
Minimizing Nonconvex Non-Separable Functions

Yaoliang Yu

Xun Zheng

Micol Marchetti-Bowick

Eric P. Xing

Machine Learning Department, Carnegie Mellon University

Abstract

Regularization has played a key role in deriving sensible estimators in high dimensional statistical inference. A substantial amount of recent works has argued for nonconvex regularizers in favor of their superior theoretical properties and excellent practical performances. In a different but analogous vein, nonconvex loss functions are promoted because of their robustness against “outliers”. However, these nonconvex formulations are computationally more challenging, especially in the presence of nonsmoothness and non-separability. To address this issue, we propose a new proximal gradient meta-algorithm by rigorously extending the *proximal average* to the nonconvex setting. We formally prove its nice convergence properties, and illustrate its effectiveness on two applications: multi-task graph-guided fused lasso and robust support vector machines. Experiments demonstrate that our method compares favorably against other alternatives.

1 Introduction

Regularization has played a major role in recent development of statistical machine learning algorithms and applications. Many regularizers, with their unique properties, have been designed. In particular, convex regularizers have been prevalent due to their computational convenience. However, the potential superiority of nonconvex regularizers has long been recognized and pursued [1–5]. Empirically, nonconvex regularizers often yield better results than their convex counterparts [6–8]. On the flip side, nonconvex regularizers are computationally more challenging, but there has been steady progress [6, 9–12]. For instance, [9, 10, 12]

are among the first to apply the proximal gradient to nonconvex regularizers; [11] extended the coordinate descent to the nonconvex and nonsmooth setting; [6, 8] employed the convex-concave procedure; and [7] applied the alternating direction method of multipliers; etc. These existing works have greatly expanded our tool sets for coping with nonconvexity, generating remarkable successes but also suffering some limitations: a). Only apply to special scenarios [9, 11]; b). No convergence result [7] or merely the weak “convergence” in terms of function values [10–12]; c). Slow convergence due to successive linearization [6, 8]; d). Incapable of handling non-separability [9–12]. In this work we propose a meta-algorithm that enjoys stronger convergence guarantees and works in broader settings.

We are interested in the general setting where the regularizer (or the loss function) is nonconvex, nonsmooth, and non-separable. For instance, the overlapping group pursuit [8] advocated a nonconvex regularizer for each *overlapping* group and achieved better estimates. The same idea can be extended to the graph-guided fused lasso [13], see Example 1 below. However, the resulting optimization is now highly non-trivial, rendering many of the existing algorithms inapplicable, hence deserving a serious investigation.

We borrow the *proximal averaging* idea of the recent work [14] and significantly extend it to the nonconvex setting, by making the following contributions:

- Rigorously addressing the multi-valuedness and non-uniqueness of the proximal map. This difficulty does not occur for convex functions but is common for nonconvex ones. It is the key to deal with *non-separable* functions where most existing works (such as [9–12, 15]) do not apply.
- Re-establishing, sometimes with essential modifications, the many key properties of the proximal average, including a complete characterization on the real line. For instance, we show in Example 3 that there are infinitely many functions that all lead to the same hard-thresholding rule, thus shedding new lights on both the statistical and algorithmic aspects of nonconvex regularizers. These theoretical developments are com-

pletely new and provide a solid ground for some ongoing work.

- Proving the convergence of the whole sequence produced by our algorithm, which is even new for the convex case. This contribution is particularly important in applications (such as biostatistics) where variable selection, if not the sole purpose, is as desirable as achieving a small prediction error.
- Experimentally validating the proposed algorithm on applications where nonconvexity and non-separability really makes a difference. The algorithm can be easily parallelized so even better efficiency can be anticipated.

To further demonstrate its flexibility, we also apply our approach to the robust support vector machines (RSVM) by swapping the role of loss and regularizer. Here, we encounter a second motivation for nonconvex functions: As shown in [16, 17], any convex loss cannot be robust against adversarial outliers. Accordingly, RSVM replaces the convex hinge loss with the nonconvex truncated hinge loss [18–20]. Through experiments we show that our algorithm is much more efficient than previous approaches such as alternating [18] and the convex-concave procedure [19–21].

We formally state our problem in Section 2. Section 3 contains all technical results that are essential for the convergence proof in Section 4. Experiments on both multi-task GFlasso and robust SVM are conducted in Section 5, and we conclude in Section 6.

2 Problem Formulation

We are interested in solving the minimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^p} \ell(\mathbf{w}) + \bar{f}(\mathbf{w}), \text{ where } \bar{f}(\mathbf{w}) = \sum_{k=1}^K \alpha_k f_k(\mathbf{w}), \quad (1)$$

and the scalars, $\alpha_k \geq 0$, $\sum_k \alpha_k = 1$, are fixed constants throughout the paper. It is clear that many statistical machine learning algorithms can be cast under our general formulation (1). For instance, take ℓ as the least squares loss and f_1 as the 1-norm (with $K = 1$) we recover lasso [22]. We can also swap the role of the “loss” ℓ and the “regularizer” \bar{f} . For instance, let f_k be the hinge loss for the k -th training data and ℓ be the squared 2-norm, we recover the support vector machines [23]. These special cases are convex problems, and have been extensively studied in the past. Instead, we will focus on the more general setting where both functions ℓ and $\{f_k\}$ are *nonconvex* and *non-smooth*. The motivation to have nonconvex functions can be diverse, and will be illustrated in Example 1 and Example 2 below. We emphasize that the function \bar{f} is *non-separable*, in the sense that its components $\{f_k\}$ have overlapping argument \mathbf{w} .

In general, (1) can be very challenging to solve, even when we lower our expectation to the convergence to some critical point. Fortunately, practical problems usually come with some structure that we can (and should) exploit. In particular, we make the following assumption:

Assumption 1 *The function ℓ has L -Lipschitz continuous gradient $\nabla \ell$, each f_k is M_k -Lipschitz continuous (all w.r.t. the Euclidean norm $\|\cdot\|$), and the proximal map $\mathbf{P}_{f_k}^\mu$ can be computed “easily” for any $\mu > 0$.*

Recall that a mapping $g : \mathbb{R}^p \rightarrow \mathbb{R}^d$ is M -Lipschitz continuous for some $M > 0$ if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$, $\|g(\mathbf{x}) - g(\mathbf{y})\| \leq M\|\mathbf{x} - \mathbf{y}\|$. The proximal map \mathbf{P}_f^μ for any function f and parameter $\mu > 0$ is defined as:

$$\mathbf{P}_f^\mu(\mathbf{w}) = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^p} \frac{1}{2\mu} \|\mathbf{z} - \mathbf{w}\|^2 + f(\mathbf{z}). \quad (2)$$

When f is the indicator function of some set C , $\mathbf{P}_f^\mu(\mathbf{w})$ simply returns the closest point in C to \mathbf{w} , namely the familiar Euclidean projection. If f is the 1-norm, then \mathbf{P}_f^μ becomes the well-known soft-shrinkage operator that is widely used in sparse methods, e.g. lasso. The parameter $\mu > 0$ in the definition (2) plays the role of step size in the algorithms we will develop, and needs to be set properly. Assumption 1 requires the proximal map to be “easily” computable, meaning roughly that its complexity should be on par with that of computing the gradient of function ℓ . This avoids the proximal map to become the bottleneck if we use a gradient-type algorithm. As we will see, Assumption 1 is quite reasonable in a number of applications.

