

A globally and linearly convergent PGM for zero-norm regularized quadratic optimization with sphere constraint

Yuqia Wu*, Shaohua Pan[†] and Shujun Bi[‡]

November 11, 2018

Abstract

This paper is concerned with the zero-norm regularized quadratic optimization with a sphere constraint, which has an important application in sparse eigenvalue problems. For this class of nonconvex and nonsmooth optimization problems, we establish the KL property of exponent 1/2 for its extended-valued objective function and develop a globally and linearly convergent proximal gradient method (PGM). Numerical experiments are included for sparse principal component analysis (PCA) with synthetic and real data to confirm the obtained theoretic results.

Keywords: KL property of exponent 1/2; zero-norm; sphere constraint; PGM

1 Introduction

Let \mathbb{S}^p denote the space of all $p \times p$ real symmetric matrices, equipped with the trace inner product and its induced Frobenius norm $\|\cdot\|_F$, and \mathbb{S}_+^p denote the set of all positive semidefinite matrices in \mathbb{S}^p . Given a $A \in \mathbb{S}^p$, we are interested in the following problem

$$\min_{x \in \mathcal{S}} \left\{ x^\top A x + \nu \|x\|_0 \right\} \quad (1)$$

where $\nu > 0$ is the regularization parameter, $\|x\|_0$ denotes the zero-norm (cardinality) of the vector x , and $\mathcal{S} := \{x \in \mathbb{R}^p \mid \|x\| = 1\}$ is the unit sphere in \mathbb{R}^p . With the indicator function of \mathcal{S} , the problem (1) can be compactly written as the following form

$$\min_{x \in \mathbb{R}^p} \left\{ \Phi_\nu(x) := x^\top A x + \nu(\|x\|_0 + \delta_{\mathcal{S}}(x)) \right\}. \quad (2)$$

It is well-known that the minimization of the structured nonconvex function $\|\cdot\|_0 + \delta_{\mathcal{S}}(\cdot)$ aims to capture a sparse unit vector. Such a nonconvex and nonsmooth problem has an

*School of Mathematics, South China University of Technology, Guangzhou.

[†]shhpan@scut.edu.cn. School of Mathematics, South China University of Technology, Guangzhou.

[‡]bishj@scut.edu.cn. School of Mathematics, South China University of Technology, Guangzhou.

important application in sparse eigenvalue problems such as the sparse PCA [9, 13, 28] or the densest subgraph finding problem (see [20, Section 4.3]). However, the nonconvexity of $\|\cdot\|_0 + \delta_S(\cdot)$ restricts the development of globally convergent methods for solving (2), although the components of its proximal mapping are available (see Section 4).

For the past several years, it has witnessed that the successful use of the Kurdyka-Łojasiewicz (KL) property in analyzing the global convergence of first-order algorithms for nonconvex and nonsmooth optimization problems (see, e.g., [3, 4, 6, 7]). In particular, the KL property of exponent 1/2 plays a key role in achieving the linear rate of convergence. As recently discussed in [16, 22, 17], for the structured semiconvex function and the primal lower nice function, their KL property of exponent 1/2 is usually weaker than the metric subregularity of their subdifferential operators [2] or the Luo-Tseng error bound [24], which are the common regularity for deriving the linear convergence of first-order algorithms (see, e.g., [11, 18, 25, 26]). Thus, a valuable research direction is to discover the class of functions that precisely possesses the KL property of exponent 1/2. We notice that some researchers have made some positive progress in this direction; for example, Liu et al. [19] established a restricted-type KL property of exponent 1/2 for the quadratic function over orthogonal constraints, and Zhang et al. [27] verified the KL property of exponent 1/2 for several classes of regularized matrix factorization functions.

The main contribution of this work is to establish the KL property of exponent 1/2 for the nonconvex composite function Φ_ν . Since Φ_ν involves two nonconvex nonsmooth functions, its (limiting) subdifferential characterization is not an easy task, let alone its KL property of exponent 1/2. We not only provide the subdifferential characterization of Φ_ν but also achieve this goal. Then, motivated by the finding that the components of the proximal mapping of $\|\cdot\|_0 + \delta_S(\cdot)$ are available, we develop a globally and linearly convergent proximal gradient method (PGM) for solving the problem (1). Numerical comparisons with the generalized power method [10] for sparse PCA on synthetic and real data confirm our theoretical results as well as validate the efficiency of the PGM.

2 Notations and preliminaries

Throughout this paper, we denote by \mathbb{X} a finite-dimensional vector space equipped with the inner product $\langle \cdot, \cdot \rangle$ and its induced norm $\|\cdot\|$, by \mathbb{R}^p the p -dimensional Euclidean space, and by \mathbb{O}^p the set of all $p \times p$ orthogonal matrices. For an extended-valued $f: \mathbb{X} \rightarrow (-\infty, +\infty]$, we say that f is proper if $\text{dom } f := \{x \in \mathbb{X} \mid f(x) < \infty\}$ is nonempty, and for any given $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}$, set $[\alpha \leq f \leq \beta] := \{x \in \mathbb{X} \mid \alpha \leq f(x) \leq \beta\}$. The notation $x' \xrightarrow[f]{} x$ to signify $x' \rightarrow x$ and $f(x') \rightarrow f(x)$. For a given $\bar{x} \in \mathbb{X}$ and $\delta > 0$, $\mathbb{B}(\bar{x}, \delta)$ denotes the closed ball centered at \bar{x} with radius δ ; and for a given set $\Omega \subseteq \mathbb{X}$, $\delta_\Omega(\cdot)$ means the indicator function of Ω . For a vector $z \in \mathbb{R}^p$, $|z|^\downarrow$ means the vector obtained by arranging the entries of $|z|$ in a decreasing order, and $|z|^{\kappa, \downarrow} \in \mathbb{R}^\kappa$ is the vector composed of the first κ entries of $|z|^\downarrow$. For $t \in \mathbb{R}$, write $t_- = \min(0, t)$ and $t_+ = \max(0, t)$. For a matrix $A \in \mathbb{S}^p$ and an index set $J \subseteq \{1, 2, \dots, p\}$, $A_{JJ} \in \mathbb{S}^{|J|}$ denotes the submatrix consisting of those entries A_{ij} with $i \in J$ and $j \in J$.

2.1 Generalized subdifferential

Definition 2.1 (see [15, Definition 8.3]) Consider a function $f: \mathbb{X} \rightarrow (-\infty, +\infty]$ and an arbitrary point $x \in \text{dom} f$. The regular subdifferential of f at x is defined as

$$\widehat{\partial}f(x) := \left\{ v \in \mathbb{X} \mid \liminf_{\substack{x' \rightarrow x \\ x' \neq x}} \frac{f(x') - f(x) - \langle v, x' - x \rangle}{\|x' - x\|} \geq 0 \right\};$$

the (limiting) subdifferential of f at x is defined as

$$\partial f(x) := \left\{ v \in \mathbb{X} \mid \exists x^k \xrightarrow{f} x \text{ and } v^k \in \widehat{\partial}f(x^k) \rightarrow v \text{ as } k \rightarrow \infty \right\};$$

and the horizon subdifferential of f at x is defined as

$$\partial^\infty f(x) := \left\{ v \in \mathbb{X} \mid \exists x^k \xrightarrow{f} x, \lambda^k \searrow 0 \text{ and } v^k \in \widehat{\partial}f(x^k) \text{ with } \lambda^k v^k \rightarrow v \text{ as } k \rightarrow \infty \right\}.$$

Remark 2.1 (i) At each $x \in \text{dom} f$, $\widehat{\partial}f(x)$ and $\partial f(x)$ are both closed with $\widehat{\partial}f(x) \subseteq \partial f(x)$, and the former is always convex but the latter is generally nonconvex. When f is convex, $\widehat{\partial}f(x) = \partial f(x)$ and is precisely the subdifferential of f at x in the sense of convex analysis.

(ii) Let $\{(x^k, v^k)\}_{k \in \mathbb{N}}$ be a sequence in $\text{gph} \partial f$ that converges to (x, v) as $k \rightarrow \infty$. By Definition 2.1, if $f(x^k) \rightarrow f(x)$ as $k \rightarrow \infty$, then $(x, v) \in \text{gph} \partial f$.

(iii) The point \bar{x} at which $0 \in \partial f(\bar{x})$ is called a (limiting) critical point of f . In the sequel, we denote by $\text{crit} f$ the set of critical points of f . By [15, Theorem 10.1], we know that a local minimizer of f is necessarily a critical point of f .

For the zero-norm function, we have the following conclusion for its subdifferential.

Lemma 2.1 Let $h(x) = \|x\|_0$ for $x \in \mathbb{R}^p$. Consider an arbitrary point $\bar{x} \in \mathbb{R}^p$. Then,

$$\widehat{\partial}h(\bar{x}) = \partial h(\bar{x}) = \left\{ \mu \in \mathbb{R}^p \mid \mu_i = 0 \text{ for } i \in \text{supp}(\bar{x}) \right\} = \partial^\infty h(\bar{x}) = [\widehat{\partial}h(\bar{x})]^\infty.$$

Proof: By [1, Theorem 1] the first two equalities hold. To establish the third one, let $\bar{v} \in \partial^\infty h(\bar{x})$. By Definition 2.1, there exist $x^k \xrightarrow{h} \bar{x}$, $\lambda^k \downarrow 0$ and $v^k \in \widehat{\partial}h(x^k)$ with $\lambda^k v^k \rightarrow \bar{v}$ as $k \rightarrow \infty$. From $x^k \xrightarrow{h} \bar{x}$, it follows that $\text{supp}(x^k) = \text{supp}(\bar{x})$ for all sufficiently large k . Along with $v^k \in \widehat{\partial}h(x^k)$ and $\lambda^k v^k \rightarrow \bar{v}$, we have $\bar{v}_i = 0$ for $i \in \text{supp}(\bar{x})$. Then,

$$\partial^\infty h(\bar{x}) \subseteq \left\{ \xi \in \mathbb{R}^n \mid \xi_i = 0 \text{ for } i \in \text{supp}(\bar{x}) \right\}. \quad (3)$$

