

Subspace Methods for Nonlinear Optimization

Xin Liu*
Zaiwen Wen†
Ya-xiang Yuan‡

Abstract. Subspace techniques such as Krylov subspace methods have been well known and extensively used in numerical linear algebra. They are also ubiquitous and becoming indispensable tools in nonlinear optimization due to their ability to handle large scale problems. There are generally two types of principals: i) the decision variable is updated in a lower dimensional subspace; ii) the objective function or constraints are approximated in a certain smaller functional subspace. The key ingredients are the constructions of suitable subspaces and subproblems according to the specific structures of the variables and functions such that either the exact or inexact solutions of subproblems are readily available and the corresponding computational cost is significantly reduced. A few relevant techniques include but not limited to direct combinations, block coordinate descent, active sets, limited-memory, Anderson acceleration, subspace correction, sampling and sketching. This paper gives a comprehensive survey on the subspace methods and their recipes in unconstrained and constrained optimization, nonlinear least squares problem, sparse and low rank optimization, linear and nonlinear eigenvalue computation, semidefinite programming, stochastic optimization and etc. In order to provide helpful guidelines, we emphasize on high level concepts for the development and implementation of practical algorithms from the subspace framework.

Key words. nonlinear optimization, subspace techniques, block coordinate descent, active sets, limited memory, Anderson acceleration, subspace correction, subsampling, sketching

AMS subject classification. 65K05, 90C30

1	Introduction	3
1.1	Overview of Subspace Techniques	4
1.2	Notation	5
1.3	Organization	5
2	General Unconstrained Optimization	5
2.1	The Line Search Methods	5
2.1.1	The Nonlinear Conjugate Gradient (CG) Method	6
2.1.2	Nesterov’s Accelerated Gradient Method	6
2.1.3	The Heavy-ball Method	7
2.1.4	A Search Direction Correction (SDC) Method	7
2.1.5	Quasi-Newton Methods	8
2.1.6	Acceleration Techniques	8

*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 (liuxin@lsec.cc.ac.cn). 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 QYZDJ-SSW-SYS010, CAS.

†Beijing International Center for Mathematical Research, Peking University, China (wenzw@pku.edu.cn). Research supported in part by NSFC grant 11831002, and by Beijing Academy of Artificial Intelligence (BAAI).

‡State Key Laboratory of Scientific and Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, China (yyx@lsec.cc.ac.cn). Research supported in part by NSFC grants 11331012 and 11688101.

36	2.1.7	Search Direction From Minimization Subproblems	9
37	2.1.8	Subspace By Coordinate Directions	10
38	2.2	Trust Region Methods	11
39	3	Nonlinear Equations and Nonlinear Least Squares Problem	13
40	3.1	General Subspace Methods	13
41	3.2	Subspace by Subsampling/Sketching	13
42	3.3	Partition of Variables	15
43	3.4	τ -steepest Descent Coordinate Subspace	15
44	4	Stochastic Optimization	15
45	4.1	Stochastic First-order Methods	16
46	4.2	Stochastic Second-Order method	17
47	5	Sparse Optimization	18
48	5.1	Basis Pursuit	18
49	5.2	Active Set Methods	18
50	6	The Domain Decomposition Methods	20
51	6.1	A Two-level Subspace Method	20
52	6.2	The Subspace Correction Method	21
53	6.3	Parallel Line Search Subspace Correction Method	21
54	7	General Constrained Optimization	22
55	7.1	Direct Subspace Techniques	23
56	7.2	Second Order Correction Steps	23
57	7.3	The Celis-Dennis-Tapia (CDT) Subproblem	24
58	7.4	Simple Bound-constrained Problems	25
59	8	Eigenvalue Computation	26
60	8.1	Classic Subspace Iteration	26
61	8.2	Polynomial Filtering	27
62	8.3	Limited Memory Methods	28
63	8.4	Augmented Rayleigh-Ritz Method	28
64	8.5	Singular Value Decomposition	29
65	8.6	Randomized SVD	31
66	8.7	Truncated Subspace Method for Tensor Train	31
67	9	Optimization with Orthogonality Constraints	34
68	9.1	Regularized Newton Type Approaches	34
69	9.2	A Structured Quasi-Newton Update with Nyström Approximation	35
70	9.3	Electronic Structure Calculations	36
71	9.3.1	The Mathematical Models	37
72	9.3.2	The Self-Consistent Field (SCF) Iteration	38
73	9.3.3	Subspace Methods For HF using Nyström Approximation	39
74	9.3.4	A Regularized Newton Type Method	39
75	9.3.5	Subspace Refinement for KSDFT	40
76	10	Semidefinite Programming (SDP)	40
77	10.1	The Maxcut SDP	40
78	10.1.1	Examples: Phase Retrieval	41

79	10.2 Community Detection	42
80	11 Low Rank Matrix Optimization	44
81	11.1 Low Rank Structure of First-order Methods	44
82	11.2 A Polynomial-filtered Subspace Method	45
83	11.3 The Polynomial-filtered Proximal Gradient Method	46
84	11.3.1 Examples: Maximal Eigenvalue and Matrix Completion	46
85	11.4 The Polynomial-filtered ADMM Method	47
86	11.4.1 Examples: 2-RDM and Cryo-EM	48
87	12 Conclusion	48
88	References	49
89		
90		

91 **1. Introduction.** Large scale optimization problems appear in a wide variety of scien-
92 tific and engineering domains. In this paper, we consider a general optimization problem

93 (1.1)
$$\min_x f(x), \text{ s. t. } x \in \mathcal{X},$$

94 where x is the decision variable, $f(x)$ is the objective function and \mathcal{X} is the feasible set. Effi-
95 cient numerical optimization algorithms have been extensively developed for (1.1) with vari-
96 ous types of objective functions and constraints [111, 88]. With the rapidly increasing prob-
97 lem scales, subspace techniques are ubiquitous and becoming indispensable tools in nonlinear
98 optimization due to their ability to handle large scale problems. For example, the Krylov sub-
99 space methods developed in the numerical linear algebraic community have been widely used
100 for the linear least squares problem and linear eigenvalue problem. The characteristics of the
101 subspaces are clear in many popular optimization algorithms such as the linear and nonlin-
102 ear conjugate gradient methods, Nesterov’s accelerated gradient method, the Quasi-Newton
103 methods and the block coordinate decent (BCD) method. The subspace correction method
104 for convex optimization can be viewed as generalizations of multigrid and domain decompo-
105 sition methods. The Anderson acceleration or the direct inversion of iterative subspace (DIIS)
106 methods have been successful in computational quantum physics and chemistry. The stochas-
107 tic gradient type methods usually take a mini-batch from a large collection samples so that
108 the computational cost of each inner iteration is small. The sketching techniques formulate a
109 reduced problem by a multiplication with random matrices with certain properties.

110 The purpose of this paper is to provide a review of the subspace methods for nonlinear
111 optimization, for their further improvement and for their future usage in even more diverse
112 and emerging fields. The subspaces techniques for (1.1) are generally divided into two cat-
113 egories. The first type is to update the decision variable in a lower dimensional subspace,
114 while the second type is to construct approximations of the objective function or constraints
115 in a certain smaller subspace of functions. Usually, there are three key steps.

- 116 • Identify a suitable subspace either for the decision variables or the functions.
- 117 • Construct a proper subproblem by various restrictions or approximations.
- 118 • Find either an exact or inexact solution of subproblems.

119 These steps are often mixed together using the specific structures of the problems case by
120 case. The essence is how to reduce the corresponding computational cost significantly.
121 During the practice in unconstrained and constrained optimization, nonlinear least squares
122 problem, sparse and low rank optimization, linear and nonlinear eigenvalue computation,
123 semidefinite programming, stochastic optimization, manifold optimization, phase retrieval,

124 variational minimization and etc, the collection of subspaces techniques is growing ever rich.
125 It includes but not limited to direct combinations, BCD, active sets, limited-memory, Ander-
126 son acceleration, subspace correction, sampling and sketching. We aim to provide helpful
127 guidelines for the development and implementation of practical algorithms using the sub-
128 space framework. Hence, only high level algorithmic ideas rather than theoretical properties
129 of the subspace techniques are covered in various contexts.

130 **1.1. Overview of Subspace Techniques.** We next summarize the concepts and
131 contexts of a few main subspace techniques.

132 **Direct Combinations.** It is a common practice to update the decision variables using a
133 combination of a few known directions which forms a subspace. The linear and nonlinear
134 conjugate gradient methods [111, 88], the Nesterov’s accelerated gradient method [84, 85],
135 the Heavy-ball method [90], the search direction correction method [126] and the momentum
136 method [47] take a linear combination of the gradient and the previous search direction. The
137 main difference is reflected in the choices of the coefficients according to different explicit
138 formulas.

139 **BCD.** The variables in many problems can be split naturally into a few blocks whose sub-
140 spaces are spanned by the coordinate directions. The Gauss-Seidel type of the BCD method
141 updates only one block by minimizing the objective function or its surrogate while all other
142 blocks are fixed at each iteration. It has been one of the core algorithmic idea in solving
143 problems with block structures, such as convex programming [77], nonlinear programming
144 [9], semidefinite programming [129, 145], compressive sensing [72, 32], etc. A proximal
145 alternating linearized minimization method is developed in [10] for solving a summation of
146 nonconvex but differentiable and nonsmooth functions. The alternating direction methods of
147 multipliers (ADMM) [11, 27, 41, 45, 55, 125] minimize the augmented Lagrangian function
148 with respect to the primal variables by BCD, then update the Lagrangian multiplier.

149 **Active Sets.** When a clear partition of variables is not available, a subset of the variables
150 can be fixed in the so-called active sets under certain mechanisms and the remaining variables
151 are determined from certain subproblems for optimization problems with bound constraints
152 or linear constraints in [17, 18, 51, 81, 82], ℓ_1 -regularized problem for sparse optimization
153 in [133, 105, 64] and general nonlinear programs in [19, 20]. In quadratic programming, the
154 inequality constraints that have zero values at the optimal solution are called active, and they
155 are replaced by equality constraints in the subproblem [111].

156 **Limited-memory.** A typical subspace is constructed from a number of history infor-
157 mation, for example, the previous iterates $\{x_k\}$, the previous gradients $\{\nabla f(x_k)\}$, the dif-
158 ferences between two consecutive iterates $\{x_k - x_{k-1}\}$, and the differences between two
159 consecutive gradients $\{\nabla f(x_k) - \nabla f(x_{k-1})\}$. After the new iterate is formed, the oldest
160 vectors in the storage are replaced by the most recent vectors if certain justification rules are
161 satisfied. Two examples are the limited memory BFGS method [111, 88], and the limited
162 memory block Krylov subspace optimization method (LMSVD) [74].

163 **Anderson Acceleration.** For a sequence $\{x_k\}$ generated by a general fixed-point iter-
164 ation, the Anderson acceleration produces a new point using a linear combination of a few
165 points in $\{x_k\}$, where the coefficients are determined from an extra linear least squares prob-
166 lem with a normalized constraint [13, 4, 123]. A few related schemes include the minimal
167 polynomial extrapolation, modified minimal polynomial extrapolation, reduced rank extrap-
168 olation, the vector Epsilon algorithm and the topological Epsilon algorithm. The Anderson
169 acceleration is also known as Anderson mixing, Pulay mixing, DIIS or the commutator DIIS
170 [92, 93, 115] in electronic structure calculation. These techniques have also been applied to
171 optimization problems in [99, 147].

172 **Subspace correction.** For variational problems, the domain decomposition methods

173 split the spatial domain into several subdomains and solve the corresponding problems on
 174 these subdomains iteratively using certain strategies. The successive subspace correction
 175 (SSC) and parallel subspace correction (PSC) methods [22, 36, 39, 38, 68, 112] are similar to
 176 the Gauss-Seidel-type and Jacobian-type BCD methods, respectively. However, the subspace
 177 correction is significantly different from BCD due to the strong connections between variables
 178 in the subdomains. The PSC methods have been studied for LASSO in [36, 39, 29] and total
 179 variation minimization in [37, 38, 39, 68].

180 **Sampling.** Assume that there are a large number of data. The general concept of sam-
 181 pling is to randomly select a small set of samples with an appropriate probability distribution
 182 with or without replacement. In the stochastic gradient descent type methods, the gradient in
 183 expectation is approximated by a sum of sample gradients over a mini-batch [47]. Random
 184 sampling is also helpful in many other contexts, for example, a greedy algorithm for a mixed
 185 integer programming in volumetric modulated arc therapy [139].

186 **Sketching.** For huge data represented in matrices, the sketching technique builds low-
 187 dimensional approximations using random linear maps [78, 136, 118]. It has been adopted
 188 for nonlinear least squares problems in [141, 103] and large scale SDP problems in [144].
 189 The Nyström approximation can be viewed as a special sketching scheme. An initial quasi-
 190 Newton matrix can be constructed if a single Hessian-matrix multiplication is affordable in
 191 [58].

192 **1.2. Notation.** Let \mathcal{S}^n be the collection of all n -by- n symmetric matrices. For any
 193 matrix $X \in \mathbb{R}^{n \times n}$, $\text{diag}(X)$ denotes a column vector consisting of all diagonal entries of X .
 194 For any vector $x \in \mathbb{R}^n$, $\text{Diag}(x)$ is an n -by- n diagonal matrix whose i -th diagonal entry is
 195 x_i . Given two matrices $A, B \in \mathbb{C}^{n \times p}$, the Frobenius inner product is defined as $\langle A, B \rangle =$
 196 $\text{tr}(A^*B)$, and the corresponding Frobenius norm is defined as $\|A\|_F = \sqrt{\text{tr}(A^*A)}$. The
 197 operation $A \odot B$ denotes the Hadamard product between two matrices A and B of the same
 198 sizes. Let e_n be a vector of all ones in \mathbb{R}^n . For any matrix $X \in \mathbb{R}^{n \times p}$, $\text{Range}(X)$ denotes the
 199 subspace spanned by the columns of X . The subscript usually denotes the iteration number,
 200 while the superscript is reserved as the index of a vector or matrix.

201 **1.3. Organization.** The rest of this paper is organized as follows. The subspace meth-
 202 ods applied in general unconstrained optimization, nonlinear equations and nonlinear least
 203 squares problem, stochastic optimization, sparse optimization, the domain decomposition,
 204 general constrained optimization, eigenvalue computation, optimization problems with or-
 205 thogonality constraints, semidefinite programming and low rank matrix optimization are dis-
 206 cussed in Sections 2 to 11, respectively. Finally, a few typical scenarios are summarized in
 207 Section 12.

208 **2. General Unconstrained Optimization.** In this section, we consider the uncon-
 209 strained optimization

$$210 \quad (2.1) \quad \min_{x \in \mathbb{R}^n} f(x),$$

211 where $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable function. The line search and trust region methods
 212 are the two main types of approaches for solving (2.1). The main difference between them
 213 is the order of determining the so-called step size and search direction. Subspace techniques
 214 have been substantially studied in [26, 48, 140, 142, 143, 87, 128, 127, 49].

215 **2.1. The Line Search Methods.** At the k -th iteration x_k , the line search methods
 216 first generate a descent search direction d_k and then search along this direction for a step size
 217 α_k such that the objective function at the next point

$$218 \quad (2.2) \quad x_{k+1} = x_k + \alpha_k d_k$$

219 is suitably reduced. The step size α_k is often selected by the monotone line search procedures
 220 with the Armijo, Goldstein or the Wolfe-Powell rules. The nonmonotone line procedures are
 221 also widely used. Interested readers are referred to [111, 88] for further information. Here,
 222 we mainly focus on generating the direction d_k in a subspace \mathfrak{S}_k , i.e.,

$$223 \quad d \in \mathfrak{S}_k.$$

224 For simplicity, we often denote $g_k = \nabla f(x_k)$.

225 **2.1.1. The Nonlinear Conjugate Gradient (CG) Method.** The nonlinear CG method
 226 is popular for solving large scale optimization problems. The search direction d_k lies in a par-
 227 ticular subspace

$$228 \quad (2.3) \quad \mathfrak{S}_k = \text{span}\{g_k, d_{k-1}\},$$

229 which is spanned by the gradient g_k and the last search direction d_{k-1} . More specifically, d_k
 230 is a linear combination of $-g_k$ and d_{k-1} with a weight β_{k-1} , i.e.,

$$231 \quad (2.4) \quad d_k = -g_k + \beta_{k-1}d_{k-1},$$

232 where $d_0 = -g_0$ and $\beta_0 = 0$. A few widely used choices for the weight β_{k-1} are

$$233 \quad \beta_{k-1} = \frac{g_k^\top g_k}{g_{k-1}^\top g_{k-1}}, \quad (\text{F-R Formula}),$$

$$234 \quad \beta_{k-1} = \frac{g_k^\top (g_k - g_{k-1})}{d_{k-1}^\top (g_k - g_{k-1})}, \quad (\text{H-S or C-W Formula}),$$

$$235 \quad \beta_{k-1} = \frac{g_k^\top (g_k - g_{k-1})}{g_{k-1}^\top g_{k-1}}, \quad (\text{PRP Formula}),$$

$$236 \quad \beta_{k-1} = -\frac{g_k^\top g_k}{d_{k-1}^\top g_{k-1}}, \quad (\text{Dixon Formula}),$$

$$237 \quad \beta_{k-1} = -\frac{g_k^\top g_k}{d_{k-1}^\top (g_k - g_{k-1})}, \quad (\text{D-Y Formula}).$$

238
 239 It is easy to observe that these formulas are equivalent in the sense that they yield the same
 240 search directions when the function $f(x)$ is quadratic with a positive definite Hessian matrix.
 241 In this case, the directions d_1, \dots, d_k are conjugate to each other with respect to the Hessian
 242 matrix. It can also be proved that the CG method has global convergence and n -step local
 243 quadratic convergence. However, for a general nonlinear function with inexact line search,
 244 the behavior of the methods with different β_k can be significantly different.

245 **2.1.2. Nesterov's Accelerated Gradient Method.** The steepest descent gradient
 246 method simply uses $d_k = -g_k$ in (2.2) for unconstrained optimization. Assume that the
 247 function $f(x)$ is convex, the optimal value f^* of (2.1) is finite and it attains at a point x^* , and
 248 the gradient $f(x)$ is Lipschitz continuous with a constant L , i.e.,

$$249 \quad \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|.$$

250 Let $\{x_k\}_{k=0}^\infty$ be a sequence generated by the gradient method with a fixed step size $\alpha_k = \frac{1}{L}$.
 251 Then it can be proved that the convergence of the objective function values is

$$252 \quad f(x_k) - f(x^*) \leq \frac{L}{2k} \|x_0 - x^*\|^2,$$

253 which is often described as a convergence rate at $O(1/k)$.

254 A natural question is whether a faster convergence rate can be achieved if only the gra-
 255 dient information is used. We now present the so-called FISTA method proposed by Beck
 256 and Teboulle [5] which is equivalent to Nesterov accelerated gradient method [84, 85]. The
 257 FISTA method first calculates a new point by an extrapolation of the previous two points,
 258 then performs a gradient step at this new point:

$$259 \quad y_k = x_{k-1} + \frac{k-2}{k+1}(x_{k-1} - x_{k-2}),$$

$$260 \quad x_k = y_k - \alpha_k \nabla f(y_k).$$

An illustration of the FISTA method is shown in Figure 2.1. Under the same assumptions as

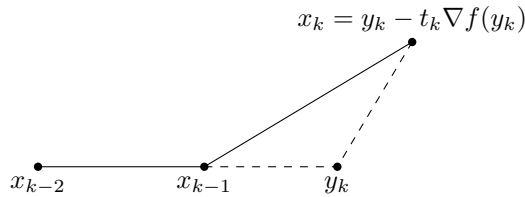


Fig. 2.1 The FISTA method

261 the gradient method, the FISTA method with a fixed step size $\alpha_k = \frac{1}{L}$ has a convergence rate
 262 of $O(1/k^2)$, i.e.,

$$264 \quad f(x_k) - f^* \leq \frac{2L}{(k+1)^2} \|x_0 - x^*\|^2.$$

265 Obviously, the FISTA method can also be interpreted as a subspace method whose subspace
 266 is

$$267 \quad (2.5) \quad \mathfrak{S}_k = \text{span}\{x_{k-1}, x_{k-2}, \nabla f(y_k)\}.$$

268 **2.1.3. The Heavy-ball Method.** The heavy-ball method [90] is also a two-step scheme: ■

$$269 \quad d_k = -g_k + \beta d_{k-1},$$

$$270 \quad x_{k+1} = x_k + \alpha d_k,$$

271 with $p_0 = 0$ and $\alpha, \beta > 0$. If $\beta \in [0, 1)$ and $\alpha \in \left(0, \frac{1-\beta}{L}\right]$ and under the same assumptions
 272 as in Sec. 2.1.2, it is established in [42] that

$$273 \quad f(\bar{x}_k) - f^* \leq \frac{1}{k+1} \left(\frac{\beta}{1-\beta} (f(x_0) - f^*) + \frac{1-\beta}{2\alpha} \|x_0 - x^*\|^2 \right),$$

274 where $\bar{x}_k = \frac{1}{1+k} \sum_{i=1}^k x_i$. We can see that the Heavy-ball method is the same as the nonlin-
 275 ear CG method (2.4) except that the parameter β is different.

276 **2.1.4. A Search Direction Correction (SDC) Method.** The search direction (2.4)
 277 can also be modified by adding a non-trivial weight to g_k . Let $d_0 = 0$. At the beginning of
 278 the $(k+1)$ -th iteration, if a descent condition

$$279 \quad (2.6) \quad \langle g_k, d_k \rangle \leq 0$$

280 holds, we update

$$281 \quad (2.7) \quad d_{k+1} = (1 - \beta_k)d_k - \gamma_k \frac{\|d_k\|}{\|g_k\|} g_k - g_k.$$

282 Then we update β_{k+1} and γ_{k+1} as follows:

$$283 \quad (2.8) \quad \beta_k = \frac{r}{l_k - 1 + r}, \quad \gamma_k = \frac{r - 3}{l_k - 1 + r},$$

284 where $r \geq 3$, $\{l_k\}$ is a sequence of parameters with of $l_1 = 1$ and $l_{k+1} = l_k + 1$. If the
285 criterion (2.6) is not met, we reset d_{k+1} , β_{k+1} and γ_{k+1} as

$$286 \quad d_{k+1} = -g_k, \beta_{k+1} = \beta_1, \gamma_{k+1} = \gamma_1, l_{k+1} = l_1.$$

287 For more details, we refer the reader to [126].

288 **2.1.5. Quasi-Newton Methods.** The search directions of the limited-memory quasi-
289 Newton methods [111, 88] also lie in subspaces. Let B_k be the limited-memory BFGS (L-
290 BFGS) matrix and H_k be its inverse matrix generated from a few most recent pairs $\{s_i, y_i\}$,
291 where

$$292 \quad s_i = x_{i+1} - x_i, \quad y_i = g_{i+1} - g_i.$$

293 Then the search direction is

$$294 \quad (2.9) \quad d_k = -B_k^{-1}g_k = -H_k g_k,$$

295 which is usually computed by the two-loop recursion. In fact, both B_k and H_k can be written
296 in a compact representation [21]. Assume that there are p pairs of vectors:

$$297 \quad (2.10) \quad U_k = [s_{k-p}, \dots, s_{k-1}] \in \mathbb{R}^{n \times p}, \quad Y_k = [y_{k-p}, \dots, y_{k-1}] \in \mathbb{R}^{n \times p}.$$

298 For a given initial matrix H_k^0 , the H_k matrix is:

$$299 \quad (2.11) \quad H_k = H_k^0 + C_k P_k C_k^\top,$$

300 where

$$301 \quad C_k := [U_k, H_k^0 Y_k] \in \mathbb{R}^{n \times 2p}, \quad D_k = \text{diag} [s_{k-p}^\top y_{k-p}, \dots, s_{k-1}^\top y_{k-1}]$$

$$302 \quad P_k := \begin{bmatrix} R_k^{-\top} (D_k + Y_k^\top H_k^0 Y_k) R_k^{-1} & -R_k^{-\top} \\ -R_k^{-1} & 0 \end{bmatrix}, (R_k)_{i,j} = \begin{cases} s_{k-p+i-1}^\top y_{k-p+j-1}, & \text{if } i \leq j, \\ 0, & \text{o.w.} \end{cases}$$

304 The initial matrix H_k^0 is usually set to be a positive scalar γ_k times the identity matrix, i.e.,
305 $\gamma_k I$. Therefore, we have

$$306 \quad d_k \in \text{span}\{g_k, s_{k-1}, \dots, s_{k-p}, y_{k-1}, \dots, y_{k-p}\}.$$

307 **2.1.6. Acceleration Techniques.** Gradient descent algorithms may converge slowly
308 after certain iterations. This issue can be resolved by using acceleration techniques such
309 as Anderson Acceleration (AA) [4, 123]. An extrapolation-based acceleration techniques
310 proposed in [99] can be applied to overcome the instability of the Anderson Acceleration. To
311 be precise, we perform linear combinations of the points x_k every $l + 2$ iterations to obtain a
312 better estimation $\tilde{x} = \sum_{i=0}^l \tilde{c}_i x_{k-l+i}$. Define the difference of $l + 2$ iteration points as

$$313 \quad U = [x_{k-l+1} - x_{k-l}, \dots, x_{k+1} - x_k].$$

314 Then the coefficients $\tilde{c} = (\tilde{c}_0, \dots, \tilde{c}_l)^\top$ is the solution of the following problem

$$315 \quad (2.12) \quad \tilde{c} = \arg \min_{c^\top e_{l+1}=1} c^\top (U^\top U + \lambda I) c,$$

316 where $\lambda > 0$ is a regularization parameter.

317 **2.1.7. Search Direction From Minimization Subproblems.** We next construct
 318 the search direction by solving a subproblem defined in a subspace \mathfrak{S}_k as

$$319 \quad (2.13) \quad \min_{d \in \mathfrak{S}_k} Q_k(d),$$

320 where $Q_k(d)$ is an approximation to $f(x_k + d)$ for d in the subspace \mathfrak{S}_k . It would be de-
 321 sirable that the approximation model $Q_k(d)$ has the following properties: (i) it is easy to be
 322 minimized in the subspace \mathfrak{S}_k ; (ii) it is a good approximation to f and the solution of the
 323 subspace subproblem will give a sufficient reduction with respect to the original objective
 324 function f .