In the nonconvex setting, we need to pay extra care to even some “obvious” properties of the proximal map (such as non-emptiness and non-uniqueness). Such technicalities, albeit important, will be postponed until Section 3. For the purpose of explaining our main idea let us pretend momentarily that \mathbf{P}_f^μ is a well-defined “function”, i.e., $\mathbf{P}_f^\mu(\mathbf{w})$ is some “closest point” to \mathbf{w} , measured by the function f . With this “simplification” we can now demonstrate how Assumption 1 is naturally satisfied in some important applications:

Example 1 (GFlasso, [13]) *The graph-guided fused lasso exploits some graph structure to improve feature selection. Given some a priori graph whose nodes correspond to the feature variables, [13] used the regularizer $\tilde{f}_{ij}(\mathbf{w}) = |w_i - w_j|$ for every edge $(i, j) \in E$, to encourage connected nodes to be selected jointly. However, as pointed out in [13], this regularizer brings a large bias as it also requires connected nodes to have similar weights, which is likely not true in general. Here we reduce the bias by proposing the nonconvex regularizer: $f_{ij} = \min\{\tilde{f}_{ij}, \tau\}$, i.e., we cap the regularizer at the threshold τ . This will allow some weights to differ significantly without getting heavily penalized. In this example ℓ is the least squares loss.*

Algorithm 1 PA-PG.

- 1: Initialize \mathbf{w}_0, μ .
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: $\mathbf{z}_t = \mathbf{w}_{t-1} - \mu \nabla \ell(\mathbf{w}_{t-1})$,
- 4: $\mathbf{w}_t \in \overline{\sum_k \alpha_k \cdot \hat{\mathbf{P}}_{f_k}^\mu(\mathbf{z}_t)}$.
- 5: **end for**

Algorithm 2 PA-APG.

- 1: Initialize $\mathbf{w}_0 = \mathbf{u}_1, \mu, \eta_1 = 1$.
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: $\mathbf{z}_t = \mathbf{u}_t - \mu \nabla \ell(\mathbf{u}_t)$,
- 4: $\mathbf{w}_t \in \sum_k \alpha_k \cdot \hat{\mathbf{P}}_{f_k}^\mu(\mathbf{z}_t)$,
- 5: $\eta_{t+1} = \frac{1 + \sqrt{1 + 4\eta_t^2}}{2}$,
- 6: $\mathbf{u}_{t+1} = \mathbf{w}_t + \frac{\eta_t - 1}{\eta_{t+1}}(\mathbf{w}_t - \mathbf{w}_{t-1})$.
- 7: **end for**

It is easily verified that each f_{ij} is $\sqrt{2}$ -Lipschitz continuous w.r.t. the Euclidean norm. Moreover, the proximal map $\mathbf{P}_{f_{ij}}^\mu$ can be computed in closed-form (see Appendix J for the detailed derivation): For $s \in \{i, j\}$, trivially $[\mathbf{P}_{f_{ij}}^\mu(\mathbf{w})]_s = w_s$, while for $\{s, t\} = \{i, j\}$,

$$[\mathbf{P}_{f_{ij}}^\mu(\mathbf{w})]_s = w_s - \text{sign}(w_s - w_t) \min\{\mu\eta, \frac{|w_i - w_j|}{2}\},$$

$$\eta = \begin{cases} 0, & |w_i - w_j| \geq 2\sqrt{\mu\tau} + ((\sqrt{\tau} - \sqrt{\mu})_+)^2 \\ 1, & |w_i - w_j| \leq 2\sqrt{\mu\tau} + ((\sqrt{\tau} - \sqrt{\mu})_+)^2 \end{cases}.$$

Roughly speaking, $\eta = 0$ iff $|w_i - w_j|$ is large, and consequently $[\mathbf{P}_{f_{ij}}^\mu(\mathbf{w})]_s = w_s$, i.e., the algorithm gives up “fusing” w_i and w_j , which can be beneficial in reducing the estimation bias when w_i and w_j are truly different.

The next example swaps the role of the loss ℓ and the regularizer f_k , demonstrating the flexibility of the general formulation (1).

Example 2 (Robust SVM, [18–20]) Support vector machine (SVM) is one of the most popular algorithms for binary classification. However, it is known not to be robust against outliers [16–18, 20]. In fact, [16] constructed an example on which all algorithms based on convex losses fail. Instead, [18, 20] proposed the (nonconvex) truncated hinge loss as a robust alternative: $f_i(\mathbf{w}) = \min\{\tau, (1 - y_i \mathbf{x}_i^\top \mathbf{w})_+\}$. It is easy to verify that f_i is $\|\mathbf{x}_i\|$ -Lipschitz continuous w.r.t. the Euclidean norm, and its proximal map can be computed as (see Appendix K for the detailed derivation):

$$\mathbf{P}_{f_i}^\mu(\mathbf{w}) = \mathbf{w} + \left[\frac{1 - y_i \mathbf{w}^\top \mathbf{x}_i}{\mathbf{x}_i^\top \mathbf{x}_i} \right]_0^{\mu\eta} \cdot y_i \mathbf{x}_i, \quad (3)$$

where $[\cdot]_0^\mu$ denotes the projection onto the interval $[0, \mu]$, and the parameter $\eta \in \{0, 1\}$ is explicitly given in (33) in the appendix. Roughly speaking, $\eta = 0$ iff the margin $y_i \mathbf{x}_i^\top \mathbf{w}$ is small, i.e., the pair (\mathbf{x}_i, y_i) is likely to be an outlier, in which case $\mathbf{P}_{f_i}^\mu(\mathbf{w}) = \mathbf{w}$, i.e. the algorithm “refuses” to update the weight. In this example ℓ is the (multiple of the) squared 2-norm.

