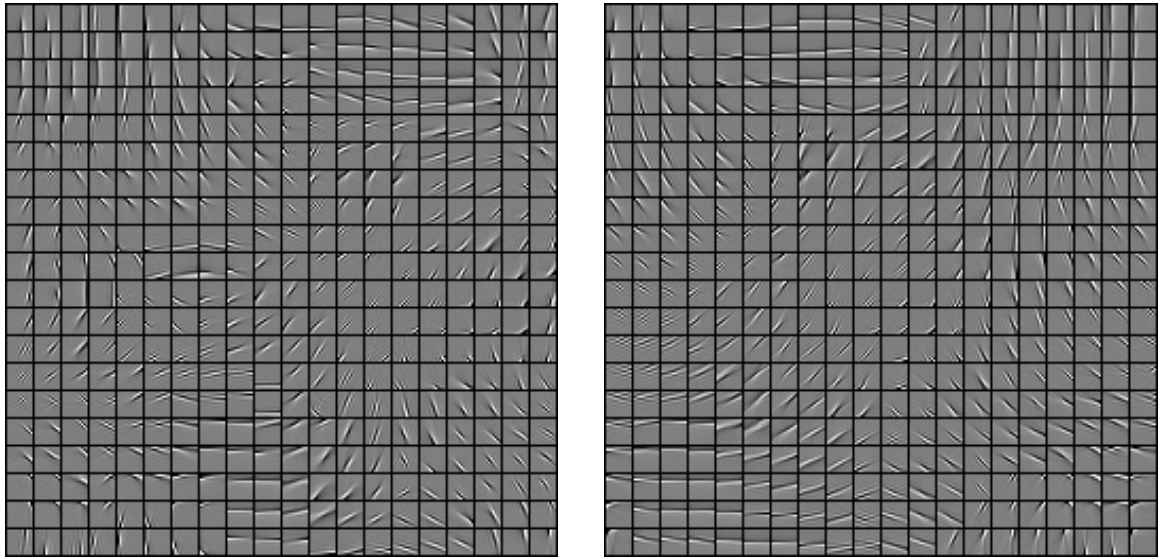


Figure 5: Topographic dictionaries with 400 elements, learned on a database of 12×12 whitened natural image patches. Left: with 3×3 cyclic overlapping groups. Right: with 4×4 cyclic overlapping groups.

As Kavukcuoglu et al. (2009); Olshausen and Field (1996), we consider a projected stochastic gradient descent algorithm for learning \mathbf{X} —that is, at iteration t , we randomly draw one signal \mathbf{y}^t from the database \mathbf{Y} , compute a sparse code $\mathbf{w}^t = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y}^t - \mathbf{X}\mathbf{w}^t\|_2^2 + \lambda \Omega(\mathbf{w})$, and update \mathbf{X} as follows: $\mathbf{X} \leftarrow \Pi_{\mathcal{C}}[\mathbf{X} - \rho(\mathbf{X}\mathbf{w}^t - \mathbf{y}^t)\mathbf{w}^{t\top}]$, where ρ is a fixed learning rate, and $\Pi_{\mathcal{C}}$ denotes the operator performing orthogonal projections onto the set \mathcal{C} . In practice, to further improve the performance, we use a mini-batch, drawing 500 signals at each iteration instead of one (see Mairal et al., 2010a). Our approach mainly differs from Kavukcuoglu et al. (2009) in the way the sparse codes \mathbf{w}^t are obtained. Whereas Kavukcuoglu et al. (2009) uses a subgradient descent algorithm to solve them, we use the proximal splitting methods presented in Section 4. The natural image patches we use are also preprocessed: They are first centered by removing their mean value (often called DC component), and whitened, as often done in the literature (Hyvärinen et al., 2001; Garrigues and Olshausen, 2010). The parameter λ is chosen such that in average $\|\mathbf{y}^i - \mathbf{X}\mathbf{w}^i\|_2 \approx 0.4\|\mathbf{y}^i\|_2$ for all new patch considered by the algorithm. Examples of obtained results are shown on Figure 5, and exhibit similarities with the topographic maps of Hyvärinen et al. (2001). Note that even though Eq. (14) is convex with respect to each variable \mathbf{X} and \mathbf{W} when one fixes the other, it is not jointly convex, and one can not guarantee our method to find a global optimum. Despite its intrinsic non-convex nature, local minima obtained with various optimization procedures have been shown to be good enough for many tasks (Elad and Aharon, 2006; Mairal et al., 2009; Kavukcuoglu et al., 2009).

5.6 Multi-Task Learning of Hierarchical Structures

As mentioned in the previous section, Jenatton et al. (2010a) have recently proposed to use a hierarchical structured norm to learn dictionaries of natural image patches. In Jenatton et al. (2010a), the dictionary elements are embedded in a *predefined* tree \mathcal{T} , via a particular instance of the structured norm Ω , which we refer to it as Ω_{tree} , and call \mathcal{G} the underlying set of groups. In this case, using

the same notation as in Section 5.5, each signal \mathbf{y}^i admits a sparse decomposition in the form of a subtree of dictionary elements.

Inspired by ideas from multi-task learning (Obozinski et al., 2010), we propose to learn the tree structure \mathcal{T} by pruning irrelevant parts of a larger initial tree \mathcal{T}_0 . We achieve this by using an additional regularization term Ω_{joint} across the different decompositions, so that subtrees of \mathcal{T}_0 will *simultaneously* be removed for all signals \mathbf{y}^i . With the notation from Section 5.5, the approach of Jenatton et al. (2010a) is then extended by the following formulation:

$$\min_{\mathbf{X} \in \mathcal{C}, \mathbf{W} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} \|\mathbf{y}^i - \mathbf{X}\mathbf{w}^i\|_2^2 + \lambda_1 \Omega_{\text{tree}}(\mathbf{w}^i) \right] + \lambda_2 \Omega_{\text{joint}}(\mathbf{W}), \quad (15)$$

where $\mathbf{W} \triangleq [\mathbf{w}^1, \dots, \mathbf{w}^n]$ is the matrix of decomposition coefficients in $\mathbb{R}^{p \times n}$. The new regularization term operates on the rows of \mathbf{W} and is defined as $\Omega_{\text{joint}}(\mathbf{W}) \triangleq \sum_{g \in \mathcal{G}} \max_{i \in \{1, \dots, n\}} |\mathbf{w}_g^i|$.²⁶ The overall penalty on \mathbf{W} , which results from the combination of Ω_{tree} and Ω_{joint} , is itself an instance of Ω with general overlapping groups, as defined in Eq (3).

