
Block Coordinate Descent Methods for Optimization under J-Orthogonality Constraints with Applications

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The J-orthogonal matrix, also referred to as the hyperbolic orthogonal matrix, is
2 a class of special orthogonal matrix in hyperbolic space, notable for its advanta-
3 geous properties. These matrices are integral to optimization under J-orthogonal
4 constraints, which have widespread applications in statistical learning and data
5 science. However, addressing these problems is generally challenging due to
6 their non-convex nature and the computational intensity of the constraints. Cur-
7 rently, algorithms for tackling these challenges are limited. This paper introduces
8 **JOB**CD, a novel Block Coordinate Descent method designed to address opti-
9 mizations with J-orthogonality constraints. We explore two specific variants of
10 **JOB**CD: one based on a Gauss-Seidel strategy (**GS-JOB**CD), the other on a
11 variance-reduced and Jacobi strategy (**VR-J-JOB**CD). Notably, leveraging the
12 parallel framework of a Jacobi strategy, **VR-J-JOB**CD integrates variance reduc-
13 tion techniques to decrease oracle complexity in the minimization of finite-sum
14 functions. For both **GS-JOB**CD and **VR-J-JOB**CD, we establish the oracle com-
15 plexity under mild conditions and strong limit-point convergence results under the
16 Kurdyka-Lojasiewicz inequality. To demonstrate the effectiveness of our method,
17 we conduct experiments on hyperbolic eigenvalue problems, hyperbolic structural
18 probe problems, and the ultrahyperbolic knowledge graph embedding problem.
19 Extensive experiments using both real-world and synthetic data demonstrate that
20 **JOB**CD consistently outperforms state-of-the-art solutions, by large margins.

21 1 Introduction

22 A matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$ is a J-orthogonal matrix if $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$, where $\mathbf{J} = \begin{bmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_{n-p} \end{bmatrix}$, and \mathbf{I}_p is a $p \times p$
23 identity matrix. Here, $\mathbf{J} \in \mathbb{R}^{n \times n}$ is the signature matrix with signature $(p, n - p)$. In this paper, we
24 mainly focus on the following optimization problem under J-orthogonality constraints:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} f(\mathbf{X}) \triangleq \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{X}), \text{ s. t. } \mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}. \quad (1)$$

25 Here, $f(\mathbf{X})$ could have a finite-sum structure, each component function $f_i(\mathbf{X})$ is assumed to be
26 differentiable, and N is the number of data points. For brevity, the J-orthogonality constraint
27 $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$ in Problem (1) is rewritten as $\mathbf{X} \in \mathcal{J}$.

28 We impose the following assumptions on Problem (1) throughout this paper. (A-i) For any matrices
29 \mathbf{X} and \mathbf{X}^+ , we assume $f_i : \mathbb{R}^{n \times n} \mapsto \mathbb{R}$ is continuously differentiable for some symmetric positive
30 semidefinite matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ that:

$$f_i(\mathbf{X}^+) \leq f_i(\mathbf{X}) + \langle \mathbf{X}^+ - \mathbf{X}, \nabla f_i(\mathbf{X}) \rangle + \frac{1}{2} \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbf{H}}^2, \quad (2)$$

31 for all $i \in [N]$, where $\|\mathbf{H}\| \leq L_f$ for some constant $L_f > 0$ and $\|\mathbf{X}\|_{\mathbf{H}}^2 \triangleq \text{vec}(\mathbf{X})^\top \mathbf{H} \text{vec}(\mathbf{X})$.
32 This further implies that: $\|\nabla f_i(\mathbf{X}) - \nabla f_i(\mathbf{X}^+)\|_{\mathbf{F}} \leq L_f \|\mathbf{X} - \mathbf{X}^+\|_{\mathbf{F}}$ for all $i \in [N]$. Import-
33 tantly, the function $f(\mathbf{X}) = \frac{1}{2} \text{tr}(\mathbf{X}^\top \mathbf{C} \mathbf{X} \mathbf{D}) = \frac{1}{2} \|\mathbf{X}\|_{\mathbf{H}}^2$ with $\mathbf{H} = \mathbf{D} \otimes \mathbf{C}$ satisfies the equality

34 $\forall \mathbf{X}, \mathbf{X}^+, f(\mathbf{X}^+) = Q(\mathbf{X}^+; \mathbf{X})$ in (2), where $\mathbf{C} \in \mathbb{R}^{n \times n}$ and $\mathbf{D} \in \mathbb{R}^{n \times n}$ are arbitrary symmetric
35 matrices. (A-ii) The function $f_i(\mathbf{X})$ is coercive for all $i \in N$, that is, $\lim_{\|\mathbf{X}\|_F \rightarrow \infty} f_i(\mathbf{X}) = \infty, \forall i$.

36 Problem (1) defines an optimization framework that is fundamental to a wide range of models in
37 statistical learning and data science, including hyperbolic eigenvalue problem [6, 43, 40], hyperbolic
38 structural probe problem [20, 7], and ultrahyperbolic knowledge graph embedding [48]. Additionally,
39 it is closely related to machine learning in hyperbolic spaces, including Lorentz model learning
40 [35, 50, 8] and ultrahyperbolic neural networks [27, 54, 42]. It also intersects with hyperbolic linear
41 algebra [3, 21], addressing problems such as the indefinite least squares problem, hyperbolic QR
42 factorization, and indefinite polar decomposition.

43 1.1 Related Work

44 ► **Block Coordinate Descent Methods.** Block Coordinate Descent (BCD) is a well-established
45 iterative algorithm that sequentially minimizes along block coordinate directions. Its simplicity
46 and efficiency have led to its widespread adoption in structured convex applications [37]. Recently,
47 BCD has gained traction in non-convex problems due to its robust optimality guarantees and/or
48 excellent empirical performance in areas including optimal transport [22], matrix optimization [12],
49 fractional minimization [52], deep neural networks [5, 53, 32], federated learning [47], black-box
50 optimization [4], and optimization with orthogonality constraints [51, 14]. To our knowledge, this is
51 the first application of BCD methods to optimization under J-orthogonality constraints, with a focus
52 on analyzing their theoretical guarantees and empirical efficacy.

53 ► **Minimizing Smooth Functions under J-Orthogonality Constraints.** The J-orthogonal matrix
54 belongs to a subset of generalized orthogonal matrices [16, 36, 23]. However, projecting onto the
55 J-orthogonality constraint poses challenges, complicating the extension of conventional optimization
56 algorithms to address optimization problems under these constraints [1, 16]. This contrasts with
57 computing orthogonal projections using methods such as polar or SVD decomposition, or approxi-
58 mating them via QR factorization. Existing methods for addressing Problem (1) can be categorized
59 into three classes. (i) CS-Decomposition Based Methods. These approaches involve parameterizing
60 four orthogonal matrices (as described in Proposition 2.2) and subsequently minimizing a smooth
61 function over these matrices in an alternating fashion. The involvement of 3×3 block matrices makes
62 the implementation of these methods very challenging. Consequently, the work of [48] focuses on
63 optimizing a reduced subspace of the CS decomposition parameters, albeit at the expense of losing
64 some degrees of freedom. (ii) Unconstrained Multiplier Correction Methods [31, 13, 14]. These
65 methods leverage the symmetry and explicit closed-form expression of the Lagrangian multiplier at
66 the first-order optimality condition. Consequently, they address an unconstrained problem, resulting
67 in efficient first-order infeasible approaches. (iii) Alternating Direction Method of Multipliers [19].
68 This method reformulates the original problem into a bilinear constrained optimization problem by
69 introducing auxiliary variables. It employs dual variables to handle bilinear constraints, iteratively
70 optimizing primal variables while keeping other primal and dual variables fixed, and using a gradient
71 ascent strategy to update the dual variables. This approach has become widely adopted for solving
72 general nonconvex and nonsmooth composite optimization problems. Notably, all the aforementioned
73 methods solely identify critical points of Problem (1).

74 ► **Finite-Sum Problems via Stochastic Gradient Descent.** The finite-sum structure is prevalent in
75 machine learning and statistical modeling, facilitating decomposition into smaller, more manageable
76 components. This property is advantageous for developing efficient algorithms for large-scale prob-
77 lems, such as Stochastic Gradient Descent (SGD). Reducing variance is crucial in SGD because it can
78 lead to more stable and faster convergence. Various techniques, such as mini-batch SGD, momentum
79 methods, and variance reduction methods like SAGA [10], SVRG [25], SARAH [34], SPIDER
80 [11, 44], SNVRG [55], and PAGE [30], have been developed to address this issue. Additionally, SGD
81 for minimizing composite functions has also been investigated by the authors [15, 24, 29].

82 1.2 Contributions

83 This paper makes the following contributions. (i) Algorithmically: We introduce the **JOB**
84 **CD** algorithm, a novel Block Coordinate Descent method specifically designed to tackle optimizations
85 constrained by J-orthogonality. We explore two specific variants of **JOB**
86 **CD**, one based on a Gauss-Seidel strategy (**GS-JOB**
87 **CD**), the other on a variance-reduced and Jacobi strategy (**VR-**
88 **J-JOB**
89 **CD**). Notably, **VR-J-JOB**
89 **CD** incorporates a variance-reduction technique into a parallel
2). (ii) Theoretically: We provide comprehensive optimality and convergence analyses for both

90 algorithms (see Sections 3 and 4). (iii) Empirically: Extensive experiments across hyperbolic
 91 eigenvalue problems, structural probe problems, and ultrahyperbolic knowledge graph embedding,
 92 using both real-world and synthetic data, consistently show the significant superiority of **JOB**
 93 over state-of-the-art solutions (see Section 5).

94 2 The Proposed **JOB** Algorithm

95 This section proposes **JOB** for solving optimization problems under J-orthogonality constraints
 96 in Problem (1), which is based on randomized block coordinate descent. Two variants of **JOB**
 97 are explored, one based on a Gauss-Seidel strategy (**GS-JOB**), the other on a variance-reduced and
 98 Jacobi strategy (**VR-JOB**).

99 **Notations.** We define $[n] \triangleq \{1, 2, \dots, n\}$. We denote $\Omega \triangleq \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{\binom{n}{2}}\}$ as all the possible
 100 combinations of the index vectors choosing 2 items from n without repetition. For any $B \in \Omega$, we
 101 define $\mathbf{U}_B \in \mathbb{R}^{n \times 2}$ as $(\mathbf{U}_B)_{ji} = 1$ if $B_i = j$, else 0 for all j and i , leading to $\mathbf{U}_B^\top \mathbf{X} = \mathbf{X}(B, :) \in$
 102 $\mathbb{R}^{2 \times n}$. We denote $\mathcal{J}_B \triangleq \{\mathbf{V} \mid \mathbf{V}^\top \mathbf{J}_{BB} \mathbf{V} = \mathbf{J}_{BB}\}$, where $\mathbf{J}_{BB} \in \mathbb{R}^{2 \times 2}$ is the sub-matrix of \mathbf{J} indexed by
 103 B . Further notations are provided in Appendix A.1.

104 2.1 Gauss-Seidel Block Coordinate Descent Algorithm

105 This subsection describes the proposed **GS-JOB** algorithm. We consider Problem (1) with $N = 1$
 106 only, without utilizing its finite-sum structure.

107 **GS-JOB** is an iterative algorithm that, in each iteration t , randomly and uniformly (with replace-
 108 ment) selects a coordinate B from the set Ω and then solves a small-sized subproblem. The row index
 109 $[n]$ of the decision variable \mathbf{X} are separated to two sets B and B^c , where $B \in \Omega$ with $|B| = 2$ is the work-
 110 ing set and $B^c = [n] \setminus B$. For simplicity, we use B instead of B^t . Following [51], we consider the follow-
 111 ing block coordinate update rule: $[\mathbf{X}^{t+1}(B, :) = \mathbf{V}\mathbf{X}^t(B, :)] \Leftrightarrow [\mathbf{X}^{t+1} = \mathbf{X}^t + \mathbf{U}_B(\mathbf{V} - \mathbf{I})\mathbf{U}_B^\top \mathbf{X}^t]$,
 112 where $\mathbf{V} \in \mathbb{R}^{2 \times 2}$ is some suitable matrix.

113 The following lemma illustrates matrix selection for enforcing J-orthogonality constraints via the
 114 update rule $\mathbf{X}^+ \leftarrow \mathcal{X}_B(\mathbf{V}) \triangleq \mathbf{X} + \mathbf{U}_B(\mathbf{V} - \mathbf{I})\mathbf{U}_B^\top \mathbf{X}$, and presents associated properties.

115 **Lemma 2.1.** (Proof in Section C.1) For any $B \in \Omega$, we define $\mathbf{X}^+ \triangleq \mathcal{X}_B(\mathbf{V}) \triangleq \mathbf{X} + \mathbf{U}_B(\mathbf{V} - \mathbf{I})\mathbf{U}_B^\top \mathbf{X}$.
 116 We have: (a) If $\mathbf{V} \in \mathcal{J}_B$ and $\mathbf{X} \in \mathcal{J}$, then $\mathbf{X}^+ \in \mathcal{J}$. (b) $\|\mathbf{X}^+ - \mathbf{X}\|_F^2 \leq \|\mathbf{X}\|_F^2 \cdot \|\mathbf{V} - \mathbf{I}\|_F^2$. (c)
 117 $\|\mathbf{X}^+ - \mathbf{X}\|_H^2 \leq \|\mathbf{V} - \mathbf{I}\|_Q^2$ for all $\mathbf{Q} \succcurlyeq \underline{\mathbf{Q}} \triangleq (\mathbf{Z}^\top \otimes \mathbf{U}_B)^\top \mathbf{H}(\mathbf{Z}^\top \otimes \mathbf{U}_B)$, $\mathbf{Z} \triangleq \mathbf{U}_B^\top \mathbf{X} \in \mathbb{R}^{k \times n}$.

118 **► The Main Algorithm.** Using the above update rule, we consider the following iterative procedure:
 119 $\mathbf{X}^{t+1} \leftarrow \mathcal{X}_B(\bar{\mathbf{V}}^t)$, where $\bar{\mathbf{V}}^t \in \arg \min_{\mathbf{V}} f(\mathcal{X}_B^t(\mathbf{V}))$. However, the resulting subproblem could be
 120 still difficult to solve. This inspires us to use sequential majorization minimization [38, 33] to address
 121 it. This technique iteratively constructs a surrogate function that upper-bounds the objective function,
 122 allowing for effective optimization and gradual reduction of the objective function. We derive:

$$\begin{aligned} f(\mathcal{X}_B^t(\mathbf{V})) &\stackrel{\textcircled{1}}{\leq} f(\mathbf{X}^t) + \frac{1}{2} \|\mathcal{X}_B^t(\mathbf{V}) - \mathbf{X}^t\|_H^2 + \langle \mathcal{X}_B^t(\mathbf{V}) - \mathbf{X}^t, \nabla f(\mathbf{X}^t) \rangle \\ &\stackrel{\textcircled{2}}{\leq} f(\mathbf{X}^t) + \frac{1}{2} \|\mathbf{V} - \mathbf{I}\|_{\mathbf{Q} + \theta \mathbf{I}}^2 + \langle \mathbf{V} - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{BB} \rangle \triangleq \mathcal{G}(\mathbf{V}; \mathbf{X}^t, B^t), \end{aligned} \quad (3)$$

123 where step $\textcircled{1}$ uses Inequality (2); step $\textcircled{2}$ uses Claim (c) of Lemma 2.1, $\theta \geq 0$ and the fact that
 124 $\langle \mathbf{U}_B(\mathbf{V} - \mathbf{I})\mathbf{U}_B^\top \mathbf{X}, \nabla f(\mathbf{X}) \rangle = \langle \mathbf{V} - \mathbf{I}, [\nabla f(\mathbf{X})\mathbf{X}^\top]_{BB} \rangle$, and the choice of $\mathbf{Q} \in \mathbb{R}^{4 \times 4}$ that:

$$\mathbf{Q} = \underline{\mathbf{Q}}, \text{ or } \mathbf{Q} = \varsigma \mathbf{I}_2, \text{ with } \|\underline{\mathbf{Q}}\| \leq \varsigma \leq L_f. \quad (4)$$

125 Therefore, the function $\mathcal{G}(\mathbf{V}; \mathbf{X}^t, B^t)$ becomes a majorization function of $f(\mathbf{X})$ at $\mathbf{X}^t \in \mathcal{J}$ for all $B^t \in$
 126 Ω . We can consider the following optimization problem to find $\bar{\mathbf{V}}^t$: $\bar{\mathbf{V}}^t \in \arg \min_{\mathbf{V}} \mathcal{G}(\mathbf{V}; \mathbf{X}^t, B^t)$.

127 We summarize the proposed **GS-JOB** in Algorithm 1.

128 Although the J-orthogonality constraint typically has a sorted diagonal with $\text{diag}(\mathbf{J}) \in \{-1, +1\}^n$,
 129 **GS-JOB** is also applicable to problems with more general constraints $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$ where
 130 $\text{diag}(\mathbf{J}) \in \{\pm 1\}^n$ is unsorted.

131 **► Solving the Small-Sized Subproblem.** We now elaborate on how to find the global opti-
 132 mal solution of Problem (6). We notice that $\mathbf{V} \in \mathcal{J}_B \triangleq \{\mathbf{V} \mid \mathbf{V}^\top \mathbf{J}_{BB} \mathbf{V} = \mathbf{J}_{BB}\}$, where
 133 $\mathbf{J}_{BB} \in \left\{ \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \right\}$. We now concentrate on the first case where $\mathbf{J}_{BB} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$. The
 134 following proposition provides a strategy to decompose any J-orthogonal matrix.

Algorithm 1: GS-JOBCD: Block Coordinate Descent Methods using a Gauss-Seidel Strategy for Solving Problem (1)

Init.: Set \mathbf{X}^0 to satisfy J-orthogonality constraints (e.g., via Hyperbolic CS Decomposition).

for t from 0 to T **do**

(S1) Choose a coordinate \mathbf{B}^t with $|\mathbf{B}^t| = 2$ from the set Ω randomly and uniformly (with replacement) for the t -th iteration. Denote $\mathbf{B} = \mathbf{B}^t$.

(S2) Choose a matrix $\mathbf{Q} \in \mathbb{R}^{4 \times 4}$ using Formula (4).

(S3) Solve the following small-size subproblem globally.

$$\bar{\mathbf{V}}^t \in \arg \min_{\mathbf{V} \in \mathcal{J}_{\mathbf{B}}} \frac{1}{2} \|\mathbf{V} - \mathbf{I}\|_{\mathbf{Q} + \theta \mathbf{I}}^2 + \langle \mathbf{V} - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} \rangle + f(\mathbf{X}^t) \quad (5)$$

$$= \arg \min_{\mathbf{V} \in \mathcal{J}_{\mathbf{B}} \in \mathbb{R}^{2 \times 2}} \frac{1}{2} \|\mathbf{V}\|_{\mathbf{Q}}^2 + \langle \mathbf{V}, \mathbf{P} \rangle + c \quad (6)$$

where $\mathbf{P} \triangleq [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} - \text{mat}(\dot{\mathbf{Q}} \text{vec}(\mathbf{I}_2))$, $\dot{\mathbf{Q}} = \mathbf{Q} + \theta \mathbf{I}$ and

$c \triangleq f(\mathbf{X}^t) - \langle \mathbf{I}_2, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} \rangle + \frac{1}{2} \|\mathbf{I}\|_{\mathbf{Q}}^2$ is a constant.

(S4) $\mathbf{X}^{t+1}(\mathbf{B}, :) = \bar{\mathbf{V}}^t \mathbf{X}^t(\mathbf{B}, :)$

end

135 **Proposition 2.2.** (Hyperbolic CS Decomposition [41]) Let \mathbf{V} be J-orthogonal with signature
 136 $(p, n - p)$. Assume that $n - p \leq p$. Then there exist vectors $\hat{c}, \hat{s} \in \mathbb{R}^{n-p}$ with $\hat{c} \odot \hat{c} -$
 137 $\hat{s} \odot \hat{s} = \mathbf{1}$, and orthogonal matrices $\mathbf{U}_1, \mathbf{V}_1 \in \mathbb{R}^{p \times p}$ and $\mathbf{U}_2, \mathbf{V}_2 \in \mathbb{R}^{(n-p) \times (n-p)}$ such that:

$$138 \mathbf{V} = \begin{bmatrix} \mathbf{U}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_2 \end{bmatrix} \begin{bmatrix} \text{Diag}(\hat{c}) & \mathbf{0} & \text{Diag}(\hat{s}) \\ \mathbf{0} & I_{p-(n-p)} & \mathbf{0} \\ \text{Diag}(\hat{s}) & \mathbf{0} & \text{Diag}(\hat{c}) \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_2^\top \end{bmatrix}.$$

139 Applying Proposition 2.2 with $n = 2, p = 1$, and $\mathbf{U}_1 = \mathbf{U}_2 = \mathbf{V}_1 = \mathbf{V}_2 = \pm \mathbf{1}$, $\hat{c}^2 - \hat{s}^2 = 1$ with
 140 $\hat{c}, \hat{s} \in \mathbb{R}$, we parametrize \mathbf{V} as: $\mathbf{V} = \begin{pmatrix} \pm 1 & \mathbf{0} \\ \mathbf{0} & \pm 1 \end{pmatrix} \cdot \begin{pmatrix} \hat{c} & \hat{s} \\ \hat{s} & \hat{c} \end{pmatrix} \cdot \begin{pmatrix} \pm 1 & \mathbf{0} \\ \mathbf{0} & \pm 1 \end{pmatrix}$, where we denote \hat{s} as $\sinh(\mu)$, \hat{c} as
 141 $\cosh(\mu)$, and \tilde{t} as $\tanh(\mu)$ for some $\mu \in \mathbb{R}$, for simplicity of notation. It is not difficult to show that
 142 Problem (6) reduces to the following one-dimensional search problem:

$$\bar{\mu} \in \min_{\mu} \frac{1}{2} \text{vec}(\mathbf{V})^\top \dot{\mathbf{Q}} \text{vec}(\mathbf{V}) + \langle \mathbf{V}, \mathbf{P} \rangle, \text{ s. t. } \mathbf{V} \in \left\{ \begin{pmatrix} \hat{c} & \hat{s} \\ \hat{s} & \hat{c} \end{pmatrix}, \begin{pmatrix} \hat{c} & -\hat{s} \\ -\hat{s} & \hat{c} \end{pmatrix}, \begin{pmatrix} -\hat{c} & -\hat{s} \\ \hat{s} & \hat{c} \end{pmatrix}, \begin{pmatrix} \hat{c} & -\hat{s} \\ \hat{s} & -\hat{c} \end{pmatrix} \right\}. \quad (7)$$

143 We apply a breakpoint search method to solve Problem (7). For simplicity, we provide an analysis
 144 only for the first case. A detailed discussion of all four cases can be found in Appendix Section B.1.
 145 For the case where $\mathbf{V} = \begin{pmatrix} \hat{c} & \hat{s} \\ \hat{s} & \hat{c} \end{pmatrix}$, Problem (7) reduces to the following problem:

$$\min_{\hat{c}, \hat{s}} a \hat{c} + b \hat{s} + c \hat{c}^2 + d \hat{c} \hat{s} + e \hat{s}^2, \quad (8)$$

146 where $a = \mathbf{P}_{11} + \mathbf{P}_{22}$, $b = \mathbf{P}_{12} + \mathbf{P}_{21}$, $c = \frac{1}{2}(\dot{\mathbf{Q}}_{11} + \dot{\mathbf{Q}}_{41} + \dot{\mathbf{Q}}_{14} + \dot{\mathbf{Q}}_{44})$, $d = \frac{1}{2}(\dot{\mathbf{Q}}_{21} + \dot{\mathbf{Q}}_{31} +$
 147 $\dot{\mathbf{Q}}_{12} + \dot{\mathbf{Q}}_{42} + \dot{\mathbf{Q}}_{13} + \dot{\mathbf{Q}}_{43} + \dot{\mathbf{Q}}_{24} + \dot{\mathbf{Q}}_{34})$, and $e = \frac{1}{2}(\dot{\mathbf{Q}}_{22} + \dot{\mathbf{Q}}_{32} + \dot{\mathbf{Q}}_{23} + \dot{\mathbf{Q}}_{33})$. Then we perform
 148 a substitution to convert Problem (8) into an equivalent problem that depends on the trigonometric
 149 functions: (i) $\hat{c}^2 = \frac{1}{1-\tilde{t}^2}$; (ii) $\hat{s}^2 = \frac{\tilde{t}^2}{1-\tilde{t}^2}$; (iii) $\tilde{t} = \frac{\hat{s}}{\hat{c}}$. The following lemma provides a characterization
 150 of the global optimal solution for Problem (8).

151 **Lemma 2.3.** (Proof in Section C.2) We let $\tilde{F}(\tilde{c}, \tilde{s}) \triangleq a \tilde{c} + b \tilde{s} + c \tilde{c}^2 + d \tilde{c} \tilde{s} + e \tilde{s}^2$. The optimal so-
 152 lution $\bar{\mu}$ to Problem (8) can be computed as: $[\cosh(\bar{\mu}), \sinh(\bar{\mu})] \in \arg \min_{[c, s]} \tilde{F}(c, s)$, s. t. $[c, s] \in$
 153 $\left\{ \left[\frac{1}{\sqrt{1-(\bar{t}_+)^2}}, \frac{\bar{t}_+}{\sqrt{1-(\bar{t}_+)^2}} \right], \left[\frac{-1}{\sqrt{1-(\bar{t}_-)^2}}, \frac{-\bar{t}_-}{\sqrt{1-(\bar{t}_-)^2}} \right] \right\}$, where $\bar{t}_+ \in \arg \min_t p(t) \triangleq \frac{a+bt}{\sqrt{1-t^2}} + \frac{w+dt}{1-t^2}$,
 154 $\bar{t}_- \in \arg \min_t \tilde{p}(t) \triangleq \frac{-a-bt}{\sqrt{1-t^2}} + \frac{w+dt}{1-t^2}$. Here $w = c + e$.

155 We now describe how to find the optimal solution \bar{t}_+ , where $\bar{t}_+ \in \arg \min_t p(t) \triangleq \frac{a+bt}{\sqrt{1-t^2}} +$
 156 $\frac{w+dt}{1-t^2}$; this strategy can naturally be extended to find \bar{t}_- . Initially, we have the following
 157 first-order optimality conditions for the problem: $0 = \nabla p(t) = [b(1-t^2) + (a+bt)t]\sqrt{1-t^2} +$
 158 $[d(1-t^2) + (w+dt)(2t)] \Leftrightarrow dt^2 + 2wt + d = -[b+at]\sqrt{1-t^2}$. Squaring both sides yields the
 159 following quartic equation: $c_4 t^4 + c_3 t^3 + c_2 t^2 + c_1 t + c_0 = 0$, where $c_4 = d^2 + a^2$, $c_3 = 4wd + 2ab$,

160 $c_2 = 4w^2 + 2d^2 - a^2 + b^2$, $c_1 = 4wd - 2ab$, $c_0 = d^2 - b^2$. This equation can be solved analytically
 161 by Lodovico Ferrari's method [46], resulting in all its real roots $\{\tilde{t}_1, \tilde{t}_2, \dots, \tilde{t}_j\}$ with $1 \leq j \leq 4$.

162 For the second and third cases, Problem (6) essentially boils down to optimization under orthogonality
 163 constraints. The work of [51] derives a breakpoint search method for finding the optimal solution for
 164 Problem (6) with $\mathbf{J}_{\text{BB}} \in \left\{ \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \right\}$ using the Givens rotation and Jacobi reflection matrices.

165 2.2 Variance-Reduced Jacobi Block Coordinate Descent Algorithm

166 This subsection proposes the **VR-J-JOBCD** algorithm, a randomized block coordinate descent
 167 method derived from **GS-JOBCD**. Importantly, by leveraging the parallel framework of a Jacobi
 168 strategy [17, 9], **VR-J-JOBCD** integrates variance reduction techniques [39, 30, 18] to decrease
 169 oracle complexity in the minimization of finite-sum functions. This makes the algorithm effective for
 170 minimizing large-scale problems under J-orthogonality constraints.

171 **Notations.** We assume n is an even number in this paper. We create $(n/2)$ pairs by non-overlapping
 172 grouping of the numbers in any arbitrary combination, with each pair containing two distinct numbers
 173 from the set $[n]$. It is not hard to verify that such grouping yields $C_J = (n!)/(2^{n/2} \frac{n!}{2})$ possible
 174 combinations. The set of these combinations is denoted as $\Upsilon \triangleq \{\tilde{\mathcal{B}}_i\}_{i=1}^{C_J} \triangleq \{\tilde{\mathcal{B}}_1, \tilde{\mathcal{B}}_2, \dots, \tilde{\mathcal{B}}_{C_J}\}^1$.

175 **► Variance Reduction Strategy.** We incorporate state-of-the-art variance reduction strategies from
 176 the literature [30, 5] into our algorithm to solve Problem (1). These methods iteratively generate a
 177 stochastic gradient estimator as follows:

$$\tilde{\mathbf{G}}^t = \begin{cases} \frac{1}{b} \sum_{i \in \mathcal{S}_+^t} \nabla f_i(\mathbf{X}^t), & \text{with probability } p; \\ \tilde{\mathbf{G}}^{t-1} + \frac{1}{b'} \sum_{i \in \mathcal{S}_*^t} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})), & \text{with probability } 1 - p. \end{cases} \quad (9)$$

178 Here, $\{\mathcal{S}_+^t, \mathcal{S}_*^t\}$ are uniform random minibatch samples with $|\mathcal{S}_+^t| = b$, $|\mathcal{S}_*^t| = b'$, and $\tilde{\mathbf{G}}^0 =$
 179 $\frac{1}{b} \sum_{i \in \mathcal{S}_+^0} \nabla f_i(\mathbf{X}^0)$. We drop the superscript t for $\{\mathcal{S}_+^t, \mathcal{S}_*^t\}$ as t can be inferred from context. We
 180 only focus on the default setting that [30, 5]: $b = N$, $b' = \sqrt{b}$ and $p = \frac{b'}{b+b'}$.

181 **► Jacobi Block Coordinate Descent Method.** The proposed algorithm is built upon the parallel
 182 framework of a Jacobi strategy. In each iteration t , we randomly and uniformly (with replacement)
 183 select a coordinate set $\mathbf{B}^t \triangleq \{\mathbf{B}_{(1)}^t, \mathbf{B}_{(2)}^t, \dots, \mathbf{B}_{(\frac{n}{2})}^t\}$ from the set Υ with $\mathbf{B}^t \in \mathbb{N}^{\frac{n}{2} \times 2}$ and $\mathbf{B}_{(i)}^t \in \mathbb{N}^2$.
 184 For all t , we have: $\mathbf{B}_{(i)}^t \cap \mathbf{B}_{(j)}^t = \emptyset$ and $\cup_{i=1}^{n/2} \mathbf{B}_{(i)}^t = [n]$. We drop the superscript t if t can be
 185 inferred from context.

186 The following lemma shows how to choose a suitable matrix \mathbf{Q} so that the Jacobi strategy can be
 187 applied.

188 **Lemma 2.4.** (Proof in Section C.3) We let $\mathbf{B}^t \triangleq \{\mathbf{B}_{(1)}^t, \mathbf{B}_{(2)}^t, \dots, \mathbf{B}_{(\frac{n}{2})}^t\} \in \Upsilon$ for all t . We let $\mathbf{Q} = \varsigma \mathbf{I}_4$,
 189 where ς is some suitable constant with $\varsigma \leq L_f$. For any $\mathbf{B}_{(i)}^t$ and $\mathbf{B}_{(j)}^t$ with $i \neq j$, their corresponding
 190 objective functions as in Equation (3) are independent.

191 We consider the following block coordinate update rule in **VR-J-JOBCD**: $\mathbf{X}^{t+1} \leftarrow \tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V};) \triangleq$
 192 $\mathbf{X}^t + [\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top] \mathbf{X}^t$. The following lemma provides properties of this rule.

193 **Lemma 2.5.** (Proof in Section C.4) We let $\mathbf{B} \in \Upsilon$, $\mathbf{V}_i \in \mathcal{J}_{\mathbf{B}_{(i)}}$, $\mathbf{X} \in \mathcal{J}$, and $i \in$
 194 $[\frac{n}{2}]$. We define $\mathbf{X}^+ \triangleq \tilde{\mathcal{X}}_{\mathbf{B}}(\mathbf{V};) \triangleq \mathbf{X} + [\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top] \mathbf{X}$. We have: (a)
 195 $\sum_{i=1}^{\frac{n}{2}} \|\mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{F}}^2 = \|\sum_{i=1}^{\frac{n}{2}} \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{F}}^2$. (b) $\|\mathbf{X}^+ - \mathbf{X}\|_{\mathbb{F}}^2 \leq \|\mathbf{X}\|_{\mathbb{F}}^2 \cdot$
 196 $\sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbb{F}}^2$. (c) $\|\mathbf{X}^+ - \mathbf{X}\|_{\mathbb{H}}^2 \leq \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbf{Q}}^2$ with $\mathbf{Q} = \varsigma \mathbf{I}_4$. (d) For all $\tilde{\mathbf{G}} \in \mathbb{R}^{n \times n}$, it
 197 follows that: $2 \sum_{i=1}^{n/2} \langle \mathbf{V}_i - \mathbf{I}_2, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}) \mathbf{X}^\top]_{\mathbf{B}_{(i)} \mathbf{B}_{(i)}} \rangle \leq \|\mathbf{X}\|_{\mathbb{F}}^2 \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbb{F}}^2 + \|[\nabla f(\mathbf{X}) -$
 198 $\tilde{\mathbf{G}}]\|_{\mathbb{F}}^2$.

199 **► The Main Algorithm.** Using the update rule above, we consider the following iterative procedure:
 200 $\mathbf{X}^{t+1} \leftarrow \tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V};)$, where $\tilde{\mathbf{V}}^t \in \arg \min_{\mathbf{V}} f(\tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V};))$. We establish the majorization function for

¹Taking $n = 4$ for example, we have: $\Upsilon = \{(1, 2), (3, 4)\}, \{(1, 3), (2, 4)\}, \{(1, 4), (2, 3)\}$.

