

Optimal Diagonal Preconditioning Beyond Worst-Case Conditioning: Theory and Practice of Omega Scaling ^{*†}

Saeed Ghadimi[‡] Woosuk L. Jung[§] Arnesh Sujana[§] David Torregrosa-Belén[¶]
Henry Wolkowicz[§]

Revising as of May 12, 2026, 8:12pm

Key words and phrases: κ , ω -condition numbers, diagonal preconditioning, iterative methods, PCG, linear systems, eigenvalues

AMS subject classifications: 15A12, 65F35, 65F08, 65G50, 49J52, 49K10, 90C32

Contents

10	1 Introduction	3
11	1.1 Main Contributions	4
12	1.2 Background; Preliminaries	5
13	1.2.1 Notation	6
14	2 κ-Optimal Diagonal Preconditioning	7
15	2.1 Optimality Conditions	7
16	2.1.1 Characterization of κ -Optimal Scaling	9
17	2.2 A Projected Subgradient Method for Minimizing $\kappa(d)$	12
18	2.2.1 Asymptotic Convergence Analysis	13
19	2.3 Avoiding Positive Homogeneity in Minimizing $\kappa(d)$	15
20	2.3.1 Nonasymptotic Convergence Rate Analysis for Minimizing $\kappa(v)$ on $\hat{\Omega}$	17
21	3 ω-Optimal Structured Preconditioning	18
22	3.1 Right- and Left-Sided Diagonal	18
23	3.1.1 Right-Sided Diagonal	19
24	3.1.2 Left-Sided Diagonal	19
25	3.2 Two-Sided Diagonal	20
26	3.2.1 Optimality Conditions for Two-Sided Problem	20
27	3.2.2 Square-Root Sinkhorn-Knopp Algorithm for Two Sided Problem	22
28	3.3 Right-, Left-Sided Block Diagonal	23
29	3.3.1 Right-Sided Block Diagonal	23
30	3.3.2 Left-Sided Block Diagonal	24

*Emails resp.: sghadimi@uwaterloo.ca, w2jung@uwaterloo.ca, a3sujana@uwaterloo.ca, david.torregrosa@ua.es, hwolkowicz@uwaterloo.ca

[†]This report is available at URL: www.math.uwaterloo.ca/~hwolkowi/henry/reports/ABSTRACTS.html A previous version of this work appeared on arXiv under the title “New Insights and Algorithms for Optimal Diagonal Preconditioning.”

[‡]Department of Management Science and Engineering, University of Waterloo, ON, Canada

[§]Department of Combinatorics and Optimization, University of Waterloo, ON, Canada

[¶]Department of Mathematics, University of Alicante, Alicante, Spain

31	4 Computational Experiments	26
32	4.1 Minimizing κ Efficiently	26
33	4.1.1 Moderate Sizes	26
34	4.1.2 Large Sizes	28
35	4.2 PCG Comparison for Solving Linear Systems	28
36	4.3 LSQR Comparison: Algorithm 3.1 vs. Two-Sided κ -optimal Scaling in [24]	29
37	4.4 From κ -optimal M to “Improve PCG” using ω -Optimal Scaling	31
38	5 Conclusion	32
39	A Assumptions and Technical Results from [19]	34
40	B An Efficient Line-Search Subgradient Method	34
41	C List of SuiteSparse Matrices Used in Experiments	35
42	Index	37
43	Bibliography	38

44 List of Algorithms

45	2.1 A Subgradient Method for Minimizing $\kappa(d)$ over Ω	12
46	2.2 A Subgradient Method for Minimizing $\kappa(v) := \kappa(\mathcal{V}(v))$ over $\hat{\Omega}$	16
47	3.1 Square-Root Sinkhorn-Knopp Algorithm for Two Sided ω -Optimal Preconditioner	22
48	B.1 An Efficient Line-Search Subgradient Method for Minimizing $\kappa(v)$	35

49 List of Tables

50	3.1 Relation of ω -Optimal Preconditioners to the Literature	18
51	4.1 Comparison of Algorithms for Min κ on SuiteSparse	27
52	4.2 Comparison of Algorithms for Min κ on Random	27
53	4.3 Comparison of Algorithms for Min κ on Large SuiteSparse	28
54	4.4 Comparison of Algorithms for Min κ on Large Random	28
55	4.5 PCG Comparison Using Preconditioners Found by Algorithm in [24] and Algorithm B.1	29
56	4.6 LSQR Comparison on Small Suitesparse Matrices: Algorithm 3.1 vs [24]	30
57	4.7 LSQR Comparison on Small Random Matrices: Algorithm 3.1 vs [24]	30
58	4.8 LSQR Comparison on Large Suitesparse Matrices: Algorithm 3.1 vs No Preconditioning	31
59	4.9 PCG tol. $1e-7$; Medium; κ -opt $\underline{\mathbf{VS}} J, \omega$ -opt of M	32
60	4.10 PCG tol. $1e-7$; Large; κ -opt $\underline{\mathbf{VS}} J, \omega$ -opt of M	32

61 List of Figures

62

Abstract

63 Preconditioning is essential in many areas of mathematics, and in particular is a fundamental tool
64 for accelerating iterative methods for solving linear systems. In this work, we study optimal diagonal
65 preconditioning under two distinct notions of conditioning: the classical worst-case κ -condition number
66 and the averaging-based ω -condition number. We observe that ω -optimal preconditioning generally
67 outperforms κ -optimal preconditioning.

68 For the κ -optimal preconditioning problem, we derive an affine-based pseudoconvex reformulation
69 with three key advantages: all stationary points are global minima, subgradients are inexpensive to

70 compute, and the optimization variable is an n -dimensional vector rather than an $n \times n$ matrix as
71 in semidefinite programming (SDP) approaches. We develop a simple and highly efficient subgradient
72 method, with convergence guarantees, for solving this pseudoconvex formulation. Numerical results
73 indicate that our approach is substantially more scalable, efficient, and accurate than existing SDP-
74 based methods, often achieving dramatic speedups.

75 For the ω -condition number, we provide explicit characterizations of optimal diagonal and block
76 diagonal preconditioners. In particular, we show that several classical preconditioners, including Jacobi
77 and row/column normalization, are ω -optimal, and that matrix balancing schemes monotonically reduce
78 ω and converge to stationary points of the two-sided problem. To the best of our knowledge, this is the
79 first unified and explicit characterization of optimality conditions for both κ and ω -based preconditioning.

80 Our numerical experiments further reveal a striking phenomenon: although κ -optimal preconditioners
81 achieve stronger reductions in the worst-case condition number, ω -optimal preconditioners are substan-
82 tially cheaper to compute and yield better performance for iterative methods such as preconditioned con-
83 jugate gradient (PCG) and least squares method (LSQR). Moreover, applying ω -optimal scaling to linear
84 systems that are already κ -optimally preconditioned leads to further improvements in PCG iterations.
85 These results suggest that the ω -condition number is more predictive of practical solver performance and
86 highlight the advantages of ω -based preconditioning in large-scale settings.

87 1 Introduction

88 Preconditioning plays a central role in accelerating iterative methods for large-scale linear systems, e.g., [10–
89 12, 24]. All of these works focus on studying the classical κ -condition number of a matrix A (the ratio of
90 largest to smallest singular values). On the other hand, another line of research has focused on studying
91 the ω -condition number, which is the ratio of the arithmetic and geometric mean of singular values, [7, 17].
92 Although κ -optimal preconditioners minimize the worst-case condition number, our results show that ω -
93 optimal preconditioners are cheaper to compute and are more predictive of PCG/LSQR performance.

94 This work is divided into three complementary parts.

95 (i): First, motivated by the recent work in [12, 24] that evaluates κ -optimal diagonal preconditioners using
96 convex semidefinite programming relaxations, we provide a different approach based on using the nonconvex
97 formulation with κ . In particular, we provide a modified model that exploits the invariance of eigenvalues
98 for a product of matrices and solves an essentially unconstrained pseudoconvex minimization problem. Our
99 approach is significantly more efficient and robust, often achieving speedups of over 100x compared with the
100 methods in [12, 24].

101 (ii): Second, we consider properties of the ω -condition number. In particular, we illustrate how to obtain
102 explicit formulae for ω -optimal diagonal and block-diagonal preconditioners. In both settings, we obtain
103 either known or new preconditioners, including the Jacobi preconditioner, column and row normalization
104 preconditioners, and those obtained via matrix balancing algorithms such as Sinkhorn-Knopp. Compared to
105 κ -optimal preconditioners, the ω -optimal counterparts are substantially cheaper to compute while retaining
106 strong practical effectiveness.

107 (iii): Finally, in our computational experiments we demonstrate the robustness and efficiency of our algo-
108 rithms for finding κ -optimal preconditioners. We also demonstrate major advantages of using ω -optimal
109 preconditioners for iteratively solving linear systems using preconditioned conjugate gradient (PCG) for pos-
110 itive definite systems and least squares method (LSQR) for general invertible systems. In general, we find
111 that reductions in the ω -condition number are more strongly correlated with improvements in LSQR and
112 PCG iteration counts than reductions in the κ -condition number.

113 For simplicity, we restrict our comparisons to diagonal and block diagonal preconditioning for large
114 scale sparse positive definite and general invertible linear systems. The positive definite case arises from
115 many diverse applications including: finite element analysis, sparse regression, Newton type optimization
116 algorithms, and also from e.g., in [12, 13]: (i) the so-called normal equations from interior point methods in
117 solving linear programs and (ii) the Hessians in minimizing logistic regression.

118 The κ - and ω -condition numbers for the positive definite case $M > 0$ are, respectively, given by ratios of

119 eigenvalues

$$\kappa(M) = \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)}; \quad \omega(M) = \frac{\text{tr}(M)/n}{\det(M)^{1/n}} = \frac{\sum_i \lambda_i(M)/n}{\prod_i (\lambda_i(M))^{1/n}}. \quad (1.1)$$

120 Preconditioning with D uses

$$Ax = b \iff D(A(D(D^{-1}x))) = Db,$$

121 and reduces a condition number of DAD . For our purposes we exploit the fact that the following functions
122 have the same values for nonsingular D :

$$\kappa(DAD) = \frac{\lambda_{\max}(DAD)}{\lambda_{\min}(DAD)} = \frac{\lambda_{\max}(ADD)}{\lambda_{\min}(ADD)}; \quad \omega(DAD) = \frac{\text{tr}(DAD)/n}{\det(DAD)^{1/n}} = \frac{\text{tr}(ADD)/n}{\det(ADD)^{1/n}}. \quad (1.2)$$

123 Our algorithms exploit the functions with argument affine in $\bar{D} = DD$ rather than quadratic in D .

124 In addition, we note that κ and ω can be considered as *worst-case* and *average-case* condition numbers,
125 respectively. In fact, it is known that the ratio of the arithmetic to geometric mean is approximately
126 the *coefficient of variation*, \mathbf{CV} ; namely, the ratio between the standard deviation and the mean of a
127 distribution; when this quantity is close to zero.¹ And, one of the big advantages of ω over κ is that finding
128 explicit formulae for optimal preconditioners by minimizing ω with special structure is often possible due to
129 the analyticity and simplicity of differentiation of trace and determinant.

130 1.1 Main Contributions

131 The main contributions of this paper consist of both theoretical and numerical analyses of the κ and ω -
132 condition numbers. First, we provide an affine based pseudoconvex reformulation of the κ -optimal diagonal
133 preconditioning problem that allows for an elegant characterization of its optimality conditions. There are
134 three advantages of this reformulation: all its stationary points are global minima, its subgradients are
135 inexpensive to compute, and its optimization variable is just an $n \times 1$ vector rather than an $n \times n$ matrix
136 as in *semidefinite programming, SDP* [24]. Our extensive numerical experiments show that subgradient
137 methods based on our reformulation can effectively converge to a κ -optimal diagonal preconditioner more
138 efficiently in time and accuracy than current SDP-based approaches in the literature [12, 24].

139 Second, we provide and exploit the ability to find explicit formulae for ω -optimal diagonal and block
140 diagonal preconditioners. We demonstrate that many popular classic preconditioners, such as the Jacobi
141 preconditioner and row/column normalization preconditioners are ω -optimal preconditioners. We also display
142 that matrix balancing schemes reduce ω at every iteration and converge to a stationary point of the two-sided
143 ω -optimal diagonal preconditioning problem. To the best of our knowledge, this is the first time that such
144 a comprehensive characterization of optimality conditions for both condition numbers is provided.

145 Finally, we conduct additional extensive numerical experiments that compare the efficiency of the κ -
146 and ω -condition numbers. Our results show that since ω -optimal preconditioners typically have closed-
147 form solutions, they are much cheaper to construct than κ -optimal preconditioners. Our results further
148 demonstrate that PCG iterations and LSQR iterations for solving linear systems are, on average, more
149 correlated with the ω -condition number than the κ -condition number.

150 We now continue with the organization of the paper and more details on the main contributions. In Sec-
151 tion 2, we present several formulations for minimizing κ along with the derivatives and optimality conditions.
152 The formulations are presented to take advantage of the affine approach and then avoid the positive homo-
153 geneity of the problem in order to improve stability of the problem. In particular, we include Theorem 2.7 that
154 presents a characterization of a κ -optimally diagonally preconditioned A that is based on the *largest/smallest*
155 eigenpairs. This leads to an efficient subgradient algorithm and we include convergence results. In Section 3,
156 we characterize ω -optimal preconditioners with special structures. A characterization for a ω -optimal diag-
157 onal preconditioner is simple as it corresponds to the classical Jacobi preconditioner for symmetric positive

¹The Young-Trent formula [33] establishes $GM/AM \approx 1 - \frac{1}{2}\mathbf{CV}^2$. Moreover, Taylor's approximation of $1/(1+x)$ around zero yields $AM/GM \approx 1 + \frac{1}{2}\mathbf{CV}^2$.

158 definite matrices and the row/column normalization preconditioner for general nonsingular matrices. We
 159 also show that the matrix balancing algorithm, Sinkhorn-Knopp, linearly converges to a stationary point of
 160 the two-sided ω -optimal diagonal preconditioning problem. Finally, we include results on extending from
 161 diagonal to block diagonal preconditioning in Section 3.3. We show that the right-sided ω -optimal block
 162 diagonal preconditioner for full column rank matrices A comes from taking Q-less QR decompositions of the
 163 blocks of A . We also give a characterization for the left-sided ω -optimal block diagonal preconditioner and
 164 show that when A is square that it is related to applying the right-sided optimal preconditioner to A^T .

165 Extensive numerical tests are then conducted in Section 4. We demonstrate that our subgradient methods
 166 are more scalable and efficient at computing κ -optimal diagonal preconditioners than existing SDP-based
 167 approaches in the literature [12,24]. Moreover, the resulting preconditioners yield more substantial improve-
 168 ments in PCG performance for solving linear systems compared to those produced by [12,24].

169 We further compare optimal diagonal preconditioning strategies based on the condition numbers κ and
 170 ω . In particular, we observe that the Sinkhorn-Knopp algorithm, which converges to a stationary point
 171 of the two-sided ω -optimal preconditioning problem, produces preconditioners that significantly outperform
 172 the two-sided κ -optimal method of [24] in terms of LSQR iteration counts for minimizing $\|b - Ax\|$. While
 173 Sinkhorn-Knopp dramatically reduces ω , the method of [24] achieves a stronger reduction in κ . This is
 174 notable as despite this difference, Sinkhorn-Knopp yields better LSQR performance on the majority of
 175 instances while also being computationally cheaper.

176 Finally, we demonstrate that applying ω -optimal diagonal scaling to positive definite linear systems that
 177 are already κ -optimally scaled leads to significant improvements in PCG performance. Overall, our results
 178 highlight two key advantages of ω -optimal preconditioners: they are inexpensive to compute and substantially
 179 reduce the iteration counts of iterative methods.

180 1.2 Background; Preliminaries

181 Given a linear system $Ax = b$ with $A \in \mathcal{M}^n$, $n \times n$ matrix, nonsingular, and $A \approx P_1 P_2$, preconditioning
 182 effectively solves the given system by solving the following system for y :

$$P_1^{-1} A P_2^{-1} y = P_1^{-1} b \text{ for } y, \quad x = P_2^{-1} y. \quad (1.3)$$

183 It is assumed that the solutions with P_1, P_2 are inexpensive.² Surveys are given in e.g., [1, 5, 25, 31] and the
 184 references therein. There is no matrix-matrix multiplication in preconditioning algorithms, such as PCG,
 185 as they do not form the products explicitly; and only matrix-vector multiplications/divisions are performed.
 186 For example, the well known *Jacobi preconditioner* uses the diagonal $d = \text{diag}(A)$, if it is not zero, and the
 187 diagonal matrices $P_2 = \text{Diag}(d)$, $P_1 = I$.

188 A discussion and references on the two condition numbers κ, ω is given recently in [17]. The matrix
 189 nearness problem $\min_P \|I - AP\|_F$ is often used to find preconditioners P . This also arises in the derivation
 190 of quasi-Newton methods. The ω -condition number is introduced in [7] as a measure for nearness to the
 191 identity using $\min_{P>0} \omega(MP)$ for finding optimal quasi-Newton updates; see also [8, 32, 34]. It is related
 192 to the measure $\text{tr} A - \log \det A$ used in the convergence proofs in [3, 4]. Then, Kaporin [18] used ω to
 193 derive new conjugate gradient convergence rate estimates and guarantees. More recently [2] presents further
 194 relationships and convergence analyses. As mentioned above, it is known that the ratio of AM to GM is
 195 approximately the *coefficient of variation*, CV , $AM/GM \approx 1 + \frac{1}{2} CV^2$ [33]. When the data is log-normal,
 196 then we get the exact relation $AM/GM = \sqrt{1 + CV^2}$. Thus by reducing CV , ω promotes clustering of
 197 eigenvalues/singular values which is essential in convergence in PCG type methods, e.g., [14, 22].

198 Note that for invertible $A \in \mathcal{M}^n$, [24] uses the former in $\sqrt{\kappa(A^T A)} = \kappa(A)$ as both involve largest to
 199 smallest singular values. In fact, as emphasized in [17], $\kappa(A) = \kappa(A^{-1})$ and it is the latter that is the
 200 operation for solving linear systems. However, these statements are not true for ω , as the numerator being
 201 the sum of singular values, the nuclear norm, is not equivalent to $\sqrt{\text{tr}(A^T A)}$ and ω is not inverse invariant.
 202 In this paper, for ω , we work with general invertible A and then use $\omega(A^T A)$, i.e., use the ratios of eigenvalues
 203 of $A^T A$.

²In some of our theoretical work for notational convenience, we use $D \leftarrow P^{-1}$, below.

204 In contrast to using κ , for preconditioners with special structure one can exploit the simple derivatives
 205 of tr, \det in ω and obtain explicit formulae for optimal preconditioners. Moreover, these often coincide
 206 with heuristics in the literature, see [17]. Thus no optimization algorithm is needed to obtain the optimal
 207 preconditioner. For example, if we restrict to a diagonal preconditioner $D = \text{diag}(d)$, then the classic Jacobi
 208 preconditioner is a multiple of the ω -optimal diagonal preconditioner. If we restrict to a diagonal structure
 209 along with a partial upper triangular part size k , then the diagonal part again corresponds to the Jacobi
 210 preconditioner, while the upper triangular part $D_{ij}, i \leq j \leq k$, yields exactly the formula for the Cholesky
 211 factorization. (See Proposition 3.9 for more details). Both these observations highlight the close connections
 212 ω -optimal preconditioners have with classical heuristic preconditioners.