To address problem (15), we use the same optimization scheme as Jenatton et al. (2010a), i.e., alternating between \mathbf{X} and \mathbf{W} , fixing one variable while optimizing with respect to the other. The task we consider is the denoising of natural image patches, with the same dataset and protocol as Jenatton et al. (2010a). We study whether learning the hierarchy of the dictionary elements improves the denoising performance, compared to standard sparse coding (i.e., when Ω_{tree} is the ℓ_1 -norm and $\lambda_2 = 0$) and the hierarchical dictionary learning of Jenatton et al. (2010a) based on predefined trees (i.e., $\lambda_2 = 0$). The dimensions of the training set — 50000 patches of size 8×8 for dictionaries with up to $p = 400$ elements — impose to handle extremely large graphs, with $|E| \approx |V| \approx 4.10^7$. Since problem (15) is too large to be solved exactly sufficiently many times to select the regularization parameters (λ_1, λ_2) rigorously, we use the following heuristics: we optimize mostly with the currently pruned tree held fixed (i.e., $\lambda_2 = 0$), and only prune the tree (i.e., $\lambda_2 > 0$) every few steps on a random subset of 10000 patches. We consider the same hierarchies as in Jenatton et al. (2010a), involving between 30 and 400 dictionary elements. The regularization parameter λ_1 is selected on the validation set of 25000 patches, for both sparse coding (Flat) and hierarchical dictionary learning (Tree). Starting from the tree giving the best performance (in this case the largest one, see Figure 6), we solve problem (15) following our heuristics, for increasing values of λ_2 . As shown in Figure 6, there is a regime where our approach performs significantly better than the two other compared methods. The standard deviation of the noise is 0.2 (the pixels have values in $[0, 1]$); no significant improvements were observed for lower levels of noise. Our experiments use the algorithm of Beck and Teboulle (2009) based on our proximal operator, with weights η_g set to 1. We present this algorithm in more details in Appendix C.

6. Conclusion

We have presented new optimization methods for solving sparse structured problems involving sums of ℓ_2 - or ℓ_∞ -norms of any (overlapping) groups of variables. Interestingly, this sheds new light on connections between sparse methods and the literature of network flow optimization. In particular, the proximal operator for the sum of ℓ_∞ -norms can be cast as a specific form of quadratic min-cost flow problem, for which we proposed an efficient and simple algorithm.

26. The simplified case where Ω_{tree} and Ω_{joint} are the ℓ_1 - and mixed ℓ_1/ℓ_2 -norms (Yuan and Lin, 2006) corresponds to Sprechmann et al. (2010).

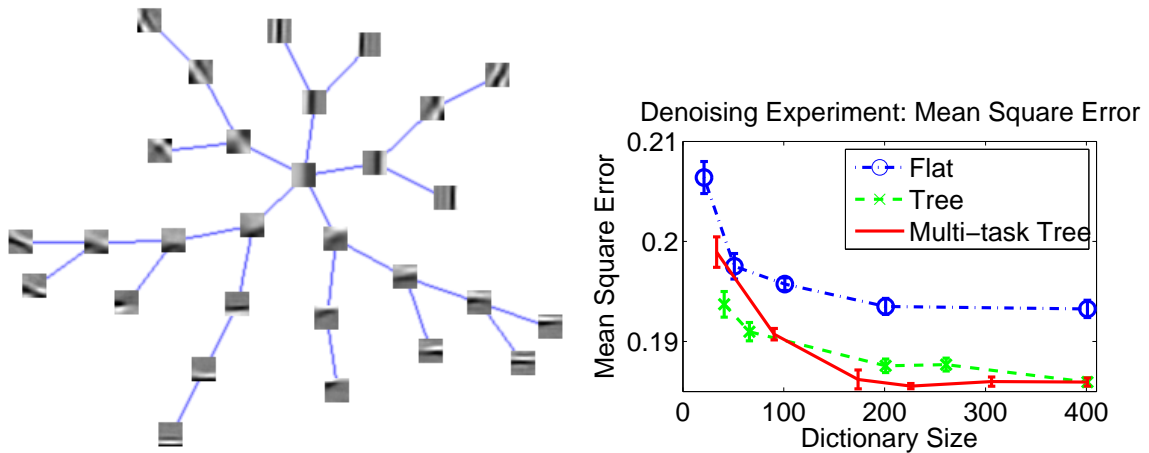


Figure 6: Left: Hierarchy obtained by pruning a larger tree of 76 elements. Right: Mean square error versus dictionary size. The error bars represent two standard deviations, based on three runs.

In addition to making it possible to resort to accelerated gradient methods, an efficient computation of the proximal operator offers more generally a certain modularity, in that it can be used as a building-block for other optimization problems. A case in point is dictionary learning where proximal problems come up and have to be solved repeatedly in an inner-loop. Interesting future work includes the computation of other structured norms such as the one introduced in Jacob et al. (2009), or total-variation based penalties, whose proximal operators are also based on minimum cost flow problems (Chambolle and Darbon, 2009). Several experiments demonstrate that our algorithm can be applied to a wide class of learning problems, which have not been addressed before with convex sparse methods.

Acknowledgments

This paper was partially supported by grants from the Agence Nationale de la Recherche (MGA Project) and from the European Research Council (SIERRA Project). In addition, Julien Mairal was supported by the NSF grant SES-0835531 and NSF award CCF-0939370. The authors would like to thank Jean Ponce for interesting discussions and suggestions for improving this manuscript, Jean-Christophe Pesquet and Patrick-Louis Combettes for pointing us to the literature of proximal splitting methods, Ryota Tomioka for his advice on using augmented Lagrangian methods.

Appendix A. Equivalence to Canonical Graphs

Formally, the notion of equivalence between graphs can be summarized by the following lemma:

Lemma 2 (Equivalence to canonical graphs.)

Let $G = (V, E, s, t)$ be the canonical graph corresponding to a group structure \mathcal{G} . Let $G' = (V, E', s, t)$ be a graph sharing the same set of vertices, source and sink as G , but with a different arc set E' . We say that G' is equivalent to G if and only if the following conditions hold:

- Arcs of E' outgoing from the source are the same as in E , with the same costs and capacities.

- Arcs of E' going to the sink are the same as in E , with the same costs and capacities.
- For every arc (g, j) in E , with (g, j) in $V_{gr} \times V_u$, there exists a unique path in E' from g to j with zero costs and infinite capacities on every arc of the path.
- Conversely, if there exists a path in E' between a vertex g in V_{gr} and a vertex j in V_u , then there exists an arc (g, j) in E .

Then, the cost of the optimal min-cost flow on G and G' are the same. Moreover, the values of the optimal flow on the arcs (j, t) , j in V_u , are the same on G and G' .

Proof. We first notice that on both G and G' , the cost of a flow on the graph only depends on the flow on the arcs (j, t) , j in V_u , which we have denoted by $\bar{\xi}$ in E .

We will prove that finding a feasible flow π on G with a cost $c(\pi)$ is equivalent to finding a feasible flow π' on G' with the same cost $c(\pi) = c(\pi')$. We now use the concept of *path flow*, which is a flow vector in G carrying the same positive value on every arc of a directed path between two nodes of G . It intuitively corresponds to sending a positive amount of flow along a path of the graph.