201 $f(\tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}:\cdot))$, as follows:

$$\begin{aligned} f(\tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}:\cdot)) &\stackrel{\textcircled{1}}{\leq} f(\mathbf{X}^t) + \langle \tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}:\cdot) - \mathbf{X}^t, \nabla f(\mathbf{X}^t) \rangle + \frac{1}{2} \|\tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}:\cdot) - \mathbf{X}^t\|_{\mathbf{H}}^2 \\ &\stackrel{\textcircled{2}}{\leq} f(\mathbf{X}^t) + \sum_{i=1}^{n/2} \{ \langle \mathbf{V}_i - \mathbf{I}_2, [\nabla f(\mathbf{X})(\mathbf{X}^t)^\top]_{\mathbf{B}(i)\mathbf{B}(i)} \rangle + \frac{1}{2} \|\mathbf{V}_i - \mathbf{I}_2\|_{(\theta+\zeta)\mathbf{I}}^2 \} \end{aligned} \quad (10)$$

202 where step $\textcircled{1}$ uses the results of telescoping Inequality (2) over i from 1 to N ; step $\textcircled{2}$ uses $\mathbf{X}^{t+1} -$
203 $\mathbf{X}^t = [\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}(i)} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}(i)}^\top] \mathbf{X}^t$, Claim (c) of Lemma 2.5, $\theta \geq 0$, and $\mathbf{Q} = \zeta \mathbf{I}$.

204 Instead of computing the exact Euclidean gradient $\nabla f(\mathbf{X}^t)$ as **GS-JOBCD**, **VR-J-JOBCD** maintains
205 and updates a recursive gradient estimator $\tilde{\mathbf{G}}^t$ using a variance-reduced strategy as in Formula (9).
206 We consider minimizing the following function instead of the one on the right-hand side of Inequality
207 (10):

$$\mathcal{T}(\mathbf{V}:\cdot; \mathbf{X}^t, \mathbf{B}^t) \triangleq f(\mathbf{X}^t) + \sum_{i=1}^{n/2} \langle \mathbf{V}_i - \mathbf{I}_2, [\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top]_{\mathbf{B}(i)\mathbf{B}(i)} \rangle + \frac{1}{2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbf{Q}}^2. \quad (11)$$

208 Here, $\mathcal{T}(\mathbf{V}:\cdot; \mathbf{X}^t, \mathbf{B}^t)$ can be termed as a stochastic majorization function of $f(\tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}:\cdot))$ at the current
209 solution \mathbf{X}^t . Therefore, we can consider the following optimization problem to find $\{\mathbf{V}:\cdot\}$ using:
210 $\tilde{\mathbf{V}}^t \in \arg \min_{\mathbf{V}:\cdot} \mathcal{T}(\mathbf{V}:\cdot; \mathbf{X}^t, \mathbf{B}^t)$, which can be decomposed into $(n/2)$ independent subproblems and
211 solved in parallel. It is important to note that each \mathbf{V}_i in Problem (12) is identical to Problem (6),
212 which can be efficiently solved in $\mathcal{O}(1)$ using the breakpoint search method, as in **GS-JOBCD**.

213 We summarize the proposed **VR-J-JOBCD** in Algorithm 2. Notably, when $N = 1$, **VR-J-JOBCD**
214 simplifies to a direct Jacobi strategy for solving Problem (1), which we refer to as **J-JOBCD**.

Algorithm 2: VR-J-JOBCD: Block Coordinate Descent Methods using a variance-reduced and Jacobi strategy for Solving Problem 1

Init.: Set \mathbf{X}^0 to satisfy J-orthogonality constraints (e.g., via Hyperbolic CS Decomposition).

for t from 0 to T **do**

(S1) Choose a coordinate \mathbf{B}^t from the set Υ randomly and uniformly (with replacement) for the t -th iteration. Denote $\mathbf{B} = \mathbf{B}^t$. In our implementation, we simply randomly permute the set $\{1, 2, \dots, n\}$ and then output the grouping $\{[1, 2], [3, 4], [5, 6], \dots, [n-1, n]\}$.

(S2) Use a variance-reduced strategy (9) to obtain $\tilde{\mathbf{G}}^t$.

(S3) Solve small-sized subproblems in parallel with $\mathbf{Q} = \zeta \mathbf{I} \in \mathbb{R}^{4 \times 4}$.

for $i = 1$ to $n/2$ **in parallel do**

$$\begin{aligned} \tilde{\mathbf{V}}_i^t &\in \arg \min_{\mathbf{V}_i \in \mathcal{J}_{\mathbf{B}(i)}} \frac{1}{2} \|\mathbf{V}_i - \mathbf{I}\|_{\mathbf{Q}}^2 + \langle \mathbf{V}_i - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}(i)\mathbf{B}(i)} \rangle + f(\mathbf{X}^t) \\ &= \arg \min_{\mathbf{V}_i \in \mathcal{J}_{\mathbf{B}(i)}} \frac{1}{2} \|\mathbf{V}_i\|_{\mathbf{Q}}^2 + \langle \mathbf{V}_i, \mathbf{P}_i \rangle \end{aligned} \quad (12)$$

where $\mathbf{P}_i \triangleq [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}(i)\mathbf{B}(i)} - \text{mat}(\ddot{\mathbf{Q}} \text{vec}(\mathbf{I}_2)) - \theta \mathbf{I}_2$, $\ddot{\mathbf{Q}} = (\zeta + \theta) \mathbf{I}$.

(S4) Update the solution \mathbf{X}^{t+1} in parallel as follows:

for $i = 1$ to $n/2$ **in parallel do**

$$\mathbf{X}^{t+1}(\mathbf{B}(i), :) = \tilde{\mathbf{V}}_i^t \mathbf{X}^t(\mathbf{B}(i), :)$$

end

215 3 Optimality Analysis

216 This section provides an optimality analysis for the proposed algorithms.

217 Initially, we define the first-order optimality condition for Problem (1). Since the matrix $\mathbf{X}^\top \mathbf{J} \mathbf{X}$
218 is symmetric, the Lagrangian multiplier Λ corresponding to the constraints $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$ is also a
219 symmetric matrix. The Lagrangian function of problem (1) is $\mathcal{L}(\mathbf{X}, \Lambda) = f(\mathbf{X}) - \frac{1}{2} \langle \Lambda, \mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J} \rangle$.

220 We obtain the following lemma for the first-order optimality condition for Problem (1).

221 **Lemma 3.1.** (Proof in Section D.1, First-Order Optimality Condition) We let $\mathcal{J} \triangleq \{\mathbf{X} \mid \mathbf{X}^\top \mathbf{J} \mathbf{X} =$
222 $\mathbf{J}\}$. We have (a) A solution $\tilde{\mathbf{X}} \in \mathcal{J}$ is a critical point of problem (1) if and only if: $\mathbf{0} =$

223 $\nabla_{\mathcal{J}} f(\tilde{\mathbf{X}}) \triangleq \nabla f(\tilde{\mathbf{X}}) - \mathbf{J}\tilde{\mathbf{X}}[\nabla f(\tilde{\mathbf{X}})]^{\top} \tilde{\mathbf{X}}\mathbf{J}$. The associated Lagrangian multiplier can be computed as
 224 $\Lambda = \mathbf{J}\tilde{\mathbf{X}}^{\top} \nabla f(\tilde{\mathbf{X}})$. (b) The critical point condition is equivalent to the requirement that the matrix
 225 $\mathbf{X}\nabla f(\tilde{\mathbf{X}})^{\top} \mathbf{J}$ is symmetric, which is expressed as $\mathbf{X}\mathbf{G}^{\top} \mathbf{J} = [\mathbf{X}\mathbf{G}^{\top} \mathbf{J}]^{\top}$.

226 **Remarks.** While our results in Lemma 3.1 show similarities to existing works focusing on problems
 227 under orthogonality constraints [45], this study marks the first investigation into the first-order
 228 optimality condition for optimization problems under J-orthogonality constraints.

229 The following definition is useful in our subsequent analysis of the proposed algorithms.

230 **Definition 3.2.** (Block Stationary Point, abbreviated as BS-point) Let $\theta > 0$. A solution $\tilde{\mathbf{X}} \in \mathcal{J}$
 231 is termed as a block stationary point if, for all $\mathbf{B} \in \Omega \triangleq \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{C_n^2}\}$, the following condition is
 232 satisfied: $\mathbf{I}_2 \in \arg \min_{\mathbf{V} \in \mathcal{J}_{\mathbf{B}}} \mathcal{G}(\mathbf{V}; \tilde{\mathbf{X}}, \mathbf{B})$.

233 The following theorem shows the relation between critical points and BS-points.

234 **Theorem 3.3.** (Proof in Section D.2) Any BS-point is a critical point, while the reverse is not
 235 necessarily true.

236 4 Convergence Analysis

237 This section provides a convergence analysis for **GS-JOBCD** and **VR-J-JOBCD**.

238 For **GS-JOBCD**, the randomness of output $(\bar{\mathbf{V}}^t, \mathbf{X}^{t+1})$ for all t are influenced by the random variable
 239 $\xi^t \triangleq (\mathbf{B}^1; \mathbf{B}^2; \dots; \mathbf{B}^t)$. For **VR-J-JOBCD**, the randomness of output $(\bar{\mathbf{V}}^t, \mathbf{X}^{t+1})$ are influenced by
 240 the random variables $\iota^t \triangleq (\mathbf{B}^1, \mathbf{S}_+^1, \mathbf{S}_*^1; \mathbf{B}^2, \mathbf{S}_+^2, \mathbf{S}_*^2; \dots; \mathbf{B}^t, \mathbf{S}_+^t, \mathbf{S}_*^t)$.

241 We denote $\bar{\mathbf{X}}$ as the global optimal solution of Problem (1). To simplify notations, we define:
 242 $u^t = \|\tilde{\mathbf{G}}^t - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2$, and $\Delta_i = f(\mathbf{X}^i) - f(\bar{\mathbf{X}})$.

243 We impose the following additional assumptions on the proposed algorithms.

244 **Assumption 4.1.** There exists constants $\{\bar{\mathbf{X}}, \bar{\mathbf{V}}\}$ that: $\|\mathbf{X}^t\|_{\mathbb{F}} \leq \bar{\mathbf{X}}$, and $\|\mathbf{V}^t\|_{\mathbb{F}} \leq \bar{\mathbf{V}}$ for all t .

245 **Assumption 4.2.** There exists a constant $\bar{\mathbf{G}}$ that: $\|\nabla f(\mathbf{X}^t)\|_{\mathbb{F}} \leq \bar{\mathbf{G}}$, and $\|\tilde{\mathbf{G}}^t\|_{\mathbb{F}} \leq \bar{\mathbf{G}}$ for all t .

246 **Assumption 4.3.** For any $\mathbf{X} \in \mathbb{R}^{n \times n}$, $\mathbb{E}_i[\|\nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] \leq \sigma^2$, where i is drawn uniformly
 247 at random from $[N]$.

248 **Remarks.** (i) Assumption 4.1 is satisfied as the function $f_i(\mathbf{X})$ is coercive for all i . (ii) Assumption
 249 4.2 imposes a bound on the (stochastic) gradient, a fairly moderate condition frequently employed in
 250 nonconvex optimization [26]. (iii) Assumption 4.3 ensures that the variance of the stochastic gradient
 251 is bounded, which is a common requirement in stochastic optimization [30, 5].

252 4.1 Global Convergence

253 We define the ϵ -BS-point as follows.

254 **Definition 4.4.** (ϵ -BS-point) Given any constant $\epsilon > 0$, a point $\tilde{\mathbf{X}}$ is called an ϵ -BS-point if: $\mathcal{E}(\tilde{\mathbf{X}}) \leq \epsilon$.

255 Here, $\mathcal{E}(\mathbf{X})$ is defined as $\mathcal{E}(\mathbf{X}) \triangleq \frac{1}{C_n^2} \sum_{i=1}^{C_n^2} \text{dist}(\mathbf{I}_2, \arg \min_{\mathbf{V}} \mathcal{G}(\mathbf{V}; \mathbf{X}, \mathcal{B}_i))^2$ for **GS-JOBCD**,
 256 while it is defined as $\mathcal{E}(\mathbf{X}) \triangleq \frac{1}{C_J} \sum_{i=1}^{C_J} \mathbb{E}_{\iota^t}[\text{dist}(\mathbf{I}_2, \arg \min_{\mathbf{V}} \mathcal{T}(\mathbf{V}; \mathbf{X}, \tilde{\mathcal{B}}_i))^2]$ for **VR-J-JOBCD**,
 257 where the expectation is with respect to the randomness inherent in the algorithm [30].

258 We have the following useful lemma for **VR-J-JOBCD**.

259 **Lemma 4.5.** (Proof in Section E.1) Suppose Assumption 4.3 holds, then the variance $\mathbb{E}_{\iota^t}[u_k]$ of the
 260 gradient estimators $\{\tilde{\mathbf{G}}^t\}$ of Algorithm 2 is bounded by: $\mathbb{E}_{\iota^t}[u^t] \leq \frac{p(N-b)}{b(N-1)} \sigma^2 + (1-p) \mathbb{E}_{\iota^{t-1}}[u^{t-1}] +$
 261 $\frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{b'} \mathbb{E}_{\iota^{t-1}}[\sum_{i=1}^{n/2} \|\mathbf{V}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2]$

262 The following two theorems establish the iteration complexity (or oracle complexity) for **GS-JOBCD**
 263 and **VR-J-JOBCD**.

264 **Theorem 4.6.** (Proof in Section E.2) **GS-JOBCD** finds an ϵ -BS-point of Problem (1) within $\mathcal{O}(\frac{\Delta_0 N}{\epsilon})$
 265 arithmetic operations.

266 **Theorem 4.7.** (Proof in Section E.3) Let $b = N$, $b' = \sqrt{N}$, and $p = \frac{b'}{b+b'}$. **VR-J-JOBCD** finds an
 267 ϵ -BS-point of Problem (1) within $\mathcal{O}(nN + \frac{\Delta_0\sqrt{N}}{\epsilon})$ arithmetic operations.

268 **Remark.** Theorems 4.6 and 4.7 demonstrate that the arithmetic operation complexity of **GS-JOBCD**
 269 is linearly dependent on N , while **VR-J-JOBCD** is linearly dependent on \sqrt{N} . Therefore, **VR-J-**
 270 **JOBCD** reduces the iteration complexity significantly.

271 4.2 Strong Convergence under KL Assumption

272 We prove algorithms achieve strong convergence based on a non-convex analysis tool called Kurdyka-
 273 Łojasiewicz inequality[2].

274 We impose the following assumption on Problem (1).

275 **Assumption 4.8.** (Kurdyka-Łojasiewicz Property). Assume that $f^\circ(\mathbf{X}) = f(\mathbf{X}) + \mathcal{I}_{\mathcal{J}}(\mathbf{X})$ is a KL
 276 function. For all $\mathbf{X} \in \text{dom } f^\circ$, there exists $\sigma \in [0, 1)$, $\eta \in (0, +\infty]$ a neighborhood Υ of \mathbf{X} and a
 277 concave and continuous function $\varphi(t) = ct^{1-\sigma}$, $c > 0$, $t \in [0, \eta)$ such that for all $\mathbf{X}' \in \Upsilon$ and satisfies
 278 $f^\circ(\mathbf{X}') \in (f^\circ(\mathbf{X}), f^\circ(\mathbf{X}) + \eta)$, the following holds: $\text{dist}(\mathbf{0}, \nabla f^\circ(\mathbf{X}'))\varphi'(f^\circ(\mathbf{X}') - f^\circ(\mathbf{X})) \geq 1$.

279 We establish strong limit-point convergence for **VR-J-JOBCD** and **GS-JOBCD**.

280 **Theorem 4.9.** (Proof in Section E.5, a Finite Length Property). The sequence $\{\mathbf{X}^t\}_{t=0}^\infty$ of **GS-**
 281 **JOBCD** has finite length property that: $\forall t, \sum_{i=1}^t \mathbb{E}_{\xi^t}[\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F] \leq \mathcal{O}(\varphi(\Delta_1)) < +\infty$, where
 282 $\varphi(\cdot)$ is the desingularization function defined in Proposition 4.8.

283 **Theorem 4.10.** (Proof in Section E.4, a Finite Length Property). Choosing $b = N$, $b' = \sqrt{N}$
 284 and $p = \frac{b'}{b+b'}$, then the sequence $\{\mathbf{X}^t\}_{t=0}^\infty$ of **VR-J-JOBCD** has finite length property that:
 285 $\forall t, \sum_{i=1}^t \mathbb{E}_{\xi^t}[\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F] \leq \mathcal{O}(\frac{\varphi(\Delta_1)}{N^{1/4}}) < +\infty$, where $\varphi(\cdot)$ is the desingularization function
 286 defined in Assumption 4.8.

287 5 Applications and Numerical Experiments

288 This section demonstrates the effectiveness and efficiency of **JOBCD** on three optimization tasks:
 289 (i) the hyperbolic eigenvalue problem, (ii) structural probe problem, and (iii) Ultra-hyperbolic
 290 Knowledge Graph Embedding problem. We provide experiments for the last problem in Section F.2.

291 **► Application to the Hyperbolic Eigenvalue Problem (HEVP).** The hyperbolic eigenvalue problem
 292 refers to the generalized eigenvalue problem in hyperbolic spaces [40]. This problem is a fundamental
 293 component in machine learning models, such as Hyperbolic PCA [43, 6]. Given a data matrix
 294 $\mathbf{D} \in \mathbb{R}^{m \times n}$ and a signature matrix \mathbf{J} with signature $(p, n - p)$, HEVP can be formulated as the
 295 following optimization problem: $\min_{\mathbf{X}} -\text{tr}(\mathbf{X}^\top \mathbf{D}^\top \mathbf{D} \mathbf{X})$, s. t. $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$.

296 **► Application to the Hyperbolic Structural Probe Problem (HSPP).** The Structure Probe (SP) is
 297 a metric learning model aimed at understanding the intrinsic semantic information of large language
 298 models [20] [7]. Given a data matrix $\mathbf{D} \in \mathbb{R}^{m \times n}$ and its associated Euclidean distance metric
 299 matrix $\mathbf{T} \in \mathbb{R}^{m \times m}$, HSPP employs a smooth homeomorphic mapping function $\varphi(\cdot)$ to project
 300 the data \mathbf{D} into ultra-hyperbolic space. Subsequently, it seeks an appropriate linear transformation
 301 $\mathbf{X} \in \mathbb{R}^{n \times n}$ constrained within a specific structure $\mathbf{X} \in \mathcal{J}$, such that the resulting transformed
 302 data $\mathbf{Q} \triangleq \varphi(\mathbf{D})\mathbf{X} \in \mathbb{R}^{m \times n}$ exhibits similarity to the original distance metric matrix \mathbf{T} under the
 303 ultra-hyperbolic geodesic distance $d_\alpha(\mathbf{Q}_{i:}, \mathbf{Q}_{j:})$, expressed as $\mathbf{T}_{i,j} \approx d_\alpha(\mathbf{Q}_{i:}, \mathbf{Q}_{j:})$ for all $i, j \in [m]$,
 304 where $\mathbf{Q}_{i:}$ is i -th row of the matrix $\mathbf{Q} \in \mathbb{R}^{m \times n}$. This can be formulated as the following optimization
 305 problem: $\min_{\mathbf{X}} \frac{1}{m^2} \sum_{i,j \in [m]} (\mathbf{T}_{i,j} - d_\alpha(\mathbf{Q}_{i:}, \mathbf{Q}_{j:}))^2$, s. t. $\mathbf{Q} \triangleq \varphi(\mathbf{D})\mathbf{X}$, $\mathbf{X} \in \mathcal{J}$. For more details
 306 on the functions $\varphi(\cdot)$ and $d_\alpha(\cdot, \cdot)$, please refer to Appendix Section F.1.

307 **► Datasets.** To generate the matrix $\mathbf{D} \in \mathbb{R}^{m \times n}$, we use 8 real-world or synthetic data sets for both
 308 HEVP and HSPP tasks: ‘Cifar’, ‘CnnCaltech’, ‘Gisette’, ‘Mnist’, ‘randn’, ‘Sector’, ‘TDT2’, ‘w1a’.
 309 We randomly extract a subset from the original data sets for the experiments.

310 **► Compared Methods.** We compare **GS-JOBCD** and **VR-J-JOBCD** with 3 state-of-the-art
 311 optimization algorithms under J-orthogonality constraints. (i) The CS Decomposition Method
 312 (CSDM) [48]. (ii) Stardard ADMM (ADMM) [19]. UMCM: Unconstrained Multiplier Correction
 313 Method [31, 13].

Table 1: Comparisons of the objectives for HEVP across all the compared methods. The time limit is set to 90s. The notation ‘(+)’ indicates that **GS-JOB**CD significantly improves upon the initial solution provided by **CSDM**. The 1st, 2nd, and 3rd best results are colored with red, green and blue, respectively. The value in (·) stands for $\sum_{ij}^n |\mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J}|_{ij}$.

dataname(m-n-p)	UMCM	ADMM	CSDM	GS-JOB	J-JOB	CSDM+GS-JOB
cifar(1000-100-50)	-1.05e+04(3.0e-09)	-1.05e+04(3.0e-09)	-5.28e+04(5.4e-09)	-1.03e+05(2.6e-08)	-1.11e+05(1.4e-07)	-1.24e+05(2.6e-08)(+)
CnnCal(2000-1000-500)	-5.89e+02(2.9e-08)	-5.89e+02(3.1e-10)	-1.11e+03(5.2e-10)	-1.07e+03(1.3e-09)	-9.16e+03(6.9e-08)	-1.15e+03(6.9e-10)(+)
gisette(3000-1000-500)	-3.22e+06(3.1e-10)	-3.22e+06(3.1e-10)	-8.53e+06(4.9e-10)	-9.49e+06(1.2e-09)	-1.36e+07(2.6e-08)	-9.65e+06(7.9e-10)(+)
mnist(1000-780-390)	-8.65e+04(4.1e-10)	-8.65e+04(4.1e-10)	-2.56e+05(5.6e-10)	-3.14e+05(1.2e-09)	-1.20e+06(4.1e-08)	-3.06e+05(7.6e-10)(+)
randn(10-10-5)	1.29e+02(9.7e-02)	1.29e+02(9.7e-02)	2.45e+02(2.3e-01)	-3.96e+01(9.7e-02)	-3.97e+02(9.7e-02)	1.55e+01(2.3e-01)(+)
randn(100-100-50)	-1.03e+04(3.0e-09)	-1.03e+04(2.5e-07)	-1.98e+04(4.4e-09)	-2.28e+04(5.6e-08)	-4.37e+04(2.6e-07)	-2.41e+04(4.2e-08)(+)
randn(1000-1000-500)	-1.16e+06(3.1e-10)	-1.16e+06(3.1e-10)	-1.93e+06(5.0e-10)	-1.22e+06(6.9e-10)	-1.04e+07(2.3e-07)	-1.95e+06(6.7e-10)(+)
sector(500-1000-500)	-3.61e+03(3.1e-10)	-3.61e+03(3.1e-10)	-7.90e+03(4.9e-10)	-9.24e+03(1.3e-09)	-1.06e+04(2.0e-08)	-8.51e+03(6.4e-10)(+)
TDT2(1000-1000-500)	-4.25e+06(3.1e-10)	-4.25e+06(3.1e-10)	-9.39e+06(4.8e-10)	-1.05e+07(1.1e-09)	-1.42e+07(2.1e-08)	-1.04e+07(6.5e-10)(+)
w1a(2470-290-145)	-3.02e+04(1.1e-04)	-3.02e+04(1.1e-04)	-5.72e+04(2.7e-05)	-9.21e+04(1.1e-04)	-9.32e+06(1.1e-04)	-7.94e+04(2.7e-05)(+)

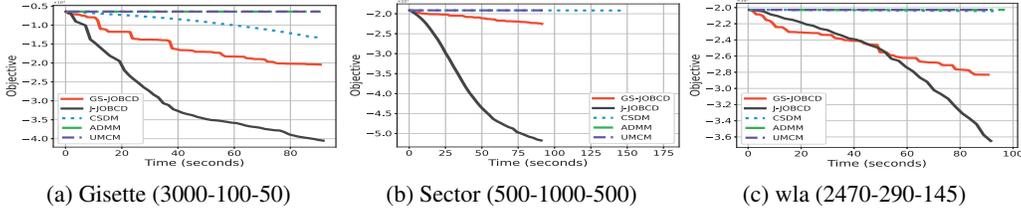


Figure 1: The convergence curve for the HEVP across various datasets with different parameters (m, n, p).

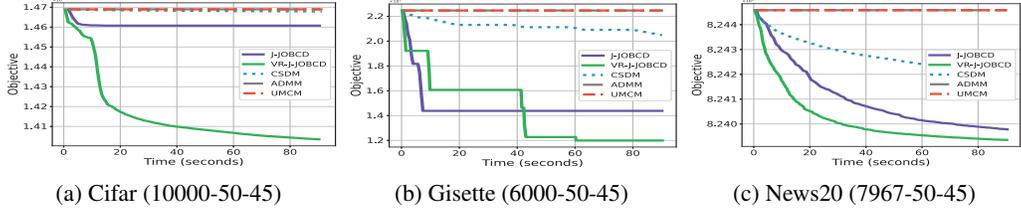


Figure 2: The convergence curve for HEVP across various datasets with different parameters (m, n, p).

314 **► Experiment Settings.** All methods are implemented using Pytorch on an Intel 2.6 GHz processor
315 with an A40 (48GB). For **HSPP**, we fix α to 1. Each method employs the same random J-orthogonal
316 matrix. The built-in solver *Admm* is used to solve the unconstrained minimization problem in **CSDM**.
317 We provide our code in the supplemental material.

318 **► Experiment Results.** Table 1 and Figure 1 display the accuracy and computational efficiency for
319 HEVP, while Figure 2 presents the results for **HSPP**, leading to the following observations: (i)
320 **GS-JOB**CD and **JJOB**CD consistently deliver better performance than the other methods. (ii) Other
321 methods frequently encounter poor local minima, whereas **GS-JOB**CD effectively escapes these
322 minima and typically achieves lower objective values, aligning with our theory that our methods
323 locate stronger stationary points. (iii) **VR-J-JOB**CD outperforms both **J-JOB**CD and **CSDM** when
324 dealing with a large dataset characterized by an infinite-sum structure.

325 6 Conclusions

326 In this paper, we propose a new approach **JOB**CD, which is based on block coordinate descent, for
327 solving the optimization problem under J-orthogonality constraints. We discuss two specific variants
328 of **JOB**CD: one based on a Gauss-Seidel strategy (**GS-JOB**CD), the other on a variance-reduced
329 Jacobi strategy. Both algorithms capitalize on specific structural characteristics of the constraints to
330 converge to more favorable stationary solutions. Notably, **VR-J-JOB**CD incorporates a variance-
331 reduction technique into a parallel framework to reduce oracle complexity in the minimization of
332 finite-sum functions. For both **GS-JOB**CD and **VR-J-JOB**CD, we establish the oracle complexity
333 under mild conditions and strong limit-point convergence results under the Kurdyka-Lojasiewicz
334 inequality. Some experiments on the hyperbolic eigenvalue problem and structural probe problem
335 show the efficiency and efficacy of the proposed methods.

References

- 336
- 337 [1] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix*
338 *manifolds*. Princeton University Press, 2008.
- 339 [2] Hédý Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal alternating
340 minimization and projection methods for nonconvex problems: An approach based on the
341 kurdyka-Łojasiewicz inequality. *Mathematics of operations research*, 35(2):438–457, 2010.
- 342 [3] Adam Bojanczyk, Nicholas J Higham, and Harikrishna Patel. Solving the indefinite least squares
343 problem by hyperbolic qr factorization. *SIAM Journal on Matrix Analysis and Applications*,
344 24(4):914–931, 2003.
- 345 [4] HanQin Cai, Yuchen Lou, Daniel McKenzie, and Wotao Yin. A zeroth-order block coordinate
346 descent algorithm for huge-scale black-box optimization. In *International Conference on*
347 *Machine Learning*, pages 1193–1203. PMLR, 2021.
- 348 [5] Xufeng Cai, Chaobing Song, Stephen Wright, and Jelena Diakonikolas. Cyclic block coordinate
349 descent with variance reduction for composite nonconvex optimization. In *International*
350 *Conference on Machine Learning*, pages 3469–3494. PMLR, 2023.
- 351 [6] Ines Chami, Albert Gu, Dat P Nguyen, and Christopher Re. Horopca: Hyperbolic dimensionality
352 reduction via horospherical projections. In *International Conference on Machine Learning*
353 *(ICML)*, volume 139, pages 1419–1429, 2021.
- 354 [7] Boli Chen, Yao Fu, Guangwei Xu, Pengjun Xie, Chuanqi Tan, Mosha Chen, and Liping Jing.
355 Probing bert in hyperbolic spaces. *ICLR*, 2021.
- 356 [8] Weize Chen, Xu Han, Yankai Lin, Hexu Zhao, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie
357 Zhou. Fully hyperbolic neural networks. *arXiv preprint arXiv:2105.14686*, 2021.
- 358 [9] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex
359 sgd. *Advances in neural information processing systems*, 32, 2019.
- 360 [10] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradi-
361 ent method with support for non-strongly convex composite objectives. *Advances in neural*
362 *information processing systems*, 27, 2014.
- 363 [11] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-
364 convex optimization via stochastic path-integrated differential estimator. *Advances in neural*
365 *information processing systems*, 31, 2018.
- 366 [12] Hamza Fawzi and Harry Goulbourne. Faster proximal algorithms for matrix optimization
367 using jacobi-based eigenvalue methods. *Advances in Neural Information Processing Systems*,
368 34:11397–11408, 2021.
- 369 [13] Bin Gao, Xin Liu, Xiaojun Chen, and Ya-xiang Yuan. A new first-order algorithmic framework
370 for optimization problems with orthogonality constraints. *SIAM Journal on Optimization*,
371 28(1):302–332, 2018.
- 372 [14] Bin Gao, Xin Liu, and Ya-xiang Yuan. Parallelizable algorithms for optimization problems with
373 orthogonality constraints. *SIAM Journal on Scientific Computing*, 41(3):A1949–A1983, 2019.
- 374 [15] Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation
375 methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-
376 2):267–305, 2016.
- 377 [16] Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- 378 [17] Eldon R Hansen. On cyclic jacobi methods. *Journal of the Society for Industrial and Applied*
379 *Mathematics*, 11(2):448–459, 1963.
- 380 [18] Vjeran Hari and Erna Begović Kovač. On the convergence of complex jacobi methods. *Linear*
381 *and multilinear algebra*, 69(3):489–514, 2021.

- 382 [19] Bingsheng He and Xiaoming Yuan. On the $\mathcal{O}(1/n)$ convergence rate of the douglas-rachford
383 alternating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.
- 384 [20] John Hewitt and Christopher D. Manning. A structural probe for finding syntax in word
385 representations. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of*
386 *the 2019 Conference of the North American Chapter of the Association for Computational*
387 *Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4129–4138, 2019.
- 388 [21] Nicholas J Higham. J-orthogonal matrices: Properties and generation. *SIAM review*, 45(3):504–
389 519, 2003.
- 390 [22] Minhui Huang, Shiqian Ma, and Lifeng Lai. A riemannian block coordinate descent method for
391 computing the projection robust wasserstein distance. In *International Conference on Machine*
392 *Learning*, pages 4446–4455. PMLR, 2021.
- 393 [23] Bo Hui and Wei-Shinn Ku. Low-rank nonnegative tensor decomposition in hyperbolic space. In
394 *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*,
395 pages 646–654, 2022.
- 396 [24] Sashank J Reddi, Suvrit Sra, Barnabas Poczos, and Alexander J Smola. Proximal stochastic
397 methods for nonsmooth nonconvex finite-sum optimization. *Advances in neural information*
398 *processing systems*, 29, 2016.
- 399 [25] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance
400 reduction. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors,
401 *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- 402 [26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua
403 Bengio and Yann LeCun, editors, *International Conference on Learning Representations (ICLR)*,
404 2015.
- 405 [27] Marc Law. Ultrahyperbolic neural networks. *Advances in Neural Information Processing*
406 *Systems*, 34:22058–22069, 2021.
- 407 [28] Marc Law and Jos Stam. Ultrahyperbolic representation learning. *Advances in neural informa-*
408 *tion processing systems*, 33:1668–1678, 2020.
- 409 [29] Qunwei Li, Yi Zhou, Yingbin Liang, and Pramod K Varshney. Convergence analysis of proximal
410 gradient with momentum for nonconvex optimization. In *International Conference on Machine*
411 *Learning*, pages 2111–2119. PMLR, 2017.
- 412 [30] Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. Page: A simple and optimal
413 probabilistic gradient estimator for nonconvex optimization. In *International conference on*
414 *machine learning*, pages 6286–6295. PMLR, 2021.
- 415 [31] Wei Liu, Yinyu Zhang, Hongqiao Yang, and Shuzhong Zhang. A class of smooth exact penalty
416 function methods for optimization problems with orthogonality constraints. *Optimization*,
417 69(3):399–426, 2020.
- 418 [32] Wan-Duo Kurt Ma, JP Lewis, and W Bastiaan Kleijn. The hsic bottleneck: Deep learning
419 without back-propagation. In *Proceedings of the AAAI conference on artificial intelligence*,
420 volume 34, pages 5085–5092, 2020.
- 421 [33] Julien Mairal. Optimization with first-order surrogate functions. In *International Conference*
422 *on Machine Learning (ICML)*, volume 28, pages 783–791, 2013.
- 423 [34] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for
424 machine learning problems using stochastic recursive gradient. In *International conference on*
425 *machine learning*, pages 2613–2621. PMLR, 2017.
- 426 [35] Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model
427 of hyperbolic geometry. In *International Conference on Machine Learning*, pages 3779–3788.
428 PMLR, 2018.