Conversely, let \bar{v} be an arbitrary vector from the set on the right hand side of (3). Take $x^k = \bar{x}$, $\lambda^k = \frac{1}{k}$ and $v^k = k\bar{v}$ for each k . Clearly, $x^k \xrightarrow{h} \bar{x}$ and $v^k \in \widehat{\partial}h(x^k)$ with $\lambda^k v^k \rightarrow \bar{v}$. So, $\bar{v} \in \partial^\infty h(\bar{x})$ and the converse inclusion in (3) holds. Thus, the third equality holds. Recall that $\widehat{\partial}h(\bar{x})$ is closed and convex. Since $0 \in \widehat{\partial}h(\bar{x})$ and $tv \in \widehat{\partial}h(\bar{x})$ for any $v \in \widehat{\partial}h(\bar{x})$ and $t \geq 0$, by [14, Theorem 8.3] we have $\widehat{\partial}h(\bar{x}) = [\widehat{\partial}h(\bar{x})]^\infty$. The last equality holds. \square

Lemma 2.2 Consider an arbitrary point $\bar{x} \in \mathcal{S}$. Then the following equalities hold:

$$\widehat{\partial}\delta_{\mathcal{S}}(\bar{x}) = \partial\delta_{\mathcal{S}}(\bar{x}) = \{\omega\bar{x} \mid \omega \in \mathbb{R}\} = \partial^{\infty}\delta_{\mathcal{S}}(\bar{x}) = [\widehat{\partial}\delta_{\mathcal{S}}(\bar{x})]^{\infty}.$$

Proof: Let $F(x) = \|x\| - 1$ for $x \in \mathbb{R}^p$. Since $\|\bar{x}\| = 1$, there exists an open neighborhood \mathcal{U} of \bar{x} such that for all $x \in \mathcal{U}$, $x \neq 0$, which implies that F is continuously differentiable in \mathcal{U} . Moreover, $F'(\bar{x})$ has full rank 1. From $\mathcal{S}^p \cap \mathcal{U} = F^{-1}(0) \cap \mathcal{U}$ and [15, Exercise 6.7],

$$\widehat{\mathcal{N}}_{F^{-1}(0) \cap \mathcal{U}}(\bar{x}) = \mathcal{N}_{F^{-1}(0) \cap \mathcal{U}}(\bar{x}) = \{\omega\bar{x} \mid \omega \in \mathbb{R}\}.$$

Notice that $\mathcal{N}_{\mathcal{S}}(\bar{x}) = \mathcal{N}_{\mathcal{S} \cap \mathcal{U}}(\bar{x})$ and $\widehat{\mathcal{N}}_{\mathcal{S}}(\bar{x}) = \widehat{\mathcal{N}}_{\mathcal{S} \cap \mathcal{U}}(\bar{x})$. The last equation implies that $\widehat{\mathcal{N}}_{\mathcal{S}}(\bar{x}) = \mathcal{N}_{\mathcal{S}}(\bar{x}) = \{\omega\bar{x} \mid \omega \in \mathbb{R}\}$. This shows that the set \mathcal{S} is regular in the sense of [15, Definition 6.4]. The desired result then follows by [15, Exercise 8.14]. \square

By Lemma 2.1 and 2.2, we can provide the following characterization for $\widehat{\partial}(\delta_{\mathcal{S}} + h)$.

Lemma 2.3 Let $h(x) = \|x\|_0$ for $x \in \mathbb{R}^p$. Consider an arbitrary point $\bar{x} \in \mathcal{S}^p$. Then,

$$\widehat{\partial}(\delta_{\mathcal{S}} + h)(\bar{x}) = \widehat{\partial}\delta_{\mathcal{S}}(\bar{x}) + \widehat{\partial}h(\bar{x}).$$

Proof: By [15, Corollary 10.9], $\widehat{\partial}(\delta_{\mathcal{S}} + h)(\bar{x}) \supseteq \widehat{\partial}\delta_{\mathcal{S}}(\bar{x}) + \widehat{\partial}h(\bar{x})$. So, it suffices to argue

$$\widehat{\partial}(\delta_{\mathcal{S}} + h)(\bar{x}) \subseteq \widehat{\partial}\delta_{\mathcal{S}}(\bar{x}) + \widehat{\partial}h(\bar{x}).$$

Write $J = \text{supp}(\bar{x})$ and $\bar{J} = \{1, \dots, p\} \setminus J$. Pick up an arbitrary $v \in \widehat{\partial}(\delta_{\mathcal{S}} + h)(\bar{x})$. We shall prove that $v \in \widehat{\partial}\delta_{\mathcal{S}}(\bar{x}) + \widehat{\partial}h(\bar{x})$, and the stated inclusion then holds. Notice that $v = (0_J; v_{\bar{J}}) + (v_J; 0_{\bar{J}})$. By Lemma 2.1, clearly, $(0_J; v_{\bar{J}}) \in \widehat{\partial}h(\bar{x})$. Thus, it suffices to prove that $(v_J; 0_{\bar{J}}) \in \widehat{\partial}\delta_{\mathcal{S}}(\bar{x})$. We proceed the arguments by the following two cases.

Case 1: $|J| = 1$. Since $|J| = 1$, by Lemma 2.2 it is immediate to have $(v_J; 0_{\bar{J}}) \in \widehat{\partial}\delta_{\mathcal{S}}(\bar{x})$.

Case 2: $|J| > 1$. Since $v \in \widehat{\partial}(\delta_{\mathcal{S}} + h)(\bar{x})$, by Definition 2.1, it follows that

$$\begin{aligned} 0 &\leq \liminf_{x' \rightarrow \bar{x}, x' \neq \bar{x}} \frac{h(x') + \delta_{\mathcal{S}}(x') - h(\bar{x}) - \delta_{\mathcal{S}}(\bar{x}) - \langle v, x' - \bar{x} \rangle}{\|x' - \bar{x}\|} \\ &\leq \liminf_{x' \in \mathcal{S}, \text{supp}(x') = J, x' \rightarrow \bar{x}, x' \neq \bar{x}} \frac{h(x') - h(\bar{x}) - \langle v, x' - \bar{x} \rangle}{\|x' - \bar{x}\|} \\ &= \liminf_{z \in \mathcal{S}^{|J|}, \text{supp}(z) = J, z \rightarrow \bar{x}_J, z \neq \bar{x}_J} \frac{-\langle v_J, z - \bar{x}_J \rangle}{\|z - \bar{x}_J\|} \\ &= \liminf_{z \in \mathcal{S}^{|J|}, \text{supp}(z) = J, z \rightarrow \bar{x}_J, z \neq \bar{x}_J} \frac{\delta_{\mathcal{S}^{|J|}}(z) - \delta_{\mathcal{S}^{|J|}}(\bar{x}_J) - \langle v_J, z - \bar{x}_J \rangle}{\|z - \bar{x}_J\|} \\ &= \liminf_{z \rightarrow \bar{x}_J, \text{supp}(z) = J, z \neq \bar{x}_J} \frac{\delta_{\mathcal{S}^{|J|}}(z) - \delta_{\mathcal{S}^{|J|}}(\bar{x}_J) - \langle v_J, z - \bar{x}_J \rangle}{\|z - \bar{x}_J\|} \\ &= \liminf_{z \rightarrow \bar{x}_J, \|z - \bar{x}_J\| \leq \delta, z \neq \bar{x}_J} \frac{\delta_{\mathcal{S}^{|J|}}(z) - \delta_{\mathcal{S}^{|J|}}(\bar{x}_J) - \langle v_J, z - \bar{x}_J \rangle}{\|z - \bar{x}_J\|} \\ &= \liminf_{z \rightarrow \bar{x}_J, z \neq \bar{x}_J} \frac{\delta_{\mathcal{S}^{|J|}}(z) - \delta_{\mathcal{S}^{|J|}}(\bar{x}_J) - \langle v_J, z - \bar{x}_J \rangle}{\|z - \bar{x}_J\|} \end{aligned}$$

where $\mathcal{S}^{|J|} = \{z \in \mathbb{R}^{|J|} \mid \|z\| = 1\}$ is the unit sphere in $\mathbb{R}^{|J|}$. The last inequality shows that $v_J \in \widehat{\partial}\delta_{\mathcal{S}^{|J|}}(\bar{x}_J)$. By Lemma 2.2, there exists some $\bar{\omega} \in \mathbb{R}$ such that $v_J = \bar{\omega}x_J$, which in turn implies that $(v_J; 0_{\bar{J}}) \in \widehat{\partial}\delta_{\mathcal{S}}(\bar{x})$. The stated inclusion follows. \square

By using Lemma 2.1-2.3, we have the subdifferential characterization for $\delta_{\mathcal{S}} + h$.

Proposition 2.1 *Let $h(x) = \|x\|_0$ for $x \in \mathbb{R}^p$. Consider an arbitrary $\bar{x} \in \mathcal{S}$. Then,*

$$\partial(\delta_{\mathcal{S}} + h)(\bar{x}) = \delta_{\mathcal{S}}(\bar{x}) + \partial h(\bar{x}).$$

Proof: Let $v \in \partial(\delta_{\mathcal{S}} + h)(\bar{x})$. Then there exist $x^k \xrightarrow{\delta_{\mathcal{S}} + h} \bar{x}$ and $v^k \in \widehat{\partial}(\delta_{\mathcal{S}} + h)(x^k)$ with $v^k \rightarrow v$ as $k \rightarrow \infty$. Write $J = \text{supp}(\bar{x})$ and $\bar{J} = \{1, \dots, p\} \setminus J$. Since $x^k \rightarrow \bar{x} \in \mathcal{S}$, we have $x^k \neq 0$ and $\text{supp}(x^k) \supseteq J$ for all sufficiently large k . Since $\delta_{\mathcal{S}}(x^k) + h(x^k) \rightarrow \delta_{\mathcal{S}}(\bar{x}) + h(\bar{x})$, we must have $x^k \in \mathcal{S}$ and $h(x^k) \rightarrow h(\bar{x})$ for all sufficiently large k . The latter, along with $\text{supp}(x^k) \supseteq J$, implies that $\text{supp}(x^k) = J$ for all sufficiently large k . By Lemma 2.3, for each k , $v^k \in \widehat{\partial}\delta_{\mathcal{S}}(x^k) + \widehat{\partial}h(x^k) = \partial\delta_{\mathcal{S}}(x^k) + \partial h(x^k)$. So, there exist $\xi^k \in \partial\delta_{\mathcal{S}}(x^k)$ and $\eta^k \in \partial h(x^k)$ such that $v^k = \xi^k + \eta^k$ for each k . By Lemma 2.2, for each sufficiently large k , there exists $\omega_k \in \mathbb{R}$ such that $\xi^k = \omega_k x^k$. While by Lemma 2.1 and $\text{supp}(x^k) = J$ for all sufficiently large k , we deduce that for each sufficiently large k ,

$$v_J^k = \xi_J^k, \quad v_{\bar{J}}^k = \eta_{\bar{J}}^k \quad \text{and} \quad \eta_J^k = 0.$$