325 It is natural to use quadratic approximations to the objective function. This leads to
 326 quadratic models in subspaces. A successive two-dimensional search algorithm is developed
 327 by Stoer and Yuan in [143] based on

$$328 \quad \min_{d \in \text{span}\{-g_k, d_{k-1}\}} Q_k(d).$$

Let the dimension $\dim(\mathfrak{S}_k) = \tau_k$ and \mathfrak{S}_k be a set generated by all linear combinations of
 vectors $p_1, p_2, \dots, p_{\tau_k} \in \mathbb{R}^n$, i.e.,

$$\mathfrak{S}_k = \text{span}\{p_1, p_2, \dots, p_{\tau_k}\}.$$

329 Define $P_k = [p_1, p_2, \dots, p_{\tau_k}]$. Then $d \in \mathfrak{S}_k$ can be represented as $d = P_k \bar{d}$ with $\bar{d} \in \mathbb{R}^{\tau_k}$.
 330 Hence, a quadratic function $Q_k(d)$ defined in the subspace can be expressed as a function \bar{Q}_k
 331 in a lower dimension space \mathbb{R}^{τ_k} in terms of $Q_k(d) = \bar{Q}_k(\bar{d})$. Since τ_k usually is quite small,
 332 the Newton method can be used to solve (2.13) efficiently.

333 We now discuss a few possible choices for the subspace \mathfrak{S}_k . A special subspace is a
 334 combination of the current gradient and the previous search directions, i.e.,

$$335 \quad (2.14) \quad \mathfrak{S}_k = \text{span}\{-g_k, s_{k-1}, \dots, s_{k-m}\}.$$

336 In this case, any vector d in the subspace \mathfrak{S}_k can be expressed as

$$337 \quad (2.15) \quad d = \alpha g_k + \sum_{i=1}^m \beta_i s_{k-i} = (-g_k, s_{k-1}, \dots, s_{k-m}) \bar{d}$$

338 where $\bar{d} = (\alpha, \beta_1, \dots, \beta_m)^\top \in \mathbb{R}^{m+1}$. All second order terms of the Taylor expansion of
 339 $f(x_k + d)$ in the subspace \mathfrak{S}_k can be approximated by secant conditions

$$340 \quad (2.16) \quad s_{k-i}^\top \nabla^2 f(x_k) s_{k-j} \approx s_{k-i}^\top y_{k-j}, \quad s_{k-i}^\top \nabla^2 f(x_k) g_k \approx y_{k-i}^\top g_k,$$

341 except $g_k^\top \nabla^2 f(x_k) g_k$. Therefore, it is reasonable to use the following quadratic model in the
 342 subspace \mathfrak{S}_k :

$$343 \quad (2.17) \quad \bar{Q}_k(\bar{d}) = (-\|g_k\|^2, g_k^\top s_{k-1}, \dots, g_k^\top s_{k-m}) \bar{d} + \frac{1}{2} \bar{d}^\top \bar{B}_k \bar{d},$$

344 where

$$345 \quad (2.18) \quad \bar{B}_k = \begin{pmatrix} \rho_k & -g_k^\top y_{k-1} & \dots & -g_k^\top y_{k-m} \\ -g_k^\top y_{k-1} & y_{k-1}^\top s_{k-1} & \dots & y_{k-m}^\top s_{k-1} \\ \vdots & \vdots & \ddots & \vdots \\ -g_k^\top y_{k-m} & y_{k-m}^\top s_{k-1} & \dots & y_{k-m}^\top s_{k-m} \end{pmatrix}$$

346 with $\rho_k \approx g_k^\top \nabla^2 f(x_k) g_k$. Hence, once a good estimation to the term $g_k^\top \nabla^2 f(x_k) g_k$ is
 347 available, we can obtain a good quadratic model in the subspace \mathfrak{S}_k .

348 There are different ways to choose ρ_k . Similar to the approach in [143], we can let

$$349 \quad (2.19) \quad \rho_k = 2 \frac{(s_{k-1}^\top g_k)^2}{s_{k-1}^\top y_{k-1}},$$

350 due to the fact that the mean value of $\cos^2(\theta)$ is $\frac{1}{2}$, which gives

$$351 \quad (2.20) \quad g_k^\top \nabla^2 f(x_k) g_k = \frac{1}{\cos^2 \theta_k} \frac{(s_{k-1}^\top \nabla^2 f(x_k) g_k)^2}{s_{k-1}^\top \nabla^2 f(x_k) s_{k-1}} \approx 2 \frac{(s_{k-1}^\top g_k)^2}{s_{k-1}^\top y_{k-1}},$$

352 where θ_k is the angle between $(\nabla^2 f(x_k))^{\frac{1}{2}} g_k$ and $(\nabla^2 f(x_k))^{\frac{1}{2}} s_{k-1}$. Another way to es-
 353 timate $g_k^\top (\nabla^2 f(x_k)) g_k$ is by replacing $\nabla^2 f(x_k)$ by a quasi-Newton matrix. We can also
 354 obtain ρ_k by computing an extra function value $f(x_k + t g_k)$ and setting

$$355 \quad (2.21) \quad \rho_k = \frac{2(f(x_k + t g_k) - f(x_k) - t \|g_k\|_2^2)}{t^2}.$$

356 By letting the second order curvature along g_k to be the average of those along s_{k-i} ($i =$
 357 $1, \dots, m$), we obtain

$$358 \quad (2.22) \quad \rho_k = \frac{\|g_k\|_2^2}{m} \sum_{i=1}^m \frac{s_{k-i}^\top y_{k-i}}{s_{k-i}^\top s_{k-i}}.$$

359 Similar to (2.14), a slightly different subspace is

$$360 \quad (2.23) \quad \mathfrak{S}_k = \text{span}\{-g_k, y_{k-1}, \dots, y_{k-m}\}.$$

361 In this case, any vector in \mathfrak{S}_k can be represented as

$$362 \quad (2.24) \quad d = \alpha g_k + \sum_{i=1}^m \beta_i y_{k-i} = W_k \bar{d}$$

363 where $W_k = [-g_k, y_{k-1}, \dots, y_{k-m}] \in \mathbb{R}^{n \times (m+1)}$. The Newton step in the subspace \mathfrak{S}_k is
 364 $W_k \bar{d}_k$ with

$$365 \quad (2.25) \quad \bar{d}_k = -[W_k^\top \nabla^2 f(x_k) W_k]^{-1} W_k^\top \nabla f(x_k).$$

366 Thus, the remaining issue is to obtain a good estimate of \bar{d}_k , using the fact that all the elements
 367 of $[W_k^\top (\nabla^2 f(x_k))^{-1} W_k]$ is known except one entry $g_k^\top \nabla^2 f(x_k)^{-1} g_k$.

368 **2.1.8. Subspace By Coordinate Directions.** We next consider subspaces spanned
 369 by coordinate directions with sparsity structures. Let g_k^i be the i -th component of the gradient
 370 g_k . The absolute values $|g_k^i|$ are sorted in the descending order such that

$$371 \quad (2.26) \quad |g_k^{i_1}| \geq |g_k^{i_2}| \geq |g_k^{i_3}| \geq \dots \geq |g_k^{i_n}|.$$

372 The subspace

$$373 \quad (2.27) \quad \mathfrak{S}_k = \text{span}\{e^{i_1}, e^{i_2}, \dots, e^{i_\tau}\}$$

374 is called as the τ -steepest coordinates subspace, where e^i is a vector of all zeros except that
 375 the i -th component is one. Then, the steepest descent direction in the subspace is sufficiently
 376 descent, namely

$$377 \quad (2.28) \quad \min_{d \in \mathfrak{S}_k} \frac{d^\top g_k}{\|d\|_2 \|g_k\|_2} \leq -\frac{\tau}{n}.$$

378 When $(g_k^{i_{\tau+1}})^2 \leq \epsilon \sum_{j=1}^{\tau} (g_k^{i_j})^2$, the following estimation can be established:

$$379 \quad (2.29) \quad \min_{d \in \mathfrak{S}_k} \frac{d^\top g_k}{\|d\|_2 \|g_k\|_2} \leq -\frac{1}{\sqrt{1 + \epsilon(n - \tau)}}.$$

Consequently, a sequential steepest coordinates search (SSCS) technique can be designed
 by augmenting the steepest coordinate directions into the subspace sequentially. For example,
 consider minimizing a convex quadratic function

$$Q(x) = g^\top x + \frac{1}{2} x^\top Bx.$$

380 At the k -th iteration of SSCS, we first compute $g_k = \nabla Q(x_k)$, then choose

$$381 \quad i_k = \arg \min_i \{|g_k^i|\}.$$

382 Let $\mathfrak{S}_k = \text{span}\{e^{i_1}, \dots, e^{i_k}\}$. The next iteration is to find

$$383 \quad x_{k+1} = \arg \min_{x \in x_k + \mathfrak{S}_k} Q(x).$$

384 **2.2. Trust Region Methods.** The trust region methods for (2.1) compute a search
 385 direction in a ball determined by a given trust region radius whose role is similar to the step
 386 size. The trust region subproblem (TRS) is normally

$$387 \quad (2.30) \quad \begin{aligned} \min_{s \in \mathbb{R}^n} \quad & Q_k(s) = g_k^\top s + \frac{1}{2} s^\top B_k s \\ \text{s. t.} \quad & \|s\|_2 \leq \Delta_k, \end{aligned}$$

388 where B_k is an approximation to the Hessian $\nabla^2 f(x_k)$ and $\Delta_k > 0$ is the trust region radius.

389 A subspace version of the trust region subproblem is suggested in [101]:

$$390 \quad \begin{aligned} \min_{s \in \mathbb{R}^n} \quad & Q_k(s) \\ \text{s. t.} \quad & \|s\|_2 \leq \Delta_k, \quad s \in \mathfrak{S}_k. \end{aligned}$$

391 The Steihaug truncated CG method [107] for solving (2.30) is essentially a subspace method.
 392 When the approximate Hessian B_k is generated by the quasi-Newton updates such as the SR1,
 393 PSB or the Broyden family [111, 88], the TRS has subspace properties. Suppose $B_1 = \sigma I$
 394 with $\sigma > 0$, let s_k be an optimal solution of TRS (2.30) and set $x_{k+1} = x_k + s_k$. Let
 395 $\mathcal{G}_k = \text{span}\{g_1, g_2, \dots, g_k\}$. Then it can be proved that $s_k \in \mathcal{G}_k$ and for any $z \in \mathcal{G}_k$,
 396 $w \in \mathcal{G}_k^\perp$, it holds

$$397 \quad (2.31) \quad B_k z \in \mathcal{G}_k, \quad B_k w = \sigma w.$$

398 Therefore, the subspace trust region algorithm generates the same sequences as the full space
 399 trust region quasi-Newton algorithm. Based on the above results, Wang and Yuan [128]

400 presented a subspace trust region quasi-Newton method for large scale unconstrained opti-
 401 mization. Similar results for the line search quasi-Newton methods were obtained by Gill
 402 and Leonard [44, 43].

403 We next discuss a special trust region subproblem which can make good use of subspace
 404 properties. If the norm $\|\cdot\|_2$ is replaced by a general norm $\|\cdot\|_W$ in (2.30), we can obtain a
 405 general TRS subproblem

$$406 \quad \min_{s \in \mathbb{R}^n} \quad g^\top s + \frac{1}{2} s^\top B s$$

$$\text{s. t.} \quad \|s\|_W \leq \Delta.$$

Here, the subscript k in g_k and B_k is omitted for simplicity. Intuitively, we should choose
 the norm $\|\cdot\|_W$ properly so that the TRS can easily be solved by using the corresponding
 subspace properties of the objective function $g^\top s + \frac{1}{2} s^\top B s$. Assume that B is a limited
 memory quasi-Newton matrix which can be expressed as

$$B = \sigma I + P D P^\top, \quad P \in \mathbb{R}^{n \times l},$$

407 where $P^\top P = I$. Let P_\perp^\top be the projection onto the space orthogonal to $\text{Range}(P)$. Then
 408 the following function

$$409 \quad (2.32) \quad \|s\|_P = \max\{\|P^\top s\|_\infty, \|P_\perp^\top s\|_2\}$$

410 is a well-defined norm, which leads to a P -norm TRS:

$$411 \quad (2.33) \quad \min_{s \in \mathbb{R}^n} \quad g^\top s + \frac{1}{2} s^\top B s$$

$$\text{s. t.} \quad \|s\|_P \leq \Delta.$$

412 Due to the definition of the $\|\cdot\|_P$, the solution s of the TRS (2.33) can be expressed by

$$413 \quad s = P s_1 + P_\perp s_2,$$

414 where s_1 and s_2 can be computed easily. In fact, s_1 is the solution of the box constrained
 415 quadratic program (QP):

$$416 \quad \min_{s \in \mathbb{R}^l} \quad s^\top (P^\top g) + \frac{1}{2} s^\top (\sigma I + D) s$$

$$\text{s. t.} \quad \|s\|_\infty \leq \Delta,$$

417 It can be verified that s_1 has a closed form solution:

$$418 \quad (s_1)_i = \begin{cases} \frac{-(P^\top g)_i}{\sigma + D_{ii}} & \text{if } |(P^\top g)_i| < (\sigma + D_{ii})\Delta, \\ \Delta \text{sign}(-(P^\top g)_i) & \text{otherwise,} \end{cases}$$

419 for $i = 1, \dots, l$. On the other hand, s_2 is solution of the 2-norm constrained special QP

$$420 \quad \min_{s \in \mathbb{R}^{n-l}} \quad s^\top (P_\perp^\top g) + \frac{1}{2} \sigma s^\top s$$

$$\text{s. t.} \quad \|s\|_2 \leq \Delta.$$

421 whose closed-form solution is

$$422 \quad s_2 = -\min\left(\frac{1}{\sigma}, \frac{\Delta}{\|P_\perp^\top g\|}\right) P_\perp^\top g.$$

423 Numerical results based on a trust region algorithm that uses the the W -norm TRS were
 424 shown in [15].

425 **3. Nonlinear Equations and Nonlinear Least Squares Problem.** In this sec-
 426 tion, we consider the system of nonlinear equations

$$427 \quad (3.1) \quad F(x) = 0, \quad x \in \mathbb{R}^n,$$

428 and nonlinear least squares problem:

$$429 \quad (3.2) \quad \min_{x \in \mathbb{R}^n} \|F(x)\|_2^2,$$

430 where $F(x) = (F^1(x), F^2(x), \dots, F^m(x))^\top \in \mathbb{R}^m$.

431 **3.1. General Subspace Methods.** Due to the special structures of nonlinear equa-
 432 tions, several implementations of Newton-like iteration schemes based on Krylov subspace
 433 projection methods are considered in [14]. Newton–Krylov methods with a global strategy
 434 restricted to a suitable Krylov subspace are developed in [7]. Because the nonlinear least
 435 squares problem (3.2) is also an unconstrained optimization problem, all the subspace tech-
 436 niques discussed in Section 2 can be applied. For example, assume that there are i_k lin-
 437 early independent vectors $\{q_k^1, q_k^2, \dots, q_k^{i_k}\}$ which spans \mathfrak{S}_k . Let $Q_k = [q_k^1, q_k^2, \dots, q_k^{i_k}]$. Then
 438 $d \in \mathfrak{S}_k$ can be expressed as $Q_k z$ with respect to a variable $z \in \mathbb{R}^{i_k}$. For (3.1), one can find a
 439 subspace step from

$$440 \quad (3.3) \quad F(x_k + Q_k z) = 0.$$

441 The linearized system for subproblem (3.3) is

$$442 \quad (3.4) \quad F(x_k) + J_k Q_k z = 0,$$

443 where J_k is the Jacobian of F at x_k . Similarly, one can construct a subspace type of the
 444 Levenberg-Marquardt method for (3.2) as

$$445 \quad \min_z \|F(x_k) + J_k Q_k z\|_2^2 + \frac{\lambda_k}{2} \|z\|_2^2,$$

446 where λ_k is a regularization parameter adjusted to ensure global convergence.

447 **3.2. Subspace by Subsampling/Sketching.** We start from solving a linear least
 448 squares problem on massive data sets:

$$449 \quad (3.5) \quad \min_x \|Ax - b\|_2^2,$$

450 where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. Although applying the direct or iterative methods to (3.5) is
 451 straightforward, it may be prohibitive for large values of m . The sketching technique chooses
 452 a matrix $W \in \mathbb{R}^{r \times m}$ with $r \ll m$ and formulates a reduced problem

$$453 \quad (3.6) \quad \min_x \|W(Ax - b)\|_2^2.$$

454 It can be proved that the solution of (3.6) can be a good approximation to that of (3.5) in
 455 certain sense if the matrix W is chosen suitably. For example, one may randomly select r
 456 rows from the identity matrix to form W so that WA is a submatrix of A . Another choice is
 457 that each element of W is sampled from an i.i.d. normal random variable with mean zero and
 458 variance $1/r$. These concepts have been extensively investigated for randomized algorithms
 459 in numerical linear algebra [78, 136].

460 For nonlinear equations, the simple sketching approach is to ignore some equations.
 461 Instead of requiring the original system (3.1), we consider

$$462 \quad (3.7) \quad F^i(x) = 0, \quad i \in I_k,$$

463 which is an incomplete set of equations. To solve the nonlinear equations (3.1) is to find a x
 464 at which F maps to the origin [141]. Let P_k^\top be a mapping from \mathbb{R}^m to a lower dimensional
 465 subspace. Solving the reduced system

$$466 \quad (3.8) \quad P_k^\top F(x) = 0$$

467 is exactly replacing $F = 0$ by requiring its mapping to the subspace spanned by P_k to be
 468 zero. Together with (3.3) yields:

$$469 \quad (3.9) \quad P_k^\top F(x_k + Q_k z) = 0,$$

470 The linearized system for subproblem (3.9) is

$$471 \quad (3.10) \quad P_k^\top [F(x_k) + J_k Q_k z] = 0.$$

472 Of course, the efficiency of such an approach depends on how to select P_k and Q_k . We can
 473 borrow ideas from subspace techniques for large scale linear systems [98]. Instead of using
 474 (3.10), we can use a subproblem of the following form:

$$475 \quad (3.11) \quad P_k^\top F(x_k) + \hat{J}_k z = 0,$$

476 where $\hat{J}_k \in \mathbb{R}^{i_k \times i_k}$ is an approximation to $P_k^\top J_k Q_k$. The reason for preferring (3.11) over
 477 (3.10) is that in (3.11) we do not need the Jacobian matrix J_k , whose size is normally signifi-
 478 cantly larger than that of \hat{J}_k .

479 Similar idea has also been studied for nonlinear least squares problems. At the k -th
 480 iteration, we consider minimizing the sum of squares of some randomly selected terms in an
 481 index set $I_k \subset \{1, \dots, m\}$ instead of all terms:

$$482 \quad (3.12) \quad \min_{x \in \mathbb{R}^n} \sum_{i \in I_k} (F^i(x))^2.$$

483 This approach works quite well on the distance geometry problem which has lots of applica-
 484 tions including protein structure prediction, where the nonlinear least squares of all the terms
 485 have lots of local minimizers [103]. Combining subspace with sketching yields a counterpart
 486 to (3.9) for nonlinear least squares:

$$487 \quad (3.13) \quad \min_{d \in \mathfrak{S}_k} \|P_k^\top F(x_k + d)\|_2^2.$$

488 A projected nonlinear least squares method is studied in [57] to fit a model ψ to (noisy)
 489 measurements y for the exponential fitting problem:

$$490 \quad (3.14) \quad \min_{x \in \mathbb{R}^n} \|\psi(x) - y\|_2^2,$$

491 where $\psi(x) \in \mathbb{R}^m$ and $n \ll m$. Since computing the Jacobian of (3.14) can be expensive,
 492 a sequence of low-dimensional surrogate problems are constructed from a sequence of sub-
 493 spaces $\{\mathcal{W}_\ell\} \subset \mathbb{R}^m$. Let $P_{\mathcal{W}_\ell}$ be an orthogonal projection onto \mathcal{W}_ℓ and W_ℓ is an orthonormal
 494 basis for \mathcal{W}_ℓ , i.e., $P_{\mathcal{W}_\ell} = W_\ell W_\ell^\top$. Then it solves the following minimization problem:

$$495 \quad \min_x \|P_{\mathcal{W}_\ell}[\psi(x) - y]\|_2^2 = \min_x \|W_\ell^\top \psi(x) - W_\ell^\top y\|_2^2.$$

496 **3.3. Partition of Variables.** We now consider the partition of variables, which is also
 497 a subspace technique for nonlinear least squares problem. Let I_k be a subset of $\{1, \dots, n\}$.
 498 The variables are partitioned into two group $x = (\bar{x}, \hat{x})$, where $\bar{x} = (x^i, i \in I_k)$ and
 499 $\hat{x} = (x^i, i \notin I_k)$. At the k -th iteration, the variables $x^i (i \notin I_k)$ are fixed and $x^i (i \in I_k)$ are
 500 free to be changed in order to obtain a better iterate point. This procedure yields a subproblem
 501 with fewer variables:

$$502 \quad (3.15) \quad \min_{\bar{x} \in \mathbb{R}^{|I_k|}} \sum_{i=1}^n (F^i(\bar{x}, \hat{x}_k))^2.$$

503 It is easy to see that partition of variables use special subspaces that spanned by coordinate
 504 directions. An obvious generalization of partition of variables is decomposition of the space
 505 which uses subspaces spanned by any given directions.

506 **3.4. τ -steepest Descent Coordinate Subspace.** The τ -steepest descent coordi-
 507 nate subspace discussed in Section 2 can also be extended to nonlinear equations and non-
 508 linear least squares. Assume that

$$509 \quad (3.16) \quad |F^{i_1}(x_k)| > \dots > |F^{i_\tau}(x_k)| > \dots$$

510 at the k -th iteration. If $F(x)$ is a monotone operator, applying the method in a subspace
 511 spanned by the coordinate directions $\{e^{i_j}, j = 1, \dots, \tau\}$ generates a system

$$512 \quad (3.17) \quad F^{i_j}(x_k) + d^\top \nabla F^{i_j}(x_k) = 0, \quad j = 1, \dots, \tau.$$

For general nonlinear functions $F(x)$, it is better to replace e^{i_j} by the steepest descent coordi-
 nate direction of the function $F^{i_j}(x)$ at x_k , i.e., substituting i_j by an index l_j such that

$$l_j = \arg \max_{t=1, \dots, n} \left| \frac{\partial F^{i_j}(x_k)}{\partial x^t} \right|.$$

513 However, it may be possible to have two different j at one l_j so that subproblem (3.17) has no
 514 solution in the subspace spanned by $\{e^{l_1}, \dots, e^{l_\tau}\}$. Further discussion on subspace methods
 515 for nonlinear equations and nonlinear least squares can be found in [141].

516 **4. Stochastic Optimization.** The supervised learning model in machine learning as-
 517 sumes that (a, b) follows a probability distribution P , where a is an input data and b is the
 518 corresponding label. Let $\phi(a, x)$ be a prediction function in a certain functional space and
 519 $\ell(\cdot, \cdot)$ represent a loss function to measure the accuracy between the prediction and the true la-
 520 bel. The task is to predict a label b from a given input a , that is, finding a function ϕ such that
 521 the expected risk $\mathbb{E}[\ell(\phi(a, x), b)]$ is minimized. In practice, the real probability distribution
 522 P is unknown, but a data set $\mathcal{D} = \{(a_1, b_1), (a_2, b_2), \dots, (a_N, b_N)\}$ is obtained by random
 523 sampling, where $(a_i, b_i) \sim P$ i.i.d. Then an empirical risk minimization is formulated as

$$524 \quad (4.1) \quad \min_x f(x) := \frac{1}{N} \sum_{i=1}^N f_i(x),$$

525 where $f_i(x) = \ell(\phi(a_i; x), b_i)$. In machine learning, a large number of problems can be
 526 expressed in the form of (4.1). For example, the function ϕ in deep learning is expressed
 527 by a deep neural network. Since the size N usually is huge, it is time consuming to use the
 528 information of all components $f_i(x)$. However, it is affordable to compute the information at
 529 a few samples so that the amount of calculation in each step is greatly reduced.

530 **4.1. Stochastic First-order Methods.** In this subsection, we briefly review a few
531 widely used stochastic first-order methods [47]. Instead of using the full gradient $\nabla f(x_k)$,
532 the stochastic gradient method (SGD) for (4.1) selects a uniformly random sample s_k from
533 $\{1, \dots, N\}$ and updates

$$534 \quad (4.2) \quad x_{k+1} = x_k - \alpha_k \nabla f_{s_k}(x_k).$$

535 A common assumption for convergence is that the expectation of the stochastic gradient is
536 equal to the full gradient, i.e.,

$$537 \quad \mathbb{E}[\nabla f_{s_k}(x_k) \mid x_k] = \nabla f(x_k).$$

538 When $f_i(x_k)$ is not differentiable, then its subgradient is used in (4.2). Note that only one
539 sample is used in (4.2). The mini-batch SGD method tries to balance between the robustness
540 of the SGD and the computational cost. It randomly selects a mini-batch $I_k \subset \{1, \dots, N\}$
541 such that $|I_k|$ is quite small, then computes

$$542 \quad (4.3) \quad x_{k+1} = x_k - \frac{\alpha_k}{|I_k|} \sum_{s_k \in I_k} \nabla f_{s_k}(x_k).$$

543 Obviously, subsampling defines a kind of subspace in terms of the component functions
544 $\{f_1(x), \dots, f_N(x)\}$. For simplicity, we next only consider extensions based on (4.2).

545 After calculating a random gradient $\nabla f_{s_k}(x_k)$ at the current point, the momentum method
546 does not directly update it to the variable x_k . It introduces a speed variable v , which represents
547 the direction and magnitude of the parameter movements. Formally, the iterative scheme is

$$548 \quad (4.4) \quad \begin{aligned} v_{k+1} &= \mu_k v_k - \alpha_k \nabla f_{s_k}(x_k), \\ x_{k+1} &= x_k + v_{k+1}. \end{aligned}$$

549 This new update direction v is a linear combination of the previous update direction v_k and
550 the gradient $\nabla f_{s_k}(x_k)$ to obtain a new v_{k+1} . When $\mu_k = 0$, the algorithm degenerates to
551 SGD. In the momentum method, the parameter μ_k is often in the range of $[0, 1)$. A typical
552 value is $\mu_k \geq 0.5$, which means that the iteration point has a large inertia and each iteration
553 will make a small correction to the previous direction.

554 Since the dimension of the variable x can be more than 10 million and the convergence
555 speed of each variable may be different, updating all components of x using a single step size
556 α_k may not be suitable. The adaptive subgradient method (AdaGrad) controls the step sizes of
557 each component separately by monitoring the accumulation of the gradients elementwisely:

$$558 \quad G_k = \sum_{i=1}^k \nabla f_{s_i}(x_i) \odot \nabla f_{s_i}(x_i),$$

559 where \odot the Hadamard product between two vectors. The AdaGrad method is

$$560 \quad (4.5) \quad \begin{aligned} x_{k+1} &= x_k - \frac{\alpha_k}{\sqrt{G_k + \epsilon} e_n} \odot \nabla f_{s_{k+1}}(x_{k+1}), \\ G_{k+1} &= G_k + \nabla f_{s_{k+1}}(x_{k+1}) \odot \nabla f_{s_{k+1}}(x_{k+1}), \end{aligned}$$

561 where the division in $\frac{\alpha_k}{\sqrt{G_k + \epsilon} e_n}$ is also performed elementwisely. Adding the term $\epsilon \mathbf{1}_n$ is to
562 prevent the division by zeros.

