
Computing D-Stationary Points of ρ -Margin Loss SVM

Lai Tian

The Chinese University of Hong Kong
tianlai@se.cuhk.edu.hk

Anthony Man-Cho So

The Chinese University of Hong Kong
manchoseo@se.cuhk.edu.hk

Abstract

This paper is concerned with the algorithmic aspects of sharper stationarity of a nonconvex, nonsmooth, Clarke irregular machine learning model. We study the SVM problem with a ρ -margin loss function, which is the margin theory generalization bound of SVM introduced in the learning theory textbook by Mohri et al. [2018], and has been extensively studied in operations research, statistics, and machine learning communities. However, due to its nonconvex, nonsmooth, and irregular nature, none of the existing optimization methods can efficiently compute a d(irectional)-stationary point, which turns out to be also a local minimum, for the ρ -margin loss SVM problem. After a detailed discussion of various nonsmooth stationarity notions, we propose a highly efficient nonconvex semi-proximal ADMM-based scheme that provably computes d-stationary points and enjoys a local linear convergence rate. We report concrete examples to demonstrate the necessity of our assumptions. Numerical results verify the effectiveness of the new algorithm and complement our theoretical results.

1 INTRODUCTION

Nonconvex nonsmooth models are ubiquitous in modern statistical and machine learning tasks. Though holding strong expressive power, these “non”-problems pose serious computational challenges as most of the tools from smooth and convex analysis are inapplicable in this “non”-scenario. To handle their algorithmic design and theoretical analysis without sacrificing rigor,

Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS) 2022, Valencia, Spain. PMLR: Volume 151. Copyright 2022 by the author(s).

we borrow ideas and tools from variational analysis [Rockafellar and Wets, 2009], in which concepts that we are familiar with branch into many and calculus rules that we take for granted may no longer hold. Even choosing a proper first-order stationarity concept for these “non”-problems becomes a highly nontrivial thing (see Definition 2.7, and [Li et al., 2020]). These obstacles could be alleviated if the optimization problem is Clarke regular (see Definition 2.4), which contains smooth, convex, or more generally, weakly convex problems as special cases.

In this paper, we study the SVM model with a ρ -margin loss, which is a nonconvex, nonsmooth, and irregular problem in the sense of Clarke. Specifically, for a given set of data points $\{(\mathbf{x}_i, y_i) : i \in [n]\} \subseteq \mathbb{R}^d \times \{+1, -1\}$, we aim to solve

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{c}{2} \|\boldsymbol{\theta}\|^2 + \sum_{i=1}^n \phi_\rho(y_i \cdot (\mathbf{x}_i^\top \boldsymbol{\theta} + b)), \quad (1.1)$$

where the regularization parameter $c > 0$, and $\phi_\rho : \mathbb{R} \rightarrow \mathbb{R}$ is the ρ -margin loss defined by

$$\phi_\rho(u) := \min \left(1, \max \left(0, 1 - \frac{u}{\rho} \right) \right).$$

We notice that Problem (1.1) (particularly its $\rho = 1$ version, which is usually termed ramp loss) has been widely recognized in operations research [Brooks, 2011, Carrizosa et al., 2014, Wang et al., 2021], statistics [Shen et al., 2003, Wu and Liu, 2007, Liu et al., 2005], and machine learning [Huang et al., 2014, Keshet and McAllester, 2011, Collobert et al., 2006b,a, Ertekin et al., 2010, Suzumura et al., 2017, Maibing and Igel, 2015] communities as it provides better robustness against data outliers than the vanilla SVM. The ρ -margin loss (see Figure 1) we used in this paper can also be seen from the learning theory textbook by Mohri et al. [2018, Corollary 5.11], in which Problem (1.1) serves as the ρ -margin generalization bound for linear hypothesis set.

However, due to its “non”-properties, especially the irregularity, Problem (1.1) poses serious computational

challenges in the optimization aspect. Indeed, Maibing and Igel [2015] showed that finding a global minimizer for Problem (1.1) is NP-hard. Thus, we must be satisfied with computing a local minimizer instead, which is still a highly nontrivial task even for a quadratic function. Recently, Ahmadi and Zhang [2022] showed that finding a local minimizer is NP-hard generally, while the NP-hardness of detecting local optimality was proven decades ago by Murty and Kabadi [1987].

1.1 Prior Arts

In the literature concerning Problem (1.1), DC (Difference of Convex) Algorithms are probably the most popular strategies [Huang et al., 2014, Keshet and McAllester, 2011, Collobert et al., 2006b, Shen et al., 2003, Wu and Liu, 2007, Liu et al., 2005, Collobert et al., 2006a, Ertekin et al., 2010] as the nonsmooth part ϕ_ρ admits the following DC-decomposition:

$$\phi_\rho(u) = \max\left(1 - \frac{u}{\rho}, 0\right) - \max\left(-\frac{u}{\rho}, 0\right).$$

However, as pointed out by Nouiehed et al. [2019], these DCA-type algorithms only compute so-called DC-critical points (see Definition 2.7), which is a fairly weak notion of stationarity (see Definition 2.7 and remark afterwards) and depend on the DC-decomposition rather than variational properties of the problem. Indeed, it is easy to construct an example that is DC-critical but not stationary in any conventional sense (see Figure 2). By considering a polyhedral reformulation of Problem (1.1), Suzumura et al. [2017] proposed a homotopy-type algorithmic framework. However, in Section 4.3, we will show that their argument cannot guarantee sharper stationarity convergence. Recently, Wang et al. [2021] reported optimality conditions for Problem (1.1) (specialized to $\rho = 1$) with so-called P-stationarity and carefully computed the closed-form solution of the ramp loss proximal mapping. However, no known algorithm can efficiently compute points satisfying such conditions.

On the other front, for max-structured DC programs, Pang et al. [2017] and Cui et al. [2018] proposed enhanced DCA schemes for computing d-stationary points (see Definition 2.7), which is arguably the sharpest stationarity type for certain structured DC programs [Nouiehed et al., 2019]. Nevertheless, as mentioned by Pang et al. [2017, Section 7], these enhanced DCAs require solving an exponentially large number of subproblems in a single iteration step, which is not even suitable for hundreds of data points. Then, here comes the main question that this paper aims to answer:

Can we compute d-stationary points efficiently for SVM with ρ -margin loss?

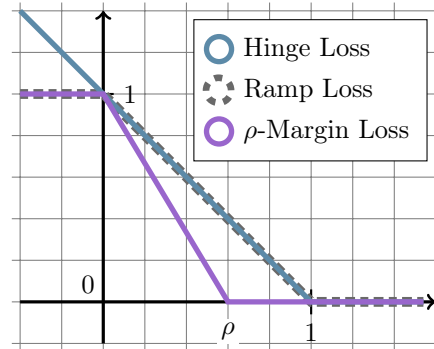


Figure 1: Hinge loss, ramp loss, and ρ -margin loss.

1.2 Contributions

We highlight the main contributions of this paper as follows:

- We propose a highly efficient *nonconvex* semi-proximal ADMM procedure, which provably computes a d-stationary point (which is also local minimal, see Proposition 3.5), of Problem (1.1).
- We show that the distance between the nonconvex ADMM generated sequence and the set of d-stationary points enjoys a local linear convergence rate to zero.
- We provide a detailed discussion on existing computing strategies and report concrete examples to demonstrate their incapability in sharp stationarity computing.

Notations. Most of the notations used in this paper are standard. Throughout this paper, scalars, vectors, and matrices are denoted by lowercase letters, boldface lowercase letters, and boldface uppercase letters, respectively; $\|\mathbf{x}\|$ is the Euclidean norm of \mathbf{x} ; $\|\mathbf{A}\|$ is the operator norm of \mathbf{A} ; for a set $S \subseteq \mathbb{R}^n$ and point \mathbf{x} , $\text{dist}(\mathbf{x}, S) := \inf_{\mathbf{v} \in S} \|\mathbf{v} - \mathbf{x}\|$; $A \oplus B$ denotes the direct sum of A and B ; $\mathbf{A} \odot \mathbf{B}$ is the Hadamard product of \mathbf{A} and \mathbf{B} ; we use $\text{prox} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ for the proximal operator, and use $\text{Prox} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ for set-valued proximal mapping (similarly, arg and Arg); $\text{cl } S$ denotes the closure of set S ; $\mathbb{B}_\epsilon(\mathbf{x}) := \{\mathbf{v} : \|\mathbf{v} - \mathbf{x}\| \leq \epsilon\}$; $\text{Co } S$ is the convex hull of set S ; $\sigma_i(\mathbf{A})$ is the i -th largest singular value; \mathbf{e}_i is the i -th column of identity matrix.

2 PRELIMINARIES

In this section, we will introduce the necessary background on variational analysis for our development. As the problem we are interested in is nonconvex and nonsmooth, conventional tools from smooth or convex analysis will be inapplicable.

Table 1: Algorithms for Problem (1.1), stationarity, computational complexity, and convergence rate.

Category	Algorithms	Stationarity	Per-Iter. Compl.	Rate
DCA (incl. CCCP)	[Huang et al., 2014] [Wu and Liu, 2007] [Brooks, 2011]	DC-critical	$\text{poly}(n, d)$	—
Nonmonotone MM	[Cui et al., 2018]	d-stationary	$O(2^n)$	—
Subgradient Method	[Davis et al., 2020] [Majewski et al., 2018]	C-stationary	$\text{poly}(n, d)$	—
Homotopy Algorithm	[Suzumura et al., 2014, 2017]	— ¹	$\text{poly}(n, d)$	—
Semi-proximal ADMM	Algorithm 1	d-stationary	$\text{poly}(n, d)$	local linear

¹ See Section 4.3.

2.1 Generalized Differentiation Theory

To tell the subtle difference between various stationary concepts and appreciate the benefit of computing sharper stationary points, we recall several notions from nonsmooth analysis in this section. We will confine the discussion to locally Lipschitz and directional differentiable functions as our Problem (3.1) is exactly in that class. To begin, let us recall the directional derivative and Clarke generalized subderivative of a locally Lipschitz function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, which will play fundamental roles in the analyses and definitions of other notions.

Definition 2.1 (directional derivative). *Given a point \mathbf{x} and direction \mathbf{d} , the directional derivative of f at \mathbf{x} in the direction \mathbf{d} is defined by*

$$f'(\mathbf{x}; \mathbf{d}) := \lim_{t \searrow 0} \frac{f(\mathbf{x} + t\mathbf{d}) - f(\mathbf{x})}{t}.$$

Definition 2.2 (Clarke subderivative). *Given a point \mathbf{x} and direction \mathbf{d} , the Clarke directional derivative of f at \mathbf{x} in the direction \mathbf{d} is defined by*

$$f^\circ(\mathbf{x}; \mathbf{d}) := \limsup_{\substack{\mathbf{x}' \rightarrow \mathbf{x} \\ t \searrow 0}} \frac{f(\mathbf{x}' + t\mathbf{d}) - f(\mathbf{x}')}{t}.$$

Using Rademacher's Theorem [Rockafellar and Wets, 2009, Theorem 9.60], the Clarke subdifferential of a locally Lipschitz function can be defined as follows [Rockafellar and Wets, 2009, Theorem 9.61]:

Definition 2.3 (Clarke subdifferential). *Given a point \mathbf{x} , the Clarke subdifferential of f at \mathbf{x} is defined by*

$$\partial_C f(\mathbf{x}) := \text{Co} \{ \mathbf{s} : \exists \mathbf{x}' \rightarrow \mathbf{x}, \nabla f(\mathbf{x}') \text{ exists, } \nabla f(\mathbf{x}') \rightarrow \mathbf{s} \}.$$

We remark that $\partial_C f(\mathbf{x})$ is convex and can be equivalently defined as the closed convex set whose support

function is $f^\circ(\mathbf{x}; \cdot)$. Now, we introduce the notion of Clarke regularity with Clarke subderivative.

Definition 2.4 (Clarke regularity [Clarke, 1990, Definition 2.3.4]). *A locally Lipschitz function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is regular at $\mathbf{x} \in \mathbb{R}^n$ if for every $\mathbf{d} \in \mathbb{R}^n$, the ordinary directional derivative $f'(\mathbf{x}; \mathbf{d})$ exists and*

$$f^\circ(\mathbf{x}; \mathbf{d}) = f'(\mathbf{x}; \mathbf{d}).$$

An important implication of regularity is the validity of various subdifferential calculus rules, which allows sharper characterization of first-order stationarity [Rockafellar and Wets, 2009, Theorem 10.6, Corollary 10.9]; see [Li et al., 2020] for a quick overview.

To characterize the local behavior of f from a variational perspective, we need the notion of Fréchet subdifferential [Rockafellar and Wets, 2009, Exercise 8.4], which has a close relation with d-stationarity (see Definition 2.7 and Proposition 3.2) and is defined as follows:

Definition 2.5 (Fréchet subdifferential). *Given a point \mathbf{x} , the Fréchet subdifferential of a locally Lipschitz and directional differentiable function f at \mathbf{x} is defined by*

$$\widehat{\partial} f(\mathbf{x}) := \{ \mathbf{s} : \mathbf{s}^\top \mathbf{d} \leq f'(\mathbf{x}; \mathbf{d}) \text{ for all } \mathbf{d} \}.$$

Though enjoying great variational properties, Fréchet subdifferential is not a robust notion as its graph is not closed (or, $\widehat{\partial} f$ is not outer semi-continuous, see Definition 2.9). Therefore, the following limiting version generalized subdifferential [Rockafellar and Wets, 2009, Theorem 8.3(b)] is sometimes more convenient:

Definition 2.6 (limiting subdifferential). *Given a point \mathbf{x} , the limiting subdifferential of a locally Lipschitz and directional differentiable function f at \mathbf{x} is defined by*

$$\partial f(\mathbf{x}) := \limsup_{\mathbf{x}' \rightarrow \mathbf{x}} \widehat{\partial} f(\mathbf{x}'),$$

where the outer limit is taken in the set-valued mapping sense (see [Rockafellar and Wets, 2009, 5(1)]).

