

---

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012

# BLOCK COORDINATE DESCENT METHODS FOR OPTIMIZATION UNDER J-ORTHOGONALITY CON- STRAINTS WITH APPLICATIONS

006      **Anonymous authors**

007      Paper under double-blind review

## ABSTRACT

013      The J-orthogonal matrix, also referred to as the hyperbolic orthogonal ma-  
014      trix, is a class of special orthogonal matrix in hyperbolic space, notable for  
015      its advantageous properties. These matrices are integral to optimization un-  
016      der J-orthogonal constraints, which have widespread applications in statisti-  
017      cal learning and data science. However, addressing these problems is gener-  
018      ally challenging due to their non-convex nature and the computational inten-  
019      sity of the constraints. Currently, algorithms for tackling these challenges  
020      are limited. This paper introduces **JOBCD**, a novel Block Coordinate Descent  
021      method designed to address optimizations with J-orthogonality constraints.  
022      We explore two specific variants of **JOBCD**: one based on a Gauss-Seidel strategy  
023      (**GS-JOBCD**), the other on a variance-reduced and Jacobi strategy (**VR-J-JOBCD**). Notably,  
024      leveraging the parallel framework of a Jacobi strategy, **VR-J-JOBCD** integrates variance reduction  
025      techniques to decrease oracle complexity in the minimization of finite-sum  
026      functions. For both **GS-JOBCD** and **VR-J-JOBCD**, we establish the  
027      oracle complexity under mild conditions and strong limit-point convergence  
028      results under the Kurdyka-Łojasiewicz inequality. To demonstrate the effec-  
029      tiveness of our method, we conduct experiments on hyperbolic eigenvalue  
030      problems, hyperbolic structural probe problems, and the ultrahyperbolic  
031      knowledge graph embedding problem. Extensive experiments using both  
032      real-world and synthetic data demonstrate that **JOBCD** consistently out-  
033      performs state-of-the-art solutions, by large margins.

## 1 INTRODUCTION

034  
035  
036      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  $\mathbf{J} = [\begin{smallmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_{n-p} \end{smallmatrix}]$ , and  $\mathbf{I}_p$  is  
037      a  $p \times p$  identity matrix. Here,  $\mathbf{J} \in \mathbb{R}^{n \times n}$  is the signature matrix with signature  $(p, n-p)$ .  
038      In this paper, we mainly focus on the following optimization problem under J-orthogonality  
039      constraints:

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

040  
041      Here,  $f(\mathbf{X})$  could have a finite-sum structure, each component function  $f_i(\mathbf{X})$  is assumed  
042      to be differentiable, and  $N$  is the number of data points. For brevity, the J-orthogonality  
043      constraint  $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$  in Problem (1) is rewritten as  $\mathbf{X} \in \mathcal{J}$ .

044  
045      We impose the following assumptions on Problem (1) throughout this paper. (A-i) For  
046      any matrices  $\mathbf{X}$  and  $\mathbf{X}^+$ , we assume  $f_i : \mathbb{R}^{n \times n} \mapsto \mathbb{R}$  is continuously differentiable for some  
047      symmetric positive semidefinite matrix  $\mathbf{H} \in \mathbb{R}^{nn \times nn}$  that:

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

048  
049      for all  $i \in [N]$ , where  $\|\mathbf{H}\| \leq L_f$  for some constant  $L_f > 0$  and  $\|\mathbf{X}\|_{\mathbf{H}}^2 \triangleq \text{vec}(\mathbf{X})^\top \mathbf{H} \text{vec}(\mathbf{X})$ .  
050      Importantly, 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  
051      the equality  $\forall \mathbf{X}, \mathbf{X}^+, f(\mathbf{X}^+) = f(\mathbf{X}) + \langle \mathbf{X}^+ - \mathbf{X}, \nabla f(\mathbf{X}) \rangle + \frac{1}{2} \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbf{H}}^2$  in (2), where  
052       $\mathbf{C} \in \mathbb{R}^{n \times n}$  and  $\mathbf{D} \in \mathbb{R}^{n \times n}$  are arbitrary symmetric matrices. (A-ii) For any matrices  $\mathbf{X}$   
053      and  $\mathbf{X}^+$ , we assume that:  $\|\nabla f_i(\mathbf{X}) - \nabla f_i(\mathbf{X}^+)\|_{\mathbf{F}} \leq L_f \|\mathbf{X} - \mathbf{X}^+\|_{\mathbf{F}}$  for all  $i \in [N]$  and

---

054  $L_f$  is mentioned in (A-i). (A-iii) The function  $f_i(\mathbf{X})$  is coercive for all  $i \in N$ , that is,  
055  $\lim_{\|\mathbf{X}\|_F \rightarrow \infty} f_i(\mathbf{X}) = \infty, \forall i$ .  
056

057 Problem (1) defines an optimization framework that is fundamental to a wide range of models  
058 in statistical learning and data science, including hyperbolic eigenvalue problem [6; 42; 39],  
059 hyperbolic structural probe problem [20; 7], and ultrahyperbolic knowledge graph embedding  
060 [49]. Additionally, it is closely related to machine learning in hyperbolic spaces, including  
061 Lorentz model learning [34; 51; 8] and ultrahyperbolic neural networks [27; 56; 41]. It also  
062 intersects with hyperbolic linear algebra [3; 21], addressing problems such as the indefinite  
063 least squares problem, hyperbolic QR factorization, and indefinite polar decomposition.  
064

## 1.1 RELATED WORK

065 **► Block Coordinate Descent Methods.** Block Coordinate Descent (BCD) is a well-  
066 established iterative algorithm that sequentially minimizes along block coordinate directions.  
067 Its simplicity and efficiency have led to its widespread adoption in structured convex appli-  
068 cations [36]. Recently, BCD has gained traction in non-convex problems due to its robust  
069 optimality guarantees and/or excellent empirical performance in areas including optimal  
070 transport [22], matrix optimization [12], fractional minimization [54], deep neural networks  
071 [5; 55; 31], federated learning[46], black-box optimization [4], and optimization with orthog-  
072 onality constraints [52; 14]. To our knowledge, this is the first application of BCD methods  
073 to optimization under J-orthogonality constraints, with a focus on analyzing their theoretical  
074 guarantees and empirical efficacy.

075 **► Minimizing Smooth Functions under J-Orthogonality Constraints.** The J-  
076 orthogonal matrix belongs to a subset of generalized orthogonal matrices [16; 35; 23]. How-  
077 ever, projecting onto the J-orthogonality constraint poses challenges, complicating the ex-  
078 tension of conventional optimization algorithms to address optimization problems under  
079 these constraints [1; 16]. This contrasts with computing orthogonal projections using meth-  
080 ods such as polar or SVD decomposition, or approximating them via QR factorization.  
081 Existing methods for addressing Problem (1) can be categorized into three classes. **(i)** CS-  
082 Decomposition Based Methods. These approaches involve parameterizing four orthogonal  
083 matrices (as described in Proposition 2.2) and subsequently minimizing a smooth function  
084 over these matrices in an alternating fashion. The involvement of  $3 \times 3$  block matrices  
085 makes the implementation of these methods very challenging. Consequently, the work of  
086 [49] focuses on optimizing a reduced subspace of the CS decomposition parameters, albeit  
087 at the expense of losing some degrees of freedom. **(ii)** Unconstrained Multiplier Correction  
088 Methods [47; 48; 13; 14]. These methods leverage the symmetry and explicit closed-form ex-  
089 pression of the Lagrangian multiplier at the first-order optimality condition. Consequently,  
090 they address an unconstrained problem, resulting in efficient first-order infeasible approaches.  
091 **(iii)** Alternating Direction Method of Multipliers [19]. This method reformulates the origi-  
092 nal problem into a bilinear constrained optimization problem by introducing auxiliary vari-  
093 ables. It employs dual variables to handle bilinear constraints, iteratively optimizing primal  
094 variables while keeping other primal and dual variables fixed, and using a gradient ascent  
095 strategy to update the dual variables. This approach has become widely adopted for solv-  
096 ing general nonconvex and nonsmooth composite optimization problems. Notably, all the  
097 aforementioned methods solely identify critical points of Problem (1).

098 **► Finite-Sum Problems via Stochastic Gradient Descent.** The finite-sum struc-  
099 ture is prevalent in machine learning and statistical modeling, facilitating decomposition  
100 into smaller, more manageable components. This property is advantageous for developing  
101 efficient algorithms for large-scale problems, such as Stochastic Gradient Descent (SGD). Re-  
102 ducing variance is crucial in SGD because it can lead to more stable and faster convergence.  
103 Various techniques, such as mini-batch SGD, momentum methods, and variance reduction  
104 methods like SAGA [10], SVRG [25], SARAH [33], SPIDER [11; 43], SNVRG [57], and  
105 PAGE [30], have been developed to address this issue. Additionally, SGD for minimizing  
106 composite functions has also been investigated by the authors [15; 24; 29].  
107

## 1.2 CONTRIBUTIONS

108 This paper makes the following contributions. **(i)** Algorithmically: We introduce the  
109 **JOB**CD algorithm, a novel Block Coordinate Descent method specifically designed to  
110 tackle optimizations constrained by J-orthogonality. We explore two specific variants of  
111

108 **JOBCD**, one based on a Gauss-Seidel strategy (**GS-JOBCD**), the other on a variance-  
 109 reduced and Jacobi strategy (**VR-J-JOBCD**). Notably, **VR-J-JOBCD** incorporates a  
 110 variance-reduction technique into a parallel framework to reduce oracle complexity in the  
 111 minimization of finite-sum functions (See Section 2). **(ii)** Theoretically: We provide com-  
 112 prehensive optimality and convergence analyses for both algorithms (see Sections 3 and 4).  
 113 **(iii)** Empirically: Extensive experiments across hyperbolic eigenvalue problems, structural  
 114 probe problems, and ultrahyperbolic knowledge graph embedding, using both real-world and  
 115 synthetic data, consistently show the significant superiority of **JOBCD** over state-of-the-art  
 116 solutions (see Section 6).

## 117 2 THE PROPOSED JOBCD ALGORITHM

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

123 **Notations.** We define  $[n] \triangleq \{1, 2, \dots, n\}$ . We denote  $\Omega \triangleq \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{C_n^2}\}$  as all the  
 124 possible combinations of the index vectors choosing 2 items from  $n$  without repetition. For  
 125 any  $B \in \Omega$ , we 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  
 126  $\mathbf{U}_B^\top \mathbf{X} = \mathbf{X}(B, :) \in \mathbb{R}^{2 \times n}$ . We denote  $\mathcal{J}_B \triangleq \{\mathbf{V} | \mathbf{V}^\top \mathbf{J}_{BB} \mathbf{V} = \mathbf{J}_{BB}\}$ , where  $\mathbf{J}_{BB} \in \mathbb{R}^{2 \times 2}$  is the  
 127 sub-matrix of  $\mathbf{J}$  indexed by  $B$ . Further notations are provided in Appendix A.1.

### 128 2.1 GAUSS-SEIDEL BLOCK COORDINATE DESCENT ALGORITHM

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

131 **GS-JOBCD** is an iterative algorithm that, in each iteration  $t$ , randomly and uniformly  
 132 (with replacement) selects a coordinate  $B^t$  from the set  $\Omega$  and then solves a small-sized  
 133 subproblem. The row index  $[n]$  of the decision variable  $\mathbf{X}^t$  are separated to two sets  $B^t$   
 134 and  $B^{t,c}$ , where  $B^t \in \Omega$  with  $|B^t| = 2$  is the working set and  $B^{t,c} = [n] \setminus B^t$ . For simplicity,  
 135 we use  $B$  instead of  $B^t$ . Following [52], we consider the following block coordinate update  
 136 rule:  $[\mathbf{X}^{t+1}(B, :) = \mathbf{V}\mathbf{X}^t(B, :)] \Leftrightarrow [\mathbf{X}^{t+1} = \mathbf{X}^t + \mathbf{U}_B(\mathbf{V} - \mathbf{I})\mathbf{U}_B^\top \mathbf{X}^t]$ , where  $\mathbf{V} \in \mathbb{R}^{2 \times 2}$  is some  
 137 suitable matrix.

138 The following lemma illustrates matrix selection for enforcing J-orthogonality constraints  
 139 via the 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.  
 140

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

145 **► The Main Algorithm.** Using the above update rule, we consider the following iterative  
 146 procedure:  $\mathbf{X}^{t+1} \Leftarrow \mathcal{X}_B^t(\bar{\mathbf{V}}^t)$ , where  $\bar{\mathbf{V}}^t \in \arg \min_{\mathbf{V}} f(\mathcal{X}_B^t(\mathbf{V}))$ . However, the resulting  
 147 subproblem could be still difficult to solve. This inspires us to use sequential majorization-  
 148 minimization [37; 32] to address it. This technique iteratively constructs a surrogate  
 149 function that upper-bounds the objective function, allowing for effective optimization and  
 150 gradual reduction of the objective function. We derive:

$$\begin{aligned} 152 \quad f(\mathcal{X}_B^t(\mathbf{V})) &\stackrel{\textcircled{1}}{\leq} f(\mathbf{X}^t) + \frac{1}{2} \|\mathcal{X}_B^t(\mathbf{V}) - \mathbf{X}^t\|_H^2 + \langle \mathcal{X}_B^t(\mathbf{V}) - \mathbf{X}^t, \nabla f(\mathbf{X}^t) \rangle \\ 153 &\stackrel{\textcircled{2}}{\leq} f(\mathbf{X}^t) + \frac{1}{2} \|\mathbf{V} - \mathbf{I}\|_{Q+\theta\mathbf{I}}^2 + \langle \mathbf{V} - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{BB} \rangle \triangleq \mathcal{G}(\mathbf{V}; \mathbf{X}^t, B^t), \end{aligned} \quad (3)$$

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

$$158 \quad \mathbf{Q} = \underline{Q}, \text{ or } \mathbf{Q} = \varsigma \mathbf{I}_4, \text{ with } \|\mathbf{Q}\| \leq \varsigma \leq L_f. \quad (4)$$

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

---

162  
 163   **Algorithm 1:** **GS-JOBCD**: Block Coordinate Descent Methods using a Gauss-Seidel  
 164   Strategy for Solving Problem (1)  
 165   **Init.:** Set  $\mathbf{X}^0$  to satisfy J-orthogonality constraints (e.g., via Hyperbolic CS  
 166   Decomposition),  $\theta$  in Inequality (3) (e.g., 1e-6).  
 167   **for**  $t$  from 0 to  $T$  **do**  
 168     (S1) Choose a coordinate  $B^t$  with  $|B^t| = 2$  from the set  $\Omega$  randomly and uniformly  
 169     (with replacement) for the  $t$ -th iteration. Denote  $B = B^t$ .  
 170     (S2) Choose a matrix  $\mathbf{Q} \in \mathbb{R}^{4 \times 4}$  using Formula (4).  
 171     (S3) Solve the following small-size subproblem globally.  
 172     
$$\bar{\mathbf{V}}^t \in \arg \min_{\mathbf{V} \in \mathcal{J}_B} \frac{1}{2} \|\mathbf{V} - \mathbf{I}\|_{\mathbf{Q} + \theta \mathbf{I}}^2 + \langle \mathbf{V} - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^T]_{BB} \rangle + f(\mathbf{X}^t) \quad (5)$$
  
 173     
$$= \arg \min_{\mathbf{V} \in \mathcal{J}_B \in \mathbb{R}^{2 \times 2}} \frac{1}{2} \|\mathbf{V}\|_{\dot{\mathbf{Q}}}^2 + \langle \mathbf{V}, \mathbf{P} \rangle + c \quad (6)$$
  
 174     where  $\mathbf{P} \triangleq [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^T]_{BB} - \text{mat}(\dot{\mathbf{Q}} \text{ vec}(\mathbf{I}_2))$ ,  $\dot{\mathbf{Q}} = \mathbf{Q} + \theta \mathbf{I}$  and  
 175      $c \triangleq f(\mathbf{X}^t) - \langle \mathbf{I}_2, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^T]_{BB} \rangle + \frac{1}{2} \|\mathbf{I}\|_{\dot{\mathbf{Q}}}^2$  is a constant.  
 176     (S4)  $\mathbf{X}^{t+1}(B, :) = \bar{\mathbf{V}}^t \mathbf{X}^t(B, :)$   
 177   **end**  
 178

---

We summarize the proposed **GS-JOBCD** in Algorithm 1.

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

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

**Proposition 2.2.** (Hyperbolic CS Decomposition [40]) Let  $\mathbf{V}$  be J-orthogonal with signature  $(p, n-p)$ . Assume that  $n-p \leq p$ . Then there exist vectors  $\dot{c}, \dot{s} \in \mathbb{R}^{n-p}$  with  $\dot{c} \odot \dot{c} - \dot{s} \odot \dot{s} = \mathbf{1}$ , 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:  $\mathbf{V} = [\begin{array}{cc} \mathbf{U}_1 & 0 \\ 0 & \mathbf{U}_2 \end{array}] [\begin{array}{ccc} \text{Diag}(\dot{c}) & 0 & \text{Diag}(\dot{s}) \\ 0 & I_{p-(n-p)} & 0 \\ \text{Diag}(\dot{s}) & 0 & \text{Diag}(\dot{c}) \end{array}] [\begin{array}{cc} \mathbf{V}_1^T & 0 \\ 0 & \mathbf{V}_2^T \end{array}]$ .

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

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

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

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

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

---

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

220  
221  
222    We now describe how to find the optimal solution  $\bar{t}_+$ , where  $\bar{t}_+ \in \arg \min_t p(t) \triangleq \frac{a+bt}{\sqrt{1-t^2}} +$   
223     $\frac{w+dt}{1-t^2}$ ; this strategy can naturally be extended to find  $\bar{t}_-$ . Initially, we have the following  
224    first-order optimality conditions for the problem:  $0 = \nabla p(t) = [b(1-t^2) + (a+bt)t]\sqrt{1-t^2} +$   
225     $[d(1-t^2) + (w+dt)(2t)] \Leftrightarrow dt^2 + 2wt + d = -[b+at]\sqrt{1-t^2}$ . Squaring both sides yields  
226    the following quartic equation:  $c_4t^4 + c_3t^3 + c_2t^2 + c_1t + c_0 = 0$ , where  $c_4 = d^2 + a^2$ ,  
227     $c_3 = 4wd + 2ab$ ,  $c_2 = 4w^2 + 2d^2 - a^2 + b^2$ ,  $c_1 = 4wd - 2ab$ ,  $c_0 = d^2 - b^2$ . This equation  
228    can be solved analytically by Lodovico Ferrari's method [45], resulting in all its real roots  
229     $\{\bar{t}_1, \bar{t}_2, \dots, \bar{t}_j\}$  with  $1 \leq j \leq 4$ .  
230

231    For the second and third cases, Problem (6) essentially boils down to optimization under  
232    orthogonality constraints. The work of [52] derives a breakpoint search method for finding  
233    the optimal solution for Problem (6) with  $\mathbf{J}_{\text{BB}} \in \{(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix}), (\begin{smallmatrix} -1 & 0 \\ 0 & -1 \end{smallmatrix})\}$  using the Givens rotation  
234    and Jacobi reflection matrices.  
235

## 236    2.2 VARIANCE-REDUCED JACOBI BLOCK COORDINATE DESCENT ALGORITHM

237  
238    This subsection proposes the **VR-J-JOBCD** algorithm, a randomized block coordinate  
239    descent method derived from **GS-JOBCD**. Importantly, by leveraging the parallel frame-  
240    work of a Jacobi strategy [17; 9], **VR-J-JOBCD** integrates variance reduction techniques  
241    [38; 30; 18] to decrease oracle complexity in the minimization of finite-sum functions. This  
242    makes the algorithm effective for minimizing large-scale problems under J-orthogonality  
243    constraints.  
244

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

249    ▶ **Variance Reduction Strategy.** We incorporate state-of-the-art variance reduction  
250    strategies from the literature [30; 5] into our algorithm to solve Problem (1). These methods  
251    iteratively generate a stochastic gradient estimator as follows:  
252

$$253 \quad \tilde{\mathbf{G}}^t = \begin{cases} \frac{1}{b} \sum_{i \in \mathbf{S}_+^t} \nabla f_i(\mathbf{X}^t), & \text{with probability } p; \\ 254 \quad \tilde{\mathbf{G}}^{t-1} + \frac{1}{b'} \sum_{i \in \mathbf{S}_*^t} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})), & \text{with probability } 1-p. \end{cases} \quad (9)$$

255  
256    Here,  $\{\mathbf{S}_+^t, \mathbf{S}_*^t\}$  are uniform random minibatch samples with  $|\mathbf{S}_+^t| = b$ ,  $|\mathbf{S}_*^t| = b'$ , and  $\tilde{\mathbf{G}}^0 =$   
257     $\frac{1}{b} \sum_{i \in \mathbf{S}_+^0} \nabla f_i(\mathbf{X}^0)$ . We drop the superscript  $t$  for  $\{\mathbf{S}_+^t, \mathbf{S}_*^t\}$  as  $t$  can be inferred from context.  
258    We only focus on the default setting that [30; 5]:  $b = N$ ,  $b' = \sqrt{b}$  and  $p = \frac{b'}{b+b'}$ .  
259

260    ▶ **Jacobi Block Coordinate Descent Method.** The proposed algorithm is built upon  
261    the parallel framework of a Jacobi strategy. In each iteration  $t$ , we randomly and uniformly  
262    (with replacement) 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  $\Upsilon$  with  
263     $\mathbf{B}^t \in \mathbb{N}^{\frac{n}{2} \times 2}$  and  $\mathbf{B}_{(i)}^t \in \mathbb{N}^2$ . 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  
264    the superscript  $t$  if  $t$  can be inferred from context.  
265

266    The following lemma shows how to choose a suitable matrix  $\mathbf{Q}$  so that the Jacobi strategy  
267    can be applied.  
268

---

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

270    **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  
 271     $\mathbf{Q} = \varsigma \mathbf{I}_4$ , where  $\varsigma$  is some suitable constant with  $\varsigma \leq L_f$ . For any  $\mathbf{B}_{(i)}^t$  and  $\mathbf{B}_{(j)}^t$  with  $i \neq j$ ,  
 272    their corresponding objective functions as in Equation (3) are independent.  
 273

274    We consider the following block coordinate update rule in **VR-J-JOBBCD**:  $\mathbf{X}^{t+1} \leftarrow \tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}_{:}) \triangleq \mathbf{X}^t + [\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top] \mathbf{X}^t$ . The following lemma provides properties  
 275    of this rule.  
 276

277    **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  
 278     $i \in [\frac{n}{2}]$ . 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:  
 279    (a)  $\sum_{i=1}^{n/2} \|\mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_F^2 = \|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_F^2$ . (b)  $\|\mathbf{X}^+ - \mathbf{X}\|_F^2 \leq$   
 280     $\|\mathbf{X}\|_F^2 \cdot \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_F^2$ . (c)  $\|\mathbf{X}^+ - \mathbf{X}\|_{\mathbf{H}}^2 \leq \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbf{Q}}^2$  with  $\mathbf{Q} = \varsigma \mathbf{I}_4$ . (d) For all  
 281     $\tilde{\mathbf{G}} \in \mathbb{R}^{n \times n}$ , it follows that:  $2 \sum_{i=1}^{n/2} \langle \mathbf{V}_i - \mathbf{I}_2, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}})\mathbf{X}^\top]_{\mathbf{B}_{(i)}\mathbf{B}_{(i)}} \rangle \leq \|\mathbf{X}\|_F^2 \sum_{i=1}^{n/2} \|\mathbf{V}_i -$   
 282     $\mathbf{I}_2\|_F^2 + \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_F^2$ .  
 283

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

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

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

289    Instead of computing the exact Euclidean gradient  $\nabla f(\mathbf{X}^t)$  as **GS-JOBBCD**, **VR-J-**  
 290    **JOBBCD** maintains and updates a recursive gradient estimator  $\tilde{\mathbf{G}}^t$  using a variance-reduced  
 291    strategy as in Formula (9). We consider minimizing the following function instead of the  
 292    one on the right-hand side of Inequality (10):  
 293

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

294    Here,  $\mathcal{T}(\mathbf{V}_{:}; \mathbf{X}^t, \mathbf{B}^t)$  can be termed as a stochastic majorization function of  $f(\tilde{\mathcal{X}}_{\mathbf{B}}^t(\mathbf{V}_{:}))$  at the  
 295    current solution  $\mathbf{X}^t$ . Therefore, we can consider the following optimization problem to find  
 296     $\{\mathbf{V}_{:}\}$  using:  $\bar{\mathbf{V}}^t \in \arg \min_{\mathbf{V}_{:}} \mathcal{T}(\mathbf{V}_{:}; \mathbf{X}^t, \mathbf{B}^t)$ , which can be decomposed into  $(n/2)$  independent  
 297    subproblems and solved in parallel. It is important to note that each  $\mathbf{V}_i$  in Problem (12) is  
 298    identical to Problem (6), which can be efficiently solved in  $\mathcal{O}(1)$  using the breakpoint search  
 299    method, as in **GS-JOBBCD**.  
 300

301    We summarize the proposed **VR-J-JOBBCD** in Algorithm 2. Notably, when  $N = 1$ , **VR-**  
 302    **J-JOBBCD** simplifies to a direct Jacobi strategy for solving Problem (1), which we refer to  
 303    as **J-JOBBCD**.  
 304

### 3 OPTIMALITY ANALYSIS

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

313    Initially, we define the first-order optimality condition for Problem (1). Since the matrix  
 314     $\mathbf{X}^\top \mathbf{J} \mathbf{X}$  is symmetric, the Lagrangian multiplier  $\Lambda$  corresponding to the constraints  $\mathbf{X}^\top \mathbf{J} \mathbf{X} =$   
 315     $\mathbf{J}$  is also a symmetric matrix. The Lagrangian function of problem (1) is  $\mathcal{L}(\mathbf{X}, \Lambda) = f(\mathbf{X}) -$   
 316     $\frac{1}{2} \langle \Lambda, \mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J} \rangle$ .  
 317

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

318    **Lemma 3.1.** (*Proof in Section D.1, First-Order Optimality Condition*) We let  $\mathcal{J} \triangleq$   
 319     $\{\mathbf{X} | \mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}\}$ . We have (a) A solution  $\check{\mathbf{X}} \in \mathcal{J}$  is a critical point of problem (1) if  
 320    and only if:  $\mathbf{0} = \nabla_{\mathcal{J}} f(\check{\mathbf{X}}) \triangleq \nabla f(\check{\mathbf{X}}) - \mathbf{J} \check{\mathbf{X}} [\nabla f(\check{\mathbf{X}})]^\top \check{\mathbf{X}} \mathbf{J}$ . The associated Lagrangian mul-  
 321    tiplier can be computed as  $\Lambda = \mathbf{J} \check{\mathbf{X}}^\top \nabla f(\check{\mathbf{X}})$ . (b) The critical point condition is equiva-  
 322    lalent to the requirement that the matrix  $\mathbf{X} \nabla f(\check{\mathbf{X}})^\top \mathbf{J}$  is symmetric, which is expressed as  
 323     $\mathbf{X} \mathbf{G}^\top \mathbf{J} = [\mathbf{X} \mathbf{G}^\top \mathbf{J}]^\top$ .  
 324

**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),  $\theta$  in Inequality (3) (e.g., 1e-6) and  $\xi$  satisfy Inequality (4).

