



Structured sparsity-inducing norms through submodular functions

Francis Bach

► To cite this version:

Francis Bach. Structured sparsity-inducing norms through submodular functions. 2010. hal-00511310v1

HAL Id: hal-00511310

<https://hal.science/hal-00511310v1>

Preprint submitted on 24 Aug 2010 (v1), last revised 12 Nov 2010 (v3)

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

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

Structured sparsity-inducing norms through submodular functions

Francis Bach
INRIA - Willow project-team
Laboratoire d'Informatique de l'Ecole Normale Supérieure
Paris, France
`francis.bach@ens.fr`

August 24, 2010

Abstract

Sparse methods for supervised learning aim at finding good linear predictors from as few variables as possible, i.e., with small cardinality of their supports. This combinatorial selection problem is often turned into a convex optimization problem by replacing the cardinality function by its convex envelope (tightest convex lower bound), in this case the ℓ_1 -norm. In this paper, we investigate more general set-functions than the cardinality, that may incorporate prior knowledge or structural constraints which are common in many applications: namely, we show that for nonincreasing submodular set-functions, the corresponding convex envelope can be obtained from its Lovász extension, a common tool in submodular analysis. This defines a family of polyhedral norms, for which we provide generic algorithmic tools (subgradients and proximal operators) and theoretical results (conditions for support recovery or high-dimensional inference). By selecting specific submodular functions, we can give a new interpretation to known norms, such as those based on rank-statistics or grouped norms with potentially overlapping groups; we also define new norms, in particular ones that can be used as non-factorial priors for supervised learning.

1 Introduction

The concept of parsimony is central in many scientific domains. In the context of statistics, signal processing or machine learning, it takes the form of variable or feature selection problems, and is commonly used in two situations: First, to make the model or the prediction more interpretable or cheaper to use, i.e., even if the underlying problem is not sparse, one looks for the best sparse approximation. Second, sparsity can also be used given prior knowledge that the model should be sparse. In these two situations, reducing parsimony to finding the model of lowest cardinality turns out to be limiting, and structured parsimony has emerged as a fruitful practical extension, with applications to image processing, text processing or bioinformatics (see, e.g., [1, 2, 3, 4, 5, 6] and Section 4).

Most of the work based on convex optimization and the design of dedicated sparsity-inducing norms has focused mainly on the specific allowed set of sparsity patterns [1, 2, 4, 6]: if $w \in \mathbb{R}^p$ denotes the predictor we aim to estimate, and $\text{Supp}(w)$ denotes its support, then these norms are designed so that penalizing with these norms only leads to supports from a given family of allowed patterns. In this paper, we instead follow the approach of [3] and consider specific penalty functions $F(\text{Supp}(w))$

of the support set, which go beyond the cardinality function, but are not limited or designed to only forbid certain sparsity patterns. As shown in Section 6.2, these may also lead to restricted sets of supports but their interpretation in terms of an *explicit* penalty on the support leads to additional insights into the behavior of structured sparsity-inducing norms (see, e.g., Section 4.1). While direct greedy approaches (forward selection) to the problem are considered in [3], we provide convex relaxations to the function $w \mapsto F(\text{Supp}(w))$, which extend the traditional link between the ℓ_1 -norm and the cardinality function.

This is done for a particular ensemble of set-functions F , namely *non-increasing submodular functions*. Submodular functions may be seen as the set-function equivalent of convex functions [7], and exhibit many interesting properties that we review in Section 2 (see [8, 9] for other applications to machine learning). This paper makes the following contributions:

- We make explicit links between submodularity and sparsity by showing that the convex envelope of the function $w \mapsto F(\text{Supp}(w))$ on the ℓ_∞ -ball may be readily obtained from the Lovász extension of the submodular function (Section 3).

- We provide generic algorithmic tools, i.e., subgradients and proximal operators (Section 5), as well as theoretical results, i.e., conditions for support recovery or high-dimensional inference (Section 6), that extend classical results for the ℓ_1 -norm and show that many norms may be tackled by the exact same analysis and algorithms.

- By selecting specific submodular functions in Section 4, we recover and give a new interpretation to known norms, such as those based on rank-statistics or grouped norms with potentially overlapping groups [1, 2], and we define new norms, in particular ones that can be used as non-factorial priors for supervised learning (Section 4). These are illustrated on simulation experiments in Section 7, where they outperform related greedy approaches [3].

Notation. For $w \in \mathbb{R}^p$, $\text{Supp}(w) \subset V = \{1, \dots, p\}$ denotes the support of w , defined as $\text{Supp}(w) = \{j \in V, w_j \neq 0\}$. For $w \in \mathbb{R}^p$ and $q \geq 1$, we denote by $\|w\|_q$ the ℓ_q -norm of w , defined as $\|w\|_q^q = \sum_{i=1}^p |w_i|^q$, and by $\|w\|_\infty = \max_{i \in \{1, \dots, p\}} |w_i|$ its ℓ_∞ -norm. We denote by $|w|$ the vector of absolute values of the components of w . Moreover, given a vector w and a matrix Q , w_A and Q_{AA} are the corresponding subvector and submatrices of w and Q . Finally, we use the common notation from submodular analysis: for $w \in \mathbb{R}^p$ and $A \subset V$, $w(A) = \sum_{k \in A} w_k$ (this defines a modular function).

2 Review of submodular function theory

Throughout this paper, we consider a *nondecreasing submodular* function F defined on the power set 2^V of $V = \{1, \dots, p\}$, i.e., such that:

$$\begin{aligned} \forall A, B \subset V, \quad F(A) + F(B) &\geq F(A \cup B) + F(A \cap B), & (\text{submodularity}) \\ \forall A, B \subset V, \quad A \subset B &\Rightarrow F(A) \leq F(B). & (\text{monotonicity}) \end{aligned}$$

Moreover, we assume (without loss of generality) that $F(\emptyset) = 0$. These set-functions are often referred to as *polymatroid set-functions* [10] or β -functions [11]. Also, without loss of generality, we may assume that F is strictly positive on singletons, i.e., for all $k \in V$, $F(\{k\}) > 0$. Indeed, if $F(\{k\}) = 0$, then by submodularity and monotonicity, if $A \ni k$, $F(A) = F(A \setminus \{k\})$ and thus we can simply consider $V \setminus \{k\}$ instead of V .

Classical examples are the cardinality function (which will lead to the ℓ_1 -norm) and, given a partition of V into $B_1 \cup \dots \cup B_k = V$, the set function $A \mapsto F(A)$ which is equal to the number of groups B_1, \dots, B_k with non empty intersection with A (which will lead to the grouped ℓ_1/ℓ_∞ -norm [1, 12]).

Lovász extension. Given any set-function F , one can define its *Lovász extension* [7] (a.k.a. *Choquet integral* [13]) $f : \mathbb{R}_+^p \rightarrow \mathbb{R}$, as follows: given $w \in \mathbb{R}_+^p$, we can order the components of w in decreasing order $w_{j_1} \geq \dots \geq w_{j_p} \geq 0$; the value $f(w)$ is then defined as:

$$f(w) = \sum_{k=1}^p w_{j_k} [F(\{j_1, \dots, j_k\}) - F(\{j_1, \dots, j_{k-1}\})]. \quad (1)$$

Note that if some of the components of w are equal, all orderings lead to the same value of $f(w)$. The Lovász extension f is always piecewise-linear, and when F is submodular, it is also convex (see, e.g., [7, 10]). Moreover, for all $\delta \in \{0, 1\}^p$, $f(\delta) = F(\text{Supp}(\delta))$: f is indeed an extension from vectors in $\{0, 1\}^p$ (which can be identified with indicator vectors of sets) to all vectors in \mathbb{R}_+^p . Moreover, it turns out that minimizing F over subsets, i.e., minimizing f over $\{0, 1\}^p$ is equivalent to minimizing f over $[0, 1]^p$ [7, 11].