- 429 [36] Vedran Novaković and Sanja Singer. A kogbetliantz-type algorithm for the hyperbolic svd.
430 *Numerical algorithms*, 90(2):523–561, 2022.
- 431 [37] Julie Nutini, Issam Laradji, and Mark Schmidt. Let’s make block coordinate descent converge
432 faster: faster greedy rules, message-passing, active-set complexity, and superlinear convergence.
433 *Journal of Machine Learning Research*, 23(131):1–74, 2022.
- 434 [38] Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Luo. A unified convergence analysis of block
435 successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization*,
436 23(2):1126–1153, 2013.
- 437 [39] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic
438 average gradient. *Mathematical Programming*, 162:83–112, 2017.
- 439 [40] Ivan Slapnicar and Ninoslav Truhar. Relative perturbation theory for hyperbolic eigenvalue
440 problem. *Linear Algebra and its Applications*, 309(1):57–72, 2000.
- 441 [41] Michael Stewart and Paul Van Dooren. On the factorization of hyperbolic and unitary trans-
442 formations into rotations. *SIAM Journal on Matrix Analysis and Applications*, 27(3):876–890,
443 2005.
- 444 [42] Puoya Tabaghi and Ivan Dokmanić. Hyperbolic distance matrices. In *Proceedings of the*
445 *26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages
446 1728–1738, 2020.
- 447 [43] Puoya Tabaghi, Michael Khazadeh, Yusu Wang, and Sivash Mirarab. Principal component
448 analysis in space forms. *ArXiv*, abs/2301.02750, 2023.
- 449 [44] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. Spiderboost and momentum:
450 Faster variance reduction algorithms. *Advances in Neural Information Processing Systems*, 32,
451 2019.
- 452 [45] Zaiwen Wen and Wotao Yin. A feasible method for optimization with orthogonality constraints.
453 *Mathematical Programming*, 142:397 – 434, 2012.
- 454 [46] WikiContributors. Quartic equation. [https://en.wikipedia.org/wiki/Quartic_](https://en.wikipedia.org/wiki/Quartic_equation)
455 [equation](https://en.wikipedia.org/wiki/Quartic_equation).
- 456 [47] Ruiyuan Wu, Anna Scaglione, Hoi-To Wai, Nurullah Karakoc, Kari Hreinsson, and Wing-Kin
457 Ma. Federated block coordinate descent scheme for learning global and personalized models. In
458 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10355–10362,
459 2021.
- 460 [48] Bo Xiong, Shichao Zhu, Mojtaba Nayyeri, Chengjin Xu, Shirui Pan, Chuan Zhou, and Steffen
461 Staab. Ultrahyperbolic knowledge graph embeddings. In *Proceedings of the 28th ACM SIGKDD*
462 *Conference on Knowledge Discovery and Data Mining*, pages 2130–2139, 2022.
- 463 [49] Bo Xiong, Shichao Zhu, Nico Potyka, Shirui Pan, Chuan Zhou, and Steffen Staab. Semi-
464 riemannian graph convolutional networks. *ArXiv*, abs/2106.03134, 2021.
- 465 [50] Tao Yu and Christopher M De Sa. Numerically accurate hyperbolic embeddings using tiling-
466 based models. *Advances in Neural Information Processing Systems*, 32, 2019.
- 467 [51] Ganzhao Yuan. A block coordinate descent method for nonsmooth composite optimization
468 under orthogonality constraints. *ArXiv*, abs/2304.03641, 2023.
- 469 [52] Ganzhao Yuan. Coordinate descent methods for fractional minimization. In *International*
470 *Conference on Machine Learning*, pages 40488–40518, 2023.
- 471 [53] Jinshan Zeng, Tim Tsz-Kit Lau, Shaobo Lin, and Yuan Yao. Global convergence of block
472 coordinate descent in deep learning. In *International conference on machine learning*, pages
473 7313–7323. PMLR, 2019.
- 474 [54] Yiding Zhang, Xiao Wang, Chuan Shi, Nian Liu, and Guojie Song. Lorentzian graph convolu-
475 tional networks. In *Proceedings of the Web Conference 2021*, pages 1249–1261, 2021.
- 476 [55] Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic nested variance reduction for nonconvex
477 optimization. *The Journal of Machine Learning Research*, 21(1):4130–4192, 2020.

Appendix

The appendix is organized as follows.

Appendix A introduces some notations, technical preliminaries, and relevant lemmas.

Appendix B concludes some additional discussions.

Appendix C presents the proofs for Section 2.

Appendix D offers the proofs for Section 3.

Appendix E contains the proofs for Section 4.

Appendix F contains several extra experiments, extensions and discussions of the proposed methods.

A Notations, Technical Preliminaries, and Relevant Lemmas

A.1 Notations

In this paper, we denote the Lowercase boldface letters represent vectors, while uppercase letters represent real-valued matrices. We use the Matlab colon notation to denote indices that describe submatrices. The following notations are used throughout this paper.

- \mathbb{N} : Set of natural numbers
- \mathbb{R} : Set of real numbers
- $[n]$: $\{1, 2, \dots, n\}$
- $\|\mathbf{x}\|$: Euclidean norm: $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$
- \mathbf{x}_i : the i -th element of vector \mathbf{x}
- $\mathbf{X}_{i,j}$ or \mathbf{X}_{ij} : the $(i^{\text{th}}, j^{\text{th}})$ element of matrix \mathbf{X}
- $\text{vec}(\mathbf{X})$: $\text{vec}(\mathbf{X}) \in \mathbb{R}^{nn \times 1}$, the vector formed by stacking the column vectors of \mathbf{X}
- $\text{mat}(\mathbf{x}) \in \mathbb{R}^{n \times n}$, Convert $\mathbf{x} \in \mathbb{R}^{nn \times 1}$ into a matrix with $\text{mat}(\text{vec}(\mathbf{X})) = \mathbf{X}$
- \mathbf{X}^T : the transpose of the matrix \mathbf{X}
- $\text{sign}(t)$: the signum function, $\text{sign}(t) = 1$ if $t \geq 0$ and $\text{sign}(t) = -1$ otherwise
- $\mathbf{X} \otimes \mathbf{Y}$: Kronecker product of \mathbf{X} and \mathbf{Y}
- $\det(\mathbf{D})$: Determinant of a square matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$
- \mathbf{C}_n^2 : the number of possible combinations choosing k items from n without repetition.
- $\mathbf{0}_{n,r}$: A zero matrix of size $n \times r$; the subscript is omitted sometimes
- \mathbf{I}_r : $\mathbf{I}_r \in \mathbb{R}^{r \times r}$, Identity matrix
- $\mathbf{X} \succeq \mathbf{0}$ (or $\succ \mathbf{0}$) : the Matrix \mathbf{X} is symmetric positive semidefinite (or definite)
- $\text{Diag}(\mathbf{x})$: Diagonal matrix with \mathbf{x} as the main diagonal entries.
- $\text{tr}(\mathbf{A})$: Sum of the elements on the main diagonal \mathbf{A} : $\text{tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}$
- $\|\mathbf{X}\|_*$: Nuclear norm: sum of the singular values of matrix \mathbf{X}
- $\|\mathbf{X}\|$: Operator/Spectral norm: the largest singular value of \mathbf{X}
- $\|\mathbf{X}\|_F$: Frobenius norm: $(\sum_{ij} \mathbf{X}_{ij}^2)^{1/2}$
- $\nabla f(\mathbf{X})$: classical (limiting) Euclidean gradient of $f(\mathbf{X})$ at \mathbf{X}
- $\nabla_{\mathcal{J}} f(\mathbf{X})$: Riemannian gradient of $f(\mathbf{X})$ at \mathbf{X}
- $\mathcal{I}_{\xi}(\mathbf{X})$: the indicator function of a set ξ with $\mathcal{I}_{\xi}(\mathbf{X}) = 0$ if $\mathbf{X} \in \xi$ and otherwise $+\infty$
- $\text{dist}(\xi, \xi')$: the distance between two sets with $\text{dist}(\xi, \xi') \triangleq \inf_{\mathbf{X} \in \xi, \mathbf{X}' \in \xi'} \|\mathbf{X} - \mathbf{X}'\|_F$
- $\mathcal{I}_{\xi}(\mathbf{x})$: the indicator function of a set ξ with $\mathcal{I}_{\xi}(\mathbf{x}) = 0$ if $\mathbf{x} \in \xi$ and otherwise $+\infty$.

517 **A.2 Relevant Lemmas**

518 **Lemma A.1.** (Lemma 6.6 of [51]) For any $\mathbf{W} \in \mathbb{R}^{n \times n}$, we have: $\sum_{i=1}^{C_n^k} \|\mathbf{W}(\mathcal{B}_i, \mathcal{B}_i)\|_{\mathbb{F}}^2 =$
 519 $\frac{k}{n} C_n^k \sum_i \mathbf{W}_{ii}^2 + C_{n-2}^{k-2} \sum_i \sum_{j,j \neq i} \mathbf{W}_{ij}^2$. Here, the set $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{C_n^k}\}$ represents all possible
 520 combinations of the index vectors choosing k items from n without repetition.

Lemma A.2. We have \mathcal{S}_+ be the set of $|\mathcal{S}_+| = b$ samples from $[N]$, drawn with replacement and uniformly at random. Then, $\forall t, \mathbf{X}^t \in \mathbb{R}^{n \times n}$, we have:

$$\mathbb{E}_{t^t} [\|\frac{1}{b} \sum_{i \in \mathcal{S}_+} \nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] = \frac{N-b}{b(N-1)} \mathbb{E}_{t^t} [\|\nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2].$$

521 *Proof.* The proof is exactly the same as in Lemma 2.8 of [5]. □

522 **Lemma A.3.** The tangent space $\mathbf{T}_{\mathbf{X}} \mathcal{J}$ of manifold constructed by $\mathbf{X}^{\top} \mathbf{J} \mathbf{X} = \mathbf{J}$, with $\mathbf{X} \in \mathbb{R}^{n \times n}$, is :

$$\mathbf{T}_{\mathbf{X}} \mathcal{J} \triangleq \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} = 0\}, \quad (13)$$

523 where $\mathbf{Y} = t\tilde{\mathbf{Y}}$ with t is a positive scalar approaching 0.

524 *Proof.* Assuming point $\mathbf{X} \in \mathbb{R}^{n \times n}$ lies on manifold \mathcal{J} , we have: $h(\mathbf{X}) = \mathbf{X}^{\top} \mathbf{J} \mathbf{X} - \mathbf{J}$. Moving
 525 along $\mathbf{Y} \in \mathbb{R}^{n \times n}$ in the tangent space of \mathbf{X} , we obtain:

$$\begin{aligned} h(\mathbf{X} + \mathbf{Y}) &= (\mathbf{X} + \mathbf{Y})^{\top} \mathbf{J} (\mathbf{X} + \mathbf{Y}) - \mathbf{J} \\ &= \mathbf{X}^{\top} \mathbf{J} \mathbf{X} + \mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{Y} - \mathbf{J} \\ &\stackrel{\textcircled{1}}{=} \mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{Y} \\ &\stackrel{\textcircled{2}}{=} t \mathbf{X}^{\top} \mathbf{J} \tilde{\mathbf{Y}} + t \tilde{\mathbf{Y}}^{\top} \mathbf{J} \mathbf{X} + t^2 \tilde{\mathbf{Y}}^{\top} \mathbf{J} \tilde{\mathbf{Y}} \end{aligned}$$

526 where step ① uses $\mathbf{X}^{\top} \mathbf{J} \mathbf{X} = \mathbf{J}$; step ② uses $\mathbf{Y} = t\tilde{\mathbf{Y}}$.

527 Since t is a positive scalar approaching 0, we can ignore the higher-order term: $t^2 \tilde{\mathbf{Y}}^{\top} \mathbf{J} \tilde{\mathbf{Y}}$. Ac-
 528 cording to the properties of the tangent space of any manifold, we have: $h(\mathbf{X} + \mathbf{Y}) = 0$. In
 529 other words, $\mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} = 0$, i.e. we obtain the defining equation for the tangent space:
 530 $\mathbf{T}_{\mathbf{X}} \mathcal{J} \triangleq \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} = 0\}$. □

531 **B Additional Discussions**

532 **B.1 On the Global Optimal Solution for Problem (7)**

533 In Section 2.1, we have demonstrated how to use the breakpoint search method to obtain an optimal
 534 solution for the case of $\mathbf{V} = (\frac{\tilde{c}}{\tilde{s}} \frac{\tilde{s}}{\tilde{c}})$ of Problem (7). Since the structure of the other three cases
 535 $\mathbf{V} \in \{(\frac{\tilde{c}}{-\tilde{s}} \frac{-\tilde{s}}{\tilde{c}}), (\frac{-\tilde{c}}{\tilde{s}} \frac{-\tilde{s}}{\tilde{c}}), (\frac{\tilde{c}}{\tilde{s}} \frac{-\tilde{s}}{-\tilde{c}})\}$ is exactly the same except for the coefficients of Problem (8), we
 536 will provide the corresponding coefficients in Problem (8): $\min_{\tilde{c}, \tilde{s}} a \tilde{c} + b \tilde{s} + c \tilde{c}^2 + d \tilde{c} \tilde{s} + e \tilde{s}^2$, and
 537 omit the specific analysis process.

538 **Case (a).** $\mathbf{V} = (\frac{\tilde{c}}{-\tilde{s}} \frac{-\tilde{s}}{\tilde{c}})$: $a = \mathbf{P}_{11} + \mathbf{P}_{22}$, $b = -\mathbf{P}_{12} - \mathbf{P}_{21}$, $c = \frac{1}{2}(\dot{\mathbf{Q}}_{11} + \dot{\mathbf{Q}}_{41} + \dot{\mathbf{Q}}_{14} + \dot{\mathbf{Q}}_{44})$,
 539 $d = -\frac{1}{2}(\dot{\mathbf{Q}}_{21} + \dot{\mathbf{Q}}_{31} + \dot{\mathbf{Q}}_{12} + \dot{\mathbf{Q}}_{42} + \dot{\mathbf{Q}}_{13} + \dot{\mathbf{Q}}_{43} + \dot{\mathbf{Q}}_{24} + \dot{\mathbf{Q}}_{34})$, and $e = \frac{1}{2}(\dot{\mathbf{Q}}_{22} + \dot{\mathbf{Q}}_{32} + \dot{\mathbf{Q}}_{23} + \dot{\mathbf{Q}}_{33})$.

540 **Case (b).** $\mathbf{V} = (\frac{-\tilde{c}}{\tilde{s}} \frac{-\tilde{s}}{\tilde{c}})$: $a = -\mathbf{P}_{11} + \mathbf{P}_{22}$, $b = -\mathbf{P}_{12} + \mathbf{P}_{21}$, $c = \frac{1}{2}(\dot{\mathbf{Q}}_{11} - \dot{\mathbf{Q}}_{41} - \dot{\mathbf{Q}}_{14} + \dot{\mathbf{Q}}_{44})$,
 541 $d = \frac{1}{2}(\dot{\mathbf{Q}}_{21} - \dot{\mathbf{Q}}_{31} + \dot{\mathbf{Q}}_{12} - \dot{\mathbf{Q}}_{42} - \dot{\mathbf{Q}}_{13} + \dot{\mathbf{Q}}_{43} - \dot{\mathbf{Q}}_{24} + \dot{\mathbf{Q}}_{34})$, and $e = \frac{1}{2}(\dot{\mathbf{Q}}_{22} - \dot{\mathbf{Q}}_{32} - \dot{\mathbf{Q}}_{23} + \dot{\mathbf{Q}}_{33})$.

542 **Case (c).** $\mathbf{V} = (\frac{\tilde{c}}{\tilde{s}} \frac{-\tilde{s}}{-\tilde{c}})$: $a = \mathbf{P}_{11} - \mathbf{P}_{22}$, $b = -\mathbf{P}_{12} + \mathbf{P}_{21}$, $c = \frac{1}{2}(\dot{\mathbf{Q}}_{11} - \dot{\mathbf{Q}}_{41} - \dot{\mathbf{Q}}_{14} + \dot{\mathbf{Q}}_{44})$,
 543 $d = \frac{1}{2}(-\dot{\mathbf{Q}}_{21} + \dot{\mathbf{Q}}_{31} - \dot{\mathbf{Q}}_{12} + \dot{\mathbf{Q}}_{42} + \dot{\mathbf{Q}}_{13} - \dot{\mathbf{Q}}_{43} + \dot{\mathbf{Q}}_{24} - \dot{\mathbf{Q}}_{34})$, and $e = \frac{1}{2}(\dot{\mathbf{Q}}_{22} - \dot{\mathbf{Q}}_{32} - \dot{\mathbf{Q}}_{23} + \dot{\mathbf{Q}}_{33})$.

544 **C Proofs for Section 2**

545 **C.1 Proof of Lemma 2.1**

546 *Proof.* Defining $\mathbf{J}_{\text{BB}} = \mathbf{J}(\mathbf{U}_{\text{B}}, \mathbf{U}_{\text{B}})$, then we have: $\mathbf{J} \mathbf{U}_{\text{B}} = \mathbf{U}_{\text{B}} \mathbf{J}_{\text{BB}}$, $\mathbf{U}_{\text{B}}^{\top} \mathbf{J} = \mathbf{J}_{\text{BB}}^{\top} \mathbf{U}_{\text{B}}^{\top}$, and $\mathbf{U}_{\text{B}}^{\top} \mathbf{J} \mathbf{U}_{\text{B}} =$
 547 \mathbf{J}_{BB} .

548 **Part (a).** For any $\mathbf{V} \in \mathbb{R}^{2 \times 2}$ and $\mathbf{B} \in \{\mathcal{B}_i\}_{i=1}^{\mathbf{C}_n^2}$, we have:

$$\begin{aligned}
& [\mathbf{X}^+]^\top \mathbf{J} \mathbf{X}^+ - \mathbf{X}^\top \mathbf{J} \mathbf{X} \\
& \stackrel{\textcircled{1}}{=} \mathbf{X}^\top \mathbf{J} \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X} + [\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X}]^\top \mathbf{J} \mathbf{X} \\
& \quad + [\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X}]^\top \mathbf{J} [\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X}] \\
& = \mathbf{X}^\top [\mathbf{J} \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top + \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2)^\top \mathbf{U}_B^\top \mathbf{J} + \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2)^\top \mathbf{U}_B^\top \mathbf{J} \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top] \mathbf{X} \\
& = \mathbf{X}^\top [\mathbf{U}_B \mathbf{J}_{BB} (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top + \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2)^\top \mathbf{J}_{BB} \mathbf{U}_B^\top + \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2)^\top \mathbf{J}_{BB} (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top] \mathbf{X} \\
& = \mathbf{X}^\top \mathbf{U}_B [\mathbf{J}_{BB} (\mathbf{V} - \mathbf{I}_2) + (\mathbf{V} - \mathbf{I}_2)^\top \mathbf{J}_{BB} + (\mathbf{V} - \mathbf{I}_2)^\top \mathbf{J}_{BB} (\mathbf{V} - \mathbf{I}_2)] \mathbf{U}_B^\top \mathbf{X} \\
& \stackrel{\textcircled{2}}{=} \mathbf{0}.
\end{aligned}$$

549 **Part (b).** Using the update rule for $\mathbf{X}^+ = \mathbf{X} + \mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X} \in \mathbb{R}^{n \times n}$, we derive:

$$\begin{aligned}
\|\mathbf{X}^+ - \mathbf{X}\|_F &= \|\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X}\|_F \\
&\stackrel{\textcircled{1}}{\leq} \|\mathbf{U}_B\|_F \cdot \|(\mathbf{V} - \mathbf{I}_2) \mathbf{U}_B^\top \mathbf{X}\|_F, \\
&\stackrel{\textcircled{2}}{\leq} \|\mathbf{U}_B\|_F \cdot \|(\mathbf{V} - \mathbf{I}_2)\|_F \cdot \|\mathbf{U}_B^\top\|_F \cdot \|\mathbf{X}\|_F, \\
&\stackrel{\textcircled{3}}{=} \|\mathbf{V} - \mathbf{I}_2\|_F \cdot \|\mathbf{X}\|_F,
\end{aligned}$$

550 where step $\textcircled{1}$ and step $\textcircled{2}$ use the norm inequality that $\|\mathbf{A} \mathbf{X}\|_F \leq \|\mathbf{A}\|_F \cdot \|\mathbf{X}\|_F$ for any \mathbf{A} and \mathbf{X} ;
551 step $\textcircled{3}$ uses $\|\mathbf{U}_B\| = \|\mathbf{U}_B^\top\| = 1$.

552 **Part (c).** We define $\mathbf{Z} \triangleq \mathbf{U}_B^\top \mathbf{X}$. We derive:

$$\begin{aligned}
\|\mathbf{X}^+ - \mathbf{X}\|_H^2 &= \|\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{Z}\|_H^2 \\
&\stackrel{\textcircled{1}}{=} \text{vec}(\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{Z})^\top \mathbf{H} \text{vec}(\mathbf{U}_B (\mathbf{V} - \mathbf{I}_2) \mathbf{Z}) \\
&\stackrel{\textcircled{2}}{=} \text{vec}(\mathbf{V} - \mathbf{I}_2)^\top (\mathbf{Z}^\top \otimes \mathbf{U}_B)^\top \mathbf{H} (\mathbf{Z}^\top \otimes \mathbf{U}_B) \text{vec}(\mathbf{V} - \mathbf{I}_2) \\
&= \|\mathbf{V} - \mathbf{I}_2\|_{(\mathbf{Z}^\top \otimes \mathbf{U}_B)^\top \mathbf{H} (\mathbf{Z}^\top \otimes \mathbf{U}_B)}^2 \\
&\stackrel{\textcircled{3}}{\leq} \|\mathbf{V} - \mathbf{I}_2\|_Q^2,
\end{aligned}$$

553 where step $\textcircled{1}$ uses $\|\mathbf{X}\|_H^2 = \text{vec}(\mathbf{X})^\top \mathbf{H} \text{vec}(\mathbf{X})$; step $\textcircled{2}$ uses $(\mathbf{Z}^\top \otimes \mathbf{R}) \text{vec}(\mathbf{U}) = \text{vec}(\mathbf{R} \mathbf{U} \mathbf{Z})$ for
554 all \mathbf{R} , \mathbf{Z} and \mathbf{U} of suitable dimensions; step $\textcircled{3}$ uses the choice of $\mathbf{Q} \succcurlyeq \underline{\mathbf{Q}} \triangleq (\mathbf{Z}^\top \otimes \mathbf{U}_B)^\top \mathbf{H} (\mathbf{Z}^\top \otimes$
555 $\mathbf{U}_B)$. \square

556 C.2 Proof of Lemma 2.3

557 *Proof.* We denote $w = c + e$. According to the properties of trigonometric functions, we have: (i)
558 $\tilde{c}^2 = \frac{1}{1-\tilde{t}^2}$; (ii) $\tilde{s}^2 = \frac{\tilde{t}^2}{1-\tilde{t}^2}$; (iii) $\tilde{t} = \frac{\tilde{s}}{\tilde{c}}$, leading to: $\tilde{c} = \frac{\pm 1}{\sqrt{1-\tilde{t}^2}}$, $\tilde{s} = \frac{\pm \tilde{t}}{\sqrt{1-\tilde{t}^2}}$ with $|\tilde{t}| < 1$.

559 We discuss two cases for Problem (8).

560 **Case (a).** $\tilde{c} = \frac{1}{\sqrt{1-\tilde{t}^2}}$, $\tilde{s} = \frac{\tilde{t}}{\sqrt{1-\tilde{t}^2}}$. Problem (8) is equivalent to the following problem: $\bar{\mu}_+ =$
561 $\arg \min_{\mu} \frac{a+\tilde{t}b}{\sqrt{1-\tilde{t}^2}} + \frac{w+\tilde{t}d}{1-\tilde{t}^2} - e$. Therefore, the optimal solution $\bar{\mu}_+$ can be computed as:

$$\cosh(\bar{\mu}_+) = \frac{1}{\sqrt{1-(\tilde{t}_+)^2}}, \text{ and } \sinh(\bar{\mu}_+) = \frac{\tilde{t}_+}{\sqrt{1-(\tilde{t}_+)^2}} \quad (14)$$

562 **Case (b).** $\tilde{c} = \frac{-1}{\sqrt{1-\tilde{t}^2}}$, $\tilde{s} = \frac{-\tilde{t}}{\sqrt{1-\tilde{t}^2}}$. Problem (8) is equivalent to the following problem: $\bar{\mu}_- =$
563 $\arg \min_{\mu} \frac{-a-\tilde{t}b}{\sqrt{1-\tilde{t}^2}} + \frac{w+\tilde{t}d}{1-\tilde{t}^2} - e$. Therefore, the optimal solution $\bar{\mu}_-$ can be computed as:

$$\cosh(\bar{\mu}_-) = \frac{-1}{\sqrt{1-(\tilde{t}_-)^2}}, \text{ and } \sinh(\bar{\mu}_-) = \frac{-\tilde{t}_-}{\sqrt{1-(\tilde{t}_-)^2}}. \quad (15)$$

564 We define the objective function as: $\check{F}(\tilde{c}, \tilde{s}) \triangleq a\tilde{c} + b\tilde{s} + c\tilde{c}^2 + d\tilde{c}\tilde{s} + e\tilde{s}^2$. In view of (14) and (15),
 565 the optimal solution pair $[\cosh(\bar{\mu}), \sinh(\bar{\mu})]$ for problem (8) can be computed as:

$$\begin{aligned} & [\cosh(\bar{\mu}), \sinh(\bar{\mu})] = \arg \min_{[c,s]} \check{F}(c, s), \\ & \text{s. t. } [c, s] \in \{[\cosh(\bar{\mu}_+), \sinh(\bar{\mu}_+)], [\cosh(\bar{\mu}_-), \sinh(\bar{\mu}_-)]\} \end{aligned}$$

566 Importantly, it is not necessary to compute the values $\bar{\mu}_+$ for (14) and $\bar{\mu}_-$ for (15).

567

□

568 C.3 Proof of Lemma 2.4

Proof. The objective function for $\mathbf{B}_{(i)}^t$ as in Equation (3) is formulated as :

$$f(\mathbf{X}^t) + \frac{1}{2} \|\mathbf{V}_i - \mathbf{I}\|_{\mathbf{Q} + \theta \mathbf{I}}^2 + \langle \mathbf{V}_i - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}_{(i)}^t, \mathbf{B}_{(i)}^t} \rangle$$

569 **Part (1).** For the part of $\frac{1}{2} \|\mathbf{V}_i - \mathbf{I}\|_{\mathbf{Q} + \theta \mathbf{I}}^2$, it is obviously irrelevant.

570 **Part (2).** For the part of $\langle \mathbf{V}_i - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}_{(i)}^t, \mathbf{B}_{(i)}^t} \rangle$, we note that $[\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}_{(i)}^t, \mathbf{B}_{(i)}^t} =$
 571 $[\nabla f(\mathbf{X}^t)](\mathbf{B}_{(i)}^t, :) [(\mathbf{X}^t)^\top](:, \mathbf{B}_{(i)}^t) = [\nabla f(\mathbf{X}^t)](\mathbf{B}_{(i)}^t, :) [(\mathbf{X}^t)(\mathbf{B}_{(i)}^t, :)]^\top$, which just use the informa-
 572 tion of block $\mathbf{B}_{(i)}^t$. The proof ends. □

573 C.4 Proof of Lemma 2.5

574 *Proof.* Part (a). For the purpose of analysis, we define the following: $\forall i \in [\frac{n}{2}], \mathbf{K}_i = \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i -$
 575 $\mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}$.

$$\begin{aligned} \|\sum_{i=1}^{\frac{n}{2}} [\mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}]\|_{\mathbb{F}}^2 & \stackrel{\textcircled{1}}{=} \left\| \begin{bmatrix} \mathbf{K}_1 \\ \mathbf{K}_2 \\ \vdots \\ \mathbf{K}_{\frac{n}{2}} \end{bmatrix} \right\|_{\mathbb{F}}^2 \\ & \stackrel{\textcircled{2}}{=} \|\mathbf{K}_1\|_{\mathbb{F}}^2 + \|\mathbf{K}_2\|_{\mathbb{F}}^2 + \cdots + \|\mathbf{K}_{\frac{n}{2}}\|_{\mathbb{F}}^2 \\ & \stackrel{\textcircled{3}}{=} \sum_{i=1}^{\frac{n}{2}} \|\mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{F}}^2 \end{aligned}$$

576 where step ① uses the definition of \mathbf{K}_i and the assumption that $\mathbf{B} \in \Upsilon$; step ② uses the definition of
 577 Squared Frobenius Norm; step ③ uses the definition of \mathbf{K}_i .

578 **Part (b).** Using the update rule for $\mathbf{X}^+ = \mathbf{X} + [\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top] \mathbf{X} \in \mathbb{R}^{n \times n}$, we have the
 579 following inequalities:

$$\|\mathbf{X}^+ - \mathbf{X}\|_{\mathbb{F}}^2 = \|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{F}}^2 \quad (16)$$

$$\stackrel{\textcircled{1}}{=} \sum_{i=1}^{n/2} \|\mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{F}}^2 \quad (17)$$

$$\stackrel{\textcircled{2}}{\leq} \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbb{F}}^2 \cdot \|\mathbf{X}\|_{\mathbb{F}}^2, \quad (18)$$

580 where step ① uses the conclusion of Part (a); step ② uses the same proof process of Part (b) of lemma
 581 2.1.

Part (c). We derive the following results:

$$\begin{aligned} \frac{1}{2} \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbb{H}}^2 & = \frac{1}{2} \|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{H}}^2 \\ & \stackrel{\textcircled{1}}{=} \frac{1}{2} \sum_{i=1}^{n/2} \|\mathbf{U}_{\mathbf{B}_{(i)}} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_{\mathbb{H}}^2 \\ & \stackrel{\textcircled{2}}{\leq} \frac{1}{2} \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbb{Q}}^2 \end{aligned}$$

582 where step ① uses the conclusion of Part (a); step ② uses the same proof process of Part (c) of lemma
583 2.1.

584 Part (d). We derive the following results:

$$\begin{aligned}
& \sum_{i=1}^{n/2} \langle \mathbf{V}_i - \mathbf{I}_2, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}})\mathbf{X}^\top]_{\mathbf{B}_i \mathbf{B}_i} \rangle \\
&= \sum_{i=1}^{n/2} \langle [\mathbf{U}_{\mathbf{B}_i}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_i}^\top]\mathbf{X}, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}})] \rangle \\
&= \langle \mathbf{X}^+ - \mathbf{X}, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}})] \rangle \\
&\stackrel{\textcircled{1}}{\leq} \frac{1}{2} \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbb{F}}^2 + \frac{1}{2} \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_{\mathbb{F}}^2 \\
&\stackrel{\textcircled{2}}{\leq} \frac{1}{2} \|\mathbf{X}\|_{\mathbb{F}}^2 \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbb{F}}^2 + \frac{1}{2} \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_{\mathbb{F}}^2 \tag{19}
\end{aligned}$$

585 where step ① uses $\forall \mathbf{A}, \mathbf{B}, \frac{1}{2} \|\mathbf{A} - \mathbf{B}\|_{\mathbb{F}}^2 = \frac{1}{2} \|\mathbf{A}\|_{\mathbb{F}}^2 + \frac{1}{2} \|\mathbf{B}\|_{\mathbb{F}}^2 - \langle \mathbf{A}, \mathbf{B} \rangle \geq 0$, with $\mathbf{A} = \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbb{F}}^2$
586 and $\mathbf{B} = \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_{\mathbb{F}}^2$; step ② uses the conclusion of Part (b). \square

587 D Proofs for Section 3

588 D.1 Proof of Lemma 3.1

589 *Proof.* We consider the Lagrangian function of problem (1):

$$\mathcal{L}(\mathbf{X}, \Lambda) = f(\mathbf{X}) - \frac{1}{2} \langle \Lambda, \mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J} \rangle. \tag{20}$$

590 Setting the gradient of $\mathcal{L}(\mathbf{X}, \Lambda)$ w.r.t. \mathbf{X} to zero yields:

$$\nabla f(\mathbf{X}) - \mathbf{J} \mathbf{X} \Lambda = \mathbf{0}. \tag{21}$$

591 **Part (a).** Multiplying both sides by \mathbf{X}^\top and using the fact that $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$, we have $\mathbf{J} \Lambda =$
592 $\mathbf{X}^\top \nabla f(\mathbf{X})$. Multiplying both sides by \mathbf{J}^\top and using $\mathbf{J}^\top \mathbf{J} = \mathbf{I}$, we have $\Lambda = \mathbf{J} \mathbf{X}^\top \nabla f(\mathbf{X})$. Since Λ
593 is symmetric, we have $\Lambda = \nabla f(\mathbf{X})^\top \mathbf{X} \mathbf{J}$. Putting this equality into Equality (21) yields the following
594 first-order optimality condition for Problem (1):

$$\nabla f(\mathbf{X}) = \mathbf{J} \mathbf{X} [\nabla f(\mathbf{X})]^\top \mathbf{X} \mathbf{J}. \tag{22}$$

595 **Part (b).** We let $\mathbf{G} = \nabla f(\mathbf{X})$. We derive the following results:

$$\begin{aligned}
\mathbf{G} = \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J} &\stackrel{\textcircled{1}}{\Rightarrow} \mathbf{J} \mathbf{X}^\top \cdot \mathbf{G} = \mathbf{J} \mathbf{X}^\top \cdot \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J} \\
&\stackrel{\textcircled{2}}{\Rightarrow} \mathbf{J} \mathbf{X}^\top \mathbf{G} = \mathbf{G}^\top \mathbf{X} \mathbf{J} \\
&\stackrel{\textcircled{3}}{\Rightarrow} \mathbf{X} (\mathbf{J} \mathbf{X}^\top \mathbf{G}) \mathbf{X}^\top = \mathbf{X} (\mathbf{G}^\top \mathbf{X} \mathbf{J}) \mathbf{X}^\top \\
&\stackrel{\textcircled{4}}{\Rightarrow} \mathbf{X} \underbrace{\mathbf{J} \mathbf{X}^\top \mathbf{G} \mathbf{X}^\top \mathbf{J}}_{\triangleq \mathbf{G}^\top} \mathbf{J} = \mathbf{J} \underbrace{\mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J}}_{\triangleq \mathbf{G}} \mathbf{X}^\top \tag{23} \\
&\stackrel{\textcircled{5}}{\Rightarrow} (\mathbf{X} \mathbf{G}^\top \mathbf{J}) \cdot \mathbf{J} \mathbf{X} = (\mathbf{J} \mathbf{G} \mathbf{X}^\top) \cdot \mathbf{J} \mathbf{X} \\
&\stackrel{\textcircled{6}}{\Rightarrow} \mathbf{X} \mathbf{G}^\top \mathbf{X} = \mathbf{J} \mathbf{G} \mathbf{J} \\
&\stackrel{\textcircled{7}}{\Rightarrow} \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J} = \mathbf{G},
\end{aligned}$$

596 where step ① uses the results of left-multiplying both sides by $\mathbf{J} \mathbf{X}^\top$; step ② uses $\mathbf{J} \cdot \mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J} \mathbf{J} = \mathbf{I}$;
597 step ③ uses the results of left-multiplying both sides by \mathbf{X} and subsequently right-multiplying them
598 by \mathbf{X}^\top ; ④ uses $\mathbf{G} = \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J}$; step ⑤ uses the results of right-multiplying both sides by $\mathbf{J} \mathbf{X}$;
599 step ⑥ uses $\mathbf{J} \mathbf{J} = \mathbf{I}$ and $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$; step ⑦ uses the results of left-multiply both sides by \mathbf{J} and
600 right-multiplied by \mathbf{J} .

