

The Proxy Step-size Technique for Regularized Optimization on the Sphere Manifold

Fang Bai, *Member, IEEE*, Adrien Bartoli

Abstract—We give an effective solution to the regularized optimization problem $g(\mathbf{x}) + h(\mathbf{x})$, where \mathbf{x} is constrained on the unit sphere $\|\mathbf{x}\|_2 = 1$. Here $g(\cdot)$ is a smooth cost with Lipschitz continuous gradient within the unit ball $\{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$ whereas $h(\cdot)$ is typically non-smooth but convex and absolutely homogeneous, *e.g.*, norm regularizers and their combinations. Our solution is based on the Riemannian proximal gradient, using an idea we call *proxy step-size* – a scalar variable which we prove is monotone with respect to the actual step-size within an interval. The proxy step-size exists ubiquitously for convex and absolutely homogeneous $h(\cdot)$, and decides the actual step-size and the tangent update in closed-form, thus the complete proximal gradient iteration. Based on these insights, we design a Riemannian proximal gradient method using the proxy step-size. We prove that our method converges to a critical point, guided by a line-search technique based on the $g(\cdot)$ cost only. The proposed method can be implemented in a couple of lines of code. We show its usefulness by applying nuclear norm, ℓ_1 norm, and nuclear-spectral norm regularization to three classical computer vision problems. The improvements are consistent and backed by numerical experiments.

Index Terms—Proxy step-size, Riemannian proximal gradient, non-smooth optimization, regularization, computer vision

1 INTRODUCTION

WE start with optimization problems on the unit sphere, *i.e.*, the sphere manifold:

$$\arg \min_{\mathbf{x} \in \mathbb{R}^n} g(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1. \quad (1)$$

The sphere constraint $\|\mathbf{x}\|_2 = 1$ has been widely used to model scale invariant mathematical structures, *e.g.*, the fundamental matrix [1] and the dual absolute quadric [2] in geometric vision (see [3] for more such applications). In statistics, the likelihood is defined up to scale [4], thus often results in problems in the form (1), *e.g.*, the spectral correspondence association [5] and the wavelet density estimation [6]. Some other well-known applications of problem (1) include the p -harmonic energy minimization [7] and discretized Bose-Einstein condensates [8].

Other than being constrained on the sphere, some applications require \mathbf{x} to possess additional structures, *e.g.*, low-rank (by reorganizing the elements of \mathbf{x} in matrix form) or sparsity (*i.e.*, number of nonzeros), to favor its physical/geometric meaning. These additional properties can be often enforced by a dedicated regularization term $h(\mathbf{x})$, leading to the following regularized optimization:

$$\arg \min_{\mathbf{x} \in \mathbb{R}^n} g(\mathbf{x}) + h(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1. \quad (2)$$

In general, the regularizer $h(\cdot)$ is non-smooth but convex and absolutely homogeneous. Typically $h(\cdot)$ are norm

functions or their combinations, in particular, ℓ_1 norm for sparsity and nuclear norm for low-rank [9].

Problem (2) is difficult, because of the entangling of the non-smooth cost $h(\cdot)$ and the non-convex manifold constraint $\|\mathbf{x}\|_2 = 1$. This inhibits the direct applicability of well-studied classical methods, *i.e.*, the Euclidean optimization techniques for non-smooth composite costs [10]–[14], and the Riemannian optimization techniques for smooth costs [8], [15], [16]. In the literature, researchers have explored several ideas to solve non-smooth optimization problems with non-convex constraints, *e.g.*, Riemannian subgradient methods [17]–[20], proximal point methods [21]–[23], operator-splitting methods [24]–[28], and more recently Riemannian proximal gradient methods [29], [30] (see Section 2 for a short review of these methods). Among them, the Riemannian proximal gradient methods show significant advantages over other methods, in terms of both convergence guarantees and convergence speed [29].

The Riemannian proximal gradient is quite a recent research topic, mainly due to Chen *et al.*'s [29] and Huang *et al.*'s [30] work on the Stiefel manifold. These methods [29], [30] are exact methods with convergence proofs to a critical point, while the others either lack such proofs or only have proofs for special forms of $h(\cdot)$ [31]. At its core, [29], [30] solve a non-smooth equation derived from the Karush–Kuhn–Tucker (KKT) system by an iterative semi-smooth Newton method (SSNM) [32]. This approach requires the generalized Clarke differential, which is difficult to obtain and only applicable to special forms, *e.g.*, the ℓ_1 norm regularization used in [29], [30], [32]. Besides, the use of the generalized Clarke differential contradicts the spirit of proximal gradient, as the proximal is intended to avoid the computation of sub-gradients. It should be noted that the unit sphere is a special case of the Stiefel manifold, thus the result in [29], [30] applies to problem (2). However, due to the usage of the generalized Clarke differential, it is hard

• Both the authors are with the ENCOV, IGT, Institut Pascal, Université Clermont Auvergne, CHU Clermont-Ferrand, France.
E-mail: fang.bai@yahoo.com; adrien.bartoli@gmail.com
Corresponding author: Fang Bai
The work was supported by the ANR project TOPACS and the CLARA project AIALO.

to implement $h(\cdot)$ even as the nuclear norm regularization, not to mention more advanced $h(\cdot)$.

In this work, we advance the proximal gradient method on the sphere manifold (PGS) by discovering a concept we call *proxy step-size*, for convex and absolute homogeneous $h(\cdot)$. Importantly, we prove that the proxy step-size is monotone with respect to the actual step-size in the working region and define one line-search iteration in closed-form. Thus the generalized Clarke differential is never required. Based on this novel insight, we control the optimization flow using the proxy step-size, and establish the convergence proof to a critical point (*i.e.*, a solution that satisfies the first-order optimality condition). Our final outputs are three PGS algorithms (of which two are accelerated algorithms with the Nesterov momentum technique) that retain the elegance of classical Euclidean proximal gradient methods. Our method is easy to implement and much faster than the SSNM based methods [29], [30]. Our main contributions are highlighted as follows:

- **Section 4.1.** We reveal the existence of the proxy step-size by exploiting the convexity and absolute homogeneity of the non-smooth cost $h(\cdot)$, and show how it decides both the actual step-size and the tangent update in closed-form.
- **Section 4.2.** We establish the monotonicity between the proxy step-size and the actual step-size, allowing one to control the actual step-size using the proxy step-size monotonically.
- **Section 4.3.** We establish the convergence proof to a critical point, using a line-search from the $g(\cdot)$ cost only, by mildly assuming $g(\cdot)$ has Lipschitz continuous gradient within the unit ball $\{\mathbf{x} : \|\mathbf{x}\| \leq 1\}$.
- **Section 5.** We present accelerated versions of our PGS algorithm using the Nesterov momentum technique [14], [33]. We empirically show that the accelerated algorithms converge much faster.
- **Section 6.** We demonstrate our algorithms with three applications, by applying nuclear norm, ℓ_1 norm and nuclear-spectral norm regularization to three well-know computer vision problems.

We start with related work in Section 2, and necessary background in Section 3. Then we formally present our proxy step-size technique in Section 4, and the accelerated version in Section 5. The three example applications and experimental results are given in Section 6 and Section 7, respectively. Section 8 concludes the paper.

2 RELATED WORK

Both optimizing a smooth cost on the manifold [8], [15], [16] and proximal gradient for non-smooth optimization in the Euclidean space [10], [11], [11], [12], [14] have been well studied in the literature. However, there exist only a few methods for solving non-smooth cost functions on the manifold. We review existing techniques in this regard.

2.1 Riemannian Subgradient Methods

Subgradient methods require one to evaluate the descent direction of the total cost directly. The descent direction, in the non-smooth setting, is characterized by the notion

of generalized Clarke gradient, which is usually difficult to calculate even numerically in practice, *e.g.*, see [34], [35] for certain types of functions. Instead, researchers seek for approximations of the subgradient. A key concept in this regard is the ϵ -subgradient [36]. Grohs *et al.* proposed two Riemannian ϵ -subgradient methods based on the line-search [17] and trust-region techniques [18], with convergence guarantee to a critical point. Hosseini *et al.* [20] generalized the idea to the ϵ -subgradient-oriented descent sequence which combines the idea of the BFGS algorithm. Hosseini *et al.* [19] gave a non-smooth Riemannian gradient sampling method with convergence analysis. Despite the hassle to handle the subgradient, the subgradient method is shown to have a slow convergence rate $O(1/\sqrt{k})$ [23], [37].

2.2 Proximal Point Methods

Ferreira *et al.* [21] proposed the proximal point algorithm on the Riemannian manifold, and Bento *et al.* [23] established the $O(1/k)$ convergence rate of the algorithm on the Hadamard manifold for convex cost functions. However since every smooth function that is geodesically convex on a compact Riemannian manifold is a constant [38], the analysis in [21], [23] does not apply to compact Riemannian manifolds (*e.g.*, the Stiefel manifold and the unit sphere). Bento *et al.* [22] gave a convergence analysis on the general Riemannian setting by assuming the cost function satisfies the Kurdyka–Lojasiewicz inequality. In terms of computation, the Riemannian proximal point algorithm [21] requires one to solve a subproblem to which an efficient solution does not exist within the current research. As a result, this line of research is largely restricted to theoretical interests at the moment.

2.3 Operator-splitting Methods

The hardness of problem (2) is caused by the composition of a non-smooth cost function and the non-convex manifold constraint. In the convex setting, the cost function and the constraint can be handled separately by the alternating direction methods of multipliers (ADMM) [39]. We recapitulate the major advancements of this technique in the non-smooth and non-convex setting. Lai *et al.* [24] explored the splitting of orthogonality constraints (SOC) method, which handles the orthogonality constraint and the cost function separately. Kovnatsky *et al.* [25] proposed the manifold ADMM (MADMM) method, which further exploits the composite structure of the smooth and non-smooth cost functions. However, both these methods, SOC and MADMM, lack convergence proofs in their original paper. A deeper exploitation in terms of the convergence study has been conducted by Wang *et al.* [28]. For special forms of $h(\cdot)$, the convergence of ADMM to a stationary point can be established, *e.g.*, see the stabilized ADMM (SADMM) for ℓ_1 norm regularization on the Stiefel manifold [31]. More recently, Chen *et al.* [26] proposed PAMAL, the proximal alternating minimized augmented Lagrangian method. The PAMAL method enjoys the sub-sequence convergence property, and is noticeably faster than the SOC method by the experiments in [26]. Another variant, termed EPALAML, was proposed by Zhu *et al.* [27], based on the proximal alternating linearized minimization (PALM) method. Both

PAMAL [26] and EPALAML [27] minimize the augmented Lagrangian function approximately with different methods.

2.4 Proximal Gradient Methods

The development of proximal gradient on the Riemannian manifold is a rather new topic. It started with the landmark paper from Chen *et al.* [29] which developed a proximal gradient on the Stiefel manifold with convergence guarantee to a critical point. Experiments in [29] show that the Riemannian proximal gradient method is more efficient than operator-splitting methods such as SOC and PAMAL. More recently, Huang *et al.* [30] proposed another formulation, by using retractions on the non-smooth cost as well to define an iteration. In [30], a convergence rate analysis was also given based on this adaptation, which is $O(1/k)$ for the case without acceleration. In addition to convergence guarantee to a critical point, Riemannian proximal gradient methods allow the possibility to design accelerated algorithms to obtain even faster convergence, see [30], [40].

However, the subproblem of each iteration in [29], [30] is solved by a SSNM which is expensive. In addition, the SSNM method requires the generalized Clarke differential which is difficult for advanced $h(\cdot)$ [41] and contradicts the spirit of proximal gradient (which aims to avoid the generalized Clarke differential).

In this work, we propose the proxy step-size technique for proximal gradient on the sphere manifold by exploiting the convexity and absolute homogeneity of $h(\cdot)$. Our method does not require the generalized Clarke differential, thus is more applicable to advanced $h(\cdot)$. In addition, our method is much faster thanks to the closed-form evaluation.

3 PRELIMINARIES

3.1 Absolute Homogeneous Function

A function $h(\cdot)$ is said to be *absolutely homogeneous* if $h(\alpha\mathbf{x}) = |\alpha|h(\mathbf{x})$ for any scalar α and vector \mathbf{x} .

Lemma 1. *If $h(\cdot)$ is convex and absolutely homogeneous, then:*

- $h(\mathbf{0}) = 0$;
- $h(\cdot)$ is even, i.e., $h(-\mathbf{x}) = h(\mathbf{x})$ for any \mathbf{x} .
- $h(\cdot)$ is non-negative, i.e., $h(\mathbf{x}) \geq 0$ for any \mathbf{x} .

Proof. The first two are obtained by setting $\alpha = 0$ and $\alpha = -1$ respectively in $h(\alpha\mathbf{x}) = |\alpha|h(\mathbf{x})$. The third is true because if $h(\cdot)$ is further convex:

$$\frac{h(\mathbf{x}) + h(-\mathbf{x})}{2} \geq h\left(\frac{\mathbf{x} + (-\mathbf{x})}{2}\right) = h(\mathbf{0}).$$

The proof is immediate by $h(-\mathbf{x}) = h(\mathbf{x})$ and $h(\mathbf{0}) = 0$. \square

3.2 Proximal Operator

For a convex but possibly non-smooth function $h(\cdot)$, the proximal operator evaluates the solution of the following convex optimization problem at a given point \mathbf{w} :

$$\begin{aligned} \text{prox}_{th}(\mathbf{w}) &= \arg \min_{\mathbf{x}} \left\{ th(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ h(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \mathbf{w}\|_2^2 \right\} \stackrel{\text{def}}{=} \mathbf{z}, \end{aligned} \quad (3)$$

where $t \geq 0$ is a given scalar. For $t = 0$, $\mathbf{z} = \mathbf{w}$.

Since the proximal is a convex problem, its solution is characterized by its first-order necessary condition:

$$\begin{aligned} \mathbf{z} = \text{prox}_{th}(\mathbf{w}) &\Leftrightarrow \mathbf{0} \in \partial h|_{\mathbf{z}} + \frac{1}{t}(\mathbf{z} - \mathbf{w}) \\ &\Leftrightarrow \mathbf{w} \in t\partial h|_{\mathbf{z}} + \mathbf{z}. \end{aligned} \quad (4)$$

The proximal satisfies firm non-expansiveness and non-expansiveness (see Appendix A). Additionally, we show the following statements hold true.

Lemma 2. *If $h(\cdot)$ is convex and absolutely homogeneous and $t \geq 0$, we have:*

- $\text{prox}_{th}(\mathbf{0}) = \mathbf{0}$;
- $\|\text{prox}_{th}(\mathbf{w})\|_2^2 \leq \langle \text{prox}_{th}(\mathbf{w}), \mathbf{w} \rangle$;
- $\|\text{prox}_{th}(\mathbf{w})\|_2 \leq \|\mathbf{w}\|_2$

Proof. See Appendix A. \square

3.3 Proximal Gradient in the Euclidean Space

Let $g(\cdot)$ be smooth, and $h(\cdot)$ be convex. A proximal gradient step in the Euclidean space is defined as:

$$\begin{aligned} \mathbf{v}_k &= \arg \min_{\mathbf{v} \in \mathbb{R}^n} \langle \nabla g|_{\mathbf{x}_k}, \mathbf{v} \rangle + \frac{1}{2t} \langle \mathbf{v}, \mathbf{v} \rangle + h(\mathbf{x}_k + \mathbf{v}) \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \mathbf{v}_k, \end{aligned} \quad (6)$$

where $t \geq 0$. Problem (6) is convex, whose solution is characterized by its first-order necessary condition:

$$\begin{aligned} \mathbf{0} &\in \nabla g|_{\mathbf{x}_k} + \frac{1}{t} \mathbf{v}_k + \partial h|_{\mathbf{x}_k + \mathbf{v}_k} \\ \mathbf{x}_k - t\nabla g|_{\mathbf{x}_k} &\in (\mathbf{x}_k + \mathbf{v}_k) + t\partial h|_{\mathbf{x}_k + \mathbf{v}_k}. \end{aligned}$$

Therefore, from equation (5), we have:

$$\mathbf{x}_k + \mathbf{v}_k = \text{prox}_{th}(\mathbf{x}_k - t\nabla g|_{\mathbf{x}_k}) = \mathbf{x}_{k+1}.$$

The update step differs from the classical gradient step for minimizing the $g(\cdot)$ cost only by a proximal operation, thus is termed as *proximal gradient*. Here t works as the *step-size* in the standard gradient descent methods.

3.4 Sphere Manifold

The unit sphere, or the sphere manifold, is an embedded manifold:

$$\mathcal{S} = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 = 1\}. \quad (8)$$

The tangent space at a point $\mathbf{x} \in \mathcal{S}$ is:

$$\mathcal{T}_{\mathbf{x}}\mathcal{S} = \{\mathbf{v} \in \mathbb{R}^n : \mathbf{x}^\top \mathbf{v} = 0\}. \quad (9)$$

Let $g : \mathcal{S} \rightarrow \mathbb{R}$ be a function defined on the manifold. The Riemannian gradient of $g(\cdot)$ at $\mathbf{x} \in \mathcal{S}$, denoted by $\text{grad } g|_{\mathbf{x}}$, is the unique tangent vector satisfying:

$$Dg(\mathbf{x})[\mathbf{v}] = \langle \text{grad } g|_{\mathbf{x}}, \mathbf{v} \rangle_{\mathbf{x}}, \quad \forall \mathbf{v} \in \mathcal{T}_{\mathbf{x}}\mathcal{S},$$

where $Dg(\mathbf{x})[\mathbf{v}]$ is the directional derivative of $g(\cdot)$ along the direction of \mathbf{v} . We shall use the induced Euclidean metric as the Riemannian metric, which means $\langle \mathbf{v}, \mathbf{v} \rangle_{\mathbf{x}} = \langle \mathbf{v}, \mathbf{v} \rangle = \mathbf{v}^\top \mathbf{v}$. Practically, the Riemannian gradient $\text{grad } g|_{\mathbf{x}}$ can be obtained by orthogonally projecting the Euclidean gradient

$\nabla g|_{\mathbf{x}}$ in \mathbb{R}^n into the tangent space at \mathbf{x} . Using the so-called orthogonal projector $\text{proj}_{\mathcal{T}_{\mathbf{x}}\mathcal{S}}(\cdot)$, we can write:

$$\text{grad } g|_{\mathbf{x}} = \text{proj}_{\mathcal{T}_{\mathbf{x}}\mathcal{S}} \nabla g|_{\mathbf{x}} = (\mathbf{I} - \mathbf{x}\mathbf{x}^\top) \nabla g|_{\mathbf{x}} \quad (10)$$

$$= \nabla g|_{\mathbf{x}} - \langle \mathbf{x}, \nabla g|_{\mathbf{x}} \rangle \mathbf{x}. \quad (11)$$

We shall use the following *retraction* to bring an increment in the tangent space back to the manifold:

$$\mathcal{R}_{\mathbf{x}}(\mathbf{v}) \stackrel{\text{def}}{=} \frac{\mathbf{x} + \mathbf{v}}{\|\mathbf{x} + \mathbf{v}\|_2} : \mathcal{T}_{\mathbf{x}}\mathcal{S} \mapsto \mathcal{S}. \quad (12)$$

For more details of these concepts, we refer to [15], [16].