Submodular polyhedron and greedy algorithm. We denote by \mathcal{P} the *submodular polyhedron* [10], defined as the set of $s \in \mathbb{R}_+^p$ such that for all $A \subset V$, $s(A) \leq F(A)$, i.e., $\mathcal{P} = \{s \in \mathbb{R}_+^p, \forall A \subset V, s(A) \leq F(A)\}$, where we use the notation $s(A) = \sum_{k \in A} s_k$. One important result in submodular analysis is that if F is a nondecreasing submodular function, then we have a representation of f as a maximum of linear functions [10, 7], i.e., for all $w \in \mathbb{R}_+^p$,

$$f(w) = \max_{s \in \mathcal{P}} w^\top s. \quad (2)$$

Instead of solving a linear program with $p + 2^p$ constraints, a solution s may be obtained by the following “greedy algorithm”: order the components of w in decreasing order $w_{j_1} \geq \dots \geq w_{j_p}$, and then take for all $k \in \{1, \dots, p\}$, $s_{j_k} = F(\{j_1, \dots, j_k\}) - F(\{j_1, \dots, j_{k-1}\})$.

Stable sets. A set A is said *stable* if it cannot be augmented without increasing F , i.e., if for all sets $B \supset A$, $B \neq A \Rightarrow F(B) > F(A)$. If F is strictly increasing, then all sets are stable. Stable sets are also sometimes referred to as *flat* or *closed* [11]. The set of stable sets is closed by intersection [11], and will correspond to the set of allowed sparsity patterns (see Section 6.2). For the cardinality function, all sets are stable.

Separable sets. A set A is separable if we can find a partition of A into $A = B_1 \cup \dots \cup B_k$ such that $F(A) = F(B_1) + \dots + F(B_k)$. A set A is inseparable if it is not separable. As shown in [11], the submodular polytope \mathcal{P} has full dimension p as soon as F is strictly positive on all singletons, and its faces are exactly the sets $\{s_k = 0\}$ for $k \in V$ and $\{s(A) = F(A)\}$ for stable *and* inseparable sets. We let denote \mathcal{T} the set of such sets. This implies that $\mathcal{P} = \{s \in \mathbb{R}_+^p, \forall A \in \mathcal{T}, s(A) \leq F(A)\}$. These stable inseparable sets will play a role when describing extreme points of unit balls of our new norms (Section 3) and for deriving concentration inequalities in Section 6.3. For the cardinality function, stable and inseparable sets are singletons.

Submodular function minimization. Submodular functions are particularly interesting because they can be minimized in polynomial time. In this paragraph, we consider a non-monotonic submodular function G (otherwise finding the minimum is trivial). Most algorithms for minimizing submodular functions rely on the following strong duality principle [11, 10]:

$$\min_{A \subset V} G(A) = \max_{s \in \mathcal{B}(G)} \sum_{k \in V} \min\{0, s_k\}, \quad (3)$$

where $\mathcal{B}(G) = \{s \in \mathbb{R}^p, \forall A \subset V, s(A) \leq G(A), s(V) = G(V)\}$ is referred to as the *base polyhedron*. Moreover, algorithms for minimizing G will usually output A and s such that $G(A) = \sum_{k \in V} \min\{0, s_k\}$ as a certificate for optimality. The two main types of algorithms are combinatorial

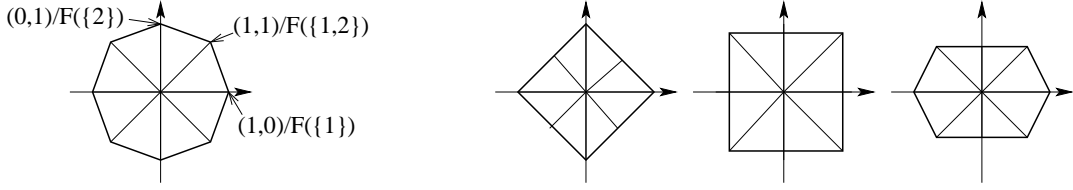


Figure 1: Polyhedral unit ball, for 4 different submodular functions (two variables), with different stable inseparable sets leading to different sets of extreme points; changing values of F may make some of the extreme points disappear. From left to right: $F(A) = |A|^{1/2}$ (all possible extreme points), $F(A) = |A|$ (leading to the ℓ_1 -norm), $F(A) = \min\{|A|, 1\}$ (leading to the ℓ_∞ -norm), $F(A) = 1_{\{A \cap \{1\} \neq \emptyset\}} + \frac{1}{2} 1_{\{A \neq \emptyset\}}$ (leading to the structured norm $\Omega(w) = |w_1| + \frac{1}{2} \|w\|_\infty$).

algorithms (that explicitly looks for A) and ones based on convex optimization (that explicitly looks for s). The first type of algorithm leads to strongly polynomial algorithms with best known complexity $O(p^6)$ [14], while the minimum point algorithm of [10] has no worst-time complexity bounds but is usually much faster in practice [10] and is based on the equivalent problem of finding the minimum-norm point in $\mathcal{B}(G)$, i.e., $\min_{s \in \mathcal{B}(G)} \|s\|_2^2$.

3 Definition and properties of structured norms

We define the function $\Omega(w) = f(|w|)$, where $|w|$ is the vector in \mathbb{R}^p composed of absolute values of w and f the Lovász extension of F . We have the following properties (see proof in appendix), which show that we indeed define a norm and that it is the desired convex envelope:

Proposition 1 (Convex envelope, dual norm) *Assume that the set function F is submodular, nondecreasing, is equal to zero for the empty set, and strictly positive for all singletons. Define the function $\Omega : w \mapsto f(|w|)$. Then, the following properties hold:*

- (i) Ω is a norm on \mathbb{R}^p ,
- (ii) Ω is the convex envelope of the function $g : w \mapsto F(\text{Supp}(w))$ on the unit ℓ_∞ -ball,
- (iii) the dual norm (see, e.g., [15]) of Ω is equal to $\Omega^*(s) = \max_{A \subset V} \frac{\|s_A\|_1}{F(A)} = \max_{A \in \mathcal{T}} \frac{\|s_A\|_1}{F(A)}$.

We provide examples of submodular set-functions and norms in Section 4, where we go from set-functions to norms, and vice-versa. From the definition of the Lovász extension in Eq. (1), we see that Ω is a polyhedral norm (i.e., its unit ball is a polyhedron). The following proposition gives the set of extreme points of the unit ball (see proof in appendix):

Proposition 2 (Extreme points of unit ball) *The extreme points of the unit ball of Ω are the vectors $\frac{1}{F(A)}s$, with $s \in \{-1, 0, 1\}^p$, $\text{Supp}(s) = A$ and A a stable inseparable set.*

This proposition shows, that depending on the number and cardinality of the inseparable stable sets, we can go from $2p$ (only singletons) to $3^p - 1$ (all possible subsets) extreme points. We show in Figure 1 examples of balls and sets of extreme points. These extreme points will play a role in concentration inequalities derived in Section 6.