Together with $v^k \rightarrow v$, it follows that $\eta^k \rightarrow (0_J; v_{\bar{J}})$, and consequently $\{\xi^k\}$ is convergent. Without loss of generality, we assume that $\xi^k \rightarrow \bar{\xi}$. Then, $\omega_k x^k \rightarrow \bar{\xi}$. From $\omega_k x_J^k \rightarrow \bar{\xi}_J$, it follows that the sequence $\{\omega_k\}$ is bounded. We assume (if necessary taking a subsequence) that $\omega^k \rightarrow \bar{\omega}$. Then $\bar{\xi} = \bar{\omega}\bar{x}$. Thus, $v = (0_J; v_{\bar{J}}) + \bar{\omega}\bar{x}$. Since $\xi^k \in \partial\delta_{\mathcal{S}}(x^k)$, $\eta^k \in \partial h(x^k)$, $\delta_{\mathcal{S}}(x^k) \rightarrow \delta_{\mathcal{S}}(\bar{x})$ and $h(x^k) \rightarrow h(\bar{x})$, we have $\bar{\omega}\bar{x} \in \partial\delta_{\mathcal{S}}(\bar{x})$ and $(0_J; v_{\bar{J}}) \in \partial h(\bar{x})$. Thus, $v \in \partial\delta_{\mathcal{S}}(\bar{x}) + \partial h(\bar{x})$. By the arbitrariness of v in $\partial(\delta_{\mathcal{S}} + h)(\bar{x})$, we obtain

$$\partial(\delta_{\mathcal{S}} + h)(\bar{x}) \subseteq \delta_{\mathcal{S}}(\bar{x}) + \partial h(\bar{x}).$$

On the other hand, by Lemma 2.3 and Lemma 2.1, the following inclusions hold:

$$\partial\delta_{\mathcal{S}}(\bar{x}) + \partial h(\bar{x}) = \widehat{\partial}\delta_{\mathcal{S}}(\bar{x}) + \widehat{\partial}h(\bar{x}) = \widehat{\partial}(\delta_{\mathcal{S}} + h)(\bar{x}) \subseteq \partial(\delta_{\mathcal{S}} + h)(\bar{x}).$$

The desired result then follows from the last two equations. The proof is completed. \square

2.2 Kurdyka-Łojasiewicz property

Definition 2.2 *Let $f: \mathbb{X} \rightarrow (-\infty, +\infty]$ be a proper function. The function f is said to have the Kurdyka-Łojasiewicz (KL) property at $\bar{x} \in \text{dom } \partial f$ if there exist $\eta \in (0, +\infty]$, a continuous concave function $\varphi: [0, \eta) \rightarrow \mathbb{R}_+$ satisfying the following conditions*

- (i) $\varphi(0) = 0$ and φ is continuously differentiable on $(0, \eta)$;
- (ii) for all $s \in (0, \eta)$, $\varphi'(s) > 0$,

and a neighborhood \mathcal{U} of \bar{x} such that for all $x \in \mathcal{U} \cap [f(\bar{x}) < f(x) < f(\bar{x}) + \eta]$,

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

If the corresponding φ can be chosen as $\varphi(s) = c\sqrt{s}$ for some $c > 0$, then f is said to have the KL property at \bar{x} with an exponent of $1/2$. If f has the KL property of exponent $1/2$ at each point of $\text{dom } \partial f$, then f is called a KL function of exponent $1/2$.

Remark 2.2 By [3, Lemma 2.1], a proper function has the KL property of exponent $1/2$ at any noncritical point. Hence, to show that it is a KL function of exponent $1/2$, it suffices to check whether it has the KL property of exponent $1/2$ at all critical points.

3 KL property of exponent $1/2$ of Φ_ν

To achieve the KL property of exponent $1/2$ of Φ_ν , we first establish this property of

$$g(z) := z^\top H z + \delta_{\mathcal{S}^m}(z) \quad \forall z \in \mathbb{R}^m \quad (4)$$

where H is an $m \times m$ real symmetric matrix and \mathcal{S}^m is the unit sphere in \mathbb{R}^m .

Lemma 3.1 For the function g in (4), it holds that $\text{crit}g = \{z \in \mathcal{S}^m \mid H z = \langle z, H z \rangle z\}$. Also, by letting H have the eigenvalue decomposition $P \Lambda P^\top$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ for $\lambda_1 \geq \dots \geq \lambda_m$ and $P \in \mathbb{O}^m$, $\text{crit}g = P W$ with $W = \{y \in \mathcal{S}^m \mid \Lambda y = \langle y, \Lambda y \rangle y\}$.

Proof: By [15, Exercise 8.8] and Lemma 2.2, it immediately follows that for any $z \in \mathbb{R}^m$,

$$\partial g(z) = 2H z + \partial \delta_{\mathcal{S}^m}(z) = 2H z + \llbracket z \rrbracket \quad (5)$$

where $\llbracket z \rrbracket$ is the subspace spanned by z . Choose an arbitrary $\bar{z} \in \text{crit}g$. Then, $0 \in \partial g(\bar{z})$. From (5), there exists $\bar{t} \in \mathbb{R}$ such that $0 = 2H \bar{z} + \bar{t} \bar{z}$. Together with $\|\bar{z}\| = 1$, we have $\bar{t} = -2\langle \bar{z}, H \bar{z} \rangle$. This shows that $\bar{z} \in \{z \in \mathcal{S}^m \mid H z = \langle z, H z \rangle z\}$. By the arbitrariness of \bar{z} ,

$$\text{crit}g \subseteq \{z \in \mathcal{S}^m \mid H z = \langle z, H z \rangle z\}.$$

The converse inclusion is immediate to check by Lemma 2.2. Thus, the first part follows. The second part is immediate by the conclusion of the first part. \square

Proposition 3.1 The function g defined by (4) is a KL function of exponent $1/2$.

Proof: Fix an arbitrary $\bar{z} \in \text{crit}g$. Let H have the eigenvalue decomposition as in Lemma 3.1. Then $\bar{y} = P^\top \bar{z} \in \text{crit}\psi$ where ψ is defined by (18) with $D = \Lambda$. By Proposition 1, there exist $\eta > 0, \delta > 0$ and $c > 0$ such that for all $y \in \mathbb{B}(\bar{y}, \delta) \cap [\psi(\bar{y}) < \psi(y) < \psi(\bar{y}) + \eta]$,

$$\text{dist}(0, \partial \psi(y)) \geq c\sqrt{\psi(y) - \psi(\bar{y})}.$$

Fix an arbitrary $z \in \mathbb{B}(\bar{z}, \delta) \cap [g(\bar{z}) < g(z) < g(\bar{z}) + \eta]$. Clearly, $z \in \mathcal{S}^m$. Write $y = P^\top z$. Then $y \in \mathcal{S}^m$ and $g(z) = \psi(y)$. Together with $g(\bar{z}) = g(\bar{y})$, it follows that

$$y \in \mathbb{B}(\bar{y}, \delta) \cap [\psi(\bar{y}) < \psi(y) < \psi(\bar{y}) + \eta].$$

In addition, from (5) and the eigenvalue decomposition of H , it is easy to check that

$$\partial g(z) = P\partial\psi(y).$$

Thus, $\text{dist}(0, \partial g(z)) = \text{dist}(0, P\partial\psi(y)) = \text{dist}(0, \partial\psi(y)) \geq c\sqrt{\psi(y) - \psi(\bar{y})}$. Together with $\psi(y) - \psi(\bar{y}) = g(z) - g(\bar{z})$, it follows that g has the KL property with exponent of $1/2$ at \bar{z} . By the arbitrariness of \bar{z} in $\text{crit}g$, g is a KL function of exponent $1/2$. \square

From Proposition 2.1 and [15, Exercise 8.8(c)], it follows that for any given $x \in \mathcal{S}$,

$$\partial\Phi_\nu(x) = 2Ax + \partial\delta_S(x) + \nu\partial\|\cdot\|_0(x). \quad (6)$$

In particular, the following conclusion holds for the set of critical points of Φ_ν .

Lemma 3.2 *For the function Φ_ν in equation (2), we have $\bar{x} \in \text{crit}\Phi_\nu$ if and only if there exists $\bar{\omega} \in \mathbb{R}$ such that $A_{JJ}\bar{x}_J = \bar{\omega}\bar{x}_J$ with $J = \text{supp}(\bar{x})$.*

Proof: Let \bar{x} be an arbitrary point in $\text{crit}\Phi_\nu$. Notice that $0 \in 2A\bar{x} + \partial\delta_S(\bar{x}) + \nu\partial\|\cdot\|_0(\bar{x})$. By Lemma 2.1 and Lemma 2.2, there exist $\tilde{\omega} \in \mathbb{R}$ and $\tilde{\xi} \in \nu\partial\|\cdot\|_0(\bar{x})$ such that

$$0 = 2A\bar{x} + \tilde{\omega}\bar{x} + \tilde{\xi}.$$

Notice that $\tilde{\xi}_J = 0$. We have $A_{JJ}\bar{x}_J + \tilde{\omega}\bar{x}_J = 0$. Thus, $\bar{\omega} = -\tilde{\omega}$ satisfies the requirement. Conversely, suppose that there exists $\bar{\omega} \in \mathbb{R}$ such that $A_{JJ}\bar{x}_J = \bar{\omega}\bar{x}_J$ with $J = \text{supp}(\bar{x})$. Then $0 \in 2A_{JJ}\bar{x}_J + \llbracket\bar{x}_J\rrbracket$. Take $\xi_J = 0$ and $\xi_{\bar{J}} = -2A_{\bar{J}J}\bar{x}_J$ where $\bar{J} = \{1, 2, \dots, p\} \setminus J$. Clearly, $0 \in 2A\bar{x} + \llbracket\bar{x}\rrbracket + \xi$, that is, $0 \in \partial\Phi_\nu(\bar{x})$. Thus, $\bar{x} \in \text{crit}\Phi_\nu$. \square

Now we are in a position to establish the KL property of exponent $1/2$ for Φ_ν .