563 The adaptive moment estimation (Adam) method combines the momentum and AdaGrad
564 together by adding a few small corrections. At each iteration, it performs:

$$565 \quad v_k = \rho_1 v_{k-1} + (1 - \rho_1) \nabla f_{s_k}(x_k),$$

$$\begin{aligned}
566 \quad G_k &= \rho_2 G_{k-1} + (1 - \rho_2) \nabla f_{s_k}(x_k) \odot \nabla f_{s_k}(x_k), \\
567 \quad \hat{v}_k &= \frac{v_k}{1 - \rho_1^k}, \\
568 \quad \hat{G}_k &= \frac{G_k}{1 - \rho_2^k}, \\
569 \quad x_{k+1} &= x_k - \frac{\alpha_k}{\sqrt{\hat{G}_k + \epsilon \epsilon_n}} \odot \hat{v}_k,
\end{aligned}$$

570 where the typical values for ρ_1 and ρ_2 are $\rho_1 = 0.9$ and $\rho_2 = 0.999$. We can see that the
571 direction v_k is a convex combination of v_{k-1} and $\nabla f_{s_k}(x_k)$, then it is corrected to \hat{v}_k . The
572 value \hat{G}_k is also obtained in a similar fashion. The main advantage of Adam is that after the
573 deviation correction, the step size of each iteration has a certain range, making the parameters
574 relatively stable.

575 The above algorithms have been implemented in mainstream deep learning frameworks,
576 which can be very convenient for training neural networks. The algorithms implemented
577 in Pytorch are AdaDelta, AdaGrad, Adam, Nesterov, RMSProp, etc. The algorithms im-
578 plemented in Tensorflow are AdaDelta, AdaGradDA, AdaGrad, ProximalAdagrad, Ftrl, Mo-
579 mentum, Adam and CenteredRMSProp, etc.

580 **4.2. Stochastic Second-Order method.** The subsampled Newton method takes an
581 additional random set $I_k^H \subset \{1, \dots, N\}$ independent to I_k and compute a search direction as

$$582 \quad \left[\frac{1}{|I_k^H|} \sum_{i \in I_k^H} \nabla^2 f_i(x) \right] d_k = -\frac{1}{|I_k|} \sum_{s_k \in I_k} \nabla f_{s_k}(x_k).$$

583 Therefore, the subspace techniques in [section 2](#) can also be adopted here.

584 Assume that the loss function is the negative logarithm probability associated with a
585 distribution with a density function $p(y|a, x)$ defined by the neural network and parameterized
586 by x . The so-called KFAC method [79] is based on the Kronecker-factored approximate
587 Fisher matrix. Take an L -layer feed-forward neural network for example. Namely, each layer
588 $j \in \{1, 2, \dots, L\}$ is given by

$$589 \quad (4.6) \quad s_j = T_j w_{j-1}, \quad w_j = \psi_j(s_j),$$

590 where $w_0 = a$ is the input of the neural network, $w_L(x) \in \mathbb{R}^m$ is the output of the neural
591 network under the input a , the constant term 1 is not considered in w_{j-1} for simplicity, T_j is
592 the weight matrix and ψ_j is the block-wise activation function. The j th diagonal block of F
593 corresponding to the parameters in the j th layer using a sample set B can be written in the
594 following way:

$$595 \quad (4.7) \quad F^j := Q_{j-1, j-1} \otimes G_{j, j},$$

596 where

$$\begin{aligned}
597 \quad Q_{j-1, j-1} &= \frac{1}{|B|} \sum_{i \in B} w_{j-1}^i (w_{j-1}^i)^\top, \\
G_{j, j} &= \frac{1}{|B|} \sum_{i \in B} \mathbb{E}_{z \sim p(z|a_i, x)} [\tilde{g}_j^i(z) \tilde{g}_j^i(z)^\top],
\end{aligned}$$

598 and $\tilde{g}_j^i(z) := \frac{\partial \ell(\phi(a_i, x), z)}{\partial s_j^i}$. Therefore, the KFAC method computes a search direction in the
599 j th layer from

$$600 \quad F^j d_k^j = -g_k^j,$$

601 where g_k^j is the corresponding subsampled gradient in the j th layer.

602 5. Sparse Optimization.

603 **5.1. Basis Pursuit.** Given a matrix $A \in \mathbb{R}^{m \times n}$ and a vector $b \in \mathbb{R}^m$ such that $m \ll$
604 n , basis pursuit is to find the sparsest signal among all solutions of the equation $Ax = b$. It
605 leads to a NP-hard problem:

$$606 \quad (5.1) \quad \min_x \|x\|_0, \quad \text{s. t. } Ax = b,$$

607 where $\|x\|_0 = |\{j \mid x_j \neq 0\}|$, i.e., the number of the nonzero elements of x . An exact
608 recovery of the sparsest signal often requires the so-called restricted isometry property (RIP),
609 i.e., there exists a constant δ_r such that

$$610 \quad (1 - \delta_r)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_r)\|x\|_2^2, \quad \text{whenever } \|x\|_0 \leq r.$$

611 The greedy pursuit methods build up an approximation in a subspace at the k -th iteration.
612 Let T_k be a subset of $\{1, \dots, n\}$, x^{T_k} be a subvector of x corresponding to the set T_k and
613 A_{T_k} be a column submatrix of A whose column indices are collected in the set T_k . Then the
614 subspace problem is

$$615 \quad x_k^{T_k} = \arg \min_x \frac{1}{2} \|A_{T_k} x - b\|_2^2.$$

616 Clearly, the solution is $x_k^{T_k} = A_{T_k}^\dagger b$ where $A_{T_k}^\dagger$ is the pseudoinverse of A_{T_k} . Since the size of
617 T_k is controlled to be small, A_{T_k} roughly has full rank column due to the RIP property. All
618 other elements of x_k whose indices are in the complementary set of T_k are set to 0.

619 We next explain the choices of T_k . Assume that the initial index set T_0 is empty. The
620 orthogonal matching pursuit (OMP) [116] computes the gradient

$$621 \quad g_k = A^\top (A_{T_k} x_k^{T_k} - b),$$

622 then selects an index such that $t_k = \arg \max_{j=1, \dots, n} |g_j|$. If multiple indices attain the
623 maximum, one can break the tie deterministically or randomly. Then the index set at the next
624 iteration is augmented as

$$625 \quad T_{k+1} = T_k \cup \{t_k\}.$$

626 The CoSaMP [83] method generates an s -sparse solution, i.e., the number of nonzero com-
627 ponents is at most s . Let $(x_k)_s$ be a truncation of x_k such that only the s largest entries in the
628 absolute values are kept and all other elements are set to 0. The support of $(x_k)_s$ is denoted as
629 $\text{supp}((x_k)_s)$. Then a gradient g_k is computed at the point $(x_k)_s$ and the set T_{k+1} is updated
630 by

$$631 \quad T_{k+1} = \text{supp}((g_k)_{2s}) \cup \text{supp}((x_k)_s).$$

632 **5.2. Active Set Methods.** Consider the ℓ_1 -regularized minimization problem

$$633 \quad (5.2) \quad \min_{x \in \mathbb{R}^n} \psi_\mu(x) := \mu \|x\|_1 + f(x),$$

634 where $\mu > 0$ and $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable. The optimality condition of
635 (5.2) is that there exists a vector

$$636 \quad (5.3) \quad (\nabla f(x))^i \begin{cases} = -\mu, & x_i > 0, \\ = +\mu, & x_i < 0, \\ \in [-\mu, \mu], & \text{otherwise.} \end{cases}$$

637 A two-stage active-set algorithm called “FPC_AS” is proposed in [133]. First, an active set
638 is identified by a first-order type method using the so-called “shrinkage” operation. Then,
639 a smooth and smaller subproblem is constructed based on this active set and solved by a
640 second-order type method. These two operations are iterated until convergence criteria are
641 satisfied. While shrinkage is very effective in obtaining a support superset, it can take a lot
642 of steps to find the corresponding values. On the other hand, if one imposes the signs of the
643 components of the variable x that are the same as those of the exact solution, problem (5.2)
644 reduces to a small smooth optimization problem, which can be relatively easily solved to
645 obtain x . Consequently, the key components are the identification of a “good” support set by
646 using shrinkage and the construction of a suitable approximate smooth optimization problem.

647 The iterative shrinkage procedure for solving (5.2) is indeed a proximal gradient method.
648 Given an initial point x_0 , the algorithm iteratively computes

$$649 \quad x_{k+1} := \arg \min_x \quad \mu \|x\|_1 + (x - x_k)^\top g_k + \frac{1}{2\alpha_k} \|x - x_k\|_2^2,$$

650 where $g_k := \nabla f(x_k)$ and $\alpha_k > 0$. A simple calculation shows that

$$651 \quad (5.4) \quad x_{k+1} = \mathcal{S}(x_k - \alpha_k g_k, \mu \alpha_k),$$

652 where for $y \in \mathbb{R}^n$ and $\nu \in \mathbb{R}$, the shrinkage operator is defined as

$$653 \quad \mathcal{S}(y, \nu) = \arg \min_x \nu \|x\|_1 + \frac{1}{2} \|x - y\|_2^2$$

$$654 \quad = \text{sign}(y) \odot \max\{|y| - \nu, \mathbf{0}\}.$$

655 Note that the scheme (5.4) first executes a gradient step with a step size α_k , then performs
656 a shrinkage. In practice, α_k can be computed by a non-monotone line search in which the
657 initial value is set to the BB step size. The convergence of (5.4) has been studied in [53]
658 under suitable conditions on α_k and the Hessian $\nabla^2 f$. An appealing feature proved in [53] is
659 that (5.4) yields the support and the signs of the minimizer x^* of (5.2) after a finite number
660 of steps under favorable conditions. For more references related to shrinkage, the reader is
661 referred to [133].

662 We now describe the subspace optimization in the second stage. For a given vector
663 $x \in \mathbb{R}^n$, the active set is denoted by $\mathcal{A}(x)$ and the inactive set (or support) is denoted by $\mathcal{I}(x)$
664 as follows

$$665 \quad (5.5) \quad \mathcal{A}(x) := \{i \in \{1, \dots, n\} \mid |x^i| = 0\} \text{ and } \mathcal{I}(x) := \{i \in \{1, \dots, n\} \mid |x^i| > 0\}.$$

666 We require that each component x^i either has the same sign as x_k^i or is zero, i.e., x is required
667 to be in the set

$$668 \quad (5.6) \quad \Omega(x_k) := \{x \in \mathbb{R}^n : \text{sign}(x_k^i) x^i \geq 0, i \in \mathcal{I}(x_k) \text{ and } x^i = 0, i \in \mathcal{A}(x_k)\}.$$

669 Then, a smooth subproblem is formulated as either an essentially unconstrained problem

$$670 \quad (5.7) \quad \min_x \quad \mu \text{sign}(x_k^{\mathcal{I}_k})^\top x^{\mathcal{I}_k} + f(x), \text{ s. t. } x^i = 0, i \in \mathcal{A}(x_k)$$

671 or a simple bound-constrained problem

$$672 \quad (5.8) \quad \min_x \quad \mu \text{sign}(x_k^{\mathcal{I}_k})^\top x^{\mathcal{I}_k} + f(x), \text{ s. t. } x \in \Omega(x_k).$$

673 Since the size of the variables in (5.7) and (5.8) is relatively small, these two problems can be
 674 solved efficiently by methods such as L-BFGS-B. If $f(x)$ is quadratic, problem (5.7) can be
 675 solved by the CG method for a system of linear equations.

676 The active set strategies have also been studied in [105, 64]. Specifically, the method
 677 in [64] solves a smooth quadratic subproblem determined by the active sets and invokes a
 678 corrective cycle that greatly improves the efficiency and robustness of the algorithm. The
 679 method is globalized by using a proximal gradient step to check the desired progress.

680 6. The Domain Decomposition Methods.

681 **6.1. A Two-level Subspace Method.** Consider an infinite dimensional minimiza-
 682 tion problem

$$683 (6.1) \quad \min_{\mathbf{x} \in \mathcal{V}} \mathcal{F}(\mathbf{x}),$$

684 where \mathcal{F} is a mapping from an infinite-dimensional space \mathcal{V} to \mathbb{R} . Many large-scale finite
 685 dimensional optimization problems arise from infinite dimensional optimization problems
 686 [28]. Since explicit solutions for these problems are usually not available, we solve the dis-
 687 cretized version of them from the “discretize-then-optimize” strategy by using the concept of
 688 multilevel optimization.

689 Let \mathcal{V}_h be a finite dimensional subset of \mathcal{V} at the grid level h , for example, a standard
 690 finite element space associated with the grid level h . For consecutive coarser levels, we
 691 choose nested spaces, so that $\mathcal{V}_1 \subset \cdots \subset \mathcal{V}_{N-1} \subset \mathcal{V}_N \subset \mathcal{V}$, where N is reserved for the
 692 index of the finest level and 1 for the coarsest level. The functional $\mathcal{F}(\mathbf{x})$ restricted on \mathcal{V}_h is
 693 constructed as $f_h(\mathbf{x}_h)$ for $\mathbf{x}_h \in \mathcal{V}_h$. Therefore, the discretization of problem (6.1) is

$$694 (6.2) \quad \min_{\mathbf{x}_h \in \mathcal{V}_h} f_h(\mathbf{x}_h).$$

695 Let $\mathbf{x}_{h,k}$ be a vector where the first subscript h denotes the discretization level of the
 696 multigrid and the second subscript k denotes the iteration count. We next briefly describe a
 697 two-level subspace method in [24] instead of simply finding a point $\mathbf{x}_{h,k+1}$ in the coarser grid
 698 space \mathcal{V}_H . We seek a point $\mathbf{x}_{h,k+1}$ in $\mathfrak{S}_{h,k} + \mathcal{V}_H$, satisfying some conditions, where $\mathfrak{S}_{h,k}$
 699 is a subspace including the descent information, such as the coordinate direction of the current
 700 iteration and the previous iterations or the gradient $\mathcal{D}_h f(\mathbf{x}_{h,k})$ in the finite space \mathcal{V}_h . Then,
 701 we solve

$$702 (6.3) \quad \mathfrak{S}_{h,k} = \text{span}\{\mathbf{x}_{h,k-1}, \mathbf{x}_{h,k}, \mathcal{D}_h f(\mathbf{x}_{h,k})\} \subseteq \mathcal{V}_h.$$

703 When $\mathbf{x}_{h,k}$ is not optimal on a coarse level $H \in \{1, 2, \dots, N\}$, we go to this level and
 704 compute a new solution $\mathbf{x}_{h,k+1}$ by

$$705 (6.4) \quad \mathbf{x}_{h,k+1} \approx \arg \min f(\mathbf{x}), \quad \text{s. t.} \quad \mathbf{x} \in \mathfrak{S}_{h,k} + \mathcal{V}_H.$$

706 Otherwise, we find a point $\mathbf{x}_{h,k+1} \in \mathcal{V}_h$ on level h .

707 The so-called full multigrid skill or mesh refinement technique can often help to generate
 708 a good initial point so that the total number of iterations is reduced. Firstly, we solve the
 709 target problem at the coarsest level which is computationally cheap. After an approximate
 710 solution \mathbf{x}_h^* at the current level is obtained, we prolongate it to the next finer level $h+1$ with
 711 interpolation as an initial point, and then apply the two-level subspace method at this new
 712 level to find a solution \mathbf{x}_{h+1}^* . This process is repeated until the finest level is reached.

713 **6.2. The Subspace Correction Method.** We next briefly review the subspace cor-
 714 rection methods [112]. Given the current point $\mathbf{x}_{h,k}$, the relaxation (or smoothing) procedure
 715 is the operation on the current level h , namely, to find a direction $d_{h,k}$ to approximate the
 716 solution of

$$717 \quad (6.5) \quad \min_{\mathbf{d} \in \mathcal{V}_h} f_h(\mathbf{x}_{h,k} + \mathbf{d}),$$

718 and the coarse grid correction procedure is the operation on the coarse level H , namely, to
 719 find a direction $\mathbf{d}_{h,k}$ to approximate the solution

$$720 \quad (6.6) \quad \min_{\mathbf{d} \in \mathcal{V}_H} f_h(\mathbf{x}_{h,k} + \mathbf{d}).$$

721 The concept of the subspace correction methods can be used to solve subproblems (6.5)
 722 and (6.6). Let $\{\phi_h^{(j)}\}_{j=1}^{n_h}$ be a basis for \mathcal{V}_h , where n_h is the dimension of \mathcal{V}_h . Denote \mathcal{V}_h
 723 as a direct sum of the one-dimensional subspaces $\mathcal{V}_h = \mathcal{V}_h^{(1)} \oplus \dots \oplus \mathcal{V}_h^{(n_h)}$. Then for each
 724 $j = 1, \dots, n_h$ in turn, we perform the following correction step for subproblem (6.5) at the
 725 k -th iteration:

$$726 \quad (6.7) \quad \begin{cases} \mathbf{d}_{h,k}^{(j)*} = \min_{\mathbf{d}_{h,k}^{(j)} \in \mathcal{V}_h^{(j)}} f_h(\mathbf{x}_{h,k} + \mathbf{d}_{h,k}^{(j)}) \\ \mathbf{x}_{h,k} = \mathbf{x}_{h,k} + \mathbf{d}_{h,k}^{(j)*}. \end{cases}$$

727 For subproblem (6.6), a similar strategy can be adopted by decompose space \mathcal{V}_H into a di-
 728 rect sum. Global convergence of this algorithm has been proved in [113] for strictly convex
 729 functions under some assumptions. The subspace correction method can be viewed as a gen-
 730 eralization of the coordinate search method or the pattern search method.

731 **6.3. Parallel Line Search Subspace Correction Method.** In this subsection, we
 732 consider the following optimization problem

$$733 \quad (6.8) \quad \min_{x \in \mathbb{R}^n} \varphi(x) := f(x) + h(x),$$

734 where $f(x)$ is differentiable convex function and $h(x)$ is a convex function that is possibly
 735 nonsmooth. The ℓ_1 -regularized minimization (LASSO) [114] and the sparse logistic regres-
 736 sion [100] are examples of (6.8). The PSC methods have been studied for LASSO in [36, 39]
 737 and total variation minimization in [37, 38, 39, 68].

738 Suppose that \mathbb{R}^n is split into p subspaces, namely,

$$739 \quad (6.9) \quad \mathbb{R}^n = X^1 + X^2 + \dots + X^p,$$

where

$$X^i = \{x \in \mathbb{R}^n \mid \text{supp}(x) \subset \mathcal{J}_i\}, \quad 1 \leq i \leq p,$$

740 such that $\mathcal{J} := \{1, \dots, n\}$ and $\mathcal{J} = \bigcup_{i=1}^p \mathcal{J}_i$. For any $i \neq j, 1 \leq i, j \leq p$, $\mathcal{J}_i \cap \mathcal{J}_j = \emptyset$ holds in
 741 a non-overlapping domain decomposition of \mathbb{R}^n . Otherwise, there exist $i, j \in \{1, \dots, p\}$ and
 742 $i \neq j$ such that $\mathcal{J}_i \cap \mathcal{J}_j \neq \emptyset$ in an overlapping domain decomposition of \mathbb{R}^n .

743 Let φ_k^i be a surrogate function of φ restricted to the i -th subspace at k -th iteration. The
 744 PSC framework for solving (6.8) is:

$$745 \quad (6.10) \quad d_k^i = \arg \min_{d^i \in X^i} \varphi_k^i(d^i), \quad i = 1, \dots, p,$$

746

$$x_{k+1} = x_k + \sum_{i=1}^p \alpha_k^i d_k^i.$$

747 The convergence can be proved if the step sizes α_k^i ($1 \leq i \leq p$) satisfy the conditions:
 748 $\sum_{i=1}^p \alpha_k^i \leq 1$ and $\alpha_k^i > 0$ ($1 \leq i \leq p$). Usually, the step size α_k^i is quite small under these
 749 conditions and convergence becomes slow. For example, the diminishing step size $\alpha_k^i = \frac{1}{p}$
 750 tends to be smaller and smaller as the number of subspaces increases.

751 A parallel subspace correction method (PSCL) with the Armijo backtracking line search
 752 for a large step size is proposed in [29]. At the k -th iteration, it chooses a surrogate functions
 753 φ_k^i and solves the subproblem (6.10) for each block, then computes a summation of the
 754 direction $d_k = \sum_{i=1}^p d_k^i$. The next iteration is

755

$$x_{k+1} = x_k + \alpha_k d_k,$$

756 where α_k satisfies the Armijo backtracking conditions. When $h(x) = 0$ and $f(x)$ is strongly
 757 convex, the surrogate function can be set to the original objective function φ . Otherwise, it
 758 can be a first-order Taylor expansion of the smooth part $f(x)$ with a proximal term and the
 759 nonsmooth part $h(x)$:

$$760 \quad (6.11) \quad \varphi_k^i(d^i) = \nabla f(x_k)^\top d^i + \frac{1}{2\lambda^i} \|d^i\|_2^2 + h(x_k + d^i), \text{ for } d^i \in X^i.$$

761 Both non-overlapping and overlapping schemes can be designed for PSCL.

762 The directions from different subproblems can be equipped with different step sizes. Let
 763 $Z_k = (d_k^1, d_k^2, \dots, d_k^p)$. The next iteration is set to

764

$$x_{k+1} = x_k + Z_k \alpha_k.$$

765 One can find α_k as an optimal solution of

766

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}^p} \varphi(x_k + Z_k \alpha).$$

767 Alternatively, we can solve the following approximation:

$$768 \quad a_k \approx \arg \min_{\alpha \in \mathbb{R}^p} \nabla f(x_k)^\top Z_k \alpha + \frac{1}{2t_k} \|Z_k \alpha\|_2^2 + h(x_k + Z_k \alpha).$$

769 The global convergence of PSCL is established by following the convergence analysis
 770 of the subspace correction methods for strongly convex problem [112], the active-set method
 771 for l_1 minimization [134] and the BCD method for nonsmooth separable minimization [119].
 772 Specifically, linear convergence rate is proved for the strongly convex case and convergence
 773 to the solution set of problem (6.8) globally is obtained for the general nonsmooth case.

774 **7. General Constrained Optimization.** In this section, we first present subspace
 775 methods for solving general equality constrained optimization problems:

$$776 \quad (7.1) \quad \begin{aligned} & \min_{x \in \mathbb{R}^n} f(x) \\ & \text{s. t. } c(x) = 0, \end{aligned}$$

777 where $c(x) = (c_1(x), \dots, c_m(x))^\top$, $f(x)$ and $c_i(x)$ are real functions defined in \mathbb{R}^n and at
 778 least one of the functions $f(x)$ and $c_i(x)$ is nonlinear. Note that inequality constraints can also
 779 be added to (7.1) but they are omitted to simplify our discussion in the first few subsections.

780 In the last subsection, we discuss methods for bound-constrained minimization problems.
 781 Problem (7.1) is often minimized by computing solutions of a sequence of subproblems which
 782 are simpler than (7.1) itself. However, they are still large-scale linear or quadratic problems
 783 because normally subproblems are also defined in the same dimensional space as the original
 784 nonlinear problem.

785 **7.1. Direct Subspace Techniques.** The sequential quadratic programming (SQP)
 786 is an important method for solving (7.1). It successively minimizes quadratic approximations
 787 to the Lagrangian function subject to the linearized constraints. Let $Q_k(d)$ be a quadratic
 788 approximation to the Lagrangian function of (7.1) at the k -th iteration:

$$789 \quad (7.2) \quad Q_k(d) = g_k^\top d + \frac{1}{2} d^\top B_k d,$$

790 where $g_k = \nabla f(x_k)$ and B_k is an approximation to the Hessian of the Lagrangian function.
 791 The search direction d_k of a line search type SQP method is obtained by solving the following
 792 QP subproblem

$$793 \quad (7.3) \quad \min_{d \in \mathbb{R}^n} Q_k(d)$$

$$794 \quad (7.4) \quad \text{s. t.} \quad c(x_k) + A_k^\top d = 0,$$

796 where $A_k = \nabla c(x_k)$. Although the SQP subproblem is simpler than (7.1), it is still difficult
 797 when the dimension n is large.

798 In general, the subspace SQP method constructs the search direction d_k by solving a QP
 799 in a subspace:

$$800 \quad (7.5) \quad \min Q_k(d)$$

$$\text{s. t.} \quad c_k + A_k^\top d = 0, \quad d \in \mathfrak{S}_k,$$

where \mathfrak{S}_k is a subspace. Lee et al. [70] considered the following choice:

$$\mathfrak{S}_k = \text{span}\{g_k, s_{k-\bar{m}}, \dots, s_{k-1}, \bar{y}_{k-\bar{m}}, \dots, \bar{y}_{k-1}, \nabla c_1(x_k), \dots, \nabla c_m(x_k)\},$$

801 where \bar{m} is the memory size of the limited memory BFGS method for constructing B_k in
 802 (7.2), and \bar{y}_i is a linear combination of y_i and $B_i s_i$. Let U_k be a matrix of linearly independent
 803 vectors in \mathfrak{S}_k . A reduced constrained version of (7.5) is

$$804 \quad (7.6) \quad \min_z (U_k^\top g_k)^\top z + \frac{1}{2} z^\top U_k^\top B_k U_k z$$

$$\text{s. t.} \quad T_k^\top (c_k + A_k^\top U_k z) = 0,$$

805 where T_k is a projection matrix so that the constraints are not over-determined.

806 **7.2. Second Order Correction Steps.** The SQP step d_k can be decomposed into
 807 two parts $d_k = h_k + v_k$ where $v_k \in \text{Range}(A_k)$ and $h_k \in (A_k^\top)^\perp$. Thus, v_k is a solution of
 808 the linearized constrained (7.4) in the range space of A_k , while h_k is the minimizer of the
 809 quadratic function $Q_k(v_k + d)$ in the null space of A_k^\top .

810 One good property of the SQP method is its superlinear convergence rate, namely when
 811 x_k is close to a Karush–Kuhn–Tucker (KKT) point x^* it holds

$$812 \quad (7.7) \quad x_k + d_k - x^* = o(\|x_k - x^*\|).$$

However, a superlinearly convergent step d_k may generate a point that seems “bad” since it may increase both the objective function and the constraint violations. Even though (7.7) holds, the Maratos effect shows that it is possible for the SQP step d_k to have both

$$f(x_k + d_k) > f(x_k), \quad \|c(x_k + d_k)\| > \|c(x_k)\|.$$

813 The second order correction step method [35, 80] solves a QP subproblem whose constraints
814 (7.4) are replaced by

$$815 \quad (7.8) \quad c(x_k + d_k) + A_k^\top (d - d_k) = 0,$$

816 because the left hand side of (7.8) is a better approximation to $c(x_k + d)$ close to the point
817 $d = d_k$. Since the modification of the constraints is a second order term, the new step can
818 be viewed as the SQP step d_k adding a second order correction step \hat{d}_k . Consequently, the
819 Maratos effect is overcome. For detailed discussions on the SQP method and the second
820 order correction step, we refer the reader to [111].