**for**  $t$  from 0 to  $T$  **do**

(S1) Choose a coordinate  $B^t$  from the set  $\Upsilon$  randomly and uniformly (with replacement) for the  $t$ -th iteration. Denote  $B = B^t$ . 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  $\mathbf{Q} = \varsigma \mathbf{I} \in \mathbb{R}^{4 \times 4}$ .

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

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

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

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

end

**Remarks.** While our results in Lemma 3.1 show similarities to existing works focusing on problems under orthogonality constraints [44], this study marks the first investigation into the first-order optimality condition for optimization problems under  $J$ -orthogonality constraints.

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

**Definition 3.2.** (Block Stationary Point, abbreviated as BS-point) Let  $\theta > 0$ . A solution  $\ddot{\mathbf{X}} \in \mathcal{J}$  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 satisfied:  $\mathbf{I}_2 \in \arg \min_{\mathbf{V} \in \mathcal{J}_{\mathbf{B}}} \mathcal{G}(\mathbf{V}; \ddot{\mathbf{X}}, \mathbf{B})$ .

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

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

## 4 CONVERGENCE ANALYSIS

This section provides a convergence analysis for **GS-JOB<sub>CD</sub>** and **VR-J-JOB<sub>CD</sub>**.

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

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

We impose the following additional assumptions on the proposed algorithms:

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

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

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

**Remarks.** (i) Assumption 4.1 is satisfied as the function  $f_i(\mathbf{X})$  is coercive for all  $i$ . (ii) Assumption 4.2 imposes a bound on the (stochastic) gradient, a fairly moderate condition

frequently employed in nonconvex optimization [26]. (*iii*) Assumption 4.3 ensures that the variance of the stochastic gradient is bounded, which is a common requirement in stochastic optimization [30; 5].

#### 4.1 GLOBAL CONVERGENCE

We define the  $\epsilon$ -BS-point as follows.

**Definition 4.4.** ( $\epsilon$ -BS-point) Given any constant  $\epsilon > 0$ , a point  $\ddot{\mathbf{X}}$  is called an  $\epsilon$ -BS-point if:  $\mathcal{E}(\ddot{\mathbf{X}}) \leq \epsilon$ . 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**, while it is defined as  $\mathcal{E}(\mathbf{X}) \triangleq \frac{1}{C_J} \sum_{i=1}^{C_J} \mathbb{E}_{t^i} [\text{dist}(\mathbf{I}_2, \arg \min_{\mathbf{V}_i} \mathcal{T}(\mathbf{V}_i; \mathbf{X}, \tilde{\mathcal{B}}_i))^2]$  for **VR-J-JOBCD**, where the expectation is with respect to the randomness inherent in the algorithm [30].

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

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

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

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

**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  $\epsilon$ -BS-point of Problem (1) within  $\mathcal{O}(nN + \frac{\Delta_0 \sqrt{N}}{\epsilon})$  arithmetic operations.

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

#### 4.2 STRONG CONVERGENCE UNDER KL ASSUMPTION

We prove algorithms achieve strong convergence based on a non-convex analysis tool called Kurdyka-ojasiewicz inequality[2]. We impose the following assumption on Problem (1).

**Assumption 4.8.** (Kurdyka-ojasiewicz Property). Assume that  $f^\circ(\mathbf{X}) = f(\mathbf{X}) + \mathcal{I}_{\mathcal{J}}(\mathbf{X})$  is a KL function. For all  $\mathbf{X} \in \text{dom } f^\circ$ , there exists  $\sigma \in [0, 1], \eta \in (0, +\infty]$  a neighborhood  $\Upsilon$  of  $\mathbf{X}$  and a 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  $f^\circ(\mathbf{X}') \in (f^\circ(\mathbf{X}), f^\circ(\mathbf{X}) + \eta)$ , the following holds:  $\text{dist}(\mathbf{0}, \nabla f^\circ(\mathbf{X}'))\varphi'(f^\circ(\mathbf{X}') - f^\circ(\mathbf{X})) \geq 1$ .

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

**Theorem 4.9.** (Proof in Section E.5, a Finite Length Property). The sequence  $\{\mathbf{X}^t\}_{t=0}^\infty$  of **GS-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  $\varphi(\cdot)$  is the desingularization function defined in Proposition 4.8.

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

### 5 DISCUSSION

► **Differences with [53].** Since both our paper and [53] adopt the block coordinate descent method as the framework, we compare and analyze this paper with [53] in Table 1. Moreover, this paper is the first to propose the first-order optimality condition, the tangent space of the optimization manifold, the optimality condition, and the convergence properties for J orthogonality constraint problem for the first time.

---

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

Table 1: Comparisons with JOBCD and [53]

	ours	[53]
Problem nature	unbounded	compact
parallelizability	✓	✗
Include Stochastic strategy	✓	✗
Convergence analysis under parallelization/Variance-Reduction strategy	✓	✗
New applications	✓	✗

► **Comparisons with GS-JOBCD and J-JOBCD.** GS-JOBCD and J-JOBCD are complementary and mutually confirmatory. (*i*) The  $\mathbf{Q}$  in GS-JOBCD (s2) has two selection methods as described in (3). However, in the (VR)-J-JOBCD (s1), only  $\mathbf{Q} = \zeta \mathbf{I} \in \mathbb{R}^{4 \times 4}$  is applicable due to the requirement for independent updates of each block in the parallelization strategy. (*ii*) In small-scale problems, the benefits of parallel strategies are limited. Therefore, the sequential GS-JOBCD generally performs better.

## 6 APPLICATIONS AND NUMERICAL EXPERIMENTS

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

► **Application to the Hyperbolic Eigenvalue Problem (HEVP).** The hyperbolic eigenvalue problem refers to the generalized eigenvalue problem in hyperbolic spaces [39]. This problem is a fundamental component in machine learning models, such as Hyperbolic PCA [42; 6]. Given a data matrix  $\mathbf{D} \in \mathbb{R}^{m \times n}$  and a signature matrix  $\mathbf{J}$  with signature  $(p, n - p)$ , HEVP can be formulated as the following optimization problem:  
 $\min_{\mathbf{X}} -\text{tr}(\mathbf{X}^\top \mathbf{D}^\dagger \mathbf{DX})$ , s.t.  $\mathbf{X}^\top \mathbf{JX} = \mathbf{J}$ .

► **Application to the Hyperbolic Structural Probe Problem (HSPP).** The Structure Probe (SP) is a metric learning model aimed at understanding the intrinsic semantic information of large language models [20] [7]. Given a data matrix  $\mathbf{D} \in \mathbb{R}^{m \times n}$  and its associated Euclidean distance metric matrix  $\mathbf{T} \in \mathbb{R}^{m \times m}$ , HSPP employs a smooth homeomorphic mapping function  $\varphi(\cdot)$  to project the data  $\mathbf{D}$  into ultra-hyperbolic space. Subsequently, it seeks an appropriate linear transformation  $\mathbf{X} \in \mathbb{R}^{n \times n}$  constrained within a specific structure  $\mathbf{X} \in \mathcal{J}$ , such that the resulting transformed 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 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  $i, j \in [m]$ , 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 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 on the functions  $\varphi(\cdot)$  and  $d_\alpha(\cdot, \cdot)$ , please refer to Appendix Section F.1.

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

► **Compared Methods.** We compare **GS-JOBCD** and **VR-J-JOBCD** with 3 state-of-the-art optimization algorithms under J-orthogonality constraints. (*i*) The CS Decomposition Method (**CSDM**) [49]. (*ii*) Standard ADMM (**ADMM**) (Appendix Section F.3) [19]. (*iii*) **UMCM**: Unconstrained Multiplier Correction Method (Appendix Section F.4) [47; 48; 13].

► **Experiment Settings.** All methods are implemented using Pytorch on an Intel 2.6 GHz processor with an A40 (48GB). For **HSPP**, we fix  $\alpha$  to 1. Each method employs the same random J-orthogonal matrix  $\mathbf{X}^0$  with  $\frac{1}{n^2} \sum_{ij}^n |\mathbf{X}^{0\top} \mathbf{JX}^0 - \mathbf{J}|_{ij} \leq 1e-9$ . For **(VR)-J-JOBCD**, we define  $\mathbf{X}$  as a high-dimensional tensor in PyTorch to achieve parallelization. The built-in solver *Adagrad* is used to solve the unconstrained minimization problem in **CSDM** and **UMCM**, the optimal learning rate is selected from the range [5e-4, 5e-3]. For more experimental details for **ADMM** and **UMCM**, please refer to Appendix Section F.5. We provide our code in the supplemental material.

Table 2: Comparisons of the objectives for HEVP across all the compared methods. The time limit is set to 90s. The notation ‘(+)’ indicates that **GS-JOBCD** significantly improves upon the initial solution provided by **CSDM**. The 1<sup>st</sup>, 2<sup>nd</sup>, and 3<sup>rd</sup> best results are colored with red, green and blue, respectively. The value in (·) stands for  $\frac{1}{n^2} \sum_{ij}^n |\mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J}|_{ij}$  and cells with this value greater than 1e-5 are highlighted in gray.

dataname(m-n-p)	UMCM	ADMM	CSDM	GS-JOBCD	J-JOBCD	CSDM+GS-JOBCD
cifar(1000-100-50)	-1.11e+04(9.9e-10)	-7.32e+03(1.0e-07)	-1.07e+04(1.4e-09)	-1.53e+04(1.3e-08)	-3.43e+04(6.7e-08)	-6.97e+04(1.6e-08)(+)
CanCaltech(2000-1000-500)	-2.62e+02(1.9e-05)	-2.51e+02(6.5e-08)	-2.80e+02(1.4e-10)	-2.51e+02(1.2e-10)	-1.96e+03(1.6e-08)	-2.87e+02(1.8e-10)(+)
gisette(3000-1000-500)	-1.02e+08(5.4e-05)	-1.49e+11(2.5e+00)	-1.77e+06(1.3e-10)	-1.41e+06(1.9e-10)	-3.58e+06(6.0e-09)	-1.81e+06(2.4e-10)(+)
mnist(1000-780-390)	-2.13e+05(1.4e-09)	-3.32e+06(1.4e-02)	-5.22e+04(1.6e-10)	-4.41e+04(2.9e-10)	-4.63e+05(2.1e-08)	-8.88e+04(5.6e-10)(+)
randn10(10-10-5)	-4.23e+01(8.4e-09)	-4.21e+01(1.1e-07)	-2.29e+02(2.7e-01)	-4.23e+01(1.5e-08)	-8.46e+01(7.5e-07)	-3.01e+02(2.7e-01)(+)
randn100(100-100-50)	-5.01e+03(1.4e-09)	-4.93e+03(1.1e-07)	-1.34e+04(3.4e-09)	-4.90e+03(2.6e-09)	-1.71e+04(1.2e-07)	-2.13e+04(5.8e-08)(+)
randn1000(1000-1000-500)	-5.64e+05(4.9e-07)	-5.63e+05(6.6e-08)	-5.46e+05(1.4e-10)	-5.01e+05(1.1e-10)	-2.84e+06(6.0e-08)	-5.52e+05(2.4e-10)(+)
sector(500-1000-500)	-1.42e+03(7.5e-10)	-1.51e+03(5.2e-10)	-1.63e+03(1.1e-10)	-1.54e+03(1.6e-10)	-2.52e+03(3.4e-09)	-1.63e+03(1.4e-10)
TDT2(1000-1000-500)	-1.05e+08(2.8e-06)	-7.34e+08(5.8e+00)	-1.93e+06(1.1e-10)	-1.81e+06(1.8e-10)	-3.21e+06(3.6e-09)	-1.93e+06(1.4e-10)
w1a(2470-290-145)	-1.46e+04(2.0e-09)	-1.39e+04(7.8e-08)	-1.45e+04(1.1e-05)	-1.38e+04(3.1e-09)	-3.60e+06(5.3e-07)	-1.63e+04(1.1e-05)(+)

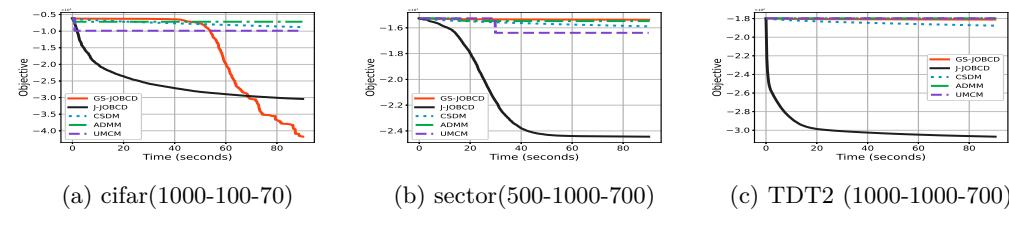


Figure 1: The convergence curve for the HEVP across various datasets with different parameters ( $m - n - p$ ).

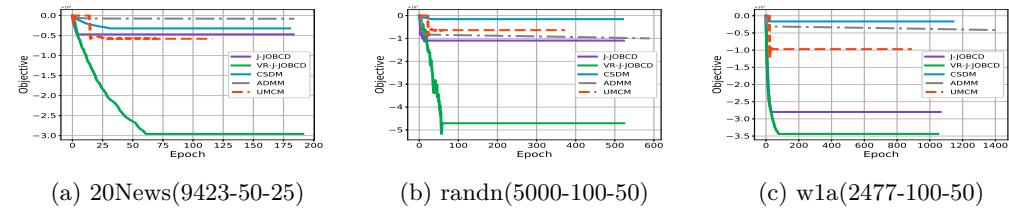


Figure 2: Comparisons of objective values ( $F(\mathbf{X}) - F^0$ ) of HSPP for all the compared methods with different parameters ( $m - n - p$ ).

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

## 7 CONCLUSIONS

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

---

540 REFERENCES  
541

- 542 [1] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on ma-*  
543 *trix manifolds*. Princeton University Press, 2008.
- 544 [2] Hédy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal  
545 alternating minimization and projection methods for nonconvex problems: An ap-  
546 proach based on the kurdyka-łojasiewicz inequality. *Mathematics of operations research*,  
547 35(2):438–457, 2010.
- 548 [3] Adam Bojanczyk, Nicholas J Higham, and Harikrishna Patel. Solving the indefinite  
549 least squares problem by hyperbolic qr factorization. *SIAM Journal on Matrix Analysis*  
550 and *Applications*, 24(4):914–931, 2003.
- 552 [4] HanQin Cai, Yuchen Lou, Daniel McKenzie, and Wotao Yin. A zeroth-order block  
553 coordinate descent algorithm for huge-scale black-box optimization. In *International*  
554 *Conference on Machine Learning*, pages 1193–1203. PMLR, 2021.
- 555 [5] Xufeng Cai, Chaobing Song, Stephen Wright, and Jelena Diakonikolas. Cyclic block  
556 coordinate descent with variance reduction for composite nonconvex optimization. In  
557 *International Conference on Machine Learning*, pages 3469–3494. PMLR, 2023.
- 559 [6] Ines Chami, Albert Gu, Dat P Nguyen, and Christopher Re. Horopca: Hyperbolic  
560 dimensionality reduction via horospherical projections. In *International Conference on*  
561 *Machine Learning (ICML)*, volume 139, pages 1419–1429, 2021.
- 562 [7] Boli Chen, Yao Fu, Guangwei Xu, Pengjun Xie, Chuanqi Tan, Mosha Chen, and Liping  
563 Jing. Probing bert in hyperbolic spaces. *ICLR*, 2021.
- 565 [8] Weize Chen, Xu Han, Yankai Lin, Hexu Zhao, Zhiyuan Liu, Peng Li, Maosong Sun, and  
566 Jie Zhou. Fully hyperbolic neural networks. *arXiv preprint arXiv:2105.14686*, 2021.
- 567 [9] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-  
568 convex sgd. *Advances in neural information processing systems*, 32, 2019.
- 570 [10] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental  
571 gradient method with support for non-strongly convex composite objectives. *Advances*  
572 *in neural information processing systems*, 27, 2014.
- 573 [11] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal  
574 non-convex optimization via stochastic path-integrated differential estimator. *Advances*  
575 *in neural information processing systems*, 31, 2018.
- 577 [12] Hamza Fawzi and Harry Goulbourne. Faster proximal algorithms for matrix optimiza-  
578 tion using jacobi-based eigenvalue methods. *Advances in Neural Information Processing*  
579 *Systems*, 34:11397–11408, 2021.
- 581 [13] Bin Gao, Xin Liu, Xiaojun Chen, and Ya-xiang Yuan. A new first-order algorithmic  
582 framework for optimization problems with orthogonality constraints. *SIAM Journal on*  
583 *Optimization*, 28(1):302–332, 2018.
- 584 [14] Bin Gao, Xin Liu, and Ya-xiang Yuan. Parallelizable algorithms for optimization  
585 problems with orthogonality constraints. *SIAM Journal on Scientific Computing*,  
586 41(3):A1949–A1983, 2019.
- 587 [15] Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approxi-  
588 mation methods for nonconvex stochastic composite optimization. *Mathematical Pro-*  
589 *gramming*, 155(1-2):267–305, 2016.
- 591 [16] Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- 592 [17] Eldon R Hansen. On cyclic jacobi methods. *Journal of the Society for Industrial and*  
593 *Applied Mathematics*, 11(2):448–459, 1963.

- 594 [18] Vjeran Hari and Erna Begović Kovač. On the convergence of complex jacobi methods.  
595 *Linear and multilinear algebra*, 69(3):489–514, 2021.  
596
- 597 [19] Bingsheng He and Xiaoming Yuan. On the  $\mathcal{O}(1/n)$  convergence rate of the  
598 douglas-rachford alternating direction method. *SIAM Journal on Numerical Analysis*,  
599 50(2):700–709, 2012.
- 600 [20] John Hewitt and Christopher D. Manning. A structural probe for finding syntax in  
601 word representations. In Jill Burstein, Christy Doran, and Thamar Solorio, editors,  
602 *Proceedings of the 2019 Conference of the North American Chapter of the Association  
603 for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages  
604 4129–4138, 2019.
- 605 [21] Nicholas J Higham. J-orthogonal matrices: Properties and generation. *SIAM review*,  
606 45(3):504–519, 2003.
- 607 [22] Minhui Huang, Shiqian Ma, and Lifeng Lai. A riemannian block coordinate descent  
608 method for computing the projection robust wasserstein distance. In *International  
609 Conference on Machine Learning*, pages 4446–4455. PMLR, 2021.
- 610 [23] Bo Hui and Wei-Shinn Ku. Low-rank nonnegative tensor decomposition in hyperbolic  
611 space. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery  
612 and Data Mining*, pages 646–654, 2022.
- 613 [24] Sashank J Reddi, Suvrit Sra, Barnabas Poczos, and Alexander J Smola. Proximal  
614 stochastic methods for nonsmooth nonconvex finite-sum optimization. *Advances in  
615 neural information processing systems*, 29, 2016.
- 616 [25] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive  
617 variance reduction. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q.  
618 Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26.  
619 Curran Associates, Inc., 2013.
- 620 [26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization.  
621 In Yoshua Bengio and Yann LeCun, editors, *International Conference on Learning  
622 Representations (ICLR)*, 2015.
- 623 [27] Marc Law. Ultrahyperbolic neural networks. *Advances in Neural Information Process-  
624 ing Systems*, 34:22058–22069, 2021.
- 625 [28] Marc Law and Jos Stam. Ultrahyperbolic representation learning. *Advances in neural  
626 information processing systems*, 33:1668–1678, 2020.
- 627 [29] Qunwei Li, Yi Zhou, Yingbin Liang, and Pramod K Varshney. Convergence analysis  
628 of proximal gradient with momentum for nonconvex optimization. In *International  
629 Conference on Machine Learning*, pages 2111–2119. PMLR, 2017.
- 630 [30] Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. Page: A simple and  
631 optimal probabilistic gradient estimator for nonconvex optimization. In *International  
632 conference on machine learning*, pages 6286–6295. PMLR, 2021.
- 633 [31] Wan-Duo Kurt Ma, JP Lewis, and W Bastiaan Kleijn. The hsic bottleneck: Deep  
634 learning without back-propagation. In *Proceedings of the AAAI conference on artificial  
635 intelligence*, volume 34, pages 5085–5092, 2020.
- 636 [32] Julien Mairal. Optimization with first-order surrogate functions. In *International  
637 Conference on Machine Learning (ICML)*, volume 28, pages 783–791, 2013.
- 638 [33] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method  
639 for machine learning problems using stochastic recursive gradient. In *International  
640 conference on machine learning*, pages 2613–2621. PMLR, 2017.