2.2 Stationarity Notions

Now, let us introduce various stationarity notions for a nonconvex nonsmooth problem.

Definition 2.7. For a locally Lipschitz and directionally differentiable DC function $f = h - g$, a point $\bar{\mathbf{x}}$ is

(a) **local minimal** if there exists $\epsilon > 0$:

$$f(\mathbf{y}) \geq f(\bar{\mathbf{x}}), \quad \forall \mathbf{y} \in \mathbb{B}_\epsilon(\bar{\mathbf{x}});$$

(b) **d(irectional)-stationary** if:

$$f'(\bar{\mathbf{x}}; \mathbf{d}) \geq 0, \quad \forall \mathbf{d} \in \mathbb{R}^n;$$

(c) **l(imiting)-stationary** if:

$$0 \in \partial f(\bar{\mathbf{x}});$$

(d) **C(larke)-stationary** if:

$$0 \in \partial_C f(\bar{\mathbf{x}}) \Leftrightarrow f^\circ(\bar{\mathbf{x}}; \mathbf{d}) \geq 0, \quad \forall \mathbf{d} \in \mathbb{R}^n;$$

(e) **DC-critical** if:

$$\partial h(\bar{\mathbf{x}}) \cap \partial g(\bar{\mathbf{x}}) \neq \emptyset.$$

By [Rockafellar and Wets, 2009, Theorem 8.6, 8.15, 8.49], the following relation holds (see also [Li et al., 2020]):

$$\{(a)\} \subseteq \{(b)\} \subseteq \{(c)\} \subseteq \{(d)\} \subseteq \{(e)\}.$$

If the problem is Clarke regular, we have the following tighter relation:

$$\{(a)\} \subseteq \{(b)\} = \{(c)\} = \{(d)\} \subseteq \{(e)\}.$$

Remark 2.8. Note that the DC-criticality in Definition 2.7 is actually dependent on the DC-decomposition $f = h - g$ of the function. How to find the best DC-decomposition is still an important open problem [Migdalas and Pardalos, 2018]. An illustration of various stationarity notions is provided in Figure 2. See its caption for more information. Besides, it is notable here that for Problem (1.1), we actually have the sharper stationarity relation $\{(a)\} = \{(b)\}$, see Proposition 3.5.

2.3 Other Useful Notions

Below we record some other notions that will be useful for our reference and development. These definitions are taken from [Rockafellar and Wets, 2009, Definition 5.4] and [Rockafellar and Wets, 2009, Definition 1.23].

Definition 2.9 (outer semi-continuity). A set-valued mapping $S : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ is outer semi-continuous (osc) at $\bar{\mathbf{x}}$ if

$$\limsup_{\mathbf{x} \rightarrow \bar{\mathbf{x}}} S(\mathbf{x}) \subset S(\bar{\mathbf{x}}).$$

Definition 2.10 (prox-boundedness). A function $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ is prox-bounded if there exists $\lambda > 0$ such that $f(\text{Prox}_{\lambda f}(\mathbf{x})) > -\infty$ for some $\mathbf{x} \in \mathbb{R}^n$.

3 MAIN RESULTS

In this section, we reformulate Problem (1.1) into a compact form and consider the computation of d-stationary points (see Definition 2.3) of the following nonconvex, nonsmooth, and irregular problem:

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} F(\mathbf{w}) := f(\mathbf{w}) + g(\mathbf{Z}\mathbf{w}), \quad (3.1)$$

where we define $\mathbf{w} = \begin{bmatrix} \boldsymbol{\theta} \\ b \end{bmatrix} \in \mathbb{R}^{d+1}$,

$$f(\mathbf{w}) := \frac{c}{2} \|\mathbf{w}\|_{\mathbf{K}}^2 = \frac{c}{2} \mathbf{w}^\top \mathbf{K} \mathbf{w}, \quad g(\mathbf{q}) := \sum_{i=1}^n \phi_\rho(q_i),$$

and the i -th row of $\mathbf{Z} \in \mathbb{R}^{n \times (d+1)}$ as

$$\mathbf{z}_i = \begin{bmatrix} y_i \mathbf{x}_i \\ y_i \end{bmatrix} \in \mathbb{R}^{d+1}, \quad \mathbf{K} = \begin{bmatrix} \mathbb{I}_{d \times d} & \\ & 0 \end{bmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}.$$

3.1 Algorithm: Semi-proximal ADMM

To begin, we make the following surjectivity assumption on \mathbf{Z} , which is crucial for the dual characterization of d-stationarity and convergence of the ADMM scheme. We will elaborate on that assumption in Section 4.4.

Assumption 3.1 (surjective). The matrix $\mathbf{Z} \in \mathbb{R}^{n \times (d+1)}$ in Problem (3.1) is surjective and $\sigma := \sigma_n(\mathbf{Z}) > 0$.

The following characterization of a d-stationary point is convenient for both analysis and computation.

Proposition 3.2 (dual d-stationarity characterization). If Assumption 3.1 holds, then a point \mathbf{w} is d-stationary of F if and only if

$$0 \in \nabla f(\mathbf{w}) + \mathbf{A}^\top \widehat{\partial} g(\mathbf{A}\mathbf{w}).$$

Proof. By [Rockafellar and Wets, 2009, Exercise 8.4], \mathbf{w} is d-stationary if and only if $0 \in \widehat{\partial} F(\mathbf{w})$. Then, the relation is immediate from [Rockafellar and Wets, 2009, Corollary 8.8(c), Exercise 10.7]. \square

The main contribution of this section is algorithmic ideas that provably computes a point \mathbf{w} satisfying the dual characterization of d-stationarity (see Proposition 3.2). Before the new development, we highlight the key difficulties below:

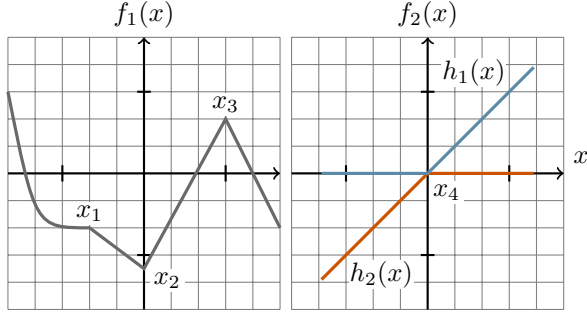


Figure 2: Illustration of various stationarities. In the left subfigure, $\{\text{local min.}\} = \{x_2\}$, $\{\text{d-stat.}\} = \{x_2\}$, $\{\text{l-stat.}\} = \{x_1, x_2\}$, and $\{\text{C-stat.}\} = \{x_1, x_2, x_3\}$. In the right subfigure, we consider a simple DC function $f_2(x) := h_1(x) + h_2(x) = \max\{x, 0\} - \max\{-x, 0\} = x$. Then, $x_4 = 0$ is DC-critical while it is obviously not stationary in any conventional (d-, l-, or even C-) sense.

- The first difficulty is the irregularity of the term g in the sense of Clarke. The lack of Clarke regularity may lead to invalidation of subdifferential calculus rules and thus difficulties in computing and characterizing sharper stationarity.
- The second difficulty, which is a consequence of the first one but specific to d-stationarity, is the lack of outer semi-continuity of the Fréchet subdifferential mapping. In other words, the graph of the set-valued mapping $\widehat{\partial}F(\cdot)$ is not closed, which implies that even if we generate a sequence satisfying $\mathbf{w}^k \rightarrow \mathbf{w}^*$, $\mathbf{v}^k \in \widehat{\partial}F(\mathbf{w}^k)$, $\mathbf{v}^k \rightarrow 0$, we cannot say $0 \in \widehat{\partial}F(\mathbf{w}^*)$, which invalidates the conventional convergence analysis that relies on the outer semi-continuity of a subdifferential mapping. An illustration of that issue with Figure 2 is that considering $x^k \nearrow x_1$, we have $\mathbf{v}^k \in \widehat{\partial}f(x^k)$, $\mathbf{v}^k \rightarrow 0$ as $k \rightarrow +\infty$. So $0 \in \partial f(x_1)$. However, it holds that $0 \notin \widehat{\partial}f(x_1) = \emptyset$.

To proceed, with the operator splitting technique, we introduce the following reformulation of Problem (3.1):

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^{d+1}, \mathbf{q} \in \mathbb{R}^n} \quad & Q(\mathbf{w}, \mathbf{q}) := f(\mathbf{w}) + g(\mathbf{q}) \\ \text{s.t.} \quad & \mathbf{Z}\mathbf{w} = \mathbf{q}. \end{aligned} \quad (\text{MSVM})$$

Then, we tackle the above reformulation with a nonconvex semi-proximal ADMM scheme. The augmented Lagrangian for Problem (MSVM) can be written as:

$$L_\beta(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) := f(\mathbf{w}) + g(\mathbf{q}) + \langle \boldsymbol{\lambda}, \mathbf{Z}\mathbf{w} - \mathbf{q} \rangle + \frac{\beta}{2} \|\mathbf{Z}\mathbf{w} - \mathbf{q}\|^2,$$

where $\beta \geq 0$ is the dual step size, $\boldsymbol{\lambda} \in \mathbb{R}^n$ is the Lagrange multiplier. To ensure convergence in such

a nonconvex scenario, we need a cautiously chosen dual step size β . The overall iteration scheme can be summarized into the following diagram:

Algorithm 1 Nonconvex Semi-proximal ADMM

Input: $\mathbf{Z} \in \mathbb{R}^{(d+1) \times n}$, choose $\gamma > 0, c > 0$ and $\beta > 1 + \max \left\{ \frac{2 \left(1 + \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 \right)}{c\sigma^2}, \frac{2}{n\sigma^2}, \frac{4}{cn\sigma^4} \right\}$.
for all $k \in \{0, 1, 2, \dots\}$ **do**

$$\mathbf{q}^{k+1} \in \text{Arg min}_{\mathbf{q}} L_\beta(\mathbf{w}^k, \mathbf{q}, \boldsymbol{\lambda}^k) + \frac{\gamma}{2} \|\mathbf{q} - \mathbf{q}^k\|^2,$$

$$\mathbf{w}^{k+1} = \arg \min_{\mathbf{w}} L_\beta(\mathbf{w}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^k),$$

$$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k + \beta(\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1}).$$

end for

In the sequel, we will detail the practical updating computation of each variable.

3.1.1 \mathbf{q}^{k+1} Updating Step

To compute \mathbf{q}^{k+1} , we note that the \mathbf{q} -subproblem can be rewritten as the following nonconvex proximal problem:

$$\mathbf{q}^{k+1} \in \text{Prox}_{\frac{1}{\beta+\gamma}g} \left(\frac{\beta}{\beta+\gamma} \left(\mathbf{Z}\mathbf{w}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k \right) + \frac{\gamma}{\beta+\gamma} \mathbf{q}^k \right),$$

where $\text{Prox}_{\frac{1}{\beta+\gamma}g} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ is a set-valued mapping due to the nonconvexity of g . In general, evaluation of a nonconvex proximal operator can be extremely difficult. Note that $g(\mathbf{v}) = \sum_{i=1}^n \phi_\rho(v_i)$ is separable and

$$\text{Prox}_{\lambda g}(\mathbf{v}) = \bigoplus_{i=1}^n \text{Prox}_{\lambda \phi_\rho}(v_i).$$

We provide the following closed-form solution for the one-dimensional separated term $\text{Prox}_{\lambda \phi_\rho}$:

Lemma 3.3 (proximal operator). *For $0 < \lambda < 2\rho^2$, it holds*

$$\text{Prox}_{\lambda \phi_\rho}(v) = \begin{cases} \{v\} & \text{for } v < -\frac{\lambda}{2\rho}, \\ \left\{ v, v + \frac{\lambda}{\rho} \right\} & \text{for } v = -\frac{\lambda}{2\rho}, \\ \left\{ v + \frac{\lambda}{\rho} \right\} & \text{for } -\frac{\lambda}{2\rho} < v \leq \rho - \frac{\lambda}{\rho}, \\ \{\rho\} & \text{for } \rho - \frac{\lambda}{\rho} < v \leq \rho, \\ \{v\} & \text{for } v > \rho. \end{cases}$$

For $0 < 2\rho^2 \leq \lambda$, we have

$$\text{Prox}_{\lambda \phi_\rho}(v) = \begin{cases} \{v\} & \text{for } v < \rho - \sqrt{2\lambda}, \\ \{v, \rho\} & \text{for } v = \rho - \sqrt{2\lambda}, \\ \{\rho\} & \text{for } \rho - \sqrt{2\lambda} < v < \rho, \\ \{v\} & \text{for } v \geq \rho. \end{cases}$$

3.1.2 w^{k+1} Updating Step

In the w -subproblem, to compute w^{k+1} , we need to solve the following linear system:

$$(c\mathbf{K} + \beta\mathbf{Z}^\top\mathbf{Z})\mathbf{w}^{k+1} = \mathbf{Z}^\top(\beta\mathbf{q}^{k+1} - \boldsymbol{\lambda}^k).$$

The following proposition certifies the existence and uniqueness of solution w^{k+1} :

Proposition 3.4. *Let $p = n\beta \cdot \|\frac{1}{n}\sum_{i=1}^n \mathbf{x}_i\|^2$. For any $c > 0, \beta > 0$, we have*

$$(c\mathbf{K} + \beta\mathbf{Z}^\top\mathbf{Z}) \succcurlyeq \mu \cdot \mathbf{I},$$

where $\mu > 0$ and defined by

$$\mu := \frac{1}{2} \left(c + n\beta + p - \sqrt{(c - n\beta)^2 + 2p(c + n\beta) + p^2} \right),$$

which reduces to $\mu = \min\{c, \beta n\}$ if the data points are centered and thus $p = 0$.

We pre-compute $(c\mathbf{K} + \beta\mathbf{Z}\mathbf{Z}^\top)^{-1}\mathbf{Z}^\top \in \mathbb{R}^{n \times (d+1)}$, which by careful matrix partition will cost $O(n^2d)$. But in every iteration, we only need to do matrix-vector product computation and never re-compute the matrix inverse again, which only costs $O(nd)$.