The two examples represent two extremes in applications: the former has a nonconvex *non-separable* regularizer while the loss is the simple least squares, and the latter has a nonconvex *non-separable* loss while the regularizer is the simple (squared) 2-norm. Of course

it is possible to have other combinations (e.g. the overlapping group lasso [8]), but for illustration purpose we shall contend ourselves with the above examples.

Having demonstrated the relevance of problem (1), we now turn to how to solve it efficiently. The main difficulty, apart from the nonconvexity, is the non-separability of the functions $\{f_k\}$: they all share the same weight \mathbf{w} . Accordingly, the coordinate descent algorithm of [11] cannot be efficiently applied. Similarly, the (block) proximal gradient (PG) algorithm, such as those in [9, 10, 12, 15], cannot be directly applied either, because we do not know how to efficiently compute the proximal map $\mathbf{P}_{f_k}^\mu$, even when we assume each $\mathbf{P}_{f_k}^\mu$ is easy to compute. Other possible algorithms include the alternating strategy [18] and the convex-concave procedure [8, 20]. However, due to successive linearizations, these algorithms can be slow, and normally would need the functions f_k to be smooth.

Our idea is to approximate the proximal map $\mathbf{P}_{\bar{f}}^\mu$ using the linearization:

$$\mathbf{P}_{\bar{f}}^\mu \approx \sum_k \alpha_k \mathbf{P}_{f_k}^\mu. \quad (4)$$

Pretending $\mathbf{P}_{f_k}^\mu(\mathbf{w})$ is a single “point”, we can plug the right-hand approximation into the PG algorithm. The resulting algorithm (PA-PG), summarized in Algorithm 1, can now be used to solve (1). It is extremely simple: alternating between a standard gradient step w.r.t. the loss ℓ and a proximal step¹ w.r.t. the regularizers $\{f_k\}$. Although the approximation (4) may seem overly naive, linearizing a nonlinear object (at least “locally”) is a ubiquitously useful technique (such as the Taylor expansion in calculus). Indeed, for convex functions ℓ and $\{f_k\}$, the recent work [14] gave a formal justification of Algorithm 1 and also the accelerated variation in Algorithm 2. In our later experiments we found Algorithm 2 to be again more effective than Algorithm 1 even for nonconvex functions.

We aim to generalize the nice results of [14] to the current nonconvex setting and demonstrate its effectiveness through experiments. This goal, as clear as it

¹The big overbar and hat notation will be understood after we present relevant technical results in Section 3.

is, is far from trivial though. The difficulties we face include: a). How should we interpret the sum in (4), considering that for nonconvex f_k , $P_{f_k}^\mu$ may no longer be single-valued? b). Does the right-hand side still correspond to a proximal map of some function? c). Can we formally justify the linearization? d). What guarantee does (4) enjoy if plugged into the proximal gradient algorithm? These questions cannot be answered by existing works, such as [14] which relies entirely on convexity, or [15] which relies entirely on separability (i.e., the functions f_k have non-overlapping arguments). A substantial technical development is needed, which we do in Section 3.

Before going into the technical details, let us point out a few practical advantages of Algorithm 1 and Algorithm 2: 1). As a general meta-algorithm, they can be used in a variety of settings. 2). For *separable* functions, they reduce to the block PG algorithm [15] while for $K = 1$ we recover the algorithms in [9, 10, 12], including the popular FISTA [24, 25] when convexity is present. 3). They can be easily parallelized when K is large (such as the RSVM example). 4). The iterates they generate converge to a critical point (globally optimal if convexity is assumed). These nice properties are not shared by too many algorithms, let alone simultaneously.

3 Technical Results

To justify our new algorithm, we need a few technical tools from variational analysis [26]. We equip \mathbb{R}^p with the usual inner product $\langle \cdot, \cdot \rangle$ and the induced Euclidean norm $\|\cdot\|$. For any *closed*² function f (not necessarily convex), its Moreau envelope (with parameter $\mu > 0$) is defined as [26]:

$$e_f^\mu(\mathbf{w}) = \inf_{\mathbf{z}} \frac{1}{2\mu} \|\mathbf{w} - \mathbf{z}\|^2 + f(\mathbf{z}), \quad (5)$$

and the proximal map is the corresponding minimizer(s), see (2). One can roughly think the envelope function e_f^μ as a “regularized” version of f . For instance, if f is the 1-norm, then e_f^μ is the celebrated Huber’s function in robust statistics [17]. Since we have stepped out of the convex domain, many “obvious” properties, such as well-definedness, smoothness, and uniqueness, can no longer be taken for granted. Fortunately, many appealing properties retain, possibly under an alternative interpretation.

It can be shown that P_f^μ is nonempty-valued iff f majorizes some quadratic function and μ is small [26]. Here, for simplicity, we assume throughout that f is bounded from below so that we need not restrict μ . Thus, $P_f^\mu : \mathbb{R}^p \rightrightarrows \mathbb{R}^p$ is nonempty-, compact-, possibly

²The function $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ is closed iff its epigraph $\{(\mathbf{x}, t) \in \mathbb{R}^p \times \mathbb{R} : f(\mathbf{x}) \leq t\}$ is a closed set.

nonconvex- and multi-valued. In fact, Lemma 1 in Appendix C showed that the proximal map $P_f^\mu(\mathbf{w})$ is single valued iff the envelope function e_f^μ is differentiable at \mathbf{w} . Both are trivially true for convex functions, but they may fail for general nonconvex functions, particularly functions that are “capped” at a certain value—our running examples in this paper. When $\mu \downarrow 0$, $e_f^\mu(\mathbf{w}) \uparrow f(\mathbf{w})$ for all \mathbf{w} [26]. As mentioned before, the proximal map is the key component of the proximal gradient algorithm.