Alternative representation. We can also give another interpretation to the norm Ω , as follows (following [4]): we consider the decomposition of w as $w = \sum_{A \subset V} w^A$, where each w^A has support included in A , i.e., $\text{Supp}(w^A) \subset A$;

Proposition 3 (Alternative representation) *The function $\Omega : w \mapsto f(|w|)$ is equal to the minimum norm of a decomposition of x into components supported by all subsets of V , i.e.,*

$$\min_{w=\sum_{A\subset V} w^A} \sum_{A\subset V} F(A) \|w^A\|_\infty.$$

Proof Simply considering convex duality, we get:

$$\begin{aligned} \min_{w=\sum_{A\subset V} w^A} \sum_{A\subset V} F(A) \|w^A\|_\infty &= \min_{w^A, A\subset V} \max_{s\in\mathbb{R}^p} \sum_{A\subset V} F(A) \|w^A\|_\infty + s^\top (w - \sum_{A\subset V} w^A) \\ &= \max_{s\in\mathbb{R}^p} \min_{w^A, A\subset V} \sum_{A\subset V} F(A) \|w^A\|_\infty + s^\top (w - \sum_{A\subset V} w^A) \\ &= \max_{s\in\mathbb{R}^p} \min_{w^A, A\subset V} \sum_{A\subset V} \left\{ F(A) \|w^A\|_\infty - s^\top w^A \right\} + s^\top w \\ &= \max_{\forall A\subset V, \|s_A\|_1 \leq F(A)} s^\top w = \max_{\Omega^*(s) \leq 1} s^\top w = \Omega(w) \end{aligned}$$

Note that it can be limited to stable inseparable sets. Following [4], it seems that it would lead to the result that the set of allowed sparsity patterns are unions of such stable and inseparable sets. This is not the case, as shown in Section 6. ■

4 Examples of nondecreasing submodular functions

We consider three main types of submodular functions with motivation for regularization for supervised learning. Some existing norms are shown to be examples of our frameworks (Section 4.1, Section 4.3), while other novel norms are designed from specific submodular functions (Section 4.2). Other examples of submodular functions, in particular in terms of matroids and entropies, may be found in [10, 8, 9] and could also lead to interesting new norms. Note that set covers, which are common examples of submodular functions are subcases of set-functions defined in Section 4.1.

4.1 Norms defined with non-overlapping or overlapping groups

We consider grouped norms defined with potentially overlapping groups [1, 2], i.e.,

$$\Omega(w) = \sum_{G\subset V} d(G) \|x_G\|_\infty,$$

where d is a nonnegative set function (with potentially $d(G) = 0$ when G should not be considered in the norm). It is a norm as soon as $\cup_{A, d(A)>0} A = V$ and it corresponds to the non-decreasing submodular function $F(A) = \sum_{G\cap A \neq \emptyset} d(G)$. In the case where ℓ_∞ -norms are replaced by ℓ_2 -norms, [2] has shown that the set of allowed sparsity patterns are unions of complements of groups G with positive weights. These sets happen to be the set of stable sets for the corresponding submodular function; thus the analysis provided in Section 6.2 extends the result of [2] to the new case of ℓ_∞ -norms. However, in our situation, we can give a reinterpretation through a submodular function that counts the number of times the support A intersects groups G with non zero weights. This goes beyond restricting the set of allowed sparsity patterns to stable sets. We show later in this section some insights gained by this reinterpretation. We now give some examples of norms, with various topologies of groups:



Figure 2: Sequence and groups: (left) groups for contiguous patterns, (right) groups for penalizing the number of jumps in the indicator vector sequence.

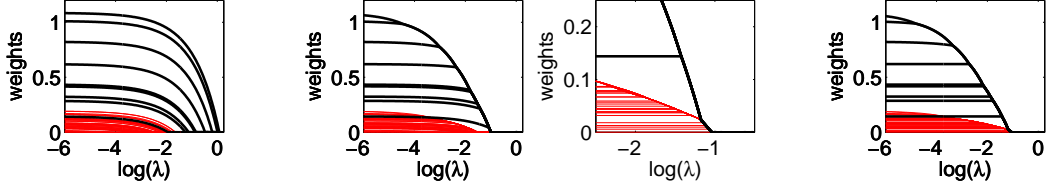


Figure 3: Regularization path for a penalized least-squares problem (black: variables that should be active, red: variables that should be left out). From left to right: ℓ_1 -norm penalization, polyhedral norm for rectangles in 2D (with zoom), mix of the two norms.

Hierarchical norms. Hierarchical norms defined on directed acyclic graphs [1, 5, 6] correspond to the set-function $F(A)$ which is the cardinality of the union of ancestors of elements in A . These have been applied to bioinformatics [5], computer vision [16] and topic models [6].

Norms defined on grids. If we assume that the p variables are organized in a 1D, 2D or 3D grid, [2] considers norms based on overlapping groups leading to stable sets equal to rectangular or convex shapes, with applications in computer vision [16]. For example, for the groups defined in the left side Figure 2 (with unit weights), we have $F(A) = p + \text{range}(A)$ if $A \neq \emptyset$ and $F(A) = 0$ if $A = \emptyset$ (the range of A is equal to $\max(A) - \min(A) + 1$). From empty sets to non-empty sets, there is a gap of $p + 1$, which is larger than differences among non-empty sets. This leads to the undesired result, which has been already observed by [2] of adding all variables in one step, rather than gradually, when the regularization parameter decreases in a regularized optimization problem. In order to counterbalance this effect, adding a constant times the cardinality function has the effect of making the first gap relatively smaller. This corresponds to adding a constant times the ℓ_1 -norm and, as shown in Figure 3, solves the problem of having all variables coming together. All patterns are then allowed, but contiguous ones are *encouraged rather than forced*.

Another interesting new norm may be defined from the groups in the right side of Figure 2. Indeed, it corresponds to the function $F(A)$ equal to $|A|$ plus the number of intervals of A . Note that this also favors contiguous patterns but is not limited to selecting a single interval (like the norm obtained from groups in the left side of Figure 2). Note that it is to be contrasted with the total variation (a.k.a. fused Lasso penalty [17]), which is a relaxation of the number of jumps in a vector w rather than in its support. In 2D or 3D, this extends to the notion of perimeter and area, but we do not pursue such extensions here.

4.2 Spectral functions of submatrices

Given a positive semidefinite matrix $Q \in \mathbb{R}^{p \times p}$ and a real-valued function h from $\mathbb{R}_+ \rightarrow \mathbb{R}$, one may define $\text{tr}[h(Q)]$ as $\sum_{i=1}^p h(\lambda_i)$ where $\lambda_1, \dots, \lambda_p$ are the (nonnegative) eigenvalues of Q [18]. We can thus define the function $F(A) = \text{tr} h(Q_{AA})$ for $A \subset V$. If Q is diagonal, then the function f is a weighted cardinal function, and it is well-known that the concavity of h is a sufficient condition for submodularity of F [10]. However, when Q is not diagonal, the concavity of h is not sufficient (as can be seen by generating random examples with $h(\lambda) = \lambda/(\lambda + 1)$).

We know however from [10] that the functions $h(\lambda) = \log(\lambda + t)$ for $t \geq 0$ lead to submodular functions; thus, since for $p \in (0, 1)$, $\lambda^p = \frac{p \sin p\pi}{\pi} \int_0^\infty \log(1 + \lambda/t) t^{p-1} dt$ (see, e.g., [19]), $h(\lambda) = \lambda^p$ for $p \in (0, 1]$ are positive linear combinations of functions that lead to nondecreasing submodular functions. Thus, they are also nondecreasing submodular functions, and, to the best of our knowledge, provide novel examples of such functions.

In the context of supervised learning from a design matrix X , we naturally use $Q = X^\top X$. If h is linear, then $F(A) = \text{tr } X_A^\top X_A$ (where X_A denotes the submatrix of X with columns in A) and we obtain a weighted cardinality function and hence a weighted ℓ_1 -norm, which is a *factorial prior*, i.e., it is a sum of terms depending on each variable independently.

In the frequentist setting, the Mallows C_L penalty [20] depends on the degrees of freedom, of the form $\text{tr } X_A^\top X_A (X_A^\top X_A + I)^{-1}$. This is a non-factorial prior but unfortunately it does not lead to a submodular function. In a Bayesian context however, it is shown by [21] that penalties of the form $\log \det(X_A^\top X_A + I)$ (which lead to submodular functions) correspond to marginal likelihoods associated to the set A and has good behavior when used within a non-convex framework.

