

A NEW FIRST-ORDER ALGORITHMIC FRAMEWORK FOR OPTIMIZATION PROBLEMS WITH ORTHOGONALITY CONSTRAINTS*

BIN GAO[†], XIN LIU[†], XIAOJUN CHEN[‡], AND YA-XIANG YUAN[§]

Abstract. In this paper, we consider a class of optimization problems with orthogonality constraints, the feasible region of which is called the Stiefel manifold. Our new framework combines a function value reduction step with a correction step. Different from the existing approaches, the function value reduction step of our algorithmic framework searches along the standard Euclidean descent directions instead of the vectors in the tangent space of the Stiefel manifold, and the correction step further reduces the function value and guarantees a symmetric dual variable at the same time. We construct two types of algorithms based on this new framework. The first type is based on gradient reduction including the gradient reflection (GR) and the gradient projection (GP) algorithms. The other one adopts a columnwise block coordinate descent (CBCD) scheme with a novel idea for solving the corresponding CBCD subproblem inexactly. We prove that both GR/GP with a fixed step size and CBCD belong to our algorithmic framework, and any clustering point of the iterates generated by the proposed framework is a first-order stationary point. Preliminary experiments illustrate that our new framework is of great potential.

Key words. orthogonality constraint, Stiefel manifold, Householder transformation, gradient projection, block coordinate descent

AMS subject classifications. 15A18, 65F15, 65K05, 90C06

DOI. 10.1137/16M1098759

1. Introduction. We consider numerical methods for solving the following matrix variable optimization problem with orthogonality constraints,

$$(1.1) \quad \begin{aligned} \min_{X \in \mathbb{R}^{n \times p}} \quad & f(X) \\ \text{s. t.} \quad & X^\top X = I_p, \end{aligned}$$

where I_p stands for the p -by- p identity matrix, $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$, with $p \leq n$, satisfying the following assumption.

Assumption 1.1 (blanket assumption).

(i) f is twice differentiable. We define ρ as

$$\rho := \sup_{X \in \tilde{\mathcal{S}}} \|\nabla^2 f(X)\|_2,$$

where $\tilde{\mathcal{S}} := \{Y \mid \|Y\|_F^2 < p + 1\}$.¹

*Received by the editors October 13, 2016; accepted for publication (in revised form) October 6, 2017; published electronically February 1, 2018.

<http://www.siam.org/journals/siopt/28-1/M109875.html>

Funding: The second author's research supported in part by NSFC grants 11622112, 11471325, 91530204, and 11688101, the National Center for Mathematics and Interdisciplinary Sciences, CAS, and Key Research Program of Frontier Sciences, CAS. The third author's research supported in part by Hong Kong Research Council Grant N.PolyU504/14. The fourth authors research supported in part by NSFC grants 11331012 and 11461161005.

[†]State Key Laboratory of Scientific and Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, and University of Chinese Academy of Sciences, China (gaobin@lsec.cc.ac.cn, liuxin@lsec.cc.ac.cn).

[‡]Department of Applied Mathematics, The Hong Kong Polytechnic University, Hong Kong, China (xiaojun.chen@polyu.edu.hk).

[§]State Key Laboratory of Scientific and Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, China (yyx@lsec.cc.ac.cn).

¹In fact, $\tilde{\mathcal{S}}$ can be defined as any given bounded open set containing $\mathcal{S}_{n,p} := \{Y \in \mathbb{R}^{n \times p} \mid Y^\top Y = I_p\}$.

- (ii) $f(X)$ can be represented as $h(X) + \text{tr}(G^\top X)$, where $G \in \mathbb{R}^{n \times p}$, and $h(X)$ is orthogonal invariant, namely, $h(XQ) = h(X)$ holds for any $Q \in \mathcal{S}_{p,p}$, and $\nabla h(X) = H(X)X$, where $H : \mathbb{R}^{n \times p} \rightarrow \mathbb{S}^n$ is a matrix function.

Here, \mathbb{S}^n refers to the set of n -by- n symmetric matrices. The feasible region of problem (1.1) can be consequently denoted as $\mathcal{S}_{n,p}$. In practice, the value of ρ is often not known and difficult to estimate. Fortunately, we can overcome this difficulty in computation as shown in section 5.1.

Optimization problems of the above type with orthogonality constraints have many applications in scientific engineering computing and data science. More specifically, they play an important role in electronic structure calculations [35, 36, 34], linear eigenvalue problems [6], low-rank correlation matrix problems [14], sparse principal component analysis [39, 8], the orthogonal Procrustes problem [27, 11], etc. For other applications, we refer the interested readers to [10, 33, 17].

Remark 1.2. If $\rho = 0$, the objective function $f(X)$ reduces to a linear function $\text{tr}(G^\top X)$. In this case, the solution of (1.1) has the closed form $X = -RQ^\top$, where RSQ^\top is the reduced singular value decomposition² of G . In this paper, this special situation will not be discussed.

Assumption 1.1 is sufficient for the global convergence of our algorithmic framework. In this paper, we will not investigate how to weaken this sufficient condition. Fortunately, many interesting problems satisfy this assumption. Here are two simple examples.

Example 1.1.

$$f(X) := \frac{1}{2} \text{tr}(X^\top AX) + \text{tr}(G^\top X),$$

where $A \in \mathbb{S}^n$. In this case

$$\nabla f(X) = AX + G.$$

We notice that if the objective function defined in Example 1.1 takes $G = 0$, the corresponding optimization problem with orthogonality constraints (1.1) reduces to the Rayleigh–Ritz minimization which is exactly the optimization model for the eigenvalue problem. However, the problem with $G \neq 0$ is difficult to solve, even if A is positive definite. Example 1.1 is a key subproblem in the trust-region method for solving optimization problems with orthogonality constraints (see [36, (4.5)]). Therefore, it is challenging and interesting to explore efficient solvers for this problem.

Example 1.2.

$$f(X) := \frac{1}{2} \text{tr}(X^\top AX) + \frac{1}{2} \sum_{i=1}^m q_i(z),$$

where $z = \text{diag}(XX^\top)$, $q_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ($i = 1, \dots, m$), and $A \in \mathbb{S}^n$. In this case,

$$\nabla f(X) = \left(A + \sum_{i=1}^m \text{Diag}(\nabla q_i(z)) \right) X.$$

This example often appears in electronic structure calculations [20], which is one of the most important topics in materials science.

²For $G \in \mathbb{R}^{n \times p}$ with $p < n$, the reduced singular value decomposition refers to RSQ^\top , where $R \in \mathbb{R}^{n \times p}$ and $Q \in \mathbb{R}^{p \times p}$ are orthogonal matrices and $S \in \mathbb{R}^{p \times p}$ is diagonal.

1.1. Overview of existing methods. In general, it is difficult to find a global solution of problem (1.1) due to the nonconvexity. In fact, finding a stationary point or a feasible point is not an easy task because it can be numerically expensive to maintain the orthogonality for large p . There are some existing infeasible methods such as the splitting method [19] or the penalty method for large-scale eigenspace computation [32]. However, the former does not guarantee global convergence, and the latter only works for a very special case. Exploring practically useful infeasible methods for optimization problems with orthogonality constraints is beyond the discussion of this paper.

Recently, some algorithms have been developed for special cases of (1.1), such as electronic structure calculations [38, 31], dominant eigenpair calculation [21, 22], computing the coupling between matrices [12]. Usually, these approaches utilize the special structures of the problems and can hardly be extended to the generic optimization problems with orthogonality constraints.

The feasible region of problem (1.1), $\mathcal{S}_{n,p}$, is usually called the Stiefel manifold [28]. Various optimization methods designed for solving optimization problems restricted on a matrix manifold can be applied to problem (1.1). For instance, gradient-based methods [23, 24, 1], conjugate gradient methods [10, 2], trust-region methods [36], Newton methods [10], quasi-Newton methods [26, 16, 15], etc. The key principle of these methods is to find a feasible point with a lower function value than that at the current iterate. In [10, 3], the authors study the geometric structure of the Stiefel manifold from the optimization point of view, and bring up a new concept, which is called “retraction,” to connect previously unrelated algorithms. A map $\mathcal{R}_X : \mathcal{T}_X \mathcal{S}_{n,p} \rightarrow \mathcal{S}_{n,p}$ is called a retraction if the following properties hold:

- (1) $\mathcal{R}_X(0_X) = X$, where 0_X is the origin of $\mathcal{T}_X \mathcal{S}_{n,p}$;
- (2) $\frac{d}{dt} \mathcal{R}_X(tZ)|_{t=0} = Z$ for all $Z \in \mathcal{T}_X \mathcal{S}_{n,p}$,

where $\mathcal{T}_X \mathcal{S}_{n,p} := \{Y \in \mathbb{R}^{n \times p} \mid Y^\top X + X^\top Y = 0\}$ is the tangent space of the Stiefel manifold $\mathcal{S}_{n,p}$ at point X . The retraction \mathcal{R}_X maps a tangent vector into the manifold, so it defines an update rule to preserve the orthogonality.

There are two major classes of retractions for optimization problems with orthogonality constraints. The first one searches along the geodesic of a manifold to find a suitable trial point. Methods in this class are called geodesic-like retractions [10, 1, 3]. Calculating geodesics involves solving ordinary differential equations which often causes computational difficulties. The authors of [24] propose a quasi-geodesic updating formula based on the Cayley transformation whose main computation is to solve an n -by- n linear system. The methods in the other major class consist of two steps, line search in the tangent space and projection back to the Stiefel manifold. Thus, they are called projection-like methods [23, 3, 4]. The orthogonal projection can be calculated by QR factorization or polar decomposition. The projection-like methods coincide with the geodesic-like methods, in the special case of $p = 1$. The above mentioned retraction-based approaches, including both geodesic-like and projection-like methods, should work with a certain line search strategy, such as the Armijo inexact line search [25, 29] or a nonmonotonic line search strategy. The line search procedure is to guarantee the global convergence, but in the meantime, it induces additional function value evaluations.

Recently, Wen and Yin [33] proposed a feasible method for optimization with orthogonality constraints. In their work, an efficient way to calculate the Cayley transformation is introduced. In each iteration, it only requires one to solve a $2p \times 2p$ linear system instead of an $n \times n$ one. Combining a curvilinear search algorithm [13] with Barzilai–Borwein (BB) [5] nonmonotonic line search [37], it achieves much

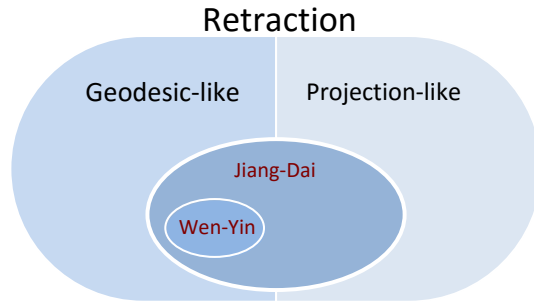


FIG. 1. Relationship among retraction-based methods.

lower computational cost than the other existing retraction-based algorithms and is illustrated to have robust numerical performance in solving a bunch of optimization problems with orthogonality constraints. Later on, Jiang and Dai [17] significantly extended the idea of [33], and found out that a large group of retraction-based methods enjoy such a reducible iterative formulation. It can be proved that all the algorithms under their framework with BB nonmonotonic line search are globally convergent to a stationary point.

In order to clarify the difference among the aforementioned retraction-based algorithms, we demonstrate their relationship through Figure 1.

It is worth mentioning that the retraction-based algorithms highly depend on the geometry of the Stiefel manifold and hence have very low compatibility with additional constraints such as nonnegative constraints or linear inequality constraints.

1.2. Contributions. In this paper, we revisit the first-order optimality condition of problem (1.1), and find that it is of the following form,

$$(1.2) \quad \begin{cases} (I_n - XX^\top)\nabla f(X) = 0, & \text{substationarity,} \\ X^\top \nabla f(X) = \nabla f(X)^\top X, & \text{symmetry,} \\ X^\top X = I_p, & \text{feasibility.} \end{cases}$$

For convenience, we call the three equalities of (1.2) substationarity,³ symmetry, and feasibility, respectively. Based on the first-order optimality condition, we propose a new algorithmic framework consisting of two main steps.

The first step is function value reduction. Namely, we find a feasible point which reduces the objective function value to a certain amount in proportion to the norm square of the projected gradient. We then propose two types of algorithms which can achieve such a requirement. Gradient reflection (GR) and gradient projection (GP) are the representatives of the first type of algorithm which uses different strategies to pull a gradient descent point back to the Stiefel manifold. The second type of algorithm employs a columnwise block coordinate descent (CBCD) iteration. A novel idea for solving the corresponding subproblem efficiently is proposed.

The second step is to find a feasible point satisfying the symmetry property. This correction step, whose main calculation is a $p \times p$ singular value decomposition, is highly dependent on Assumption 1.1. The correction step can be viewed as a rotation of the trial point obtained in the first step. In the special cases when $p = 1$ or $G = 0$, the symmetry of (1.2) always holds and hence this step can be waived.

³Here, substationarity stands for the stationarity of the gradient of the objective function in the null space of X^\top .

It is worth mentioning that GR and GP iterations belong to particular retractions if the correction step is not necessary, but either CBCD or GR/GP with correction step is not a retraction-based iteration. According to the construction way, the proposed algorithmic framework is expected to be compatible with additional nonmanifold constraints. Our framework exposes the essential mechanism of the gradient methods for optimization problems with orthogonality constraints, with which the global convergence of gradient-based algorithms with fixed step sizes can be established. Moreover, the numerical experiments for solving a class of generic quadratic minimization problems and the instances arising from electronic structure calculations show that our new algorithmic framework performs robustly and more efficiently than the existing algorithms.

Finally, the global convergence of CBCD is of great potential itself, as this is the first convergence result for the BCD method for nonconvex optimization problems with coupled constraints.

