
Scale Invariant Power Iteration

Cheolmin Kim
Northwestern University
Evanston, IL
cheolminkim2019@u.northwestern.edu

Youngseok Kim
University of Chicago
Chicago, IL
youngseok@uchicago.edu

Diego Klabjan
Northwestern University
Evanston, IL
d-klabjan@northwestern.edu

Abstract

Power iteration has been generalized to solve many interesting problems in machine learning and statistics. Despite its striking success, theoretical understanding of when and how such an algorithm enjoys good convergence property is limited. In this work, we introduce a new class of optimization problems called scale invariant problems and prove that they can be efficiently solved by scale invariant power iteration (SCI-PI) with a generalized convergence guarantee of power iteration. By deriving that a stationary point is an eigenvector of the Hessian evaluated at the point, we show that scale invariant problems indeed resemble the leading eigenvector problem near a local optimum. Also, based on a novel reformulation, we geometrically derive SCI-PI which has a general form of power iteration. The convergence analysis shows that SCI-PI attains local linear convergence with a rate being proportional to the top two eigenvalues of the Hessian at the optimum. Moreover, we discuss some extended settings of scale invariant problems and provide similar convergence results for them. In numerical experiments, we introduce applications to independent component analysis, Gaussian mixtures, and non-negative matrix factorization. Experimental results demonstrate that SCI-PI is competitive to state-of-the-art benchmark algorithms and often yield better solutions.

1 Introduction

We consider a generalization of power iteration for finding the leading eigenvector of a matrix A . In power iteration, the update rule $x_{k+1} \leftarrow Ax_k / \|Ax_k\|$ is repeatedly applied until some stopping criterion is satisfied. Since no hyperparameter is required, this update rule is very practical yet attains global linear convergence with the rate of $|\lambda_2|/|\lambda_1|$ where $|\lambda_i|$ is the i^{th} largest absolute eigenvalue of A . This convergence result is analogous to that of gradient descent for convex optimization. Therefore, many variants including coordinate-wise [13], accelerated [23], stochastic [17], stochastic variance-reduced (VR) [19, 20], and stochastic VR heavy ball [10] power iterations have been developed, drawing a parallel literature to gradient descent for convex optimization. Also, a general form of power iteration has been used to solve

$$\text{maximize } f(x) \quad \text{subject to } x \in \partial\mathcal{B}_d \triangleq \{x \in \mathbb{R}^d : \|x\| = 1\} \quad (1)$$

in many applications such as sparse principal component analysis (PCA) [8, 16], L_1 -norm kernel PCA [9], phase synchronization [15], and the Burer-Monteiro factorization of semi-definite programs [3]. (All norms are 2-norms unless indicated otherwise.) Nevertheless, theoretical understanding of when such an algorithm enjoys the attractive convergence property of power iteration is limited.

While convex f is considered in [8], only global sublinear convergence is shown, not generalizing the appealing linear convergence property of power iteration.

In this work, we introduce a new class of optimization problems called *scale invariant problems* and show that they can be efficiently solved by *scale invariant power iteration* (SCI-PI) with a generalized convergence guarantee of power iteration. Scale invariant problems consider *scale invariant functions* in (1). We say that $f(x)$ is scale invariant, which is rigorously defined later, if its geometric surface is invariant under constant multiplication of x . Many important optimization problems in statistics and machine learning can be formulated as scale invariant problems, for instance, L_p -norm kernel PCA and maximum likelihood estimation of mixture proportions, to name a few. Moreover, as studied herein, independent component analysis (ICA), non-negative matrix factorization (NMF), and Gaussian mixture models (GMM) can be formulated as extended settings of scale invariant problems.

Derivatives of scale invariant functions have the interesting relation that $\nabla^2 f(x)x = k\nabla f(x)$ holds for some k . Using the KKT condition, we derive an eigenvector property stating that any stationary point x^* satisfying $\nabla f(x^*) = \lambda^*x^*$ for some λ^* is an eigenvector of $\nabla^2 f(x^*)$. Due to the eigenvector property, scale invariant problems can be locally seen as the leading eigenvector problem. Therefore, we can expect that a simple update rule like power iteration would efficiently solve scale invariant problems near a local optimum x^* . Another interesting property of scale invariant problems is that by swapping the objective function and the constraint, a geometrically interpretable dual problem with the goal of finding the closest point w to the origin from the constraint $f(w) = 1$ is obtained. By mapping an iterate x_k to the dual space, taking a descent step in the dual space and mapping it back to the original space, we geometrically derive SCI-PI, which replaces Ax_k with $\nabla f(x_k)$ in power iteration. We show that SCI-PI converges to a local maximum x^* at a linear rate when initialized close to it. The convergence rate is proportional to $\bar{\lambda}_2 / \lambda^*$ where $\bar{\lambda}_2$ is the spectral norm of $\nabla^2 f(x^*)(I - x^*(x^*)^T)$ and λ^* is the Lagrange multiplier corresponding to x^* , generalizing the convergence rate of power iteration. Moreover, under some mild conditions, we provide an explicit expression regarding the initial condition on $\|x_0 - x^*\|$ to ensure linear convergence.

In the extended settings, we discuss three variants of scale invariant problems. In the first setting, f is replaced with a sum of scale invariant functions. This setting covers a Kurtosis-based ICA and can be solved by SCI-PI with similar convergence guarantees. We also consider a block version of scale invariant problems which covers NMF and the Burer-Monteiro factorization of semi-definite programs. To solve block scale invariant problems, we present a block version of SCI-PI and show that it attains linear convergence in a two-block case. Lastly, we consider partially scale invariant problems which include general mixture problems such as GMM. To solve partially scale invariant problems, we present an alternative algorithm based on SCI-PI and the gradient method and prove its local linear convergence. In numerical experiments, we benchmark the proposed algorithms against state-of-the-art methods for KL-NMF, GMM and ICA. The experimental results show that our algorithms are computationally competitive and result in better solutions in select cases.

Our work has the following contributions.

1. We introduce scale invariant problems which cover interesting examples in statistics and machine learning yet can be efficiently solved by SCI-PI due to the eigenvector property.
2. We present a geometric derivation of SCI-PI using a dual reformulation and provide a convergence analysis for it. We show that SCI-PI converges to a local maximum x^* at a linear rate when initialized close to x^* , generalizing the attractive convergence property of power iteration. Moreover, we introduce 3 extended settings of scale invariant problems together with their convergence analyses.
3. We report numerical experiments including a novel reformulation of KL-NMF to extended settings of scale invariant problems. The experimental results demonstrate that SCI-PI are not only computationally competitive to state-of-the-art methods but also often yield better solutions.

The paper is organized as follows. In Section 2, we define scale invariance and present interesting properties of scale invariant problems including an eigenvector property and a dual formulation. We then provide a geometric derivation of SCI-PI and a convergence analysis in Section 3. The extended settings are discussed in Section 4 and we report the numerical experiments in Section 5. All proofs are deferred to Supplementary Material.

2 Scale Invariant Problems

Before presenting properties of scale invariant problems, we first define scale invariant functions.

Definition 2.1. We say that a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is multiplicatively scale invariant if it satisfies

$$f(cx) = u(c)f(x) \quad (2)$$

and $f : \mathbb{R}^d \setminus \{0\} \rightarrow \mathbb{R}$ is additively scale invariant if it satisfies

$$f(cx) = f(x) + v(c) \quad (3)$$

for some even functions $u : \mathbb{R} \rightarrow \mathbb{R}^+$ with $u(0) = 0$ and $v : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R}$ with $v(1) = 0$.

The following proposition characterizes the exact form of u and v for continuous f .

Proposition 2.2. If a continuous function $f \neq 0$ satisfies (2) with a multiplicative factor u , then we have $u(c) = |c|^p$ for some $p > 0$. Also, if a continuous function f satisfies (3) with an additive factor v , then we have $v(c) = \log_a |c|$ for some a where $0 < a$, $a \neq 1$.

Next, we establish derivative-based properties of scale invariant functions.

Proposition 2.3. Suppose f is twice differentiable. If f satisfies (2) with $u(c) = |c|^p$, we have

$$c\nabla f(cx) = |c|^p \nabla f(x), \quad \nabla f(x)^T x = pf(x), \quad \nabla^2 f(x)x = (p-1)\nabla f(x). \quad (4)$$

Also, if f satisfies (3) with $v(c) = \log_a |c|$, we have

$$c\nabla f(cx) = \nabla f(x), \quad \nabla f(x)^T x = \log^{-1}(a), \quad \nabla^2 f(x)x = -\nabla f(x). \quad (5)$$

The interesting relation that $\nabla^2 f(x) = k\nabla f(x)$ holds for some k is presented in Proposition 2.3. Using the first-order optimality conditions, we derive an eigenvector property as follows.

Proposition 2.4. Suppose that f is twice differentiable and let (λ^*, x^*) be a stationary point of (1) such that $\nabla f(x^*) = \lambda^* x^*$. If f satisfies (2) with $u(c) = |c|^p$, then we have $\nabla^2 f(x^*)x^* = (p-1)\lambda^* x^*$. Also, if f satisfies (3) with $v(c) = \log_a |c|$, then we have $\nabla^2 f(x^*)x^* = -\lambda^* x^*$. In both cases, x^* is an eigenvector of $\nabla^2 f(x^*)$. Moreover, if λ^* is greater than the remaining eigenvalues of $\nabla^2 f(x^*)$, then x^* is a local maximum to (1).

Proposition 2.4 states that a stationary point x^* is an eigenvector of $\nabla^2 f(x^*)$ and becomes a local maximum if the Lagrange multiplier λ^* is greater than the remaining eigenvalues of $\nabla^2 f(x^*)$. Due to this property, scale invariant problems can be considered as a generalization of the leading eigenvector problem. Next, we introduce a dual formulation of scale invariant problems.

Proposition 2.5. Suppose that a continuous function f is either multiplicatively scale invariant such that $f(x^*) > 0$ or additively scale invariant with an additive factor $u(c) = \log_a |c|$ with $a > 1$. Then, solving (1) is equivalent to solving the following optimization problem

$$\text{minimize } \|w\| \quad \text{subject to } f(w) = 1. \quad (6)$$

Note that a dual reformulation for a multiplicatively scale invariant f with $f(x^*) < 0$ or an additively scale invariant f with $0 < a < 1$ can be obtained by replacing $f(w) = 1$ with $f(w) = -1$ in (6). The dual formulation (6) has a nice geometric interpretation that an optimal solution w^* is the closest point to the origin from $\{w : f(w) = 1\}$. This understanding is used to derive SCI-PI in Section 3.

L_p -norm kernel PCA, estimation of mixture proportions, and KL-divergence NMF are all cases of scale invariant problems. The details of these cases are provided in Appendix A.1.

3 Scale Invariant Power Iteration

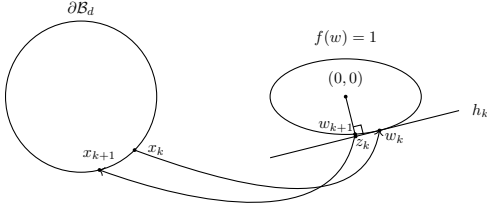
In this section, we provide a geometric derivation of SCI-PI to find a local optimal solution of (1). The algorithm is developed using the geometric interpretation of the dual formulation (6) as illustrated in Figure 1. Starting with an iterate $x_k \in \partial\mathcal{B}$, we obtain a dual iterate w_k by projecting x_k to the constraint $f(w) = 1$. Given w_k , we identify the hyperplane h_k which the current iterate w_k lies on and is tangent to $f(w) = 1$. After identifying the equation of h_k , we find the closest point z_k to

the origin from h_k and obtain a new dual iterate w_{k+1} by projecting z_k to the constraint $f(w) = 1$. Finally, we obtain a new primal iterate x_{k+1} by mapping w_{k+1} back to the set $\partial\mathcal{B}_d$.

Now, we develop an algorithm based on the above idea. For derivation of the algorithm, we assume that an objective function f is continuous and satisfies either (2) with $u(c) = |c|^p$ where $p > 0$ and $f(x) > 0$ for all $x \in \partial\mathcal{B}$ or (3) with $v(c) = \log_a|c|$ where $1 < a$. Under these conditions, a scalar mapping from x_k to w_k can be well defined as $w_k = x_k/f(x_k)^{1/p}$ or $w_k = a^{1-f(x_k)}x_k$, respectively. Let $w_k = c_k x_k$. Since w_k is on the constraint $f(w) = 1$, the tangent vector of the hyperplane h_k is $\nabla f(w_k)$. Therefore, we can write down the equation of the hyperplane h_k as $\{w : \nabla f(w_k)^T(w - w_k) = 0\}$. Note that z_k is a scalar multiple of $\nabla f(w_k)$ where the scalar can be determined from the requirement that z_k is on h_k . Since w_{k+1} is the projection of z_k , it must be a scalar multiple of the tangent vector $y_k = \nabla f(w_k)$. Therefore, we can write w_{k+1} as $w_{k+1} = d_k y_k$. Finally, by projecting w_{k+1} to $\partial\mathcal{B}$, we obtain

$$x_{k+1} = \frac{w_{k+1}}{\|w_{k+1}\|} = \frac{d_k y_k}{\|d_k y_k\|} = \frac{y_k}{\|y_k\|} = \frac{\nabla f(w_k)}{\|\nabla f(w_k)\|} = \frac{\nabla f(c_k x_k)}{\|\nabla f(c_k x_k)\|} = \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$$

where the last equality follows from Proposition 2.3. Summarizing all the above, we obtain SCI-PI presented in Algorithm 1.



Algorithm 1 SCI-PI

Input: initial point x_0
for $k = 0, 1, \dots, T-1$ **do**
 $x_{k+1} \leftarrow \nabla f(x_k) / \|\nabla f(x_k)\|$
end for
Output: x_T

Figure 1: Geometric derivation of SCI-PI

Next, we provide a convergence analysis of SCI-PI. For convex f , global sublinear convergence has been addressed in [8]. If f is not convex, $\{f(x_k)\}_{k=0,1,\dots}$ is no longer increasing, making it hard to analyze global convergence. Assuming that an initial point x_0 is close to a local maximum x^* , we study local convergence of SCI-PI. Note that $\|x^* - x_0\|^2 = 2(1 - x_0^T x^*)$ since $\|x^*\| = \|x_0\| = 1$.

Theorem 3.1. *Let f be a scale invariant, twice continuously differentiable function on an open set containing $\partial\mathcal{B}_d$ and let x^* be a local maximum satisfying $\nabla f(x^*) = \lambda^* x^*$ and $\lambda^* > \bar{\lambda}_2 = \max_{2 \leq i \leq d} |\lambda_i|$ where (λ_i, v_i) is an eigen-pair of $\nabla^2 f(x^*)$ with $x^* = v_1$. Then, there exists some $\delta > 0$ such that under the initial condition $1 - x_0^T x^* < \delta$, the sequence of iterates $\{x_k\}_{k=0,1,\dots}$ generated by SCI-PI satisfies*

$$1 - (x_k^T x^*)^2 \leq \prod_{t=0}^{k-1} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma_t \right)^2 (1 - (x_0^T x^*)^2) \quad \text{and} \quad \lim_{k \rightarrow \infty} \gamma_k = 0.$$

Moreover, if $\nabla_i f = \partial f / \partial x_i$ has a continuous Hessian H_i on an open set containing $\mathcal{B}_{d,\infty} \triangleq \{x \in \mathbb{R}^d : \|x\|_\infty \leq 1\}$, we can explicitly write δ as

$$\delta(\lambda^*, \bar{\lambda}_1, \bar{\lambda}_2, M) = \min \left\{ \left(\frac{\lambda^*}{M + \bar{\lambda}_1} \right)^2, \left(\frac{\lambda^* - \bar{\lambda}_2}{2M + \bar{\lambda}_1} \right)^2, 1 \right\}$$

where

$$\bar{\lambda}_1 = |\lambda_1|, \quad M = \max_{x \in \partial\mathcal{B}_d, y \in \mathcal{B}_{d,\infty}} \sqrt{\sum_{i=1}^d (x^T G_i(y) x)^2}, \quad G_i(y) = \sum_{j=1}^d v_{i,j} H_j(y).$$

Theorem 3.1 presents a local convergence result of SCI-PI with $\lambda^*/\bar{\lambda}_2$ generalizing the convergence rate of power iteration. Note that Theorem 3.1 requires $\lambda^* > \bar{\lambda}_2$ while it is sufficient to have $\lambda^* > \lambda_i$ for $2 \leq i \leq d$ to ensure local optimality. However, by adding $\sigma\|x\|^2$ for some $\sigma > 0$ to the objective function f , we can always enforce $\lambda^* > \bar{\lambda}_2$. On the other hand, by adding $\sigma\|x\|^2$ for some $\sigma < 0$, we may improve the convergence rate by increasing the relative gap between λ^* and $\bar{\lambda}_2$.