Proposition 3.2 *The function Φ_ν in equation (2) is a KL function of exponent $1/2$.*

Proof: Fix an arbitrary $\bar{x} \in \text{crit}\Phi_\nu$ and write $J = \text{supp}(\bar{x})$. Let g be the function defined by (4) with $H = A_{JJ}$. By Proposition 3.1, g is a KL function of exponent $1/2$. So, there exist $\delta_1 > 0, \eta_1 > 0$ and $c > 0$ such that for all $z \in \mathbb{B}(\bar{x}_J, \delta_1) \cap [g(\bar{x}_J) < g(z) < g(\bar{x}_J) + \eta_1]$

$$\text{dist}(0, \partial g(z)) \geq c\sqrt{g(z) - g(\bar{x}_J)}. \quad (7)$$

By the continuity, for any $\eta_2 \in (0, \frac{\nu}{3})$ there exist $\delta_2 > 0$ such that for all $x \in \mathbb{B}(\bar{x}, \delta_2)$,

$$|x^\top Ax - \bar{x}^\top A\bar{x}| < \eta_2. \quad (8)$$

Take $\delta = \min(\delta_1, \delta_2)$ and $\eta = \min(\eta_1, \eta_2, \frac{\nu}{3})$. Let x be an arbitrary point from the set $\mathbb{B}(\bar{x}, \delta) \cap [\Phi_\nu(\bar{x}) < \Phi_\nu(x) < \Phi_\nu(\bar{x}) + \eta]$. Clearly, $x \in \mathcal{S}$. Moreover, we necessarily have

$$x^\top Ax > \bar{x}^\top A\bar{x}. \quad (9)$$

(If not, we have $x^\top Ax \leq \bar{x}^\top A\bar{x}$, which together with $\Phi_\nu(\bar{x}) < \Phi_\nu(x) < \Phi_\nu(\bar{x}) + \eta$ and (8) yields that $\|\bar{x}\|_0 + 1 \leq \|x\|_0 < \|\bar{x}\|_0 + \frac{1}{\nu}(\eta + \eta_2) < \|\bar{x}\|_0 + 1$, which is impossible.) Then, combining (9) with $\Phi_\nu(x) < \Phi_\nu(\bar{x}) + \eta$, we deduce that $\|x\|_0 \leq \|\bar{x}\|_0$. In addition,

by reducing δ if necessary, we also have $\|x\|_0 \geq \|\bar{x}\|_0$. Thus, $\|x\|_0 = \|\bar{x}\|_0$. Notice that $\text{supp}(x) \supseteq \text{supp}(\bar{x})$ (if necessary shrinking the value of δ). Therefore, it holds that $\text{supp}(x) = \text{supp}(\bar{x}) = J$. From the expression of $\partial\Phi_\nu(x)$ in equation (6), it follows that

$$\begin{aligned} \text{dist}(0, \partial\Phi_\nu(x)) &= \min_{\zeta \in \partial\delta_S(x), \xi \in \nu\partial\|\cdot\|_0(x)} \|2Ax + \zeta + \xi\| \\ &= \min_{w \in \mathbb{R}, \xi \in \partial\|\cdot\|_0(x)} \|2Ax + wx + \xi\| \\ &= \min_{w \in \mathbb{R}} \|2A_J x_J + w x_J\| = \text{dist}(0, \partial g(x_J)) \end{aligned} \quad (10)$$

where the second equality is using Lemma 2.2, the third one is due to Lemma 2.1 and $\text{supp}(x) = \text{supp}(\bar{x}) = J$, and the last one is by the definition of g and equation (5). In addition, from $\text{supp}(x) = \text{supp}(\bar{x}) = J$ and the expressions of Φ_ν and g , it follows that

$$\begin{aligned} \Phi_\nu(x) - \Phi_\nu(\bar{x}) &= x^\top Ax - \bar{x}^\top A\bar{x} + \nu(\|x\|_0 - \|\bar{x}\|_0) = x^\top Ax - \bar{x}^\top A\bar{x} \\ &= x_J^\top A_J x_J - \bar{x}_J^\top A_J \bar{x}_J = g(x_J) - g(\bar{x}_J) \end{aligned}$$

which, along with $x \in [\Phi_\nu(\bar{x}) < \Phi_\nu(x) < \Phi_\nu(\bar{x}) + \eta]$ and $\eta \leq \eta_1$, implies that

$$x_J \in [g(\bar{x}_J) < g(x_J) < g(\bar{x}_J) + \eta_1].$$

Notice that $\|x_J - \bar{x}_J\| = \|x - \bar{x}\| \leq \delta < \delta_1$. Thus, from (7) and inequality (10),

$$\text{dist}(0, \partial\Phi_\nu(x)) = \text{dist}(0, \partial g(x_J)) \geq c\sqrt{g(x_J) - g(\bar{x}_J)} = c\sqrt{\Phi_\nu(x) - \Phi_\nu(\bar{x})}.$$

By the arbitrariness of x in $\mathbb{B}(\bar{x}, \delta) \cap [\Phi_\nu(\bar{x}) < \Phi_\nu(x) < \Phi_\nu(\bar{x}) + \eta]$, this shows that f has the KL property of exponent $1/2$ at \bar{x} . By the arbitrariness of \bar{x} in $\text{crit}\Phi_\nu$, the function Φ_ν is a KL function of exponent $1/2$. The proof is then completed. \square

4 Globally and linearly convergent PGM

The proximal mapping of $\|\cdot\|_0 + \delta_S(\cdot)$ is multivalued instead of single-valued. Now we show that the components of its proximal mapping are accessible, that is, for any given $z \in \mathbb{R}^p$, one may obtain a global optimal solution of the following problem:

$$\min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - z\|^2 + \nu \|x\|_0 : \|x\| = 1 \right\}. \quad (11)$$

Lemma 4.1 *Let Q be a $p \times p$ signed permutation matrix such that $|z|^\downarrow = Qz$. If x^* is a global optimal solution of (11), then Qx^* is a global optimal solution of the problem*

$$\min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - |z|^\downarrow\|^2 + \nu \|x\|_0 : \|x\| = 1 \right\}. \quad (12)$$

Conversely, if x^ is globally optimal to (12), then $Q^\top x^*$ is globally optimal to (11).*

By Lemma 4.1 it suffices to argue that the optimal solutions of (12) are available.

Proposition 4.1 For each $\kappa \in \{1, 2, \dots, p\}$, write $\vartheta(\kappa) := \||z|^{\kappa, \downarrow}\| - \||z|^{\kappa-1, \downarrow}\|$ where $|z|^{0, \downarrow}$ is stipulated to be 0. Then, the following assertions hold.

- (i) $\vartheta(1) \geq \vartheta(2) \geq \dots \geq \vartheta(p)$.
- (ii) If $\nu \leq \vartheta(p)$, $x^* = \frac{1}{\||z|^\downarrow}\|} |z|^\downarrow$ is a global optimal solution of (12); if $\nu \geq \vartheta(1)$, then $x^* = \frac{1}{\||z|^{1, \downarrow}\|} (|z|^{1, \downarrow}; 0)$ is globally optimal to (12); and if there exists $l \in \{1, \dots, p\}$ such that $\nu \in [\vartheta(l+1), \vartheta(l)]$, then $x^* = \frac{1}{\||z|^{l, \downarrow}\|} (|z|^{l, \downarrow}; 0)$ is globally optimal to (12).

Proof: (i) According to the definition of $\vartheta(\cdot)$, it is immediate to obtain that

$$\vartheta(\kappa) = \frac{\||z|^{\kappa, \downarrow}\|^2 - \||z|^{\kappa-1, \downarrow}\|^2}{\||z|^{\kappa, \downarrow}\| + \||z|^{\kappa-1, \downarrow}\|}. \quad (13)$$

For each $\kappa \in \{1, 2, \dots, p-1\}$, by the definition of $|z|^{\kappa, \downarrow}$, it is immediate to obtain that

$$\begin{aligned} \||z|^{\kappa, \downarrow}\|^2 - \||z|^{\kappa-1, \downarrow}\|^2 &= (|z|_{\kappa}^\downarrow)^2 \geq (|z|_{\kappa+1}^\downarrow)^2 = \||z|^{\kappa+1, \downarrow}\|^2 - \||z|^{\kappa, \downarrow}\|^2, \\ \||z|^{\kappa, \downarrow}\| + \||z|^{\kappa-1, \downarrow}\| &\leq \||z|^{\kappa+1, \downarrow}\| + \||z|^{\kappa, \downarrow}\|. \end{aligned}$$

Together with (13), we have $\vartheta(\kappa) \geq \vartheta(\kappa+1)$, and the desired result follows.

(ii) Let v^* denote the optimal value of (12). One can check that $v^* = \min_{\kappa \in \{1, \dots, p\}} \phi(\kappa)$ with

$$\phi(\kappa) := \min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - |z|^\downarrow\|^2 : \|x\|_0 \leq \kappa, \|x\| = 1 \right\}. \quad (14)$$

By Lemma 2 in Appendix B, it follows that $\phi(\kappa) = \frac{1}{2}(1 + \|z\|^2 - 2\||z|^{\kappa, \downarrow}\|)$. Then,

$$\Delta\phi(\kappa) := \phi(\kappa) - \phi(\kappa+1) = \||z|^{\kappa+1, \downarrow}\| - \||z|^{\kappa, \downarrow}\| = \vartheta(\kappa+1).$$

When $\nu \leq \vartheta(p)$, by part (i) we have $\nu \leq \vartheta(\kappa)$ for all $\kappa \in \{1, \dots, p\}$. Along with the last equation, $\Delta\phi(\kappa) \leq 0$ for $\kappa = 1, \dots, p-1$, which implies that $\phi(1) + \nu \geq \phi(2) + 2\nu \geq \dots \geq \phi(p) + p\nu$. Hence, $v^* = \phi(p) + p\nu$, and the corresponding optimal solution is $x^* = \frac{1}{\||z|^\downarrow}\|} |z|^\downarrow$ by Lemma 2 in Appendix B. Using the similar arguments, one may obtain the rest conclusions. \square

Motivated by Proposition 4.1, we use the following PGM to solve the problem (11).