1.3. Organization. The rest of this paper is organized as follows. In section 2, we study the first-order optimality condition of problem (1.1), and provide a new first-order framework. The main step of our new framework is only required to meet a condition for sufficient function value reduction. We then develop two types of algorithms, in section 3, to achieve this requirement and form three concrete algorithms under the scheme of the new framework, namely, GR, GP, and CBCD, respectively. Global convergence of our new algorithmic framework is established in section 4. In section 5, we demonstrate the efficiency of our algorithmic framework in solving a class of general quadratic minimization problems and the energy minimization problem arising from the electronic structure calculations. We show the great potential of our proposed approach in solving large-scale problems. Finally, concluding remarks are given in the last section.

1.4. Notation. The Euclidean inner product of two matrices $X, Y \in \mathbb{R}^{n \times p}$ is defined as $\langle X, Y \rangle = \text{tr}(X^T Y)$, where $\text{tr}(A)$ is the trace of a matrix $A \in \mathbb{R}^{p \times p}$. $\|\cdot\|_2$ and $\|\cdot\|_F$ represent the 2-norm and the Frobenius norm, respectively. The notations $\text{diag}(A)$ and $\text{Diag}(x)$ stand for the vector formed by the diagonal entries of matrix A , and the diagonal matrix with the entries of $x \in \mathbb{R}^n$ to be its diagonal, respectively. X^\dagger refers to the pseudoinverse of X . We denote the smallest positive eigenvalue and the smallest eigenvalue in magnitude of A by $\lambda_{\min}^+(A)$ and $\lambda_{|\min|}(A)$, respectively. The i th column of matrix $X \in \mathbb{R}^{n \times p}$ is denoted by X_i . $X_{\bar{i}} \in \mathbb{R}^{n \times (p-1)}$ denotes the matrix X with its i th column removed, i.e., $X_{\bar{i}} = [X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p]$. We use $X_{i,v} \in \mathbb{R}^{n \times p}$ to denote X with its i th column replaced by a given vector v , i.e., $X_{i,v} = [X_1, \dots, X_{i-1}, v, X_{i+1}, \dots, X_p]$. Finally, $\mathcal{B}(C, r)$ is the ball defined as $\{X \in \mathbb{R}^{m_1 \times m_2} \mid \|X - C\|_F \leq r\}$, where $C \in \mathbb{R}^{m_1 \times m_2}$ is the center and r is the radius. $\mathbf{qr}(X)$ is the Q matrix of the reduced QR decomposition⁴ of X . $\mathcal{P}_{\mathcal{S}_{n,p}}(X)$ denotes the projection⁵ of X to the Stiefel manifold $\mathcal{S}_{n,p}$. Finally, $\mathbf{rand}(n, p)$ and $\mathbf{randn}(n, p)$ represent $n \times p$ randomly generated matrices under independently and identically distributed (i.i.d.) uniform distribution in $[0, 1]$ and i.i.d. standard Gaussian distribution, respectively.

2. A new first-order framework. In this section, we first give a new presentation of the first-order optimality condition of the optimization problem with

⁴ $Q \in \mathbb{R}^{n \times p}$ is the Q matrix of the reduced QR decomposition of $X \in \mathbb{R}^{n \times p}$, if $X = QR$, $Q \in \mathbb{R}^{n \times p}$ is orthogonal, and $R \in \mathbb{R}^{p \times p}$ is an upper triangle matrix.

⁵ $\mathcal{P}_{\mathcal{S}_{n,p}}(X) = \hat{U}\hat{V}^T$, where $\hat{U}\hat{\Sigma}\hat{V}^T$ is the reduced singular value decomposition of X .

orthogonality constraints (1.1), which motivates our new framework. The details of the new framework will also be presented.

2.1. Optimality condition. The first-order optimality condition of problem (1.1) can be interpreted as follows.

DEFINITION 2.1. *Given a point $X \in \mathbb{R}^{n \times p}$, if the relationship*

$$(2.1) \quad \begin{cases} \operatorname{tr}(Y^\top \nabla f(X)) \geq 0, \\ X^\top X = I_p, \end{cases}$$

holds for any $Y \in \mathcal{T}_X \mathcal{S}_{n,p}$, we call X a first-order stationary point of (1.1). The set containing all the first-order stationary points is denoted as Ω_{FON} .

Since condition (2.1) cannot be verified numerically, we show the following equivalent result.

LEMMA 2.2. *A point X is a first-order stationary point if and only if equalities (1.2) hold.*

Proof. We notice that any $Y \in \mathcal{T}_X \mathcal{S}_{n,p}$ can be uniquely decomposed as $Y = XS + K$,⁶ where $S \in \mathbb{R}^{p \times p}$ is a skew matrix (i.e., $S^\top + S = 0$) and $K \in \mathbb{R}^{n \times p}$ satisfies $K^\top X = 0$, which is equivalent to $K = (I_n - XX^\top)K$. Likewise, any matrix of the form $XS + K$ lies in $\mathcal{T}_X \mathcal{S}_{n,p}$.

Since S and K are arbitrary, condition (2.1) is equivalent to the following relationships,

$$(2.2) \quad \operatorname{tr}(S^\top X^\top \nabla f(X)) \geq 0 \quad \forall S \in \mathbb{R}^{p \times p} \text{ and } S^\top + S = I_p,$$

$$(2.3) \quad \begin{aligned} \operatorname{tr}(K^\top \nabla f(X)) &\geq 0 \quad \forall K \in \mathbb{R}^{n \times p} \text{ and } K^\top X = 0, \\ X^\top X &= I_p. \end{aligned}$$

By using (2.2) and the skew symmetry of $Q^\top - Q$, where $Q := X^\top \nabla f(X)$, we obtain

$$(2.4) \quad \operatorname{tr}((Q - Q^\top)Q) \geq 0.$$

It then follows from (2.4) that

$$\begin{aligned} 0 &\leq \operatorname{tr}((Q - Q^\top)Q) + \operatorname{tr}((Q - Q^\top)Q) = \operatorname{tr}(QQ - Q^\top Q) + \operatorname{tr}(Q^\top(Q^\top - Q)) \\ &= \operatorname{tr}(QQ - Q^\top Q) + \operatorname{tr}((Q^\top - Q)Q^\top) = \operatorname{tr}(QQ - Q^\top Q + Q^\top Q^\top - QQ^\top) \\ &= \operatorname{tr}((Q - Q^\top)(Q - Q^\top)) = -\operatorname{tr}((Q - Q^\top)^\top(Q - Q^\top)) \leq 0. \end{aligned}$$

This implies $Q = Q^\top$. On the other hand, if $X^\top \nabla f(X)$ is symmetric, the equality $\operatorname{tr}(S^\top X^\top \nabla f(X)) = 0$ holds for any skew symmetric matrix S . Hence, (2.2) is equivalent to the symmetry of $X^\top \nabla f(X)$.

Following from the property $K = (I_n - XX^\top)K$ and the arbitrariness of K , we can easily obtain the equivalence between (2.3) and $(I_n - XX^\top)\nabla f(X) = 0$. This completes the proof. \square

⁶This is because $Y = XX^\top Y + (I_n - XX^\top)Y$ for any Y . It is not difficult to verify that $(I_n - XX^\top)Y$ satisfies $Y^\top(I_n - XX^\top)^\top X = 0$. $X^\top Y$ is a skew matrix due to the fact that $Y \in \mathcal{T}_X \mathcal{S}_{n,p} := \{Y \in \mathbb{R}^{n \times p} \mid Y^\top X + X^\top Y = 0\}$.

Remark 2.3. It is very easy to check that our first-order optimality condition (1.2) in the Euclidean space is exactly the same as the one in the tangent space:

$$\begin{cases} \nabla f(X) - X\nabla f(X)^\top X = 0, \\ X^\top X = I_p, \end{cases}$$

which is stated in [33]. Moreover, it actually holds that

$$(2.5) \quad \begin{aligned} \|\nabla f(X) - X\nabla f(X)^\top X\|_F^2 &= \|\nabla f(X) - XX^\top \nabla f(X)\|_F^2 \\ &\quad + \|X^\top \nabla f(X) - \nabla f(X)^\top X\|_F^2. \end{aligned}$$

2.2. Correction step and algorithm framework. We notice that there are three properties, substationarity, symmetry, and feasibility in our first-order optimality condition (1.2) of problem (1.1). Motivated by the relationship (2.5), to make the gradient in the tangent space equal to zero, we can adopt the following two step procedure. From the current iterate, we first find a trial point which reduces the function value in proportion to the norm square of the projected gradient. Based on this trial point, we then find the next iterate which makes the symmetry property hold without increasing the function value. Then we repeat the procedure until convergence. In these two steps, the feasibility holds all the time. The details of these two steps are described in the following.

Suppose the current iteration point is X^k . In the first step, we find an intermediate point $\bar{X} \in \mathcal{S}_{n,p}$, which satisfies sufficient function value reduction, i.e.,

$$(2.6) \quad f(X^k) - f(\bar{X}) \geq C_1 \cdot \left\| (I_n - X^k X^{k\top}) \nabla f(X^k) \right\|_F^2,$$

where $C_1 > 0$ is a positive constant. The right-hand side of (2.6) measures the square of the Frobenius norm of the projected gradient at X^k in the Euclidean space.

Although the intermediate point $\bar{X} \in \mathcal{S}_{n,p}$ satisfies (2.6), it does not satisfy the symmetry property in (1.2). In the second part, we consider constructing a correction step which makes the symmetry property hold without increasing the function value.

Resulting from Assumption 1.1, it holds that

$$(2.7) \quad \bar{X}^\top \nabla f(\bar{X}) = \bar{X}^\top H(\bar{X})\bar{X} + \bar{X}^\top G.$$

The term $\bar{X}^\top H(\bar{X})\bar{X}$ is symmetric. Hence, the next iterate X^{k+1} can take \bar{X} , if $\bar{X}^\top G$ is symmetric. Otherwise, it suffices to find a point X^{k+1} satisfying the symmetry property $X^{k+1\top} G = G^\top X^{k+1}$. To achieve this, we use the rotation correction, namely, $X^{k+1} = -\bar{X}UT^\top$, where U and T come from the singular value decomposition of a $p \times p$ matrix

$$(2.8) \quad \bar{X}^\top G = U\Lambda T^\top.$$

The motivation of this correction step is to find a $p \times p$ orthogonal matrix Q^* which minimizes $f(\bar{X}Q^*)$, and then set $X^{k+1} = \bar{X}Q^*$. By recalling Assumption 1.1, we obtain $Q^* = -UT^\top$, which is the global minimizer of

$$\min_{Q \in \mathcal{S}_{p,p}} \text{tr} \left((\bar{X}Q)^\top G \right).$$

Therefore, we set the next iterate as

$$(2.9) \quad X^{k+1} = \begin{cases} \bar{X} & \text{if } \bar{X}^\top G = G^\top \bar{X}, \\ -\bar{X}UT^\top & \text{otherwise.} \end{cases}$$

We can then establish the following properties of such a correction step X^{k+1} .

LEMMA 2.4. Suppose $\bar{X} \in \mathcal{S}_{n,p}$. Let $\{X^{k+1}\}$ be calculated by (2.9), where U and T are determined by (2.8). Then, it holds that $X^{k+1} \in \mathcal{S}_{n,p}$ and $X^{k+1\top} \nabla f(X^{k+1})$ is symmetric. Furthermore, we have

$$(2.10) \quad 8\theta (f(\bar{X}) - f(X^{k+1})) \geq \|\bar{X}^\top \nabla f(\bar{X}) - \nabla f(\bar{X})^\top \bar{X}\|_F^2,$$

where

$$(2.11) \quad \theta := \|G\|_2.$$

Proof. The orthogonality of X^{k+1} and the symmetry of $X^{k+1\top} \nabla f(X^{k+1})$ can be directly derived by formula (2.9). Next, we prove inequality (2.10). If $\theta = 0$, which means $\nabla f(X) = H(X)X$, then the symmetry of $\bar{X}^\top \nabla f(\bar{X})$ implies (2.10) immediately. On the other hand, according to Assumption 1.1, we have

$$(2.12) \quad \begin{aligned} f(\bar{X}) - f(X^{k+1}) &= h(\bar{X}) + \text{tr}(G^\top \bar{X}) - h(X^{k+1}) - \text{tr}(G^\top X^{k+1}) \\ &= \text{tr}(G^\top \bar{X} - G^\top X^{k+1}) \\ &= \text{tr}(U\Lambda T^\top + \Lambda) = \text{tr}(B + \Lambda), \end{aligned}$$

where $B = (\Lambda T^\top U + U^\top T \Lambda)/2$.

On the other hand,

$$(2.13) \quad \begin{aligned} \|\bar{X}^\top \nabla f(\bar{X}) - \nabla f(\bar{X})^\top \bar{X}\|_F^2 &= \|\bar{X}^\top G - G^\top \bar{X}\|_F^2 \\ &= \|U\Lambda T^\top - T\Lambda U^\top\|_F^2 = 2\text{tr}(\Lambda^2) - 2\text{tr}(\Lambda T^\top U \Lambda T^\top U) \\ &= 4\text{tr}(\Lambda^2 - B^2), \end{aligned}$$

where the last equality uses the fact that

$$\text{tr}(B^2) = \frac{1}{2}\text{tr}(\Lambda^2) + \frac{1}{2}\text{tr}(\Lambda T^\top U \Lambda T^\top U).$$

Moreover, we have

$$(2.14) \quad \begin{aligned} \text{tr}(\Lambda^2 - B^2) &\leq \sum_{i=1}^p (\Lambda_{ii}^2 - B_i^\top B_i) \leq \sum_{i=1}^p (\Lambda_{ii}^2 - B_{ii}^2) = \sum_{i=1}^p (\Lambda_{ii} - B_{ii})(\Lambda_{ii} + B_{ii}) \\ &\leq \sum_{i=1}^p 2\Lambda_{ii}(\Lambda_{ii} + B_{ii}) \leq 2\|\Lambda\|_2 \cdot \sum_{i=1}^p (\Lambda_{ii} + B_{ii}) = 2\|\Lambda\|_2 \cdot \text{tr}(\Lambda + B) \\ &\leq 2\|G\|_2 \cdot \text{tr}(\Lambda + B) = 2\theta \cdot \text{tr}(\Lambda + B). \end{aligned}$$

Here the third inequality uses the fact that

$$|B_{ii}| = \Lambda_{ii} \cdot |T_i^\top U_i| \leq \Lambda_{ii}.$$

Combining (2.12)–(2.14), we complete the proof. \square

We can adopt

$$(2.15) \quad c(X) := (I_n - XX^\top) \nabla f(X),$$

because the symmetry and the feasibility of (1.2) hold at each iteration. The complete framework can be described as the following.