4 Extended Settings

4.1 Sum of Scale Invariant Functions

Consider a sum of scale invariant functions having the form of $f(x) = \sum_{i=1}^m g_i(x) + \sum_{j=1}^n h_j(x)$ where g_i is a multiplicatively scale invariant function with $u(c) = |c|^{p_i}$ and h_j is an additively scale invariant function with $v(c) = \log_{a_j} |c|$. Note that this does not imply that f is scale invariant in general. However, by Proposition 2.3, the gradient of f has the form of

$$\nabla f(x) = \sum_{i=1}^m \nabla g_i(x) + \sum_{j=1}^n \nabla h_j(x) = \left[\sum_{i=1}^m \left(\frac{1}{p_i - 1} \right) \nabla^2 g_i(x) - \sum_{j=1}^n \nabla^2 h_j(x) \right] x = F(x)x,$$

therefore a stationary point x^* satisfying $\nabla f(x^*) = \lambda^* x^*$ is an eigenvector of $F(x)$. We present a local convergence analysis of SCI-PI for a sum of scale invariant functions as follows.

Theorem 4.1. *Let f be a sum of scale invariant functions and twice continuously differentiable on an open set containing $\partial\mathcal{B}_d$ and let x^* be a local maximum x^* satisfying $\nabla f(x^*) = \lambda^* x^*$ and $\lambda^* > \bar{\lambda}_2 = \|\nabla^2 f(x^*)(I - x^*(x^*)^T)\|$. Then, there exists some $\delta > 0$ such that under the initial condition $1 - x_0^T x^* < \delta$, the sequence of iterates $\{x_k\}_{k=0,1,\dots}$ generated by SCI-PI satisfies*

$$1 - (x_k^T x^*)^2 \leq \prod_{t=0}^{k-1} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma_t \right)^2 (1 - (x_0^T x^*)^2) \text{ and } \lim_{k \rightarrow \infty} \gamma_k = 0.$$

Moreover, if $\nabla_i f = \partial f / \partial x_i$ has a continuous Hessian H_i on an open set containing $\mathcal{B}_{d,\infty}$,

$$\delta(\lambda^*, \bar{\lambda}_1, \bar{\lambda}_2, M) = \min \left\{ \left(\frac{\lambda^*}{M + \bar{\lambda}_1} \right)^2, \left(\frac{\lambda^* - \bar{\lambda}_2}{2M + \bar{\lambda}_1 + \bar{\lambda}_2} \right)^2, 1, \right\}$$

where

$$\bar{\lambda}_1 = \sqrt{2} \cdot \|\nabla^2 f(x^*)x^*\|, \quad M = \max_{x \in \partial\mathcal{B}_d, y \in \mathcal{B}_{d,\infty}} \sqrt{\sum_{i=1}^d (x^T G_i(y)x)^2}, \quad G_i(y) = \sum_{j=1}^d v_{i,j} H_j(y).$$

Note that $\bar{\lambda}_1$ has the additional $\sqrt{2}$ factor which comes from the fact that x^* is not necessarily an eigenvector of $\nabla^2 f(x^*)$. Nonetheless, the asymptotic convergence rate in Theorem 4.1 provides a generalization of the convergence rate in Theorem 3.1.

4.2 Block Scale Invariant Problems

Next, consider a class of optimization problems having the form of: $\max f(x, y)$ subject to $x \in \partial\mathcal{B}_{d_1}$, $y \in \partial\mathcal{B}_{d_2}$ where $f : \mathbb{R}^{d_1+d_2} \rightarrow \mathbb{R}$ is scale invariant in x for fixed y and vice versa. We derive the following alternating maximization algorithm called *block SCI-PI* as

$$x_{k+1} \leftarrow \nabla_x f(x, y_k) / \|\nabla_x f(x, y_k)\|, \quad y_{k+1} \leftarrow \nabla_y f(x_k, y) / \|\nabla_y f(x_k, y)\|. \quad (7)$$

Theorem 4.2. *Suppose that f is twice continuously differentiable on an open set containing $\partial\mathcal{B}_{d_1} \times \partial\mathcal{B}_{d_2}$ and let (x^*, y^*) be a local maximum satisfying $\nabla_x f(x^*, y^*) = \lambda^* x^*$, $\lambda^* > \bar{\lambda}_2 = \max_{2 \leq i \leq d_1} |\lambda_i|$, $\nabla_y f(x^*, y^*) = s^* y^*$, $s^* > \bar{s}_2 = \max_{2 \leq i \leq d_2} |s_i|$ where (λ_i, v_i) and (s_i, u_i) are eigen-pairs of $\nabla_x^2 f(x^*, y^*)$ and $\nabla_y^2 f(x^*, y^*)$, respectively such that $x^* = v_1$ and $y^* = u_1$. If $\nu^2 = \|\nabla_{yx} f(x^*, y^*)\|^2 < (\lambda^* - \bar{\lambda}_2)(s^* - \bar{s}_2)$ holds, then for the sequence of iterates $\{(x_k, y_k)\}_{k=0,1,\dots}$ generated by (7), there exists some $\delta > 0$ such that if $\max\{1 - x_0^T x^*, 1 - y_0^T y^*\} < \delta$, then we have $\|\Delta_k\| \leq \prod_{t=0}^{k-1} (\rho + \gamma_t) \|\Delta_0\|$ for some sequence γ_k such that $\lim_{k \rightarrow \infty} \gamma_k = 0$ where*

$$\Delta_k = \begin{bmatrix} \sqrt{1 - (x_k^T x^*)^2} \\ \sqrt{1 - (y_k^T y^*)^2} \end{bmatrix}, \quad \rho = \frac{1}{2} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \frac{\bar{s}_2}{s^*} + \sqrt{\left(\frac{\bar{\lambda}_2}{\lambda^*} - \frac{\bar{s}_2}{s^*} \right)^2 + \frac{4\nu^2}{\lambda^* s^*}} \right) < 1.$$

If x and y are independent ($\nu = 0$), we have $\rho = \max\{\bar{\lambda}_2/\lambda^*, \bar{s}_2/s^*\}$. Otherwise, ρ increases as ν increases. Note the result of Theorem 3.1 can be restored by dropping x or y in Theorem 4.2 and the algorithm and convergence analysis can be easily generalized to more than two blocks.

4.3 Partially Scale Invariant Problems

Lastly, we consider a class of optimization problems of the form: $\max f(x, y)$ subject to $x \in \partial\mathcal{B}_{d_1}$ where $f(x, y) : \mathbb{R}^{d_1+d_2} \rightarrow \mathbb{R}$ is a scale invariant function in x for each $y \in \mathbb{R}^{d_2}$. This problem has the form of (1) with respect to x once y is fixed. Also, by fixing x , we obtain an unconstrained optimization problem with respect to y . Using SCI-PI and the gradient method, an alternative maximization algorithm is derived as

$$x_{k+1} \leftarrow \nabla_x f(x_k, y_k) / \|\nabla_x f(x_k, y_k)\|, \quad y_{k+1} \leftarrow y_k + \alpha \nabla_y f(x_k, y_k). \quad (8)$$

While the gradient method is used in (8), any method for unconstrained optimization can replace it.

Theorem 4.3. *Suppose that $f(x, y)$ is scale invariant in x for each $y \in \mathbb{R}^{d_2}$, μ -strongly concave in y with an L -Lipschitz continuous $\nabla_y f(x, y)$ for each $x \in \partial\mathcal{B}_{d_1}$, and three-times continuously differentiable on an open set containing $\partial\mathcal{B}_{d_1} \times \mathbb{R}^{d_2}$. Let (x^*, y^*) be a local maximum satisfying $\nabla f(x^*) = \lambda^* x^*$ and $\lambda^* > \bar{\lambda}_2 = \max_{2 \leq i \leq d} |\lambda_i|$ where (λ_i, v_i) is an eigen-pair of $\nabla^2 f(x^*)$ with $x^* = v_1$. If $\nu^2 = \|\nabla_{yx}^2 f(x^*, y^*)\|^2 < \mu(\lambda^* - \bar{\lambda}_2)$ holds, then for the sequence of iterates $\{(x_k, y_k)\}_{k=0,1,\dots}$ generated by (8) with $\alpha = 2/(L + \mu)$, there exists some $\delta > 0$ such that if $\max\{1 - x_0^T x^*, \|y - y^*\|\} < \delta$, then we have $\|\Delta_k\| \leq \prod_{t=0}^{k-1} (\rho + \gamma_t) \|\Delta_0\|$ for some sequence γ_k such that $\lim_{k \rightarrow \infty} \gamma_k = 0$ where*

$$\Delta_k = \left[\begin{array}{c} \sqrt{1 - (x_k^T x^*)^2} \\ \|y_k - y^*\| \end{array} \right], \quad \rho = \frac{1}{2} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \frac{L - \mu}{L + \mu} + \sqrt{\left(\frac{\bar{\lambda}_2}{\lambda^*} - \frac{L - \mu}{L + \mu} \right)^2 + \frac{8\nu^2}{\lambda^*(L + \mu)}} \right) < 1.$$

As in the result of Theorem 4.2, the rate ρ increases as ν increases and is equal to $\max\{\bar{\lambda}_2/\lambda^*, (L - \mu)/(L + \mu)\}$ when $\nu = 0$. Also, by dropping y , we can restore the convergence result of Theorem 3.1.

5 Numerical Experiments

We tested the proposed algorithms on real-world data sets. All experiments were implemented on a standard laptop (2.6 GHz Intel Core i7 processor and 16GM memory) using the Julia programming language. Let us emphasize that scale invariant problems frequently appear in many important applications in statistics and machine learning. We select 3 important applications, KL-NMF, GMM and ICA. A description of the data sets is provided in Supplementary Material. All of them are standard sets used in prior works on these 3 problems.

KL-NMF: The KL-divergence NMF (KL-NMF) subproblem can be solved via SCI-PI (see Supplementary Material A.1). Our focus is to compare this algorithm with other famous alternating minimization algorithms listed below, updating H and W alternatively. To lighten the notation, let \odot , \oslash and $(\cdot)^{\odot 2}$ denote element-wise product, division and square, respectively. We let $z = V \oslash (Wh)$ and $\mathbb{1}_n$ denote a vector of ones.

- Projected gradient descent (PGD): It iterates $h^{\text{new}} \leftarrow h - \eta \odot W^T(z - \mathbb{1}_n)$ followed by projection onto the simplex, where $\eta \propto h$ is an appropriate learning rate [14].
- Multiplicative update (MU): A famous multiplicative update algorithm is originally suggested by [12], which iterates $h^{\text{new}} \leftarrow h \odot (W^T z) \oslash (W^T \mathbb{1}_n)$ and is learning rate free.
- Our method (SCI-PI): It iterates $h^{\text{new}} \leftarrow h \odot (\sigma + W^T z)^{\odot 2}$ and rescales h , where σ is a shift parameter. We simply use $\sigma = 1$ for preconditioning.
- Sequential quadratic programming (MIXSQP): Solving each subproblem via a convex solver `mixsqp` [11]. This algorithm performs sequential non-negative least squares.

To study the convergence rate for KL-NMF subproblems, we use four simulated data sets exhibited in [11]. We study MU, PGD and SCI-PI since they have the same order of computational complexity per iteration, but omit MIXSQP since it is a second-order method which cannot be directly compared. For PGD, the learning rate is optimized by grid search. The stopping criterion is $\|f(x_k) - f^*\| \leq 10^{-6} f^*$ where f^* is the solution obtained by MIXSQP after extensive computation time. The average runtime for aforementioned 3 methods are 33, 33 and 30 seconds for 10,000 iterations, respectively. The

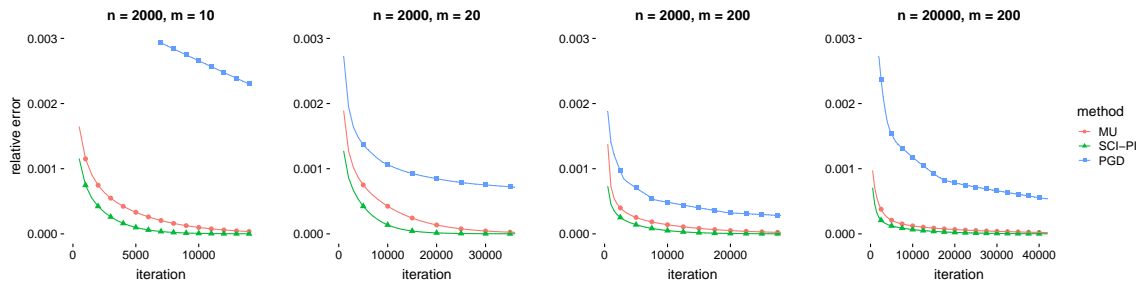


Figure 2: Convergence of 3 algorithms for the KL-NMF subproblem. n, m : the number of samples/features of the data matrix.

result is shown in Figure 2¹. It shows that SCI-PI outperforms the other 2 for all simulated data sets. Also, all methods seems to exhibit linear convergence.

Next, we test the 4 algorithms on 4 real-world data sets for 3 different purposes: 287 waving tree (WT) images for image reconstruction, two bag-of-words data sets from the KOS blog and NIPS full papers for topic modeling, and a Wikipedia (WIKI) vote network for graph clustering. We estimate $k = 20$ factors. At each iteration, all 4 algorithms solve m subproblems simultaneously.

The result is summarized in Figure 3². The convergence plots are based on the average relative errors over 10 repeated runs with random initializations. The result shows that SCI-PI is an overall winner, showing faster convergence rates. The stopping criterion is the same as above. To assess overall performance when initialized differently, we select KOS and WIKI and run MU, PGD, SCI-PI, and MIXSQP 10 times¹. The 3 algorithms except MIXSQP have (approximately) the same computational cost per iteration, take runtime of 391, 396, 408 seconds for KOS data and 372, 390, 418 seconds for WIKI data, respectively for 200 iterations. MIXSQP has a larger per iteration cost. After 400 seconds, SCI-PI achieves lowest objective values in all cases but one for each data set (38 out of 40 in total). Thus it clearly outperforms other methods and also achieves the lowest variance. Unlike the other 3 algorithms, SCI-PI is not an ascent algorithm but an eigenvalue-based fixed-point algorithm. We observe that sometimes SCI-PI converges to a better solution due to this fact. Admittedly, this can be potentially dangerous but for the KL-NMF problem its performance turns out to be stable.

GMM: GMM fits a mixture of Gaussian distributions to the underlying data. Let $L_{ik} = \mathcal{N}(x_i; \mu_k, \Sigma_k)$ where i is the sample index and k the cluster index and let π be the actual mixture proportion vector. GMM fits into our restricted scale invariant setting (Sec. 4.3) with reparametrization, but the gradient update for μ_k, Σ_k is replaced by the exact coordinate ascent step. The EM and SCI-PI updates for π can be written respectively as

$$r = \mathbb{1} \odot (L\pi), \quad \pi_k^{\text{new}} \propto \pi \odot (L^T r) \quad (\text{EM}), \quad \pi_k^{\text{new}} \propto \pi \odot (\alpha + L^T r)^{\odot 2}, \quad (\text{SCI-PI}). \quad (9)$$

We compare SCI-PI and EM for different real-world data sets. All the algorithms initialize from the same standard Gaussian random variable, repeatedly for 10 times. The result is summarized in the left panel in Figure 4. The stopping criterion is $\|x_{k+1} - x_k\| < 10^{-8}$. In some cases, SCI-PI achieves much larger objective values even if initialized the same. In many cases the 2 algorithms exhibit the same performance. This is because estimation of μ_k 's and Σ_k 's are usually harder than estimation of π , and EM and SCI-PI have the same updates for μ and Σ . For a few cases EM outperforms SCI-PI. Let us mention that SCI-PI and EM have the same order of computational complexity and require 591 and 590 seconds of total computation time, respectively.

ICA: We implement SCI-PI on the Kurtosis-based ICA problem [6] and compare it with the benchmark algorithm FastICA [5], which is the most popular algorithm. Given a pre-processed³ data matrix $W \in \mathbb{R}^{n \times d}$, we seek to maximize an approximated negative entropy $f(x) = \sum_{i=1}^n [(w_i^T x)^4 - 3]^2$

¹For each evaluation, we randomly draw 10 initial points and report the averaged relative errors with respect to f^* . The initial input for the KL-NMF problem is a one-step MU update of a $\text{Unif}(0, 1)$ random matrix.

²In all plots we do not show the first few iterations. The initial random solutions have the gap of approximately 50% which drops to a few percent after 10 iterates where the plots start.

³A centered matrix $\widetilde{W} = n^{1/2}UDV^T$ is pre-processed by $W = \widetilde{W}VD^{-1}V^T$ so that $W^T W = nVV^T$.