3.2 Convergence Analysis

Following the idea of [Cui et al., 2020, Proposition 4.1], for our specific MSVM problem, we claim that d-stationarity is necessary and sufficient for local optimality.

Proposition 3.5. *A point w is local minimal for Problem (3.1) if and only if w is d-stationary.*

It is well-known that the first-order condition, e.g., d-stationarity, is only necessary for local optimality (see [Rockafellar and Wets, 2009, Theorem 10.1]). The sufficient optimality conditions usually require some sort of convexity [Rockafellar and Wets, 2009, Theorem 8.15] or second-order information [Rockafellar and Wets, 2009, Theorem 13.24]. Luckily, as shown in Proposition 3.5, for the ρ -margin loss SVM problem, d-stationarity is sufficient for local minimality, which is the tightest stationarity that we can hope for without knowing more global information.

We notice that recently [Suzumura et al., 2017, Theorem 4] report a KKT-type necessary and sufficient condition for local minimality of Problem (3.1). However, it turns out that their argument is flawed. We will provide a detailed discussion on that issue and a concrete counterexample to [Suzumura et al., 2017, Theorem 4] in Section 4.3.

In view of the above, an efficient algorithm that can provably generate d-stationary points for Problem (3.1)

is highly desirable. We will show that the new nonconvex ADMM scheme in Algorithm 1 is capable of doing so with a local linear convergence rate.

Theorem 3.6. *Let $\gamma > 0, c > 0$, and*

$$\beta > 1 + \max \left\{ \frac{2 \left(1 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right\|^2 \right)}{c\sigma^2}, \frac{2}{n\sigma^2}, \frac{4}{cn\sigma^4} \right\}.$$

For the sequence $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1, the following holds:

- When $k \rightarrow \infty$, we have sequential convergence $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \rightarrow (\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ and \mathbf{w}^* is a d-stationary point of Problem (3.1).
- For any $0 \leq T < \infty$, we have

$$\min_{k \in [T]} \text{dist} \left(0, \partial L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \right) \leq \frac{C_1}{\sqrt{T}},$$

where $C_1 := \tau_2 \cdot \sqrt{\frac{L_\beta(\mathbf{w}^0, \mathbf{q}^0, \boldsymbol{\lambda}^0) - F^*}{\tau_1}}$. See Section 5 for definitions of τ_1, τ_2 .

- Let \mathcal{S}_d be the set of d-stationary points of Problem (3.1). There exist $C_2 < \infty, \rho \in [0, 1)$ and finite \bar{k} such that for all $k \geq \bar{k}$:

$$\text{dist}(\mathbf{w}^k, \mathcal{S}_d) \leq C_2 \rho^k.$$

The constants C_2 and ρ can be determined from [Han et al., 2018, Theorem 2] and Hoffman's error bound for linear system [Dontchev and Rockafellar, 2009, Lemma 3C.4]. We will sketch the main idea of proof in Section 5 and defer the formal argument to Appendix C.

4 DISCUSSION

In this section, we will first construct concrete examples that lead to failure of computing sharp stationary points by the subgradient, DCA, and homotopy methods. Then, we will explain the subtleties of the technical Assumption 3.1 and its consequences on optimality conditions and algorithmic convergence.

4.1 Subgradient Method

We consider the subgradient method with step size $t_k < \frac{1}{c}$ and iteration scheme:

$$\mathbf{w}^{k+1} := \mathbf{w}^k - t_k \mathbf{v}_k, \quad \mathbf{v}^k \in \partial F(\mathbf{w}^k).$$

Then, we introduce the following set:

$$S := \{\mathbf{w} \in \mathbb{R}^{d+1} : \mathbf{Z}\mathbf{K}\mathbf{w} < 0, w_{d+1} = 0\}.$$

Assuming $\mathbf{w}_k \in S$, we have $\partial F(\mathbf{w}^k) = c \cdot \mathbf{K}\mathbf{w}^k$, which implies that

$$\mathbf{Z}\mathbf{K}\mathbf{w}^{k+1} = \mathbf{Z}(1 - t_k \cdot c)\mathbf{K}\mathbf{w}^k < 0.$$

Thus, the sequence generated by the subgradient method cannot escape from the set S in finite time. Considering any accumulation point \mathbf{w}^* of the sequence $\{\mathbf{w}^k\}_k$, by $\mathbf{w}^* \in \text{cl}S$, we will show that if $0 \in \partial F(\mathbf{w}^*)$ then $\mathbf{K}\mathbf{w}^* = 0$ and $0 \notin \widehat{\partial}F(\mathbf{w}^*)$. So, any accumulation point \mathbf{w}^* cannot be locally minimal due to [Rockafellar and Wets, 2009, Theorem 10.1]. To this end, note that

$$0 \in (\mathbf{w}^*)^\top \partial F(\mathbf{w}^*) \subseteq c \cdot \|\mathbf{w}^*\|_{\mathbf{K}}^2 + (\mathbf{Z}\mathbf{w}^*)^\top \partial g(\mathbf{Z}\mathbf{w}^*).$$

By $\mathbf{w}^* \in \text{cl}S$, it holds that $\mathbf{Z}\mathbf{w}^* \leq 0$ and thus $\partial g(\mathbf{Z}\mathbf{w}^*) \subseteq \{0, -1\}^n$. This gives $-c \cdot \|\mathbf{w}^*\|_{\mathbf{K}}^2 \geq 0$ and thus $\|\mathbf{w}^*\|_{\mathbf{K}} = 0$ and $\mathbf{w}^* = 0$. Hence, $0 \notin \widehat{\partial}F(\mathbf{w}^*) = \mathbf{Z}^\top \widehat{\partial}g(\mathbf{0}) = \emptyset$.

4.2 DCA-type Algorithms

For DCA-type algorithms, the failure case is similar to that for the subgradient method. Consider the DC decomposition of $\phi_\rho(\mathbf{z}_i^\top \mathbf{w}) = \phi_{\rho,1}(\mathbf{z}_i^\top \mathbf{w}) - \phi_{\rho,2}(\mathbf{z}_i^\top \mathbf{w})$, where $\phi_{\rho,1}(\mathbf{z}_i^\top \mathbf{w}) = \max\left(1 - \frac{\mathbf{z}_i^\top \mathbf{w}}{\rho}, 0\right)$, $\phi_{\rho,2}(\mathbf{z}_i^\top \mathbf{w}) = \max\left(-\frac{\mathbf{z}_i^\top \mathbf{w}}{\rho}, 0\right)$. Then, for $\mathbf{w} = \mathbf{0}$, we have

$$0 \in \partial \phi_{\rho,1}(0) \cap \phi_{\rho,2}(0),$$

which, by summing up, implies that $\mathbf{w} = 0$ is DC-critical for Problem (3.1). However, as shown in Section 4.1, $\mathbf{w} = 0$ cannot be a d-stationary point.

4.3 Homotopy Algorithm

Recently, Suzumura et al. [2017, Theorem 4] reported a KKT-type necessary and sufficient condition for local minimality of Problem (3.1). However, the following example shows that their KKT-type condition is in fact not necessary:

Example 4.1. Let $d = 1, n = 6, c = 1, \rho = 1$. We consider the following data points $\{(x_i, y_i) : i \in [n]\} := \{(1, +1), (1, +1), (-1, -1), (-1, -1), (0, +1), (0, -1)\} \subseteq \mathbb{R}^d \times \{+1, -1\}$. We claim that $\mathbf{w} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ is a global minimizer for Problem (3.1) on $\{(x_i, y_i) : i \in [n]\}$. However, \mathbf{w} does not satisfy the condition in [Suzumura et al., 2017, Theorem 4]. The proof of that claim is deferred to Appendix D due to space limitation.

The reason for the invalidation of [Suzumura et al., 2017, Theorem 4] is that, in the proof of [Suzumura et al., 2017, Theorem 3], the failure of a primal-dual pair $(f_{\mathcal{P}}^*, \boldsymbol{\alpha})$ being a KKT point does not imply the sub-optimality of primal solution $f_{\mathcal{P}}^*$.

4.4 On Assumption 3.1

The surjectivity assumption is crucial in our subdifferential d-stationary characterization (Proposition 3.2), without which efficient computation of d-stationary points would be quite hard if not impossible. An exception is the nonmonotone MM algorithm [Pang et al., 2017, Cui et al., 2018], which computes d-stationary points using the primal directional derivative characterization Definition 2.7(b) without Assumption 3.1. However, as pointed out in [Pang et al., 2017, Section 7], these DCA-type algorithms may need to solve an exponentially large number of subproblems in a single iteration step. We do not notice any algorithm that can efficiently compute d-stationary points for Problem (1.1) without the surjectivity assumption.

On the other front, the surjectivity assumption is necessary for a nonconvex ADMM scheme to converge in general. Concrete examples can be constructed [Li and Pong, 2015, Example 7] in theory, and oscillating behavior can be observed in numerical experiments. Besides, we notice that in the nonconvex ADMM literature, the surjectivity assumption on mapping \mathbf{Z} seems pervasive and necessary [Jiang et al., 2019, Wang et al., 2019, Li and Pong, 2015].

5 PROOF SKETCH

In this section, we lay out a sketch for the proof of the main result Theorem 3.6. The formal argument is deferred to Appendix C.

5.1 Limiting Stationary Convergence

The first step is to prove the sequential convergence of $\{\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k\}_k$. To this end, we will use the augmented Lagrangian $L_\beta(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda})$ as a Lyapunov function. We first show that the sequence $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ is bounded and $L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) > -\infty$ for any $k \in \mathbb{N}$.

Lemma 5.1 (bounded and proper). For $\beta > \frac{2}{c\sigma^2}, \gamma > 0, c > 0$, and $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1, there exists an $R > 0$ such that $\|(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\| \leq R$ and $L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) > -\infty$ for any $k \in \mathbb{N}$.

Then, we show that the update rules in Algorithm 1 satisfy the following sufficient decrease and limiting safeguard properties:

Lemma 5.2 (sufficient decrease). For $\gamma > 0, c > 0, \beta > 1 + \max\left\{\frac{2(1 + \|\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i\|^2)}{c\sigma^2}, \frac{2}{n\sigma^2}, \frac{4}{cn\sigma^4}\right\}$, and $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1, there exists a constant $\tau_1 > 0$ only depending on $\{\sigma, c, n, \gamma, \beta\}$, such

that

$$\begin{aligned} & L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \\ & \leq -\tau_1 \|(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - (\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\|^2. \end{aligned}$$

Lemma 5.3 (limiting safeguard). *There exists a $\tau_2 > 0$ such that*

$$\begin{aligned} & \text{dist}\left(0, \partial L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1})\right) \\ & \leq \tau_2 \|(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - (\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\|. \end{aligned}$$

Finally, by showing that the Lyapunov function satisfies the KL property (see Definition C.1 in Appendix C), and with the general convergence result in [Attouch et al., 2013, Theorem 2.9], we establish the sequential convergence of Algorithm 1:

Lemma 5.4 (limiting stationarity convergence). *Let $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ be the sequence generated by Algorithm 1. As $k \rightarrow +\infty$, we have $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \rightarrow (\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ with \mathbf{w}^* is limiting-stationary, that is,*

$$0 \in c\mathbf{K}\mathbf{w}^* + \mathbf{Z}^\top \partial g(\mathbf{Z}\mathbf{w}^*).$$

5.2 Stationarity Refinement

For the second step, we will refine the stationary characterization of the limiting point $(\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ by showing that \mathbf{w}^* satisfies the following sharper condition:

$$0 \in c\mathbf{K}\mathbf{w}^* + \mathbf{Z}^\top \widehat{\partial} g(\mathbf{Z}\mathbf{w}^*),$$

which by Proposition 3.2 indicates that \mathbf{w}^* is d-stationary. The crux of the proof is to bypass the lack of outer semi-continuity of the subdifferential mapping $\widehat{\partial} g$ in the optimality condition. Our strategy is to wrap $\widehat{\partial} g : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ up with the nonconvex proximal mapping $\text{Prox}_{\lambda\phi_g} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$, whose outer semi-continuity can be obtained from the prox-boundedness of g (see Definition 2.10). See Appendix C for details.

5.3 Local Linear Convergence

Third, we will establish the local linear convergence rate of Algorithm 1 with a manifold identification argument. The key observation from the analytic solution of $\text{Prox}_{\lambda\phi_\rho}$ (see Lemma 3.3) is that for any $\lambda > 0, \rho > 0, q \in \mathbb{R}$, and $q \in \text{Prox}_{\lambda\phi_\rho}(v)$, where $v \in \mathbb{R}$ is arbitrary, we have $|q| \geq \min\left\{\frac{\lambda}{2\rho}, \rho\right\}$. This uniform lower bound allows us to find a \bar{k} such that for any $k \geq \bar{k}$, the iteration $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)$ from Algorithm 1 can be viewed as running Algorithm 1 on a special convex piecewise-linear quadratic (PLQ) problem [Rockafellar and Wets, 2009, Definition 10.20]. Moreover, the KKT mapping of this convex PLQ problem is polyhedral, which implies that its inverse mapping satisfies the so-called

calmness condition (see [Han et al., 2018, Definition 2]). By [Han et al., 2018, Theorem 2], the sequence $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)$ converges linearly to the set of KKT points of the convex problem. After showing that all these KKT points are indeed d-stationary for Problem (3.1), the proof is complete.

6 NUMERICAL RESULTS

In this section, we present numerical results to demonstrate the effectiveness of the new algorithmic scheme. All simulations are implemented² with MATLAB R2021a on a Windows 10 PC with Intel Core i7-11700K, and 32 GB RAM. We will show the numerical performance of the subgradient method, DCA, SpADMM in Algorithm 1, and a hybrid method (Algorithm 3) that runs DCA until convergence and then starts SpADMM from that DC-critical point (see Appendix E).