213 1.2.1 Notation

214 We work in *real* Euclidean vector spaces. We let \mathbb{S}^n denote the space of symmetric matrices with the
 215 trace inner product and corresponding Frobenius norm; $\mathbb{S}_+^n, \mathbb{S}_{++}^n$ are the cones of positive semidefinite, and
 216 definite, symmetric matrices of order n , respectively; denoted $S \geq 0, S > 0$, respectively. We let \mathcal{M}^n denote
 217 the space of square matrices of order n , also equipped with the trace inner product and Frobenius norm. We
 218 only work with real matrices in this paper. Throughout, we let $A \in \mathcal{M}^n$ be a nonsingular matrix. Hence,
 219 $M = A^T A > 0$.

220 In the literature, for $A \in \mathcal{M}^n$, the *classical condition number*, $\kappa(A) = \|A\| \|A^{-1}\| = \frac{\max \sigma_i(A)}{\min \sigma_i(A)}$, i.e., the ratio
 221 of largest to smallest singular values. For $B \in \mathcal{M}^n$ with real eigenvalues we let

$$\lambda_i = \lambda_i(B) : \quad \lambda_{\max} = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{\min} = \lambda_n,$$

222 be the eigenvalues in nonincreasing order, and thus define: $\lambda_1(B) = \max_i \lambda_i(B)$; $\lambda_n(B) = \min_i \lambda_i(B)$. We
 223 let $I_n \in \mathcal{M}^n$ be the identity matrix and let $e_n \in \mathbb{R}^n$ be the vector of ones. We often use the following matrix
 224 for a basis for the orthogonal complement e_n^\perp ,

$$V = \frac{1}{\sqrt{2}} \begin{bmatrix} I_{n-1} \\ -e_{n-1}^T \end{bmatrix}, \quad \text{range}(V) = e_n^\perp. \quad (1.4)$$

225 We use $X \bullet Y$ to denote the *Hadamard product* (elementwise product) of two (compatible) matrices. For a
 226 vector $d \in \mathbb{R}^n$, $\text{Diag}(d) \in \mathbb{S}^n$ denotes the diagonal matrix formed using d . The adjoint linear transformation
 227 for $S \in \mathbb{S}^n$ is denoted $\text{Diag}^*(S) = \text{diag}(S) \in \mathbb{R}^n$. For two compatible functions f and g , we let $f \circ g$ denote
 228 the composite function.

229 The classical Fenchel subdifferential of a convex function h is defined as

$$\partial h(x) := \{v : \langle v, y - x \rangle \leq h(y) - h(x), \quad \forall y\}. \quad (1.5)$$

230 The Clarke subdifferential of a locally Lipschitz, not necessarily convex, function h is defined as

$$\partial_C h(x) = \text{conv} \{s : \exists x^k \rightarrow x, \nabla h(x^k) \text{ exists, and } \nabla h(x^k) \rightarrow s\}, \quad (1.6)$$

231 where conv denotes the convex hull of a set. It is well known that the Fenchel subdifferential and the Clarke
 232 subdifferential coincide for convex functions. Further notation is introduced below as needed.

233 We now recall the well-known facts about the largest and smallest eigenvalues of symmetric matrices and
 234 include a proof for completeness as it is useful further below.

235 **Lemma 1.1.** *The maximum and minimum eigenvalue functions on \mathbb{S}^n , $\lambda_{\max}, \lambda_{\min} : \mathbb{S}^n \rightarrow \mathbb{R}$, are convex,*
 236 *concave, respectively.*

237 *Proof.* Letting $B \in \mathbb{S}^n$ and using the Rayleigh quotient, we get $\lambda_{\max}(B) = \max\{x^T B x : \|x\| = 1\}$, which is a
 238 maximum of linear (so convex) functions in B , and therefore is convex. We get the concavity part similarly
 239 from $\lambda_{\min}(B) = \min\{x^T B x : \|x\| = 1\}$.

240

241

242 Now let

$$M \in \mathbb{S}_{++}^n, d \in \mathbb{R}_{++}^n, D = \text{Diag}(d).$$

243 We note the following useful relations that follow by the commutativity property for eigenvalues, the fact
244 that similarity transforms preserve the spectrum, and our definitions of κ, ω that use eigenvalues for products
245 of positive definite matrices:

246 2 κ -Optimal Diagonal Preconditioning

247 Preconditioning (or scaling) plays a fundamental role in numerical linear algebra and optimization; see, for
248 example, the surveys [1, 5, 31]). In this section, we study the problem of computing a **κ -optimal diagonal
249 preconditioner** for positive definite linear systems

$$Mx = b, \quad M = A^T A, \quad A \in \mathcal{M}^n \text{ nonsingular.}$$

250 Our objective is to find a diagonal matrix $D > 0$ that minimizes the classical condition number $\kappa(D^{1/2}MD^{1/2})$.³

251 Rather than viewing this as a nonlinear problem in $D^{1/2}$, we exploit the eigenvalue invariance in (1.2) to
252 reformulate to the simpler affine mapping MD and still minimize $\kappa(D^{1/2}MD^{1/2})$. The reformulation yields
253 a pseudoconvex minimization problem with several advantages:

- 254 • all stationary points are global minimizers,
- 255 • subgradients are easy to compute,
- 256 • the composite functions are used to avoid homogeneity and design very efficient and scalable algorithms.

257 Note that [12] also considers minimizing $\kappa(D^{1/2}MD^{1/2})$, although they aim to find a diagonal precondi-
258 tioner that is a convex combination of a small basis or list of diagonal preconditioners.

259 2.1 Optimality Conditions

260 Let $M \in \mathbb{S}_{++}^n$. We denote the quadratic type scaling as follows:

$$\mathcal{D}_q : \mathbb{R}_{++}^n \rightarrow \mathbb{S}^n, \quad \mathcal{D}_q(d) = \text{Diag}(d)^{1/2} M \text{Diag}(d)^{1/2}.$$

261 By abuse of notation, we let the argument determine the function,

$$\kappa(d) = \kappa(\mathcal{D}_q(d)) = (\kappa \circ \mathcal{D}_q)(d).$$

262 Similarly,

$$\lambda_{\max}(d) = \lambda_{\max}(\mathcal{D}_q(d)) = (\lambda_{\max} \circ \mathcal{D}_q)(d), \quad \lambda_{\min}(d) = \lambda_{\min}(\mathcal{D}_q(d)) = (\lambda_{\min} \circ \mathcal{D}_q)(d).$$

263 In the literature, e.g., [11, 12, 24], κ -optimal diagonal preconditioning refers to solving

$$d^* \in \text{argmin} \left\{ \kappa(\mathcal{D}_q(d)) = \frac{\lambda_{\max}(\mathcal{D}_q(d))}{\lambda_{\min}(\mathcal{D}_q(d))} : d \in \mathbb{R}_{++}^n \right\}. \quad (2.1)$$

264 Here we restrict to $d \in \mathbb{R}_{++}^n$ as we are taking square roots, see also Remark 2.4, below.

³Note that we are using $\kappa(D^{1/2}MD^{1/2})$ and not the inverse $D^{-1/2}$ or D^{-1} as is customary, see e.g., (1.3). This simplifies the derivatives in the optimization problems. Moreover, the equivalent expression for $\kappa(D^{1/2}MD^{1/2})$ follows from the invariance of the spectrum after commuting matrices.

265

We show below, that for our purposes, we can use the equivalent *linear in d* transformation

$$\mathcal{D} : \mathbb{R}_{++}^n \rightarrow \mathbb{S}^n, \quad \mathcal{D}(d) = M \text{Diag}(d). \quad (2.2)$$

266

267

268

With $D = \text{Diag}(d)$, we first note the relationships in the eigenpairs of $\mathcal{D}_q(d) = D^{1/2}MD^{1/2}$ and the nonsymmetric but similar $\mathcal{D}(d) = MD$, and $\mathcal{D}(d)^T = DM$; as well we have the convexity and concavity of the maximum and minimum eigenvalues of these nonsymmetric matrices MD .⁴

269

270

Lemma 2.1. *Let $M, D = \text{Diag}(d) \in \mathbb{S}_{++}^n$, and $(u_i, \lambda_i), i = 1, \dots, n$, be orthogonal eigenpairs of $\mathcal{D}_q(d)$, i.e., $\langle u_i, u_j \rangle = 0, \forall i \neq j$. Define*

$$x_i = D^{-1/2}u_i, \forall i, \quad X = [x_1 \quad \dots \quad x_n], \quad Y = X^{-T}, \quad \Lambda = \text{Diag}(\lambda).$$

271

272

273

274

Then X and Y are matrices of right and left eigenvectors of $\mathcal{D}(d)$, respectively, with corresponding matrix of eigenvalues Λ ; and Y is given by $Y = DX(X^TDX)^{-1}$.

Moreover, $\lambda_{\max}(d) = \lambda_1(d)$ (respectively, $\lambda_{\min}(d) = \lambda_n(d)$) is a convex (respectively, concave) function of $d \in \mathbb{R}_{++}^n$.

275

276

Proof. Let $U := [u_1 \dots u_n]$ be the corresponding matrix with *orthogonal*, not necessarily orthonormal, eigenvectors. We have $X = D^{-1/2}U$ and thus $U^TU = X^TDX$. This brings us to

$$U^{-T} = U(X^TDX)^{-1}.$$

277

278

Notice that X^TDX is a diagonal matrix with its diagonal entries given by $\|u_i\|^2$, which are strictly positive as eigenvectors are nonzero. Thus the inverse is well-defined. Now, observe that

$$\begin{aligned} Y &= X^{-T} = D^{1/2}U^{-T} \\ &= D^{1/2}U(X^TDX)^{-1} \\ &= DX(X^TDX)^{-1} \end{aligned}$$

279

as desired. Now,

$$\begin{aligned} D^{1/2}MD^{1/2}U = U\Lambda &\implies MDD^{-1/2}U = D^{-1/2}U\Lambda \\ &\implies MDX = X\Lambda \\ &\implies X^TDM = \Lambda X^T \\ &\implies DMX^{-T} = X^{-T}\Lambda \\ &\implies DMY = Y\Lambda. \end{aligned}$$

280

281

282

The second and last equation show that X and Y are right, and left, eigenvectors of MD , respectively. Moreover, we note that $M > 0$ has a positive definite square root and $\lambda_i(M \text{Diag}(d)) = \lambda_i(M^{1/2} \text{Diag}(d)M^{1/2})$. Therefore,

$$\max \lambda_i(\mathcal{D}(d)) = \max \lambda_i(M^{1/2} \text{Diag}(d)M^{1/2}) = \max_{\|x\|=1} x^T M^{1/2} \text{Diag}(d)M^{1/2} x.$$

283

284

As in the proof of Lemma 1.1, the latter is a maximum of linear functions in d and therefore is convex. The concavity result follows similarly. ■

285

286

287

Corollary 2.2. *Let $M \in \mathbb{S}_{++}^n, d \in \mathbb{R}_{++}^n$. Then*

$$\lambda_{\max}(d) = \lambda_{\max}(\mathcal{D}_q(d)) = \lambda_{\max}(\mathcal{D}(d)), \quad \lambda_{\min}(d) = \lambda_{\min}(\mathcal{D}_q(d)) = \lambda_{\min}(\mathcal{D}(d));$$

288

and

$$\kappa(d) = \frac{\lambda_{\max}(\mathcal{D}_q(d))}{\lambda_{\min}(\mathcal{D}_q(d))} = \frac{\lambda_{\max}(\mathcal{D}(d))}{\lambda_{\min}(\mathcal{D}(d))},$$

289

where recall $\mathcal{D}(d)$ is as in (2.2).

⁴These are a special case of so-called *K*-pd matrices in the literature, i.e., a product of two symmetric matrices where at least one is positive definite. The product of two positive definite matrices $P = AB$ arises in the optimality conditions in semidefinite programming. One of the difficulties is that the symmetrization of P is *not* necessarily positive definite, see e.g., [29].

290 *Proof.* The proof follows immediately from Lemma 2.1.

291

292

293 This allows us to consider the equivalent simplified view of finding a κ -optimal diagonal preconditioner,
 294 i.e., we need to solve the fractional pseudoconvex⁵ minimization problem

$$\bar{d} \in \operatorname{argmin} \left\{ \kappa(d) = \frac{\lambda_{\max}(\mathcal{D}(d))}{\lambda_{\min}(\mathcal{D}(d))} : d \in \mathbb{R}_{++}^n \right\}. \quad (2.3)$$

295 The following Lemma 2.3 shows that for $M > 0$, the matrix MD having all positive eigenvalues is
 296 equivalent to the positive definiteness of D . Combining this observation with Lemma 2.1, we obtain that
 297 $\kappa(d)$ is the ratio of a convex function and a positive concave function on its domain, and is therefore
 298 pseudoconvex.

299 **Lemma 2.3.** *Let $M \in \mathbb{S}_{++}^n, d \in \mathbb{R}^n, D = \operatorname{Diag}(d)$. Then*

$$\lambda_i(MD) > 0, \forall i \iff d \in \mathbb{R}_{++}^n.$$

300 *Proof.* We note that the eigenvalues of MD are the same as the eigenvalues of $M^{1/2}DM^{1/2}$. Sylvester's
 301 Lemma of inertia implies that the number of negative eigenvalues of $M^{1/2}DM^{1/2}$ is the same as that of D
 302 and so it is the same as the number of negative elements in d .

303

304

305 **Remark 2.4.** *Note that the two equivalent problems (2.1) and (2.3) are both essentially unconstrained, and*
 306 *the optima are attained and characterized as stationary points in \mathbb{R}_{++}^n . This is not obvious but follows since*
 307 *we have the ratios of convex and concave functions using $\lambda_{\max}, \lambda_{\min}$ and this is over the open cone constraint*
 308 *$d \in \mathbb{R}_{++}^n$. If we add the constraints $d^T e_n = n, d > 0$, or equivalently that $d = e_n + Vv > 0$, with $V \in \mathbb{R}^{n \times (n-1)}$*
 309 *a matrix whose columns contain a basis for e_n^\perp as done above, then we have a bounded problem in $v \in \mathbb{R}^{n-1}$*
 310 *and the spectral radius $\lambda_{\max}(MD) \leq \|M\| \|D\|$ is bounded above. Bounded below away from 0 follows from*
 311 *applying the greedy solution to the knapsack problem with constraints $\sum_i d_i = n, d \geq 0$. We get*

$$\lambda_{\max}(d) \geq \operatorname{tr}(MD)/n = \sum_i M_{ii}d_i/n \geq \min_i M_{ii} > 0.$$

312 *Therefore, with the denominator going to 0, we have κ going to ∞ as $d = e_n + Vv$ approaches the boundary*
 313 *of the simplex. That is, minimizing κ provides a self-barrier function for this boundary.*

314 *Moreover, $\alpha > 0 \implies \kappa(d) = \kappa(\alpha d)$, i.e., we have positive homogeneity of degree zero. Therefore,*
 315 *the optimal d is not unique. By pseudoconvexity, the optimal set is a convex set. This set can be large,*
 316 *see Proposition 2.8.*

317 2.1.1 Characterization of κ -Optimal Scaling

318 Using eigenvalue derivative formulas in Lemmas 2.5 and 2.6, we compute the gradient of $\kappa(d)$ in the smooth
 319 setting. When the largest and smallest eigenvalues of $\mathcal{D}(d)$ are simple, the gradient takes an explicit form
 320 in terms of the corresponding eigenvectors. The main result Theorem 2.7 provides a characterization for
 321 κ -optimality.

⁵A function $f : X \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be pseudoconvex if for all $x, y \in X, \nabla f(x) \cdot (y - x) \geq 0 \implies f(y) \geq f(x)$. A key property of pseudoconvex functions is that every stationary point is a global minimizer.

322 **Lemma 2.5** ([30, (3.1)]). Let $B : \mathbb{R} \rightarrow \mathcal{M}^n$ be differentiable with derivative $\dot{B} = \dot{B}(t)$, and, at $t = \bar{t}$, let
 323 $B(\bar{t})$ be diagonalizable with real eigenpairs (ignoring the argument \bar{t})

$$BX = X\Lambda, B^T Y = Y\Lambda, \quad \Lambda = \text{Diag}(\lambda).$$

324 And make the choice $Y = X^{-T}$. Let $1 \leq k \leq n$ and λ_k be a singleton eigenvalue with right and left eigenvector
 325 x_k, y_k , respectively, taken from the corresponding columns of X, Y , respectively. Then the derivative of
 326 eigenvalues is given by

$$\dot{\lambda}_k(\bar{t}) = y_k^T(\bar{t}) \dot{B}(\bar{t}) x_k(\bar{t}) = \text{tr} \dot{B} x_k y_k^T. \quad (2.4)$$

327 *Proof.* The proof is in [30, Pg 303].⁷

328 ■

329

330 **Lemma 2.6.** The Fréchet derivative of the linear transformation \mathcal{D} at d acting on Δd is (simply)

$$\mathcal{D}'(d)(\Delta d) = M \text{Diag}(\Delta d).$$

331 Now, let λ be a singleton eigenvalue of $\mathcal{D}(d)$ with a right eigenvector x . Then the gradient of the composite
 332 function at d is

$$\nabla \lambda(\mathcal{D}(d)) = \frac{\lambda}{x^T D x} x \bullet x.$$

333 *Proof.* The derivative of the linear transformation is clear.

334 Suppose as above that we have a linear transformation $\mathcal{D}(d)$ with singleton eigenvalue λ . Let $y =$
 335 $Dx/x^T Dx$ be the corresponding left eigenvector given in Lemma 2.1. Then the derivative of the composite
 336 function at d acting on Δd is

$$\begin{aligned} \langle \nabla(\lambda \circ \mathcal{D})(d), \Delta d \rangle &= \lambda'(\mathcal{D}'(d)(\Delta d)) = y^T (\mathcal{D}'(d)(\Delta d)) x \\ &= \frac{1}{x^T D x} x^T (DM \text{Diag}(\Delta d)) x \\ &= \frac{\lambda}{x^T D x} x^T \text{Diag}(\Delta d) x = \frac{\lambda}{x^T D x} \langle x \bullet x, \Delta d \rangle. \end{aligned}$$

337 ■

338

339 Theorem 2.7 shows that at optimal preconditioning, the extreme eigenvectors must “spread mass equally”
 340 across coordinates. This balance condition is what characterizes optimal diagonal preconditioners.

341 **Theorem 2.7** (Characterization of κ -optimal diagonal preconditioning). Consider $M \in \mathbb{S}_{++}^n, d \in \mathbb{R}_{++}^n$ be
 342 given. Let $D = \text{Diag}(d)$ and let (λ_1, x_1) (resp. (λ_n, x_n)) be a singleton maximum (resp. minimum) eigenpair
 343 of $\mathcal{D}(d)$. Then the gradient of the composite function $\kappa(d)$ is

$$\nabla \kappa(d) = \kappa(d) \left(\frac{1}{x_1^T D x_1} (x_1 \bullet x_1) - \frac{1}{x_n^T D x_n} (x_n \bullet x_n) \right),$$

344 where recall that $\kappa(d)$ is as in Corollary 2.2. Hence, M is κ -optimally diagonally preconditioned if, and only
 345 if, the orthonormal eigenvector pair satisfies

$$x_1 \bullet x_1 = x_n \bullet x_n. \quad (2.5)$$

⁶If Y is an invertible matrix of right eigenvectors not chosen using X^{-T} , then the normalization scaling needs to be added explicitly as $\dot{\lambda}_k(\bar{t}) = y_k^T(\bar{t}) \dot{B}(\bar{t}) x_k(\bar{t}) / (y_k(\bar{t})^T x(\bar{t}))$.

⁷The derivative of eigenvectors is included in the reference along with a useful normalization that allows for stability of the evaluations of the eigenvectors.