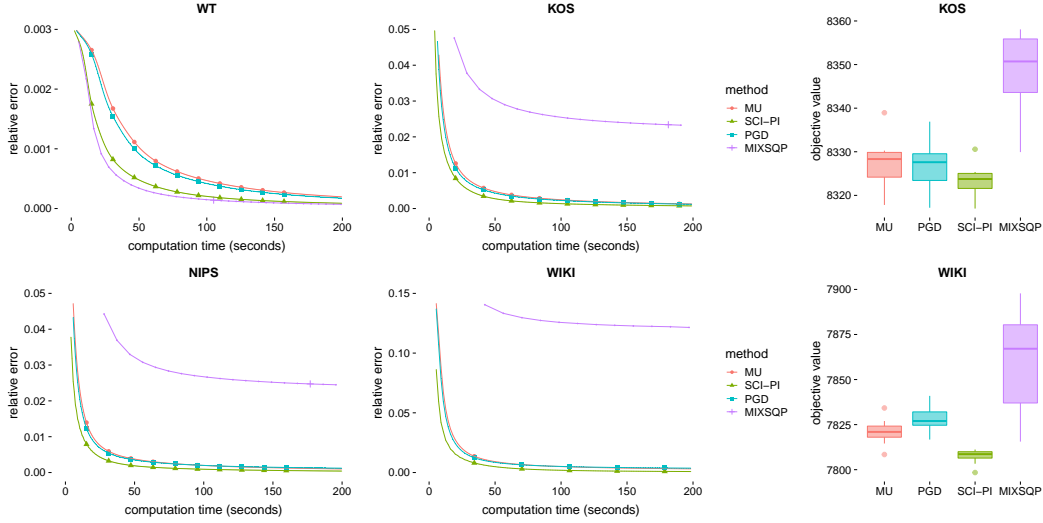


Figure 3: (Left and center) Convergence of the 4 NMF algorithms. (Right) Boxplots containing 10 objective values achieved after 400 seconds.

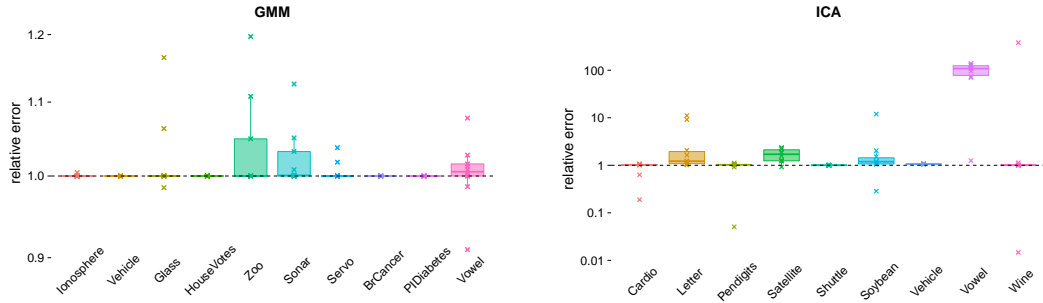


Figure 4: The relative error f_{SCI-PI}^*/f_{EM}^* for GMM (Left) and $f_{SCI-PI}^*/f_{FastICA}^*$ for ICA (Right).

subject to $x \in \partial\mathcal{B}_d$, for maximizing Kurtosis-based non-Gaussianity [7]. This problem fits into the sum of scale invariant setting (Sec. 4.1). SCI-PI iterates $x_{k+1} \leftarrow W^T[(Wx_k)^{\odot 4} - 3\mathbb{1}_n] \odot (Wx_k)^{\odot 3}$ and FastICA iterates $x_{k+1} \leftarrow W^T(Wx_k)^{\odot 3} - 3(\mathbb{1}^T(Wx_k)^{\odot 2})x_k$, both followed by normalization.

We compare SCI-PI and FastICA for different real-world data sets. The majority of data points (81 out of 100 in total) show that SCI-PI tends to find a better solution with a larger objective value, but in a few cases SCI-PI converges to a sub-optimal point. Both algorithms are fixed-point based and thus have no guarantee of global convergence but overall SCI-PI outperforms FastICA. SCI-PI and FastICA have the same order of computational complexity and require 11 and 12 seconds of total computation time, respectively.

6 Final Remarks

In this paper, we propose a new class of optimization problems called the scale invariant problems, together with a generic solver SCI-PI, which is indeed an eigenvalue-based fixed-point iteration. We showed that SCI-PI directly generalizes power iteration and enjoys similar properties, for instance, that SCI-PI has local linear convergence under mild conditions and its convergence rate is determined by eigenvalues of the Hessian matrix at a solution. Also, we extend scale invariant problems to problems with more general settings. Although scale invariance is a rather restrictive assumption, we show by experiments that SCI-PI can be a competitive option for numerous important problems such

as KL-NMF, GMM and ICA. Finding more examples and extending SCI-PI further to a more general setting is a promising direction for future studies.

References

- [1] Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [2] Sébastien Bubeck. Convex Optimization: Algorithms and Complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–357, 2015.
- [3] Murat A Erdogdu, Asuman Ozdaglar, Pablo A Parrilo, and Nuri Denizcan Vanli. Convergence Rate of Block-Coordinate Maximization Burer-Monteiro Method for Solving Large SDPs. *arXiv preprint arXiv:1807.04428*, 2018.
- [4] Cédric Févotte and Jérôme Idier. Algorithms for Nonnegative Matrix Factorization with the β -divergence. *Neural Computation*, 23(9):2421–2456, 2011.
- [5] Aapo Hyvarinen. Fast ICA for Noisy Data using Gaussian Moments. In *Proceedings of the 1999 IEEE International Symposium on Circuits and Systems VLSI*, volume 5, pages 57–61. IEEE, 1999.
- [6] Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. *Independent Component Analysis*, volume 46. John Wiley & Sons, 2004.
- [7] Aapo Hyvärinen and Erkki Oja. Independent Component Analysis: Algorithms and Applications. *Neural Networks*, 13(4-5):411–430, 2000.
- [8] Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized Power Method for Sparse Principal Component Analysis. *Journal of Machine Learning Research*, 11(Feb):517–553, 2010.
- [9] Cheolmin Kim and Diego Klabjan. A Simple and Fast Algorithm for L1-norm Kernel PCA. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [10] Cheolmin Kim and Diego Klabjan. Stochastic Variance-reduced Heavy Ball Power Iteration. *arXiv preprint arXiv:1901.08179*, 2019.
- [11] Youngseok Kim, Peter Carbonetto, Matthew Stephens, and Mihai Anitescu. A Fast Algorithm for Maximum Likelihood Estimation of Mixture Proportions Using Sequential Quadratic Programming. *arXiv preprint arXiv:1806.01412*, 2018.
- [12] Daniel D Lee and H Sebastian Seung. Algorithms for Non-negative Matrix Factorization. In *Advances in Neural Information Processing Systems*, pages 556–562, 2001.
- [13] Qi Lei, Kai Zhong, and Inderjit S Dhillon. Coordinate-wise Power method. In *Advances in Neural Information Processing Systems*, pages 2064–2072, 2016.
- [14] Chih-Jen Lin. Projected Gradient Methods for Non-negative Matrix Factorization. *Neural Computation*, 19(10):2756–2779, 2007.
- [15] Huikang Liu, Man-Chung Yue, and Anthony Man-Cho So. On the Estimation Performance and Convergence Rate of the Generalized Power Method for Phase Synchronization. *SIAM Journal on Optimization*, 27(4):2426–2446, 2017.
- [16] Ronny Luss and Marc Teboulle. Conditional Gradient Algorithms for Rank-One Matrix Approximations with a Sparsity Constraint. *SIAM REVIEW*, 55(1):65–98, 2013.
- [17] Erkki Oja. Simplified Neuron Model as a Principal Component Analyzer. *Journal of Mathematical Biology*, 15(3):267–273, 1982.
- [18] Prasanna K Sahoo and Palaniappan Kannappan. *Introduction to Functional Equations*. Chapman and Hall/CRC, 2011.
- [19] Ohad Shamir. A Stochastic PCA and SVD Algorithm with an Exponential Convergence Rate. In *International Conference on Machine Learning*, pages 144–152, 2015.
- [20] Ohad Shamir. Fast Stochastic Algorithms for SVD and PCA: Convergence Properties and Convexity. In *International Conference on Machine Learning*, pages 248–256, 2016.

- [21] Rong Wang, Feiping Nie, Xiaojun Yang, Feifei Gao, and Minli Yao. Robust 2DPCA With Non-greedy ℓ_1 -Norm Maximization for Image Analysis. *IEEE Transactions on Cybernetics*, 45(5):1108–1112, 2015.
- [22] Yu-Xiong Wang and Yu-Jin Zhang. Nonnegative Matrix Factorization: A Comprehensive Review. *IEEE Transactions on Knowledge and Data Engineering*, 25(6):1336–1353, 2013.
- [23] Peng Xu, Bryan He, Christopher De Sa, Ioannis Mitliagkas, and Chris Re. Accelerated Stochastic Power Iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 58–67, 2018.

A Supplementary Material

A.1 Examples

We introduce two immediate applications, L_p -norm kernel PCA and the mixture model, which have been intensively studied over the past few decades in the field of statistics and machine learning.

Example A.1 (L_p -norm kernel PCA). For $p > 0$, L_p -norm kernel PCA [21] is defined as

$$\text{maximize } f_p(x) = n^{-1} \sum_{i=1}^n \|\Phi(a_i)^T x\|_p^p \quad \text{subject to } x \in \mathcal{B}_d, \quad (10)$$

which satisfies property (2) with $u(c) = |c|^p$. The example includes the standard L_2 -norm PCA.

Example A.2 (Estimation of Mixture Proportions). Given a design matrix $L \in \mathbb{R}^{n \times d}$ satisfying $L_{jk} \geq 0$, the problem of estimating mixture proportions seeks to find a vector π of mixture proportions on the probability simplex $\mathcal{S}^d = \{\pi : \sum_{k=1}^d \pi_k = 1, \pi \geq 0\}$ that maximizes the log-likelihood $\sum_{j=1}^n \log \left(\sum_{k=1}^d L_{jk} \pi_k \right)$. We reformulate the problem by reparametrizing x_j by π_k^2 and obtain

$$\text{maximize } f_0(x) = n^{-1} \sum_{j=1}^n \log \left(\sum_{k=1}^d L_{jk} x_k^2 \right) \quad \text{subject to } x \in \mathcal{B}_d, \quad (11)$$

which now satisfies property (3) with $v(c) = 2 \log |c|$.

The reformulation idea in Example A.2 implies that any simplex-constrained problem with scale invariant f can be reformulated to a scale invariant problem. Example A.2 has a direct application to general mixture models, including the GMM [1]. The same optimization problem also appears in the Kullback-Leibler (KL) divergence NMF problem. In what follows, we show that the KL divergence NMF subproblem is indeed a scale invariant problem.

Example A.3 (KL-NMF). The KL-NMF problem [4, 12, 22] is defined as

$$\begin{aligned} \text{minimize } & D_{KL}(V \| WH) \triangleq \sum_{i,j} \left[V_{ij} \log \frac{V_{ij}}{\sum_k W_{ik} H_{kj}} - V_{ij} + \sum_k W_{ik} H_{kj} \right] \\ \text{subject to } & W_{ik} \geq 0, H_{kj} \geq 0, i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, K. \end{aligned} \quad (12)$$

Many popular algorithms for KL-NMF are based on alternate minimization of W and H . We consider a subproblem given $W \geq 0$ and $j \in \{1, \dots, m\}$:

$$\text{minimize } f_{KL}(h) = \sum_i \left[v_i \log \frac{v_i}{\sum_k W_{ik} h_k} - v_i + \sum_k W_{ik} h_k \right] \quad \text{subject to } h_k \geq 0, \quad (13)$$

where we let $v_i = V_{ij}$ and $h_k = H_{kj}$, as the objective is decomposed into m separate subproblems. Problem (13) can be reformulated to a scale invariant problem as follows.

Lemma A.4. The KL-NMF subproblem (13) is equivalent to the following scale invariant problem:

$$\text{maximize } - \sum_i v_i \log \sum_k W_{ik} \bar{h}_k \quad \text{subject to } \sum_k \bar{h}_k = 1, \bar{h}_k \geq 0, \quad (14)$$

with the relationship $(\sum_i v_i) \bar{h}_k = (\sum_i W_{ik}) h_k$.

A.2 Description of Data Sets

Table 1: A brief summary of data sets used for KL-NMF

Name	# of samples	# of features	# of nonzeros	Sparsity
WIKI	8,274	8,297	104,000	0.999
NIPS	1,500	12,419	280,000	0.985
KOS	3,430	6,906	950,000	0.960
WT	287	19,200	5,510,000	0.000

For KL divergence nonnegative matrix factorization (Section 5), we used 4 public real data sets available online⁴. Waving Trees (WT) has 287 images, each having 160×120 pixels. KOS and NIPS

Table 2: A brief summary of data sets used for GMM

Name	# of classes	# of samples	Dimension
Sonar	2	208	60
Ionosphere	2	351	34
HouseVotes84	2	435	16
BrCancer	2	699	10
PIDiabetes	2	768	8
Vehicle	4	846	18
Glass	6	214	9
Zoo	7	101	16
Vowel	11	990	10
Servo	51	167	4

are sparse, large matrices implemented for topic modeling. WIKI is a large binary matrix having values 0 or 1 representing the adjacency matrix of a directed graph.

For GMM (Section 5), we used 10 public real data sets. We used all small/moderate data sets provided by the `mlbench` package in R. We select data sets for multi-class classification problems and run EM and SCI-PI for given number of classes without class labels. The sample size varies from 101 to 990, the dimension varies from 2 to 60, and the number of classes varies from 2 to 51. If missing data exists, we simply replace it by 0 since our main focus is to solve the optimization problem better.

Table 3: A brief summary of data sets used for ICA

Name	# of samples	# of features
Wine	178	14
Soybean	683	35
Vehicle	846	18
Vowel	990	10
Cardio	2,126	22
Satellite	6,435	37
Pendigits	10,992	17
Letter	20,000	16
Shuttle	58,000	9

For ICA, we used 9 public data sets from the UCI Machine Learning repository⁵. The sample size varies from 178 to 58,000 and the dimension varies from 9 to 37.

A.3 More on Nonnegative Matrix Factorization

We first draw averaged convergence plots for the 4 real world data sets in Figure ?? . For the WT data set, MIXSQP exhibits a best convergence. Also, the convergence of SCI-PI is much faster than those of MU and PGD. For the other 3 data sets, MIXSQP sometimes converges to suboptimal points. Also, SCI-PI exhibits fastest convergence.

Next, we design a simple simulation study to evaluate the performance of block SCI-PI on KL-NMF problems. To this end, we sample a data matrix V independently from a single “zero-inflated” Poisson distribution (ZIP):

$$V_{ij} \sim \pi_0 \delta_0 + (1 - \pi_0) \text{Poisson}(l) \quad (15)$$

where π_0 is the proportion of zero inflation and l is the mean parameter of the Poisson distribution. Although this data generating distribution does not always reflect empirical distributions of real-world data sets, our focus here is to understand the behavior of SCI-PI compared to the other two methods, MU and PGD. Let n and m be the row and column lengths of $V \in \mathbb{R}_+^{n \times m}$, K be the number of non-negative factors and s be the number of zero entries in V . We set $n = 500$, $m = 300$, $K = 10$,

⁴These 4 data sets are retrieved from <https://www.microsoft.com/en-us/research/project>, <https://archive.ics.uci.edu/ml/datasets/bag+of+words>, and <https://snap.stanford.edu/data/wiki-Vote.html>

⁵<https://archive.ics.uci.edu/ml/index.php>

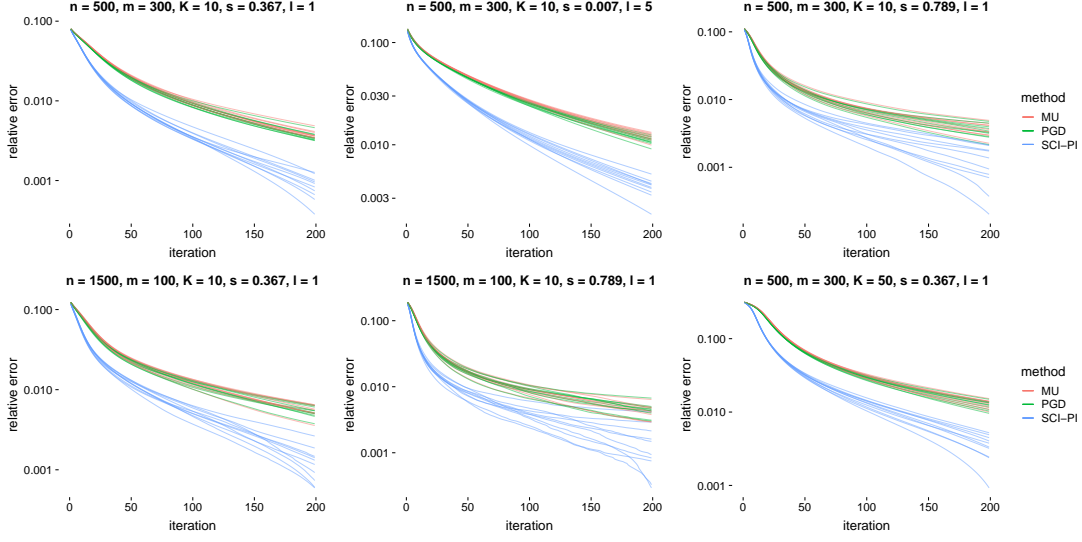


Figure 5: Convergence plots of 3 methods for KL-NMF on 6 synthetic data sets. We draw 10 convergence plots for each method differently initialized at random.