This highlights the need for non-factorial priors which are sub-linear functions of the eigenvalues of $X_A^\top X_A$, which is exactly what nondecreasing submodular function of submatrices are. We do not pursue the extensive evaluation of non-factorial convex priors in this paper but provide in simulations examples with $F(A) = \text{tr}(X_A^\top X_A)^{1/2}$ (which is equal to the trace norm of X_A [15]).

4.3 Functions of cardinality

For $F(A) = h(|A|)$ where h is nondecreasing, such that $h(0) = 0$ and concave, then, from Eq. (1), $\Omega(x)$ defined from the rank statistics of $|x| \in \mathbb{R}_+^p$, i.e., if $|x_{(1)}| \geq |x_{(2)}| \geq \dots \geq |x_{(p)}|$, then $\Omega(x) = \sum_{k=1}^p [h(k) - h(k-1)] |x_{(k)}|$. This includes the sum of the q largest elements, and might lead to interesting new norms for unstructured variable selection but this is not pursued here. However, the algorithms and analysis presented in Section 5 and Section 6 apply to this case.

5 Convex analysis and optimization

In this section we provide analytic and algorithmic tools related to optimization problems based on the regularization by our novel sparsity-inducing norms. Note that since these norms are polyhedral norms with unit balls having potentially exponential number of vertices or faces, regular linear programming toolboxes may not be used.

Subgradient. From $\Omega(w) = \max_{s \in \mathcal{P}} s^\top |w|$ and the greedy algorithm¹ presented in Section 2, one can easily get in *polynomial time* one subgradient as one of the maximizers s . This allows to use subgradient descent [22], with, as shown in Figure 4, slow convergence.

Proximal operator. Given regularized problems of the form $\min_{w \in \mathbb{R}^p} L(w) + \lambda \Omega(w)$, where L is differentiable with Lipschitz-continuous gradient, *proximal methods* have been shown to be particularly efficient first-order methods (see, e.g., [23]). In this paper we consider the methods “ISTA” and its accelerated variants “FISTA” [23], which are compared in Figure 4.

To apply these methods, it suffices to be able to solve efficiently problems of the form: $\min_{w \in \mathbb{R}^p} \frac{1}{2} \|w - z\|_2^2 + \lambda \Omega(w)$. In the case of the ℓ_1 -norm, this reduces to soft thresholding of z , the following proposition (see proof in the appendix) shows that it is equivalent to submodular function minimization,

¹The greedy algorithm to find extreme points of the submodular polyhedron should not be confused with the greedy algorithm (e.g., forward selection) that we consider in Section 7.

and thus that we can use existing algorithms (in simulations, we have use the minimum-norm point algorithm, which is empirically faster than algorithms with complexity bounds [10]):

Proposition 4 (Proximal operator) *Let $z \in \mathbb{R}^p$ and $\lambda > 0$, minimizing $\frac{1}{2}\|w - z\|_2^2 + \lambda\Omega(w)$ is equivalent to finding the minimum of the submodular function $A \mapsto \lambda F(A) - |z|(A)$.*

In the proof of the last proposition, it is shown how a solution for one problem may be obtained from a solution to the other problem. Note that using the minimum-norm point algorithm leads to a *generic* algorithms that can be applied to *any* submodular functions F , and that it may be rather inefficient for simpler subcases (e.g., the ℓ_1 -norm, grouped norms with non-overlapping groups [24], or hierarchical norms [6]).

Note that similar links between convex optimization and minimization of submodular functions have been considered (see, e.g., [25]). However, these are dedicated to symmetric submodular functions (such as the ones obtained from graph cuts) and are thus not directly applicable to our situation of non-increasing submodular functions.

6 Sparsity-inducing properties

In this section, we consider a fixed design matrix $X \in \mathbb{R}^{n \times p}$ and $y \in \mathbb{R}^n$ a vector of random responses. Given $\lambda > 0$, we define \hat{w} as a minimizer of the regularized least-squares cost:

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

We study the sparsity-inducing properties of solutions of Eq. (6), i.e., we determine in Section 6.2 which patterns are allowed and in Section 6.3 which sufficient conditions lead to correct estimation. Like regular analysis of sparsity-inducing norms, the analysis provided in this section relies heavily on decomposability properties of our norm Ω

6.1 Decomposability

For a subset J of V , we denote by $F_J : 2^J \rightarrow \mathbb{R}$ the *restriction* of F to J , defined for $A \subset J$ by $F_J(A) = F(A)$, and by $F^J : 2^{J^c} \rightarrow \mathbb{R}$ the *contraction* of F by J , defined for $A \subset J^c$ by $F^J(A) = F(A \cup J) - F(A)$. These two functions are submodular and nondecreasing as soon as F is (see, e.g., [10]).

We denote by Ω_J the norm on \mathbb{R}^J defined through the submodular function F_J , and Ω^J the pseudo-norm defined on \mathbb{R}^{J^c} defined through F^J (as shown in Proposition 5, it is a norm only when J is a stable set). Note that Ω_{J^c} (a norm on J^c) is in general different from Ω^J . Moreover, $\Omega_J(x_J)$ is actually equal to $\Omega(\tilde{x})$ where $\tilde{x}_J = x_J$ and $\tilde{x}_{J^c} = 0$, i.e., it is the restriction of Ω to J .

We can now prove the following decomposition property, which shows that under certain circumstances, we can decompose the norm Ω on subsets J and their complements:

Proposition 5 (Decomposition) *Given $J \subset V$ and Ω_J and Ω^J defined as above, we have:*

- (i) $\forall w \in \mathbb{R}^p$, $\Omega(w) \geq \Omega_J(w_J) + \Omega^J(w_{J^c})$,
- (ii) $\forall w \in \mathbb{R}^p$, and $\forall J \subset V$, if $\min_{j \in J} |w_j| \geq \max_{j \in J^c} |w_j|$, then $\Omega(w) = \Omega_J(w_J) + \Omega^J(w_{J^c})$,
- (iii) Ω^J is a norm on \mathbb{R}^{J^c} if and only if J is a stable set.

6.2 Sparsity patterns

In this section, we do not make any assumptions regarding the correct specification of the linear model (see proof in the appendix):

Proposition 6 (Stable sparsity patterns) *Assume $y \in \mathbb{R}^n$ has an absolutely continuous density with respect to the Lebesgue measure and that $X^\top X$ is invertible. Then the minimizer \hat{w} of Eq. (6) is unique and, with probability one, its support $\text{Supp}(\hat{w})$ is a stable set.*

For simplicity, we have assumed invertibility of $X^\top X$, which forbids the high-dimensional situation $p \geq n$, but we could extend to the type of assumptions used in [2]. However, the theoretical analysis we now provide is applicable to high-dimensional inference.

6.3 High-dimensional inference

We now assume that the linear model is well-specified and extend results from [26] for sufficient support recovery conditions and from [27] for estimation consistency. As seen in Proposition 5, the norm Ω is decomposable and we use this property extensively in this section. We denote by $\rho(J) = \min_{B \subset J^c} \frac{F(B \cup J) - F(J)}{F(B)}$; by submodularity and monotonicity of F , $\rho(J)$ is always between zero and one, and, as soon as J is stable it is strictly positive (for the ℓ_1 -norm, $\rho(J) = 1$). Moreover, we denote by $c(J) = \sup_{w \in \mathbb{R}^p} \frac{\Omega_J(w_J)}{\|w_J\|_2}$, the equivalence constant between the norm Ω_J and the ℓ_2 -norm. We always have $c(J) \leq |J|^{1/2} \max_{k \in V} F(\{k\})$ (with equality for the ℓ_1 -norm).

The following propositions allow us to get back and extend well-known results for the ℓ_1 -norm, i.e., Propositions 7 and 9 extend results based on support recovery conditions [26]; while Propositions 8 and 9 extend results based on restricted eigenvalue conditions [27]. We can also get back results for the ℓ_1/ℓ_∞ -norm [12]. As shown in the appendix, proof techniques are similar and are adapted through the decomposition properties from Proposition 5.