Algorithm 1: First-order framework for optimization problems with orthogonality constraints.

```

1 Set tolerance  $\epsilon > 0$ ; Initialize:  $X^0 \in \mathcal{S}_{n,p}$ ; Set  $k := 0$ 
2 while  $\|c(X^k)\|_F > \epsilon$  do
3   Based on  $X^k$ , find a feasible point  $\bar{X}$  satisfying (2.6);
4   Based on  $\bar{X}$ , calculate a feasible point  $X^{k+1}$  by (2.9);
5   Set  $k := k + 1$ .
6 Return  $X^k$ .
```

3. Algorithms for finding \bar{X} from X^k . In section 2, we propose a new algorithmic framework; however, how to find a point \bar{X} satisfying sufficient function value reduction (2.6) is still open. In this section, we introduce two types of algorithms to achieve step 3 in Algorithm 1. The first type of algorithm is based on gradient descent in the Euclidean space which will be introduced in the first two subsections. The second type of algorithm adopts a columnwise coordinate descent idea, and it will be introduced in the third subsection. In the last subsection, we list the computational cost per iteration of some existing algorithms and our new proposed algorithms.

3.1. Gradient-type methods. An intuitive idea to reduce the function value in the Euclidean space is to take the gradient descent direction. Unfortunately, a trial point obtained by a gradient descent step from the current iterate may violate the orthogonality constraint. Therefore, in this section we discuss two concrete strategies to pull the trial point back to the Stiefel manifold. Each of them can be used in step 3 in Algorithm 1.

Both strategies are based on the following observation.

LEMMA 3.1. *For any $Y \in \mathcal{B}_{X,\tau} := \mathcal{B}(X - \tau \nabla f(X), \tau \|\nabla f(X)\|_F)$, where $\tau \in (0, \rho^{-1})$, it holds that*

$$(3.1) \quad f(X) - f(Y) \geq \frac{1 - \rho\tau}{2\tau} \cdot \|X - Y\|_F^2.$$

Proof. For any $Y \in \mathcal{B}_{X,\tau}$, we can derive

$$\langle Y - X, Y - X + 2\tau \nabla f(X) \rangle \leq 0,$$

which implies

$$\begin{aligned} f(Y) &\leq f(X) + \langle Y - X, \nabla f(X) \rangle + \frac{\rho}{2} \|Y - X\|_F^2 \\ &= f(X) + \frac{1}{2\tau} \cdot \langle Y - X, Y - X + 2\tau \nabla f(X) \rangle - \frac{\tau^{-1} - \rho}{2} \cdot \|Y - X\|_F^2 \\ &\leq f(X) - \frac{\tau^{-1} - \rho}{2} \cdot \|Y - X\|_F^2. \end{aligned}$$

This completes the proof. \square

We illustrate the relationship among the feasible region, current iterate, gradient step, and the auxiliary ball $\mathcal{B}_{X,\tau}$ in Figure 2.

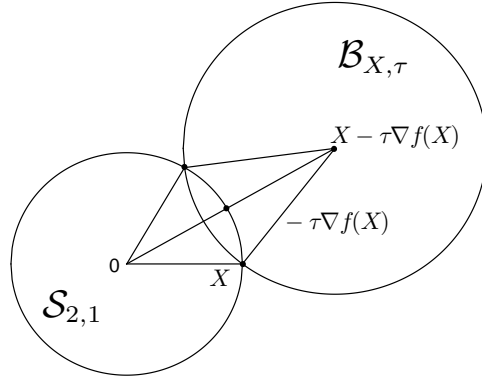


FIG. 2. Gradient-type method.

3.1.1. Gradient reflection. The first possible choice of a feasible trial point is to take the reflection point of the current iterate X^k reflecting on the null space of $X^k - \tau \nabla f(X^k)$. This point can be actually calculated by the Householder transformation

$$(3.2) \quad \mathbf{GR} : \quad \begin{cases} V = X^k - \tau \nabla f(X^k) & \text{for a fixed chosen } \tau \in (0, \rho^{-1}), \\ \bar{X}_{\text{GR}} = (-I_n + 2V(V^\top V)^\dagger V^\top)X^k. \end{cases}$$

Since \bar{X}_{GR} is the reflection point, we call Algorithm 1, using (3.2), to get $\bar{X} := \bar{X}_{\text{GR}}$ in step 3, the GR.

Next, we show the intermediate point \bar{X}_{GR} defined in (3.2) is feasible and achieves sufficient function value reduction (2.6).

LEMMA 3.2. *Let $X^k \in \mathcal{S}_{n,p}$ and \bar{X}_{GR} be defined by (3.2). Then it holds that $\bar{X}_{\text{GR}} \in \mathcal{S}_{n,p}$ and*

$$(3.3) \quad f(X^k) - f(\bar{X}_{\text{GR}}) \geq \frac{2(\tau^{-1} - \rho)}{(\tau^{-1} + \rho + \theta)^2} \cdot \left\| (I_n - X^k X^{k\top}) \nabla f(X^k) \right\|_{\text{F}}^2,$$

where $\tau \in (0, \rho^{-1})$, ρ , and θ are defined in Assumption 1.1 and equality (2.11), respectively.

Proof. With a slight abuse of notation, we omit the superscript k and use X to denote X^k in this proof.

First, by simple calculation, we have

$$\bar{X}_{\text{GR}}^\top \bar{X}_{\text{GR}} = X^\top (-I_n + 2V(V^\top V)^\dagger V^\top)^\top (-I_n + 2V(V^\top V)^\dagger V^\top) X = I_p,$$

which implies $\bar{X}_{\text{GR}} \in \mathcal{S}_{n,p}$.

Let RSQ^\top be the reduced singular value decomposition of V . If $S = 0$, we have $X = \tau \nabla f(X)$, which implies that $(I_n - XX^\top) \nabla f(X) = 0$, and inequality (3.3) holds immediately. Now we consider the case that $S \neq 0$. We have

$$(3.4) \quad \begin{aligned} & \|\bar{X}_{\text{GR}} - X\|_{\text{F}} \\ &= 2 \|(I_n - V(V^\top V)^\dagger V^\top) X\|_{\text{F}} = 2 \|(I_n - RR^\top) X\|_{\text{F}} = 2\sqrt{p - \|R^\top X\|_{\text{F}}^2} \\ &= 2 \|(I_n - XX^\top) R\|_{\text{F}} \geq 2 \|(I_n - XX^\top) VQS^\dagger\|_{\text{F}} \geq 2 \|(I_n - XX^\top) VQ\|_{\text{F}} \cdot \lambda_{\min}^+(S^\dagger) \\ &= 2 \|(I_n - XX^\top) V\|_{\text{F}} / \|S\|_2 = 2\tau \|c(X)\|_{\text{F}} / \|V\|_2 \geq \frac{2}{\tau^{-1} + \rho + \theta} \|c(X)\|_{\text{F}}, \end{aligned}$$

where λ_{\min}^+ and $c(X)$ are defined in section 1.4 and equality (2.15), respectively. Here, the second inequality uses the fact that all the entries of the j th column of VQ are zero for any j satisfying $S_{jj} = 0$ which is implied by the equality $RS = VQ$. The last inequality of (3.4) results from

$$(3.5) \quad \begin{aligned} \|\nabla f(X)\|_2 &\leq \|H(X)X\|_2 + \|G\|_2 \leq \rho + \theta, \\ \|V\|_2 &\leq \|X\|_2 + \tau\|\nabla f(X)\|_2 \leq 1 + \tau(\rho + \theta). \end{aligned}$$

Substituting inequality (3.4) into (3.1) in Lemma 3.1 with $Y = \bar{X}_{\text{GR}}$, we arrive at

$$(3.6) \quad \begin{aligned} f(X) - f(\bar{X}_{\text{GR}}) &\geq \frac{4}{(\tau^{-1} + \rho + \theta)^2} \cdot \frac{\tau^{-1} - \rho}{2} \cdot \|(I_n - XX^\top) \nabla f(X)\|_{\text{F}}^2 \\ &= \frac{2(\tau^{-1} - \rho)}{(\tau^{-1} + \rho + \theta)^2} \cdot \|(I_n - XX^\top) \nabla f(X)\|_{\text{F}}^2, \end{aligned}$$

which completes the proof. \square

3.1.2. Gradient projection. Another possible choice of a feasible trial point is to directly take the projection of $X^k - \tau\nabla f(X^k)$ onto the Stiefel manifold, which can be calculated by,

$$(3.7) \quad \mathbf{GP} : \quad \begin{cases} V = X^k - \tau\nabla f(X^k) & \text{for a fixed chosen } \tau \in (0, \rho^{-1}), \\ \bar{X}_{\text{GP}} = \mathcal{P}_{\mathcal{S}_{n,p}}(V). \end{cases}$$

We call Algorithm 1 using (3.7) to get $\bar{X} := \bar{X}_{\text{GP}}$ in step 3 the GP. We can similarly prove the feasibility of \bar{X}_{GP} and show sufficient function value reduction can be achieved.

LEMMA 3.3. *Let $X^k \in \mathcal{S}_{n,p}$ and \bar{X}_{GP} be defined by (3.7). Then it holds that $\bar{X}_{\text{GP}} \in \mathcal{S}_{n,p}$ and*

$$(3.8) \quad f(X^k) - f(\bar{X}_{\text{GP}}) \geq \frac{\tau^{-1} - \rho}{2(\tau^{-1} + \rho + \theta)^2} \cdot \|(I_n - X^k X^{k\top}) \nabla f(X^k)\|_{\text{F}}^2,$$

where $\tau \in (0, \rho^{-1})$, ρ , and θ are defined in Assumption 1.1 and equality (2.11), respectively.

Proof. The first part of the argument can be derived in the same manner as Lemma 3.2. Here, we just focus on the proof of (3.8).

By using the singular value decomposition $V = RSQ^\top$ and the first two equalities of (3.4), we arrive at

$$\begin{aligned} \|\bar{X}_{\text{GP}} - X^k\|_{\text{F}}^2 &- \frac{1}{4} \|\bar{X}_{\text{GR}} - X\|_{\text{F}}^2 \\ &= \|RQ^\top - X^k\|_{\text{F}}^2 - \|(I_n - RR^\top)X^k\|_{\text{F}}^2 \\ &= \text{tr}(I_p) - 2\text{tr}(QR^\top X^k) + \text{tr}(I_p) - \text{tr}(I_p) + 2\text{tr}(X^{k\top}RR^\top X^k) - \|R^\top X^k\|_{\text{F}}^2 \\ &= p - 2\text{tr}(QR^\top X^k) + \|R^\top X^k\|_{\text{F}}^2 = \|Q^\top - R^\top X^k\|_{\text{F}}^2 \geq 0, \end{aligned}$$

which implies

$$\|\bar{X}_{\text{GP}} - X^k\|_{\text{F}} \geq \frac{1}{2} \|\bar{X}_{\text{GR}} - X^k\|_{\text{F}} \geq \frac{1}{\tau^{-1} + \rho + \theta} \|c(X)\|_{\text{F}}.$$

Then, we can obtain (3.8) along the lines of the proof of inequality (3.6), and then complete the proof. \square

3.2. CBCD method. Another popular first-order method is block coordinate descent. For optimization problems with orthogonality constraints, a natural way to build up blocks is to partition the variables by the columns. On the other hand, the convergence of block coordinate descent with blocks coupled in nonconvex constraints cannot be guaranteed by existing results. Therefore, it is worthwhile studying the CBCD for optimization problems with orthogonality constraints. In this subsection, we consider Algorithm 1 using CBCD in step 3, discuss the way to solve the subproblem efficiently, and prove that such an approach belongs to Algorithm 1.

Once we fix the values of $p - 1$ columns of X and only leave the i th column as variable, we arrive at the following subproblem:

$$(3.9) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f_{i,X}(x) \\ \text{s. t.} \quad & \|x\|_2 = 1, \\ & X_{\bar{i}}^\top x = 0, \end{aligned}$$

where $f_{i,X}(x) := f(X_{i,x})$, $X_{i,x}$, and $X_{\bar{i}}$ are defined in subsection 1.4.

Suppose we can obtain the solution of the above subproblem or find a feasible point x^+ with sufficient function value reduction compared to $f_{i,X}(X_i)$. Then we can use this feasible point to update our iterate columnwisely in a Gauss–Seidel manner. More specifically, if X is the current iterate, the trial point \bar{X} can be calculated by the following CBCD scheme.

Algorithm 2: Columnwise block coordinate descent.

- 1 Set $W^0 = X$, $i := 1$;
 - 2 **while** $i \leq p$ **do**
 - 3 Solve the subproblem (3.9) with X replaced by W^{i-1} , and obtain feasible point x^+ satisfying the following sufficient function value reduction and asymptotic small step size safeguard

$$(3.10) \quad f_{i,W^{i-1}}(X_i) - f_{i,W^{i-1}}(x^+) \geq k_1 \|X_i - x^+\|_2^2,$$

$$(3.11) \quad \|X_i - x^+\|_2 \geq k_2 \|(I_n - W^{i-1}W^{i-1\top})\nabla f_{i,W^{i-1}}(X_i)\|_2;$$

Set $W^i = W_{i,x^+}^{i-1}$, $i := i + 1$;
 - 4 **Return** $\bar{X} = W^p$.
-

Remark 3.4. Algorithm 2 actually provides a cyclic CBCD scheme, i.e., the columns are updated in a cyclic order. We can similarly implement the greedy order, stochastic order (sampling with replacement), or randomly permuted order (sampling without replacement) which often appear in classical block coordinate descent algorithms. However, as we will show in section 5, these strategies will not help to improve the performance of the cyclic CBCD. Therefore, we omit the detailed descriptions and analyses of these strategies.