In the nonsmooth and nonconvex setting, the usual gradient or subgradient no longer applies to characterize critical points. Instead, we are forced to “localize”. We first define the regular (or Fréchet) subdifferential $\hat{\partial}f(\mathbf{w})$ at \mathbf{w} , as the collection of vectors \mathbf{v} such that

$$\forall \mathbf{z}, f(\mathbf{z}) \geq f(\mathbf{w}) + \langle \mathbf{z} - \mathbf{w}, \mathbf{v} \rangle + o(\|\mathbf{z} - \mathbf{w}\|),$$

where the little-o term signifies a local neighborhood. Since $\hat{\partial}f$ can be empty even for Lipschitz functions (e.g. $-|\cdot|$ at the origin), we take its “closure” to avoid this degeneracy, arriving at the subdifferential $\partial f(\mathbf{w})$:

$$\{\mathbf{v} : \exists \mathbf{w}_n \rightarrow \mathbf{w}, f(\mathbf{w}_n) \rightarrow f(\mathbf{w}), \mathbf{v}_n \in \hat{\partial}f(\mathbf{w}_n), \mathbf{v}_n \rightarrow \mathbf{v}\}.$$

Clearly, $\hat{\partial}f(\mathbf{w}) \subseteq \partial f(\mathbf{w})$ for all \mathbf{w} . Pleasantly, if f is (resp. continuously) differentiable at \mathbf{w} , then $\hat{\partial}f(\mathbf{w})$ (resp. $\partial f(\mathbf{w})$) coincides with the usual derivative. From the definition it follows that if \mathbf{w} is a local minimizer, then $0 \in \hat{\partial}f(\mathbf{w}) \subseteq \partial f(\mathbf{w})$, which generalizes the familiar Fermat’s rule. We will be interested in finding (asymptotically) some \mathbf{w} so that $0 \in \partial f(\mathbf{w})$, i.e., the critical points of f .

We reassure ourselves some nice properties of the Moreau envelope and the proximal map in Appendix A. In particular, we proved that any Moreau envelope is concave after subtracting the function $\frac{1}{2\mu} \|\cdot\|^2$. The converse, which we prove next, will allow us to “average” functions in a somewhat peculiar but computationally appealing way.

Let SCV_μ be the set of *finite-valued* μ -semiconcave functions, that is, functions $f : \mathbb{R}^p \rightarrow \mathbb{R}$ such that $f - \frac{1}{2\mu} \|\cdot\|^2$ is concave. For conciseness, denote CPB the class of closed, proper, and bounded from below functions. The next result, whose proof is deferred to Appendix B, significantly extends [14, Proposition 2]:

Proposition 1 *Fix $\mu > 0$ and $f \in \text{CPB}$. Then $f = e_g^\mu$ for some function $g \in \text{CPB}$ iff $f \in \text{SCV}_\mu$. Moreover, the Moreau envelope map $e^\mu : \text{CPB} \rightarrow \text{SCV}_\mu$ that sends $f \in \text{CPB}$ to e_f^μ is increasing, and concave on any convex subset of CPB (under the pointwise order).*

Note that in the nonconvex setting, the Moreau envelope (for any fixed $\mu > 0$) is no longer injective (see

Example 3 below). It is clear that SCV_μ is a *convex set*. Therefore we can average μ -semiconcave functions and still be able to find a pre-image under the Moreau envelope map, thanks to Proposition 1. This leads to the generalized notion of the proximal average. Specifically, recall that $\bar{f} = \sum_k \alpha_k f_k$ with each $f_k \in \text{CPB}$, i.e. \bar{f} is the convex combination of the component functions $\{f_k\}$ under the weight $\{\alpha_k\}$. Note that we always assume $\bar{f} \in \text{CPB}$.

Definition 1 (Proximal average) *The proximal average A^μ is any function g such that $e_g^\mu = \sum_{k=1}^K \alpha_k e_{f_k}^\mu$. In face of non-uniqueness, we will always pick $A^\mu = -e_M^\mu$, where $M := -\sum_k \alpha_k e_{f_k}^\mu$.*

The main idea behind this definition is to find some function whose Moreau envelope is simply the average $\sum_k \alpha_k e_{f_k}^\mu$. Indeed, the existence of such a function follows from the surjectivity of e^μ , which we proved in Proposition 1. However, unlike the convex case, the proximal average for nonconvex functions need *not* be unique, and for concreteness we have picked a convenient representative in Definition 1. We remark that [27] used a slightly different definition for the sake of pursuing smoothness; for us the current definition is more useful. Note that for any function f , the so-called μ -proximal hull $h_f^\mu := -e_{(-e_f^\mu)}^\mu$ has the same Moreau envelope as f but need not coincide with f [26] (while for convex functions f , always $h_f^\mu = f$). One easily verifies, through the proximal hull, that our particular choice in Definition 1 is indeed legitimate (and will prove convenient later, see Proposition 5).

To facilitate our discussions, let us first prove some results that are interesting in their own rights. For any multi-valued map $P : \mathbb{R}^p \rightrightarrows \mathbb{R}^p$, we define its closure $\bar{P} : \mathbb{R}^p \rightrightarrows \mathbb{R}^p$, $\mathbf{w} \mapsto \{\mathbf{z} : \exists \{\mathbf{w}_n, \mathbf{z}_n\} \rightarrow (\mathbf{w}, \mathbf{z}), \mathbf{z}_n \in P(\mathbf{w}_n)\}$, i.e., the graph of the closure \bar{P} is simply the closure of the graph of P . By ‘‘closing’’ a map we gain some continuity property. Also define $\hat{P}(\mathbf{w}) = P(\mathbf{w})$ at points \mathbf{w} where $P(\mathbf{w})$ is single-valued and empty otherwise.

Definition 2 (Extremal proximal maps) *Define the limiting proximal map $L_f^\mu = \overline{\hat{P}_f^\mu}$, and the hull proximal map $H_f^\mu : \mathbb{R}^p \rightrightarrows \mathbb{R}^p$, $\mathbf{w} \mapsto \text{conv}(P_{f_k}^\mu(\mathbf{w}))$, where conv denotes the convex hull.*

Thanks to item iv) of Proposition 7 (in Appendix A) and [26, Theorem 1.25], we know $\emptyset \neq L_f^\mu(\mathbf{w}) \subseteq P_f^\mu(\mathbf{w}) \subseteq H_f^\mu(\mathbf{w})$ for all \mathbf{w} . The inclusion can be strict, as we will see shortly. It may help to keep in mind that L_f^μ and H_f^μ are the smallest and biggest proximal map ‘‘compatible’’ with the function f , respectively. In Appendix C we prove many new structural properties of these different notions of proximal maps.

The next result characterizes exactly when Moreau envelopes coincide (proof in Appendix D).