We first consider a toy setting with multi-dimensional Gaussian inputs and random binary labels. Let $\mathbf{x}_i \sim N(0, \mathbf{I}_d)$ and $y_i \sim 2 \cdot \text{Bernoulli}(\frac{1}{2}) - 1$ for any $i \in [n]$, where $n = 50, d = 100$. We run subGD (subgradient), DCA, SpADMM (Algorithm 1) and Hybrid Algorithm 3 with the same Gaussian initial point \mathbf{w}^0 for 150 iteration steps and record the objective function value and the augmented Lagrangian function value of Algorithm 1. We report the difference between the current objective function value and the best objective value among all algorithms in Figure 3. It is notable that the reflection-like curve for augmented Lagrangian of Algorithm 1 is due to the fact that the augmented Lagrangian value could be smaller than the best objective function value. From Figure 3, we highlight two interesting observations:

- It is easy to see in Figure 3 that both Algorithm 1 and Algorithm 3 outperform subGD and DCA, which is because the semi-proximal ADMM scheme escapes from non-locally minimal stationary points and will eventually converge only to a locally minimal one.
- When running the Hybrid Algorithm 3, we can see it decreases the objective value even further from the DC-critical point that DCA converges to. In other words, SpADMM escaped from the DC-critical point that is not locally minimal. This provides empirical evidence that computing a sharper kind of stationarity does not only have advantages in theory but also in numerical performance.

We also conduct simulations on real-world datasets, which are all from the LIBSVM datasets collection.

²See <https://github.com/icety3/rho-marginSVM>

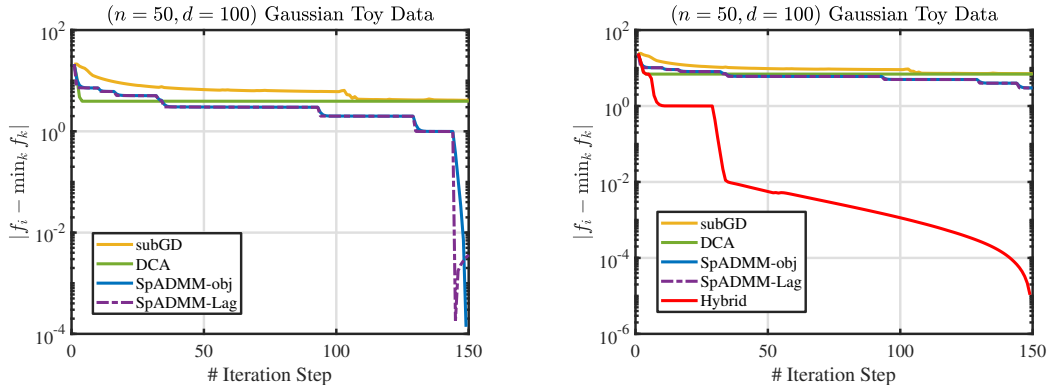


Figure 3: Performance of subGD (subgradient), DCA, SpADMM (Algorithm 1), and Hybrid Algorithm 3 on Gaussian Toy Data. The min in log-scaled y -axis is taken over all aforementioned algorithms.

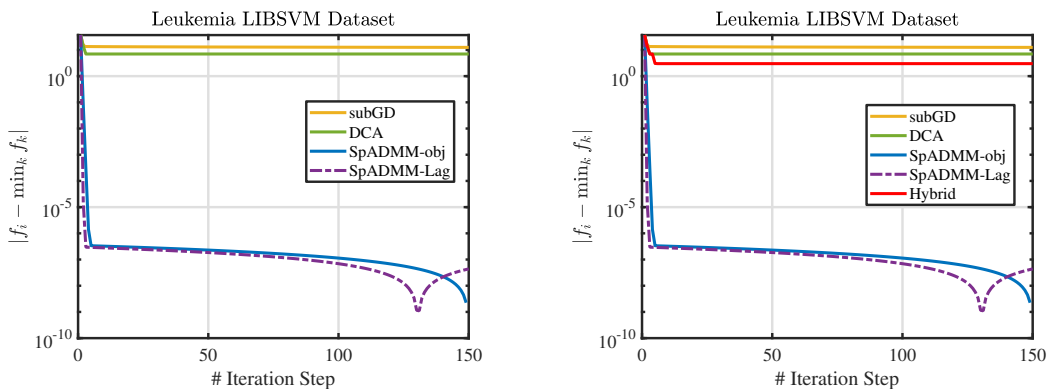


Figure 4: Performance of subGD (subgradient), DCA, SpADMM (Algorithm 1), and Hybrid Algorithm 3 on Real-world Data. The min in log-scaled y -axis is taken over all aforementioned algorithms.

The results are reported in Figure 4 (see also Figure 5 in Appendix F.1). It is easy to see that our observations from toy data remain true for real-world data. Besides, we note that the Hybrid Algorithm 3 does not necessarily outperform SpADMM in Algorithm 1 (see the results for Leukemia in Figure 4 and Duke breast-cancer in Appendix F.1, Figure 5).

7 LIMITATIONS & REMARKS

In this paper, we proposed a highly efficient nonconvex semi-proximal ADMM-based scheme that provably computes d -stationary points for the ρ -margin loss SVM problem, which is nonconvex, nonsmooth, and Clarke irregular. We further show that the local convergence rate of the nonconvex ADMM is linear. Our development certainly has its limitations and might inspire further investigation for new tools and ideas for nonconvex nonsmooth models in machine learning.

- One limitation is of course the surjectivity Assumption 3.1 in our analysis. Without this assumption, we are not aware of any algorithm that can efficiently compute d -stationary points for Prob-

lem (1.1). One might guess that Assumption 3.1 can be relaxed with a better analysis of the ADMM scheme. But we think a more intriguing direction is to find a better computable d -stationary characterization (e.g., Proposition 3.2). In [Pang et al., 2017], their enhanced DCA algorithm works on a primal directional derivative condition, which needs less regularity conditions to hold but causes serious computational trouble in the updating steps. Our dual subdifferential condition decomposes the non-smooth parts into separable pieces and leads to efficient computing with extra cost on more assumptions.

- The other limitation is that our scheme is required to solve a nonconvex subproblem (i.e., to compute $\text{Prox}_{\lambda g}$) to optimality, which is crucial to do a limiting argument without osc in $\widehat{\partial}g$. For Problem (1.1), we do these computations by hand and the optimal solution is carefully computed as a closed-form (see Lemma 3.3). However, for a general nonconvex problem, it is unrealistic to compute its proximal mapping. Finding a computational cheaper mechanism to do limiting argument for $\widehat{\partial}g$ would be interesting.

References

- Amir Ali Ahmadi and Jeffrey Zhang. On the complexity of finding a local minimizer of a quadratic function over a polytope. *Mathematical Programming*, 2022.
- Hedy Attouch, Jérôme Bolte, and Benar Fux Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized gauss-seidel methods. *Mathematical Programming*, 137(1):91–129, 2013.
- J. Paul Brooks. Support vector machines with the ramp loss and the hard margin loss. *Operations Research*, 59(2):467–479, 2011.
- Emilio Carrizosa, Amaya Nogales-Gómez, and Dolores Romero Morales. Heuristic approaches for support vector machines with the ramp loss. *Optimization Letters*, 8(3):1125–1135, 2014.
- Augustin-Louis Cauchy. *Exercices de Mathématiques*, volume 3. De Bure frères, 1828.
- Frank H. Clarke. *Optimization and Nonsmooth Analysis*. SIAM, 1990.
- Ronan Collobert, Fabian Sinz, Jason Weston, and Léon Bottou. Trading convexity for scalability. In *International Conference on Machine Learning*, pages 201–208. PMLR, 2006a.
- Ronan Collobert, Fabian Sinz, Jason Weston, Léon Bottou, and Thorsten Joachims. Large scale transductive SVMs. *Journal of Machine Learning Research*, 7(8), 2006b.
- Ying Cui, Jong-Shi Pang, and Bodhisattva Sen. Composite difference-max programs for modern statistical estimation problems. *SIAM Journal on Optimization*, 28(4):3344–3374, 2018.
- Ying Cui, Tsung-Hui Chang, Mingyi Hong, and Jong-Shi Pang. A study of piecewise linear-quadratic programs. *Journal of Optimization Theory and Applications*, 186(2):523–553, 2020.
- Damek Davis, Dmitriy Drusvyatskiy, Sham Kakade, and Jason D. Lee. Stochastic subgradient method converges on tame functions. *Foundations of Computational Mathematics*, 20(1):119–154, 2020.
- Asen L. Dontchev and R. Tyrrell Rockafellar. *Implicit Functions and Solution Mappings*, volume 543. Springer, 2009.
- Seyda Ertekin, Leon Bottou, and C. Lee Giles. Nonconvex online support vector machines. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(2):368–381, 2010.
- Francisco Facchinei and Jong-Shi Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer Science & Business Media, 2007.
- Deren Han, Defeng Sun, and Liwei Zhang. Linear rate convergence of the alternating direction method of multipliers for convex composite programming. *Mathematics of Operations Research*, 43(2):622–637, 2018.
- D. Hertz, C.S. Adjiman, and C.A. Floudas. Two results on bounding the roots of interval polynomials. *Computers & Chemical Engineering*, 23(9):1333–1339, 1999.
- Xiaolin Huang, Lei Shi, and Johan A.K. Suykens. Ramp loss linear programming support vector machine. *Journal of Machine Learning Research*, 15(1):2185–2211, 2014.
- Bo Jiang, Tianyi Lin, Shiqian Ma, and Shuzhong Zhang. Structured nonconvex and nonsmooth optimization: algorithms and iteration complexity analysis. *Computational Optimization and Applications*, 72(1):115–157, 2019.
- Joseph Keshet and David McAllester. Generalization bounds and consistency for latent structural probit and ramp loss. In *Advances in Neural Information Processing Systems*, volume 24, 2011.
- Krzysztof Kurdyka. On gradients of functions definable in o-minimal structures. In *Annales de l’institut Fourier*, volume 48, pages 769–783, 1998.
- Guoyin Li and Ting Kei Pong. Global convergence of splitting methods for nonconvex composite optimization. *SIAM Journal on Optimization*, 25(4):2434–2460, 2015.
- Jiajin Li, Anthony Man-Cho So, and Wing-Kin Ma. Understanding notions of stationarity in non-smooth optimization. *IEEE Signal Processing Magazine*, 37(5):18–31, 2020.
- Yufeng Liu, Xiaotong Shen, and Hani Doss. Multi-category ψ -learning and support vector machine: computational tools. *Journal of Computational and Graphical Statistics*, 14(1):219–236, 2005.
- Søren Frejstrup Maibing and Christian Igel. Computational complexity of linear large margin classification with ramp loss. In *Artificial Intelligence and Statistics*, pages 259–267. PMLR, 2015.
- Szymon Majewski, Błażej Miasojedow, and Eric Moulines. Analysis of nonsmooth stochastic approximation: the differential inclusion approach. *arXiv preprint arXiv:1805.01916*, 2018.
- Athanasios Migdalas and P.M. Pardalos. A note on open problems and challenges in optimization theory and algorithms. In *Open Problems in Optimization and Data Analysis*, pages 1–8. Springer, 2018.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT press, 2018.

- Katta G. Murty and Santosh N. Kabadi. Some NP-complete problems in quadratic and nonlinear programming. *Mathematical Programming*, 39(2):117–129, 1987.
- Yurii Nesterov. *Lectures on Convex Optimization*, volume 137 of *Springer Optimization and Its Applications*. Springer Nature, second edition, 2018.
- Maher Nouiehed, Jong-Shi Pang, and Meisam Razaviyayn. On the pervasiveness of difference-convexity in optimization and statistics. *Mathematical Programming*, 174(1):195–222, 2019.
- Jong-Shi Pang, Meisam Razaviyayn, and Alberth Alvarado. Computing B-stationary points of nonsmooth DC programs. *Mathematics of Operations Research*, 42(1):95–118, 2017.
- Stephen M Robinson. Some continuity properties of polyhedral multifunctions. In *Mathematical Programming at Oberwolfach*, pages 206–214. Springer, 1981.
- R. Tyrrell Rockafellar and Roger J.-B. Wets. *Variational Analysis*, volume 317. Springer Science & Business Media, 2009.
- Xiaotong Shen, George C. Tseng, Xuegong Zhang, and Wing Hung Wong. On ψ -learning. *Journal of the American Statistical Association*, 98(463):724–734, 2003.
- Shinya Suzumura, Kohei Ogawa, Masashi Sugiyama, and Ichiro Takeuchi. Outlier path: a homotopy algorithm for robust SVM. In *International Conference on Machine Learning*, pages 1098–1106. PMLR, 2014.
- Shinya Suzumura, Kohei Ogawa, Masashi Sugiyama, Masayuki Karasuyama, and Ichiro Takeuchi. Homotopy continuation approaches for robust SV classification and regression. *Machine Learning*, 106(7):1009–1038, 2017.
- Huajun Wang, Yuanhai Shao, and Naihua Xiu. Proximal operator and optimality conditions for ramp loss SVM. *Optimization Letters*, pages 1–16, 2021.
- Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of ADMM in nonconvex nonsmooth optimization. *Journal of Scientific Computing*, 78(1):29–63, 2019.
- Yichao Wu and Yufeng Liu. Robust truncated hinge loss support vector machines. *Journal of the American Statistical Association*, 102(479):974–983, 2007.

Supplementary Material: Computing D-Stationary Points of ρ -Margin Loss SVM

A Optimality Condition

For preparation, the optimality conditions for every update iteration of Algorithm 1 can be written as:

$$\begin{cases} 0 \in \widehat{\partial}g(\mathbf{q}^{k+1}) + \beta \left(\mathbf{q}^{k+1} - \mathbf{Z}\mathbf{w}^k - \frac{1}{\beta}\boldsymbol{\lambda}^k \right) + \gamma(\mathbf{q}^{k+1} - \mathbf{q}^k), & \text{(A.1a)} \\ 0 = c\mathbf{K}\mathbf{w}^{k+1} + \beta\mathbf{Z}^\top \left(\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1} + \frac{1}{\beta}\boldsymbol{\lambda}^k \right), & \text{(A.1b)} \\ \boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k + \beta(\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1}). & \text{(A.1c)} \end{cases}$$

B Technical Lemmas

Lemma B.1. For any $\mathbf{A} \in \mathbb{R}^{n \times (d+1)}$, and any $\mathbf{p} \in \mathbb{R}^n$, it holds $\|\mathbf{A}^\top \mathbf{p}\| \geq \sigma_{\min}(\mathbf{A}) \cdot \|\mathbf{p}\|$.