Algorithm 1 Proximal gradient method for (11)

Initialization: Select $\nu > 0$, $\tau = 2\gamma\|A\|$ for $\gamma > 1$ and an initial $x^0 \in \mathcal{S}$. Set $k = 0$.

while the stopping conditions are not satisfied **do**

$$x^{k+1} \in \arg \min_{x \in \mathcal{S}} \left\{ 2\langle Ax^k, x - x^k \rangle + \frac{\tau}{2} \|x - x^k\|^2 + \nu \|x\|_0 \right\}. \quad (15)$$

end while

Remark 4.1 Since \mathcal{S} is compact and the zero-norm is lower semicontinuous, from Weierstrass' theorem it follows that x^{k+1} in (15) is well defined. Notice that for any $x \in \mathbb{R}^p$,

$$x^\top Ax \leq (x^k)^\top Ax^k + 2\langle Ax^k, x - x^k \rangle + \|A\| \|x - x^k\|^2.$$

So, each iterate of Algorithm 1 is minimizing the upper approximation of (2).

Lemma 4.2 Let $\{x^k\}_{k \in \mathbb{N}}$ be the sequence given by Algorithm 1. Then, for each $k \in \mathbb{N}$,

$$\Phi_\nu(x^{k+1}) \leq \Phi_\nu(x^k) - (\gamma - 1)\|A\| \|x^{k+1} - x^k\|^2, \quad (16)$$

and hence $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\|^2 < \infty$, which implies that $\lim_{k \rightarrow \infty} \|x^{k+1} - x^k\| = 0$.

Proof: From the optimality of x^{k+1} and the feasibility of x^k to the subproblem (15),

$$2\langle Ax^k, x^{k+1} - x^k \rangle + \frac{\tau}{2} \|x^{k+1} - x^k\|^2 + \nu \|x^{k+1}\|_0 \leq \nu \|x^k\|_0.$$

Together with the descent lemma (see [5, Appendix A.24]), it then follows that

$$\begin{aligned} \Phi_\nu(x^{k+1}) &= \langle x^{k+1}, Ax^{k+1} \rangle + \nu \|x^{k+1}\|_0 + \delta_{\mathcal{S}}(x^{k+1}) \\ &\leq \langle x^k, Ax^k \rangle + \langle 2Ax^k, x^{k+1} - x^k \rangle + \|A\| \|x^{k+1} - x^k\|^2 + \nu \|x^{k+1}\|_0 \\ &\leq \langle x^k, Ax^k \rangle + \nu \|x^k\|_0 - (\gamma - 1)\|A\| \|x^{k+1} - x^k\|^2 \\ &= \Phi_\nu(x^k) - (\gamma - 1)\|A\| \|x^{k+1} - x^k\|^2 \end{aligned}$$

where the last equality is using $x^k \in \mathcal{S}$. This implies the desired result. \square

Lemma 4.3 Let $\{x^k\}_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1. For each $k \in \mathbb{N}$, we have $u^k := 2A(x^k - x^{k-1}) - \tau(x^k - x^{k-1}) \in \partial\Phi_\nu(x^k)$ and $\|u^k\| \leq 2(\gamma + 1)\|A\| \|x^k - x^{k-1}\|$.

Proof: Fix an arbitrary $k \in \mathbb{N}$. From the optimality of x^k to (15), it follows that

$$0 \in 2Ax^{k-1} + \tau(x^k - x^{k-1}) + \nu\partial(\|\cdot\|_0 + \delta_{\mathcal{S}}(\cdot))(x^k),$$

which is equivalent to saying that

$$2A(x^k - x^{k-1}) - \tau(x^k - x^{k-1}) \in 2Ax^k + \nu\partial(\|\cdot\|_0 + \delta_{\mathcal{S}}(\cdot))(x^k) = \partial\Phi_\nu(x^k).$$

That is, $u^k \in \partial\Phi_\nu(x^k)$. The desired result then follows by using $\tau = 2\gamma\|A\|$. \square

By using Lemma 4.2, we can obtain the following weak convergence result.

Proposition 4.2 Let $\{x^k\}_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1 and denote by $\varpi(x^0)$ the set of limit points of $\{x^k\}_{k \in \mathbb{N}}$. Then, the following assertions holds:

- (i) $\emptyset \neq \varpi(x^0) \subseteq \text{crit}\Phi_\nu$;
- (ii) $\lim_{k \rightarrow \infty} \text{dist}(x^k, \varpi(x^0)) = 0$;

(iii) $\varpi(x^0)$ is a nonempty, compact and connected set;

(iv) The function Φ_ν is finite and constant on $\varpi(x^0)$.

Proof: (i) Since $\{x^k\}_{k \in \mathbb{N}} \subseteq \mathcal{S}$, $\varpi(x^0) \neq \emptyset$. Let x^* be an arbitrary point from $\varpi(x^0)$. There exists a subsequence $\{x^{k_j}\} \rightarrow x^*$ as $j \rightarrow \infty$. By Lemma 4.3, for each $j \in \mathbb{N}$,

$$u^{k_j} := 2A(x^{k_j} - x^{k_j-1}) - \tau(x^{k_j} - x^{k_j-1}) \in \partial\Phi_\nu(x^{k_j}).$$

From Lemma 4.2, we know that $\lim_{j \rightarrow \infty} \|x^{k_j} - x^{k_j-1}\| = 0$. Then $\lim_{j \rightarrow \infty} u^{k_j} = 0$ and $\lim_{j \rightarrow \infty} x^{k_j-1} = x^*$. We next argue that $\lim_{j \rightarrow \infty} \Phi_\nu(x^{k_j}) = \Phi_\nu(x^*)$. By Algorithm 1,

$$\begin{aligned} \Phi_\nu(x^{k_j}) &\leq (x^{k_j-1})^\top A x^{k_j-1} + 2\langle A x^{k_j-1}, x^{k_j} - x^{k_j-1} \rangle + \frac{\tau}{2} \|x^{k_j} - x^{k_j-1}\|^2 + \nu \|x^{k_j}\|_0 \\ &\leq (x^{k_j-1})^\top A x^{k_j-1} + 2\langle A x^{k_j-1}, x^* - x^{k_j-1} \rangle + \frac{\tau}{2} \|x^* - x^{k_j-1}\|^2 + \nu \|x^*\|_0, \end{aligned} \quad (17)$$

which by $\lim_{j \rightarrow \infty} x^{k_j-1} = x^*$ implies that $\limsup_{j \rightarrow \infty} \Phi_\nu(x^{k_j}) \leq \Phi_\nu(x^*)$. In addition, by the lower semicontinuity of Φ_ν , $\liminf_{j \rightarrow \infty} \Phi_\nu(x^{k_j}) \geq \Phi_\nu(x^*)$. Thus, we obtain $\lim_{j \rightarrow \infty} \Phi_\nu(x^{k_j}) = \Phi_\nu(x^*)$. Together with $\lim_{j \rightarrow \infty} u^{k_j} = 0$, by Remark 2.1(ii) we have $0 \in \partial\Phi_\nu(x^*)$. By the arbitrariness of x^* in $\varpi(x^0)$, the inclusion $\varpi(x^0) \subset \text{crit}\Phi_\nu$ follows.

(ii) Suppose on the contrary that there exists a subsequence $\{x^{k_j}\}$ such that

$$\text{dist}(x^{k_j}, \varpi(x^0)) \geq \epsilon \text{ for some } \epsilon > 0.$$

Since the subsequence $\{x^{k_j}\}$ is bounded, we assume (if necessary taking a subsequence) that $\{x^{k_j}\}$ is convergent, say, $\lim_{j \rightarrow \infty} x^{k_j} = \tilde{x}$. By the continuity of the distance function,

$$\text{dist}(\tilde{x}, \varpi(x^0)) = \lim_{j \rightarrow \infty} \text{dist}(x^{k_j}, \varpi(x^0)) \geq \epsilon.$$

On the other hand, $\tilde{x} \in \varpi(x^0)$ by part (i). Thus, we obtain a contradiction.

(iii) Notice that $\emptyset \neq \varpi(x^0)$. The set $\varpi(x^0)$ is compact and nonempty. By following the same arguments as those for [6, Lemma 5(iii)], the set $\varpi(x^0)$ is connected.

(iv) By Lemma 4.2, the sequence $\{\Phi_\nu(x^k)\}_{k \in \mathbb{N}}$ is decreasing. Since $\inf_{x \in \mathbb{R}^p} \Phi_\nu(x) > -\infty$, it follows that $\{\Phi_\nu(x^k)\}_{k \in \mathbb{N}}$ is convergent, and denote its limit by ω^* . Let x^* be an arbitrary point from $\varpi(x^0)$. Then there exists a subsequence $x^{k_j} \rightarrow x^*$ as $j \rightarrow \infty$. From the proof of by part (i) we know that $\lim_{j \rightarrow \infty} \Phi_\nu(x^{k_j}) = \Phi_\nu(x^*)$. From the convergence of $\{\Phi_\nu(x^k)\}_{k \in \mathbb{N}}$, it follows that $\Phi_\nu(x^*) = \omega^*$. By the arbitrariness of $x^* \in \varpi(x^0)$, the function Φ_ν is finite and constant on $\varpi(x^0)$. The proof is completed. \square

Now by invoking Lemma 4.2-4.3, Proposition 4.2 and Proposition 3.2 and following the similar arguments as those for [3, Theorem 3.2 & 3.4] or [6, Theorem 1], we can obtain the following global convergence and linear rate of convergence.