Before claiming that Algorithm 2 can find \bar{X} in step 3 of Algorithm 1, we need to answer two questions: can we cheaply calculate a solution or feasible point achieving sufficient function value reduction and asymptotic small step size safeguard (3.10)–(3.11)? Does Algorithm 2 provide a feasible point of problem (1.1) satisfying (2.6)? We answer these two questions in the following two subsections.

3.2.1. Solving the CBCD subproblem. In this subsection, we discuss how to obtain a feasible trial point of subproblem (3.9) efficiently. We notice that the second constraint of (3.9) restricts the variable x lying in the null space of $X_{\bar{i}}$. Hence, we can use the variable change $x = (I_n - X_{\bar{i}}X_{\bar{i}}^\top)x$ to reduce this constraint.

First, the fact $X_{\bar{i}}^\top x = 0$ holds if and only if $x = (I_n - X_{\bar{i}}X_{\bar{i}}^\top)x$. Hence, subproblem (3.9) is equivalent to the following problem,

$$(3.12) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f_{i,X} \left((I_n - X_{\bar{i}}X_{\bar{i}}^\top) x \right) \\ \text{s. t.} \quad & \left\| (I_n - X_{\bar{i}}X_{\bar{i}}^\top) x \right\|_2 = 1. \end{aligned}$$

Furthermore, problem (3.12) is equivalent to a well-posed problem if it is restricted to the null space of $X_{\bar{i}}$. More specifically, we have the following proposition.

PROPOSITION 3.5. *The equivalence between problem (3.9) restricted to a subspace \mathcal{D} and the following sphere constrained problem*

$$(3.13) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & q_i(x) := f_{i,X} \left((I_n - X_{\bar{i}}X_{\bar{i}}^\top) x \right) \\ \text{s. t.} \quad & \|x\|_2 = 1, \\ & x \in \mathcal{D} \end{aligned}$$

holds, if $X_{\bar{i}}^\top x = 0$ holds for any $x \in \mathcal{D}$.

Proof. For any $x \in \mathcal{D}$, it holds that $x = (I_n - X_{\bar{i}}X_{\bar{i}}^\top)x$ which implies the equivalence of problems (3.13) and (3.12) restricted to the subspace \mathcal{D} . By using the equivalence between problems (3.12) and (3.9), we complete the proof. \square

Proposition 3.5 tells us that we can calculate a feasible point of subproblem (3.9) with sufficient function value reduction through solving problem (3.13) if we can find a suitable subspace \mathcal{D} .

We notice that both X_i and $\nabla q_i(X_i) = (I_n - X_{\bar{i}}X_{\bar{i}}^\top)\nabla f_{i,X}((I_n - X_{\bar{i}}X_{\bar{i}}^\top)X_i)$ lie in the null space of $X_{\bar{i}}$. Therefore, any point in the subspace $\text{span}\{X_i, \nabla q_i(X_i)\}$ satisfies the orthogonality. In other words, $\text{span}\{X_i, \nabla q_i(X_i)\}$ is a qualified choice of orthogonal subspace \mathcal{D} in Proposition 3.5. Considering that subproblem (3.13) with $\mathcal{D} = \text{span}\{X_i, \nabla q_i(X_i)\}$ is a special case of the original optimization problem with orthogonality constraints (1.1) with $n = 2$ and $p = 1$, we recommend using the GR step (3.2) or the GP step (3.7) introduced in subsection 3.1 to calculate x^+ .

It can be verified that the GR as well as the GP step satisfy sufficient function value reduction (3.10) and asymptotic small step size safeguard (3.11).

LEMMA 3.6. *Let $x^+ = (-1 + 2v(v^\top v)^{-1}v^\top)X_i$ or $x^+ = (v^\top v)^{-\frac{1}{2}}v$, where $v = X_i - \tau \cdot \nabla q_i(X_i)$, $\tau \in (0, \rho^{-1})$. Then x^+ satisfies the constraints of (3.9) and conditions (3.10) and (3.11).*

The proof of Lemma 3.6 directly follows from Lemmas 3.1, 3.2, 3.3, and the fact that $I_n - XX^\top = (I_n - X_iX_i^\top)(I_n - X_{\bar{i}}X_{\bar{i}}^\top)$, and hence it is omitted here.

Remark 3.7. If $f_{i,X}$ is quadratic, subproblem (3.13) restricted to the subspace $\text{span}\{X_i, \nabla q_i(X_i)\}$ is equivalent to finding the roots of a quartic equation, which can be calculated in a closed form. In this case, the global minimizer of subproblem (3.13) restricted to the subspace $\text{span}\{X_i, \nabla q_i(X_i)\}$ can be an alternative option of x^+ .

3.2.2. Sufficient function value reduction. In this subsection, we show that \bar{X} calculated by Algorithm 2 is a feasible point of problem (1.1) and satisfies sufficient function value reduction (2.6).

LEMMA 3.8. Let $X \in \mathcal{S}_{n,p}$ and \bar{X} be calculated by Algorithm 2. Then it holds that $\bar{X} \in \mathcal{S}_{n,p}$ and

$$(3.14) \quad f(X) - f(\bar{X}) \geq \frac{k_1 k_2^2}{(1 + (p-1)k_2((1 + \sqrt{2})\rho + \sqrt{2}\theta))^2} \cdot \|(I_n - XX^\top) \nabla f(X)\|_F^2.$$

Proof. The feasibility of \bar{X} directly follows from the cyclic Gauss–Seidel-type update and the constraints of subproblem (3.9).

Now, we prove the second part. First, we have

$$(3.15) \quad \begin{aligned} f(X) - f(\bar{X}) &= f(W^0) - f(W^p) = \sum_{i=1}^p (f(W^{i-1}) - f(W^i)) \\ &= \sum_{i=1}^p (f_{i,W^{i-1}}(W_i^{i-1}) - f_{i,W^{i-1}}(W_i^i)) \end{aligned}$$

and

$$(3.16) \quad \begin{aligned} f_{i,W^{i-1}}(W_i^{i-1}) - f_{i,W^{i-1}}(W_i^i) \\ \geq k_1 k_2^2 \left\| (I_n - W^{i-1}W^{i-1\top}) \nabla f_{i,W^{i-1}}(W_i^{i-1}) \right\|_2^2. \end{aligned}$$

By using the Lipschitz continuity and the boundedness of gradient (3.5), we have

$$(3.17) \quad \begin{aligned} &\left\| (I_n - W^{i-1}W^{i-1\top}) \nabla f_{i,W^{i-1}}(W_i^{i-1}) \right\|_2 \\ &\geq \left\| (I_n - W^0W^{0\top}) \nabla f_{i,W^{i-1}}(W_i^{i-1}) \right\|_2 \\ &\quad - \left\| (W^{i-1}W^{i-1\top} - W^0W^{0\top}) \nabla f_{i,W^{i-1}}(W_i^{i-1}) \right\|_2 \\ &\geq \left\| (I_n - XX^\top) \nabla f_{i,W^{i-1}}(X_i) \right\|_2 \\ &\quad - \sum_{j=1}^{i-1} \left\| W^j W^{j\top} - W^{j-1} W^{j-1\top} \right\|_F \cdot \left\| \nabla f_{i,W^{i-1}}(W_i^{i-1}) \right\|_2 \\ &\geq \left\| (I_n - XX^\top) \nabla f_{i,W^0}(X_i) \right\|_2 - \left\| (I_n - XX^\top) (\nabla f_{i,W^0}(X_i) - \nabla f_{i,W^{i-1}}(X_i)) \right\|_2 \\ &\quad - (\rho + \theta) \sum_{j=1}^{i-1} \left(\sqrt{2 - 2(W_j^j{}^\top W_j^{j-1})^2} \right) \\ &\geq \left\| (I_n - XX^\top) \nabla f_{i,X}(X_i) \right\|_2 - \left\| \nabla f_{i,W^0}(X_i) - \nabla f_{i,W^{i-1}}(X_i) \right\|_2 \\ &\quad - \sqrt{2}(\rho + \theta) \cdot \sum_{j=1}^{i-1} \left(\sqrt{2 - 2(W_j^j{}^\top W_j^{j-1})^2} \right) \\ &\geq \left\| (I_n - XX^\top) \nabla f_{i,X}(X_i) \right\|_2 - \rho \cdot \|W^0 - W^{i-1}\|_2 - \sqrt{2}(\rho + \theta) \sum_{j=1}^{i-1} \|W_j^j - W_j^{j-1}\|_2 \\ &\geq \left\| (I_n - XX^\top) \nabla f_{i,X}(X_i) \right\|_2 \\ &\quad - \left((1 + \sqrt{2})\rho + \sqrt{2}\theta \right) \sqrt{k_1}^{-1} \sum_{j=1}^{i-1} \sqrt{f_{j,W^{j-1}}(W_j^{j-1}) - f_{j,W^{j-1}}(W_j^j)}, \end{aligned}$$

where the third last inequality uses the facts $|W_i^{i-1\top} W_i^i| \leq 1$ and $\sqrt{2-2\delta^2} \leq 2\sqrt{1-\delta}$ ($\forall |\delta| \leq 1$). Together with (3.16), we have

$$(3.18) \quad \begin{aligned} \sqrt{f_{i,W^{i-1}}(W_i^{i-1}) - f_{i,W^{i-1}}(W_i^i)} &\geq \sqrt{k_1 k_2} \left\| (I_n - W^{i-1} W^{i-1\top}) \nabla f_{i,W^{i-1}}(W_i^{i-1}) \right\|_2 \\ &\geq \sqrt{k_1 k_2} \left\| (I_n - X X^\top) \nabla f_{i,X}(X_i) \right\|_2 \\ &\quad - k_2 \left((1 + \sqrt{2}) \rho + \sqrt{2}\theta \right) \sum_{j=1}^{i-1} \sqrt{f_{j,W^{j-1}}(W_j^{j-1}) - f_{j,W^{j-1}}(W_j^j)}. \end{aligned}$$

Let $\delta_j := \sqrt{f_{j,W^{j-1}}(W_j^{j-1}) - f_{j,W^{j-1}}(W_j^j)}$, $c := k_2((1 + \sqrt{2})\rho + \sqrt{2}\theta)$, substituting relationship (3.18) into the fact that

$$\left(\delta_i + c \sum_{j=1}^{i-1} \delta_j \right)^2 \leq (1 + (i-1)c) \delta_i^2 + \sum_{j=1}^{i-1} c(1 + (i-1)c) \delta_j^2,$$

we obtain

$$(3.19) \quad \begin{aligned} (1 + (i-1)c) \delta_i^2 + \sum_{j=1}^{i-1} c(1 + (i-1)c) \delta_j^2 \\ \geq k_1 k_2^2 \left\| (I_n - X X^\top) \nabla f_{i,X}(X_i) \right\|_2^2. \end{aligned}$$

Summing up inequality (3.19) from $i = 1$ to p , and recalling (3.15), we arrive at

$$(3.20) \quad \begin{aligned} (1 + (p-1)k_2 \left((1 + \sqrt{2}) \rho + \sqrt{2}\theta \right))^2 \sum_{i=1}^p \left(\sqrt{f_{i,W^{i-1}}(W_i^{i-1}) - f_{i,W^{i-1}}(W_i^i)} \right)^2 \\ \geq \sum_{i=1}^p k_1 k_2^2 \cdot \left\| (I_n - X X^\top) \nabla f_{i,X}(X_i) \right\|_2^2 = k_1 k_2^2 \cdot \left\| (I_n - X X^\top) \nabla f(X) \right\|_F^2, \end{aligned}$$

which implies

$$f(X) - f(\bar{X}) \geq \frac{k_1 k_2^2}{(1 + (p-1)k_2 \left((1 + \sqrt{2}) \rho + \sqrt{2}\theta \right))^2} \cdot \left\| (I_n - X X^\top) \nabla f(X) \right\|_F^2.$$

This completes the proof. \square

A byproduct of the proof of Lemma 3.8 is the following asymptotic small step size safeguard property of CBCD.

COROLLARY 3.9. *Let $X \in \mathcal{S}_{n,p}$ and \bar{X} be calculated by Algorithm 2. Then it holds that*

$$(3.21) \quad \|X - \bar{X}\|_F \geq \frac{k_2}{1 + (p-1)k_2 \left((1 + \sqrt{2}) \rho + \sqrt{2}\theta \right)} \cdot \left\| (I_n - X X^\top) \nabla f(X) \right\|_F.$$

Proof. Using condition (3.11), the second last inequality of (3.17), and following along the same lines of inequalities (3.18) and (3.20), we can immediately obtain the desired result. \square

TABLE 1
Comparison on computational cost.