Proposition 2 *Fix $\mu > 0$. For any $f \in \text{CPB}$, there exist $h_f^\mu, \ell_f^\mu \in \text{CPB}$ such that for any $g \in \text{CPB}$, $e_g^\mu = e_f^\mu + c$ for some constant c iff $h_f^\mu \leq g - c \leq \ell_f^\mu$ iff $P_{\ell_f^\mu}^\mu(\mathbf{w}) \subseteq P_g^\mu(\mathbf{w}) \subseteq P_{h_f^\mu}^\mu(\mathbf{w})$ for all \mathbf{w} .*

In fact, ℓ_f^μ is the restriction of h_f^μ onto some closed set. Their explicit forms can be found in the proof. It is also true that $P_{\ell_f^\mu}^\mu = L_f^\mu$ on the real line. Using Proposition 2 we can easily characterize when the proximal average is unique (essentially our particular choice in Definition 1 plays the role of h_f^μ). It also leads to the following result that completely characterizes proximal maps on the real line (proof in Appendix D):

Proposition 3 *The map $P : \mathbb{R} \rightrightarrows \mathbb{R}$ is a proximal map iff it is (nonempty) compact-valued, monotone, and has a closed graph. Moreover, there is a unique function (up to addition of a constant) f such that $P_f = P$ iff P is also convex-valued.*

Thus, both the SCAD [2] and the MC+ [3] thresholding rules correspond to a *unique* regularization function. In contrast, there are infinitely many different regularizers that all lead to the hard thresholding rule, see Example 3. Importantly, Proposition 3 allows us to directly design the proximal map (thresholding rule), without even the need to refer to the regularizer f ! We now come to the main result for justifying Algorithm 1. Recall that the main property of the proximal average, as seen from its definition, is that its Moreau envelope is the convex combination of the Moreau envelopes of the component functions. We wish to say something similar for its proximal map. Indeed, this is possible after an appropriate modification (proof in Appendix E):

Proposition 4 *For all \mathbf{w} , $\emptyset \neq \overline{\sum_k \alpha_k \hat{P}_{f_k}^\mu(\mathbf{w})} \subseteq [P_{A^\mu}^\mu(\mathbf{w}) \cap (\sum_{k=1}^K \alpha_k P_{f_k}^\mu(\mathbf{w}))]$.*

Recall that the middle term is exactly the approximation we employed in Section 2. Unlike the convex case, we have to replace the simpler average $\sum_{k=1}^K \alpha_k P_{f_k}^\mu$ with the slightly more complicated closure $\overline{\sum_k \alpha_k \hat{P}_{f_k}^\mu}$, due to the possible multi-valuedness of $P_{f_k}^\mu$. Indeed, some element in $\sum_{k=1}^K \alpha_k P_{f_k}^\mu(\mathbf{w})$ may not be in $P_{A^\mu}^\mu(\mathbf{w})$, which itself may change if we use a different proximal average. Proposition 4 avoids such pathology by always picking a common element. Moreover, in Section 4 we prove Algorithm 1 converges to a critical point of the proximal average, which itself may not even be unique! This ambiguity is resolved using Proposition 4, which guarantees all realizations

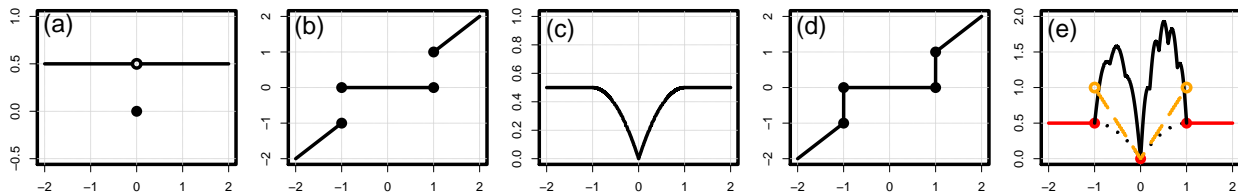


Figure 1: (a): the ℓ_0 function; (b): hard-thresholding operator; (c): proximal hull $h_{|\cdot|}$; (d): proximal map of $h_{|\cdot|}$; (e): solid: another (continuous) function that has (b) as proximal map, dashed (+ red solid): function g . See Example 3 for the explanations and formulas. Throughout $\mu = 1, \lambda = 1$.

of the proximal average coincide along the trajectory of Algorithm 1.

On the real line, thanks to Proposition 3, $\overline{\sum_k \alpha_k \hat{P}_{f_k}^\mu} = \sum_k \alpha_k L_{f_k}^\mu$, with the latter being readily available. For our Example 1 and Example 2, the computations also reduce to the real line (see Appendix J and Appendix K, respectively), hence can be easily addressed.

Let us now demonstrate some pathology of the proximal map, using a familiar nonconvex function.

Example 3 Consider the cardinality function on the real line $|x|_0 = \frac{\lambda^2}{2} \mathbf{1}_{x \neq 0}$. Its proximal map (with $\mu = 1$) is the well-known hard-thresholding operator:

$$P_{|\cdot|_0}(x) = \begin{cases} x, & |x| > \lambda \\ \{0, x\}, & |x| = \lambda \\ 0, & |x| < \lambda \end{cases} \quad (6)$$

In the literature, the above proximal map at $|x| = \lambda$ is usually set to 0. Mathematically, this is not precise and possibly confusing: If $P_{|\cdot|_0}$ was single-valued at $|x| = \lambda$, then $e_{|\cdot|_0}(x) = \frac{1}{2} \min\{\lambda^2, x^2\}$ would be continuously differentiable, which is not true. Interestingly, without the tools we have developed so far, [1] noticed that the functions

$$h(x) := \frac{1}{2}(\lambda^2 - (|x| - \lambda)^2 \mathbf{1}_{|x| \leq \lambda}) \quad (7)$$

$$g(x) := \lambda|x| \mathbf{1}_{|x| < \lambda} + \frac{\lambda^2}{2} \mathbf{1}_{|x| \geq \lambda} \quad (8)$$

also have (6) as their (limiting) proximal map. We verify that the former is exactly the proximal hull $h_{|\cdot|_0}$. Applying Proposition 2 we know any function $f \geq h_{|\cdot|_0}$ with equality on the closed set $[-\infty, -\lambda] \cup [\lambda, \infty] \cup \{0\}$ will have (6) as its limiting proximal map. Furthermore, if f is strictly larger than $h_{|\cdot|_0}$ on $]-\lambda, \lambda[$, such as the function g above, then according to Lemma 7 (in Appendix C) it has the same proximal map as the cardinality function! See Figure 1 for the illustrations.

Statistically, one is often interested in the hard-thresholding rule (6), rather than the cardinality function itself [1, 2, 10]. Figure 1 shows that there are actually infinitely many functions that all yield the same proximal map (6). This observation suggests that we

should not base our algorithm on any particular function form but on the proximal map directly (which is less ambiguous). In this sense the proximal gradient algorithm seems to be a well fit. Similar conclusions have been made in [10]. We point out that the lessons we learned from this example extend to most nonconvex regularizers therefore deserve some attention.