Proof. Let the thin-SVD of \mathbf{A} be $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$. We compute

$$\|\mathbf{A}^\top \mathbf{p}\|^2 = \|\boldsymbol{\Sigma}\mathbf{U}^\top \mathbf{p}\|^2 = \sum_{i=1}^n \sigma_i^2(\mathbf{A}) \cdot (\mathbf{u}_i^\top \mathbf{p})^2 \geq \sigma_{\min}^2(\mathbf{A}) \cdot \|\mathbf{U}^\top \mathbf{p}\|^2 = \sigma_{\min}^2(\mathbf{A}) \cdot \|\mathbf{p}\|^2,$$

which completes the proof. □

Lemma B.2. For $c_1, c_2, c_3 > 0$, it holds

$$\frac{1}{2} \left(c_1 + c_2 c_3^2 + c_2 - \sqrt{(c_1 + c_2 c_3^2 + c_2)^2 - 4c_1 c_2} \right) = \min_{0 \leq t \leq 1} Q(t) := c_1 \cdot t^2 + c_2 \cdot (c_3 \cdot t - \sqrt{1-t^2})^2$$

Proof. Note that

$$Q(t) = (c_1 + c_2 c_3^2 - c_2) \cdot t^2 - 2c_2 c_3 \cdot t \sqrt{1-t^2} + c_2.$$

As mapping $t^2 \mapsto t$ is bijective in $[0, 1]$, we have

$$\min_{0 \leq t \leq 1} Q(t) = \min_{0 \leq t \leq 1} P(t) := (c_1 - c_2 + c_2 c_3^2) \cdot t - 2c_2 c_3 \cdot \sqrt{t} \cdot \sqrt{1-t} + c_2.$$

It is elementary to see $x \mapsto -\sqrt{x} \cdot \sqrt{1-x}$ is convex on $[0, 1]$. Thus, $P(t)$ is convex on $[0, 1]$. By the first-order optimality condition and the fact $\frac{1}{2} (c_4 - \sqrt{1+c_4^2}) < \frac{1}{2} (c_4 - |c_4|) = \min\{c_4, 0\}$, it holds for $c_4 \in \mathbb{R}$ that

$$\min_{0 \leq t \leq 1} c_4 \cdot t - \sqrt{t} \cdot \sqrt{1-t} = \frac{1}{2} \left(c_4 - \sqrt{1+c_4^2} \right) \quad \text{with} \quad t^* = \frac{c_4^2 - c_4 \sqrt{1+c_4^2} + 1}{2c_4^2 + 2} \in (0, 1).$$

The proof completes with simple algebraic manipulation. □

Proposition 3.4. Let $p = n\beta \cdot \|\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i\|^2$. For any $c > 0, \beta > 0$, we have

$$(c\mathbf{K} + \beta\mathbf{Z}^\top \mathbf{Z}) \succcurlyeq \mu \cdot \mathbf{I},$$

where $\mu > 0$ and defined by

$$\mu := \frac{1}{2} \left(c + n\beta + p - \sqrt{(c - n\beta)^2 + 2p(c + n\beta) + p^2} \right),$$

which reduces to $\mu = \min\{c, \beta n\}$ if the data points are centered and thus $p = 0$.

Proof. Let data points $\mathbf{X} \in \mathbb{R}^{n \times d}$ and labels $\mathbf{y} \in \{+1, -1\}^n$. We define $\widetilde{\mathbf{X}} := \mathbf{X} \odot \mathbf{y} \mathbf{1}_d^\top$, $\mathbf{v}_1 \in \mathbb{R}^d$, $v_2 \in \mathbb{R}$, and $\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ v_2 \end{bmatrix} \in \mathbb{R}^{d+1}$. Thus, $\mathbf{Z} = \begin{bmatrix} \widetilde{\mathbf{X}} & \mathbf{y} \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$. Then, for any $\mathbf{v} : \|\mathbf{v}\| = 1$, we compute

$$\begin{aligned}
 & \mathbf{v}^\top (c\mathbf{K} + \beta\mathbf{Z}^\top \mathbf{Z}) \mathbf{v} \\
 &= c \cdot \|\mathbf{v}_1\|^2 + \beta \cdot \|\mathbf{Z}\mathbf{v}\|^2 \\
 &= c \cdot \|\mathbf{v}_1\|^2 + \beta \cdot \|\widetilde{\mathbf{X}}\mathbf{v}_1 + v_2 \cdot \mathbf{y}\|^2 \\
 &= c \cdot \|\mathbf{v}_1\|^2 + \beta \cdot \left\| \left(\mathbf{I} - \frac{1}{n} \mathbf{y}\mathbf{y}^\top \right) \widetilde{\mathbf{X}}\mathbf{v}_1 + \frac{1}{n} \mathbf{y}\mathbf{y}^\top \widetilde{\mathbf{X}}\mathbf{v}_1 + v_2 \cdot \mathbf{y} \right\|^2 \\
 &\geq c \cdot \|\mathbf{v}_1\|^2 + \beta \cdot \left\| \frac{1}{n} \mathbf{y}\mathbf{y}^\top \widetilde{\mathbf{X}}\mathbf{v}_1 + v_2 \cdot \mathbf{y} \right\|^2 \quad (\text{Pythagorean}) \\
 &= c \cdot \|\mathbf{v}_1\|^2 + n\beta \cdot \left(\left(\frac{1}{n} \mathbf{1}_n^\top \mathbf{X} \right) \mathbf{v}_1 + v_2 \right)^2. \quad (\natural)
 \end{aligned}$$

Let $\bar{\mathbf{x}} := \frac{1}{n} \mathbf{1}_n^\top \mathbf{X}$, $\theta := \frac{\bar{\mathbf{x}}^\top \mathbf{v}_1}{\|\mathbf{v}_1\|}$, $\|\mathbf{v}_1\| = t$ with $0 \leq t \leq 1$. Without loss of generality, as $\bar{\mathbf{x}}$ is arbitrary, we assume $\theta \geq 0$ and $v_2 = -\sqrt{1-t^2} \leq 0$ (otherwise, the value of Equation (\natural) only increases). We continue with

$$\text{Equation } (\natural) \geq \min_{0 \leq t \leq 1} c \cdot t^2 + n\beta \cdot (\theta \cdot t - \sqrt{1-t^2})^2.$$

By Lemma B.2, we have

$$\text{Equation } (\natural) \geq \frac{1}{2} \left(c + n\beta + n\beta\theta^2 - \sqrt{(c - n\beta)^2 + n\beta\theta^2(2c + 2n\beta + n\beta\theta^2)} \right).$$

Note that $n\beta\theta^2 \leq n\beta\|\bar{\mathbf{x}}\|^2 =: p$ and the right-hand side of above lower bound is monotone decreasing with respect to θ^2 . We have Equation $(\natural) \geq \mu$. To see $\mu > 0$, note the following identity

$$\mu = \frac{1}{2} \left((c + n\beta + p) - \sqrt{(c + n\beta + p)^2 - 4cn\beta} \right) > 0,$$

which completes the proof. \square

Lemma 5.2 (sufficient decrease). *For $\gamma > 0, c > 0, \beta > 1 + \max \left\{ \frac{2(1 + \|\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i\|^2)}{c\sigma^2}, \frac{2}{n\sigma^2}, \frac{4}{cn\sigma^4} \right\}$, and $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1, there exists a constant $\tau_1 > 0$ only depending on $\{\sigma, c, n, \gamma, \beta\}$, such that*

$$\begin{aligned}
 & L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \\
 &\leq -\tau_1 \|(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - (\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\|^2.
 \end{aligned}$$

Proof. Note that

$$\begin{aligned}
 L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) &\leq \underbrace{L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^k)}_{T_1} + \\
 &\quad \underbrace{L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^k) - L_\beta(\mathbf{w}^k, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^k)}_{T_2} + \\
 &\quad \underbrace{L_\beta(\mathbf{w}^k, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^k) - L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)}_{T_3}.
 \end{aligned}$$

For T_1 , a direct computation and Equation (A.1c) gives

$$T_1 = \langle \boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k, \mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1} \rangle = \frac{1}{\beta} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|^2.$$

Using Equation (A.1b) and Equation (A.1c), we have

$$\mathbf{Z}^\top \boldsymbol{\lambda}^{k+1} = \mathbf{Z}^\top (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k) + \mathbf{Z}^\top \boldsymbol{\lambda}^k = \beta \mathbf{Z}^\top (\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1}) + \mathbf{Z}^\top \boldsymbol{\lambda}^k = -\mathbf{K}\mathbf{w}^{k+1}.$$

As \mathbf{Z} is surjective, it holds $\sigma := \sigma_{\min}(\mathbf{Z}) > 0$. Combining above computation and Lemma B.1, we have

$$T_1 = \frac{1}{\beta} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|^2 \leq \frac{1}{\sigma^2 \beta} \|\mathbf{Z}^\top (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)\|^2 = \frac{1}{\sigma^2 \beta} \|\mathbf{K}(\mathbf{w}^{k+1} - \mathbf{w}^k)\|^2 \leq \frac{1}{\sigma^2 \beta} \|\mathbf{w}^{k+1} - \mathbf{w}^k\|^2.$$

For T_2 , we first note the strong convexity of the \mathbf{w} -subproblem, which is due to the positive definite Hessian matrix $c\mathbf{K} + \beta\mathbf{Z}^\top \mathbf{Z} \succcurlyeq \mu \cdot \mathbf{I}$ in Proposition 3.4. Then, by [Nesterov, 2018, Theorem 2.1.8], we have

$$T_2 \leq -\mu \|\mathbf{w}^{k+1} - \mathbf{w}^k\|^2 \leq -\frac{\mu}{2} \|\mathbf{w}^{k+1} - \mathbf{w}^k\|^2 - \frac{\sigma^2 \mu}{2} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|^2,$$

where the last inequality is due to the primal-dual control:

$$\sigma \cdot \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\| \leq \|\mathbf{Z}^\top (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)\| = \|\mathbf{K}(\mathbf{w}^{k+1} - \mathbf{w}^k)\| \leq \|\mathbf{w}^{k+1} - \mathbf{w}^k\|.$$

For T_3 , by optimality condition Equation (A.1a), it holds

$$T_3 \leq -\frac{\gamma}{2} \|\mathbf{q}^{k+1} - \mathbf{q}^k\|^2.$$

In summary, we bound

$$T_1 + T_2 + T_3 \leq \left(\frac{1}{\sigma^2 \beta} - \frac{\mu}{2} \right) \|\mathbf{w}^{k+1} - \mathbf{w}^k\|^2 - \frac{\sigma^2 \mu}{2} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|^2 - \frac{\gamma}{2} \|\mathbf{q}^{k+1} - \mathbf{q}^k\|^2.$$

Finally, we choose parameters as

$$\beta > \frac{2}{\sigma^2 \mu}, \quad \gamma > 0, \quad \text{and } \tau_1 = \min \left\{ \left(\frac{1}{\sigma^2 \beta} - \frac{\mu}{2} \right), \frac{\sigma^2 \mu}{2}, \frac{\gamma}{2} \right\} > 0.$$

To get an explicit lower estimation for β , we plug in μ in Proposition 3.4. After algebraic simplification on condition $\beta > \frac{2}{\sigma^2 \mu}$, we have

$$\beta^2 (\theta^2 n \sigma^2 + n \sigma^2) + \beta c \sigma^2 - 4 > \beta \sigma^2 \sqrt{(c + \beta \theta^2 n + \beta n)^2 - 4 \beta c n}.$$

Then, it suffices to ensure $\beta^2 (\theta^2 n \sigma^2 + n \sigma^2) + \beta c \sigma^2 - 4 > 0$ and

$$\left(\beta^2 (\theta^2 n \sigma^2 + n \sigma^2) + \beta c \sigma^2 - 4 \right)^2 - \left(\beta \sigma^2 \sqrt{(c + \beta \theta^2 n + \beta n)^2 - 4 \beta c n} \right)^2 > 0.$$

Solving the quadratic inequality and simplifying the cubic one, we get $\beta > \max\{\text{Rt}_3(\#), \nu\}$, where

$$\nu := \frac{1}{2} \sqrt{\frac{c^2}{n^2 (\theta^2 + 1)^2} + \frac{16}{n (\theta^2 + 1) \sigma^2}} - \frac{c}{2n (\theta^2 + 1)},$$

and $\text{Rt}_3(\#)$ is the largest positive root of the following cubic equation with respect to t :

$$c n \sigma^4 t^3 - 2 n \sigma^2 \left(1 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right\|^2 \right) t^2 - 2 \sigma^2 c t + 4 = 0. \quad (\#)$$

Then, we use the Cauchy's bound Cauchy [1828] (see also [Hertz et al., 1999, Theorem 2.1]) for real zeros of polynomials to get

$$\text{Rt}_3(\#) \leq 1 + \max \left\{ \frac{2 \left(1 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right\|^2 \right)}{c \sigma^2}, \frac{2}{n \sigma^2}, \frac{4}{c n \sigma^4} \right\}.$$

Finally, note that $\nu \leq \text{Rt}_3(\#)$ and the proof completes. \square

Lemma 5.1 (bounded and proper). *For $\beta > \frac{2}{c\sigma^2}, \gamma > 0, c > 0$, and $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1, there exists an $R > 0$ such that $\|(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\| \leq R$ and $L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) > -\infty$ for any $k \in \mathbb{N}$.*

Proof. By Lemma 5.2, we have for any $k \in \mathbb{N}$:

$$\begin{aligned} L_\beta(\mathbf{w}^0, \mathbf{q}^0, \boldsymbol{\lambda}^0) &\geq L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \\ &= \frac{c}{2} \|\mathbf{w}^k\|_{\mathbf{K}}^2 + g(\mathbf{q}^k) + \frac{\beta}{2} \|\mathbf{Z}\mathbf{w}^k - \mathbf{q}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k\|^2 - \frac{1}{2\beta} \|\boldsymbol{\lambda}^k\|^2 \\ &\geq \frac{c}{2} \|\mathbf{w}^k\|_{\mathbf{K}}^2 + g(\mathbf{q}^k) + \frac{\beta}{2} \|\mathbf{Z}\mathbf{w}^k - \mathbf{q}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k\|^2 - \frac{1}{2\beta\sigma^2} \|\mathbf{K}\mathbf{w}^k\|^2 \\ &\geq \frac{c}{4} \|\mathbf{w}^k\|_{\mathbf{K}}^2 + g(\mathbf{q}^k) + \frac{\beta}{2} \|\mathbf{Z}\mathbf{w}^k - \mathbf{q}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k\|^2 \geq 0, \end{aligned}$$

where the second inequality is due to $\sigma \cdot \|\boldsymbol{\lambda}^k\| \leq \|\mathbf{Z}^\top \boldsymbol{\lambda}^k\| = \|\mathbf{K}\mathbf{w}^k\|$, the final inequality is by facts that \mathbf{K} is idempotent and $\left(\frac{c}{2} - \frac{1}{2\beta\sigma^2}\right) > \frac{c}{4}$. Therefore, the augmented Lagrangian is lower bounded by 0 for any $k \in \mathbb{N}$ and the properness of $L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)$ for any $k \in \mathbb{N}$ directly follows.