3.5 Proximal Gradient on the Sphere Manifold

Formulation in the tangent space. Inspired by the proximal gradient in the Euclidean space, researchers have tried to formulate the update vector \mathbf{v}_k in the tangent space of \mathbf{x}_k as an extension to the manifold setting [29], [30]. In this work, we propose the following:

$$\mathbf{v}_k = \arg \min_{\mathbf{v} \in \mathcal{T}_{\mathbf{x}_k}\mathcal{S}} \langle \text{grad } g|_{\mathbf{x}_k}, \mathbf{v} \rangle + \frac{1}{2t} \langle \mathbf{v}, \mathbf{v} \rangle + h(\mathbf{x}_k + \mathbf{v}) \quad (13)$$

$$\mathbf{x}_{k+1} = \mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k). \quad (14)$$

The subproblem (13) was first proposed by Chen *et al.* [29] in their work on the Stiefel manifold. However in Chen *et al.*'s work, the update equation (14) is $\mathbf{x}_{k+1} = \mathcal{R}_{\mathbf{x}_k}(\alpha_k \mathbf{v}_k)$ with α_k acting as another step-size to control the length of \mathbf{v}_k similar to Riemannian subgradient methods. We shall see that this α_k is not required, as in accordance with equation (7) in the Euclidean proximal gradient.

The solution from the KKT system. Problem (13) is convex, thus its solution is uniquely characterized by the KKT system. We take the tangent space constraint explicitly, and write the Lagrange function as:

$$\mathcal{L}(\mathbf{v}, \mu) = \langle \text{grad } g|_{\mathbf{x}_k}, \mathbf{v} \rangle + \frac{1}{2t} \langle \mathbf{v}, \mathbf{v} \rangle + h(\mathbf{x}_k + \mathbf{v}) + \mu \mathbf{x}_k^\top \mathbf{v}.$$

The KKT system is $\mathbf{0} \in \partial \mathcal{L}_{\mathbf{v}}, \mathbf{x}_k^\top \mathbf{v}_k = 0$. With some trivial calculations, the KKT system is reduced to the following (see Appendix B for details):

$$\begin{cases} \mathbf{v}_k = \text{prox}_{th}((1 - \mu t) \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k}) - \mathbf{x}_k & (15a) \\ \mathbf{x}_k^\top \text{prox}_{th}((1 - \mu t) \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k}) = 1. & (15b) \end{cases}$$

This KKT system can be solved by the SSNM method [32], by solving the non-smooth equation (15b) first to obtain the Lagrange multiplier μ and then computing the tangent update \mathbf{v}_k by equation (15a). This approach has been discussed in [32] and used by Chen *et al.* in [29].

Limitations of the existing solution. While subproblem (13) seems a trivial extension from the Euclidean case (6), the existing solution is not as elegant as its Euclidean counterpart. First, to apply the SSNM method, as $h(\cdot)$ is non-smooth, the generalized Clarke differential of $h(\cdot)$ is required. This contradicts the spirit of proximal gradient, as the proximal operator is typically used to avoid the generalized differential. Second, the generalized Clarke differential is usually difficult to obtain and only applicable for special forms, *e.g.*, the ℓ_1 norm used in [29], [30]. If $h(\cdot)$ is the nuclear norm or more advanced functions, the SSNM is hard

to implement. Third, solving an inner loop by an iterative method (like SSNM) can degenerate the numerical accuracy and even the overall convergence.

In this work, we propose proxy step-size, an effective technique to handle subproblem (13). Instead of solving (13) directly, we aim to generate valid solutions to problem (13) in closed-form controlled by the proxy step-size. With this new technique, the generalized differential is never required, and we show that proximal gradient on the sphere can be elegantly formulated as a proximal gradient step with respect to the proxy step-size followed by normalizations.

One iteration in Chen *et al.* [29]. After solving \mathbf{v}_k , a line-search process is used to ensure the descent of the total cost $f(\cdot) = g(\cdot) + h(\cdot)$. In general, they propose:

- given t , solve \mathbf{v}_k using the SSNM method;
- set $\alpha_k \leftarrow 1$ and shrink α_k until the following line-search criterion is met:

$$f(\mathcal{R}_{\mathbf{x}_k}(\alpha \mathbf{v}_k)) \leq f(\mathbf{x}_k) - \frac{\alpha_k}{2t} \langle \mathbf{v}_k, \mathbf{v}_k \rangle; \quad (16)$$

- set $\mathbf{x}_{k+1} \leftarrow \mathcal{R}_{\mathbf{x}_k}(\alpha_k \mathbf{v}_k)$.

The validity of the line-search (16) is proved in [29].

Convex toolbox for problem (13). Problem (13) is convex, thus a naive idea is to simply call a convex toolbox. However, the computation is prohibitive for practical usage for high dimensional \mathbf{x} . Needless to say this is only one iteration, which we want to solve as efficiently as possible. Mathematically, this solution is not elegant.

4 THE PROXY STEP-SIZE TECHNIQUE

We now formally present our proxy step-size technique to solve one proximal gradient iteration, by assuming $h(\cdot)$ to be absolutely homogeneous. We term the proposed algorithm based on the proxy step-size technique as PGS (short for Proximal Gradient on the Sphere manifold).

4.1 Proxy Step-size

Lemma 3. *If $h(\cdot)$ is convex and absolutely homogeneous and $t \geq 0$, then $\text{prox}_{th}(\alpha \mathbf{w}) = \alpha \text{prox}_{\frac{t}{|\alpha|}h}(\mathbf{w})$ for any scalar α .*

Proof. See Appendix C. □

Now we introduce $t' = \frac{t}{1 - \mu t}$, which we term *proxy step-size*. We show that the update equations from \mathbf{x}_k to \mathbf{x}_{k+1} are completely determined by the proxy step-size t' . To that end, we rewrite equation (15a) as:

$$\begin{aligned} \mathbf{x}_k + \mathbf{v}_k &= \text{prox}_{th}((1 - \mu t) \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k}) \\ &= \text{prox}_{th} \left((1 - \mu t) \left(\mathbf{x}_k - \frac{t}{1 - \mu t} \text{grad } g|_{\mathbf{x}_k} \right) \right) \\ &= (1 - \mu t) \text{prox}_{\frac{t}{|1 - \mu t|}h} \left(\mathbf{x}_k - \frac{t}{1 - \mu t} \text{grad } g|_{\mathbf{x}_k} \right) \\ &= \frac{t}{t'} \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k}), \end{aligned}$$

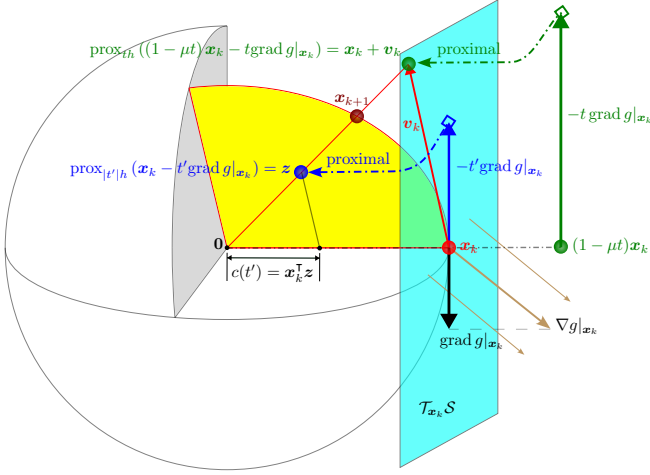


Fig. 1: The proxy step-size technique. In the KKT system (15), given t , the essence of the non-smooth equation (15b) is to decide a proper Lagrange multiplier μ that lands the point $\text{prox}_{t'h}((1 - \mu)t \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k})$ into the tangent plane at \mathbf{x}_k . Such a process is difficult as it is hard to solve equation (15b) for advanced $h(\cdot)$. Instead, we propose to use the proxy step-size t' to generate valid solutions to the KKT system (15). In specific, we first move in the tangent plane at \mathbf{x}_k by $-t' \text{grad } g|_{\mathbf{x}_k}$, and then apply the proximal to reach the point $\mathbf{z} = \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k})$. We see the point $\text{prox}_{t'h}((1 - \mu)t \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k})$ is the intersection of the line $(\mathbf{0}, \mathbf{z})$ and the tangent plane at \mathbf{x}_k , given as $\frac{1}{\mathbf{x}_k^\top \mathbf{z}} \mathbf{z}$.

where the third equality stems from Lemma 3. By $\mathbf{x}_k^\top \mathbf{v}_k = 0$, we obtain $t = \frac{t'}{\mathbf{x}_k^\top \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k})}$. Therefore the KKT system (15) can be rewritten as:

$$\begin{cases} \mathbf{v}_k = \frac{\text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k})}{\mathbf{x}_k^\top \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k})} - \mathbf{x}_k & (17a) \\ t = \frac{t'}{\mathbf{x}_k^\top \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k})} \stackrel{\text{def}}{=} \phi(t'). & (17b) \end{cases}$$

The new KKT system (17) can be considered as a reparameterization of the previous KKT system (15), using t and t' . However, in the new KKT system, both \mathbf{v}_k and t are completely decided by the proxy step-size t' .

Therefore, problem (13) can be solved in closed-form with respect to a given proxy step-size t' as follow:

$$\begin{cases} \mathbf{z} = \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k}) \\ t = \frac{1}{\mathbf{x}_k^\top \mathbf{z}} t' \\ \mathbf{v}_k = \frac{1}{\mathbf{x}_k^\top \mathbf{z}} \mathbf{z} - \mathbf{x}_k. \end{cases} \quad (18)$$

Note that t and \mathbf{v}_k computed from t' satisfy the KKT system (15), thus they are optimal for problem (13). An illustration of the proxy step-size technique is given in Fig. 1.

To design iterations based on the proxy step-size entirely, we need to reveal the relation between the proxy step-size t' and the actual step-size t , and design a line-search process to govern the convergence.

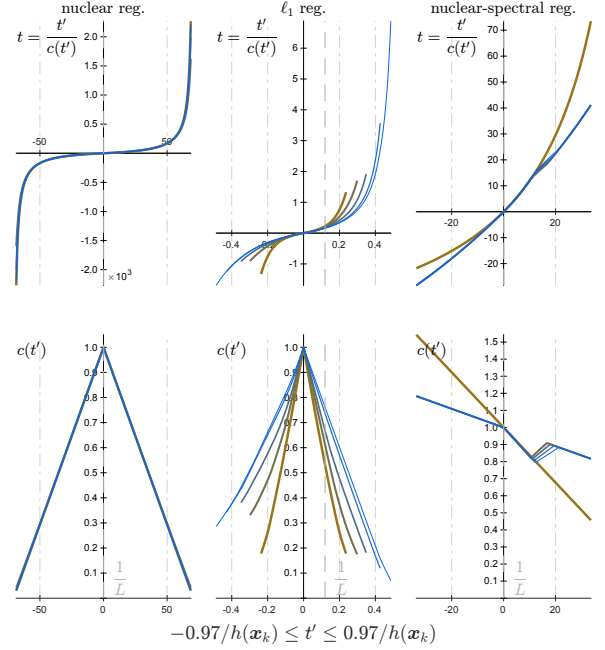


Fig. 2: Numerical examples for Proposition 2 and Theorem 1. If $h(\cdot)$ is convex and absolutely homogeneous, then for $|t'| < 1/h(\mathbf{x}_k)$, we have $c(t') = \mathbf{x}_k^\top \mathbf{z} > 0$ and thus the mapping $t = \phi(t') = t'/c(t')$ is monotonically increasing within this interval. The reference proxy step-size $1/L$ to be described in Section 4.4.1 is plotted as vertical dashed lines. These examples are obtained from the first 5 iterations of: 1) fundamental matrix estimation with nuclear norm regularization, 2) correspondence association with ℓ_1 norm regularization, and 3) self-calibration with nuclear-spectral norm regularization to be presented in Section 6.

4.2 The Mapping between Proxy Step-size and Actual Step-size

We denote $t = \phi(t') = t'/c(t')$ the mapping between t' and t defined by equation (17b), where we introduce $c(t') = \mathbf{x}_k^\top \mathbf{z} = \mathbf{x}_k^\top \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k})$.

Lemma 4. Given arbitrary $t'_1 \neq 0, t'_2 \neq 0, t'_1 \neq t'_2$, it can be shown that:

$$\epsilon(t'_1, t'_2) = \left(\frac{1}{\phi(t'_1)} - \frac{1}{\phi(t'_2)} \right) \left(\frac{1}{t'_1} - \frac{1}{t'_2} \right) \geq 0, \quad (19)$$

which means $1/\phi(t')$ is monotonically increasing with respect to $1/t'$.

Proof. See Appendix D. \square

Lemma 4 states the monotonicity between $1/t'$ and $1/t$. To establish the monotonicity between t' and t explicitly, it suffices to identify an interval where $c(t') > 0$.

Proposition 1. If $c(t') > 0$ for $t' \in [l, u]$, then $\phi(t')$ is monotonically increasing with respect to t' for $t' \in [l, u]$.

Proof. Note that $\phi(t') = t'/c(t')$. We thus have:

$$t'_1 t'_2 \phi(t'_1) \phi(t'_2) = (t'_1 t'_2)^2 / (c(t'_1) c(t'_2)) > 0 \Leftrightarrow c(t'_1) c(t'_2) > 0$$

Therefore if $c(t'_1) c(t'_2) > 0$, inequality (19) is equivalent to:

$$(\phi(t'_1) - \phi(t'_2))(t'_1 - t'_2) \geq 0. \quad (20)$$

We see $c(t') > 0$ is sufficient for $c(t'_1)c(t'_2) > 0$. \square

Now we characterize an interval in which $c(t') > 0$. For this purpose, we first prove the following global inequality for convex and absolutely homogeneous $h(\cdot)$:

Lemma 5. *Let $h(\cdot)$ be convex and absolutely homogeneous. Then for any t, \mathbf{x} and \mathbf{w} , we have:*

$$\langle \mathbf{w} - \text{prox}_{|t|h}(\mathbf{w}), \mathbf{x} \rangle \leq |t| h(\mathbf{x}). \quad (21)$$

The inequality is tight for $\mathbf{x} = \text{prox}_{|t|h}(\mathbf{w})$.

Proof. See Appendix E. \square

Proposition 2. *If $h(\cdot)$ is convex and absolutely homogeneous, then $c(t') \geq 1 - |t'| h(\mathbf{x}_k)$.*

Proof. In Lemma 5, we set $\mathbf{w} = \mathbf{x}_k - t' \text{grad} g|_{\mathbf{x}_k}$ and $\mathbf{x} = \mathbf{x}_k$, which yields:

$$\begin{aligned} \langle \mathbf{x}_k - t' \text{grad} g|_{\mathbf{x}_k} - \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad} g|_{\mathbf{x}_k}), \mathbf{x}_k \rangle \\ \leq |t'| h(\mathbf{x}_k). \end{aligned} \quad (22)$$

We reorganize inequality (22) to complete the proof:

$$c(t') = \mathbf{x}_k^\top \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad} g|_{\mathbf{x}_k}) \geq 1 - |t'| h(\mathbf{x}_k),$$

where we use $\langle \mathbf{x}_k, \mathbf{x}_k \rangle = 1$ and $\langle \text{grad} g|_{\mathbf{x}_k}, \mathbf{x}_k \rangle = 0$. \square

Theorem 1. *If $h(\cdot)$ is convex and absolutely homogeneous, then the mapping $t = \phi(t')$ from proxy step-size t' to actual step-size t is monotonically increasing for $|t'| < 1/h(\mathbf{x}_k)$, and $\phi(0) = 0$.*

Proof. By Proposition 2, we see $c(t') > 0$ for $|t'| < 1/h(\mathbf{x}_k)$. Then by Proposition 1, $\phi(t')$ is monotonically increasing with respect to t' for $|t'| < 1/h(\mathbf{x}_k)$. By definition of the proximal operator in equation (3), we observe:

$$\lim_{t' \rightarrow 0} \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad} g|_{\mathbf{x}_k}) = \mathbf{x}_k, \quad (23)$$

from which we conclude $c(0) = 1$, thus $\phi(0) = 0$. \square

Numerical examples to these results are given in Fig. 2. In practice, we require $t > 0$, thus we use $0 < t' < 1/h(\mathbf{x}_k)$. Theorem 1 suggests that we can control the step-size t by the proxy step-size t' within this interval.