821 We now examine the second order correction step from subspace point of views. It can
822 be verified that the second order correction step \hat{d}_k is a solution of

$$823 \quad \min_{d \in \mathbb{R}^n} Q_k(d_k + d) \\ \text{s. t.} \quad c(x_k + d_k) + A_k^\top d = 0.$$

824 Compute the QR factorization

$$825 \quad A_k = [Y_k, Z_k] \begin{bmatrix} R_k \\ 0 \end{bmatrix}$$

826 and assume that R_k is nonsingular. Therefore, the second order correction step can be repre-
827 sented as $\hat{d}_k = \hat{v}_k + \hat{h}_k$, where $\hat{v}_k = -Y_k R_k^{-T} c(x_k + d_k)$ and \hat{h}_k is the minimizer of

$$828 \quad (7.9) \quad \min_{h \in \text{Null}(A_k^\top)} Q(d_k + \hat{v}_k + h).$$

829 Since d_k is the SQP step, it follows that $g_k + B_k d_k \in \text{Range}(A_k)$, which implies that the
830 minimization problem (7.9) is equivalent to

$$831 \quad (7.10) \quad \min_{h \in \text{Null}(A_k^\top)} \frac{1}{2} (\hat{v}_k + h)^\top B_k (\hat{v}_k + h).$$

832 Examining the SQP method from the perspective of subspace enables us to get more
833 insights. If $Y_k^\top B_k Z_k = 0$, it holds $\hat{h}_k = 0$, which means that the second order correction
834 step $\hat{d}_k \in \text{Range}(A_k)$ is also a range space step. Hence, the second order correction uses
835 two range space steps and one null space step. Note that a range space step is fast since it is
836 a Newton step, while a null space step is normally slower because B_k is often approximated
837 by a quasi-Newton approximation to the Hessian of the Lagrangian function. Intuitively, it
838 might be better to have two slower steps with one fast step. Therefore, it might be reasonable
839 to study a correction step $\hat{d}_k \in \text{Null}(A_k^\top)$ in a modified SQP method.

840 **7.3. The Celis-Dennis-Tapia (CDT) Subproblem.** The CDT subproblem [23] is
841 often needed in some trust region algorithms for constrained optimization. It has two trust
842 region ball constraints:

$$843 \quad (7.11) \quad \min_{d \in \mathbb{R}^n} Q_k(d) = g_k^\top d + \frac{1}{2} d^\top B_k d \\ \text{s. t.} \quad \|c_k + A_k^\top d\|_2 \leq \xi_k, \quad \|d\|_2 \leq \Delta_k,$$

844 where ξ_k and Δ_k are both trust region radii. Let $S_k = \text{span}\{Z_k\}$, $Z_k^\top Z_k = I$, $\text{span}\{A_k, g_k\} \subset$
 845 S_k and $B_k u = \sigma u$, $\forall u \in S_k^\perp$. It is shown in [50] that the CDT subproblem is equivalent to

$$846 \quad \min_{\bar{d} \in \mathbb{R}^r} \quad \bar{Q}_k(\bar{d}) = \bar{g}_k^\top \bar{d} + \frac{1}{2} \bar{d}^\top \bar{B}_k \bar{d}$$

$$\text{s. t.} \quad \|c_k + \bar{A}_k^\top \bar{d}\|_2 \leq \xi_k, \quad \|\bar{d}\|_2 \leq \Delta_k,$$

847 where $\bar{g}_k = Z_k^\top g_k$, $\bar{B}_k = Z_k^\top B_k Z_k$ and $\bar{A}_k = Z_k^\top A_k$. Consequently, a subspace version of
 848 the Powell-Yuan trust algorithm [91] was developed in [50].

849 **7.4. Simple Bound-constrained Problems.** We now consider the optimization
 850 problems with simple bound-constraints:

$$851 \quad (7.12) \quad \min_{x \in \mathbb{R}^n} \quad f(x)$$

$$\text{s. t.} \quad l \leq x \leq u,$$

852 where l and u are two given vectors in \mathbb{R}^n . In this subsection, the superscript of a vector
 853 denotes its indices, for example, x^i is the i th component of x .

854 A subspace adaptation of the Coleman-Li trust region and interior method is proposed in
 855 [12]. The affine scaling matrices D_k and C_k are defined from examining the KKT conditions
 856 of (7.12) as:

$$857 \quad D_k = D(x_k) = \text{diag}(|v(x_k)|^{-1/2}), \quad C_k = D_k \text{diag}(g_k) J_k^v D_k$$

858 where $J^v(x)$ is a diagonal matrix whose diagonal elements equal to zero or ± 1 , and

$$859 \quad v^i = \begin{cases} x^i - u^i, & \text{if } g^i < 0 \text{ and } u^i < +\infty, \\ x^i - l^i, & \text{if } g^i \geq 0 \text{ and } l^i > -\infty, \\ -1, & \text{if } g^i < 0 \text{ and } u^i = +\infty, \\ +1, & \text{if } g^i \geq 0 \text{ and } l^i = -\infty. \end{cases}$$

860 Let H_k be an approximation to the Hessian matrix $\nabla^2 f(x_k)$ and define

$$861 \quad \hat{g}_k = D_k^{-1} g_k, \quad \hat{M}_k = D_k^{-1} H_k D_k^{-1} + \text{diag}(g_k) J_k^v.$$

862 Then the subspace trust region subproblem is

$$863 \quad (7.13) \quad \min_s \quad g_k^\top s + \frac{1}{2} s^\top (H_k + C_k) s$$

$$\text{s. t.} \quad \|D_k s\|_2 \leq \Delta_k, \quad s \in \mathfrak{S}_k.$$

864 If the matrix \hat{M}_k is positive definite, the subspace is taken as

$$865 \quad \mathfrak{S}_k = \text{span}\{D_k^{-2} g_k, w_k\},$$

866 where w_k is either $\hat{s}_k^N = -\hat{M}_k^{-1} \hat{g}_k$ or its inexact version. Otherwise, \mathfrak{S}_k is set to

$$867 \quad \text{span}\{D_k^{-2} \text{sign}(g_k)\} \text{ or } \text{span}\{D_k^{-2} \text{sign}(g_k), w_k\},$$

868 where \hat{w}_k is a vector of nonpositive curvature of \hat{M}_k .

869 A subspace limited memory quasi-Newton method is developed by Ni and Yuan in [87].
870 There are three types of search directions: a subspace quasi-Newton direction, subspace gra-
871 dient and modified gradient directions. The limited memory quasi-Newton method is used
872 to update the variables with indices outside of the active set, while the projected gradient
873 method is used to update the active variables. An active set algorithm is designed in [52].
874 The algorithm consists of a nonmonotone gradient projection step, an unconstrained opti-
875 mization step, and a set of rules for branching between the two steps. After a suitable active
876 set is detected, some components of variables are fixed and the method is switched to the
877 unconstrained optimization algorithm in a lower-dimensional space.

878 **8. Eigenvalue Computation.** The eigenvalue decomposition (EVD) and singular value
879 decomposition (SVD) are fundamental computational tools with extraordinarily wide-ranging
880 applications in science and engineering. For example, most algorithms in high dimensional-
881 ity reduction, such as the principal component analyses (PCA), the multidimensional scaling,
882 Isomap and etc, use them to transform the data into a meaningful representation of reduced
883 dimensionality. More recently, identifying dominant eigenvalue or singular value decom-
884 positions of a sequence of closely related matrices has become an indispensable algorithmic
885 component for many first-order optimization methods for various convex and nonconvex opti-
886 mization problems, such as semidefinite programming, low-rank matrix computation, robust
887 principal component analysis, sparse principal component analysis, sparse inverse covari-
888 ance matrix estimation, nearest correlation matrix estimation and the self-consistent iteration
889 in electronic structure calculation. The computational cost of these decompositions is a major
890 bottleneck which significantly affects the overall efficiency of these algorithms.

891 For a given real symmetric matrix $A \in \mathbb{R}^{n \times n}$, let $\lambda_1, \lambda_2, \dots, \lambda_n$ be the eigenvalues
892 of A sorted in a descending order: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$, and $q_1, \dots, q_n \in \mathbb{R}^n$ be the
893 corresponding eigenvectors such that $Aq_i = \lambda_i q_i$, $\|q_i\|_2 = 1$, $i = 1, \dots, n$ and $q_i^\top q_j = 0$
894 for $i \neq j$. The eigenvalue decomposition of A is defined as $A = Q_n \Lambda_n Q_n^\top$, where, for any
895 integer $i \in [1, n]$,

$$896 \quad (8.1) \quad Q_i = [q_1, q_2, \dots, q_i] \in \mathbb{R}^{n \times i}, \quad \Lambda_i = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_i) \in \mathbb{R}^{i \times i},$$

897 and $\text{Diag}(\cdot)$ denotes a diagonal matrix with its arguments on the diagonal. For simplicity,
898 we also write $A = Q\Lambda Q^\top$ where $Q = Q_n$ and $\Lambda = \Lambda_n$. Without loss of generality, we
899 assume for convenience that A is positive definite (after a shift if necessary). Our task is to
900 compute p largest eigenpairs (Q_p, Λ_p) for some $p \ll n$ where by definition $AQ_p = Q_p \Lambda_p$
901 and $Q_p^\top Q_p = I \in \mathbb{R}^{p \times p}$. Replacing A by a suitable function of A , say $\lambda_1 I - A$, one can also
902 in principle apply the same algorithms to finding p smallest eigenpairs as well.

903 We next describe the Rayleigh-Ritz (RR) step which is to extract approximate eigenpairs,
904 called Ritz-pairs, from the range space $\mathcal{R}(Z)$ spanned by a given matrix $Z \in \mathbb{R}^{n \times m}$. This
905 procedure is widely used as an important component for an approximation to a desired m -
906 dimensional eigenspace of A . It consists of the following four steps.

- 907 (i) Given $Z \in \mathbb{R}^{n \times m}$, orthonormalize Z to obtain $U \in \text{orth}(Z)$, where $\text{orth}(Z)$ is the
908 set of orthonormal bases for the range space of Z
- 909 (ii) Compute $H = U^\top A U \in \mathbb{R}^{m \times m}$, the projection of A onto the range space of U .
- 910 (iii) Compute the eigenvalue decomposition $H = V^\top \Sigma V$, where $V^\top V = I$ and Σ is
911 diagonal.
- 912 (iv) Assemble the Ritz pairs (Y, Σ) where $Y = UV \in \mathbb{R}^{n \times m}$ satisfies $Y^\top Y = I$.

913 The RR procedure is denoted as a map $(Y, \Sigma) = \text{RR}(A, Z)$ where the output (Y, Σ) is a Ritz
914 pair block.

915 **8.1. Classic Subspace Iteration.** The simple (simultaneous) subspace iteration (SSI)
916 method [95, 96, 108, 110] is an extension of the power method which computes a single eigen-

917 pair corresponding to the largest eigenvalue in magnitude. Starting from an initial matrix U ,
 918 SSI repeatedly performs the matrix-matrix multiplications AU , followed by an orthogonal-
 919 ization and RR projection, i.e.,

$$920 \quad (8.2) \quad Z = \text{orth}(AU), \quad U = \text{RR}(A, Z).$$

921 The major purpose of orthogonalization is to guarantee the full-rankness of the matrix Z
 922 since AU may lose rank numerically. The so-called deflation can be executed after each
 923 RR projection to fix the numerically converged eigenvectors since the convergence rates for
 924 different eigenpairs are not the same. Moreover, q extra vectors, often called guard vectors,
 925 are added to U to accelerate convergence. Although the iteration cost is increased at the initial
 926 stage, the overall performance may be better.

927 Due to fast memory access and highly parallelizable computation on modern computer
 928 architectures, simultaneous matrix-block multiplications have advantages over individual matrix-
 929 vector multiplications. Whenever there is a gap between the p -th and the $(p + 1)$ -th eigen-
 930 values of A , the SSI method is ensured to converge to the largest p eigenpairs from any
 931 generic starting point. However, the convergence speed of the SSI method depends critically
 932 on eigenvalue distributions. It can be intolerably slow if the eigenvalue distributions are not
 933 favorable.

934 **8.2. Polynomial Filtering.** The idea of polynomial filtering is originated from a well-
 935 known fact that polynomials are able to manipulate the eigenvalues of any symmetric matrix
 936 A while keeping its eigenvectors unchanged. Due to the eigenvalue decomposition (8.1), it
 937 holds that

$$938 \quad (8.3) \quad \rho(A) = Q\rho(\Lambda)Q^T = \sum_{i=1}^n \rho(\lambda_i)q_iq_i^T,$$

939 where $\rho(\Lambda) = \text{diag}(\rho(\lambda_1), \rho(\lambda_2), \dots, \rho(\lambda_n))$. Ideally, the eigenvalue distribution $\rho(A)$ is
 940 more favorable than the original one.

941 The convergence of the desired eigen-space of SSI is determined by the gap of the eigen-
 942 values, which can be very slow if the gap is nearly zero. Polynomial filtering has been used
 943 to manipulate the gap in eigenvalue computation through various ways [97, 109, 150, 34] in
 944 order to obtain a faster convergence. One popular choice of $\rho(t)$ is the Chebyshev polynomial
 945 of the first kind, which can be written as

$$946 \quad (8.4) \quad T_d(t) = \begin{cases} \cos(d \arccos t) & |t| \leq 1, \\ \frac{1}{2}((t - \sqrt{t^2 - 1})^d + (t + \sqrt{t^2 - 1})^d) & |t| > 1, \end{cases}$$

947 where d is the degree of the polynomial. Since Chebyshev polynomials grow pretty fast
 948 outside the interval $[-1, 1]$, they can help to suppress all unwanted eigenvalues in this interval
 949 efficiently. For these eigenvalues in a general interval $[a, b]$, the polynomial can be chosen as

$$950 \quad (8.5) \quad \rho(t) = T_d\left(\frac{t - (b+a)/2}{(b-a)/2}\right).$$

951 From an initial matrix U , the polynomial-filtered subspace iteration is simply

$$952 \quad (8.6) \quad Z = \text{orth}(\rho(A)U), \quad U = \text{RR}(A, Z).$$

953 **8.3. Limited Memory Methods.** Finding a p -dimensional eigenspace associated with
 954 p largest eigenvalues of A is equivalent to solving a trace maximization problem with orthog-
 955 onality constraints:

$$956 \quad (8.7) \quad \max_{X \in \mathbb{R}^{n \times p}} \operatorname{tr}(X^\top AX), \text{ s. t. } X^\top X = I.$$

957 The first-order optimality conditions of (8.7) are

$$958 \quad AX = X\Lambda, \quad X^\top X = I,$$

959 where $\Lambda = X^\top AX \in \mathbb{R}^{p \times p}$ is the matrix of Lagrangian multipliers. Once the matrix Λ
 960 is diagonalized, the matrix pair (Λ, X) provides p eigenpairs of A . When maximization is
 961 replaced by minimization, (8.7) computes an eigenspace associated with p smallest eigenval-
 962 ues. A few block algorithms have been designed based on solving (8.7), including the locally
 963 optimal block preconditioned conjugate gradient method (LOBPCG) [65] and the limited
 964 memory block Krylov subspace optimization method (LMSVD) [74]. At each iteration, these
 965 methods in fact solve a subspace trace maximization problem of the form

$$966 \quad (8.8) \quad Y = \arg \max_{X \in \mathbb{R}^{n \times p}} \{ \operatorname{tr}(X^\top AX) : X^\top X = I, X \in \mathfrak{S} \}.$$

967 Obviously, the closed-form solution of (8.8) can be obtained by using the RR procedure.

968 The subspace \mathfrak{S} is varied from method to method. In LOBPCG, \mathfrak{S} is the span of the two
 969 most recent iterations X_{i-1} and X_i , and the residual $AX_i - X_i\Lambda_i$ at X_i , which is essentially
 970 equivalent to

$$971 \quad (8.9) \quad \mathfrak{S} = \operatorname{span} \{ X_{i-1}, X_i, AX_i \}.$$

972 The term AX_i can be pre-multiplied by a pre-conditioning matrix if it is available. The
 973 LMSVD method constructs the subspace \mathfrak{S} as a limited memory of the current i -th iterate
 974 and the previous t iterates; i.e.,

$$975 \quad (8.10) \quad \mathfrak{S} = \operatorname{span} \{ X_i, X_{i-1}, \dots, X_{i-t} \}.$$

976 In general, the subspace \mathfrak{S} should be constructed such that the computational cost of solving
 977 the subproblem (8.8) is relatively small.

978 **8.4. Augmented Rayleigh-Ritz Method.** We next introduce the augmented Rayleigh-
 979 Ritz (ARR) procedure. It is easy to see that the RR map $(Y, \Sigma) = \operatorname{RR}(A, Z)$ is equivalent to
 980 solving the trace-maximization subproblem (8.8) with the subspace $\mathfrak{S} = \mathcal{R}(Z)$, while requir-
 981 ing $Y^\top AY$ to be a diagonal matrix Σ . For a fixed number p , the larger the subspace $\mathcal{R}(Z)$
 982 is, the greater chance there is to extract better Ritz pairs. The augmentation of the subspaces
 983 in LOGPCG and LMSVD is the main reason why they generally achieve faster convergence
 984 than the classic SSI.

985 The augmentation in ARR is based on a block Krylov subspace structure, i.e., for some
 986 integer $t \geq 0$,

$$987 \quad (8.11) \quad \mathfrak{S} = \operatorname{span} \{ X, AX, A^2X, \dots, A^tX \}.$$

988 Then the optimal solution of the trace maximization problem (8.8), restricted in the sub-
 989 space \mathfrak{S} in (8.11), is computed via the RR procedure using $(\hat{Y}, \hat{\Sigma}) = \operatorname{RR}(A, K_t)$, where
 990 $K_t = [X, AX, A^2X, \dots, A^tX]$. Finally, the p leading Ritz pairs (Y, Σ) is extracted from
 991 $(\hat{Y}, \hat{\Sigma})$. This augmented RR procedure is simply referred as ARR. It looks identical to a

992 block Lanczos algorithm. However, a fundamental dissimilarity is that the ARR is primarily
 993 developed to compute a relatively large number of eigenpairs by using only a few augmenta-
 994 tion blocks.

995 We next describe an ‘‘Arrabit’’ algorithmic framework with two main steps at each outer
 996 iteration: a subspace update (SU) step and an ARR projection step, for computing a subset
 997 of eigenpairs of large matrices. The goal of the subspace update step is finding a matrix
 998 $X \in \mathbb{R}^{n \times p}$ so that its column space is a good approximation to the p -dimensional eigenspace
 999 spanned by p desired eigenvectors. Once X is obtained, the projection step aims to extract
 1000 from X a set of approximate eigenpairs that are optimal in certain sense. The SU step is
 1001 often performed on a transformed matrix $\rho(A)$, where $\rho(t) : \mathbb{R} \rightarrow \mathbb{R}$ is a suitable polynomial
 1002 function. For a reasonable choice $X \in \mathbb{R}^{n \times p}$, it follows from (8.3) that $\rho(A)X \approx Q_p Q_p^T X$
 1003 would be an approximate basis for the desired eigenspace. The analysis of ARR in [135,
 1004 Corollary 4.6] shows that the convergence rate of SSI is improved from $|\rho(\lambda_{p+1})/\rho(\lambda_p)|$ for
 1005 RR ($t = 0$) to $|\rho(\lambda_{(t+1)p+1})/\rho(\lambda_p)|$ for ARR ($t > 0$). Therefore, a significant improvement
 1006 is possible with a suitably chosen polynomial $\rho(\cdot)$ such that $|\rho(\lambda_{(t+1)p+1})| \ll |\rho(\lambda_{p+1})|$.

1007 In principle, the SU step can be fulfilled by many kinds of updating schemes without
 1008 explicit orthogonalizations. The Gauss-Newton (GN) algorithm in [75] solves the nonlinear
 1009 least squares problem:

$$1010 \quad \min \|XX^\top - A\|_F^2.$$

1011 For any full-rank matrix $X \in \mathbb{R}^{n \times p}$, it takes the simple closed form

$$1012 \quad X_+ = X + \alpha \left(I - \frac{1}{2} X(X^\top X)^{-1} X^\top \right) (AX(X^\top X)^{-1} - X),$$

1013 where the parameter $\alpha > 0$ is a step size. The classic power iteration can be modified without
 1014 orthogonalization at each step. For $X = [x^1 \ x^2 \ \cdots \ x^m] \in \mathbb{R}^{n \times m}$, the power iteration is
 1015 applied individually to all columns of the iterate matrix X , i.e.,

$$1016 \quad x^i = \rho(A)x^i \quad \text{and} \quad x^i = \frac{x^i}{\|x^i\|_2}, \quad i = 1, 2, \dots, m.$$

1017 This scheme is called a multi-power method.

1018 **8.5. Singular Value Decomposition.** Computing the singular value decomposition
 1019 of a real symmetric matrix $A \in \mathbb{R}^{m \times n}$ is equivalent to finding the eigenvalue decomposition
 1020 of AA^\top . Although the methods in the previous subsections can be applied to AA^\top directly,
 1021 the efficiency can be improved when some operations are performed carefully. We first state
 1022 the abstract form of the LMSVD method [74], then describe a few implementation details.

1023 There are two main steps. For a chosen subspace \mathfrak{S}_i with a block Krylov subspace
 1024 structure, an intermediate iterate is computed from

$$1025 \quad (8.12) \quad \hat{X}_i := \arg \max_{X \in \mathbb{R}^{m \times p}} \|A^\top X\|_F^2, \text{ s. t. } X^\top X = I, X \in \mathfrak{S}_i.$$

1026 The next iterate X_{i+1} is generated from a SSI step on \hat{X}_i , i.e.,

$$1027 \quad (8.13) \quad X_{i+1} \in \text{orth} \left(AA^\top \hat{X}_i \right).$$

1028 We collect a limited memory of the last a few iterates in (8.10) into a matrix

$$1029 \quad (8.14) \quad \mathbf{X} = \mathbf{X}_i^t := [X_i, X_{i-1}, \dots, X_{i-t}] \in \mathbb{R}^{m \times q}$$

1030 where $q = (t + 1)p$ is the total number of columns in \mathbf{X}_i^t . For simplicity of notation, the su-
 1031 perscript and subscript of \mathbf{X}_i^t are dropped whenever no confusion would arise. The collection
 1032 matrix \mathbf{X} is written in boldfaces to differentiate it from its blocks. Similarly, a collection of
 1033 matrix-vector multiplications from the SSI steps are saved in

$$1034 \quad (8.15) \quad \mathbf{Y} = \mathbf{Y}_i^t := A^\top \mathbf{X}_i^t := [A^\top X_i, A^\top X_{i-1}, \dots, A^\top X_{i-t}] \in \mathbb{R}^{m \times q}.$$

1035 Assume that \mathbf{X} has a full rank and this assumption will be relaxed later. A stable ap-
 1036 proach for solving (8.12) is to find an orthonormal basis for \mathfrak{S}_i , say,

$$1037 \quad \mathbf{Q} = \mathbf{Q}_i^t \in \text{orth}(\mathbf{X}_i^t).$$

1038 Note that $X \in \mathfrak{S}_i$ if and only if $X = \mathbf{Q}V$ for some $V \in \mathbb{R}^{q \times p}$. The generalized eigenvalue
 1039 problem (8.12) is converted into an equivalent eigenvalue problem

$$1040 \quad (8.16) \quad \max_{V \in \mathbb{R}^{q \times p}} \|\mathbf{R}V\|_{\mathbb{F}}^2, \text{ s. t. } V^\top V = I,$$

1041 where

$$1042 \quad (8.17) \quad \mathbf{R} = \mathbf{R}_i^t := A^\top \mathbf{Q}_i^t.$$

1043 The matrix product \mathbf{R} in (8.17) can be computed from historical information without any
 1044 additional computation involving the matrix A . Since $\mathbf{Q} \in \text{orth}(\mathbf{X})$ and \mathbf{X} has a full rank,
 1045 there exists a nonsingular matrix $C \in \mathbb{R}^{q \times q}$ such that $\mathbf{X} = \mathbf{Q}C$. Therefore, $\mathbf{Q} = \mathbf{X}C^{-1}$, and
 1046 \mathbf{R} in (8.17) can be assembled as

$$1047 \quad (8.18) \quad \mathbf{R} = A^\top \mathbf{Q} = (A^\top \mathbf{X})C^{-1} = \mathbf{Y}C^{-1},$$

1048 where $\mathbf{Y} = A^\top \mathbf{X}$ is accessible from our limited memory. Once \mathbf{R} is available, a solution \hat{V}
 1049 to (8.16) can be computed from the p leading eigenvectors of the $q \times q$ matrix $\mathbf{R}^\top \mathbf{R}$. The
 1050 matrix product can then be calculated as

$$1051 \quad (8.19) \quad AA^\top \hat{X}_i = \mathbf{R} \hat{V} = \mathbf{Y}C^{-1} \hat{V}.$$

1052 We now explain how to efficiently and stably compute \mathbf{Q} and \mathbf{R} when the matrix \mathbf{X} is
 1053 numerically rank deficient. Since each block itself in \mathbf{X} is orthonormal, keeping the latest
 1054 block X_i intact and projecting the rest of the blocks onto the null space of X_i^\top yields

$$1055 \quad (8.20) \quad \mathbf{P}_X = \mathbf{P}_i^X := (I - X_i X_i^\top) [X_{i-1} \ \dots \ X_{i-p}].$$

1056 An orthonormalization of \mathbf{P}_X is performed via the eigenvalue decomposition of its Gram
 1057 matrix

$$1058 \quad (8.21) \quad \mathbf{P}_X^\top \mathbf{P}_X = U_X \Lambda_X U_X^\top,$$

1059 where U_X is orthogonal and Λ_X is diagonal. If Λ_X is invertible, it holds

$$1060 \quad (8.22) \quad \mathbf{Q} = \mathbf{Q}_i^t := \left[X_i, \mathbf{P}_X U_X \Lambda_X^{-\frac{1}{2}} \right] \in \text{orth}(\mathbf{X}_i^t).$$

1061 The above procedure can be stabilized by deleting the columns of \mathbf{P}_X whose Euclidean
 1062 norms are below a threshold or deleting the small eigenvalues in Λ_X and the corresponding
 1063 columns in U_X . The same notations are still used for \mathbf{P}_X , U_X and Λ_X after these possible
 1064 deletions. Therefore, a stable construction of \mathbf{Q} is still provided by formula (8.22) and the
 1065 corresponding \mathbf{R} matrix can be formulated as

$$1066 \quad (8.23) \quad \mathbf{R} = \mathbf{R}_i^t := \left[Y_i, \mathbf{P}_Y U_X \Lambda_X^{-\frac{1}{2}} \right],$$

1067 where $\mathbf{P}_Y = \mathbf{P}_i^Y := A^\top \mathbf{P}_X$ before the stabilization procedure but some of the columns of
 1068 \mathbf{P}_Y may have been deleted due to the stabilization steps. Therefore, the \mathbf{R} matrix in (8.23) is
 1069 well defined as is the \mathbf{Q} matrix in (8.22) after the numerical rank deficiency is removed. ■

1070 **8.6. Randomized SVD.** Given an $m \times n$ matrix A and an integer $p < \min(m, n)$, we
 1071 want to find an orthonormal $m \times p$ matrix Q such that

$$1072 \quad A \approx QQ^T A.$$

1073 A prototype randomized SVD in [54] is essentially one step of the Power method using an
 1074 initial random input. We select an oversampling parameter $l \geq 2$ and an exponent t (for
 1075 example, $t = 1$ or $t = 2$), then perform the following steps.