- 
- 648 [34] Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz  
649 model of hyperbolic geometry. In *International Conference on Machine Learning*, pages  
650 3779–3788. PMLR, 2018.
- 651 [35] Vedran Novaković and Sanja Singer. A kogbetliantz-type algorithm for the hyperbolic  
652 svd. *Numerical algorithms*, 90(2):523–561, 2022.
- 654 [36] Julie Nutini, Issam Laradji, and Mark Schmidt. Let’s make block coordinate descent  
655 converge faster: faster greedy rules, message-passing, active-set complexity, and super-  
656 linear convergence. *Journal of Machine Learning Research*, 23(131):1–74, 2022.
- 657 [37] Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Luo. A unified convergence analysis  
658 of block successive minimization methods for nonsmooth optimization. *SIAM Journal  
659 on Optimization*, 23(2):1126–1153, 2013.
- 661 [38] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the  
662 stochastic average gradient. *Mathematical Programming*, 162:83–112, 2017.
- 663 [39] Ivan Slapnicar and Ninoslav Truhar. Relative perturbation theory for hyperbolic eigen-  
664 value problem. *Linear Algebra and its Applications*, 309(1):57–72, 2000.
- 666 [40] Michael Stewart and Paul Van Dooren. On the factorization of hyperbolic and unitary  
667 transformations into rotations. *SIAM Journal on Matrix Analysis and Applications*,  
668 27(3):876–890, 2005.
- 669 [41] Puoya Tabaghi and Ivan Dokmanić. Hyperbolic distance matrices. In *Proceedings of the  
670 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*,  
671 pages 1728–1738, 2020.
- 673 [42] Puoya Tabaghi, Michael Khanzadeh, Yusu Wang, and Sivash Mirarab. Principal com-  
674 ponent analysis in space forms. *ArXiv*, abs/2301.02750, 2023.
- 675 [43] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. Spiderboost and  
676 momentum: Faster variance reduction algorithms. *Advances in Neural Information  
677 Processing Systems*, 32, 2019.
- 679 [44] Zaiwen Wen and Wotao Yin. A feasible method for optimization with orthogonality  
680 constraints. *Mathematical Programming*, 142:397 – 434, 2012.
- 681 [45] WikiContributors. Quartic equation. [https://en.wikipedia.org/wiki/Quartic\\_equation](https://en.wikipedia.org/wiki/Quartic_equation).
- 684 [46] Ruiyuan Wu, Anna Scaglione, Hoi-To Wai, Nurullah Karakoc, Kari Hreinsson, and  
685 Wing-Kin Ma. Federated block coordinate descent scheme for learning global and  
686 personalized models. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
687 volume 35, pages 10355–10362, 2021.
- 688 [47] Nachuan Xiao, Xin Liu, and Kim-Chuan Toh. Dissolving constraints for riemannian  
689 optimization. *Mathematics of Operations Research*, 49(1):366–397, 2024.
- 691 [48] Nachuan Xiao, Xin Liu, and Ya-xiang Yuan. A class of smooth exact penalty func-  
692 tion methods for optimization problems with orthogonality constraints. *Optimization  
693 Methods and Software*, 37(4):1205–1241, 2022.
- 694 [49] Bo Xiong, Shichao Zhu, Mojtaba Nayyeri, Chengjin Xu, Shirui Pan, Chuan Zhou, and  
695 Steffen Staab. Ultrahyperbolic knowledge graph embeddings. In *Proceedings of the 28th  
696 ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2130–2139,  
697 2022.
- 698 [50] Bo Xiong, Shichao Zhu, Nico Potyka, Shirui Pan, Chuan Zhou, and Steffen Staab.  
699 Semi-riemannian graph convolutional networks. *ArXiv*, abs/2106.03134, 2021.
- 700 [51] Tao Yu and Christopher M De Sa. Numerically accurate hyperbolic embeddings using  
701 tiling-based models. *Advances in Neural Information Processing Systems*, 32, 2019.

- 
- 702 [52] Ganzhao Yuan. A block coordinate descent method for nonsmooth composite optimization  
703 under orthogonality constraints. *ArXiv*, abs/2304.03641, 2023.  
704
- 705 [53] Ganzhao Yuan. Coordinate descent methods for dc minimization: Optimality conditions  
706 and global convergence. In *Proceedings of the AAAI Conference on Artificial*  
707 *Intelligence*, volume 37, pages 11034–11042, 2023.
- 708 [54] Ganzhao Yuan. Coordinate descent methods for fractional minimization. In *International*  
709 *Conference on Machine Learning*, pages 40488–40518, 2023.  
710
- 711 [55] Jinshan Zeng, Tim Tszi-Kit Lau, Shaobo Lin, and Yuan Yao. Global convergence of  
712 block coordinate descent in deep learning. In *International conference on machine*  
713 *learning*, pages 7313–7323. PMLR, 2019.
- 714 [56] Yiding Zhang, Xiao Wang, Chuan Shi, Nian Liu, and Guojie Song. Lorentzian graph  
715 convolutional networks. In *Proceedings of the Web Conference 2021*, pages 1249–1261,  
716 2021.
- 717 [57] Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic nested variance reduction for  
718 nonconvex optimization. *The Journal of Machine Learning Research*, 21(1):4130–4192,  
719 2020.

720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

---

# 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  $\mathbf{X}$
- $\|\mathbf{X}\|_{\text{F}}$  : 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  $\xi$  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  $\text{dist}(\xi, \xi') \triangleq \inf_{\mathbf{X} \in \xi, \mathbf{X}' \in \xi'} \|\mathbf{X} - \mathbf{X}'\|_{\text{F}}$
- $\mathcal{I}_{\xi}(\mathbf{x})$  : the indicator function of a set  $\xi$  with  $\mathcal{I}_{\xi}(\mathbf{x}) = 0$  if  $\mathbf{x} \in \xi$  and otherwise  $+\infty$ .

---

810    A.2 RELEVANT LEMMAS  
811

812    **Lemma A.1.** (Lemma 6.6 of [52]) For any  $\mathbf{W} \in \mathbb{R}^{n \times n}$ , we have:  $\sum_{i=1}^{C_n^k} \|\mathbf{W}(\mathcal{B}_i, \mathcal{B}_i)\|_F^2 =$   
813     $\frac{k}{n} C_n^k \sum_i \mathbf{W}_{ii}^2 + C_n^{k-2} \sum_i \sum_{j,j \neq i} \mathbf{W}_{ij}^2$ . Here, the set  $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{C_n^k}\}$  represents all possible  
814    combinations of the index vectors choosing  $k$  items from  $n$  without repetition.  
815

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

$$818 \quad \mathbb{E}_{t \in \mathbf{S}_+} [\|\frac{1}{b} \sum_{i \in \mathbf{S}_+} \nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_F^2] = \frac{N-b}{b(N-1)} \mathbb{E}_{t \in \mathbf{S}_+} [\|\nabla f_i(\mathbf{X}^t) - \nabla f(\mathbf{X}^t)\|_F^2].$$

821    *Proof.* The proof is exactly the same as in Lemma 2.8 of [5].  $\square$   
822

824    **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  $\mathbf{X} \in$   
825     $\mathbb{R}^{n \times n}$ , is :

$$826 \quad \mathbf{T}_{\mathbf{X}}\mathcal{J} \triangleq \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \mathbf{X}^\top \mathbf{J} \mathbf{Y} + \mathbf{Y}^\top \mathbf{J} \mathbf{X} = 0\}, \quad (13)$$

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

831    *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  
832    along  $\mathbf{Y} \in \mathbb{R}^{n \times n}$  in the tangent space of  $\mathbf{X}$ , we obtain:

$$\begin{aligned} 834 \quad h(\mathbf{X} + \mathbf{Y}) &= (\mathbf{X} + \mathbf{Y})^\top \mathbf{J}(\mathbf{X} + \mathbf{Y}) - \mathbf{J} \\ 835 &= \mathbf{X}^\top \mathbf{J} \mathbf{X} + \mathbf{X}^\top \mathbf{J} \mathbf{Y} + \mathbf{Y}^\top \mathbf{J} \mathbf{X} + \mathbf{Y}^\top \mathbf{J} \mathbf{Y} - \mathbf{J} \\ 836 &\stackrel{\textcircled{1}}{=} \mathbf{X}^\top \mathbf{J} \mathbf{Y} + \mathbf{Y}^\top \mathbf{J} \mathbf{X} + \mathbf{Y}^\top \mathbf{J} \mathbf{Y} \\ 837 &\stackrel{\textcircled{2}}{=} t\mathbf{X}^\top \mathbf{J}\tilde{\mathbf{Y}} + t\tilde{\mathbf{Y}}^\top \mathbf{J} \mathbf{X} + t^2\tilde{\mathbf{Y}}^\top \mathbf{J}\tilde{\mathbf{Y}} \end{aligned}$$

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

842    Since  $t$  is a positive scalar approaching 0, we can ignore the higher-order term:  $t^2\tilde{\mathbf{Y}}^\top \mathbf{J}\tilde{\mathbf{Y}}$ .  
843    According to the properties of the tangent space of any manifold, we have:  $h(\mathbf{X} + \mathbf{Y}) = 0$ ,  
844    In 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  
845    space:  $\mathbf{T}_{\mathbf{X}}\mathcal{J} \triangleq \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \mathbf{X}^\top \mathbf{J} \mathbf{Y} + \mathbf{Y}^\top \mathbf{J} \mathbf{X} = 0\}$ .  $\square$

847    B ADDITIONAL DISCUSSIONS  
848

849    B.1 ON THE GLOBAL OPTIMAL SOLUTION FOR PROBLEM (7)  
850

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

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

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

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

---

864 C PROOFS FOR SECTION 2
865

## 866 C.1 PROOF OF LEMMA 2.1 867

868 *Proof.* Defining  $\mathbf{J}_{\mathbf{BB}} = \mathbf{J}(\mathbf{U}_B, \mathbf{U}_B)$ , then we have:  $\mathbf{J}\mathbf{U}_B = \mathbf{U}_B\mathbf{J}_{\mathbf{BB}}$ ,  $\mathbf{U}_B^\top \mathbf{J} = \mathbf{J}_{\mathbf{BB}}\mathbf{U}_B^\top$ , and
869  $\mathbf{U}_B^\top \mathbf{J}\mathbf{U}_B = \mathbf{J}_{\mathbf{BB}}$ .

870 **Part (a).** For any  $\mathbf{V} \in \mathbb{R}^{2 \times 2}$  and  $B \in \{\mathcal{B}_i\}_{i=1}^{\mathbf{C}_n^2}$ , we have:
871

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

883 **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:
884

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

885 where step ① and step ② use the norm inequality that  $\|\mathbf{AX}\|_F \leq \|\mathbf{A}\|_F \cdot \|\mathbf{X}\|_F$  for any  $\mathbf{A}$  and  $\mathbf{X}$ ; step ③ uses  $\|\mathbf{U}_B\| = \|\mathbf{U}_B^\top\| = 1$ .
886

887 **Part (c).** We define  $\mathbf{Z} \triangleq \mathbf{U}_B^\top \mathbf{X}$ . We derive:
888

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

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

## 906 C.2 PROOF OF LEMMA 2.3 907

908 *Proof.* We denote  $w = c + e$ . According to the properties of trigonometric functions, we
909 have: (i)  $\tilde{c}^2 = \frac{1}{1-\tilde{t}^2}$ ; (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
910  $|\tilde{t}| < 1$ .
911

912 We discuss two cases for Problem (8).
913

914 **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:
915  $\bar{\mu}_+ = \arg \min_{\mu} \frac{a+\tilde{t}b}{\sqrt{1-\tilde{t}^2}} + \frac{w+\tilde{t}d}{1-\tilde{t}^2} - e$ . Therefore, the optimal solution  $\bar{\mu}_+$  can be computed as:
916

917 
$$\cosh(\bar{\mu}_+) = \frac{1}{\sqrt{1-(\bar{\mu}_+)^2}}, \text{ and } \sinh(\bar{\mu}_+) = \frac{\bar{\mu}_+}{\sqrt{1-(\bar{\mu}_+)^2}} \quad (14)$$

918   **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:  
919    $\bar{\mu}_- = \arg \min_{\mu} \frac{-a-\tilde{t}b}{\sqrt{1-\tilde{t}^2}} + \frac{w+\tilde{t}d}{1-\tilde{t}^2} - e$ . Therefore, the optimal solution  $\bar{\mu}_-$  can be computed as:  
920  
921

$$922 \quad \cosh(\bar{\mu}_-) = \frac{-1}{\sqrt{1-(\bar{\mu}_-)^2}}, \text{ and } \sinh(\bar{\mu}_-) = \frac{-\bar{\mu}_-}{\sqrt{1-(\bar{\mu}_-)^2}}. \quad (15)$$

923   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  
924   (15), the optimal solution pair  $[\cosh(\bar{\mu}_-), \sinh(\bar{\mu}_-)]$  for problem (8) can be computed as:

$$925 \quad [\cosh(\bar{\mu}_-), \sinh(\bar{\mu}_-)] = \arg \min_{[c,s]} \check{F}(c,s), \\ 926 \quad \text{s. t. } [c,s] \in \{[\cosh(\bar{\mu}_+), \sinh(\bar{\mu}_+)], [\cosh(\bar{\mu}_-), \sinh(\bar{\mu}_-)]\}$$

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

### 934   C.3 PROOF OF LEMMA 2.4

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

$$936 \quad f(\mathbf{X}^t) + \frac{1}{2} \|\mathbf{V}_i - \mathbf{I}\|_{\mathbf{Q}+\theta\mathbf{I}}^2 + \langle \mathbf{V}_i - \mathbf{I}, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}_{(i)}^t \mathbf{B}_{(i)}^t} \rangle$$

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

938   **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  
939    $[\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{B}_{(i)}^t \mathbf{B}_{(i)}^t} = [\nabla f(\mathbf{X}^t)](\mathbf{B}_{(i)}^t, :)[(\mathbf{X}^t)^\top](; \mathbf{B}_{(i)}^t) = [\nabla f(\mathbf{X}^t)](\mathbf{B}_{(i)}^t, :)[(\mathbf{X}^t)(\mathbf{B}_{(i)}^t, :)^\top]$ ,  
940   which just use the information of block  $\mathbf{B}_{(i)}^t$ . The proof ends.  $\square$

### 941   C.4 PROOF OF LEMMA 2.5

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

$$944 \quad \begin{aligned} \|\sum_{i=1}^{\frac{n}{2}} [\mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}]\|_F^2 &\stackrel{\textcircled{1}}{=} \left\| \begin{bmatrix} \mathbf{K}_1 \\ \mathbf{K}_2 \\ \vdots \\ \mathbf{K}_{\frac{n}{2}} \end{bmatrix} \right\|_F^2 \\ 945 &\stackrel{\textcircled{2}}{=} \|\mathbf{K}_1\|_F^2 + \|\mathbf{K}_2\|_F^2 + \cdots + \|\mathbf{K}_{\frac{n}{2}}\|_F^2 \\ 946 &\stackrel{\textcircled{3}}{=} \sum_{i=1}^{\frac{n}{2}} [\|\mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top \mathbf{X}\|_F^2] \end{aligned}$$

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

949   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  
950   have the following inequalities:

$$951 \quad \|\mathbf{X}^+ - \mathbf{X}\|_F^2 = \|[\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top] \mathbf{X}\|_F^2 \quad (16)$$

$$952 \quad \stackrel{\textcircled{1}}{=} \sum_{i=1}^{n/2} \|[\mathbf{U}_{\mathbf{B}_{(i)}}(\mathbf{V}_i - \mathbf{I}_2)\mathbf{U}_{\mathbf{B}_{(i)}}^\top] \mathbf{X}\|_F^2 \quad (17)$$

$$953 \quad \stackrel{\textcircled{2}}{\leq} \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_F^2 \cdot \|\mathbf{X}\|_F^2, \quad (18)$$

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

972 Part (**c**). We derive the following results:  
973

$$\begin{aligned} 974 \quad \frac{1}{2} \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbf{H}}^2 &= \frac{1}{2} \|[\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}(i)} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}(i)}^\top] \mathbf{X}\|_{\mathbf{H}}^2 \\ 975 \\ 976 &\stackrel{\textcircled{1}}{=} \frac{1}{2} \sum_{i=1}^{n/2} \|[\mathbf{U}_{\mathbf{B}(i)} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}(i)}^\top] \mathbf{X}\|_{\mathbf{H}}^2 \\ 977 \\ 978 &\stackrel{\textcircled{2}}{\leq} \frac{1}{2} \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbf{Q}}^2 \end{aligned}$$

979 where step ① uses the conclusion of Part (**a**); step ② uses the same proof process of Part  
980 (**c**) of lemma 2.1.  
981

982 Part (**d**). We derive the following results:  
983

$$\begin{aligned} 983 \quad &\sum_{i=1}^{n/2} \langle \mathbf{V}_i - \mathbf{I}_2, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}) \mathbf{X}^\top]_{\mathbf{B}_i \mathbf{B}_i} \rangle \\ 984 \\ 985 &= \sum_{i=1}^{n/2} \langle [\mathbf{U}_{\mathbf{B}(i)} (\mathbf{V}_i - \mathbf{I}_2) \mathbf{U}_{\mathbf{B}(i)}^\top] \mathbf{X}, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}})] \rangle \\ 986 \\ 987 &= \langle \mathbf{X}^+ - \mathbf{X}, [(\nabla f(\mathbf{X}) - \tilde{\mathbf{G}})] \rangle \\ 988 &\stackrel{\textcircled{1}}{\leq} \frac{1}{2} \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbf{F}}^2 + \frac{1}{2} \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_{\mathbf{F}}^2 \\ 989 \\ 990 &\stackrel{\textcircled{2}}{\leq} \frac{1}{2} \|\mathbf{X}\|_{\mathbf{F}}^2 \sum_{i=1}^{n/2} \|\mathbf{V}_i - \mathbf{I}_2\|_{\mathbf{F}}^2 + \frac{1}{2} \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_{\mathbf{F}}^2 \end{aligned} \tag{19}$$

991 where step ① uses  $\forall \mathbf{A}, \mathbf{B}, \frac{1}{2} \|\mathbf{A} - \mathbf{B}\|_{\mathbf{F}}^2 = \frac{1}{2} \|\mathbf{A}\|_{\mathbf{F}}^2 + \frac{1}{2} \|\mathbf{B}\|_{\mathbf{F}}^2 - \langle \mathbf{A}, \mathbf{B} \rangle \geq 0$ , with  $\mathbf{A} = \|\mathbf{X}^+ - \mathbf{X}\|_{\mathbf{F}}^2$   
992 and  $\mathbf{B} = \|[\nabla f(\mathbf{X}) - \tilde{\mathbf{G}}]\|_{\mathbf{F}}^2$ ; step ② uses the conclusion of Part (**b**).  $\square$   
993

## 994 D PROOFS FOR SECTION 3

### 995 D.1 PROOF OF LEMMA 3.1

996 *Proof.* We consider the Lagrangian function of problem (1):  
997

$$1000 \quad \mathcal{L}(\mathbf{X}, \Lambda) = f(\mathbf{X}) - \frac{1}{2} \langle \Lambda, \mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J} \rangle. \tag{20}$$

1001 Setting the gradient of  $\mathcal{L}(\mathbf{X}, \Lambda)$  w.r.t.  $\mathbf{X}$  to zero yields:  
1002

$$1003 \quad \nabla f(\mathbf{X}) - \mathbf{J} \mathbf{X} \Lambda = \mathbf{0}. \tag{21}$$

1004 **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  
1005  $\mathbf{J} \Lambda = \mathbf{X}^\top \nabla f(\mathbf{X})$ . Multiplying both sides by  $\mathbf{J}^\top$  and using  $\mathbf{J}^\top \mathbf{J} = \mathbf{I}$ , we have  $\Lambda = \mathbf{J} \mathbf{X}^\top \nabla f(\mathbf{X})$ .  
1006 Since  $\Lambda$  is symmetric, we have  $\Lambda = \nabla f(\mathbf{X})^\top \mathbf{J} \mathbf{X}$ . Putting this equality into Equality (21)  
1007 yields the following first-order optimality condition for Problem (1):

$$1008 \quad \nabla f(\mathbf{X}) = \mathbf{J} \mathbf{X} [\nabla f(\mathbf{X})]^\top \mathbf{X} \mathbf{J}. \tag{22}$$

1009 **Part (b).** We let  $\mathbf{G} = \nabla f(\mathbf{X})$ . We derive the following results:  
1010

$$\begin{aligned} 1011 \quad \mathbf{G} = \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J} &\stackrel{\textcircled{1}}{\Rightarrow} \mathbf{J} \mathbf{X}^\top \cdot \mathbf{G} = \mathbf{J} \mathbf{X}^\top \cdot \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J} \\ 1012 &\stackrel{\textcircled{2}}{\Rightarrow} \mathbf{J} \mathbf{X}^\top \mathbf{G} = \mathbf{G}^\top \mathbf{X} \mathbf{J} \\ 1013 &\stackrel{\textcircled{3}}{\Rightarrow} \mathbf{X} (\mathbf{J} \mathbf{X}^\top \mathbf{G}) \mathbf{X}^\top = \mathbf{X} (\mathbf{G}^\top \mathbf{X} \mathbf{J}) \mathbf{X}^\top \\ 1014 &\stackrel{\textcircled{4}}{\Rightarrow} \mathbf{X} \underbrace{\mathbf{J} \mathbf{X}^\top \mathbf{G} \mathbf{X}^\top \mathbf{J} \mathbf{J}}_{\triangleq \mathbf{G}^\top} = \mathbf{J} \underbrace{\mathbf{J} \mathbf{X}^\top \mathbf{X} \mathbf{J}}_{\triangleq \mathbf{G}} \mathbf{X}^\top \\ 1015 &\stackrel{\textcircled{5}}{\Rightarrow} (\mathbf{X} \mathbf{G}^\top \mathbf{J}) \cdot \mathbf{J} \mathbf{X} = (\mathbf{J} \mathbf{G} \mathbf{X}^\top) \cdot \mathbf{J} \mathbf{X} \\ 1016 &\stackrel{\textcircled{6}}{\Rightarrow} \mathbf{X} \mathbf{G}^\top \mathbf{X} = \mathbf{J} \mathbf{G} \mathbf{J} \\ 1017 &\stackrel{\textcircled{7}}{\Rightarrow} \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J} = \mathbf{G}, \end{aligned} \tag{23}$$

1022 where step ① uses the results of left-multiplying both sides by  $\mathbf{J} \mathbf{X}^\top$ ; step ② uses  $\mathbf{J} \cdot \mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J} \mathbf{J} = \mathbf{I}$ ;  
1023 step ③ uses the results of left-multiplying both sides by  $\mathbf{X}$  and subsequently right-  
1024 multiplying them by  $\mathbf{X}^\top$ ; ④ uses  $\mathbf{G} = \mathbf{J} \mathbf{X} \mathbf{G}^\top \mathbf{X} \mathbf{J}$ ; step ⑤ uses the the results of right-  
1025 multiplying both sides by  $\mathbf{J} \mathbf{X}$ ; 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 right-multiplied by  $\mathbf{J}$ .