Proposition 3. *If $h(\cdot)$ is convex and absolutely homogeneous and if $\|\text{grad} g|_{\mathbf{x}_k}\|_2$ is bounded, then for any $t' \neq 0$, we have $t = \phi(t') \neq 0$.*

Proof. It can be shown (see Appendix F for details) that:

$$|t| \geq \frac{1}{\sqrt{(1/t')^2 + \|\text{grad} g|_{\mathbf{x}_k}\|_2^2}}.$$

Thus if $\|\text{grad} g|_{\mathbf{x}_k}\|_2$ is bounded, then $t \neq 0$ for $t' \neq 0$. \square

We notice that $\|\text{grad} g|_{\mathbf{x}_k}\|_2 \leq \|\mathbf{I} - \mathbf{x}_k \mathbf{x}_k^\top\| \|\nabla g|_{\mathbf{x}_k}\|_2 = \|\nabla g|_{\mathbf{x}_k}\|_2$, thus it suffices to have a bounded $\|\nabla g|_{\mathbf{x}_k}\|_2$. In particular, if the Euclidean gradient ∇g is Lipschitz continuous on the sphere, $\|\nabla g|_{\mathbf{x}_k}\|_2$ is bounded.

Algorithm 1: One line-search iteration.

```

1 function ( $\mathbf{v}_k, t, t'$ )  $\leftarrow$  lineSearch ( $\mathbf{x}_k, t'_{\max}$ )
2    $t' \leftarrow \min \{t'_{\max}, 1/h(\mathbf{x}_k)\}$ 
3    $\mathbf{z} \leftarrow \text{prox}_{|t'|h}(\mathbf{x}_k - t' \text{grad} g|_{\mathbf{x}_k})$ 
4    $\mathbf{v}_k \leftarrow \frac{1}{\mathbf{x}_k^\top \mathbf{z}} \mathbf{z} - \mathbf{x}_k$ 
5    $t \leftarrow \frac{1}{\mathbf{x}_k^\top \mathbf{z}} t'$ 
6    $Q_L = g(\mathbf{x}_k) + \langle \text{grad} g|_{\mathbf{x}_k}, \mathbf{v}_k \rangle + \frac{1}{2t} \langle \mathbf{v}_k, \mathbf{v}_k \rangle$ 
7   if  $g(\mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k)) \leq Q_L$  then return ( $\mathbf{v}_k, t, t'$ )
8   else  $t' \leftarrow 0.8t'$ , goto step 3
9 end

```

4.3 Line-search and Convergence

To establish the convergence proof, we make the following mild assumption on the cost $g(\cdot)$.

Assumption 1. *We assume the pullback function $g(\mathcal{R}_{\mathbf{x}}(\mathbf{v}))$ with $\mathbf{v} \in \mathcal{T}_{\mathbf{x}}\mathcal{S}$ satisfies:*

$$g(\mathcal{R}_{\mathbf{x}}(\mathbf{v})) \leq g(\mathbf{x}) + \langle \text{grad} g(\mathbf{x}), \mathbf{v} \rangle + \frac{L}{2} \langle \mathbf{v}, \mathbf{v} \rangle, \quad (24)$$

for a (known or unknown) constant $L > 0$.

Assumption 1 holds if the ambient Euclidean space function $g(\cdot)$ on \mathbb{R}^n has Lipschitz continuous gradient ∇g within the convex hull of \mathcal{S} (Lemma 3 in [42]), *i.e.*, within the unit ball $\{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$. In other words, this assumption is satisfied if $g(\cdot)$ is not changing radically within the ball (like going to infinity at some point). Thus in practice, this assumption is hardly violated.

Based on Assumption 1, we propose the following line-search.

Line-search criterion. Search for a proxy step-size t' , such that the actual step-size t and the tangent update \mathbf{v}_k solved from equation (18) satisfy:

$$g(\mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k)) \leq g(\mathbf{x}_k) + \langle \text{grad} g|_{\mathbf{x}_k}, \mathbf{v}_k \rangle + \frac{1}{2t} \langle \mathbf{v}_k, \mathbf{v}_k \rangle. \quad (25)$$

By Assumption 1, the inequality (25) is satisfied for any $0 \leq t \leq 1/L$. This line-search criterion is in the same spirit as the one used in the classical Euclidean proximal gradient literature, where the Euclidean gradient $\nabla g|_{\mathbf{x}_k}$ is used instead of the Riemannian gradient here.

Line-search process. We start with an initial proxy step-size $t' < 1/h(\mathbf{x}_k)$ and reduce t' until $t = \phi(t')$ satisfies $0 \leq t \leq 1/L$. This process is well-defined by the monotonicity of $t = \phi(t')$, as proved by Theorem 1. One line-search iteration is given in Algorithm 1.

We now show the line-search criterion guarantees a descent for the total cost $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$ at each iteration. See the PGS curve in Fig. 4 for an illustration.

Lemma 6. *For retraction (12), if $h(\cdot)$ is absolutely homogeneous and $\mathbf{v}_k \in \mathcal{T}_{\mathbf{x}_k}\mathcal{S}$, we have $h(\mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k)) \leq h(\mathbf{x}_k + \mathbf{v}_k)$.*

Proof. See Appendix G. \square

Theorem 2. *Let $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$. If the line-search criterion (25) holds then:*

$$f(\mathbf{x}_{k+1}) = f(\mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k)) \leq f(\mathbf{x}_k) - \frac{1}{2t} \langle \mathbf{v}_k, \mathbf{v}_k \rangle. \quad (26)$$

Proof. The first-order necessary optimality condition of problem (13) states that:

$$\mathbf{0} \in \text{grad } g|_{\mathbf{x}_k} + \frac{1}{t} \mathbf{v}_k + \text{proj}_{\mathcal{T}_{\mathbf{x}_k} \mathcal{S}} \partial h|_{\mathbf{x}_k + \mathbf{v}_k} \quad (27)$$

$$\Leftrightarrow -\text{grad } g|_{\mathbf{x}_k} - \frac{1}{t} \mathbf{v}_k \in \text{proj}_{\mathcal{T}_{\mathbf{x}_k} \mathcal{S}} \partial h|_{\mathbf{x}_k + \mathbf{v}_k}. \quad (28)$$

Because $\mathbf{v}_k \in \mathcal{T}_{\mathbf{x}_k} \mathcal{S}$ then $\langle \mathbf{z}, \mathbf{v}_k \rangle = \langle \text{proj}_{\mathcal{T}_{\mathbf{x}_k} \mathcal{S}} \mathbf{z}, \mathbf{v}_k \rangle$ for any $\mathbf{z} \in \mathbb{R}^n$. By the convexity of $h(\cdot)$ at $\mathbf{x}_k + \mathbf{v}_k$, we obtain:

$$h(\mathbf{x}_k + \mathbf{v}_k) \leq h(\mathbf{x}_k) + \langle \partial h|_{\mathbf{x}_k + \mathbf{v}_k}, \mathbf{v}_k \rangle \quad (29)$$

$$= h(\mathbf{x}_k) + \langle \text{proj}_{\mathcal{T}_{\mathbf{x}_k} \mathcal{S}} \partial h|_{\mathbf{x}_k + \mathbf{v}_k}, \mathbf{v}_k \rangle \quad (30)$$

$$\stackrel{(28)}{=} h(\mathbf{x}_k) + \langle -\text{grad } g|_{\mathbf{x}_k} - \frac{1}{t} \mathbf{v}_k, \mathbf{v}_k \rangle. \quad (31)$$

By considering Lemma 6, we have the following inequality:

$$h(\mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k)) \leq h(\mathbf{x}_k) + \langle -\text{grad } g|_{\mathbf{x}_k} - \frac{1}{t} \mathbf{v}_k, \mathbf{v}_k \rangle. \quad (32)$$

Summing together inequalities (25) and (32), we obtain inequality (26). \square

We denote t'_k ($k = 0, 1, \dots, K-1$) to be the proxy step-size satisfying the line-search criterion (25) at iteration k . In the line-search process, we have:

$$0 < t'_k \leq \min\{t'_{\max}, c/h(\mathbf{x}_k)\}, \quad \text{with } 0 < c < 1.$$

We further denote the corresponding actual step-sizes as t_k where $t_k = \phi(t'_k) = t'_k/c(t'_k)$. From Theorem 1, we know $t_k \geq 0$. By Proposition 3, we know if $t'_k \neq 0$ then $t_k \neq 0$ as $\|\text{grad } g|_{\mathbf{x}_k}\|_2$ is bounded by the assumption of Lipschitz continuous gradient ∇g within the unit ball. Therefore we conclude $t_k > 0$. To summarize, we have:

$$0 < t_{\min} \leq t_k \leq t_{\max},$$

where we denote $t_{\min} = \min\{t_k\}$, $t_{\max} = \max\{t_k\}$.

At last, we show the iterations guided by the line-search process converge to a critical point of problem (2).

Proposition 4. *Assume $f(\mathbf{x})$ is bounded from below on \mathcal{S} , i.e., the problem is well-posed. If $g(\cdot)$ has Lipschitz continuous gradient ∇g within the unit ball $\{\mathbf{x} : \|\mathbf{x}\| \leq 1\}$, the line-search iterations converge to a critical point of problem (2).*

Proof. Let f^* be the optimal value, i.e., $f^* \leq f(\mathbf{x})$ for any \mathbf{x} on \mathcal{S} . Given K iterations, inequality (26) implies:

$$\begin{aligned} f(\mathbf{x}_0) - f^* &\geq f(\mathbf{x}_0) - f(\mathbf{x}_K) = \sum_{k=0}^{K-1} (f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})) \\ &\geq \sum_{k=0}^{K-1} \frac{1}{2t_k} \langle \mathbf{v}_k, \mathbf{v}_k \rangle = \sum_{k=0}^{K-1} \frac{t_k}{2} \langle \mathbf{v}_k/t_k, \mathbf{v}_k/t_k \rangle. \end{aligned}$$

Considering $0 < t_{\min} \leq t_k \leq t_{\max}$, we have:

$$\begin{aligned} f(\mathbf{x}_0) - f^* &\geq \frac{1}{2t_{\max}} \sum_{k=0}^{K-1} \|\mathbf{v}_k\|_2^2, \\ f(\mathbf{x}_0) - f^* &\geq \frac{t_{\min}}{2} \sum_{k=0}^{K-1} \|\mathbf{v}_k/t_k\|_2^2. \end{aligned}$$

Denote $\epsilon_0 = f(\mathbf{x}_0) - f^*$. Taking K to the infinity, we obtain:

$$\lim_{K \rightarrow \infty} \sum_{k=0}^{K-1} \|\mathbf{v}_k\|_2^2 \leq 2\epsilon_0 t_{\max}, \quad \lim_{K \rightarrow \infty} \sum_{k=0}^{K-1} \left\| \frac{\mathbf{v}_k}{t_k} \right\|_2^2 \leq \frac{2\epsilon_0}{t_{\min}}.$$

Since $f(\mathbf{x})$ is bounded from below on \mathcal{S} , ϵ_0 is a non-negative constant. Thus the right side of each inequality is bounded. Noting the left side is the summation of an infinite non-negative sequence, we have $\lim_{k \rightarrow \infty} \|\mathbf{v}_k\|_2 = 0$, $\lim_{k \rightarrow \infty} \|\mathbf{v}_k/t_k\|_2 = 0$, which means $\lim_{k \rightarrow \infty} \mathbf{v}_k = \mathbf{0}$, $\lim_{k \rightarrow \infty} \mathbf{v}_k/t_k = \mathbf{0}$. Lastly, we notice if $\mathbf{v}_k = \mathbf{0}$ and $\mathbf{v}_k/t_k = \mathbf{0}$, the first-order necessary optimality condition of problem (13), i.e., equation (27), becomes:

$$\mathbf{0} \in \text{grad } g|_{\mathbf{x}_k} + \text{proj}_{\mathcal{T}_{\mathbf{x}_k} \mathcal{S}} \partial h|_{\mathbf{x}_k}, \quad (33)$$

which is exactly the first-order necessary optimality condition of problem (2) [29], [43]. \square

Both $\|\mathbf{v}_k\|$ and $\|\mathbf{v}_k/t_k\|$ have linear convergence rates:

$$\min_{k=1,2,\dots,K} \|\mathbf{v}_k\|_2^2 \leq \frac{2\epsilon_0 t_{\max}}{K}, \quad \min_{k=1,2,\dots,K} \left\| \frac{\mathbf{v}_k}{t_k} \right\|_2^2 \leq \frac{\epsilon_0}{t_{\min} K},$$

with the constant decided by line-search strategies.

4.4 Algorithm

The final algorithm is to repeat the line-search until convergence. The pseudocode is provided in Algorithm 3, which unifies the PGS method in this section and the accelerated methods (A-PGS, AM-PGS) to be discussed shortly. Here we give additional components to complete the PGS method, before moving to its accelerated versions.

4.4.1 Maximum proxy step-size t'_{\max}

While the constant L in Assumption 1 (i.e., inequality (24)) exists ubiquitously, it is often not clear how to obtain L in closed-form except for certain types of $g(\cdot)$ e.g., $g(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$ where $L = 2\sigma_{\max}(\mathbf{A})$ (see Section 6.1). Importantly, the existence of L is used to establish proofs, but its actual value is never required to be known explicitly. Instead, to complete Algorithm 1, we need to decide t'_{\max} .

Setting t'_{\max} from known L . In line-search criterion (25), L establishes a lower-bound for the search, where if step-size $t < 1/L$ then the total cost is guaranteed to descend (Theorem 2). While we work with proxy step-size t' and control step-size t by $t = \phi(t')$ in its monotone region, the value $1/L$ provides a good reference for the maximum proxy step-size t'_{\max} . Practically, if L is known ahead, we recommend using $t'_{\max} = 1/L$.

Setting t'_{\max} from line-search. Nonetheless, if L is unknown (which is the usual case), t'_{\max} can be decided effectively from line-search. Such a line-search process is described in Algorithm 2, which is well-defined owing to the existence of L . Numerical examples are provided in Fig. 3. We compare the searched t'_{\max} and the known reference proxy step-size $t' = 1/L$ using their ratio, and see that a proper t'_{\max} close to $1/L$ can be found cheaply within 5 - 10 iterations.

Adaptive maximum proxy step-size t'_{\max} . We observe that typically the corresponding $t = \phi(t')$ is slightly larger than t' for t' around $1/L$, as $c(t')$ is slightly below

Algorithm 2: Search for proxy step-size t'_{\max} .

```

1 function  $t'_{\max} \leftarrow \text{searchMaxProxyStepsize}(\mathbf{x}_0)$ 
2    $found \leftarrow false, ub \leftarrow 0.7/h(\mathbf{x}_0), t' \leftarrow ub$ 
3    $\mathbf{z} \leftarrow \text{prox}_{|t'|h}(\mathbf{x}_0 - t' \text{grad } g|_{\mathbf{x}_0})$ 
4    $\mathbf{v}_0 \leftarrow \frac{1}{\mathbf{x}_0^\top \mathbf{z}} \mathbf{z} - \mathbf{x}_0$ 
5    $t \leftarrow \frac{1}{\mathbf{x}_0^\top \mathbf{z}} t'$ 
6    $Q_L = g(\mathbf{x}_0) + \langle \text{grad } g|_{\mathbf{x}_0}, \mathbf{v}_0 \rangle + \frac{1}{2t} \langle \mathbf{v}_0, \mathbf{v}_0 \rangle$ 
7   if  $g(\mathcal{R}_{\mathbf{x}_0}(\mathbf{v}_0)) \leq Q_L$  then
8     if  $t' = ub$  then return  $t'_{\max} \leftarrow ub$ 
9      $found \leftarrow true, t' \leftarrow \min\{2t', ub\}$ , goto step 3
10  else if  $found = true$  then return  $t'_{\max} \leftarrow 0.5t'$ 
11  else  $t' \leftarrow 0.1t'$ , goto step 3
12 end

```

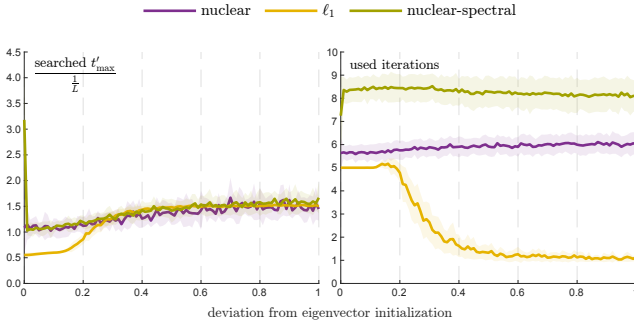


Fig. 3: Examples of line-searched maximum proxy step-size t'_{\max} , using different initializations \mathbf{x}_0 (detailed in Section 7.1.1). On the left we report the ratio $t'_{\max}/(1/L)$, and on the right the used iterations.

1 as seen in Fig. 2. Thus setting $t' = 1/L$ does not necessarily guarantee a line-search success from criterion (25) established on t . On the other hand, criterion (25) can also be satisfied for some t greater than $1/L$. This motivates us to use an adaptive t'_{\max} , where we set t'_{\max} to the working t' obtained at the previous iteration. This choice is controlled by the *AdaptiveMaxProxyStepsize* flag in Algorithm 3. We shall see such a strategy is useful to reduce subsequent total line-search iterations if t'_{\max} is initially obtained from the line-search in Algorithm 2 (see results in Fig. 5).

4.4.2 Stop criteria for convergence

In general, we propose to monitor at least $\|\mathbf{v}_k\|_2$ and $\|\mathbf{v}_k/t_k\|_2$ as the stop criteria for convergence. Numerical examples are provided in Fig. 4. While other options are also possible, *e.g.*, by monitoring $f(\cdot)$, these two indicators are important for the following reasons.

Convergence of estimates. The convergence of the estimates can be determined from the distance of \mathbf{x}_k and \mathbf{x}_{k+1} , *e.g.*, by using the chordal distance $d = \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2$ or the angle $\theta = \arccos(\mathbf{x}_{k+1}^\top \mathbf{x}_k)$. Here we propose to check the length of the tangent vector $\|\mathbf{v}_k\|_2$. In fact, $\mathbf{v}_k = \mathcal{R}_{\mathbf{x}_k}^{-1}(\mathbf{x}_{k+1})$ with $\mathcal{R}_{\cdot}^{-1}(\cdot)$ to be defined in equation (36), thus measuring the distance of \mathbf{x}_k and \mathbf{x}_{k-1} .