346 Equivalently, after a permutation of the elements to account for the sign,

$$x_1 = \begin{pmatrix} u \\ v \end{pmatrix}, x_n = \begin{pmatrix} u \\ -v \end{pmatrix}, \|u\| = \|v\|. \quad (2.6)$$

347 In the nonsmooth case, we can choose the normalized x_1 , respectively x_n , in the eigenspace of $\lambda_{\max}(\mathcal{D}(d))$,
348 respectively $\lambda_{\min}(\mathcal{D}(d))$.

349 *Proof.* From Lemma 2.6, we can find the gradient as follows:

$$\begin{aligned} \nabla \kappa(d) &= \frac{1}{\lambda_n^2} \left(\lambda_n \dot{\lambda}_1 - \lambda_1 \dot{\lambda}_n \right) = \frac{1}{\lambda_n^2} \left(\frac{\lambda_n \lambda_1}{x_1^T D x_1} x_1 \bullet x_1 - \frac{\lambda_1 \lambda_n}{x_n^T D x_n} x_n \bullet x_n \right) \\ &= \kappa(d) \left((x_1 \bullet x_1) / (x_1^T D x_1) - (x_n \bullet x_n) / (x_n^T D x_n) \right). \end{aligned}$$

350 The characterization for κ -optimally diagonally preconditioned matrix A follows from solving $\nabla \kappa(e_n) = 0$.

351 ■

352

353 Because $\kappa(d)$ is positively homogeneous of degree zero, the minimizer is not unique. Proposition 2.8
354 shows that this nonuniqueness is structural rather than pathological, i.e., when the eigenspace for the largest
355 and smallest eigenvalues is unique of dimension 2, the orthogonal complement is unique. When the extreme
356 eigenvectors do not uniquely determine the eigenspace of the orthogonal complement, entire continua of opti-
357 mal diagonal scalings exist. This observation motivates the reformulation in Section 2.3, where homogeneity
358 is removed.

359 **Proposition 2.8** (Nonuniqueness of κ -optimal diagonal preconditioning). *Let $M > 0$ be a κ -optimal diagonal*
360 *preconditioned matrix as given in Theorem 2.7 with $x_i, \lambda_i, i = 1, n$, being two, singleton, eigenpairs that satisfy*
361 *the optimality conditions in Theorem 2.7. Thus we have*

$$\lambda_1 > \lambda_2 \geq \dots \geq \lambda_{n-1} > \lambda_n > 0, \quad \lambda = (\lambda_i) \in \mathbb{R}_{++}^n, \quad \Lambda = \text{Diag}(\lambda),$$

362 and we let $Q = [x_1 \quad \bar{Q} \quad x_n]$ be an orthogonal matrix and $M = Q\Lambda Q^T$ from the spectral theorem. Let V be
363 the basis for e_n^\perp defined in (1.4). Suppose in addition that the lack of strict complementarity condition holds,
364 i.e., there exists v such that

$$0 \neq v \in \{u \in \mathbb{R}^{n-1} : 0 = Vu \bullet x_1 = Vu \bullet x_n\}. \quad (2.7)$$

365 Then with $\Delta D = \text{Diag}(Vv)$, there exists $\epsilon > 0$ such that

$$\kappa(M) = \kappa((I + t\Delta D)M(I + t\Delta D)), \quad \forall |t| \leq \epsilon.$$

366 *Proof.* The result follows from expanding

$$(I + \epsilon\Delta D)M(I + \epsilon\Delta D) = M + O(\epsilon)$$

367 and using the continuity of eigenvalues and the fact that the optimality conditions imply

$$\begin{aligned} Vv \bullet x_i = 0, i = 1, n &\iff \text{Diag}(Vv)x_i = 0, i = 1, n \\ &\iff \Delta D x_i = 0, i = 1, n. \end{aligned}$$

368 Specifically, let

$$M_2 = \bar{Q} \text{Diag}((\lambda_2, \dots, \lambda_{n-1})^T) \bar{Q}^T, \quad M_1 = M - M_2.$$

369 The above equation then implies $\Delta D M_1 = M_1 \Delta D = 0$. Moreover, let $D = (I + \epsilon\Delta D)$, then

$$DMD = M_1 + DM_2D = M + \epsilon\Delta D M_2 + \epsilon M_2 \Delta D + \epsilon^2 \Delta D M_2 \Delta D.$$

370 Since $\text{range}(\Delta D) \subseteq \text{null}(M_1) = \text{range}(M_2)$, the range of the perturbation of M above is restricted to the
371 eigenspace of M_2 , i.e., to the span($\{x_2, x_3, \dots, x_{n-1}\}$). This means that after the perturbation, with $\epsilon > 0$
372 sufficiently small, λ_1 and λ_n remain the largest and smallest eigenvalues of DMD respectively. As stated
373 above, we are using the continuity of eigenvalues and the orthogonality $M_1 M_2 = 0$ that arises using the
374 spectral theorem and $M_1, M_2 \geq 0$.

375

376

377 Note that if condition (2.7) holds, then we get a nonsingleton set in \mathbb{R}^{n-1} of solutions. Moreover, the
 378 structure of V and support of v implies that (2.7) restricts the support of the eigenspace $\text{span}(\{x_1, x_n\})$ and
 379 we can permute to get the support in the last entries.⁸

380 2.2 A Projected Subgradient Method for Minimizing $\kappa(d)$

381 In Theorem 2.7, we derived the gradient and optimality conditions for minimizing the pseudoconvex function
 382 $\kappa(d)$ in (2.3), over the open set $d > 0$. Convergence of subgradient methods for pseudoconvex minimization
 383 typically requires optimizing over a closed set. Hence, we consider the following optimization problem:

$$\bar{d} \in \operatorname{argmin} \left\{ \kappa(d) = \frac{\lambda_{\max}(\mathcal{D}(d))}{\lambda_{\min}(\mathcal{D}(d))} : d \in \Omega := \{d : d \geq \delta e_n\} \right\}, \quad 0 < \delta \ll 1. \quad (2.8)$$

384 We now exploit the pseudoconvex structure that guarantees all stationary points are global minimizers
 385 making subgradient methods natural. The search directions used in Algorithm 2.1 are shown to lie in
 386 the quasisubdifferential of κ , allowing us to directly invoke classical convergence results for quasiconvex
 387 minimization. We first treat the unconstrained formulation in Algorithm 2.1 and address homogeneity in
 388 Section 2.3.

Algorithm 2.1 A Subgradient Method for Minimizing $\kappa(d)$ over Ω

Inputs: symmetric positive definite matrix $M > 0$; sequence of positive stepsizes $\{t_k\} \rightarrow 0$ with $\sum_{k=1}^{\infty} t_k = \infty$, scalar $\delta \in (0, 1)$; tolerance tol ; rule for the stopping criterion, stopcrit .

set $k \leftarrow 0$;

set $\text{stopcrit} \leftarrow \infty$;

set $d_1 = e_n \in \mathbb{R}^n$;

1: **while** $\text{stopcrit} > \text{tol}$ **do** (main outer loop)

2: set $k \leftarrow k + 1$;

3: compute min eigenpair (λ_n^k, x_n^k) and max eigenpair (λ_1^k, x_1^k) of $M \operatorname{Diag}(d_k)$;

4: compute direction

$$s_k = \frac{\lambda_1^k}{\lambda_n^k} \left(\frac{1}{\langle x_1^k, d_k \bullet x_1^k \rangle} (x_1^k \bullet x_1^k) - \frac{1}{\langle x_n^k, d_k \bullet x_n^k \rangle} (x_n^k \bullet x_n^k) \right); \quad (2.9)$$

5: perform projected gradient step

$$d_{k+1} = \max \left\{ d_k - t_k \frac{s_k}{\|s_k\|}, \delta e_n \right\}; \quad (2.10)$$

6: update stopcrit ;

7: **end while**(main outer loop)

Output: $\hat{D} := \operatorname{Diag}(d_{k+1})$.

389 **Remark 2.9.** For Algorithm 2.1, we use a popular choice for stepsize sequence $\{t_k\}$ in subgradient methods,
 390 $t_k = 1/k$, where k is the iteration index [19]. We note that since Algorithm 2.1 is a subgradient method, in
 391 general it is not a descent method. Finally, the algorithm is general in that it does not specify the stopping
 392 rule stopcrit . There are several possible rules that the user can employ in practice for updating stopcrit in
 393 step 6. For example, at every iteration the user can take stopcrit as $\|s_k\|$ or $|\kappa(d_{k+1}) - \kappa(d_k)|$.

⁸Necessity of (2.7) for nonuniqueness is still an open problem.

394 We now briefly describe each of the steps of Algorithm 2.1. First, step 3 computes a maximal and
395 minimal eigenpair of $M \text{Diag}(d_k)$ to construct the search direction s_k . It should be noted that the user
396 never has to form the matrix $M \text{Diag}(d_k)$ to compute the eigenpairs. Instead, the user only has to construct
397 subroutines which compute $M \text{Diag}(d_k) * x$ and $(M \text{Diag}(d_k)) \backslash x$ efficiently, where $x \in \mathbb{R}^n$. Second, it will
398 be shown in Section 2.2.1 that the direction s_k computed in step 4 lies in the quasisubdifferential (also to
399 be defined in Section 2.2.1) of $\kappa(d_k)$. Finally, using this direction, step 5 performs the projected gradient
400 update $d_{k+1} = \Pi_\Omega(d_k - t_k \frac{s_k}{\|s_k\|})$, where $\Pi_\Omega(x) := \operatorname{argmin}_{y \in \Omega} \|x - y\|$ denotes the projection onto Ω .

401 2.2.1 Asymptotic Convergence Analysis

402 Our convergence analysis of Algorithm 2.1 relies mostly on [19] where efficient subgradient methods for
403 minimizing quasiconvex functions are presented. Under various assumptions in addition to quasiconvexity,
404 the author establishes asymptotic convergence of subgradient methods with decaying stepsizes. Since our
405 function $\kappa(d)$ is pseudoconvex, it is quasiconvex. For a more detailed discussion related to the assumptions
406 of Kiwiel [19] and how our set-up satisfies these assumptions, the reader is referred to Section A. In the
407 remaining part of this subsection, we show that Algorithm 2.1 asymptotically converges by showing that it
408 is an instance of the subgradient framework proposed by Kiwiel.

409 First, we need to define a special subdifferential called the quasisubdifferential for a quasiconvex func-
410 tion f . Define the strict sublevel set or inner slice of f as:

$$411 \quad \bar{S}(x) := \{y \in \operatorname{int} \bar{D} : f(y) < f(x)\}, \quad (2.11)$$

411 where $\operatorname{int} \bar{D}$ denotes the interior of the domain of f . The quasisubdifferential of a quasiconvex function f
412 relative to the above sublevel set is defined as

$$413 \quad \bar{\partial}^\circ f(x) := \{g : \langle g, y - x \rangle < 0, \quad \forall y \in \bar{S}(x)\}. \quad (2.12)$$

413 To minimize a quasiconvex function f over a closed convex set X , Kiwiel proposes in [19] the following basic
414 subgradient algorithm:

$$415 \quad x_{k+1} := \Pi_X(x_k - t_k \hat{g}_k), \quad \hat{g}_k := g_k / \|g_k\|, \quad g_k \in \bar{\partial}^\circ f(x_k), \quad k = 1, 2, \dots, \quad x_1 \in X,$$

415 where $t_k > 0$ are the stepsizes.

416 Hence, we need to characterize vectors that lie in the quasisubdifferential of $\kappa(d)$ to show that Algo-
417 rithm 2.1 is an instance of Kiwiel's subgradient framework. The result below presents one characterization
418 of vectors that lie in the quasisubdifferential of fractional programs like the one in our set-up (2.8).

419 **Lemma 2.10.** *Suppose $f(x) := a(x)/b(x)$ is well-defined for all $x \in X$ where $a(x)$ is a convex function
420 that is positive on X , $b(x)$ is a concave function that is positive on X . If for $x \in X$, $B := 1/b(x)$ and
421 $A := a(x)/b^2(x)$, then*

$$422 \quad B[\partial a(x)] + A[\partial(-b)(x)] \subseteq \bar{\partial}^\circ f(x).$$

422 *Proof.* Let $x \in X$ and consider B and A as in the assumptions of the lemma. It follows from the fact that a
423 and $-b$ are convex functions, B and A are positive scalars, and Fenchel subdifferential calculus rules that

$$424 \quad B[\partial a(x)] + A[\partial(-b)(x)] = \partial(Ba)(x) + \partial(-Ab)(x) \subseteq \partial(Ba - Ab)(x).$$

424 Now, let $g \in B[\partial a(x)] + A[\partial(-b)(x)] \subseteq \partial(Ba - Ab)(x)$, and suppose $y \in \bar{S}(x)$, i.e., y is in the interior of the
425 domain of f and $f(y) < f(x)$. It then follows that

$$\begin{aligned} 426 \quad \langle g, y - x \rangle &\leq Ba(y) - Ab(y) - Ba(x) + Ab(x) \\ &= \frac{a(y)}{b(x)} - \frac{a(x)b(y)}{b^2(x)} < 0, \end{aligned}$$

426 where the first inequality follows from the definition of Fenchel subdifferential, the equality follows from the
427 definitions of A and B , and the last inequality follows from the fact that $a(y)/b(y) < a(x)/b(x)$ and $b(y)$ and
428 $b(x)$ are positive. It then follows from the definition of quasisubdifferential in (2.12) that $g \in \bar{\partial}^\circ f(x)$.

429

430

431 Corollary 2.11 constructs a vector that lies in the quasisubdifferential of $\kappa(d)$.

432 **Corollary 2.11.** *Let $d \in \Omega$ and let (λ_1, x_1) and (λ_n, x_n) be a maximal and minimal eigenpair of $\mathcal{D}(d)$,*
 433 *respectively. Then, it holds that*

$$\frac{\lambda_1}{\lambda_n} \left(\frac{1}{x_1^T(d \bullet x_1)} (x_1 \bullet x_1) - \frac{1}{x_n^T(d \bullet x_n)} (x_n \bullet x_n) \right) \in \bar{\partial}^\circ(\kappa(d)).$$

434 *Proof.* It follows from Lemma 2.1 that $\lambda_{\max}(\mathcal{D}(d))$ and $\lambda_{\min}(\mathcal{D}(d))$ are convex and concave functions of d ,
 435 respectively. Also, it is easy to see that $\lambda_{\max}(\mathcal{D}(d))$ and $\lambda_{\min}(\mathcal{D}(d))$ are positive on Ω . Hence, it follows from
 436 Lemma 2.10 that

$$\frac{1}{\lambda_{\min}(\mathcal{D}(d))} \partial \lambda_{\max}(\mathcal{D}(d)) + \frac{\lambda_{\max}(\mathcal{D}(d))}{\lambda_{\min}^2(\mathcal{D}(d))} \partial(-\lambda_{\min}(\mathcal{D}(d))) \subseteq \bar{\partial}^\circ(\kappa(d)) \quad (2.13)$$

437 holds for $d \in \Omega$. Moreover, it follows from Lemma 2.6, the definition of Clarke subdifferential in (1.6), and the
 438 fact that the Fenchel and Clarke subdifferentials coincide for convex functions that the following inclusions
 439 hold.

$$\frac{\lambda_1}{x_1^T(d \bullet x_1)} (x_1 \bullet x_1) \in \partial \lambda_{\max}(\mathcal{D}(d)), \quad -\frac{\lambda_n}{x_n^T(d \bullet x_n)} (x_n \bullet x_n) \in \partial(-\lambda_{\min}(\mathcal{D}(d))), \quad (2.14)$$

440 where here (λ_1, x_1) and (λ_n, x_n) are maximal and minimal eigenpairs of $\mathcal{D}(d)$, respectively. The result then
 441 follows from (2.13) and (2.14).

442

443

444 **Remark 2.12.** *It follows from Corollary 2.11 that the update rule (2.10) in step 5 is of the form*

$$d_{k+1} = \Pi_\Omega \left(d_k - t_k \frac{s_k}{\|s_k\|} \right), \text{ where } s_k \in \bar{\partial}^\circ(\kappa(d_k)).$$

445 We are now ready to present the main theorem that shows Algorithm 2.1 converges asymptotically.

446 **Theorem 2.13** (Asymptotic convergence). *Let $\{d_k\}$ be generated by Algorithm 2.1. Let*

$$\kappa_* := \min\{\kappa(d) : d \in \Omega\}; \quad \kappa_*^k = \min\{\kappa(d_j), j = 1 \dots, k\}.$$

447 *Then*

$$\liminf_{k \rightarrow \infty} \kappa(d_k) = \kappa_*; \quad \text{and } \kappa_*^k \downarrow \kappa_*.$$

448 *Proof.* It follows from the last remark in Lemma 2.1 and the definitions of $\kappa(d)$ and Ω in (2.8), that
 449 $\lambda_{\max}(\mathcal{D}(d))$ (resp. $\lambda_{\min}(\mathcal{D}(d))$) is a convex (resp. concave) function that is positive on Ω . It then fol-
 450 lows from this observation, Lemma A.2, and the definition of $\kappa(d)$, that assumptions **A1-A4** in Section A
 451 hold for the minimization problem in (2.8). Clearly, also assumption **A5** in Section A also holds since Ω in
 452 (2.8) is a closed set and the intersection of Ω with the interior of the domain of $\kappa(d)$ is clearly nonempty. The
 453 result of the theorem then immediately follows from this observation, Remark 2.12, the facts that $t_k \rightarrow 0$,
 454 and $\sum_{k=1}^{\infty} t_k = \infty$, and Theorem 1 in [19].

455

456

2.3 Avoiding Positive Homogeneity in Minimizing $\kappa(d)$

As mentioned in Remark 2.4, $\kappa(d)$ is a positively homogenous function (of degree zero). Thus there are multiple optimal solutions and our problem is *Hadamard ill-posed*. From a computational stability perspective, it is more efficient to minimize an equivalent smaller dimensional formulation that is not positively homogeneous. With this in mind, we let $M > 0$ and $e_n \in \mathbb{R}^n$ be the vector of all ones. We consider the function

$$\mathcal{V}(v) := M \text{Diag}(e_n + Vv), \quad (2.15)$$

where $v \in \mathbb{R}^{n-1}$ and V is as in (1.4). In this subsection, we consider the following formulation

$$\min \left\{ \kappa(v) := \frac{\lambda_{\max}(\mathcal{V}(v))}{\lambda_{\min}(\mathcal{V}(v))} : e_n + Vv \geq \hat{\delta}e_n, v \in \mathbb{R}^{n-1} \right\}, \quad (2.16)$$

where $\hat{\delta} \in (0, 1)$ is a scalar. It is easy to see that with the choice of V in (1.4), the above (2.16) can be rewritten as

$$\min \left\{ \kappa(v) : v \in \hat{\Omega} \right\}, \quad (2.17)$$

where

$$\hat{\Omega} := \left\{ v \in \mathbb{R}^{n-1} : \sum_{i=1}^{n-1} v_i \leq -\sqrt{2}(\hat{\delta} - 1), \quad v_i \geq \sqrt{2}(\hat{\delta} - 1), \quad i = 1, \dots, n-1 \right\}. \quad (2.18)$$

Note that the set $\hat{\Omega}$ is convex and compact. Projecting onto it is also easy since it is just a simplex. Also, since $\hat{\Omega}$ is bounded, we will be able to get nonasymptotic convergence guarantees for the subgradient method that we propose to minimize (2.17).

The following lemma will be useful for developing our subgradient method. It gives useful characterizations of the derivatives of $\mathcal{V}(v)$ and $\kappa(v)$ in the smooth setting when the maximum and minimum eigenvalues of $\mathcal{V}(v)$ have multiplicity one.