To provide a strong convergence guarantee for Algorithm 1, we will (and perhaps should) restrict the (nonconvex and nonsmooth) functions under our consideration, for otherwise they can behave very pathologically. To do so we recall some notions from semi-algebraic geometry [28]. A set $A \subseteq \mathbb{R}^p$ is semi-algebraic if it is the finite unions of finite intersections of the sets $\{\mathbf{w} \in \mathbb{R}^p : p_0(\mathbf{w}) = 0, p_1(\mathbf{w}) < 0\}$, where p_0, p_1 are polynomials with real coefficients. For instance, hyperplanes, halfspaces, spheres, ellipsoids, the positive semi-definite cone, are all semi-algebraic. The most striking property of semi-algebraic sets is that their intersection with any line is the union of finitely many points and open intervals (due to the fact that any polynomial admits only finitely many roots). Thus, for instance, the set of all natural numbers is not semi-algebraic. A function $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ is semi-algebraic iff its graph $\{(\mathbf{w}, f(\mathbf{w})) : \mathbf{w} \in \text{dom } f\}$ is a semi-algebraic set. For instance, all power functions with rational exponent and all polyhedral functions are semi-algebraic. On one hand, semi-algebraic functions are extremely well-structured, allowing one to prove many strong results; on the other hand, they appear very naturally in various applications, such as the ones we consider here: All functions appear in our experiments are semi-algebraic. Note that the exponential function and power functions with irrational exponent are not semi-algebraic. They belong to the more general class of definable functions³. For brevity, we omit the relevant definitions here but refer to the nice article [28]. Since all of our results extend to definable functions (w.r.t. some order-minimal structure that contains all semi-algebraic functions), we will use the term definable freely below; for all practical purposes, one can simply think of “definable” as semi-algebraic.

³In case one wonders, there do exist non-definable (even convex) functions, which oscillate infinitely often. Such functions are unlikely to be useful in applications though.

Definable functions are closed under all familiar algebraic operations. For instance, the sum, product, and composition of definable functions are definable, respectively. So is the scalar multiple and the inverse. Moreover, the following result is useful to us (proof in Appendix F).

Proposition 5 *The function f is definable iff $e_f^\mu(\mathbf{w})$, as a function of (\mathbf{w}, μ) , is definable on $\mathbb{R}^p \times \mathbb{R}_{++}$. In particular, the proximal average (cf. Definition 1) of definable functions is definable.*

As we mentioned before, when $\mu \downarrow 0$, $e_f^\mu \uparrow f$ pointwise [26]. Under the Lipschitz assumption, we can strengthen the convergence to be uniform (Proof in Appendix G):

Proposition 6 *Under Assumption 1 we have $0 \leq \bar{f} - A^\mu \leq \bar{f} - \sum_k \alpha_k e_{f_k}^\mu \leq \frac{\mu}{2} \sum_{k=1}^K \alpha_k M_k^2$.*

Similar as the convex case, we see that the proximal average A^μ is a better under-approximation to \bar{f} than the average of Moreau envelopes, i.e. $\sum_k \alpha_k e_{f_k}^\mu$.

4 Theoretical Justification

Given our development in the previous section, it is now clear that Algorithm 1 aims at solving the approximate problem:

$$\min_{\mathbf{w}} \ell(\mathbf{w}) + A^\mu(\mathbf{w}). \quad (9)$$

The next important pieces are to show a). Algorithm 1 converges for the approximate problem (9); b). The approximate problem (9) is reasonably ‘‘close’’ to the original problem (1). Indeed, for the first piece, we have the following result (proof in Appendix H):

Theorem 1 *Let Assumption 1 hold and the functions ℓ and $\{f_k\}$ be definable. Choose $\mu < 1/L$, then Algorithm 1 converges to a critical point of (9), provided that the iterates are bounded.*

The last assumption is trivially met if, say the objective in (9) has bounded sublevel sets. To appreciate the significance of Theorem 1, let us consider a simple example: Assume say both 1 and -1 are critical points of our problem, then any limit point of the iterate sequence $\{1, -1, 1, -1, \dots\}$ is indeed critical, but the whole sequence does not converge at all. This behavior can happen for the coordinate descent algorithm [11] or the convex-concave procedures (cccp) [8, 19, 20], but is eliminated for our algorithm, thanks to Theorem 1.

To fulfill our second piece, we need a notion of approximate minimizer. We call \mathbf{w} an ϵ -local minimizer of f if there exists some neighborhood \mathcal{N} of \mathbf{w} such that for all $\mathbf{z} \in \mathcal{N}$, $f(\mathbf{w}) \leq f(\mathbf{z}) + \epsilon$. Of course, when $\epsilon = 0$,

we retrieve the usual notion of local minimizer. Then we have (proof in Appendix I):

Theorem 2 *Let Assumption 1 hold. Fix the accuracy $\epsilon > 0$ and choose $\mu < \min\{1/L, 2\epsilon / \sum_k \alpha_k M_k^2\}$. If Algorithm 1 converges to an ϵ -local minimizer of (9), $\tilde{\mathbf{w}}$, then $\tilde{\mathbf{w}}$ is also a (2ϵ) -local minimizer of (1). Same is true if $\tilde{\mathbf{w}}$ is in fact an ϵ -global minimizer.*

It is possible to prove that locally Algorithm 1 converges at a rate no slower than sublinear (when all functions are semi-algebraic). Moreover, if we let $\mu \downarrow 0$, we can prove that the iterates converge to a critical point of the original problem (1). In experiments, we found that a relatively small μ already yields satisfying results, therefore we omit the rather technical discussions.

5 Experiments

We evaluate our algorithm on two application domains: truncated GFlasso (Example 1) and robust SVM (Example 2). We demonstrate the benefits of nonconvex formulations by comparing with the convex counterparts, and we verify the effectiveness of our proposed algorithm against alternative optimization methods such as alternating coordinate descent (`alter`) [18] and the convex-concave procedure (`cccp`) [19]. We found that Algorithm 2 is always faster than Algorithm 1 in all our experiments so only the former (denoted as `proxavg`) is included here.

5.1 Multi-task GFlasso (Example 1)

Formally, the multi-task graph-guided fused lasso model is given by:

$$\begin{aligned} \frac{1}{2} \|Y - XW\|_F^2 + \lambda \sum_{j=1}^q \phi(\mathbf{w}_j) \\ + \gamma \sum_{(j,k) \in E} \omega_{jk} \psi(\mathbf{w}_j - \text{sign}(\omega_{jk}) \mathbf{w}_k), \end{aligned} \quad (10)$$

where $\mathbf{w}_j \in \mathbb{R}^p$ is the j -th column of W . Here ϕ is a regularizer that encourages sparsity among the elements of \mathbf{w}_j , and ψ is a regularizer that encourages fusion between the elements of \mathbf{w}_j and \mathbf{w}_k when output variables j and k are connected in the graph. In our experiment we use $\phi(\mathbf{u}) = \|\mathbf{u}\|_1$ and $\psi(\mathbf{u}) = \sum_i \min\{|u_i|, \tau\}$. If $\tau = 0$, we recover the multi-task lasso while if $\tau = \infty$, we recover the GFlasso which is convex. For any $\tau > 0$, the fusion regularizer is non-separable and nonconvex.