$\pi_0 = 0$ and $l = 1$ as a default and change some parameters to understand how the algorithms work for different settings.

Figure 5 summarizes the result. We conclude that SCI-PI tends to perform better in comparison to MU and PGD when V is denser ((1,1) vs. (1,3) and (2,1) vs. (2,2) in Figure 5), when K is larger ((1,1) vs. (2,3) in Figure 5) and when V is more uniformly distributed ((1,1) vs (1,2) and (1,1) vs (2,1) in Figure 5).

A.4 Proofs of Results in Section 2

Proof of Proposition 2.2. We first consider the multiplicative scale invariant case. Let x be a point such that $f(x) \neq 0$. Then, we have

$$f(rsx) = u(rs)f(x) = u(r)u(s)f(x),$$

resulting in

$$u(rs) = u(r)u(s)$$

for all $r, s \in \mathbb{R}$. Letting $g(r) = \ln(u(e^r))$, we have

$$g(r+s) = \ln(u(e^{r+s})) = \ln(u(e^r e^s)) = \ln(u(e^r)u(e^s)) = \ln(u(e^r)) + \ln(u(e^s)) = g(r) + g(s),$$

implying that g satisfies the Cauchy functional equation. Since f is continuous, so is u and thus g . Therefore, by [18], we have

$$g(r) = rg(1) \tag{16}$$

for all $r \geq 0$. From the definition of g and (16), we have

$$u(e^r) = e^{g(r)} = (e^r)^{g(1)}. \tag{17}$$

Representing $r > 0$ as $r = e^{\ln(r)}$ and using (17), we obtain

$$u(r) = u\left(e^{\ln(r)}\right) = r^{g(1)} = r^{\ln(u(e))} = r^p.$$

If $p = \ln(u(e)) < 0$, then since $f(x) \neq 0$, we have

$$\lim_{r \rightarrow 0^+} f(rx) = \lim_{r \rightarrow 0^+} u(r)f(x) = f(x) \cdot \lim_{r \rightarrow 0^+} r^p = f(x) \cdot \infty \neq f(0) < \infty,$$

contradicting the fact that f is continuous at 0. If $p = 0$, we get $u(r) = 1$, which contradicts $u(0) = 0$. Therefore, we must have $p > 0$. From u being an even function, we finally have

$$u(r) = |r|^p$$

for $r \in \mathbb{R}$.

Consider now the additive scale invariant case. For any $x \in \text{dom}(f)$, we have

$$f(rsx) = f(x) + v(rs) = f(x) + v(r) + v(s),$$

resulting in

$$v(rs) = v(r) + v(s)$$

for all $r, s \in \mathbb{R}$. Letting $g(r) = v(e^r)$, we have

$$g(r + s) = v(e^{r+s}) = v(e^r e^s) = v(e^r) + v(e^s) = g(r) + g(s).$$

Since g is continuous and satisfies the Cauchy functional equation, we have

$$g(r) = rg(1)$$

for all $r \geq 0$. For $r > 0$, letting $r = e^{\ln(r)}$, we have

$$v(r) = v(e^{\ln(r)}) = g(\ln(r)) = g(1)\ln(r) = v(e)\ln(r) = \log_a(r)$$

where $a = e^{\frac{1}{v(e)}}$ satisfying $0 < a, a \neq 1$. From the fact that v is an even function, we finally have

$$v(r) = \log_a|r|$$

for $r \in \mathbb{R} \setminus \{0\}$. □

Proof of Proposition 2.3. Without loss of generality, we can represent a scale-invariant function f as

$$f(cx) = u(c)f(x) + v(c) \tag{18}$$

since we can restore a multiplicatively or additively scale-invariant function by setting $v(c) = 0$ or $u(c) = 1$, respectively. By differentiating (18) with respect to x , we have

$$\nabla f(cx) = \frac{u(c)}{c} \nabla f(x).$$

On the other hand, by differentiating (18) with respect to c , we have

$$\nabla f(cx)^T x = u'(c)f(x) + v'(c). \tag{19}$$

By differentiating (19) with respect to x , we obtain

$$c\nabla^2 f(cx)x + \nabla f(cx) = u'(c)\nabla f(x). \tag{20}$$

Plugging $c = 1$ into (19) and (20) completes the proof. □

Proof of Proposition 2.4. Consider the Lagrangian function

$$L(x, \lambda) = f(x) + \frac{\lambda}{2} (1 - \|x\|^2)$$

and a stationary point (λ^*, x^*) satisfying

$$\nabla f(x^*) = \lambda^* x^*, \quad \|x^*\| = 1.$$

If f is multiplicative scale invariant with the degree of p , by Proposition 2.3, we have

$$\nabla^2 f(x^*)x^* = (p-1)\nabla f(x^*) = (p-1)\lambda^* x^*.$$

Also, by Proposition 2.3, if f is additive scale invariant f , we have

$$\nabla^2 f(x^*)x^* = -\nabla f(x^*) = -\lambda^* x^*.$$

Therefore, in both cases, a stationary point x^* is an eigenvector of $\nabla^2 f(x^*)$.

Next, suppose that λ^* is greater than the remaining eigenvalues of $\nabla^2 f(x^*)$. Since we have

$$d^T \nabla_{xx}^2 L(x^*, \lambda^*) d = d^T (\nabla^2 f(x^*) - \lambda^* I) d = d^T \nabla^2 f(x^*) d - \lambda^* \|d\|^2 < 0.$$

for any d satisfying $d^T x^* = 0$, the second-order sufficient condition is satisfied. Therefore, x^* is a local maximum. □

Proof of Proposition 2.5. First, we consider the case where an objective function f is multiplicative scale invariant with a multiplicative factor $u(c) = |c|^p$ where $p > 0$. Let w^* be an optimal solution to (6). From $f(0) = 0$ and $f(w^*) = 1$, we have $w^* \neq 0$, leading to

$$\|w^*\| > 0, \quad f\left(\frac{w^*}{\|w^*\|}\right) = \frac{1}{\|w^*\|^p} > 0.$$

Suppose an optimal solution to (1) is y with

$$f(y) > f\left(\frac{w^*}{\|w^*\|}\right) > 0. \quad (21)$$

Letting

$$\hat{y} = \frac{y}{f(y)^{1/p}},$$

we have

$$f(\hat{y}) = 1, \quad y = \frac{\hat{y}}{\|\hat{y}\|}.$$

Since $f(\hat{y}) = f(w^*) = 1$, we have

$$f(y) = f\left(\frac{\hat{y}}{\|\hat{y}\|}\right) = \frac{1}{\|\hat{y}\|^{1/p}}, \quad f\left(\frac{w^*}{\|w^*\|}\right) = \frac{1}{\|w^*\|^{1/p}}. \quad (22)$$

From (21) and (22), we have $\|\hat{y}\| < \|w^*\|$, contradicting that w^* is an optimal solution to (6).

On the other hand, let x^* be an optimal solution to (1) with $f(x^*) > 0$. Suppose that an optimal solution to (6) is z with

$$\|z\| < \left\| \frac{x^*}{f(x^*)^{1/p}} \right\|. \quad (23)$$

Letting

$$\hat{z} = \frac{z}{\|z\|},$$

we have

$$\|\hat{z}\| = 1, \quad z = \frac{\hat{z}}{f(\hat{z})^{1/p}}.$$

Since $\|\hat{z}\| = \|x^*\| = 1$, we have

$$\|z\| = \left\| \frac{\hat{z}}{f(\hat{z})^{1/p}} \right\| = \frac{1}{f(\hat{z})^{1/p}}, \quad \left\| \frac{x^*}{f(x^*)^{1/p}} \right\| = \frac{1}{f(x^*)^{1/p}}. \quad (24)$$

From (23) and (24), we have

$$f(x^*) < f(\hat{z})$$

due to $p > 0$, contradicting the assumption that x^* is an optimal solution to (1).

Next, let f be an additively scale invariant function with an additive factor $v(c) = \log_a |c|$ where $a > 1$. In the same way as above, let w^* be an optimal solution to (6) and suppose that an optimal solution of (1) is y with

$$f(y) > f\left(\frac{w^*}{\|w^*\|}\right). \quad (25)$$

Letting

$$\hat{y} = a^{1-f(y)} y,$$

we have

$$f(\hat{y}) = 1, \quad y = \frac{\hat{y}}{\|\hat{y}\|}.$$

Since $f(\hat{y}) = f(w^*) = 1$, we have

$$f(y) = f(\hat{y}) - \log_a \|\hat{y}\| = 1 - \log_a \|\hat{y}\|, \quad f\left(\frac{w^*}{\|w^*\|}\right) = 1 - \log_a \|w^*\|. \quad (26)$$

From (25) and (26), we have

$$\|\hat{y}\| < \|w^*\|$$

due to $a > 1$, contradicting the fact that w^* is an optimal solution to (6).

On the other hand, let x^* be an optimal solution to (1) and suppose that an optimal solution to (6) is z with

$$\|z\| < \|a^{1-f(x^*)} x^*\|. \quad (27)$$

Letting

$$\hat{z} = \frac{z}{\|z\|},$$

we have

$$\|\hat{z}\| = 1, \quad z = a^{1-f(\hat{z})} \hat{z}.$$

Since $\|\hat{z}\| = \|x^*\| = 1$, we have

$$\|z\| = a^{1-f(\hat{z})}, \quad \|a^{1-f(x^*)} x^*\| = a^{1-f(x^*)}. \quad (28)$$

From (27) and (28), we have

$$f(x^*) < f(\hat{z})$$

due to $a > 1$, contradicting the assumption that x^* is an optimal solution to (1). \square

Proof of Lemma A.4. Since a log-linear function is concave, (13) is a convex problem in h . Consider the Lagrangian of the original problem

$$\mathcal{L}(h, \lambda) = f_{KL}(h) - \sum_k \lambda_k h_k \quad (29)$$

where $\lambda \geq 0$. By the first-order KKT conditions, we must have

$$\nabla_k f_{KL}(h^*) = \lambda_k^*, \quad \lambda_k^* h_k^* = 0, \quad \forall k = 1, \dots, K \quad (30)$$

at an optimal solution (h^*, λ^*) . Since (30) implies $\sum_k h_k^* \lambda_k^* = 0$, we have

$$\sum_k h_k^* \lambda_k^* = \sum_k h_k^* \nabla_k f_{KL}(h^*) = - \sum_{i,k} \frac{v_i W_{ik} h_k^*}{\sum_{k'} W_{ik'} h_{k'}^*} + \sum_{i,k} W_{ik} h_k^* = \sum_i v_i - \sum_{i,k} W_{ik} h_k^* = 0,$$

resulting in

$$\sum_i v_i - \sum_{i,k} W_{ik} h_k^* = 0. \quad (31)$$

We can show that

$$\text{minimize } f_{SCI}(h) \triangleq \sum_i v_i \log \frac{v_i}{\sum_k W_{ik} h_k} \quad \text{subject to } \sum_i v_i = \sum_{i,k} W_{ik} h_k, \quad h_k \geq 0. \quad (32)$$

is equivalent to the original subproblem (13), due to the following.

1. It always satisfies $f_{SCI}^* \geq f_{KL}^*$ since (32) has an additional constraint $\sum_i v_i = \sum_{i,k} W_{ik} h_k$ compared to (13).
2. A solution h^* of (13) is a feasible point of (32) since we have shown that $\sum_i v_i = \sum_{i,k} W_{ik} h_k^*$. This implies $f_{KL}^* \geq f_{SCI}^*$.

Now we can reparametrize h by \bar{h} so that $\sum_i v_i = \sum_{i,k} W_{ik} h_k$ if and only if $\sum_k \bar{h}_k = 1$, which yields the relationship between two variables $\bar{h}_k = h_k (\sum_i W_{ik}) / (\sum_i v_i)$. This completes the proof. \square

A.5 Proofs of Theorem 3.1 and Theorem 4.1

On several occasions, we use if $x \in \partial B_d$, $y \in \partial B_d$, then

$$\|x - y\|^2 = \|x\|^2 + \|y\|^2 - 2x^T y = 2(1 - x^T y).$$

Note that if $x^T y \geq 0$, then

$$\sqrt{1 - (x^T y)^2} = \sqrt{(1 - x^T y)(1 + x^T y)} \geq \sqrt{1 - x^T y} = \frac{\|x - y\|}{\sqrt{2}}.$$

By Cauchy-Schwarz, we also have

$$\sqrt{1 - (x^T y)^2} = \sqrt{(1 - x^T y)(1 + x^T y)} \leq \sqrt{2}\sqrt{1 - x^T y} = \|x - y\|.$$

Lemma A.5. *Let $\{v_1, \dots, v_d\}$ be an orthogonal basis in \mathbb{R}^d and $\{x_k\}_{k=0,1,\dots}$ be the sequence of iterates generated by SCI-PI. If for every $x \in \partial B_d$ we have*

$$\nabla f(x)^T v_1 = \lambda^* + \alpha(x), \quad \sum_{i=2}^d (\nabla f(x)^T v_i)^2 \leq (\bar{\lambda}_2 \|x - x^*\| + \beta(x))^2 \quad (33)$$

where

$$\alpha(x) = o(\sqrt{\|x - x^*\|}), \quad \beta(x) = o(\|x - x^*\|),$$

then there exists some $\delta > 0$ such that under the initial condition $1 - x_0^T x^* < \delta$, we have

$$1 - (x_k^T x^*)^2 \leq \prod_{t=0}^{k-1} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma_t \right)^2 (1 - (x_0^T x^*)^2) \quad \text{and} \quad \lim_{k \rightarrow \infty} \gamma_k = 0.$$

Proof. By (33) for every $x \in \partial B_d$, we have

$$\frac{\sum_{i=2}^d (\nabla f(x)^T v_i)^2}{(\nabla f(x)^T v_1)^2} \leq \left(\frac{\bar{\lambda}_2 \|x - x^*\| + \alpha(x)}{\lambda^* + \beta(\sqrt{x})} \right)^2.$$

From

$$\frac{\bar{\lambda}_2 \|x - x^*\| + \alpha(x)}{\lambda^* + \beta(x)} = \frac{\bar{\lambda}_2}{\lambda^*} \|x - x^*\| + \theta(x)$$

where $\theta(x) = o(\|x - x^*\|)$, we have

$$\frac{\sum_{i=2}^d (\nabla f(x)^T v_i)^2}{(\nabla f(x)^T v_1)^2} \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \frac{\theta(x)}{\|x - x^*\|} \right)^2 \|x - x^*\|^2. \quad (34)$$

Using

$$\epsilon(x) = \frac{\theta(x)}{\|x - x^*\|}, \quad \|x - x^*\|^2 = \left(1 + \frac{1 - x^T x^*}{1 + x^T x^*} \right) (1 - (x^T x^*)^2), \quad (35)$$

we can further represent (34) as

$$\begin{aligned} \frac{\sum_{i=2}^d (\nabla f(x)^T v_i)^2}{(\nabla f(x)^T v_1)^2} &\leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \epsilon(x) \right)^2 \left(1 + \frac{1 - x^T x^*}{1 + x^T x^*} \right) (1 - (x^T x^*)^2) \\ &= \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma(x) \right)^2 (1 - (x^T x^*)^2) \end{aligned} \quad (36)$$

where

$$\gamma(x) = \frac{\bar{\lambda}_2}{\lambda^*} \left(\frac{1 - x^T x^*}{1 + x^T x^* + \sqrt{2(1 + x^T x^*)}} \right) + \epsilon(x) \sqrt{1 + \frac{1 - x^T x^*}{1 + x^T x^*}}. \quad (37)$$

From (33), there exists some $\delta_1 > 0$ such that if $1 - x^T x^* < \delta_1$, then

$$\nabla f(x)^T v_1 > 0. \quad (38)$$

Also, by (35), for any $\bar{\gamma} > 0$ satisfying

$$\frac{\bar{\lambda}_2}{\lambda^*} + \bar{\gamma} < 1, \quad (39)$$

there exists some constant $\delta_2 > 0$ such that if $1 - x^T x^* < \delta_2$, then

$$|\epsilon(x)| \leq \frac{\bar{\gamma}}{4}. \quad (40)$$

Let $\delta = \min\{\delta_1, \delta_2, \frac{\lambda^*}{\lambda_2} \bar{\gamma}, 1\}$.