601 Given Equality (23), we conclude that the critical point condition is equivalent to the requirement
602 that the matrix $\mathbf{X} \nabla f(\tilde{\mathbf{X}})^\top \mathbf{J}$ is symmetric, which is expressed as $\mathbf{X} \mathbf{G}^\top \mathbf{J} = [\mathbf{X} \mathbf{G}^\top \mathbf{J}]^\top$. \square

603 D.2 Proof of Theorem 3.3

604 *Proof.* We use $\check{\mathbf{X}}$ and $\check{\check{\mathbf{X}}}$ to denote any BS-point and critical point, respectively.

605 For all $\mathbf{B} \in \Omega \triangleq \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{C_n^2}\}$, we have:

$$\mathbf{I}_2 \in \arg \min_{\mathbf{V} \in \mathcal{J}_{\mathbf{B}}} \mathcal{G}(\mathbf{V}; \ddot{\mathbf{X}}, \mathbf{B}).$$

606 where $\mathcal{G}(\mathbf{V}; \mathbf{X}, \mathbf{B}) \triangleq f(\mathbf{X}) + \frac{1}{2} \|\mathbf{V} - \mathbf{I}_2\|_{\mathbf{Q} + \theta \mathbf{I}}^2 + \langle \mathbf{V} - \mathbf{I}, [\nabla f(\mathbf{X})(\mathbf{X})^\top]_{\mathbf{BB}} \rangle$.

607 The Euclidean gradient of $\mathcal{G}(\mathbf{V}; \ddot{\mathbf{X}}, \mathbf{B})$ can be computed as:

$$\ddot{\mathbf{G}} \triangleq \text{mat}((\mathbf{Q} + \theta \mathbf{I}_2) \text{vec}(\mathbf{V} - \mathbf{I}_2)) + [\nabla f(\ddot{\mathbf{X}})(\ddot{\mathbf{X}})^\top]_{\mathbf{BB}}. \quad (24)$$

608 Given Lemma 3.1, we set the Riemannian gradient of $\mathcal{G}(\mathbf{V}; \ddot{\mathbf{X}}, \mathbf{B})$ w.r.t. \mathbf{V} to zero, leading to the
609 following first-order optimality condition:

$$\mathbf{0} = \nabla_{\mathcal{J}} \mathcal{G}(\mathbf{V}; \ddot{\mathbf{X}}, \mathbf{B}) = \ddot{\mathbf{G}} - \mathbf{U}_{\mathbf{B}}^\top \mathbf{J} \mathbf{V} \ddot{\mathbf{G}}^\top \mathbf{V} \mathbf{J} \mathbf{U}_{\mathbf{B}}. \quad (25)$$

610 Letting $\mathbf{V} = \mathbf{I}_2$, and using the definition of $\ddot{\mathbf{G}}$, we have:

$$\begin{aligned} & \mathbf{0}_{2,2} = [\nabla f(\mathbf{X})(\mathbf{X})^\top]_{\mathbf{BB}} - \mathbf{J}_{\mathbf{BB}} \ddot{\mathbf{G}}^\top \mathbf{J}_{\mathbf{BB}}, \quad \forall \mathbf{B} \in \{\mathcal{B}_i\}_{i=1}^{C_n^2} \\ \Rightarrow & \mathbf{0}_{2,2} = \mathbf{U}_{\mathbf{B}}^\top [\nabla f(\ddot{\mathbf{X}}) \ddot{\mathbf{X}}^\top]_{\mathbf{U}_{\mathbf{B}}} - \mathbf{J}_{\mathbf{BB}} \mathbf{U}_{\mathbf{B}}^\top [\ddot{\mathbf{X}} \nabla f(\ddot{\mathbf{X}})^\top]_{\mathbf{U}_{\mathbf{B}}} \mathbf{J}_{\mathbf{BB}}, \quad \forall \mathbf{B} \in \{\mathcal{B}_i\}_{i=1}^{C_n^2} \\ \stackrel{\textcircled{1}}{\Rightarrow} & \mathbf{0}_{2,2} = \mathbf{U}_{\mathbf{B}}^\top [\nabla f(\ddot{\mathbf{X}}) \ddot{\mathbf{X}}^\top]_{\mathbf{U}_{\mathbf{B}}} - \mathbf{U}_{\mathbf{B}}^\top \mathbf{J} [\ddot{\mathbf{X}} \nabla f(\ddot{\mathbf{X}})^\top]_{\mathbf{J}} \mathbf{U}_{\mathbf{B}}, \quad \forall \mathbf{B} \in \{\mathcal{B}_i\}_{i=1}^{C_n^2} \\ \stackrel{\textcircled{2}}{\Rightarrow} & \mathbf{0}_{n,n} = [\nabla f(\ddot{\mathbf{X}}) \ddot{\mathbf{X}}^\top] - \mathbf{J} [\ddot{\mathbf{X}} \nabla f(\ddot{\mathbf{X}})^\top]_{\mathbf{J}}, \\ \stackrel{\textcircled{3}}{\Rightarrow} & [\mathbf{J} \nabla f(\ddot{\mathbf{X}}) \ddot{\mathbf{X}}^\top] = [\mathbf{J} \nabla f(\ddot{\mathbf{X}}) \ddot{\mathbf{X}}^\top]^\top, \end{aligned}$$

611 where step $\textcircled{1}$ uses $\mathbf{U}_{\mathbf{B}}^\top \mathbf{J} = \mathbf{J}_{\mathbf{BB}} \mathbf{U}_{\mathbf{B}}^\top$ and $\mathbf{J} \mathbf{U}_{\mathbf{B}} = \mathbf{U}_{\mathbf{B}} \mathbf{J}_{\mathbf{BB}}$; step $\textcircled{2}$ uses the the following results for any
612 $\mathbf{W} \in \mathbb{R}^{n \times n}$:

$$(\forall \mathbf{B} \in \{\mathcal{B}_i\}_{i=1}^{C_n^2}, \mathbf{0}_{2,2} = \mathbf{U}_{\mathbf{B}}^\top \mathbf{W} \mathbf{U}_{\mathbf{B}} = \mathbf{W}_{\mathbf{BB}}) \Rightarrow (\mathbf{W} = \mathbf{0}_{n,n}); \quad (26)$$

613 step $\textcircled{3}$ uses the fact that both sides are left-multiplied by \mathbf{J} . We conclude that the matrix $\mathbf{J} \nabla f(\ddot{\mathbf{X}}) \ddot{\mathbf{X}}^\top$
614 is symmetric. Using Claim (b) of Lemma 3.1, we conclude that $\ddot{\mathbf{X}}$ is also a critical point.

615 Notably, the condition in Equation (25) is a necessary but not sufficient condition. This is because
616 BS-point is the global minimum of Problem: $\arg \min_{\mathbf{V} \in \mathcal{J}_{\mathbf{B}}} \mathcal{G}(\mathbf{V}; \ddot{\mathbf{X}}, \mathbf{B})$, according to Definition
617 3.2. \square

618 E Proofs for Section 4

619 E.1 Proof of Lemma 4.5

620 *Proof.* By the definition of $\tilde{\mathbf{G}}^t$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}^t} [\|\tilde{\mathbf{G}}^t - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] \\ \stackrel{\textcircled{1}}{=} & p \mathbb{E}_{\mathbf{X}^t} [\|\frac{1}{b} \sum_{i=1}^b \nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] + \\ & (1-p) \mathbb{E}_{\mathbf{X}^t} [\|\tilde{\mathbf{G}}^{t-1} + \frac{1}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] \\ \stackrel{\textcircled{2}}{=} & p \mathbb{E}_{\mathbf{X}^t} [\|\frac{1}{b} \sum_{i=1}^b \nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] + (1-p) \mathbb{E}_{\mathbf{X}^{t-1}} [\|\tilde{\mathbf{G}}^{t-1} - \nabla f(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] \\ & + (1-p) \mathbb{E}_{\mathbf{X}^t} [\|\frac{1}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})) - \nabla f(\mathbf{X}^t) + \nabla f(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] \end{aligned}$$

621 where step ① uses formula (9); step ② uses that $\tilde{\mathbf{G}}^{t-1} - \nabla f(\mathbf{X}^{t-1})$ is measurable w.r.t. t^{t-1} and
 622 $\mathbb{E}_{\iota^t}[\|\frac{1}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})) - \nabla f(\mathbf{X}^t) + \nabla f(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] = 0$. We further have

$$\begin{aligned}
 & \mathbb{E}_{\iota^t}[\|\tilde{\mathbf{G}}^t - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] \\
 \stackrel{\textcircled{1}}{\leq} & p\mathbb{E}_{\iota^t}[\|\frac{1}{b} \sum_{i=1}^b \nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] + (1-p)\mathbb{E}_{\iota^{t-1}}[\|\tilde{\mathbf{G}}^{t-1} - \nabla f(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] \\
 & + (1-p)\mathbb{E}_{\iota^t}[\|\frac{1}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1}))\|_{\mathbb{F}}^2] \\
 \stackrel{\textcircled{2}}{\leq} & \frac{p(N-b)}{b(N-1)}\mathbb{E}_{\iota^t}[\|\nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}^2] + (1-p)\mathbb{E}_{\iota^{t-1}}[\|\tilde{\mathbf{G}}^{t-1} - \nabla f(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] \\
 & + \frac{1-p}{b'}\mathbb{E}_{\iota^{t-1}}[\|\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] \\
 \stackrel{\textcircled{3}}{\leq} & \frac{p(N-b)}{b(N-1)}\sigma^2 + (1-p)\mathbb{E}_{\iota^{t-1}}[\|\tilde{\mathbf{G}}^{t-1} - \nabla f(\mathbf{X}^{t-1})\|_{\mathbb{F}}^2] \\
 & + \frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{b'}\mathbb{E}_{\iota^{t-1}}[\sum_{i=1}^{n/2} \|\mathbf{V}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2] \tag{27}
 \end{aligned}$$

623 where step ① uses that for any random variable \mathbf{X} , $\mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2] \leq \mathbb{E}[\mathbf{X}^2]$; step ② uses lemma
 624 A.2; step ③ uses assumption 4.3, Inequality (2) and Part (b) of lemma 2.5. \square

625 E.2 Proof of theorem 4.6

626 *Proof.* For simplicity, we use \mathbf{B} instead of \mathbf{B}^t . We will show that the following inequality holds :

$$\frac{\theta}{2} \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2 \leq f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}). \tag{28}$$

Since $\bar{\mathbf{V}}^t$ is the global optimal solution of Problem (5), we have:

$$\mathcal{G}(\bar{\mathbf{V}}^t; \mathbf{X}^t, \mathbf{B}) \leq \mathcal{G}(\mathbf{V}; \mathbf{X}^t, \mathbf{B}), \mathbf{V} \in \mathcal{J}_{\mathbf{B}}$$

627 Letting $\mathbf{V} = \mathbf{I}_2$, we have: $\mathcal{G}(\bar{\mathbf{V}}^t; \mathbf{X}^t, \mathbf{B}) \leq \mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})$. We further obtain:

$$\frac{1}{2} \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbf{Q} + \theta \mathbf{I}}^2 + \langle \bar{\mathbf{V}}^t - \mathbf{I}_2, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} \rangle \leq 0. \tag{29}$$

628 Using Inequality (2) with $N = 1$ and Part (c) of Lemma 2.1, we have:

$$f(\mathbf{X}^{t+1}) \leq f(\mathbf{X}^t) + \langle \bar{\mathbf{V}}^t - \mathbf{I}_2, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} \rangle + \frac{1}{2} \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbf{Q}}^2. \tag{30}$$

629 Adding Inequality (29) and (30) together, we obtain the inequality in (28). Using the result of Part (b)
 630 in Lemma 2.1 that $\frac{\|\mathbf{X}^t - \mathbf{X}\|_{\mathbb{F}}^2}{\|\mathbf{X}\|_{\mathbb{F}}^2} \leq \|\mathbf{V} - \mathbf{I}_2\|_{\mathbb{F}}^2$, we have the following sufficient decrease condition:

$$f(\mathbf{X}^{t+1}) - f(\mathbf{X}^t) \leq -\frac{\theta}{2} \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2 \leq -\frac{\theta}{2} \frac{\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2}{\|\mathbf{X}^t\|_{\mathbb{F}}^2} \tag{31}$$

We now prove the global convergence. Taking the expectation for Inequality (31), we obtain a lower bound on the expected progress made by each iteration for Algorithm 1:

$$\mathbb{E}_{\xi^{t+1}}[f(\mathbf{X}^{t+1})] - \mathbb{E}_{\xi^t}[f(\mathbf{X}^t)] \leq -\mathbb{E}_{\xi^t}[\frac{\theta}{2} \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2].$$

Summing up the inequality above over $t = 0, 1, \dots, T$, we have:

$$\mathbb{E}_{\xi^t}[\frac{\theta}{2} \sum_{t=0}^T \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \leq f(\mathbf{X}^0) - \mathbb{E}_{\xi^{T+1}}[f(\mathbf{X}^{T+1})] \leq f(\mathbf{X}^0) - f(\bar{\mathbf{X}}).$$

631 As a result, there exists an index \bar{t} with $0 \leq \bar{t} \leq T$ such that

$$\mathbb{E}_{\xi^{\bar{t}}}[\|\bar{\mathbf{V}}^{\bar{t}} - \mathbf{I}_2\|_{\mathbb{F}}^2] \leq \frac{2}{\theta(T+1)} [f(\mathbf{X}^0) - f(\bar{\mathbf{X}})]. \tag{32}$$

632 Furthermore, for any t , we have:

$$\mathcal{E}(\mathbf{X}^t) \triangleq \frac{1}{C_n^2} \sum_{i=1}^{C_n^2} \text{dist}(\mathbf{I}_2, \arg \min_{\mathbf{V}} \mathcal{G}(\mathbf{V}; \mathbf{X}^t, \mathcal{B}_i))^2 = \mathbb{E}_{\xi^t}[\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \tag{33}$$

633 Combining Inequality (32) and equality (33), we have the following result:

$$\mathbb{E}_{\xi^t}[\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2] = \mathcal{E}(\mathbf{X}^{\bar{t}}) \leq \frac{2(f(\mathbf{X}^0) - f(\bar{\mathbf{X}}))}{\theta(T+1)}. \tag{34}$$

We will give the arithmetic operations of **GS-JOBCD**. By the chosen parameters and Inequality (34), we have

$$\mathcal{E}(\mathbf{X}^{\bar{t}}) \leq \frac{2(f(\mathbf{X}^0) - f(\bar{\mathbf{X}}))}{\theta(T+1)} \leq \epsilon.$$

We define $\Delta_0 = f(\mathbf{X}_0) - f(\bar{\mathbf{X}})$ and set $T + 1 = \frac{2\Delta_0}{\epsilon\theta}$. Denoting m_t to be the number of arithmetic operations at t -th iteration, we have for $t \geq 1$:

$$\mathbb{E}_{\xi^t} [m_t] = \mathcal{O}(2N).$$

Then we have for $t \geq 1$, the total number of arithmetic operations M^T in T iterations to obtain ϵ -BS-point is

$$\mathbb{E}_{\xi^T} [M^T] = \mathbb{E}_{\xi^t} [\sum_{t=0}^T m_t] = 2(T+1)N = \mathcal{O}((T+1)N).$$

634 We have $(T+1)N = N \frac{2\Delta_0}{\epsilon\theta} = \mathcal{O}(\frac{\Delta_0 N}{\epsilon})$. □

635 E.3 Proof of Theorem 4.7

Proof. For simplicity, we use \mathbf{B} instead of \mathbf{B}^t . Defining $\bar{\mathbf{V}}_i^t$ as the global optimal solution of $\arg \min_{\mathbf{V}_i} \mathcal{T}(\mathbf{V}_i; \mathbf{X}^t, \mathbf{B})$, we have:

$$\mathcal{T}(\bar{\mathbf{V}}_i^t; \mathbf{X}^t, \mathbf{B}) \leq \mathcal{T}(\mathbf{V}_i; \mathbf{X}^t, \mathbf{B}), \forall i, \mathbf{V}_i \in \mathcal{J}_{\mathbf{B}(i)}$$

636 Letting $\mathbf{V}_i = \mathbf{I}_2, \forall i$, we have: $\mathcal{T}(\bar{\mathbf{V}}_i^t; \mathbf{X}^t, \mathbf{B}) \leq \mathcal{T}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})$. We further obtain:

$$\frac{1}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{(\zeta+\theta)\mathbf{I}}^2 + \sum_{i=1}^{n/2} \langle \bar{\mathbf{V}}_i^t - \mathbf{I}_2, [\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top]_{\mathbf{B}(i)\mathbf{B}(i)} \rangle \leq 0. \quad (35)$$

637 Using the results of telescoping Inequality (2) over i from 1 to N with Part (c) of Lemma 2.5, we
638 have:

$$f(\mathbf{X}^{t+1}) \leq f(\mathbf{X}^t) + \sum_{i=1}^{n/2} \langle \bar{\mathbf{V}}_i^t - \mathbf{I}_2, [\nabla f(\mathbf{X})\mathbf{X}^\top]_{\mathbf{B}(i)\mathbf{B}(i)} \rangle + \frac{1}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\zeta\mathbf{I}}^2. \quad (36)$$

639 Adding inequality (35), and (36) together, we obtain the inequality in (37).

$$\begin{aligned} & \frac{\theta}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2 \\ & \leq f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}) + \sum_{i=1}^{n/2} \langle \bar{\mathbf{V}}_i^t - \mathbf{I}_2, [(\nabla f(\mathbf{X}^t) - \tilde{\mathbf{G}}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}(i)\mathbf{B}(i)} \rangle \\ & \stackrel{\textcircled{1}}{\leq} f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}) + \frac{1}{2} \|\mathbf{X}^t\|_{\mathbb{F}}^2 \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2 + \frac{1}{2} \|[\nabla f(\mathbf{X}^t) - \tilde{\mathbf{G}}^t]\|_{\mathbb{F}}^2 \end{aligned} \quad (37)$$

640 where step $\textcircled{1}$ uses Part (d) of Lemma 2.5.

641 Taking expectation on both sides of inequality (37) with respect to all randomness of the algorithm,
642 and adding the inequality in Lemma 4.5 $\times \frac{1}{2p}$ to (37), we have:

$$\begin{aligned} & \left(\frac{\theta - \bar{\mathbf{X}}^2}{2} - \frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{2pb'} \right) \mathbb{E}_{\iota^t} [\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \\ & \leq \mathbb{E}_{\iota^t} [f(\mathbf{X}^t)] - \mathbb{E}_{\iota^{t+1}} [f(\mathbf{X}^{t+1})] + \frac{(N-b)}{2b(N-1)} \sigma^2 + \frac{1-p}{2p} (\mathbb{E}_{\iota^t} [u^t] - \mathbb{E}_{\iota^{t+1}} [u^{t+1}]) \end{aligned} \quad (38)$$

643 Summing up the inequality above over $t = 0, 1, \dots, T$, we have:

$$\begin{aligned} & \left(\frac{\theta - \bar{\mathbf{X}}^2}{2} - \frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{2pb'} \right) \mathbb{E}_{\iota^T} [\sum_{t=0}^T \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \\ & \leq f(\mathbf{X}^0) - \mathbb{E}_{\iota^T} [f(\mathbf{X}^T)] + \frac{(T+1)(N-b)}{2b(N-1)} \sigma^2 + \frac{1-p}{2p} (u^0 - \mathbb{E}_{\iota^{T+1}} [u^{T+1}]) \\ & \leq f(\mathbf{X}^0) - f(\bar{\mathbf{X}}) + \frac{(T+1)(N-b)}{2b(N-1)} \sigma^2 + \frac{1-p}{2p} (u^0 - \mathbb{E}_{\iota^{T+1}} [u^{T+1}]) \end{aligned} \quad (39)$$

644 As a result, there exists an index \bar{t} with $0 \leq \bar{t} \leq T$ such that

$$\begin{aligned} & \left(\frac{\theta - \bar{\mathbf{X}}^2}{2} - \frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{2pb'} \right) (T+1) \mathbb{E}_{\iota^{\bar{t}}} [\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{\bar{t}} - \mathbf{I}_2\|_{\mathbb{F}}^2] \\ & \leq f(\mathbf{X}^0) - f(\bar{\mathbf{X}}) + \frac{(T+1)(N-b)}{2b(N-1)} \sigma^2 + \frac{1-p}{2p} (u^0 - \mathbb{E}_{\iota^{T+1}} [u^{T+1}]) \end{aligned} \quad (40)$$

645 Defining $\varpi = \frac{\theta - \bar{\mathbf{X}}^2}{2} - \frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{2pb'}$, furthermore, for any t and $\forall i$, we have:

$$\mathcal{E}(\mathbf{X}^t) = \frac{1}{C_J} \sum_{i=1}^{C_J} \mathbb{E}_{\iota^t} [\text{dist}(\mathbf{I}_2, \arg \min_{\mathbf{V}_i} \mathcal{T}(\mathbf{V}_i; \mathbf{X}^t, \tilde{\mathbf{B}}_i))^2] = \mathbb{E}_{\iota^t} [\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \quad (41)$$

646 Combining inequality (40) and (41), we have the following result:

$$\mathcal{E}(\mathbf{X}^{\bar{t}}) \leq \frac{1}{(T+1)\varpi} (f(\mathbf{X}^0) - f(\bar{\mathbf{X}}) + \frac{(T+1)(N-b)}{2b(N-1)} \sigma^2 + \frac{1-p}{2p} (u^0 - \mathbb{E}_{\mathcal{L}^{T+1}}[u^{T+1}])) \quad (42)$$

By the chosen parameters and Inequality (42), we have

$$\mathcal{E}(\mathbf{X}^{\bar{t}}) \leq \frac{1}{(T+1)\varpi} (f(\mathbf{X}^0) - f(\bar{\mathbf{X}}) + \frac{(T+1)(N-b)}{2b(N-1)} \sigma^2 + \frac{1-p}{2p} (u^0 - \mathbb{E}_{\mathcal{L}^{T+1}}[u^{T+1}])) \leq \epsilon.$$

We define $\Delta_0 = f(\mathbf{X}_0) - f(\bar{\mathbf{X}})$ and set $T+1 = \frac{\Delta_0}{\epsilon\varpi}$. Denoting m_t^i to be the number of arithmetic operations to update the i -th block at t -th iteration, we have for $t \geq 1$

$$\mathbb{E}_{\mathcal{L}^t}[m_t^i] = \mathcal{O}(2(pb + (1-p)b')).$$

Letting m_t be the number of arithmetic operations in the t -th iteration, we have for $t \geq 1$

$$\mathbb{E}_{\mathcal{L}^t}[m_t] = \mathbb{E}_{\mathcal{L}^t}[\sum_{i=1}^{n/2} m_t^i] = \mathcal{O}((pb + (1-p)b')n/2 \times 2) = \mathcal{O}(n(pb + (1-p)b')).$$

Hence, the total number of arithmetic operations M^T in T iterations to obtain ϵ -BS-point is

$$\mathbb{E}_{\mathcal{L}^T}[M] = \mathbb{E}_{\mathcal{L}^T}[\sum_{t=0}^T m_t] = \mathcal{O}(bn) + \mathbb{E}_{\mathcal{L}^T}[\sum_{t=1}^T m_t] = \mathcal{O}(bn + Tn(pb + (1-p)b')).$$

Since $b = N, b' = \sqrt{b}$ and $p = \frac{b'}{b+b'}, \varpi = \frac{\theta - \bar{\mathbf{X}}^2}{2} - \frac{L_f^2 \bar{\mathbf{X}}^2 (1-p)}{2pb'} = \frac{1}{2}(\theta - \bar{\mathbf{X}}^2 - L_f^2 \bar{\mathbf{X}}^2)$, we have

$$nT(pb + (1-p)b') = n \frac{\Delta_0}{\epsilon(\theta - \bar{\mathbf{X}}^2 - L_f^2 \bar{\mathbf{X}}^2)} \frac{2bb'}{b+b'} \leq \frac{n\Delta_0}{\epsilon(\theta - \bar{\mathbf{X}}^2 - L_f^2 \bar{\mathbf{X}}^2)} 2b' = \mathcal{O}\left(\frac{\Delta_0 \sqrt{N}}{\epsilon}\right).$$

647

□

648 E.4 Proof of Theorem 4.10

649 *Proof.* For simplicity, we use \mathbf{B} instead of \mathbf{B}^t . We notice that the Riemannian gradient of $\mathcal{T}(\mathbf{V};; \mathbf{X}^t, \mathbf{B})$
 650 at the point $\mathbf{V}_i = \mathbf{I}_2, \forall i$. Defining $\mathbf{G} = \tilde{\mathbf{G}}^t[\mathbf{X}^t]^\top$ and using $\mathbf{J}\mathbf{U}_{\mathbf{B}} = \mathbf{U}_{\mathbf{B}}\mathbf{J}_{\mathbf{B}\mathbf{B}}, \mathbf{U}_{\mathbf{B}}^\top \mathbf{J} = \mathbf{J}_{\mathbf{B}\mathbf{B}} \mathbf{U}_{\mathbf{B}}^\top$, we
 651 have:

$$\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{V};; \mathbf{I}_2; \mathbf{X}^t, \mathbf{B}) = \sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}(i)}^\top \mathbf{G} \mathbf{U}_{\mathbf{B}(i)} - \mathbf{U}_{\mathbf{B}(i)}^\top \mathbf{J} \mathbf{G}^\top \mathbf{J} \mathbf{U}_{\mathbf{B}(i)} \quad (43)$$

652 Then, we prove the following important lemmas.

653 **Lemma E.1.** *We have the following result for VR-J-JOBCD: $\mathbb{E}_{\mathcal{L}^{t+1}}[\|\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^{t+1}\|_{\mathbb{F}}] \leq p\mathbb{E}_{\mathcal{L}^t}[\sqrt{u^t}] +$*
 654 *$L_f \mathbb{E}_{\mathcal{L}^{t+1}}[\|\mathbf{X}^t - \mathbf{X}^{t+1}\|_{\mathbb{F}}]$*

655 *Proof.* By the definition of $\tilde{\mathbf{G}}^t$, with the choice of $b = N, b' = \sqrt{b}$ and $p = \frac{b'}{b+b'}$, we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{L}^{t+1}}[\|\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^{t+1}\|_{\mathbb{F}}] \\ \stackrel{\textcircled{1}}{=} & \mathbb{E}_{\mathcal{L}^{t+1}}[\|\tilde{\mathbf{G}}^t - \frac{p}{b} \sum_{i=1}^b \nabla f_i(\mathbf{X}^{t+1}) - \frac{1-p}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^{t+1}) - \nabla f_i(\mathbf{X}^t)) - (1-p)\tilde{\mathbf{G}}^t\|_{\mathbb{F}}] \\ = & \mathbb{E}_{\mathcal{L}^{t+1}}[\|p\tilde{\mathbf{G}}^t - \frac{p}{b} \sum_{i=1}^b \nabla f_i(\mathbf{X}^{t+1}) - \frac{1-p}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^{t+1}) - \nabla f_i(\mathbf{X}^t))\|_{\mathbb{F}}] \\ \stackrel{\textcircled{2}}{\leq} & p\mathbb{E}_{\mathcal{L}^{t+1}}[\|\tilde{\mathbf{G}}^t - \nabla f(\mathbf{X}^{t+1})\|_{\mathbb{F}}] + \frac{1-p}{b'} \mathbb{E}_{\mathcal{L}^t}[\|\sum_{i=1}^{b'} \nabla f_i(\mathbf{X}^{t+1}) - \nabla f_i(\mathbf{X}^t)\|_{\mathbb{F}}] \\ \stackrel{\textcircled{3}}{\leq} & p\mathbb{E}_{\mathcal{L}^t}[\|\tilde{\mathbf{G}}^t - \nabla f(\mathbf{X}^t)\|_{\mathbb{F}}] + p\mathbb{E}_{\mathcal{L}^{t+1}}[\|\nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}^{t+1})\|_{\mathbb{F}}] \\ & + \frac{1-p}{b'} \mathbb{E}_{\mathcal{L}^{t+1}}[\|\sum_{i=1}^{b'} \nabla f_i(\mathbf{X}^{t+1}) - \nabla f_i(\mathbf{X}^t)\|_{\mathbb{F}}] \\ \stackrel{\textcircled{4}}{\leq} & p\mathbb{E}_{\mathcal{L}^t}[\sqrt{u^t}] + p\mathbb{E}_{\mathcal{L}^{t+1}}[\|\nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}^{t+1})\|_{\mathbb{F}}] + (1-p)\mathbb{E}_{\mathcal{L}^{t+1}}[\|\nabla f_i(\mathbf{X}^{t+1}) - \nabla f_i(\mathbf{X}^t)\|_{\mathbb{F}}] \\ \stackrel{\textcircled{5}}{\leq} & p\mathbb{E}_{\mathcal{L}^t}[\sqrt{u^t}] + L_f \mathbb{E}_{\mathcal{L}^{t+1}}[\|\mathbf{X}^t - \mathbf{X}^{t+1}\|_{\mathbb{F}}] \end{aligned}$$

656 where step ① uses formula (9); step ② uses norm inequality and $\frac{1}{b} \sum_{i=1}^b \nabla f_i(\mathbf{X}^{t+1}) = \nabla f(\mathbf{X}^{t+1})$
 657 with $b = N$ and norm inequality; step ③ uses triangle inequality that $\|\mathbf{A} - \mathbf{B}\|_{\mathbb{F}} \leq \|\mathbf{A} - \mathbf{C}\|_{\mathbb{F}} +$
 658 $\|\mathbf{C} - \mathbf{B}\|_{\mathbb{F}}$, for any \mathbf{A}, \mathbf{B} and \mathbf{C} ; step ④ the definition of u^t ; step ⑤ uses Inequality (2) and the results
 659 of telescoping it over i from 1 to N . □

660 **Lemma E.2.** (*Riemannian gradient Lower Bound for the Iterates Gap*) We de-
661 fine $\phi \triangleq (3\bar{X} + \overline{VX})\bar{G} + (1 + \bar{V}^2 + \frac{n}{2}(\bar{X}^2 + \bar{V}^2\bar{X}^2))L_f + (1 + \bar{V}^2)\theta$. It holds that:
662 $\mathbb{E}_{\iota^{t+1}}[\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}\mathcal{T}(\mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1}))] \leq \phi \cdot \mathbb{E}_{\iota^t}[\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_F] + \frac{np\sqrt{u^t}}{2}(\bar{X} + \bar{V}^2\bar{X})$.

663 *Proof.* For notation simplicity, we define:

$$\Omega_{i0} \triangleq \mathbf{U}_{\mathbf{B}^{(i)}}^\top [\tilde{\mathbf{G}}^{t+1}][\mathbf{X}^{t+1}]^\top \mathbf{U}_{\mathbf{B}^{(i)}}, \forall i \quad (44)$$

$$\Omega_{i1} \triangleq \mathbf{U}_{\mathbf{B}^{(i)}}^\top [\tilde{\mathbf{G}}^{t+1}][\mathbf{X}^t]^\top \mathbf{U}_{\mathbf{B}^{(i)}}, \forall i, \quad (45)$$

$$\Omega_{i2} \triangleq \mathbf{U}_{\mathbf{B}^{(i)}}^\top [\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^{t+1}][\mathbf{X}^t]^\top \mathbf{U}_{\mathbf{B}^{(i)}}, \forall i. \quad (46)$$

664 First, using the optimality of $\bar{\mathbf{V}}_i^t, i \in \{1, \dots, \frac{n}{2}\}$ for the subproblem, we have:

$$\mathbf{0}_{2,2} = \tilde{\mathbf{G}}_i - \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t \tilde{\mathbf{G}}_i^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} \quad (47)$$

$$\text{where } \tilde{\mathbf{G}}_i = \underbrace{\text{mat}((\mathbf{Q} + \theta \mathbf{I}_2) \text{vec}(\bar{\mathbf{V}}_i^t - \mathbf{I}_2))}_{\triangleq \Upsilon_{i1}} + \underbrace{\mathbf{U}_{\mathbf{B}^{(i)}}^\top \tilde{\mathbf{G}}^t (\mathbf{X}^t)^\top \mathbf{U}_{\mathbf{B}^{(i)}}}_{\triangleq \Upsilon_{i2}}. \quad (48)$$

665 Using the relation that $\tilde{\mathbf{G}}_i = \Upsilon_{i1} + \Upsilon_{i2}$, we obtain the following results from the above equality:

$$\begin{aligned} \mathbf{0}_{2,2} &= (\Upsilon_{i1} + \Upsilon_{i2}) - \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Upsilon_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} \\ &\stackrel{\textcircled{1}}{=} \Upsilon_{i1} + \Omega_{i1} + \Omega_{i2} - \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Omega_{i1} + \Omega_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} \\ &\Rightarrow \Omega_{i1} = \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Omega_{i1} + \Omega_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} - \Upsilon_{i1} - \Omega_{i2}, \end{aligned} \quad (49)$$

666 where step $\textcircled{1}$ uses $\Upsilon_{i2} = \Omega_{i1} + \Omega_{i2}$. Then we derive the following results:

$$\begin{aligned} &\mathbb{E}_{\iota^{t+1}}[\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}\mathcal{T}(\mathbf{V}; = \mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1}))] = \mathbb{E}_{\iota^{t+1}}[\|\nabla_{\mathcal{J}}\mathcal{T}(\mathbf{V}; = \mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1})\|_F] \\ &\stackrel{\textcircled{1}}{=} \mathbb{E}_{\iota^{t+1}}[\|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}^{(i)}}^\top (\tilde{\mathbf{G}}^{t+1}[\mathbf{X}^{t+1}]^\top - \mathbf{J}\mathbf{X}^{t+1}[\tilde{\mathbf{G}}^{t+1}]^\top \mathbf{J}) \mathbf{U}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{2}}{=} \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}^{(i)}}^\top (\tilde{\mathbf{G}}^{t+1}[\mathbf{X}^{t+1}]^\top - \mathbf{J}\mathbf{X}^{t+1}[\tilde{\mathbf{G}}^{t+1}]^\top \mathbf{J}) \mathbf{U}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{3}}{=} \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i0}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{4}}{=} \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} (\Omega_{i0} - \Omega_{i1}) + \Omega_{i1} - (\mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i0}^\top \mathbf{J}_{\mathbf{B}^{(i)}} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}) - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{5}}{\leq} \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_F] + \mathbb{E}_{\iota^{t+1}}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i0}^\top \mathbf{J}_{\mathbf{B}^{(i)}} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\quad + \mathbb{E}_{\iota^{t+1}}[\|\sum_{i=1}^{n/2} \Omega_{i1} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{6}}{\leq} \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_F] + \mathbb{E}_{\iota^{t+1}}[\|\sum_{i=1}^{n/2} \Omega_{i0}^\top - \Omega_{i1}^\top\|_F] + \mathbb{E}_{\iota^{t+1}}[\|\sum_{i=1}^{n/2} \Omega_{i1} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{7}}{\leq} 2\mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_F] + \mathbb{E}_{\iota^{t+1}}[\|\sum_{i=1}^{n/2} \Omega_{i1} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{8}}{=} 2\mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_F] \\ &\quad + \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Omega_{i1} + \Omega_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} - \Upsilon_{i1} - \Omega_{i2} - \mathbf{J}_{\mathbf{B}^{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}^{(i)}}\|_F] \\ &\stackrel{\textcircled{9}}{\leq} 2\mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_F] + \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t \Upsilon_{i1}^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} - \Upsilon_{i1}\|_F] + \\ &\quad \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Omega_{i1}^\top \bar{\mathbf{V}}_i^t - \Omega_{i1}^\top\|_F] + \mathbb{E}_{\iota^t}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}^{(i)}} \bar{\mathbf{V}}_i^t \Omega_{i2}^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}^{(i)}} - \Omega_{i2}\|_F] \end{aligned} \quad (50)$$

667 where step $\textcircled{1}$ uses Equality (43); step $\textcircled{2}$ uses the fact that both the working set \mathbf{B}^t and \mathbf{B}^{t+1} are selected
668 randomly and uniformly; step $\textcircled{3}$ uses the definition of Ω_{i0} in (44); step $\textcircled{4}$ uses $-\Omega_{i1} + \Omega_{i1} = \mathbf{0}$
669 and $-\Omega_{i1}^\top + \Omega_{i1}^\top = \mathbf{0}$; step $\textcircled{5}$ uses the norm inequality; step $\textcircled{6}$ uses the norm inequality; step $\textcircled{7}$ uses
670 the norm inequality; step $\textcircled{8}$ uses Equality (49); step $\textcircled{9}$ uses the norm inequality. We now establish
671 individual bounds for each term for Inequality (50).