- 1076 • Generate an $n \times (p + l)$ Gaussian matrix Ω .
- 1077 • Compute $Y = (AA^T)^t A\Omega$ by the multiplications of A and A^T alternatively.
- 1078 • Construct a matrix $Q = \text{orth}(Y)$ by the QR factorization.
- 1079 • Form the matrix $B = Q^T A$.
- 1080 • Calculate an SVD of B to obtain $B = \tilde{U}\Sigma V^T$, and set $U = Q\tilde{U}$.

1081 Consequently, we have the approximation $A \approx U\Sigma V^T$. For the eigenvalue computation, we
 1082 can simply run the SSI (8.2) for only one step with an Gaussian matrix U . Assume that the
 1083 computation is performed in exact arithmetic. It is shown in [54] that

$$1084 \quad \mathbb{E}\|A - QQ^T A\|_2 \leq \left[1 + \frac{4\sqrt{p+l}}{l-1}\right] \sigma_{p+1},$$

1085 where the expectation is taken with respect to the random matrix Ω and σ_{p+1} is the $(p+1)$ -th
 1086 largest singular value of A .

1087 Suppose that a low rank approximation of A with a target rank r is needed. A sketching
 1088 method is further developed in [118] for selected p and ℓ . Again, we draw independent
 1089 Gaussian matrix $\Omega \in \mathbb{R}^{n \times p}$ and $\Psi \in \mathbb{R}^{\ell \times m}$, and compute the matrix-matrix multiplications:

$$1090 \quad Y = A\Omega, \quad W = \Psi A,$$

1091 Then an approximation \hat{A} is computed:

- 1092 • Calculate an orthogonal-triangular factorization $Y = QR$ where $Q \in \mathbb{R}^{m \times p}$.
- 1093 • Compute a least-squares problem to derive $X = (\Psi Q)^\dagger W \in \mathbb{R}^{p \times n}$
- 1094 • Assemble the rank- p approximation $\hat{A} = QX$

1095 Assume that $p = 2r + 1$ and $\ell = 4r + 2$. It is established that

$$1096 \quad \mathbb{E}\|A - \hat{A}\|_F \leq 2 \min_{\text{rank}(Z) \leq r} \|A - Z\|_F.$$

1097 **8.7. Truncated Subspace Method for Tensor Train.** In this subsection, we con-
 1098 sider the trace maximization problem (8.7) whose dimension reaches the magnitude of $O(10^{42})$.
 1099 Due to the scale of data storage, a tensor train (TT) format is used to express data matrices
 1100 and eigenvectors in [148]. The corresponding eigenvalue problem can be solved based on the
 1101 subspace algorithm and the alternating direction method with suitable truncations.

1102 The goal is to express a vector $x \in \mathbb{R}^n$ as a tensor $\mathbf{x} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ for some positive
 1103 integers n_1, \dots, n_d such that $n = n_1 n_2 \dots n_d$ using a collection of three-dimensional tensor
 1104 cores $\mathbf{X}_\mu \in \mathbb{R}^{r_{\mu-1} \times r_\mu \times n_\mu}$ with fixed dimensions r_μ , $\mu = 1, \dots, d$ and $r_0 = r_d = 1$. A
 1105 tensor \mathbf{x} is stored in the TT format if its elements can be written as

$$1106 \quad x_{i_1 i_2 \dots i_d} = X_1(i_1)X_2(i_2) \dots X_d(i_d),$$

1107 where $X_\mu(i_\mu) \in \mathbb{R}^{r_{\mu-1} \times r_\mu}$ is the i_μ -th slice of \mathbf{X}_μ for $i_\mu = 1, 2, \dots, n_\mu$. The values r_μ are
 1108 often equal to a constant r , which is then called the TT-rank. Consequently, storing a vector
 1109 $x \in \mathbb{R}^{n_1}$ only needs $\mathcal{O}(dn_1 r^2)$ entries if the corresponding tensor \mathbf{x} has a TT format. The
 1110 representation of \mathbf{x} is shown as graphs in Figure 8.1.

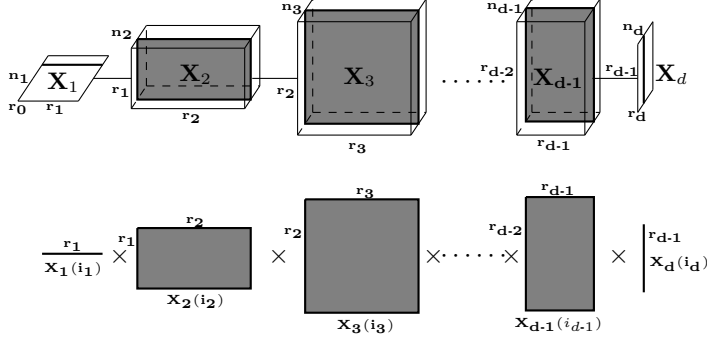


Fig. 8.1 The first row is a TT format of \mathbf{u} with cores \mathbf{X}_μ , $\mu = 1, 2, \dots, d$. The second row is a representation of its elements $x_{i_1 i_2 \dots i_d}$.

1111 There are several ways to express a matrix $X \in \mathbb{R}^{n \times p}$ with $p \ll n$ in the TT format. A
 1112 direct way is to store each column of X as tensors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ in the TT format separately.
 1113 Another economic choice is that these p tensors share all except one core. Let the shared
 1114 cores be \mathbf{X}_i , $i \neq \mu$ and the μ -th core of \mathbf{x}_i be $\mathbf{X}_{\mu,i}$, for $i = 1, 2, \dots, p$. Then the $i_1 i_2 \dots i_d$
 1115 component of \mathbf{x}_j is

$$1116 \quad (8.24) \quad X(i_1, \dots, i_\mu, \dots, i_d; j) = X_1(i_1) \cdots X_{\mu,j}(i_\mu) \cdots X_d(i_d).$$

1117 The above scheme generates a block- μ TT (μ -BTT) format, which is depicted in **Figure 8.2**.
 1118 Similarly, a matrix $A \in \mathbb{R}^{n \times n}$ is in an operator TT format \mathbf{A} if the components of A can be
 1119 assembled as

$$1120 \quad (8.25) \quad A_{i_1 i_2 \dots i_d, j_1 j_2 \dots j_d} = A_1(i_1, j_1) A_2(i_2, j_2) \cdots A_d(i_d, j_d),$$

1121 where $A_\mu(i_\mu, j_\mu) \in \mathbb{R}^{r_{\mu-1} \times r_\mu}$ for $i_\mu, j_\mu \in \{1, \dots, n_\mu\}$.

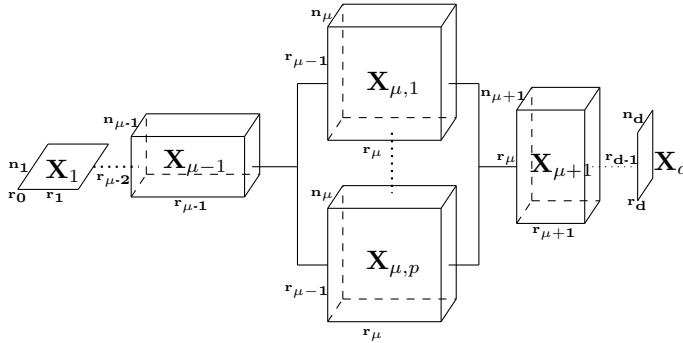


Fig. 8.2 Demonstration of a μ -BTT format.

1122 Assume that the matrix A itself can be written in the operator TT format. Let $X \in \mathbb{R}^{n \times p}$
 1123 with $n = n_1 n_2 \dots n_d$ whose BTT format is \mathbf{X} , and $\mathbf{T}_{n,r,p}$ be the set of the BTT formats
 1124 whose TT-ranks are no more than r . Then the eigenvalue problem in the block BTT format is

$$1125 \quad (8.26) \quad \min_{X \in \mathbb{R}^{n \times p}} \text{tr}(\mathbf{X}^\top \mathbf{A} \mathbf{X}), \quad \text{s. t.} \quad \mathbf{X}^\top \mathbf{X} = I_p \text{ and } \mathbf{X} \in \mathbf{T}_{n,r,p},$$

1126 where $\mathbf{X} \in \mathbf{T}_{\mathbf{n},r,p}$ means that all calculations are performed in the BTT format. Since
 1127 the TT-ranks increase dramatically after operations such as the addition and matrix-vector
 1128 multiplication in the TT formats, the computational cost and the storage becomes more and
 1129 more expensive as the TT-ranks increase. Therefore, the subspace methods in [subsection 8.3](#)
 1130 can only be applied with projections to $\mathbf{T}_{\mathbf{n},r,p}$ at some suitable places so that the overall
 1131 computational cost is still tractable.

1132 In our truncated subspace optimization methods, solving the subproblem (8.8) is split
 1133 into a few steps. First, the subspace \mathfrak{S}_k is modified with truncations so that the computation of
 1134 the coefficient matrix $U^\top AU$ in the RR procedure is affordable. Let $\mathcal{P}_{\mathbf{T}}(\mathbf{X})$ be the truncation
 1135 of \mathbf{X} to the BTT format $\mathbf{T}_{\mathbf{n},r,p}$. One can choose either the following subspace

$$1136 \quad (8.27) \quad \mathfrak{S}_k^{\mathbf{T}} = \text{span}\{\mathcal{P}_{\mathbf{T}}(\mathbf{A}\mathbf{X}_k), \mathbf{X}_k, \mathbf{X}_{k-1}\},$$

1137 or a subspace similar to that of LOBPCG with two truncations as

$$1138 \quad (8.28) \quad \mathfrak{S}_k^{\mathbf{T}} = \text{span}\{\mathbf{X}_k, \mathcal{P}_{\mathbf{T}}(\mathbf{R}_k), \mathcal{P}_{\mathbf{T}}(\mathbf{P}_k)\},$$

1139 where the conjugate gradient direction is $\mathbf{P}_k = X_k - X_{k-1}$ and the residual vector is $\mathbf{R}_k =$
 1140 $\mathbf{A}\mathbf{X}_k - X_k\Lambda_k$.

1141 Consequently, the subspace problem in the BTT format is

$$1142 \quad (8.29) \quad \mathbf{Y}_{k+1} := \arg \min_{\mathbf{X} \in \mathbb{R}^{n \times p}} \text{tr}(\mathbf{X}^\top \mathbf{A}\mathbf{X}), \text{ s. t. } \mathbf{X}^\top \mathbf{X} = I_p, \mathbf{X} \in \mathfrak{S}_k^{\mathbf{T}},$$

1143 which is equivalent to a generalized eigenvalue decomposition problem:

$$1144 \quad (8.30) \quad \min_{V \in \mathbb{R}^{q \times p}} \text{tr}(V^\top (S^\top \mathbf{A}S)V), \text{ s. t. } V^\top S^\top S V = I_p.$$

1145 Note that $\mathbf{Y}_{k+1} \notin \mathbf{T}_{\mathbf{n},r,p}$ because the rank of \mathbf{Y}_{k+1} is larger than r due to several additions
 1146 between the BTT formats. Since \mathbf{Y}_{k+1} is a linear combination of the BTT formats in $\mathfrak{S}_k^{\mathbf{T}}$,
 1147 problem (8.29) still can be solved easily but only the coefficients of the linear combinations
 1148 are stored.

1149 We next project \mathbf{Y}_{k+1} to the required space $\mathbf{T}_{\mathbf{n},r,p}$ as

$$1150 \quad (8.31) \quad \mathbf{X}_{k+1} = \arg \min_{\mathbf{X} \in \mathbb{R}^{n \times p}} \|\mathbf{X} - \mathbf{Y}_{k+1}\|_F^2, \text{ s. t. } \mathbf{X}^\top \mathbf{X} = I_p, \mathbf{X} \in \mathbf{T}_{\mathbf{n},r,p}.$$

1151 This problem can be solved by using the alternating minimization scheme. By fixing all
 1152 except the μ th core, we obtain

$$1153 \quad (8.32) \quad \min_V \|\mathcal{X}_{\neq \mu}^\top V - \text{vec}(\mathbf{Y}_{k+1})\|_F^2, \text{ s. t. } V^\top \mathcal{X}_{\neq \mu}^\top \mathcal{X}_{\neq \mu} V = I_p,$$

1154 where

$$1155 \quad \mathcal{X}_{\neq \mu} := (\mathbf{X}_{\geq \mu+1} \otimes I_{n_\mu} \otimes \mathbf{X}_{\leq \mu-1}),$$

1156 and

$$1157 \quad X_{\leq \mu} = [X_1(i_1)X_2(i_2) \cdots X_\mu(i_\mu)] \in \mathbb{R}^{n_1 n_2 \cdots n_\mu \times r_\mu},$$

$$1158 \quad X_{\geq \mu} = [X_\mu(i_\mu)X_{\mu+1}(i_{\mu+1}) \cdots X_d(i_d)]^\top \in \mathbb{R}^{n_\mu n_{\mu+1} \cdots n_d \times r_{\mu-1}}.$$

1159 Therefore, after imposing orthogonality on $\mathcal{X}_{\neq \mu}$, (8.32) is reformulated as

$$1160 \quad (8.33) \quad \min_V \|V - \mathcal{X}_{\neq \mu}^\top \text{vec}(\mathbf{Y}_{k+1})\|_F^2, \text{ s. t. } V^\top V = I_p,$$

1161 whose optimal solution can be computed by the p -dominant SVD of $\mathcal{X}_{\neq \mu}^\top \text{vec}(\mathbf{Y}_{k+1})$.

1162 **9. Optimization with Orthogonality Constraints.** In this section, we consider the
 1163 optimization problem with orthogonality constraints [132, 59, 2]:

$$1164 \quad (9.1) \quad \min_{X \in \mathbb{C}^{n \times p}} f(X) \quad \text{s. t.} \quad X^* X = I_p,$$

1165 where $f(X) : \mathbb{C}^{n \times p} \rightarrow \mathbb{R}$ is a \mathbb{R} -differentiable function [67]. The set $\text{St}(n, p) := \{X \in$
 1166 $\mathbb{C}^{n \times p} : X^* X = I_p\}$ is called the Stiefel manifold. Obviously, the eigenvalue problem in sec-
 1167 tion 8 is a special case of (9.1). Other important applications include the density functional
 1168 theory [131], Bose-Einstein condensates [137], low rank nearest correlation matrix comple-
 1169 tion [121], and etc. Although (9.1) can be treated from the perspective of general nonlinear
 1170 programming, the intrinsic structure of the Stiefel manifold enables us to develop more effi-
 1171 cient algorithms. In fact, it can be solved by the Riemannian gradient descent, Riemannian
 1172 conjugate gradient, proximal Riemannian gradient methods [40, 104, 2, 59]. The Riemann-
 1173 nian Newton, trust-region, adaptive regularized Newton methods [120, 1, 2, 59] can used
 1174 when the Hessian information is available. Otherwise, the quasi-Newton types methods are
 1175 good alternatives [62, 61, 58].

1176 The tangent space is $T_X := \{\xi \in \mathbb{C}^{n \times p} : X^* \xi + \xi^* X = 0\}$. The operator $\text{Proj}_X(Z) :=$
 1177 $Z - X \text{sym}(X^* Z)$ is the projection of Z onto the tangent space T_X and $\text{sym}(A) := (A +$
 1178 $A^*)/2$. The symbols $\nabla f(X)$ ($\nabla^2 f(X)$) and $\text{grad } f(X)$ ($\text{Hess } f(X)$) denote the Euclidean
 1179 and Riemannian gradient (Hessian) of f at X . Using the real part of the Frobenius inner
 1180 product $\Re \langle A, B \rangle$ as the Euclidean metric, the Riemannian Hessian $\text{Hess } f(X)$ [31, 3] can be
 1181 written as

$$1182 \quad (9.2) \quad \text{Hess } f(X)[\xi] = \text{Proj}_X(\nabla^2 f(X)[\xi] - \xi \text{sym}(X^* \nabla f(X))),$$

1183 where ξ is any tangent vector in T_X . A retraction R is a smooth mapping from the tangent
 1184 bundle to the manifold. Moreover, the restriction R_X of R to T_X has to satisfy $R_X(0_X) = X$
 1185 and $DR_X(0_X) = \text{id}_{T_X}$, where id_{T_X} is the identity mapping on T_X .

1186 **9.1. Regularized Newton Type Approaches.** We now describe an adaptively reg-
 1187 ularized Riemannian Newton type method with a subspace refinement procedure [59, 58].
 1188 Note that the Riemannian Hessian-vector multiplication (9.2) involves the Euclidean Hessian
 1189 and gradient with simple structures. We construct a second-order Taylor approximation in the
 1190 Euclidean space rather than the Riemannian space at the k -th iteration:

$$1191 \quad (9.3) \quad m_k(X) := \Re \langle \nabla f(X_k), X - X_k \rangle + \frac{1}{2} \Re \langle \mathcal{B}_k[X - X_k], X - X_k \rangle + \frac{\tau_k}{2} \|X - X_k\|_F^2,$$

1192 where \mathcal{B}_k is either $\nabla^2 f(X_k)$ or its approximation based on whether $\nabla^2 f(X_k)$ is affordable
 1193 or not, and τ_k is a regularization parameter to control the distance between X and X_k . Then
 1194 the subproblem is

$$1195 \quad (9.4) \quad \min_{X \in \mathbb{C}^{n \times p}} m_k(X) \quad \text{s. t.} \quad X^* X = I.$$

1196 After obtaining an approximate solution Z_k of (9.4), we calculate a ratio between the pre-
 1197 dicted reduction and the actual reduction, then use the ratio to decide whether X_{k+1} is set to
 1198 Z_k or X_k and to adjust the parameter τ_k similar to the trust region methods.

1199 In particular, the model (9.4) can be minimized by using a modified CG method to solve
 1200 a single Riemannian Newton system:

$$1201 \quad (9.5) \quad \text{grad } m_k(X_k) + \text{Hess } m_k(X_k)[\xi] = 0.$$

1202 A simple calculation yields:

$$1203 \quad (9.6) \quad \text{Hess } m_k(X_k)[\xi] = \text{Proj}_{X_k}(\mathcal{B}_k[\xi] - \xi \text{sym}((X_k)^* \nabla f(X_k)) + \tau_k \xi), \quad \xi \in T_{X_k}.$$

1204 Hence, the regularization term shifts the spectrum of the Riemannian Hessian by τ_k . The
 1205 modified CG method is a direct adaption of the truncated CG method for solving the classic
 1206 trust region subproblem, see [88, Chapter 5] and [2, Chapter 7] for a comparison. It is ter-
 1207 minated when either the residual becomes small or a negative curvature is detected since the
 1208 Hessian may be indefinite. During the process, two different vectors s_k and d_k are generated,
 1209 where the vector d_k represents the negative curvature direction and s_k corresponds to the con-
 1210 jugate direction from the CG iteration. The direction d_k is zero unless a negative curvature is
 1211 detected. Therefore, a possible choice of the search direction ξ_k is

$$1212 \quad (9.7) \quad \xi_k = \begin{cases} s_k + \tau_k d_k & \text{if } d_k \neq 0, \\ s_k & \text{if } d_k = 0, \end{cases} \quad \text{with} \quad \tau_k := \frac{\langle d_k, \text{grad } m_k(X_k) \rangle}{\langle d_k, \text{Hess } m_k(X_k)[d_k] \rangle}.$$

1213 Once the direction ξ_k is computed, a trial point Z_k is searched along ξ_k followed by a retrac-
 1214 tion, i.e.,

$$1215 \quad (9.8) \quad Z_k = R_{X_k}(\alpha_k \xi_k).$$

1216 The step size $\alpha_k = \alpha_0 \delta^h$ is chosen by the Armijo rule such that h is the smallest integer
 1217 satisfying

$$1218 \quad (9.9) \quad m_k(R_{X_k}(\alpha_0 \delta^h \xi_k)) \leq \rho \alpha_0 \delta^h \langle \text{grad } m_k(X_k), \xi_k \rangle,$$

1219 where $\rho, \delta \in (0, 1)$ and $\alpha_0 \in (0, 1]$ are given constants.

1220 The performance of the Newton-type method may be seriously deteriorated when the
 1221 Hessian is close to be singular. One reason is that the Riemannian Newton direction is nearly
 1222 parallel to the negative gradient direction. Consequently, the next iteration X_{k+1} very likely
 1223 belongs to the subspace $\text{span}\{X_k, \text{grad } f(X_k)\}$, which is similar to the Riemannian gradient
 1224 approach. To overcome the numerical difficulty, we can further solve (9.1) in a restricted
 1225 subspace. Specifically, a q -dimensional subspace \mathfrak{S}_k is constructed with an orthogonal basis
 1226 $Q_k \in \mathbb{C}^{n \times q}$ ($p \leq q \leq n$). Then the representation of any point X in the subspace \mathfrak{S}_k is

$$1227 \quad X = Q_k M$$

1228 for some $M \in \mathbb{C}^{q \times p}$. In a similar fashion to these constructions for the linear eigen-
 1229 value problems in section 8, the subspace \mathfrak{S}_k can be built by using the history information
 1230 $\{X_k, X_{k-1}, \dots\}$, $\{\text{grad } f(X_k), \text{grad } f(X_{k-1}), \dots\}$ and other useful information. Once a
 1231 subspace \mathfrak{S}_k is given, (9.1) with an additional constraint $X \in \mathfrak{S}_k$ becomes

$$1232 \quad (9.10) \quad \min_{M \in \mathbb{C}^{q \times p}} f(Q_k M) \quad \text{s. t.} \quad M^* M = I_p.$$

1233 Suppose that M_k is an inexact solution of the problem (9.10) from existing optimization
 1234 methods on manifold. Then $X_{k+1} = Q_k M_k$ is a better point than X_k . For extremely difficult
 1235 problems, one may alternate between the Newton type method and the subspace refinement
 1236 procedure for a few cycles.

1237 **9.2. A Structured Quasi-Newton Update with Nyström Approximation.** The
 1238 secant condition in the classical quasi-Newton methods for constructing the quasi-Newton
 1239 matrix \mathcal{B}_k

$$1240 \quad (9.11) \quad \mathcal{B}_k[S_k] = \nabla f(X_k) - \nabla f(X_{k-1}),$$

1241 where

$$1242 \quad S_k := X_k - X_{k-1}.$$

1243 Assume that the Euclidean Hessian $\nabla^2 f(X)$ is a summation of a relatively cheap part $\mathcal{H}^c(X)$
1244 and a relatively expensive or even inaccessible part $\mathcal{H}^e(X)$, i.e.,

$$1245 \quad (9.12) \quad \nabla^2 f(X) = \mathcal{H}^c(X) + \mathcal{H}^e(X).$$

1246 Then it is reasonable to keep the cheaper part $\mathcal{H}^c(X)$ but approximate $\mathcal{H}^e(X)$ using the
1247 quasi-Newton update \mathcal{E}_k . It yields an approximation \mathcal{B}_k to the Hessian $\nabla^2 f(X_k)$ as

$$1248 \quad (9.13) \quad \mathcal{B}_k = \mathcal{H}^c(X_k) + \mathcal{E}_k,$$

1249 Plugging (9.13) into (9.11) gives the following revised secant condition

$$1250 \quad (9.14) \quad \mathcal{E}_k[S_k] = Y_k,$$

1251 where

$$1252 \quad (9.15) \quad Y_k := \nabla f(X_k) - \nabla f(X_{k-1}) - \mathcal{H}^c(X_k)[S_k].$$

1253 A good initial matrix \mathcal{E}_k^0 to \mathcal{E}_k is important to ensure the convergence speed of the limited-
1254 memory quasi-Newton method. We assume that a known matrix $\hat{\mathcal{E}}_k^0$ can approximate the ex-
1255 pensive part of the Hessian $\mathcal{H}^e(X_k)$ well, a very limited number of matrix-matrix products
1256 involving $\hat{\mathcal{E}}_k^0$ is affordable but many of them are still prohibitive. We next use the Nyström
1257 approximation [117] to construct a low rank matrix. Let Ω be a matrix whose columns consti-
1258 tute an orthogonal basis of a well-chosen subspace \mathfrak{S} and denote $W = \hat{\mathcal{E}}_k^0[\Omega]$. The Nyström
1259 approximation is

$$1260 \quad (9.16) \quad \mathcal{E}_k^0[U] := W(W^*\Omega)^\dagger W^*U,$$

where $U \in \mathbb{C}^{n \times p}$ is any direction. When the dimension of the subspace \mathfrak{S} is small enough,
the rank of $W(W^*\Omega)^\dagger W^*$ is also small so that the computational cost of $\mathcal{E}_k^0[U]$ is significantly
cheaper than the original $\hat{\mathcal{E}}_k^0[U]$. Suppose the subspace \mathfrak{S} is chosen as

$$\text{span}\{X_{k-1}, X_k\},$$

1261 which contains the element S_k . If $\hat{\mathcal{E}}_k^0[UV] = \hat{\mathcal{E}}_k^0[U]V$ for any matrices U, V with proper
1262 dimension (this condition is satisfied when $\hat{\mathcal{E}}_k^0$ is a matrix), then the secant condition still
1263 holds at \mathcal{E}_k^0 , i.e.,

$$1264 \quad \mathcal{E}_k^0[S_k] = Y_k.$$

1265 The subspace \mathfrak{S} can also be defined as

$$1266 \quad (9.17) \quad \text{span}\{X_{k-1}, X_k, \hat{\mathcal{E}}_k^0[X_k]\} \quad \text{or} \quad \text{span}\{X_{k-h}, \dots, X_{k-1}, X_k\}$$

1267 with small memory length h . Consequently, we obtain a limited-memory Nyström approxi-
1268 mation.