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

## D.2 PROOF OF THEOREM 3.3

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

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

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

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

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

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

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 following first-order optimality condition:

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

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

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

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

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

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

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

## E PROOFS FOR SECTION 4

### E.1 PROOF OF LEMMA 4.5

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

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

1080 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.  $\iota^{t-1}$   
1081 and  $\mathbb{E}_{\iota^t}[\|\frac{1}{b'} \sum_{i=1}^{b'} (\nabla f_i(\mathbf{X}^t) - \nabla f_i(\mathbf{X}^{t-1})) - \nabla f(\mathbf{X}^t) + \nabla f(\mathbf{X}^{t-1})\|_F^2] = 0$ . We further have  
1082

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

1095 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  
1096 lemma A.2; step ③ uses assumption 4.3, Inequality (2) and Part (b) of lemma 2.5.  $\square$   
1097

## 1098 E.2 PROOF OF THEOREM 4.6

1100 *Proof.* For simplicity, we use  $\mathbf{B}$  instead of  $\mathbf{B}^t$ . We will show that the following inequality  
1101 holds :

$$1102 \frac{\theta}{2}\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F^2 \leq f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}). \tag{28}$$

1103 Since  $\bar{\mathbf{V}}^t$  is the global optimal solution of Problem (5), we have:

$$1105 \mathcal{G}(\bar{\mathbf{V}}^t; \mathbf{X}^t, \mathbf{B}) \leq \mathcal{G}(\mathbf{V}; \mathbf{X}^t, \mathbf{B}), \mathbf{V} \in \mathcal{J}_{\mathbf{B}}$$

1106 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:  
1107

$$1108 \frac{1}{2}\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbf{Q}+\theta\mathbf{I}}^2 + \langle \bar{\mathbf{V}}^t - \mathbf{I}_2, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} \rangle \leq 0. \tag{29}$$

1109 Using Inequality (2) with  $N = 1$  and Part (c) of Lemma 2.1, we have:  
1110

$$1111 f(\mathbf{X}^{t+1}) \leq f(\mathbf{X}^t) + \langle \bar{\mathbf{V}}^t - \mathbf{I}_2, [\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top]_{\mathbf{BB}} \rangle + \frac{1}{2}\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbf{Q}}^2. \tag{30}$$

1112 Adding Inequality (29) and (30) together, we obtain the inequality in (28). Using the  
1113 result of Part (b) in Lemma 2.1 that  $\frac{\|\mathbf{X}^+ - \mathbf{X}\|_F^2}{\|\mathbf{X}\|_F^2} \leq \|\mathbf{V} - \mathbf{I}_2\|_F^2$ , we have the following sufficient  
1114 decrease condition:  
1115

$$1116 f(\mathbf{X}^{t+1}) - f(\mathbf{X}^t) \leq -\frac{\theta}{2}\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F^2 \leq -\frac{\theta}{2}\frac{\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2}{\|\mathbf{X}^t\|_F^2} \tag{31}$$

1117 We now prove the global convergence. Taking the expectation for Inequality (31), we obtain  
1118 a lower bound on the expected progress made by each iteration for Algorithm 1:  
1119

$$1121 \mathbb{E}_{\xi^{t+1}}[f(\mathbf{X}^{t+1})] - \mathbb{E}_{\xi^t}[f(\mathbf{X}^t)] \leq -\mathbb{E}_{\xi^t}[\frac{\theta}{2}\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F^2].$$

1122 Summing up the inequality above over  $t = 0, 1, \dots, T$ , we have:  
1123

$$1124 \mathbb{E}_{\xi^T}[\frac{\theta}{2} \sum_{t=0}^T \|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F^2] \leq f(\mathbf{X}^0) - \mathbb{E}_{\xi^{T+1}}[f(\mathbf{X}^{T+1})] \leq f(\mathbf{X}^0) - f(\bar{\mathbf{X}}).$$

1125 As a result, there exists an index  $\bar{t}$  with  $0 \leq \bar{t} \leq T$  such that  
1126

$$1127 \mathbb{E}_{\xi^{\bar{t}}}[\|\bar{\mathbf{V}}^{\bar{t}} - \mathbf{I}_2\|_F^2] \leq \frac{2}{\theta(T+1)}[f(\mathbf{X}^0) - f(\bar{\mathbf{X}})]. \tag{32}$$

1128 Furthermore, for any  $t$ , we have:  
1129

$$1130 \mathcal{E}(\mathbf{X}^t) \triangleq \frac{1}{C_n^2} \sum_{i=1}^{C_n^2} \text{dist}(\mathbf{I}_2, \arg \min_{\mathbf{V}} \mathcal{G}(\mathbf{V}; \mathbf{X}^t, \mathcal{B}_i))^2 = \mathbb{E}_{\xi^t}[\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F^2] \tag{33}$$

1131 Combining Inequality (32) and equality (33), we have the following result:  
1132

$$1133 \mathbb{E}_{\xi^t}[\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F^2] = \mathcal{E}(\mathbf{X}^{\bar{t}}) \leq \frac{2(f(\mathbf{X}^0) - f(\bar{\mathbf{X}}))}{\theta(T+1)}. \tag{34}$$

We will give the arithmetic operations of **GS-JOB****CD**. 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  $\epsilon$ -BS-point is

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

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

### 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:

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

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:

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

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

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

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

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

where step ① uses Part (**d**) of Lemma 2.5.

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

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

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

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

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

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

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

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

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

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

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

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

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  $t \geq 1$

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

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

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

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

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

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

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

□

#### E.4 PROOF OF THEOREM 4.10

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

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

Then, we prove the following important lemmas.

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

*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

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

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

---

1242   **Lemma E.2.** (*Riemannian gradient Lower Bound for the Iterates Gap*) We de-  
 1243    fine  $\phi \triangleq (3\bar{X} + \bar{V}\bar{X})\bar{G} + (1 + \bar{V}^2 + \frac{n}{2}(\bar{X}^2 + \bar{V}^2\bar{X}^2))L_f + (1 + \bar{V}^2)\theta$ . It holds that:  
 1244  
 1245  $\mathbb{E}_{t+1}[\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}\mathcal{T}(\mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1}))] \leq \phi \cdot \mathbb{E}_{t+1}[\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbf{F}}] + \frac{np\sqrt{u^t}}{2}(\bar{X} + \bar{V}^2\bar{X})$ .

1246  
 1247   *Proof.* For notation simplicity, we define:  
 1248

1249    $\Omega_{i0} \triangleq \mathbf{U}_{\mathbf{B}_{(i)}}^\top [\tilde{\mathbf{G}}^{t+1}] [\mathbf{X}^{t+1}]^\top \mathbf{U}_{\mathbf{B}_{(i)}}, \forall i$    (44)

1250    $\Omega_{i1} \triangleq \mathbf{U}_{\mathbf{B}_{(i)}}^\top [\tilde{\mathbf{G}}^{t+1}] [\mathbf{X}^t]^\top \mathbf{U}_{\mathbf{B}_{(i)}}, \forall i$ ,   (45)

1251    $\Omega_{i2} \triangleq \mathbf{U}_{\mathbf{B}_{(i)}}^\top [\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^{t+1}] [\mathbf{X}^t]^\top \mathbf{U}_{\mathbf{B}_{(i)}}, \forall i$ .   (46)

1252  
 1253   First, using the optimality of  $\bar{\mathbf{V}}_i^t, i \in \{1, \dots, \frac{n}{2}\}$  for the subproblem, we have:  
 1254

1255    $\mathbf{0}_{2,2} = \tilde{\mathbf{G}}_i - \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t \tilde{\mathbf{G}}_i^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}}$    (47)

1256   where  $\tilde{\mathbf{G}}_i = \underbrace{\text{mat}((\mathbf{Q} + \theta\mathbf{I}_2) \text{vec}(\bar{\mathbf{V}}_i^t - \mathbf{I}_2))}_{\triangleq \Upsilon_{i1}} + \underbrace{\mathbf{U}_{\mathbf{B}_{(i)}}^\top \tilde{\mathbf{G}}^t (\mathbf{X}^t)^\top \mathbf{U}_{\mathbf{B}_{(i)}}}_{\triangleq \Upsilon_{i2}}$ .   (48)

1257  
 1258   Using the relation that  $\tilde{\mathbf{G}}_i = \Upsilon_{i1} + \Upsilon_{i2}$ , we obtain the following results from the above  
 1259   equality:  
 1260

1261    $\mathbf{0}_{2,2} = (\Upsilon_{i1} + \Upsilon_{i2}) - \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Upsilon_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}}$   
 1262  
 1263    $\stackrel{\textcircled{1}}{\Rightarrow} \mathbf{0}_{2,2} = \Upsilon_{i1} + \Omega_{i1} + \Omega_{i2} - \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Omega_{i1} + \Omega_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}}$   
 1264    $\Rightarrow \Omega_{i1} = \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Omega_{i1} + \Omega_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}} - \Upsilon_{i1} - \Omega_{i2}$ ,   (49)

1265  
 1266   where step ① uses  $\Upsilon_{i2} = \Omega_{i1} + \Omega_{i2}$ . Then we derive the following results:  
 1267

1268    $\mathbb{E}_{t+1}[\text{dist}(\mathbf{0}, \nabla_{\mathcal{J}}\mathcal{T}(\mathbf{V}_: = \mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1}))] = \mathbb{E}_{t+1}[\|\nabla_{\mathcal{J}}\mathcal{T}(\mathbf{V}_: = \mathbf{I}_2; \mathbf{X}^{t+1}, \mathbf{B}^{t+1})\|_{\mathbf{F}}]$   
 1269  
 1270    $\stackrel{\textcircled{1}}{=} \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}^{t+1}}^\top (\tilde{\mathbf{G}}^{t+1} [\mathbf{X}^{t+1}]^\top - \mathbf{J} \mathbf{X}^{t+1} [\tilde{\mathbf{G}}^{t+1}]^\top \mathbf{J}) \mathbf{U}_{\mathbf{B}_{(i)}^{t+1}}\|_{\mathbf{F}}]$   
 1271  
 1272    $\stackrel{\textcircled{2}}{=} \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \mathbf{U}_{\mathbf{B}_{(i)}}^\top (\tilde{\mathbf{G}}^{t+1} [\mathbf{X}^{t+1}]^\top - \mathbf{J} \mathbf{X}^{t+1} [\tilde{\mathbf{G}}^{t+1}]^\top \mathbf{J}) \mathbf{U}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1273  
 1274    $\stackrel{\textcircled{3}}{=} \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i0}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1275  
 1276    $\stackrel{\textcircled{4}}{=} \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} (\Omega_{i0} - \Omega_{i1}) + \Omega_{i1} - (\mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i0}^\top \mathbf{J}_{\mathbf{B}_{(i)}} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}) - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1277  
 1278    $\stackrel{\textcircled{5}}{\leq} \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbf{F}}] + \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i0}^\top \mathbf{J}_{\mathbf{B}_{(i)}} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1279  
 1280    $+ \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i1} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1281  
 1282    $\stackrel{\textcircled{6}}{\leq} \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbf{F}}] + \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0}^\top - \Omega_{i1}^\top\|_{\mathbf{F}}] + \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i1} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1283  
 1284    $\stackrel{\textcircled{7}}{\leq} 2\mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbf{F}}] + \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i1} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1285  
 1286    $\stackrel{\textcircled{8}}{=} 2\mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbf{F}}]$   
 1287  
 1288    $+ \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t (\Upsilon_{i1} + \Omega_{i1} + \Omega_{i2})^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}} - \Upsilon_{i1} - \Omega_{i2} - \mathbf{J}_{\mathbf{B}_{(i)}} \Omega_{i1}^\top \mathbf{J}_{\mathbf{B}_{(i)}}\|_{\mathbf{F}}]$   
 1289  
 1290    $\stackrel{\textcircled{9}}{\leq} 2\mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \Omega_{i0} - \Omega_{i1}\|_{\mathbf{F}}] + \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t \Upsilon_{i1}^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}} - \Upsilon_{i1}\|_{\mathbf{F}}] +$   
 1291  
 1292    $\mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \bar{\mathbf{V}}_i^t \Omega_{i1}^\top \bar{\mathbf{V}}_i^t - \Omega_{i1}^\top\|_{\mathbf{F}}] + \mathbb{E}_{t+1}[\|\sum_{i=1}^{n/2} \mathbf{J}_{\mathbf{B}_{(i)}} \bar{\mathbf{V}}_i^t \Omega_{i2}^\top \bar{\mathbf{V}}_i^t \mathbf{J}_{\mathbf{B}_{(i)}} - \Omega_{i2}\|_{\mathbf{F}}]$    (50)

1293   where step ① uses Equality (43); step ② uses the fact that both the working set  $\mathbf{B}^t$  and  
 1294    $\mathbf{B}^{t+1}$  are selected randomly and uniformly; step ③ uses the definition of  $\Omega_{i0}$  in (44); step ④  
 1295   uses  $-\Omega_{i1} + \Omega_{i1} = \mathbf{0}$  and  $-\Omega_{i1}^\top + \Omega_{i1}^\top = \mathbf{0}$ ; step ⑤ uses the norm inequality; step ⑥ uses the  
 1296   norm inequality; step ⑦ uses the norm inequality; step ⑧ uses Equality (49); step ⑨ uses  
 1297   the norm inequality. We now establish individual bounds for each term for Inequality (50).

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

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

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

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

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

where step ① uses the triangle inequality; step ② uses the inequality  $\|\mathbf{XY}\|_F \leq \|\mathbf{X}\|_F \|\mathbf{Y}\|_F$  for all  $\mathbf{X}$  and  $\mathbf{Y}$  and  $\forall t, \|\mathbf{V}^t\|_F \leq \bar{V}$ ; step ③ uses the definition of  $\Upsilon_{i1}$ ; step ④ uses the choice of  $\mathbf{Q} \preceq L_f \mathbf{I}$  and the norm inequality.

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

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

where step 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 ② uses the norm inequality and  $\forall t, \|\mathbf{V}^t\|_F \leq \bar{V}$ ; step ③ uses the fact that  $\|\Omega_{i1}\|_F = \|\mathbf{U}_{B(i)}^\top \tilde{\mathbf{G}}^t [\mathbf{X}^t]^\top \mathbf{U}_{B(i)}\|_F \leq \bar{X} \|\tilde{\mathbf{G}}^t\|_F, \forall i$  which can be derived using the norm inequality ; step ④ uses the fact that  $\forall \mathbf{X}, \|\tilde{\mathbf{G}}^t\|_F \leq \bar{G}$ .

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

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

where step ① uses the triangle inequality; step ② uses the norm inequality and  $\forall t, \|\mathbf{V}^t\|_{\mathbb{F}} \leq \bar{V}$ ; step ③ uses the definition of  $\forall i, \Omega_{i2} = \mathbf{U}_{B(i)}^\top [\tilde{\mathbf{G}}^t - \tilde{\mathbf{G}}^t][\mathbf{X}^t]^\top \mathbf{U}_{B(i)}$  in (46); step ④ uses the norm inequality and  $\forall t, \|\mathbf{X}^t\|_{\mathbb{F}} \leq \bar{X}$ ; step ⑤ uses Lemma E.1; step ⑥ uses Part (b) in Lemma 2.5 and  $\forall t, \|\mathbf{X}^t\|_{\mathbb{F}} \leq \bar{X}$ .

In view of( 51), (52), (53), (54), and (50), we have:

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

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  $c_4 = \frac{n}{2}(\bar{X}^2 + \bar{V}^2 \bar{X}^2)L_f$ .  $\square$

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

*Proof.* We have the following inequalities:

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

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

We Consider  $\|\nabla f(\mathbf{X}^t)(\mathbf{X}^t)^\top - \mathbf{J}\mathbf{X}^t(\nabla f(\mathbf{X}^t))^\top \mathbf{J}\|_{\mathbb{F}}$ :

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

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

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

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 definition of  $\nabla_{\mathcal{J}} \mathcal{T}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})$ .  $\square$

We now present the following useful lemma.

**Lemma E.4.** *We define  $\mathbf{T}_{\mathbf{X}} \mathcal{J} \triangleq \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \mathcal{A}_X(\mathbf{Y}) = \mathbf{0}\}$  and  $\mathcal{A}_X(\mathbf{Y}) \triangleq \mathbf{X}^\top \mathbf{J}\mathbf{Y} + \mathbf{Y}^\top \mathbf{J}\mathbf{X}$ . 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:*

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

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

1404 *Proof.* We note that  $\bar{\mathbf{Y}} = \arg \min_{\mathbf{Y} \in \mathbf{T}_{\mathbf{X}} \mathcal{J}} \frac{1}{2} \|\mathbf{Y} - \mathbf{G}\|_F^2 = \arg \min_{\mathbf{Y}} \frac{1}{2} \|\mathbf{Y} - \mathbf{G}\|_F^2$ , s.t.  
1405  $\mathbf{X}^\top \mathbf{JY} + \mathbf{Y}^\top \mathbf{JX} = \mathbf{0}$ . Introducing a multiplier  $\Lambda \in \mathbb{R}^{n \times n}$  for the linear con-  
1406 straints  $\mathbf{X}^\top \mathbf{JY} + \mathbf{Y}^\top \mathbf{JX} = \mathbf{0}$ , we have following Lagrangian function:  $\tilde{\mathcal{L}}(\mathbf{Y}; \Lambda) =$   
1407  $\frac{1}{2} \|\mathbf{Y} - \mathbf{G}\|_F^2 + \langle \mathbf{X}^\top \mathbf{JY} + \mathbf{Y}^\top \mathbf{JX}, \Lambda \rangle$ . We naturally derive the following first-order op-  
1408 timality condition:  $\mathbf{Y} - \mathbf{G} + \mathbf{JX}\Lambda = \mathbf{0}$ ,  $\mathbf{X}^\top \mathbf{JY} + \mathbf{Y}^\top \mathbf{JX} = \mathbf{0}$ . Incorporating the term  
1409  $\mathbf{Y} = \mathbf{G} - \mathbf{JX}\Lambda$  into  $\mathbf{X}^\top \mathbf{JY} + \mathbf{Y}^\top \mathbf{JX} = \mathbf{0}$ , we obtain:  
1410

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

1411 Any  $\Lambda$  satisfying formula (55) is a feasible point, so we can easily find :  
1412

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

1413 where step ① uses the fact that any matrix  $\mathbf{X}$  satisfying the  $\mathbf{J}$ -orthogonality constraint has  
1414 a determinant of 1 or -1, thus  $\text{inv}(\mathbf{X})$  exists; step ② multiply both sides of the equation by  
1415  $\mathbf{XJ}$ ; step ③ uses  $\mathbf{X}^\top \mathbf{JX} = \mathbf{J}$  and  $\mathbf{JJ} = \mathbf{I}$ ; step ④ multiply both sides of the equation by  $\mathbf{J}$   
1416 and uses  $\mathbf{JJ} = \mathbf{I}$ ; step ⑤ uses the fact that  $\Lambda$  is a symmetric matrix.  
1417

1418 Therefore, a feasible solution  $\mathbf{Y}$  can be computed as  $\mathbf{Y} = \mathbf{G} - \mathbf{JX}\Lambda = \mathbf{G} - \mathbf{JXG}^\top \mathbf{XJ}$ . Since  
1419  $\bar{\mathbf{Y}}$  is the optimal solution, there must be  $h(\bar{\mathbf{Y}}) \leq h(\mathbf{G} - \mathbf{JXG}^\top \mathbf{XJ})$ .  $\square$   
1420

1421 We now present the proof of this lemma.  
1422

1423 **Lemma E.5.** *For any  $\mathbf{X} \in \mathbb{R}^{n \times n}$ , it holds that  $\text{dist}(\mathbf{0}, \nabla f^\circ(\mathbf{X})) \leq \text{dist}(\mathbf{0}, \nabla_{\mathcal{J}} f(\mathbf{X}))$ .*  
1424

1425 *Proof.* For the purpose of analysis, we define the nearest  $\mathbf{J}$  orthogonal matrix to an arbitrary  
1426 matrix  $\mathbf{Y} \in \mathbb{R}^{n \times n}$  is given by  $\mathcal{P}_{\mathcal{J}}(\mathbf{Y})$ . Similarly, we have  $\mathcal{P}_{\mathbf{T}_{\mathbf{X}} \mathcal{J}}(\nabla f(\mathbf{X}))$  for projecting  
1427 gradient  $\nabla f(\mathbf{X})$  into space  $\mathbf{T}_{\mathbf{X}} \mathcal{J}$ .  
1428

1429 We recall that the following first-order optimality conditions are equivalent for all  $\mathbf{X} \in \mathbb{R}^{n \times n}$   
1430 :  
1431

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

1432 Therefore, we derive the following results:  
1433

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

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

1434 We let  $\mathbf{G} \in \nabla f(\mathbf{X})$  and obtain the following results from the above equality:  
1435

$$\text{dist}(\mathbf{0}, \nabla f^\circ(\mathbf{X})) \stackrel{\textcircled{1}}{\leq} \|\mathbf{G} - \mathbf{JXG}^\top \mathbf{XJ}\|_F, \quad (60)$$

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

1436 where step ① uses Lemma E.4; step ② uses  $\nabla_{\mathcal{J}} f(\mathbf{X}) = \mathbf{G} - \mathbf{JXG}^\top \mathbf{XJ}$  with  $\mathbf{G} \in \nabla f(\mathbf{X})$ .  $\square$   
1437

1438 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  
1439 that:  
1440

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

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

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

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:

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

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

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

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

$$(67)$$

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

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

where step ① uses the sufficient descent condition as shown in Theorem 4.7; step ② uses Inequality (64) and (62) with  $\mathbf{X}' = \mathbf{X}^t$  and  $\mathbf{X} = \bar{\mathbf{X}}$ ; step ③ uses lemma E.3 ; step ④ uses Lemma E.2 ; step ⑤ uses  $\forall x_i \in \mathbb{R}, \frac{x_1 + \dots + x_n}{n} \leq \sqrt{\frac{x_1^2 + \dots + x_n^2}{n}}$  ; step ⑥ applies the inequality that  $\forall \theta' > 0, a, b, ab \leq \frac{\theta' a^2}{2} + \frac{b^2}{2\theta'}$  with  $a = \sqrt{\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^{t-1} - \mathbf{I}_2\|_{\mathbb{F}}^2}, b = \mathcal{E}^t \gamma \phi \sqrt{\frac{n}{2}}, a = \sqrt{\mathbb{E}_{\ell^t}[u^t]}, b = \mathcal{E}^t (2\bar{X}^2 + \gamma \frac{np}{2} \bar{X} + \gamma \frac{np}{2} \bar{V}^2 \bar{X})$ ; step ⑦ denote  $\mathfrak{A}^2 \triangleq \frac{(2\bar{X}^2 + \gamma \frac{np}{2} \bar{X} + \gamma \frac{np}{2} \bar{V}^2 \bar{X})^2}{2\theta} + \frac{n\gamma^2\phi^2}{4\theta'}$ . To simplify the formula, we define  $\mathfrak{N}^t = \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^t - \mathbf{I}_2\|_{\mathbb{F}}^2$ .

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

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

To recursively eliminate term  $\sqrt{(\theta+1)\mathbb{E}_{t^t}[u^t]}$ , we take the root of both sides of the Inequality in Lemma 4.5:

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

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

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

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:

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

Rearranging terms, we have:

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

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

$$\begin{aligned} &\mathbb{E}_{t^T}[\sqrt{\theta-\bar{X}^2}\sqrt{\mathbf{N}^T}] + \mathbb{E}_{t^{T-1}}[\frac{\sqrt{\theta-\bar{X}^2}}{2}\sum_{t=i}^{T-1}\sqrt{\mathbf{N}^t}] \\ &\leq \mathfrak{A}\sum_{t=i}^T\mathcal{E}^t + \frac{\sqrt{1-p}\sqrt{(\theta+1)}}{1-\sqrt{1-p}}(\sqrt{\mathbb{E}_{t^{i-1}}[u^{i-1}]} - \sqrt{\mathbb{E}_{t^T}[u^T]}) + \\ &\quad \frac{(T-i+1)\sqrt{\theta+1}}{1-\sqrt{1-p}}\sqrt{\frac{p(N-b)}{b(N-1)}\sigma^2} + \frac{\sqrt{\theta-\bar{X}^2}}{2}\mathbb{E}_{t^{i-1}}[\sqrt{\mathbf{N}^{i-1}}] \\ &\stackrel{\textcircled{1}}{\leq} \mathfrak{A}\sum_{t=i}^T\mathcal{E}^t + \frac{\sqrt{1-p}\sqrt{(\theta+1)}}{1-\sqrt{1-p}}\sqrt{\mathbb{E}_{t^{i-1}}[u^{i-1}]} + \frac{(T-i+1)\sqrt{\theta+1}}{1-\sqrt{1-p}}\sqrt{\frac{p(N-b)}{b(N-1)}\sigma^2} + \frac{\sqrt{\theta-\bar{X}^2}}{2}\mathbb{E}_{t^{i-1}}[\sqrt{\mathbf{N}^{i-1}}] \end{aligned}$$

where step ① uses the fact that  $\mathbb{E}_{t^T}[u^T] \geq 0$ .