We next argue that if $1 - x_k^T x^* < \delta$, then we have

$$x_{k+1}^T x^* > 0, \quad 1 - (x_{k+1}^T x^*)^2 \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma_k\right)^2 (1 - (x_k^T x^*)^2) \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \bar{\gamma}\right)^2 (1 - (x_k^T x^*)^2). \quad (41)$$

From $\delta < 1$, we have $x_k^T x^* > 0$. Also, from $1 - x_k^T x^* < \delta_1$, $x^* = v_1$ and the update rule of SCI-PI, we obtain

$$x_{k+1}^T x^* = \frac{\nabla f(x_k)^T x^*}{\|\nabla f(x_k)\|} = \frac{\nabla f(x_k)^T v_1}{\|\nabla f(x_k)\|} > 0$$

due to (38). From $|x_{k+1}^T v_1| \leq \|x_{k+1}\| \|v_1\| = 1$, we have

$$1 - (x_{k+1}^T x^*)^2 \leq \frac{1 - (x_{k+1}^T v_1)^2}{(x_{k+1}^T v_1)^2}.$$

Also, from the update rule of SCI-PI and the fact that $\{v_1, \dots, v_d\}$ forms an orthogonal basis in \mathbb{R}^d implying $\nabla f(x_k) = \sum_{i=1}^d (\nabla f(x_k)^T v_i) v_i$ and $\|\nabla f(x_k)\|^2 = \sum_{i=1}^d (\nabla f(x_k)^T v_i)^2$, we obtain

$$\frac{1 - (x_{k+1}^T v_1)^2}{(x_{k+1}^T v_1)^2} = \frac{\|\nabla f(x_k)\|^2 - (\nabla f(x_k)^T v_1)^2}{(\nabla f(x_k)^T v_1)^2} = \frac{\sum_{i=2}^d (\nabla f(x_k)^T v_i)^2}{(\nabla f(x_k)^T v_1)^2},$$

resulting in

$$1 - (x_{k+1}^T x^*)^2 \leq \frac{\sum_{i=2}^d (\nabla f(x_k)^T v_i)^2}{(\nabla f(x_k)^T v_1)^2}.$$

Using $x_k^T x^* > 0$ and $1 - x_k^T x^* < \min\{\delta_2, \frac{\lambda^*}{\lambda_2} \bar{\gamma}\}$ in (37) for iteration k , we have

$$\gamma_k = \gamma(x_k) = \frac{\bar{\lambda}_2}{\lambda^*} \left(\frac{1 - x_k^T x^*}{1 + x_k^T x^* + \sqrt{2(1 + x_k^T x^*)}} \right) + \epsilon(x_k) \sqrt{1 + \frac{1 - x_k^T x^*}{1 + x_k^T x^*}} \leq \frac{\bar{\gamma}}{2} + \frac{\bar{\gamma}}{2} = \bar{\gamma},$$

which from (36) leads to

$$1 - (x_{k+1}^T x^*)^2 \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma_k\right)^2 (1 - (x_k^T x^*)^2) \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \bar{\gamma}\right)^2 (1 - (x_k^T x^*)^2) < \delta. \quad (42)$$

Next, using mathematical induction, we show that if

$$1 - x_0^T x^* < \delta, \quad (43)$$

then

$$1 - x_k^T x^* < \delta \quad (44)$$

for all $k \geq 0$.

By (43), we have $1 - x_0^T x^* < \delta$, which shows the base case.

Suppose that $1 - x_k^T x^* < \delta$ holds. Then, we have (41). Also, from $\delta < 1$, we have $x_k^T x^* > 0$. From $x_{k+1}^T x^* > 0$, $x_k^T x^* > 0$, and $1 - (x_{k+1}^T x^*)^2 \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} + \bar{\gamma}\right)^2 (1 - (x_k^T x^*)^2) < 1 - (x_k^T x^*)^2$, we have

$$1 - x_{k+1}^T x^* < 1 - x_k^T x^* < \delta.$$

This completes the induction proof.

Since (44) holds for all $k \geq 0$, we can repeatedly apply (41) to obtain

$$1 - (x_k^T x^*)^2 \leq \prod_{t=0}^{k-1} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \gamma_t\right)^2 (1 - (x_0^T x^*)^2) < \left(\frac{\bar{\lambda}_2}{\lambda^*} + \bar{\gamma}\right)^{2k} (1 - (x_0^T x^*)^2). \quad (45)$$

From (45), we obtain $(x_k^T x^*)^2 \rightarrow 1$. Since $x_k^T x^* > 0$ for all $k \geq 0$ by (44), we have $x_k \rightarrow x^*$, and thus $\lim_{k \rightarrow \infty} \gamma_k = 0$ by (37). With (45), this gives the desired result. \square

Lemma A.6. *Let $\{v_1, \dots, v_d\}$ be an orthogonal basis in \mathbb{R}^d . If $x^* = v_1$ and a sequence of iterates $\{x_k\}_{k=0,1,\dots}$ generated by SCI-PI satisfies*

$$\nabla f(x_k)^T v_1 \geq A - B(1 - x_k^T x^*) - C\sqrt{1 - x_k^T x^*} \quad (46)$$

and

$$\sum_{i=2}^d (\nabla f(x_k)^T v_i)^2 \leq \left(D\sqrt{1 - (x_k^T x^*)^2} + E\sqrt{2(1 - x_k^T x^*)} + \frac{F}{2}\|x_k - x^*\|^2\right)^2 \quad (47)$$

where $A > 0$ and B, C, D, E, F are non-negative real numbers such that $B + C > 0$ and $\frac{D+E}{A} < 1$, then under the initial condition

$$1 - x_0^T x^* < \delta$$

where

$$\delta = \min \left\{ \left(\frac{A}{B+C}\right)^2, \left(\frac{A-D-E}{B+C+E+F}\right)^2, 1 \right\} \quad (48)$$

we have

$$1 - (x_k^T x^*)^2 \leq \prod_{t=0}^{k-1} \left(\frac{D+E}{A} + \gamma_t\right)^2 (1 - (x_0^T x^*)^2) \quad \text{and} \quad \lim_{k \rightarrow \infty} \gamma_k = 0.$$

Proof. We first show that if $1 - x_k^T x^* < \delta$, then we have

$$x_{k+1}^T x^* > 0, \quad 1 - (x_{k+1}^T x^*)^2 < \left(\frac{D+E}{A} + \gamma_k\right)^2 (1 - (x_k^T x^*)^2), \quad \frac{D+E}{A} + \gamma_k < 1 \quad (49)$$

for all $k \geq 0$ where

$$\gamma_k = \frac{(A(E+F) + (B+C)(D+E))\sqrt{1 - x_k^T x^*}}{A(A - (B+C)\sqrt{1 - x_k^T x^*})}. \quad (50)$$

Since $\sqrt{1 - x_k^T x^*} \geq 1 - x_k^T x^*$ holds from $0 < 1 - x_k^T x^* < 1$, using the update rule of SCI-PI, (46), and $\delta \leq \left(\frac{A}{B+C}\right)^2$, we have

$$x_{k+1}^T x^* = \frac{\nabla f(x_k)^T v_1}{\|\nabla f(x_k)\|} \geq \frac{A - B(1 - x_k^T x^*) - C\sqrt{1 - x_k^T x^*}}{\|\nabla f(x_k)\|} \geq \frac{A - (B+C)\sqrt{1 - x_k^T x^*}}{\|\nabla f(x_k)\|} > 0. \quad (51)$$

From $x^* = v_1$, $|x_{k+1}^T v_1| \leq \|x_{k+1}\| \|v_1\| = 1$, the update rule of SCI-PI and the fact that $\{v_1, \dots, v_d\}$ forms an orthogonal basis in \mathbb{R}^d implying $\nabla f(x_k) = \sum_{i=1}^d (\nabla f(x_k)^T v_i) v_i$ and $\|\nabla f(x_k)\|^2 = \sum_{i=1}^d (\nabla f(x_k)^T v_i)^2$, we have

$$1 - (x_{k+1}^T x^*)^2 \leq \frac{1 - (x_{k+1}^T v_1)^2}{(x_{k+1}^T v_1)^2} = \frac{\|\nabla f(x_k)\|^2 - (\nabla f(x_k)^T v_1)^2}{(\nabla f(x_k)^T v_1)^2} = \frac{\sum_{i=2}^d (\nabla f(x_k)^T v_i)^2}{(\nabla f(x_k)^T v_1)^2}. \quad (52)$$

By (51), we have

$$A - B(1 - x_k^T x^*) - C\sqrt{1 - x_k^T x^*} > 0.$$

Therefore, by plugging (46) and (47) into (52) and using that $x_k^T x^* > 0$, we have

$$\begin{aligned} 1 - (x_{k+1}^T x^*)^2 &\leq \left(\frac{D\sqrt{1 - (x_k^T x^*)^2} + E\sqrt{2(1 - x_k^T x^*)} + \frac{F}{2}\|x_k - x^*\|^2}{A - B(1 - x_k^T x^*) - C\sqrt{1 - x_k^T x^*}} \right)^2 \\ &= \left(\frac{D + E\sqrt{1 + \frac{1 - x_k^T x^*}{1 + x_k^T x^*}} + F\sqrt{\frac{1 - x_k^T x^*}{1 + x_k^T x^*}}}{A - B(1 - x_k^T x^*) - C\sqrt{1 - x_k^T x^*}} \right)^2 (1 - (x_k^T x^*)^2) \\ &\leq \left(\frac{D + E\left(1 + \sqrt{1 - x_k^T x^*}\right) + F\sqrt{1 - x_k^T x^*}}{A - (B + C)\sqrt{1 - x_k^T x^*}} \right)^2 (1 - (x_k^T x^*)^2) \\ &= \left(\frac{D + E}{A} + \gamma_k \right)^2 (1 - (x_k^T x^*)^2). \end{aligned} \quad (53)$$

In the above, we use the fact that $\sqrt{1 + x} \leq 1 + \sqrt{x}$ for $x \geq 0$ to derive the second inequality.

Moreover, from

$$\sqrt{1 - x_k^T x^*} < \sqrt{\delta} \leq \frac{A - D - E}{B + C + E + F},$$

we have

$$\gamma_k < 1 - \frac{D + E}{A}.$$

Next, using mathematical induction, we show that if

$$1 - x_0^T x^* < \delta, \quad (54)$$

then we have

$$1 - x_k^T x^* < \delta \quad (55)$$

for all $k \geq 0$.

By (54), we have $1 - x_0^T x^* < \delta$, which proves the base case.

Suppose that we have $1 - x_k^T x^* < \delta$. Then, we have (49). Also, from $\delta < 1$, we have $x_k^T x^* > 0$. From $x_{k+1}^T x^* > 0$, $x_k^T x^* > 0$, and that $1 - (x_{k+1}^T x^*)^2 < 1 - (x_k^T x^*)^2$, we have

$$1 - x_{k+1}^T x^* < 1 - x_k^T x^* < \delta.$$

This completes the induction proof.

Since (55) holds for all $k \geq 0$, by repeatedly applying (49), we obtain

$$1 - (x_k^T x^*)^2 \leq \prod_{t=0}^{k-1} \left(\frac{D + E}{A} + \gamma_t \right)^2 (1 - (x_0^T x^*)^2). \quad (56)$$

Since $\frac{D+E}{A} + \gamma_k < 1$ for all $k \geq 0$, $1 - (x_k^T x^*)^2$ is monotone decreasing, and so is $1 - x_k^T x^*$ by non-negativity. Since γ_k is a monotone increasing function of $1 - x_k^T x^*$, we have $\gamma_{k+1} \leq \gamma_k$ for all $k \geq 0$, resulting in

$$\prod_{t=0}^{k-1} \left(\frac{D+E}{A} + \gamma_t \right) \leq \left(\frac{D+E}{A} + \gamma_0 \right)^{2k}.$$

Since $\frac{D+E}{A} + \gamma_0 < 1$ by (49), we have $(x_k^T x^*)^2 \rightarrow 1$. Due to $x_k^T x^* > 0$ for all $k \geq 0$, this implies $x_k \rightarrow x^*$, and thus $\lim_{k \rightarrow \infty} \gamma_k = 0$ due to (50). With (56), this gives the desired result. \square

Proof of Theorem 3.1. Since $\nabla^2 f(x^*)$ is real and symmetric, without loss of generality, we assume that $\{v_1, \dots, v_d\}$ form an orthogonal basis in \mathbb{R}^d .

Since f is twice continuously differentiable on an open set containing $\partial\mathcal{B}_d$, for $x \in \partial\mathcal{B}_d$, using the Taylor expansion of $\nabla f(x)^T v_i$ at x^* , we have

$$\nabla f(x)^T v_i = \nabla f(x^*)^T v_i + (x - x^*)^T \nabla^2 f(x^*) v_i + R_i(x) \quad (57)$$

where

$$R_i(x) = o(\|x - x^*\|). \quad (58)$$

From $\nabla f(x^*) = \lambda^* x^*$ and $x^* = v_1$, we have

$$\begin{aligned} \nabla f(x)^T v_1 &= \nabla f(x^*)^T x^* + (x - x^*)^T \nabla^2 f(x^*) x^* + R_1(x) \\ &= \lambda^* - \lambda_1(1 - x^T x^*) + R_1(x) \\ &= \lambda^* + \alpha(x) \end{aligned} \quad (59)$$

where

$$\alpha(x) = -\lambda_1(1 - x^T x^*) + R_1(x) = o(\|x - x^*\|)$$

due to $R_1(x) = o(\|x - x^*\|)$ and $1 - x^T x^* = o(\|x - x^*\|)$.

On the other hand, for $2 \leq i \leq d$, due to $\nabla f(x^*) = \lambda^* x^*$, we have

$$\nabla f(x^*)^T v_i = \lambda^* (x^*)^T v_i = 0. \quad (60)$$

From (57), this results in

$$\nabla f(x)^T v_i = \lambda_i x^T v_i + R_i(x). \quad (61)$$

Let $\bar{R}_2(x) = \max_{2 \leq i \leq d} |R_i(x)|$. Note that $\bar{R}_2(x) = o(\|x - x^*\|)$. By (61), we obtain

$$\begin{aligned} \sum_{i=2}^d (\nabla f(x)^T v_i)^2 &= \sum_{i=2}^d \left[\lambda_i^2 (x^T v_i)^2 + 2\lambda_i (x^T v_i) R_i(x) + (R_i(x))^2 \right] \\ &\leq \bar{\lambda}_2^2 \sum_{i=2}^d (x^T v_i)^2 + 2\bar{\lambda}_2 \bar{R}_2(x) \sum_{i=2}^d |x^T v_i| + d (\bar{R}_2(x))^2. \end{aligned} \quad (62)$$

From $x \in \partial\mathcal{B}_d$, $x^* = v_1$, and the fact that $\{v_1, \dots, v_d\}$ forms an orthogonal basis in \mathbb{R}^d , we have

$$\sum_{i=2}^d (x^T v_i)^2 = 1 - (x^T v_1)^2 = 1 - (x^T x^*)^2 = (1 - x^T x^*)(1 + x^T x^*) \leq 2(1 - x^T x^*) = \|x - x^*\|^2.$$

Also, by the Cauchy Schwartz inequality, we have

$$\sum_{i=2}^d |x^T v_i| \leq \sqrt{d} \sqrt{\sum_{i=2}^d (x^T v_i)^2} \leq \sqrt{d} \|x - x^*\|.$$

Therefore, we obtain from (62) that

$$\begin{aligned} \sum_{i=2}^d (\nabla f(x)^T v_i)^2 &\leq \bar{\lambda}_2^2 \|x - x^*\|^2 + 2\bar{\lambda}_2 \bar{R}_2(x) \sqrt{d} \|x - x^*\| + d (\bar{R}_2(x))^2 \\ &= (\bar{\lambda}_2 \|x - x^*\| + \beta(x))^2 \end{aligned} \quad (63)$$

where

$$\beta(x) = \sqrt{d}\bar{R}_2(x) = o(\|x - x^*\|).$$

By (59), (63), and Lemma A.5, we obtain the first part of the desired result.

Next, we consider the case where $\nabla_i f$ has a continuous Hessian H_i . From $\nabla_i f(x)$ being twice continuously differentiable in \mathcal{B}_∞ , we have

$$\nabla_i f(x_k) = \nabla_i f(x^*) + \nabla \nabla_i f(x^*)(x_k - x^*) + \frac{1}{2}(x_k - x^*)^T H_i(\hat{x}_k^i)(x_k - x^*) \quad (64)$$

where

$$\hat{x}_k^i \in \mathcal{N}(x_k, x^*) \triangleq \{x : x_s = t_s x_s^* + (1 - t_s)x_{k,s}, 0 \leq t_s \leq 1, s = 1, \dots, d\}.$$

In the above, x_s^* and $x_{k,s}$ denote the s^{th} coordinates of x^* and x_k , respectively.