Lemma 2.14 (Derivatives of $\mathcal{V}(v), \kappa(v)$). *Let $M > 0, v \in \mathbb{R}^{n-1}$ be given and set $w = e_n + Vv \in \mathbb{R}_{++}^n$ and $D = \text{Diag}(w)$. Also, let (λ_1, x_1) and (λ_n, x_n) be maximal and minimal eigenpairs of $\mathcal{V}(v)$, respectively, where λ_1 and λ_n have multiplicity one and x_1 and x_n are assumed to be normalized. Then the following hold:*

1 The derivative at v acting on $\Delta v \in \mathbb{R}^{n-1}$ is

$$\dot{\mathcal{V}}(v)(\Delta v) := \mathcal{V}'(v)(\Delta v) = M \text{Diag}(V\Delta v).$$

2 The gradient of the composite function $\kappa(v) := \kappa(\mathcal{V}(v))$ is

$$\nabla \kappa(v) = \kappa(v) V^T \left(\frac{1}{x_1^T(w \bullet x_1)} (x_1 \bullet x_1) - \frac{1}{x_n^T(w \bullet x_n)} (x_n \bullet x_n) \right). \quad (2.19)$$

Proof. 1 The proof follows from the function being affine.

2 For a singleton eigenvalue λ with normalized right eigenvector x , we have the left eigenvector $y = Dx/(x^T Dx)$ as in Lemma 2.1, which satisfies $x^T y = 1$. Then

$$\begin{aligned} \langle \nabla(\lambda \circ \mathcal{V})(v), \Delta v \rangle &= y^T \dot{\mathcal{V}}(v)(\Delta v) x = y^T (M \text{Diag}(V\Delta v)) x \\ &= \langle V^T \text{diag}(\lambda D^{-1} y x^T), \Delta v \rangle \\ &= \frac{1}{x^T D x} \langle V^T \text{diag}(\lambda D^{-1} D x x^T), \Delta v \rangle \\ &= \frac{1}{x^T(w \bullet x)} \langle V^T \text{diag}(\lambda x x^T), \Delta v \rangle = \frac{1}{x^T(w \bullet x)} \langle \lambda V^T(x \bullet x), \Delta v \rangle. \end{aligned}$$

Therefore the gradient of $\kappa(v) := \kappa(\mathcal{V}(v))$ is:

$$\begin{aligned} \nabla \kappa(v) &= \frac{\lambda_n \lambda_1 - \lambda_1 \lambda_n}{\lambda_n^2} = \frac{1}{\lambda_n^2} \left(\frac{\lambda_n}{x_1^T(w \bullet x_1)} \lambda_1 V^T(x_1 \bullet x_1) - \frac{\lambda_1}{x_n^T(w \bullet x_n)} \lambda_n V^T(x_n \bullet x_n) \right) \\ &= \kappa(v) V^T \left(\frac{1}{x_1^T(w \bullet x_1)} (x_1 \bullet x_1) - \frac{1}{x_n^T(w \bullet x_n)} (x_n \bullet x_n) \right). \end{aligned} \quad (2.20)$$

482

483

484 The following remark shows that, as in Lemma 2.14, $\nabla\kappa(v)$ lies in the quasisubdifferential of $\kappa(v)$. ■

485 **Remark 2.15.** Let $M > 0$ and $v \in \mathbb{R}^{n-1}$. Also, suppose that $w = e_n + Vv$ and let (x_i, λ_i) , $i = 1, n$, be
 486 eigenpairs of $\mathcal{V}(v)$, where $\|x_i\| = 1$. It then follows from Lemma 2.10 and a similar argument as in the proof
 487 of Corollary 2.11 that

$$\kappa(v)V^T \left(\frac{1}{x_1^T(w \bullet x_1)}(x_1 \bullet x_1) - \frac{1}{x_n^T(w \bullet x_n)}(x_n \bullet x_n) \right) \in \bar{\partial}^\circ(\kappa(v)) \quad (2.21)$$

488 where recall that $\kappa(v) := \kappa(\mathcal{V}(v))$.

489 We now present our subgradient algorithm for minimizing (2.17).

Algorithm 2.2 A Subgradient Method for Minimizing $\kappa(v) := \kappa(\mathcal{V}(v))$ over $\hat{\Omega}$

Inputs: symmetric positive definite matrix $M > 0$; $V \in \mathbb{R}^{n \times (n-1)}$ a basis matrix for the orthogonal complement e^\perp ; sequence of stepsizes $\{t_k\} = 1/\sqrt{k}$; scalar $\hat{\delta} \in (0, 1)$; a tolerance $\text{tol} > 0$; a rule for the stopping criterion, stopcrit.

set $k \leftarrow 0$;

set stopcrit $\leftarrow \infty$;

set $v_1 = 0 \in \mathbb{R}^{n-1}$ and $w_1 = e_n \in \mathbb{R}_+^n$;

1: **while** stopcrit $>$ tol **do** (main outer loop)

2: set $k \leftarrow k + 1$.

3: compute min eigenpair (λ_n^k, x_n^k) and max eigenpair (λ_1^k, x_1^k) of $M \text{Diag}(w_k)$;

4: compute direction

$$g_k = \frac{\lambda_1^k}{\lambda_n^k} V^T \left(\frac{1}{\langle x_1^k, w_k \bullet x_1^k \rangle} (x_1^k \bullet x_1^k) - \frac{1}{\langle x_n^k, w_k \bullet x_n^k \rangle} (x_n^k \bullet x_n^k) \right) \quad (2.22)$$

5: perform projected gradient step

$$v_{k+1} = \Pi_{\hat{\Omega}} \left(v_k - t_k \frac{g_k}{\|g_k\|} \right) \quad (2.23)$$

where $\hat{\Omega}$ is as in (2.18) and set

$$w_{k+1} = e + Vv_{k+1}; \quad (2.24)$$

6: update stopcrit;

7: **end while**(main outer loop)

Output: $\hat{D} := \text{Diag}(w_{k+1})$.

490 Several remarks about Algorithm 2.2 are now given. First, the matrix V does not need to be stored in
 491 memory and is actually not needed as input. All the user needs to input is a subroutine that outputs Vv
 492 given a vector $v \in \mathbb{R}^{n-1}$. Likewise, the matrix $M \text{Diag}(w_k)$ does not need to be stored. Second, it follows
 493 from Remark 2.15 that $g_k \in \bar{\partial}^\circ(\kappa(v_k))$, which is defined in (2.12). Finally, the projected gradient step in
 494 (2.23) can be performed very efficiently since projecting onto $\hat{\Omega}$ just involves computing the root of a simple
 495 equation.

496 **2.3.1 Nonasymptotic Convergence Rate Analysis for Minimizing $\kappa(v)$ on $\hat{\Omega}$**

497 In this subsection, we show a nonasymptotic convergence rate for Algorithm 2.2 for minimizing $\kappa(v)$ over $\hat{\Omega}$.
 498 More specifically, we show that

$$\min_{1 \leq k \leq K} \kappa(v_k) - \kappa_* \leq \epsilon$$

499 holds for $K = \mathcal{O}(1/\epsilon^2)$, where $\kappa_* = \min_{v \in \hat{\Omega}} \kappa(v)$. Our nonasymptotic convergence analysis of Algorithm 2.2
 500 relies mostly on the results from [19] and [15]. The function $\kappa(v)$ is pseudoconvex since it is the ratio of a
 501 convex function and a positive concave function. Also, observe that the set $\hat{\Omega}$ is just a box so it is a convex
 502 compact set that is easy to project onto. The only other assumption that needs to be verified to be able to
 503 show a nonasymptotic convergence rate of Algorithm 2.2 is that the function $\kappa(v)$ is Lipschitz continuous
 504 on $\hat{\Omega}$. This result is proved in the following proposition.

505 **Proposition 2.16.** *The function $\kappa(v)$ is Lipschitz continuous on $\hat{\Omega}$, i.e.,*

$$\|\kappa(v_1) - \kappa(v_2)\| \leq L\|v_1 - v_2\|, \quad \forall v_1, v_2 \in \hat{\Omega}$$

506 where $L > 0$ and $\hat{\Omega}$ is as in (2.18).

507 *Proof.* First, it is easy to see that $\lambda_{\max}(\mathcal{V}(v))$ (resp. $\lambda_{\min}(\mathcal{V}(v))$) is a convex (resp. concave) function. It then
 508 follows from this observation, the fact $\hat{\Omega}$ is a compact set that is contained in the domain of both functions,
 509 and Proposition A.48(b) in [23] that $\lambda_{\max}(\mathcal{V}(v))$ (resp. $\lambda_{\min}(\mathcal{V}(v))$) is L_1 -Lipschitz (resp. L_2 -Lipschitz) on
 510 $\hat{\Omega}$. Consider now the composite function $h(v) = g(\lambda_{\min}(\mathcal{V}(v)))$ where $g(y) = 1/y$. The above conclusion and
 511 the fact that $g(y) = 1/y$ is Lipschitz continuous on any positive open interval then imply that h is L_3 -Lipschitz
 512 continuous on $\hat{\Omega}$, where $L_3 > 0$.

513 We now show that $\kappa(v)$ is Lipschitz continuous. First, it holds that $|\lambda_{\max}(v)| \leq M_1$ and $|h(v)| \leq M_2$ for
 514 $v \in \hat{\Omega}$ since a Lipschitz function on a compact set is bounded. Then, for any $v_1 \in \hat{\Omega}$ and $v_2 \in \hat{\Omega}$, the following
 515 holds:

$$\begin{aligned} \|\kappa(\mathcal{V}(v_1)) - \kappa(\mathcal{V}(v_2))\| &= \|\lambda_{\max}(\mathcal{V}(v_1))h(v_1) - \lambda_{\max}(\mathcal{V}(v_2))h(v_2)\| \\ &= \|\lambda_{\max}(\mathcal{V}(v_1))h(v_1) - \lambda_{\max}(\mathcal{V}(v_1))h(v_2) + \lambda_{\max}(\mathcal{V}(v_1))h(v_2) - \lambda_{\max}(v_2)h(v_2)\| \\ &\leq M_1\|h(v_1) - h(v_2)\| + M_2\|\lambda_{\max}(\mathcal{V}(v_1)) - \lambda_{\max}(\mathcal{V}(v_2))\| \\ &\leq M_1L_3\|v_1 - v_2\| + M_2L_1\|v_1 - v_2\| = (M_1L_3 + M_2L_1)\|v_1 - v_2\|, \end{aligned}$$

516 which immediately implies the statement of the proposition with $L = M_1L_3 + M_2L_1$.

517 ■

518

519 We now state Theorem 2.17, which displays that Algorithm 2.2 is able to find an ϵ -approximate optimal
 520 solution of (2.17) with a sublinear $\mathcal{O}(1/\epsilon^2)$ rate of convergence.

521 **Theorem 2.17.** *Let $\epsilon > 0$ be a given tolerance and suppose that $K = \mathcal{O}(1/\epsilon^2)$. It then holds that*

$$\min_{1 \leq k \leq K} \kappa(v_k) - \kappa_* \leq \epsilon, \tag{2.25}$$

522 where $\kappa_* = \min_{v \in \hat{\Omega}} \kappa(v)$ and $\kappa(v) := \kappa(\mathcal{V}(v))$.

523 *Proof.* It is immediate to see that $\hat{\Omega}$ is a convex compact set and that $\kappa(v)$ is a quasiconvex function since it is
 524 the ratio of a convex and a positive concave function. It follows from Proposition 2.16 that $\kappa(v)$ is Lipschitz
 525 continuous on $\hat{\Omega}$. Finally, it follows from Remark 2.15 that the update (2.23) in step 5 of Algorithm 2.2 is
 526 of the form $v_{k+1} = \Pi_{\hat{\Omega}}\left(v_k - t_k \frac{g_k}{\|g_k\|}\right)$ where $g_k \in \bar{\partial}^\circ(\kappa(v_k))$. Hence, it follows from these observations, the
 527 fact that $\{t_k\} = 1/\sqrt{k}$, and Theorem 3.2(ii) in [15] with $s = 1/2$, $\delta = \epsilon$, and $p = 1$ that the statement of
 528 Theorem 2.17 holds.

531 3 ω -Optimal Structured Preconditioning

532 This section focuses on obtaining ω -optimal preconditioners of special structure for finding least squares
 533 solutions of $Ax = b$, where A might not be symmetric nor square. In what follows, A will be assumed to
 534 have full column rank, so that the matrix $M := A^T A$ is positive definite. Recall that for $M > 0$, we define
 535 $\omega(M)$ as

$$\omega(M) = \frac{\text{tr}(M)/n}{\det(M)^{1/n}},$$

536 i.e., it is the ratio of the arithmetic to geometric means of the eigenvalues of M . Unlike the classical condition
 537 number κ , the quantity ω is smooth, and closely related to eigenvalue clustering. These properties make it
 538 particularly well suited to deriving explicit solutions to optimality conditions for structured preconditioners.

539 This section is divided into three main subsections.

- 540 • Section 3.1 treats one-sided diagonal preconditioning, distinguishing between right- and left-sided scal-
 541 ings.
- 542 • Section 3.2 addresses two-sided diagonal scaling, establishing optimality conditions and linking them
 543 to matrix balancing.
- 544 • Section 3.3 extends the analysis to block-diagonal preconditioners, including triangular block structures.

545 A unifying theme of this section is that many classical preconditioners—Jacobi scaling, row and column
 546 normalization, and Sinkhorn–Knopp balancing—arise naturally as exact solutions or stationary points of
 547 ω -optimal problems. Table 3.1 summarizes several of these connections.

Table 3.1: Relation of ω -Optimal Preconditioners to the Literature

Special Structure	Preconditioner	Location	Literature
Right Diagonal	Column Normalization	Proposition 3.1	[7, 8]
Diagonal for Symmetric $M > 0$	Jacobi Preconditioner	Corollary 3.2	[16]
Left Diagonal	Row Normalization	Theorem 3.3	
Two-Sided Diagonal	Sinkhorn-Knopp/Matrix Balancing	Theorem 3.6	[20, 27, 28]
Right Block-Diagonal	Block QR	Corollary 3.10	[8, 17]
Left Block-Diagonal	Characterization	Theorem 3.11	

548 3.1 Right- and Left-Sided Diagonal

549 Suppose that $\text{diag}(A)$ has no zero entries. The classical Jacobi preconditioner [16], [26, Sect. 10.2] is defined
 550 by:

$$P^{-1}A = P^{-1}b, \quad P = \text{Diag}(\text{diag}(A)). \quad (3.1)$$

551 We note that this can be found by using: a splitting $A = M - N$ so that $A = \text{Diag}(\text{diag}(A)) - N$; or
 552 the following variational problem that implicitly finds the best approximation of the identity using diagonal
 553 matrices:

$$\min_d \|A - \text{Diag}(d)\|_F, \quad \text{Diag}(d)^{-1}A \approx I.$$

554 In what follows, we reinterpret such classical constructions through the lens of the ω -condition number.

555 **3.1.1 Right-Sided Diagonal**

556 The following result shows that a ω -optimal right-sided diagonal scaling of an overdetermined full rank
 557 matrix A is formed using (\pm) column normalization, see [7].

558 **Proposition 3.1.** *Let A be $m \times n$ full column rank, and $D := \text{Diag}(\text{diag}(A^T A))^{-1}$. Then $AD^{1/2}$ is a*
 559 *ω -optimal right-sided diagonal scaling i.e.,*

$$D \in \operatorname{argmin} \left\{ \omega(D^{1/2} A^T A D^{1/2}) : \text{diag}(D) \in \mathbb{R}_{++}^n \right\}.$$

560 *Proof.* The proof of the result can be found in [7, Proposition 2.1(v)] and a modified proof in [17, Proposition
 561 2.1(3)].

562 ■

563

564 As an immediate specialization, we obtain the following classical result.

565 **Corollary 3.2.** *Let $M \in \mathbb{S}_{++}^n$. Then the Jacobi preconditioner, $D = \text{Diag}(\text{diag}(M))^{-1}$, is the ω -optimal*
 566 *diagonal scaling, i.e.,*

$$D \in \operatorname{argmin} \left\{ \omega(D^{1/2} M D^{1/2}) : \text{diag}(D) \in \mathbb{R}_{++}^n \right\}.$$

567 *Proof.* The proof is immediate from Proposition 3.1 with $A^T A$ replaced by $M \in \mathbb{S}_{++}^n$.

568 ■

569

570 Thus, in contrast to its heuristic origins, Jacobi scaling is optimal in the precise sense of minimizing ω .

571 **3.1.2 Left-Sided Diagonal**

572 In this section we let $A \in \mathcal{M}^n$ be invertible, and we consider the optimal left-sided diagonal preconditioning
 573 problem:

$$\begin{aligned} \min \quad & \omega((\text{Diag}(c)^{1/2} A)^T \text{Diag}(c)^{1/2} A) \\ \text{s.t.} \quad & c \in \mathbb{R}_{++}^n. \end{aligned} \tag{3.2}$$

574 For simplicity, let $C := \text{Diag}(c)$. We can characterize the ω -optimal left-sided diagonal preconditioner using
 575 the above results in Section 3.1.1 for the right preconditioner.

576 **Theorem 3.3.** *Let $A \in \mathcal{M}^n$ be nonsingular. Then a ω -optimal left-sided diagonal preconditioner minimizing*
 577 *(3.2) is given by $C = \text{Diag}(\text{diag}(AA^T))^{-1}$.*

578 *Proof.* Note that

$$\omega((\text{Diag}(c)^{1/2} A)^T \text{Diag}(c)^{1/2} A) = \omega(A^T \text{Diag}(c) A) = \omega(\text{Diag}(c)^{1/2} A A^T \text{Diag}(c)^{1/2}).$$

579 Hence, Proposition 3.1 and the above equivalence implies that

$$\begin{aligned} C = \text{Diag}(\text{diag}(AA^T))^{-1} & \in \operatorname{argmin} \left\{ \omega(C^{1/2} A A^T C^{1/2}) : C = \text{diag}(c) \in \mathbb{S}_{++}^n \right\} \\ & = \operatorname{argmin} \left\{ \omega((\text{Diag}(c)^{1/2} A)^T \text{Diag}(c)^{1/2} A) : C = \text{diag}(c) \in \mathbb{S}_{++}^n \right\}. \end{aligned}$$

580 ■

581

3.2 Two-Sided Diagonal

We now consider the problem of finding the two-sided ω -optimal diagonal scaling for nonsingular matrices $A \in \mathcal{M}^n$. This subsection is broken up into two smaller parts. The first part derives the optimality conditions for the two-sided ω -optimal scaling problem. The second part presents an iterative matrix balancing scheme that reduces ω at every iteration and whose output is proved to be a stationary point of the two-sided ω -optimal scaling problem.

3.2.1 Optimality Conditions for Two-Sided Problem

This part presents the formulation of the two-sided ω -optimal diagonal scaling problem and derives its optimality conditions. Namely, we consider

$$\begin{aligned} \min \quad & \omega((\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})^T (\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})) \\ \text{s.t.} \quad & c, d \in \mathbb{R}_{++}^n, \end{aligned} \quad (3.3)$$

where $A \in \mathcal{M}^n$ is nonsingular. For simplicity, let $C := \text{Diag}(c)$ and $D := \text{Diag}(d)$. We also define

$$\mathcal{Q}(c, d) := A^T \text{Diag}(c) A \text{Diag}(d), \quad f(c, d) := \frac{1}{n} \text{tr}(\mathcal{Q}(c, d)), \quad g(c, d) := \det(\mathcal{Q}(c, d))^{1/n}. \quad (3.4)$$

Therefore we have the following equivalent problem to (3.3):

$$\begin{aligned} \min \quad & \omega_{\mathcal{Q}}(c, d) := \frac{f(c, d)}{g(c, d)} \\ \text{s.t.} \quad & c, d \in \mathbb{R}_{++}^n. \end{aligned} \quad (3.5)$$

Finally, for notational convenience, we define $\text{inv} : \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ by

$$\text{inv}(c, d) = \begin{pmatrix} \text{diag}(\text{Diag}(c)^{-1}) \\ \text{diag}(\text{Diag}(d)^{-1}) \end{pmatrix}. \quad (3.6)$$

The following result provides the gradient of $\omega_{\mathcal{Q}}(c, d)$.