1269 **9.3. Electronic Structure Calculations.** The density functional theory (DFT) in
1270 electronic structure calculation is an important source of optimization problems with orthog-
1271 onality constraints. By abuse of notation, we refer to Kohn-Sham (KS) equations with lo-
1272 cal or semi-local exchange-correlation functionals as KSDFT, and KS equations with hybrid
1273 functionals as HF (short for Hartree-Fock). The KS/HF equations try to identify orthogo-
1274 nal eigenvectors to satisfy the nonlinear eigenvalue problems, while the KS/HF minimization
1275 problem minimizes the KS/HF total energy functionals under the orthogonality constraints.
1276 These two problems are connected by the optimality conditions. ■

1277 **9.3.1. The Mathematical Models.** The wave functions of p occupied states can be
 1278 expressed as $X = [x_1, \dots, x_p] \in \mathbb{C}^{n \times p}$ with $X^*X = I_p$ after some suitable discretization.
 1279 The KS total energy functional is defined as

(9.18)

$$1280 E_{\text{ks}}(X) := \frac{1}{4} \text{tr}(X^* L X) + \frac{1}{2} \text{tr}(X^* V_{\text{ion}} X) + \frac{1}{2} \sum_l \sum_i \zeta_l |x_i^* w_l|^2 + \frac{1}{4} \rho^\top L^\dagger \rho + \frac{1}{2} e_n^\top \epsilon_{\text{xc}}(\rho),$$

1281 where L is a discretized Laplacian operator, the charge density is $\rho(X) = \text{diag}(X X^*)$, V_{ion}
 1282 is the constant ionic pseudopotentials, w_l represents a discretized pseudopotential reference
 1283 projection function, ζ_l is a constant whose value is ± 1 , and ϵ_{xc} is related to the exchange
 1284 correlation energy. The Fock exchange operator $\mathcal{V}(\cdot) : \mathbb{C}^{n \times n} \rightarrow \mathbb{C}^{n \times n}$ is usually a fourth-
 1285 order tensor [69] which satisfies the following properties: (i) $\langle \mathcal{V}(D_1), D_2 \rangle = \langle \mathcal{V}(D_2), D_1 \rangle$
 1286 for any $D_1, D_2 \in \mathbb{C}^{n \times n}$; (ii) $\mathcal{V}(D)$ is Hermitian if D is Hermitian. Then the Fock exchange
 1287 energy is

$$1288 (9.19) \quad E_{\text{f}}(X) := \frac{1}{4} \langle \mathcal{V}(X X^*) X, X \rangle = \frac{1}{4} \langle \mathcal{V}(X X^*), X X^* \rangle.$$

1289 Therefore, the total energy minimization problem can be formulated as

$$1290 (9.20) \quad \min_{X \in \mathbb{C}^{n \times p}} E(X), \quad \text{s. t.} \quad X^* X = I_p,$$

1291 where $E(X)$ is $E_{\text{ks}}(X)$ in KSDFT and

$$1292 E_{\text{hf}}(X) := E_{\text{ks}}(X) + E_{\text{f}}(X)$$

1293 in HF. Computing $E_{\text{f}}(X)$ is very expensive since a multiplication between an $n \times n \times n \times n$
 1294 fourth-order tensor and an n -by- n matrix is needed in $\mathcal{V}(\cdot)$.

1295 Denote the KS Hamiltonian $H_{\text{ks}}(X)$ as

$$1296 (9.21) \quad H_{\text{ks}}(X) := \frac{1}{2} L + V_{\text{ion}} + \sum_l \zeta_l w_l w_l^* + \text{Diag}((\Re L^\dagger) \rho) + \text{Diag}(\mu_{\text{xc}}(\rho)^* e_n),$$

1297 where $\mu_{\text{xc}}(\rho) = \frac{\partial \epsilon_{\text{xc}}(\rho)}{\partial \rho}$. Since $H_{\text{ks}}(X)$ is essentially determined by the charge density $\rho(X)$,
 1298 it is often written as $H_{\text{ks}}(\rho)$. The HF Hamiltonian is

$$1299 (9.22) \quad H_{\text{hf}}(X) := H_{\text{ks}}(X) + \mathcal{V}(X X^*).$$

1300 A detailed calculation shows that the Euclidean gradient of $E_{\text{ks}}(X)$ is

$$1301 (9.23) \quad \nabla E_{\text{ks}}(X) = H_{\text{ks}}(X) X.$$

1302 The gradient of $E_{\text{f}}(X)$ is $\nabla E_{\text{f}}(X) = \mathcal{V}(X X^*) X$. Assume that $\epsilon_{\text{xc}}(\rho(X))$ is twice differen-
 1303 tiable with respect to $\rho(X)$, the Hessian of $E_{\text{ks}}(X)$ is

$$1304 (9.24) \quad \nabla^2 E_{\text{ks}}(X)[U] = H_{\text{ks}}(X) U + \mathcal{R}(X)[U],$$

1305 where $U \in \mathbb{C}^{n \times p}$ and $\mathcal{R}(X)[U] := \text{Diag} \left((\Re L^\dagger + \frac{\partial^2 \epsilon_{\text{xc}}}{\partial \rho^2} e_n) (\bar{X} \odot U + X \odot \bar{U}) e_n \right) X$. The
 1306 Hessian of $E_{\text{f}}(X)$ is

$$1307 (9.25) \quad \nabla^2 E_{\text{f}}(X)[U] = \mathcal{V}(X X^*) U + \mathcal{V}(X U^* + U X^*) X.$$

1308 **9.3.2. The Self-Consistent Field (SCF) Iteration.** The first-order optimality con-
 1309 ditions for the total energy minimization problem are

$$1310 \quad (9.26) \quad H(X)X = X\Lambda, \quad X^*X = I_p,$$

1311 where $X \in \mathbb{C}^{n \times p}$, Λ is a diagonal matrix and H represents H_{ks} in (9.21) or H_{hf} in (9.22). For
 1312 KSDFT, one of the most popular methods is the SCF iteration. At the k -th iteration, we first
 1313 fix the Hamiltonian to be $H_{\text{ks}}(\tilde{\rho}_k)$ for a given $\tilde{\rho}_k$ and solve the following linear eigenvalue
 1314 problem

$$1315 \quad (9.27) \quad H_{\text{ks}}(\tilde{\rho}_k)X = X\Lambda, \quad X^*X = I_p.$$

1316 The eigenvectors corresponding to the p smallest eigenvalues of $H_{\text{ks}}(\rho_k)$ is denoted as X_{k+1} ,
 1317 which leads to a new charge density $\rho_{k+1} = \rho(X_{k+1})$. It is then mixed with charge densities
 1318 from previous steps to produce the new charge density $\tilde{\rho}_{k+1}$ in order to accelerate the con-
 1319 vergence instead of using ρ_{k+1} directly. This procedure is repeated until self-consistency is
 1320 reached.

1321 A particular charge mixing scheme is the direct inversion of iterative subspace (DIIS) or
 1322 the Pulay mixing [92, 93, 115]. Choose an integer m with $m \leq k$. Let

$$1323 \quad W = (\Delta\rho_k, \Delta\rho_{k-1}, \dots, \Delta\rho_{k-m+1}), \quad \Delta\rho_j = \rho_j - \rho_{j-1}.$$

1324 The Pulay mixing generates the charge density $\tilde{\rho}_k$ by a linear combination of the previously
 1325 charge densities

$$1326 \quad \tilde{\rho}_k = \sum_{j=0}^{m-1} c_j \rho_{k-j},$$

1327 where $c = (c_0, c_1, \dots, c_{m-1})$ is the solution to the minimization problem:

$$1328 \quad \min_c \quad \|Wc\|^2, \quad \text{s. t.} \quad c^\top e_m = 1.$$

1329 Other types of mixing includes Broyden mixing, Kerker mixing and Anderson mixing, etc.
 1330 Charge mixing is widely used for improving the convergence of SCF even though its conver-
 1331 gence property is still not clear in few cases.

1332 In HF, the SCF method at the k -th iteration solves:

$$1333 \quad \tilde{H}_k X = X\Lambda, \quad X^*X = I_p,$$

1334 where \tilde{H}_k is formed from certain mixing schemes. Note that the Hamiltonian (9.22) can be
 1335 written as $H_{\text{hf}}(D)$ with respect to the density matrix $D = XX^*$. In the commutator DIIS
 1336 (C-DIIS) method [92, 93], the residual W_j is defined as the commutator between $H_{\text{hf}}(D_j)$
 1337 and D_j , i.e.,

$$1338 \quad (9.28) \quad W_j = H_{\text{hf}}(D_j)D_j - D_j H_{\text{hf}}(D_j).$$

1339 We next solve the following minimization to obtain a coefficient c :

$$1340 \quad \min_c \quad \left\| \sum_{j=0}^{m-1} c_j W_j \right\|_F^2, \quad \text{s. t.} \quad c^\top e_m = 1.$$

1341 Then, a new Hamiltonian matrix is obtained $\tilde{H}_k = \sum_{j=0}^{m-1} c_j H_{k-j}$. Since an explicit storage
 1342 of the density matrix can be prohibitive, the projected C-DIIS in [60] uses projections of the
 1343 density and commutator matrices so that the sizes are much smaller.

1344 **9.3.3. Subspace Methods For HF using Nyström Approximation.** Note that
 1345 the most expensive part in HF is the evaluation of $E_f(X)$ and the related derivatives. We
 1346 apply the limited-memory Nyström technique to approximate $\mathcal{V}(X_k X_k^*)$ by $\hat{\mathcal{V}}(X_k X_k^*)$. Let
 1347 $Z = \mathcal{V}(X_k X_k^*) \Omega$ where Ω is an orthogonal basis of the subspace such as

$$1348 \quad \text{span}\{X_k\}, \text{span}\{X_{k-1}, X_k\} \text{ or } \text{span}\{X_{k-1}, X_k, \mathcal{V}(X_k X_k^*)X_k\}.$$

1349 Then the low rank approximation

$$1350 \quad (9.29) \quad \hat{\mathcal{V}}(X_k X_k^*) := Z(Z^* \Omega)^\dagger Z^*$$

1351 is able to reduce the computational cost significantly. Note that the adaptive compression
 1352 method in [73] compresses the operator $\mathcal{V}(X_k X_k^*)$ on the subspace $\text{span}\{X_k\}$. Conse-
 1353 quently, we can keep the easier parts E_{ks} but approximate $E_f(X)$ by using (9.29). Hence, a
 1354 new subproblem is formulated as

$$1355 \quad (9.30) \quad \min_{X \in \mathbb{C}^{n \times p}} E_{\text{ks}}(X) + \frac{1}{4} \left\langle \hat{\mathcal{V}}(X_k X_k^*) X, X \right\rangle \quad \text{s. t.} \quad X^* X = I_p.$$

1356 The subproblem (9.30) can be solved by the SCF iteration, the Riemannian gradient method
 1357 or the modified CG method based on the following linear equation

$$1358 \quad \text{Proj}_{X_k} \left(\nabla^2 E_{\text{ks}}(X_k)[\xi] + \frac{1}{2} \hat{\mathcal{V}}(X_k X_k^*) \xi - \xi \text{sym}(X_k^* \nabla f(X_k)) \right) = -\text{grad} E_{\text{hf}}(X_k).$$

1359 **9.3.4. A Regularized Newton Type Method.** Computing the p -smallest eigenpairs
 1360 of $H_{\text{ks}}(\tilde{\rho})$ is equivalent to a trace minimization problem

$$1361 \quad (9.31) \quad \min_{X \in \mathbb{C}^{n \times p}} q(X) := \frac{1}{2} \text{tr}(X^* H_{\text{ks}}(\tilde{\rho}) X) \quad \text{s. t.} \quad X^* X = I_p.$$

1362 Note that $q(X)$ is a second-order approximation to the total energy $E_{\text{ks}}(X)$ without consid-
 1363 ering the second term in the Hessian (9.24). Hence, the SCF method may not converge if this
 1364 second term dominates. The regularized Newton in (9.1) can be applied to solve both KSDFT
 1365 and HF with convergence guarantees. We next explain a particular version in [138] whose
 1366 subproblem is

$$1367 \quad (9.32) \quad \min_{X \in \mathbb{C}^{n \times p}} \frac{1}{2} \text{tr}(X^* H_{\text{ks}}(\tilde{\rho}) X) + \frac{\tau_k}{4} \|X X^\top - X_k X_k^\top\|_F^2 \quad \text{s. t.} \quad X^* X = I_p.$$

1368 Since X_k and X are orthonormal matrices, we have

$$1369 \quad \begin{aligned} \|X X^\top - X_k X_k^\top\|_F^2 &= \text{tr}((X X^\top - X_k X_k^\top)(X X^\top - X_k X_k^\top)) \\ &= 2p - 2\text{tr}(X^\top X_k X_k^\top X). \end{aligned}$$

1370 Therefore, (9.32) is a linear eigenvalue problem:

$$1371 \quad \begin{aligned} (H_{\text{ks}}(\tilde{\rho}) - \tau_k X_k X_k^\top) X &= X \Lambda, \\ X^\top X &= I_p. \end{aligned}$$

1372 **9.3.5. Subspace Refinement for KSDFT.** The direct minimization method in [138]
 1373 is a kind of subspace refinement procedure using

$$1374 \quad Y = [X_k, P_k, R_k],$$

1375 where $P_k = X_k - X_{k-1}$ and $R_k = H_{\text{ks}}(X_k)X_k - X_k\Lambda_k$. Then the variable X can be
 1376 expressed as $X = YG$ where $G \in \mathbb{C}^{3p \times p}$. The total energy minimization problem (9.20)
 1377 becomes:

$$1378 \quad \min_G E_{\text{ks}}(YG), \text{ s. t. } G^*Y^*YG = I_p,$$

1379 whose first-order optimality condition is a generalized linear eigenvalue problem:

$$1380 \quad (Y^*H_{\text{ks}}(YG)Y)G = Y^*YG\Omega, \quad G^*Y^*YG = I_p.$$

1381 The subspace refinement method may help when the regularized Newton method does
 1382 not perform well. Note that the total energy minimization problem (9.20) is not necessary
 1383 equivalent to a nonlinear eigenvalue problem (9.26) for finding the p smallest eigenvalues of
 1384 $H(X)$. Although an intermediate iterate X is orthogonal and contains eigenvectors of $H(X)$,
 1385 these eigenvectors are not necessary the eigenvectors corresponding to the p smallest eigen-
 1386 values. Hence, we can form a subspace which contains these possible target eigenvectors. In
 1387 particular, we first compute the first γp smallest eigenvalues for some small integer γ . Their
 1388 corresponding eigenvectors of $H(X_k)$, denoted by Γ_k , are put in a subspace as

$$1389 \quad (9.33) \quad \text{span}\{X_{k-1}, X_k, \text{grad} E(X_k), \Gamma_k\}.$$

1390 Numerical experience shows that the refinement scheme in subsection 9.1 with this subspace
 1391 is likely escape a stagnated point.

1392 **10. Semidefinite Programming (SDP).** In this section, we present two specialized
 1393 subspace methods for solving the maxcut SDP and the maxcut SDP with nonnegative con-
 1394 straints from community detection.

1395 **10.1. The Maxcut SDP.** The maxcut problem partition the vertices of a graph into
 1396 two sets so that the sum of the weights of the edges connecting vertices in one set with these
 1397 in the other set is maximized. The corresponding SDP relaxation [46, 16, 56, 8] is

$$1398 \quad (10.1) \quad \begin{aligned} & \min_{X \in \mathcal{S}^n} \langle C, X \rangle \\ & \text{s. t. } X^{ii} = 1, \quad i = 1, \dots, n, \\ & \quad X \succeq 0. \end{aligned}$$

1399 We first describe a second-order cone program (SOCP) restriction for the SDP pro-
 1400 blem (10.1) by fixing all except one row and column of the matrix X . For any integer
 1401 $i \in \{1, \dots, n\}$, the complement of the set $\{i\}$ is $i^c = \{1, \dots, n\} \setminus \{i\}$. Let $B = X^{i^c, i^c}$
 1402 be the submatrix of X after deleting its i -th row and column, and $y = X^{i^c, i}$ be the i th col-
 1403 umn of the matrix X without the element $X^{i,i}$. Since $X_{ii} = 1$, the variable X of (10.1) can
 1404 be written as

$$1405 \quad X := \begin{pmatrix} 1 & y^\top \\ y & B \end{pmatrix} := \begin{pmatrix} 1 & y^\top \\ y & X^{i^c, i^c} \end{pmatrix}$$

1406 without loss of generality. Suppose that the submatrix B is fixed. It then follows from the
 1407 Schur complement theorem that $X \succeq 0$ is equivalent to

$$1408 \quad \xi - y^\top B^{-1}y \geq 0.$$

1409 In order to maintain the strict positive definiteness of X , we require $1 - y^\top B^{-1}y \geq \nu$ for a
 1410 small constant $\nu > 0$. Therefore, the SDP problem (10.1) is reduced to a SOCP:

$$1411 \quad (10.2) \quad \begin{aligned} & \min_{y \in \mathbb{R}^{n-1}} \hat{c}^\top y \\ & \text{s. t.} \quad 1 - y^\top B^\dagger y \geq \nu, \quad y \in \text{Range}(B), \end{aligned}$$

1412 where $\hat{c} := 2C^{i^c, i}$. If $\gamma := \hat{c}^\top B \hat{c} > 0$, an explicit solution of (10.2) is given by

$$1413 \quad (10.3) \quad y = -\sqrt{\frac{1-\nu}{\gamma}} B \hat{c}.$$

1414 Otherwise, the solution is $y = 0$.

1415 We next describe the RBR method [130]. Starting from a positive definite feasible so-
 1416 lution X_1 , it updates one row/column of X at each of the inner steps. The operations from
 1417 the first row to the last row is called a cycle. At the first step of the k -th cycle, the matrix B
 1418 is set to $X_k^{1^c, 1^c}$ and y is computed by (10.3). Then the first row/column of X_k is substituted
 1419 by $X_k^{1^c, 1} := y$. Other rows/columns are updated in a similar fashion until all row/column
 1420 are updated. Then we set $X_{k+1} := X_k$ and this procedure is repeated until certain stopping
 1421 criteria are met.

1422 The RBR method can also be derived from the logarithmic barrier problem

$$1423 \quad (10.4) \quad \begin{aligned} & \min_{X \in \mathcal{S}^n} \phi_\sigma(X) := \langle C, X \rangle - \sigma \log \det X \\ & \text{s. t.} \quad X^{ii} = 1, \forall i = 1, \dots, n, \quad X \succ 0. \end{aligned}$$

1424 Fixing the block $B = X^{i^c, i^c}$ gives

$$1425 \quad \det(X) = \det(B)(1 - (X^{i^c, i})^\top B^{-1} X^{i^c, i}).$$

1426 Therefore, the RBR subproblem for (10.4) is

$$1427 \quad (10.5) \quad \min_{y \in \mathbb{R}^{n-1}} \hat{c}^\top y - \sigma \log(1 - y^\top B^{-1}y).$$

1428 If $\gamma := \hat{c}^\top B \hat{c} > 0$, the solution of problem (10.5) is

$$1429 \quad (10.6) \quad y = -\frac{\sqrt{\sigma^2 + \gamma} - \sigma}{\gamma} B \hat{c}.$$

1430 Consequently, the subproblem (10.2) has the same solution as (10.5) if $\nu = 2\sigma \frac{\sqrt{\sigma^2 + \gamma} - \sigma}{\gamma}$.

1431 **10.1.1. Examples: Phase Retrieval.** Given a matrix $A \in \mathbb{C}^{m \times n}$ and a vector $b \in$
 1432 \mathbb{R}^m , the phase retrieval problem can be formulated as a feasibility problem:

$$1433 \quad \text{find } x, \text{ s. t. } |Ax| = b.$$

1434 An equivalent model in [122] is

$$1435 \quad \begin{aligned} & \min_{x \in \mathbb{C}^n, y \in \mathbb{R}^m} \frac{1}{2} \|Ax - y\|_2^2 \\ & \text{s. t.} \quad |y| = b, \end{aligned}$$

1436 which can be further reformulated as

$$1437 \quad (10.7) \quad \min_{x \in \mathbb{C}^n, u \in \mathbb{C}^m} \frac{1}{2} \|Ax - \text{diag}(b)u\|_2^2$$

$$\text{s. t.} \quad |u^i| = 1, i = 1, \dots, m.$$

1438 By fixing the variable u , it becomes a least squares problem with respect to x , whose explicit
1439 solution is $x = A^\dagger \text{diag}(b)u$. Substituting x back to (10.7) yields a general maxcut problem:

$$1440 \quad \min_{u \in \mathbb{C}^m} u^* M u$$

$$\text{s. t.} \quad |u^i| = 1, i = 1, \dots, m,$$

1441 where $M = \text{diag}(b)(I - AA^\dagger)\text{diag}(b)$ is positive semidefinite. Hence, the corresponding
1442 SDP relaxation is

$$1443 \quad \min_{U \in \mathcal{S}^m} \text{tr}(UM)$$

$$\text{s. t.} \quad U^{ii} = 1, i = 1, \dots, m, U \succeq 0.$$

1444 The above problem can be further solved by the RBR method.

1445 **10.2. Community Detection.** Suppose that the nodes $[n] = \{1, \dots, n\}$ of a network
1446 can be partitioned into $r \geq 2$ disjoint sets $\{K_1, \dots, K_r\}$. A binary matrix X is called a
1447 partition matrix if $X^{ij} = 1$ for $i, j \in K_t, t \in \{1, \dots, r\}$ and otherwise $X^{ij} = 0$. Let A be the
1448 adjacency matrix and d be the degree vector, where $d_i = \sum_j A^{ij}, i \in [n]$. Define the matrix

$$1449 \quad (10.8) \quad C = -(A - \lambda dd^\top),$$

1450 where $\lambda = 1/\|d\|_1$. A popular method for the community detection problem is to maximize
1451 the modularity [86] as:

$$1452 \quad (10.9) \quad \min_X \langle C, X \rangle \text{ s. t. } X \in \mathcal{P}_n^r,$$

1453 where \mathcal{P}_n^r is the set of all partition matrices of n nodes with no more than r subsets. Since
1454 the modularity optimization (10.9) is NP-hard, a SDP relaxation proposed in [25] is:

$$1455 \quad (10.10) \quad \min_{X \in \mathbb{R}^{n \times n}} \langle C, X \rangle$$

$$\text{s. t.} \quad X^{ii} = 1, i = 1, \dots, n,$$

$$0 \leq X^{ij} \leq 1, \forall i, j,$$

$$X \succeq 0.$$

1456 The RBR method in subsection 10.1 can not be applied to (10.10) directly due to the compo-
1457 nentwise constraints $0 \leq X^{ij} \leq 1$.

1458 Note that the true partition matrix X^* can be decomposed as $X^* = \Phi^*(\Phi^*)^\top$, where
1459 $\Phi^* \in \{0, 1\}^{n \times r}$ is the true assignment matrix. This decomposition is unique up to a permu-
1460 tation of the columns of Φ^* . The structures of Φ^* leads to a new relaxation of the original
1461 partition matrix X [146]. Define a matrix

$$1462 \quad U = [u^1, \dots, u^n]^\top \in \mathbb{R}^{n \times r}.$$

1463 We can consider a decomposition $X = UU^\top$. The constraints $X^{ii} = 1$ and $\Phi^* \in \{0, 1\}^{n \times r}$
1464 imply that

$$1465 \quad \|u^i\|_2 = 1, \quad U \geq 0, \quad \|u^i\|_0 \leq p,$$

1466 where the cardinality constraints are added to impose sparsity of the solution. Therefore, an
 1467 alternative relaxation to (10.9) is

$$\begin{aligned}
 & \min_{U \in \mathbb{R}^{n \times r}} \langle C, UU^\top \rangle \\
 & \text{s. t.} \quad \|u^i\|_2 = 1, i = 1, \dots, n, \\
 & \quad \quad \|u^i\|_0 \leq p, i = 1, \dots, n, \\
 & \quad \quad U \geq 0.
 \end{aligned}
 \tag{10.11}$$

1469 Although (10.11) is still NP-hard, it enables us to develop a computationally efficient
 1470 RBR method. The feasible set for each block u^i is

$$\mathcal{U} = \{u \in \mathbb{R}^r \mid \|u\|_2 = 1, \quad u \geq 0, \quad \|u\|_0 \leq p\}.$$

1472 Then, problem (10.11) can be rewritten as

$$\min_{U \in \mathbb{R}^{n \times r}} f(U) \equiv \langle C, UU^\top \rangle, \quad \text{s. t.} \quad u^i \in \mathcal{U}.
 \tag{10.12}$$

For the i -th subproblem, we fix all except the i -th row of U and formulate the subproblem as

$$u^i = \arg \min_{x \in \mathcal{U}} f(u^1, \dots, u^{i-1}, x, u^{i+1}, \dots, u^n) + \frac{\sigma}{2} \|x - \bar{u}^i\|^2,$$

1474 where the last part in the objective function is the proximal term and $\sigma > 0$ is a parameter.
 1475 Note that the quadratic term $\|x\|^2$ is eliminated due to the constraint $\|u\|_2 = 1$. Therefore,
 1476 the subproblem becomes

$$u^i = \arg \min_{x \in \mathcal{U}} b^\top x,
 \tag{10.13}$$

1478 where $b = 2C^{i,i^c}U^{-i} - \sigma\bar{u}^i$, and C^{i,i^c} is the i -th row of C without the i -th component, U^{-i}
 1479 is the matrix U without the i -th row. Define $b_+ = \max\{b, 0\}$, $b_- = \max\{-b, 0\}$, where the
 1480 max is taken component-wisely. Then the closed-form solution of (10.13) is given by

$$u = \begin{cases} \frac{b_-^p}{\|b_-^p\|}, & \text{if } b_- \neq 0, \\ e^{j_0}, \text{ with } j_0 = \arg \min_j b^j, & \text{otherwise,} \end{cases}
 \tag{10.14}$$

1482 where b_-^p is obtained by keeping the largest p components in b_- and letting the others be
 1483 zero, and when $\|b_-\|_0 \leq p$, $b_-^p = b_-$. Then the RBR method goes over all rows of U by
 1484 using (10.14).

1485 We next briefly describe the parallelization of the RBR method on a shared memory
 1486 computer with many threads. The variable U is stored in the shared memory so that it can be
 1487 accessed by all threads. Even when some row u^i is updating in a thread, the other threads can
 1488 still access U whenever necessary. In the sequential RBR method, the main cost of updating
 1489 one row u^i is the computation of $b = 2C^{i,i^c}U^{-i} - \sigma\bar{u}^i$, where \bar{u}^i and U are the current
 1490 iterates. The definition of C in (10.8) gives

$$b^\top = -2A^{i,i^c}U^{-i} + 2\lambda d^i (d^{i^c})^\top U^{-i} - \sigma\bar{u}^i,
 \tag{10.15}$$

1492 where A^{i,i^c} is the i -th row of A without the i -th component. The parallel RBR method is
 1493 outlined in Figure 10.1 where many threads are working simultaneously. The vector $d^\top U$
 1494 and matrix U are stored in the shared memory and all threads can access and update them.