Proposition 7 (Support recovery) *Assume that $y = Xw^* + \sigma\varepsilon$, where ε is a standard multivariate normal vector. Let $Q = \frac{1}{n}X^\top X \in \mathbb{R}^{p \times p}$. Denote by J the smallest stable set containing the support $\text{Supp}(w^*)$ of w^* . Define $\nu = \min_{j, w_j^* \neq 0} |w_j^*| > 0$, assume $\kappa = \lambda_{\min}(Q_{JJ}) > 0$ and that for $\eta > 0$,*

$$(\Omega^J)^*[(\Omega_J(Q_{JJ}^{-1}Q_{Jj}))_{j \in J^c}] \leq 1 - \eta. \quad (5)$$

Then, if $\lambda \leq \frac{\kappa\nu}{2c(J)}$, the minimizer \hat{w} is unique and has support equal to J , with probability larger than $P(\Omega^(z) \leq \frac{\lambda\eta\rho(J)}{2\sigma\sqrt{n}})$, where z is multivariate normal with covariance matrix Q .*

Proposition 8 (Consistency) *Assume that $y = Xw^* + \sigma\varepsilon$, where ε is a standard multivariate normal vector. Let $Q = \frac{1}{n}X^\top X \in \mathbb{R}^{p \times p}$. Denote by J the smallest stable set containing the support $\text{Supp}(w^*)$ of w^* . Assume that for all Δ such that $\Omega^J(\Delta_{J^c}) \leq 3\Omega_J(\Delta_J)$, then $\Delta^\top Q \Delta \geq \kappa \|\Delta_J\|_2^2$. Then, we have $\Omega(\Delta) \leq \frac{24c(J)^2\lambda}{\kappa\rho(J)^2}$ and $\Delta^\top Q \Delta \leq \frac{36c(J)^2\lambda^2}{\kappa\rho(J)^2}$, with probability larger than $P(\Omega^*(z) \leq \frac{\lambda\rho(J)}{2\sigma\sqrt{n}})$ where z is multivariate normal with covariance matrix Q .*

Proposition 9 (Concentration inequalities) *Let z be a random normal variable with covariance matrix Q . Let \mathcal{T} be the set of stable inseparable sets. Then*

$$P(\Omega^*(z) > t) \leq \sum_{A \in \mathcal{T}} 2^{|A|} \exp\left(-\frac{t^2 F(A)}{2 \times 1^\top Q_{AA} 1}\right).$$

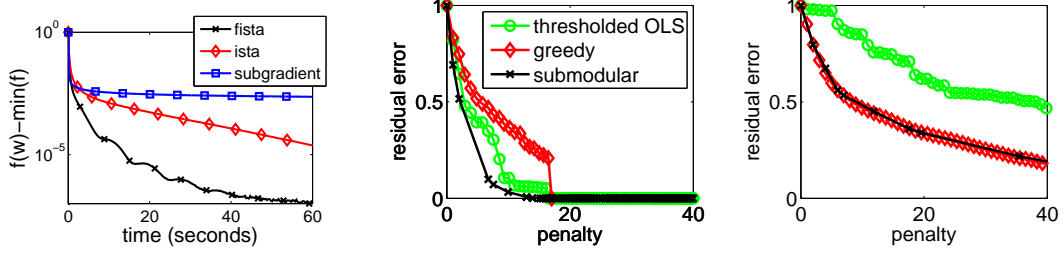


Figure 4: (Left) Comparison of iterative optimization algorithms (value of objective function vs. running time). (Middle/Right) Best subset selection with structured penalty: (middle) high-dimensional case ($p = 120, n = 20, k = 40$), (right) lower-dimensional case ($p = 120, n = 120, k = 40$).

7 Experiments

In this section, we provide illustrations on toy examples of some of the results presented in the paper. We consider the regularized least-squares problem of Eq. (6), with data generated as follows: given p, n, k , the design matrix $X \in \mathbb{R}^{n \times p}$ is a matrix of i.i.d. standard Gaussian components, normalized to have unit ℓ_2 -norm columns. A set J of cardinality k is chosen at random and the weights w_J^* are sampled from a standard multivariate Gaussian distribution and $w_{J^c}^* = 0$. We then generate $y = Xw^* + n^{-1/2}\|Xw^*\|_2 \times \varepsilon$ where ε is a standard Gaussian vector (this corresponds to a unit signal-to-noise ratio).

Proximal methods vs. subgradient descent. For the submodular function $F(A) = |A|^{1/2}$ (a simple submodular function beyond the cardinality) we compare three optimization algorithms described in Section 5, subgradient descent and two proximal methods, ISTA and its accelerated version FISTA [23], for $p = n = 1000, k = 100$ and $\lambda = 0.1$. Other settings and other set-functions would lead to similar results than the ones presented in Figure 4: FISTA is faster than ISTA, and much faster than subgradient descent.

Relaxation of combinatorial optimization problem. We compare three strategies for solving the combinatorial optimization problem $\min_{w \in \mathbb{R}^p} \frac{1}{2n}\|y - Xw\|_2^2 + \lambda F(\text{Supp}(w))$ with $F(A) = \text{tr}(X_A^\top X_A)^{1/2}$, the approach based on our sparsity-inducing norms, the simpler greedy (forward selection) approach proposed in [3], and by thresholding the ordinary least-squares estimate. For all methods, we try all possible regularization parameters. We see in the left plots of Figure 4 that for hard cases (middle plot) convex optimization techniques perform better than other approaches for the subset selection problem, while for easier cases where there are more observations (right plot), it does as well as greedy approaches.

Non factorial priors for variable selection. We perform similar experiments, but now focus on the predictive performance and compare our new norm with $F(A) = \text{tr}(X_A^\top X_A)^{1/2}$, with greedy approaches [3] and to regularization by ℓ_1 or ℓ_2 norms. As shown in Table 1, the new norm based on non-factorial priors is more robust than other approaches to low number of observations n and to lack of sparsity k .

8 Conclusions

We have presented a family of sparsity-inducing norms dedicated to incorporating prior knowledge or structural constraints on the support of linear predictors. We have provided a set of common algorithms and theoretical results, as well as simulations on synthetic examples illustrating the good behavior of these norms. Several avenues are worth investigating: first, we could follow current practice in sparse methods, e.g., by considering related adapted concave penalties to enhance

p	n	k	submodular	ℓ_2	ℓ_1	greedy
120	40	80	3.5	1.8	4.4	30.1
120	40	40	8.9	7.8	10.4	31.6
120	40	20	13.0	14.6	13.9	42.9
120	40	10	17.3	23.6	17.3	37.0
120	20	80	1.0	1.5	7.4	28.9
120	20	40	3.4	2.5	9.9	28.6
120	20	20	7.4	7.4	11.5	38.5
120	20	10	20.9	19.8	24.8	41.7

Table 1: Normalized mean-square errors $\frac{1}{n}\|X\hat{w} - Xw^*\|_2^2$ (multiplied by 100) with optimal regularization parameters (averaged over 5 replications).

sparsity-inducing norms, or by extending some of the concepts for norms of matrices, with potential applications in matrix factorization or multi-task learning (see, e.g., [28] for application of submodular functions to dictionary learning). Second, links between submodularity and sparsity could be studied further, in particular by considering submodular relaxations of other combinatorial functions, e.g., by exploring relationships with other polyhedral norms such as total variations, which are known to be similarly associated with non monotonic submodular set-functions such as graph cuts [29, 30, 17, 25].