672 For the first term $2\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbb{F}}]$ in (50):

$$\begin{aligned}
2\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbb{F}}] &= 2\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbb{B}(i)}^{\top} [\tilde{\mathbf{G}}^t] [\mathbf{X}^t - \mathbf{X}^t]^{\top} \mathbf{U}_{\mathbb{B}(i)}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{1}}{=} 2\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} [\tilde{\mathbf{G}}^t] [\mathbf{U}_{\mathbb{B}(i)} (\bar{\mathbf{V}}_i^t - \mathbf{I}_2) \mathbf{U}_{\mathbb{B}(i)}^{\top} \mathbf{X}^t]^{\top}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{2}}{\leq} 2\bar{\mathbf{X}}\bar{\mathbf{G}}\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{3}}{\leq} 2\bar{\mathbf{X}}\bar{\mathbf{G}}\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \tag{51}
\end{aligned}$$

673 where step $\textcircled{1}$ uses $[\mathbf{X}^t - \mathbf{X}^t]_{\mathbb{B}(i)} = \mathbf{U}_{\mathbb{B}(i)} (\bar{\mathbf{V}}_i^t - \mathbf{I}_2) \mathbf{U}_{\mathbb{B}(i)}^{\top} \mathbf{X}^t$; step $\textcircled{2}$ uses the inequality $\|\mathbf{X}\mathbf{Y}\|_{\mathbb{F}} \leq$
674 $\|\mathbf{X}\|_{\mathbb{F}}\|\mathbf{Y}\|_{\mathbb{F}}$ for all \mathbf{X} and \mathbf{Y} repeatedly and the fact that $\forall t, \|\tilde{\mathbf{G}}^t\|_{\mathbb{F}} \leq \bar{\mathbf{G}}$ and $\forall t, \|\mathbf{X}^t\|_{\mathbb{F}} \leq \bar{\mathbf{X}}$; step $\textcircled{3}$
675 uses the norm inequality.

676 For the second term $\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbb{B}(i)} \bar{\mathbf{V}}_i^t \Upsilon_{i1}^{\top} \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbb{B}(i)} - \Upsilon_{i1}\|_{\mathbb{F}}]$ in (50):

$$\begin{aligned}
&\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbb{B}(i)} \bar{\mathbf{V}}_i^t \Upsilon_{i1}^{\top} \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbb{B}(i)} - \Upsilon_{i1}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Upsilon_{i1}^{\top} \bar{\mathbf{V}}_i^t\|_{\mathbb{F}}] + \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \Upsilon_{i1}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{2}}{\leq} (1 + \bar{\mathbf{V}}^2) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \Upsilon_{i1}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{3}}{\leq} (1 + \bar{\mathbf{V}}^2) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \text{mat}((\mathbf{Q} + \theta \mathbf{I}_2) \text{vec}(\bar{\mathbf{V}}_i^t - \mathbf{I}_2))\|_{\mathbb{F}}] \\
&\leq (1 + \bar{\mathbf{V}}^2) \|\mathbf{Q} + \theta \mathbf{I}_2\|_{\mathbb{F}} \cdot \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{4}}{\leq} (1 + \bar{\mathbf{V}}^2) (L_f + \theta) \cdot \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \tag{52}
\end{aligned}$$

677 where step $\textcircled{1}$ uses the triangle inequality; step $\textcircled{2}$ uses the inequality $\|\mathbf{X}\mathbf{Y}\|_{\mathbb{F}} \leq \|\mathbf{X}\|_{\mathbb{F}}\|\mathbf{Y}\|_{\mathbb{F}}$ for all
678 \mathbf{X} and \mathbf{Y} and $\forall t, \|\bar{\mathbf{V}}_i^t\|_{\mathbb{F}} \leq \bar{\mathbf{V}}$; step $\textcircled{3}$ uses the definition of Υ_{i1} ; step $\textcircled{4}$ uses the choice of $\mathbf{Q} \preceq L_f \mathbf{I}$
679 and the norm inequality.

680 For the third term $\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Omega_{i1}^{\top} \bar{\mathbf{V}}_i^t - \Omega_{i1}^{\top}\|_{\mathbb{F}}]$ in (50), we have:

$$\begin{aligned}
&\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Omega_{i1}^{\top} \bar{\mathbf{V}}_i^t - \Omega_{i1}^{\top}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{1}}{=} \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Omega_{i1}^{\top} (\bar{\mathbf{V}}_i^t - \mathbf{I}_2) + (\bar{\mathbf{V}}_i^t - \mathbf{I}_2) \Omega_{i1}^{\top}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{2}}{\leq} (1 + \bar{\mathbf{V}}) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \|\Omega_{i1}\|_{\mathbb{F}} \cdot \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{3}}{\leq} (\bar{\mathbf{X}} + \bar{\mathbf{V}}\bar{\mathbf{X}}) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \|\tilde{\mathbf{G}}^t\|_{\mathbb{F}} \cdot \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{4}}{\leq} (\bar{\mathbf{X}} + \bar{\mathbf{V}}\bar{\mathbf{X}}) \bar{\mathbf{G}} \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \tag{53}
\end{aligned}$$

681 where step $\textcircled{1}$ uses the fact that $-\bar{\mathbf{V}}_i^t \Omega_{i1}^{\top} \mathbf{I}_2 + \bar{\mathbf{V}}_i^t \Omega_{i1}^{\top} = \mathbf{0}$; step $\textcircled{2}$ uses the norm inequality and
682 $\forall t, \|\bar{\mathbf{V}}_i^t\|_{\mathbb{F}} \leq \bar{\mathbf{V}}$; step $\textcircled{3}$ uses the fact that $\|\Omega_{i1}\|_{\mathbb{F}} = \|\mathbf{U}_{\mathbb{B}(i)}^{\top} \tilde{\mathbf{G}}^t [\mathbf{X}^t]^{\top} \mathbf{U}_{\mathbb{B}(i)}\|_{\mathbb{F}} \leq \bar{\mathbf{X}} \|\tilde{\mathbf{G}}^t\|_{\mathbb{F}}$, $\forall i$ which
683 can be derived using the norm inequality; step $\textcircled{4}$ uses the fact that $\forall \mathbf{X}, \|\tilde{\mathbf{G}}^t\|_{\mathbb{F}} \leq \bar{\mathbf{G}}$.

684 For the fourth term $\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbb{B}(i)} \bar{\mathbf{V}}_i^t \Omega_{i2}^{\top} \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbb{B}(i)} - \Omega_{i2}\|_{\mathbb{F}}]$ in (50), we have:

$$\begin{aligned}
&\mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbb{B}(i)} \bar{\mathbf{V}}_i^t \Omega_{i2}^{\top} \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbb{B}(i)} - \Omega_{i2}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Omega_{i2}^{\top} \bar{\mathbf{V}}_i^t\|_{\mathbb{F}}] + \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \Omega_{i2}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{2}}{\leq} (1 + \bar{\mathbf{V}}^2) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \Omega_{i2}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{3}}{=} (1 + \bar{\mathbf{V}}^2) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbb{B}(i)}^{\top} [\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^t] [\mathbf{X}^t]^{\top} \mathbf{U}_{\mathbb{B}(i)}\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{4}}{\leq} \frac{n}{2} (\bar{\mathbf{X}} + \bar{\mathbf{V}}^2 \bar{\mathbf{X}}) \mathbb{E}_{\iota^t} [\|\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^t\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{5}}{\leq} \frac{n}{2} (\bar{\mathbf{X}} + \bar{\mathbf{V}}^2 \bar{\mathbf{X}}) (p \mathbb{E}_{\iota^t} [\sqrt{u^t}] + L_f \mathbb{E}_{\iota^t} [\|\mathbf{X}^t - \mathbf{X}^t\|_{\mathbb{F}}]) \\
&\stackrel{\textcircled{6}}{\leq} \frac{np}{2} (\bar{\mathbf{X}} + \bar{\mathbf{V}}^2 \bar{\mathbf{X}}) \mathbb{E}_{\iota^t} [\sqrt{u^t}] + \frac{nL_f}{2} (\bar{\mathbf{X}}^2 + \bar{\mathbf{V}}^2 \bar{\mathbf{X}}^2) \mathbb{E}_{\iota^t} [\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}] \tag{54}
\end{aligned}$$

685 where step ① uses the triangle inequality; step ② uses the norm inequality and $\forall t, \|\mathbf{V}^t\|_F \leq \bar{V}$; step ③
686 uses the definition of $\forall i, \Omega_{i2} = \mathbf{U}_{\mathbf{B}(i)}^\top [\tilde{\mathbf{G}}^t - \hat{\mathbf{G}}^t] [\mathbf{X}^t]^\top \mathbf{U}_{\mathbf{B}(i)}$ in (46); step ④ uses the norm inequality
687 and $\forall t, \|\mathbf{X}^t\|_F \leq \bar{X}$; step ⑤ uses Lemma E.1; step ⑥ uses Part (b) in Lemma 2.5 and $\forall t, \|\mathbf{X}^t\|_F \leq \bar{X}$.
688 In view of (51), (52), (53), (54), and (50), we have:

$$\begin{aligned} & \mathbb{E}_{\iota^{t+1}} [\|\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1})\|_F] \\ & \leq \frac{n\rho}{2} (\bar{X} + \bar{V}^2 \bar{X}) \mathbb{E}_{\iota^t} [\sqrt{u^t}] + (c_1 + c_2 + c_3 + c_4) \cdot \mathbb{E}_{\iota^t} [\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_F] \\ & = \frac{n\rho}{2} (\bar{X} + \bar{V}^2 \bar{X}) \mathbb{E}_{\iota^t} [\sqrt{u^t}] + \phi \mathbb{E}_{\iota^t} [\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_F] \end{aligned}$$

689 where $c_1 = 2\bar{X}\bar{G}$, $c_2 = (1 + \bar{V}^2)(L_f + \theta)$, $c_3 = (\bar{X} + \bar{V}\bar{X})\bar{G}$, and $c_4 = \frac{n}{2}(\bar{X}^2 + \bar{V}^2\bar{X}^2)L_f$. \square

690 **Lemma E.3.** *We have the following results: $\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}} f(\mathbf{X}^t)) \leq \gamma \cdot \|\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_F +$
691 $2\bar{X}^2 \sqrt{\mathbb{E}_{\iota^t}[u^t]}$ with $\gamma \triangleq \bar{X} \sqrt{C_n^2}$.*

692 *Proof.* We have the following inequalities:

$$\begin{aligned} \|\nabla_{\mathcal{J}} f(\mathbf{X}^t)\|_F & \stackrel{\textcircled{1}}{=} \|\nabla f(\mathbf{X}^t) - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{X}^t \mathbf{J}\|_F \\ & \stackrel{\textcircled{2}}{=} \|\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top \mathbf{J}\mathbf{X}^t \mathbf{J} - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{J}\mathbf{J}\mathbf{X}^t \mathbf{J}\|_F \\ & \stackrel{\textcircled{3}}{\leq} \|\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{J}\|_F \|\mathbf{J}\mathbf{X}^t \mathbf{J}\|_F \\ & \stackrel{\textcircled{4}}{\leq} \bar{X} \|\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{J}\|_F \end{aligned}$$

693 where step ① uses the definition of $\nabla_{\mathcal{J}} f(\mathbf{X}^t)$; step ② uses $\mathbf{J}\mathbf{J} = \mathbf{I}$ and $\mathbf{X}^\top \mathbf{J}\mathbf{X} = \mathbf{J} \Rightarrow \mathbf{X}^\top \mathbf{J}\mathbf{X}\mathbf{J} =$
694 $\mathbf{J}\mathbf{J} = \mathbf{I}$; step ③ uses the norm inequality and; step ④ uses $\forall t, \|\mathbf{X}^t\|_F \leq \bar{X}$.

695 We Consider $\|\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{J}\|_F$:

$$\begin{aligned} & \|\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{J}\|_F \\ & \stackrel{\textcircled{1}}{\leq} \|\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\tilde{\mathbf{G}}^t)^\top \mathbf{J}\|_F + \|(\nabla f(\mathbf{X}^t) - \tilde{\mathbf{G}}^t)(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t) - \tilde{\mathbf{G}}^t)^\top \mathbf{J}\|_F \\ & \stackrel{\textcircled{2}}{\leq} \|\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\tilde{\mathbf{G}}^t)^\top \mathbf{J}\|_F + \|\nabla f(\mathbf{X}^t) - \tilde{\mathbf{G}}^t\|_F \cdot \|\mathbf{X}^t\|_F + \|\mathbf{X}^t\|_F \cdot \|\nabla f(\mathbf{X}^t) - \tilde{\mathbf{G}}^t\|_F \\ & \stackrel{\textcircled{3}}{\leq} \|\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\tilde{\mathbf{G}}^t)^\top \mathbf{J}\|_F + 2\bar{X} \sqrt{\mathbb{E}_{\iota^t}[u^t]} \end{aligned}$$

696 where step ① uses $\forall \mathbf{A}, \mathbf{B}, \|\mathbf{A}\|_F - \|\mathbf{B}\|_F \leq \|\mathbf{A} - \mathbf{B}\|_F$; step ② uses the norm inequality; step ③
697 uses $\forall t, \|\mathbf{X}^t\|_F \leq \bar{X}$. Thus,

$$\begin{aligned} \|\nabla_{\mathcal{J}} f(\mathbf{X}^t)\|_F & \leq \bar{X} \|\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\tilde{\mathbf{G}}^t)^\top \mathbf{J}\|_F + 2\bar{X}^2 \sqrt{\mathbb{E}_{\iota^t}[u^t]} \\ & \stackrel{\textcircled{1}}{\leq} \bar{X} \sqrt{C_n^2} \cdot \|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}(i)}^\top [\tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\tilde{\mathbf{G}}^t)^\top \mathbf{J}\mathbf{U}_{\mathbf{B}(i)}]\|_F + 2\bar{X}^2 \sqrt{\mathbb{E}_{\iota^t}[u^t]} \\ & \stackrel{\textcircled{2}}{=} \bar{X} \sqrt{C_n^2} \cdot \|\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_F + 2\bar{X}^2 \sqrt{\mathbb{E}_{\iota^t}[u^t]} \end{aligned}$$

698 where step ① uses Lemma A.1 with $\mathbf{W} = \tilde{\mathbf{G}}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\tilde{\mathbf{G}}^t)^\top \mathbf{J}$ and $k = 2$; step ② uses the
699 definition of $\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})$. \square

700 We now present the following useful lemma.

Lemma E.4. *We define $\mathbf{T}_{\mathbf{X}} \mathcal{J} \triangleq \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \mathcal{A}_{\mathbf{X}}(\mathbf{Y}) = \mathbf{0}\}$ and $\mathcal{A}_{\mathbf{X}}(\mathbf{Y}) \triangleq \mathbf{X}^\top \mathbf{J}\mathbf{Y} + \mathbf{Y}^\top \mathbf{J}\mathbf{X}$.
For any $\mathbf{G} \in \mathbb{R}^{n \times n}$ and $\mathbf{X}^\top \mathbf{J}\mathbf{X} = \mathbf{J}$, the unique minimizer of the following optimization problem:*

$$\bar{\mathbf{Y}} = \arg \min_{\mathbf{Y} \in \mathbf{T}_{\mathbf{X}} \mathcal{J}} h(\mathbf{Y}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{G}\|_F^2,$$

701 *satisfy $h(\bar{\mathbf{Y}}) \leq h(\mathbf{G} - \mathbf{J}\mathbf{X}\mathbf{G}^\top \mathbf{X}\mathbf{J})$.*

702 *Proof.* We note that $\bar{\mathbf{Y}} = \arg \min_{\mathbf{Y} \in \mathbf{T}_{\mathbf{X}} \mathcal{J}} \frac{1}{2} \|\mathbf{Y} - \mathbf{G}\|_F^2 = \arg \min_{\mathbf{Y}} \frac{1}{2} \|\mathbf{Y} - \mathbf{G}\|_F^2,$
703 s.t. $\mathbf{X}^\top \mathbf{J}\mathbf{Y} + \mathbf{Y}^\top \mathbf{J}\mathbf{X} = \mathbf{0}$. Introducing a multiplier $\boldsymbol{\Lambda} \in \mathbb{R}^{n \times n}$ for the linear con-
704 straints $\mathbf{X}^\top \mathbf{J}\mathbf{Y} + \mathbf{Y}^\top \mathbf{J}\mathbf{X} = \mathbf{0}$, we have following Lagrangian function: $\tilde{\mathcal{L}}(\mathbf{Y}; \boldsymbol{\Lambda}) =$

705 $\frac{1}{2}\|\mathbf{Y} - \mathbf{G}\|_{\mathbb{F}}^2 + \langle \mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X}, \mathbf{\Lambda} \rangle$. We naturally derive the following first-order optimality
706 condition: $\mathbf{Y} - \mathbf{G} + \mathbf{J} \mathbf{X} \mathbf{\Lambda} = \mathbf{0}$, $\mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} = \mathbf{0}$. Incorporating the term $\mathbf{Y} = \mathbf{G} - \mathbf{J} \mathbf{X} \mathbf{\Lambda}$
707 into $\mathbf{X}^{\top} \mathbf{J} \mathbf{Y} + \mathbf{Y}^{\top} \mathbf{J} \mathbf{X} = \mathbf{0}$, we obtain:

$$\mathbf{X}^{\top} \mathbf{X} \mathbf{\Lambda} + \mathbf{\Lambda}^{\top} \mathbf{X}^{\top} \mathbf{X} = \mathbf{G}^{\top} \mathbf{J} \mathbf{X} + \mathbf{X}^{\top} \mathbf{J} \mathbf{G} \quad (55)$$

708 Any $\mathbf{\Lambda}$ satisfying formula (55) is a feasible point, so we can easily find :

$$\begin{aligned} & \mathbf{X}^{\top} \mathbf{X} \mathbf{\Lambda} = \mathbf{X}^{\top} \mathbf{J} \mathbf{G} \\ \stackrel{\textcircled{1}}{\Rightarrow} & \mathbf{X} \mathbf{\Lambda} = \mathbf{J} \mathbf{G} \\ \stackrel{\textcircled{2}}{\Rightarrow} & \mathbf{X}^{\top} \mathbf{J} \mathbf{X} \mathbf{\Lambda} = \mathbf{X}^{\top} \mathbf{J} \mathbf{J} \mathbf{G} \\ \stackrel{\textcircled{3}}{\Rightarrow} & \mathbf{J} \mathbf{\Lambda} = \mathbf{X}^{\top} \mathbf{G} \\ \stackrel{\textcircled{4}}{\Rightarrow} & \mathbf{\Lambda} = \mathbf{J} \mathbf{X}^{\top} \mathbf{G} \\ \stackrel{\textcircled{5}}{\Rightarrow} & \mathbf{\Lambda} = \mathbf{G}^{\top} \mathbf{X} \mathbf{J} \end{aligned} \quad (56)$$

709 where step $\textcircled{1}$ uses the fact that any matrix \mathbf{X} satisfying the J-orthogonality constraint has a determinant
710 of 1 or -1, thus $\text{inv}(\mathbf{X})$ exists; step $\textcircled{2}$ multiply both sides of the equation by $\mathbf{X} \mathbf{J}$; step $\textcircled{3}$ uses
711 $\mathbf{X}^{\top} \mathbf{J} \mathbf{X} = \mathbf{J}$ and $\mathbf{J} \mathbf{J} = \mathbf{I}$; step $\textcircled{4}$ multiply both sides of the equation by \mathbf{J} and uses $\mathbf{J} \mathbf{J} = \mathbf{I}$; step $\textcircled{5}$
712 uses the fact that $\mathbf{\Lambda}$ is a symmetric matrix.

713 Therefore, a feasible solution \mathbf{Y} can be computed as $\mathbf{Y} = \mathbf{G} - \mathbf{J} \mathbf{X} \mathbf{\Lambda} = \mathbf{G} - \mathbf{J} \mathbf{X} \mathbf{G}^{\top} \mathbf{X} \mathbf{J}$. Since $\bar{\mathbf{Y}}$
714 is the optimal solution, there must be $h(\bar{\mathbf{Y}}) \leq h(\mathbf{G} - \mathbf{J} \mathbf{X} \mathbf{G}^{\top} \mathbf{X} \mathbf{J})$. \square

715 We now present the proof of this lemma.

716 **Lemma E.5.** For any $\mathbf{X} \in \mathbb{R}^{n \times n}$, it holds that $\text{dist}(\mathbf{0}, \nabla f^{\circ}(\mathbf{X})) \leq \text{dist}(\mathbf{0}, \nabla_{\mathcal{J}} f(\mathbf{X}))$.

717 *Proof.* For the purpose of analysis, we define the nearest J orthogonal matrix to an arbitrary matrix
718 $\mathbf{Y} \in \mathbb{R}^{n \times n}$ is given by $\mathcal{P}_{\mathcal{J}}(\mathbf{X})$. Similarly, we have $\mathcal{P}_{\mathbf{T}_{\mathbf{X}} \mathcal{J}}(\nabla f(\mathbf{X}))$ for projecting gradient $\nabla f(\mathbf{X})$
719 into space $\mathbf{T}_{\mathbf{X}} \mathcal{J}$.

720 We recall that the following first-order optimality conditions are equivalent for all $\mathbf{X} \in \mathbb{R}^{n \times n}$:

$$(\mathbf{0} \in \nabla f^{\circ}(\mathbf{X})) \Leftrightarrow (\mathbf{0} \in \mathcal{P}_{\mathbf{T}_{\mathbf{X}} \mathcal{J}}(\nabla f(\mathbf{X}))). \quad (57)$$

721 Therefore, we derive the following results:

$$\text{dist}(\mathbf{0}, \nabla f^{\circ}(\mathbf{X})) = \inf_{\mathbf{Y} \in \nabla f^{\circ}(\mathbf{X})} \|\mathbf{Y}\|_{\mathbb{F}} \quad (58)$$

$$= \inf_{\mathbf{Y} \in \mathcal{P}_{(\mathbf{T}_{\mathbf{X}} \mathcal{J})}(\nabla f(\mathbf{X}))} \|\mathbf{Y}\|_{\mathbb{F}} \quad (59)$$

722 We let $\mathbf{G} \in \nabla f(\mathbf{X})$ and obtain the following results from the above equality:

$$\text{dist}(\mathbf{0}, \nabla f^{\circ}(\mathbf{X})) \stackrel{\textcircled{1}}{\leq} \|\mathbf{G} - \mathbf{J} \mathbf{X} \mathbf{G}^{\top} \mathbf{X} \mathbf{J}\|_{\mathbb{F}}, \quad (60)$$

$$\stackrel{\textcircled{2}}{=} \|\nabla_{\mathcal{J}} f(\mathbf{X})\|_{\mathbb{F}} \triangleq \text{dist}(\mathbf{0}, \nabla_{\mathcal{J}} f(\mathbf{X})). \quad (61)$$

723 where step $\textcircled{1}$ uses Lemma E.4; step $\textcircled{2}$ uses $\nabla_{\mathcal{J}} f(\mathbf{X}) = \mathbf{G} - \mathbf{J} \mathbf{X} \mathbf{G}^{\top} \mathbf{X} \mathbf{J}$ with $\mathbf{G} \in \nabla f(\mathbf{X})$. \square

724 First of all, since $f^{\circ}(\mathbf{X}) \triangleq f(\mathbf{X}) + \mathcal{I}_{\mathcal{J}}(\mathbf{X})$ is a KL function, we have from Proposition 4.8 that:

$$\begin{aligned} \frac{1}{\varphi'(f^{\circ}(\mathbf{X}') - f^{\circ}(\mathbf{X}))} & \leq \text{dist}(\mathbf{0}, \nabla f^{\circ}(\mathbf{X}')) \\ & \stackrel{\textcircled{1}}{=} \|\nabla_{\mathcal{J}} f(\mathbf{X}')\|_{\mathbb{F}}, \end{aligned} \quad (62)$$

725 where step $\textcircled{1}$ uses Lemma E.5. Here, $\varphi(\cdot)$ is some certain concave desingularization function. Since
726 $\varphi(\cdot)$ is concave, we have:

$$\forall \Delta \in \mathbb{R}, \Delta^+ \in \mathbb{R}, \varphi(\Delta^+) + (\Delta - \Delta^+) \varphi'(\Delta) \leq \varphi(\Delta). \quad (63)$$

727 Applying the inequality above with $\Delta = f(\mathbf{X}^t) - f(\bar{\mathbf{X}})$ and $\Delta^+ = f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}})$, we have:

$$\begin{aligned} & (f(\mathbf{X}^t) - f(\mathbf{X}^{t+1})) \varphi'(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) \\ & \leq \varphi(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}})) \triangleq \mathcal{E}^t. \end{aligned} \quad (64)$$

728 With the sufficient descent condition as shown in Theorem 4.7, we derive the following inequalities:

$$\begin{aligned} & \mathbb{E}_{\iota^t} [\frac{\theta}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \\ \leq & \mathbb{E}_{\iota^t} [f(\mathbf{X}^t) - f(\mathbf{X}^{t+1})] + \frac{1}{2} \mathbb{E}_{\iota^t} [\|\mathbf{X}^t\|_{\mathbb{F}}^2] \mathbb{E}_{\iota^t} [\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \quad (65) \end{aligned}$$

$$\begin{aligned} \stackrel{\textcircled{1}}{\Rightarrow} & \mathbb{E}_{\iota^t} [\frac{\theta - \bar{\mathbf{X}}^2}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \leq \mathbb{E}_{\iota^t} [f(\mathbf{X}^t) - f(\mathbf{X}^{t+1})] + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \quad (66) \\ & \quad \quad \quad (67) \end{aligned}$$

729 where step ① uses $\forall t, \|\mathbf{X}^t\|_{\mathbb{F}} \leq \bar{\mathbf{X}}$.

$$\begin{aligned} & \mathbb{E}_{\iota^t} [\frac{\theta - \bar{\mathbf{X}}^2}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2] \\ \stackrel{\textcircled{1}}{\leq} & \mathbb{E}_{\iota^t} [\frac{\mathcal{E}^t}{\varphi'(f(\mathbf{X}^t) - f(\mathbf{X}))}] + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \\ \stackrel{\textcircled{2}}{\leq} & \mathbb{E}_{\iota^t} [\mathcal{E}^t \|\nabla_{\mathcal{J}} f(\mathbf{X}^t)\|_{\mathbb{F}}] + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \\ \stackrel{\textcircled{3}}{\leq} & \mathbb{E}_{\iota^t} [\mathcal{E}^t \gamma \|\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_{\mathbb{F}}] + 2\mathcal{E}^t \bar{\mathbf{X}}^2 \sqrt{\mathbb{E}_{\iota^t} [u^t]} + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \\ \stackrel{\textcircled{4}}{\leq} & \mathbb{E}_{\iota^t} [\mathcal{E}^t \gamma \phi \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}] + \mathcal{E}^t \gamma \frac{np}{2} (\bar{\mathbf{X}} + \bar{\mathbf{V}}^2 \bar{\mathbf{X}}) \sqrt{\mathbb{E}_{\iota^t} [u^t]} \\ & + 2\mathcal{E}^t \bar{\mathbf{X}}^2 \sqrt{\mathbb{E}_{\iota^t} [u^t]} + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \\ \stackrel{\textcircled{5}}{\leq} & \mathbb{E}_{\iota^t} [\mathcal{E}^t \gamma \phi \sqrt{\frac{n}{2}} \sqrt{\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2}] \\ & + \mathcal{E}^t (2\bar{\mathbf{X}}^2 + \gamma \frac{np}{2} \bar{\mathbf{X}} + \gamma \frac{np}{2} \bar{\mathbf{V}}^2 \bar{\mathbf{X}}) \sqrt{\mathbb{E}_{\iota^t} [u^t]} + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \\ \stackrel{\textcircled{6}}{\leq} & \mathbb{E}_{\iota^t} [\frac{n\mathcal{E}^t \gamma^2 \phi^2}{4\theta'} + \frac{\theta'}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2 + \frac{\bar{\theta} \mathbb{E}_{\iota^t} [u^t]}{2}] \\ & + \frac{\mathcal{E}^t (2\bar{\mathbf{X}}^2 + \gamma \frac{np}{2} \bar{\mathbf{X}} + \gamma \frac{np}{2} \bar{\mathbf{V}}^2 \bar{\mathbf{X}})^2}{2\theta} + \frac{1}{2} \mathbb{E}_{\iota^t} [u^t] \\ \stackrel{\textcircled{7}}{=} & \mathbb{E}_{\iota^t} [\mathcal{E}^t \mathfrak{Q}^2 + \frac{\theta'}{2} \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2] + \frac{\bar{\theta} + 1}{2} \mathbb{E}_{\iota^t} [u^t] \quad (68) \end{aligned}$$

730 where step ① uses the sufficient descent condition as shown in Theorem 4.7; step

731 ② uses Inequality (64) and (62) with $\mathbf{X}' = \mathbf{X}^t$ and $\mathbf{X} = \bar{\mathbf{X}}$; step ③ uses lemma

732 E.3 ; step ④ uses Lemma E.2 ; step ⑤ uses $\forall x_i \in \mathbb{R}, \frac{x_1 + \dots + x_n}{n} \leq \sqrt{\frac{x_1^2 + \dots + x_n^2}{n}}$

733 ; step ⑥ applies the inequality that $\forall \theta' > 0, a, b, ab \leq \frac{\theta' a^2}{2} + \frac{b^2}{2\theta'}$ with

734 $a = \sqrt{\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2}, b = \mathcal{E}^t \gamma \phi \sqrt{\frac{n}{2}}; a = \sqrt{\mathbb{E}_{\iota^t} [u^t]}, b = \mathcal{E}^t (2\bar{\mathbf{X}}^2 + \gamma \frac{np}{2} \bar{\mathbf{X}} + \gamma \frac{np}{2} \bar{\mathbf{V}}^2 \bar{\mathbf{X}});$

735 step ⑦ denote $\mathfrak{Q}^2 \triangleq \frac{(2\bar{\mathbf{X}}^2 + \gamma \frac{np}{2} \bar{\mathbf{X}} + \gamma \frac{np}{2} \bar{\mathbf{V}}^2 \bar{\mathbf{X}})^2}{2\theta} + \frac{n\gamma^2 \phi^2}{4\theta'}$. To simplify the formula, we define

736 $\aleph^t = \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2$.