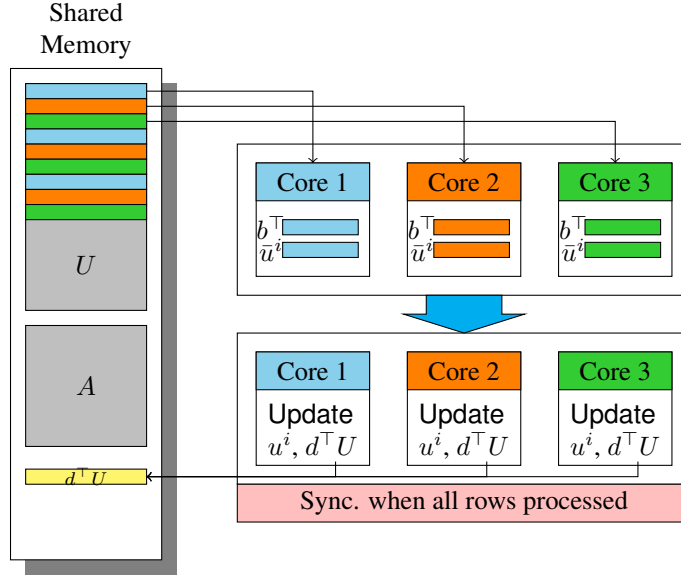


Fig. 10.1 An illustration of the asynchronous parallel proximal RBR method

1495 Every thread picks up their own row u^i at a time and then reads U and the vector $d^\top U$. Then,
 1496 a private copy of b^\top is computed. Thereafter, the variable u^i is updated and $d^\top U$ is set to
 1497 $d^\top U \leftarrow d^\top U + d^i(u^i - \bar{u}^i)$ in the shared memory. It immediately proceeds to another row
 1498 without waiting for other threads to finish their tasks. Therefore, when a thread is updating
 1499 its variables, other blocks of variables $u^j, j \neq i$ are not necessarily the most new version.
 1500 Moreover, if this thread is reading some row u^j or the vector $d^\top U$ from memory and another
 1501 thread is just modifying them, the data of u^i will be partially updated. Since the memory
 1502 locking is removed, the parallel RBR method may be able to provide near-linear speedups.
 1503 See also the HOGWILD! [94] and CYCLADES [89] for the asynchronous methods.

1504 **11. Low Rank Matrix Optimization.** Optimization problems whose variable is re-
 1505 lated to low-rank matrices arise in many applications, for example, semidefinite programming
 1506 (SDP), matrix completion, robust principle component analysis, control and systems theory,
 1507 model reduction [76], phase retrieval, blind deconvolution, data mining, pattern recognitions
 1508 [33], latent semantic indexing, collaborative prediction and low-dimensional embedding.

1509 **11.1. Low Rank Structure of First-order Methods.** A common feature of many
 1510 first-order methods for the low rank matrix optimization problems is that the next iterate x_{k+1}
 1511 is defined by the current iterate x_k and a partial eigenvalue decomposition of certain matrix.
 1512 They can be unified as the following fixed-point iteration scheme [71]:

1513 (11.1)
$$x_{k+1} = \mathcal{T}(x_k, \Psi(\mathcal{B}(x_k))), \quad x_k \in \mathcal{D},$$

1514 where $\mathcal{B} : \mathcal{D} \rightarrow \mathcal{S}^n$ is a bounded mapping from a given Euclidean space \mathcal{D} to the n -
 1515 dimensional symmetric matrix space \mathcal{S}^n , and \mathcal{T} is a general mapping from $\mathcal{D} \times \mathcal{S}^n$ to \mathcal{D} .
 1516 The spectral operator $\Psi : \mathcal{S}^n \rightarrow \mathcal{S}^n$ is given by

1517 (11.2)
$$\Psi(X) = V \text{Diag}(\psi(\lambda(X))) V^\top,$$

1518 where $X = V\text{Diag}(\lambda_1, \dots, \lambda_n)V^\top$ is the eigenvalue decomposition of X with eigenvalues
 1519 in descending order $\lambda_1 \geq \lambda_2 \cdots \geq \lambda_n$, $\lambda(X) = (\lambda_1, \dots, \lambda_n)^T$, the operator $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$
 1520 is a vector-valued *symmetric mapping*, i.e., $\psi(P\lambda) = P\psi(\lambda)$ for any permutation matrix P .
 1521 The orthogonal projection of a symmetric matrix X on to a given $\text{Range}(Q)$ with $Q^\top Q = I$
 1522 is defined as:

$$1523 \quad (11.3) \quad \mathcal{P}_Q(X) := \arg \min_{Y \in \mathcal{S}^n, \text{Range}(Y) = \text{Range}(Q)} \|Y - X\|_F^2 = QQ^\top XQQ^\top.$$

1524 The operator Ψ has the low-rank property at X if there exists an orthogonal matrix $V_{\mathcal{I}} \in$
 1525 $\mathbb{R}^{n \times p}$ ($p \ll n$) that span a p -dimensional eigen-space corresponding to $\lambda_i(X)$, $i \in \mathcal{I}$, such
 1526 that $\Psi(X) = \Phi(\mathcal{P}_{V_{\mathcal{I}}}(X))$, where Φ is either the same as Ψ or a different spectral operator
 1527 induced by ϕ , and \mathcal{I} is an index set depending on X . The low-rank property ensures that the
 1528 full eigenvalue decomposition is not needed.

1529 The scheme (11.1) is time-consuming for large scale problems since first-order methods
 1530 often take thousands of iterations to converge and each iteration requires at least one full
 1531 or partial eigenvalue decomposition for evaluating Ψ . However, $\Psi(\mathcal{B}(x_k))$ often lives in a
 1532 low-dimensional eigenspace in practice. A common practice is to use inexact method such as
 1533 the Lanczos method, LOBPCG, and randomized methods with early stopping rules [149, 6,
 1534 106]. The so-called subspace method performs refinement on a low-dimensional subspace for
 1535 univariate maximal eigenvalue optimization problem [66, 102, 63] and in the SCF iteration
 1536 for KSDFT [151]. In the rest of this section, we present approaches [71] which integrate
 1537 eigenvalue computation coherently with the underlying optimization methods.

1538 **11.2. A Polynomial-filtered Subspace Method.** We now describe a general sub-
 1539 space framework for the scheme (11.1) using Chebyshev polynomials $\rho_k(\cdot)$ defined in (8.5).
 1540 Assume that x^* is a limit point of the fixed-point iteration (11.1) and the low-rank property
 1541 holds for every $\mathcal{B}(x_k)$ in (11.1). Consequently, the scheme (11.1) is equivalent to

$$1542 \quad (11.4) \quad x_{k+1} = \mathcal{T}(x_k, \Phi(\mathcal{P}_{V_{\mathcal{I}_k}}(\mathcal{B}(x_k))))),$$

1543 where $V_{\mathcal{I}_k}$ is determined by $\mathcal{B}(x_k)$. Although the exact subspace $V_{\mathcal{I}_k}$ usually is unknown,
 1544 it can be approximated by an estimated subspace U_k so that the computational cost of Ψ
 1545 is significantly reduced. After the next point x_{k+1} is formed, a polynomial filter step is
 1546 performed in order to extract a new subspace U_{k+1} based on U_k . Therefore, combining the
 1547 two steps (8.6) and (11.4) together gives

$$1548 \quad (11.5) \quad x_{k+1} = \mathcal{T}(x_k, \Phi(\mathcal{P}_{U_k}(\mathcal{B}(x_k))))),$$

$$1549 \quad (11.6) \quad U_{k+1} = \text{orth}(\rho_{k+1}^{q_{k+1}}(\mathcal{B}(x_{k+1}))U_k),$$

1551 where q_k is a small number (e.g. 1 to 3) of the polynomial filter $\rho_k(\cdot)$ applied to U_k . The
 1552 Chebyshev polynomials are suitable when the targeted eigenvalues are located within an inter-
 1553 val, for example, finding a few largest/smallest eigenvalues in magnitude or all positive/neg-
 1554 ative eigenvalues.

1555 The main feature is that the exact subspace $V_{\mathcal{I}_k}$ is substituted by its approximation U_k
 1556 in (11.5). The principle angle between the true and extracted subspace is controlled by the
 1557 polynomial degree. Then the error between one exact and inexact iteration is bounded. When
 1558 the initial space is not orthogonal to the target space, the convergence of (11.5)-(11.6) is
 1559 established under mild assumptions. In fact, the subspace often becomes more and more
 1560 accurate so that the warm start property is helpful, i.e., the subspace of the current iteration
 1561 can be refined from the previous one.

1562 **11.3. The Polynomial-filtered Proximal Gradient Method.** We next show how
 1563 to apply the subspace update (11.5) and (11.6) to the proximal gradient method on a set of
 1564 composite optimization problems

$$1565 \quad (11.7) \quad \min h(x) := F(x) + R(x),$$

1566 where $F(x) = f \circ \lambda(\mathcal{B}(x))$ with $\mathcal{B}(x) = G + \mathcal{A}^*(x)$ and $R(x)$ is a regularization term with
 1567 simple structures but need not be smooth. Here G is a known matrix in \mathcal{S}^n , the linear operator
 1568 \mathcal{A} and its adjoint operator \mathcal{A}^* are defined as

$$1569 \quad (11.8) \quad \mathcal{A}(X) = [\langle A_1, X \rangle, \dots, \langle A_m, X \rangle]^T, \quad \mathcal{A}^*(x) = \sum_{i=1}^m x_i A_i,$$

1570 for given symmetric matrices $A_i \in \mathcal{S}^n$. The function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is smooth and *absolutely*
 1571 *symmetric*, i.e., $f(x) = f(Px)$ for all $x \in \mathbb{R}^n$ and any permutation matrix $P \in \mathbb{R}^{n \times n}$.

1572 Let Ψ be a spectral operator induced by $\psi = \nabla f$. It can be verified that the gradient of
 1573 F in (11.7) is

$$1574 \quad (11.9) \quad \nabla F(x) = \mathcal{A}(\Psi(\mathcal{B}(x))).$$

1575 The proximal operator is defined by

$$1576 \quad (11.10) \quad \text{prox}_{tR}(x) = \arg \min_u R(u) + \frac{1}{2t} \|u - x\|_2^2.$$

1577 Consequently, the proximal gradient method is

$$1578 \quad (11.11) \quad x_{k+1} = \text{prox}_{\tau_k R}(x_k - \tau_k \mathcal{A}(\Psi(\mathcal{B}(x_k))))),$$

1579 where τ_k is the step size. Therefore, the iteration (11.11) is a special case of (11.1) with

$$1580 \quad \begin{aligned} \mathcal{T}(x, X) &= \text{prox}_{\tau_k R}(x - \mathcal{A}(X)), \\ \Psi(X) &= V \text{Diag}(\nabla f(\lambda(X))) V^\top. \end{aligned}$$

1581 Assume that the low-rank property holds at every iteration. The corresponding polynomial-
 1582 filtered method can be written as

$$1583 \quad (11.12) \quad x_{k+1} = \text{prox}_{\tau_k R}(x_k - \tau_k \mathcal{A}(\Phi(\mathcal{P}_{U_k}(\mathcal{B}(x_k))))),$$

$$1584 \quad (11.13) \quad U_{k+1} = \text{orth}(\rho_{k+1}^{q_{k+1}}(\mathcal{B}(x_{k+1}))U_k).$$

1586 **11.3.1. Examples: Maximal Eigenvalue and Matrix Completion.** Consider the
 1587 maximal eigenvalue optimization problem:

$$1588 \quad (11.14) \quad \min_x F(x) + R(x) := \lambda_1(\mathcal{B}(x)) + R(x),$$

1589 where $\mathcal{B}(x) = G + \mathcal{A}^*(x)$. Certain specific formulations of phase recovery and blind decon-
 1590 volution are special case of (11.14). The subgradient of $F(x)$ is

$$1591 \quad \partial F(x) = \{\mathcal{A}(U_1 S U_1^T) \mid S \succeq 0, \text{tr}(S) = 1\},$$

1592 where $U_1 \in \mathbb{R}^{n \times r_1}$ is the subspace spanned by eigenvectors of $\lambda_1(\mathcal{B}(x))$ with multiplicity
 1593 r_1 . For simplicity, we assume $r_1 = 1$ and $\lambda_1(\mathcal{B}(x)) > 0$, which means that $\partial F(x)$ has only
 1594 one element and the function $F(x)$ is differentiable. Then the polynomial-filtered method is

$$1595 \quad x_{k+1} = \text{prox}_{\tau R}(x_k - \tau \mathcal{A}(u_1 u_1^T)),$$

1596 where u_1 is the eigenvector of $\lambda_1(\mathcal{B}(x_k))$. Hence, we have

$$1597 \quad \mathcal{T}(x, W) = \text{prox}_{\tau R}(x - \tau \mathcal{A}(W)), \quad \Psi(X) = u_1 u_1^T.$$

1598 In addition, $\Psi(\cdot)$ satisfies the low-rank property around x^* with $\mathcal{I} = \{1\}$ and

$$1599 \quad (\psi(\lambda))_i = (\phi(\lambda))_i = \begin{cases} 1, & i = 1, \\ 0, & \text{otherwise.} \end{cases}$$

1600 Another example is the penalized formulation of the matrix completion problem:

$$1601 \quad (11.15) \quad \min \|X\|_* + \frac{1}{2\mu} \|\mathcal{P}_\Omega(X - M)\|_F^2,$$

1602 where Ω is a given index set of the true matrix M , and $\mathcal{P}_\Omega : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ denotes
 1603 the projection operator onto the sparse matrix space with non-zero entries on Ω . Problem
 1604 (11.15) can be solved by the proximal gradient method. At the k -th iteration, the main cost
 1605 is to compute the truncated SVD of a matrix. Although (11.15) is not a direct special case of
 1606 (11.7), we can still insert the polynomial filter into the proximal gradient method to reduce
 1607 the cost of SVD.

1608 **11.4. The Polynomial-filtered ADMM Method.** Consider the standard SDP:

$$1609 \quad (11.16) \quad \begin{aligned} & \min \quad \langle C, X \rangle, \\ & \text{s. t.} \quad \mathcal{A}X = b, \\ & \quad \quad X \succeq 0, \end{aligned}$$

1610 where C , \mathcal{A} and b are given, the linear operator \mathcal{A} and its adjoint are defined in (11.8).
 1611 Note that the ADMM on the dual problem of (11.16) is equivalent to the Douglas-Rachford
 1612 Splitting (DRS) method [30] on the primal SDP (11.16). Define $F(X) = 1_{\{X \succeq 0\}}(X)$ and
 1613 $G(X) = 1_{\{\mathcal{A}X = b\}}(X) + \langle C, X \rangle$, where $1_\Omega(X)$ is the indicator function on a set Ω . The
 1614 proximal operators $\text{prox}_{tF}(Z)$ and $\text{prox}_{tG}(Y)$ can be computed explicitly as

$$1615 \quad \begin{aligned} \text{prox}_{tF}(Z) &= \mathcal{P}_+(Z), \\ \text{prox}_{tG}(Y) &= (Y + tC) - \mathcal{A}^*(\mathcal{A}\mathcal{A}^*)^{-1}(\mathcal{A}Y + t\mathcal{A}C - b), \end{aligned}$$

1617 where $\mathcal{P}_+(Z)$ is the projection operator onto the positive semi-definite cone. Hence, DRS
 1618 can be formulated as

$$1619 \quad Z_{k+1} = T_{\text{DRS}}(Z_k) \triangleq \text{prox}_{tG}(2\text{prox}_{tF}(Z_k) - Z_k) - \text{prox}_{tF}(Z_k) + Z_k,$$

1620 which is also a special case of (11.1) with

$$1621 \quad \begin{aligned} \mathcal{T}(x, X) &= \text{prox}_{tG}(2X - x) - X + x, \\ \Psi(X) &= \mathcal{P}_+(X). \end{aligned}$$

1622 Note that $\mathcal{P}_+(X)$ is a spectral operator induced by ψ with the form

$$1623 \quad (\psi(\lambda))_i = \max\{\lambda_i, 0\}, \quad \forall i.$$

1624 It can be verified that $\Psi(X) = \Psi(\mathcal{P}_{V_{\mathcal{I}}}(X))$, where \mathcal{I} contains all indices of the positive
 1625 eigenvalues $\lambda_i(X)$. The operator $\Psi(X)$ satisfies the low-rank property if X only has a few
 1626 positive eigenvalues. Hence, the polynomial-filtered method method can be written as

$$1627 \quad (11.17) \quad Z_{k+1} = \text{prox}_{tG}(2\mathcal{P}_+(\mathcal{P}_{U_k}(Z_k)) - Z_k) - \mathcal{P}_+(\mathcal{P}_{U_k}(Z_k)) + Z_k,$$

$$1628 \quad (11.18) \quad U_{k+1} = \text{orth}(\rho_{k+1}^{q_{k+1}}(Z_{k+1})U_k).$$

1629 **11.4.1. Examples: 2-RDM and Cryo-EM.** The two-body reduced density matrix
 1630 (2-RDM) problem can be formulated as a standard SDP. It has a block diagonal structure
 1631 with respect to the variable X , where each block is a low rank matrix. Hence, the polynomial
 1632 filters can be applied to each block to reduce the cost. As an extension, we can plug poly-
 1633 nomial filters into multi-block ADMM for the nonlinear SDPs from the weighted LS model
 1634 with spectral norm constraints and least unsquared deviations (LUD) model in orientation de-
 1635 termination of cryo-EM images [124]. For these examples we only introduce the formulation
 1636 of the corresponding model. The details of the multi-block ADMM can be found in [124].

1637 Suppose K is a given integer and S and W are two known matrices, the weighted LS
 1638 model with spectral norm constraints is

$$\begin{aligned}
 & \max \quad \langle W \odot S, G \rangle, \\
 & \text{s.t.} \quad G_{ii} = I_2, \\
 & \quad \quad G \succeq 0, \\
 & \quad \quad \|G\|_2 \leq \alpha K,
 \end{aligned}
 \tag{11.19}$$

1640 where $G = (G_{ij})_{i,j=1,\dots,K} \in \mathcal{S}^{2K}$ is the variable, with each block G_{ij} being a 2-by-2 small
 1641 matrix, and $\|\cdot\|_2$ is the spectral norm. A three-block ADMM is introduced to solve (11.19).
 1642 The cost of the projection onto the semidefinite cone can be reduced by the polynomial filters.

1643 The semidefinite relaxation of the LUD problem is

$$\begin{aligned}
 & \min \quad \sum_{1 \leq i < j \leq K} \|c_{ij} - G_{ij}c_{ji}\|_2, \\
 & \text{s.t.} \quad G_{ii} = I_2, \\
 & \quad \quad G \succeq 0, \\
 & \quad \quad \|G\|_2 \leq \alpha K,
 \end{aligned}
 \tag{11.20}$$

1645 where G , G_{ij} , K are defined the same in (11.19), and $c_{ij} \in \mathbb{R}^2$ are known vectors. The
 1646 spectral norm constraint in (11.20) is optional. A four-block ADMM is proposed to solve
 1647 (11.20). Similarly, the polynomial filters can be inserted into the ADMM update to reduce
 1648 the computational cost.

1649 **12. Conclusion.** In this paper, we provide a comprehensive survey on various sub-
 1650 space techniques for nonlinear optimization. The main idea of subspace algorithms aims
 1651 to conquer large scale nonlinear problems by performing iterations in a lower dimensional
 1652 subspace. We next summarize a few typical scenarios as follows.

- 1653 • Find a linear combination of several known directions. Examples are the linear and
 1654 nonlinear conjugate gradient methods, the Nesterov's accelerated gradient method,
 1655 the Heavy-ball method and the momentum method.
- 1656 • Keep the objective function and constraints, but add an extra restriction in a cer-
 1657 tain subspace. Examples are OMP, CoSaMP, LOBPCG, LMSVD, Arrabit, subspace
 1658 refinement and multilevel methods.
- 1659 • Approximate the objective objective function but keep the constraints. Examples are
 1660 BCD, RBR, trust region with subspaces and parallel subspace correction.
- 1661 • Approximate the objective objective function and design new constraints. Examples
 1662 are trust region with subspaces and FPC_AS.
- 1663 • Add a postprocess procedure after the subspace problem is solved. An example is
 1664 the truncated subspace method for tensor train.
- 1665 • Use subspace techniques to approximate the objective functions. Examples are sam-
 1666 pling, sketching and Nyström approximation.
- 1667 • Integrate the optimization method and subspace update in one framework. An ex-
 1668 ample is the polynomial-filtered subspace method for low-rank matrix optimization.

1669 The competitive performance of the methods adopting the above mentioned subspace
 1670 techniques in the related examples implies that the subspace methods are very promising
 1671 tools for large scale optimization problems. In fact, how to choose subspaces, how to con-
 1672 struct subproblems, and how to solve them efficiently are the key questions of designing a
 1673 successful subspace method. A good tradeoff between the simplicity of subproblems and the
 1674 computational cost has to be made carefully. We are confident that many future directions are
 1675 worth to be pursued from the point view of subspaces.

1676

REFERENCES

- 1677 [1] P.-A. ABSIL, C. G. BAKER, AND K. A. GALLIVAN, *Trust-region methods on Riemannian manifolds*,
 1678 Found. Comput. Math., 7 (2007), pp. 303–330.
- 1679 [2] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton
 1680 University Press, Princeton, NJ, 2008.
- 1681 [3] P.-A. ABSIL, R. MAHONY, AND J. TRUMPF, *An extrinsic look at the Riemannian Hessian*, in Geometric
 1682 science of information, Springer, 2013, pp. 361–368.
- 1683 [4] D. G. ANDERSON, *Iterative procedures for nonlinear integral equations*, Journal of the ACM (JACM), 12
 1684 (1965), pp. 547–560.
- 1685 [5] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*,
 1686 SIAM Journal on Imaging Sciences, 2 (2009), pp. 183–202.
- 1687 [6] S. BECKER, V. CEVHER, AND A. KYRILLIDIS, *Randomized low-memory singular value projection*, arXiv
 1688 preprint arXiv:1303.0167, (2013).
- 1689 [7] S. BELLAVIA AND B. MORINI, *A globally convergent Newton-GMRES subspace method for systems of*
 1690 *nonlinear equations*, SIAM J. Sci. Comput., 23 (2001), pp. 940–960.
- 1691 [8] S. J. BENSON, Y. YE, AND X. ZHANG, *Solving large-scale sparse semidefinite programs for combinatorial*
 1692 *optimization*, SIAM J. Optim., 10 (2000), pp. 443–461.
- 1693 [9] D. P. BERTSEKAS, *Nonlinear Programming*, Athena Scientific, September 1999.
- 1694 [10] J. BOLTE, S. SABACH, AND M. TEOULLE, *Proximal alternating linearized minimization for nonconvex*
 1695 *and nonsmooth problems*, Math. Program., 146 (2014), pp. 459–494.
- 1696 [11] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, AND J. ECKSTEIN, *Distributed optimization and statistical*
 1697 *learning via the alternating direction method of multipliers*, Foundations and Trends® in Machine
 1698 Learning, 3 (2011), pp. 1–122.
- 1699 [12] M. A. BRANCH, T. F. COLEMAN, AND Y. LI, *A subspace, interior, and conjugate gradient method for*
 1700 *large-scale bound-constrained minimization problems*, SIAM J. Sci. Comput., 21 (1999), pp. 1–23.
- 1701 [13] C. BREZINSKI, M. REDIVO-ZAGLIA, AND Y. SAAD, *Shanks sequence transformations and Anderson*
 1702 *acceleration*, SIAM Rev., 60 (2018), pp. 646–669.
- 1703 [14] P. N. BROWN AND Y. SAAD, *Hybrid Krylov methods for nonlinear systems of equations*, SIAM J. Sci.
 1704 Statist. Comput., 11 (1990), pp. 450–481.
- 1705 [15] O. BURDAKOV, L. GONG, S. ZIKRIN, AND Y.-X. YUAN, *On efficiently combining limited-memory and*
 1706 *trust-region techniques*, Math. Program. Comput., 9 (2017), pp. 101–134.
- 1707 [16] S. BURER AND R. D. C. MONTEIRO, *A projected gradient algorithm for solving the maxcut SDP relax-*
 1708 *ation*, Optim. Methods Softw., 15 (2001), pp. 175–200.
- 1709 [17] J. V. BURKE AND J. J. MORÉ, *On the identification of active constraints*, SIAM J. Numer. Anal., 25 (1988),
 1710 pp. 1197–1211.
- 1711 [18] ———, *Exposing constraints*, SIAM J. Optim., 4 (1994), pp. 573–595.
- 1712 [19] R. H. BYRD, N. I. M. GOULD, J. NOCEDAL, AND R. A. WALTZ, *An algorithm for nonlinear optimization*
 1713 *using linear programming and equality constrained subproblems*, Math. Program., 100 (2004), pp. 27–
 1714 48.
- 1715 [20] ———, *On the convergence of successive linear-quadratic programming algorithms*, SIAM J. Optim., 16
 1716 (2005), pp. 471–489.
- 1717 [21] R. H. BYRD, J. NOCEDAL, AND R. B. SCHNABEL, *Representations of quasi-Newton matrices and their*
 1718 *use in limited memory methods*, Math. Programming, 63 (1994), pp. 129–156.
- 1719 [22] C. CARSTENSEN, *Domain decomposition for a non-smooth convex minimization problem and its applica-*
 1720 *tion to plasticity*, Numerical linear algebra with applications, 4 (1997), pp. 177–190.
- 1721 [23] M. R. CELIS, J. E. DENNIS, AND R. A. TAPIA, *A trust region strategy for nonlinear equality constrained*
 1722 *optimization*, in Numerical optimization, 1984 (Boulder, Colo., 1984), SIAM, Philadelphia, PA, 1985,
 1723 pp. 71–82.
- 1724 [24] C. CHEN, Z. WEN, AND Y.-X. YUAN, *A general two-level subspace method for nonlinear optimization*, J.
 1725 Comput. Math., 36 (2018), pp. 881–902.
- 1726 [25] Y. CHEN, X. LI, AND J. XU, *Convexified modularity maximization for degree-corrected stochastic block*