Optimality as critical points. If $\mathbf{v}_k/t_k = \mathbf{0}$ and $\mathbf{v}_k = \mathbf{0}$, equation (27) becomes equation (33), thus \mathbf{x}_k admits a critical point of problem (2) by satisfying the first-order

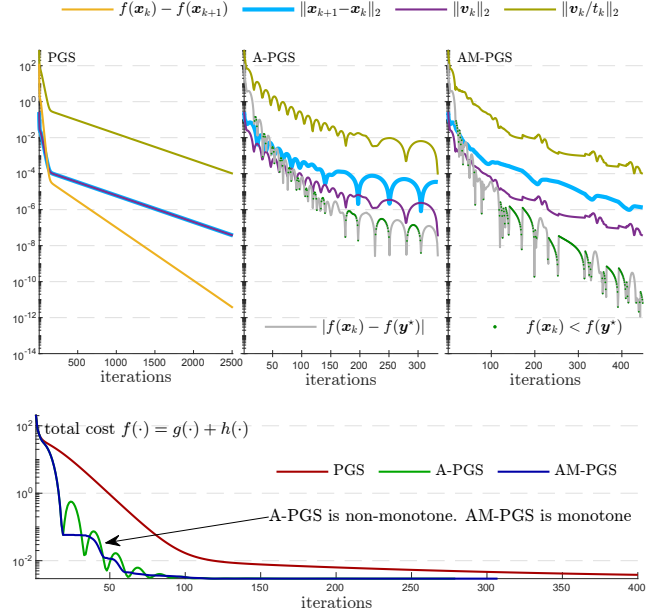


Fig. 4: The convergence behavior of the PGS, A-PGS and AM-PGS methods. The example is drawn from self-calibration with nuclear norm regularization. For each method, we report the convergence of the estimate as $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2$, of the tangent vector as $\|\mathbf{v}_k\|_2$ and of the first-order optimality as $\|\mathbf{v}_k/t_k\|_2$. The progression of the total cost $f(\mathbf{x}_k)$ for each method is plotted in the second figure. In the first figure, for PGS, we report the convergence of the total cost as $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})$ which is always positive from Theorem 2. For A-PGS and AM-PGS, we report the absolute value $|f(\mathbf{x}_k) - f(\mathbf{y}^*)|$ and mark down cases where $f(\mathbf{x}_k) < f(\mathbf{y}^*)$.

optimality condition. For a deeper understanding, we see from the KKT system (17):

$$\begin{aligned} \frac{\mathbf{v}_k}{t_k} &= \frac{1}{t'_k} (\mathbf{I} - \mathbf{x}_k \mathbf{x}_k^\top) \text{prox}_{|t'_k|h}(\mathbf{x}_k - t'_k \text{grad } g|_{\mathbf{x}_k}) \\ &= \text{proj}_{\mathcal{T}_{\mathbf{x}_k} \mathcal{S}} \text{prox}_h\left(\frac{1}{t'_k} \mathbf{x}_k - \text{grad } g|_{\mathbf{x}_k}\right). \end{aligned}$$

Thus \mathbf{v}_k/t_k is to project the proximal of $\frac{1}{t'_k} \mathbf{x}_k - \text{grad } g|_{\mathbf{x}_k}$ into the tangent space at \mathbf{x}_k . A critical point is where this tangent projection goes to zero.

4.4.3 Initialization \mathbf{x}_0

Typically, it is a good idea to initialize the regularized problem (2) from the solution of the original problem (1). However this choice is problem dependent and should be discussed specifically according to the problem at hand. For example, we initialize the regularized problems for fundamental matrix estimation and correspondence association with the solution of the Rayleigh quotient optimization (the original problem of these instances). However, for self-calibration, the regularized problems are initialized from the canonical DAQ after quasi-calibration, as the solution space of the original problem is likely to be ambiguous due to the critical motion sequence. Details are given in Section 6, and numerical validations provided in Section 7.1.

5 ACCELERATION USING THE NESTEROV MOMENTUM TECHNIQUE

The conventional PGS method evaluates the gradient and proximal at the current estimate \mathbf{x}_k .

For the Nesterov momentum technique, the gradient and the proximal are instead evaluated at an auxiliary state \mathbf{y}_k defined as a linear combination of the current estimate \mathbf{x}_k and the previous estimate \mathbf{x}_{k-1} [12], [14], [33]. In the Euclidean case, one accelerated iteration is defined as:

$$\begin{cases} \mathbf{x}_{k+1} = \mathbf{y}_k + \mathbf{v}_k \\ \mathbf{y}_{k+1} = \mathbf{x}_{k+1} + \frac{1-\alpha_k}{\alpha_{k+1}} (\mathbf{x}_k - \mathbf{x}_{k+1}), \quad \mathbf{y}_0 = \mathbf{x}_0, \end{cases} \quad (34)$$

where \mathbf{v}_k is obtained by evaluating iteration (6) at \mathbf{y}_k as:

$$\mathbf{v}_k = \arg \min_{\mathbf{v} \in \mathbb{R}^n} \langle \nabla g|_{\mathbf{y}_k}, \mathbf{v} \rangle + \frac{1}{2t} \langle \mathbf{v}, \mathbf{v} \rangle + h(\mathbf{y}_k + \mathbf{v}),$$

and the scalar sequence α_k ($k = 1, 2, \dots$) is defined as:

$$\alpha_{k+1} = \frac{1 + \sqrt{1 + 4\alpha_k^2}}{2}, \quad \alpha_0 = 1. \quad (35)$$

The iteration defined above was first proposed by Nesterov for smooth optimization [33], and later extended to non-smooth composite optimization in [12], [14]. The original proof shows that the accelerated iteration attains quadratic convergence for convex cost functions. These results are also valid if the cost function is locally convex around a local minimum. Although Riemannian manifolds are typically non-convex, it has been shown that the Nesterov sequence can attain quadratic convergence rate for geodesically convex optimization problems on the manifold [44], [45].

To extend the above result to the sphere manifold, we need to evaluate the difference between \mathbf{x}_{k+1} and \mathbf{x}_k on the sphere. Inspired by [40], we define this difference as a vector $\Delta \mathbf{v}$ in the tangent space of \mathbf{x}_{k+1} , thus the first summation can be extended by the retraction at \mathbf{x}_{k+1} . Such $\Delta \mathbf{v}$ must satisfy $\mathcal{R}_{\mathbf{x}_{k+1}}(\Delta \mathbf{v}) = \mathbf{x}_k$. Abusing notations, we define the inverse of the retraction, $\Delta \mathbf{v} = \mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{x}_k)$, which can be calculated in closed-form as:

$$\mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{x}_k) \stackrel{\text{def}}{=} \frac{1}{\mathbf{x}_k^\top \mathbf{x}_{k+1}} \mathbf{x}_k - \mathbf{x}_{k+1}. \quad (36)$$

It can be easily verified that $\mathcal{R}_{\mathbf{x}_{k+1}}(\mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{x}_k)) = \mathbf{x}_k$ and $\mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{x}_k) \in \mathcal{T}_{\mathbf{x}_{k+1}} \mathcal{S}$ for any $\mathbf{x}_k, \mathbf{x}_{k+1} \in \mathcal{S}$.

Overall we extend the Nesterov sequence to the sphere manifold as:

$$\begin{cases} \mathbf{x}_{k+1} = \mathcal{R}_{\mathbf{y}_k}(\mathbf{v}_k) \\ \mathbf{y}_{k+1} = \mathcal{R}_{\mathbf{x}_{k+1}} \left(\frac{1-\alpha_k}{\alpha_{k+1}} \mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{x}_k) \right), \quad \mathbf{y}_0 = \mathbf{x}_0. \end{cases} \quad (37)$$

The extension is to replace Euclidean addition and subtraction with Riemannian retraction $\mathcal{R}(\cdot)$ and its inverse $\mathcal{R}^{-1}(\cdot)$. Here \mathbf{x}_k is obtained by evaluating iteration (13) at \mathbf{y}_k as:

$$\mathbf{v}_k = \arg \min_{\mathbf{v} \in \mathcal{T}_{\mathbf{y}_k} \mathcal{S}} \langle \text{grad } g|_{\mathbf{y}_k}, \mathbf{v} \rangle + \frac{1}{2t} \langle \mathbf{v}, \mathbf{v} \rangle + h(\mathbf{y}_k + \mathbf{v}),$$

and the scalar α_k is defined as in equation (35). The tangent update \mathbf{v}_k can be solved in closed-form using the proxy step-size technique proposed in Section 4.

The estimates \mathbf{x}_k ($k = 1, 2, \dots$) generated from equation (37) do not guarantee the monotonicity of the total cost $f(\mathbf{x}_k) = g(\mathbf{x}_k) + h(\mathbf{x}_k)$ [46]. On the manifold setting, this can sometimes lead to divergence if \mathbf{v}_k is computed based on the SSNM method [40]. To detect and recover from potential failures, the authors in [40] introduced a safeguard by monitoring the progression of the cost function within several iterations. One potential reason for the divergence is the inexact computation of each iteration [46]:

“In our case, where the denoising subproblems are not solved exactly, monotonicity becomes an important issue. It might happen that due to the inexact computations of the denoising subproblems, the algorithm might become extremely non-monotone and in fact can even diverge!”

Practically in our experiments to be presented in Section 7, we did not observe the divergence of iterations by using the proxy step-size to obtain the tangent update \mathbf{v}_k . This may be due to the fact that we solve each iteration exactly (in closed-form) while the SSNM method used in [29], [30], [40] is iterative thus incurring inexact solutions.

Nonetheless, we propose a monotone algorithm for the sphere manifold based on Beck *et al.*'s [46] Euclidean version which has been proved to retain the quadratic convergence rate. Beck *et al.*'s [46] monotone algorithm is defined as:

$$\begin{cases} \mathbf{y}^* = \mathbf{y}_k + \mathbf{v}_k \\ \mathbf{x}_{k+1} = \begin{cases} \mathbf{y}^* & \text{if } f(\mathbf{y}^*) < f(\mathbf{x}_k) \\ \mathbf{x}_k & \text{otherwise} \end{cases} \\ \mathbf{y}_{k+1} = \begin{cases} \mathbf{y}^* + \frac{1-\alpha_k}{\alpha_{k+1}} (\mathbf{x}_k - \mathbf{y}^*) & \text{if } f(\mathbf{y}^*) < f(\mathbf{x}_k) \\ \mathbf{x}_k + \frac{\alpha_k}{\alpha_{k+1}} (\mathbf{y}^* - \mathbf{x}_k) & \text{otherwise.} \end{cases} \end{cases}$$

The above iteration ensures the monotonicity of the total cost $f(\cdot)$ by leveraging between the new estimate \mathbf{y}^* and the previous estimate \mathbf{x}_k . We extend Beck *et al.*'s monotone algorithm to the sphere manifold as follow:

$$\begin{cases} \mathbf{y}^* = \mathcal{R}_{\mathbf{y}_k}(\mathbf{v}_k) \\ \mathbf{x}_{k+1} = \begin{cases} \mathbf{y}^* & \text{if } f(\mathbf{y}^*) < f(\mathbf{x}_k) \\ \mathbf{x}_k & \text{otherwise} \end{cases} \\ \mathbf{y}_{k+1} = \begin{cases} \mathcal{R}_{\mathbf{y}^*} \left(\frac{1-\alpha_k}{\alpha_{k+1}} \mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{x}_k) \right) & \text{if } f(\mathbf{y}^*) < f(\mathbf{x}_k) \\ \mathcal{R}_{\mathbf{x}_k} \left(\frac{\alpha_k}{\alpha_{k+1}} \mathcal{R}_{\mathbf{x}_{k+1}}^{-1}(\mathbf{y}^*) \right) & \text{otherwise.} \end{cases} \end{cases}$$

The accelerated PGS (A-PGS), *i.e.*, Nesterov sequence, and the accelerated monotone PGS (AM-PGS), *i.e.*, Beck's sequence, are implemented as pseudocode in Algorithm 3. An illustration of the convergence is given in Fig. 4.

6 APPLICATIONS

6.1 Rayleigh Quotient Optimization

We consider $g(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ with \mathbf{A} being a symmetric matrix. The Euclidean gradient at $\mathbf{x} \in \mathbb{R}^n$ is $\nabla g|_{\mathbf{x}} = 2\mathbf{A}\mathbf{x}$. The Riemannian gradient at $\mathbf{x} \in \mathcal{S}$ is:

$$\begin{aligned} \text{grad } g(\mathbf{x}) &= \nabla g|_{\mathbf{x}} - \langle \mathbf{x}, \nabla g|_{\mathbf{x}} \rangle \mathbf{x} \\ &= 2\mathbf{A}\mathbf{x} - 2 \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle \mathbf{x} = 2\mathbf{A}\mathbf{x} - 2(\mathbf{x}^\top \mathbf{A}\mathbf{x})\mathbf{x}. \end{aligned} \quad (38)$$

Algorithm 3: A unified implementation of the PGS, A-PGS and AM-PGS methods.

```

input : method = PGS or A-PGS or AM-PGS
input :  $\mathbf{x}_0$ 
1 if known Lipschitz constant  $L$  then  $t'_{\max} \leftarrow 1/L$ 
2 else  $t'_{\max} \leftarrow \text{searchMaxProxyStepsize}(\mathbf{x}_0)$ 
3  $\mathbf{y}_0 \leftarrow \mathbf{x}_0, \alpha_0 \leftarrow 1, k \leftarrow 0$ 
4 while  $k < \text{maxIterations}$  do
5    $(\mathbf{v}_k, t, t') \leftarrow \text{lineSearch}(\mathbf{y}_k, t'_{\max})$ 
6    $\mathbf{y}^* \leftarrow \mathcal{R}_{\mathbf{y}_k}(\mathbf{v}_k)$ 
7   if AdaptiveMaxProxyStepsize then  $t'_{\max} \leftarrow t'$ 
8   if method = PGS then
9      $\mathbf{x}_{k+1} \leftarrow \mathbf{y}^*$ 
10     $\mathbf{y}_{k+1} \leftarrow \mathbf{y}^*$ 
11  else if method = A-PGS or AM-PGS then
12     $\alpha_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4\alpha_k^2}}{2}$ 
13     $\mathbf{x}_{k+1} \leftarrow \mathbf{y}^*$ 
14     $\mathbf{y}_{k+1} \leftarrow \mathcal{R}_{\mathbf{y}^*} \left( \frac{1 - \alpha_k}{\alpha_{k+1}} \mathcal{R}_{\mathbf{y}^*}^{-1}(\mathbf{x}_k) \right)$ 
15    if method = AM-PGS then
16      if  $g(\mathbf{y}^*) + h(\mathbf{y}^*) > g(\mathbf{x}_k) + h(\mathbf{x}_k)$  then
17         $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k$ 
18         $\mathbf{y}_{k+1} \leftarrow \mathcal{R}_{\mathbf{x}_k} \left( \frac{\alpha_k}{\alpha_{k+1}} \mathcal{R}_{\mathbf{x}_k}^{-1}(\mathbf{y}^*) \right)$ 
19      end
20    end
21  end
22  if  $\|\mathbf{v}_k\|_2 < 1e - 5$  and  $\|\mathbf{v}_k/t\|_2 < 1e - 3$  then
23    return  $\mathbf{x}_{k+1}$ 
24  else  $k \leftarrow k + 1$ 
25 end

```

Using retraction (12), the Lipschitz-type constant of $g(\mathcal{R}_{\mathbf{x}}(\mathbf{v}))$ is $L = 2\sigma_{\max}(\mathbf{A})$ (see Appendix H) where $\sigma_{\max}(\mathbf{A})$ denotes the largest singular value of \mathbf{A} .

Minimizing $g(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ on the sphere manifold is called Rayleigh quotient optimization:

$$\mathbf{x}_0 = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^\top \mathbf{A} \mathbf{x} \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1. \quad (39)$$

The solution to this problem is given in closed-form, which is the bottom eigenvector of \mathbf{A} (i.e., the eigenvector associated with the smallest eigenvalue), which often gives as a good initialization for its regularized versions.

6.2 Fundamental Matrix Estimation

6.2.1 Problem Statement

The fundamental matrix is a key algebraic model of the two-view geometry [1], [3], [47]. We denote $\mathbf{p}_i \leftrightarrow \mathbf{p}'_i$ ($i = [1 : m]$) the homogeneous coordinates of corresponding points in two images. In the noise-free case, the epipolar constraint holds as:

$$\mathbf{p}'_i{}^\top \mathbf{F} \mathbf{p}_i = 0, \quad i = [1 : m]. \quad (40)$$

The matrix $\mathbf{F} \in \mathbb{R}^{3 \times 3}$ is called the *fundamental matrix*, defined up to scale, thus we seek for \mathbf{F} on the unit sphere such that $\|\mathbf{F}\|_F = 1$. Importantly, $\text{rank}(\mathbf{F}) = 2$ is required. A related concept is the essential matrix for normalized calibration, for which we refer to a recent work [48].

To estimate \mathbf{F} , we can formulate a cost function based on the algebraic error as:

$$\begin{aligned} \varphi(\mathbf{F}) &= \frac{1}{m} \sum_{i=1}^m (\mathbf{p}'_i{}^\top \mathbf{F} \mathbf{p}_i)^2 = \frac{1}{m} \sum_{i=1}^m \left((\mathbf{p}'_i{}^\top \otimes \mathbf{p}'_i{}^\top) \text{vec}(\mathbf{F}) \right)^2 \\ &= \frac{1}{m} \|\mathbf{H} \text{vec}(\mathbf{F})\|_2^2, \end{aligned}$$

where the i -th row of \mathbf{H} is $\mathbf{p}'_i{}^\top \otimes \mathbf{p}'_i{}^\top$. We denote $\mathbf{x} = \text{vec}(\mathbf{F})$ and $\mathbf{F} = \text{mat}(\mathbf{x})$, where $\text{vec}(\cdot)$ is the standard column-wise matrix vectorization and $\text{mat}(\cdot)$ is its inverse operation. Defining $\mathbf{A} = \frac{1}{m} \mathbf{H}^\top \mathbf{H}$, we see the fundamental matrix problem is an instance of problem (39).