737 Multiplying both sides by 2 and taking the square root of both sides, we have:

$$\begin{aligned} \mathbb{E}_{\iota^t} [\sqrt{\theta - \bar{\mathbf{X}}^2} \sqrt{\aleph^t}] & \leq \sqrt{\mathbb{E}_{\iota^t} [\mathcal{E}^t \mathfrak{Q}^2 + \theta' \aleph^{t-1}] + (\bar{\theta} + 1) \mathbb{E}_{\iota^t} [u^t]} \\ & \leq \sqrt{\mathbb{E}_{\iota^t} [\mathcal{E}^t \mathfrak{Q}^2]} + \mathbb{E}_{\iota^{t-1}} [\sqrt{\theta' \aleph^{t-1}}] + \sqrt{(\bar{\theta} + 1) \mathbb{E}_{\iota^t} [u^t]} \\ & \leq \mathcal{E}^t \mathfrak{Q} + \sqrt{\theta'} \mathbb{E}_{\iota^{t-1}} [\sqrt{\aleph^{t-1}}] + \sqrt{(\bar{\theta} + 1) \mathbb{E}_{\iota^t} [u^t]} \quad (69) \end{aligned}$$

738 To recursively eliminate term $\sqrt{(\bar{\theta} + 1) \mathbb{E}_{\iota^t} [u^t]}$, we take the root of both sides of the Inequality in

739 Lemma 4.5:

$$\begin{aligned} \sqrt{\mathbb{E}_{\iota^t} [u^t]} & \leq \sqrt{\frac{p(N-b)}{b(N-1)} \sigma^2} + \sqrt{(1-p) \mathbb{E}_{\iota^{t-1}} [u^{t-1}]} + \sqrt{\frac{L_j^2 \bar{\mathbf{X}}^2 (1-p)}{b'}} \mathbb{E}_{\iota^{t-1}} [\aleph^{t-1}] \\ & \leq \sqrt{\frac{p(N-b)}{b(N-1)} \sigma^2} + \sqrt{(1-p)} \sqrt{\mathbb{E}_{\iota^{t-1}} [u^{t-1}]} + \sqrt{\frac{L_j^2 \bar{\mathbf{X}}^2 (1-p)}{b'}} \sqrt{\mathbb{E}_{\iota^{t-1}} [\aleph^{t-1}]} \quad (70) \end{aligned}$$

740 Adding Inequality $\frac{\sqrt{\theta+1}}{1-\sqrt{1-p}} \times (70)$ to (69)

$$\begin{aligned} \mathbb{E}_{\iota^t}[\sqrt{\theta - \bar{X}^2} \sqrt{\aleph^t}] &\leq \mathcal{E}^t \mathfrak{A} + (\sqrt{\theta'} + \sqrt{\frac{L_f^2 \bar{X}^2 (1-p)}{b'}} \frac{\sqrt{\theta+1}}{1-\sqrt{1-p}}) \mathbb{E}_{\iota^{t-1}}[\sqrt{\aleph^{t-1}}] + \\ &\quad \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} (\sqrt{\mathbb{E}_{\iota^{t-1}}[u^{t-1}]} - \sqrt{\mathbb{E}_{\iota^t}[u^t]}) + \frac{\sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 \end{aligned} \quad (71)$$

741 With the choice $\sqrt{\theta'} = \frac{\sqrt{\theta - \bar{X}^2}}{2} - \sqrt{\frac{L_f^2 \bar{X}^2 (1-p)}{b'}} \frac{\sqrt{\theta+1}}{1-\sqrt{1-p}}$, we have:

$$\begin{aligned} \mathbb{E}_{\iota^t}[\sqrt{\theta - \bar{X}^2} \sqrt{\aleph^t}] &\leq \mathcal{E}^t \mathfrak{A} + (\frac{\sqrt{\theta - \bar{X}^2}}{2}) \mathbb{E}_{\iota^{t-1}}[\sqrt{\aleph^{t-1}}] + \\ &\quad \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} (\sqrt{\mathbb{E}_{\iota^{t-1}}[u^{t-1}]} - \sqrt{\mathbb{E}_{\iota^t}[u^t]}) + \frac{\sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 \end{aligned} \quad (72)$$

742 Rearranging terms, we have:

$$\begin{aligned} &\mathbb{E}_{\iota^t}[\sqrt{\theta - \bar{X}^2} \sqrt{\aleph^t}] - \mathbb{E}_{\iota^{t-1}}[\frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\aleph^{t-1}}] \\ &\leq \mathcal{E}^t \mathfrak{A} + \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} (\sqrt{\mathbb{E}_{\iota^{t-1}}[u^{t-1}]} - \sqrt{\mathbb{E}_{\iota^t}[u^t]}) + \frac{\sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 \end{aligned} \quad (73)$$

743 Summing the inequality above over $t = 1, 2, \dots, T$, we have:

$$\begin{aligned} &\mathbb{E}_{\iota^T}[\sqrt{\theta - \bar{X}^2} \sqrt{\aleph^T}] + \mathbb{E}_{\iota^{T-1}}[\frac{\sqrt{\theta - \bar{X}^2}}{2} \sum_{t=1}^{T-1} \sqrt{\aleph^t}] \\ &\leq \mathfrak{A} \sum_{t=1}^T \mathcal{E}^t + \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} (\sqrt{\mathbb{E}_{\iota^0}[u^0]} - \sqrt{\mathbb{E}_{\iota^T}[u^T]}) + \frac{T \sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\aleph^0} \\ &\leq \mathfrak{A} \sum_{t=1}^T \mathcal{E}^t + \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} \sqrt{\frac{N-b}{b(N-1)}} \sigma^2 + \frac{T \sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 + \mathbb{E}_{\iota^t}[\frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\aleph^0}] \\ &\leq \mathfrak{A} \sum_{t=1}^T \mathcal{E}^t + \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} \sqrt{\frac{N-b}{b(N-1)}} \sigma^2 + \frac{T \sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2} (\bar{V} + \sqrt{2})^2} \end{aligned}$$

744 where step ① uses the fact that $\mathbb{E}_{\iota^T}[u^T] \geq 0$ and $\mathbb{E}_{\iota^0}[u^0] \leq \frac{N-b}{b(N-1)} \sigma^2$; step ② uses $\forall t, \|\mathbf{V}\|_F \leq \bar{V}$,

745 then, $\|\mathbf{V}_i - \mathbf{I}_2\|_F^2 \leq (\|\mathbf{V}_i\|_F + \|\mathbf{I}_2\|_F)^2 \leq (\bar{X} + \sqrt{2})^2$ and $\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^0 - \mathbf{I}_2\|_F^2 \leq \frac{n}{2} (\bar{V} + \sqrt{2})^2$.

746 Define $\mathfrak{C} = \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} \sqrt{\frac{N-b}{b(N-1)}} \sigma^2 + \frac{T \sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2$ and rearrange terms, we have:

$$\mathbb{E}_{\iota^t}[\frac{\theta - \bar{X}^2}{2} \sum_{t=1}^T \sqrt{\aleph^t}] \leq \mathfrak{A} \sum_{t=1}^T \mathcal{E}^t + \mathfrak{C} + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2} (\bar{V} + \sqrt{2})^2} \quad (74)$$

747 Considering $\mathfrak{A} \sum_{t=1}^T \mathcal{E}^t$, we have:

$$\begin{aligned} \mathfrak{A} \sum_{t=1}^T \mathcal{E}^t &\stackrel{\text{①}}{=} \mathfrak{A} \sum_{t=1}^T \varphi(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}})) \\ &\stackrel{\text{②}}{=} \mathfrak{A} [\varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{T+1}) - f(\bar{\mathbf{X}}))] \\ &\stackrel{\text{③}}{\leq} \mathfrak{A} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) \end{aligned} \quad (75)$$

748 where step ① uses the definition of \mathcal{E}^i in (64); step ② uses a basic recursive reduction; step ③ uses the
749 fact the desingularization function $\varphi(\cdot)$ is positive. Combining Inequality (74) and (75), we obtain :

$$\mathbb{E}_{\iota^t}[\frac{\theta - \bar{X}^2}{2} \sum_{t=1}^T \sqrt{\aleph^t}] \leq \mathfrak{A} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \mathfrak{C} + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2} (\bar{V} + \sqrt{2})^2}$$

750 Using the inequality that $\frac{\|\mathbf{X}^+ - \mathbf{X}\|_F^2}{\bar{X}^2} \leq \frac{\|\mathbf{X}^+ - \mathbf{X}\|_F^2}{\|\mathbf{X}\|_F^2} \leq \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i - \mathbf{I}_2\|_F^2$ as shown in Part (b) in Lemma
751 2.5, we have:

$$\mathbb{E}_{\iota^t}[\frac{\theta - \bar{X}^2}{2\bar{X}} \sum_{t=1}^T \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F] \leq \mathfrak{A} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \mathfrak{C} + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2} (\bar{V} + \sqrt{2})^2}$$

Since $b = N, b' = \sqrt{b}$ and $p = \frac{b'}{b+b'}$, $\mathfrak{C} = \frac{\sqrt{1-p} \sqrt{(\theta+1)}}{1-\sqrt{1-p}} \sqrt{\frac{N-b}{b(N-1)}} \sigma^2 + \frac{T \sqrt{\theta+1}}{1-\sqrt{1-p}} \sqrt{\frac{p(N-b)}{b(N-1)}} \sigma^2 = 0$,
we can get the expression for C:

$$\mathbb{E}_{\iota^t}[\sum_{j=1}^t \|\mathbf{X}^{j+1} - \mathbf{X}^j\|_F] \leq C$$

$$C \triangleq \frac{2\bar{X}}{\theta - \bar{X}^2} (\mathfrak{A}\varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2}(\bar{\mathbf{V}} + \sqrt{2})^2})$$

753 Considering that: $\sqrt{\theta'} = \frac{\sqrt{\theta - \bar{X}^2}}{2} - \sqrt{\frac{L_f^2 \bar{X}^2 (1-p)}{b'}} \frac{\sqrt{\theta+1}}{1-\sqrt{1-p}} = \frac{\sqrt{\theta - \bar{X}^2}}{2} - \sqrt{L_f^2 \bar{X}^2 (1+\bar{\theta})} ((1+N^{\frac{1}{2}})^{\frac{1}{2}} +$
 754 $N^{\frac{1}{4}}) = \mathcal{O}(N^{\frac{1}{4}})$, we have: $\mathfrak{A} = \sqrt{\frac{(2\bar{X}^2 + \gamma \frac{np}{2} \bar{X} + \gamma \frac{np}{2} \bar{\mathbf{V}}^2 \bar{X})^2}{2\theta} + \frac{n\gamma^2 \phi^2}{4\theta'}} = \mathcal{O}(\frac{1}{N^{1/4}})$. Finally, we have
 755 $C = \mathcal{O}(\frac{\varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}}))}{N^{1/4}})$ \square

756 E.5 Proof of Theorem 4.9

757 *Proof.* For simplicity, we use \mathbf{B} instead of \mathbf{B}^t . Initially, we prove the following important lemmas.

758 **Lemma E.6.** (Riemannian gradient Lower Bound for the Iterates Gap) We define $\phi \triangleq (3\bar{X} + \bar{\mathbf{V}}\bar{X})\bar{\mathbf{G}} +$
 759 $(1 + \bar{X}^2 + \bar{\mathbf{V}}^2 + \bar{\mathbf{V}}^2 \bar{X}^2)L_f + (1 + \bar{\mathbf{V}}^2)\theta$. It holds that: $\mathbb{E}_{\xi^{t+1}}[\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1}))] \leq$
 760 $\phi \cdot \mathbb{E}_{\xi^t}[\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}]$.

761 *Proof.* The proof process is exactly the same as in lemma E.2 and will not be repeated here. \square

762 The following lemma is useful to outline the relation of $\|\nabla_{\mathcal{J}}f(\mathbf{X}^t)\|_{\mathbb{F}}$ and $\|\nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_{\mathbb{F}}$.

763 **Lemma E.7.** We have the following results:

$$764 \text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}f(\mathbf{X}^t)) \leq \gamma \cdot \mathbb{E}_{\xi^{t-1}}[\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B}))] \text{ with } \gamma \triangleq \bar{X}\sqrt{C_n^2}.$$

765 *Proof.* We have the following inequalities:

$$\begin{aligned} \|\nabla_{\mathcal{J}}f(\mathbf{X}^t)\|_{\mathbb{F}}^2 &\stackrel{\textcircled{1}}{=} \|\mathbf{G}^t - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^\top \mathbf{X}^t \mathbf{J}\|_{\mathbb{F}}^2 \\ &\stackrel{\textcircled{2}}{=} \|\mathbf{G}^t(\mathbf{X}^t)^\top \mathbf{J}\mathbf{X}^t \mathbf{J} - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^\top \mathbf{J}\mathbf{J}\mathbf{X}^t \mathbf{J}\|_{\mathbb{F}}^2 \\ &\stackrel{\textcircled{3}}{\leq} \|\mathbf{G}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^\top \mathbf{J}\|_{\mathbb{F}}^2 \|\mathbf{J}\mathbf{X}^t \mathbf{J}\|_{\mathbb{F}}^2 \\ &\stackrel{\textcircled{4}}{\leq} \bar{X}^2 \|\mathbf{W}\|_{\mathbb{F}}^2, \text{ with } \mathbf{W} \triangleq \mathbf{G}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^\top \mathbf{J} \\ &\stackrel{\textcircled{5}}{\leq} \bar{X}^2 C_n^2 \cdot \mathbb{E}_{\xi^{t-1}}[\|\mathbf{U}_B^\top [\mathbf{G}^t(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^\top \mathbf{J}] \mathbf{U}_B\|_{\mathbb{F}}^2] \\ &\stackrel{\textcircled{6}}{=} \bar{X}^2 C_n^2 \cdot \mathbb{E}_{\xi^{t-1}}[\|\nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_{\mathbb{F}}^2] \end{aligned}$$

766 where step ① uses the definition of $\nabla_{\mathcal{J}}f(\mathbf{X}^t)$; step ② uses $\mathbf{J}\mathbf{J} = \mathbf{I}$ and $\mathbf{X}^\top \mathbf{J}\mathbf{X} = \mathbf{J} \Rightarrow \mathbf{X}^\top \mathbf{J}\mathbf{X}\mathbf{J} =$
 767 $\mathbf{J}\mathbf{J} = \mathbf{I}$; step ③ uses the norm inequality and ; step ④ uses the definition of $\mathbf{W} \triangleq \mathbf{G}^t(\mathbf{X}^t)^\top -$
 768 $\mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^\top \mathbf{J}$ and $\forall t, \|\mathbf{X}^t\|_{\mathbb{F}} \leq \bar{X}$; step ⑤ uses Lemma (A.1) with $k = 2$; step ⑥ uses the definition
 769 of $\nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})$. Taking the square root of both sides, we finish the proof of this lemma. \square

770 Finally, we obtain our main convergence results. First of all, since $f^\circ(\mathbf{X}) \triangleq f(\mathbf{X}) + \mathcal{I}_{\mathcal{J}}(\mathbf{X})$ is a KL
 771 function, we have from Proposition 4.8 that:

$$\frac{1}{\varphi'(f^\circ(\mathbf{X}') - f^\circ(\mathbf{X}))} \leq \text{dist}(\mathbf{0}, \nabla f^\circ(\mathbf{X}')) \stackrel{\textcircled{1}}{\leq} \|\nabla_{\mathcal{J}}f(\mathbf{X}')\|_{\mathbb{F}}, \quad (76)$$

where step ① uses Lemma E.5. Here, $\varphi(\cdot)$ is some certain concave desingularization function. Since $\varphi(\cdot)$ is concave, we have:

$$\forall \Delta \in \mathbb{R}, \Delta^+ \in \mathbb{R}, \varphi(\Delta^+) + (\Delta - \Delta^+)\varphi'(\Delta) \leq \varphi(\Delta).$$

772 Applying the inequality above with $\Delta = f(\mathbf{X}^t) - f(\bar{\mathbf{X}})$ and $\Delta^+ = f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}})$, we have:

$$\begin{aligned} &(f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}))\varphi'(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) \\ &\leq \varphi(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}})) \triangleq \mathcal{E}^t. \end{aligned} \quad (77)$$

773 We derive the following inequalities:

$$\begin{aligned}
\mathbb{E}_{\xi^t} \left[\frac{\theta}{2} \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}^2 \right] &\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\xi^t} [f(\mathbf{X}^t) - f(\mathbf{X}^{t+1})] \\
&\stackrel{\textcircled{2}}{\leq} \mathbb{E}_{\xi^t} \left[\frac{\mathcal{E}^t}{\varphi'(f(\mathbf{X}^t) - f(\bar{\mathbf{X}}))} \right] \\
&\stackrel{\textcircled{3}}{\leq} \mathbb{E}_{\xi^t} [\mathcal{E}^t \|\nabla_{\mathcal{J}} f(\mathbf{X}^t)\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{4}}{\leq} \mathbb{E}_{\xi^t} [\mathcal{E}^t \gamma \|\nabla_{\mathcal{J}} \mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{5}}{\leq} \mathbb{E}_{\xi^{t-1}} [\mathcal{E}^t \gamma \phi \|\bar{\mathbf{V}}^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\stackrel{\textcircled{6}}{\leq} \mathbb{E}_{\xi^{t-1}} \left[\frac{\theta'}{2} \|\bar{\mathbf{V}}^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2 + \frac{(\mathcal{E}^t \gamma \phi)^2}{2\theta'} \right], \forall \theta' > 0,
\end{aligned}$$

774 where step ① uses the sufficient descent condition as shown in Theorem 4.6; step ② uses Inequality
775 (77); step ③ uses Inequality (76) with $\mathbf{X}' = \mathbf{X}^t$ and $\mathbf{X} = \bar{\mathbf{X}}$; step ④ uses Lemma E.7; step ⑤ uses
776 Lemma E.6; step ⑥ applies the inequality that $\forall \theta' > 0, a, b, ab \leq \frac{\theta' a^2}{2} + \frac{b^2}{2\theta'}$ with $a = \|\bar{\mathbf{V}}^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}$
777 and $b = \mathcal{E}^t \gamma \phi$.

778 Multiplying both sides by 2 and taking the square root of both sides, we have:

$$\begin{aligned}
\sqrt{\theta} \mathbb{E}_{\xi^t} [\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}] &\leq \sqrt{\frac{(\mathcal{E}^t \gamma \phi)^2}{\theta'} + \theta' \mathbb{E}_{\xi^{t-1}} [\|\bar{\mathbf{V}}^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2]}, \forall \theta' > 0 \\
&\stackrel{\textcircled{1}}{\leq} \sqrt{\theta'} \mathbb{E}_{\xi^{t-1}} [\|\bar{\mathbf{V}}^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}] + \frac{\mathcal{E}^t \gamma \phi}{\sqrt{\theta'}}, \forall \theta' > 0,
\end{aligned}$$

779 where step ① uses the inequality that $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ for all $a \geq 0$ and $b \geq 0$. Summing the
780 inequality above over $i = 1, 2, \dots, t$, we have:

$$\begin{aligned}
&\sqrt{\theta} \mathbb{E}_{\xi^t} [\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}] - \sqrt{\theta'} \mathbb{E}_{\xi^0} [\|\bar{\mathbf{V}}^0 - \mathbf{I}_2\|_{\mathbb{F}}] + \sum_{i=1}^{t-1} (\sqrt{\theta} - \sqrt{\theta'}) \mathbb{E}_{\xi^i} [\|\bar{\mathbf{V}}^i - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\leq \frac{\gamma \phi}{\sqrt{\theta'}} \sum_{i=1}^t \mathcal{E}^i \\
&\stackrel{\textcircled{1}}{=} \frac{\gamma \phi}{\sqrt{\theta'}} \sum_{i=1}^t \varphi(f(\mathbf{X}^i) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{i+1}) - f(\bar{\mathbf{X}})) \\
&\stackrel{\textcircled{2}}{=} \frac{\gamma \phi}{\sqrt{\theta'}} [\varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}}))] \\
&\stackrel{\textcircled{3}}{\leq} \frac{\gamma \phi}{\sqrt{\theta'}} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})),
\end{aligned}$$

781 where step ① uses the definition of \mathcal{E}^i in (77); step ② uses a basic recursive reduction; step ③ uses
782 the fact the desingularization function $\varphi(\cdot)$ is positive. With the choice $\theta' = \frac{\theta}{4}$, we have:

$$\begin{aligned}
&\sqrt{\theta} \mathbb{E}_{\xi^t} [\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}] + \frac{\sqrt{\theta}}{2} \sum_{i=1}^{t-1} \mathbb{E}_{\xi^i} [\|\bar{\mathbf{V}}^i - \mathbf{I}_2\|_{\mathbb{F}}] \\
&\leq \frac{2\gamma \phi}{\sqrt{\theta}} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \frac{\sqrt{\theta}}{2} \mathbb{E}_{\xi^0} [\|\bar{\mathbf{V}}^0 - \mathbf{I}_2\|_{\mathbb{F}}] \tag{78}
\end{aligned}$$

$$\stackrel{\textcircled{1}}{\leq} \frac{2\gamma \phi}{\sqrt{\theta}} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \frac{\sqrt{\theta}}{2} (\bar{\mathbf{V}} + \sqrt{2}), \tag{79}$$

783 where step ① uses $\forall t, \|\mathbf{V}\|_{\mathbb{F}} \leq \bar{\mathbf{V}}$, then, $\|\mathbf{V} - \mathbf{I}_2\|_{\mathbb{F}} \leq \|\mathbf{V}\|_{\mathbb{F}} + \|\mathbf{I}\|_{\mathbb{F}} \leq \bar{\mathbf{V}} + \sqrt{2}$. Finally, we obtain
784 from Inequality (79):

$$\begin{aligned}
&\frac{1}{2} \sum_{i=1}^t \mathbb{E}_{\xi^i} [\|\bar{\mathbf{V}}^i - \mathbf{I}_2\|_{\mathbb{F}}] \leq \frac{2\gamma \phi}{\theta} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \frac{1}{2} (\bar{\mathbf{V}} + \sqrt{2}) \\
&\stackrel{\textcircled{1}}{\Rightarrow} \frac{1}{2} \sum_{i=1}^t \mathbb{E}_{\xi^i} [\|\mathbf{X}^{i+1} - \mathbf{X}^i\|_{\mathbb{F}}] \leq \left(\frac{2\bar{\mathbf{X}}\gamma \phi}{\theta} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \frac{\bar{\mathbf{X}}}{2} (\bar{\mathbf{V}} + \sqrt{2}) \right)
\end{aligned}$$

where step ① uses the inequality that $\frac{\|\mathbf{X}^{i+1} - \mathbf{X}^i\|_{\mathbb{F}}}{\bar{\mathbf{X}}} \leq \|\bar{\mathbf{V}}^i - \mathbf{I}_2\|_{\mathbb{F}}$ as shown in Part (b) in Lemma 2.1.
Finally, we can get the expression for C :

$$C \triangleq \frac{4\bar{\mathbf{X}}\gamma \phi}{\theta} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \bar{\mathbf{X}} (\bar{\mathbf{V}} + \sqrt{2}) = n\mathcal{O}(\varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})))$$

785

□

786 F Additional Experiment Details and Results

787 F.1 Additional Details for Hyperbolic Structural Probe Problem

788 To begin with, we give the definition of the Ultrahyperbolic manifold $\mathbb{U}_\alpha^{p,q}$, which will be used in
789 Ultra-hyperbolic geodesic distance $\mathbf{d}_\alpha(\mathbf{x}, \mathbf{y})$ and Diffeomorphism $\varphi(\cdot)$.

790 ► **Ultrahyperbolic manifold.** Vectors in an ultrahyperbolic manifold is defined as $\mathbb{U}_\alpha^{p,q} = \{\mathbf{x} =$
791 $(x_1, x_2, \dots, x_{p+q})^\top \in \mathbb{R}^{p,q} : \|\mathbf{x}\|_q^2 = -\alpha^2\}$ [48], where α is a non-negative real number denoting
792 the radius of curvature. $\|\mathbf{x}\|_q^2 = \langle \mathbf{x}, \mathbf{x} \rangle_q, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{p,q}, \langle \mathbf{x}, \mathbf{y} \rangle_q = \sum_{i=1}^p \mathbf{x}_i \mathbf{y}_i - \sum_{j=p+1}^{p+q} \mathbf{x}_j \mathbf{y}_j$ is
793 a norm of the induced scalar product. The hyperbolic and spherical manifolds can be defined as
794 $\mathbb{H}_\alpha = \mathbb{U}_\alpha^{p,1}, \mathbb{S}_\alpha = \mathbb{U}_\alpha^{0,q}$.

795 ► **Ultra-hyperbolic geodesic distance.** The ultra-hyperbolic geodesic distance [27][28] $\mathbf{d}_\alpha(\cdot, \cdot)$ is
796 formulated: $\forall \mathbf{x} \in \mathbb{U}_\alpha^{p,q}, \mathbf{y} \in \mathbb{U}_\alpha^{p,q}$ and $\alpha > 0, \mathbf{d}_\alpha(\mathbf{x}, \mathbf{y}) = \begin{cases} \alpha \cosh^{-1}(|\frac{\langle \mathbf{x}, \mathbf{y} \rangle_q}{\alpha^2}|) & \text{if } |\frac{\langle \mathbf{x}, \mathbf{y} \rangle_q}{\alpha^2}| \geq 1 \\ \alpha \cos^{-1}(|\frac{\langle \mathbf{x}, \mathbf{y} \rangle_q}{\alpha^2}|) & \text{otherwise.} \end{cases}$

797 ► **Diffeomorphism.** [Theorem 1 Diffeomorphism of [49]]: Any vector $\mathbf{x} \in \mathbb{R}^p \times \mathbb{R}_*^q$ can be mapped
798 into $\mathbb{U}_\alpha^{p,q}$ by a double projection $\varphi = \phi^{-1} \circ \phi$, with $\psi(\mathbf{x}) = \begin{pmatrix} \mathbf{s} \\ \alpha \frac{\mathbf{t}}{\|\mathbf{t}\|} \end{pmatrix}, \psi^{-1}(\mathbf{z}) = \begin{pmatrix} \mathbf{v} \\ \frac{\sqrt{\alpha^2 + \|\mathbf{v}\|^2}}{\alpha} \mathbf{u} \end{pmatrix},$
799 where $\mathbf{x} = \begin{pmatrix} \mathbf{s} \\ \mathbf{t} \end{pmatrix} \in \mathbb{U}_\alpha^{p,q}$ with $\mathbf{s} \in \mathbb{R}^p$ and $\mathbf{t} \in \mathbb{R}_*^q, \mathbf{z} = \begin{pmatrix} \mathbf{v} \\ \mathbf{u} \end{pmatrix} \in \mathbb{R}^p \times \mathbb{S}_\alpha^q$ with $\mathbf{v} \in \mathbb{R}^p$ and $\mathbf{u} \in \mathbb{S}_\alpha^q$.

800 F.2 Additional application: Ultra-hyperbolic Knowledge Graph Embedding

801 The J orthogonal matrix can be used as an isometric linear operator in the Ultrahyperbolic manifold,
802 [48] et al. extended the knowledge graph model from hyperbolic space to Ultra-hyperbolic space
803 (named as **UltraE**) by this property. The **UltraE** model is formulated as follows:

$$\min_{\mathbf{R}, \mathbf{E}, \mathbf{b}} \mathcal{L}(\mathbf{R}, \mathbf{E}, \mathbf{b}) \triangleq -\frac{1}{N} \sum_{(h,r,t) \in \Delta} (\log s(h, r, t) + \sum_{(h',r',t') \in \Delta'_{(h,r,t)}} \log(1 - s(h', r', t'))) \\ \text{s.t. } \begin{cases} s(h, r, t) = \sigma(-d_\alpha^2(\mathbf{R}_r \mathbf{E}_h, \mathbf{E}_t) + \mathbf{b}_h + \mathbf{b}_t + \delta) \\ \mathbf{R}_r^\top \mathbf{J} \mathbf{R}_r = \mathbf{J} \end{cases}$$

804 where $\mathbf{E} \in \mathbb{R}^{n_e \times n}$ with $\mathbf{E}_h = \mathbf{E}(h, \cdot) \in \mathbb{U}_\alpha^{p,q}, \mathbf{b} \in \mathbb{R}^{n_r}$ with $\mathbf{b}_h = \mathbf{b}(r) \in \mathbb{R}, \mathbf{R} \in \mathbb{R}^{n_r \times n \times n}$ with
805 $\mathbf{R}_r = \mathbf{R}(r, \cdot, \cdot) \in \mathbb{R}^{n \times n}$ and $\mathbf{J} = \begin{bmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_q \end{bmatrix}; \Delta \in \mathbb{N}^{N \times 3}$ is the set of positive triplets, $\Delta'_{(h,r,t)} \in$
806 $\mathbb{N}^{N \times k \times 3}$ denotes the set of negative triples constructed by corrupting $(h, r, t); \delta$ is a global margin
807 hyper-parameter, $\sigma(\cdot)$ is the sigmoid function, n_e represents the number of entities and n_r represents
808 the number of relations; $d_\alpha(\cdot)$ stands for the Ultra-hyperbolic geodesic distance (refer to F.1).

809 ► **Experiment Details.** We selected a batch of **FB15K** and **WN18RR** respectively as the data set for
810 the Ultra-hyperbolic Knowledge Graph Embedding problem, (training set size, test set size, number
811 of entities, number of relations) are (719,308,135,22) and (545,233,208,5) respectively. $n = 36,$
812 $p = 18, \delta = 5, \alpha = 1$ and $k = 50$. In order to highlight the difference between J orthogonal
813 optimization, in the **UltraE** model, all entities and biases of the optimization algorithm are optimized
814 using **ADMM** by **Pytorch**, $lr = 5e - 4$. We use the **Adagrad** optimizer in Pytorch to optimize the
815 J-orthogonality constraint variable in the **CS** model.

816 F.3 Experiment result

817 ► **Hyperbolic Eigenvalue Problem.** Table 2 and Figure 3, 4, 5 are supplementary experiments for
818 HEVP. Several conclusions can be drawn. (i) **GS-JOB**CD often greatly improves upon **UMCM**,
819 **ADMM** and **CSDM**. This is because our methods find stronger stationary points than them. (ii)
820 **J-JOB**CD is a parallel version of **GS-JOB**CD and thus exhibits significantly faster convergence. (iii)
821 The proposed methods generally give the best performance.

822 ► **Hyperbolic Structural Probe Problem.** Table 3 and Figure 6, 7 are supplementary experiments
823 for HSPP. Several conclusions can be drawn. (i) **J-JOB**CD often greatly improves upon **UMCM**,
824 **ADMM** and **CSDM** (ii) **VR-J-JOB**CD is a reduced variance version of **J-JOB**CD and thus exhibits

825 significantly faster convergence for problems with large samples. (iii) The proposed methods generally
826 give the best performance.

827 ► **Ultra-hyperbolic Knowledge Graph Embedding Problem.** Figure 8, 9, 10 and 11 are supplement-
828 ary experiments for **UltraE**. Several conclusions can be drawn. (i) In terms of Epoch performance,
829 **J-JOBCD** and **VR-J-JOBCD** often greatly improves upon **CSDM**, thus they show better MRR and
830 hits results. (ii) In models with limited sample sizes, the computational efficiency of **VR-J-JOBCD**
831 is inferior to that of **J-JOBCD**. This discrepancy arises because each iteration in **VR-J-JOBCD**
832 necessitates two instances of backpropagation, thus consuming substantial computational resources.
833 (iii) The proposed methods generally give the best performance.

Table 2: The convergence curve of the compared methods for solving HEVP. (+) indicates that after the convergence of the CSDM, UMCM and ADMM, utilizing the GS-JOB CD for optimization markedly enhances the objective value. The 1st, 2nd, and 3rd best results are colored with red, green and blue, respectively. (n, p) represents the dimension and p-value of the J orthogonal matrix (square matrix). The value in () stands for $\sum_{ij} |X^T J X - J|_{ij}$.