Then, similar to the proof of Proposition 3.4, we can prove $\frac{c}{4}\mathbf{K} + \frac{\beta}{2}\mathbf{Z}^\top\mathbf{Z} \succ 0$. Thus the level-set of L_β with respect to \mathbf{w} is bounded and we have $\{\mathbf{w}^k\}_k$ is bounded. Again, by Equations (A.1b) and (A.1c) and surjective \mathbf{Z} , we have $\|\mathbf{w}^k\| \geq \|\mathbf{K}\mathbf{w}^k\| = \|\mathbf{Z}^\top \boldsymbol{\lambda}^k\| \geq \sigma \|\boldsymbol{\lambda}^k\|$, and $\{\boldsymbol{\lambda}^k\}_k$ is bounded. By Equation (A.1c), $\|\mathbf{q}^k\| = \|\mathbf{Z}\mathbf{w}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k - \frac{1}{\beta} \boldsymbol{\lambda}^{k+1}\| \leq \|\mathbf{Z}\| \|\mathbf{w}^k\| + \frac{1}{\beta} \|\boldsymbol{\lambda}^k\| + \frac{1}{\beta} \|\boldsymbol{\lambda}^{k+1}\|$, and $\{\mathbf{q}^k\}_k$ is bounded. \square

Lemma 5.3 (limiting safeguard). *There exists a $\tau_2 > 0$ such that*

$$\begin{aligned} &\text{dist}\left(0, \partial L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1})\right) \\ &\leq \tau_2 \|(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) - (\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\|. \end{aligned}$$

Proof. To prove the safeguard property, we need to characterize the limiting subdifferential ∂L_β . Due to the irregularity of $g(\mathbf{q})$, a rigorous argument needs to apply a series of subdifferential calculus rules in a carefully chosen order. To this end, by grouping smooth terms in L_β together, by [Rockafellar and Wets, 2009, Exercise 8.8(c)], calmness of L_β from local Lipschitz, [Rockafellar and Wets, 2009, Proposition 8.32], and then with [Rockafellar and Wets, 2009, Proposition 10.5], we have

$$\partial L_\beta(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) = \begin{bmatrix} c\mathbf{K}\mathbf{w} + \mathbf{Z}^\top \boldsymbol{\lambda} + \beta\mathbf{Z}^\top(\mathbf{Z}\mathbf{w} - \mathbf{q}) \\ -\boldsymbol{\lambda} + \beta(\mathbf{q} - \mathbf{Z}\mathbf{w}) + \partial g(\mathbf{q}) \\ \mathbf{Z}\mathbf{w} - \mathbf{q} \end{bmatrix}.$$

By Cauchy-Schwarz, it holds

$$\begin{aligned} \text{dist}\left(0, \partial L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1})\right) &\leq \|c\mathbf{K}\mathbf{w}^{k+1} + \mathbf{Z}^\top \boldsymbol{\lambda}^{k+1} + \beta\mathbf{Z}^\top(\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1})\| + \\ &\quad \text{dist}\left(0, -\boldsymbol{\lambda}^{k+1} + \beta(\mathbf{q}^{k+1} - \mathbf{Z}\mathbf{w}^{k+1}) + \partial g(\mathbf{q}^{k+1})\right) + \|\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1}\|. \end{aligned}$$

By Equation (A.1b), we have

$$\|c\mathbf{K}\mathbf{w}^{k+1} + \mathbf{Z}^\top \boldsymbol{\lambda}^{k+1} + \beta\mathbf{Z}^\top(\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1})\| = \|\mathbf{Z}^\top(\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)\| \leq \|\mathbf{Z}\| \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|.$$

Meanwhile, by Equation (A.1a) and triangle inequality, we compute

$$\text{dist}\left(0, -\boldsymbol{\lambda}^{k+1} + \beta(\mathbf{q}^{k+1} - \mathbf{Z}\mathbf{w}^{k+1}) + \partial g(\mathbf{q}^{k+1})\right) \leq \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\| + \gamma \|\mathbf{q}^{k+1} - \mathbf{q}^k\| + \beta \|\mathbf{Z}\| \|\mathbf{w}^{k+1} - \mathbf{w}^k\|.$$

Finally, using Equation (A.1c), summing up, we get

$$\text{dist}\left(0, \partial L_\beta(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1})\right) \leq \beta \|\mathbf{Z}\| \|\mathbf{w}^{k+1} - \mathbf{w}^k\| + \gamma \|\mathbf{q}^{k+1} - \mathbf{q}^k\| + \left(\frac{\beta+1}{\beta} + \|\mathbf{Z}\|\right) \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|,$$

which combining with Cauchy-Schwarz and the choosing

$$\tau_2 := \sqrt{3} \max \left\{ \beta \|\mathbf{Z}\|, \gamma, \left(\frac{\beta+1}{\beta} + \|\mathbf{Z}\| \right) \right\},$$

completes the proof. \square

Lemma 5.4 (limiting stationarity convergence). *Let $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ be the sequence generated by Algorithm 1. As $k \rightarrow +\infty$, we have $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \rightarrow (\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ with \mathbf{w}^* is limiting-stationary, that is,*

$$0 \in c\mathbf{K}\mathbf{w}^* + \mathbf{Z}^\top \partial g(\mathbf{Z}\mathbf{w}^*).$$

Proof. The proof is by verifying the conditions in [Attouch et al., 2013, Theorem 2.9]. For H1, it directly holds from Lemma 5.2. H2 is from Lemma 5.3. H3 is due to Lemma 5.1 and continuity of L_β . We still need to show L_β is a KL function. To see this, we simply notice that $h_1(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) := \frac{c}{2} \|\mathbf{w}\|_{\mathbf{K}}^2 + \frac{\beta}{2} \|\mathbf{Z}\mathbf{w} - \mathbf{q} + \frac{1}{\beta} \boldsymbol{\lambda}\|^2 - \frac{1}{2\beta} \|\boldsymbol{\lambda}\|^2$ is a polynomial function and $h_2(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) := g(\mathbf{q})$ has a piecewise linear graph. Hence the sum $L_\beta(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) = h_1(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) + h_2(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda})$ is semi-algebraic and thus KL Kurdyka [1998]. The proof completes by using [Attouch et al., 2013, Theorem 2.9]. \square

Lemma 3.3 (proximal operator). *For $0 < \lambda < 2\rho^2$, it holds*

$$\text{Prox}_{\lambda\phi_\rho}(v) = \begin{cases} \{v\} & \text{for } v < -\frac{\lambda}{2\rho}, \\ \{v, v + \frac{\lambda}{\rho}\} & \text{for } v = -\frac{\lambda}{2\rho}, \\ \{v + \frac{\lambda}{\rho}\} & \text{for } -\frac{\lambda}{2\rho} < v \leq \rho - \frac{\lambda}{\rho}, \\ \{\rho\} & \text{for } \rho - \frac{\lambda}{\rho} < v \leq \rho, \\ \{v\} & \text{for } v > \rho. \end{cases}$$

For $0 < 2\rho^2 \leq \lambda$, we have

$$\text{Prox}_{\lambda\phi_\rho}(v) = \begin{cases} \{v\} & \text{for } v < \rho - \sqrt{2\lambda}, \\ \{v, \rho\} & \text{for } v = \rho - \sqrt{2\lambda}, \\ \{\rho\} & \text{for } \rho - \sqrt{2\lambda} < v < \rho, \\ \{v\} & \text{for } v \geq \rho. \end{cases}$$

Proof. As the computation of the proximal operator of the margin-loss Φ_ρ is quite tedious and complicated, we will deduce its explicit expression by calling on existing result Wang et al. [2021]. By definition, we have

$$\text{Prox}_{\lambda\phi_\rho}(v) = \text{Arg min}_x \left\{ \frac{1}{2\lambda} (x - v)^2 + \phi_\rho(x) \right\}.$$

Let $y = 1 - \frac{x}{\rho}$ and using identity $\max(a, \min(b, c)) = \min(b, \max(a, c))$, we have

$$\text{Prox}_{\lambda\phi_\rho}(v) = \rho - \rho \cdot \text{Arg min}_y \left\{ \frac{\rho^2}{2\lambda} \left(y + \frac{v}{\rho} - 1 \right)^2 + \max(0, \min(1, y)) \right\} = \rho - \rho \cdot \text{Prox}_{\frac{\lambda}{\rho^2} \ell_r} \left(1 - \frac{v}{\rho} \right),$$

where ℓ_r is defined in [Wang et al., 2021, Equation (2)]. Then, the claim follows from [Wang et al., 2021, Proposition 1], [Wang et al., 2021, Proposition 2]. \square

Proposition B.3. *g in Equation (MSVM) is prox-bounded for all $\lambda_g > 0$.*

Proof. Trivial as g is lower-bounded by 0 globally. \square

Proposition 3.5. *A point \mathbf{w} is local minimal for Problem (3.1) if and only if \mathbf{w} is d-stationary.*

Proof. The ‘‘only if’’ part is directly from [Rockafellar and Wets, 2009, Theorem 10.1]. For the ‘‘if’’ part, let \mathbf{w}^* be a d-stationary point of F with $F'(\mathbf{w}^*; \mathbf{w} - \mathbf{w}^*) \geq 0, \forall \mathbf{w} \in \mathbb{R}^{d+1}$ (see Definition 2.7(b)). As $g(\mathbf{Z}\mathbf{w})$ is piecewise

affine with respect to \mathbf{w} , by the min-max representation and with [Facchinei and Pang, 2007, Equation (4.2.7)], there exists $\delta > 0$ such that for all $\|\mathbf{w} - \mathbf{w}^*\| \leq \delta$, we have

$$g(\mathbf{Z}\mathbf{w}) = g(\mathbf{Z}\mathbf{w}^*) + g'(\mathbf{Z}\mathbf{w}^*; \mathbf{Z}\mathbf{w} - \mathbf{Z}\mathbf{w}^*).$$

Then, we compute

$$\begin{aligned} F(\mathbf{w}) &= \frac{c}{2} \|\mathbf{w}^* + (\mathbf{w} - \mathbf{w}^*)\|_{\mathbf{K}}^2 + \sum_{i=1}^n \phi_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}) \\ &\geq \frac{c}{2} \|\mathbf{w}^*\|_{\mathbf{K}}^2 + c(\mathbf{w} - \mathbf{w}^*)^{\top} \mathbf{K} \mathbf{w}^* + \sum_{i=1}^n \phi_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}^*) + \sum_{i=1}^n \phi_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}) - \phi_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}^*) \\ &= \frac{c}{2} \|\mathbf{w}^*\|_{\mathbf{K}}^2 + \sum_{i=1}^n \phi_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}^*) + c(\mathbf{w} - \mathbf{w}^*)^{\top} \mathbf{K} \mathbf{w}^* + \sum_{i=1}^n \phi'_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}^*; \mathbf{z}_i^{\top} \mathbf{w} - \mathbf{z}_i^{\top} \mathbf{w}^*) \\ &= \frac{c}{2} \|\mathbf{w}^*\|_{\mathbf{K}}^2 + \sum_{i=1}^n \phi_{\rho}(\mathbf{z}_i^{\top} \mathbf{w}^*) + F'(\mathbf{w}^*; \mathbf{w} - \mathbf{w}^*) \geq F(\mathbf{w}^*), \end{aligned}$$

which completes the proof. \square

C Proof of Theorem 3.6

We recall the definition of KL property in [Attouch et al., 2013, Definition 2.4] as follows:

Definition C.1 (KL property). *The function $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is said to have the Kurdyka–Łojasiewicz (KL) property at $\bar{\mathbf{x}} \in \text{dom } \partial f$ if there exists $\eta \in (0, +\infty]$, a neighborhood U of $\bar{\mathbf{x}}$ and a continuous concave function $\phi : [0, \eta) \rightarrow \mathbb{R}_+$ such that*

- $\phi(0) = 0$;
- ϕ is C^1 on $(0, \eta)$;
- for all $s \in (0, \eta)$, $\phi'(s) > 0$;
- for all \mathbf{x} in $U \cap [f(\bar{\mathbf{x}}) < f < f(\bar{\mathbf{x}}) + \eta]$, the KL inequality holds

$$\phi'(f(\mathbf{x}) - f(\bar{\mathbf{x}})) \cdot \text{dist}(0, \partial f(\mathbf{x})) \geq 1.$$

Now, we are ready to prove the Theorem 3.6.