Update schemes	Computational cost	
	First τ	Subsequent τ
Geodesic-like algorithms		
$Y_{\text{geoc}}(\tau; X)$ [1]	$O(n^3)$	$O(n^3)$
$Y_{\text{qgeo}}(\tau; X)$ [24]	$O(n^3)$	$O(n^3)$
$Y_{\text{geoe}}(\tau; X)$ [10]	$10np^2 + 2np + O(p^3)$	$4np^2 + O(p^3)$
$Y_{\text{wy}}(\tau; X)$ [33]	$7np^2 + 2np + O(p^3)$	$4np^2 + np + O(p^3)$
Projection-like algorithms		
$Y_{\text{qr}}(\tau; X)$ [3]	$6np^2 + 3np + O(p^3)$	$2np^2 + 2np$
$Y_{\text{pd}}(\tau; X)$ [3]	$7np^2 + 4np + O(p^3)$	$2np^2 + 2np + O(p^3)$
$Y_{\text{pj}}(\tau; X)$ [23]	$7np^2 + 4np + O(p^3)$	$3np^2 + 3np + O(p^3)$
$Y_{\text{jd}}(\tau; X)$ [17]	$7np^2 + 3np + O(p^3)$	$2np^2 + 3np + O(p^3)$
Our algorithms		
GR	$9np^2 + 4np + O(p^3)$	
GP	$7np^2 + 3np + O(p^3)$	
CBCD-GR	$4np^2 + 8np + O(p^3)$	
CBCD-GP	$4np^2 + 5np + O(p^3)$	

3.3. Computational cost. In this subsection, we compare the computational cost per iteration among the existing algorithms and our GR, GP, and CBCD algorithms. First of all, we clarify the computational cost of the basic linear algebra operations as the following. Given $A \in \mathbb{R}^{n \times n}$, $B_1, B_2 \in \mathbb{R}^{n \times p}$, $S_1, S_2 \in \mathbb{R}^{p \times p}$, and $x \in \mathbb{R}^n$, calculating matrix-matrix products $B_1^\top B_2$, $B_1^\top B_1$, $B_1 S_1$, and $S_1 S_2$ we need $2np^2$, $np^2 + np$, $2np^2$, and $2p^3$ flops, respectively. Computing A^{-1} and S^{-1} we need $8n^3/3$ and $8p^3/3$ flops, respectively, and computing Ax needs $2n^2$ flops. Calculating the full singular value decomposition of $S \in \mathbb{R}^{p \times p}$ to a fixed precision costs $O(p^3)$ flops [30]. We assume $\nabla f(X)$ is already assembled and hence the computation of $\nabla f(X)$ is not counted in the computational cost per iteration. The other settings are similar to [17, Table 1]. We illustrate the comparison result as follows.

In Table 1, the two columns “first τ ” and “subsequent τ ” refer to the computational cost for the first trial point, and for subsequent trial points, respectively. However, the additional function evaluations have not been counted yet. For our GP, GR, and CBCD algorithms, line search is waived, as GR and GP converge with a fixed step size and the subproblem of CBCD only needs to be solved inexactly by one iteration. Hence, our computational cost per iteration is much cheaper than the retraction-based algorithms in general. Nevertheless, we have to point out that the computation time does not only depend on the flops count, but also an efficient use of the BLAS.

Moreover, CBCD-GR or CBCD-GP refer to the CBCD (using Algorithm 2 in step 3 of Algorithm 1) with the GR or GP updating formula used once in step 3 of Algorithm 2. We notice that the calculation of $\nabla f_{i,X}((I_n - W_i^{i-1} W_i^{i-1 \top}) X_i)$ is waived because it is equal to $\nabla f_{i,X}(X_i)$ which is implied by $W_i^{i-1 \top} X_i = 0$. If $f_{i,X}(X_i)$ ($i = 1, \dots, p$) are quadratic, and we solve the subproblem (3.13) restricted to the subspace $\text{span}\{X_i, \nabla q_i(X_i)\}$ to global optimality in Algorithm 2, the corresponding computational cost is $12np^2 + 3np + O(p^3)$.

4. Convergence analysis. In this section, we establish the global convergence of our new algorithmic framework, Algorithm 1. First, the function value convergence is shown.

LEMMA 4.1. *Let $\{X^k\}$ be the iterate sequence generated by Algorithm 1 initiated from a point $X^0 \in \mathcal{S}_{n,p}$, then $\{f(X^k)\}$ converges.*

Proof. According to the construction of \bar{X} in step 3 of Algorithm 1 and Lemma 2.4, we obtain

$$\begin{aligned}
 (4.1) \quad & f(X^k) - f(X^{k+1}) \\
 &= f(X^k) - f(\bar{X}) + f(\bar{X}) - f(X^{k+1}) \\
 &\geq C_1 \cdot \left\| \nabla f(X^k) - X^k X^k{}^\top \nabla f(X^k) \right\|_{\mathbb{F}}^2 + \frac{1}{8\theta + 1} \cdot \left\| \bar{X}^\top \nabla f(\bar{X}) - \nabla f(\bar{X})^\top \bar{X} \right\|_{\mathbb{F}}^2 \\
 &\geq C_1 \cdot \left\| \nabla f(X^k) - X^k \nabla f(X^k)^\top X^k \right\|_{\mathbb{F}}^2.
 \end{aligned}$$

Hence, $\{f(X^k)\}$ is monotonically decreasing. Since $\mathcal{S}_{n,p}$ is a compact set, $\{f(X^k)\}$ is bounded below so that we can conclude that $\{f(X^k)\}$ converges. \square

Next, we show the iterate subsequence convergence.

THEOREM 4.2. *Let $\{X^k\}$ be the sequence generated by Algorithm 1 initiated from a point $X^0 \in \mathcal{S}_{n,p}$. Then there exists a convergent subsequence of $\{X^k\}$. Moreover, each accumulation point X^* of $\{X^k\}$ satisfies the first-order optimality condition of (1.1) defined by Definition 2.1.*

Proof. Notice that $\{X^k\}$ is bounded due to the feasibility of each iterate X^k , hence it has a convergent subsequence. Let X^* be an accumulation point of $\{X^k\}$. Due to the feasibility of X^k , X^* satisfies the feasibility equality of condition (1.2).

Recalling inequality (4.1) in the proof of Lemma 4.1 and the boundedness of $\{f(X^k)\}$, we can conclude that

$$\lim_{k \rightarrow +\infty} \left\| \nabla f(X^k) - X^k \nabla f(X^k)^\top X^k \right\|_{\mathbb{F}} = 0,$$

which implies $\left\| \nabla f(X^*) - X^* \nabla f(X^*)^\top X^* \right\|_{\mathbb{F}} = 0$. Hence, the first two equalities of condition (1.2) are satisfied at X^* as well. Combining with Lemma 2.2, we complete the proof. \square

Lemma 4.1 and Theorem 4.2 guarantee that $f(X)$ is a constant on the accumulation point set of $\{X^k\}$. We denote this constant as f^* , and

$$(4.2) \quad \Omega_{FON}^{f^*} = \Omega_{FON} \cap \{X \mid f(X) = f^*\},$$

where Ω_{FON} is defined in Definition 2.1.

Next, we show that the distance between X^k and $\Omega_{FON}^{f^*}$ goes to zero.

COROLLARY 4.3. *Let $\{X^k\}$ be the sequence generated by Algorithm 1 initiated from a point $X^0 \in \mathcal{S}_{n,p}$, then it holds that*

$$(4.3) \quad f(X^k) \geq f^* \quad \forall k = 1, \dots,$$

and

$$(4.4) \quad \lim_{k \rightarrow \infty} \text{dist}(X^k, \Omega_{FON}^{f^*}) = 0.$$

Proof. Since $\{f(X^k)\}$ is nonincreasing, relationship (4.3) holds. Now, we assume statement (4.4) is not true. Then there exist $\delta > 0$ and a subsequence of $\{X^k\}$, denoted as $\{X^{k_j}\}$, such that

$$(4.5) \quad \text{dist} \left(X^{k_j}, \Omega_{FON}^{f^*} \right) \geq \delta.$$

Since $\{X^{k_j}\}$ is bounded, there exists a convergent subsequence of $\{X^{k_j}\}$ and any accumulation point shall also satisfy the first-order optimality condition, which contradicts (4.5). \square

5. Numerical experiments. In this section, we report the numerical performance of the algorithms based on Algorithm 1. Two types of testing problems are chosen based on Examples 1.1 and 1.2. All experiments are performed in MATLAB R2016a under a Windows 10 operating system on a Dell Optiplex 9020 personal computer with an Intel Core i7-4790 CPU at 3.6 GHz \times 2 and 8 GB of RAM.

5.1. Implementation details. In Lemmas 3.2 and 3.3, we show that GR and GP satisfy sufficient function value reduction (3.1) if the fixed step size τ is smaller than ρ^{-1} . However, to obtain a good estimation of ρ is often intractable, and ρ^{-1} can be very small, which leads to slow convergence. In practice, we can use an alternating BB step size introduced in [7], which has been already adopted in the retraction-based algorithm in [33]. More specifically, the updating rule for τ can be described as follows:

$$(5.1) \quad \tau := \begin{cases} \tau^{\text{BB1}} & \text{for odd } k, \\ \tau^{\text{BB2}} & \text{for even } k, \end{cases}$$

where

$$\begin{aligned} \tau^{\text{BB1}} &:= \frac{\langle J^{k-1}, J^{k-1} \rangle}{|\langle J^{k-1}, K^{k-1} \rangle|}, & \tau^{\text{BB2}} &:= \frac{|\langle J^{k-1}, K^{k-1} \rangle|}{\langle K^{k-1}, K^{k-1} \rangle}, \\ J^{k-1} &= X^k - X^{k-1}, & K^{k-1} &= c(X^k) - c(X^{k-1}). \end{aligned}$$

We call GR and GP with step size τ defined by (5.1) as GR-BB and GP-BB, respectively. In contrast, GR and GP with a fixed step size are called GR-F and GP-F, respectively.

CBCD will only be tested in solving quadratic problem (1.1). Therefore, in each inner iteration, the subproblem (3.13) restricted to the 2-dimensional subspace span $\{X_i, \nabla q_i(X_i)\}$ can be solved to the global optimality. However, in each outer iteration, the column updating order $\{j_1, j_2, \dots, j_p\}$ determines the way to classify different types of algorithms. Usually, there are four orders:

- (a) *cyclic type*: $j_i = i$ for $i = 1, 2, \dots, p$;
- (b) *random 1*: $j_i = \lceil p \cdot \text{rand}(1, 1) \rceil$ for $i = 1, 2, \dots, p$ (sampling with replacement);
- (c) *random 2*: $\{j_1, j_2, \dots, j_p\}$ is a random permutation of $\{1, 2, \dots, p\}$ (sampling without replacement);
- (d) *greedy type*: for $i = 1, 2, \dots, p$,

$$j_i := \arg \max_{j=1, \dots, p} \left\| \left(I_n - W^{i-1} W^{i-1 \top} \right) \nabla f_j(X_j) \right\|_2.$$

The corresponding CBCD are denoted as CBCD-C, CBCD-R1, CBCD-R2, and CBCD-G, respectively.

We have already shown that any iterate generated by any algorithm based on our new framework satisfies the symmetry and feasibility in (1.2). Hence, for the stopping criterion, we only need to check the projected gradient, $\|(I_n - XX^\top)\nabla f(X)\|_F$. More specifically, the stopping criterion can be described as follows:

$$(5.2) \quad \|(I_n - XX^\top)\nabla f(X)\|_F < \epsilon \left\| \nabla f(X^0) - X^0 \nabla f(X^0)^\top X^0 \right\|_F,$$

where $\epsilon > 0$ is a small number. The right-hand side of (5.2) is to match the scale of the initial projected gradient. On the other hand, convergence of first-order methods may slow down as the iterates approach a stationary point, so it is critical to detect the slowdown and stop properly. It is usually beneficial to have flexible stopping rules for identifying the situation that the algorithm gets trapped in a certain region. As suggested in [33], we use the following rule based on the relative error in addition:

$$(5.3) \quad \text{tol}_k^x := \frac{\|X^k - X^{k+1}\|_F}{\sqrt{n}} < \epsilon_x \quad \text{and} \quad \text{tol}_k^f := \frac{|f(X^k) - f(X^{k+1})|}{|f(X^k)| + 1} < \epsilon_f,$$

$$(5.4) \quad \begin{aligned} & \text{mean} \left(\left[\text{tol}_{k-\min\{k, T\}+1}^x, \dots, \text{tol}_k^x \right] \right) < 10\epsilon_x \\ & \text{and} \quad \text{mean} \left(\left[\text{tol}_{k-\min\{k, T\}+1}^f, \dots, \text{tol}_k^f \right] \right) < 10\epsilon_f. \end{aligned}$$

We terminate the algorithm when one of the above three criteria (5.2)–(5.4) or a maximum iteration number `MaxIter` is reached. Unless otherwise specified, the default tolerance parameters are chosen as $\epsilon = 10^{-5}$, $\epsilon_x = 10^{-6}$, $\epsilon_f = 10^{-10}$, $T = 5$, and `MaxIter` = 3000.

5.2. Testing problems. In this subsection, we introduce two types of testing problems.

The first type of testing problem is based on Example 1.1. We consider the following quadratic minimization problems with orthogonality constraints,

$$(5.5) \quad \begin{aligned} & \min_{X \in \mathbb{R}^{n \times p}} \quad \frac{1}{2} \text{tr}(X^\top AX) + \text{tr}(G^\top X) \\ & \text{s.t.} \quad \quad \quad X^\top X = I_p, \end{aligned}$$

where the matrices $A \in \mathbb{R}^{n \times n}$ and $G \in \mathbb{R}^{n \times p}$ are randomly generated by

$$(5.6) \quad A := P\Lambda P^\top,$$

$$(5.7) \quad G := \alpha \cdot QD,$$

where the matrices $P = \mathbf{qr}(\mathbf{rand}(n, n)) \in \mathbb{R}^{n \times n}$, $\tilde{Q} = \mathbf{rand}(n, p) \in \mathbb{R}^{n \times p}$, $Q \in \mathbb{R}^{n \times p}$, and $Q_i = \tilde{Q}_i / \|\tilde{Q}_i\|_2$ ($i = 1, 2, \dots, p$), and matrices $\Lambda \in \mathbb{R}^{n \times n}$ and $D \in \mathbb{R}^{p \times p}$ are diagonal matrices with

$$(5.8) \quad \Lambda_{ii} := \begin{cases} \beta^{1-i} & \text{if } \omega_i < \xi, \\ -\beta^{1-i} & \text{otherwise} \end{cases} \quad \text{for all } i = 1, 2, \dots, n,$$

$$(5.9) \quad D_{jj} := \zeta^{j-1} \quad \text{for all } j = 1, 2, \dots, p,$$

where $\omega_i \in [0, 1]$ ($i = 1, 2, \dots, n$) are randomly generated numbers. Here, $n \times p$ is the variable size; $\beta \geq 1$ is a parameter determining the decay of eigenvalues of A ; $\zeta \geq 1$ is a parameter referring to the growth rate of the column's norm of G . The parameter $\alpha > 0$ represents the scale difference between the quadratic term and the linear term. When α is large, the linear term dominates the objective. The parameter $\xi \in [0, 1]$ is

for determining the definiteness of A . Once $\xi = 1$, matrix A is positive definite, while $\xi = 0$ means A is negative definite. In contrast, unless specifically mentioned, the default setting of these parameters are $n = 3000$, $p = 60$, $\alpha = 1$, $\beta = 1.01$, $\zeta = 1.2$, $\xi = 1$. The initial point is chosen as $X^0 = \mathbf{qr}(\mathbf{rand}(n, p)) \in \mathbb{R}^{n \times p}$.