Proposition 3.4. *The gradient of $\omega_{\mathcal{Q}}(c, d)$ is given by*

$$\nabla \omega_{\mathcal{Q}}(c, d) = \frac{1}{ng(c, d)} \begin{pmatrix} \text{diag}(A \text{Diag}(d) A^T) \\ \text{diag}(A^T \text{Diag}(c) A) \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{Q}}(c, d) \text{inv}(c, d), \quad (3.7)$$

where $g(c, d)$ is as in (3.4), $\omega_{\mathcal{Q}}(c, d)$ is defined in (3.5), and $\text{inv}(c, d)$ is as in (3.6).

Proof. To compute $\nabla \omega_{\mathcal{Q}}(c, d)$, we have to compute $\nabla f(c, d)$ and $\nabla g(c, d)$. It is easy to see from the definition of $f(c, d)$ in (3.4) that $\nabla f(c, d)$ can be computed as

$$\nabla f(c, d) = \frac{1}{n} \begin{pmatrix} \text{diag}(A \text{Diag}(d) A^T) \\ \text{diag}(A^T \text{Diag}(c) A) \end{pmatrix}. \quad (3.8)$$

To compute $\nabla g(c, d)$, observe that it follows from the definition of $g(c, d)$ in (3.4) and the formula for the gradient of the determinant that

$$\begin{aligned} \langle \nabla g(c, d), (\Delta c, \Delta d) \rangle &= \det(A^T A)^{1/n} \langle \nabla \det(CD)^{1/n}, (\Delta c, \Delta d) \rangle \\ &= \det(A^T A)^{1/n} \left\langle \frac{1}{n} \det(CD)^{\frac{1}{n}-1} \text{adj}(CD), (C\Delta D + \Delta CD) \right\rangle \\ &= \frac{1}{n} g(c, d) \left\langle (CD)^{-1}, (C\Delta D + \Delta CD) \right\rangle \\ &= \frac{1}{n} g(c, d) \text{tr}(C^{-1} \Delta C + D^{-1} \Delta D). \end{aligned}$$

601 Hence, it follows from the above equations

$$\nabla g(c, d) = \frac{1}{n} g(c, d) \begin{pmatrix} \text{diag}(\text{Diag}(c)^{-1}) \\ \text{diag}(\text{Diag}(d)^{-1}) \end{pmatrix} = \frac{1}{n} g(c, d) \text{inv}(c, d). \quad (3.9)$$

602 It then follows from (3.8), (3.9) and the quotient rule that $\nabla \omega_{\mathcal{Q}}(c, d)$ can be computed as

$$\begin{aligned} \nabla \omega_{\mathcal{Q}}(c, d) &= \frac{1}{g(c, d)^2} (g(c, d) \nabla f(c, d) - f(c, d) \nabla g(c, d)) = \frac{1}{g(c, d)} \nabla f(c, d) - \frac{\omega_{\mathcal{Q}}(c, d)}{g(c, d)} \nabla g(c, d) \\ &= \frac{1}{ng(c, d)} \begin{pmatrix} \text{diag}(A \text{Diag}(d) A^T) \\ \text{diag}(A^T \text{Diag}(c) A) \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{Q}}(c, d) \text{inv}(c, d). \end{aligned}$$

603

604

605 **Remark 3.5.** *It is easy to see that $\nabla \omega_{\mathcal{Q}}(c, d)$ can be equivalently written in the useful form:*

$$\nabla \omega_{\mathcal{Q}}(c, d) = \frac{1}{n} \begin{bmatrix} 0 & A \bullet A \\ (A \bullet A)^T & 0 \end{bmatrix} \begin{pmatrix} c \\ d \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{Q}}(c, d) \text{inv}(c, d).$$

606 We now state our main theorem of this subsection on a necessary condition for a vector to be globally
607 optimal for (3.3). This relates to optimal matrix balancing discussed below.

608 **Theorem 3.6.** *Let $A \in \mathbb{R}^{n \times n}$ be nonsingular. Let $c^*, d^* \in \mathbb{R}_{++}^n$ be a global optimal solution pair of (3.3).
609 Let $e_{2n} \in \mathbb{R}^{2n}$ be the vector of all ones, $C^* = \text{Diag}(c^*)$, and $D^* = \text{Diag}(d^*)$. Then:*

$$\begin{pmatrix} \text{diag}(C^* A D^* A^T) \\ \text{diag}(D^* A^T C^* A) \end{pmatrix} = \alpha e_{2n}, \quad \text{for some } \alpha > 0. \quad (3.10)$$

610 Hence, if c^*, d^* is a global optimal solution of (3.3), then all row and column of $\sqrt{C^*} A \sqrt{D^*}$ have identical
611 Euclidean norms, i.e., optimal two-sided ω -scaling produces a balanced matrix.

612 *Proof.* For $c^*, d^* \in \mathbb{R}_{++}^n$ to be a global optimal solution of (3.3), it must satisfy the necessary condition
613 $\nabla \omega_{\mathcal{Q}}(c^*, d^*) = 0$. Thus, it follows from (3.7) that c^*, d^* must satisfy

$$0 = \frac{1}{ng(c^*, d^*)} \begin{pmatrix} \text{diag}(A D^* A^T) \\ \text{diag}(A^T C^* A) \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{M}}(c^*, d^*) \text{inv}(c^*, d^*).$$

614 It follows from taking the Hadamard product with the vector $\begin{pmatrix} c^* \\ d^* \end{pmatrix}$ on both sides of the above equation as
615 well as multiplying both sides by $ng(c^*, d^*)$ that the following necessary condition must hold

$$f(c^*, d^*) e_{2n} = \begin{pmatrix} \text{diag}(C^* A D^* A^T) \\ \text{diag}(D^* A^T C^* A) \end{pmatrix}. \quad (3.11)$$

616 Observe that the definition of $f(c, d)$ in (3.4) implies that if c^*, d^* satisfies the above relation then any
617 positive scalar multiple of c^*, d^* also satisfies (3.11). It then follows from this observation and (3.11) that a
618 globally optimal solution c^*, d^* must satisfy the necessary condition (3.10) for some scalar $\alpha > 0$. The last
619 observation of the theorem then immediately follows from the condition in (3.10).

620

621

622 3.2.2 Square-Root Sinkhorn-Knopp Algorithm for Two Sided Problem

623 Let a nonsingular matrix $A \in \mathcal{M}^n$ be given. We have seen above that the ω -optimal right-sided preconditioner
624 is equivalent to the column normalization preconditioner. Similarly, the ω -optimal left-sided preconditioner
625 is equivalent to the row normalization preconditioner. Therefore, we can alternate between ω -optimal right-
626 and left-sided preconditioners and decrease the value of ω at each iteration. This yields the following com-
627 putationally efficient matrix balancing scheme where we alternate between column and row normalization.
628 This algorithm can be seen as equivalent to the Square-Root Sinkhorn-Knopp algorithm and is related to
629 other balancing methods in the literature [20, 21, 27, 28].

630 We prove below in Theorem 3.8 that the Square-Root Sinkhorn-Knopp algorithm, namely Algorithm 3.1,
631 converges and the outputted two-sided diagonal preconditioner satisfies the necessary condition in (3.10) for
632 global optimality of a two-sided ω -optimal preconditioner. This continues the theme that the ω -measure
633 provides justification for heuristics used in the literature.

Algorithm 3.1 Square-Root Sinkhorn-Knopp Algorithm for Two Sided ω -Optimal Preconditioner

Inputs: A square nonsingular matrix $A \in \mathcal{M}^n$, a tolerance $\text{tol} > 0$, and a rule for the stopping criterion,
stopcrit.

set stopcrit $\leftarrow \infty$;

set $A_0 \leftarrow A$;

set $k \leftarrow 0$;

1: **while** stopcrit $>$ tol **do** (main outer loop)

2: set $k \leftarrow k + 1$;

3: compute $(\sqrt{d_k})_i = 1./\|(A_{k-1})_{:,i}\|$, $i = 1, \dots, n$, where $(A_{k-1})_{:,i}$ is the i -th column of A_{k-1} ;

4: set $\tilde{A}_k \leftarrow A_{k-1} \text{Diag}(\sqrt{d_k})$;

5: compute $\sqrt{c_k} \in \mathbb{R}^n$ with $(\sqrt{c_k})_i = 1./\|(\tilde{A}_k)_{i,:}\|$, $i = 1, \dots, n$, where $(\tilde{A}_k)_{i,:}$ is the i -th row of \tilde{A}_k ;

6: update stopcrit;

7: set $A_k \leftarrow \text{Diag}(\sqrt{c_k})\tilde{A}_k$;

8: **end while**(main outer loop)

Output:

$$\hat{C}^{1/2} := \text{Diag} \left(\prod_{j=1}^k \sqrt{c_j} \right) \text{ and } \hat{D}^{1/2} := \text{Diag} \left(\prod_{j=1}^k \sqrt{d_j} \right).$$

634 There several viable options for the choice of stopcrit in Algorithm 3.1. One natural choice is

$$\text{stopcrit} = \max \left(\max_i \|(A_k)_{i,:} - 1\|, \max_j \|(A_k)_{:,j} - 1\| \right).$$

635 It is shown in Theorem 3.8 that the sequence A_k converges to a matrix that is balanced, i.e., one that has
636 column and row norms all equal to one. The remark below also discusses that the ω values of the iterates
637 are monotonically decreasing.

638 **Remark 3.7.** At each iteration, \tilde{A}_k (resp. A_k) is computed using the ω -optimal right-sided (resp. left-sided)
639 preconditioner. Hence, the ω values of the iterates are monotonically decreasing, i.e., $\omega(A_k A_k^T) \leq \omega(\tilde{A}_k \tilde{A}_k^T) \leq$
640 $\omega(A_{k-1} A_{k-1}^T)$ for all $k \geq 1$. The output of Algorithm 3.1 thus satisfies $\omega((\hat{C}^{1/2} \hat{A} \hat{D}^{1/2})^T (\hat{C}^{1/2} \hat{A} \hat{D}^{1/2})) \leq$
641 $\omega(A^T A)$.

642 Theorem 3.8 now provides a convergence guarantee for Algorithm 3.1 and shows that the necessary
643 condition (3.10) for a two-sided ω -optimal diagonal preconditioner holds in the limit for Algorithm 3.1.

644 **Theorem 3.8.** Let A be such that $A \bullet A$ has total support.⁹ Then the sequence A_k converges linearly to a
645 matrix $\bar{A} := \bar{C}^{1/2} A \bar{D}^{1/2}$ that has column and row norms all equal to 1. That is,

$$\lim_{k \rightarrow \infty} \text{Diag} \left(\prod_{j=1}^k \sqrt{c_j} \right) A \text{Diag} \left(\prod_{j=1}^k \sqrt{d_j} \right) = \bar{C}^{1/2} A \bar{D}^{1/2}$$

646 with a linear rate of convergence. Hence, it follows from Theorem 3.6 that limiting matrices \bar{C} and \bar{D} satisfy
647 the necessary condition (3.10) for a two-sided ω -optimal diagonal preconditioner.

648 *Proof.* Let $\tilde{P}_k = \tilde{A}_k \bullet \tilde{A}_k$ and let $P_k = A_k \bullet A_k$. It follows that $\tilde{P}_k = (A_{k-1} \bullet A_{k-1}) \text{Diag}(d_k) = P_{k-1} \text{Diag}(d_k)$
649 where $(d_k)_i = 1./\|(A_{k-1})_{:,i}\|^2$. Note $\|(A_{k-1})_{:,i}\|^2$ is exactly the i -th column sum of P_{k-1} so $(d_k)_i =$
650 $1./\sum_{l=1}^n (P_{k-1})_{l,i}$. Moreover, also observe that $P_k = \text{Diag}(c_k)(\tilde{A}_k \bullet \tilde{A}_k) = \text{Diag}(c_k)\tilde{P}_k$. Note $(c_k)_i =$
651 $1./\sum_{l=1}^n (\tilde{P}_k)_{i,l}$. Hence, the sequences of iterates \tilde{P}_k and P_k can be viewed as the iterates of Sinkhorn-
652 Knopp Algorithm applied to $A \bullet A$. By [20], we know then that the sequence P_k converges linearly to a
653 matrix \bar{P} where $\bar{P} := \bar{C}(A \bullet A)\bar{D}$ is a doubly stochastic matrix. The sequences of iterates \tilde{A}_k and A_k are
654 just the square roots of the iterates of the Sinkhorn-Knopp algorithm applied to $A \bullet A$. Therefore, it follows
655 from the above argument that the sequence A_k converges linearly to a matrix $\bar{A} := \bar{C}^{1/2} A \bar{D}^{1/2}$ that has
656 column and row norms all equal to 1. The last conclusion of the theorem follows from this observation and
657 the necessary condition (3.10) in Theorem 3.6. ■

658
659

660 3.3 Right-, Left-Sided Block Diagonal

661 We now show that the above results extend directly to finding *block diagonal* preconditioners for least
662 squares solutions of full rank possibly overdetermined linear systems $Ax = b, A \in \mathbb{R}^{m \times n}, m \geq n$. The
663 preconditioner is extended from diagonal to block diagonal (block upper-triangular). This relates to sparse
664 QR preconditioning and extends the results for ω -optimal preconditioning with partial Cholesky structure
665 in [17, Theorem 2.7] and the block diagonal ω -optimal preconditioner in [8, Prop. 3 part 3].

666 3.3.1 Right-Sided Block Diagonal

667 We let the linear transformation $B = \text{Blkdiag}(\mathcal{B}) \in \mathcal{M}^n$ denote the block diagonal matrix with blocks formed
668 from the set of square matrices $\mathcal{B} = \{B_i\}_{i=1}^k$ of order $n_i, \sum_{i=1}^k n_i = n$. For our application we restrict the
669 blocks to be of upper-triangular structure denoted $\mathcal{R} = \{R_i\}_{i=1}^k$, and see that the best, with respect to the
670 ω measure, comes from Q-less QR decompositions of the corresponding blocks of A .

671 First we choose the number of blocks k and the block sizes $n_i, \sum_{i=1}^k n_i = n$. Thus we get the block
672 structure

$$A = [A_1 \ A_2 \ \dots \ A_k], \quad A/\text{Blkdiag}(\mathcal{R}) = [A_1/R_1 \ A_2/R_2 \ \dots \ A_k/R_k],$$

673 where we use the MATLAB notation that $/$ denotes matrix division $A/R = AR^{-1}$. Optimally, we want the
674 sizes of the blocks chosen so that the matrices $A_i^T A_i$ are chordal so that there is no loss of sparsity. Moreover,
675 we can allow a preliminary permutation of the columns $A \leftarrow AP$ so that we can increase the effect of the
676 preconditioner.

677 To extend the diagonal preconditioner results we solve:

$$\min_{\mathcal{B}} \{\omega((AB)^T(AB)) : B = \text{Blkdiag}(\mathcal{B})\}.$$

678 The basic result we use follows.

⁹A nonnegative square matrix B is said to have total support if $B \neq 0$ and all its nonzero elements lie on a positive diagonal. In this case, a diagonal of a matrix is just a collection of elements with one element from each row and one element from each column of B . More specifically, it will consist of $b_{i,\sigma(i)}$ for $i = 1, 2, \dots, n$ for a permutation σ and where b_{ij} denotes the (i, j) entry of B .

679 **Proposition 3.9** ([8, Prop. 3 part 3], [9, Prop. 2.2]). *Let*

$$A = [A_1 \ A_2 \ \dots \ A_k], \quad A_i \in \mathbb{R}^{m \times n_i}, \forall i,$$

680 *be a full rank $m \times n$ matrix, $m \geq n$. Then an optimal block diagonal scaling*

$$\mathcal{B} := \{B_1, B_2, \dots, B_k\}, \quad B_i \in \mathbb{R}^{n_i \times n_i}, \quad B := \text{Blkdiag}(\mathcal{B}),$$

681 *that minimizes the measure ω , i.e.,*

$$\min \{\omega((AB)^T(AB)) : B = \text{Blkdiag}(\mathcal{B})\},$$

682 *is found by satisfying the factorization*

$$B_i B_i^T = (A_i^T A_i)^{-1}, \quad i = 1, \dots, k.$$

683 **Corollary 3.10.** *Let A, \mathcal{B}, B be as in Proposition 3.9.*

684 *(i) If we restrict the diagonal blocks B_i to be diagonal themselves, then we get the optimal diagonal preconditioner in Proposition 3.1.*

686 *(ii) If we restrict the diagonal blocks B_i to be upper-triangular, then we get $B_i = R_i^{-1}, A_i = Q_i R_i$ from the QR decomposition of A_i (up to sign of rows of R).*

688 3.3.2 Left-Sided Block Diagonal

689 We continue with A full column rank but with left-sided preconditioning. We present a nonlinear equation that characterizes optimality.

691 **Theorem 3.11.** *Let $A \in \mathbb{R}^{m \times n}$ be full column rank with block structure given by*

$$A = \begin{bmatrix} \bar{A}_1 \\ \bar{A}_2 \\ \dots \\ \bar{A}_\ell \end{bmatrix}, \quad \bar{A}_j \in \mathbb{R}^{m_j \times n}, \text{rank}(\bar{A}_j) = m_j, \forall j.$$

692 *Then the ω -optimal left-sided block diagonal preconditioner*

$$E := \text{Blkdiag}(E_1, E_2, \dots, E_\ell), \quad E_j \in \mathbb{R}^{m_j \times m_j},$$

693 *that minimizes the measure ω , i.e.,*

$$\min \omega((EA)^T EA),$$

694 *is characterized by:*

$$\bar{A}_j \bar{A}_j^T = \frac{1}{n} \left(\text{tr} \sum_{t=1}^{\ell} (\bar{A}_t^T (E_t^T E_t) \bar{A}_t) \right) \bar{A}_j (A^T (E^T E) A)^{-1} \bar{A}_j^T, \quad \forall j = 1, \dots, \ell.$$

695 *Proof.* To begin, we substitute $K := E^T E$. That is, we solve

$$\min \{\omega(A^T K A) : K = \text{Blkdiag}(K_1, K_2, \dots, K_\ell), \quad K_j \in \mathbb{R}^{m_j \times m_j}\}.$$

696 Notice that $A^T K A = \sum_{j=1}^{\ell} \bar{A}_j^T K_j \bar{A}_j$. Define the linear transformation $\mathcal{G}(K) = \sum_{j=1}^{\ell} \bar{A}_j^T K_j \bar{A}_j$. This yields

$$\omega_{\mathcal{G}}(K) := \frac{\text{tr}(A^T K A)/n}{\det(A^T K A)^{1/n}} = \frac{\sum_{j=1}^{\ell} \text{tr}(\bar{A}_j^T K_j \bar{A}_j)/n}{\det(A^T K A)^{1/n}}.$$

697 Hence, for convenience and for the remainder of the proof define

$$f(K) := \frac{1}{n} \sum_{j=1}^{\ell} \text{tr}(\bar{A}_j^T K_j \bar{A}_j) \quad \text{and} \quad g(K) := \det(A^T K A)^{1/n}.$$

698 The partial derivative of g with respect to the block K_j is given by

$$\begin{aligned} \langle \nabla_{K_j} g(K), \Delta K_j \rangle &= \langle \nabla \det(\mathcal{G}(K))^{1/n}, \Delta K_j \rangle \\ &= \left\langle \frac{1}{n} \det(\mathcal{G}(K))^{\frac{1}{n}-1} \text{adj}(\mathcal{G}(K)), \bar{A}_j^T \Delta K_j \bar{A}_j \right\rangle \\ &= \frac{1}{n} g(K) \langle \mathcal{G}(K)^{-1}, \bar{A}_j^T \Delta K_j \bar{A}_j \rangle \\ &= \frac{1}{n} g(K) \langle \bar{A}_j \mathcal{G}(K)^{-1} \bar{A}_j^T, \Delta K_j \rangle. \end{aligned}$$

699 Therefore, we have

$$\nabla_{K_j} f(K) = \frac{1}{n} \bar{A}_j \bar{A}_j^T \quad \text{and} \quad \nabla_{K_j} g(K) = \frac{1}{n} g(K) \bar{A}_j (A^T K A)^{-1} \bar{A}_j.$$

700 These in turn lead to the partial derivative of $\omega_{\mathcal{G}}$ with respect to K_j given by

$$\begin{aligned} \nabla_{K_j} \omega_{\mathcal{G}}(K) &= \frac{1}{g(K)^2} (g(K) \nabla_{K_j} f(K) - f(K) \nabla_{K_j} g(K)) \\ &= \frac{1}{n g(K)} (\bar{A}_j \bar{A}_j^T - f(K) \bar{A}_j (A^T K A)^{-1} \bar{A}_j^T) \end{aligned}$$

701 Therefore

$$\nabla \omega_{\mathcal{G}}(K) = 0 \quad \iff \quad \bar{A}_j \bar{A}_j^T = \frac{1}{n} \sum_{t=1}^{\ell} \text{tr}(\bar{A}_t \bar{A}_t^T K_t) \bar{A}_j (A^T K A)^{-1} \bar{A}_j^T, \quad \forall j.$$

702 We can then recover the E_j from K using factorizations of the blocks $K_j = E_j^T E_j$.

703

704

705 Finding an explicit solution for E in Theorem 3.11 is an open question. However, if we assume further
706 that A is a square matrix, the following Corollary 3.12 connects the result from Corollary 3.10 to this left
707 sided preconditioning.