- 1727 *models*, Ann. Statist., 46 (2018), pp. 1573–1602.
- 1728 [26] A. CONN, N. GOULD, A. SARTENAER, AND P. TOINT, *On iterated-subspace methods for nonlinear opt-*
1729 *imization*, in Linear and Nonlinear Conjugate Gradient-Related Methods, J. Adams and J. Nazareth,
1730 eds., 1996, pp. 50–79.
- 1731 [27] W. DENG, M.-J. LAI, Z. PENG, AND W. YIN, *Parallel multi-block ADMM with $o(1/k)$ convergence*, J.
1732 Sci. Comput., 71 (2017), pp. 712–736.
- 1733 [28] E. D. DOLAN, J. J. MORÉ, AND T. S. MUNSON, *Benchmarking optimization software with cops 3.0*, tech.
1734 rep., Mathematics and Computer Science Division, Argonne National Laboratory, February 2004.
- 1735 [29] Q. DONG, X. LIU, Z.-W. WEN, AND Y.-X. YUAN, *A parallel line search subspace correction method for*
1736 *composite convex optimization*, J. Oper. Res. Soc. China, 3 (2015), pp. 163–187.
- 1737 [30] J. DOUGLAS AND H. H. RACHFORD, *On the numerical solution of heat conduction problems in two and*
1738 *three space variables*, Transactions of the American mathematical Society, 82 (1956), pp. 421–439.
- 1739 [31] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*,
1740 SIAM J. Matrix Anal. Appl., 20 (1999), pp. 303–353.
- 1741 [32] M. ELAD, B. MATALON, AND M. ZIBULEVSKY, *Coordinate and subspace optimization methods for linear*
1742 *least squares with non-quadratic regularization*, Appl. Comput. Harmon. Anal., 23 (2007), pp. 346–
1743 367.
- 1744 [33] L. ELDÉN, *Matrix Methods in Data Mining and Pattern Recognition (Fundamentals of Algorithms)*, Society
1745 for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2007.
- 1746 [34] H.-R. FANG AND Y. SAAD, *A filtered Lanczos procedure for extreme and interior eigenvalue problems*,
1747 SIAM J. Sci. Comput., 34 (2012), pp. A2220–A2246.
- 1748 [35] R. FLETCHER, *Second order corrections for nondifferentiable optimization*, in Numerical analysis (Dundee,
1749 1981), vol. 912 of Lecture Notes in Math., Springer, Berlin-New York, 1982, pp. 85–114.
- 1750 [36] M. FORNASIER, *Domain decomposition methods for linear inverse problems with sparsity constraints*, In-
1751 verse Problems, 23 (2007), p. 2505.
- 1752 [37] M. FORNASIER, Y. KIM, A. LANGER, AND C.-B. SCHÖNLIEB, *Wavelet decomposition method for l_2/tv -*
1753 *image deblurring*, SIAM Journal on Imaging Sciences, 5 (2012), pp. 857–885.
- 1754 [38] M. FORNASIER, A. LANGER, AND C.-B. SCHÖNLIEB, *A convergent overlapping domain decomposition*
1755 *method for total variation minimization*, Numerische Mathematik, 116 (2010), pp. 645–685.
- 1756 [39] M. FORNASIER AND C.-B. SCHÖNLIEB, *Subspace correction methods for total variation and l_1 -*
1757 *minimization*, SIAM Journal on Numerical Analysis, 47 (2009), pp. 3397–3428.
- 1758 [40] D. GABAY, *Minimizing a differentiable function over a differential manifold*, J. Optim. Theory Appl., 37
1759 (1982), pp. 177–219.
- 1760 [41] D. GABAY AND B. MERCIER, *A dual algorithm for the solution of nonlinear variational problems via finite*
1761 *element approximation*, Computers & Mathematics with Applications, 2 (1976), pp. 17–40.
- 1762 [42] E. GHADIMI, H. R. FEYZMAHDAVIAN, AND M. JOHANSSON, *Global convergence of the heavy-ball*
1763 *method for convex optimization*, in 2015 European Control Conference (ECC), 2015, pp. 310–315.
- 1764 [43] P. E. GILL AND M. W. LEONARD, *Reduced-Hessian quasi-Newton methods for unconstrained optimiza-*
1765 *tion*, SIAM J. Optim., 12 (2001), pp. 209–237.
- 1766 [44] ———, *Limited-memory reduced-Hessian methods for large-scale unconstrained optimization*, SIAM J.
1767 Optim., 14 (2003), pp. 380–401.
- 1768 [45] R. GLOWINSKI AND A. MARROCCO, *Sur l’approximation, par éléments finis d’ordre un, et la résolution,*
1769 *par pénalisation-dualité, d’une classe de problèmes de Dirichlet non linéaires*, Rev. Française Automat.
1770 Informat. Recherche Opérationnelle Sér. Rouge Anal. Numér., 9 (1975), pp. 41–76.
- 1771 [46] M. X. GOEMANS AND D. P. WILLIAMSON, *Improved approximation algorithms for maximum cut and*
1772 *satisfiability problems using semidefinite programming*, J. Assoc. Comput. Mach., 42 (1995), pp. 1115–
1773 1145.
- 1774 [47] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep Learning*, MIT Press, 2016. [http://www.](http://www.deeplearningbook.org)
1775 [deeplearningbook.org](http://www.deeplearningbook.org).
- 1776 [48] N. GOULD, D. ORBAN, AND P. TOINT, *Numerical methods for large-scale nonlinear optimization*, Acta
1777 Numer., 14 (2005), pp. 299–361.
- 1778 [49] G. N. GRAPIGLIA, J. YUAN, AND Y.-X. YUAN, *A subspace version of the powell–yuan trust-region al-*
1779 *gorithm for equality constrained optimization*, Journal of the Operations Research Society of China, 1
1780 (2013), pp. 425–451.
- 1781 [50] G. N. GRAPIGLIA, Y. YUAN, AND Y.-X. YUAN, *A subspace version of the powell–yuan trust-region al-*
1782 *gorithm for equality constrained optimization*, J. Operations Research Society of China, 1 (2013),
1783 pp. 425–451.
- 1784 [51] W. W. HAGER AND H. ZHANG, *A new active set algorithm for box constrained optimization*, SIAM J.
1785 Optim., 17 (2006), pp. 526–557.
- 1786 [52] ———, *A new active set algorithm for box constrained optimization*, SIAM J. Optim., 17 (2006), pp. 526–
1787 557.
- 1788 [53] E. T. HALE, W. YIN, AND Y. ZHANG, *Fixed-point continuation for l_1 -minimization: methodology and*

- 1789 convergence, *SIAM J. Optim.*, 19 (2008), pp. 1107–1130.
- 1790 [54] N. HALKO, P. G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: probabilistic*
1791 *algorithms for constructing approximate matrix decompositions*, *SIAM Rev.*, 53 (2011), pp. 217–288.
- 1792 [55] B. HE, H.-K. XU, AND X. YUAN, *On the proximal Jacobian decomposition of ALM for multiple-block*
1793 *separable convex minimization problems and its relationship to ADMM*, *J. Sci. Comput.*, 66 (2016),
1794 pp. 1204–1217.
- 1795 [56] C. HELMBERG AND F. RENDL, *A spectral bundle method for semidefinite programming*, *SIAM J. Optim.*,
1796 10 (2000), pp. 673–696.
- 1797 [57] J. M. HOKANSON, *Projected nonlinear least squares for exponential fitting*, *SIAM J. Sci. Comput.*, 39
1798 (2017), pp. A3107–A3128.
- 1799 [58] J. HU, B. JIANG, L. LIN, Z. WEN, AND Y.-X. YUAN, *Structured quasi-Newton methods for optimization*
1800 *with orthogonality constraints*, *SIAM J. Sci. Comput.*, 41 (2019), pp. A2239–A2269.
- 1801 [59] J. HU, A. MILZAREK, Z. WEN, AND Y. YUAN, *Adaptive quadratically regularized Newton method for*
1802 *Riemannian optimization*, *SIAM J. Matrix Anal. Appl.*, 39 (2018), pp. 1181–1207.
- 1803 [60] W. HU, L. LIN, AND C. YANG, *Projected commutator DIIS method for accelerating hybrid functional*
1804 *electronic structure calculations*, *J. Chem. Theory Comput.*, (2017).
- 1805 [61] W. HUANG, P.-A. ABSIL, AND K. GALLIVAN, *A Riemannian BFGS method without differentiated retraction*
1806 *for nonconvex optimization problems*, *SIAM J. Optim.*, 28 (2018), pp. 470–495.
- 1807 [62] W. HUANG, K. A. GALLIVAN, AND P.-A. ABSIL, *A Broyden class of quasi-Newton methods for Riemannian*
1808 *optimization*, *SIAM J. Optim.*, 25 (2015), pp. 1660–1685.
- 1809 [63] F. KANGAL, K. MEERBERGEN, E. MENGI, AND W. MICHIELS, *A subspace method for large-scale eigen-*
1810 *value optimization*, *SIAM Journal on Matrix Analysis and Applications*, 39 (2018), pp. 48–82.
- 1811 [64] N. KESKAR, J. NOCEDAL, F. ÖZTOPRAK, AND A. WÄCHTER, *A second-order method for convex ℓ_1 -*
1812 *regularized optimization with active-set prediction*, *Optim. Methods Softw.*, 31 (2016), pp. 605–621.
- 1813 [65] A. V. KNYAZEV, *Toward the optimal preconditioned eigensolver: locally optimal block preconditioned con-*
1814 *jugate gradient method*, *SIAM J. Sci. Comput.*, 23 (2001), pp. 517–541. Copper Mountain Conference
1815 (2000).
- 1816 [66] D. KRESSNER, D. LU, AND B. VANDEREYCKEN, *Subspace acceleration for the crawford number and re-*
1817 *lated eigenvalue optimization problems*, *SIAM Journal on Matrix Analysis and Applications*, 39 (2018),
1818 pp. 961–982.
- 1819 [67] K. KREUTZ-DELGADO, *The complex gradient operator and the CR-calculus*, 2009.
1820 <http://arxiv.org/abs/0906.4835>.
- 1821 [68] A. LANGER, S. OSHER, AND C.-B. SCHÖNLIEB, *Bregmanized domain decomposition for image restora-*
1822 *tion*, *Journal of Scientific Computing*, 54 (2013), pp. 549–576.
- 1823 [69] C. LE BRIS, *Computational chemistry from the perspective of numerical analysis*, *Acta Numer.*, 14 (2005),
1824 pp. 363–444.
- 1825 [70] J. H. LEE, Y. M. JUNG, Y.-X. YUAN, AND S. YUN, *A subspace SQP method for equality constrained*
1826 *optimization*, *Comput. Optim. Appl.*, 74 (2019), pp. 177–194.
- 1827 [71] Y. LI, H. LIU, Z. WEN, AND Y. YUAN, *Low-rank matrix optimization using polynomial-filtered subspace*
1828 *extraction*.
- 1829 [72] Y. LI AND S. OSHER, *Coordinate descent optimization for ℓ^1 minimization with application to compressed*
1830 *sensing; a greedy algorithm*, *Inverse Probl. Imaging*, 3 (2009), pp. 487–503.
- 1831 [73] L. LIN, *Adaptively compressed exchange operator*, *J. Chem. Theory Comput.*, 12 (2016), pp. 2242–2249.
- 1832 [74] X. LIU, Z. WEN, AND Y. ZHANG, *Limited memory block Krylov subspace optimization for computing*
1833 *dominant singular value decompositions*, *SIAM J. Sci. Comput.*, 35 (2013), pp. A1641–A1668.
- 1834 [75] ———, *An efficient Gauss-Newton algorithm for symmetric low-rank product matrix approximations*, *SIAM*
1835 *J. Optim.*, 25 (2015), pp. 1571–1608.
- 1836 [76] Z. LIU AND L. VANDENBERGHE, *Interior-point method for nuclear norm approximation with application*
1837 *to system identification*, *SIAM Journal on Matrix Analysis and Applications*, 31 (2009), pp. 1235–1256.
- 1838 [77] Z. Q. LUO AND P. TSENG, *On the convergence of the coordinate descent method for convex differentiable*
1839 *minimization*, *J. Optim. Theory Appl.*, 72 (1992), pp. 7–35.
- 1840 [78] M. W. MAHONEY, *Randomized algorithms for matrices and data*, *Foundations and Trends[®] in Machine*
1841 *Learning*, 3 (2011), pp. 123–224.
- 1842 [79] J. MARTENS AND R. GROSSE, *Optimizing neural networks with kronecker-factored approximate curvature*,
1843 in *International Conference on Machine Learning*, 2015, pp. 2408–2417.
- 1844 [80] D. Q. MAYNE AND E. POLAK, *A superlinearly convergent algorithm for constrained optimization prob-*
1845 *lems*, *Math. Programming Stud.*, (1982), pp. 45–61.
- 1846 [81] J. J. MORÉ AND G. TORALDO, *Algorithms for bound constrained quadratic programming problems*, *Nu-*
1847 *mer. Math.*, 55 (1989), pp. 377–400.
- 1848 [82] ———, *On the solution of large quadratic programming problems with bound constraints*, *SIAM J. Optim.*,
1849 1 (1991), pp. 93–113.
- 1850 [83] D. NEEDELL AND J. A. TROPP, *CoSaMP: iterative signal recovery from incomplete and inaccurate sam-*

- 1851 *ples*, Appl. Comput. Harmon. Anal., 26 (2009), pp. 301–321.
- 1852 [84] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science &
- 1853 Business Media, 2013.
- 1854 [85] Y. E. NESTEROV, *A method for solving the convex programming problem with convergence rate $o(1/k^2)$* , in
- 1855 Dokl. Akad. Nauk SSSR, vol. 269, 1983, pp. 543–547.
- 1856 [86] M. E. NEWMAN, *Modularity and community structure in networks*, Proceedings of the national academy of
- 1857 sciences, 103 (2006), pp. 8577–8582.
- 1858 [87] Q. NI AND Y. YUAN, *A subspace limited memory quasi-Newton algorithm for large-scale nonlinear bound*
- 1859 *constrained optimization*, Math. Comp., 66 (1997), pp. 1509–1520.
- 1860 [88] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer Series in Operations Research and
- 1861 Financial Engineering, Springer, New York, second ed., 2006.
- 1862 [89] X. PAN, M. LAM, S. TU, D. PAPALIOPOULOS, C. ZHANG, M. I. JORDAN, K. RAMCHANDRAN, C. RE,
- 1863 AND B. RECHT, *Cyclades: Conflict-free asynchronous machine learning*, in Advances in Neural Infor-
- 1864 mation Processing Systems, 2016, p. 2576–2584.
- 1865 [90] B. POLYAK, *Some methods of speeding up the convergence of iteration methods*, USSR Computational
- 1866 Mathematics and Mathematical Physics, 4 (1964), pp. 1 – 17.
- 1867 [91] M. J. D. POWELL AND Y. YUAN, *A trust region algorithm for equality constrained optimization*, Math.
- 1868 Programming, 49 (1990/91), pp. 189–211.
- 1869 [92] P. PULAY, *Convergence acceleration of iterative sequences. the case of SCF iteration*, Chemical Physics
- 1870 Letters, 73 (1980), pp. 393–398.
- 1871 [93] P. PULAY, *Improved SCF convergence acceleration*, Journal of Computational Chemistry, 3 (1982),
- 1872 pp. 556–560.
- 1873 [94] B. RECHT, C. RE, S. WRIGHT, AND F. NIU, *Hogwild: A lock-free approach to parallelizing stochastic*
- 1874 *gradient descent*, in Advances in Neural Information Processing Systems, 2011, pp. 693–701.
- 1875 [95] H. RUTISHAUSER, *Computational aspects of F. L. Bauer’s simultaneous iteration method*, Numer. Math.,
- 1876 13 (1969), pp. 4–13.
- 1877 [96] H. RUTISHAUSER, *Simultaneous iteration method for symmetric matrices*, Numer. Math., 16 (1970),
- 1878 pp. 205–223.
- 1879 [97] Y. SAAD, *Chebyshev acceleration techniques for solving nonsymmetric eigenvalue problems*, Mathematics
- 1880 of Computation, 42 (1984), pp. 567–588.
- 1881 [98] Y. SAAD, *Iterative methods for sparse linear systems*, Society for Industrial and Applied Mathematics,
- 1882 Philadelphia, PA, second ed., 2003.
- 1883 [99] D. SCIEUR, A. D’ASPREMONT, AND F. BACH, *Regularized nonlinear acceleration*, in Advances In Neural
- 1884 Information Processing Systems, 2016, pp. 712–720.
- 1885 [100] S. K. SHEVADE AND S. S. KEERTHI, *A simple and efficient algorithm for gene selection using sparse*
- 1886 *logistic regression*, Bioinformatics, 19 (2003), pp. 2246–2253.
- 1887 [101] G. A. SHULTZ, R. B. SCHNABEL, AND R. H. BYRD, *A family of trust-region-based algorithms for uncon-*
- 1888 *strained minimization with strong global convergence properties*, SIAM J. Numer. Anal., 22 (1985),
- 1889 pp. 47–67.
- 1890 [102] P. SIRKOVIC AND D. KRESSNER, *Subspace acceleration for large-scale parameter-dependent Hermitian*
- 1891 *eigenproblems*, SIAM Journal on Matrix Analysis and Applications, 37 (2016), pp. 695–718.
- 1892 [103] A. SIT, Z. WU, AND Y. YUAN, *A geometric buildup algorithm for the solution of the distance geometry*
- 1893 *problem using least-squares approximation*, Bull. Math. Biol., 71 (2009), pp. 1914–1933.
- 1894 [104] S. T. SMITH, *Optimization techniques on Riemannian manifolds*, Fields Institute Communications, 3 (1994).
- 1895 [105] S. SOLNTSEV, J. NOCEDAL, AND R. H. BYRD, *An algorithm for quadratic ℓ_1 -regularized optimization*
- 1896 *with a flexible active-set strategy*, Optim. Methods Softw., 30 (2015), pp. 1213–1237.
- 1897 [106] M. SOLTANI AND C. HEGDE, *Fast low-rank matrix estimation without the condition number*, arXiv preprint
- 1898 arXiv:1712.03281, (2017).
- 1899 [107] T. STEIHAUG, *The conjugate gradient method and trust regions in large scale optimization*, SIAM J. Numer.
- 1900 Anal., 20 (1983), pp. 626–637.
- 1901 [108] G. W. STEWART, *Simultaneous iteration for computing invariant subspaces of non-Hermitian matrices*,
- 1902 Numer. Math., 25 (1975/76), pp. 123–136.
- 1903 [109] ———, *Matrix algorithms Vol. II: Eigensystems*, Society for Industrial and Applied Mathematics (SIAM),
- 1904 Philadelphia, PA, 2001.
- 1905 [110] W. J. STEWART AND A. JENNINGS, *A simultaneous iteration algorithm for real matrices*, ACM Trans.
- 1906 Math. Software, 7 (1981), pp. 184–198.
- 1907 [111] W. SUN AND Y. YUAN, *Optimization theory and methods: nonlinear programming*, vol. 1, Springer Science
- 1908 & Business Media, 2006.
- 1909 [112] X.-C. TAI AND J. XU, *Global and uniform convergence of subspace correction methods for some convex*
- 1910 *optimization problems*, Mathematics of Computation, 71 (2002), pp. 105–124.
- 1911 [113] ———, *Global and uniform convergence of subspace correction methods for some convex optimization*
- 1912 *problems*, Math. Comp., 71 (2002), pp. 105–124.

- 1913 [114] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society.
 1914 Series B (Methodological), (1996), pp. 267–288.
- 1915 [115] A. TOTH, J. A. ELLIS, T. EVANS, S. HAMILTON, C. KELLEY, R. PAWLOWSKI, AND S. SLATTERY, *Local*
 1916 *improvement results for Anderson acceleration with inaccurate function evaluations*, SIAM Journal on
 1917 Scientific Computing, 39 (2017), pp. S47–S65.
- 1918 [116] J. A. TROPP AND A. C. GILBERT, *Signal recovery from random measurements via orthogonal matching*
 1919 *pursuit*, IEEE Trans. Inform. Theory, 53 (2007), pp. 4655–4666.
- 1920 [117] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Fixed-rank approximation of a positive-*
 1921 *semidefinite matrix from streaming data*, in Advances in Neural Information Processing Systems, 2017,
 1922 pp. 1225–1234.
- 1923 [118] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for low-rank*
 1924 *matrix approximation*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 1454–1485.
- 1925 [119] P. TSENG AND S. YUN, *A coordinate gradient descent method for nonsmooth separable minimization*,
 1926 Mathematical Programming, 117 (2009), pp. 387–423.
- 1927 [120] C. UDRISTE, *Convex functions and optimization methods on Riemannian manifolds*, vol. 297, Springer
 1928 Science & Business Media, 1994.
- 1929 [121] B. VANDEREYCKEN, *Low-rank matrix completion by Riemannian optimization*, SIAM Journal on Opti-
 1930 mization, 23 (2013), pp. 1214–1236.
- 1931 [122] I. WALDSPURGER, A. D’ASPREMONT, AND S. MALLAT, *Phase recovery, MaxCut and complex semidefi-*
 1932 *nite programming*, Math. Program., 149 (2015), pp. 47–81.
- 1933 [123] H. F. WALKER AND P. NI, *Anderson acceleration for fixed-point iterations*, SIAM Journal on Numerical
 1934 Analysis, 49 (2011), pp. 1715–1735.
- 1935 [124] L. WANG, A. SINGER, AND Z. WEN, *Orientation determination of cryo-em images using least unsquared*
 1936 *deviations*, SIAM journal on imaging sciences, 6 (2013), pp. 2450–2483.
- 1937 [125] X. WANG, M. HONG, S. MA, AND Z.-Q. LUO, *Solving multiple-block separable convex minimization*
 1938 *problems using two-block alternating direction method of multipliers*, arXiv preprint arXiv:1308.5294,
 1939 (2013).
- 1940 [126] Y. WANG, Z. JIA, AND Z. WEN, *The search direction correction makes first-order methods faster*, (2019).
 1941 Arxiv 1905.06507.
- 1942 [127] Z. WANG, Z. WEN, AND Y. YUAN, *A subspace trust region method for large scale unconstrained opti-*
 1943 *mization*, in Numerical Linear Algebra and Optimization, Y.Yuan, ed., 2004, pp. 265–274.
- 1944 [128] Z.-H. WANG AND Y.-X. YUAN, *A subspace implementation of quasi-Newton trust region methods for*
 1945 *unconstrained optimization*, Numer. Math., 104 (2006), pp. 241–269.
- 1946 [129] Z. WEN, D. GOLDFARB, AND K. SCHEINBERG, *Block coordinate descent methods for semidefinite pro-*
 1947 *gramming*, Handbook on Semidefinite, Cone and Polynomial Optimization, (2011).
- 1948 [130] ———, *Block coordinate descent methods for semidefinite programming*, in Handbook on semidefinite,
 1949 conic and polynomial optimization, vol. 166 of Internat. Ser. Oper. Res. Management Sci., Springer,
 1950 New York, 2012, pp. 533–564.
- 1951 [131] Z. WEN, A. MILZAREK, M. ULBRICH, AND H. ZHANG, *Adaptive regularized self-consistent field iteration*
 1952 *with exact Hessian for electronic structure calculation*, SIAM J. Sci. Comput., 35 (2013), pp. A1299–
 1953 A1324.
- 1954 [132] Z. WEN AND W. YIN, *A feasible method for optimization with orthogonality constraints*, Math. Program.,
 1955 142 (2013), pp. 397–434.
- 1956 [133] Z. WEN, W. YIN, D. GOLDFARB, AND Y. ZHANG, *A fast algorithm for sparse reconstruction based on*
 1957 *shrinkage, subspace optimization and continuation*, SIAM Journal on Scientific Computing, 32 (2010),
 1958 pp. 1832–1857.
- 1959 [134] Z. WEN, W. YIN, H. ZHANG, AND D. GOLDFARB, *On the convergence of an active-set method for l_1*
 1960 *minimization*, Optimization Methods and Software, 27 (2012), pp. 1127–1146.
- 1961 [135] Z. WEN AND Y. ZHANG, *Accelerating convergence by augmented rayleigh–ritz projections for large-scale*
 1962 *eigenpair computation*, SIAM Journal on Matrix Analysis and Applications, 38 (2017), pp. 273–296.
- 1963 [136] D. P. WOODRUFF, *Sketching as a tool for numerical linear algebra*, Found. Trends Theor. Comput. Sci., 10
 1964 (2014), pp. iv+157.
- 1965 [137] X. WU, Z. WEN, AND W. BAO, *A regularized newton method for computing ground states of Bose-Einstein*
 1966 *condensates*, arXiv preprint arXiv:1504.02891, (2015).
- 1967 [138] C. YANG, J. C. MEZA, AND L.-W. WANG, *A trust region direct constrained minimization algorithm for*
 1968 *the Kohn-Sham equation*, SIAM J. Sci. Comput., 29 (2007), pp. 1854–1875.
- 1969 [139] Y. YANG, B. DONG, AND Z. WEN, *Randomized algorithms for high quality treatment planning in volu-*
 1970 *metric modulated arc therapy*, Inverse Problems, 33 (2017), pp. 025007, 22.
- 1971 [140] Y.-X. YUAN, *Subspace techniques for nonlinear optimization*, in Some topics in industrial and applied
 1972 mathematics, vol. 8 of Ser. Contemp. Appl. Math. CAM, Higher Ed. Press, Beijing, 2007, pp. 206–
 1973 218.
- 1974 [141] Y.-X. YUAN, *Subspace methods for large scale nonlinear equations and nonlinear least squares*, Optim.

- 1975 Eng., 10 (2009), pp. 207–218.
- 1976 [142] Y.-X. YUAN, *A review on subspace methods for nonlinear optimization*, in Proceedings of the International
 1977 Congress of Mathematicians—Seoul 2014. Vol. IV, Kyung Moon Sa, Seoul, 2014, pp. 807–827.
- 1978 [143] Y.-X. YUAN AND J. STOER, *A subspace study on conjugate gradient algorithms*, Z. Angew. Math. Mech.,
 1979 75 (1995), pp. 69–77.
- 1980 [144] A. YURTSEVER, J. A. TROPP, O. FERCOQ, M. UDELL, AND V. CEVHER, *Scalable semidefinite program-*
 1981 *ming*, 2019. arXiv:1912.02949.
- 1982 [145] J. ZHANG, H. LIU, Z. WEN, AND S. ZHANG, *A sparse completely positive relaxation of the modularity*
 1983 *maximization for community detection*, SIAM J. Sci. Comput., 40 (2018), pp. A3091–A3120.
- 1984 [146] ———, *A sparse completely positive relaxation of the modularity maximization for community detection*,
 1985 SIAM J. Sci. Comput., 40 (2018), pp. A3091–A3120.
- 1986 [147] J. ZHANG, B. O’DONOGHUE, AND S. BOYD, *Globally convergent Type-I Anderson acceleration for non-*
 1987 *smooth fixed-point iterations*, (2018). arXiv:1808.03971.
- 1988 [148] J. ZHANG, Z. WEN, AND Y. ZHANG, *Subspace methods with local refinements for eigenvalue computation*
 1989 *using low-rank tensor-train format*, Journal of Scientific Computing, 70 (2017), pp. 478–499.
- 1990 [149] T. ZHOU AND D. TAO, *Godec: Randomized low-rank & sparse matrix decomposition in noisy case*, in
 1991 International conference on machine learning, Omnipress, 2011.
- 1992 [150] Y. ZHOU AND Y. SAAD, *A Chebyshev–Davidson algorithm for large symmetric eigenproblems*, SIAM J.
 1993 Matrix Anal. and Appl., 29 (2007), pp. 954–971.
- 1994 [151] Y. ZHOU, Y. SAAD, M. L. TIAGO, AND J. R. CHELIKOWSKY, *Self-consistent-field calculations using*
 1995 *Chebyshev-filtered subspace iteration*, Journal of Computational Physics, 219 (2006), pp. 172–184.