The second type of testing problem is a special case of Example 1.2. It is called Kohn–Sham total energy minimization which comes from electronic structure calculations [18]. The original Kohn–Sham equations are the Euler–Lagrange equations for the continuous total energy minimization problem. Under the planewave discretization scheme, the Kohn–Sham total energy can be transformed into a finite-dimensional approximation as follows:

$$(5.10) \quad E_{\text{total}}(X) = \text{tr} \left[X^\top \left(\frac{1}{2}L + V_{\text{ion}} \right) X \right] + \frac{1}{2} \varrho(X)^\top L^\dagger \varrho(X) + \varrho(X)^\top \epsilon_{\text{xc}}(\varrho(X)),$$

where $\varrho(X) := \text{diag}(XX^\top)$ denotes the charge density, and L is a finite-dimensional representation of the Laplacian operator in the planewave basis. The discretized local ionic potential can be represented by a diagonal matrix V_{ion} . And the matrix L^\dagger which is the discrete form of the Hartree potential corresponds to the pseudoinverse of L . The exchange correlation function ϵ_{xc} is used to model the nonclassical and quantum interaction between electrons. We aim to solve the following total energy minimization problem,

$$(5.11) \quad \min_{X^\top X = I_p} E_{\text{total}}(X).$$

It is not difficult to verify that the gradient of the energy function is $H(X)X$, where $H(X) = L/2 + V_{\text{ion}} + \text{Diag}(L^\dagger \varrho(X)) + \text{Diag}(\mu_{\text{xc}}(\varrho(X)))$ is the Kohn–Sham Hamiltonian and $\mu_{\text{xc}}(\varrho(X)) = d\epsilon_{\text{xc}}/d\varrho(X)$.

5.3. Default settings of our algorithms. In this subsection, we determine the default settings for our GR, GP, and CBCD algorithms by numerical experiments.

We first compare the performance of GR-F and GP-F with different fixed step sizes for choosing a proper value of the step size. The parameter p in the test is chosen as $0.1n$, the parameter ζ is 1.01, and the other parameters take their default values. We will compare four measurements: CPU time in seconds, total number of iterations, KKT violation, and function value variance, which is defined in the following. Suppose f_{\min} is the smallest absolute function value of those obtained by all solvers in the comparison, then for f_s the function value returned by solver s , the function value variance is defined as

$$(5.12) \quad \frac{|f_s - f_{\min}|}{1 + |f_{\min}|} + \text{eps},$$

where $\text{eps} = 2.2204e-16$ is the machine precision in MATLAB. Here, we add eps to the relative variance of function value, which is the first part of (5.12), in order to plot the variance of the function value with a logarithmic scale for the y-axis. Since both retraction-based approaches and our new framework are feasible methods, we do not report the feasibility violation $\|I - X^\top X\|_F$. The results with respect to the above four measurements are demonstrated in subfigures of Figure 3(a)–(d), respectively. We choose different τ 's ranging from $0.1\rho^{-1}$ to ρ^{-1} .

From Figure 3, we observe that $\tau = 1/3\rho$ and $1/\rho$ are the best choices for GR-F and GP-F, respectively, in this testing problem. Hence, we choose them as step sizes in the comparison with GR-BB and GP-BB.

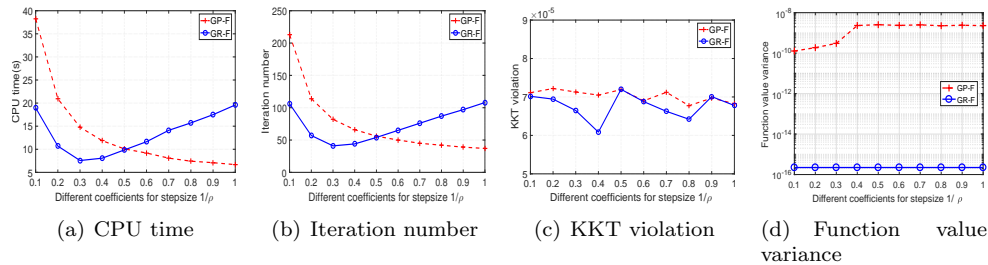


FIG. 3. Performance of GR-F and GP-F with different step sizes.

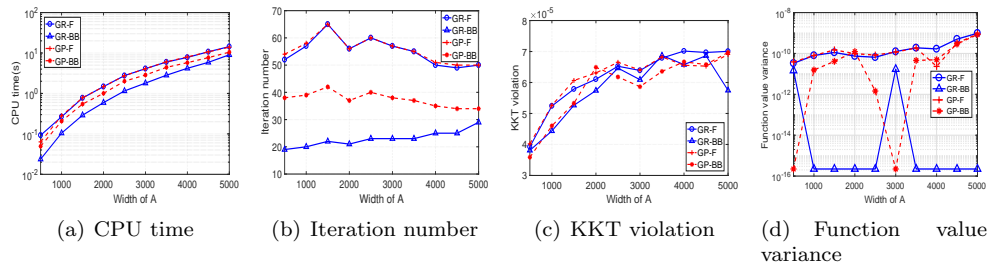


FIG. 4. Performance of gradient-based algorithms.

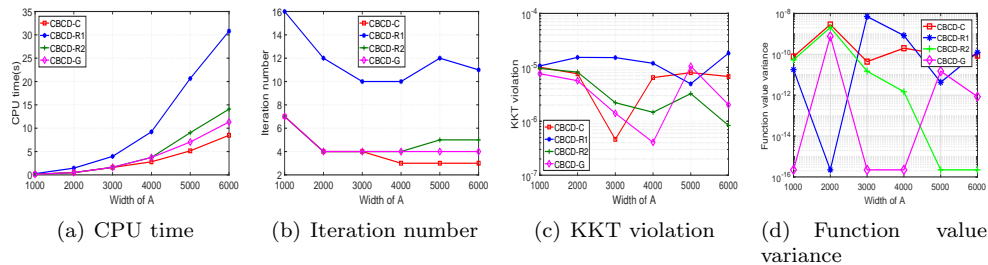


FIG. 5. Performance of CBCD with different types of choosing working index.

Next, we perform on a set of testing problems with ten randomly generated matrices with size n ranging from 500 to 5000, and the width of variable p is still $10\%n$. The parameter ζ is 1.01, and the other parameters take their default values. Numerical results of this test are illustrated in Figure 4.

From Figure 4, we notice that GR-BB and GP-BB require a much smaller number of iterations and less CPU time than GR-F and GP-F, and also achieve the same first-order stationary point with comparable KKT violation. Moreover, GR-BB outperforms GP-BB in terms of CPU time and iteration number in most cases. Thus, we choose GR-BB to represent the gradient-based class of algorithms in the following comparison in subsection 5.4.

We next compare the performance among CBCD variations corresponding to different updating orders. In this comparison, we run CBCD-C, CBCD-R1, CBCD-R2, and CBCD-RG to solve the testing problems with n ranging from 1000 to 6000, $p = 2\%n$, and other parameters taking their default values. The numerical results are presented in Figure 5.

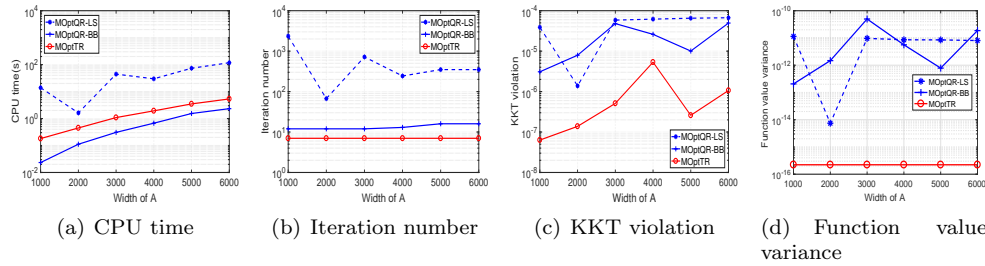


FIG. 6. Performance of MOptQR with different types of step sizes.

From Figure 5, we can see that CBCD-C, CBCD-G, and CBCD-R2 have a similar performance with respect to CPU time and iteration number, and are better than CBCD-R1. Among CBCD-C, CBCD-G, and CBCD-R2, CBCD-C performs slightly better and it is easy to implement. Therefore, we will use CBCD-C to represent the CBCD class of algorithms in the following tests.

5.4. Performance comparison on random problems. In this subsection, we compare the performance of our algorithms GR-BB and CBCD-C with two state-of-the-art solvers in solving a large variety of problems (5.5). We first choose the solver OptM⁷ based on the algorithm in [33]. For the other existing solver for comparison, we intend to choose one from MOptQR-LS (manifold QR method with line search⁸ [3]), MOptQR-BB (for fair comparison, we implement the same alternating BB step size strategy as GR-BB for the manifold QR method), and MOptTR (manifold trust-region method [3]). We compare MOptQR-LS, MOptQR-BB, and MOptTR to solve the problem (5.5) with default settings. The result is illustrated in Figure 6.

We can learn from Figure 6 that MOptQR-BB outperforms the other two methods in the testing problems and, hence, we will choose MOptQR-BB to be the other solver to compare with our algorithms. By abuse of notation, we use MOptQR to denote MOptQR-BB hereafter.

In the following experiments, we only compare the performance between GR-BB, CBCD-C, OptM, and MOptQR. We will set the same stopping criteria as introduced in subsection 5.1, and the tolerance takes its default value. We design six groups of testing problems, in each of which there is only one parameter varying with all the others fixed. More specifically, we describe the varying parameters of each group as follows:

- Number of rows of the variable, $n = 1000j$ for $j = 1, 2, 3, 4, 5, 6$.
- Number of columns of the variable, $p = 20j$ for $j = 1, 2, 3, 4, 5, 6$.
- Decay of the eigenvalues of A , $\beta = 1.01 + 0.03j$ for $j = 0, 1, 2, 3, 4, 5, 6, 7, 8$.
- Difference between column norms of G , $\zeta = 1.01 + 0.03j$ for $j = 0, 1, 2, 3, 4, 5, 6, 7, 8$.
- The dominance of the linear term, $\alpha = 10^{-2}, 10^{-1}, 1, 10, 10^2$.
- The definiteness of A , $\xi = 0.2(j - 1)$ for $j = 1, 2, 3, 4, 5, 6$.

The meanings of these parameters refer to equalities (5.6)–(5.9). The linear eigenvalue problem, i.e., problem (5.5) with $\alpha = 0$, is not in our testing problems. The step sizes in our new proposed GR and GP need to be tuned for different problems, and hence the algorithms do not become practically useful. On the other hand, there are a bunch

⁷ Available from <http://optman.blogs.rice.edu>

⁸ Available from <http://www.manopt.org>

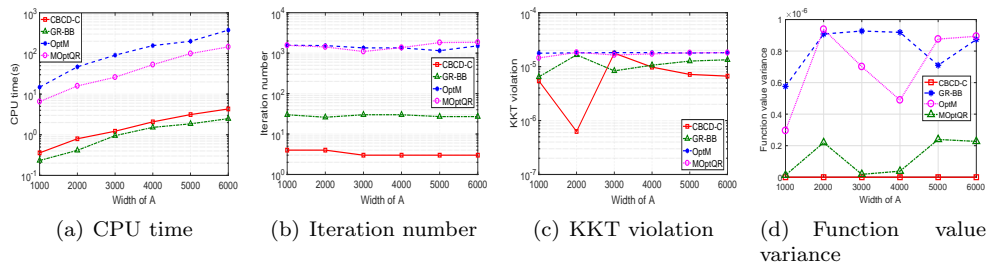


FIG. 7. Comparison with varying matrix dimension n .

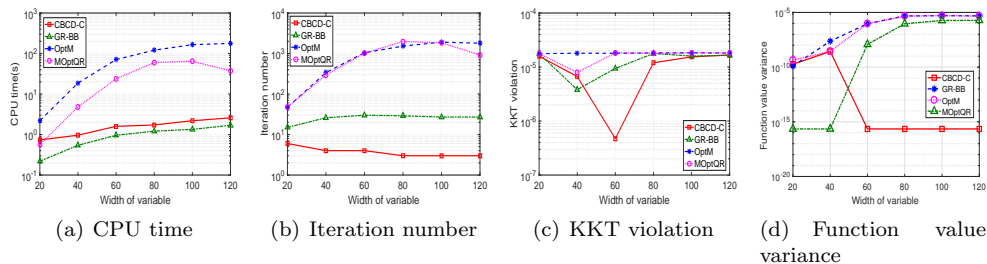


FIG. 8. Comparison with varying width of variable p .

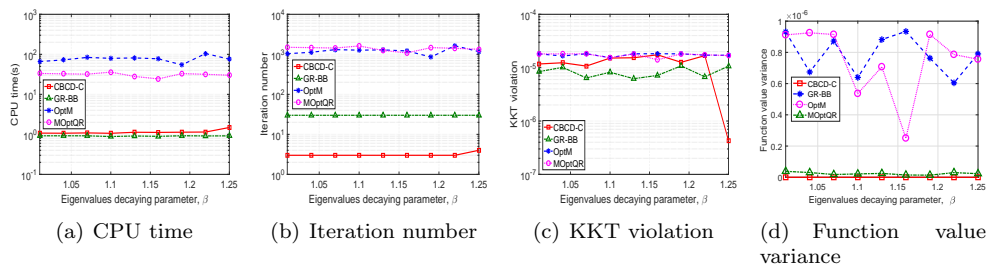


FIG. 9. Comparison with varying decay parameter β .