A Properties of the norm

A.1 Proof of Proposition 1

(i) Ω is positively homogeneous by definition of the Lovász extension in Eq. (1), convex because of the representation in Eq. (2) as the maximum of $s^\top |w|$ for some $s \in \mathcal{P} \subset \mathbb{R}_+^p$, and it is a norm as soon as $\Omega(w) = 0$ implies that $w = 0$, which is true since $\Omega(w) \geq \min_k F(\{k\}) \|w\|_\infty$. (ii) We denote by g^* the Fenchel conjugate of g on the domain $\{w \in \mathbb{R}^p, \|w\|_\infty \leq 1\}$, and g^{**} its bidual [15]. By definition of the Fenchel conjugate, we have:

$$\begin{aligned} g^*(s) &= \max_{\|w\|_\infty \leq 1} w^\top s - g(w) \\ &= \max_{\delta \in \{0,1\}^p} \max_{\|w\|_\infty \leq 1} (\delta \circ w)^\top s - f(\delta) \\ &= \max_{\delta \in \{0,1\}^p} \delta^\top |s| - f(\delta) \\ &= \max_{\delta \in [0,1]^p} \delta^\top |s| - f(\delta) \text{ because } F - |s| \text{ is submodular} \end{aligned}$$

Thus, for all w such that $\|w\|_\infty \leq 1$,

$$\begin{aligned} g^{**}(w) &= \max_{s \in \mathbb{R}^p} s^\top w - g^*(s) \\ &= \max_{s \in \mathbb{R}^p} \min_{\delta \in [0,1]^p} s^\top w - \delta^\top |s| + f(\delta) \\ &= \min_{\delta \in [0,1]^p} \max_{s \in \mathbb{R}^p} s^\top w - \delta^\top |s| + f(\delta) \text{ by strong duality and Slater's condition [15]} \\ &= \min_{\delta \in [0,1]^p, \delta \geq |w|} f(\delta) = f(|w|) \text{ because } F \text{ is nonincreasing.} \end{aligned}$$

Note that F non-increasing implies that f is non-increasing with respect to all of its components. (ii) We have $\Omega(w) = f(|w|) = \max_{s \in \mathcal{P}} s^\top |w| = \max_{|s| \in \mathcal{P}} s^\top w = \max_{\|s_A\|_1 \leq F(A), A \subset V} s^\top w = \max_{\max_{A \subset V} \frac{\|s_A\|_1}{F(A)} \leq 1} s^\top w$, which implies the desired result. Note that the maximization may indeed be limited to the stable inseparable sets $A \in \mathcal{T}$.

A.2 Proof of Proposition 2

We have seen in Section 2 that for $A \in \mathcal{T}$ (set of stable inseparable sets), then $\{x(A) = F(A)\}$ is a face of \mathcal{P} (and those sets are the only ones for which this happens). We get to the desired result by considering potential different signs.

B Convex optimization results

We first prove an additional result related to decomposition of subdifferentials. Note that the exact subdifferential for the non-zero components of w is rather complicated when w has components with equal magnitude. If this is not the case, i.e., $|w_{j_1}| > \dots > |w_{j_k}| > 0$, where $k = |J|$, then the subdifferential $\partial\Omega_J(w_J)$ is reduced to a point s such that $s_{j_k} = F(\{j_1, \dots, j_k\}) - F(\{j_1, \dots, j_{k-1}\})$. For more details on the subdifferential for nonzero components, see [10].

Lemma 1 (Decomposition of subdifferential) *Let $w \in \mathbb{R}^p$, with support $J = \text{Supp}(w)$ and with H equal to the smallest stable set containing J . The subdifferential $\partial\Omega(w)$ at w , can then be decomposed as follows on $\mathbb{R}^V = \mathbb{R}^J \times \mathbb{R}^{H \setminus J} \times \mathbb{R}^{H^c}$: $\partial\Omega(w) = \partial\Omega_J(w_J) \times \{0\} \times \{s_{H^c}, (\Omega^H)^*(s_{H^c}) \leq 1\}$.*

Proof For all sufficient small $\Delta \in \mathbb{R}^p$, the components in $(w + \Delta)_J$ have all greater absolute values than the ones in $(w + \Delta)_{J^c}$. Thus, from Proposition 5, $\Omega(w + \Delta) = \Omega_J(w_J + \Delta_J) + \Omega^J(\Delta_{J^c}) = \Omega_J(w_J + \Delta_J) + \Omega^H(\Delta_{H^c})$, and thus the subdifferential decomposes as $\partial\Omega_J(w_J) \times \{0\} \times \partial\Omega^H(0)$. The subdifferential of a norm at zero is exactly the unit ball of the dual norm, which leads to the desired result. \blacksquare

B.1 Proof of Proposition 4

Following [6], without loss of generality, we assume that z has nonnegative components. We have by convex duality (which is applicable here because of Slater's condition):

$$\begin{aligned} \min_{w \in \mathbb{R}^p} \frac{1}{2} \|w - z\|_2^2 + \lambda \Omega(w) &= \min_{w \in \mathbb{R}^p} \max_{\Omega^*(s) \leq 1} \frac{1}{2} \|w - z\|_2^2 + \lambda s^\top w \\ &= \max_{\Omega^*(s) \leq 1} \min_{w \in \mathbb{R}^p} \frac{1}{2} \|w - z\|_2^2 + \lambda s^\top w \\ &= \max_{\Omega^*(s) \leq 1} \frac{1}{2} \|z\|_2^2 - \frac{1}{2} \|\lambda s - z\|_2^2, \end{aligned}$$

where the (unique) optimal w is obtained from the optimal s by $w = z - \lambda s$. s is defined constrained to satisfy $\Omega^*(s) \leq 1$, which is equivalent to $|s| \in \mathcal{P}$. Since z has nonnegative components, the minimum restricted to $|s| \in \mathcal{P}$ is the same as the minimum restricted to $s \in \mathcal{P}$, and also the same as the one restricted to the submodular polyhedron without constraints on positivity, i.e., our problem reduces to $\min_{\forall A \subset V, s(A) \leq F(A)} \|s - z/\lambda\|_2^2$, which is also equivalent to

$$\min_{\forall A \subset V, t(A) \leq F(A) - \lambda^{-1}z(A)} \|t\|_2^2,$$

which is exactly (up to the constraints $s(V) = F(V) - \lambda^{-1}z(V)$) the minimum-norm point problem for the submodular function $G : A \mapsto F(A) - \lambda^{-1}z(A)$. We can thus follow [10]: if (t, A) is a primal-dual optimal pair for the submodular function G (from Eq. (3)), then we can obtain s as $\lambda^{-1}z$ plus the negative part of t . From s we then get w through $w = z - \lambda s$.

C Sparse estimation

In this section, we consider a design $X \in \mathbb{R}^{n \times p}$ be a fixed design and $y \in \mathbb{R}^n$ a set of random responses. Given $\lambda > 0$, we define \hat{w} as a minimizer of the regularized least-squares cost:

$$\min_{w \in \mathbb{R}^p} \frac{1}{2n} \|y - Xw\|_2^2 + \lambda \Omega(w). \quad (6)$$

C.1 Proof of Proposition 5

(i) for $s \in \mathbb{R}_+^p$, if $\forall B \subset J$, $s(B) \leq F(B)$ and $\forall C \subset J^c$, $s(C) \leq F(C \cap J) - F(J)$, then $\forall A \subset V$, $s(A) = s(A \cap J) + s(A \cap J^c) \leq F(A \cap J) + F(A \cap J^c) - F(J) \leq F(A)$ by submodularity. This

implies that the desired result by considering the representation of the Lovász extension in Eq. (1) and the fact that we have just prove that \mathcal{P} contains the product of the two submodular polyhedra associated to F^J and F_J .

(ii) This is immediate from the expression of the Lovász extension in Eq. (1). Indeed, the order within J and the one within J^c do not interact. Note that this case includes cases where some of the components of $|w_J|$ are equal to some of $|w_{J^c}|$.

(iii) Ω^J corresponds to the submodular function obtained as the contraction of F by J . It is thus a norm as soon as F^J is positive on all singletons, which is itself equivalent to the stability of J . The equivalence of being a norm with stability of the set J is then straightforward.