Since  $b = N, b' = \sqrt{b}, p = \frac{b'}{b+b'}$ , we have  $\frac{(T+i-1)\sqrt{\theta+1}}{1-\sqrt{1-p}}\sqrt{\frac{p(N-b)}{b(N-1)}\sigma^2} = 0$ . Rearranging terms, we have:

$$\mathbb{E}_{t^T}[\frac{\theta-\bar{X}^2}{2}\sum_{t=i}^T\sqrt{\mathbf{N}^t}] \leq \mathfrak{A}\sum_{t=i}^T\mathcal{E}^t + \frac{\sqrt{1-p}\sqrt{(\theta+1)}}{1-\sqrt{1-p}}\sqrt{\mathbb{E}_{t^{i-1}}[u^{i-1}]} + \frac{\sqrt{\theta-\bar{X}^2}}{2}\mathbb{E}_{t^{i-1}}[\sqrt{\mathbf{N}^{i-1}}] \quad (74)$$

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

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

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

$$\begin{aligned} \mathbb{E}_{t^T}[\frac{\theta-\bar{X}^2}{2}\sum_{t=i}^T\sqrt{\mathbf{N}^t}] &\leq \mathfrak{A}\varphi(f(\mathbf{X}^i) - f(\bar{\mathbf{X}})) + \\ &\quad \frac{\sqrt{1-p}\sqrt{(\theta+1)}}{1-\sqrt{1-p}}\sqrt{\mathbb{E}_{t^{i-1}}[u^{i-1}]} + \frac{\sqrt{\theta-\bar{X}^2}}{2}\mathbb{E}_{t^i}[\sqrt{\mathbf{N}^{i-1}}] \end{aligned} \quad (76)$$

Using  $\forall t, \|\mathbf{V}^t\|_{\text{F}} \leq \bar{V}$ , we have the fact that  $\|\mathbf{V}_i - \mathbf{I}_2\|_{\text{F}}^2 \leq (\|\mathbf{V}_i\|_{\text{F}} + \|\mathbf{I}_2\|_{\text{F}})^2 \leq (\bar{X} + \sqrt{2})^2$  and  $\sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i^0 - \mathbf{I}_2\|_{\text{F}}^2 \leq \frac{n}{2}(\bar{V} + \sqrt{2})^2$ . Using the inequality that  $\frac{\|\mathbf{x}^+ - \mathbf{x}\|_{\text{F}}^2}{\bar{X}^2} \leq \frac{\|\mathbf{x}^+ - \mathbf{x}\|_{\text{F}}^2}{\|\mathbf{x}\|_{\text{F}}^2} \leq \sum_{i=1}^{n/2} \|\bar{\mathbf{V}}_i - \mathbf{I}_2\|_{\text{F}}^2$  as shown in Part (b) in Lemma 2.5 and letting  $i = 1$ , we have:

$$\begin{aligned} \mathbb{E}_{t^0} \left[ \frac{\theta - \bar{X}}{2\bar{X}} \sum_{t=1}^T \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\text{F}} \right] &\leq \mathfrak{A} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \\ &\quad \frac{\sqrt{1-p}\sqrt{(\bar{\theta}+1)}}{1-\sqrt{1-p}} \sqrt{\mathbb{E}_{t^0}[u^0]} + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2}(\bar{V} + \sqrt{2})^2} \end{aligned}$$

Since  $\mathbb{E}_{t^0}[u^0] \leq \frac{N-b}{b(N-1)}\sigma^2 = 0$ , we have:

$$\mathbb{E}_{t^0} \left[ \frac{\theta - \bar{X}}{2\bar{X}} \sum_{t=1}^T \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\text{F}} \right] \leq \mathfrak{A} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \frac{\sqrt{\theta - \bar{X}^2}}{2} \sqrt{\frac{n}{2}(\bar{V} + \sqrt{2})^2}$$

We can get the expression for C:

$$\mathbb{E}_{t^0} \left[ \sum_{j=1}^t \|\mathbf{X}^{j+1} - \mathbf{X}^j\|_{\text{F}} \right] \leq C$$

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

## E.5 PROOF OF THEOREM 4.9

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

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

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

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

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

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

*Proof.* We have the following inequalities:

$$\begin{aligned} \|\nabla_{\mathcal{J}}f(\mathbf{X}^t)\|_{\text{F}}^2 &\stackrel{\textcircled{1}}{=} \|\mathbf{G}^t - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^{\top} \mathbf{X}^t \mathbf{J}\|_{\text{F}}^2 \\ &\stackrel{\textcircled{2}}{=} \|\mathbf{G}^t(\mathbf{X}^t)^{\top} \mathbf{J}\mathbf{X}^t \mathbf{J} - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^{\top} \mathbf{J}\mathbf{J}\mathbf{X}^t \mathbf{J}\|_{\text{F}}^2 \\ &\stackrel{\textcircled{3}}{\leq} \|\mathbf{G}^t(\mathbf{X}^t)^{\top} - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^{\top} \mathbf{J}\|_{\text{F}}^2 \|\mathbf{J}\mathbf{X}^t \mathbf{J}\|_{\text{F}}^2 \\ &\stackrel{\textcircled{4}}{\leq} \|\mathbf{X}^t\|_{\text{F}}^2 \|\mathbf{W}\|_{\text{F}}^2, \text{ with } \mathbf{W} \triangleq \mathbf{G}^t(\mathbf{X}^t)^{\top} - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^{\top} \mathbf{J} \\ &\stackrel{\textcircled{5}}{\leq} \|\mathbf{X}^t\|_{\text{F}}^2 C_n^2 \cdot \mathbb{E}_{\xi^{t-1}}[\|\mathbf{U}_{\mathbf{B}}^{\top}[\mathbf{G}^t(\mathbf{X}^t)^{\top} - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^{\top} \mathbf{J}] \mathbf{U}_{\mathbf{B}}\|_{\text{F}}^2] \\ &\stackrel{\textcircled{6}}{=} \|\mathbf{X}^t\|_{\text{F}}^2 C_n^2 \cdot \mathbb{E}_{\xi^{t-1}}[\|\nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_{\text{F}}^2] \\ &\stackrel{\textcircled{7}}{\leq} \bar{X}^2 C_n^2 \cdot \mathbb{E}_{\xi^{t-1}}[\|\nabla_{\mathcal{J}}\mathcal{G}(\mathbf{I}_2; \mathbf{X}^t, \mathbf{B})\|_{\text{F}}^2] \end{aligned} \tag{77}$$

where step  $\textcircled{1}$  uses the definition of  $\nabla_{\mathcal{J}}f(\mathbf{X}^t)$ ; step  $\textcircled{2}$  uses  $\mathbf{J}\mathbf{J} = \mathbf{I}$  and  $\mathbf{X}^{\top} \mathbf{J}\mathbf{X} = \mathbf{J} \Rightarrow \mathbf{X}^{\top} \mathbf{J}\mathbf{X}\mathbf{J} = \mathbf{J}\mathbf{J} = \mathbf{I}$ ; step  $\textcircled{3}$  uses the norm inequality and ; step  $\textcircled{4}$  uses the definition of  $\mathbf{W} \triangleq \mathbf{G}^t(\mathbf{X}^t)^{\top} - \mathbf{J}\mathbf{X}^t(\mathbf{G}^t)^{\top} \mathbf{J}$ ; step  $\textcircled{5}$  uses Lemma (A.1) with  $k = 2$ ; step  $\textcircled{6}$  uses the definition 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; step  $\textcircled{7}$  uses  $\forall t, \|\mathbf{X}^t\|_{\text{F}} \leq \bar{X}$ .  $\square$

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

$$\frac{1}{\varphi'(f^\circ(\mathbf{X}') - f^\circ(\mathbf{X}))} \leq \text{dist}(0, \nabla f^\circ(\mathbf{X}')) \stackrel{\textcircled{1}}{\leq} \|\nabla_J f(\mathbf{X}')\|_F, \quad (78)$$

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

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

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:

$$\begin{aligned} & (f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}))\varphi'(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) \\ & \leq \varphi(f(\mathbf{X}^t) - f(\bar{\mathbf{X}})) - \varphi(f(\mathbf{X}^{t+1}) - f(\bar{\mathbf{X}})) \triangleq \mathcal{E}^t. \end{aligned} \quad (79)$$

We derive the following inequalities:

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

where step ① uses the sufficient descent condition as shown in Theorem 4.6; step ② uses Inequality (79); step ③ uses Inequality (78) with  $\mathbf{X}' = \mathbf{X}^t$  and  $\mathbf{X} = \bar{\mathbf{X}}$ ; step ④ uses Lemma E.7; step ⑤ uses 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  $a = \|\bar{\mathbf{V}}^{t-1} - \mathbf{I}_2\|_F$  and  $b = \mathcal{E}^t \gamma \phi$ .

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

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

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 inequality above over  $t = i, 2, \dots, T$ , we have:

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

where step ① uses the definition of  $\mathcal{E}^i$  in (79); step ② uses a basic recursive reduction; step ③ uses the fact the desingularization function  $\varphi(\cdot)$  is positive. With the choice  $\theta' = \frac{\theta}{4}$ , we have:

$$\begin{aligned} & \sqrt{\theta} \mathbb{E}_{\xi^T} [\|\bar{\mathbf{V}}^T - \mathbf{I}_2\|_F] + \frac{\sqrt{\theta}}{2} \sum_{t=i}^{T-1} \mathbb{E}_{\xi^t} [\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_F] \\ & \leq \frac{2\gamma\phi}{\sqrt{\theta}} \varphi(f(\mathbf{X}^i) - f(\bar{\mathbf{X}})) + \frac{\sqrt{\theta}}{2} \mathbb{E}_{\xi^{i-1}} [\|\bar{\mathbf{V}}^{i-1} - \mathbf{I}_2\|_F] \end{aligned} \quad (80)$$

1674 We obtain from Inequality (80):  
1675

$$\begin{aligned} 1676 \quad & \frac{1}{2} \sum_{t=i}^T \mathbb{E}_{\xi^t} [\|\bar{\mathbf{V}}^t - \mathbf{I}_2\|_{\mathbb{F}}] \leq \frac{2\gamma\phi}{\theta} \varphi(f(\mathbf{X}^i) - f(\bar{\mathbf{X}})) + \frac{1}{2} \mathbb{E}_{\xi^{i-1}} [\|\bar{\mathbf{V}}^{i-1} - \mathbf{I}_2\|_{\mathbb{F}}] \quad (81) \\ 1677 \quad \stackrel{\textcircled{1}}{\Rightarrow} \quad & \frac{1}{2} \sum_{t=i}^T \mathbb{E}_{\xi^t} [\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}] \leq (\frac{2\bar{X}\gamma\phi}{\theta} \varphi(f(\mathbf{X}^i) - f(\bar{\mathbf{X}})) + \frac{\bar{X}}{2}(\bar{V} + \sqrt{2})) \end{aligned}$$

1679 where step ① uses  $\forall t, \|\mathbf{V}\|_{\mathbb{F}} \leq \bar{V}$ , then  $\|\mathbf{V} - \mathbf{I}_2\|_{\mathbb{F}} \leq \|\mathbf{V}\|_{\mathbb{F}} + \|\mathbf{I}_2\|_{\mathbb{F}} \leq \bar{V} + \sqrt{2}$  and the inequality  
1680 that  $\frac{\|\mathbf{X}^{i+1} - \mathbf{X}^i\|_{\mathbb{F}}}{\bar{X}} \leq \|\bar{\mathbf{V}}^i - \mathbf{I}_2\|_{\mathbb{F}}$  as shown in Part (b) in Lemma 2.1. Finally, let  $i = 1$  we can  
1681 get:  
1682

$$\sum_{t=1}^T \mathbb{E}_{\xi^t} [\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}] \leq C$$

1684 where  $C \triangleq \frac{4\bar{X}\gamma\phi}{\theta} \varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})) + \bar{X}(\bar{V} + \sqrt{2}) \leq \mathcal{O}(\varphi(f(\mathbf{X}^1) - f(\bar{\mathbf{X}})))$ .  $\square$   
1685

## F ADDITIONAL EXPERIMENT DETAILS AND RESULTS

### F.1 ADDITIONAL DETAILS FOR HYPERBOLIC STRUCTURAL PROBE PROBLEM

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

► **Ultrahyperbolic manifold.** Vectors in an ultrahyperbolic manifold is defined as  $\mathbb{U}_{\alpha}^{p,q} = \{\mathbf{x} = (x_1, x_2, \dots, x_{p+q})^\top \in \mathbb{R}^{p,q} : \|\mathbf{x}\|_q^2 = -\alpha^2\}$ [49], where  $\alpha$  is a non-negative real number denoting 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 a norm of the induced scalar product. The hyperbolic and spherical manifolds can be defined as : $\mathbb{H}_{\alpha} = \mathbb{U}_{\alpha}^{p,1}$ ,  $\mathbb{S}_{\alpha} = \mathbb{U}_{\alpha}^{0,q}$ .

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

► **Diffeomorphism.** [Theorem 1 Diffeomorphism of [50]]: Any vector  $\mathbf{x} \in \mathbb{R}^p \times \mathbb{R}_*^q$  can be mapped into  $\mathbb{U}_{\alpha}^{p,q}$  by a double projection  $\varphi = \phi^{-1} \circ \phi$ , with  $\psi(\mathbf{x}) = (\frac{\mathbf{s}}{\alpha \|\mathbf{t}\|})$ ,  $\psi^{-1}(\mathbf{z}) = (\frac{\mathbf{v}}{\sqrt{\alpha^2 + \|\mathbf{v}\|^2}} \mathbf{u})$ , where  $\mathbf{x} = (\frac{\mathbf{s}}{\mathbf{t}}) \in \mathbb{U}_{\alpha}^{p,q}$  with  $\mathbf{s} \in \mathbb{R}^p$  and  $\mathbf{t} \in \mathbb{R}_*^q \cdot \mathbf{z} = (\frac{\mathbf{v}}{\mathbf{u}}) \in \mathbb{R}^p \times \mathbb{S}_{\alpha}^q$  with  $\mathbf{v} \in \mathbb{R}^p$  and  $\mathbf{u} \in \mathbb{S}_{\alpha}^q$ .

### F.2 ADDITIONAL APPLICATION: ULTRA-HYPERBOLIC KNOWLEDGE GRAPH EMBEDDING

The J orthogonal matrix can be used as an isometric linear operator in the Ultrahyperbolic manifold, [49] et al. extended the knowledge graph model from hyperbolic space to Ultra-hyperbolic space (named as **UltraE**) by this property. The **UltraE** model is formulated as follows:

$$\begin{aligned} 1717 \quad & \min_{\mathbf{R}, \mathbf{E}, \mathbf{b}} \mathcal{L}(\mathbf{R}, \mathbf{E}, \mathbf{b}) \triangleq -\frac{1}{N} \sum_{(h, r, t) \in \Delta} (\log s(h, r, t) + \sum_{(h', r', t') \in \Delta'_{(h, r, t)}} \log(1 - s(h', r', t'))) \\ 1718 \quad & \text{s.t. } \begin{cases} s(h, r, t) = \sigma(-d_{\alpha}^2(\mathbf{R}_r \mathbf{E}_h, \mathbf{E}_t) + \mathbf{b}_h + \mathbf{b}_t + \delta) \\ \mathbf{R}_r^\top \mathbf{J} \mathbf{R}_r = \mathbf{J} \end{cases} \end{aligned}$$

1722 where  $\mathbf{E} \in \mathbb{R}^{n_e \times n}$  with  $\mathbf{E}_h = \mathbf{E}(h, :) \in \mathbb{U}_{\alpha}^{p,q}$ ,  $\mathbf{b} \in \mathbb{R}^{n_r}$  with  $\mathbf{b}_h = \mathbf{b}(r) \in \mathbb{R}$ ,  $\mathbf{R} \in \mathbb{R}^{n_r \times n \times n}$   
1723 with  $\mathbf{R}_r = \mathbf{R}(r, :, :) \in \mathbb{R}^{n \times n}$  and  $\mathbf{J} = [\begin{smallmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_q \end{smallmatrix}]$ ;  $\Delta \in \mathbb{N}^{N \times 3}$  is the set of positive triplets,  
1724  $\Delta'_{(h, r, t)} \in \mathbb{N}^{N \times k \times 3}$  denotes the set of negative triples constructed by corrupting  $(h, r, t)$ ;  $\delta$   
1725 is a global margin hyper-parameter,  $\sigma(\cdot)$  is the sigmoid function,  $n_e$  represents the number  
1726 of entities and  $n_r$  represents the number of relations;  $d_{\alpha}(\cdot)$  stands for the Ultra-hyperbolic  
1727 geodesic distance (refer to F.1).

---

1728 **► Experiment Details.** We selected a batch of **FB15K** and **WN18RR** respectively  
 1729 as the data set for the Ultra-hyperbolic Knowledge Graph Embedding problem, (training  
 1730 set size, test set size, number of entities, number of relations) are (719,308,135,22) and  
 1731 (545,233,208,5) respectively.  $n = 36$ ,  $p = 18$ ,  $\delta = 5$ ,  $\alpha = 1$  and  $k = 50$ . In order to  
 1732 highlight the difference between J orthogonal optimization, in the **UltraE** model, all entities  
 1733 and biases of the optimization algorithm are optimized using **ADMM** by **Pytorch**,  $lr =$   
 1734  $5e-4$ . We use the **Adagrad** optimizer in Pytorch to optimize the J-orthogonality constraint  
 1735 variable in the **CS** model.  
 1736

### F.3 IMPLEMENTATION OF ADMM ALGORITHM FOR PROBLEM (1)

1738 We consider the following smooth J-orthogonality constraint problem:  
 1739  $\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} f(\mathbf{X})$ , s.t.  $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$ . Defining  $\mathbf{Y} = \mathbf{J} \mathbf{X} \in \mathbb{R}^{n \times n}$ , we have:  $\min_{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}} f(\mathbf{X})$   
 1740 s.t.  $\mathbf{X}^\top \mathbf{Y} = \mathbf{J}$ ,  $\mathbf{Y} = \mathbf{J} \mathbf{X}$ . Introducing Lagrange multipliers  $\mathbf{Z} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  
 1741  $\beta \in \mathbb{R}$ , we have the following Lagrange function:  
 1742

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}; \mathbf{Z}, \mathbf{W}) = f(\mathbf{X}) + \langle \mathbf{X}^\top \mathbf{Y} - \mathbf{J}, \mathbf{Z} \rangle + \langle \mathbf{J} \mathbf{X} - \mathbf{Y}, \mathbf{W} \rangle + \frac{\beta}{2} \|\mathbf{X}^\top \mathbf{Y} - \mathbf{J}\|_F^2 + \frac{\beta}{2} \|\mathbf{J} \mathbf{X} - \mathbf{Y}\|_F^2$$

1743 Supposing  $f(\mathbf{X})$  is  $l$ -Lipschitz gradient continuous:  $f(\mathbf{X}) \leq f(\mathbf{X}^t) + \langle \mathbf{X} - \mathbf{X}^t, \nabla f(\mathbf{X}^t) \rangle +$   
 1744  $\frac{l}{2} \|\mathbf{X} - \mathbf{X}^t\|_F^2$ , We get the following majorization function of  $\mathcal{L}(\mathbf{X}, \mathbf{Y}; \mathbf{Z}, \mathbf{W})$  at  $(\mathbf{X}^t, \mathbf{Y}; \mathbf{Z}, \mathbf{W})$ :  
 1745

$$\begin{aligned} \mathcal{L}(\mathbf{X}, \mathbf{Y}; \mathbf{Z}, \mathbf{W}) &\leq f(\mathbf{X}^t) + \langle \mathbf{X} - \mathbf{X}^t, \nabla f(\mathbf{X}^t) \rangle + \frac{l}{2} \|\mathbf{X} - \mathbf{X}^t\|_F^2 + \langle \mathbf{X}^\top \mathbf{Y} - \mathbf{J}, \mathbf{Z} \rangle + \\ &\quad \langle \mathbf{J} \mathbf{X} - \mathbf{Y}, \mathbf{W} \rangle + \frac{\beta}{2} \|\mathbf{X}^\top \mathbf{Y} - \mathbf{J}\|_F^2 + \frac{\beta}{2} \|\mathbf{J} \mathbf{X} - \mathbf{Y}\|_F^2 \end{aligned}$$

1746 We solve the following subproblem to update  $\mathbf{X}^{t+1}$  and  $\mathbf{Y}^{t+1}$  alternately:  
 1747

$$\begin{aligned} \mathbf{X}^{t+1} &= \arg \min_{\mathbf{X}} \mathcal{L}_{\mathbf{X}}(\mathbf{X}, \mathbf{Y}^t; \mathbf{Z}^t, \mathbf{W}^t) \triangleq \langle \mathbf{X} - \mathbf{X}^t, \nabla f(\mathbf{X}^t) \rangle + \frac{l}{2} \|\mathbf{X} - \mathbf{X}^t\|_F^2 + \\ &\quad \langle \mathbf{X}^\top \mathbf{Y} - \mathbf{J}, \mathbf{Z}^t \rangle + \langle \mathbf{J} \mathbf{X} - \mathbf{Y}, \mathbf{W}^t \rangle + \frac{\beta}{2} \|\mathbf{X}^\top \mathbf{Y} - \mathbf{J}\|_F^2 + \frac{\beta}{2} \|\mathbf{J} \mathbf{X} - \mathbf{Y}\|_F^2 \\ \mathbf{Y}^{t+1} &= \arg \min_{\mathbf{Y}} \mathcal{L}_{\mathbf{Y}}(\mathbf{X}^{t+1}, \mathbf{Y}; \mathbf{Z}^t, \mathbf{W}^t) \triangleq \langle \mathbf{X}^\top \mathbf{Y} - \mathbf{J}, \mathbf{Z}^t \rangle + \langle \mathbf{J} \mathbf{X}^{t+1} - \mathbf{Y}, \mathbf{W}^t \rangle \\ &\quad + \frac{\beta}{2} \|\mathbf{X}^{t+1} - \mathbf{Y}\|_F^2 + \frac{\beta}{2} \|\mathbf{J} \mathbf{X}^{t+1} - \mathbf{Y}\|_F^2 \\ \mathbf{Z}^{t+1} &= \mathbf{Z}^t + \beta \cdot (\mathbf{X}^{t+1} - \mathbf{Y}^{t+1} - \mathbf{J}) \\ \mathbf{W}^{t+1} &= \mathbf{W}^t + \beta \cdot (\mathbf{J} \mathbf{X}^{t+1} - \mathbf{Y}^{t+1}) \end{aligned}$$

1748 Considering first-order optimality conditions for functions  $\mathcal{L}_{\mathbf{X}}(\mathbf{X}, \mathbf{Y}^t; \mathbf{Z}^t, \mathbf{W}^t)$  and  
 1749  $\mathcal{L}_{\mathbf{Y}}(\mathbf{X}^{t+1}, \mathbf{Y}; \mathbf{Z}^t, \mathbf{W}^t)$ , we can get the updated formula for  $\mathbf{X}^{t+1}$  and  $\mathbf{Y}^{t+1}$ :  
 1750

$$\begin{aligned} \mathbf{X}^{t+1} &= -(l\mathbf{I} + \beta(\mathbf{Y}^t \mathbf{Y}^{t\top} + \mathbf{I}))^{-1}(\nabla f(\mathbf{X}^{t\top}) - l\mathbf{X}^t + \mathbf{Y}^t \mathbf{X}^{t\top} + \mathbf{J} \mathbf{W}^t - \beta \mathbf{Y}^t \mathbf{J} - \beta \mathbf{J} \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= -(\beta(\mathbf{X}^{t+1} \mathbf{X}^{t+1\top} + \mathbf{I}))^{-1}(\mathbf{X}^{t+1} \mathbf{Z}^t - \mathbf{J} \mathbf{W}^t - \beta \mathbf{X}^{t+1} \mathbf{J} - \beta \mathbf{J} \mathbf{X}^{t+1}) \end{aligned}$$

### F.4 IMPLEMENTATION OF UMCM ALGORITHM FOR PROBLEM (1)

1751 We consider the following smooth J-Orthogonality constraint problem:  
 1752  $\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} f(\mathbf{X})$ , s.t.  $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$ . We consider the Lagrangian function of the above  
 1753 J-Orthogonality constraint problem with  $\Lambda \in \mathbb{R}^{n \times n}$ :  
 1754