The solution $\mathbf{F}_0 = \text{mat}(\mathbf{x}_0)$ solved from problem (39) is usually of rank 3. A remedy is to subsequently round the solution \mathbf{F}_0 using the rank-2 approximation via the Singular Value Decomposition (SVD). This two-stage solution is the eight-point algorithm. Instead, we give a low-rank solution using the nuclear norm regularization.

6.2.2 Nuclear Norm Regularization

The nuclear norm of a matrix $\|\mathbf{X}\|_*$, defined as the summation of its singular values, is the tightest convex envelop of the rank function within the unit ball $\{\mathbf{X} \in \mathbb{R}^{m \times n} : \|\mathbf{X}\|_2 \leq 1\}$ [9]. Nuclear norm regularization has been widely used as a technique to promote low-rank in Euclidean optimization problems [13], while a direct deployment to the manifold setting seems to be obscure with the results in [29], [32], mainly due to the challenge incurred in evaluating the generalized Jacobian matrix of a non-smooth function. In contrast, our technique can handle nuclear norm regularization with no effort. We apply nuclear norm regularization to problem (39) as:

$$\arg \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \lambda \|\text{mat}(\mathbf{x})\|_* \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1, \quad (41)$$

where $\lambda \geq 0$ is a given constant controlling the strength of regularization.

In problem (41), the regularization term is $h(\mathbf{x}) = \lambda \|\text{mat}(\mathbf{x})\|_*$, which is convex but non-smooth. Besides, $h(\cdot)$ is absolutely homogeneous, i.e., $h(\alpha \mathbf{x}) = |\alpha| h(\mathbf{x})$, because:

$$\lambda \|\text{mat}(\alpha \mathbf{x})\|_* = \lambda \|\alpha \text{mat}(\mathbf{x})\|_* = \lambda |\alpha| \|\text{mat}(\mathbf{x})\|_*. \quad (42)$$

The first equality holds because $\text{mat}(\mathbf{x})$ is a linear operator, and the second because norms are absolutely homogeneous.

Let $\mathbf{X} = \text{mat}(\mathbf{x})$ and denote its SVD as $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$. For the nuclear norm function $\lambda \|\mathbf{X}\|_*$, the proximal is given in closed-form [49], where $\text{prox}_{t\lambda \|\cdot\|_*}(\mathbf{X}) = \mathbf{U} (\mathbf{\Sigma} - t\lambda \mathbf{I})_+ \mathbf{V}^\top$. Here $\mathbf{A}_+ \stackrel{\text{def}}{=} \max\{\mathbf{A}, \mathbf{0}\}$ proceeds element-wise. Thus we obtain the proximal of $h(\cdot)$ as:

$$\text{prox}_{th}(\mathbf{x}) = \text{vec} \left(\mathbf{U} (\mathbf{\Sigma} - t\lambda \mathbf{I})_+ \mathbf{V}^\top \right). \quad (43)$$

6.3 Correspondence Association

6.3.1 Problem Statement

The correspondence association problem using pairwise constraints can be formulated as a Rayleigh quotient optimization as well [5]. We denote the association hypothesis that a point i in the point-cloud \mathcal{Q} is matched with a point

i' in the point-cloud \mathcal{Q}' as $h_{ii'}$. The correspondence problem is to estimate the likelihood $p(h_{ii'})$ ($i \in \mathcal{Q}$, $i' \in \mathcal{Q}'$) of all possible association hypotheses, collected as components of the state vector \mathbf{x} .

To that end, we can design an adjacency matrix M from the pairwise consistency of hypotheses [5]. A common practice is based on the change of distance:

$$M(h_{ii'}, h_{jj'}) = \begin{cases} 4.5 - \frac{(d_{ij} - d_{i'j'})^2}{2\delta_d^2} & \text{if } |d_{ij} - d_{i'j'}| < 3\delta_d \\ 0 & \text{otherwise,} \end{cases} \quad (44)$$

where d_{ij} is the Euclidean distance between the points i and j in \mathcal{Q} , $d_{i'j'}$ the distance between the points i' and j' in \mathcal{Q}' , and δ_d a tuning parameter.

The resulting problem is formalized as maximizing the overall consistency $\mathbf{x}^T M \mathbf{x}$ on the unit sphere, as an instance of problem (39) by letting $\mathbf{A} = -M$. The estimate of \mathbf{x} is further used to decide the final correspondences, based on various assumptions, *e.g.*, one point in \mathcal{Q} can only be matched with one point in \mathcal{Q}' [5].

The match hypotheses solved this way are dense, while many of them present with contradictions or low probabilities. It is thus favorable to have a sparse \mathbf{x} , where some unlikely hypotheses and contradictions are pruned away. We give such a sparse solution by ℓ_1 norm regularization.

6.3.2 ℓ_1 Norm Regularization

It has been known that ℓ_1 norm regularization can favor sparsity in Euclidean [12] and manifold optimization [29], [32]. We apply ℓ_1 norm regularization to problem (39) as:

$$\arg \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda \|\mathbf{x}\|_{\ell_1} \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1, \quad (45)$$

where $\lambda \geq 0$ is a given constant controlling the strength of regularization. In this case, $h(\mathbf{x}) = \lambda \|\mathbf{x}\|_{\ell_1}$, which is convex and absolutely homogeneous. The proximal of $h(\mathbf{x})$, often termed soft shrinkage operator, is given element-wise as:

$$\text{prox}_{t h}(\mathbf{x})_i = \text{sgn}(\mathbf{x}_i)(|\mathbf{x}_i| - t\lambda)_+. \quad (46)$$

6.4 Camera Self-calibration

6.4.1 Problem Statement

In projective reconstruction, we obtain a set of projective cameras $\tilde{P}_i \in \mathbb{R}^{3 \times 4}$:

$$\tilde{P}_i \propto P_i H^{-1} \quad (i \in [1 : n]), \quad (47)$$

which differ from the Euclidean cameras $P_i \in \mathbb{R}^{3 \times 4}$ by a common projective transformation $H \in \mathbb{R}^{4 \times 4}$. The camera self-calibration problem is to infer H from \tilde{P}_i [3].

The key algebraic model to this task is the *Dual Absolute Quadric (DAQ)*, a rank-3 symmetric matrix in $\mathbb{R}^{4 \times 4}$ defined up to scale [2]. In specific, the DAQ in the Euclidean space, termed the canonical DAQ, takes the form $\Omega_\infty^* = \text{diag}(1, 1, 1, 0)$. We denote the DAQ in the projective space (where \tilde{P}_i is defined) by $Q_\infty^* = H \Omega_\infty^* H^T$. The image of the DAQ, denoted by ω_i^* , is invariant under H :

$$\omega_i^* \propto \tilde{P}_i Q_\infty^* \tilde{P}_i^T \propto P_i \Omega_\infty^* P_i^T \propto K_i K_i^T, \quad (48)$$

where $K_i \in \mathbb{R}^{3 \times 3}$ is the intrinsic matrix of P_i . At its core, the self-calibration problem is to estimate Q_∞^* from equation (48) using various constraints on K_i .

Here we consider a linear approach developed for cameras with varying focal lengths [50]. In this case, we have $K_i K_i^T \propto \text{diag}(f_i^2, f_i^2, 1)$. Based on equation (48), we have:

$$\begin{cases} \tilde{a}_i^T Q_\infty^* \tilde{a}_i = \tilde{b}_i^T Q_\infty^* \tilde{b}_i \\ \tilde{a}_i^T Q_\infty^* \tilde{b}_i = 0, \tilde{a}_i^T Q_\infty^* \tilde{c}_i = 0, \tilde{b}_i^T Q_\infty^* \tilde{c}_i = 0, \end{cases} \quad (49a)$$

$$\tilde{a}_i^T Q_\infty^* \tilde{b}_i = 0, \tilde{a}_i^T Q_\infty^* \tilde{c}_i = 0, \tilde{b}_i^T Q_\infty^* \tilde{c}_i = 0, \quad (49b)$$

where $\tilde{P}_i^T = [\tilde{a}_i \ \tilde{b}_i \ \tilde{c}_i]$. Equation (49) is linear in Q_∞^* thus can be rewritten as:

$$M_i \text{vec}_t(Q_\infty^*) = M_i \mathbf{x} = \mathbf{0}, \quad (50)$$

where we have defined $\text{vec}_t(Q_\infty^*) = \mathbf{x} \in \mathbb{R}^{10}$ comprising of the upper triangular elements of Q_∞^* , and $\text{mat}_t(\mathbf{x}) = Q_\infty^*$ its inverse operation. Since Q_∞^* is defined up to scale so is \mathbf{x} and we minimize the cost $\varphi(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \|M_i \mathbf{x}\|_2^2$ on the unit sphere. Upon defining $\mathbf{A} = \frac{1}{n} \sum_{i=1}^n M_i^T M_i$, we see it is another instance of problem (39).

Just as the case of the fundamental matrix estimation, the DAQ estimated from solving problem (39) is typically of rank 4 instead of rank 3, thus an SVD based rounding process is used subsequently.

Once obtaining a rank-3 estimate of Q_∞^* , we can recover H up to a similarity transformation. This is usually done by the eigen decomposition of Q_∞^* . Let $Q_\infty^* = U \Lambda U^T$, with $\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3, 0)$. We then set $H \propto U \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \sqrt{\lambda_3}, 1)$. The estimate of H determines camera P_i , and K_i afterwards by decomposing P_i .

6.4.2 Nuclear Norm Regularization

One issue regarding self-calibration is the critical motion sequences (CMS) [51]. The CMSs are camera configurations where self-calibration is ambiguous due to the lack of sufficient constraints. Among which, we consider the artificial CMS that can be resolved by enforcing the rank deficiency of the DAQ during the estimation rather than a posteriori. For the linear self-calibration described in equation (49), one such CMS is that - all the cameras' principal axes intersect at a fixed point, *i.e.*, all the cameras look towards a common point. In this case, there only exist two rank deficient solutions [51]: one rank-3 solution (desired) and one rank-1 solution (undesired).

In analogy to the fundamental matrix estimation, we use nuclear norm regularization to promote low rank. This case is similar to what we have discussed in Section 6.2.2. We omit the details as they can be readily derived by replacing $\text{mat}(\cdot)$ with $\text{mat}_t(\cdot)$ and $\text{vec}(\cdot)$ with $\text{vec}_t(\cdot)$ in problem (41).

6.4.3 Nuclear-Spectral Norm Regularization

The aforementioned nuclear norm regularization resolves the CMS only partly due to the existence of the rank-1 solution. We propose to avoid the rank-1 solution by additionally including a spectral norm to penalize the largest singular value. The nuclear-spectral norm regularizer is:

$$h(\mathbf{x}) = \lambda_1 \|\text{mat}_t(\mathbf{x})\|_* + \lambda_2 \|\text{mat}_t(\mathbf{x})\|_2, \quad (51)$$

where λ_1 and λ_2 control the regularization strength of each part. We apply this regularizer to problem (39) as:

$$\begin{aligned} \arg \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda_1 \|\text{mat}_t(\mathbf{x})\|_* + \lambda_2 \|\text{mat}_t(\mathbf{x})\|_2 \\ \text{s.t.} \quad \|\mathbf{x}\|_2 = 1. \end{aligned} \quad (52)$$

Denote the SVD of $\text{mat}_t(\mathbf{x})$ as $\text{mat}_t(\mathbf{x}) = \mathbf{U}\Sigma\mathbf{V}^\top$. Since both the nuclear and spectral norm are orthogonal invariant, we can derive the proximal of $h(\cdot)$ in a similar manner as the derivation used for the nuclear norm. The proximal of $h(\cdot)$ in equation (51) is given as follow:

$$\text{prox}_{th}(\mathbf{x}) = \text{vec}_t\left(\mathbf{U}(\Sigma - t\lambda_1\mathbf{I} - t\lambda_2\mathbf{E}_1)_+\mathbf{V}^\top\right), \quad (53)$$

where $\mathbf{E}_1 = \text{diag}(1, 0, \dots, 0)$ is a diagonal matrix where the top-left element is 1 and the rests are all-zeros.

7 EXPERIMENTAL RESULTS

7.1 The Proposed PGS Algorithm

We first provide an evaluation of the PGS methods, *i.e.*, PGS, A-PGS and AM-PGS methods in Algorithm 3.

7.1.1 Experiment Setup

Numerical instances. We experiment with the Rayleigh quotient optimization, and draw numerical examples from different applications which essentially form different \mathbf{A} matrices in problem (39):

- nuclear norm reg. — fundamental matrix estimation with nuclear norm regularization,
- ℓ_1 norm reg. — correspondence association with ℓ_1 norm regularization,
- nuclear-spectral norm reg. — self-calibration with nuclear-spectral norm regularization.

In this section, we distinguish these instances by the type of regularization used, *i.e.*, nuclear, ℓ_1 and nuclear-spectral.

Initialization. The numerical examples used in this section are special forms of the Rayleigh quotient optimization problem (39), and its optimal solution is known to be the bottom eigenvector of matrix \mathbf{A} which we denote by $\xi_{\mathbf{A}}$. To examine the convergence behavior with respect to different initializations, we create a range of initial values \mathbf{x}_0 by adding independent zero-mean Gaussian noise element-wisely to $\xi_{\mathbf{A}}$. For the k -th vector element, we let:

$$\xi[k] \leftarrow \xi_{\mathbf{A}}[k] + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \delta_{\text{init}}), \quad (54)$$

and normalize ξ to get the initial value $\mathbf{x}_0 = \xi / \|\xi\|_2$. Intuitively, δ_{init} controls the deviation from the eigenvector initialization $\xi_{\mathbf{A}}$, and if δ_{init} is large, the above process simulates random initialization.

Proxy step-size strategies. We examine the following proxy step-size strategies.

- LipschitzFixed — $t'_{\max} = 1/L$ initially, and t'_{\max} is kept fixed in the following iterations;
- LipschitzAdaptive — $t'_{\max} = 1/L$ initially, and at each iteration t'_{\max} is updated to the previous working proxy step-size t' ;
- SearchedFixed — t'_{\max} is obtained from Algorithm 2 initially, and t'_{\max} is kept fixed in the following iterations;
- SearchedAdaptive — t'_{\max} is obtained from Algorithm 2, and at each iteration t'_{\max} is updated to the previous working proxy step-size t' .

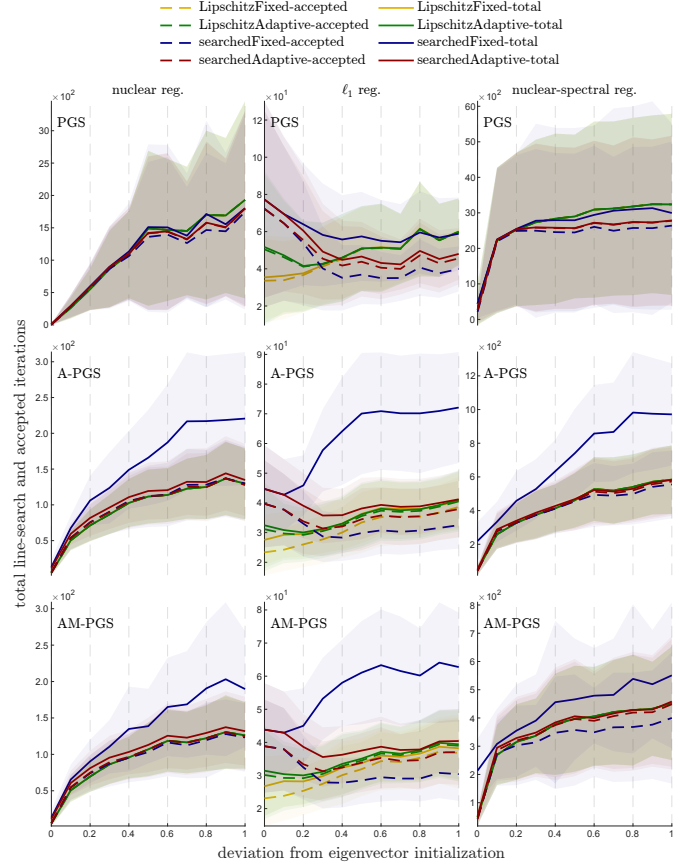


Fig. 5: The convergence of the PGS, A-PGS and AM-PGS methods by using different proxy step-size strategies. Results are reported as the total line-search and the accepted iterations with respect to different initialization.

7.1.2 Results

We report the results with a 20 run Monte-Carlo simulation for each PGS method and each proxy step-size strategy with respect to different initializations.

Convergence by used iterations. We examine the convergence by the used iterations in Fig. 5. In particular, we distinguish the total line-search iterations (the overall iterations processed in Algorithm 1) and the accepted iterations (the iterations used in Algorithm 3). The final verdicts is as follows: a) In general, we observe no significant differences for the accepted iterations across different proxy step-size strategies. b) However, the proxy step-size strategy SearchedFixed is not recommended, as it often leads to line-search failures especially when the initialization is bad, as reflected by the number of total line-search iterations. Therefore, if t'_{\max} is obtained from line-search in Algorithm 2, we recommend at each iteration updating t'_{\max} to the previous working proxy step-size t' . c) If t'_{\max} is initialized as $1/L$ from the Lipschitz constant L , both strategies LipschitzFixed and LipschitzAdaptive give similar results. d) We observe accelerated methods A-PGS and AM-PGS converge much faster than the unaccelerated PGS method, thus are generally recommended.