C.2 Proof of Proposition 6

Let $Q = \frac{1}{n}X^\top X \in \mathbb{R}^{p \times p}$ and $r = \frac{1}{n}X^\top y \in \mathbb{R}^p$. The unicity of the minimizer \hat{w} is a consequence of the invertibility of $Q = \frac{1}{n}X^\top X$. Let $J \subset V$. We will show that if $\text{Supp}(\hat{w}) = J$, then \hat{w}_J is an affine function of r (and hence y), of the form $\hat{w}_J = (Q_{JJ}^{-1} - A_{JJ})r_J + b_J$, where $0 \preceq A_{JJ} \preceq Q_{JJ}^{-1}$ and (A_{JJ}, b_J) belongs to a finite set independent of r .

If J is not a stable set, then, by Proposition 1, this will implies that there exists $j \in J^c$ such that $Q_{jJ}[(Q_{JJ}^{-1} - A_{JJ})r_J + b_J] - r_j = 0$, i.e.,

$$0 = Q_{jJ}[(Q_{JJ}^{-1} - A_{JJ})r_J + b_J] - r_j = \frac{1}{n}[Q_{jJ}Q_{JJ}^{-1}X_J^\top - X_j^\top - Q_{jJ}A_{JJ}X_J^\top]y + Q_{jJ}b_J.$$

The row vector $Q_{jJ}Q_{JJ}^{-1}X_J^\top - X_j^\top - Q_{jJ}A_{JJ}X_J^\top$ cannot be equal to zero, otherwise,

$$0 = \frac{1}{n}[Q_{jJ}Q_{JJ}^{-1}X_J^\top - X_j^\top - Q_{jJ}A_{JJ}X_J^\top]X_j = Q_{jJ}Q_{JJ}^{-1}Q_{Jj} - Q_{jj} - Q_{jJ}A_{JJ}Q_{Jj} \leq Q_{jJ}Q_{JJ}^{-1}Q_{Jj} - Q_{jj}$$

which is a contradiction because of the invertibility of Q and the Schur complement lemma [31] (which implies that the previous quantity must be strictly negative). Thus, we have shown that if $\text{Supp}(\hat{w}) = J$ and J is not a stable subset, then for a finite number of non zero $(c, d) \in \mathbb{R}^n \times \mathbb{R}$, then $c^\top y$ is constant. This occurs with probability zero.

What remains to be shown is the affine representation of \hat{w}_J when the support is given; it is essentially equivalent to showing that the path is piecewise affine, which is not surprising for a polyhedral norm [32]. We use the representation $\Omega_J(w_J) = \max_{z \in B} z^\top w_J$ where B is the finite set of z such that $|z|$ in an extreme point of the submodular polyhedron associated with Ω_J .

Necessary optimality conditions [33] for such the problem in Eq. (6) is the existence of $\eta_z \geq 0$ (for each $z \in B$) such that (1) $\sum_{z \in B} \eta_z = 1$, (2) $\eta_z = 0$ if z is not a maximizer of $\max_{z \in B} z^\top w_J$, and (3) w_J is a minimizer of $\frac{1}{2}w_J^\top Q_{JJ}w_J - r_J^\top w_J + \lambda w_J^\top \sum_{z \in A} \eta_z z$, i.e., $Q_{JJ}w_J + \lambda \sum_{z \in A} \eta_z z = r_J$. Moreover, by Carathéodory's theorem [33], the number k of non-zero η may be taken to be less than $|J| + 1$.

This thus implies that, if consider the vector $\zeta \in \mathbb{R}^k$ of non-zero η , and the matrix $Z \in \mathbb{R}^{|J| \times k}$ of corresponding z 's, then we have

$$\begin{aligned} Q_{JJ}w_J + \lambda Z\zeta &= r_J \\ \zeta^\top 1 &= 1 \end{aligned}$$

$$\exists c \in \mathbb{R} \text{ such that } Z^\top w_J = c1.$$

In matrix form, this can be written as:

$$\begin{pmatrix} Q_{JJ} & \lambda Z & 0 \\ \lambda Z^\top & 0 & -\lambda 1 \\ 0 & -\lambda 1^\top & 0 \end{pmatrix} \begin{pmatrix} w_J \\ \zeta \\ c \end{pmatrix} = \begin{pmatrix} r_J \\ 0 \\ -\lambda \end{pmatrix}.$$

It is then a simple linear algebra exercise to show that if $k \leq |J| + 1$, then w_J is of the desired form.

C.3 Proof of Proposition 7

Let $q = \frac{1}{n}X^\top \varepsilon \in \mathbb{R}^p$, which is normal with mean zero and covariance matrix $\sigma^2 Q/n$. We have $\Omega(x) \geq \Omega_J(x_J) + \Omega^J(x_{J^c}) \geq \Omega_J(x_J) + \rho(J)\Omega_{J^c}(x_{J^c}) \geq \rho(J)\Omega(x)$. Thus, if we assume $\Omega^*(q) \leq \lambda\rho(J)\eta/2$, then $\Omega_J^*(q_J) \leq \lambda/2$ and $(\Omega^J)^*(q_{J^c}) \leq \lambda\eta/2$, which implies $(\Omega^J)^*(q_{J^c} - Q_{J^c J}Q_{JJ}^{-1}q_J) \leq \lambda\eta/2$, because the norm Ω^* is increasing with respect to each of its components.

We denote by \tilde{w} the unique (because Q_{JJ} is invertible) minimum of $\frac{1}{2n}\|y - Xw\|_2^2 + \lambda\Omega(w)$, subject to $w_{J^c} = 0$. \tilde{w}_J is defined through $Q_{JJ}(\tilde{w}_J - w_J^*) - q_J = -\lambda s_J$ where $s_J \in \partial\Omega_J(\tilde{w}_J)$ (which implies that $\Omega_J^*(s_J) \leq 1$), i.e., $\tilde{w}_J - w_J^* = Q_{JJ}^{-1}(q_J - \lambda s_J)$. We have:

$$\begin{aligned} \|\tilde{w}_J - w_J^*\|_\infty &\leq \max_{j \in J} |\delta_j^\top Q_{JJ}^{-1}(q_J - \lambda s_J)| \\ &\leq \max_{j \in J} \Omega_J(Q_{JJ}^{-1}\delta_j)\Omega_J^*(q_J - \lambda s_J) \\ &\leq \max_{j \in J} c(J)\|Q_{JJ}^{-1}\delta_j\|_2[\Omega_J^*(q_J) + \lambda\Omega_J^*(s_J)] \leq \frac{3}{2}\lambda c(J)\kappa^{-1} \end{aligned}$$

Thus if $2\lambda c(J)\kappa^{-1} \leq \nu$, then $\text{Supp}(\tilde{w}) \supset \text{Supp}(w^*)$.

We now show that since we have $(\Omega^J)^*(q_{J^c} - Q_{J^c J}Q_{JJ}^{-1}q_J) \leq \lambda\eta/2$, \tilde{w} is the unique minimizer of Eq. (6). For that it suffices to show that $(\Omega^J)^*(Q_{J^c J}(\tilde{w}_J - w_J^*) - q_{J^c}) < \lambda$. We have:

$$\begin{aligned} (\Omega^J)^*(Q_{J^c J}(\tilde{w}_J - w_J^*) - q_{J^c}) &= (\Omega^J)^*(Q_{J^c J}Q_{JJ}^{-1}(q_J - \lambda s_J) - q_{J^c}) \\ &\leq (\Omega^J)^*(Q_{J^c J}Q_{JJ}^{-1}q_J - q_{J^c}) + \lambda(\Omega^J)^*(Q_{J^c J}Q_{JJ}^{-1}s_J) \\ &\leq (\Omega^J)^*(Q_{J^c J}Q_{JJ}^{-1}q_J - q_{J^c}) + \lambda(\Omega^J)^*[(\Omega_J(Q_{JJ}^{-1}Q_{Jj}))_{j \in J^c}] \\ &\leq \lambda\eta/2 + \lambda(1 - \eta) < \lambda \end{aligned}$$

which leads to the desired result.