According to the definition of graph equivalence introduced in the Lemma, it is easy to show that there is a bijection between the arcs in E , and the paths in E' with positive capacities on every arc. Given now a feasible flow π in G , we build a feasible flow π' on G' which is a *sum* of path flows. More precisely, for every arc a in E , we consider its equivalent path in E' , with a path flow carrying the same amount of flow as a . Therefore, each arc a' in E' has a total amount of flow that is equal to the sum of the flows carried by the path flows going over a' . It is also easy to show that this construction builds a flow on G' (capacity and conservation constraints are satisfied) and that this flow π' has the same cost as π , that is, $c(\pi) = c(\pi')$.

Conversely, given a flow π' on G' , we use a classical path flow decomposition (see Bertsekas, 1998, Proposition 1.1), saying that there exists a decomposition of π' as a sum of path flows in E' . Using the bijection described above, we know that each path in the previous sums corresponds to a unique arc in E . We now build a flow π in G , by associating to each path flow in the decomposition of π' , an arc in E carrying the same amount of flow. The flow of every other arc in E is set to zero. It is also easy to show that this builds a valid flow in G that has the same cost as π' . ■

Appendix B. Convergence Analysis

We show in this section the correctness of Algorithm 1 for computing the proximal operator, and of Algorithm 2 for computing the dual norm Ω^* .

B.1 Computation of the Proximal Operator

We first prove that our algorithm converges and that it finds the optimal solution of the proximal problem. This requires that we introduce the optimality conditions for problem (6) derived from Jenatton et al. (2010a, 2011) since our convergence proof essentially checks that these conditions are satisfied upon termination of the algorithm.

Lemma 3 (Optimality conditions of the problem (6) from Jenatton et al. 2010a, 2011)

The primal-dual variables (\mathbf{w}, ξ) are respectively solutions of the primal (4) and dual problems (6)

Recall that we assume (cf. Section 3.3) that the scalars \mathbf{u}_j are all non negative, and that we add non-negativity constraints on ξ . With the optimality conditions of Lemma 3 in hand, we can show our first convergence result.

Proposition 1 (Convergence of Algorithm 1)

Algorithm 1 converges in a finite and polynomial number of operations.

Proof. Our algorithm splits recursively the graph into disjoint parts and processes each part recursively. The processing of one part requires an orthogonal projection onto an ℓ_1 -ball and a max-flow algorithm, which can both be computed in polynomial time. To prove that the procedure converges, it is sufficient to show that when the procedure `computeFlow` is called for a graph (V, E, s, t) and computes a cut (V^+, V^-) , then the components V^+ and V^- are both non-empty.

Suppose for instance that $V^- = \emptyset$. In this case, the capacity of the min-cut is equal to $\sum_{j \in V_u} \gamma_j$, and the value of the max-flow is $\sum_{j \in V_u} \bar{\xi}_j$. Using the classical max-flow/min-cut theorem (Ford and Fulkerson, 1956), we have equality between these two terms. Since, by definition of both γ and $\bar{\xi}$, we have for all j in V_u , $\bar{\xi}_j \leq \gamma_j$, we obtain a contradiction with the existence of j in V_u such that $\bar{\xi}_j \neq \gamma_j$.

Conversely, suppose now that $V^+ = \emptyset$. Then, the value of the max-flow is still $\sum_{j \in V_u} \bar{\xi}_j$, and the value of the min-cut is $\lambda \sum_{g \in V_{gr}} \eta_g$. Using again the max-flow/min-cut theorem, we have that $\sum_{j \in V_u} \bar{\xi}_j = \lambda \sum_{g \in V_{gr}} \eta_g$. Moreover, by definition of γ , we also have $\sum_{j \in V_u} \bar{\xi}_j \leq \sum_{j \in V_u} \gamma_j \leq \lambda \sum_{g \in V_{gr}} \eta_g$, leading to a contradiction with the existence of j in V_u satisfying $\bar{\xi}_j \neq \gamma_j$. We remind the reader of the fact that such a $j \in V_u$ exists since the cut is only computed when the current estimate $\bar{\xi}$ is not optimal yet. This proof holds for any graph that is equivalent to the canonical one. ■

After proving the convergence, we prove that the algorithm is correct with the next proposition.

Proposition 2 (Correctness of Algorithm 1)

Algorithm 1 solves the proximal problem of Eq. (4).

Proof. For a group structure \mathcal{G} , we first prove the correctness of our algorithm if the graph used is its associated canonical graph that we denote $G_0 = (V_0, E_0, s, t)$. We proceed by induction on the number of nodes of the graph. The induction hypothesis $\mathcal{H}(k)$ is the following:

For all canonical graphs $G = (V = V_u \cup V_{gr}, E, s, t)$ associated with a group structure \mathcal{G}_V with weights $(\eta_g)_{g \in \mathcal{G}_V}$ such that $|V| \leq k$, `computeFlow` (V, E) solves the following optimization problem:

$$\min_{(\xi_j^g)} \sum_{j \in V_u, g \in V_{gr}} \frac{1}{2} (\mathbf{u}_j - \sum_{g \in V_{gr}} \xi_j^g)^2 \quad \text{s.t.} \quad \forall g \in V_{gr}, \sum_{j \in V_u} \xi_j^g \leq \lambda \eta_g \quad \text{and} \quad \xi_j^g = 0, \forall j \notin g. \quad (16)$$

Since $\mathcal{G}_{V_0} = \mathcal{G}$, it is sufficient to show that $\mathcal{H}(|V_0|)$ to prove the proposition.

We initialize the induction by $\mathcal{H}(2)$, corresponding to the simplest canonical graph, for which $|V_{gr}| = |V_u| = 1$. Simple algebra shows that $\mathcal{H}(2)$ is indeed correct.

We now suppose that $\mathcal{H}(k')$ is true for all $k' < k$ and consider a graph $G = (V, E, s, t)$, $|V| = k$. The first step of the algorithm computes the variable $(\gamma_j)_{j \in V_u}$ by a projection on the ℓ_1 -ball. This is

itself an instance of the dual formulation of Eq. (6) in a simple case, with one group containing all variables. We can therefore use Lemma 3 to characterize the optimality of $(\gamma_j)_{j \in V_u}$, which yields

$$\begin{cases} \sum_{j \in V_u} (\mathbf{u}_j - \gamma_j) \gamma_j = (\max_{j \in V_u} |\mathbf{u}_j - \gamma_j|) \sum_{j \in V_u} \gamma_j \text{ and } \sum_{j \in V_u} \gamma_j = \lambda \sum_{g \in V_{gr}} \eta_g, \\ \text{or } \mathbf{u}_j - \gamma_j = 0, \forall j \in V_u. \end{cases} \quad (17)$$

The algorithm then computes a max-flow, using the scalars γ_j as capacities, and we now have two possible situations:

1. If $\bar{\xi}_j = \gamma_j$ for all j in V_u , the algorithm stops; we write $\mathbf{w}_j = \mathbf{u}_j - \bar{\xi}_j$ for j in V_u , and using Eq. (17), we obtain

$$\begin{cases} \sum_{j \in V_u} \mathbf{w}_j \bar{\xi}_j = (\max_{j \in V_u} |\mathbf{w}_j|) \sum_{j \in V_u} \bar{\xi}_j \text{ and } \sum_{j \in V_u} \bar{\xi}_j = \lambda \sum_{g \in V_{gr}} \eta_g, \\ \text{or } \mathbf{w}_j = 0, \forall j \in V_u. \end{cases} \quad (18)$$

We can rewrite the condition above as

$$\sum_{g \in V_{gr}} \sum_{j \in g} \mathbf{w}_j \xi_j^g = \sum_{g \in V_{gr}} (\max_{j \in V_u} |\mathbf{w}_j|) \sum_{j \in V_u} \xi_j^g.$$

Since all the quantities in the previous sum are positive, this can only hold if for all $g \in V_{gr}$,

$$\sum_{j \in V_u} \mathbf{w}_j \xi_j^g = (\max_{j \in V_u} |\mathbf{w}_j|) \sum_{j \in V_u} \xi_j^g.$$

Moreover, by definition of the max flow and the optimality conditions, we have

$$\forall g \in V_{gr}, \sum_{j \in V_u} \xi_j^g \leq \lambda \eta_g, \text{ and } \sum_{j \in V_u} \bar{\xi}_j = \lambda \sum_{g \in V_{gr}} \eta_g,$$

which leads to

$$\forall g \in V_{gr}, \sum_{j \in V_u} \xi_j^g = \lambda \eta_g.$$

By Lemma 3, we have shown that the problem (16) is solved.

2. Let us now consider the case where there exists j in V_u such that $\bar{\xi}_j \neq \gamma_j$. The algorithm splits the vertex set V into two parts V^+ and V^- , which we have proven to be non-empty in the proof of Proposition 1. The next step of the algorithm removes all edges between V^+ and V^- (see Figure 7). Processing (V^+, E^+) and (V^-, E^-) independently, it updates the value of the flow matrix ξ_j^g , $j \in V_u$, $g \in V_{gr}$, and the corresponding flow vector $\bar{\xi}_j$, $j \in V_u$. As for V , we denote by $V_u^+ \triangleq V^+ \cap V_u$, $V_u^- \triangleq V^- \cap V_u$ and $V_{gr}^+ \triangleq V^+ \cap V_{gr}$, $V_{gr}^- \triangleq V^- \cap V_{gr}$.

Then, we notice that (V^+, E^+, s, t) and (V^-, E^-, s, t) are respective canonical graphs for the group structures $\mathcal{G}_{V^+} \triangleq \{g \cap V_u^+ \mid g \in V_{gr}\}$, and $\mathcal{G}_{V^-} \triangleq \{g \cap V_u^- \mid g \in V_{gr}\}$.

Writing $\mathbf{w}_j = \mathbf{u}_j - \bar{\xi}_j$ for j in V_u , and using the induction hypotheses $\mathcal{H}(|V^+|)$ and $\mathcal{H}(|V^-|)$, we now have the following optimality conditions deriving from Lemma 3 applied on Eq. (16) respectively for the graphs (V^+, E^+) and (V^-, E^-) :

$$\forall g \in V_{gr}^+, g' \triangleq g \cap V_u^+, \begin{cases} \mathbf{w}_{g'}^\top \xi_{g'}^g = \|\mathbf{w}_{g'}\|_\infty \sum_{j \in g'} \xi_j^g \text{ and } \sum_{j \in g'} \xi_j^g = \lambda \eta_g, \\ \text{or } \mathbf{w}_{g'} = 0, \end{cases} \quad (19)$$

and

$$\forall g \in V_{gr}^-, g' \triangleq g \cap V_u^-, \begin{cases} \mathbf{w}_{g'}^\top \xi_{g'}^g = \|\mathbf{w}_{g'}\|_\infty \sum_{j \in g'} \xi_j^g & \text{and } \sum_{j \in g'} \xi_j^g = \lambda \eta_g, \\ \text{or } \mathbf{w}_{g'} = 0. \end{cases} \quad (20)$$

We will now combine Eq. (19) and Eq. (20) into optimality conditions for Eq. (16). We first notice that $g \cap V_u^+ = g$ since there are no arcs between V^+ and V^- in E (see the properties of the cuts discussed before this proposition). It is therefore possible to replace g' by g in Eq. (19). We will show that it is possible to do the same in Eq. (20), so that combining these two equations yield the optimality conditions of Eq. (16).

More precisely, we will show that for all $g \in V_{gr}^-$ and $j \in g \cap V_u^+$, $|\mathbf{w}_j| \leq \max_{l \in g \cap V_u^-} |\mathbf{w}_l|$, in which case g' can be replaced by g in Eq. (20). This result is relatively intuitive: (s, V^+) and (V^-, t) being an (s, t) -cut, all arcs between s and V^- are saturated, while there are unsaturated arcs between s and V^+ ; one therefore expects the residuals $\mathbf{u}_j - \bar{\xi}_j$ to decrease on the V^+ side, while increasing on the V^- side. The proof is nonetheless a bit technical.

Let us show first that for all g in V_{gr}^+ , $\|\mathbf{w}_g\|_\infty \leq \max_{j \in V_u} |\mathbf{u}_j - \gamma_j|$. We split the set V^+ into disjoint parts:

$$\begin{aligned} V_{gr}^{++} &\triangleq \{g \in V_{gr}^+ \text{ s.t. } \|\mathbf{w}_g\|_\infty \leq \max_{j \in V_u} |\mathbf{u}_j - \gamma_j|\}, \\ V_u^{++} &\triangleq \{j \in V_u^+ \text{ s.t. } \exists g \in V_{gr}^{++}, j \in g\}, \\ V_{gr}^{+-} &\triangleq V_{gr}^+ \setminus V_{gr}^{++} = \{g \in V_{gr}^+ \text{ s.t. } \|\mathbf{w}_g\|_\infty > \max_{j \in V_u} |\mathbf{u}_j - \gamma_j|\}, \\ V_u^{+-} &\triangleq V_u^+ \setminus V_u^{++}. \end{aligned}$$

As previously, we denote $V^{+-} \triangleq V_u^{+-} \cup V_{gr}^{+-}$ and $V^{++} \triangleq V_u^{++} \cup V_{gr}^{++}$. We want to show that V_{gr}^{+-} is necessarily empty. We reason by contradiction and assume that $V_{gr}^{+-} \neq \emptyset$.