dataname	(m-n-p)	UMCM	ADMM	CSDM	GS-JOB CD		J-JOB CD	UMCM+GS-JOB CD	ADMM+GS-JOB CD	CSDM+GS-JOB CD
					time limit:30s	time limit:30s				
CnnCaltech	cfar (1000-1000-50)	-1.05e+04(3.0e-09)	-1.05e+04(3.0e-09)	-5.28e+04(4.8e-09)	-7.76e+04(1.6e-08)	-1.19e+05(8.1e-08)	-7.96e+04(1.1e-08)(+)	-5.86e+04(8.4e-09)(+)	-8.50e+04(1.2e-08)(+)	-8.50e+04(1.2e-08)(+)
	giettete (2000-1000-500)	-5.89e+02(2.9e-08)	-5.89e+02(3.1e-10)	-7.86e+02(3.8e-10)	-7.68e+02(3.2e-10)	-3.90e+03(2.3e-08)	-6.71e+02(2.9e-08)(+)	-6.73e+02(3.2e-10)(+)	-8.49e+02(3.3e-10)(+)	-8.49e+02(3.3e-10)(+)
	mimst (3000-1000-500)	-3.22e+06(3.1e-10)	-3.22e+06(3.1e-10)	-5.16e+06(3.9e-10)	-7.13e+06(5.9e-10)	-1.15e+07(1.5e-08)	-4.63e+06(3.7e-10)(+)	-4.74e+06(3.7e-10)(+)	-6.10e+06(4.6e-10)(+)	-6.10e+06(4.6e-10)(+)
	randn10 (10-10-5)	-8.65e+04(4.1e-10)	-8.65e+04(4.1e-10)	-1.63e+05(4.9e-10)	-2.22e+05(6.6e-10)	-6.99e+05(2.1e-08)	-1.59e+05(4.7e-10)(+)	-1.48e+05(4.7e-10)(+)	-2.04e+05(5.7e-10)(+)	-2.04e+05(5.7e-10)(+)
	randn100 (100-100-50)	-1.03e+04(3.0e-09)	-1.03e+04(2.5e-07)	-1.98e+04(5.1e-09)	-1.49e+04(2.7e-08)	-3.44e+04(1.3e-07)	-1.33e+04(1.4e-08)(+)	-1.31e+04(2.5e-07)(+)	-2.20e+04(1.8e-08)(+)	-2.20e+04(1.8e-08)(+)
	randn1000 (1000-1000-500)	-1.16e+06(3.1e-10)	-1.16e+06(3.1e-10)	-1.47e+06(3.9e-10)	-1.22e+06(6.4e-10)	-4.83e+06(7.3e-08)	-1.18e+06(3.5e-10)(+)	-1.18e+06(3.5e-10)(+)	-1.49e+06(4.7e-10)(+)	-1.49e+06(4.7e-10)(+)
	sector (500-1000-500)	-3.61e+03(3.1e-10)	-3.61e+03(3.1e-10)	-5.35e+03(2.9e-10)	-6.68e+03(5.8e-10)	-1.07e+04(1.2e-08)	-4.73e+03(3.7e-10)(+)	-4.85e+03(3.6e-10)(+)	-6.47e+03(4.6e-10)(+)	-6.47e+03(4.6e-10)(+)
	TD2T (10-10-7)	-4.25e+06(2.2e-10)	-4.25e+06(2.2e-10)	-6.37e+06(4.0e-10)	-8.20e+06(6.0e-10)	-1.32e+07(1.2e-08)	-5.67e+06(3.7e-10)(+)	-5.93e+06(3.7e-10)(+)	-7.85e+06(4.8e-10)(+)	-7.85e+06(4.8e-10)(+)
	w1a (2470-290-145)	-3.02e+04(1.1e-04)	-3.02e+04(1.1e-04)	-5.42e+04(3.9e-09)	-5.74e+04(1.1e-04)	-6.73e+05(1.1e-04)	-4.76e+04(1.1e-04)(+)	-4.57e+04(1.1e-04)(+)	-6.32e+04(4.4e-04)(+)	-6.32e+04(4.4e-04)(+)
	w1a (10-10-9)	-1.29e+02(9.7e-02)	-1.29e+02(9.7e-02)	-3.03e+02(2.3e-01)	-3.96e+01(9.7e-02)	-2.98e+02(9.7e-02)	-7.05e+01(9.7e-02)(+)	-1.75e+01(9.7e-02)(+)	-2.29e+02(2.3e-01)(+)	-2.29e+02(2.3e-01)(+)
CnnCaltech	cfar (1000-1000-70)	-7.32e+03(1.9e-09)	-7.32e+03(1.9e-09)	-3.29e+04(3.2e-09)	-6.01e+04(1.5e-08)	-1.12e+05(7.4e-08)	-4.84e+04(1.0e-08)(+)	-4.30e+04(7.6e-09)(+)	-7.52e+04(1.3e-08)(+)	-7.52e+04(1.3e-08)(+)
	giettete (2000-1000-500)	-4.33e+02(2.1e-08)	-4.33e+02(2.2e-10)	-5.43e+02(2.5e-10)	-5.02e+02(2.6e-10)	-2.88e+03(1.6e-08)	-4.86e+02(2.1e-08)(+)	-4.85e+02(2.1e-08)(+)	-5.98e+02(2.3e-10)(+)	-5.98e+02(2.3e-10)(+)
	mimst (1000-780-500)	-2.45e+06(2.2e-10)	-2.45e+06(2.2e-10)	-3.59e+06(2.5e-10)	-5.02e+06(4.2e-10)	-9.17e+06(1.0e-08)	-3.15e+06(2.5e-10)(+)	-3.25e+06(2.5e-10)(+)	-4.25e+06(2.7e-10)(+)	-4.25e+06(2.7e-10)(+)
	randn10 (10-10-7)	-7.05e+04(3.1e-10)	-7.05e+04(3.1e-10)	-1.12e+05(3.6e-10)	-1.18e+05(3.5e-10)	-6.28e+05(1.9e-08)	-1.14e+05(3.5e-10)(+)	-1.21e+05(3.5e-10)(+)	-1.59e+05(4.3e-10)(+)	-1.59e+05(4.3e-10)(+)
	randn100 (100-100-50)	-1.03e+04(3.0e-09)	-1.03e+04(2.5e-07)	-1.98e+04(5.1e-09)	-1.41e+04(2.7e-08)	-3.44e+04(1.3e-07)	-1.33e+04(1.4e-08)(+)	-1.31e+04(2.5e-07)(+)	-2.20e+04(1.8e-08)(+)	-2.20e+04(1.8e-08)(+)
	randn1000 (1000-1000-500)	-1.16e+06(3.1e-10)	-1.16e+06(3.1e-10)	-1.47e+06(3.9e-10)	-1.22e+06(6.4e-10)	-4.83e+06(7.3e-08)	-1.18e+06(3.5e-10)(+)	-1.18e+06(3.5e-10)(+)	-1.49e+06(4.7e-10)(+)	-1.49e+06(4.7e-10)(+)
	sector (500-1000-500)	-3.61e+03(3.1e-10)	-3.61e+03(3.1e-10)	-5.35e+03(2.9e-10)	-6.68e+03(5.8e-10)	-1.07e+04(1.2e-08)	-4.73e+03(3.7e-10)(+)	-4.85e+03(3.6e-10)(+)	-6.47e+03(4.6e-10)(+)	-6.47e+03(4.6e-10)(+)
	TD2T (10-10-7)	-4.25e+06(2.2e-10)	-4.25e+06(2.2e-10)	-6.37e+06(4.0e-10)	-8.20e+06(6.0e-10)	-1.32e+07(1.2e-08)	-5.67e+06(3.7e-10)(+)	-5.93e+06(3.7e-10)(+)	-7.85e+06(4.8e-10)(+)	-7.85e+06(4.8e-10)(+)
	w1a (2470-290-145)	-3.02e+04(1.1e-04)	-3.02e+04(1.1e-04)	-5.42e+04(3.9e-09)	-5.74e+04(1.1e-04)	-6.73e+05(1.1e-04)	-4.76e+04(1.1e-04)(+)	-4.57e+04(1.1e-04)(+)	-6.32e+04(4.4e-04)(+)	-6.32e+04(4.4e-04)(+)
	w1a (10-10-9)	-1.29e+02(9.7e-02)	-1.29e+02(9.7e-02)	-3.03e+02(2.3e-01)	-3.96e+01(9.7e-02)	-2.98e+02(9.7e-02)	-7.05e+01(9.7e-02)(+)	-1.75e+01(9.7e-02)(+)	-2.29e+02(2.3e-01)(+)	-2.29e+02(2.3e-01)(+)
CnnCaltech	cfar (1000-1000-70)	-7.32e+03(1.9e-09)	-7.32e+03(1.9e-09)	-3.29e+04(3.2e-09)	-6.01e+04(1.5e-08)	-1.12e+05(7.4e-08)	-4.84e+04(1.0e-08)(+)	-4.30e+04(7.6e-09)(+)	-7.52e+04(1.3e-08)(+)	-7.52e+04(1.3e-08)(+)
	giettete (2000-1000-500)	-4.33e+02(2.1e-08)	-4.33e+02(2.2e-10)	-5.43e+02(2.5e-10)	-5.02e+02(2.6e-10)	-2.88e+03(1.6e-08)	-4.86e+02(2.1e-08)(+)	-4.85e+02(2.1e-08)(+)	-5.98e+02(2.3e-10)(+)	-5.98e+02(2.3e-10)(+)
	mimst (1000-780-500)	-2.45e+06(2.2e-10)	-2.45e+06(2.2e-10)	-3.59e+06(2.5e-10)	-5.02e+06(4.2e-10)	-9.17e+06(1.0e-08)	-3.15e+06(2.5e-10)(+)	-3.25e+06(2.5e-10)(+)	-4.25e+06(2.7e-10)(+)	-4.25e+06(2.7e-10)(+)
	randn10 (10-10-7)	-7.05e+04(3.1e-10)	-7.05e+04(3.1e-10)	-1.12e+05(3.6e-10)	-1.18e+05(3.5e-10)	-6.28e+05(1.9e-08)	-1.14e+05(3.5e-10)(+)	-1.21e+05(3.5e-10)(+)	-1.59e+05(4.3e-10)(+)	-1.59e+05(4.3e-10)(+)
	randn100 (100-100-50)	-1.03e+04(3.0e-09)	-1.03e+04(2.5e-07)	-1.98e+04(5.1e-09)	-1.41e+04(2.7e-08)	-3.44e+04(1.3e-07)	-1.33e+04(1.4e-08)(+)	-1.31e+04(2.5e-07)(+)	-2.20e+04(1.8e-08)(+)	-2.20e+04(1.8e-08)(+)
	randn1000 (1000-1000-500)	-1.16e+06(3.1e-10)	-1.16e+06(3.1e-10)	-1.47e+06(3.9e-10)	-1.22e+06(6.4e-10)	-4.83e+06(7.3e-08)	-1.18e+06(3.5e-10)(+)	-1.18e+06(3.5e-10)(+)	-1.49e+06(4.7e-10)(+)	-1.49e+06(4.7e-10)(+)
	sector (500-1000-500)	-3.61e+03(3.1e-10)	-3.61e+03(3.1e-10)	-5.35e+03(2.9e-10)	-6.68e+03(5.8e-10)	-1.07e+04(1.2e-08)	-4.73e+03(3.7e-10)(+)	-4.85e+03(3.6e-10)(+)	-6.47e+03(4.6e-10)(+)	-6.47e+03(4.6e-10)(+)
	TD2T (10-10-7)	-4.25e+06(2.2e-10)	-4.25e+06(2.2e-10)	-6.37e+06(4.0e-10)	-8.20e+06(6.0e-10)	-1.32e+07(1.2e-08)	-5.67e+06(3.7e-10)(+)	-5.93e+06(3.7e-10)(+)	-7.85e+06(4.8e-10)(+)	-7.85e+06(4.8e-10)(+)
	w1a (2470-290-145)	-3.02e+04(1.1e-04)	-3.02e+04(1.1e-04)	-5.42e+04(3.9e-09)	-5.74e+04(1.1e-04)	-6.73e+05(1.1e-04)	-4.76e+04(1.1e-04)(+)	-4.57e+04(1.1e-04)(+)	-6.32e+04(4.4e-04)(+)	-6.32e+04(4.4e-04)(+)
	w1a (10-10-9)	-1.29e+02(9.7e-02)	-1.29e+02(9.7e-02)	-3.03e+02(2.3e-01)	-3.96e+01(9.7e-02)	-2.98e+02(9.7e-02)	-7.05e+01(9.7e-02)(+)	-1.75e+01(9.7e-02)(+)	-2.29e+02(2.3e-01)(+)	-2.29e+02(2.3e-01)(+)
CnnCaltech	cfar (1000-1000-900)	-6.42e+03(1.2e-09)	-6.42e+03(1.2e-09)	-1.19e+04(1.2e-09)	-3.14e+04(1.4e-08)	-5.05e+04(5.6e-08)	-2.64e+04(6.3e-09)(+)	-3.35e+04(7.0e-09)(+)	-3.89e+04(6.7e-09)(+)	-3.89e+04(6.7e-09)(+)
	giettete (2000-1000-900)	-3.10e+02(1.1e-08)	-3.10e+02(1.2e-10)	-3.41e+02(1.3e-10)	-3.64e+02(1.7e-10)	-1.65e+03(8.0e-09)	-3.29e+02(1.1e-08)(+)	-3.32e+02(1.1e-08)(+)	-3.61e+02(1.4e-10)(+)	-3.61e+02(1.4e-10)(+)
	mimst (3000-1000-900)	-1.74e+06(1.2e-10)	-1.74e+06(1.2e-10)	-2.05e+06(1.2e-10)	-2.57e+06(1.8e-10)	-6.46e+06(8.0e-09)	-2.00e+06(1.3e-10)(+)	-1.99e+06(1.2e-10)(+)	-2.33e+06(1.4e-10)(+)	-2.33e+06(1.4e-10)(+)
	randn10 (10-10-9)	-5.33e+02(1.7e-01)	-5.33e+02(1.7e-01)	-6.12e+02(2.2e-10)	-1.02e+05(2.9e-10)	-1.70e+05(2.5e-08)	-6.75e+04(2.1e-10)(+)	-6.58e+04(2.1e-10)(+)	-8.02e+04(2.4e-10)(+)	-8.02e+04(2.4e-10)(+)
	randn100 (100-100-90)	-6.14e+03(1.1e-07)	-6.14e+03(1.2e-09)	-8.14e+03(1.2e-09)	-1.03e+04(1.5e-08)	-1.31e+04(5.8e-08)	-2.21e+03(1.7e-01)(+)	-3.46e+02(1.7e-01)(+)	-1.54e+02(1.3e-01)(+)	-1.54e+02(1.3e-01)(+)
	randn1000 (1000-1000-900)	-6.33e+05(1.2e-10)	-6.33e+05(1.2e-10)	-6.84e+05(1.3e-10)	-6.46e+05(1.9e-10)	-1.87e+06(1.8e-08)	-6.39e+05(1.3e-10)(+)	-6.39e+05(1.3e-10)(+)	-6.90e+05(1.5e-10)(+)	-6.90e+05(1.5e-10)(+)
	sector (500-1000-900)	-1.92e+03(1.2e-10)	-1.92e+03(1.2e-10)	-2.18e+03(1.2e-10)	-2.50e+03(1.6e-10)	-5.84e+03(6.5e-09)	-2.12e+03(1.2e-10)(+)	-2.13e+03(1.2e-10)(+)	-2.34e+03(1.4e-10)(+)	-2.34e+03(1.4e-10)(+)
	TD2T (1000-1000-900)	-2.26e+06(1.2e-10)	-2.26e+06(1.2e-10)	-2.58e+06(1.6e-10)	-2.90e+06(1.6e-10)	-7.40e+06(6.6e-09)	-2.53e+06(1.2e-10)(+)	-2.51e+06(1.2e-10)(+)	-2.83e+06(1.3e-10)(+)	-2.83e+06(1.3e-10)(+)
	w1a (2470-290-250)	-2.03e+04(5.4e-10)	-2.03e+04(5.4e-10)	-2.41e+04(5.9e-10)	-2.74e+04(7.9e-09)	-2.59e+05(3.3e-08)	-2.82e+04(1.2e-09)(+)	-3.17e+04(1.4e-09)(+)	-3.78e+04(1.5e-09)(+)	-3.78e+04(1.5e-09)(+)
	w1a (10-10-9)	-1.29e+02(9.7e-02)	-1.29e+02(9.7e-02)	-3.03e+02(2.3e-01)	-3.96e+01(9.7e-02)	-2.98e+02(9.7e-02)	-7.05e+01(9.7e-02)(+)	-1.75e+01(9.7e-02)(+)	-2.29e+02(2.3e-01)(+)	-2.29e+02(2.3e-01)(+)
CnnCaltech	cfar (1000-1000-50)	-1.05e+04(3.0e-09)	-1.05e+04(3.0e-09)	-5.28e+04(4.8e-09)	-7.76e+04(1.6e-08)	-1.19e+05(8.1e-08)	-7.96e+04(1.1e-08)(+)	-5.86e+04(8.4e-09)(+)	-8.50e+04(1.2e-08)(+)	-8.50e+04(1.2e-08)(+)
	giettete (2000-1000-500)	-5.89e+02(2.9e-08)	-5.89e+02(3.1e-10)	-7.86e+02(3.8e-10)	-7.68e+02(3.2e-10)	-3.90e+03(2.3e-08)	-6.71e+02(2.9e-08)(+)	-6.73e+02(3.2e-10)(+)	-8.49e+02(3.3e-10)(+)	-8.49e+02(3.3e-10)(+)
	mimst (3000-1000-500)	-3.22e+06(3.1e-10)	-3.22e+06(3.1e-10)	-5.16e+06(3.9e-10)	-7.13e+06(5.9e-10)	-1.15e+07(1.5e-08)	-4.63e+06(3.7e-10)(+)	-4.74e+06(3.7e-10)(+)	-6.10e+06(4.6e-10)(+)	-6.10e+06(4.6e-10)(+)
	randn10 (10-10-5)	-8.65e+04(4.1e-10)	-8.65e+04(4.1e-10)	-1.63e+05(4.9e-10)	-2.22e+05(6.6e-10)	-6.99e+05(2.1e-08)	-1.59e+05(4.7e-10)(+)	-1.48e+05(4.7e-10)(+)	-2.04e+05(5.7e-10)(+)	-2.04e+05(5.7e-10)(+)
	randn100 (100-100-50)	-1.03e+04(3.0e-09)	-1.03e+04(2.5e-07)	-1.98e+04(5.1e-09)	-1.49e+04(2.7e-08)	-3.44e+04(1.3e-07)	-1.33e+04(1.4e-08)(+)	-1.31e+04(2.5e-07)(+)	-2.20e+04(1.8e-08)(+)	-2.20e+04(1.8e-08)(+)
	randn1000 (1000-1000-500)	-1.16e+06(3.1e-10)	-1.16e+06(3.1e-10)	-1.47e+06(3.9e-10)	-1.22e+06(6.4e-10)	-4.83e+06(7.3e-08)	-1.18e+06(3.5e-10)(+)	-1.18e+06(3.5e-10)(+)	-1.49e+06(4.7e-10)(+)	-1.49e+06(4.7e-10)(+)
	sector (500-1000-500)	-3.61e+03(3.1e-10)	-3.61e+03(3.1e-10)	-5.35e+03(2.9e-10)	-6.68e+03(5.8e-10)	-1.07e+04(1.2e-08)	-4.73e+03(3.7e-10)(+)	-4.85e+03(3.6e-10)(+)	-6.47e+03(4.6e-10)(+)	-6.47e+03(4.6e-10)(+)
	TD2T (10-10-7)	-4.25e+06(2.2e-10)	-4.25e+06(2.2e-10)	-6.37e+06(4.0e-10)	-8.20e+06(6.0e-10)	-1.32e+07(1.2e-08)	-5.67e+06(3.7e-10)(+)	-5.93e+06(3.7e-10)(+)	-7.85e+06(4.8e-10)(+)	-7.85e+06(4.8e-10)(+)
	w1a (2470-290-145)	-3.02e+04(1.1e-04)	-3.02e+04(1.1e-04)	-5.42e+04(3.9e-09)	-5.74e+04(1.1e-04)	-6.73e+05(1.1e-04)	-4.76e+04(1.1e-04)(+)	-4.57e+04(1.1e-04)(+)	-6.32e+04(4.4e-04)(+)	-6.32e+04(4.4e-04)(+)
	w1a (10-10-9)	-1.29e+02(9.7e-02)	-1.29e+02(9.7e-02)	-3.03e+02(2.3e-01)	-3.96e+01(9.7e-02)	-2.98e+02(9.7e-02)	-7.05e+01(9.7e-02)(+)	-1.75e+01(9.7e-02)(+)	-2.29e+02(2.3e-01)(+)	-2.29e+02(2.3e-01)(+)
CnnCaltech	cfar (1000-1000-70)	-7.32e+03(1.9e-09)	-7.32e+03(1.9e-09)	-3.29e+04(3.2e-09)	-6.01e+04(1.5e-08)	-1.12e+05(7.4e-08)	-4.84e+04(1.0e-08)(+)	-4.30e+04(7.		

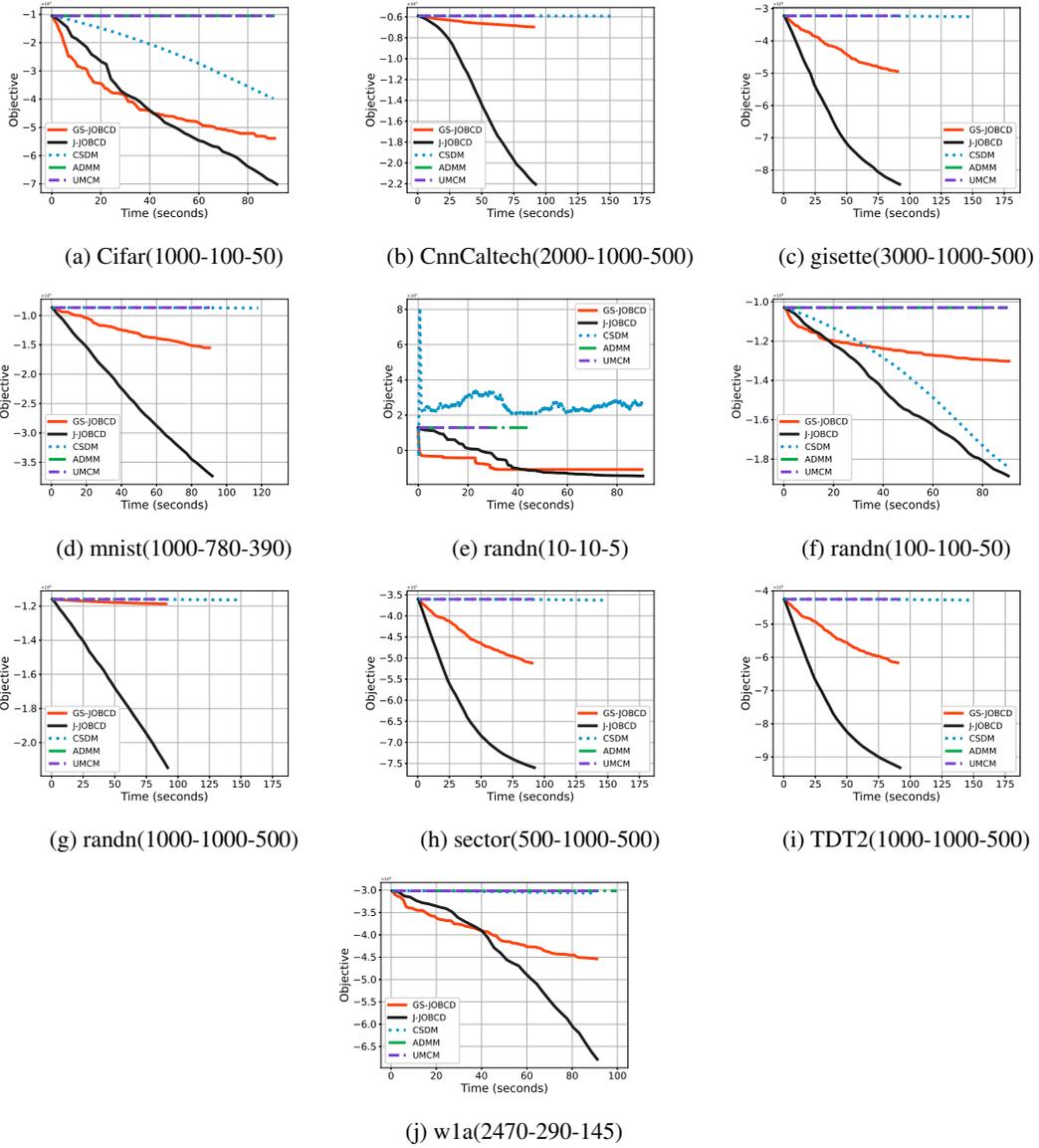


Figure 3: The convergence curve of the compared methods for solving HEVP with varying (m, n, p) .

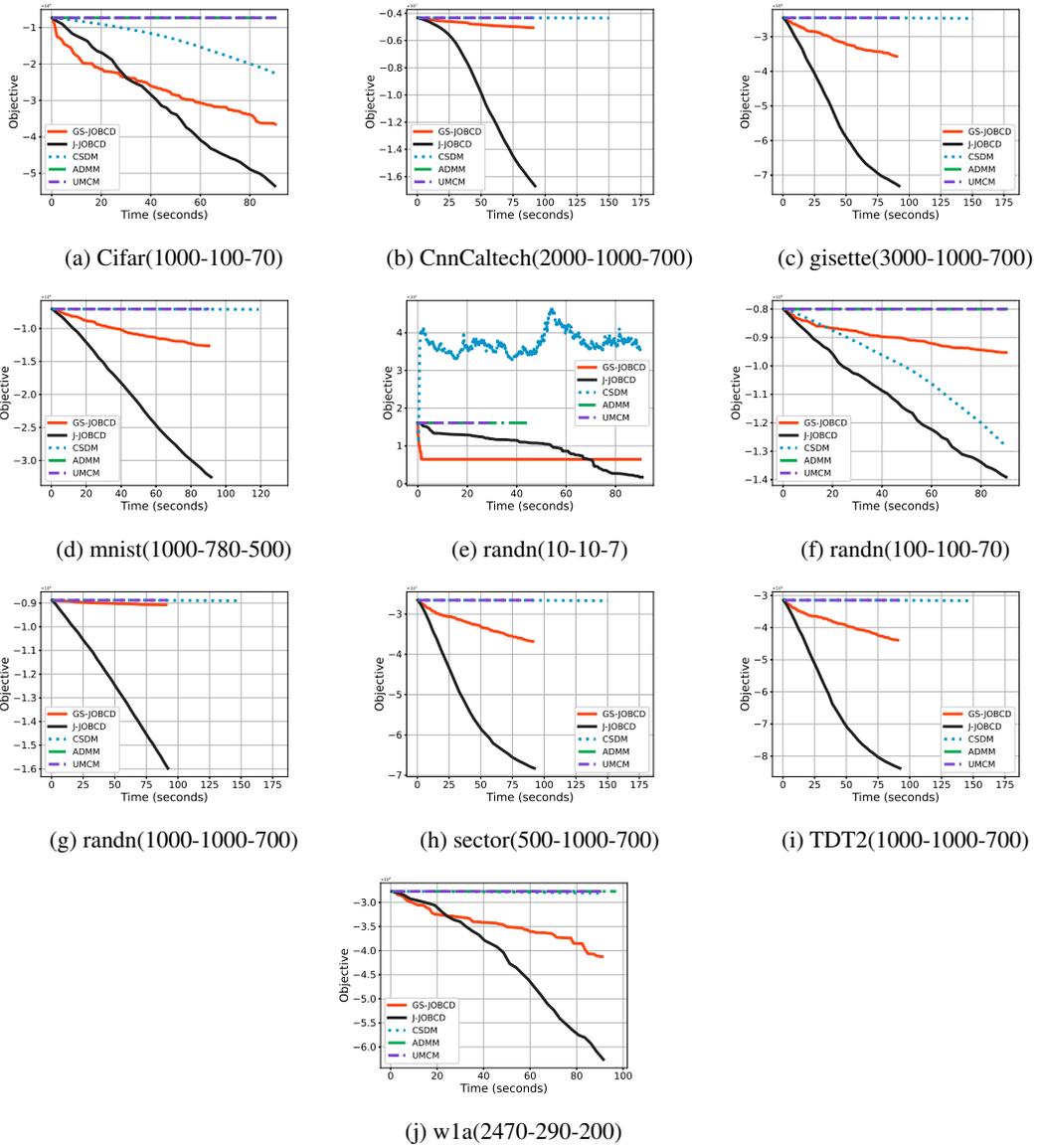


Figure 4: The convergence curve of the compared methods for solving HEVP with varying (m, n, p) .

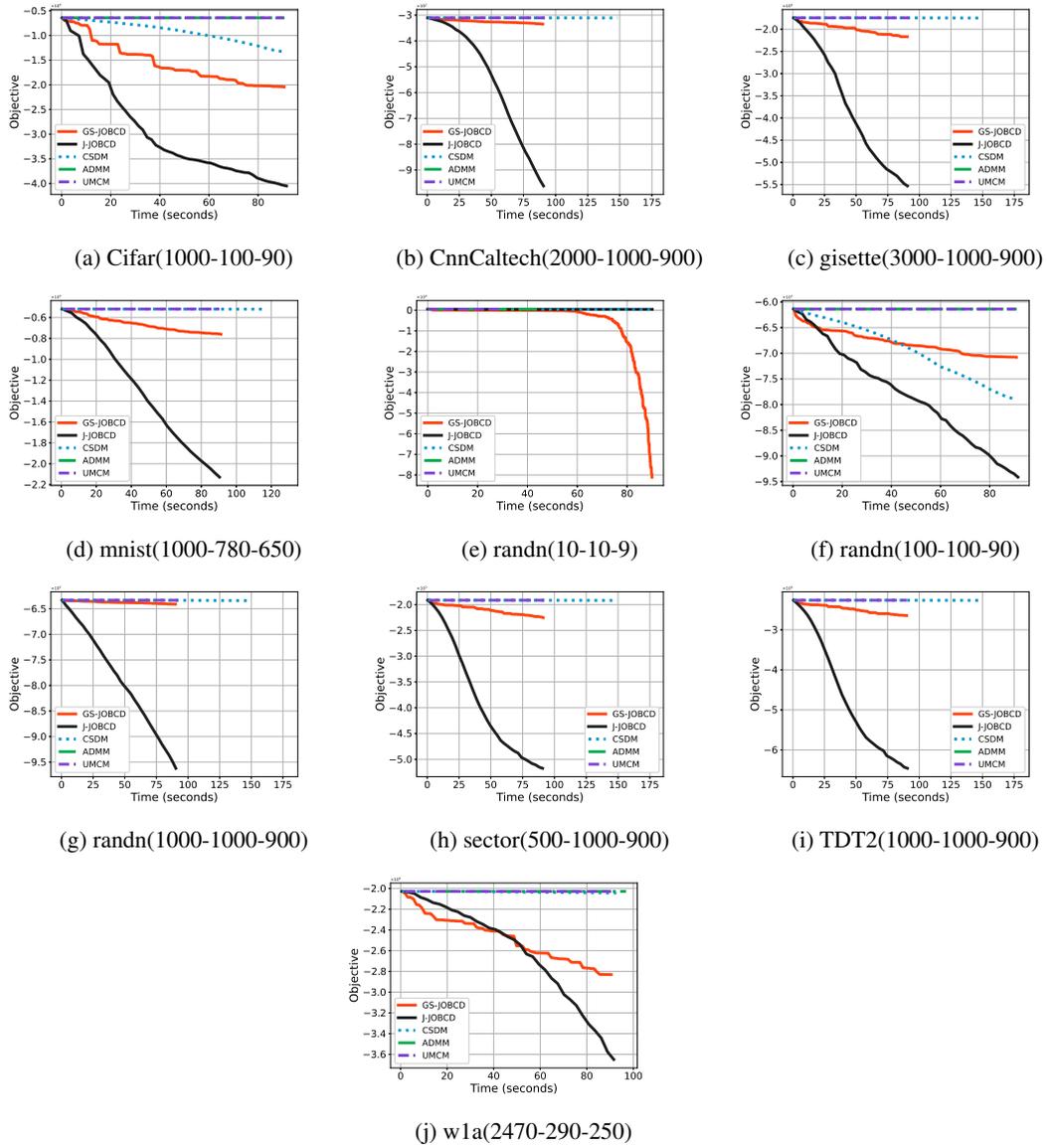


Figure 5: The convergence curve of the compared methods for solving HEVP with varying (m, n, p) .

Table 3: The convergence curve of the compared methods for solving HSPP. (+) indicates that after the convergence of the **CSDM**, utilizing the **J-OB**CD for optimization markedly enhances the objective value. The 1st, 2nd, and 3rd best results are colored with **red**, **green** and **blue**, respectively. (n, p) represents the dimension and p-value of the **J** orthogonal matrix (square matrix). The value in $()$ stands for $\sum_{ij}^n |X^T J X - J|_{ij}$.