C.4 Proof of Proposition 8

Like for the proof of Proposition 7, we have $\Omega(x) \geq \Omega_J(x_J) + \Omega^J(x_{J^c}) \geq \Omega_J(x_J) + \rho(J)\Omega_{J^c}(x_{J^c}) \geq \rho(J)\Omega(x)$. Thus, if we assume $\Omega^*(q) \leq \lambda\rho(J)/2$, then $\Omega_J^*(q_J) \leq \lambda/2$ and $(\Omega^J)^*(q_{J^c}) \leq \lambda/2$. Let $\Delta = \hat{w} - w^*$.

We follow the proof from [27] by using the decomposition property of the norm Ω . We have, by optimality of \hat{w} :

$$\frac{1}{2}\Delta^\top Q\Delta + \lambda\Omega(w^* + \Delta) + q^\top \Delta \leq \lambda\Omega(w^* + \Delta) + q^\top \Delta \leq \lambda\Omega(w^*)$$

Using the decomposition property,

$$\begin{aligned} \lambda\Omega_J((w^* + \Delta)_J) + \lambda\Omega^J((w^* + \Delta)_{J^c}) + q_J^\top \Delta_J + q_{J^c}^\top \Delta_{J^c} &\leq \lambda\Omega_J(w_J^*) \\ \lambda\Omega^J(\Delta_{J^c}) &\leq \lambda\Omega_J(w_J^*) - \lambda\Omega_J(w_J^* + \Delta_J) + \Omega_J^*(q_J)\Omega_J(\Delta_J) + (\Omega^J)^*(q_{J^c})\Omega^J(\Delta_{J^c}) \\ (\lambda - (\Omega^J)^*(q_{J^c}))\Omega^J(\Delta_{J^c}) &\leq (\lambda + \Omega_J^*(q_J))\Omega_J(\Delta_J) \end{aligned}$$

Thus $\Omega^J(\Delta_{J^c}) \leq 3\Omega_J(\Delta_J)$, which implies $\Delta^\top Q \Delta \geq \kappa \|\Delta_J\|_2^2$ (we have assumed a restricted eigenvalue condition). Moreover, we have:

$$\begin{aligned} \Delta^\top Q \Delta &= \Delta^\top (Q \Delta) \leq \Omega(\Delta) \Omega^*(Q \Delta) \\ &\leq \Omega(\Delta) (\Omega^*(q) + \lambda) \leq \frac{3\lambda}{2} \Omega(\Delta) \text{ by optimality of } \hat{w} \\ \Omega(\Delta) &\leq \Omega_J(\Delta_J) + \rho(J)^{-1} \Omega^J(\Delta_{J^c}) \\ &\leq \Omega_J(\Delta_J) (3 + \frac{1}{\rho(J)}) \leq \frac{4}{\rho(J)} \Omega_J(\Delta_J) \end{aligned}$$

This implies that $\frac{\kappa}{c(J)^2} \Omega_J(\Delta_J)^2 \leq \kappa \|\Delta_J\|_2^2 \leq \Delta^\top Q \Delta \leq \frac{6\lambda}{\rho(J)} \Omega_J(\Delta_J)$, and thus $\Omega_J(\Delta_J) \leq \frac{6c(J)^2 \lambda}{\kappa \rho(J)}$, which leads to the desired result, given the previous inequalities.

C.5 Proof of Proposition 9

We have $\Omega^*(z) = \max_{\Omega(w) \leq 1} w^\top z$; the maximum can be taken over the set of extreme points of the unit ball, which leads to the desired result given Proposition 2.

References

- [1] P. Zhao, G. Rocha, and B. Yu. Grouped and hierarchical model selection through composite absolute penalties. *Annals of Statistics*, 37(6A):3468–3497, 2009.
- [2] R. Jenatton, J.Y. Audibert, and F. Bach. Structured variable selection with sparsity-inducing norms. Technical report, arXiv:0904.3523, 2009.
- [3] J. Huang, T. Zhang, and D. Metaxas. Learning with structured sparsity. In *Proc. ICML*, 2009.
- [4] L. Jacob, G. Obozinski, and J.-P. Vert. Group Lasso with overlaps and graph Lasso. In *Proc. ICML*, 2009.
- [5] S. Kim and E. Xing. Tree-guided group Lasso for multi-task regression with structured sparsity. In *Proc. ICML*, 2010.
- [6] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for sparse hierarchical dictionary learning. In *Proc. ICML*, 2010.
- [7] L. Lovász. Submodular functions and convexity. *Mathematical programming: the state of the art, Bonn*, pages 235–257, 1982.
- [8] A. Krause and C. Guestrin. Near-optimal nonmyopic value of information in graphical models. In *Proc. UAI*, 2005.
- [9] Y. Kawahara, K. Nagano, K. Tsuda, and J.A. Bilmes. Submodularity cuts and applications. In *Adv. NIPS 22*, 2009.
- [10] S. Fujishige. *Submodular Functions and Optimization*. Elsevier, 2005.
- [11] J. Edmonds. Submodular functions, matroids, and certain polyhedra. In *Combinatorial optimization - Eureka, you shrink!*, pages 11–26. Springer, 2003.
- [12] S. Negahban and M. J. Wainwright. Joint support recovery under high-dimensional scaling: Benefits and perils of ℓ_1 - ℓ_∞ -regularization. In *Adv. NIPS*, 2008.

- [13] G. Choquet. Theory of capacities. *Ann. Inst. Fourier*, 5(131-295):54, 1953.
- [14] J.B. Orlin. A faster strongly polynomial time algorithm for submodular function minimization. *Mathematical Programming*, 118(2):237–251, 2009.
- [15] S. P. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [16] R. Jenatton, G. Obozinski, and F. Bach. Structured sparse principal component analysis. In *Proc. AISTATS*, 2009.
- [17] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused Lasso. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, pages 91–108, 2005.
- [18] R. A. Horn and C. R. Johnson. *Matrix analysis*. Cambridge Univ. Press, 1990.
- [19] T. Ando. Concavity of certain maps on positive definite matrices and applications to hadamard products. *Linear Algebra and its Applications*, 26:203–241, 1979.
- [20] C. L. Mallows. Some comments on C_p . *Technometrics*, 15(4):661–675, 1973.
- [21] D. Wipf and S. Nagarajan. Sparse estimation using general likelihoods and non-factorial priors. In *Adv. NIPS 22*, 2009.
- [22] D. Bertsekas. *Nonlinear programming*. Athena Scientific, 1995.
- [23] 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.
- [24] A. Quattoni, X. Carreras, M. Collins, and T. Darrell. An efficient projection for ℓ_1 - ℓ_∞ regularization. In *Proc. ICML*, 2009.
- [25] A. Chambolle and J. Darbon. On total variation minimization and surface evolution using parametric maximum flows. *International journal of computer vision*, 84(3):288–307, 2009.
- [26] P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2563, 2006.
- [27] P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37(4):1705–1732, 2009.
- [28] A. Krause and V. Cevher. Submodular dictionary selection for sparse representation. In *Proc. ICML*, 2010.
- [29] Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Trans. PAMI*, 23(11):1222–1239, 2001.
- [30] A. Chambolle. Total variation minimization and a class of binary MRF models. In *Energy Minimization Methods in Computer Vision and Pattern Recognition*, pages 136–152. Springer, 2005.
- [31] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, 1996.
- [32] S. Rosset and J. Zhu. Piecewise linear regularized solution paths. *Ann. Statist.*, 35(3):1012–1030, 2007.
- [33] J. M. Borwein and A. S. Lewis. *Convex Analysis and Nonlinear Optimization: Theory and Examples*. Springer, 2006.