According to the definition of the different sets above, we observe that no arcs are going from V^{++} to V^{+-} , that is, for all g in V_{gr}^{++} , $g \cap V_u^{+-} = \emptyset$. We observe as well that the flow from V_{gr}^{+-} to V_u^{++} is the null flow, because optimality conditions (19) imply that for a group g only nodes $j \in g$ such that $\mathbf{w}_j = \|\mathbf{w}_g\|_\infty$ receive some flow, which excludes nodes in V_u^{++} provided $V_{gr}^{+-} \neq \emptyset$; Combining this fact and the inequality $\sum_{g \in V_{gr}^+} \lambda \eta_g \geq \sum_{j \in V_u^+} \gamma_j$ (which is a direct consequence of the minimum (s, t) -cut), we have as well

$$\sum_{g \in V_{gr}^{+-}} \lambda \eta_g \geq \sum_{j \in V_u^{+-}} \gamma_j.$$

Let $j \in V_u^{+-}$, if $\bar{\xi}_j \neq 0$ then for some $g \in V_{gr}^{+-}$ such that j receives some flow from g , which from the optimality conditions (19) implies $\mathbf{w}_j = \|\mathbf{w}_g\|_\infty$; by definition of V_{gr}^{+-} , $\|\mathbf{w}_g\|_\infty > \mathbf{u}_j - \gamma_j$. But since at the optimum, $\mathbf{w}_j = \mathbf{u}_j - \bar{\xi}_j$, this implies that $\bar{\xi}_j < \gamma_j$, and in turn that $\sum_{j \in V_u^{+-}} \bar{\xi}_j = \lambda \sum_{g \in V_{gr}^{+-}} \eta_g$. Finally,

$$\lambda \sum_{g \in V_{gr}^{+-}} \eta_g = \sum_{j \in V_u^{+-}, \bar{\xi}_j \neq 0} \bar{\xi}_j < \sum_{j \in V_u^{+-}} \gamma_j$$

and this is a contradiction.

We now have that for all g in V_{gr}^+ , $\|\mathbf{w}_g\|_\infty \leq \max_{j \in V_u} |\mathbf{u}_j - \gamma_j|$. The proof showing that for all g in V_{gr}^- , $\|\mathbf{w}_g\|_\infty \geq \max_{j \in V_u} |\mathbf{u}_j - \gamma_j|$, uses the same kind of decomposition for V^- , and follows along similar arguments. We will therefore not detail it.

To summarize, we have shown that for all $g \in V_{gr}^-$ and $j \in g \cap V_u^+$, $|\mathbf{w}_j| \leq \max_{l \in g \cap V_u^-} |\mathbf{w}_l|$. Since there is no flow from V^- to V^+ , i.e., $\xi_j^g = 0$ for g in V_{gr}^- and j in V_u^+ , we can now replace the definition of g' in Eq. (20) by $g' \triangleq g \cap V_u$, the combination of Eq. (19) and Eq. (20) gives us optimality conditions for Eq. (16).

The proposition being proved for the canonical graph, we extend it now for an equivalent graph in the sense of Lemma 2. First, we observe that the algorithm gives the same values of γ for two equivalent graphs. Then, it is easy to see that the value $\bar{\xi}$ given by the max-flow, and the chosen (s, t) -cut is the same, which is enough to conclude that the algorithm performs the same steps for two equivalent graphs. ■

B.2 Computation of the Dual Norm Ω^*

As for the proximal operator, the computation of dual norm Ω^* can itself be shown to solve another network flow problem, based on the following variational formulation, which extends a previous result from Jenatton et al. (2009):

Lemma 4 (Dual formulation of the dual-norm Ω^* .)

Let $\kappa \in \mathbb{R}^p$. We have

$$\Omega^*(\kappa) = \min_{\xi \in \mathbb{R}^{p \times |\mathcal{G}|}, \tau \in \mathbb{R}} \tau \quad \text{s.t.} \quad \sum_{g \in \mathcal{G}} \xi^g = \kappa, \text{ and } \forall g \in \mathcal{G}, \|\xi^g\|_1 \leq \tau \eta_g \text{ with } \xi_j^g = 0 \text{ if } j \notin g.$$

Proof. By definition of $\Omega^*(\kappa)$, we have

$$\Omega^*(\kappa) \triangleq \max_{\Omega(\mathbf{z}) \leq 1} \mathbf{z}^\top \kappa.$$

By introducing the primal variables $(\alpha_g)_{g \in \mathcal{G}} \in \mathbb{R}^{|\mathcal{G}|}$, we can rewrite the previous maximization problem as

$$\Omega^*(\kappa) = \max_{\sum_{g \in \mathcal{G}} \eta_g \alpha_g \leq 1} \kappa^\top \mathbf{z}, \quad \text{s.t.} \quad \forall g \in \mathcal{G}, \|\mathbf{z}_g\|_\infty \leq \alpha_g,$$

with the additional $|\mathcal{G}|$ conic constraints $\|\mathbf{z}_g\|_\infty \leq \alpha_g$. This primal problem is convex and satisfies Slater's conditions for generalized conic inequalities, which implies that strong duality holds (Boyd and Vandenberghe, 2004). We now consider the Lagrangian \mathcal{L} defined as

$$\mathcal{L}(\mathbf{z}, \alpha_g, \tau, \gamma_g, \xi) = \kappa^\top \mathbf{z} + \tau(1 - \sum_{g \in \mathcal{G}} \eta_g \alpha_g) + \sum_{g \in \mathcal{G}} \begin{pmatrix} \alpha_g \\ \mathbf{z}_g \end{pmatrix}^\top \begin{pmatrix} \gamma_g \\ \xi_g^g \end{pmatrix},$$

with the dual variables $\{\tau, (\gamma_g)_{g \in \mathcal{G}}, \xi\} \in \mathbb{R}_+ \times \mathbb{R}^{|\mathcal{G}|} \times \mathbb{R}^{p \times |\mathcal{G}|}$ such that for all $g \in \mathcal{G}$, $\xi_j^g = 0$ if $j \notin g$ and $\|\xi^g\|_1 \leq \gamma_g$. The dual function is obtained by taking the derivatives of \mathcal{L} with respect to the primal variables \mathbf{z} and $(\alpha_g)_{g \in \mathcal{G}}$ and equating them to zero, which leads to

$$\begin{aligned} \forall j \in \{1, \dots, p\}, \quad \kappa_j + \sum_{g \in \mathcal{G}} \xi_j^g &= 0 \\ \forall g \in \mathcal{G}, \quad \tau \eta_g - \gamma_g &= 0. \end{aligned}$$

After simplifying the Lagrangian and flipping the sign of ξ , the dual problem then reduces to

$$\min_{\xi \in \mathbb{R}^{p \times |\mathcal{G}|}, \tau \in \mathbb{R}} \tau \quad \text{s.t.} \quad \begin{cases} \forall j \in \{1, \dots, p\}, \kappa_j = \sum_{g \in \mathcal{G}} \xi_j^g \text{ and } \xi_j^g = 0 \text{ if } j \notin g, \\ \forall g \in \mathcal{G}, \|\xi^g\|_1 \leq \tau \eta_g, \end{cases}$$

which is the desired result. \blacksquare

We now prove that Algorithm 2 is correct.

Proposition 3 (Convergence and correctness of Algorithm 2)

Algorithm 2 computes the value of the dual norm of Eq. (9) in a finite and polynomial number of operations.