Theorem 4.1 *Let $\{x^k\}_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1. Then,*

(i) *the sequence $\{x^k\}_{k \in \mathbb{N}}$ has a finite length, i.e., $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| < \infty$.*

(ii) *The sequence $\{x^k\}_{k \in \mathbb{N}}$ converges to a critical point x^* of Φ_ν ;*

(iii) *There exist a constant $c_1 > 0$ and $\varrho \in [0, 1)$ such that $\|x^k - x^*\| \leq c_1 \varrho^k$.*

5 Numerical experiments

In this section, we apply Algorithm 1 for seeking the sparse principal components (PCs) of a given sample covariance matrix $\Sigma = X^T X \in \mathbb{S}_+^p$, i.e., solve the problem (1) for $A = -\Sigma$ with Algorithm 1. The sample matrix $X \in \mathbb{R}^{n \times p}$ comes from synthetic or real data. We compare the performance of Algorithm 1 with that of the generalized power method (GPower $_{\ell_0}$) proposed in [10] for solving an equivalent reformulation of (1).

Unless otherwise stated, we always choose $\nu = 0.1\|A\|$, $\gamma = 1.00005$ and the largest eigenvector of Σ as the starting point x^0 for Algorithm 1, and use the default parameter and starting point in the code for GPower $_{\ell_0}$ ¹. Our code can be found in the Website². All numerical results are computed by a desktop computer running on 64-bit Windows Operating System with an Intel(R) Core(TM) i7-8700 CPU 3.20GHz and 8.00 GB RAM.

5.1 Convergence performance of PGM

We test the convergence performance of Algorithm 1 with a synthetic sample covariance matrix $\Sigma \in \mathbb{S}_+^p$, generated in the procedure proposed by Shen and Huang [23]. We first generate a $p \times p$ matrix Z in Matlab command `randn(p, p)`, decompose $(Z + Z^T)/2$ into $U\Lambda U^T$, replace the first m columns of $U \in \mathbb{O}^p$ with pre-specified m sparse orthonormal vectors, and finally synthesize Σ through the eigenvalue decomposition $\Sigma = UDU^T$. A data matrix $X \in \mathbb{R}^{n \times p}$ is generated by drawing n samples from a zero-mean normal distribution with covariance matrix Σ . We consider the setting with $p = 500$ and $m = 2$, the pre-specified m sparse orthonormal vectors are defined as follows:

$$v_{1i} = \begin{cases} \frac{1}{\sqrt{10}} & \text{for } i = 1, \dots, 10, \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad v_{2i} = \begin{cases} \frac{1}{\sqrt{10}} & \text{for } i = 11, \dots, 20, \\ 0 & \text{otherwise,} \end{cases}$$

and $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ with $\lambda_1 = 400$, $\lambda_2 = 300$ and $\lambda_j = 1$ for $j = 3, \dots, p$.

Figure 5.1 plots the iteration error and the objective value error curves yielded by Algorithm 1 for solving (1) with $\nu = 10^{-3}\|A\|$ and $A = -\Sigma$, where Σ is generated randomly as above with $n = 50$. We see that as k increases, the iteration error $\|x^k - x^*\|$ and the objective value error $\Phi_\nu(x^k) - \Phi_\nu(x^*)$ decreases with a linear rate, where x^* is the final output of Algorithm 1. This coincides with the results of Theorem 4.1.

¹<http://www.montefiore.ulg.ac.be/~journee/>

²https://github.com/SoilWu/Code_of_PGM

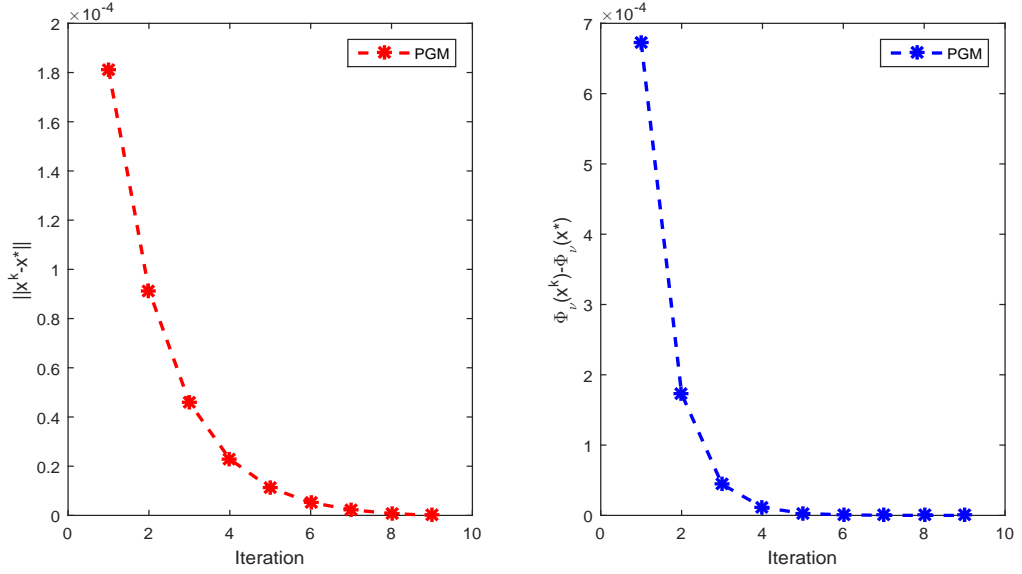


Figure 1: Iteration error and objective value error curves yielded by Algorithm 1

5.2 Recoverability for sparse principal components

We test the recoverability of Algorithm 1 for sparse PCs with a synthetic sample covariance $\Sigma \in \mathbb{S}_+^p$ which is generated in the same way as Subsection 5.1 describes. Figure 5.2 plots the successful recovery curve of Algorithm 1 and GPower_{ℓ_0} under different sample sizes. Among others, the left subfigure in Figure 5.2 plots the successful recovery curve for v_1 , and the right one does for v_2 . For each sample size n , we generate 500 data matrices $X \in \mathbb{R}^{n \times p}$ to formulate $A = -\Sigma$, and then solve the 500 problems with Algorithm 1 to obtain its successful recovery ratio under the sample size n . We say the model underlying the data to be successfully recovered when both $|v^\top z_1|$ and $|v^\top z_2|$ are greater than 0.99, where $z_1, z_2 \in \mathbb{R}^p$ are the unit-norm sparse loading vector computed by the solvers. We see that as the sample size increases, the recoverability of two solvers become higher. When the sample size is small, Algorithm 1 has a little higher recoverability (especially for v_2) than GPower_{ℓ_0} does; when the sample size is large, their recoverability is almost the same. This shows that Algorithm 1 is effective for the recovery of sparse PCs.

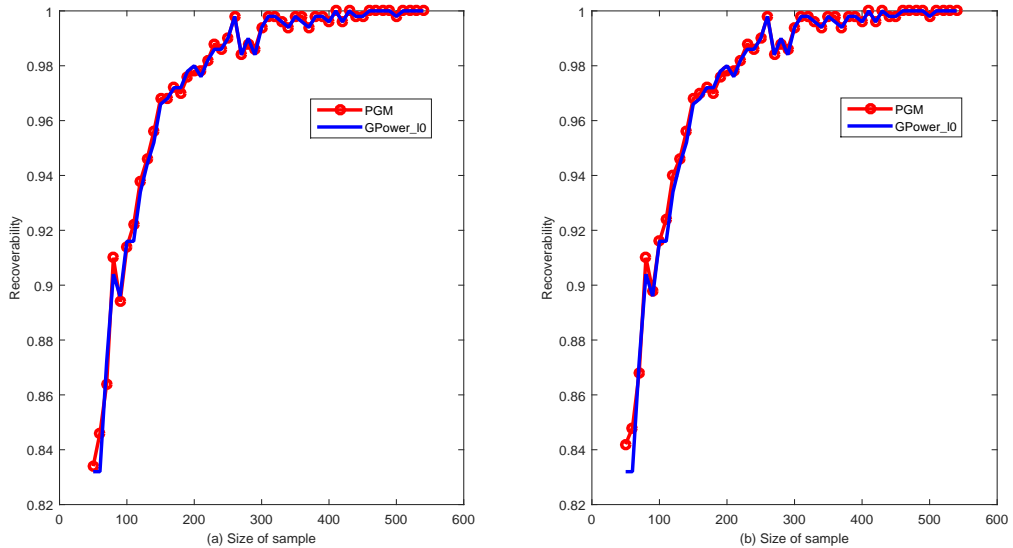


Figure 2: Recoverability comparisons for sparse PCs of PGM and GPower_{ℓ_0}

5.3 Pitprops data set

The PitProps data set [8], consisting of 180 observations for 13 measured variables, is a classical example to illustrate the difficulty of interpreting PCs. We use Algorithm 1 and GPower_{ℓ_0} to compute the first six sparse PCs of the data set. To find the first six sparse PCs rather than the first one, a common way is to use the deflation method [12] for PCA, i.e., the contribution of the previously found PCs is removed from the covariance matrix and then solve (1) with the new one $A = -\Sigma'$, where $\Sigma' = (I_{p \times p} - xx^T)\Sigma(I_{p \times p} - xx^T)$ and x is the optimal solution of (1) with the previous A . Let \hat{Z} be the matrix consisting of sparse PCs yielded by the solvers. As pointed out by Zou [28], when \hat{Z} is correlated, $\text{tr}(\hat{Z}^T \hat{Z})$ is too optimistic for representing the total variance explained by \hat{Z} . So, we follow the method in [28, Section 3.4] to compute the total adjusted variance.

Table 5.3 reports the first six PCs computed by Algorithm 1 and GPower_{ℓ_0} , and the final row reports the cumulative adjusted covariance explained by these sparse PCs. The parameter γ of GPower_{ℓ_0} is chosen to be $0.15 \max_i \|A_i\|^2$, which is much better than its default one. We see that the first sparse PC yielded by Algorithm 1 and GPower_{ℓ_0} have the same cardinality and variance, and the first six sparse PCs yielded by Algorithm 1 have a little higher total adjusted variance than those given by GPower_{ℓ_0} do. Of course, the “cardinality” row shows that the latter has less nonzero loading that the former does.