708 **Corollary 3.12.** *Let $A \in \mathcal{M}^n$ be nonsingular. Let the block of rows, $\bar{A}_j \in \mathbb{R}^{n_j \times n}$, be given by $A^T =$
709 $[\bar{A}_1^T \bar{A}_2^T \dots \bar{A}_\ell^T]$, $\sum_j n_j = n$. Define $\mathcal{E} := \{E_1, E_2, \dots, E_k\}$ by $E_i := \bar{R}_i^{-T}$, where $\bar{A}_i^T := \bar{Q}_i \bar{R}_i$ is the QR
710 decomposition of \bar{A}_i^T . Then, \mathcal{E} is the optimal lower triangular block diagonal scaling for*

$$\min \{ \omega((EA)^T(EA)) : E = \text{Blkdiag}(E_1, E_2, \dots, E_k) \}.$$

711 *Proof.* Since A and E are square matrices, we observe

$$\omega((EA)^T(EA)) = \omega(A^T E^T E A) = \omega(E A A^T E^T),$$

712 and thus Corollary 3.10 implies $E_i = \bar{R}_i^{-T}$, where $\bar{A}_i \bar{A}_i^T := \bar{R}_i^T \bar{R}_i$ as defined above.

713

714

715 **Remark 3.13.** *The result in Corollary 3.12 agrees with our characterization in Theorem 3.11. Indeed, by*
 716 *defining $E_i := \bar{R}_i^{-T}$ with $\bar{A}_i \bar{A}_i^T = \bar{R}_i^T \bar{R}_i$, we get that $\text{tr}(\bar{A}_i \bar{A}_i^T E_i^T E_i) = \text{tr}(\bar{R}_i^T \bar{R}_i \bar{R}_i^{-1} \bar{R}_i^{-T}) = \text{tr}(I_{n_i}) = n_i$.*
 717 *Thus,*

$$\sum_{t=1}^{\ell} \text{tr}(\bar{A}_t \bar{A}_t^T E_t^T E_t) = \sum_{t=1}^{\ell} n_t = n.$$

718 *Now, define $\mathcal{I}_j := \begin{bmatrix} 0 & I_{n_j} & 0 \end{bmatrix} \in \mathbb{R}^{n_j \times n}$, with the identity $I_{n_j} \in \mathcal{M}^{n_j}$ of order n_j , appropriately located such*
 719 *that $\mathcal{I}_j A = \bar{A}_j$. Then,*

$$\begin{aligned} \frac{1}{n} \left(\text{tr} \sum_{t=1}^{\ell} (\bar{A}_t^T E_t^T E_t \bar{A}_t) \right) \bar{A}_j (A^T E^T E A)^{-1} \bar{A}_j^T &= \bar{A}_j (A^T E^T E A)^{-1} \bar{A}_j^T \\ &= \mathcal{I}_j A A^{-1} E^{-1} E^{-T} A^{-T} A^T \mathcal{I}_j^T \\ &= \mathcal{I}_j E^{-1} E^{-T} \mathcal{I}_j^T \\ &= E_j^{-1} E_j^{-T} = \bar{R}_j^T \bar{R}_j = \bar{A}_j \bar{A}_j^T. \end{aligned}$$

720 4 Computational Experiments

721 Our computational experiments are organized into four parts, each designed to evaluate different aspects of
 722 the proposed methods.

723 In Section 4.1, we compare the computational efficiency of our subgradient-based algorithms with existing
 724 SDP-based approaches for computing approximate κ -optimal diagonal preconditioners.

725 In Section 4.2, we assess the impact of these preconditioners on PCG performance for solving symmetric
 726 positive definite systems.

727 In Section 4.3, we study two-sided preconditioning for LSQR and compare ω -optimal scaling with κ -optimal
 728 approaches.

729 Finally, in Section 4.4, we investigate whether applying ω -optimal scaling to κ -optimally preconditioned
 730 systems can further improve PCG performance.

731 All experiments in Sections 4.1 to 4.4 are run on a 2023 Macbook Pro with an 8-core CPU and 128
 732 GB of memory using MATLAB 2025b. The latest codes are available at this [clickable-link](https://github.com/asujanani6/Optimal_Preconditioning) or with URL
 733 https://github.com/asujanani6/Optimal_Preconditioning.

734 4.1 Minimizing κ Efficiently

735 In this subsection, we study the practical performance of our subgradient methods, Algorithms B.1 and 2.2,
 736 for computing approximate κ -optimal diagonal preconditioners. We first compare the proposed methods with
 737 the SDP-based approaches of [12, 24] on moderate-scale instances. We then investigate the scalability of the
 738 proposed methods on substantially larger problems for which SDP-based approaches become impractical.

739 Throughout this subsection, condition numbers are evaluated using MATLAB's EIGS routine with tol-
 740 erance 10^{-10} . CPU times include all preprocessing and algorithmic overhead. Unless otherwise stated, the
 741 default parameter settings from [12, 24] are used. Algorithm 2.2 uses tolerance 10^{-4} , maximum iteration limit
 742 500, sets $\text{stopcrit} = 2|\kappa(d_{k+1}) - \kappa(d_k)| / |\kappa(d_{k+1}) + \kappa(d_k)|$ at every iteration, and uses $\hat{\delta} = 10^{-3}$. Algorithm B.1
 743 sets a maximum of 80 iterations and uses the same tolerance 10^{-4} .

744 4.1.1 Moderate Sizes

745 In this subsection, we compare our Algorithms B.1 and 2.2, with the SDP-based approaches of [12, 24] on
 746 moderate-scale SuiteSparse and randomly generated instances. The precise details of our line-search variant
 747 Algorithm B.1 are given in Section B.

748 Our methods yield an optimal diagonal matrix D that minimizes $\kappa(D^{1/2} M D^{1/2})$ for given $M > 0$.
 749 Convergence is proved for our Algorithm 2.2 that often achieved stronger reductions in κ than Algorithm B.1;

our line-search variant Algorithm B.1 was often slightly faster in practice than Algorithm 2.2. In contrast, the method in [12] solves an SDP over a prescribed subspace of diagonal preconditioners, while [24] computes a right-scaling E for M by minimizing $\kappa((BE)^T(BE))$, with $M = B^T B$ for given nonsingular B . Hence, when implementing this latter method, we take B as a Cholesky factorization of M .

Table 4.1 reports results on eighteen matrices from the SuiteSparse Matrix Collection [6], while Table 4.2 considers randomly generated instances produced using MATLAB's SPRANDSYM routine.

The proposed methods consistently achieved substantially larger reductions in κ than the SDP-based approaches. On the SuiteSparse instances, Algorithm 2.2 and Algorithm B.1 reduced κ on average by 79.7% and 68.5%, respectively, compared to 30.3% and 21.1% for [24] and [12]. Moreover, both proposed methods were substantially faster, with Algorithm B.1 often providing the best runtime performance.

For the random instances, a time limit of 3600 seconds was imposed on [12] due to its substantially higher computational cost. Across all tested instances, Algorithm 2.2 and Algorithm B.1 consistently reduced κ more effectively and more efficiently than the competing methods.

Table 4.1: Comparison of Algorithms for Min κ on SuiteSparse

Dim & Density, & $\kappa(M)$			% Reduction in κ				CPU Time (seconds)				CPU Ratios	
n	density	$\kappa(M)$	[24]	[12]	Algorithm 2.2	Algorithm B.1	[24]	[12]	Algorithm 2.2	Algorithm B.1	[24]/Algorithm 2.2	[24]/Algorithm B.1
1074	1.1e-02	2.6e+07	0.0e+00	0.0e+00	9.9e+01	9.9e+01	48.203	51.700	8.469	5.644	5.7	8.5
2003	2.1e-02	1.1e+10	0.0e+00	0.0e+00	9.8e+01	6.8e+01	330.887	324.334	12.968	0.120	25.5	2752.3
3600	2.1e-03	1.8e+07	0.0e+00	0.0e+00	6.0e+01	6.0e+01	510.244	1102.270	24.644	9.538	20.7	53.5
3134	4.6e-03	2.6e+12	-9.7e-04	0.0e+00	7.5e+01	7.0e+01	178.360	1347.551	2.542	0.410	70.2	435.3
3562	1.3e-02	1.9e+11	0.0e+00	0.0e+00	8.6e+01	7.8e+01	1621.648	948.059	36.940	1.694	43.9	957.4
1922	8.2e-03	1.7e+08	0.0e+00	0.0e+00	9.6e+01	8.8e+01	182.146	448.809	4.524	0.405	40.3	450.0
4410	1.1e-02	9.5e+08	0.0e+00	8.7e+01	9.6e+01	8.3e+01	3170.182	2602.683	4.481	2.784	707.5	1138.7
588	6.2e-02	2.8e+04	9.9e+01	0.0e+00	9.0e+01	9.4e+01	19.195	7.152	0.463	1.379	41.5	13.9
494	6.8e-03	2.4e+06	9.0e+01	0.0e+00	9.1e+01	3.3e+01	9.222	2.959	0.832	0.116	11.1	79.8
662	5.6e-03	7.9e+05	9.5e+01	0.0e+00	7.4e+01	4.5e+01	16.293	5.947	0.324	0.203	50.2	80.3
1824	1.2e-02	1.9e+06	0.0e+00	0.0e+00	8.7e+01	7.1e+01	182.642	165.698	3.468	1.958	52.7	93.3
2146	1.6e-02	1.7e+03	0.0e+00	5.3e+01	5.5e+01	4.9e+01	428.711	137.240	2.604	16.508	164.6	26.0
2910	2.1e-02	6.0e+06	0.0e+00	0.0e+00	9.1e+01	7.5e+01	619.758	613.440	8.306	3.731	74.6	166.1
237	1.8e-02	2.0e+07	2.1e+01	5.5e+00	7.9e+01	6.5e+01	2.178	0.620	0.276	0.056	7.9	38.8
957	4.5e-03	5.1e+09	2.8e+01	3.5e+01	6.0e+01	6.5e+01	10.107	8.676	6.419	0.145	1.6	69.7
100	5.9e-02	1.6e+03	3.8e+01	3.1e+01	3.5e+01	3.5e+01	0.800	0.297	0.036	0.970	22.2	0.8
468	2.4e-02	1.1e+04	8.3e+01	7.9e+01	7.7e+01	7.1e+01	15.428	2.957	1.081	1.298	14.3	11.9
729	8.7e-03	2.4e+09	9.2e+01	8.9e+01	8.6e+01	8.4e+01	39.553	6.439	0.987	1.108	40.1	35.7

Table 4.2: Comparison of Algorithms for Min κ on Random

Dim & Density, & $\kappa(A)$			% Reduction in κ				CPU Time (seconds)				CPU Ratios	
n	density	$\kappa(A)$	[24]	[12]	Algorithm 2.2	Algorithm B.1	[24]	[12]	Algorithm 2.2	Algorithm B.1	[24]/Algorithm 2.2	[24]/Algorithm B.1
2000	1.0e-03	2.0e+06	-4.8e+01	4.8e+01	5.3e+01	6.6e+01	14.498	76.446	4.470	1.154	3.2	12.6
2500	8.1e-04	2.5e+06	-5.0e+01	2.9e+01	5.2e+01	5.2e+01	21.500	229.744	9.068	1.233	2.4	17.4
3000	6.7e-04	3.0e+06	-1.8e+01	0.0e+00	1.9e+01	1.3e+01	32.757	358.866	5.011	0.476	6.5	68.8
3500	5.8e-04	3.5e+06	-1.9e+01	2.5e+01	4.3e+01	4.9e+01	43.032	565.527	8.757	1.961	4.9	21.9
4000	5.0e-04	4.0e+06	-4.4e+01	0.0e+00	2.8e+01	6.3e+00	56.246	631.303	11.281	0.519	5.0	108.4
4500	4.5e-04	4.5e+06	-3.0e+01	0.0e+00	3.3e+01	4.2e+01	69.000	888.172	3.774	2.202	18.3	31.3
5000	4.0e-04	5.0e+06	-4.4e+01	0.0e+00	3.9e+01	3.9e+01	85.069	2009.496	13.994	2.341	6.1	36.3
5500	3.7e-04	5.5e+06	-3.5e+01	1.4e+01	2.8e+01	2.8e+01	104.179	2512.849	2.763	2.614	37.7	39.8
6000	3.3e-04	6.0e+06	-4.3e+01	0.0e+00	3.0e+01	3.4e+01	127.060	3680.562	5.383	2.748	23.6	46.2
6500	3.1e-04	6.5e+06	-2.9e+01	1.9e+01	2.7e+01	3.2e+01	146.476	3755.564	4.694	3.124	31.2	46.9
7000	2.9e-04	7.0e+06	-2.9e+01	2.4e+01	3.3e+01	3.0e+01	168.726	3838.763	11.441	3.243	14.7	52.0
7500	2.7e-04	7.5e+06	-4.1e+01	2.5e+01	3.2e+01	2.8e+01	193.993	3879.417	12.929	3.649	15.0	53.2
8000	2.5e-04	8.0e+06	-1.9e+01	3.0e+01	3.2e+01	2.7e+01	215.347	3790.944	16.234	3.645	13.3	59.1
8500	2.4e-04	8.5e+06	-2.8e+01	3.6e+01	2.5e+01	2.6e+01	246.213	3929.044	8.086	4.716	30.5	52.2
9000	2.2e-04	9.0e+06	-1.5e+01	-9.9e+00	9.1e+00	3.2e+00	271.863	3903.914	7.354	0.931	37.0	292.1
9500	2.1e-04	9.5e+06	-4.2e+01	3.6e+01	2.6e+01	2.3e+01	323.748	4147.347	13.329	5.681	24.3	57.0
10000	2.0e-04	1.0e+07	-3.2e+01	1.2e+01	2.2e+01	1.8e+01	341.504	4304.042	8.240	5.232	41.4	65.3
10500	1.9e-04	1.0e+07	-1.2e+01	0.0e+00	2.7e+01	2.1e+01	384.690	4302.243	20.255	6.627	19.0	58.0
11000	1.8e-04	1.1e+07	-4.1e+01	3.0e+00	2.2e+01	1.9e+01	463.158	4273.218	11.945	6.085	38.8	76.1
11500	1.7e-04	1.1e+07	-3.4e+01	0.0e+00	2.3e+01	2.0e+01	487.216	4000.550	9.915	6.705	49.1	72.7
12000	1.7e-04	1.2e+07	-3.0e+01	0.0e+00	2.2e+01	1.9e+01	511.686	4581.828	10.963	6.709	46.7	76.3
12500	1.6e-04	1.3e+07	-3.1e+01	0.0e+00	4.0e+00	3.2e+00	546.856	4059.532	3.496	2.080	156.4	263.0
13000	1.5e-04	1.3e+07	-2.9e+01	0.0e+00	2.4e+01	1.8e+01	593.252	4586.045	40.099	7.383	14.8	80.4
13500	1.5e-04	1.4e+07	-3.0e+01	0.0e+00	2.2e+01	1.8e+01	665.554	5145.811	15.941	8.300	41.8	80.2
14000	1.4e-04	1.4e+07	-3.0e+01	0.0e+00	2.2e+01	1.7e+01	719.336	4122.784	14.430	7.864	49.9	91.5
14500	1.4e-04	1.4e+07	-3.4e+01	0.0e+00	2.1e+01	1.2e+01	901.989	4606.777	24.572	6.888	36.7	130.9
15000	1.3e-04	1.5e+07	-4.1e+01	0.0e+00	1.6e+01	1.1e+01	920.641	5115.021	14.827	6.354	62.1	144.9

763 **4.1.2 Large Sizes**

764 Tables 4.3 and 4.4 illustrate the scalability of Algorithms B.1 and 2.2 on substantially larger instances for
 765 which the SDP-based methods of [12,24] became computationally impractical. Table 4.3 considers matrices
 766 from the SuiteSparse Matrix Collection, while Table 4.4 considers randomly generated matrices with very
 767 large condition numbers.

768 On the SuiteSparse instances, Algorithm 2.2 and Algorithm B.1 reduced κ by averages of 45.7% and
 769 41.9%, respectively, while often requiring less than one minute. On the random instances, which satisfy
 770 $\kappa(M) > 10^{11}$, the corresponding average reductions were 11.6% and 4.4%.

Table 4.3: Comparison of Algorithms for Min κ on Large SuiteSparse

Dim & Density, & $\kappa(M)$			% Reduction in κ		CPU Time	
n	density	$\kappa(M)$	Algorithm 2.2	Algorithm B.1	Algorithm 2.2	Algorithm B.1
14822	3.3e-03	2.0e+06	3.9e+01	3.0e+01	7.850	283.641
15439	1.1e-03	4.4e+12	8.0e+01	7.2e+01	8.734	2.185
15439	6.5e-05	6.1e+09	4.2e+01	5.4e+01	2.283	2.725
17361	1.1e-03	2.5e+06	1.3e+01	1.2e+01	1.097	9.086
17361	3.4e-03	1.1e+09	5.6e+01	3.0e+01	19.675	9.306
23052	2.2e-03	7.4e+11	9.0e+01	7.2e+01	503.668	6.933
30401	5.1e-04	5.8e+03	9.1e+01	9.3e+01	204.718	14.951
36441	4.3e-04	2.6e+03	8.5e+01	9.2e+01	459.869	59.632
40806	1.2e-04	8.1e+01	4.5e-02	1.1e-01	29.901	16.301
48962	2.1e-04	1.6e+05	6.2e+00	4.1e+00	22.109	43.123
150102	3.2e-05	1.3e+07	4.6e-03	1.5e+00	53.450	864.384

Table 4.4: Comparison of Algorithms for Min κ on Large Random

Dim & Density, & $\kappa(M)$			% Reduction in κ		CPU Time	
n	density	$\kappa(M)$	Algorithm 2.2	Algorithm B.1	Algorithm 2.2	Algorithm B.1
50000	3.8e-05	1.0e+11	1.7e+01	7.7e+00	244.577	74.798
60000	3.2e-05	1.2e+11	1.5e+01	6.5e+00	400.467	103.404
70000	2.7e-05	1.4e+11	1.4e+01	5.6e+00	492.698	123.284
80000	2.4e-05	1.6e+11	1.1e+01	5.0e+00	436.111	198.894
90000	2.1e-05	1.8e+11	1.0e+01	4.3e+00	550.245	223.299
100000	1.9e-05	2.0e+11	1.2e+01	4.0e+00	986.187	250.367
110000	1.7e-05	2.2e+11	1.1e+01	3.7e+00	1162.321	293.284
120000	1.6e-05	2.4e+11	1.1e+01	3.4e+00	1172.119	332.887
130000	1.5e-05	2.6e+11	8.5e+00	3.1e+00	1274.404	386.352
140000	1.4e-05	2.8e+11	9.1e+00	2.9e+00	1500.467	480.042
150000	1.3e-05	3.0e+11	8.7e+00	2.7e+00	1677.360	498.512

771 In summary, the results show that the proposed subgradient-based methods not only scale substantially
 772 better than existing SDP-based approaches, but also achieve significantly larger reductions in κ at a fraction
 773 of the computational cost.