Proof. The convergence of the algorithm only requires to show that the cardinality of V in the different calls of the function `computeFlow` strictly decreases. Similar arguments to those used in the proof of Proposition 1 can show that each part of the cuts (V^+, V^-) are both non-empty. The algorithm thus requires a finite number of calls to a max-flow algorithm and converges in a finite and polynomial number of operations.

Let us now prove that the algorithm is correct for a canonical graph. We proceed again by induction on the number of nodes of the graph. More precisely, we consider the induction hypothesis $\mathcal{H}'(k)$ defined as:

for all canonical graphs $G = (V, E, s, t)$ associated with a group structure \mathcal{G}_V and such that $|V| \leq k$, $\text{dualNormAux}(V = V_u \cup V_{gr}, E)$ solves the following optimization problem:

$$\min_{\xi, \tau} \tau \quad \text{s.t.} \quad \forall j \in V_u, \kappa_j = \sum_{g \in V_{gr}} \xi_j^g, \text{ and } \forall g \in V_{gr}, \sum_{j \in V_u} \xi_j^g \leq \tau \eta_g \text{ with } \xi_j^g = 0 \text{ if } j \notin g. \quad (21)$$

We first initialize the induction by $\mathcal{H}'(2)$ (i.e., with the simplest canonical graph, such that $|V_{gr}| = |V_u| = 1$). Simple algebra shows that $\mathcal{H}'(2)$ is indeed correct.

We next consider a canonical graph $G = (V, E, s, t)$ such that $|V| = k$, and suppose that $\mathcal{H}'(k-1)$ is true. After the max-flow step, we have two possible cases to discuss:

1. If $\bar{\xi}_j = \gamma_j$ for all j in V_u , the algorithm stops. We know that any scalar τ such that the constraints of Eq. (21) are all satisfied necessarily verifies $\sum_{g \in V_{gr}} \tau \eta_g \geq \sum_{j \in V_u} \kappa_j$. We have indeed that $\sum_{g \in V_{gr}} \tau \eta_g$ is the value of an (s, t) -cut in the graph, and $\sum_{j \in V_u} \kappa_j$ is the value of the max-flow, and the inequality follows from the max-flow/min-cut theorem (Ford and Fulkerson, 1956). This gives a lower-bound on τ . Since this bound is reached, τ is necessarily optimal.
2. We now consider the case where there exists j in V_u such that $\bar{\xi}_j \neq \kappa_j$, meaning that for the given value of τ , the constraint set of Eq. (21) is not feasible for ξ , and that the value of τ should necessarily increase. The algorithm splits the vertex set V into two non-empty parts V^+ and V^- and we remark that there are no arcs going from V^+ to V^- , and no flow going from V^- to V^+ . Since the arcs going from s to V^- are saturated, we have that $\sum_{g \in V_{gr}^-} \tau \eta_g \leq \sum_{j \in V_u^-} \kappa_j$. Let us now consider τ^* the solution of Eq. (21). Using the induction hypothesis $\mathcal{H}'(|V^-|)$, the algorithm computes a new value τ' that solves Eq. (21) when replacing V by V^- and this new value satisfies the following inequality $\sum_{g \in V_{gr}^-} \tau' \eta_g \geq \sum_{j \in V_u^-} \kappa_j$. The value of τ' has therefore increased and the updated flow ξ now satisfies the constraints of Eq. (21) and therefore $\tau' \geq \tau^*$. Since there are no arcs going from V^+ to V^- , τ^* is feasible for Eq. (21) when replacing V by V^- and we have that $\tau^* \geq \tau'$ and then $\tau' = \tau^*$.

To prove that the result holds for any equivalent graph, similar arguments to those used in the proof of Proposition 1 can be exploited, showing that the algorithm computes the same values of τ and same (s, t) -cuts at each step. ■

Appendix C. Algorithm FISTA with duality gap

In this section, we describe in details the algorithm FISTA (Beck and Teboulle, 2009) when applied to solve problem (2), with a duality gap as the stopping criterion. The algorithm, as implemented in the experiments, is summarized in Algorithm 3.

Without loss of generality, let us assume we are looking for models of the form $\mathbf{X}\mathbf{w}$, for some matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ (typically, a linear model where \mathbf{X} is the design matrix composed of n observations in \mathbb{R}^p). Thus, we can consider the following primal problem

$$\min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{X}\mathbf{w}) + \lambda\Omega(\mathbf{w}), \quad (22)$$

in place of problem (2). Based on Fenchel duality arguments (Borwein and Lewis, 2006),

$$f(\mathbf{X}\mathbf{w}) + \lambda\Omega(\mathbf{w}) + f^*(-\boldsymbol{\kappa}), \text{ for } \mathbf{w} \in \mathbb{R}^p, \boldsymbol{\kappa} \in \mathbb{R}^n \text{ and } \Omega^*(\mathbf{X}^\top \boldsymbol{\kappa}) \leq \lambda,$$

is a duality gap for problem (22), where $f^*(\boldsymbol{\kappa}) \triangleq \sup_{\mathbf{z}} [\mathbf{z}^\top \boldsymbol{\kappa} - f(\mathbf{z})]$ is the Fenchel conjugate of f (Borwein and Lewis, 2006). Given a primal variable \mathbf{w} , a good dual candidate $\boldsymbol{\kappa}$ can be obtained by looking at the conditions that have to be satisfied by the pair $(\mathbf{w}, \boldsymbol{\kappa})$ at optimality (Borwein and Lewis, 2006). In particular, the dual variable $\boldsymbol{\kappa}$ is chosen to be

$$\boldsymbol{\kappa} = -\rho^{-1} \nabla f(\mathbf{X}\mathbf{w}), \text{ with } \rho \triangleq \max \{ \lambda^{-1} \Omega^*(\mathbf{X}^\top \nabla f(\mathbf{X}\mathbf{w})), 1 \}.$$

Consequently, computing the duality gap requires evaluating the dual norm Ω^* , as explained in Algorithm 2. We sum up the computation of the duality gap in Algorithm 3. Moreover, we refer to the proximal operator associated with $\lambda\Omega$ as $\text{prox}_{\lambda\Omega}$.²⁷

In our experiment, we choose the line-search parameter ν to be equal to 1.5.

Appendix D. Speed comparison of Algorithm 1 with parametric max-flow algorithms

As shown by Hochbaum and Hong (1995), min-cost flow problems, and in particular, the dual problem of (4), can be reduced to a specific *parametric max-flow* problem. We thus compare our approach (ProxFlow) with the efficient parametric max-flow algorithm proposed by Gallo et al. (1989) and a simplified version of the latter proposed by Babenko and Goldberg (2006). We refer to these two algorithms as GGT and SIMP respectively. The benchmark is established on the same datasets as those already used in the experimental section of the paper, namely: (1) three datasets built from overcomplete bases of discrete cosine transforms (DCT), with respectively 10^4 , 10^5 and 10^6 variables, and (2) images used for the background subtraction task, composed of 57600 pixels. For GGT and SIMP, we use the paraF software which is a C++ parametric max-flow implementation available at <http://www.avglab.com/andrew/soft.html>. Experiments were conducted on

27. As a brief reminder, it is defined as the function that maps the vector \mathbf{u} in \mathbb{R}^p to the (unique, by strong convexity) solution of Eq. (4).