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

1755 Setting the gradient of  $\mathcal{L}(\mathbf{X}, \Lambda)$  w.r.t.  $\mathbf{X}$  to zero yields:  $\nabla f(\mathbf{X}) - \mathbf{J} \mathbf{X} \Lambda = 0$ .  
 1756

1757 Multiplying both sides by  $\mathbf{X}^\top$  and using the fact that  $\mathbf{X}^\top \mathbf{J} \mathbf{X} = \mathbf{J}$ , we have  $\mathbf{J} \Lambda = \mathbf{X}^\top \nabla f(\mathbf{X})$ .  
 1758 Multiplying both sides by  $\mathbf{J}^\top$  and using  $\mathbf{J}^\top \mathbf{J} = \mathbf{I}$ , we have  $\Lambda = \mathbf{J} \mathbf{X}^\top \nabla f(\mathbf{X})$ . Thus, we obtain  
 1759 equivalent unconstrained optimization problem:  
 1760

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} f(\mathbf{X}) - \frac{1}{2} \langle \mathbf{J} \mathbf{X}^\top \nabla f(\mathbf{X}), \mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J} \rangle.$$

---

1782  
 1783   **Algorithm 3:** Alternating Direction Method of Multipliers for Problem (1)  
 1784   1: Input:  $\mathbf{X}^0 \in \mathbb{R}^{n \times n}$ ,  $\mathbf{Y}^0 = \mathbf{X}^0$ ,  $\mathbf{Z}^0 \in \mathbb{R}^{n \times n}$ ,  $\mathbf{W}^0 \in \mathbb{R}^{n \times n}$  and positive constants  $l, \beta$ .  
 1785   **while** not converged **do**  
 1786     2: Update  $\mathbf{X}^{t+1}$  and  $\mathbf{Y}^{t+1}$ .  
 1787        $\mathbf{X}^{t+1} = -(l\mathbf{I} + \beta(\mathbf{Y}^t \mathbf{Y}^{t\top} + \mathbf{I}))^{-1}(\nabla f(\mathbf{X}^{t\top}) - l\mathbf{X}^t + \mathbf{Y}^t \mathbf{X}^{t\top} + \mathbf{J}\mathbf{W}^t - \beta\mathbf{Y}^t \mathbf{J} - \beta\mathbf{J}\mathbf{Y}^t)$   
 1788        $\mathbf{Y}^{t+1} = -(\beta(\mathbf{X}^{t+1} \mathbf{X}^{t+1\top} + \mathbf{I}))^{-1}(\mathbf{X}^{t+1} \mathbf{Z}^t - \mathbf{J}\mathbf{W}^t - \beta\mathbf{X}^{t+1} \mathbf{J} - \beta\mathbf{J}\mathbf{X}^{t+1})$   
 1789     3: Update  $\mathbf{Z}^{t+1}$  and  $\mathbf{W}^{t+1}$ .  
 1790        $\mathbf{Z}^{t+1} = \mathbf{Z}^t + \beta \cdot (\mathbf{X}^{t+1\top} \mathbf{Y}^{t+1} - \mathbf{J})$   
 1791        $\mathbf{W}^{t+1} = \mathbf{W}^t + \beta \cdot (\mathbf{J}\mathbf{X}^{t+1} - \mathbf{Y}^{t+1})$   
 1792     4: if  $t \% 5 == 0$  and  $\beta \leq 10^9$ :  $\beta = 2 * \beta$ .  
 1793     5:  $t = t + 1$   
 1794   **end**


---

1795   With the quadratic term  $\frac{\beta}{2} \|\mathbf{X}^\top \mathbf{J}\mathbf{X} - \mathbf{J}\|_F^2$ , we get the objective function of UMCM as follows:  
 1796

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} f(\mathbf{X}) - \frac{1}{2} \langle \mathbf{J}\mathbf{X}^\top \nabla f(\mathbf{X}), \mathbf{X}^\top \mathbf{J}\mathbf{X} - \mathbf{J} \rangle + \frac{\beta}{2} \|\mathbf{X}^\top \mathbf{J}\mathbf{X} - \mathbf{J}\|_F^2.$$

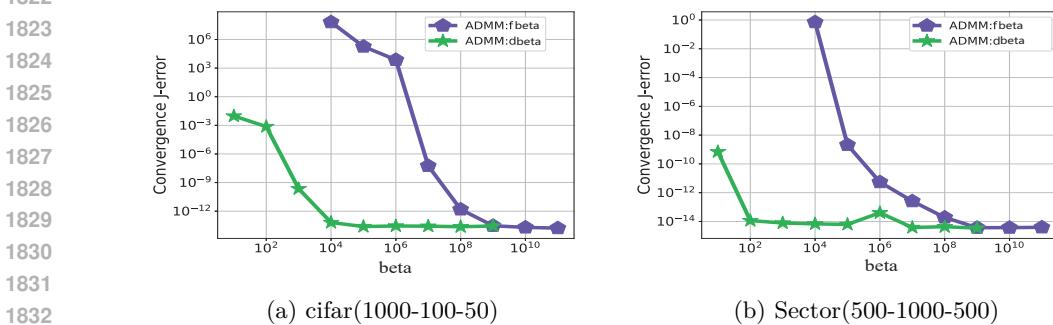
1804   Finally, we solve it by gradient-based approach. Exactly, we use the Adagrad optimizer  
 1805   built into PYTORCH in the our paper.  
 1806

## F.5 THE SELECTION OF PARAMETER $\beta$ IN ADMM AND UMCM

1809   In the ADMM algorithm,  $\beta$  is an important parameter that balances the constraint adherence  
 1810   and the optimization objective. We offer two methods for choosing  $\beta$ : one is a fixed  
 1811    $\beta$  that remains constant throughout the entire ADMM iteration process, and the other is  
 1812   a dynamic  $\beta$  that increases every specified number of iterations until it reaches an upper  
 1813   limit.

datasetname	(m-n-p)	fixed $\beta: 1e5$	fixed $\beta: 1e7$	fixed $\beta: 1e9$	dynamic $\beta: 1e3$	dynamic $\beta: 1e5$	dynamic $\beta: 1e7$
cifar	(1000-100-50)	-9.40e+09(1.5e+06)	-1.29e+07(9.2e-07)	-6.63e+03(2.9e-14)	-5.23e+06(3.6e-10)	-7.71e+03(2.3e-14)	-6.52e+03(2.9e-14)
mnist	(1000-750-300)	-9.26e+11(5.8e+05)	-6.31e+04(5.5e-11)	-3.95e+04(2.5e-14)	-1.36e+07(1.5e+02)	-9.88e+04(1.1e-14)	-3.97e+04(9.7e-15)
randn1000	(1000-1000-500)	-1.09e+06(2.1e-07)	-5.03e+05(1.2e-11)	-5.04e+05(1.1e-14)	-4.18e+06(1.7e-09)	-5.06e+05(7.5e-15)	-5.01e+05(8.7e-15)

1817   Table 3: Supplementary experiments for **HEVP** (limited to 30s). The data in the cell  
 1818   stand for the convergence values of **ADMM** under different  $\beta$  Settings and the value in  
 1819   parentheses represents  $\frac{1}{n^2} \sum_{ij}^n |\mathbf{X}^\top \mathbf{J}\mathbf{X} - \mathbf{J}|_{ij}$ . The values with an error of less than 1e-13  
 1820   are colored with red.  
 1821



1833   Figure 3: Supplementary experiments for **HEVP** (limited to 30s). The X-axis stands  
 1834   for the initial value of  $\beta$ , and the Y-axis stands for the convergence error of **ADMM**:  
 1835    $\frac{1}{n^2} \sum_{ij}^n |\mathbf{X}^\top \mathbf{J}\mathbf{X} - \mathbf{J}|_{ij}$ .

---

1836 Table 3 and Figure 3 show the objective function values and constraint violation for solving  
1837 the HEVP problem with different  $\beta$  values when ADMM converges. Generally, the dynamic  
1838  $\beta$  setting performs better than the fixed  $\beta$  setting. However, in the dynamic  $\beta$  initial setting,  
1839 a smaller  $\beta$  can lead to larger constraint violation.

1840 To effectively compare with feasible methods such as CSDM and JOBCD, we chose the  
1841 ADMM algorithm with a dynamic  $\beta$  setting, initializing  $\beta$  from range  $[1e2, 1e5]$ . The se-  
1842 lection of parameter  $\beta$  in the UMCM algorithm shares similar characteristics with that  
1843 in ADMM. We also use a dynamic setting, with the initial value of  $\beta$  chosen from range  
1844  $[1e2, 1e5]$ .

1845

## 1846 F.6 EXPERIMENT RESULT

1847

1848 ▶ **Hyperbolic Eigenvalue Problem.** Table 4 and Figure 4, 5, 6 are supplementary  
1849 experiments for HEVP. Several conclusions can be drawn. (i) **GS-JOBCD** often greatly  
1850 improves upon **UMCM**, **ADMM** and **CSDM**. This is because our methods find stronger  
1851 stationary points than them. (ii) **J-JOBCD** is a parallel version of **GS-JOBCD** and thus  
1852 exhibits significantly faster convergence. (iii) The proposed methods generally give the best  
1853 performance.

1854 ▶ **Hyperbolic Structural Probe Problem.** Table 5 and Figure 7, 8 are supplementary  
1855 experiments for HSPP. Several conclusions can be drawn. (i) **J-JOBCD** often greatly  
1856 improves upon **UMCM**, **ADMM** and **CSDM** (ii) **VR-J-JOBCD** is a reduced variance  
1857 version of **J-JOBCD** and thus exhibits significantly faster convergence for problems with  
1858 large samples. (iii) The proposed methods generally give the best performance.

1859 ▶ **Ultra-hyperbolic Knowledge Graph Embedding Problem.** Figure 9, 10, 11 and  
1860 12 are supplementary experiments for **UltraE**. Several conclusions can be drawn. (i) In  
1861 terms of Epoch performance, **J-JOBCD** and **VR-J-JOBCD** often greatly improves upon  
1862 **CSDM**, thus they show better MRR and hits results. (ii) In models with limited sample  
1863 sizes, the computational efficiency of **VR-J-JOBCD** is inferior to that of **J-JOBCD**. This  
1864 discrepancy arises because each iteration in **VR-J-JOBCD** necessitates two instances of  
1865 backpropagation, thus consuming substantial computational resources. (iii) The proposed  
1866 methods generally give the best performance.

1867

1868

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

## F.6.1 HYPERBOLIC EIGENVALUE PROBLEM

Table 4: The convergence curve of the compared methods for solving HEVP. (+) indicates that after the convergence of the **CSDM**, **UMCM** and **ADMM**, utilizing the **GS-JOBED** for optimization markedly enhances the objective value. The 1<sup>st</sup>, 2<sup>nd</sup>, and 3<sup>rd</sup> best results are colored with red, green and blue, respectively. ( $n, p$ ) represents the dimension and p-value of the  $\mathbf{J}$  orthogonal matrix (square matrix). The value in () stands for  $\frac{1}{n^2} \sum_{ij}^n |\mathbf{X}^\top \mathbf{J} \mathbf{X}| - \mathbf{J}_{ij}$  and cells with this value greater than 1e-7 are highlighted in gray.

dataname	(m-n-p)	UMCM	ADMM	CSDM	GS-JOBED	J-JOBED	UMCM+GS-JOBED	ADMM+GS-JOBED	CSDM+GS-JOBED
cifar	(1000-100-50)	-1.11e+04(5.3e-10)	-9.71e+04(1.2e-07)	1.38e+04(1.8e-09)	-6.33e+03(2.0e-09)	1.11e+04(1.8e-08)	2.31e+04(6.2e-09)(+)	-9.71e+04(4.1e-07)(+)	-2.81e+04(6.6e-09)(+)
CnCaltech	(2000-1000-500)	-2.02e+04(3.5e-08)	-2.02e+04(2.1e-07)	-2.02e+04(2.9e-09)	-2.51e+04(2.9e-07)	-2.02e+04(2.1e-08)	-2.02e+04(1.3e-07)(+)	-2.02e+04(2.1e-07)(+)	-2.02e+04(1.3e-07)(+)
gsette	(1000-100-500)	-1.09e+04(7.7e-04)	-2.84e+04(1.7e-04)	-5.66e+04(6.0e-10)	-2.76e+04(6.0e-09)	-1.09e+04(5.3e-07)(+)	-2.76e+04(6.0e-09)(+)	-1.09e+04(5.3e-07)(+)	-1.09e+04(5.3e-07)(+)
mnist	(1000-780-390)	-2.13e+03(5.1e-09)	-2.06e+03(1.7e-01)	-1.51e+04(1.0e-01)	-2.06e+03(5.7e-09)	-2.13e+03(5.1e-09)	-2.13e+03(5.2e-09)	-2.13e+03(5.2e-09)	-2.13e+03(5.2e-09)
randn10	(10-10-5)	-4.23e+04(5.8e-09)	-4.21e+04(1.1e-07)	-2.29e+04(2.2e-08)	-1.23e+01(1.5e-08)	-7.85e+01(3.7e-07)	-4.24e+04(1.2e-07)(+)	-4.22e+04(1.2e-07)(+)	-1.06e+11(7.0e-01)(+)
randn100	(100-100-50)	-5.01e+03(1.1e-09)	-4.93e+03(1.1e-07)	-7.29e+03(2.0e-09)	-4.89e+03(2.0e-09)	-1.27e+03(4.5e-08)	-5.01e+03(1.6e-09)	-4.93e+03(1.6e-09)	-5.01e+03(1.6e-09)
randn1000	(1000-1000-500)	-5.64e+04(5.7e-07)	-5.64e+04(5.0e-06)	-5.15e+04(5.1e-09)	-5.01e+04(5.1e-09)	-1.42e+04(1.6e-08)	-5.64e+04(5.0e-06)	-5.63e+04(5.0e-06)	-5.16e+04(5.1e-10)(+)
sector	(1000-1000-500)	-1.42e+03(2.6e-07)	-1.76e+03(3.0e-07)	-1.76e+03(3.0e-07)	-1.76e+03(3.0e-07)	-1.76e+03(3.0e-07)	-1.76e+03(3.0e-07)	-1.76e+03(3.0e-07)	-1.76e+03(3.0e-07)
TDT2	(1000-1000-500)	-1.06e+04(1.1e-07)	-1.97e+04(0.1e-10)	-8.10e+03(9.9e-11)	-1.80e+03(6.6e-10)	-1.06e+04(1.1e-07)	-1.06e+04(1.1e-07)	-1.06e+04(1.1e-07)	-1.98e+03(1.2e-10)
wla	(2470-290-145)	-3.72e+04(1.6e-09)	-3.72e+04(1.6e-08)	-3.60e+04(1.3e-09)	-1.38e+04(8.4e-10)	-3.97e+04(3.0e-08)	-1.48e+04(3.9e-09)	-1.73e+04(8.4e-08)(+)	-3.10e+04(1.3e-05)(+)
cifar	(1000-100-70)	-9.88e+03(1.6e-09)	-7.04e+04(1.1e-07)	1.03e+04(1.8e-09)	-6.30e+03(1.3e-09)	-6.30e+03(1.3e-09)	-7.20e+04(8.5e-09)	-7.04e+04(1.1e-07)(+)	-2.76e+04(8.5e-09)(+)
CnCaltech	(2000-1000-700)	-2.58e+04(2.5e-07)	-2.57e+04(2.1e-07)	-2.50e+04(2.7e-07)	-2.50e+04(2.7e-07)	-1.01e+04(3.5e-09)	-2.61e+04(2.1e-07)	-2.51e+04(2.1e-07)	-2.59e+04(2.1e-10)(+)
gsette	(3000-1000-700)	-6.86e+04(7.1e-04)	-1.11e+05(1.8e-01)	-1.49e+06(6.0e-11)	-1.40e+06(7.6e-09)	-2.00e+06(6.1e-09)	-3.45e+04(10.1e-09)	-1.52e+06(6.1e-10)	-1.52e+06(6.1e-10)
mnist	(1000-780-650)	-4.97e+04(2.3e-08)	-4.22e+04(1.8e-09)	-4.21e+04(1.1e-07)	-2.98e+04(2.1e-07)	-3.98e+04(1.8e-07)	-1.03e+04(3.7e-07)	-1.03e+04(3.7e-07)	-1.03e+04(3.7e-07)
randn10	(10-10-7)	-4.98e+03(1.1e-09)	-4.92e+03(3.9e-08)	-4.26e+03(3.1e-09)	-4.89e+03(2.5e-09)	-4.89e+03(2.5e-09)	-4.89e+03(2.5e-09)	-4.92e+03(3.9e-08)	-4.99e+03(9.4e-09)
randn100	(100-100-70)	-5.52e+03(1.7e-07)	-5.52e+03(3.0e-06)	-5.09e+03(5.9e-11)	-5.01e+05(6.6e-10)	-1.31e+06(6.1e-08)	-5.53e+03(1.7e-07)	-5.52e+03(1.7e-07)	-5.10e+05(1.2e-10)
randn1000	(1000-1000-700)	-1.64e+04(3.2e-08)	-1.78e+04(6.2e-04)	-1.62e+04(3.8e-11)	-1.80e+04(6.2e-04)	-1.53e+04(3.7e-11)	-1.75e+04(3.3e-08)	-1.77e+04(3.6e-04)	-1.63e+04(3.1e-11)
sector	(500-1000-500)	-1.40e+03(2.0e-08)	-1.63e+03(3.0e-07)	-1.56e+03(3.0e-08)	-1.80e+03(3.0e-08)	-1.53e+03(3.2e-10)	-1.67e+03(2.1e-08)	-1.63e+03(3.2e-07)	-1.57e+03(8.1e-11)
TDT2	(1000-1000-500)	-1.29e+07(1.1e-04)	-1.44e+11(1.9e-01)	-1.83e+06(5.6e-11)	-1.80e+06(3.0e-10)	-1.29e+07(1.1e-04)	-1.50e+11(1.9e-01)(+)	-1.85e+06(8.1e-11)(+)	-1.52e+04(1.2e-09)(+)
wla	(2470-290-250)	-1.40e+04(2.2e-09)	-1.72e+04(1.6e-09)	-1.40e+04(4.6e-09)	-1.40e+04(4.6e-09)	-1.01e+04(4.6e-09)	-1.44e+04(4.4e-09)	-1.53e+04(4.6e-08)	-1.01e+04(4.4e-09)(+)
cifar	(1000-100-20)	-9.11e+04(6.2e-10)	-9.71e+04(4.1e-09)	-9.71e+04(2.8e-09)	-6.50e+03(3.6e-09)	-1.02e+04(3.6e-08)	-2.96e+04(7.6e-09)(+)	-9.61e+04(1.3e-07)(+)	-4.36e+04(8.4e-09)(+)
CnCaltech	(2000-1000-500)	-2.52e+04(2.8e-06)	-2.50e+04(2.5e-07)	-2.50e+04(2.6e-07)	-2.50e+04(2.6e-07)	-1.01e+04(3.5e-09)	-2.75e+04(2.1e-07)	-2.53e+04(2.8e-08)(+)	-2.76e+04(2.1e-10)
gsette	(3000-1000-500)	-1.30e+07(2.1e-04)	-2.82e+04(10.3e-06)	-1.43e+06(5.9e-11)	-1.40e+06(7.6e-09)	-1.67e+06(6.1e-09)	-1.36e+07(2.1e-04)(+)	-1.52e+06(6.1e-10)	-1.52e+06(6.1e-10)
mnist	(1000-780-650)	-4.97e+04(2.3e-08)	-5.82e+04(7.9e-05)	-3.95e+04(8.4e-05)	-3.98e+04(8.4e-05)	-1.16e+04(5.6e-09)	-5.52e+04(2.3e-08)(+)	-1.25e+04(5.6e-09)	-4.39e+04(1.5e-10)
randn10	(10-10-5)	-4.21e+04(9.0e-09)	-4.21e+04(1.1e-07)	-2.29e+04(2.7e-06)	-1.23e+01(1.5e-08)	-8.32e+04(1.7e-07)	-4.24e+04(1.2e-08)	-4.22e+04(1.2e-07)	-2.80e+03(2.7e-01)
randn100	(100-100-50)	-5.01e+03(6.1e-09)	-4.93e+03(1.0e-09)	-4.17e+03(2.6e-09)	-1.40e+04(1.3e-10)	-1.40e+04(1.3e-10)	-1.96e+04(2.5e-09)	-1.40e+04(2.4e-09)	-1.12e+04(2.4e-08)
randn1000	(1000-1000-500)	-5.64e+04(5.0e-07)	-5.64e+04(5.0e-08)	-5.24e+04(1.8e-09)	-5.01e+04(5.9e-11)	-2.96e+04(3.8e-09)	-5.64e+04(5.1e-08)	-5.27e+04(5.1e-08)	-5.27e+04(5.1e-10)
sector	(100-100-21-18)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)	-1.64e+04(3.6e-09)
TDT2	(1000-1000-500)	-1.05e+08(1.8e-06)	-4.03e+06(15.4e-09)	-2.06e+06(12.1e-09)	-1.72e+06(12.1e-09)	-9.40e+04(1.5e-07)	-1.49e+04(4.6e-09)	-1.80e+04(4.8e-08)	-1.81e+04(2.4e-05)
wla	(2470-290-145)	-1.72e+06(2.0e-09)	-1.72e+06(4.8e-08)	-1.66e+06(2.2e-09)	-1.38e+04(1.5e-09)	-9.40e+04(1.5e-07)	-1.80e+04(4.6e-09)	-1.80e+04(4.8e-08)	-1.81e+04(2.4e-05)
cifar	(1000-100-70)	-9.88e+03(1.6e-09)	-7.04e+04(1.1e-07)	-1.04e+04(1.8e-09)	-1.04e+04(1.8e-09)	-1.04e+04(1.8e-09)	-2.71e+04(8.8e-09)	-7.04e+04(1.1e-07)(+)	-3.89e+04(9.4e-09)(+)
CnCaltech	(2000-1000-700)	-2.58e+04(6.6e-10)	-2.50e+04(5.0e-08)	-2.50e+04(2.1e-09)	-1.47e+04(3.0e-08)	-2.50e+04(2.1e-09)	-2.50e+04(5.1e-08)	-2.68e+04(2.1e-10)	-2.76e+04(2.1e-10)
gsette	(3000-1000-700)	-6.86e+04(7.1e-04)	-1.09e+05(1.8e-01)	-1.46e+06(6.1e-11)	-1.41e+06(7.6e-11)	-3.02e+06(3.6e-09)	-1.10e+08(1.1e-04)	-3.10e+06(1.1e-04)	-1.70e+06(1.8e-10)
mnist	(1000-780-390)	-2.13e+03(5.1e-09)	-2.06e+03(1.5e-09)	-1.91e+04(4.7e-10)	-1.11e+04(4.1e-09)	-3.28e+04(1.5e-08)	-2.14e+04(2.0e-09)	-5.82e+04(7.4e-01)(+)	-4.82e+04(3.5e-10)
randn10	(10-10-5)	-4.23e+04(9.0e-09)	-4.21e+04(1.1e-07)	-2.29e+04(2.7e-06)	-1.23e+01(1.5e-08)	-8.32e+04(1.7e-07)	-4.24e+04(1.2e-08)	-4.22e+04(1.2e-07)	-2.80e+03(2.7e-01)
randn100	(100-100-50)	-5.01e+03(1.1e-09)	-4.93e+03(1.0e-09)	-4.17e+03(2.6e-09)	-1.40e+04(1.3e-10)	-1.40e+04(1.3e-10)	-1.96e+04(2.5e-09)	-1.40e+04(2.4e-09)	-1.13e+04(2.4e-08)
randn1000	(1000-1000-500)	-5.64e+04(5.0e-07)	-5.64e+04(5.0e-08)	-5.24e+04(1.8e-09)	-5.01e+04(5.9e-11)	-2.96e+04(3.8e-09)	-5.64e+04(5.1e-08)	-5.27e+04(5.1e-08)	-5.27e+04(5.1e-10)
sector	(500-1000-500)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)	-1.64e+03(3.6e-09)
TDT2	(1000-1000-500)	-1.26e+07(3.1e-08)	-1.26e+07(4.1e-08)	-1.26e+07(4.1e-08)	-1.26e+07(4.1e-08)	-1.26e+07(4.1e-08)	-1.26e+07(4.1e-08)	-1.26e+07(4.1e-08)	-1.26e+07(4.1e-08)
wla	(2470-290-250)	-1.40e+04(2.1e-09)	-1.53e+04(6.3e-08)	-1.40e+04(4.6e-08)	-1.40e+04(4.6e-08)	-1.40e+04(4.6e-08)	-1.40e+04(4.6e-08)	-1.40e+04(4.6e-08)	-1.40e+04(4.6e-08)
cifar	(1000-100-50)	-7.23e+04(1.0e-07)	-1.07e+04(4.1e-09)	-3.43e+04(6.7e-08)	-3.43e+04(6.7e-08)	-5.76e+04(1.1e-08)	-7.62e+04(1.0e-07)	-4.97e+04(1.6e-08)	-4.97e+04(1.6e-08)
CnCaltech	(2000-1000-500)	-2.62e+04(2.1e-10)	-2.51e+04(2.5e-08)	-2.80e+04(2.1e-10)	-1.96e+03(1.6e-08)	-2.51e+04(2.6e-08)	-2.80e+04(2.1e-10)	-2.87e+04(2.1e-10)	-2.87e+04(2.1e-10)
gsette	(3000-1000-500)	-1.02e+08(4.5e-05)	-1.49e+11(2.1e-05)	-1.77e+11(2.1e-05)	-1.41e+06(1.6e-09)	-3.58e+06(6.0e-09)	-3.50e+05(12.5e-09)	-1.81e+06(2.4e-10)	-1.81e+06(2.4e-10)
mnist	(1000-780-390)	-2.13e+03(1.4e-09)	-3.32e+04(1.0e-04)	-2.22e+04(1.6e-04)	-4.41e+04(2.9e-11)	-1.63e+04(2.5e-08)	-2.15e+05(2.0e-09)	-3.16e+05(6.4e-02)(+)	-8.85e+04(5.6e-10)
randn10	(10-10-5)	-4.23e+04(8.4e-09)	-4.21e+04(1.4e-07)	-2.29e+04(2.7e-06)	-1.23e+01(1.5e-08)	-8.46e+04(1.7e-07)	-4.24e+04(1.2e-08)	-4.22e+04(1.2e-07)	-2.80e+03(2.7e-01)
randn100	(100-100-50)	-5.01e+03(1.4e-09)	-4.93e+03(9.7e-08)	-4.09e+03(8.7e-08)	-4.89e+03(8.7e-08)	-1.62e+04(1.3e-07)	-5.25e+04(2.3e-08)	-5.25e+04(2.3e-08)	-5.25e+04(2.3e-08)
randn1000	(1000-1000-500)	-5.64e+04(5.0e-07)	-5.64e+04(5.0e-08)	-5.24e+04(1.8e-09)	-5.01e+04(5.9e-11)	-2.96e+04(3.8e-09)	-5.64e+04(5.1e-08)	-5.27e+04(5.1e-08)	-5.27e+04(5.1e-10)
sector	(500-1000-500)	-1.42e+03(7.5e-10)	-1.51e+03(9.7e-08)	-8.40e+03(5.3e-08)	-5.23e+03(15.7e-08)	-5.01e+03(6.8e-09)	-1.63e+03(3.4e-08)	-1.63e+03(3.4e-08)	-1.63e+03(3.4e-10)
TDT2	(1000-1000-500)	-6.98e+04(1.0e-07)	-7.26e+04(3.8e-08)	-7.26e+04(3.8e-08)	-3.07e+04(1.0e-09)	-3.07e+04(1.0e-09)	-3.07e+04(1.0e-09)	-3.07e+04(1.0e-09)	-3.07e+04(1.0e-09)
wla	(2470-290-200)	-1.40e+06(1.6e-09)	-1.39e+04(6.2e-08)	-1.40e+04(4.6e-08)	-1.40e+04(4.6e-08)	-1.38e+04(4.6e-08)	-1.38e+04(4.6e-08)	-1.38e+04(4.6e-08)	-1.38e+04(4.6e-08)
cifar	(1000-100-70)	-7.78e+03(1.6e-09)	-7.17e+04(9.1e-08)	-3.08e+04(1.0e-08)	-3.08e+04(1.0e-08)	-3.08e+04(1.0e-08)	-5.02e+04(9.3e-08)(+)	-3.70e+04(1.1e-08)(+)	-2.87e+04(1.1e-08)(+)
CnCaltech	(2000-1000-700)	-2.58e+04(2.6e-10)	-2.50e+04(2.5e-08)	-2.50e+04(2.6e-10)	-2.50e+04(2.6e-10)	-1.91e+04(2.1e-08)	-2.58e+04(2.7e-08)	-2.58e+04(2.1e-10)	-2.58e+04(2.1e-10)
gsette	(3000-1000-700)	-1.25e+08(4.8e-05)	-1.38e+04(2.1e-05)	-1.77e+04(2.1e-05)	-1.41e+06(1.6e-09)	-3.58e+06(6.0e-09)	-3.50e+05(12.5e-09)	-1.81e+06(2.4e-10)	-1.81e+06(2.4e-10)
mnist	(1000-780-500)	-1.73e+03(5.8e-10)	-1.80e+03(5.3e-09)	-1.40e+03(5.3e-09)	-1.40e+03(5.3e-				