For each $1 \leq i \leq d$, we have

$$\frac{1}{2} \sum_{j=1}^d v_{i,j} (x_k - x^*)^T H_j(\hat{x}_k^j)(x_k - x^*) = \frac{1}{2}(x_k - x^*)^T G_i(\hat{x}_k^j)(x_k - x^*).$$

Since

$$\begin{aligned} \left| (x_k - x^*)^T G_i(\hat{x}_k^j)(x_k - x^*) \right| &= \|x_k - x^*\|^2 \left| \left[\frac{x_k - x^*}{\|x_k - x^*\|} \right]^T G_i(\hat{x}_k^j) \left[\frac{x_k - x^*}{\|x_k - x^*\|} \right] \right| \\ &\leq \|x_k - x^*\|^2 \max_{x \in \partial \mathcal{B}_d} |x^T G_i(\hat{x}_k^j)x| \\ &\leq \|x_k - x^*\|^2 \max_{x \in \partial \mathcal{B}_d, y \in \mathcal{B}_\infty} |x^T G_i(y)x| \\ &\leq \|x_k - x^*\|^2 \max_{x \in \partial \mathcal{B}_d, y \in \mathcal{B}_\infty} \sqrt{\sum_{i=1}^d (x^T G_i(y)x)^2} \\ &= M \|x_k - x^*\|^2, \end{aligned}$$

we obtain

$$\frac{1}{2} \left| \sum_{j=1}^d v_{i,j} (x_k - x^*)^T H_j(\hat{x}_k^j)(x_k - x^*) \right| \leq \frac{1}{2} M \|x_k - x^*\|^2. \quad (65)$$

From (64), (65) and that $x^* = v_1$, we have

$$\nabla f(x_k)^T v_1 \geq \nabla f(x^*)^T x^* + (x_k - x^*)^T \nabla^2 f(x^*) x^* - \frac{M}{2} \|x_k - x^*\|^2,$$

resulting in

$$\nabla f(x_k)^T v_1 \geq \lambda^* - (M + |\lambda_1|)(1 - x_k^T x^*). \quad (66)$$

For $2 \leq i \leq d$, we have

$$\begin{aligned} \nabla f(x_k)^T v_i &= \nabla f(x^*)^T v_i + (x_k - x^*)^T \nabla^2 f(x^*) v_i + \frac{1}{2}(x_k - x^*)^T G_i(\hat{x}_k^j)(x_k - x^*) \\ &= \lambda_i x_k^T v_i + \frac{1}{2}(x_k - x^*)^T G_i(\hat{x}_k^j)(x_k - x^*). \end{aligned} \quad (67)$$

Since

$$\begin{aligned} \sum_{i=2}^d \left[(x_k - x^*)^T G_i(\hat{x}_k^j)(x_k - x^*) \right]^2 &\leq \|x_k - x^*\|^4 \sum_{i=2}^d \left(\left[\frac{x_k - x^*}{\|x_k - x^*\|} \right]^T G_i(\hat{x}_k^j) \left[\frac{x_k - x^*}{\|x_k - x^*\|} \right] \right)^2 \\ &\leq \|x_k - x^*\|^4 \max_{x \in \partial \mathcal{B}_d} \sum_{i=2}^d (x^T G_i(\hat{x}_k^j)x)^2 \\ &\leq \|x_k - x^*\|^4 \max_{x \in \partial \mathcal{B}_d, y \in \mathcal{B}_\infty} \sum_{i=2}^d (x^T G_i(y)x)^2 \\ &\leq \|x_k - x^*\|^4 \max_{x \in \partial \mathcal{B}_d, y \in \mathcal{B}_\infty} \sum_{i=1}^d (x^T G_i(y)x)^2 \\ &= M^2 \|x_k - x^*\|^4, \end{aligned}$$

using (67) and the Cauchy-Schwartz inequality, we obtain

$$\begin{aligned}
\sum_{i=2}^d (\nabla f(x_k)^T v_i)^2 &\leq \sum_{i=2}^d \left(|\lambda_i| |x_k^T v_i| + \frac{1}{2} (x_k - x^*)^T G_i(\hat{x}_k^j) (x_k - x^*) \right)^2 \\
&\leq \bar{\lambda}_2^2 \sum_{i=2}^d (x_k^T v_i)^2 + \bar{\lambda}_2 \sum_{i=2}^d |(x_k - x^*)^T G_i(\hat{x}_k^j) (x_k - x^*)| |x_k^T v_i| + \frac{M^2}{4} \|x_k - x^*\|^4 \\
&\leq \bar{\lambda}_2^2 \sum_{i=2}^d (x_k^T v_i)^2 + \bar{\lambda}_2 M \|x_k - x^*\|^2 \sqrt{\sum_{i=2}^d (x_k^T v_i)^2} + \frac{M^2}{4} \|x_k - x^*\|^4 \\
&= \left(\bar{\lambda}_2 \sqrt{1 - (x_k^T x^*)^2} + \frac{M}{2} \|x_k - x^*\|^2 \right)^2. \tag{68}
\end{aligned}$$

Using (66), (68), and Lemma A.6 with

$$A = \lambda^*, \quad B = M + |\lambda_1|, \quad C = 0, \quad D = \bar{\lambda}_2, \quad E = 0, \quad F = M,$$

we obtain the desired result. \square

Proof of Theorem 4.1. Let $\{v_1, \dots, v_d\}$ be a set of eigenvectors of $F(x^*)$. Without loss of generality, we assume $x^* = v_1$. Also, since $F(x^*)$ is real and symmetric, we assume that $\{v_1, \dots, v_d\}$ forms an orthogonal basis in \mathbb{R}^d

Since f is twice continuously differentiable on an open set containing $\partial\mathcal{B}_d$, for $x \in \partial\mathcal{B}_d$, using the Taylor expansion of $\nabla f(x)^T v_i$ at x^* , we have

$$\nabla f(x)^T v_i = \nabla f(x^*)^T v_i + (x - x^*)^T \nabla^2 f(x^*) v_i + R_i(x) \quad (69)$$

where $R_i(x) = o(\|x - x^*\|)$. Using (69) with $i = 1$ and $\nabla f(x^*) = \lambda^* x^*$, we obtain

$$\begin{aligned} \nabla f(x)^T v_1 &= \lambda^* (x^*)^T v_1 + (x - x^*)^T \nabla^2 f(x^*) v_1 + R_1(x) \\ &= \lambda^* + \alpha(x) \end{aligned} \quad (70)$$

where

$$\alpha(x) = (x - x^*)^T \nabla^2 f(x^*) v_1 + R_1(x) = o(\sqrt{\|x - x^*\|})$$

due to $(x - x^*)^T \nabla^2 f(x^*) v_1 = o(\sqrt{\|x - x^*\|})$ and $R_1(x) = o(\|x - x^*\|)$.

Again, using (69) and $\nabla f(x^*) = \lambda^* x^*$ for $2 \leq i \leq d$, we have

$$\begin{aligned} \nabla f(x)^T v_i &= \lambda^* (x^*)^T v_i + (x - x^*)^T \nabla^2 f(x^*) v_i + R_i(x) \\ &= (x - x^*)^T \nabla^2 f(x^*) v_i + R_i(x), \end{aligned}$$

resulting in

$$\sum_{i=2}^d (\nabla f(x)^T v_i)^2 = \sum_{i=2}^d ((x - x^*)^T \nabla^2 f(x^*) v_i + R_i(x))^2. \quad (71)$$

Let $\bar{R}_2(x) = \max_{2 \leq i \leq d} |R_i(x)|$. Note that $\bar{R}_2(x) = o(\|x - x^*\|)$.

From $x^* = v_1$ and the fact that $\{v_1, \dots, v_d\}$ forms an orthogonal basis in \mathbb{R}^d , we have

$$\begin{aligned} \sum_{i=2}^d ((x - x^*)^T \nabla^2 f(x^*) v_i)^2 &= \|\nabla^2 f(x^*) (x - x^*)\|_2^2 - ((x - x^*)^T \nabla^2 f(x^*) v_1)^2 \\ &= (x - x^*)^T \nabla^2 f(x^*) (I - x^* (x^*)^T) \nabla^2 f(x^*) (x - x^*) \\ &= (x - x^*)^T \nabla^2 f(x^*) (I - x^* (x^*)^T)^2 \nabla^2 f(x^*) (x - x^*). \end{aligned}$$

Since

$$\begin{aligned} \|\nabla^2 f(x^*) (I - x^* (x^*)^T)^2 \nabla^2 f(x^*)\| &= \|(I - x^* (x^*)^T) \nabla^2 f(x^*)\|^2 \\ &= \|\nabla^2 f(x^*) (I - x^* (x^*)^T)\|^2, \end{aligned}$$

we have

$$\sum_{i=2}^d ((x - x^*)^T \nabla^2 f(x^*) v_i)^2 \leq \bar{\lambda}_2^2 \|x - x^*\|^2. \quad (72)$$

Also, from (72) and the Cauchy-Schwartz inequality, we obtain

$$\sum_{i=2}^d (x - x^*)^T \nabla^2 f(x^*) v_i \leq \sum_{i=2}^d |(x - x^*)^T \nabla^2 f(x^*) v_i| \leq \bar{\lambda}_2 \sqrt{d} \|x - x^*\|. \quad (73)$$

Using (72) and (73) for (71), we obtain

$$\begin{aligned} \sum_{i=2}^d ((x - x^*)^T \nabla^2 f(x^*) v_i)^2 &\leq \bar{\lambda}_2^2 \|x - x^*\|^2 + 2\bar{R}_2(x) \sum_{i=2}^d (x - x^*)^T \nabla^2 f(x^*) v_i + d(\bar{R}_2(x))^2 \\ &\leq \bar{\lambda}_2^2 \|x - x^*\|^2 + 2\bar{\lambda}_2 \bar{R}_2(x) \sqrt{d} \|x - x^*\| + d(\bar{R}_2(x))^2, \end{aligned}$$

resulting in

$$\sum_{i=2}^d (\nabla f(x)^T v_i)^2 \leq (\bar{\lambda}_2 \|x - x^*\|^2 + \beta(x))^2 \quad (74)$$

where

$$\beta(x) = \sqrt{d} \bar{R}_2(x) = o(\|x - x^*\|).$$

By (70), (74), and Lemma A.5, we obtain the first part of the desired result.

Next, we assume that $\nabla_i f$ has a continuous Hessian H_i . By the Taylor theorem, we have

$$\nabla_i f(x_k) = \nabla_i f(x^*) + \nabla \nabla_i f(x^*)(x_k - x^*) + \frac{1}{2} (x_k - x^*)^T H_i(\hat{x}_k^i) (x_k - x^*) \quad (75)$$

where $\hat{x}_k^i \in \mathcal{N}(x_k, x^*)$.

Using the triangle and Cauchy-Schwartz inequalities and $\|H_i\| \leq M$, we obtain

$$\begin{aligned} \left| \sum_{j=1}^d v_{i,j} (x_k - x^*)^T H_j(\hat{x}_k^j) (x_k - x^*) \right| &\leq \sum_{j=1}^d |v_{i,j}| \left| (x_k - x^*)^T H_j(\hat{x}_k^j) (x_k - x^*) \right| \\ &\leq M \|x_k - x^*\|^2. \end{aligned} \quad (76)$$

From (75), (76) and that $x^* = v_1$, we have

$$\nabla f(x_k)^T v_1 \geq \nabla f(x^*)^T x^* + (x_k - x^*)^T \nabla^2 f(x^*) x^* - \frac{M}{2} \|x_k - x^*\|^2$$

resulting in

$$\begin{aligned} \nabla f(x_k)^T v_1 &\geq \lambda^* - \|\nabla^2 f(x^*) x^*\| \sqrt{2(1 - x_k^T x^*)} - M(1 - x_k^T x^*) \\ &= \lambda^* - \bar{\lambda}_1 \sqrt{(1 - x_k^T x^*)} - M(1 - x_k^T x^*) \end{aligned} \quad (77)$$

For $2 \leq i \leq d$, we have

$$\begin{aligned} \nabla f(x_k)^T v_i &\leq \nabla f(x^*)^T v_i + (x_k - x^*)^T \nabla^2 f(x^*) v_i + \frac{M}{2} \|x_k - x^*\|^2 \\ &= \lambda^*(x^*)^T v_i + (x_k - x^*)^T \nabla^2 f(x^*) v_i + \frac{M}{2} \|x_k - x^*\|^2 \\ &= (x_k - x^*)^T \nabla^2 f(x^*) v_i + \frac{M}{2} \|x_k - x^*\|^2. \end{aligned} \quad (78)$$

Using (78), (72) and (73), we obtain

$$\begin{aligned} \sum_{i=2}^d (\nabla f(x_k)^T v_i)^2 &\leq \sum_{i=2}^d \left((x_k - x^*)^T \nabla^2 f(x^*) v_i + \frac{M}{2} \|x_k - x^*\|^2 \right)^2 \\ &\leq \left(\bar{\lambda}_2 \|x_k - x^*\| + \frac{M}{2} \|x_k - x^*\|^2 \right)^2. \end{aligned} \quad (79)$$

Using (77), (79), and Lemma A.6 with

$$a = \lambda^*, b = M, c = \bar{\lambda}_1, d = 0, e = \bar{\lambda}_2, f = M,$$

we obtain the desired result. \square

A.6 Proofs of Theorem 4.2 and Theorem 4.3

Lemma A.7. Suppose that $f(w, z)$ is scale invariant in $w \in \mathbb{R}^{d_w}$ for each $z \in \mathbb{R}^{d_z}$ and twice continuously differentiable on an open set containing $\partial\mathcal{B}_{d_w} \times \partial\mathcal{B}_{d_z}$. Let (w^*, z^*) be a point satisfying

$$\nabla_w f(w^*, z^*) = \lambda_w^* w^*, \quad \lambda_w^* > \bar{\lambda}_2^w = \max_{2 \leq i \leq d_w} |\lambda_i^w|, \quad w^* = v_1^w$$

where (λ_i^w, v_i^w) is an eigen-pair of $\nabla_{ww}^2 f(w^*, z^*)$. Then, for any $w \in \partial\mathcal{B}_{d_w}$ and $z \in \partial\mathcal{B}_{d_z}$, we have

$$\nabla_w f(w, z)^T v_1^w = \lambda_w^* + (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) w^* + \alpha^w(w, z)$$

and

$$\sum_{i=2}^{d_w} (\nabla_w f(w, z)^T v_i^w)^2 \leq \left(\bar{\lambda}_2^w \sqrt{1 - (w^T w^*)^2} + \nu^{wz} \|z - z^*\| + \beta^w(w, z) \right)^2$$

where

$$\alpha^w(w, z) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\| \right), \quad \beta^w(w, z) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\| \right).$$

Therefore, we have

$$1 - \frac{(\nabla_w f(w, z)^T w^*)^2}{\|\nabla_w f(w, z)\|^2} \leq \left(\frac{\bar{\lambda}_2^w}{\lambda_w^*} \sqrt{1 - (w^T w^*)^2} + \frac{\nu^{wz}}{\lambda_w^*} \|z - z^*\| + \theta^w(w, z) \right)^2$$

where

$$\nu^{wz} = \|\nabla_{wz}^2 f(w^*, z^*)\|, \quad \theta^w(w, z) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\| \right).$$

Proof. Since $\nabla_{ww}^2 f(w^*, z^*)$ is real and symmetric, without loss of generality, we assume that $\{v_1^w, \dots, v_{d_w}^w\}$ forms an orthogonal basis in \mathbb{R}^{d_w} .