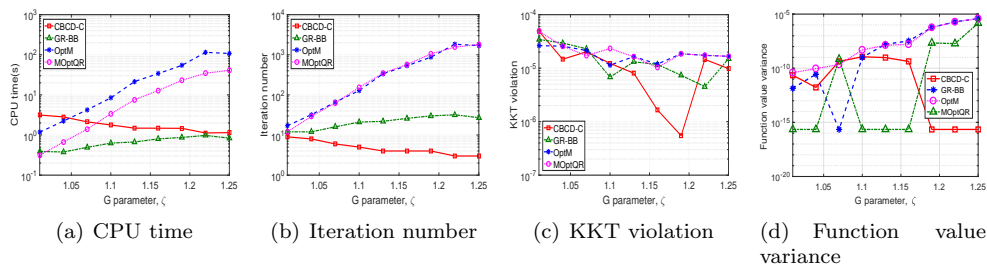


FIG. 10. Comparison with varying G parameter ζ .

of efficient solvers particularly designed for the linear eigenvalue problem which can hardly be beaten by general solvers for optimization problems with orthogonality constraints. The numerical results of the above six groups of testing problems are given in Figures 7 to 12, respectively.

From the above figures, we have the following observations. All solvers reach the same function value from the same initial point. They achieve comparable KKT

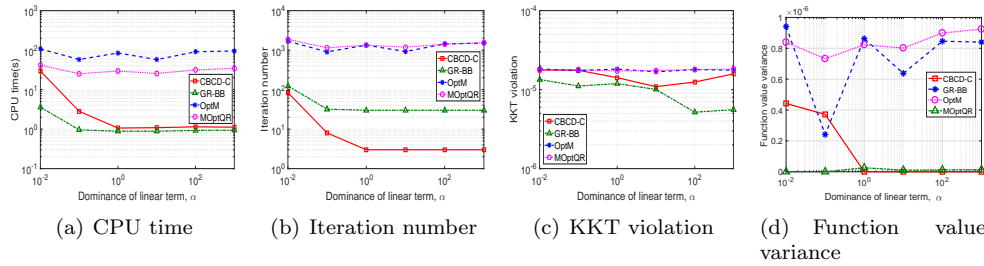


FIG. 11. Comparison with varying dominance of linear term α .

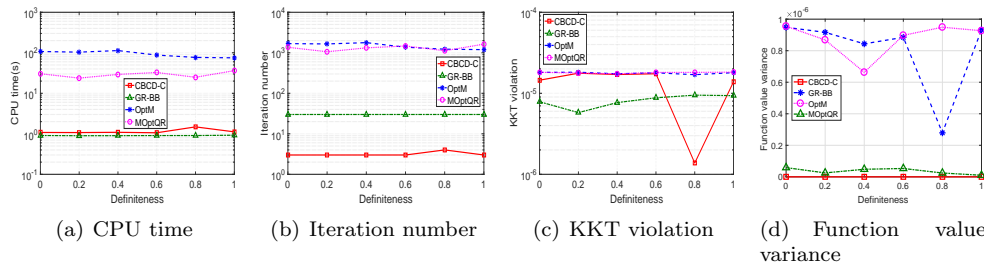


FIG. 12. Comparison with varying nonnegativity of A, ξ .

violation with magnitude around 10^{-5} . Moreover, GR-BB and CBCD-C usually have lower KKT violation than the other two in most experiments. Among the four algorithms, CBCD-C has the lowest iteration number in all the tests, while GR-BB has the least CPU time. Except for very extreme cases, CBCD-C performs the second best in terms of CPU-time.

Finally, we select all the testing problems with options in bold in the above description, and put them into a performance profile experiment [9]. There are altogether $6 \times 6 \times 3 \times 3 \times 3 = 2916$ randomly generated problems. The performance profile can eliminate the influence of a small number of difficult problems and the sensitivity of results associated with the different criteria, and also provide a way to visualize the expected performance difference among many solvers. We describe the key parameters of such a test as the following. For problem m and solver s , we denote $t_{m,s}$ to represent the CPU time or iteration number. Performance ratio is defined as $r_{m,s} := t_{m,s} / \min_s \{t_{m,s}\}$. If solver s fails to solve problem m , the ratio $r_{m,s}$ will be set to infinity or some sufficiently large number. Finally, the overall performance of solver s is defined by

$$\pi_s(\omega) := \frac{\text{number of problems where } r_{m,s} \leq \omega}{\text{total number of problems}}.$$

It means the percentage of testing problems that can be solved in $\omega \min_s t_{m,s}$ seconds (or iterations). Of course, the closer π_s is to 1, the better performance solver s has. The performance profile results with respect to CPU time and iteration number are given in Figure 13.

We observe that GR-BB performs best and CBCD-C performs the second best among all four algorithms in solving these 2916 testing problems in CPU time, and meanwhile CBCD-C requires the smallest iteration number. In addition, we also provide the average KKT violation and feasibility over these 2916 random problems in Table 2.

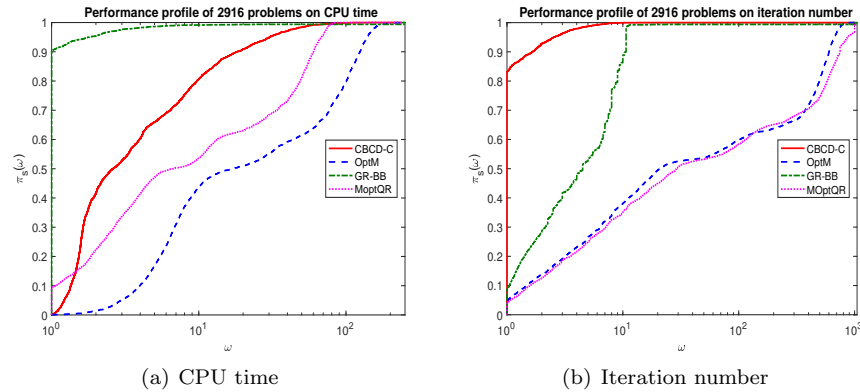


FIG. 13. Performance profile.

TABLE 2
Average KKT, feasibility violation, and function value.

	CBCD-C	OptM	GR-BB	MOptQR
KKT violation	1.6075e-05	2.1730e-05	1.9501e-05	2.5072e-05
Function value variance	6.5780e-06	8.1754e-06	3.0417e-06	7.9584e-06

Table 2 shows all solvers achieve a comparable average KKT violation, feasibility, and function value variance. Here, the function value variance of solver s in solving problem m is in the same manner as (5.12). More specifically,

$$z_{m,s} := \frac{|f_{m,s} - \min_s \{f_{m,s}\}|}{1 + |\min_s \{f_{m,s}\}|}.$$

5.5. Global property of CBCD. An interesting observation of all the experiments introduced above is that all solvers reach the same function value when they converge from a randomly generated initial guess, although our problem (1.1) is non-convex. Therefore we design a new experiment as the following. We construct the following problem

$$\begin{aligned} \min_{X \in \mathbb{R}^{3 \times 2}} & \frac{1}{2} \text{tr} \left((X - X^*)^\top A (X - X^*) \right) \\ \text{s. t.} & \quad X^\top X = I_2, \end{aligned}$$

where

$$A = \begin{bmatrix} 13/2 & 2 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

For this special problem, we can verify that

$$X^* = \begin{bmatrix} 3/5 & 0 \\ 4/5 & 0 \\ 0 & 1 \end{bmatrix}$$

TABLE 3
Test results with initial points near X^I .

Testing methods	X^*	X^I	X^{II}	X^{III}	Success rate
CBCD-C	1000	0	0	0	100%
GR-BB	0	1000	0	0	0%
OptM	0	1000	0	0	0%
MOptQR	0	1000	0	0	0%

TABLE 4
Test results with initial points near X^{II} .

Testing methods	X^*	X^I	X^{II}	X^{III}	Success rate
CBCD-C	1000	0	0	0	100%
GR-BB	0	1000	0	0	0%
OptM	729	0	271	0	72.9%
MOptQR	78	0	922	0	7.8%

TABLE 5
Test results with initial points near X^{III} .

Testing methods	X^*	X^I	X^{II}	X^{III}	Success rate
CBCD-C	1000	0	0	0	100%
GR-BB	1000	0	0	0	100%
OptM	338	28	634	0	33.8%
MOptQR	5	5	990	0	0.5%

is the unique global minimizer, while

$$X^I = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad X^{II} = \begin{bmatrix} 3/5 & 0 \\ 4/5 & 0 \\ 0 & -1 \end{bmatrix}, \quad X^{III} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & -1 \end{bmatrix}$$

are the first-order stationary points. X^I is a local minimizer, while the other two are saddle points. Then we set the initial guess from the neighborhoods of the three stationary points, and run GR-BB, CBCD-C, OptM, and MOptQR to see the different performances. More specifically,

$$\begin{aligned} X^0 &:= \mathcal{P}_{S_{n,p}}(X^i + \mu \cdot \mathbf{randn}(3, 2)) \quad \text{for } i = I, II, III, \\ X^0 &:= \mathcal{P}_{S_{n,p}}(\mathbf{randn}(3, 2)), \end{aligned}$$

where $\mu > 0$ controls the distance between X^0 and X^i for $i = I, II, III$. We set $\mu = 10^{-4}$ and compare all solvers with these four types of initial points. With repeating each test 1000 times, we record the number of each solution for every solver, and the success rates are presented in Tables 3, 4, 5, and 6.

It can be observed from the above tables that the four algorithms are not necessarily convergent to same stationary points. In our tests, CBCD-C can always find the global minimizer. We are not sure whether it is a coincidence or CBCD-C has the

TABLE 6
Test results with random initial guesses.

Testing methods	X^*	X^I	X^{II}	X^{III}	Success rate
CBCD-C	1000	0	0	0	100%
GR-BB	656	344	0	0	65.6%
OptM	864	136	0	0	86.4%
MOptQR	774	226	0	0	77.4%

nice property of converging to a global minimizer with great probability. The random initialization does increase the chance to find a global minimizer for the other three algorithms.

5.6. Kohn–Sham total energy minimization. In the end of this section, we compare GR-BB with the state-of-the-art solvers in solving Kohn–Sham total energy minimization. Our test is based on the best MATLAB platform, to the best of our knowledge, for electronic structure calculations, KSSOLV [34]. KSSOLV has a friendly interface and allows researchers to investigate their own algorithms easily for different steps in electronic structure calculations. Currently, the most widely used algorithm for (5.11) is the self-consistent field (SCF) iteration, which is provided in KSSOLV. This is an iterative method for solving the nonlinear eigenvalue problem (KKT system of (5.11) briefly). Other methods focusing on discretized Kohn–Sham total energy minimization including direct constrained minimization [35] and its improved version, trust-region direct constrained minimization (TRDCM) [36] are also integrated in KSSOLV. TRDCM combines the trust-region and the subspace strategies to this special optimization problem with orthogonality constraints, and its trust-region subproblems restricted to a subspace are solved by SCF. GR-BB and MOptQR are selected in this comparison as general solvers for optimization problems with orthogonality constraints.

We select 18 testing problems with respect to different molecules, which are assembled in KSSOLV. We run methods SCF and TRDCM with $\epsilon = 10^{-5}$, $\text{MaxIter} = 200$, and other parameters taking their default values, while GR-BB and MOptQR improve their stopping criteria with $\epsilon = 10^{-5}$, $\epsilon_x = 10^{-9}$, $\epsilon_f = 10^{-13}$, $\text{MaxIter} = 1000$ to get a comparable solution with other methods. It is worth mentioning that here the symmetry of (1.2) is already achieved, since the total energy function is homogeneous and hence without a linear term. The stopping rule is set as $\|(I_n - XX^\top)H(X)X\|_F < \epsilon$. For all of the testing algorithms, we set the same initial guess X^0 by using the function “getX0,” which is provided by KSSOLV. The numerical results are illustrated in Tables 7 and 8.

Here, “ E_{tot} ,” “KKT violation,” “Iteration,” and “CPU time(s)” represent the total energy function value, the value of $\|(I_n - XX^\top)H(X)X\|_F$, the number of iterations and the total running time in seconds, respectively. From the tables, we observe that GR-BB outperforms the other algorithms, even the heuristic ones, in most cases, and it obtains a comparable total energy function value and a lower KKT violation. In particular, in the large size problem “ctube661,” GR-BB achieves the same total energy function value and same magnitude KKT violation, but requires much less CPU time than the others.

TABLE 7
The results in total energy minimization.