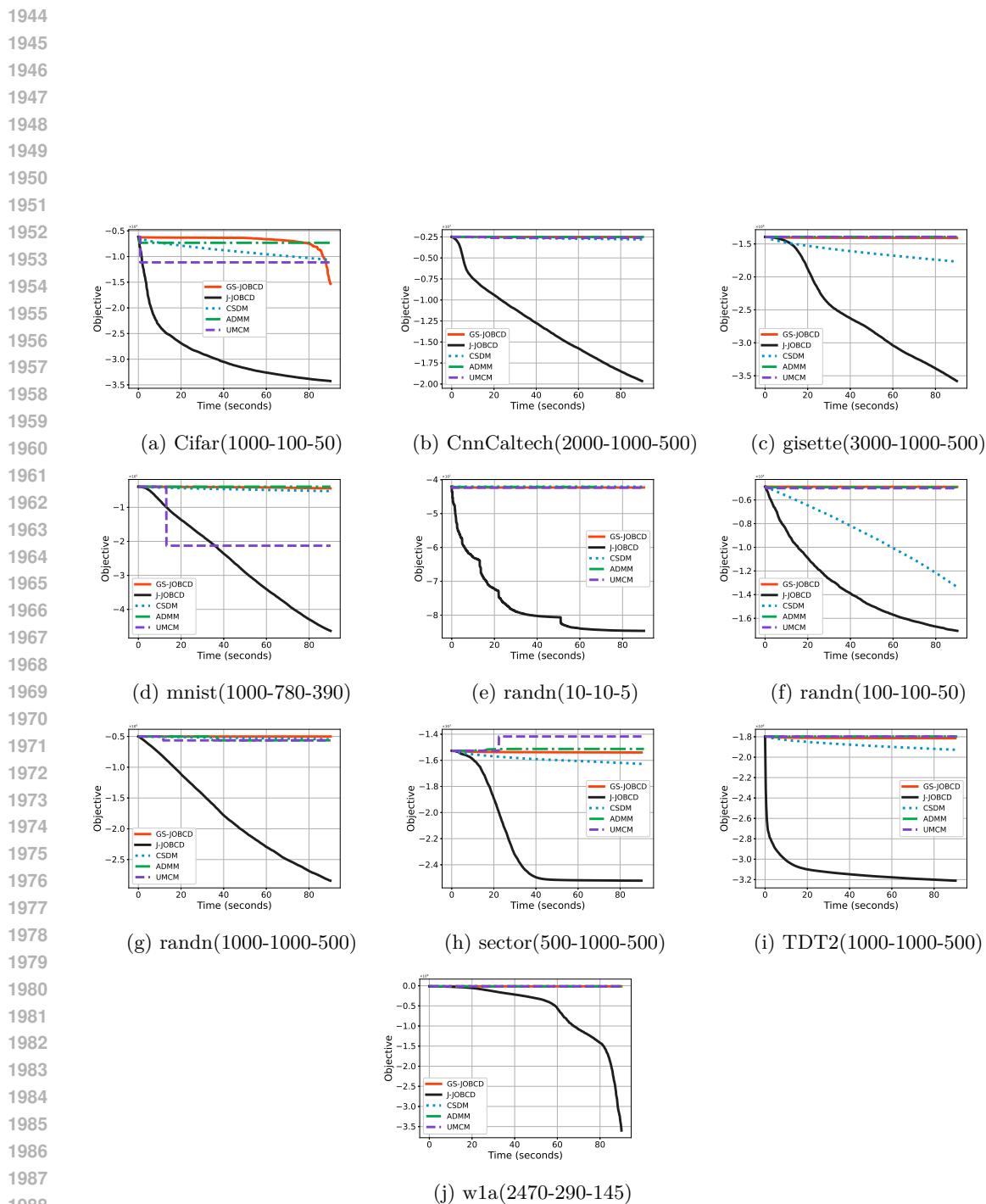


Figure 4: The convergence curve of the compared methods for solving HEVP with varying  $(m, n, p)$ .

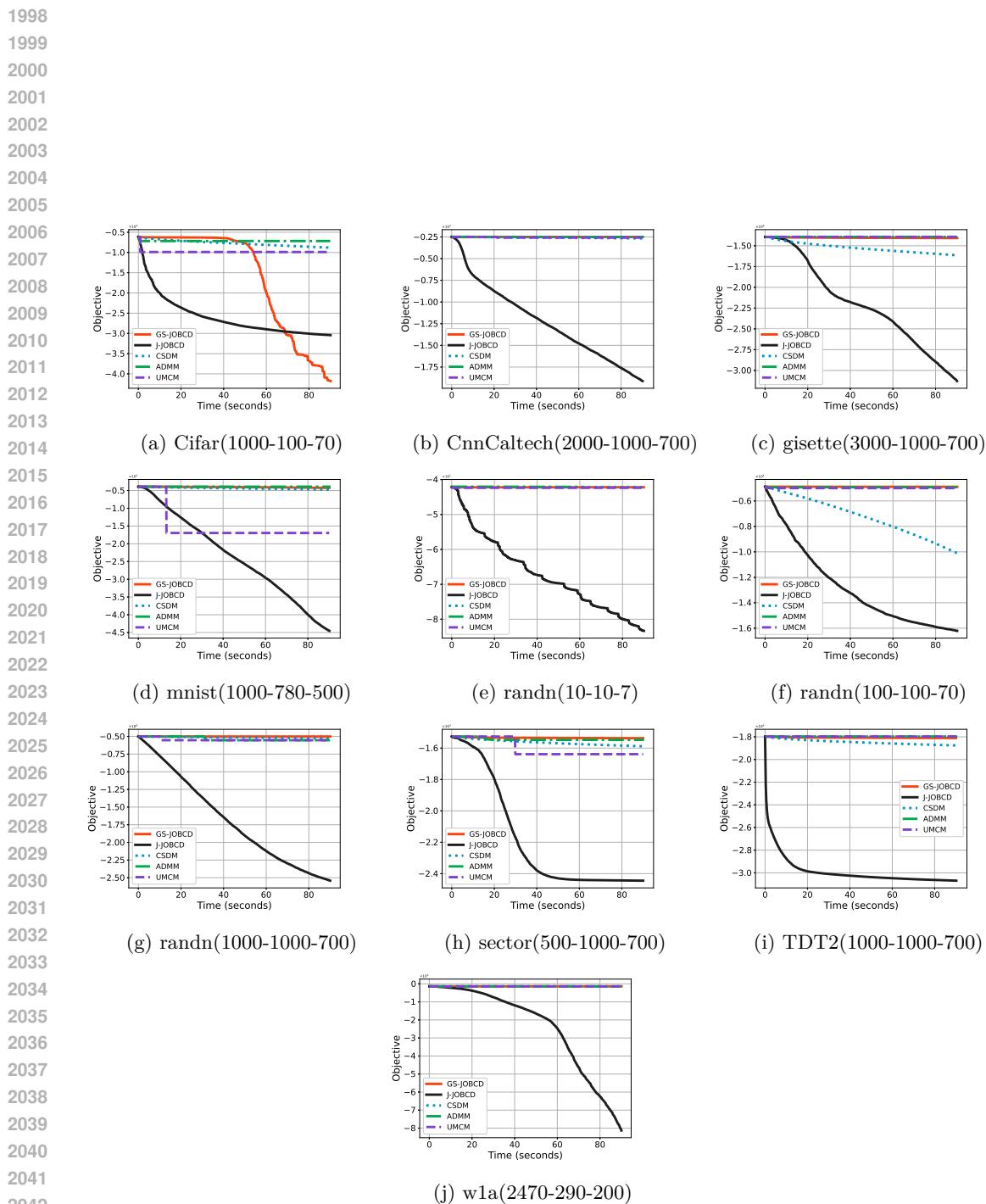


Figure 5: The convergence curve of the compared methods for solving HEVP with varying  $(m, n, p)$ .

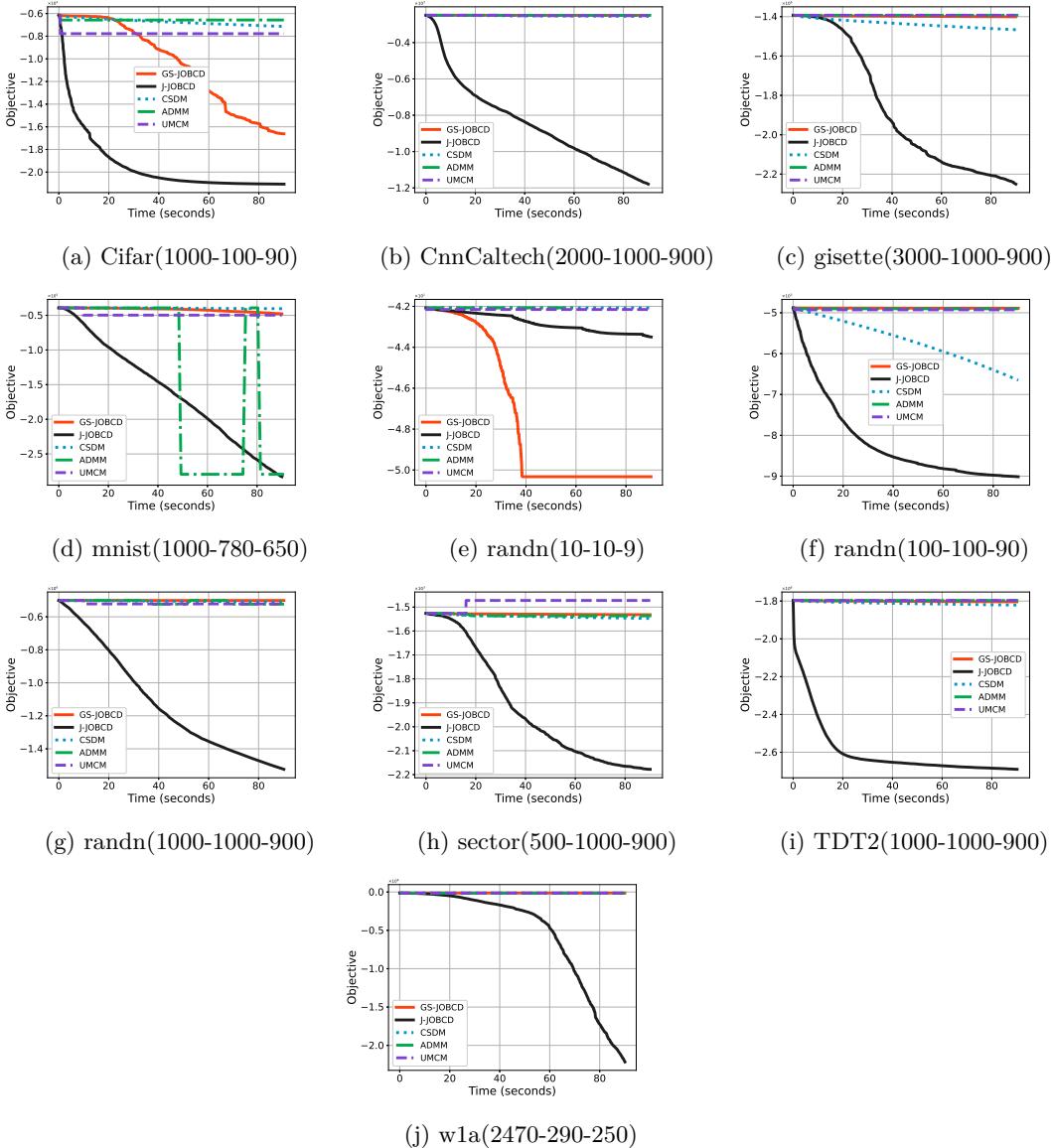


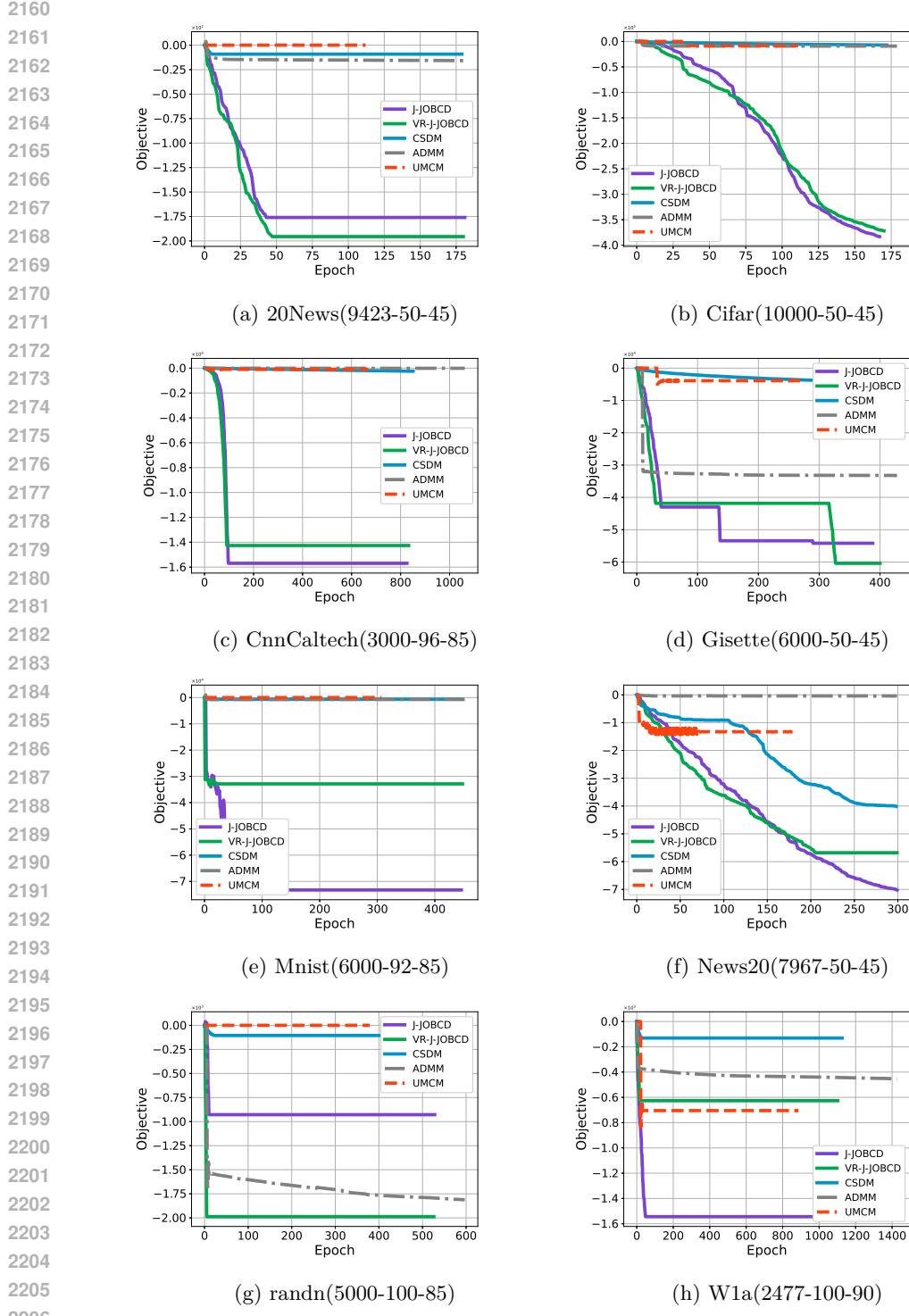
Figure 6: The convergence curve of the compared methods for solving HEVP with varying  $(m, n, p)$ .