Table 1: Total covariance explained by sparse PCs of Algorithm 1 and GPower_{ℓ_0}

Variable	Algorithm 1						GPower_{ℓ_0}					
	sPC1	sPC2	sPC3	sPC4	sPC5	sPC6	sPC1	sPC2	sPC3	sPC4	sPC5	sPC6
topdiam	0.424	0	-0.275	0	0	0	0.420	0	0.277	0	0	0
length	0.430	0	-0.280	0	0	0	0.424	0	0.282	0	0	0
moist	0	0.684	0	0	0	0	0	0.710	0	0	0	0
testsg	0	0.678	0	0	0	0	0	0.705	0	0	0	0
ovensg	0	0	0.502	0	0	0.700	0	0	-0.514	0	0	0.727
ringtop	0.268	0	0.435	0.338	0	0	0.282	0	-0.424	0.384	0	0
ringbut	0.403	0	0.352	0	0	0	0.411	0	-0.358	0	0	0
bowmax	0.313	0	0	-0.327	0.313	0	0.311	0	0	-0.453	0	0
bowdist	0.379	0	0	0	0	0	0.374	0	0	0	0	0
whorls	0.399	0	0	-0.350	0	-0.238	0.399	0	0	0	0	0
clear	0	0	0	0.351	0.889	0	0	0	0	0	1	0
knots	0	0.275	0	0.730	-0.334	0	0	0	0	0.804	0	0
diaknot	0	0	-0.531	0	0	0.674	0	0	0.522	0	0	0.686
cardinality	7	3	6	5	3	3	7	2	6	3	1	2
variance(%)	30.7	15.1	14.4	10.2	7.97	6.31	30.7	14.8	14.4	9.32	7.69	6.10
adj_variance(%)	30.7	14.8	14.2	7.11	7.87	8.01	30.7	13.9	14.2	7.98	7.43	5.92
CA_variance(%)	30.7	45.5	59.7	66.8	74.7	80.7	30.7	44.6	58.9	66.8	74.2	80.2

5.4 Gene expression data

Gene expression data results from DNA microarrays and provides the expression level of thousand of genes across several hundreds of experiments. Since the number of genes is larger than the sample size, it is necessary to reduce the dimension of the data. It is well known that the PCA is usually a linear combination of all genes, and it has a difficulty in reducing dimensionality and hence in explaining. While the sparse PCs have fewer linear terms than the PCs do, it is very convenient for them to interpret. We generate a sample covariance $\Sigma \in \mathbb{S}_+^p$ by a breast cancer data set from [21] with 1213 observations and 52 samples, and apply Algorithm 1 and GPower_{ℓ_0} to the problem (1) with $A = -\Sigma$.

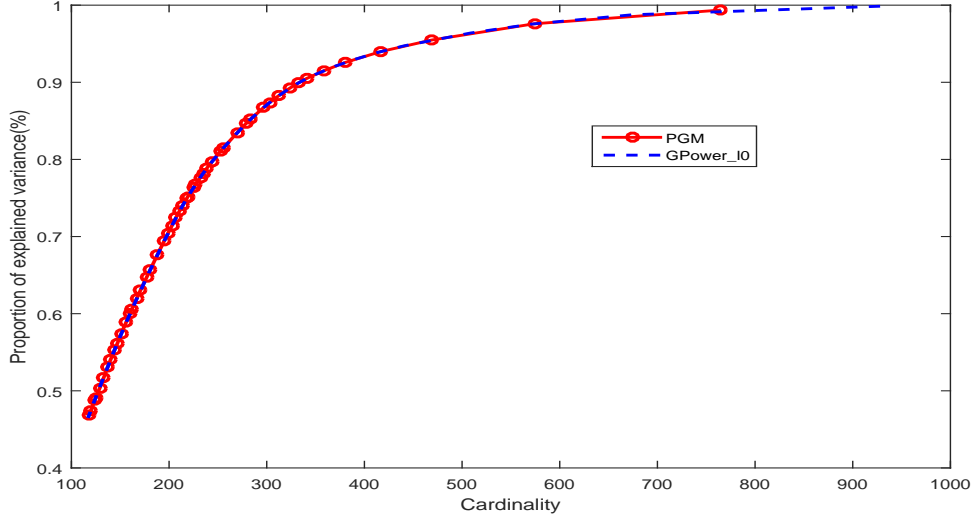


Figure 3: Proportion of explained covariance of the first sparse PC under different sparsity

Figure 5.4 plots the proportion of explained covariance against the sparsity of the first sparse PC, where the x -axis represents the cardinality of the first sparse PC, while y -axis represents the proportion of the variance of the first sparse PC to the variance of the first PC. We see that under the same sparsity, the sparse PCs yielded by Algorithm 1 almost explain the same variance as those yielded by GPower_{ℓ_0} do.

From the previous numerical results, we conclude that Algorithm 1 is comparable with GPower_{ℓ_0} in terms of the recoverability on the sparse PCs and the proportion of explained covariance of the sparse PCs. In addition, the computation work in each iteration of Algorithm 1 is similar to that of GPower_{ℓ_0} . As far as we know, there is still lack of convergence guarantee for the iterate sequence $\{x^k\}$ of GPower_{ℓ_0} , and only a weak convergence result similar to part (ii) of Proposition 4.2 is provided in [10].

6 Conclusions

We have established the KL property of exponent $1/2$ for the zero-norm regularized quadratic function over the unit sphere. To the best of our knowledge, there is no work on the KL property of exponent $1/2$ for such a composite function which involves the sum of two nonconvex nonsmooth functions. By using this crucial property and the finding that the global optimal solutions of (11) are accessible, we also develop a globally and linearly convergent PGM. Numerical comparison with GPower_{ℓ_0} further confirms the obtained theoretical results and the efficiency of the proposed PGM.

References

- [1] Y. HAI LE, *Generalized subdifferentials of the rank function*, Optimization Letters, 7(2013): 731-743.
- [2] F. J. ARAGÓN ARTACHO AND M. H. GEOFFROY, *Characterization of metric regularity of subdifferential*, Journal of Convex Analysis, 15(2008): 365-380.
- [3] H. ATTOUCH, J. BOLTE, P. REDONT AND A. SOUBEYRAN, *Proximal alternating minimization and projection methods for nonconvex problems: an approach based on the Kerdyka-Łojasiewicz inequality*, Mathematics of Operations Research, 35(2010): 438-457.
- [4] H. ATTOUCH, J. BOLTE AND B. F. SVAITER, *Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods*, Mathematical Programming, 137(2013): 91-129.
- [5] D. P. BERTSEKAS, *Nonlinear Programming*, 2nd edition, Athena Scientific, Belmont, Massachusetts, 1999.

- [6] J. BOLTE, S. SABACH AND M. TEBoulLE, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Mathematical Programming, 146(2014): 459-494.
- [7] J. BOLTE, T. P. NGUYEN, J. PEYPOUQUET AND B. W. SUTER, *From error bounds to the complexity of first-order descent methods for convex functions*, Mathematical Programming, 165(2017): 471-507.
- [8] J. JEFFERS, *Two case studies in the application of principal components*, Applied Statistics, 16(3): 225-236, 1967.
- [9] I. T. JOLLIFFE, N. T. TRENDAFILOV AND M. UDDIN, *A modified principal component technique based on the lasso*, Journal of Computational and Graphical Statistics, 12(2003): 531-547.
- [10] M. JOURNÉE, Y. NESTEROV, P. RICHTÁRIK AND R. SEPULCHRE, *Generalized power method for sparse principal component analysis*, Journal of Machine Learning Research, 11(2010): 517-553.
- [11] F. J. LUQUE, *Asymptotic convergence analysis of the proximal point algorithm*, SIAM Journal on Control and Optimization, 22(1984): 277-293.
- [12] L. MACKEY, *Deflation methods for sparse pca*. In Proceedings of the 22nd Annual Conference on Neural Information Processing Systems (NIPS'08), 2008.
- [13] B. MOGHADDAM, Y. WEISS AND S. AVIDAN, *Generalized spectral bounds for sparse lda*, In Proceedings of the 23rd International Conference on Machine Learning (ICML'06), pages 641-648, 2006.
- [14] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, 1970.
- [15] R. T. ROCKAFELLAR AND R. J-B. WETS, *Variational Analysis*, Springer, 1998.
- [16] G. Y. LI AND T. K. PONG, *Calculus of the exponent of Kurdyka-Łojasiewicz inequality and its applications to linear convergence of first-order methods*, Foundations of Computational Mathematics, DOI 10.1007/s10208-017-9366-8.
- [17] S. H. PAN, D. D. ZHANG AND Y. L. LIU, *Metric subregularity of subdifferential and KL Property of exponent 1/2*, submitted to Journal of Optimization Theory and Applications.
- [18] Z. Q. LUO AND P. TSENG, *Error bounds and convergence analysis of matrix splitting algorithms for the affine variational inequality problem*, SIAM Journal on Optimization, 1(1992):43-54.
- [19] H. K. LIU, W. J. WU AND A.M.-C. SO, *QUADRATIC OPTIMIZATION WITH ORTHOGONALITY CONSTRAINTS: EXPLICIT ŁOJASIEWICZ EXPONENT AND LINEAR CONVERGENCE OF LINE-SEARCH METHODS*, in Proceeding of Journal of 46(2016): 1158-1167.

- [20] X. T. YUAN AND T. ZHANG, *Truncated power method for sparse eigenvalue problems*, *Journal of Machine Learning Research*, 14(2013): 899-925.
- [21] L. J. VAN'T VEER, H. DAI, M. J. VAN DE VIJVER, Y. D. HE, A. A. HART AND ET AL., *Gene expression profiling predicts clinical outcome of breast cancer*, *Nature*, 415(2002): 530-536.
- [22] X. F. WANG, J. J. YE, X. M. YUAN, S. Z. ZENG AND J. ZHANG, *Perturbation techniques for convergence analysis of proximal gradient method and other first-order algorithms via variational analysis*, arXiv:1810.10051.
- [23] H. SHEN AND J. Z. HUANG, *Sparse principal component analysis via regularized low rank matrix approximation* *Journal of Multivariate Analysis*, 99(2008): 1015-1034.
- [24] P. TSENG AND S. YUN, *A coordinate gradient descent method for nonsmooth separable minimization*, *Mathematical Programming*, 117(2009): 387-423.
- [25] B. WEN, X. J. CHEN AND T. K. PONG, *Linear convergence of proximal gradient algorithm with extrapolation for a class of nonconvex nonsmooth minimization problems*, *SIAM Journal on Optimization*, 27(2017): 124-145.
- [26] Z. R. ZHOU AND A.M.-C. SO, *A unified approach to error bounds for structured convex optimization problems*, *Mathematical Programming*, 165(2017): 689-728.
- [27] Q. ZHANG, C. H. CHEN, H. K. LIU AND A.M.-C. SO, *A unified approach to error bounds for structured convex optimization problems*, *Mathematical Programming*, 165(2017): 689-728.
- [28] H. ZOU, T. HASTIE AND R. TIBSHIRANI, *Sparse principal component analysis*, *Journal of Computational and Graphical Statistics*, 15(2006): 265-286.