Algorithm 3 FISTA procedure to solve problem (22).

- 1: **Inputs:** initial $\mathbf{w}_{(0)} \in \mathbb{R}^p$, Ω , $\lambda > 0$, $\varepsilon_{\text{gap}} > 0$ (precision for the duality gap).
- 2: **Parameters:** $\nu > 1$, $L_0 > 0$.
- 3: **Outputs:** solution \mathbf{w} .
- 4: **Initialization:** $\mathbf{y}_{(1)} = \mathbf{w}_{(0)}$, $t_1 = 1$, $k = 1$.
- 5: **while** { computeDualityGap($\mathbf{w}_{(k-1)}$) $> \varepsilon_{\text{gap}}$ } **do**
- 6: Find the smallest integer $s_k \geq 0$ such that
- 7: $f(\text{prox}_{[\lambda\Omega]}(\mathbf{y}_{(k)})) \leq f(\mathbf{y}_{(k)}) + \Delta_{(k)}^\top \nabla f(\mathbf{y}_{(k)}) + \frac{\tilde{L}}{2} \|\Delta_{(k)}\|_2^2$,
- 8: with $\tilde{L} \triangleq L_k \nu^{s_k}$ and $\Delta_{(k)} \triangleq \mathbf{y}_{(k)} - \text{prox}_{[\lambda\Omega]}(\mathbf{y}_{(k)})$.
- 9: $L_k \leftarrow L_{k-1} \nu^{s_k}$.
- 10: $\mathbf{w}_{(k)} \leftarrow \text{prox}_{[\lambda\Omega]}(\mathbf{y}_{(k)})$.
- 11: $t_{k+1} \leftarrow (1 + \sqrt{1 + t_k^2})/2$.
- 12: $\mathbf{y}_{(k+1)} \leftarrow \mathbf{w}_{(k)} + \frac{t_k - 1}{t_{k+1}} (\mathbf{w}_{(k)} - \mathbf{w}_{(k-1)})$.
- 13: $k \leftarrow k + 1$.
- 14: **end while**
- 15: **Return:** $\mathbf{w} \leftarrow \mathbf{w}_{(k-1)}$.

Procedure computeDualityGap(\mathbf{w})

- 1: $\kappa \leftarrow -\rho^{-1} \nabla f(\mathbf{X}\mathbf{w})$, with $\rho \triangleq \max \{ \lambda^{-1} \Omega^*(\mathbf{X}^\top \nabla f(\mathbf{X}\mathbf{w})), 1 \}$.
 - 2: **Return:** $f(\mathbf{X}\mathbf{w}) + \lambda \Omega(\mathbf{w}) + f^*(-\kappa)$.
-

a single-core 2.33 Ghz. We report in the following table the average execution time in seconds of each algorithm for 5 runs, as well as the statistics of the corresponding problems:

Number of variables p	10000	100000	1000000	57600
$ V $	20000	200000	2000000	57600
$ E $	110000	500000	11000000	579632
ProxFlow (in sec.)	0.4	3.1	113.0	1.7
GGT (in sec.)	2.4	26.0	525.0	16.7
SIMP (in sec.)	1.2	13.1	284.0	8.31

Although we provide the speed comparison for a single value of λ (the one used in the corresponding experiments of the paper), we observed that our approach consistently outperforms GGT and SIMP for values of λ corresponding to different regularization regimes.

References

- R. K. Ahuja, T. L. Magnanti, and J. Orlin. *Network Flows*. Prentice Hall, 1993.
- A. Argyriou, T. Evgeniou, and M. Pontil. Convex multi-task feature learning. *Machine Learning*, 73(3):243–272, 2008.
- M. Babenko and A.V. Goldberg. Experimental evaluation of a parametric flow algorithm. Technical report, Microsoft Research, 2006. MSR-TR-2006-77.

- F. Bach. High-dimensional non-linear variable selection through hierarchical kernel learning. Technical report, arXiv:0909.0844, 2009.
- F. Bach. Structured sparsity-inducing norms through submodular functions. In *Advances in Neural Information Processing Systems*, 2010.
- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Convex optimization with sparsity-inducing norms. In S. Sra, S. Nowozin, and S. J. Wright, editors, *Optimization for Machine Learning*. MIT Press, 2011. To appear.
- R. G. Baraniuk, V. Cevher, M. Duarte, and C. Hegde. Model-based compressive sensing. *IEEE Transactions on Information Theory*, 56(4):1982–2001, 2010.
- A. Barron, J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE Transactions on Information Theory*, 44(6):2743–2760, 1998.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- D. P. Bertsekas. *Network Optimization: Continuous and Discrete Models*. Athena Scientific, 1998.
- D. P. Bertsekas. *Nonlinear programming*. Athena Scientific Belmont, 1999.
- D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and distributed computation: Numerical Methods*. Prentice Hall Inc., 1989.
- P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37(4):1705–1732, 2009.
- J. Bien, Y. Xu, and M. W. Mahoney. CUR from a sparse optimization viewpoint. In *Advances in Neural Information Processing Systems*, 2010.
- J. M. Borwein and A. S. Lewis. *Convex analysis and nonlinear optimization: Theory and examples*. Springer, 2006.
- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
- S. P. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- Y. Boykov and V. Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(9):1124–1137, 2004.
- P. Brucker. An $O(n)$ algorithm for quadratic knapsack problems. *Operations Research Letters*, 3: 163–166, 1984.
- T. T. Cai. Adaptive wavelet estimation: a block thresholding and oracle inequality approach. *Annals of statistics*, pages 898–924, 1999.

- V. Cehver, M. F. Duarte, C. Hedge, and R. G. Baraniuk. Sparse signal recovery using Markov random fields. In *Advances in Neural Information Processing Systems*, 2008.
- A. Chambolle and J. Darbon. On total variation minimization and surface evolution using parametric maximal flows. *International Journal of Computer Vision*, 84(3), 2009.
- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20:33–61, 1999.
- X. Chen, Q. Lin, S. Kim, J. Pena, J. G. Carbonell, and E. P. Xing. An efficient proximal-gradient method for single and multi-task regression with structured sparsity. Technical report, 2010. ArXiv:1005.4717v1.
- B. V. Cherkassky and A. V. Goldberg. On implementing the push-relabel method for the maximum flow problem. *Algorithmica*, 19(4):390–410, 1997.
- P. L. Combettes and J.-C. Pesquet. A proximal decomposition method for solving convex variational inverse problems. *Inverse Problems*, 24(27), 2008. Art. 065014.
- P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*. Springer, 2010.
- D. L. Donoho and I. M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association*, 90(432):1200–1224, 1995.
- J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32(2):407–499, 2004.
- M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image Processing*, 54(12):3736–3745, December 2006.
- L. R. Ford and D. R. Fulkerson. Maximal flow through a network. *Canadian Journal of Mathematics*, 8(3):399–404, 1956.
- J. Friedman, T. Hastie, and R. Tibshirani. A note on the group Lasso and a sparse group Lasso. Technical report, Preprint arXiv:1001.0736, 2010.
- S. Fujishige. *Submodular functions and optimization*. Elsevier, 2005.
- G. Gallo, M. E. Grigoriadis, and R. E. Tarjan. A fast parametric maximum flow algorithm and applications. *SIAM J. Comput.*, 18:30–55, 1989.
- P. Garrigues and B. Olshausen. Group sparse coding with a Laplacian scale mixture prior. In *Advances in Neural Information Processing Systems*, 2010.
- A. V. Goldberg and R. E. Tarjan. A new approach to the maximum flow problem. In *Proc. of ACM Symposium on Theory of Computing*, pages 136–146, 1986.