---

## F.6.2 HYPERBOLIC STRUCTURAL PROBE PROBLEM

Table 5: The convergence curve of the compared methods for solving HSPP. (+) indicates that after the convergence of the **CSDM**, utilizing the **J-OBCD** for optimization markedly enhances the objective value. The 1<sup>st</sup>, 2<sup>nd</sup>, and 3<sup>rd</sup> 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  $\frac{1}{n^2} \sum_{ij}^n |\mathbf{X}^\top \mathbf{J} \mathbf{X} - \mathbf{J}|_{ij}$  and cells with this value greater than 1e-7 are highlighted in gray.

	datasetname	(m-n-p)	ADMM	UMCM	CSDM	J-OBCD	VR-J-OBCD	CSDM+J-OBCD
time limit=30s								
20News	(9423-50-25)	<b>-7.56e+00(7.6e-06)</b>	-5.59e+01(3.0e-04)	-3.21e+01(1.0e-04)	-2.52e+02(2.4e-08)	<b>-2.96e+02(3.5e-08)</b>	-3.21e+01(1.0e-04)	
Cifar	(10000-50-25)	-4.81e+04(1.3e-05)	-1.00e+04(4.1e-04)	-5.15e+03(2.4e-04)	-5.26e+04(1.6e-08)	-3.19e+04(1.7e-08)	-5.15e+03(2.4e-04)(+)	
cnnCaltech	(3000-96-48)	-9.35e+00(4.0e-06)	-1.90e+02(5.6e-06)	<b>-2.11e+02(7.1e-06)</b>	-1.65e+04(2.8e-08)	-1.34e+04(2.3e-08)	-2.11e+02(7.1e-06)(+)	
gissette	(6000-50-25)	-4.50e+04(1.1e-05)	-4.50e+03(3.1e-05)	-1.67e+04(9.5e-06)	-2.27e+03(6.3e-06)	-3.08e+03(1.1e-08)	-5.11e+04(1.7e-08)	-5.00e+03(3.3e-04)(+)
Mnist	(6000-92-46)	<b>-4.90e+04(6.8e-06)</b>	-1.76e+04(1.1e-05)	-1.67e+04(9.5e-06)	-5.00e+03(3.3e-04)	-5.81e+04(1.6e-08)	-2.27e+03(6.3e-06)	
news20	(7967-50-25)	-3.91e-03(7.6e-06)	-1.25e+00(4.2e-04)	-1.06e+00(1.0e-04)	<b>-7.03e+00(6.2e-08)</b>	-6.92e+00(5.6e-08)	-1.06e+00(1.0e-04)(+)	
randn5000	(5000-100-50)	-8.82e+02(1.9e-06)	-6.37e+02(4.2e-04)	-1.56e+02(2.6e-06)	-3.53e+03(9.1e-09)	<b>-4.71e+03(1.9e-08)</b>	-1.56e+02(2.6e-06)	
w1a	(2477-100-75)	-3.38e+02(3.6e-06)	-9.70e+02(3.1e-04)	-1.65e+02(5.9e-06)	-2.97e+03(9.0e-09)	<b>-3.44e+03(2.1e-08)</b>	-1.65e+02(5.9e-06)	
20News	(9423-50-35)	-1.01e+01(7.6e-05)	-6.95e+01(1.5e-03)	-1.32e+01(7.0e-06)	-1.50e+02(8.1e-08)	<b>-1.95e+02(1.2e-08)</b>	-1.32e+01(7.0e-06)	
Cifar	(10000-50-35)	<b>-6.60e+04(1.2e-05)</b>	-1.11e+04(4.1e-04)	-4.50e+03(1.9e-04)	-2.99e+04(1.5e-08)	-3.78e+04(1.7e-08)	-4.50e+03(1.9e-04)(+)	
cnnCaltech	(3000-96-70)	-1.19e+01(4.0e-06)	-1.79e+02(7.2e-06)	-1.33e+02(4.9e-06)	-1.43e+04(2.3e-08)	<b>-1.79e+04(2.7e-08)</b>	-1.33e+02(4.9e-06)(+)	
gissette	(6000-50-35)	-2.46e+04(1.0e-05)	-4.87e+03(3.1e-05)	-3.87e+03(8.7e-05)	-5.87e+04(3.0e-08)	<b>-6.76e+04(3.0e-08)</b>	-3.87e+03(8.7e-05)(+)	
Mnist	(6000-92-70)	<b>-7.45e+04(6.3e-06)</b>	-1.71e+04(3.2e-04)	-2.15e+02(3.0e-06)	-9.38e+03(1.1e-08)	-5.00e-01(3.3e-09)	-2.15e+03(3.0e-06)	
news20	(7967-50-35)	-1.56e-02(7.6e-06)	-1.34e+00(3.0e-04)	-9.14e-01(1.2e-04)	-9.70e+00(6.1e-08)	<b>-1.06e+01(6.2e-08)</b>	-9.14e-01(1.2e-04)(+)	
randn5000	(5000-100-75)	-1.57e+03(5.4e-06)	-6.27e+02(8.9e-04)	-9.15e+01(5.4e-06)	-1.63e+03(1.7e-08)	-4.35e+01(7.1e-10)	-9.15e+01(5.4e-06)	
w1a	(2477-100-75)	-6.77e+02(3.6e-06)	-1.06e+03(6.1e-06)	-1.80e+02(3.5e-05)	-1.80e+03(5.1e-09)	<b>-1.80e+03(5.1e-09)</b>	-1.90e+02(4.5e-05)	
20News	(9423-50-45)	<b>-1.51e+01(7.6e-06)</b>	-8.44e+00(1.1e-03)	-9.15e+00(1.4e-05)	-3.84e+01(1.1e-09)	<b>-1.96e+02(1.1e-08)</b>	-9.15e+00(1.4e-05)	
Cifar	(10000-50-45)	<b>-8.45e+04(1.4e-05)</b>	-8.15e+03(3.2e-04)	-3.75e+03(1.9e-04)	-2.69e+03(3.2e-08)	<b>-5.00e+03(3.9e-09)</b>	-3.75e+03(1.9e-04)(+)	
cnnCaltech	(3000-96-85)	-7.38e+00(4.0e-06)	-9.28e+00(1.8e-06)	-8.77e+01(5.4e-05)	-1.79e+04(3.4e-08)	-1.43e+04(3.1e-08)	-8.77e+01(5.4e-05)(+)	
gissette	(6000-50-45)	-3.29e+04(4.9e-06)	-3.88e+03(3.1e-05)	-2.65e+03(3.3e-06)	-5.02e+04(1.5e-08)	-4.18e+04(1.6e-08)	-2.65e+03(3.3e-04)(+)	
Mnist	(6000-92-85)	<b>-8.69e+04(7.2e-06)</b>	-2.49e+03(1.0e-03)	-6.49e+02(3.0e-06)	-3.13e+04(1.1e-08)	<b>-3.10e+04(1.3e-08)</b>	-6.49e+02(3.0e-06)	
news20	(7967-50-45)	-3.91e-02(7.6e-06)	-1.19e+00(3.0e-04)	-8.98e-01(6.0e-05)	-4.00e+00(1.5e-07)	-3.60e+00(4.7e-07)	-8.98e-01(6.0e-05)(+)	
randn5000	(5000-100-85)	-1.66e+03(2.7e-06)	-3.77e+02(9.2e-04)	-2.90e+02(9.8e-06)	-2.68e+03(2.3e-09)	-1.99e+03(2.0e-10)	-1.05e+02(9.8e-06)	
w1a	(2477-100-90)	-4.26e+02(3.7e-06)	-7.05e+02(3.1e-04)	-1.31e+02(3.0e-06)	-6.27e+02(2.0e-09)	-6.26e+02(2.4e-09)	-1.31e+02(3.0e-06)	
time limit=60s								
20News	(9423-50-25)	-7.77e+00(7.6e-06)	<b>-5.79e+01(4.3e-05)</b>	-3.21e+01(1.0e-04)	-3.57e+02(4.6e-08)	<b>-2.96e+02(3.5e-08)</b>	-3.21e+01(1.0e-04)	
Cifar	(10000-50-25)	<b>-4.88e+04(1.1e-05)</b>	-1.01e+04(2.2e-05)	-7.96e+03(3.6e-04)	-8.47e+04(4.1e-08)	-5.92e+04(3.6e-08)	-7.96e+03(3.6e-04)(+)	
cnnCaltech	(3000-96-48)	-1.19e+01(4.0e-06)	-1.89e+02(5.6e-06)	<b>-3.58e+02(7.6e-06)</b>	-8.79e+03(2.1e-08)	<b>-1.34e+04(2.3e-08)</b>	-3.58e+02(7.6e-06)(+)	
gissette	(6000-50-25)	-1.81e+04(8.9e-06)	-4.49e+03(1.4e-05)	-7.58e+03(4.9e-04)	-6.78e+04(2.0e-08)	-5.11e+04(1.7e-08)	-7.58e+03(4.9e-04)(+)	
Mnist	(6000-92-46)	<b>-5.04e+04(6.8e-06)</b>	-1.67e+04(8.7e-06)	-2.27e+03(6.3e-06)	-4.25e+04(5.2e-08)	-2.28e+04(3.8e-08)	-2.27e+03(6.3e-06)	
news20	(7967-50-25)	-3.91e-03(7.6e-06)	-1.21e+00(1.0e-05)	-1.21e+00(1.0e-04)	-1.06e+01(1.1e-07)	<b>-1.17e+01(1.1e-07)</b>	-1.21e+00(1.0e-04)(+)	
randn5000	(5000-100-50)	-9.31e+02(9.7e-06)	-6.37e+02(2.4e-04)	-1.56e+02(2.6e-06)	-8.38e+00(6.9e-11)	<b>-4.71e+03(1.9e-08)</b>	-1.56e+02(2.6e-06)	
w1a	(2477-100-50)	-3.71e+02(3.5e-06)	-9.70e+02(3.1e-04)	-1.65e+02(5.9e-06)	-2.19e+03(4.8e-09)	<b>-3.44e+03(2.1e-08)</b>	-1.65e+02(5.9e-06)	
20News	(9423-50-35)	-1.03e+01(7.6e-06)	<b>-6.95e+01(1.5e-05)</b>	-1.32e+01(7.0e-06)	-4.52e+01(2.5e-08)	<b>-1.95e+02(1.2e-08)</b>	-1.32e+01(7.0e-06)	
Cifar	(10000-50-35)	-6.68e+04(1.0e-05)	-1.10e+04(1.1e-05)	-6.80e+03(3.2e-04)	-9.09e+04(4.0e-08)	<b>-6.64e+04(3.4e-08)</b>	-6.80e+03(3.2e-04)(+)	
cnnCaltech	(3000-96-70)	-1.46e+01(4.0e-06)	-1.79e+02(2.0e-06)	<b>-2.29e+02(4.6e-06)</b>	-1.87e+04(2.5e-08)	<b>-1.79e+04(2.7e-08)</b>	-2.29e+02(4.6e-06)(+)	
gissette	(6000-50-35)	<b>-2.52e+04(1.0e-05)</b>	-4.89e+03(1.3e-05)	-5.53e+03(9.5e-05)	-7.67e+04(3.6e-08)	-6.76e+04(3.6e-08)	-5.53e+03(9.5e-05)(+)	
Mnist	(6000-92-70)	<b>-7.71e+04(6.1e-06)</b>	-1.71e+04(3.2e-04)	-2.15e+03(3.0e-06)	-9.38e+03(1.1e-08)	-5.00e-01(3.3e-09)	-2.15e+03(3.0e-06)	
news20	(7967-50-35)	-1.56e-02(7.7e-06)	-1.42e+00(1.3e-05)	-1.26e+00(1.0e-04)	-1.46e-01(1.0e-07)	<b>-1.54e+01(1.1e-07)</b>	-1.26e+00(1.0e-04)(+)	
randn5000	(5000-100-75)	-1.65e+03(6.0e-06)	-6.27e+02(8.9e-04)	-9.15e+01(5.4e-06)	-2.01e+03(1.7e-08)	-4.35e+01(7.1e-10)	-9.15e+01(5.4e-06)	
w1a	(2477-100-75)	-7.26e+02(3.2e-06)	-1.06e+03(6.1e-06)	-1.80e+02(4.5e-05)	-1.58e+03(3.0e-09)	<b>-1.80e+03(5.1e-09)</b>	-1.90e+02(4.5e-05)	
20News	(9423-50-45)	-1.54e+01(7.6e-06)	<b>-8.44e+00(1.1e-03)</b>	-9.15e+00(1.4e-05)	-8.68e+00(5.4e-10)	<b>-1.96e+02(1.1e-08)</b>	-9.15e+00(1.4e-05)	
Cifar	(10000-50-45)	<b>-8.59e+04(1.3e-05)</b>	-8.08e+03(3.1e-05)	-5.64e+03(2.9e-04)	-9.56e+03(5.0e-07)	<b>-2.65e+03(2.9e-04)</b>	-5.64e+03(2.9e-04)(+)	
cnnCaltech	(3000-96-85)	-9.45e+00(4.0e-06)	-9.26e+00(1.8e-06)	-1.67e+02(5.5e-05)	<b>-1.65e+04(2.8e-08)</b>	-1.43e+04(3.1e-08)	-1.67e+02(5.5e-05)(+)	
gissette	(6000-50-45)	-3.32e+04(4.9e-06)	-3.87e+03(3.1e-05)	-3.81e+03(3.0e-06)	-4.50e+04(7.2e-07)	-4.18e+04(4.1e-07)	-3.81e+03(3.0e-04)(+)	
Mnist	(6000-92-85)	<b>-9.06e+04(6.8e-06)</b>	-2.49e+03(1.1e-03)	-6.49e+02(3.0e-06)	-1.63e+04(1.2e-08)	-3.10e+04(1.3e-08)	-6.49e+02(3.0e-06)	
news20	(7967-50-45)	-3.91e-02(7.6e-06)	-1.33e+00(1.1e-05)	-3.20e+00(8.0e-05)	-6.17e+00(5.0e-07)	<b>-5.50e+00(6.0e-07)</b>	-3.20e+00(8.0e-05)(+)	
randn5000	(5000-100-85)	-1.76e+03(3.9e-06)	-3.77e+02(9.2e-04)	-1.05e+02(9.8e-06)	-3.87e+03(1.3e-09)	-1.99e+03(2.0e-10)	-1.05e+02(9.8e-06)	
w1a	(2477-100-90)	-4.38e+02(3.6e-06)	-7.05e+02(3.1e-04)	-1.31e+02(3.0e-06)	-3.07e+01(9.3e-11)	<b>-6.26e+02(2.4e-09)</b>	-1.31e+02(3.0e-06)	
time limit=90s								
20News	(9423-50-25)	-7.96e+00(7.6e-06)	<b>-5.79e+01(4.3e-05)</b>	-3.21e+01(1.0e-04)	-4.70e+01(2.9e-09)	<b>-2.96e+02(3.5e-08)</b>	-3.21e+01(1.0e-04)	
Cifar	(10000-50-25)	<b>-4.93e+04(2.6e-05)</b>	-1.01e+04(9.8e-06)	-1.05e+04(4.9e-06)	-1.07e+05(6.5e-07)	-8.44e+04(5.9e-08)	-1.05e+04(4.9e-04)(+)	
cnnCaltech	(3000-96-48)	-1.55e+01(4.0e-06)	-1.89e+02(5.3e-06)	<b>-4.99e+02(8.3e-06)</b>	-1.58e+04(2.5e-08)	-1.34e+04(2.3e-08)	-4.99e+02(8.3e-06)(+)	
gissette	(6000-50-25)	-1.87e+04(9.3e-06)	-4.51e+03(1.4e-05)	-9.52e+03(6.1e-04)	-6.07e+04(1.4e-08)	-5.11e+04(1.7e-08)	-9.52e+03(6.1e-04)(+)	
Mnist	(6000-92-46)	<b>-5.24e+04(6.4e-06)</b>	-1.67e+04(8.7e-06)	-2.27e+03(6.3e-06)	-1.35e+04(2.3e-08)	-2.28e+04(3.8e-08)	-2.27e+03(6.3e-06)	
news20	(7967-50-25)	-3.91e-03(7.6e-06)	-1.21e+00(1.0e-05)	-1.34e+02(2.0e-05)	-1.37e+01(1.8e-07)	<b>-1.43e+01(1.6e-07)</b>	-1.34e+00(2.4e-05)(+)	
randn5000	(5000-100-50)	-9.98e+02(9.7e-06)	-6.37e+02(2.4e-04)	-1.56e+00(2.6e-06)	-1.09e+03(3.6e-09)	<b>-4.71e+03(1.9e-08)</b>	-1.56e+02(2.6e-06)	
w1a	(2477-100-75)	-4.16e+02(3.4e-06)	-9.70e+02(3.1e-04)	-1.65e+02(5.9e-06)	-2.80e+03(7.4e-09)	<b>-3.44e+03(2.1e-08)</b>	-1.65e+02(5.9e-06)	
20News	(9423-50-35)	-1.04e+01(7.6e-06)	-6.95e+01(1.5e-05)	<b>-1.32e+01(7.0e-06)</b>	-1.94e+01(2.1e-08)	<b>-1.95e+02(1.2e-08)</b>	-1.32e+01(7.0e-06)	
Cifar	(10000-50-35)	<b>-6.83e+04(1.2e-05)</b>	-1.10e+04(1.2e-05)	-8.44e+03(4.1e-04)	-8.76e+04(5.4e-08)	<b>-9.79e+04(5.4e-08)</b>	-8.83e+03(4.1e-04)(+)	
cnnCaltech	(3000-96-70)	-1.79e+01(4.0e-06)	-1.79e+02(7.3e-06)	<b>-3.02e+02(5.8e-06)</b>	-1.26e+04(2.3e-08)	<b>-1.79e+04(2.7e-08)</b>	-3.02e+02(5.8e-06)(+)	
gissette	(6000-50-35)	<b>-2.55e+04(9.7e-06)</b>	-4.89e+03(3.1e-05)	-6.74e+03(9.6e-05)	<b>-7.15e+04(3.3e-08)</b>	-6.76e+04(3.0e-08)	-6.74e+03(9.6e-05)(+)	
Mnist	(6000-92-70)	<b>-8.05e+04(41.6e-06)</b>	-1.71e+04(3.2e-04)	-2.15e+03(3.0e-06)	-9.38e+03(1.1e-08)	-5.00e-01(3.3e-09)	-2.15e+03(3.0e-06)	
news20	(7967-50-35)	-1.56e-02(7.6e-06)	-1.42e+00(1.4e-05)	-1.56e+00(1.0e-04)	-1.98e+01(1.8e-07)	<b>-1.95e+01(1.8e-07)</b>	-1.56e+00(1.0e-04)(+)	
randn5000	(5000-100-75)	-1.72e+03(3.6e-06)	-6.27e+02(8.9e-04)	-9.15e+01(5.4e-06)	-3.25e+00(6.9e-10)	<b>-4.35e+01(7.1e-10)</b>	-9.15e+01(5.4e-06)	
w1a	(2477-100-75)	-8.30e+02(3.3e-06)	-1.06e+03(6.8e-06)	-1.90e+02(4.5e-05)	-1.12e+03(1.9e-09)	<b>-1.80e+03(5.1e-09)</b>	-1.90e+02(4.5e-05)	
20News	(9423-50-45)	-1.58e+01(7.6e-06)	-8.44e+					



2207  
 2208  
 Figure 7: Comparisons of objective values ( $F(\mathbf{X}) - F^0$ ) of HSPP for all the compared  
 2209  
 2210  
 2211  
 2212  
 2213  
 methods by epochs with different parameters ( $m - n - p$ ).

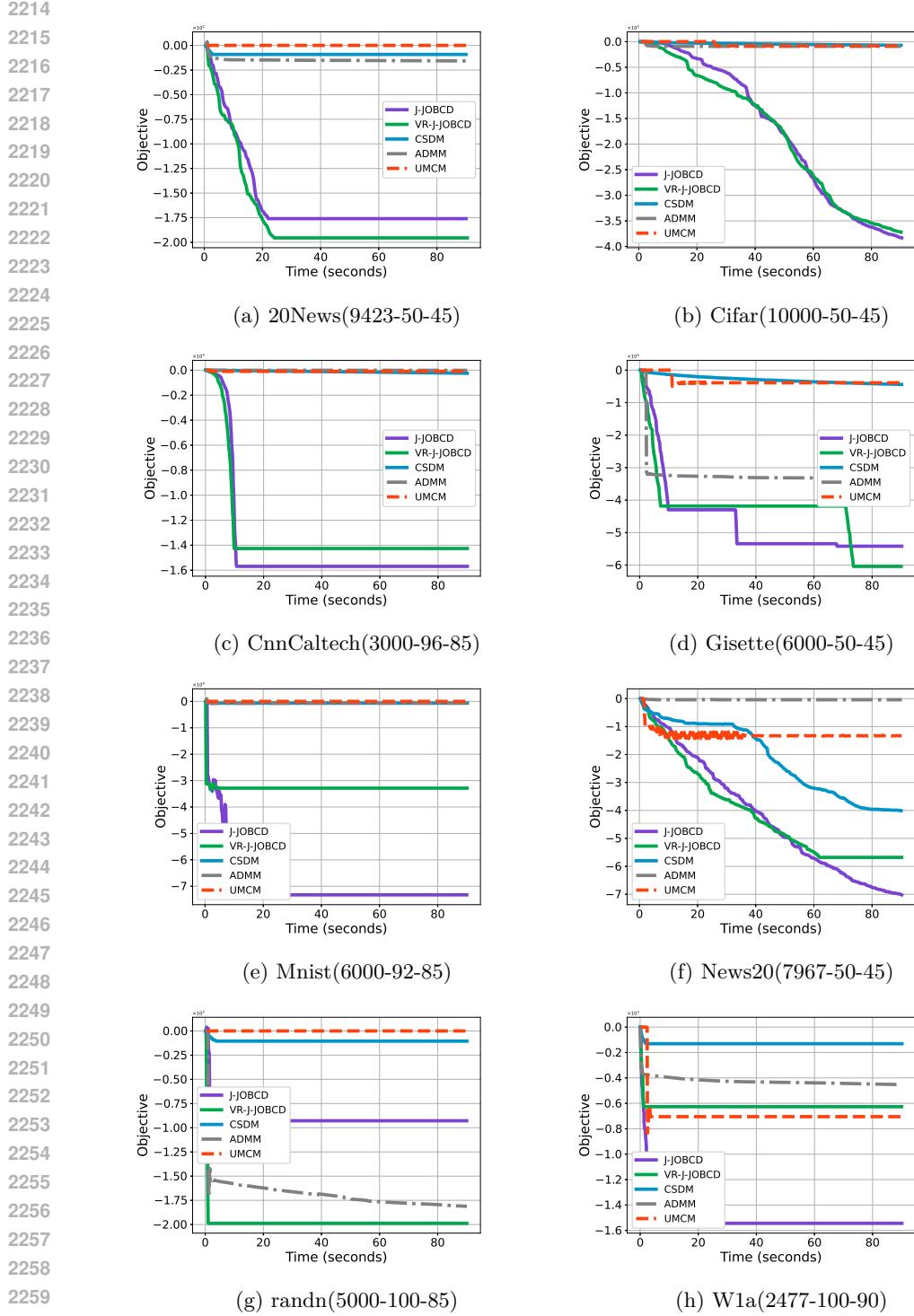


Figure 8: Comparisons of objective values ( $F(\mathbf{X}) - F^0$ ) of HSPP for all the compared methods by time with different parameters ( $m - n - p$ ).

---

### F.6.3 ULTRA-HYPERBOLIC KNOWLEDGE GRAPH EMBEDDING

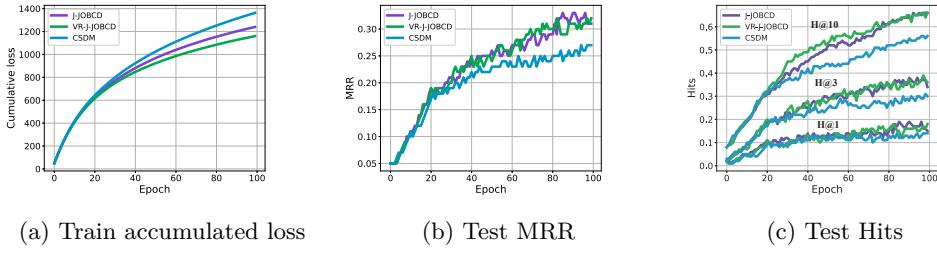


Figure 9: Epoch performance of CS, J-JOBCD, and VR-J-JOBCD in training UltraE on FB15k.

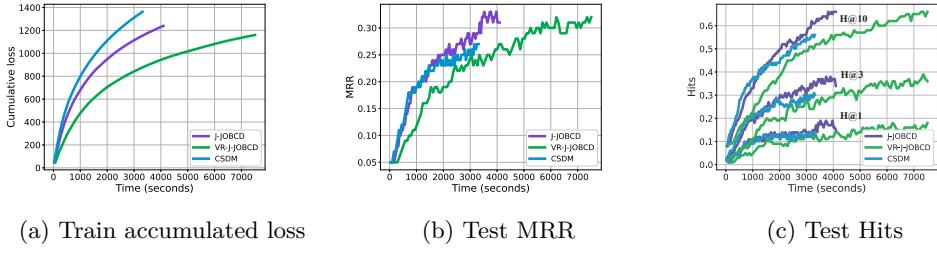


Figure 10: Time performance of CS, J-JOBCD, and VR-J-JOBCD in training UltraE on FB15k.

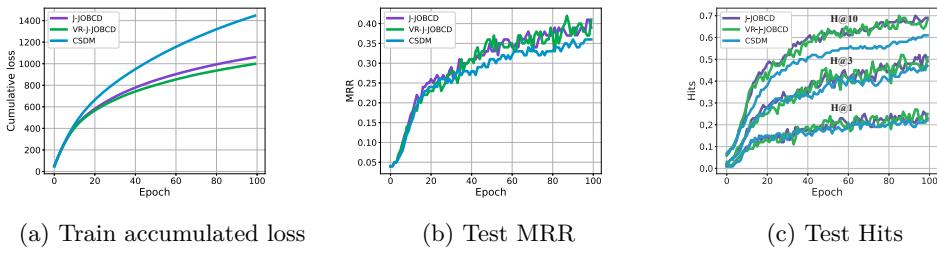


Figure 11: Epoch performance of CSDM, J-JOBCD, and VR-J-JOBCD in training UltraE on WN18RR.

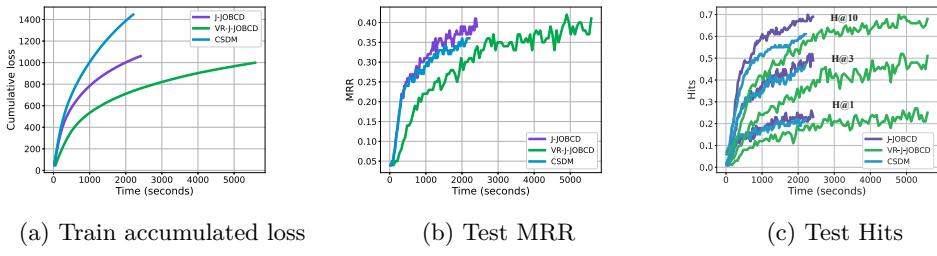


Figure 12: Time performance of CSDM, J-JOBCD, and VR-J-JOBCD in training UltraE on WN18RR.