774 **4.2 PCG Comparison for Solving Linear Systems**

775 In this subsection, we evaluate the effect of preconditioning on PCG performance for solving symmetric
 776 positive definite systems $Mx = b$. We compare the diagonal preconditioners computed by Algorithm B.1
 777 with those obtained using the κ -optimal scaling method from [24]. PCG is run with a relative residual
 778 tolerance of 10^{-4} and a maximum of 5×10^6 iterations. The results are reported in Table 4.5. All test
 779 matrices are randomly generated using MATLAB's `sprandsym` function with varying dimensions, densities,
 780 and κ values.

781 Table 4.5 shows that, over the 47 test instances, PCG requires on average 442,147 iterations with the
 782 preconditioner from Algorithm B.1, compared to 549,002 iterations with the preconditioner from [24]. Av-
 783 eraged over all test instances, PCG with Algorithm B.1 requires approximately 20% fewer iterations while
 784 reducing total CPU time by nearly an order of magnitude compared with the method of [24].

Table 4.5: PCG Comparison Using Preconditioners Found by Algorithm in [24] and Algorithm B.1

Dim, Density, & $\kappa(A)$			% Reduction in κ		PCG Iterations		PCG CPU Time		Total CPU Time	
n	density	$\kappa(M)$	[24]	Algorithm B.1	[24]	Algorithm B.1	[24]	Algorithm B.1	[24]	Algorithm B.1
1000	3.3e-03	1.0e+09	-2.6e+01	3.6e+01	115264	99023	1.14	0.99	5.33	1.29
1300	2.5e-03	1.3e+09	-2.6e+01	6.0e+01	153760	101197	1.89	1.23	8.41	1.49
1600	2.0e-03	1.6e+09	-4.3e+01	-1.0e-13	180929	179469	2.67	2.67	12.17	2.71
1900	1.7e-03	1.9e+09	-3.5e+01	6.9e+01	202276	121629	3.43	2.07	16.17	2.87
2200	1.4e-03	2.2e+09	-4.9e+01	7.1e+01	245521	126567	4.82	2.50	21.54	3.39
2500	1.2e-03	2.5e+09	-2.3e+01	2.7e+01	252757	221181	5.43	4.69	29.06	5.07
2800	1.1e-03	2.8e+09	-4.2e+01	1.9e+01	289351	248376	6.69	5.74	33.60	6.12
3100	9.9e-04	3.1e+09	-2.9e+01	6.6e+01	306195	174755	8.02	4.56	45.28	6.01
3400	9.0e-04	3.4e+09	-3.3e+01	5.4e+01	325838	215488	9.07	6.05	54.35	7.46
3700	8.3e-04	3.7e+09	-2.5e+01	4.4e+01	330583	256787	9.89	7.67	63.61	9.21
4000	7.6e-04	4.0e+09	-4.4e+01	5.7e+01	366269	234646	11.60	7.40	65.77	9.21
4300	7.1e-04	4.3e+09	-2.9e+01	5.4e+01	378370	252209	13.49	9.00	88.27	10.97
4600	6.6e-04	4.6e+09	-4.2e+01	4.1e-13	408179	384608	15.52	14.65	89.55	14.74
4900	6.2e-04	4.9e+09	-1.1e+01	8.7e+00	399371	382842	15.87	15.31	108.88	15.75
5200	5.8e-04	5.2e+09	-3.7e+01	-3.7e-14	458497	418282	19.04	17.32	115.01	17.43
5500	5.5e-04	5.5e+09	-3.3e+01	3.9e+01	460157	340511	19.97	14.79	128.37	16.99
5800	5.2e-04	5.8e+09	-3.2e+01	4.6e+01	469674	336957	21.25	15.32	143.80	17.73
6100	4.9e-04	6.1e+09	-5.0e+01	4.3e+01	506736	353612	23.89	16.68	157.42	19.34
6400	4.7e-04	6.4e+09	-2.2e+01	4.3e+01	492739	361695	23.98	17.59	166.18	20.56
6700	4.4e-04	6.7e+09	-4.0e+01	4.1e+01	526574	378065	25.65	18.54	189.81	21.40
7000	4.2e-04	7.0e+09	-2.9e+01	3.9e+01	533440	394622	28.16	20.13	202.50	22.93
7300	4.1e-04	7.3e+09	-5.3e+01	2.7e+00	572826	511612	32.32	27.95	234.05	28.40
7600	3.9e-04	7.6e+09	-2.9e+01	-3.8e-14	555453	528584	31.27	29.77	232.71	29.89
7900	3.7e-04	7.9e+09	-2.2e+01	3.5e+01	569653	434589	33.24	25.44	274.74	29.14
8200	3.6e-04	8.2e+09	-6.6e+01	3.5e+01	614529	448278	37.09	26.97	323.07	30.42
8500	3.5e-04	8.5e+09	-2.7e+01	3.4e+01	573590	458996	35.49	28.28	292.58	32.18
8800	3.3e-04	8.8e+09	-3.6e+01	1.6e+01	625430	532398	40.04	34.12	322.79	35.85
9100	3.2e-04	9.1e+09	-1.8e+01	1.5e+01	597815	541439	39.21	35.60	324.73	37.50
9400	3.1e-04	9.4e+09	-2.0e+01	1.1e+01	635767	565196	42.77	37.92	374.03	39.77
9700	3.0e-04	9.7e+09	-2.9e+01	3.1e+01	639683	508698	44.28	35.17	391.15	40.05
10000	2.9e-04	1.0e+10	-4.5e+01	3.1e+01	687294	521357	48.91	37.36	454.89	41.62
10300	2.8e-04	1.0e+10	-2.5e+01	2.9e+01	667536	538749	49.08	39.58	506.66	44.78
10600	2.8e-04	1.1e+10	-4.4e+01	2.9e+01	725311	548238	54.63	41.32	510.40	46.14
10900	2.7e-04	1.1e+10	-3.9e+01	2.6e+01	725604	566231	55.80	43.64	491.97	49.26
11200	2.6e-04	1.1e+10	-2.4e+01	2.8e+01	693639	569480	55.03	44.93	498.31	49.79
11500	2.5e-04	1.1e+10	-8.9e+00	2.4e+01	673910	597322	54.77	48.20	521.32	53.85
11800	2.5e-04	1.2e+10	-3.9e+01	1.7e+01	757749	627850	62.92	52.02	581.98	55.67
12100	2.4e-04	1.2e+10	-2.9e+01	2.6e+01	748235	601448	63.64	51.45	640.96	57.59
12400	2.3e-04	1.2e+10	-4.3e+01	2.5e+01	777435	614047	68.12	54.42	653.76	60.37
12700	2.3e-04	1.3e+10	-3.3e+01	2.5e+01	775871	620193	70.07	56.35	686.85	63.40
13000	2.2e-04	1.3e+10	-3.2e+01	1.1e+01	771101	681834	73.02	65.22	703.04	68.91
13300	2.2e-04	1.3e+10	-2.7e+01	2.4e+01	766793	648685	75.93	64.69	727.64	71.98
13600	2.1e-04	1.4e+10	-3.5e+01	1.0e+01	814120	708483	82.37	72.15	752.62	75.16
13900	2.1e-04	1.4e+10	-5.5e+01	2.3e+01	863835	662119	89.35	68.57	824.58	76.07
14200	2.0e-04	1.4e+10	-4.4e+01	2.3e+01	858674	673736	90.12	71.18	887.64	78.25
14500	2.0e-04	1.5e+10	-4.3e+01	9.4e+00	872710	741010	94.29	80.77	869.02	84.41
14800	2.0e-04	1.5e+10	-3.0e+01	8.5e+00	836054	747797	91.59	82.12	874.94	84.96

4.3 LSQR Comparison: Algorithm 3.1 vs. Two-Sided κ -optimal Scaling in [24]

In Section 4.3, we compare our two-sided Algorithm 3.1, which decreases ω at every iteration, with the two-sided κ -optimal preconditioner of [24] for solving $\min_x \|b - Ax\|$ with LSQR.

The results are reported in Tables 4.6 to 4.8. Table 4.6 and Table 4.7 consider SuiteSparse and randomly generated matrices, respectively, with $n \leq 300$, since the method of [24] could not handle larger instances. In these experiments, LSQR is run with tolerance 10^{-8} and a maximum of 5000 iterations. Table 4.8 considers SuiteSparse matrices with $n \geq 2000$ and compares LSQR with and without the preconditioner from Algorithm 3.1. Here, LSQR is run with tolerance 10^{-6} and a maximum of 10000 iterations.

In all tables, $S := (\hat{C}^{1/2} A \hat{D}^{1/2})^T (\hat{C}^{1/2} A \hat{D}^{1/2})$ and $T := (D_1^{1/2} A D_2^{-1/2})^T (D_1^{1/2} A D_2^{-1/2})$ where $\hat{C}^{1/2}$, $\hat{D}^{1/2}$ and $D_1^{1/2}$, $D_2^{-1/2}$ are the preconditioners computed by Algorithm 3.1 and [24], respectively. Reported CPU times include both preconditioning and LSQR runtimes.

Table 4.6: LSQR Comparison on Small Suitesparse Matrices: Algorithm 3.1 vs [24]

Dim & Density, & $\kappa(A)$			Ratio of Condition Numbers		CPU Time for Prec		LSQR Iterations		Total CPU Time	
n	density	$\kappa(A)$	$\kappa(S)/\kappa(T)$	$\omega(S)/\omega(T)$	Algorithm 3.1	[24]	Algorithm 3.1	[24]	Algorithm 3.1	[24]
300	6.6e-03	5.2e+03	1.3e+00	1.0e+00	0.006	23.358	5	4	0.015	23.36
130	6.1e-02	6.1e+10	1.1e-01	8.2e-01	0.002	27.943	9	35	0.005	27.95
200	2.0e-02	2.4e+03	1.0e+00	1.0e+00	0.000	19.308	750	756	0.022	19.32
104	9.2e-02	5.5e+03	1.6e+00	9.7e-01	0.001	10.056	301	288	0.004	10.06
216	9.3e-02	3.3e+04	8.9e-01	4.6e-01	0.002	309.028	56	175	0.004	309.03
115	9.3e-02	3.7e+02	1.5e+00	6.7e-01	0.001	14.324	267	363	0.004	14.33
185	2.8e-02	1.8e+05	5.7e+00	7.5e-01	0.003	40.382	242	328	0.006	40.39
216	1.7e-02	1.0e+02	1.3e+00	9.0e-01	0.001	17.672	183	215	0.004	17.67
207	1.3e-02	1.4e+08	3.9e-01	3.4e-01	0.003	19.862	122	709	0.006	19.87
137	2.1e-02	1.8e+04	1.2e+00	7.9e-01	0.002	12.505	61	95	0.003	12.51
225	2.6e-02	7.1e+06	2.5e+00	8.1e-01	0.003	146.868	78	86	0.004	146.87
198	1.4e-01	3.0e+03	1.3e+00	5.8e-01	0.002	88.701	631	973	0.010	88.71
265	2.5e-02	1.4e+03	1.6e+00	5.0e-01	0.002	38.714	735	1379	0.010	38.73
100	4.0e-02	1.5e+04	1.0e+00	9.7e-01	0.000	8.944	114	132	0.002	8.95
105	8.0e-02	7.2e+02	1.5e+00	6.6e-01	0.000	14.751	163	227	0.002	14.75
135	3.6e-02	9.2e+05	1.0e+00	2.3e-01	0.000	17.484	164	850	0.002	17.49
120	6.0e-02	4.3e+08	3.1e-03	2.7e-02	0.001	24.967	77	845	0.003	24.97
100	7.1e-02	2.4e+12	5.8e-02	2.4e-01	0.001	28.979	19	108	0.003	28.98
136	2.6e-02	2.5e+05	1.9e+00	8.0e-01	0.001	19.051	132	175	0.003	19.05
100	4.0e-02	1.3e+04	2.3e+00	9.3e-01	0.000	10.800	339	385	0.003	10.80
300	3.5e-02	8.5e+05	2.2e+00	6.0e-01	0.002	581.259	1195	1983	0.016	581.28
132	2.4e-02	4.2e+11	2.1e+00	1.0e-01	0.000	22.553	93	1356	0.002	22.56

Table 4.7: LSQR Comparison on Small Random Matrices: Algorithm 3.1 vs [24]

Dim & Density, & $\kappa(A)$			Ratio of Condition Numbers		CPU Time for Prec		LSQR Iterations		Total CPU Time	
n	density	$\kappa(A)$	$\kappa(S)/\kappa(T)$	$\omega(S)/\omega(T)$	Algorithm 3.1	[24]	Algorithm 3.1	[24]	Algorithm 3.1	[24]
50	2.5e-01	1.0e+03	1.2e+00	9.4e-01	0.007	16.583	18	18	0.02	16.59
60	2.4e-01	1.0e+03	5.2e+00	5.1e-01	0.001	11.914	46	67	0.01	11.92
70	2.3e-01	1.1e+03	1.4e+00	8.3e-01	0.001	17.909	22	27	0.00	17.91
80	2.2e-01	1.1e+03	2.0e+00	2.6e-01	0.002	13.848	51	191	0.00	13.85
90	2.2e-01	1.2e+03	2.2e+00	5.0e-01	0.002	18.674	43	90	0.00	18.68
100	2.1e-01	1.2e+03	5.2e+00	4.9e-01	0.003	27.400	52	72	0.00	27.40
110	2.1e-01	1.3e+03	2.2e+00	6.4e-01	0.003	45.059	52	52	0.00	45.06
120	2.0e-01	1.4e+03	2.0e+00	7.0e-01	0.002	37.248	63	57	0.00	37.25
130	2.0e-01	1.5e+03	7.3e+00	4.0e-01	0.002	42.239	96	140	0.00	42.24
140	1.9e-01	1.6e+03	1.1e+01	5.4e-01	0.002	55.521	54	83	0.00	55.52
150	1.9e-01	1.7e+03	2.4e+00	8.0e-01	0.001	102.427	46	35	0.00	102.43
160	1.9e-01	1.8e+03	6.2e+00	5.3e-01	0.002	90.146	66	81	0.00	90.15
170	1.9e-01	1.9e+03	2.4e+01	5.1e-01	0.003	101.889	84	111	0.01	101.89
180	1.8e-01	2.1e+03	9.3e+00	4.3e-01	0.001	128.101	90	109	0.00	128.10
190	1.8e-01	2.2e+03	8.1e+00	3.4e-01	0.001	116.720	113	219	0.00	116.72
200	1.8e-01	2.5e+03	9.1e+00	4.7e-01	0.001	147.824	94	124	0.00	147.83
210	1.8e-01	2.7e+03	1.2e+01	5.2e-01	0.002	191.344	121	113	0.00	191.35
220	1.8e-01	3.1e+03	4.6e+00	5.3e-01	0.002	293.214	94	75	0.00	293.22
230	1.7e-01	3.5e+03	2.2e+00	3.8e-01	0.001	175.652	105	508	0.00	175.66
240	1.7e-01	4.0e+03	6.1e+00	4.5e-01	0.001	287.572	129	116	0.00	287.57
250	1.7e-01	4.8e+03	1.1e+01	5.4e-01	0.002	433.302	107	92	0.00	433.30
260	1.7e-01	5.9e+03	4.0e+00	5.1e-01	0.002	384.679	139	166	0.01	384.68
270	1.7e-01	7.8e+03	3.6e+00	6.3e-01	0.002	431.242	121	75	0.00	431.24
280	1.7e-01	1.1e+04	1.5e+01	3.9e-01	0.002	498.242	153	298	0.01	498.25
290	1.7e-01	2.0e+04	1.6e+01	4.5e-01	0.002	750.635	178	128	0.01	750.64
300	1.7e-01	1.0e+05	1.0e+01	2.8e-01	0.002	751.572	197	264	0.01	751.58

Table 4.8: LSQR Comparison on Large Suitesparse Matrices: Algorithm 3.1 vs No Preconditioning

Dim & Density, & $\omega(A)$			Ratio of Omega & Prec CPU Time		LSQR Iter		Total CPU Time	
n	density	$\omega(A)$	$\omega(S)/\omega(A)$	Algorithm 3.1 CPU	No Prec	Algorithm 3.1	No Prec	Algorithm 3.1
14214	1.3e-03	5.1e+01	3.2e-02	1.0e-01	100000	8963	3.63e+01	3.43e+00
4119	2.1e-03	4.1e+05	3.6e-06	1.7e-02	5461	946	4.98e-01	1.04e-01
7479	1.2e-03	2.8e+04	4.8e-05	3.0e-02	9393	1758	1.85e+00	3.81e-01
9271	1.4e-03	8.2e+06	1.7e-07	4.1e-02	1798	377	4.44e-01	1.37e-01
4562	6.3e-03	2.3e+02	7.4e-03	1.3e-02	6497	1070	1.11e+00	2.00e-01
3072	1.3e-02	1.9e+00	1.0e+00	9.8e-03	1358	1357	5.43e-01	6.06e-01
11341	7.6e-04	6.0e+01	5.0e-02	3.9e-02	26647	3050	7.55e+00	9.09e-01
2048	2.9e-03	2.8e+00	9.3e-01	1.6e-03	4730	4187	2.10e-01	1.93e-01
8192	7.3e-04	2.8e+00	9.4e-01	6.1e-03	19233	16963	3.15e+00	2.79e+00
14734	4.4e-04	2.3e+00	7.2e-01	4.2e-02	10631	6883	3.43e+00	2.26e+00
2283	9.2e-03	3.6e+01	5.6e-02	9.1e-03	17604	1305	1.15e+00	9.74e-02
2000	2.0e-03	2.5e+00	1.0e+00	7.9e-04	1578	1578	6.33e-02	6.54e-02
2000	5.9e-03	4.2e+00	5.9e-01	3.5e-03	9674	3842	4.76e-01	1.95e-01
5000	2.4e-03	4.2e+00	5.7e-01	6.9e-03	29362	7483	4.56e+00	1.15e+00
7000	1.7e-03	4.3e+00	5.7e-01	9.3e-03	50422	13747	9.93e+00	2.82e+00
7000	1.7e-03	4.2e+00	5.7e-01	9.4e-03	48882	11272	9.71e+00	2.23e+00
3175	8.4e-03	3.0e+01	7.0e-02	9.3e-03	16004	2157	2.06e+00	2.90e-01
5952	6.3e-04	3.7e+02	4.6e-03	1.8e-02	13984	5192	1.55e+00	5.96e-01
13694	3.9e-04	1.8e+02	7.1e-03	8.0e-02	100000	2551	3.08e+01	8.70e-01
5832	9.0e-03	1.7e+00	1.0e+00	1.0e-02	2397	2390	6.25e-01	6.36e-01
9801	9.1e-04	1.5e+00	1.0e+00	2.6e-03	3132	3036	7.93e-01	7.72e-01
4000	5.5e-04	1.3e+09	7.5e-10	7.3e-03	7194	5	5.05e-01	8.63e-03
5940	2.4e-03	6.6e+00	4.6e-01	3.6e-02	100000	44539	1.74e+01	7.78e+00
4326	3.3e-03	1.1e+05	1.7e-05	1.4e-02	34107	2486	4.85e+00	3.69e-01
9800	6.0e-04	2.6e+00	9.9e-01	6.4e-03	18382	18177	4.37e+00	4.29e+00

796 The results in Tables 4.6 and 4.7 show that although the two-sided κ -optimal preconditioner in [24]
797 often achieves a stronger reduction in κ , Algorithm 3.1 consistently yields lower ω values at a significantly
798 lower computational cost. Importantly, reductions in ω correlate more strongly with LSQR iteration counts,
799 resulting in faster convergence on the majority of tested instances.