Theorem 3.6. *Let $\gamma > 0, c > 0$, and*

$$\beta > 1 + \max \left\{ \frac{2 \left(1 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right\|^2 \right)}{c\sigma^2}, \frac{2}{n\sigma^2}, \frac{4}{cn\sigma^4} \right\}.$$

For the sequence $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1, the following holds:

- When $k \rightarrow \infty$, we have sequential convergence $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \rightarrow (\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ and \mathbf{w}^* is a d -stationary point of Problem (3.1).
- For any $0 \leq T < \infty$, we have

$$\min_{k \in [T]} \text{dist} \left(0, \partial L_{\beta}(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \right) \leq \frac{C_1}{\sqrt{T}},$$

where $C_1 := \tau_2 \cdot \sqrt{\frac{L_{\beta}(\mathbf{w}^0, \mathbf{q}^0, \boldsymbol{\lambda}^0) - F^*}{\tau_1}}$. See Section 5 for definitions of τ_1, τ_2 .

- Let \mathcal{S}_d be the set of d -stationary points of Problem (3.1). There exist $C_2 < \infty$, $\rho \in [0, 1)$ and finite \bar{k} such that for all $k \geq \bar{k}$:

$$\text{dist}(\mathbf{w}^k, \mathcal{S}_d) \leq C_2 \rho^k.$$

Proof. The proof is divided into several steps.

Step 1 (limiting stationary convergence). We first prove the sequence $\{(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k)\}_k$ generated by Algorithm 1 will converge to a limiting stationary point $(\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ of the augmented Lagrangian L_β , which implies \mathbf{w}^* is a limiting stationary point of Problem (3.1). The claim direct follows from Lemma 5.4.

The global rate is a direct corollary of Lemma 5.2 and Lemma 5.3. Specifically,

$$\begin{aligned} \min_{k \in [T]} \text{dist} \left(0, \partial L_\beta(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \right) &\leq \tau_2 \cdot \min_{k \in [T]} \left\| (\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) - (\mathbf{w}^{k-1}, \mathbf{q}^{k-1}, \boldsymbol{\lambda}^{k-1}) \right\| && \text{(Lemma 5.3)} \\ &\leq \tau_2 \cdot \sqrt{\frac{1}{T} \sum_{k=1}^T \left\| (\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) - (\mathbf{w}^{k-1}, \mathbf{q}^{k-1}, \boldsymbol{\lambda}^{k-1}) \right\|^2} \\ &\leq \tau_2 \cdot \sqrt{\frac{L_\beta(\mathbf{w}^0, \mathbf{q}^0, \boldsymbol{\lambda}^0) - L_\beta(\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)}{\tau_1 T}} = \frac{C_2}{\sqrt{T}}. && \text{(Lemma 5.2)} \end{aligned}$$

Step 2 (refine stationarity characterization). Then, we will refine the characterization of $(\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ by showing \mathbf{w}^* is indeed a d-stationary point. We first rewrite the optimality condition for \mathbf{q} -subproblem as

$$\mathbf{q}^{k+1} \in \text{Prox}_{\frac{1}{\beta+\gamma}g} \left(\frac{\beta}{\beta+\gamma} \left(\mathbf{Z}\mathbf{w}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k \right) + \frac{\gamma}{\beta+\gamma} \mathbf{q}^k \right).$$

Note that, by Equation (A.1c), we observe the following identify:

$$\frac{\beta}{\beta+\gamma} \left(\mathbf{Z}\mathbf{w}^k + \frac{1}{\beta} \boldsymbol{\lambda}^k \right) + \frac{\gamma}{\beta+\gamma} \mathbf{q}^k = \mathbf{q}^{k+1} + \frac{\gamma}{\beta+\gamma} (\mathbf{q}^k - \mathbf{q}^{k+1}) + \frac{\beta}{\beta+\gamma} \mathbf{Z}(\mathbf{w}^k - \mathbf{w}^{k+1}) + \frac{1}{\beta+\gamma} \boldsymbol{\lambda}^{k+1}.$$

Let $k \rightarrow +\infty$. By the convergence of $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \rightarrow (\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ in Step 1, global prox-boundedness of g from Proposition B.3, and outer semi-continuous of proximal mapping $\text{Prox}_{\frac{1}{\beta+\gamma}g}$ (see [Rockafellar and Wets, 2009, Example 5.23]), we have

$$\begin{aligned} \mathbf{q}^* &= \lim_{k \rightarrow \infty} \mathbf{q}^k \in \limsup_{k \rightarrow \infty} \text{Prox}_{\frac{1}{\beta+\gamma}g} \left(\mathbf{q}^{k+1} + \frac{\gamma}{\beta+\gamma} (\mathbf{q}^k - \mathbf{q}^{k+1}) + \frac{\beta}{\beta+\gamma} \mathbf{Z}(\mathbf{w}^k - \mathbf{w}^{k+1}) + \frac{1}{\beta+\gamma} \boldsymbol{\lambda}^{k+1} \right) \\ &\subseteq \text{Prox}_{\frac{1}{\beta+\gamma}g} \left(\mathbf{q}^* + \frac{1}{\beta+\gamma} \boldsymbol{\lambda}^* \right). \end{aligned}$$

By Fermat's rule [Rockafellar and Wets, 2009, Theorem 10.1] on proximal mapping $\text{Prox}_{\frac{1}{\beta+\gamma}g}$, it holds $\boldsymbol{\lambda}^* \in \widehat{\partial}g(\mathbf{q}^*)$. Using Equations (A.1b) and (A.1c) and limiting argument, we also have $\mathbf{Z}\mathbf{w}^* = \mathbf{q}^*, 0 = \mathbf{K}\mathbf{w}^* + \mathbf{Z}^\top \boldsymbol{\lambda}^*$, which indicates $0 \in \mathbf{K}\mathbf{w}^* + \mathbf{Z}^\top \widehat{\partial}g(\mathbf{Z}\mathbf{w}^*)$. As \mathbf{Z} is surjective, by chain rule [Rockafellar and Wets, 2009, Exercise 10.7] and sum rule [Rockafellar and Wets, 2009, Exercise 8.8(c)], we have $0 \in \widehat{\partial}F(\mathbf{w}^*)$. As F is locally Lipschitz, by [Rockafellar and Wets, 2009, Theorem 9.13], \mathbf{w}^* is a d-stationary point.

Step 3 (linear convergence). By inspecting the analytic expression of proximal operator in Lemma 3.3, it is easy to check that, for any $\lambda > 0, \rho > 0, q \in \mathbb{R}$ and $q \in \text{Prox}_{\lambda\phi_\rho}(v)$, where $v \in \mathbb{R}$ is arbitrary, we have

$$|q| \geq \min \left\{ \frac{\lambda}{2\rho}, \rho \right\}.$$

Then, for any $k \in \mathbb{N}$, there exists a constant τ_3 with $\min \left\{ \frac{1}{2\rho(\beta+\gamma)}, \rho \right\} > \tau_3 > 0$ such that

$$\min_{i \in [n]} |q_i^k| > \tau_3 > 0, \quad \text{and} \quad \min_{i \in [n]} |q_i^*| = \min_{i \in [n]} |\mathbf{z}_i^\top \mathbf{w}^*| > \tau_3 > 0.$$

By the convergence of $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) \rightarrow (\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*)$ in Lemma 5.4, there exists \bar{k} such that for all $k \geq \bar{k}$, we have

$$\max_{i \in [n]} |q_i^k - q_i^*| \leq \tau_3, \quad \text{implying} \quad \min_{i \in [n]} q_i^k \cdot q_i^* > \frac{\tau_3^2}{2} > 0, \quad \text{and} \quad \text{sgn}(q_i^k) = \text{sgn}(q_i^*), \quad \forall i \in [n].$$

Let the partition $[n] = \mathcal{I}_<^k \sqcup \mathcal{I}_=^k \sqcup \mathcal{I}_>^k$, where $k \in \mathbb{N} \cup \{*\}$ and

$$\mathcal{I}_<^k := \{i \in [n] : q_i^k < 0\}, \quad \mathcal{I}_=^k := \{i \in [n] : q_i^k = 0\}, \quad \mathcal{I}_>^k := \{i \in [n] : q_i^k > 0\}.$$

Then, we define the following ghost convex problem:

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \widehat{F}(\mathbf{w}) := \frac{c}{2} \|\mathbf{w}\|_{\mathbf{K}}^2 + \sum_{i \in \mathcal{I}_>^k} \widehat{\phi}_\rho(\mathbf{z}_i^\top \mathbf{w}), \quad \text{where} \quad \widehat{\phi}_\rho(u) = \max\left(0, 1 - \frac{u}{\rho}\right). \quad (\text{CPLQ})$$

Let $\hat{g}(\mathbf{q}) = \sum_{i \in \mathcal{I}_>^k} \widehat{\phi}_\rho(\mathbf{z}_i^\top \mathbf{w})$. Using the same operator splitting technique $\mathbf{Z}\mathbf{w} = \mathbf{q} \in \mathbb{R}^n$, the augmented Lagrangian can be written as

$$\widehat{L}_\beta(\mathbf{w}, \mathbf{q}, \boldsymbol{\lambda}) := \widehat{F}(\mathbf{w}) + \langle \boldsymbol{\lambda}, \mathbf{Z}\mathbf{w} - \mathbf{q} \rangle + \frac{\beta}{2} \|\mathbf{Z}\mathbf{w} - \mathbf{q}\|^2.$$

Setting $\hat{\mathbf{w}}^0 = \mathbf{w}^{\bar{k}}, \hat{\mathbf{q}}^0 = \mathbf{q}^{\bar{k}}, \hat{\boldsymbol{\lambda}}^0 = \boldsymbol{\lambda}^{\bar{k}}, (\beta, \gamma) = (\beta, \gamma)$ in Algorithm 1, we have the following ghost convex semi-proximal ADMM algorithm:

Algorithm 2 Ghost Semi-proximal ADMM for Problem (CPLQ)

- 1: Set $\hat{\mathbf{w}}^0 = \mathbf{w}^{\bar{k}}, \hat{\mathbf{q}}^0 = \mathbf{q}^{\bar{k}}, \hat{\boldsymbol{\lambda}}^0 = \boldsymbol{\lambda}^{\bar{k}}, (\beta, \gamma) = (\beta, \gamma)$ in Algorithm 1.
- 2: **for** $t \in \{0, 1, 2, \dots\}$ **do**
- 3:

$$\begin{aligned} \hat{\mathbf{q}}^{t+1} &\in \arg \min_{\hat{\mathbf{q}}} \hat{g}(\hat{\mathbf{q}}) + \frac{\beta}{2} \|\hat{\mathbf{q}} - \mathbf{Z}\hat{\mathbf{w}}^t - \frac{1}{\beta} \hat{\boldsymbol{\lambda}}^t\|^2 + \frac{\gamma}{2} \|\hat{\mathbf{q}} - \hat{\mathbf{q}}^t\|^2 \\ \hat{\mathbf{w}}^{t+1} &= \arg \min_{\hat{\mathbf{w}}} \hat{f}(\hat{\mathbf{w}}) + \frac{\beta}{2} \|\mathbf{Z}\hat{\mathbf{w}} - \hat{\mathbf{q}}^{t+1} + \frac{1}{\beta} \hat{\boldsymbol{\lambda}}^t\|^2 \\ \hat{\boldsymbol{\lambda}}^{t+1} &= \hat{\boldsymbol{\lambda}}^t + \beta(\mathbf{Z}\hat{\mathbf{w}}^{t+1} - \hat{\mathbf{q}}^{t+1}). \end{aligned}$$

- 4: **end for**
-

We denote the sequence generated by Algorithm 2 as $\{(\hat{\mathbf{w}}^t, \hat{\mathbf{q}}^t, \hat{\boldsymbol{\lambda}}^t)\}_t$. As \hat{g} is a convex piecewise linear-quadratic function, by [Rockafellar and Wets, 2009, Theorem 12.30], the KKT mapping of \widehat{F} is piecewise polyhedral. Then, by the celebrated results of Robinson [1981], we have the inverse KKT mapping is calm (see also [Dontchev and Rockafellar, 2009, Proposition 3H.1]). Besides, by the strong convexity of \mathbf{w} -subproblem in the proof of Lemma 5.2, we have $c\mathbf{K} + \beta\mathbf{Z}^\top \mathbf{Z} \succ \mathbf{0}$. Using [Han et al., 2018, Theorem 2], there exists $\rho \in (0, 1)$ and constant $\tau_4 \geq 0$ such that

$$\text{dist}\left((\hat{\mathbf{w}}^k, \hat{\mathbf{q}}^k, \hat{\boldsymbol{\lambda}}^k), \mathcal{S}_{\text{KKT}}\right) \leq \tau_4 \cdot \rho^k,$$

where \mathcal{S}_{KKT} is the set KKT points for Equation (CPLQ).