By Taylor expansion of $\nabla_w f(w, z)^T v_i^w$ at (w^*, z^*) , we have

$$\nabla_w f(w, z)^T v_i^w = \nabla_x f(w^*, z^*)^T v_i^w + \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix}^T \begin{bmatrix} \nabla_{ww}^2 f(w^*, z^*) \\ \nabla_{zw}^2 f(w^*, z^*) \end{bmatrix} v_i^w + R_i^w(w, z)$$

where

$$R_i^w(w, z) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\| \right).$$

Using $\nabla_w f(w^*, z^*) = \lambda_w^* w^*$ and $w^* = v_1^w$, we have

$$\nabla_w f(w^*, z^*)^T v_1^w = \lambda_w^*, \quad (w - w^*)^T \nabla_{ww}^2 f(w^*, z^*) v_1^w = -\lambda_1^w (1 - w^T w^*).$$

Therefore, we obtain

$$\nabla_w f(w, z)^T v_1^w = \lambda_w^* + (w - w^*)^T \nabla_{zw}^2 f(w^*, z^*) w^* + \alpha^w(w, z) \quad (80)$$

where

$$\alpha^w(w, z) = R_1^w(w, z) - \lambda_1^w (1 - w^T w^*) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\| \right).$$

In the same way, for $2 \leq i \leq d_w$, we have

$$\nabla_w f(w^*, z^*)^T v_i^w = \lambda_i^w (w^*)^T v_i^w = 0, \quad (w - w^*)^T \nabla_{ww}^2 f(w^*, z^*) v_i^w = \lambda_i^w w^T v_i^w,$$

resulting in

$$\nabla_w f(w, z)^T v_i^w = \lambda_i^w w^T v_i^w + (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w + R_i^w(w, z). \quad (81)$$

From (81), we obtain

$$\begin{aligned} \sum_{i=2}^{d_w} (\nabla_w f(w, z)^T v_i^w)^2 &= \sum_{i=2}^{d_w} (\lambda_i^w)^2 (w^T v_i^w)^2 + 2 \sum_{i=2}^{d_w} \lambda_i^w (w^T v_i^w) (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w \\ &\quad + 2 \sum_{i=2}^{d_w} \lambda_i^w (w^T v_i^w) R_i^w(w, z) + 2 \sum_{i=2}^{d_w} (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w R_i^w(w, z) \\ &\quad + \sum_{i=2}^{d_w} ((z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w)^2 + \sum_{i=2}^{d_w} (R_i^w(w, z))^2. \end{aligned}$$

Since $\{v_1^w, \dots, v_{d_w}^w\}$ forms an orthogonal basis in \mathbb{R}^{d_w} , with $w^* = v_1^w$ and $\|w\|^2 = 1$, we have

$$\sum_{i=2}^{d_w} (\lambda_i^w)^2 (w^T v_i^w)^2 \leq (\bar{\lambda}_2^w)^2 (1 - (w^T w^*)^2)$$

and

$$\sum_{i=2}^{d_w} ((z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w)^2 \leq \|(z - z^*)^T \nabla_{zw}^2 f(w^*, z^*)\|^2 \leq (\nu^{wz})^2 \|z - z^*\|^2.$$

Let $\bar{R}_2^w(w, z) = \max_{2 \leq i \leq d_w} |R_i^w(w, z)|$. Note that

$$\bar{R}_2^w(w, z) = o\left(\left\|\begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix}\right\|\right).$$

Using the Cauchy-Shwartz inequality, we have

$$\sum_{i=2}^{d_w} \lambda_i^w (w^T v_i^w) (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w \leq \bar{\lambda}_2^w \nu^{wz} \|z - z^*\| \sqrt{1 - (w^T w^*)^2}.$$

Also, we have

$$\sum_{i=2}^{d_w} \lambda_i^w (w^T v_i^w) R_i^w(w, z) \leq \bar{\lambda}_2^w \bar{R}_2^w(w, z) \sqrt{d_w} \sqrt{1 - (w^T w^*)^2}$$

and

$$\sum_{i=2}^{d_w} R_i^w(w, z) (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) v_i^w \leq \nu^{wz} \bar{R}_2^w(w, z) \sqrt{d_w} \|z - z^*\|.$$

Therefore, we obtain

$$\sum_{i=2}^{d_w} (\nabla_w f(w, z)^T v_i^w)^2 \leq \left(\bar{\lambda}_2^w \sqrt{1 - (w^T w^*)^2} + \nu^{wz} \|z - z^*\| + \beta^w(w, z) \right)^2 \quad (82)$$

where

$$\beta^w(w, z) = \bar{R}_2^w(w, z) \sqrt{d_w} = o\left(\left\|\begin{bmatrix} x_k - x^* \\ y_k - y^* \end{bmatrix}\right\|\right).$$

Since $\{v_1^w, \dots, v_{d_w}^w\}$ forms an orthogonal basis in \mathbb{R}^{d_w} and $|w^T w^*| \leq \|w\| \|w^*\| = 1$, we have

$$1 - \frac{(\nabla_w f(w, z)^T w^*)^2}{\|\nabla_w f(w, z)\|^2} \leq \frac{\sum_{i=2}^{d_w} (\nabla_w f(w, z)^T v_i^w)^2}{(\nabla_w f(w, z)^T v_1^w)^2}.$$

Using (80) and (82), we have

$$\frac{\sum_{i=2}^{d_w} (\nabla_w f(w, z)^T v_i^w)^2}{(\nabla_w f(w, z)^T v_1^w)^2} \leq \left(\frac{\bar{\lambda}_2^w}{\lambda_1^w} \sqrt{1 - (w^T w^*)^2} + \frac{\nu^{wz}}{\lambda_1^w} \|z - z^*\| + \theta^w(w, z) \right)^2$$

where

$$\begin{aligned} \theta^w(w, z) &= \frac{\beta^w(w, z)}{\lambda_w^*} - \left(\frac{\bar{\lambda}_2^w \sqrt{1 - (w^T w^*)^2} + \nu^{wz} \|z - z^*\| + \sqrt{d_w} \beta^w(w, z)}{\lambda_w^*} \right) \\ &\quad \times \left(\frac{(z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) w^* + \beta^w(w, z)}{\lambda_w^* + (z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) w^* + \beta^w(w, z)} \right). \end{aligned}$$

Since

$$|(z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) w^*| \leq \nu^{wz} \|z - z^*\|,$$

we have

$$|(z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) w^*| \sqrt{1 - (w^T w^*)^2} \leq \frac{1}{2} (1 - (w^T w^*)^2) + \frac{1}{2} (\nu^{wz})^2 \|z - z^*\|^2$$

and

$$\nu^{wz} |(z - z^*)^T \nabla_{zw}^2 f(w^*, z^*) w^*| \|z - z^*\| \leq (\nu^{wz})^2 \|z - z^*\|^2.$$

From

$$1 - (w^T w^*)^2 = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\|\right), \quad \|z - z^*\|^2 = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\|\right),$$

we finally obtain

$$\theta^w(w, z) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\|\right).$$

This completes the proof. \square

Lemma A.8. *Suppose that $f(w, z)$ is μ -strongly concave in $z \in \mathbb{R}^{d_z}$ with an L -Lipschitz continuous $\nabla_z f(w, z)$ for each $w \in \partial \mathcal{B}_{d_w}$ and three-times continuously differentiable with respect to x and y on an open set containing $\partial \mathcal{B}_{d_w}$ and \mathbb{R}^{d_z} , respectively. Let (w^*, z^*) be a point such that $\nabla_z f(w^*, z^*) = 0$. Then, for any $w \in \partial \mathcal{B}_{d_w}$ and $z \in \partial \mathcal{B}_{d_z}$, with $\alpha = 2/(L + \mu)$, we have*

$$\|z + \alpha \nabla_z f(w, z) - z^*\| \leq \left(\frac{2\nu^{zw}}{L + \mu} \right) \|w - w^*\| + \left(\frac{L - \mu}{L + \mu} \right) \|z - z^*\| + \theta^z(w, z) \quad (83)$$

where

$$\nu^{zw} = \|\nabla_{zw}^2 f(w^*, z^*)\|, \quad \theta^z(w, z) = o\left(\left\| \begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix} \right\|\right).$$

Proof. Let $\nabla_{z,i} f$ be the i^{th} coordinate of $\nabla_z f$ and

$$H_{z,i} = \begin{bmatrix} H_{z,i}^{ww} & H_{z,i}^{wz} \\ H_{z,i}^{zw} & H_{z,i}^{zz} \end{bmatrix}$$

be the Hessian of $\nabla_{z,i} f$. By Taylor expansion of $\nabla_{z,i} f(w, z)$ at (w^*, z) , we have

$$\nabla_{z,i} f(w, z) = \nabla_{z,i} f(w^*, z) + \nabla_{zw,i}^2 f(w^*, z)^T (w - w^*) + R_i^z(w, z) \quad (84)$$

where $\nabla_{zw,i}^2 f(w^*, z) = \nabla_w \nabla_{z,i} f(w^*, z)$ denotes the i^{th} column of $\nabla_{zw}^2 f(w^*, z)$ and

$$R_i^z(w, z) = \frac{1}{2} (w - w^*)^T H_{z,i}^{ww}(\hat{w}^i, z) (w - w^*), \quad \hat{w}^i \in \mathcal{N}(w, w^*). \quad (85)$$

Also, from f being three-times continuously differentiable, we have

$$\nabla_{zw,i}^2 f(w^*, z) = \nabla_{zw,i}^2 f(w^*, z^*) + H_{z,i}^{wz}(w^*, \hat{z}^i)(z - z^*), \quad \hat{z}^i \in \mathcal{N}(z, z^*). \quad (86)$$

Since

$$\begin{aligned} |(z - z^*)^T H_{z,i}^{zw}(w^*, \hat{z}^i)(w - w^*)| &\leq \|H_{z,i}^{zw}(w^*, \hat{z}^i)\| \|w - w^*\| \|z - z^*\| \\ &\leq \frac{1}{2} \|H_{z,i}^{zw}(w^*, \hat{z}^i)\| (\|w - w^*\|^2 + \|z - z^*\|^2), \end{aligned}$$

we have

$$(z - z^*)^T H_{z,i}^{wz}(w^*, z^i)(w - w^*) = o\left(\left\|\begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix}\right\|\right). \quad (87)$$

By (84), (85), (86), and (87), we have

$$\nabla_z f(w, z) = \nabla_z f(w^*, z) + \nabla_{zw}^2 f(w^*, z^*)(w - w^*) + \bar{R}^z(w, z) \quad (88)$$

where

$$\bar{R}_i^z(w, z) = R_i^z(w, z) + (z - z^*)^T H_{z,i}^{zw}(w^*, z^i)(w - w^*) = o\left(\left\|\begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix}\right\|\right).$$

Using (88), we have

$$z + \alpha \nabla_z f(w, z) - z^* = z - z^* + \alpha \nabla_z f(w^*, z) + \alpha \nabla_{zw}^2 f(w^*, z^*)(w - w^*) + \bar{R}^z(w, z),$$

resulting in

$$\|z + \alpha \nabla_z f(w, z) - z^*\| \leq \|z - z^* + \alpha \nabla_z f(w^*, z)\| + \alpha \|\nabla_{zw}^2 f(w^*, z^*)(w - w^*)\| + \|\bar{R}^z(w, z)\|. \quad (89)$$

Since $-f(w^*, z)$ is μ -strongly convex in z with an L -Lipschitz continuous gradient $-\nabla_z f(w^*, z)$, by theory of convex optimization (see the proof of Theorem 3.12 on page 270 in [2]), we have

$$\|z - z^* + \alpha \nabla_z f(w^*, z)\| \leq \left(\frac{L - \mu}{L + \mu}\right) \|z - z^*\| \quad (90)$$

due to $\alpha = 2/(L + \mu)$. Also, we have

$$\alpha \|\nabla_{zw}^2 f(w^*, z^*)(w - w^*)\| \leq \left(\frac{2\nu^{zw}}{L + \mu}\right) \|w - w^*\|. \quad (91)$$

Plugging (90), (91) into (89), we finally obtain

$$\|z - z^* + \alpha \nabla_z f(w^*, z)\| \leq \left(\frac{L - \mu}{L + \mu}\right) \|z - z^*\| + \left(\frac{2\nu^{zw}}{L + \mu}\right) \|w - w^*\| + \theta^z(w, z)$$

where

$$\theta^z(w, z) = \|\bar{R}^z(w, z)\| = o\left(\left\|\begin{bmatrix} w - w^* \\ z - z^* \end{bmatrix}\right\|\right).$$

□

Lemma A.9. Let M be a 2×2 matrix such that

$$M = \begin{bmatrix} a & e/b \\ e/c & d \end{bmatrix}$$

for some $a > 0, b > 0, c > 0, d \geq 0, e \geq 0$ and let ρ be the largest absolute eigenvalue of M . Then, there exists a sequence ω_t such that

$$\|M^k\| = \prod_{t=0}^{k-1} (\rho + \omega_t) \quad \text{and} \quad \lim_{t \rightarrow \infty} \omega_t = 0.$$

Proof. The characteristic equation reads

$$\det(M - \lambda I) = \lambda^2 - \lambda(a + d) + ad - \frac{e^2}{bc} = 0$$

with the discriminant of

$$(a - d)^2 + \frac{4e^2}{bc} \geq 0.$$

Thus, all eigenvalues are real.

First, we consider the case when $\det(M - \lambda I) = 0$ has a double root. We obtain the condition for a double root as

$$(a - d)^2 + \frac{4e^2}{bc} = 0.$$

Since $b > 0$ and $c > 0$, this implies

$$a = d, \quad e = 0.$$

Therefore, $M = aI$ and $\rho = a$. From $M^k = a^k I$, we have

$$\|M^k\| = \sqrt{a^{2k}} = \rho^k,$$

resulting in

$$\omega_k = \frac{\|M^{k+1}\|}{\|M^k\|} - \rho = \rho - \rho = 0$$

for all $k \geq 0$.

Next, we consider the case when M has two distinct eigenvalues λ_1 and λ_2 . Since $a + d > 0$, we have $\lambda_1 + \lambda_2 > 0$. Without loss of generality, assume $\lambda_1 > \lambda_2$. Then, $\rho = \lambda_1$. Let v_1 and v_2 be corresponding eigenvectors of λ_1 and λ_2 , respectively. Since v_1 and v_2 are linearly independent we can represent each column of M as a linear combination of v_1 and v_2 as

$$M = [\alpha_1 v_1 + \beta_1 v_2 \quad \alpha_2 v_1 + \beta_2 v_2].$$

By repeatedly multiplying M , we obtain

$$M^k = [\alpha_1 \lambda_1^{k-1} v_1 + \beta_1 \lambda_2^{k-1} v_2 \quad \alpha_2 \lambda_1^{k-1} v_1 + \beta_2 \lambda_2^{k-1} v_2].$$

Let $C^k = (M^k)^T M^k$. Then, we have

$$\begin{aligned} C_{11}^k &= \alpha_1^2 \lambda_1^{2(k-1)} + \beta_1^2 \lambda_2^{2(k-1)} + 2\alpha_1 \beta_1 (\lambda_1 \lambda_2)^{k-1} v_1^T v_2 \\ C_{22}^k &= \alpha_2^2 \lambda_1^{2(k-1)} + \beta_2^2 \lambda_2^{2(k-1)} + 2\alpha_2 \beta_2 (\lambda_1 \lambda_2)^{k-1} v_1^T v_2 \end{aligned}$$

and

$$C_{12}^k = \alpha_1 \alpha_2 \lambda_1^{2(k-1)} + \beta_1 \beta_2 \lambda_2^{2(k-1)} + (\alpha_1 \beta_2 + \alpha_2 \beta_1) (\lambda_1 \lambda_2)^{k-1} v_1^T v_2, \quad C_{21}^k = C_{12}^k.$$

Since

$$C_{11}^k \geq \alpha_1^2 \lambda_1^{2(k-1)} + \beta_1^2 \lambda_2^{2(k-1)} - 2\alpha_1 \beta_1 (\lambda_1 \lambda_2)^{k-1} = (\alpha_1 \lambda_1^{k-1} - \beta_1 \lambda_2^{k-1})^2 \geq 0$$

and

$$C_{22}^k \geq \alpha_2^2 \lambda_1^{2(k-1)} + \beta_2^2 \lambda_2^{2(k-1)} - 2\alpha_2 \beta_2 (\lambda_1 \lambda_2)^{k-1} = (\alpha_2 \lambda_1^{k-1} - \beta_2 \lambda_2^{k-1})^2 \geq 0,$$

we have

$$\|M^k\| = \sqrt{\frac{1}{2} \left[C_{11}^k + C_{22}^k + \sqrt{(C_{11}^k - C_{22}^k)^2 + 4(C_{12}^k)^2} \right]},$$

leading to

$$\frac{\|M^{k+1}\|}{\|M^k\|} = \sqrt{\frac{C_{11}^{k+1} + C_{22}^{k+1} + \sqrt{(C_{11}^{k+1} - C_{22}^{k+1})^2 + 4(C_{12}^{k+1})^2}}{C_{11}^k + C_{22}^k + \sqrt{(C_{11}^k - C_{22}^k)^2 + 4(C_{12}^k)^2}}}.$$

From

$$\lim_{k \rightarrow \infty} \frac{C_{11}^k}{\lambda_1^{2(k-1)}} = \alpha_1^2, \quad \lim_{k \rightarrow \infty} \frac{C_{22}^k}{\lambda_1^{2(k-1)}} = \alpha_2^2, \quad \lim_{k \rightarrow \infty} \frac{C_{12}^k}{\lambda_1^{2(k-1)}} = \lim_{k \rightarrow \infty} \frac{C_{21}^k}{\lambda_1^{2(k-1)}} = \alpha_1 \alpha_2,$$

we obtain

$$\lim_{k \rightarrow \infty} \frac{\|M^{k+1}\|}{\|M^k\|} = \sqrt{\lambda_1^2} = \rho.$$

From

$$\lim_{k \rightarrow \infty} \omega_k = \lim_{k \rightarrow \infty} \frac{\|M^{k+1}\|}{\|M^k\|} - \rho = \rho - \rho = 0,$$

we obtain the desired result. \square

Proof of Theorem 4.2. From Lemma A.7 with $w = x_k, z = y_k$, we have

$$1 - \frac{(\nabla_x f(x_k, y_k)^T x^*)^2}{\|\nabla_x f(x_k, y_k)\|^2} \leq \left(\frac{\bar{\lambda}_2}{\lambda^*} \sqrt{1 - (x_k^T x^*)^2} + \frac{\nu}{\lambda^*} \|y_k - y^*\| + \theta^x(x_k, y_k) \right)^2.$$