Variation of the optimal costs. We compare the optimal costs of each Monte-Carlo run in Fig. 6. In specific, in the

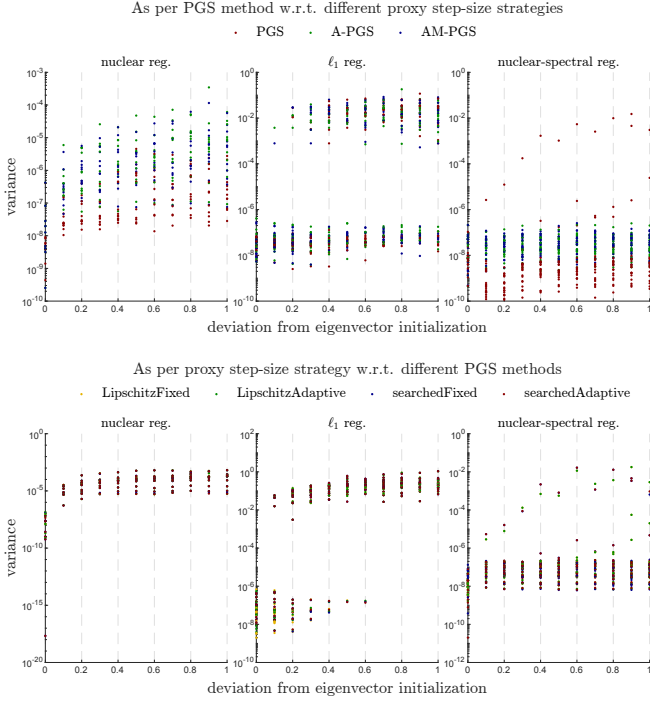


Fig. 6: The variation of the optimal costs for each case.

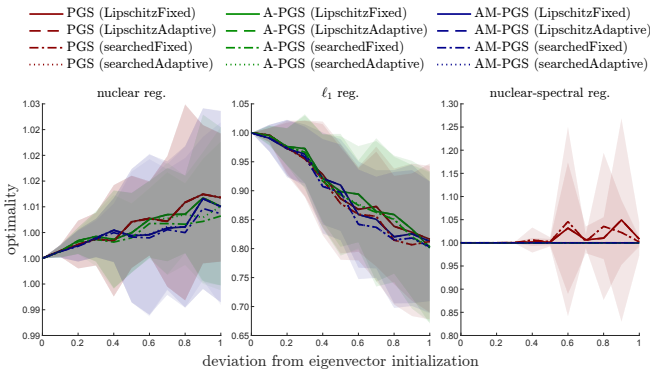


Fig. 7: The impact of different initializations.

top figure, for each PGS method (*i.e.*, PGS/A-PGS/AM-PGS), we compare the costs obtained from different proxy step-size strategies. Likewise, in the bottom figure, for each proxy step-size strategy, we compare the costs obtained from different PGS methods. The difference is evaluated as the variation of the optimal costs. From Fig. 6, we see that the optimal costs are mostly the same if δ_{init} in equation (54) is small, or otherwise stated if the initialization x_0 is close to the eigenvector initialization ξ_A where we simulate good initializations. As δ_{init} grows where we simulate bad initializations, the differences grow as different PGS methods and proxy step-size strategies can lead to the convergence to different local minima. This phenomenon is extremely clear for the ℓ_1 norm regularized instances, where x is valued mostly below 0.5 (see the example in Fig. 10) and $\delta_{\text{init}} > 0.5$ sets x_0 almost to random.

Optimality of the optimal costs. We evaluate the optimal cost obtained from the initialization x_0 by comparing with the optimal cost obtained from the eigenvector initialization

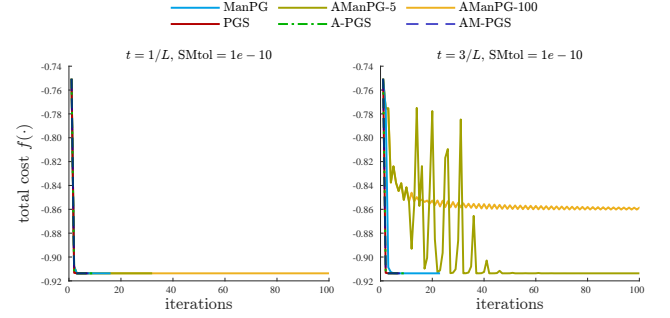


Fig. 8: The convergence of ManPG and AManPG using the line-search criterion (16) in comparison to the proposed one.

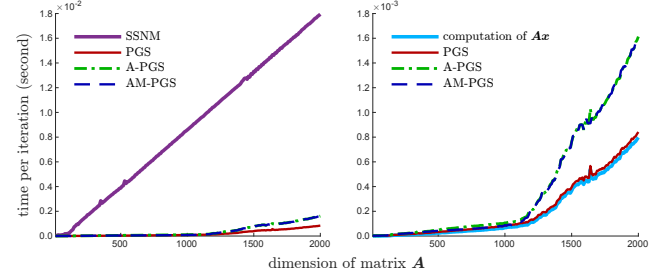


Fig. 9: The computational time per iteration by using the SSNM method [32] (used in [29], [30], [40]) and the proposed proxy step-size technique (used in the PGS, A-PGS and AM-PGS methods). The matrix-vector multiplication Ax is computed by the level-2 BLAS routine “dsymv”.

ξ_A . For each Monte-Carlo run, we define the optimality:

$$\text{optimality} = \frac{\text{optimal cost initialized from } x_0}{\text{optimal cost initialized from } \xi_A},$$

to benchmark the influence of different initializations. If the optimality metric is close to 1, then the optimal cost is close to the one obtained from the eigenvector initialization, and otherwise if this metric deviates from 1 then the computed solution is considered to be suboptimal. Intuitively, the optimality curve defines the robustness against bad initializations. The statistics for each tested case are plotted in Fig. 7. It is worth noting that the costs of the ℓ_1 norm regularized instances are negative (as $A = -M$), thus the optimality metric is below 1. For the problem instances used in this paper, the PGS, A-PGS and AM-PGS methods are in general robust to a large range of initializations, while the ℓ_1 norm regularized instances are more sensitive to initializations.

7.2 Comparison with ManPG [29] and AManPG [40]

We compare with the ManPG method [29] and its accelerated version AManPG [40]. Due to the difficulties of implementing the generalized Clarke differential for nuclear norm regularization, the comparison is only performed by applying ℓ_1 regularization to problem (39). We use the C++ implementation of ManPG and AManPG released in [40]. AManPG uses a safeguard mechanism every N iterations, where we set $N = 5$ and $N = 100$ and term the resulting methods as AManPG-5 and AManPG-100. The error-tolerance of the SSNM method [32] is set to $1e - 10$.

We implemented our methods PGS, A-PGS, and AM-PGS in C++ as well for a fair comparison.

	e_{dist} (pixels)↓					e_{rep} (pixels)↓				
	8pt	PGS ₅	PGS ₁₀	PGS	Gp	8pt	PGS ₅	PGS ₁₀	PGS	Gp
Chapel(0,1)	0.386	0.380	0.380	0.379	0.376	0.260	0.256	0.256	0.255	0.254
Keble(0,3)	0.248	0.247	0.247	0.247	0.247	0.175	0.175	0.175	0.175	0.175
Desktop(C,D)	0.574	0.336	0.303	0.288	0.268	0.406	0.238	0.214	0.204	0.190
Library(1,3)	0.428	0.409	0.405	0.405	0.400	0.302	0.289	0.286	0.286	0.282
Merton1(1,3)	0.308	0.295	0.291	0.288	0.277	0.217	0.209	0.205	0.203	0.196
Merton2(1,3)	0.596	0.528	0.498	0.472	0.404	0.421	0.373	0.352	0.334	0.286
Arch	0.304	0.299	0.299	0.299	0.298	0.215	0.211	0.211	0.211	0.211
Yard	0.433	0.429	0.429	0.428	0.426	0.306	0.303	0.303	0.302	0.301
Slate	0.246	0.184	0.177	0.170	0.163	0.129	0.097	0.093	0.089	0.085
Ben1	0.203	0.144	0.128	0.102	0.101	0.142	0.101	0.089	0.071	0.070
Ben2	0.139	0.086	0.063	0.050	0.048	0.097	0.060	0.044	0.035	0.033

TABLE 1: Fundamental matrix estimation with nuclear norm regularization (with data in [52]).

7.2.1 Line-search Criteria and Convergence

Theoretical justification. In ManPG’s line-search criterion (16), we first need to assign t and then perform line-search for α_k to ensure the descent of the total cost $f(\cdot)$. From Theorem 2 of our work, we see that if $t \leq 1/L$ with L being the Lipschitz constant, it suffices to set $\alpha_k = 1$ in the ManPG’s line-search criterion (16). The authors in [29], [40] assume known Lipschitz constant L and suggest to use $t = 1/L$ as a reference. Although for arbitrary $t > 0$, the existence of α_k for criterion (16) is proved in [29], we observe a proper t in ManPG/AManPG is required.

Numerical validation. On the left of Fig. 8, we see by setting $t = 1/L$, ManPG’s line-search (16) works almost the same way as the proposed PGS line-search. On the right of Fig. 8, we run the same numerical instance again by setting $t = 3/L$ in ManPG and AManPG. In this case, ManPG converges slower as seen from its curve being slightly shifted right. With some fluctuations, the accelerated method AManPG-5 manages to converge while the convergence is even slower than ManPG. AManPG-10 simply does not converge.

In practice, if the Lipschitz constant L is unknown, it is expected that an approximated L may cause a lot of trouble in ManPG’s line-search criterion (16) as used in [29], [40]. In the proposed line-search, this is never a problem. Intuitively, we fix $\alpha_k = 1$ in criterion (16) and use proxy step-size t' to find a working t from the line-search criterion (25). This process is well-defined by the Lipschitz type assumption (*i.e.*, Assumption 1), and ensures the descent of the total cost $f(\cdot)$ by Theorem 2.

7.2.2 Computational Complexity

The proposed PGS methods are much faster than ManPG and AManPG. Both ManPG and AManPG rely on the SSNM method [32] to solve the non-smooth KKT system, thus we compare the computation time per iteration of the proposed proxy step-size technique with that of the SSNM method. We report the timing statistics per iteration in Fig. 9, with respect to the dimension of the \mathbf{A} matrix in problem (39). It is clearly seen that the proposed proxy step-size technique is substantially faster than the SSNM method. An ablation study shows that the computation time per iteration of the PGS, A-PGS and AM-PGS methods is mostly decided by the matrix-vector multiplication $\mathbf{A}\mathbf{x}_k$ used in evaluating the Euclidean gradient $\nabla g|_{\mathbf{x}} = 2\mathbf{A}\mathbf{x}_k$ and the cost function $g(\mathbf{x}_k) = \mathbf{x}_k^T \mathbf{A}\mathbf{x}_k$. The A-PGS and AM-PGS methods have the same per iteration complexity, while being twice more

expensive than the PGS method due to an extra evaluation at the auxiliary state \mathbf{y}_k .

7.3 Fundamental Matrix Estimation

We find $\lambda = 0.01$ works well in general, after normalizing the image points [3]. For this problem, we find that the PGS algorithm converges mostly within 20 iterations. For most of the cases, 10 iterations or even 5 are sufficient to reduce the last singular value σ_3 close enough to zero. Therefore aside from the PGS method (with full convergence), aiming for an efficient engineering design, we propose the following two truncated PGS algorithms:

- PGS₅. 5 PGS iterations, followed by a rank-2 rounding by setting $\sigma_3 = 0$.
- PGS₁₀. 10 PGS iterations, followed by a rank-2 rounding by setting $\sigma_3 = 0$.

We use two benchmark algorithms: a) the normalized eight point algorithm (denoted as 8pt), b) the global polynomial optimization [52] (denoted as Gp) with a formulation based on $\det(\mathbf{F}) = 0$.

We run the 8pt, PGS₅, PGS₁₀, PGS and Gp methods on a list of standard benchmarks, and report in Table 1 a) e_{dist} the distance between the epipolar line and the corresponding image feature point [53], and b) e_{rep} the reprojection error of the triangulated 3D points. Overall, the e_{dist} and e_{rep} statistics decrease consistently over the 8pt, PGS₅, PGS₁₀, PGS, and Gp methods. Table 1 shows that the PGS₅ and PGS₁₀ methods, with a close performance towards the PGS method, consistently outperform the 8pt method, and they can give almost the same accuracy as the global method Gp.

7.4 Correspondence Association

The regularization strength is set as $\lambda = -\sigma_{\min}(\mathbf{A})/(\sqrt{n} - 1) = \sigma_{\max}(\mathbf{M})/(\sqrt{n} - 1)$ with n the dimension of \mathbf{A} . This choice is motivated by the canonical basis vector \mathbf{e}_k^T . Since the diagonal elements of \mathbf{A} are zero by construction, the total cost at \mathbf{e}_k is $f(\mathbf{e}_k) = \mathbf{e}_k^T \mathbf{A} \mathbf{e}_k + \lambda \|\mathbf{e}_k\|_{\ell_1} = \lambda$. Noting that $1 \leq \|\mathbf{x}\|_{\ell_1} \leq \sqrt{n}$, we thus have the following relation:

$$\begin{aligned} f(\mathbf{x}) < f(\mathbf{e}_k) &\Leftrightarrow \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda \|\mathbf{x}\|_{\ell_1} \leq \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda \sqrt{n} < \lambda \\ &\Leftrightarrow \lambda \leq -\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\sqrt{n} - 1} \leq -\frac{\sigma_{\min}(\mathbf{A})}{\sqrt{n} - 1}, \end{aligned}$$

where the used λ is chosen as the largest possible value.

Our experimental setup is similar to the one used in Section 5.1 of [5]. We simulate a 2-dimensional point-cloud

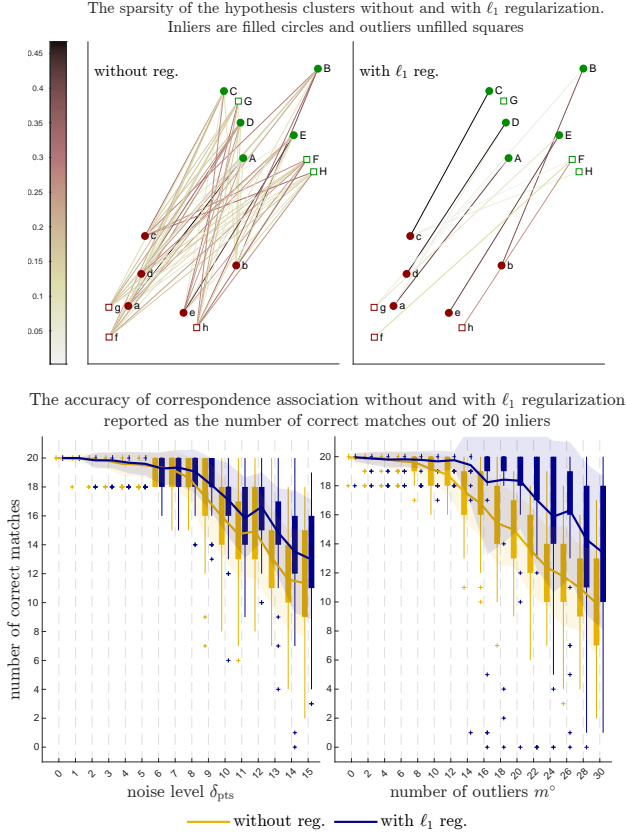


Fig. 10: Correspondence association with ℓ_1 norm regularization. With ℓ_1 regularization, high probability hypotheses are enhanced while low probability hypotheses are trimmed off, resulting in a sparser and more consistent cluster of association hypotheses. Therefore the correct correspondences (*i.e.*, inliers) can be identified more robustly.

\mathcal{Q} of m points, with $x - y$ coordinates uniformly distributed in $[0, 256\sqrt{m}/10]$. Then we add zero-mean Gaussian noise with standard deviation δ_{pts} to each point in \mathcal{Q} , and rotate and translate the whole \mathcal{Q} to obtain \mathcal{Q}' . We generate m^o outliers in both \mathcal{Q} and \mathcal{Q}' uniformly in the same region. We use $\delta_d = 5$ in equation (44) as in [5].

The results are reported in Fig. 10. In the first figure, we give an illustration (using $m = 5$) that with ℓ_1 norm regularization, many unlikely hypotheses are trimmed off, thus yielding a sparse hypothesis cluster. In the second figure, we set $m = 20$ and report the number of correct matches with respect to different noise levels and with respect to different outlier ratios. For each case, we use a 50 run Monte-Carlo simulation. It is clear that the association accuracy is consistently improved by using ℓ_1 norm regularization.

7.5 Linear Self-calibration

7.5.1 Implementation

We find normalization is essential to obtain stable self-calibration results. The key points of our implementation are sketched as follows:

- 1) Image point normalization. Let (p_x, p_y) be the principal point of the camera. If this is unknown, we approximate $(p_x, p_y) = (I_x/2, I_y/2)$, where I_x and

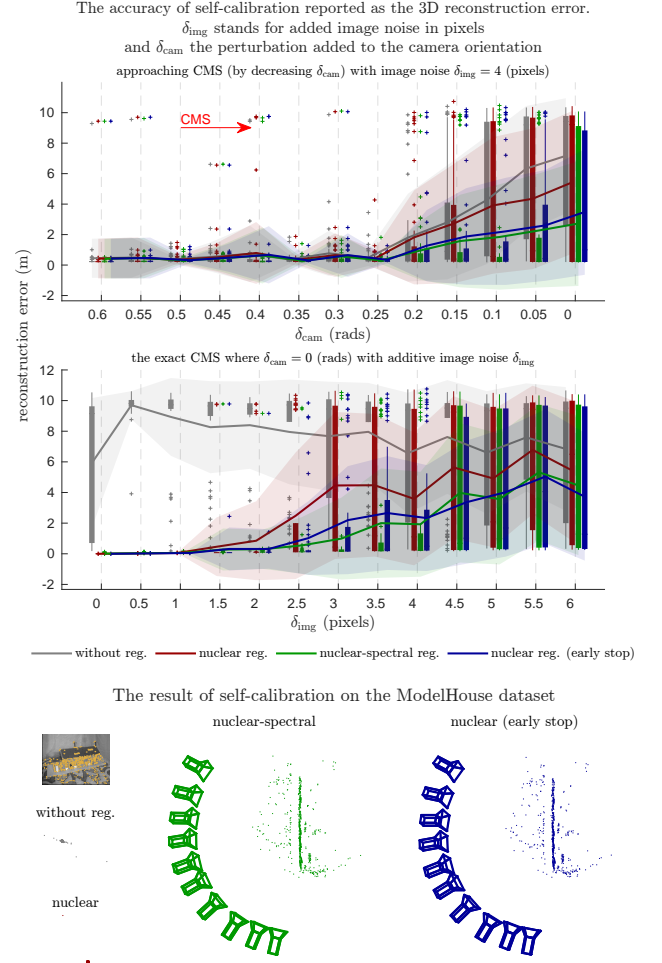


Fig. 11: Self-calibration with the nuclear norm and the nuclear-spectral norm regularization. We show the results with the critical motion sequence (CMS) where the camera principal axes pass approximately through the geometric center of the observed object (case R4 of Table 1 in [51]).