We compare different methods (corresponding to different τ 's) on a synthetic data in which pairs of correlated output variables, \mathbf{y}_j and \mathbf{y}_k , have similar weights, \mathbf{w}_j and \mathbf{w}_k . We choose a block correlation graph for concreteness and generate the data as follows. First, partition the output variables $\mathbf{y}_1, \dots, \mathbf{y}_q$ into disjoint

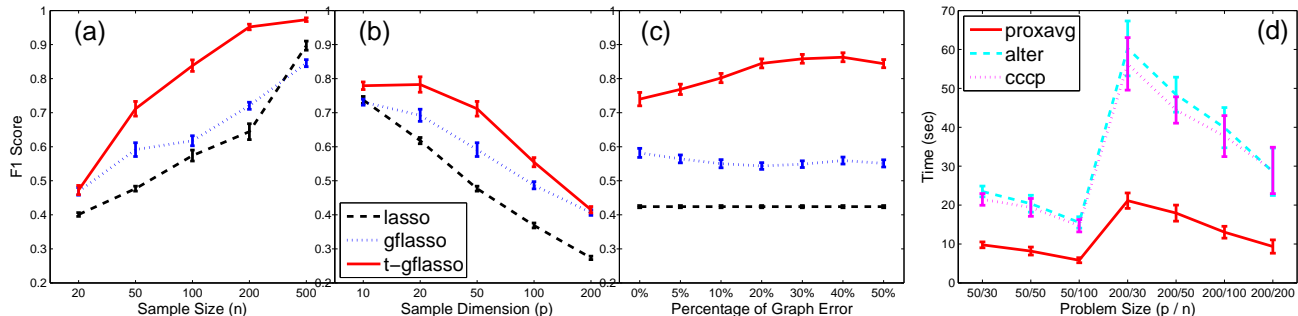


Figure 2: The F1 score under (a): varying sample size (n); (b): varying dimensionality (p); (c): varying graph error. (d): Speed comparison between our algorithm (proxavg), alternating, and cccp.

ccat	proxavg	alter	cccp	svm	sparse
training time (sec)	25.34	897.98	1151.26	69.85	110.21
test accuracy (%)	91.12	84.77	87.02	69.96	89.20
detected outlier (%)	9.94	20.12	18.45	-	10.47

Table 1: Results on the CCAT data set.

groups. Next, generate the sparsity pattern for the weight matrix W by assigning the same (randomly chosen) set of active input features to all the output variables in each group. The nonzero entries of W are drawn from Uniform(0.4, 0.8), but all variables in the same group are given the same weight for each feature. Draw $X \sim \mathcal{N}(0, I)$ and $Y \sim \mathcal{N}(XW, \sigma^2 I)$. Finally, generate the correlation graph E over the output variables by thresholding the sample covariance matrix at some value ν . We also test the robustness of the algorithms by randomly changing certain percentage of the edge weights.

A series of experimental results are shown in Figure 2. Regularization parameters λ , γ , and τ (whenever applicable) are selected by optimizing the prediction error on a held-out set. We observe that a). The nonconvex fusion regularizer ($0 < \tau < \infty$) consistently outperforms lasso ($\tau = 0$) and GFlasso ($\tau = \infty$) in terms of both feature selection (F1 score) and prediction error (not shown); b). Our algorithm is several times faster than both **alter** and **cccp**.

5.2 Robust SVM (Example 2)

We conducted experiments on two benchmark data sets. The Long-Servedio [16] dataset is a well-crafted synthetic data for testing robustness against label noise and delicate leverage points. We generate 10,000 training examples (each with 21 features) and randomly flip 10% labels. The other real dataset CCAT from RCV1 [29] contains 23,149 training examples (each with 47,152 features) and 781,265 test examples. Similarly, we randomly flip 10% of the labels in the training set, and scale the corresponding features by 10. We average the results with 10 repetitions and report them in Table 1. For both SVM, **alter** [18],

and **cccp** [19], we use the state-of-the-art LIBLINEAR solver [30]. Instead of the (squared) 2-norm regularizer, our algorithm extends easily to the 1-norm regularizer hence we also include **sparse** to further demonstrate the flexibility. In contrast, LIBLINEAR cannot deal with the 1-norm regularizer.

We confirmed that SVM fails miserably on the Long-Servedio dataset (achieving 72.14% prediction error), while all other solvers (aimed for the robust SVM) achieve nearly perfect results and identify the correct amount of outliers. Our algorithm is fastest but the margin is small (due to the small size of the dataset). For the CCAT dataset, our algorithm not only achieves superior prediction accuracy but also much pronounced efficiency. Again, SVM severely suffers from outliers while **alter** and **cccp** are slow due to their sequential nature: multiple calls of the SVM solver can only be executed consecutively. Interestingly, with small sacrifice in accuracy and training time, **sparse**, using 1-norm regularizer in SVM, learns a model with only 4.8% nonzero entries, whereas the models learned by all other methods are at least 10 times denser. This could hugely reduce the test time—a critical requirement in some financial applications.

6 Conclusions

We successfully extended the proximal average proximal gradient algorithm into the nonconvex setting, through a careful examination of the now multi-valued proximal map. We proved that the whole sequence of iterates converges to a critical point. Experimentally, the proposed algorithm has shown much promise, and naive parallelizability makes it even more favorable. We intend to strengthen the convergence guarantee and develop a fully distributed implementation.