Then, we will show for any $k \in (\mathbb{N} \setminus [\bar{k} - 1]) \cup \{*\}$, $(\mathbf{w}^k, \mathbf{q}^k, \boldsymbol{\lambda}^k) = (\hat{\mathbf{w}}^{k-\bar{k}}, \hat{\mathbf{q}}^{k-\bar{k}}, \hat{\boldsymbol{\lambda}}^{k-\bar{k}})$, which directly indicates the linear convergence. The proof is by induction on k . For $k = \bar{k}$, it is by definition the claim holds. Then, we assume for $\{k, \dots, k\}$ the claim holds. Then we need to show $(\mathbf{w}^{k+1}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^{k+1}) = (\hat{\mathbf{w}}^{k-\bar{k}+1}, \hat{\mathbf{q}}^{k-\bar{k}+1}, \hat{\boldsymbol{\lambda}}^{k-\bar{k}+1})$. For $\mathbf{q}^{k+1} = \hat{\mathbf{q}}^{k-\bar{k}+1}$, we consider the optimality condition Equation (A.1a) of the \mathbf{q} -subproblem in Algorithm 1. By $\min_{i \in [n]} |q_i^{k+1}| > \tau_3$, [Rockafellar and Wets, 2009, Proposition 10.5], [Rockafellar and Wets, 2009, Theorem 10.49], we have

$$\widehat{\partial}g(\mathbf{q}^{k+1}) = \bigoplus_{i=1}^n \partial\phi_\rho(q_i^{k+1}) = \bigoplus_{i=1}^n \delta_{\mathcal{I}_>^{k+1}}(i) \cdot \partial\widehat{\phi}_\rho(q_i^{k+1}) = \partial\hat{g}(\mathbf{q}^{k+1}),$$

where $\delta_C(i) := 1$ if $i \in C$ otherwise $\delta_C(i) := 0$, and we use $\mathcal{I}_>^k = \mathcal{I}_>^*$ for all $k \geq \bar{k}$. Then, it is clear that \mathbf{q}^{k+1} satisfies the optimality condition of the \mathbf{q} -subproblem of Algorithm 2, by whose strong convexity, we have $\mathbf{q}^{k+1} = \hat{\mathbf{q}}^{k+1}$. Similarly, we can show by examining the optimality condition that $\mathbf{w}^{k+1} = \hat{\mathbf{w}}^{k+1}$ and $\boldsymbol{\lambda}^{k+1} = \hat{\boldsymbol{\lambda}}^{k+1}$. Taking limit we have $(\mathbf{w}^*, \mathbf{q}^*, \boldsymbol{\lambda}^*) = (\hat{\mathbf{w}}^*, \hat{\mathbf{q}}^*, \hat{\boldsymbol{\lambda}}^*)$. The final building block is to show if $(\hat{\mathbf{w}}^*, \hat{\mathbf{q}}^*, \hat{\boldsymbol{\lambda}}^*) \in \mathcal{S}_{\text{KKT}}$ then

$\hat{\mathbf{w}}^* \in \mathcal{S}_d$, which indicates $\text{dist}(\hat{\mathbf{w}}^k, \mathcal{S}_d) \leq \text{dist}(\hat{\mathbf{w}}^k, \Pi_1(\mathcal{S}_{\text{KKT}})) \leq \text{dist}((\hat{\mathbf{w}}^k, \hat{\mathbf{q}}^k, \hat{\boldsymbol{\lambda}}^k), \mathcal{S}_{\text{KKT}})$. To this end, for any KKT point $(\hat{\mathbf{w}}^*, \hat{\mathbf{q}}^*, \hat{\boldsymbol{\lambda}}^*)$, we write down the KKT system of Problem (CPLQ) as

$$\begin{cases} 0 &= c\mathbf{K}\hat{\mathbf{w}}^* + \mathbf{Z}^\top \hat{\boldsymbol{\lambda}}^*, \\ 0 &\in -\hat{\boldsymbol{\lambda}}^* + \partial \hat{g}(\hat{\mathbf{q}}^*), \\ 0 &= \mathbf{Z}\hat{\mathbf{w}}^* - \hat{\mathbf{q}}^*. \end{cases}$$

Then we have $0 \in c\mathbf{K}\hat{\mathbf{w}}^* + \mathbf{Z}^\top \partial \hat{g}(\mathbf{Z}\hat{\mathbf{w}}^*)$, where \hat{g} is convex. Therefore, the claim follows from $\min_{i \in [n]} |q_i^*| > \tau_3$, [Rockafellar and Wets, 2009, Proposition 10.5], [Rockafellar and Wets, 2009, Theorem 10.49], which implies $\partial \hat{g}(\mathbf{Z}\hat{\mathbf{w}}^*) = \partial \hat{g}(\hat{\mathbf{q}}^*) = \partial \hat{g}(\hat{\mathbf{q}}^*) = \partial \hat{g}(\mathbf{Z}\hat{\mathbf{w}}^*)$. Thus we have $\hat{\mathbf{w}}^* \in \mathcal{S}_d$, which completes the proof. \square

D More Example

Counterexample. Let $d = 1, n = 6, c = 1, \rho = 1$. We consider the following data points multiset:

$$\mathbb{S} := \{(x_i, y_i) : i \in [n]\} = \{(1, +1), (1, +1), (-1, -1), (-1, -1), (0, +1), (0, -1)\} \subseteq \mathbb{R}^d \times \{+1, -1\}.$$

We claim that $\mathbf{w} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ is a global minimizer for Problem (3.1) on $\{(x_i, y_i) : i \in [n]\}$. However, \mathbf{w} does not satisfy the condition in [Suzumura et al., 2017, Theorem 4]. For the proof, we reorganized the data points into

$$\mathbf{X} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & -1 \\ 1 & -1 \\ 0 & 1 \\ 0 & -1 \end{bmatrix}, \quad \text{with } \text{rank}(\mathbf{Z}) = d + 1 = 2 < n.$$

Let $\phi_1(t) = \min\{1, \max\{1 - t, 0\}\}$. Then, Problem (1.1) on \mathbb{S} reads

$$\nu := \min_{\theta \in \mathbb{R}, b \in \mathbb{R}} f(\theta, b) := \frac{\theta^2}{2} + 2\phi_1(\theta + b) + 2\phi_1(\theta - b) + \phi_1(b) + \phi_1(-b). \quad (\text{CE})$$

To get an optimizer of Problem (CE), we consider the following 5 cases:

- $b = 0$
 - $\theta = 0$: $\nu_1 := \inf_{\theta} f(\theta, b) = 6$.
 - $\theta \neq 0$: $\nu_2 := \inf_{\theta} f(\theta, b) = 2 + \left(\inf_{\theta} \frac{\theta^2}{2} + 4\phi_1(\theta)\right) = \frac{5}{2}$ when $\theta = 1$.
- $b \neq 0$
 - $|\theta| = |b|$: by symmetry of b , we assume $\theta = b$. Then, Problem (CE) can be rewritten as

$$\min_{b \neq 0} \frac{b^2}{2} + 2\phi_1(2b) + 2 + \phi_1(b) + \phi_1(-b).$$

By monotonicity of ϕ_1 , we can assume $b > 0$. Let the convex hinge loss be $h(t) = \max\{1 - t, 0\}$. Then, we get the following convex reformulation

$$\nu_3 := \frac{7}{2} = \inf_{b > 0} \frac{b^2}{2} + 2h(2b) + 2 + h(b) + 1.$$

- $|\theta| \neq |b|$: by monotonicity of g and symmetry of b , we assume $w \geq 0$ and $b > 0$.
 - ★ $|\theta| > |b|$:

$$\nu_4 := \frac{5}{2} = \inf_{\theta \geq 0, b > 0} \frac{\theta^2}{2} + 2h(\theta + b) + 2h(\theta - b) + h(b) + 1.$$

★ $|\theta| < |b|$:

$$\nu_5 := 3 = \inf_{\theta \geq 0, b > 0} \frac{\theta^2}{2} + 2h(\theta + b) + 2 + h(b) + 1.$$

In summary, we claim $\theta^* = 1, b^* = 0$ is an optimal solution of Problem (CE) as $f(1, 0) = \frac{5}{2} = \min\{\nu_i\}_{[5]} \leq \nu$. In that case, let $\mathbf{q} := \mathbf{y} \mathbb{1}_d \odot (\mathbf{X}\theta^* - b^*\mathbf{y}) = [1 \ 1 \ 1 \ 1 \ 0 \ 0]^\top$. On the other hand, in [Suzumura et al., 2017, Theorem 4], the KKT-type condition requires $q_i \neq 0$ for all $i \in [n]$, which is not satisfied by a global minimizer $\theta^* = 1, b^* = 0$. Thus, [Suzumura et al., 2017, Theorem 4] is not necessary. Besides, we note that above example provide an explicit case that if \mathbf{Z} is not surjective, then the dual characterization in Proposition 3.2 may not hold. In the language of variational analysis, that is because $0 \in \widehat{\partial}(f(\mathbf{Z}\cdot))(\mathbf{w}^*)$ whereas $\widehat{\partial}(f)(\mathbf{Z}\mathbf{w}^*) = \emptyset$.

E DCA-SpADMM Hybrid Algorithm

The hybrid method (Algorithm 3) runs DCA until convergence and then starts SpADMM from that DC-critical point. The detailed hybrid scheme can be summarized as follows:

Algorithm 3 Hybrid Semi-proximal ADMM and DCA for Problem (3.1)

Input: $\mathbf{Z} \in \mathbb{R}^{(d+1) \times n}, \mathbf{w}^0$, choose $\gamma > 0, c > 0$ and $\beta > 1 + \max\left\{\frac{2(1 + \|\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i\|^2)}{c\sigma^2}, \frac{2}{n\sigma^2}, \frac{4}{cn\sigma^4}\right\}$.

1: Set $r = 0$.

2: **for all** $k \in \{0, 1, 2, \dots\}$ **do**

3: **if** $r = 0$ **then**

5: Set $\mathbf{g}^k = -\frac{1}{\rho} \mathbf{Z}^\top \text{sgn}\left(-\frac{1}{\rho} \mathbf{Z}^\top \mathbf{w}^k\right)$.

6: Compute

$$\mathbf{w}^{k+1} \in \text{Arg min}_w \frac{c}{2} \|\mathbf{w}\|_{\mathbf{K}}^2 + \sum_{i=1}^n \max\left(1 - \frac{\mathbf{z}_i^\top \mathbf{w}}{\rho}, 0\right) - \mathbf{w}^\top \mathbf{g}^k.$$

7: **if** $\|\mathbf{w}^{k+1} - \mathbf{w}^k\| \leq 10^{-5}$ **then**

8: Set $r = 1, \mathbf{q}^{k+1} = \mathbf{Z}\mathbf{w}^{k+1}, \boldsymbol{\lambda}^{k+1} = 0$.

9: **end if**

10: **else**

$$\mathbf{q}^{k+1} \in \text{Arg min}_q L_\beta(\mathbf{w}^k, \mathbf{q}, \boldsymbol{\lambda}^k) + \frac{\gamma}{2} \|\mathbf{q} - \mathbf{q}^k\|^2,$$

$$\mathbf{w}^{k+1} = \arg \min_w L_\beta(\mathbf{w}, \mathbf{q}^{k+1}, \boldsymbol{\lambda}^k),$$

$$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k + \beta(\mathbf{Z}\mathbf{w}^{k+1} - \mathbf{q}^{k+1}).$$

11: **end if**

12: **end for**

F More Numerical Results

F.1 Objective Value Curve

In this section, we provide extra numerical results in the setting of Section 6.

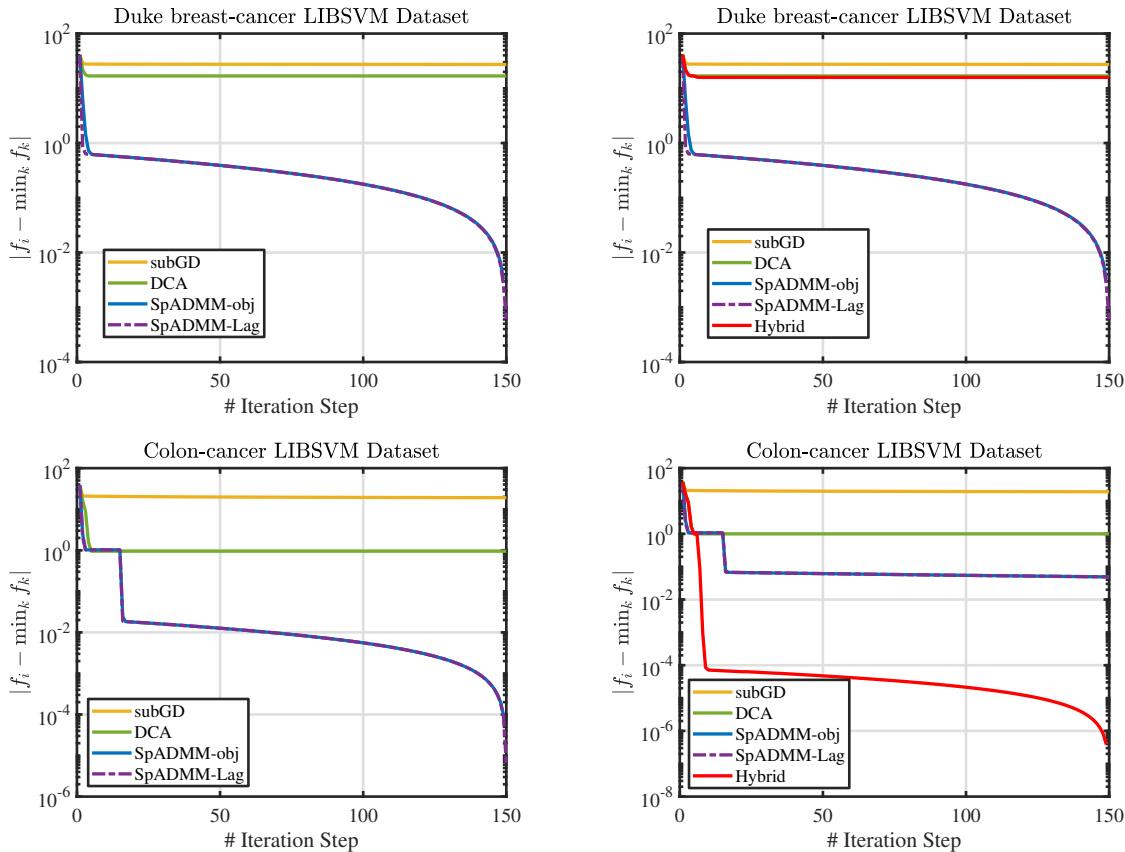


Figure 5: Performance of subGD (subgradient), DCA, SpADMM (Algorithm 1), and Hybrid Algorithm 3 on Real-world Data. The min in log-scaled y -axis is taken over all aforementioned algorithms.

F.2 Wall-Clock Time

To demonstrate the efficiency of various algorithms in the numerical experiments, we record the wall-clock time of running these algorithms for 10^2 iterative steps in Table 2.

Table 2: Wall-Clock Time of subGD (subgradient), DCA, SpADMM (Algorithm 1), and Hybrid Algorithm 3.

10^2 steps	Wall-Clock Time (s)			
	Toy	Leukemia	Duke	Colon
subGD	0.0173	24.0051	23.9875	1.7640
DCA	19.6547	60.6936	72.4713	42.5783
SpADMM	0.0138	1.1031	1.1090	0.0828
Hybrid	0.7206	2.9015	3.9245	1.9675