Appendix A

Let $D = \text{diag}(d_1, d_2, \dots, d_p)$ with $d_1 \geq \dots \geq d_p$. Consider the following function

$$\psi(x) := x^\top D x + \delta_S(x) \quad \forall x \in \mathbb{R}^p. \quad (18)$$

In this part, we shall study the KL property of exponent 1/2 of ψ . First, by Lemma 3.1 it is immediate to obtain the following characterization for the critical set of ψ .

Lemma 1 *Let ψ be the function in (18). Then $\text{crit } \psi = \{x \in \mathcal{S} \mid D x = \langle x, D x \rangle x\}$. In particular, for each $x \in \text{crit } \psi$, it holds that $d_i = \langle x, D x \rangle$ with $i \in \text{supp}(x)$.*

Proposition 1 *The function ψ defined by (18) is a KL function of exponent 1/2.*

Proof: For any $z \in \text{dom } \partial\psi$, from equation (5) it immediately follows that

$$\begin{aligned}
\text{dist}(0, \partial\psi(z))^2 &= \min_{u \in \partial\psi(z)} \|u\|^2 = \min_{w \in \mathbb{R}} \|2Dz + wz\|^2 \\
&= \min_{w \in \mathbb{R}} \left\{ 4\langle z, D^\top Dz \rangle + w^2 + 4w\langle z, Dz \rangle \right\} \\
&= 4\langle z, D^\top Dz \rangle - 4(\langle z, Dz \rangle)^2 \\
&= 4\|Dz - \langle z, Dz \rangle z\|^2.
\end{aligned} \tag{19}$$

Now fix an arbitrary $\bar{x} \in \text{crit } \psi$. From Lemma 1 it immediately follows that

$$0 = -D\bar{x} + \langle \bar{x}, D\bar{x} \rangle \bar{x}. \tag{20}$$

We next proceed the arguments by two cases as will be shown below.

Case 1: $d_1 = \dots = d_p = \gamma$ for some $\gamma \in \mathbb{R}$. Choose an arbitrary $\eta > 0$ and an arbitrary $\delta > 0$. Fix an arbitrary $x \in \mathbb{B}(\bar{x}, \delta) \cap [\psi(\bar{x}) < \psi(x) < \psi(\bar{x}) + \eta]$. Clearly, $x \in \mathcal{S}$ and $\langle x, Dx \rangle = \gamma$. Combining equality (20) and equation (19) yields that

$$\text{dist}(0, \partial\psi(x)) = 4\|Dx - \langle x, Dx \rangle x - (D\bar{x} - \langle \bar{x}, D\bar{x} \rangle \bar{x})\| = 0.$$

In addition, $\psi(x) - \psi(\bar{x}) = \langle x, Dx \rangle - \langle \bar{x}, D\bar{x} \rangle = \gamma - \gamma = 0$. This means that

$$\text{dist}(0, \partial\psi(x)) = \sqrt{\psi(x) - \psi(\bar{x})},$$

and consequently, the function ψ has the KL property of exponent 1/2 at \bar{x} .

Case 2: there exist $i \neq j \in \{1, 2, \dots, p\}$ such that $d_i \neq d_j$. Write

$$I := \{i \in \{1, \dots, p\} \mid \bar{x}_i = 0\} \quad \text{and} \quad J := \{i \in \{1, \dots, p\} \mid \bar{x}_i \neq 0\}.$$

By Lemma 1, we know that $d_i = \langle \bar{x}, D\bar{x} \rangle$ for all $i \in J$. This means that there must exist an index $\kappa \in I$ such that $d_\kappa \neq \langle \bar{x}, D\bar{x} \rangle$. Write $I_1 := \{i \in I \mid d_i \neq \langle \bar{x}, D\bar{x} \rangle\}$. By the continuity of the function $\langle \cdot, D\cdot \rangle$, there exists $\delta > 0$ such that for all $z \in \mathbb{B}(\bar{x}, \delta) \cap \mathcal{S}$,

$$\frac{1}{2}|d_j - \langle \bar{x}, D\bar{x} \rangle| \leq |d_j - \langle z, Dz \rangle| \leq \frac{3}{2}|d_j - \langle \bar{x}, D\bar{x} \rangle| \quad \forall j \in I_1. \tag{21}$$

Choose an arbitrary $\eta > 0$. Fix an arbitrary $x \in \mathbb{B}(\bar{x}, \delta) \cap [\psi(\bar{x}) < \psi(x) < \psi(\bar{x}) + \eta]$. Clearly, $x \in \mathcal{S}$. From equation (19), it follows that

$$\begin{aligned}
\frac{1}{4}\text{dist}(0, \partial\psi(x))^2 &= \|Dx - \langle x, Dx \rangle x\|^2 \\
&= \sum_{j \in I} (d_j - \langle x, Dx \rangle)^2 x_j^2 + \sum_{j \in J} (d_j - \langle x, Dx \rangle)^2 x_j^2 \\
&= \sum_{j \in I} (d_j - \langle x, Dx \rangle)^2 x_j^2 + \sum_{j \in J} (\langle \bar{x}, D\bar{x} \rangle - \langle x, Dx \rangle)^2 x_j^2 \\
&\geq \sum_{j \in I_1} (d_j - \langle x, Dx \rangle)^2 x_j^2 \geq \frac{1}{4} \sum_{j \in I_1} (d_j - \langle \bar{x}, D\bar{x} \rangle)^2 x_j^2
\end{aligned} \tag{22}$$

where the third equality is using Lemma 1, the first inequality is by the definition of I_1 , and the last inequality is due to (21). On the other hand, by the definition of ψ ,

$$\begin{aligned}
\psi(x) - \psi(\bar{x}) &= \langle x, Dx \rangle - \langle \bar{x}, D\bar{x} \rangle = \sum_{j \in I} d_j x_j^2 + \sum_{j \in J} d_j x_j^2 - \langle \bar{x}, D\bar{x} \rangle \|x\|^2 \\
&= \sum_{j \in I} (d_j - \langle \bar{x}, D\bar{x} \rangle) x_j^2 + \sum_{j \in J} (d_j - \langle \bar{x}, D\bar{x} \rangle) x_j^2 \\
&= \sum_{j \in I} (d_j - \langle \bar{x}, D\bar{x} \rangle) x_j^2 = \sum_{j \in I_1} (d_j - \langle \bar{x}, D\bar{x} \rangle) x_j^2 \\
&\leq \sum_{j \in I_1} |\langle \bar{x}, D\bar{x} \rangle - d_j| x_j^2 \leq \max_{j \in I_1} |d_j - \langle \bar{x}, D\bar{x} \rangle| \|x_{I_1}\|^2
\end{aligned} \tag{23}$$

where the fourth equality is due to Lemma 1, the fifth one is by the definition of I_1 , and the inequality is since $\psi(x) - \psi(\bar{x}) > 0$. From the above inequalities (22) and (23),

$$\begin{aligned}
\text{dist}(0, \partial\psi(x)) &\geq \sqrt{\sum_{j \in I_1} (d_j - \langle \bar{x}, D\bar{x} \rangle)^2 x_j^2} \geq \min_{j \in I_1} |d_j - \langle \bar{x}, D\bar{x} \rangle| \|x_{I_1}\| \\
&\geq \frac{\min_{j \in I_1} |d_j - \langle \bar{x}, D\bar{x} \rangle|}{\sqrt{\max_{j \in I_1} |d_j - \langle \bar{x}, D\bar{x} \rangle|}} \sqrt{\psi(x) - \psi(\bar{x})}.
\end{aligned}$$

That is, there exists a constant c (only dependent on \bar{x}) such that

$$\text{dist}(0, \partial\psi(x)) \geq c \sqrt{\psi(x) - \psi(\bar{x})}.$$

By the arbitrariness of x , the function ψ has the KL property with exponent $1/2$ at \bar{x} . From the arbitrariness of \bar{x} in $\text{crit}\psi$, ψ is a KL function of exponent $1/2$. \square

Appendix B

In this part, we consider the following problems where κ is a given positive integer:

$$\left\{ \begin{array}{l} \min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - z\|^2 : \|x\| = 1, \|x\|_0 \leq \kappa \right\}, \\ \min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - |z|^{\downarrow}\|^2 : \|x\| = 1, \|x\|_0 \leq \kappa \right\}. \end{array} \right. \tag{24a}$$

$$\left\{ \begin{array}{l} \min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - z\|^2 : \|x\| = 1, \|x\|_0 \leq \kappa \right\}, \\ \min_{x \in \mathbb{R}^p} \left\{ \frac{1}{2} \|x - |z|^{\downarrow}\|^2 : \|x\| = 1, \|x\|_0 \leq \kappa \right\}. \end{array} \right. \tag{24b}$$

By the definition of the zero-norm, it is not difficult to obtain the following result.

Lemma 2 *Let $z \in \mathbb{R}^p$ be a given vector, and let Q be a $p \times p$ signed permutation matrix such that $|z|^{\downarrow} = Qz$. Then, the following assertions hold.*

- (i) *If x^* is a global optimal solution of (24b), then $Q^{\mathbb{T}}x^*$ is globally optimal to (24a); conversely, if x^* is globally optimal to (24a), then Qx^* is globally optimal to (24b).*
- (ii) *The vector $x^* = \frac{1}{\| |z|^{\kappa, \downarrow} \|} (|z|^{\kappa, \downarrow}; \mathbf{0}_{p-\kappa})$ is a global optimal solution of (24b), and consequently $\frac{1}{\| |z|^{\kappa, \downarrow} \|} Q^{\mathbb{T}} (|z|^{\kappa, \downarrow}; \mathbf{0}_{p-\kappa})$ is globally optimal to the problem (24a).*