- D. Goldfarg and S. Ma. Fast multiple splitting algorithms for convex optimization. Technical report, 2009. Preprint arXiv:0912.4570v1.
- H. Groenevelt. Two algorithms for maximizing a separable concave function over a polymatroid feasible region. *European Journal of Operations Research*, pages 227–236, 1991.
- L. He and L. Carin. Exploiting structure in wavelet-based Bayesian compressive sensing. *IEEE Transactions on Signal Processing*, 57(9):3488–3497, 2009.
- D. S. Hochbaum and S. P. Hong. About strongly polynomial time algorithms for quadratic optimization over submodular constraints. *Mathematical Programming*, 69(1):269–309, 1995.
- H. Hoefling. A path algorithm for the fused lasso signal approximator. *Journal of Computational and Graphical Statistics*, 19(4):984–1006, 2010.
- R. A. Horn and C. R. Johnson. *Matrix analysis*. Cambridge University Press, 1990.
- J. Huang and T. Zhang. The benefit of group sparsity. *Annals of Statistics*, 38(4):1978–2004, 2010.
- J. Huang, Z. Zhang, and D. Metaxas. Learning with structured sparsity. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2009.
- A. Hyvärinen, P. Hoyer, and M. Inki. Topographic independent component analysis. *Neural Computation*, 13(7):1527–1558, 2001.
- L. Jacob, G. Obozinski, and J.-P. Vert. Group Lasso with overlap and graph Lasso. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2009.
- R. Jenatton, J.-Y. Audibert, and F. Bach. Structured variable selection with sparsity-inducing norms. Technical report, 2009. Preprint arXiv:0904.3523v3.
- R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for sparse hierarchical dictionary learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010a.
- R. Jenatton, G. Obozinski, and F. Bach. Structured sparse principal component analysis. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010b.
- R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for hierarchical sparse coding. *Journal of Machine Learning Research*, 12:2297–2334, 2011.
- K. Kavukcuoglu, M. Ranzato, R. Fergus, and Y. LeCun. Learning invariant features through topographic filter maps. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- S. Kim and E. P. Xing. Tree-guided group Lasso for multi-task regression with structured sparsity. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010.
- N. Maculan and J. R. G. Galdino de Paula. A linear-time median-finding algorithm for projecting a vector on the simplex of R^n . *Operations Research Letters*, 8(4):219–222, 1989.

- M. W. Mahoney and P. Drineas. CUR matrix decompositions for improved data analysis. *Proceedings of the National Academy of Sciences*, 106(3):697, 2009.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2009.
- J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19–60, 2010a.
- J. Mairal, R. Jenatton, G. Obozinski, and F. Bach. Network flow algorithms for structured sparsity. In *Advances in Neural Information Processing Systems*, 2010b.
- S. Mallat. *A Wavelet Tour of Signal Processing, Second Edition*. Academic Press, New York, September 1999.
- C. A. Micchelli, J. M. Morales, and M. Pontil. A family of penalty functions for structured sparsity. In *Advances in Neural Information Processing Systems*, 2010.
- J. J. Moreau. Fonctions convexes duales et points proximaux dans un espace Hilbertien. *Compte Rendus de l'Académie des Sciences, Paris, Série A, Mathématiques*, 255:2897–2899, 1962.
- D. Needell and J. A. Tropp. CoSaMP: Iterative signal recovery from incomplete and inaccurate samples. *Applied and Computational Harmonic Analysis*, 26(3):301–321, 2009.
- S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems*, 2009.
- Y. Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152, 2005.
- Y. Nesterov. Gradient methods for minimizing composite objective function. Technical report, Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, 2007.
- G. Obozinski, B. Taskar, and M. I. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, 20(2):231–252, 2010.
- B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381:607–609, 1996.
- Z. Qin and D. Goldfarb. Structured sparsity via alternating directions methods. Technical report, 2011. preprint ArXiv:1105.0728.
- A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet. SimpleMKL. *Journal of Machine Learning Research*, 9:2491–2521, 2008.
- V. Roth and B. Fischer. The group-Lasso for generalized linear models: uniqueness of solutions and efficient algorithms. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.

- M. Schmidt and K. Murphy. Convex structure learning in log-linear models: Beyond pairwise potentials. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010.
- M. Schmidt, N. Le Roux, and F. Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. In *Advances in Neural Information Processing Systems*, 2011. to appear, preprint ArXiv:1109.2415v1.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- P. Sprechmann, I. Ramirez, G. Sapiro, and Y. C. Eldar. Collaborative hierarchical sparse modeling. Technical report, 2010. Preprint arXiv:1003.0400v1.
- M. Stojnic, F. Parvaresh, and B. Hassibi. On the reconstruction of block-sparse signals with an optimal number of measurements. *IEEE Transactions on Signal Processing*, 57(8):3075–3085, 2009.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society: Series B*, 58(1):267–288, 1996.
- R. Tomioka, T. Suzuki, and M. Sugiyama. Augmented Lagrangian methods for learning, selecting and combining features. In S. Sra, S. Nowozin, and S. J. Wright, editors, *Optimization for Machine Learning*. MIT Press, 2011. To appear.
- P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of optimization theory and applications*, 109(3):475–494, 2001.
- B. A. Turlach, W. N. Venables, and S. J. Wright. Simultaneous variable selection. *Technometrics*, 47(3):349–363, 2005.
- M. J. Wainwright. Sharp thresholds for noisy and high-dimensional recovery of sparsity using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Transactions on Information Theory*, 55(5): 2183–2202, May 2009.
- J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):210–227, 2009a.
- S. Wright, R. Nowak, and M. Figueiredo. Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing*, 57(7):2479–2493, 2009b.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B*, 68:49–67, 2006.
- X. Zhang, M. Burger, and S. Osher. A unified primal-dual algorithm framework based on Bregman iteration. *Journal of Scientific Computing*, 46(1):20–46, 2011.
- P. Zhao, G. Rocha, and B. Yu. The composite absolute penalties family for grouped and hierarchical variable selection. *Annals of Statistics*, 37(6A):3468–3497, 2009.