I_y are the width and height of the image. The average distance of all image points to (p_x, p_y) is denoted by s . We normalize all image points by a common transformation T^{-1} :

$$T = \begin{pmatrix} s & 0 & p_x \\ 0 & s & p_y \\ 0 & 0 & 1 \end{pmatrix}, \quad T^{-1} = \begin{pmatrix} 1/s & 0 & -p_x/s \\ 0 & 1/s & -p_y/s \\ 0 & 0 & 1 \end{pmatrix}.$$

- 2) Projective reconstruction using projective bundle-adjustment from the normalized image points.
- 3) Quasi-Euclidean rectification [54]. We approximate the intrinsic matrix of each camera computed in step 2) as $K_i = \text{diag}(f_i, f_i, 1)$, with $f_i = 2\sqrt{m_x^2 + m_y^2}$ where m_x and m_y are the maximum range in the x - and y - coordinates. Using this approximate K_i , we compute an approximate estimate of Q_∞^* from the DAQ constraint (48) *i.e.*, $\tilde{P}_i Q_\infty^* \tilde{P}_i^\top \propto K_i K_i^\top$, and rectify the projective reconstruction approximately.
- 4) Linear self-calibration from the quasi-Euclidean rectification in step 3). We initialize the PGS methods from the canonical DAQ $\Omega_\infty^* = \text{diag}(1, 1, 1, 0)$.

- 5) Transforming cameras P_i obtained from step 4) by T . Lastly, TP_i are the final estimate of Euclidean cameras.

We use $\lambda = 0.01$ for the nuclear norm regularization, and $\lambda_1 = 0.01$, $\lambda_2 = 2\lambda_1$ for the nuclear-spectral norm regularization.

7.5.2 Simulated Data

We simulate a scene comprising: a) 50 points spreading randomly in a diameter of 3 meters; b) 7 cameras circularly distributed 30 meters away from the point-cloud. All cameras are oriented towards the centroid of the point-cloud, thus the simulated scenario is an artificial CMS whose ambiguity can be removed using the rank deficiency of the DAQ. We evaluate the performance of each method with respect to the perturbation of camera orientations δ_{cam} and the noise of image points δ_{img} . We use the 3D reconstruction error as the evaluation metric, which is computed by the similarity Procrustes analysis between the ground-truth point-cloud and the estimated point-cloud.

We conduct two sets of experiments. First, we use the fixed image noise $\delta_{\text{img}} = 4$ pixels and decrease δ_{cam} to gradually bring the camera configuration to the CMS. Second, we set the camera configuration to the exact CMS where $\delta_{\text{cam}} = 0$ and then test the performance with respect to different image noise δ_{img} . We report the reconstruction error with a 50 run Monte-Carlo simulation in Fig. 11.

As shown in Fig. 11, when facing the CMS, the classical method without regularization fails, and the method with nuclear-spectral norm regularization is more robust than the one with nuclear norm regularization *e.g.*, in case of the exact CMS where $\delta_{\text{cam}} = 0$ and $\delta_{\text{img}} > 3$. It is interesting to see that the nuclear-norm regularization works well for less noisy scenarios of the CMS, *e.g.*, when $\delta_{\text{img}} < 2$. For these cases, it seems that the solution of the nuclear-norm regularized problem is well-trapped at the local minimum (the rank-3 DAQ), while the gradient is not large enough to go to the global minimum (the rank-1 DAQ). To examine this hypothesis, we use an early-stop trick (by setting the maximum PGS iterations to 1000) in the nuclear norm regularization, and observe that with the early-stop trick the nuclear-norm regularized method mostly performs well.

7.5.3 Real Data

A qualitative example of the CMS is given in Fig. 11 using a real dataset called ModelHouse where all cameras look towards the geometric center of a model house¹. In this scenario, matrix A has two eigenvalues close to zero. The classic method fails because the solution space is ambiguous. The nuclear norm regularized method fails by converging to the rank-1 solution. The nuclear-spectral norm regularized method can find the correct rank-3 solution thus recover the correct Euclidean geometry. The nuclear norm regularized method with the early stop trick also works. Intuitively, the spectral norm in the nuclear-spectral norm regularization prevents the algorithm from gliding to the rank-1 solution, and the early stop trick has the similar functionality.

1. <https://www.robots.ox.ac.uk/~vgg/data/mview/>

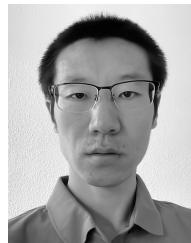
8 CONCLUSION

We have proposed the proxy step-size technique, and presented an effective solution to problem (2) for convex and absolutely homogeneous $h(\cdot)$. The proposed solution is: exact (satisfying the first-order necessary optimality condition), elegant (simple and in closed-form), and easily applicable (to nuclear norm regularization etc.). Future work includes extending the proxy-step size technique to the oblique and the Stiefel manifolds, and analyzing the convergence rate of the accelerated methods in Algorithm 3.

REFERENCES

- [1] O. Faugeras and B. Mourrain, "On the geometry and algebra of the point and line correspondences between n images," in *International Conference on Computer Vision*, 1995.
- [2] B. Triggs, "Autocalibration and the absolute quadric," in *Computer Vision and Pattern Recognition*, 1997.
- [3] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge University Press, ISBN: 0521540518, 2004.
- [4] A. Etz, "Introduction to the concept of likelihood and its applications," *Advances in Methods and Practices in Psychological Science*, vol. 1, no. 1, pp. 60–69, 2018.
- [5] M. Leordeanu and M. Hebert, "A spectral technique for correspondence problems using pairwise constraints," in *Computer Vision, IEEE International Conference on*, vol. 2. IEEE Computer Society, 2005, pp. 1482–1489.
- [6] A. M. Peter and A. Rangarajan, "Maximum likelihood wavelet density estimation with applications to image and shape matching," *IEEE Transactions on Image Processing*, vol. 17, no. 4, pp. 458–468, 2008.
- [7] R. Lai, Z. Wen, W. Yin, X. Gu, and L. M. Lui, "Folding-free global conformal mapping for genus-0 surfaces by harmonic energy minimization," *Journal of Scientific Computing*, vol. 58, no. 3, pp. 705–725, 2014.
- [8] J. Hu, X. Liu, Z.-W. Wen, and Y.-X. Yuan, "A brief introduction to manifold optimization," *Journal of the Operations Research Society of China*, vol. 8, no. 2, pp. 199–248, 2020.
- [9] M. Fazel, "Matrix rank minimization with applications," Ph.D. dissertation, Stanford University, 2002.
- [10] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media, 2003, vol. 87.
- [11] A. Beck, *First-order methods in optimization*. SIAM, 2017.
- [12] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [13] M. Fazel, T. K. Pong, D. Sun, and P. Tseng, "Hankel matrix rank minimization with applications to system identification and realization," *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 3, pp. 946–977, 2013.
- [14] Y. Nesterov, "Gradient methods for minimizing composite functions," *Mathematical Programming*, vol. 140, no. 1, pp. 125–161, 2013.
- [15] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- [16] N. Boumal, "An introduction to optimization on smooth manifolds," *Available online*, May, 2020.
- [17] P. Grohs and S. Hosseini, " ε -subgradient algorithms for locally Lipschitz functions on Riemannian manifolds," *Advances in Computational Mathematics*, vol. 42, no. 2, pp. 333–360, 2016.
- [18] —, "Nonsmooth trust region algorithms for locally Lipschitz functions on Riemannian manifolds," *IMA Journal of Numerical Analysis*, vol. 36, no. 3, pp. 1167–1192, 2016.
- [19] S. Hosseini and A. Uschmajew, "A Riemannian gradient sampling algorithm for nonsmooth optimization on manifolds," *SIAM Journal on Optimization*, vol. 27, no. 1, pp. 173–189, 2017.
- [20] S. Hosseini, W. Huang, and R. Yousefpour, "Line search algorithms for locally Lipschitz functions on Riemannian manifolds," *SIAM Journal on Optimization*, vol. 28, no. 1, pp. 596–619, 2018.
- [21] O. Ferreira and P. Oliveira, "Proximal point algorithm on Riemannian manifolds," *Optimization*, vol. 51, no. 2, pp. 257–270, 2002.

- [22] G. de Carvalho Bento, J. X. da Cruz Neto, and P. R. Oliveira, "A new approach to the proximal point method: convergence on general Riemannian manifolds," *Journal of Optimization Theory and Applications*, vol. 168, no. 3, pp. 743–755, 2016.
- [23] G. C. Bento, O. P. Ferreira, and J. G. Melo, "Iteration-complexity of gradient, subgradient and proximal point methods on Riemannian manifolds," *Journal of Optimization Theory and Applications*, vol. 173, no. 2, pp. 548–562, 2017.
- [24] R. Lai and S. Osher, "A splitting method for orthogonality constrained problems," *Journal of Scientific Computing*, vol. 58, no. 2, pp. 431–449, 2014.
- [25] A. Kovnatsky, K. Glashoff, and M. M. Bronstein, "MADMM: a generic algorithm for non-smooth optimization on manifolds," in *European Conference on Computer Vision*. Springer, 2016, pp. 680–696.
- [26] W. Chen, H. Ji, and Y. You, "An augmented Lagrangian method for ℓ_1 -regularized optimization problems with orthogonality constraints," *SIAM Journal on Scientific Computing*, vol. 38, no. 4, pp. B570–B592, 2016.
- [27] H. Zhu, X. Zhang, D. Chu, and L.-Z. Liao, "Nonconvex and non-smooth optimization with generalized orthogonality constraints: An approximate augmented Lagrangian method," *Journal of Scientific Computing*, vol. 72, no. 1, pp. 331–372, 2017.
- [28] Y. Wang, W. Yin, and J. Zeng, "Global convergence of ADMM in nonconvex nonsmooth optimization," *Journal of Scientific Computing*, vol. 78, no. 1, pp. 29–63, 2019.
- [29] S. Chen, S. Ma, A. Man-Cho So, and T. Zhang, "Proximal gradient method for nonsmooth optimization over the Stiefel manifold," *SIAM Journal on Optimization*, vol. 30, no. 1, pp. 210–239, 2020.
- [30] W. Huang and K. Wei, "Riemannian proximal gradient methods," *Mathematical Programming*, pp. 1–43, 2021.
- [31] M. Tan, Z. Hu, Y. Yan, J. Cao, D. Gong, and Q. Wu, "Learning sparse pca with stabilized admm method on stiefel manifold," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 3, pp. 1078–1088, 2019.
- [32] X. Xiao, Y. Li, Z. Wen, and L. Zhang, "A regularized semi-smooth Newton method with projection steps for composite convex programs," *Journal of Scientific Computing*, vol. 76, no. 1, pp. 364–389, 2018.
- [33] Y. E. Nesterov, "A method for solving the convex programming problem with convergence rate $o(1/k^2)$," *Dokl. akad. nauk Sssr*, vol. 269, pp. 543–547, 1983.
- [34] G. Dirr, U. Helmke, and C. Lageman, "Nonsmooth Riemannian optimization with applications to sphere packing and grasping," in *Lagrangian and Hamiltonian methods for nonlinear control 2006*. Springer, 2007, pp. 29–45.
- [35] P. B. Borckmans, S. E. Selvan, N. Boumal, and P.-A. Absil, "A Riemannian subgradient algorithm for economic dispatch with valve-point effect," *Journal of computational and applied mathematics*, vol. 255, pp. 848–866, 2014.
- [36] A. Goldstein, "Optimization of Lipschitz continuous functions," *Mathematical Programming*, vol. 13, no. 1, pp. 14–22, 1977.
- [37] H. Zhang and S. Sra, "First-order methods for geodesically convex optimization," in *Conference on Learning Theory*. PMLR, 2016, pp. 1617–1638.
- [38] R. L. Bishop and B. O'Neill, "Manifolds of negative curvature," *Transactions of the American Mathematical Society*, vol. 145, pp. 1–49, 1969.
- [39] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Machine Learning*, vol. 3, no. 1, pp. 1–122, 2010.
- [40] W. Huang and K. Wei, "An extension of fast iterative shrinkage-thresholding algorithm to riemannian optimization for sparse principal component analysis," *Numerical Linear Algebra with Applications*, p. e2409, 2021.
- [41] G. A. Watson, "Characterization of the subdifferential of some matrix norms," *Linear algebra and its applications*, vol. 170, no. 0, pp. 33–45, 1992.
- [42] N. Boumal, P.-A. Absil, and C. Cartis, "Global rates of convergence for nonconvex optimization on manifolds," *IMA Journal of Numerical Analysis*, vol. 39, no. 1, pp. 1–33, 2019.
- [43] W. H. Yang, L.-H. Zhang, and R. Song, "Optimality conditions for the nonlinear programming problems on Riemannian manifolds," *Pacific Journal of Optimization*, vol. 10, no. 2, pp. 415–434, 2014.
- [44] Y. Liu, F. Shang, J. Cheng, H. Cheng, and L. Jiao, "Accelerated first-order methods for geodesically convex optimization on Riemannian manifolds." in *NIPS*, 2017, pp. 4868–4877.
- [45] H. Zhang and S. Sra, "An estimate sequence for geodesically convex optimization," in *Conference On Learning Theory*. PMLR, 2018, pp. 1703–1723.
- [46] A. Beck and M. Teboulle, "Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems," *IEEE transactions on image processing*, vol. 18, no. 11, pp. 2419–2434, 2009.
- [47] Y. Zheng, S. Sugimoto, and M. Okutomi, "A practical rank-constrained eight-point algorithm for fundamental matrix estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1546–1553.
- [48] J. Zhao, "An efficient solution to non-minimal case essential matrix estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [49] N. Parikh and S. Boyd, "Proximal algorithms," *Foundations and Trends in optimization*, vol. 1, no. 3, pp. 127–239, 2014.
- [50] M. Pollefeys, R. Koch, and L. Van Gool, "Self-calibration and metric reconstruction inspite of varying and unknown intrinsic camera parameters," *International Journal of Computer Vision*, vol. 32, no. 1, pp. 7–25, 1999.
- [51] P. Gurdjos, A. Bartoli, and P. Sturm, "Is dual linear self-calibration artificially ambiguous?" in *International Conference on Computer Vision*, 2009.
- [52] F. Bugarin, A. Bartoli, D. Henrion, J.-B. Lasserre, J.-J. Orteu, and T. Sentenac, "Rank-constrained fundamental matrix estimation by polynomial global optimization versus the eight-point algorithm," *Journal of Mathematical Imaging and Vision*, vol. 53, no. 1, pp. 42–60, 2015.
- [53] Z. Zhang, "Determining the epipolar geometry and its uncertainty: A review," *International Journal of Computer Vision*, vol. 27, no. 2, pp. 161–195, 1998.
- [54] P. A. Beardsley, A. Zisserman, and D. W. Murray, "Sequential updating of projective and affine structure from motion," *International journal of computer vision*, vol. 23, no. 3, pp. 235–259, 1997.



Fang Bai Fang Bai was born in Ningxia Province, China, in 1988. He received the B.Sc. degree in computer science and technology from Nankai University, China, in 2010, and the Ph.D. degree in robotics from University of Technology Sydney, Australia, in 2020. His research has been focused on mathematical abstractions in robotics and computer vision. He has conducted several fundamental breakthroughs on related topics as the first-author, e.g., the cycle based pose graph optimization, the equation to predict the change of optimal values, and the closed-form solution for template-free deformable Procrustes analysis.



Adrien Bartoli Adrien Bartoli has held the position of Professor of Computer Science at Université Clermont Auvergne since fall 2009 and has been a member of Institut Universitaire de France since 2016. He is currently on leave as research scientist at the University Hospital of Clermont-Ferrand and as Chief Scientific Officer at SurgAR. He leads the Endoscopy and Computer Vision (EnCoV) research group at the University and Hospital of Clermont-Ferrand. His main research interests are in computer vision, including image registration and Shape-from-X for deformable environments, and their application to computer-aided medical interventions.