datasetname	(m-n-p)	ADMM	UMCM	CSDM	J-OB	VR-J-OB	CSDM+J-OB
time limit=30s							
20News	(9423-50-25)	6.47e+04(2.4e-03)	6.47e+04(2.4e-03)	6.45e+04(2.0e-03)	6.40e+04(2.4e-03)	6.26e+04(2.4e-03)	6.31e+04(2.0e-03)(+)
Cifar	(10000-50-25)	1.49e+07(2.4e-03)	1.49e+07(2.4e-03)	1.49e+07(2.0e-03)	1.46e+07(2.4e-03)	1.45e+07(2.5e-03)	1.46e+07(2.0e-03)(+)
cmnCaltech	(3000-96-48)	3.08e+04(3.5e-09)	3.08e+04(3.5e-09)	3.03e+04(2.3e-05)	1.64e+04(4.0e-08)	1.64e+04(4.7e-08)	1.74e+04(2.3e-05)(+)
E2006	(5000-100-50)	5.54e+04(3.5e-09)	5.54e+04(3.5e-09)	5.54e+04(1.9e-05)	5.53e+04(9.0e-08)	5.53e+04(1.2e-07)	5.53e+04(1.9e-05)(+)
gisette	(6000-50-25)	1.61e+05(2.4e-03)	1.61e+05(2.4e-03)	1.52e+05(1.6e-03)	9.06e+04(2.4e-03)	7.48e+04(2.4e-03)	8.65e+04(1.6e-03)(+)
Mnist	(6000-92-46)	6.94e+06(3.6e-09)	6.94e+06(3.6e-09)	6.93e+06(2.1e-05)	6.50e+06(4.0e-05)	6.48e+06(8.7e-05)	6.92e+06(2.1e-05)
news20	(7967-50-25)	6.12e+04(2.4e-03)	6.12e+04(2.4e-03)	6.12e+04(2.0e-03)	6.12e+04(2.4e-03)	6.12e+04(2.4e-03)	6.12e+04(2.0e-03)
randn5000	(5000-100-50)	1.84e+06(3.5e-09)	1.84e+06(3.5e-09)	1.84e+06(2.0e-05)	1.77e+06(1.4e-07)	1.56e+06(9.9e-05)	1.84e+06(2.0e-05)
randn10000	(10000-50-25)	2.66e+06(2.4e-03)	2.66e+06(2.4e-03)	2.66e+06(2.5e-03)	2.64e+06(2.4e-03)	2.34e+06(2.5e-03)	2.66e+06(2.5e-03)
w1a	(2477-100-50)	2.71e+04(3.5e-09)	2.71e+04(3.5e-09)	2.47e+04(1.6e-05)	2.05e+04(3.6e-08)	1.92e+04(5.1e-08)	2.47e+04(1.6e-05)
20News	(9423-50-35)	7.89e+04(3.7e-09)	7.89e+04(3.7e-09)	7.85e+04(2.4e-05)	7.83e+04(8.6e-09)	7.71e+04(2.7e-08)	7.64e+04(2.4e-05)(+)
Cifar	(10000-50-35)	1.48e+07(5.5e-03)	1.48e+07(5.5e-03)	1.48e+07(6.2e-03)	1.43e+07(5.5e-03)	1.43e+07(5.7e-03)	1.43e+07(6.2e-03)(+)
cmnCaltech	(3000-96-70)	3.75e+04(2.0e-09)	3.75e+04(2.0e-09)	3.71e+04(1.3e-05)	2.40e+04(2.5e-08)	1.87e+04(3.4e-08)	2.15e+04(1.3e-05)(+)
E2006	(5000-100-75)	6.79e+04(1.0e-03)	6.79e+04(6.3e-04)	6.79e+04(5.8e-04)	6.65e+04(6.3e-04)	6.34e+04(6.3e-04)	6.69e+04(5.8e-04)(+)
gisette	(6000-50-35)	1.89e+05(5.5e-03)	1.89e+05(5.5e-03)	1.83e+05(4.7e-03)	1.24e+05(5.5e-03)	1.03e+05(5.5e-03)	1.17e+05(4.7e-03)(+)
Mnist	(6000-92-70)	6.75e+06(2.0e-09)	6.75e+06(2.0e-09)	6.74e+06(1.4e-05)	6.30e+06(7.7e-05)	6.20e+06(3.7e-04)	6.24e+06(7.0e-05)(+)
news20	(7967-50-35)	7.26e+04(5.5e-03)	7.26e+04(5.5e-03)	7.26e+04(4.9e-03)	7.26e+04(5.5e-03)	7.25e+04(5.5e-03)	7.26e+04(4.9e-03)
randn5000	(5000-100-75)	1.75e+06(6.3e-04)	1.75e+06(6.3e-04)	1.75e+06(6.1e-04)	1.44e+06(6.6e-04)	1.70e+06(6.3e-04)	1.75e+06(6.1e-04)
randn10000	(10000-50-36)	2.56e+06(3.7e-09)	2.56e+06(3.7e-09)	2.56e+06(2.0e-05)	2.56e+06(5.9e-09)	2.54e+06(2.2e-08)	2.56e+06(2.0e-05)(+)
w1a	(2477-100-75)	3.36e+04(6.3e-04)	3.36e+04(6.3e-04)	3.16e+04(5.5e-04)	2.53e+04(6.3e-04)	2.48e+04(6.3e-04)	2.90e+04(5.5e-04)(+)
20News	(9423-50-45)	8.79e+04(4.8e-03)	8.79e+04(4.8e-03)	8.73e+04(4.6e-03)	8.62e+04(4.8e-03)	8.60e+04(4.8e-03)	8.68e+04(4.6e-03)(+)
Cifar	(10000-50-45)	1.47e+07(4.8e-03)	1.47e+07(4.8e-03)	1.47e+07(4.6e-03)	1.46e+07(4.8e-03)	1.41e+07(7.2e-03)	1.46e+07(4.6e-03)(+)
cmnCaltech	(3000-96-85)	4.13e+04(9.2e-04)	4.14e+04(4.4e-04)	4.11e+04(6.5e-04)	2.59e+04(4.4e-04)	2.25e+04(4.4e-04)	2.65e+04(6.5e-04)(+)
E2006	(5000-100-90)	7.44e+04(1.2e-09)	7.44e+04(1.2e-09)	7.44e+04(7.5e-06)	6.95e+04(6.5e-08)	6.92e+04(3.8e-07)	6.93e+04(7.5e-06)(+)
gisette	(6000-50-45)	2.25e+05(4.8e-03)	2.25e+05(4.8e-03)	2.10e+05(5.2e-03)	1.42e+05(4.8e-03)	1.50e+05(4.8e-03)	1.54e+05(5.2e-03)(+)
Mnist	(6000-92-85)	6.73e+06(1.0e-03)	6.73e+06(1.0e-03)	6.73e+06(1.1e-03)	6.67e+06(1.0e-03)	6.22e+06(1.1e-03)	6.61e+06(1.1e-03)(+)
news20	(7967-50-45)	8.24e+04(4.8e-03)	8.24e+04(4.8e-03)	8.24e+04(4.5e-03)	8.24e+04(4.8e-03)	8.24e+04(4.8e-03)	8.24e+04(4.5e-03)
randn5000	(5000-100-85)	1.72e+06(1.3e-03)	1.72e+06(1.3e-03)	1.72e+06(1.4e-03)	1.70e+06(1.3e-03)	1.56e+06(1.3e-03)	1.72e+06(1.4e-03)
randn10000	(10000-50-45)	2.54e+06(3.3e-09)	2.54e+06(3.3e-09)	2.54e+06(1.9e-05)	2.53e+06(4.9e-09)	2.46e+06(3.3e-08)	2.54e+06(1.9e-05)
w1a	(2477-100-90)	4.10e+04(1.2e-09)	4.10e+04(1.2e-09)	3.76e+04(6.4e-06)	3.77e+04(4.2e-09)	3.27e+04(1.1e-08)	3.48e+04(6.4e-06)(+)
time limit=60s							
20News	(9423-50-25)	6.47e+04(2.4e-03)	6.47e+04(2.4e-03)	6.42e+04(1.9e-03)	6.32e+04(2.4e-03)	6.35e+04(2.4e-03)	6.29e+04(1.9e-03)(+)
Cifar	(10000-50-25)	1.49e+07(2.4e-03)	1.49e+07(2.4e-03)	1.49e+07(1.7e-03)	1.45e+07(2.5e-03)	1.45e+07(2.5e-03)	1.45e+07(1.7e-03)(+)
cmnCaltech	(3000-96-48)	3.08e+04(3.5e-09)	3.08e+04(3.5e-09)	3.00e+04(2.4e-05)	1.92e+04(3.8e-08)	1.52e+04(3.8e-08)	1.67e+04(2.4e-05)(+)
E2006	(5000-100-50)	5.54e+04(3.5e-09)	5.54e+04(3.5e-09)	5.54e+04(2.3e-05)	5.51e+04(1.6e-07)	5.23e+04(3.9e-07)	5.51e+04(2.3e-05)(+)
gisette	(6000-50-25)	1.61e+05(2.4e-03)	1.61e+05(2.4e-03)	1.46e+05(2.1e-03)	9.76e+04(2.4e-03)	7.20e+04(2.4e-03)	8.40e+04(2.1e-03)(+)
Mnist	(6000-92-46)	6.94e+06(3.6e-09)	6.94e+06(3.6e-09)	6.93e+06(2.1e-05)	6.47e+06(1.4e-04)	6.45e+06(3.0e-04)	6.93e+06(2.1e-05)
news20	(7967-50-25)	6.12e+04(2.4e-03)	6.12e+04(2.4e-03)	6.12e+04(2.0e-03)	6.12e+04(2.4e-03)	6.12e+04(2.4e-03)	6.12e+04(2.0e-03)
randn5000	(5000-100-50)	1.84e+06(3.5e-09)	1.84e+06(3.5e-09)	1.84e+06(1.9e-05)	1.59e+06(8.3e-05)	1.51e+06(4.9e-04)	1.84e+06(1.9e-05)
randn10000	(10000-50-25)	2.66e+06(2.4e-03)	2.66e+06(2.4e-03)	2.66e+06(2.5e-03)	2.58e+06(2.4e-03)	2.37e+06(2.4e-03)	2.66e+06(2.5e-03)
w1a	(2477-100-50)	2.71e+04(3.5e-09)	2.71e+04(3.5e-09)	2.41e+04(1.5e-05)	2.33e+04(1.3e-08)	1.89e+04(6.0e-08)	2.41e+04(1.5e-05)
20News	(9423-50-35)	7.89e+04(3.7e-09)	7.89e+04(3.7e-09)	7.84e+04(2.1e-05)	7.63e+04(3.3e-08)	7.60e+04(5.3e-08)	7.84e+04(2.1e-05)
Cifar	(10000-50-35)	1.48e+07(5.5e-03)	1.48e+07(5.5e-03)	1.48e+07(6.5e-03)	1.43e+07(5.6e-03)	1.42e+07(5.7e-03)	1.43e+07(6.5e-03)(+)
cmnCaltech	(3000-96-70)	3.75e+04(2.0e-09)	3.75e+04(2.0e-09)	3.69e+04(1.1e-05)	1.93e+04(3.2e-08)	1.89e+04(3.6e-08)	2.19e+04(1.1e-05)(+)
E2006	(5000-100-75)	6.79e+04(1.0e-03)	6.79e+04(6.3e-04)	6.79e+04(5.6e-04)	6.48e+04(6.3e-04)	6.36e+04(6.3e-04)	6.39e+04(5.6e-04)(+)
gisette	(6000-50-35)	1.89e+05(5.5e-03)	1.89e+05(5.5e-03)	1.80e+05(4.4e-03)	1.14e+05(5.5e-03)	9.97e+04(5.5e-03)	1.24e+05(4.4e-03)(+)
Mnist	(6000-92-70)	6.75e+06(2.0e-09)	6.75e+06(2.0e-09)	6.74e+06(1.3e-05)	6.21e+06(5.4e-04)	6.17e+06(2.0e-04)	6.21e+06(3.4e-04)(+)
news20	(7967-50-35)	7.26e+04(5.5e-03)	7.26e+04(5.5e-03)	7.26e+04(4.8e-03)	7.25e+04(5.5e-03)	7.25e+04(5.5e-03)	7.25e+04(4.8e-03)(+)
randn5000	(5000-100-75)	1.75e+06(6.3e-04)	1.75e+06(6.3e-04)	1.75e+06(6.1e-04)	1.44e+06(7.2e-04)	1.35e+06(9.9e-04)	1.75e+06(6.1e-04)
randn10000	(10000-50-36)	2.56e+06(3.7e-09)	2.56e+06(3.7e-09)	2.56e+06(2.0e-05)	2.56e+06(5.4e-09)	2.41e+06(2.6e-07)	2.56e+06(2.0e-05)
w1a	(2477-100-75)	3.36e+04(6.3e-04)	3.36e+04(6.3e-04)	3.10e+04(5.4e-04)	3.32e+04(6.3e-04)	2.57e+04(6.3e-04)	3.10e+04(5.4e-04)
20News	(9423-50-45)	8.79e+04(4.8e-03)	8.79e+04(4.8e-03)	8.71e+04(4.5e-03)	8.73e+04(4.8e-03)	8.57e+04(4.8e-03)	8.66e+04(4.5e-03)(+)
Cifar	(10000-50-45)	1.47e+07(4.8e-03)	1.47e+07(4.8e-03)	1.47e+07(4.5e-03)	1.46e+07(4.8e-03)	1.46e+07(4.8e-03)	1.46e+07(4.5e-03)(+)
cmnCaltech	(3000-96-85)	4.13e+04(9.2e-04)	4.14e+04(4.4e-04)	4.09e+04(7.2e-04)	2.69e+04(4.4e-04)	2.44e+04(4.4e-04)	2.40e+04(7.2e-04)(+)
E2006	(5000-100-90)	7.44e+04(1.2e-09)	7.44e+04(1.2e-09)	7.44e+04(6.1e-06)	6.92e+04(8.2e-08)	6.90e+04(1.3e-07)	6.94e+04(6.1e-06)(+)
gisette	(6000-50-45)	2.25e+05(4.8e-03)	2.25e+05(4.8e-03)	2.10e+05(5.2e-03)	1.63e+05(4.8e-03)	1.37e+05(4.8e-03)	2.10e+05(5.2e-03)(+)
Mnist	(6000-92-85)	6.73e+06(1.0e-03)	6.73e+06(1.0e-03)	6.73e+06(1.1e-03)	6.65e+06(1.0e-03)	6.12e+06(1.1e-03)	6.73e+06(1.1e-03)
news20	(7967-50-45)	8.24e+04(4.8e-03)	8.24e+04(4.8e-03)	8.24e+04(4.4e-03)	8.24e+04(4.8e-03)	8.24e+04(4.8e-03)	8.24e+04(4.4e-03)
randn5000	(5000-100-85)	1.72e+06(1.3e-03)	1.72e+06(1.3e-03)	1.72e+06(1.4e-03)	1.65e+06(1.3e-03)	1.56e+06(1.3e-03)	1.72e+06(1.4e-03)
randn10000	(10000-50-45)	2.54e+06(3.3e-09)	2.54e+06(3.3e-09)	2.54e+06(1.9e-05)	2.53e+06(9.0e-09)	2.34e+06(2.5e-07)	2.54e+06(1.9e-05)
w1a	(2477-100-90)	4.10e+04(1.2e-09)	4.10e+04(1.2e-09)	3.64e+04(6.0e-06)	3.45e+04(1.0e-08)	3.16e+04(1.0e-08)	3.34e+04(6.0e-06)(+)
time limit=90s							
20News	(9423-50-25)	6.47e+04(2.4e-03)	6.47e+04(2.4e-03)	6.41e+04(1.8e-03)	6.31e+04(2.4e-03)	6.25e+04(2.4e-03)	6.41e+04(1.8e-03)
Cifar	(10000-50-25)	1.49e+07(2.4e-03)	1.49e+07(2.4e-03)	1.49e+07(2.0e-03)	1.45e+07(2.6e-03)	1.45e+07(2.6e-03)	1.45e+07(2.0e-03)(+)
cmnCaltech	(3000-96-48)	3.08e+04(3.5e-09)	3.08e+04(3.5e-09)	2.98e+04(2.2e-05)	1.89e+04(4.0e-08)	1.52e+04(4.5e-08)	1.74e+04(2.2e-05)(+)
E2006	(5000-100-50)	5.54e+04(3.5e-09)	5.54e+04(3.5e-09)	5.54e+04(2.1e-05)	5.34e+04(2.7e-07)	5.22e+04(3.0e-07)	5.30e+04(2.1e-05)(+)
gisette	(6000-50-25)	1.61e+05(2.4e-03)	1.61e+05(2.4e-03)	1.42e+05(2.3e-03)	9.50e+04(2.4e-03)	8.14e+04(2.4e-03)	7.44e+04(2.3e-03)(+)
Mnist	(6000-92-46)	6.94e+06(3.6e-09)	6.94e+06(3.6e-09)	6.93e+06(2.0e-05)	6.45e+06(3.6e-04)	6.42e+06(1.4e-03)	6.93e+06(2.0e-05)
news20	(7967-50-25)	6.12e+04(2.4e-03)	6.12e+04(2.4e-03)	6.12e+04(1.9e-03)	6.12e+04(2.4e-03)	6.12e+04(2.4e-03)	6.12e+04(1.9e-03)
randn5000	(5000-100-50)	1.84e+06(3.5e-09)	1.84e+06(3.5e-09)	1.84e+06(1.9e-05)	1.63e+06(1.4e-05)	1.52e+06(6.9e-04)	1.84e+06(1.9e-05)
randn10000	(10000-50-25)	2.66e+06(2.4e-03)	2.66e+06(2.4e-03)	2.66e+06(2.5e-03)	2.65e+06(2.4e-03)	2.25e+06(2.7e-03)	2.66e+06(2.5e-03)
w1a	(2477-100-50)	2.71e+04(3.5e-09)	2.71e+04(3.5e-09)	2.34e+04(1.5e-05)	2.24e+04(1.9e-08)	1.50e+04(3.4e-08)	1.90e+04(1.5e-05)(+)
20News	(9423-50-35)	7.89e+04(3.7e-09)	7.89e+04(3.7e-09)	7.84e+04(2.3e-05)	7.66e+04(4.5e-08)	7.57e+04(1.6e-07)	7.70e+04(2.3e-05)(+)
Cifar	(10000-50-35)	1.48e+07(5.5e-03)	1.48e+07(5.5e-03)	1.48e+07(7.0e-03)	1.43e+07(5.6e-03)	1.42	

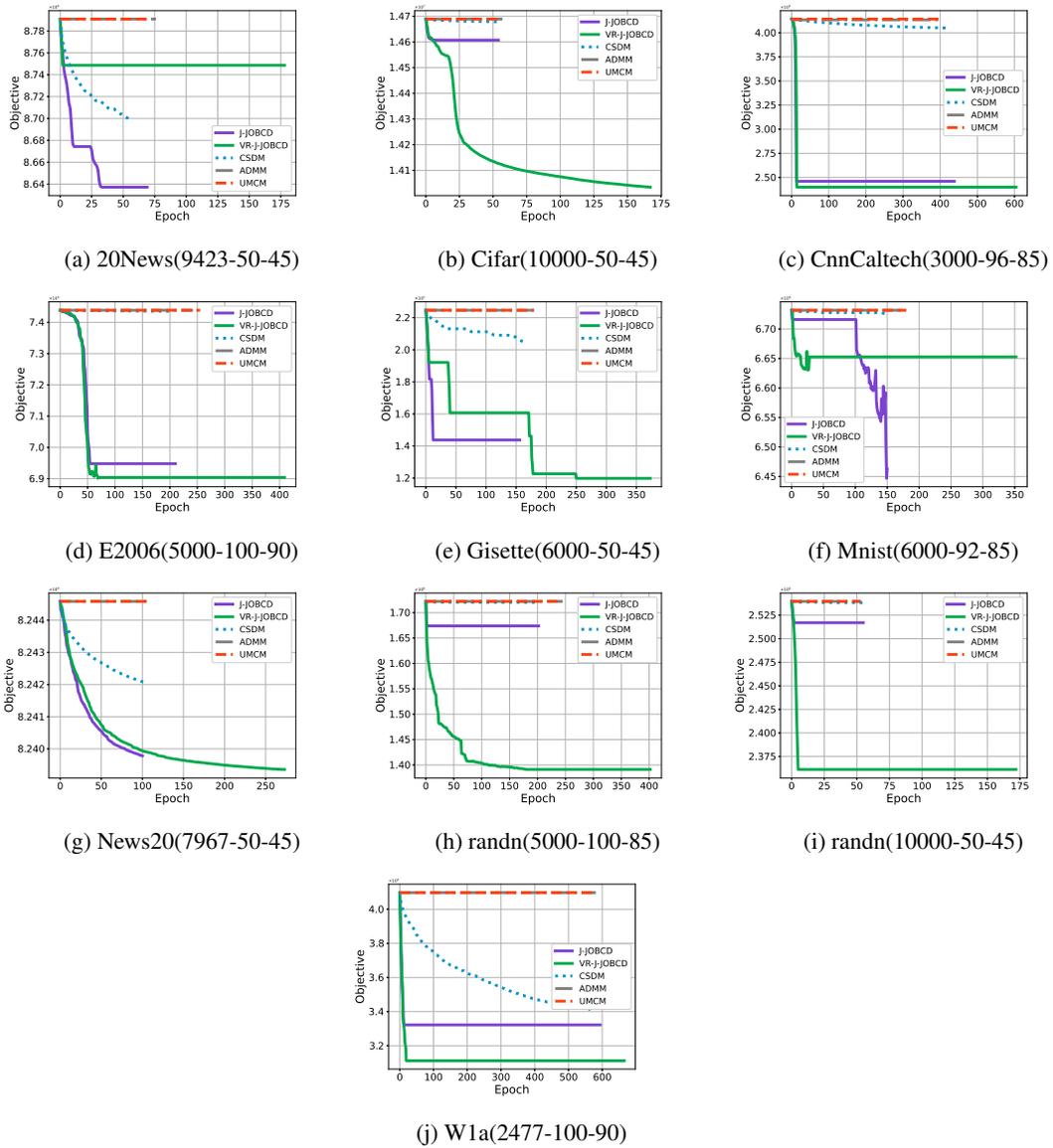


Figure 6: The convergence curve of the compared methods for solving HSPP by epochs with varying (m, n, p) .

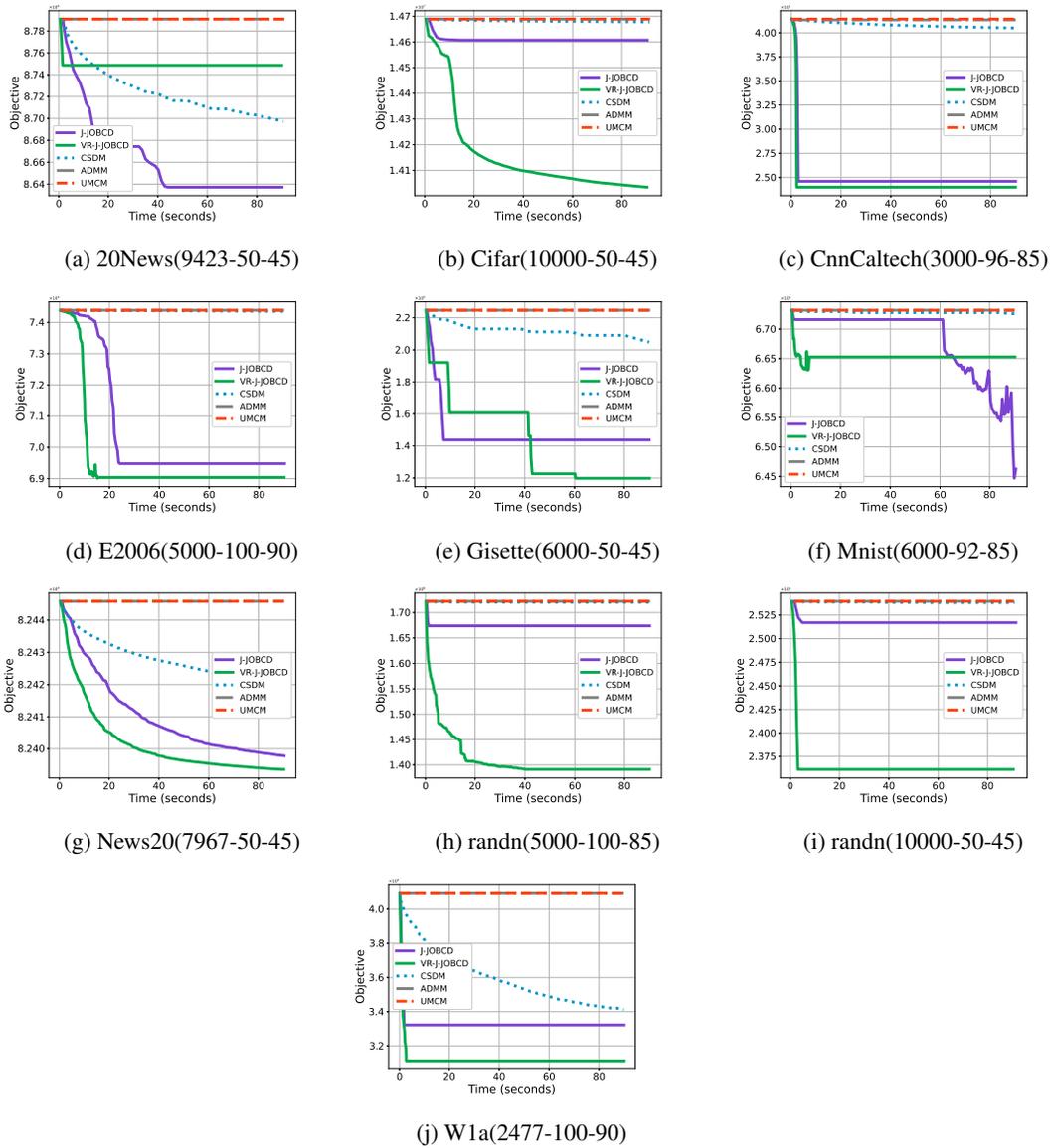


Figure 7: The convergence curve of the compared methods for solving HSPP by time with varying (m, n, p) .

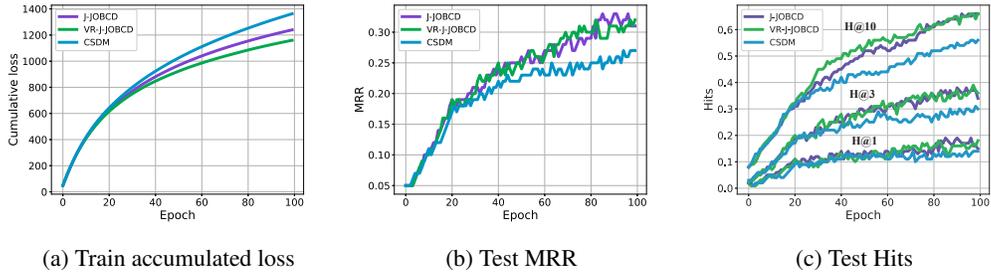


Figure 8: Epoch performance of CS, J-JOBBCD, and VR-J-JOBBCD in training UltraE on FB15k.

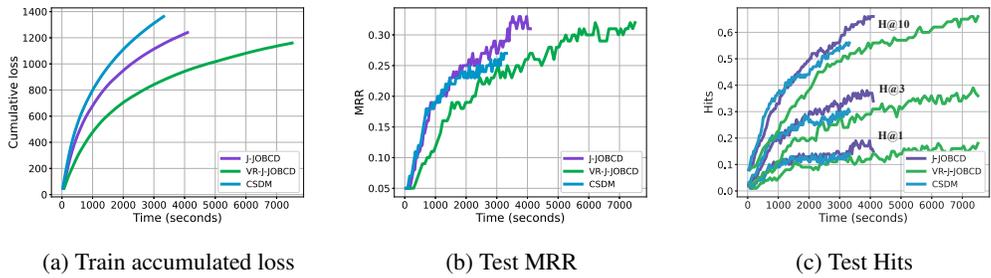


Figure 9: Time performance of CS, J-JOBBCD, and VR-J-JOBBCD in training UltraE on FB15k.

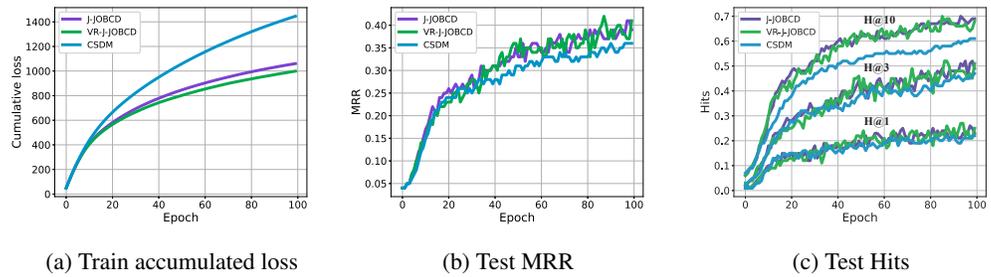


Figure 10: Epoch performance of CSDM, J-JOBBCD, and VR-J-JOBBCD in training UltraE on WN18RR.

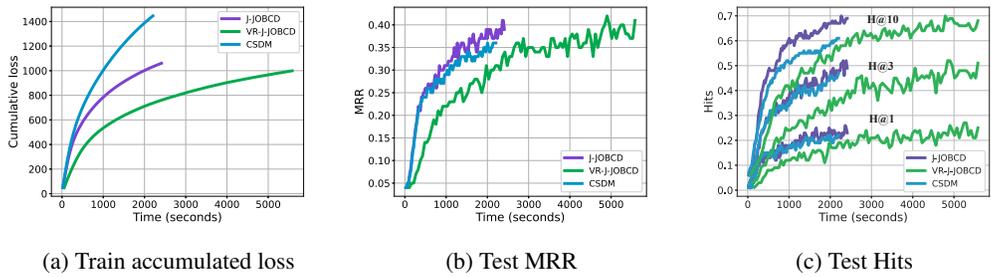


Figure 11: Time performance of CSDM, J-JOBBCD, and VR-J-JOBBCD in training UltraE on WN18RR.

837 **NeurIPS Paper Checklist**

838 **1. Claims**

839 Question: Do the main claims made in the abstract and introduction accurately reflect the
840 paper's contributions and scope?

841 Answer: [Yes]

842 Justification: In the abstract, we highlighted our contributions, including algorithm develop-
843 ment, theoretical analysis, and empirical study.

844 Guidelines:

- 845 • The answer NA means that the abstract and introduction do not include the claims
846 made in the paper.
- 847 • The abstract and/or introduction should clearly state the claims made, including the
848 contributions made in the paper and important assumptions and limitations. A No or
849 NA answer to this question will not be perceived well by the reviewers.
- 850 • The claims made should match theoretical and experimental results, and reflect how
851 much the results can be expected to generalize to other settings.
- 852 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
853 are not attained by the paper.

854 **2. Limitations**

855 Question: Does the paper discuss the limitations of the work performed by the authors?

856 Answer: [Yes]

857 Justification: Please refer to the assumptions made for the optimization problem outlined in
858 the introduction and Section 4.

859 Guidelines:

- 860 • The answer NA means that the paper has no limitation while the answer No means that
861 the paper has limitations, but those are not discussed in the paper.
- 862 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 863 • The paper should point out any strong assumptions and how robust the results are to
864 violations of these assumptions (e.g., independence assumptions, noiseless settings,
865 model well-specification, asymptotic approximations only holding locally). The authors
866 should reflect on how these assumptions might be violated in practice and what the
867 implications would be.
- 868 • The authors should reflect on the scope of the claims made, e.g., if the approach was
869 only tested on a few datasets or with a few runs. In general, empirical results often
870 depend on implicit assumptions, which should be articulated.
- 871 • The authors should reflect on the factors that influence the performance of the approach.
872 For example, a facial recognition algorithm may perform poorly when image resolution
873 is low or images are taken in low lighting. Or a speech-to-text system might not be
874 used reliably to provide closed captions for online lectures because it fails to handle
875 technical jargon.
- 876 • The authors should discuss the computational efficiency of the proposed algorithms
877 and how they scale with dataset size.
- 878 • If applicable, the authors should discuss possible limitations of their approach to
879 address problems of privacy and fairness.
- 880 • While the authors might fear that complete honesty about limitations might be used by
881 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
882 limitations that aren't acknowledged in the paper. The authors should use their best
883 judgment and recognize that individual actions in favor of transparency play an impor-
884 tant role in developing norms that preserve the integrity of the community. Reviewers
885 will be specifically instructed to not penalize honesty concerning limitations.

886 **3. Theory Assumptions and Proofs**

887 Question: For each theoretical result, does the paper provide the full set of assumptions and
888 a complete (and correct) proof?

889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942

Answer: [\[Yes\]](#)

Justification: We have added a hyperlink before each theoretical result, which points to the complete proof located in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We have provided sufficient details for reproducing the results of the paper, such as parameter settings, runtime environments, and dataset descriptions.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

943 Question: Does the paper provide open access to the data and code, with sufficient instruc-
944 tions to faithfully reproduce the main experimental results, as described in supplemental
945 material?

946 Answer: [Yes]

947 Justification: We have included all the code and data in the supplemental materials.

948 Guidelines:

- 949 • The answer NA means that paper does not include experiments requiring code.
- 950 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
951 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 952 • While we encourage the release of code and data, we understand that this might not be
953 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
954 including code, unless this is central to the contribution (e.g., for a new open-source
955 benchmark).
- 956 • The instructions should contain the exact command and environment needed to run to
957 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
958 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 959 • The authors should provide instructions on data access and preparation, including how
960 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 961 • The authors should provide scripts to reproduce all experimental results for the new
962 proposed method and baselines. If only a subset of experiments are reproducible, they
963 should state which ones are omitted from the script and why.
- 964 • At submission time, to preserve anonymity, the authors should release anonymized
965 versions (if applicable).
- 966 • Providing as much information as possible in supplemental material (appended to the
967 paper) is recommended, but including URLs to data and code is permitted.

968 6. Experimental Setting/Details

969 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
970 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
971 results?

972 Answer: [Yes]

973 Justification: We have provided sufficient details for solving the optimization problem,
974 encompassing hyperparameter settings and dataset generation.

975 Guidelines:

- 976 • The answer NA means that the paper does not include experiments.
- 977 • The experimental setting should be presented in the core of the paper to a level of detail
978 that is necessary to appreciate the results and make sense of them.
- 979 • The full details can be provided either with the code, in appendix, or as supplemental
980 material.

981 7. Experiment Statistical Significance

982 Question: Does the paper report error bars suitably and correctly defined or other appropriate
983 information about the statistical significance of the experiments?

984 Answer: [No]

985 Justification: For simplicity, we only demonstrate the convergence behavior of the objective
986 function by varying the time or iterations. Our methods exhibit clear advantages over the
987 compared methods. Such results have demonstrated significance in the experiments.

988 Guidelines:

- 989 • The answer NA means that the paper does not include experiments.
- 990 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
991 dence intervals, or statistical significance tests, at least for the experiments that support
992 the main claims of the paper.

- 993 • The factors of variability that the error bars are capturing should be clearly stated (for
994 example, train/test split, initialization, random drawing of some parameter, or overall
995 run with given experimental conditions).
- 996 • The method for calculating the error bars should be explained (closed form formula,
997 call to a library function, bootstrap, etc.)
- 998 • The assumptions made should be given (e.g., Normally distributed errors).
- 999 • It should be clear whether the error bar is the standard deviation or the standard error
1000 of the mean.
- 1001 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1002 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1003 of Normality of errors is not verified.
- 1004 • For asymmetric distributions, the authors should be careful not to show in tables or
1005 figures symmetric error bars that would yield results that are out of range (e.g. negative
1006 error rates).
- 1007 • If error bars are reported in tables or plots, The authors should explain in the text how
1008 they were calculated and reference the corresponding figures or tables in the text.

1009 8. Experiments Compute Resources

1010 Question: For each experiment, does the paper provide sufficient information on the com-
1011 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1012 the experiments?

1013 Answer: [Yes]

1014 Justification:

1015 Guidelines: We have outlined the types of compute workers, detailing CPU and memory
1016 specifications.

- 1017 • The answer NA means that the paper does not include experiments.
- 1018 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1019 or cloud provider, including relevant memory and storage.
- 1020 • The paper should provide the amount of compute required for each of the individual
1021 experimental runs as well as estimate the total compute.
- 1022 • The paper should disclose whether the full research project required more compute
1023 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1024 didn't make it into the paper).

1025 9. Code Of Ethics

1026 Question: Does the research conducted in the paper conform, in every respect, with the
1027 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1028 Answer: [Yes]

1029 Justification: Our research aligns with the ethical guidelines outlined by NeurIPS.

1030 Guidelines:

- 1031 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1032 • If the authors answer No, they should explain the special circumstances that require a
1033 deviation from the Code of Ethics.
- 1034 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1035 eration due to laws or regulations in their jurisdiction).

1036 10. Broader Impacts

1037 Question: Does the paper discuss both potential positive societal impacts and negative
1038 societal impacts of the work performed?

1039 Answer: [NA]

1040 Justification: The paper addresses theoretical questions on algorithm complexity, which, to
1041 the best of our knowledge, pose no negative social impact.

1042 Guidelines:

- 1043 • The answer NA means that there is no societal impact of the work performed.

- 1044 • If the authors answer NA or No, they should explain why their work has no societal
1045 impact or why the paper does not address societal impact.
- 1046 • Examples of negative societal impacts include potential malicious or unintended uses
1047 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1048 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1049 groups), privacy considerations, and security considerations.
- 1050 • The conference expects that many papers will be foundational research and not tied
1051 to particular applications, let alone deployments. However, if there is a direct path to
1052 any negative applications, the authors should point it out. For example, it is legitimate
1053 to point out that an improvement in the quality of generative models could be used to
1054 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1055 that a generic algorithm for optimizing neural networks could enable people to train
1056 models that generate Deepfakes faster.
- 1057 • The authors should consider possible harms that could arise when the technology is
1058 being used as intended and functioning correctly, harms that could arise when the
1059 technology is being used as intended but gives incorrect results, and harms following
1060 from (intentional or unintentional) misuse of the technology.
- 1061 • If there are negative societal impacts, the authors could also discuss possible mitigation
1062 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1063 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1064 feedback over time, improving the efficiency and accessibility of ML).

1065 11. Safeguards

1066 Question: Does the paper describe safeguards that have been put in place for responsible
1067 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1068 image generators, or scraped datasets)?

1069 Answer: [NA]

1070 Justification: The paper poses no such risks.

1071 Guidelines:

- 1072 • The answer NA means that the paper poses no such risks.
- 1073 • Released models that have a high risk for misuse or dual-use should be released with
1074 necessary safeguards to allow for controlled use of the model, for example by requiring
1075 that users adhere to usage guidelines or restrictions to access the model or implementing
1076 safety filters.
- 1077 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1078 should describe how they avoided releasing unsafe images.
- 1079 • We recognize that providing effective safeguards is challenging, and many papers do
1080 not require this, but we encourage authors to take this into account and make a best
1081 faith effort.

1082 12. Licenses for existing assets

1083 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1084 the paper, properly credited and are the license and terms of use explicitly mentioned and
1085 properly respected?

1086 Answer: [Yes]

1087 Justification: The dataset used in the experiments is published on an open site without
1088 license.

1089 Guidelines:

- 1090 • The answer NA means that the paper does not use existing assets.
- 1091 • The authors should cite the original paper that produced the code package or dataset.
- 1092 • The authors should state which version of the asset is used and, if possible, include a
1093 URL.
- 1094 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1095 • For scraped data from a particular source (e.g., website), the copyright and terms of
1096 service of that source should be provided.

- 1097
- 1098
- 1099
- 1100
- 1101
- 1102
- 1103
- 1104
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

1105 13. **New Assets**

1106 Question: Are new assets introduced in the paper well documented and is the documentation
1107 provided alongside the assets?

1108 Answer: [NA]

1109 Justification: The experiments do not involve new datasets.

1110 Guidelines:

- 1111
- 1112
- 1113
- 1114
- 1115
- 1116
- 1117
- 1118
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

1119 14. **Crowdsourcing and Research with Human Subjects**

1120 Question: For crowdsourcing experiments and research with human subjects, does the paper
1121 include the full text of instructions given to participants and screenshots, if applicable, as
1122 well as details about compensation (if any)?

1123 Answer: [NA]

1124 Justification: No crowdsourcing or human object is involved.

1125 Guidelines:

- 1126
- 1127
- 1128
- 1129
- 1130
- 1131
- 1132
- 1133
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

1134 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1135 Subjects**

1136 Question: Does the paper describe potential risks incurred by study participants, whether
1137 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1138 approvals (or an equivalent approval/review based on the requirements of your country or
1139 institution) were obtained?

1140 Answer: [NA]

1141 Justification: No crowdsourcing or human object is involved.

1142 Guidelines:

- 1143
- 1144
- 1145
- 1146
- 1147
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

1148
1149
1150
1151
1152

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.