Solver	E_{tot}	KKT violation	Iteration	CPU time(s)
al, $n = 16879$, $p = 12$				
SCF	-1.5799906179e+01	8.68e-03	200	2509.48
TRDCM	-1.5803817595e+01	8.15e-06	184	1595.83
MOptQR	-1.5802118775e+01	8.42e-03	1000	2017.61
GR-BB	-1.5802922328e+01	2.05e-03	1000	2070.80
alanine, $n = 12671$, $p = 18$				
SCF	-6.1161921213e+01	9.70e-07	15	204.20
TRDCM	-6.1161921213e+01	5.91e-06	16	147.84
MOptQR	-6.1161921213e+01	8.14e-06	65	142.70
GR-BB	-6.1161921212e+01	9.78e-06	63	142.36
benzene, $n = 8407$, $p = 15$				
SCF	-3.7225751363e+01	7.85e-07	12	85.52
TRDCM	-3.7225751363e+01	7.33e-06	14	71.13
MOptQR	-3.7225751363e+01	8.38e-06	127	154.06
GR-BB	-3.7225751362e+01	9.69e-06	50	60.38
c2h6, $n = 2103$, $p = 7$				
SCF	-1.4420491322e+01	1.12e-06	11	10.09
TRDCM	-1.4420491322e+01	5.00e-06	12	7.61
MOptQR	-1.4420491322e+01	5.56e-06	49	8.53
GR-BB	-1.4420491321e+01	9.84e-06	43	7.58
c12h26, $n = 5709$, $p = 37$				
SCF	-8.1536091936e+01	1.52e-06	16	288.09
TRDCM	-8.1536091937e+01	9.48e-06	15	171.38
MOptQR	-8.1536091935e+01	9.51e-06	442	1296.05
GR-BB	-8.1536091936e+01	8.85e-06	50	157.02
co2, $n = 2103$, $p = 8$				
SCF	-3.5124395801e+01	1.50e-06	11	11.92
TRDCM	-3.5124395801e+01	7.63e-06	13	8.72
MOptQR	-3.5124395800e+01	9.03e-06	39	7.53
GR-BB	-3.5124395801e+01	6.94e-06	39	7.52
ctube661, $n = 12599$, $p = 48$				
SCF	-1.3463843176e+02	2.80e-06	13	532.25
TRDCM	-1.3463843176e+02	5.77e-06	22	787.58
MOptQR	-1.3463843177e+02	5.06e-06	533	3817.95
GR-BB	-1.3463843176e+02	9.27e-06	68	493.53
glutamine, $n = 16517$, $p = 29$				
SCF	-9.1839425243e + 01	2.88e-06	17	616.73
TRDCM	-9.1839425244e + 01	8.49e-06	15	479.34
MOptQR	-9.1839425243e + 01	7.26e-06	87	570.86
GR-BB	-9.1839425243e + 01	9.76e-06	75	499.92
graphene16, $n = 3071$, $p = 37$				
SCF	-9.3873673630e+01	5.28e-03	200	2008.61
TRDCM	-9.4046217545e+01	6.12e-06	43	313.88
MOptQR	-9.4046217540e+01	9.56e-06	693	1110.39
GR-BB	-9.4046217543+01	8.35e-06	321	513.45

TABLE 8
The results in total energy minimization.

Solver	E_{tot}	KKT violation	Iteration	CPU time(s)
graphene30, $n = 12279$, $p = 67$				
SCF	-1.7358503892e+02	3.18e-03	200	15344.80
TRDCM	-1.7359510505e+02	9.77e-06	62	3768.22
MOptQR	-1.6908746446e+02	3.87e+00	1000	11930.80
GR-BB	-1.7359510453e+02	1.97e-04	1000	12027.63
h2o, $n = 2103$, $p = 4$				
SCF	-1.6440507246e+01	7.78e-07	9	5.48
TRDCM	-1.6440507246e+01	8.22e-06	11	4.55
MOptQR	-1.6440507245e+01	8.43e-06	44	5.13
GR-BB	-1.6440507245e+01	9.89e-06	42	4.53
hnco, $n = 2103$, $p = 8$				
SCF	-1.6440507246e+01	7.08e-07	9	5.52
TRDCM	-1.6440507246e+01	9.64e-06	11	4.27
MOptQR	-1.6440507245e+01	9.20e-06	82	10.41
GR-BB	-1.6440507246e+01	8.64e-06	40	5.11
nic, $n = 251$, $p = 7$				
SCF	-2.3543529955e+01	1.10e-06	12	3.13
TRDCM	-2.3543529955e+01	9.33e-06	49	5.50
MOptQR	-2.3543529955e+01	8.26e-06	100	2.84
GR-BB	-2.3543529955e+01	9.56e-06	39	0.88
pentacene, $n = 44791$, $p = 51$				
SCF	-1.3189029495e+02	9.83e-07	15	2448.72
TRDCM	-1.3189029495e+02	9.67e-06	23	2706.14
MOptQR	-1.3189029495e+02	7.02e-06	355	9145.66
GR-BB	-1.3189029495e+02	9.54e-06	100	2606.81
ptnio, $n = 4609$, $p = 43$				
SCF	-2.2678884273e+02	8.25e-07	70	1079.14
TRDCM	-2.2678882962e+02	2.93e-04	200	1957.89
MOptQR	-2.2678884235e+02	2.33e-05	1000	2281.22
GR-BB	-2.2678884272e+02	9.68e-06	512	1159.91
qdot, $n = 2103$, $p = 8$				
SCF	2.7702342351e+01	3.91e-02	200	175.16
TRDCM	2.7699896368e+01	2.72e-03	200	104.80
MOptQR	3.1736592205e+01	3.96e+00	1000	135.88
GR-BB	2.7700280932e+01	7.90e-04	1000	138.98
si2h4, $n = 2103$, $p = 6$				
SCF	-6.3009750460e+00	4.98e-07	13	12.42
TRDCM	-6.3009750459e+00	7.39e-06	16	9.09
MOptQR	-6.3009750460e+00	3.83e-06	75	11.67
GR-BB	-6.3009750457e+00	6.58e-06	58	8.97
sih4, $n = 2103$, $p = 4$				
SCF	-6.1769279851e+00	8.83e-07	10	5.80
TRDCM	-6.1769279850e+00	9.59e-06	10	4.50
MOptQR	-6.1769279851e+00	3.76e-06	42	5.14
GR-BB	-6.1769279850e+00	9.03e-06	36	4.41

6. Conclusion. In this paper, we propose a new first-order algorithmic framework, Algorithm 1, for optimization problems with orthogonality constraints (1.1). This algorithmic framework consists of two steps. In the first step, we choose a function value reduction approach to reduce the function value and keep the feasibility at the same time, and hence the calculation related to the tangent space of the Stiefel manifold can be waived. Second, a correction step is employed to guarantee that any accumulation point of the iterates is a first-order stationary point. Moreover, for some special cases, the correction step can be waived. We introduce two classes of approaches. The difference of them is in the first step. We first put forward a gradient-based scheme, whose global convergence can be guaranteed by a fixed step size and hence line search is no longer needed. We recommend two particular algorithms, GR and GP, from this class. The second class of algorithms is called CBCD, in which the columnwise block coordinate update is conducted in a Gauss–Seidel manner. We also propose novel ideas to solve the columnwise subproblem efficiently and guarantee the global convergence. Preliminary experiments on two large classes of testing problems including Kohn–Sham total energy minimization arising from electronic structure calculations illustrate that our new algorithms have great potential.

However, how to design second-order methods to further enhance the performance and obtain local minimizers is still under investigation. Global optimality under some random assumptions is an attractive topic for future work. How to design Jacobian-type CBCD methods is very important for the parallelization, as low scalability is an inevitable bottleneck of existing approaches for solving optimization problems with orthogonality constraints.

Acknowledgments. The authors would like to thank Zhaosong Lu, Ting Kei Pong, Zaiwen Wen, and Zaikun Zhang for the insightful discussions.

REFERENCES

- [1] T. E. ABRUDAN, J. ERIKSSON, AND V. KOIVUNEN, *Steepest descent algorithms for optimization under unitary matrix constraint*, IEEE Trans. Signal Process., 56 (2008), pp. 1134–1147.
- [2] T. E. ABRUDAN, J. ERIKSSON, AND V. KOIVUNEN, *Conjugate gradient algorithm for optimization under unitary matrix constraint*, Signal Process., 89 (2009), pp. 1704–1714.
- [3] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, Princeton, NJ, 2009.
- [4] P.-A. ABSIL AND J. MALICK, *Projection-like retractions on matrix manifolds*, SIAM J. Optim., 22 (2012), pp. 135–158.
- [5] J. BARZILAI AND J. M. BORWEIN, *Two-point step size gradient methods*, IMA J. Numer. Anal., 8 (1988), pp. 141–148.
- [6] A. CABOUSSAT, R. GLOWINSKI, AND V. PONS, *An augmented Lagrangian approach to the numerical solution of a non-smooth eigenvalue problem*, J. Numer. Math., 17 (2009), pp. 3–26.
- [7] Y.-H. DAI AND R. FLETCHER, *Projected Barzilai-Borwein methods for large-scale box-constrained quadratic programming*, Numer. Math., 100 (2005), pp. 21–47.
- [8] A. D’ASPREMONT, L. EL GHAOU, M. I. JORDAN, AND G. R. G. LANCKRIET, *A direct formulation for sparse PCA using semidefinite programming*, SIAM Rev., 49 (2007), pp. 434–448.
- [9] E. D. DOLAN AND J. J. MORÉ, *Benchmarking optimization software with performance profiles*, Math. Program., 91 (2002), pp. 201–213.
- [10] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*, SIAM J. Matrix Anal. Appl., 20 (1998), pp. 303–353.
- [11] L. ELDÉN AND H. PARK, *A Procrustes problem on the Stiefel manifold*, Numer. Math., 82 (1999), pp. 599–619.
- [12] C. FRAIKIN, Y. NESTEROV, AND P. VAN DOOREN, *A gradient-type algorithm optimizing the coupling between matrices*, Linear Algebra Appl., 429 (2008), pp. 1229–1242.
- [13] D. GOLDFARB, Z. WEN, AND W. YIN, *A curvilinear search method for p-harmonic flows on spheres*, SIAM J. Imaging Sci., 2 (2009), pp. 84–109.

- [14] I. GRUBIŠIĆ AND R. PIETERSZ, *Efficient rank reduction of correlation matrices*, Linear Algebra Appl., 422 (2007), pp. 629–653.
- [15] W. HUANG, P.-A. ABSIL AND K. A. GALLIVAN, *A Riemannian BFGS method for nonconvex optimization problems*, in Numerical Mathematics and Advanced Applications, Lecture Notes Comput. Sci. Eng. 112, Springer, Cham, Switzerland, 2016, pp. 627–634, https://doi.org/10.1007/978-3-319-39929-4_60.
- [16] W. HUANG, K. A. GALLIVAN AND P.-A. ABSIL, *A Broyden class of quasi-Newton methods for Riemannian optimization*, SIAM J. Optim., 25 (2015), pp. 1660–1685.
- [17] B. JIANG AND Y.-H. DAI, *A framework of constraint preserving update schemes for optimization on Stiefel manifold*, Math. Program., 153 (2015), pp. 535–575.
- [18] W. KOHN AND L. J. SHAM, *Self-consistent equations including exchange and correlation effects*, Phys. Rev. (2), 140 (1965), pp. A1133–1138.
- [19] R. LAI AND S. OSHER, *A splitting method for orthogonality constrained problems*, J. Sci. Comput., 58 (2014), pp. 431–449.
- [20] X. LIU, X. WANG, Z. WEN, AND Y. YUAN, *On the convergence of the self-consistent field iteration in Kohn–Sham density functional theory*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 546–558.
- [21] X. LIU, Z. WEN, AND Y. ZHANG, *Limited memory block Krylov subspace optimization for computing dominant singular value decompositions*, SIAM J. Sci. Comput., 35 (2013), pp. A1641–A1668.
- [22] X. LIU, Z. WEN, AND Y. ZHANG, *An efficient Gauss–Newton algorithm for symmetric low-rank product matrix approximations*, SIAM J. Optim., 25 (2015), pp. 1571–1608.
- [23] J. H. MANTON, *Optimization algorithms exploiting unitary constraints*, IEEE Trans. Signal Process., 50 (2002), pp. 635–650.
- [24] Y. NISHIMORI AND S. AKAHO, *Learning algorithms utilizing quasi-geodesic flows on the Stiefel manifold*, Neurocomputing, 67 (2005), pp. 106–135.
- [25] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, Springer, Berlin, 2006.
- [26] B. SAVAS AND L.-H. LIM, *Quasi-Newton methods on Grassmannians and multilinear approximations of tensors*, SIAM J. Sci. Comput., 32 (2010), pp. 3352–3393.
- [27] P. H. SCHÖNEMANN, *A generalized solution of the orthogonal Procrustes problem*, Psychometrika, 31 (1966), pp. 1–10.
- [28] E. STIEFEL, *Richtungsfelder und fernparallelismus in n -dimensionalen mannigfaltigkeiten*, Comment. Math. Helv., 8 (1935), pp. 305–353.
- [29] W. SUN AND Y.-X. YUAN, *Optimization Theory and Methods: Nonlinear Programming*, Vol. 1, Springer, New York, 2006.
- [30] L. N. TREFETHEN AND D. BAU, *Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- [31] M. ULBRICH, Z. WEN, C. YANG, D. KLÖCKNER, AND Z. LU, *A proximal gradient method for ensemble density functional theory*, SIAM J. Sci. Comput., 37 (2015), pp. A1975–A2002.
- [32] Z. WEN, C. YANG, X. LIU, AND Y. ZHANG, *Trace-penalty minimization for large-scale eigenspace computation*, J. Sci. Comput., 66 (2016), pp. 1175–1203.
- [33] Z. WEN AND W. YIN, *A feasible method for optimization with orthogonality constraints*, Math. Program., 142 (2013), pp. 397–434.
- [34] C. YANG, J. C. MEZA, B. LEE, AND L.-W. WANG, *KSSOLV—a MATLAB toolbox for solving the Kohn–Sham equations*, ACM Trans. Math. Software, 36 (2009), 10.
- [35] C. YANG, J. C. MEZA, AND L.-W. WANG, *A constrained optimization algorithm for total energy minimization in electronic structure calculations*, J. Comput. Phys., 217 (2006), pp. 709–721.
- [36] C. YANG, J. C. MEZA, AND L.-W. WANG, *A trust region direct constrained minimization algorithm for the Kohn–Sham equation*, SIAM J. Sci. Comput., 29 (2007), pp. 1854–1875.
- [37] H. ZHANG AND W. W. HAGER, *A nonmonotone line search technique and its application to unconstrained optimization*, SIAM J. Optim., 14 (2004), pp. 1043–1056.
- [38] X. ZHANG, J. ZHU, Z. WEN, AND A. ZHOU, *Gradient type optimization methods for electronic structure calculations*, SIAM J. Sci. Comput., 36 (2014), pp. C265–C289.
- [39] H. ZOU, T. HASTIE, AND R. TIBSHIRANI, *Sparse principal component analysis*, J. Comput. Graph. Statist., 15 (2006), pp. 265–286.