Since

$$x_{k+1} = \frac{\nabla_x f(x_k, y_k)}{\|\nabla_x f(x_k, y_k)\|},$$

we obtain

$$\sqrt{1 - (x_{k+1}^T x^*)^2} \leq \frac{\bar{\lambda}_2}{\lambda^*} \sqrt{1 - (x_k^T x^*)^2} + \frac{\nu}{\lambda^*} \|y_k - y^*\| + \theta^x(x_k, y_k).$$

Using

$$\|y_k - y^*\| = \sqrt{2(1 - y_k^T y^*)} = \left(1 + \frac{1 - y_k^T y^*}{1 + y_k^T y^* + \sqrt{2(1 + y_k^T y^*)}} \right) \sqrt{1 - (y_k^T y^*)^2},$$

we have

$$\sqrt{1 - (x_{k+1}^T x^*)^2} \leq \frac{\bar{\lambda}_2}{\lambda^*} \sqrt{1 - (x_k^T x^*)^2} + \frac{\nu}{\lambda^*} \sqrt{1 - (y_k^T y^*)^2} + \bar{\theta}^x(x_k, y_k) \quad (92)$$

where

$$\bar{\theta}^x(x_k, y_k) = \theta^x(x_k, y_k) + \left(\frac{1 - y_k^T y^*}{1 + y_k^T y^* + \sqrt{2(1 + y_k^T y^*)}} \right) \sqrt{1 - (y_k^T y^*)^2} = o\left(\left\| \begin{bmatrix} x_k - x^* \\ y_k - y^* \end{bmatrix} \right\|\right).$$

Using Lemma A.7 for $w = y_k, z = x_k$ and the definition of y_{k+1} , we have

$$\sqrt{1 - (y_{k+1}^T y^*)^2} \leq \frac{\nu}{s^*} \sqrt{1 - (x_k^T x^*)^2} + \frac{\bar{s}_2}{s^*} \sqrt{1 - (y_k^T y^*)^2} + \bar{\theta}^y(x_k, y_k) \quad (93)$$

where

$$\bar{\theta}^y(x_k, y_k) = \theta^y(x_k, y_k) + \left(\frac{1 - x_k^T x^*}{1 + x_k^T x^* + \sqrt{2(1 + x_k^T x^*)}} \right) \sqrt{1 - (x_k^T x^*)^2} = o\left(\left\| \begin{bmatrix} x_k - x^* \\ y_k - y^* \end{bmatrix} \right\|\right).$$

Combining (92) and (93), we obtain

$$\begin{bmatrix} \sqrt{1 - (x_{k+1}^T x^*)^2} \\ \sqrt{1 - (y_{k+1}^T y^*)^2} \end{bmatrix} \leq \begin{bmatrix} \frac{\bar{\lambda}_2}{\lambda^*} & \frac{\nu}{\lambda^*} \\ \frac{\nu}{s^*} & \frac{\bar{s}_2}{s^*} \end{bmatrix} \begin{bmatrix} \sqrt{1 - (x_k^T x^*)^2} \\ \sqrt{1 - (y_k^T y^*)^2} \end{bmatrix} + \begin{bmatrix} \bar{\theta}^x(x_k, y_k) \\ \bar{\theta}^y(x_k, y_k) \end{bmatrix} \quad (94)$$

$$\leq (M + N(x_k, y_k)) \begin{bmatrix} \sqrt{1 - (x_k^T x^*)^2} \\ \sqrt{1 - (y_k^T y^*)^2} \end{bmatrix} \quad (95)$$

where

$$M = \begin{bmatrix} \frac{\bar{\lambda}_2}{\lambda^*} & \frac{\nu}{\lambda^*} \\ \frac{\nu}{s^*} & \frac{\bar{s}_2}{s^*} \end{bmatrix}, \quad N(x, y) = \frac{\epsilon(x, y)}{\sqrt{2 - x^T x^* - y^T y^*}} \begin{bmatrix} \sqrt{\frac{1 - x^T x^*}{1 + x^T x^*}} & \sqrt{\frac{1 - y^T y^*}{1 + y^T y^*}} \\ \sqrt{\frac{1 - x^T x^*}{1 + x^T x^*}} & \sqrt{\frac{1 - y^T y^*}{1 + y^T y^*}} \end{bmatrix}$$

and

$$\epsilon(x, y) = \frac{\max\{\bar{\theta}^x(x, y), \bar{\theta}^y(x, y)\}}{\sqrt{2 - x^T x^* - y^T y^*}}.$$

Note that the spectral radius ρ of M satisfies

$$\rho = \frac{1}{2} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \frac{\bar{s}_2}{s^*} + \sqrt{\left(\frac{\bar{\lambda}_2}{\lambda^*} - \frac{\bar{s}_2}{s^*} \right)^2 + \frac{4\nu^2}{\lambda^* s^*}} \right) < 1.$$

due to $\nu^2 < (\lambda^* - \bar{\lambda}_2)(s^* - \bar{s}_2)$. Also, for $i, j = 1, 2$, we have

$$\lim_{(x,y) \rightarrow (x^*, y^*)} N_{ij}(x, y) = 0.$$

By Lemma A.9, there exists a sequence ω_t such that

$$\|M^k\| = \prod_{t=0}^{k-1} (\rho + \omega_t) \quad \text{and} \quad \lim_{t \rightarrow \infty} \omega_t = 0.$$

Let

$$\tau = \min\{k : \|M^k\| < 1\}, \quad \bar{\rho} = \frac{\|M^\tau\| + 1}{2}, \quad \rho_{\max} = \max_{1 \leq k \leq \tau} \|M^k\|.$$

We first show that $(x_{n\tau}, y_{n\tau}) \rightarrow (x^*, y^*)$ as $n \rightarrow \infty$. By Lemma A.7, we have

$$\begin{aligned} \nabla_x f(x, y)^T v_1 &= \lambda^* + (y - y^*)^T \nabla_{yx}^2 f(x^*, y^*) x^* + \alpha^x(x, y) \\ \nabla_y f(x, y)^T u_1 &= s^* + (x - x^*)^T \nabla_{xy}^2 f(x^*, y^*) y^* + \alpha^y(x, y) \end{aligned}$$

where

$$\alpha^x(x, y) = o\left(\left\| \begin{bmatrix} x - x^* \\ y - y^* \end{bmatrix} \right\|\right), \quad \alpha^y(x, y) = o\left(\left\| \begin{bmatrix} x - x^* \\ y - y^* \end{bmatrix} \right\|\right).$$

Therefore, there exists some $\delta_1 > 0$ such that if

$$x^T x^* > 0, \quad y^T y^* > 0, \quad \left\| \begin{bmatrix} \sqrt{1 - (x^T x^*)^2} \\ \sqrt{1 - (y^T y^*)^2} \end{bmatrix} \right\| < \delta_1,$$

then

$$\nabla_x f(x, y)^T v_1 > 0, \quad \nabla_y f(x, y)^T u_1 > 0. \quad (96)$$

Since $N_{ij}(x, y) \rightarrow 0$ as $(x, y) \rightarrow (x^*, y^*)$ for $i, j = 1, 2$, there exists some $\delta_2 > 0$ such that if

$$x^T x^* > 0, \quad y^T y^* > 0, \quad \left\| \begin{bmatrix} \sqrt{1 - (x^T x^*)^2} \\ \sqrt{1 - (y^T y^*)^2} \end{bmatrix} \right\| < \delta_2,$$

then we have

$$\left\| \prod_{l=0}^{\tau-1} (M + N(\phi(x, y, l))) \right\| < \bar{\rho}, \quad \max_{0 < m \leq \tau} \left\| \prod_{l=0}^{m-1} (M + N(\phi(x, y, l))) \right\| < 1 + \rho_{\max} \quad (97)$$

where $\phi(x, y, l)$ denotes the vector after l iterations of the algorithm starting with (x, y) .

To see this, let us define

$$g(x, y, m) = \left\| \prod_{l=0}^{m-1} (M + N(\phi(x, y, l))) \right\|.$$

By (95) and (96), if $x \rightarrow x^*$ and $y \rightarrow y^*$, then for any $0 \leq l \leq \tau$, we have

$$\phi(x, y, l) \rightarrow (x^*, y^*),$$

resulting in

$$g(x, y, m) \rightarrow \|M^m\|.$$

Therefore, by the definition of $\bar{\rho}$ and τ , there exists some $\delta_{2,\tau} > 0$ such that $g(x, y, \tau) < \bar{\rho}$. Also, for each $1 \leq m < \tau$, there exists some $\delta_{2,m} > 0$ such that $g(x, y, m) < 1 + \rho_{\max}$. Taking the minimum of $\delta_{2,m}$ for $1 \leq m \leq \tau$, we obtain δ_2 satisfying (97).

Let

$$\bar{\delta} = \min \left\{ \delta_1, \frac{\delta_1}{1 + \rho_{\max}}, \delta_2 \right\}, \quad N_k = N(x_k, y_k).$$

By mathematical induction, we show that for any $n \geq 0$, if

$$x_{n\tau}^T x^* > 0, \quad y_{n\tau}^T y^* > 0, \quad \Delta_{n\tau} < \bar{\delta}, \quad (98)$$

then for $0 \leq m \leq \tau$, we have

$$x_{n\tau+m}^T x^* > 0, \quad y_{n\tau+m}^T y^* > 0, \quad \Delta_{n\tau+m} \leq (1 + \rho_{\max}) \Delta_{n\tau} < \delta_1. \quad (99)$$

By (98), it is obvious that we have (99) for $m = 0$. This proves the base case.

Suppose that we have (99) for $0 \leq m < \tau$. Then, by the definition of δ_1 , we have

$$\begin{aligned} x_{n\tau+m+1}^T x^* &= x_{n\tau+m+1}^T v_1 = \frac{\nabla_x f(x_{n\tau+m}, y_{n\tau+m})^T v_1}{\|\nabla_x f(x_{n\tau+m}, y_{n\tau+m})\|} > 0, \\ y_{n\tau+m+1}^T y^* &= y_{n\tau+m+1}^T u_1 = \frac{\nabla_y f(x_{n\tau+m}, y_{n\tau+m})^T u_1}{\|\nabla_y f(x_{n\tau+m}, y_{n\tau+m})\|} > 0. \end{aligned}$$

Also, by (95), (98) and (97), we have

$$\Delta_{n\tau+m+1} \leq \left\| \prod_{l=0}^m (M + N_{n\tau+l}) \right\| \Delta_{n\tau} \leq (1 + \rho_{\max}) \Delta_{n\tau} < \delta_1.$$

This completes the induction proof.

Suppose that (x_0, y_0) satisfies $\max\{1 - x_0^T x^*, 1 - y_0^T y^*\} < \delta$ where $\delta = \min\{\bar{\delta}, 1\}$. Since $\delta < 1$, we have

$$x_0^T x^* > 0, \quad y_0^T y^* > 0, \quad \Delta_0 < \bar{\delta}. \quad (100)$$

Next, we show

$$x_{n\tau}^T x^* > 0, \quad y_{n\tau}^T y^* > 0, \quad \Delta_{n\tau} \leq \bar{\rho}^n \Delta_0. \quad (101)$$

For $n = 0$, we have (101) by (100). This proves the base case.

Suppose that we have (101) for n . Then, since $\Delta_{n\tau} \leq \bar{\rho}^n \Delta_0 < \bar{\delta}$, by (98) and (99), we have

$$x_{(n+1)\tau}^T x^* > 0, \quad y_{(n+1)\tau}^T y^* > 0.$$

Moreover, using (95) and (97), we have

$$\Delta_{(n+1)\tau} \leq \left\| \prod_{l=0}^{\tau-1} (M + N_{n\tau+l}) \right\| \Delta_{n\tau} \leq \bar{\rho} \Delta_{n\tau} < \bar{\rho}^{n+1} \Delta_0,$$

which completes the induction proof. Therefore, by (101), $(x_{n\tau}, y_{n\tau}) \rightarrow (x^*, y^*)$ as $n \rightarrow \infty$.

Furthermore, due to (99), we have $(x_{n\tau+m}, y_{n\tau+m}) \rightarrow (x^*, y^*)$ for every $0 < m \leq \tau$, indicating that $(x_k, y_k) \rightarrow (x^*, y^*)$. This in turn implies that $N_k \rightarrow 0$. Letting

$$\eta_k = \frac{\|\prod_{t=0}^k (M + N_t)\|}{\|\prod_{t=0}^{k-1} (M + N_t)\|} - \frac{\|M^{k+1}\|}{\|M^k\|}, \quad \gamma_k = \omega_k + \eta_k,$$

we have

$$\left\| \prod_{t=0}^{k-1} (M + N_t) \right\| = \prod_{t=0}^{k-1} (\rho + \omega_t + \eta_t) = \prod_{t=0}^{k-1} (\rho + \gamma_t). \quad (102)$$

Since $\eta_k \rightarrow 0$ as $N_k \rightarrow 0$, we have $\lim \gamma_k = 0$. This concludes the proof. \square

Proof of Theorem 4.3. Using Lemma A.7 for $w = x_k$, $z = y_k$ and the definition of x_{k+1} , we have

$$\sqrt{1 - (x_{k+1}^T x^*)^2} \leq \frac{\bar{\lambda}_2}{\lambda^*} \sqrt{1 - (x_k^T x^*)^2} + \frac{\nu}{\lambda^*} \|y_k - y^*\| + \theta^x(x_k, y_k). \quad (103)$$

By Lemma A.8 with $w = x_k, z = y_k$, we also have

$$\|y_{k+1} - y^*\| \leq \left(\frac{2\nu}{L+\mu}\right) \|x_k - x^*\| + \left(\frac{L-\mu}{L+\mu}\right) \|y_k - y^*\| + \theta^y(x_k, y_k). \quad (104)$$

Using

$$\bar{\theta}^y(x_k, y_k) = \theta^y(x_k, y_k) + \left(\frac{1 - x_k^T x^*}{1 + x_k^T x^* + \sqrt{2(1 + x_k^T x^*)}}\right) \sqrt{1 - (x_k^T x^*)^2} = o\left(\left\|\begin{bmatrix} x_k - x^* \\ y_k - y^* \end{bmatrix}\right\|\right),$$

we can write (104) as

$$\|y_{k+1} - y^*\| \leq \left(\frac{2\nu}{L+\mu}\right) \sqrt{1 - (x_k^T x^*)^2} + \left(\frac{L-\mu}{L+\mu}\right) \|y_k - y^*\| + \bar{\theta}^y(x_k, y_k). \quad (105)$$

Combining (103) and (105), we obtain

$$\begin{bmatrix} \sqrt{1 - (x_{k+1}^T x^*)^2} \\ \|y_{k+1} - y^*\| \end{bmatrix} \leq \begin{bmatrix} \frac{\bar{\lambda}_2}{\lambda^*} & \frac{\nu}{\lambda^*} \\ \frac{2\nu}{L+\mu} & \frac{L-\mu}{L+\mu} \end{bmatrix} \begin{bmatrix} \sqrt{1 - (x_k^T x^*)^2} \\ \|y_k - y^*\| \end{bmatrix} + \begin{bmatrix} \theta^x(x_k, y_k) \\ \bar{\theta}^y(x_k, y_k) \end{bmatrix} \quad (106)$$

$$\leq (M + N(x_k, y_k)) \begin{bmatrix} \sqrt{1 - (x_k^T x^*)^2} \\ \|y_k - y^*\| \end{bmatrix} \quad (107)$$

where

$$M = \begin{bmatrix} \frac{\bar{\lambda}_2}{\lambda^*} & \frac{\nu}{\lambda^*} \\ \frac{2\nu}{L+\mu} & \frac{L-\mu}{L+\mu} \end{bmatrix}, \quad N(x, y) = \frac{\epsilon(x, y)}{\sqrt{1 - x^T x^* + \|y - y^*\|^2}} \begin{bmatrix} \sqrt{\frac{1 - x^T x^*}{1 + x^T x^*}} \|y - y^*\| \\ \sqrt{\frac{1 - x^T x^*}{1 + x^T x^*}} \|y - y^*\| \end{bmatrix}$$

and

$$\epsilon(x, y) = \frac{\max\{\theta^x(x, y), \bar{\theta}^y(x, y)\}}{\sqrt{1 - x^T x^* + \|y - y^*\|^2}}.$$

Since $\nu^2 < \mu(\lambda^* - \bar{\lambda}_2)$, the spectral radius ρ of M satisfies

$$\rho = \frac{1}{2} \left(\frac{\bar{\lambda}_2}{\lambda^*} + \frac{L-\mu}{L+\mu} + \sqrt{\left(\frac{\bar{\lambda}_2}{\lambda^*} - \frac{L-\mu}{L+\mu}\right)^2 + \frac{8\nu^2}{\lambda^*(L+\mu)}} \right) < 1.$$

The rest of the proof is the same as the steps taken in the proof of Theorem 4.2. \square