References

- [1] Anestis Antoniadis. “Wavelets in Statistics: A Review.” *Journal of the Italian Statistical Association*, vol. 2 (1997), pp. 97–130 (cit. on pp. 1, 6, 18).
- [2] Jianqing Fan and Runze Li. “Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties.” *Journal of the American Statistical Association*, vol. 96, no. 456 (2001), pp. 1348–1360 (cit. on pp. 1, 5, 6, 16).
- [3] Cun-Hui Zhang. “Nearly unbiased variable selection under minimax concave penalty.” *Annals of Statistics*, vol. 38, no. 2 (2010), pp. 894–942 (cit. on pp. 1, 5, 16).
- [4] Cun-Hui Zhang and Tong Zhang. “A general theory of concave regularization for high-dimensional sparse estimation problems.” *Statistical Science*, vol. 27, no. 4 (2012), pp. 576–593 (cit. on p. 1).
- [5] Zhaoran Wang, Han Liu, and Tong Zhang. “Optimal computational and statistical rates of convergence for sparse nonconvex learning problems.” *The Annals of Statistics*, vol. 42, no. 6 (2014), pp. 2164–2201 (cit. on p. 1).
- [6] Gilles Gasso, Alain Rakotomamonjy, and Stéphane Canu. “Recovering Sparse Signals With a Certain Family of Nonconvex Penalties and DC Programming.” *IEEE Transactions on Signal Processing*, vol. 57, no. 12 (2009), pp. 4686–4698 (cit. on p. 1).
- [7] Rick Chartrand. “Nonconvex splitting for regularized low-rank + sparse decomposition.” *IEEE Transactions on Signal Processing*, vol. 60, no. 11 (2012), pp. 5810–5819 (cit. on p. 1).
- [8] Yunzhang Zhu, Xiaotong Shen, and Wei Pan. “Simultaneous grouping pursuit and feature selection over an undirected graph.” *Journal of the American Statistical Association*, vol. 108, no. 502 (2013), pp. 713–725 (cit. on pp. 1, 3, 7).
- [9] Thomas Blumensath and Mike E. Davies. “Iterative thresholding for sparse approximations.” *Journal of Fourier Analysis and Applications*, vol. 14 (2008), pp. 629–654 (cit. on pp. 1, 3, 4).
- [10] Yiyuan She. “Thresholding-based iterative selection procedures for model selection and shrinkage.” *Electronic Journal of Statistics*, vol. 3 (2009), pp. 384–415 (cit. on pp. 1, 3, 4, 6).
- [11] Rahul Mazumder, Jerome H. Friedman, and Trevor Hastie. “SparseNet: Coordinate Descent With Nonconvex Penalties.” *Journal of the American Statistical Association*, vol. 106, no. 495 (2011), pp. 1125–1138 (cit. on pp. 1, 3, 7).
- [12] Pinghua Gong, Changshui Zhang, Zhaosong Lu, Jianhua Z. Huang, and Jieping Ye. “A General Iterative Shrinkage and Thresholding Algorithm for Non-convex Regularized Optimization Problems.” In: *ICML*. 2013 (cit. on pp. 1, 3, 4).
- [13] Seyoung Kim and Eric P. Xing. “Statistical Estimation of Correlated Genome Associations to a Quantitative Trait Network.” *PLoS Genetics*, vol. 5, no. 8 (2009), pp. 1–18 (cit. on pp. 1, 2).
- [14] Yaoliang Yu. “Better Approximation and Faster Algorithm Using the Proximal Average.” In: *NIPS*. 2013 (cit. on pp. 1, 3, 4, 11, 17, 18, 20).
- [15] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. “Proximal alternating linearized minimization for nonconvex and nonsmooth problems.” *Mathematical Programming, Series A*, vol. 146 (2014), pp. 459–494 (cit. on pp. 1, 3, 4, 18).
- [16] Philip M. Long and Rocco A. Servedio. “Random Classification Noise Defeats All Convex Potential Boosters.” *Machine Learning*, vol. 78, no. 3 (2010), pp. 287–304 (cit. on pp. 2, 3, 8).
- [17] Yaoliang Yu, Özlem Aslan, and Dale Schuurmans. “A Polynomial-time Form of Robust Regression.” In: *NIPS*. 2012 (cit. on pp. 2–4).
- [18] Linli Xu, Koby Crammer, and Dale Schuurmans. “Robust support vector machine training via convex outlier ablation.” In: *AAAI*. 2006 (cit. on pp. 2, 3, 7, 8).
- [19] Ronan Collobert, Fabian Sinz, Jason Weston, and Léon Bottou. “Trading Convexity for Scalability.” In: *ICML*. 2006, pp. 201–208 (cit. on pp. 2, 3, 7, 8).
- [20] Yichao Wu and Yufeng Liu. “Robust truncated hinge loss support vector machines.” *Journal of the American Statistical Association*, vol. 102, no. 479 (2007), pp. 974–983 (cit. on pp. 2, 3, 7).
- [21] Yufeng Liu, Xiaotong Shen, and Hani Doss. “Multicategory ψ -learning and support vector machine: computational tools.” *Journal of Computational and Graphical Statistics*, vol. 14, no. 1 (2005), pp. 219–236 (cit. on p. 2).
- [22] Robert Tibshirani. “Regression Shrinkage and Selection Via the Lasso.” *Journal of the Royal Statistical Society, Series B*, vol. 58 (1996), pp. 267–288 (cit. on p. 2).
- [23] Corinna Cortes and Vladimir Vapnik. “Support-Vector Networks.” *Machine Learning*, vol. 20 (1995), pp. 273–297 (cit. on p. 2).
- [24] Amir Beck and Marc Teboulle. “A fast iterative shrinkage-thresholding algorithm for linear inverse problems.” *SIAM Journal on Imaging Sciences*, vol. 2, no. 1 (2009), pp. 183–202 (cit. on p. 4).
- [25] Yurii Nesterov. “Gradient methods for minimizing composite functions.” *Mathematical Programming, Series B*, vol. 140 (2013), pp. 125–161 (cit. on p. 4).
- [26] Ralph Tyrell Rockafellar and Roger J-B Wets. *Variational Analysis*. Springer, 1998 (cit. on pp. 4, 5, 7, 10–16, 18).
- [27] Warren L. Hare. “A proximal average for nonconvex functions: A proximal stability perspective.” *SIAM Journal on Optimization*, vol. 20, no. 2 (2009), pp. 650–666 (cit. on p. 5).
- [28] Lou van den Dries and Chris Miller. “Geometric Categories and O-minimal Structures.” *Duke Mathematical Journal*, vol. 84, no. 2 (1996), pp. 497–540 (cit. on pp. 6, 17).
- [29] David D. Lewis, Yiming Yang, Tony G. Rose, Fan Li, G. Dietterich, and Fan Li. “RCV1: A new benchmark collection for text categorization research.” *Journal of Machine Learning Research*, vol. 5 (2004), pp. 361–397 (cit. on p. 8).
- [30] Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S. Sathya Keerthi, and S. Sundararajan. “A Dual Coordinate Descent Method for Large-scale Linear SVM.” In: *ICML*. 2008, pp. 408–415 (cit. on p. 8).