APPENDIX A PROPERTIES OF PROXIMAL

A.1 Convexity

If $h(\cdot)$ is convex, then for any \mathbf{x} and \mathbf{x}_0 , we have:

$$h(\mathbf{x}) \geq h(\mathbf{x}_0) + \langle \partial h|_{\mathbf{x}_0}, \mathbf{x} - \mathbf{x}_0 \rangle.$$

A.2 Firm non-expansiveness and non-expansiveness

For a convex $h(\cdot)$, let $\mathbf{z}_1 = \text{prox}_{th}(\mathbf{w}_1)$ and $\mathbf{z}_2 = \text{prox}_{th}(\mathbf{w}_2)$ with $t \geq 0$. Then the following holds:

- Firm non-expansiveness:

$$\|\mathbf{z}_1 - \mathbf{z}_2\|_2^2 \leq \langle \mathbf{z}_1 - \mathbf{z}_2, \mathbf{w}_1 - \mathbf{w}_2 \rangle; \quad (55)$$

- Non-expansiveness:

$$\|\mathbf{z}_1 - \mathbf{z}_2\|_2^2 \leq \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2. \quad (56)$$

Proof. The solution of the proximal is characterized by its first-order necessary condition:

$$\mathbf{0} \in \partial h|_{\mathbf{z}_1} + \frac{1}{t}(\mathbf{z}_1 - \mathbf{w}_1) \Leftrightarrow \frac{1}{t}(\mathbf{w}_1 - \mathbf{z}_1) \in \partial h|_{\mathbf{z}_1}, \quad (57)$$

$$\mathbf{0} \in \partial h|_{\mathbf{z}_2} + \frac{1}{t}(\mathbf{z}_2 - \mathbf{w}_2) \Leftrightarrow \frac{1}{t}(\mathbf{w}_2 - \mathbf{z}_2) \in \partial h|_{\mathbf{z}_2}. \quad (58)$$

By the convexity of $h(\cdot)$, we have:

$$\begin{cases} h(\mathbf{z}_1) \geq h(\mathbf{z}_2) + \langle \partial h|_{\mathbf{z}_2}, \mathbf{z}_1 - \mathbf{z}_2 \rangle & (59a) \\ h(\mathbf{z}_2) \geq h(\mathbf{z}_1) + \langle \partial h|_{\mathbf{z}_1}, \mathbf{z}_2 - \mathbf{z}_1 \rangle. & (59b) \end{cases}$$

Substituting equations (57) and (58) into (59), we have:

$$\begin{cases} h(\mathbf{z}_1) \geq h(\mathbf{z}_2) + \langle \frac{1}{t}(\mathbf{w}_2 - \mathbf{z}_2), \mathbf{z}_1 - \mathbf{z}_2 \rangle & (60a) \\ h(\mathbf{z}_2) \geq h(\mathbf{z}_1) + \langle \frac{1}{t}(\mathbf{w}_1 - \mathbf{z}_1), \mathbf{z}_2 - \mathbf{z}_1 \rangle. & (60b) \end{cases}$$

Summing together inequalities (60a) and (60b), we have:

$$\langle \mathbf{w}_2 - \mathbf{z}_2 - \mathbf{w}_1 + \mathbf{z}_1, \mathbf{z}_1 - \mathbf{z}_2 \rangle \leq 0.$$

Expanding the above, we obtain the firm non-expansiveness (55). By the Cauchy–Schwarz inequality, we obtain:

$$|\langle \mathbf{z}_1 - \mathbf{z}_2, \mathbf{w}_1 - \mathbf{w}_2 \rangle| \leq \|\mathbf{z}_1 - \mathbf{z}_2\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2.$$

After canceling $\|\mathbf{z}_1 - \mathbf{z}_2\|_2$ in the firm non-expansiveness (55), we obtain non-expansiveness (56). \square

A.3 Proof of Lemma 2

By definition of the proximal, we write:

$$\text{prox}_{th}(\mathbf{0}) = \arg \min_{\mathbf{x}} \left\{ th(\mathbf{x}) + \frac{1}{2} \|\mathbf{x}\|^2 \right\}, \quad t \geq 0 \quad (61)$$

Since $h(\cdot)$ is convex and absolutely homogeneous, from Lemma 1, we have $h(\mathbf{x}) \geq 0$ and $h(\mathbf{0}) = 0$. Thus, we observe that the optimal cost of problem (61) is 0 which is attained at $\mathbf{x} = \mathbf{0}$ as $h(\mathbf{0}) = 0$ and $\|\mathbf{x}\| = 0$. This concludes $\text{prox}_{th}(\mathbf{0}) = \mathbf{0}$ if $h(\cdot)$ is convex and absolutely homogeneous. The inequalities $\|\text{prox}_{th}(\mathbf{w})\|_2^2 \leq \langle \text{prox}_{th}(\mathbf{w}), \mathbf{w} \rangle$ and $\|\text{prox}_{th}(\mathbf{w})\|_2 \leq \|\mathbf{w}\|_2$ follow from firm non-expansiveness (55) and non-expansiveness (56), by setting $\mathbf{w}_1 = \mathbf{w}$ and $\mathbf{w}_2 = \mathbf{0}$.

APPENDIX B DERIVATION OF KKT SYSTEM (15)

We write $\mathbf{0} \in \partial \mathcal{L}_{v_k}$ as follow:

$$\begin{aligned} \mathbf{0} &\in \text{grad } g|_{\mathbf{x}_k} + \frac{1}{t} \mathbf{v}_k + \partial h|_{\mathbf{x}_k + \mathbf{v}_k} + \mu \mathbf{x}_k \\ &\Leftrightarrow (1 - \mu t) \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k} \in t \partial h|_{\mathbf{x}_k + \mathbf{v}_k} + \mathbf{x}_k + \mathbf{v}_k \\ &\Leftrightarrow \mathbf{x}_k + \mathbf{v}_k = \text{prox}_{th}((1 - \mu t) \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k}). \end{aligned}$$

Therefore:

$$\mathbf{v}_k = \text{prox}_{th}((1 - \mu t) \mathbf{x}_k - t \text{grad } g|_{\mathbf{x}_k}) - \mathbf{x}_k,$$

which is equation (15a). Considering $\mathbf{x}_k^\top \mathbf{v}_k = 0$ with equation (15a), we obtain equation (15b).

APPENDIX C PROOF OF LEMMA 3

By definition of the proximal, we write:

$$\text{prox}_{th}(\alpha \mathbf{w}) = \arg \min_{\mathbf{x}} \left\{ h(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \alpha \mathbf{w}\|^2 \right\} \stackrel{\text{def}}{=} \mathbf{z}.$$

Since $h(\cdot)$ is absolutely homogeneous, $h(\mathbf{x}) = |\alpha| h(\frac{1}{\alpha} \mathbf{x})$. The above equation can thus be written as:

$$\begin{aligned} \mathbf{z} &= \arg \min_{\mathbf{x}} \left\{ |\alpha| h\left(\frac{1}{\alpha} \mathbf{x}\right) + \frac{\alpha^2}{2t} \left\| \frac{1}{\alpha} \mathbf{x} - \mathbf{w} \right\|^2 \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ h\left(\frac{1}{\alpha} \mathbf{x}\right) + \frac{1}{2 \frac{t}{|\alpha|}} \left\| \frac{1}{\alpha} \mathbf{x} - \mathbf{w} \right\|^2 \right\} \\ &\Leftrightarrow \frac{1}{\alpha} \mathbf{z} = \text{prox}_{\frac{t}{|\alpha|} h}(\mathbf{w}). \end{aligned}$$

Therefore $\mathbf{z} = \text{prox}_{th}(\alpha \mathbf{w}) = \alpha \text{prox}_{\frac{t}{|\alpha|} h}(\mathbf{w})$.

APPENDIX D PROOF OF LEMMA 4

From equation (17b), we obtain:

$$\begin{aligned} \frac{1}{\phi(t')} &= \frac{1}{t} = \frac{1}{t'} \mathbf{x}_k^\top \text{prox}_{|t'|h} \left(t' \left(\frac{1}{t'} \mathbf{x}_k - \text{grad } g|_{\mathbf{x}_k} \right) \right) \\ &= \mathbf{x}_k^\top \text{prox}_h \left(\frac{1}{t'} \mathbf{x}_k - \text{grad } g|_{\mathbf{x}_k} \right). \end{aligned}$$

The last equality is due to Lemma 3. We denote $\mathbf{w}_1 = \frac{1}{t'} \mathbf{x}_k - \text{grad } g|_{\mathbf{x}_k}$ and $\mathbf{w}_2 = \frac{1}{t'} \mathbf{x}_k - \text{grad } g|_{\mathbf{x}_k}$. Then $1/t'_1 = \mathbf{x}_k^\top \mathbf{w}_1$ and $1/t'_2 = \mathbf{x}_k^\top \mathbf{w}_2$. Moreover:

$$\mathbf{x}_k \mathbf{x}_k^\top (\mathbf{w}_1 - \mathbf{w}_2) = \frac{1}{t'_1} \mathbf{x}_k - \frac{1}{t'_2} \mathbf{x}_k = \mathbf{w}_1 - \mathbf{w}_2. \quad (62)$$

The proof is given by summarizing the above facts as:

$$\begin{aligned} \epsilon(t'_1, t'_2) &= \langle \mathbf{x}_k^\top \text{prox}_h(\mathbf{w}_1) - \mathbf{x}_k^\top \text{prox}_h(\mathbf{w}_2), \mathbf{x}_k^\top (\mathbf{w}_1 - \mathbf{w}_2) \rangle \\ &= \langle \text{prox}_h(\mathbf{w}_1) - \text{prox}_h(\mathbf{w}_2), \mathbf{x}_k \mathbf{x}_k^\top (\mathbf{w}_1 - \mathbf{w}_2) \rangle \\ &= \langle \text{prox}_h(\mathbf{w}_1) - \text{prox}_h(\mathbf{w}_2), \mathbf{w}_1 - \mathbf{w}_2 \rangle \\ &\geq \|\text{prox}_h(\mathbf{w}_1) - \text{prox}_h(\mathbf{w}_2)\|_2^2, \end{aligned}$$

where the inequality is due to the firm non-expansiveness of $\text{prox}_h(\cdot)$.

APPENDIX E**PROOF OF LEMMA 5****E.1 A General Inequality by Convexity**

Lemma 7. Let $h(\cdot)$ be convex. Then for any t , \mathbf{x} and \mathbf{w} , we have:

$$\begin{aligned} \langle \mathbf{w} - \text{prox}_{|t|h}(\mathbf{w}), \text{prox}_{|t|h}(\mathbf{w}) - \mathbf{x} \rangle \\ \geq |t| \left(h(\text{prox}_{|t|h}(\mathbf{w})) - h(\mathbf{x}) \right). \end{aligned} \quad (63)$$

Proof. Let $\mathbf{z} = \text{prox}_{|t|h}(\mathbf{w})$. By definition of the proximal, we write:

$$\mathbf{z} = \arg \min_{\mathbf{y}} \left\{ h(\mathbf{y}) + \frac{1}{2|t|} \|\mathbf{y} - \mathbf{w}\|^2 \right\}. \quad (64)$$

The first-order necessary condition of problem (64) states:

$$-\frac{1}{|t|}(\mathbf{z} - \mathbf{w}) \in \partial h|_{\mathbf{z}}. \quad (65)$$

By the convexity of $h(\cdot)$ at \mathbf{z} , the following inequality holds for any \mathbf{x} :

$$h(\mathbf{x}) \geq h(\mathbf{z}) + \langle \partial h|_{\mathbf{z}}, \mathbf{x} - \mathbf{z} \rangle.$$

Therefore for any \mathbf{x} , we have the following inequality:

$$h(\mathbf{x}) \geq h(\mathbf{z}) - \frac{1}{|t|} \langle \mathbf{z} - \mathbf{w}, \mathbf{x} - \mathbf{z} \rangle.$$

Reorganizing this inequality, we obtain inequality (63). \square

E.2 Proof of Lemma 5

Following Lemma 7, in inequality (63), we let $\mathbf{x} = \alpha \text{prox}_{|t|h}(\mathbf{w})$:

$$\begin{aligned} \langle \mathbf{w} - \text{prox}_{|t|h}(\mathbf{w}), (1 - \alpha) \text{prox}_{|t|h}(\mathbf{w}) \rangle \\ \geq |t| (1 - |\alpha|) h(\text{prox}_{|t|h}(\mathbf{w})). \end{aligned} \quad (66)$$

If $0 \leq \alpha < 1$, then $1 - |\alpha| = 1 - \alpha > 0$:

$$\langle \mathbf{w} - \text{prox}_{|t|h}(\mathbf{w}), \text{prox}_{|t|h}(\mathbf{w}) \rangle \geq |t| h(\text{prox}_{|t|h}(\mathbf{w})). \quad (67)$$

If $\alpha > 1$, then $1 - |\alpha| = 1 - \alpha < 0$:

$$\langle \mathbf{w} - \text{prox}_{|t|h}(\mathbf{w}), \text{prox}_{|t|h}(\mathbf{w}) \rangle \leq |t| h(\text{prox}_{|t|h}(\mathbf{w})). \quad (68)$$

By considering inequalities (67) and (68) together, we obtain:

$$\langle \mathbf{w} - \text{prox}_{|t|h}(\mathbf{w}), \text{prox}_{|t|h}(\mathbf{w}) \rangle = |t| h(\text{prox}_{|t|h}(\mathbf{w})). \quad (69)$$

Subtracting equation (69) in inequality (63), we obtain inequality (21).

APPENDIX F**PROOF OF PROPOSITION 3**

It can be shown that $t = t'/c(t')$ satisfies:

$$\begin{aligned} |t| &= \frac{|t'|}{\left| \mathbf{x}_k^\top \text{prox}_{|t'h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k}) \right|} \\ &\geq \frac{|t'|}{\|\mathbf{x}_k\|_2 \left\| \text{prox}_{|t'h}(\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k}) \right\|_2} \quad (\text{Cauchy-Schwarz}) \\ &\geq \frac{|t'|}{\|\mathbf{x}_k\|_2 \|\mathbf{x}_k - t' \text{grad } g|_{\mathbf{x}_k}\|_2} \quad (\text{Lemma 2}) \\ &= \frac{|t'|}{\|\mathbf{x}_k\|_2 \sqrt{\|\mathbf{x}_k\|_2^2 + \|t' \text{grad } g|_{\mathbf{x}_k}\|_2^2}} \quad (\text{since } \mathbf{x}_k^\top \text{grad } g|_{\mathbf{x}_k} = 0) \\ &= \frac{1}{\sqrt{(1/t')^2 + \|\text{grad } g|_{\mathbf{x}_k}\|_2^2}} \quad (\text{since } \|\mathbf{x}_k\|_2 = 1). \end{aligned}$$

APPENDIX G**PROOF OF LEMMA 6**

Since $\langle \mathbf{x}_k, \mathbf{v}_k \rangle = 0$ as $\mathbf{v}_k \in \mathcal{T}_{\mathbf{x}_k} \mathcal{S}$, the norm satisfies:

$$\begin{aligned} \|\mathbf{x}_k + \mathbf{v}_k\|_2 &= \sqrt{\|\mathbf{x}_k\|_2^2 + \|\mathbf{v}_k\|_2^2 + 2\langle \mathbf{x}_k, \mathbf{v}_k \rangle} \\ &= \sqrt{\|\mathbf{x}_k\|_2^2 + \|\mathbf{v}_k\|_2^2} \geq \sqrt{\|\mathbf{x}_k\|_2^2} = 1. \end{aligned}$$

By retraction (12), since $h(\cdot)$ is absolutely homogeneous, we write:

$$\begin{aligned} h(\mathcal{R}_{\mathbf{x}_k}(\mathbf{v}_k)) &= h\left(\frac{\mathbf{x}_k + \mathbf{v}_k}{\|\mathbf{x}_k + \mathbf{v}_k\|_2}\right) \\ &= \frac{1}{\|\mathbf{x}_k + \mathbf{v}_k\|_2} h(\mathbf{x}_k + \mathbf{v}_k) \leq h(\mathbf{x}_k + \mathbf{v}_k). \end{aligned}$$

APPENDIX H**LIPSCHITZ-TYPE CONSTANT L**

It is well-known in the Rayleigh quotient literature that:

$$\max_{\mathbf{v}} \frac{\mathbf{v}^\top \mathbf{A} \mathbf{v}}{\mathbf{v}^\top \mathbf{v}} = \sigma_{\max}(\mathbf{A}),$$

where $\sigma_{\max}(\mathbf{A})$ is the largest singular value of \mathbf{A} . Therefore for any \mathbf{v} , we have:

$$\frac{\mathbf{v}^\top \mathbf{A} \mathbf{v}}{\mathbf{v}^\top \mathbf{v}} \leq \sigma_{\max}(\mathbf{A}) \Leftrightarrow \mathbf{v}^\top \mathbf{A} \mathbf{v} \leq \sigma_{\max}(\mathbf{A}) \|\mathbf{v}\|_2^2. \quad (70)$$

From equation (38) and the fact that $\mathbf{x}^\top \mathbf{v} = 0$, we have:

$$\langle \text{grad } g(\mathbf{x}), \mathbf{v} \rangle = \langle 2\mathbf{A}\mathbf{x} - 2(\mathbf{x}^\top \mathbf{A}\mathbf{x})\mathbf{x}, \mathbf{v} \rangle = 2\mathbf{x}^\top \mathbf{A} \mathbf{v}. \quad (71)$$

Recall that $\frac{1}{\|\mathbf{x}+\mathbf{v}\|_2^2} = \frac{1}{1+\|\mathbf{v}\|_2^2} \leq 1$. For $g(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ and retraction (12), we write

$$\begin{aligned} g(\mathcal{R}_{\mathbf{x}}(\mathbf{v})) &= \left(\frac{\mathbf{x} + \mathbf{v}}{\|\mathbf{x} + \mathbf{v}\|_2} \right)^\top \mathbf{A} \left(\frac{\mathbf{x} + \mathbf{v}}{\|\mathbf{x} + \mathbf{v}\|_2} \right) \\ &= \frac{1}{\|\mathbf{x} + \mathbf{v}\|_2^2} \left(\mathbf{x}^\top \mathbf{A} \mathbf{x} + 2\mathbf{x}^\top \mathbf{A} \mathbf{v} + \mathbf{v}^\top \mathbf{A} \mathbf{v} \right) \\ &\leq \mathbf{x}^\top \mathbf{A} \mathbf{x} + 2\mathbf{x}^\top \mathbf{A} \mathbf{v} + \mathbf{v}^\top \mathbf{A} \mathbf{v} \\ &\leq g(\mathbf{x}) + \langle \text{grad } g(\mathbf{x}), \mathbf{v} \rangle + \frac{2\sigma_{\max}(\mathbf{A})}{2} \|\mathbf{v}\|_2^2. \end{aligned}$$

Therefore we obtain $L = 2\sigma_{\max}(\mathbf{A})$.