800 In summary, these experiments reinforce that ω -optimal preconditioning offers a more practical and com-
801 putationally efficient alternative to κ -optimal scaling for LSQR, particularly in large-scale and nonsymmetric
802 settings.

803 4.4 From κ -optimal M to “Improve PCG” using ω -Optimal Scaling

804 We study the effect of applying ω -optimal (Jacobi) diagonal scaling to a κ -optimally scaled matrix $M > 0$.¹⁰
805 To this end, we construct the matrix M using Theorem 2.7: the extreme eigenpairs satisfy (2.5) and (2.6),
806 while the remaining eigenvalues are spread uniformly in (λ_n, λ_1) with orthonormal eigenvectors chosen in
807 the orthogonal complement of $\text{span}\{x_1, x_n\}$. The remaining eigenvectors are generated via a sparse QR
808 factorization, so the density of M cannot be prescribed exactly in advance. A more detailed description of
809 the construction is provided in our code.

810 After constructing M , we apply the ω -optimal (Jacobi) scaling $J = D^{1/2}MD^{1/2}$ where $D = \text{Diag}(1./\text{diag}(M))$,
811 and compare PCG applied to $Mx = b$ and to the scaled system $J(D^{-1/2}x) = D^{1/2}b$. Iteration counts and
812 runtimes are averaged over five initial points. The results are reported in Tables 4.9 and 4.10.

¹⁰For the large problems in Table 4.10 we used the fastlinux server, cpu155.math.private, a Dell PowerEdge R650 with two Intel Xeon Gold 6334 8-core 3.6 GHz (Ice Lake) processors and 256 GB RAM.

Table 4.9: PCG tol. $1e-7$; Medium; κ -opt **VS** J, ω -opt of M

Dim & Density		Ratios in conds (J, M)		J : ω -opt of M		Ratios M/J	
n	density	$\kappa(J)/\kappa(M)$	$\omega(J)/\omega(M)$	iters	cpu	iters	cpu
5000	9.40e-03	2.98e+00	7.368e-01	19.4	3.927e-03	25.2	13.9
10000	7.44e-03	3.31e+00	7.362e-01	18.8	6.698e-03	36.6	30.5
15000	6.02e-03	2.64e+00	7.365e-01	17.0	1.418e-02	41.6	39.0
20000	4.57e-03	3.56e+00	7.363e-01	20.6	2.861e-02	41.6	41.1
25000	3.47e-03	3.71e+00	7.370e-01	20.6	8.006e-02	30.9	29.1
30000	2.42e-03	3.52e+00	7.362e-01	20.2	4.327e-02	48.0	43.0
35000	1.56e-03	2.77e+00	7.371e-01	16.0	4.174e-02	36.5	36.1
40000	8.93e-04	2.89e+00	7.364e-01	18.0	2.535e-02	46.7	38.9
45000	4.15e-04	3.93e+00	7.370e-01	21.8	1.820e-02	28.5	27.0
50000	4.12e-04	4.12e+00	7.361e-01	24.0	2.679e-02	49.5	43.7

Table 4.10: PCG tol. $1e-7$; Large; κ -opt **VS** J, ω -opt of M

Dim & Density		Ratios in conds (J, M)		J : ω -opt of M		Ratios M/J	
n	density	$\kappa(J)/\kappa(M)$	$\omega(J)/\omega(M)$	iters	cpu	iters	cpu
50000	9.05e-03	3.26e+00	7.366e-01	15.6	3.332e-01	48.7	42.5
52000	8.17e-03	4.34e+00	7.362e-01	20.8	4.245e-01	51.7	46.9
54000	7.49e-03	3.55e+00	7.365e-01	17.0	3.578e-01	47.5	42.4
56000	6.71e-03	3.73e+00	7.363e-01	18.0	3.672e-01	57.1	51.0
58000	5.98e-03	3.33e+00	7.369e-01	16.8	3.367e-01	38.5	34.1
60000	5.24e-03	3.90e+00	7.362e-01	17.6	3.316e-01	62.1	55.4
62000	4.63e-03	4.43e+00	7.371e-01	18.0	3.233e-01	33.2	29.7
64000	4.01e-03	4.37e+00	7.365e-01	22.0	3.582e-01	38.9	35.5
66000	3.46e-03	3.24e+00	7.369e-01	16.0	2.457e-01	38.9	34.4
68000	2.87e-03	4.31e+00	7.361e-01	20.8	2.756e-01	59.8	54.0
70000	2.40e-03	4.25e+00	7.363e-01	20.4	2.445e-01	46.5	42.0
72000	1.96e-03	4.05e+00	7.369e-01	17.8	1.883e-01	36.4	32.1
74000	1.54e-03	3.68e+00	7.362e-01	17.2	1.438e-01	64.8	56.9
76000	1.22e-03	3.95e+00	7.366e-01	18.0	1.194e-01	43.6	38.4
78000	9.02e-04	2.98e+00	7.361e-01	16.4	8.235e-02	72.1	62.2
80000	6.44e-04	3.50e+00	7.366e-01	17.0	6.009e-02	43.4	36.6
82000	4.20e-04	3.81e+00	7.366e-01	19.2	4.669e-02	40.2	33.9
84000	2.46e-04	4.12e+00	7.363e-01	20.2	3.227e-02	46.2	37.8
86000	1.20e-04	2.92e+00	7.370e-01	17.0	1.930e-02	35.8	27.8
88000	4.12e-05	3.25e+00	7.366e-01	21.6	1.801e-02	36.7	23.2
90000	1.14e-05	1.81e+00	7.366e-01	18.8	1.303e-02	38.2	21.9

813 The experiments show that surprisingly, despite increasing the κ -condition number, the ω -optimal (Ja-
814 cobi) scaling consistently yields dramatic improvements in PCG convergence. Across all medium and large
815 instances presented in Tables 4.9 and 4.10, PCG requires between 96–98.6% fewer iterations and correspond-
816 ingly less CPU time after ω -optimal scaling. This improvement is especially pronounced on the large-scale
817 instances presented in Table 4.10.

818 In summary, our computational experiments in Section 4 provide compelling evidence that the ω -condition
819 number is a stronger predictor of iterative solver performance than κ , and that ω -optimal preconditioners
820 strike a superior balance between computational cost and practical effectiveness.

821 5 Conclusion

822 In this paper, we studied optimal diagonal preconditioning through the lens of two condition numbers: the
823 classical κ and the averaging-based ω -condition number. On the theoretical side, we introduced an affine-

824 based pseudoconvex reformulation of the κ -optimal preconditioning problem, yielding simple optimality
825 conditions and enabling efficient optimization over an n -dimensional vector. Moreover, since the κ -condition
826 number is positively homogeneous of degree zero, the associated minimization problem is Hadamard ill-
827 posed. To address this, we introduce a reformulation that removes this homogeneity, leading to improved
828 computational stability in practice. Building on this formulation, we develop a highly efficient subgradient
829 method with convergence guarantees that significantly outperforms existing SDP-based approaches in both
830 scalability and accuracy. For example, Table 4.1 shows that Algorithm 2.2 (resp. Algorithm B.1) is, on
831 average, 77.48 (resp. 356.22) times faster than [24] when minimizing κ on relatively small SuiteSparse
832 matrices. Our methods also substantially outperform [12] in terms of runtime. In addition to these speedups,
833 our subgradient methods consistently achieve significantly greater reductions in κ than both [24] and [12],
834 often reducing it by more than half. Finally, our algorithms scale to large problem instances with hundreds
835 of thousands of variables, solving them within minutes, whereas existing methods [12, 24] struggle to handle
836 such instances within a reasonable time.

837 In parallel, we provided explicit and unified characterizations of ω -optimal diagonal and block-diagonal
838 preconditioners, showing that many classical schemes such as Jacobi scaling, row/column normalization, and
839 matrix balancing arise naturally as ω -optimal solutions or stationary points of the ω -optimal preconditioning
840 problem. These results offer a new perspective on widely used preconditioning techniques and, to the best
841 of our knowledge, constitute the first comprehensive comparison of κ - and ω -based optimality conditions.

842 Our numerical results further reveal a clear and practically important message: while κ -optimal precondi-
843 tioners reduce the worst-case condition number more aggressively, ω -optimal preconditioners are substantially
844 cheaper to compute and more strongly correlated with the performance of iterative methods such as PCG
845 and LSQR. Moreover, applying the ω -optimal diagonal scaling to linear systems that are already κ -optimally
846 preconditioned leads to further significant improvements in PCG’s performance.

847 Overall, our findings suggest that ω -based preconditioning provides a computationally efficient and prac-
848 tically superior alternative to κ -based approaches for large-scale problems, and highlight the importance of
849 moving beyond worst-case conditioning when designing preconditioners.

850 Statements and Declarations

851 **Funding.** The first author acknowledges the support of the Natural Sciences and Engineering Research
852 Council of Canada (NSERC) [Application ID: RGPIN-2021-02644]. The second and fifth authors acknowl-
853 edge support from Natural Sciences and Engineering Research Council of Canada (NSERC) [Work Order
854 Number: 50503-11728 2925 105]. The third author was partially supported by NSERC [Application ID:
855 RGPIN-2021-02644], NSERC [Work Order Number: 50503-11728 2925 105], the Department of Combina-
856 torics at the University of Waterloo, and Professor Walaa Moursi. The fourth author was partially supported
857 by grant PID2022-136399NB-C21 funded by ERDF/EU and by MICIU/AEI/10.13039/501100011033; by
858 Centro de Modelamiento Matemático (CMM) BASAL fund FB210005 for center of excellence; and Fondecyt
859 Postdoctorado 3250039 from ANID Chile.

860 **Competing Interests.** The authors have no relevant financial or non-financial interests to disclose.

861 **Data Availability.** The datasets and matrices generated and analyzed during the current study are publicly
862 available in the GitHub repository: https://github.com/asujanani6/Optimal_Preconditioning.

863 **Code Availability.** The full code used to generate the results and experiments in this paper is publicly
864 available in the GitHub repository: https://github.com/asujanani6/Optimal_Preconditioning.

865 A Assumptions and Technical Results from [19]

866 Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}} := \mathbb{R} \cup \{-\infty, \infty\}$ be an extended real-valued function with *domain* \bar{D} . The following
867 minimization problem is considered in [19]:

$$f_* := \inf \{f(x) : x \in X\}, \quad (\text{A.1})$$

868 with the following assumptions on f and X :

869 **A1** : $\text{int } \bar{D}$ is nonempty and convex.

870 **A2** : $\bar{\lim}_{t \downarrow 0} f(x + t(y - x)) \leq f(x), \forall x \in \bar{D}, y \in \text{int } \bar{D}$.

871 **A3** : f is upper semicontinuous (usc) on $\text{int } \bar{D}$, i.e., $f(x) = \lim_{\epsilon \downarrow 0} \sup_{B(x, \epsilon)} f, \forall x \in \text{int } \bar{D}$ where $B(x, \epsilon) :=$
872 $\{y : \|y - x\| \leq \epsilon\}$.

873 **A4** : f is quasiconvex on $\text{int } \bar{D}$, i.e., the set $\{x \in \text{int } \bar{D} : f(x) \leq \alpha\}$ is convex $\forall \alpha \in \mathbb{R}$.

874 **A5** : The constraint set $X \subset \mathbb{R}^n$ is closed convex, and $X \cap \text{int } \bar{D} \neq \emptyset$.

875 The following lemma from [19] discusses that fractional programs with certain structures are quasiconvex
876 and also provides one characterization of their quasisubgradients.

877 **Lemma A.1** ([19, Lemma 4]). *Suppose $f(x) = a(x)/b(x)$ for all $x \in \text{int } \bar{D}$, where a is a convex function, b
878 is finite and positive on $\text{int } \bar{D}$, $\text{int } \bar{D}$ is convex, and one of the following conditions holds:*

879 (a) b is affine;

880 (b) a is nonnegative on $\text{int } \bar{D}$ and b is concave;

881 (c) a is nonpositive on $\text{int } \bar{D}$ and b is convex.

882 *Then f is quasiconvex on $\text{int } \bar{D}$ and for each $x \in \text{int } \bar{D}$, if $\alpha := f(x)$ is finite then $a - \alpha b$ is convex and
883 $\partial[a - \alpha b](x) \subset \bar{\partial}^\circ f(x)$.*

884 The following lemma can be found in [19] and is useful for showing that assumptions (A1)-(A4) hold for
885 our set-up in (2.8) since $\kappa(\mathcal{D}(d))$ is the ratio of a convex and concave function that is positive on Ω .

886 **Lemma A.2** ([19, Remark 2]). *Suppose a and $\tilde{b} := -b$ are proper convex functions on $\text{int } \bar{D}^a \cap \text{int } \bar{D}^b$, a is
887 nonnegative on the (open and convex set) $C := \{y \in \text{int } \bar{D}^a \cap \text{int } \bar{D}^b : b(y) > 0\} \neq \emptyset$, and*

$$f(x) := \begin{cases} a(x)/b(x), & \text{if } x \in C, \\ \sup_{y \in C} \bar{\lim}_{t \downarrow 0} f(x + t(y - x)), & \text{if } x \in \text{bd } C, \\ \infty, & \text{if } x \notin \text{cl } C \end{cases}$$

888 *where $\text{cl } C$ denotes the closure, $\text{bd } C$ denotes the boundary of C , and \bar{D}^a and \bar{D}^b denote the domains of a
889 and b , respectively. Then assumptions **A1-A4** hold with $\text{int } \bar{D} = C$.*

890 B An Efficient Line-Search Subgradient Method

891 This section presents an efficient line-search subgradient method for minimizing $\kappa(v) := \kappa(\mathcal{V}(v))$ where $\mathcal{V}(v)$
892 is as in (2.15).

893 Let $w = e + Vv \in \mathbb{R}_{++}^n$ correspond to the current iterate v . Also, let σ be a small scalar between 0 and
894 1 and let $\Delta w = V\Delta v$ where $\Delta v = -\nabla \kappa(v)$ where $\kappa(v)$ is as in (2.16). From Lemma 2.3, we can use a ratio
895 test and guarantee positive definiteness from $w + tV\Delta v = w + t\Delta w > \sigma w$, or equivalently $t\Delta w > (\sigma - 1)w$.

896 Therefore, we conclude that the maximum step with safeguarding, so as not to get too close to the boundary,
897 is

$$t_{\max}(\sigma) \cong t_{\max} := \min \left\{ \frac{(\sigma - 1)w_i}{\Delta w_i} : \Delta w_i < 0 \right\} > 0. \quad (\text{B.1})$$

898 Recall the gradient $\nabla\kappa(v) = \nabla(\kappa \circ \mathcal{V}(v))$ in Lemma 2.14. Our line search in Algorithm B.1 maintains:

899 **LineSearch B.1** (for Algorithm B.1). *Backtrack and maintain: (i) $t \leq t_{\max}$ from (B.1); (ii) monotonic*
900 *nonincrease of the objective $\kappa(v + t\Delta v)$; and nonpositivity of the directional derivative $\nabla\kappa(v + t\Delta v)^T \Delta v \leq 0$.*

901 Algorithm B.1 is presented below.

Algorithm B.1 An Efficient Line-Search Subgradient Method for Minimizing $\kappa(v)$.

Inputs: A symmetric positive definite matrix $A > 0$, a max iteration count `mainmaxiter`, stopping tolerances
`maintoler` > 0 and `linesrchtoler` > 0 , and a small scalar $0 < \sigma \ll 1$.

set $k \leftarrow 0$, $v_1 = 0 \in \mathbb{R}^{n-1}$, $w_1 = e_n \in \mathbb{R}^n$, and $t_k \leftarrow \infty$;

compute a min eigenpair (λ_n^1, x_n^1) and max eigenpair (λ_1^1, x_1^1) of A ;

compute $\kappa(v_1)$ and $\nabla\kappa(v_1)$ according to (2.16) and (2.19), respectively;

set `relnormg` = $\|\nabla\kappa(v_1)\|^2 / (1 + \kappa(v_1)^2)$;

1: **while** `relnormg` $>$ `maintoler` & $t_k >$ `linesrchtoler` & $k \leq$ `mainmaxiter` **do** (main outer loop)

2: set $k \leftarrow k + 1$;

3: set $\Delta v_k \leftarrow -\nabla\kappa(v_k)$;

4: use (B.1) with $\Delta w := V\Delta v_k$ and $w := w_k$ to evaluate $t_{\max}(\sigma)$;

5: perform LineSearch B.1 to find t_k ;

6: set $v_{k+1} = v_k + t_k \Delta v_k$ and $w_{k+1} = e_n + Vv_{k+1}$;

7: compute min eigenpair $(\lambda_n^{k+1}, x_n^{k+1})$ and max eigenpair $(\lambda_1^{k+1}, x_1^{k+1})$ of $A \text{Diag}(w_{k+1})$;

8: compute $\kappa(v_{k+1})$ and $\nabla\kappa(v_{k+1})$ according to (2.16) and (2.19), respectively;

9: update `relnormg` = $\|\nabla\kappa(v_{k+1})\|^2 / (1 + \kappa(v_{k+1})^2)$;

10: **end while**(main outer loop)

Output: $\hat{D} := \text{Diag}(w_{k+1})$.

902 In practice, we take `linesearchtol` to be 10^{-7} and `maintoler` to be 10^{-4} .

903 C List of SuiteSparse Matrices Used in Experiments

904 Table 4.1 considers the following matrices:

905 “bcsstk08.mat”, “bcsstk13.mat”, “bcsstk21.mat”, “bcsstk23.mat”, “bcsstk24.mat”

906 “bcsstk26.mat”, “bcsstk28.mat”, “bcsstk34.mat”, “494_bus.mat”, “662_bus.mat”

907 “nasa1824.mat”, “nasa2146.mat”, “nasa2910.mat”, “nos1.mat”, “nos2.mat”, “nos4.mat”

908 “nos5.mat”, “nos7.mat”

909 Table 4.3 considers the following matrices:

910 “Pres_Poisson.mat”, “bcsstk25.mat”, “bcsstm25.mat”, “gyro_m.mat”

911 “gyro.mat”, “bcsstk36.mat”, “wathen100.mat”, “wathen120.mat”, “minsurfo.mat”

912 “gridgena.mat”, “G2_circuit.mat”

913 Table 4.6 considers the following matrices:

914 “08blocks.mat”, “arc130.mat”, “bwm200.mat”, “ck104.mat”, “ex1.mat”

915 “football.mat”, “gre_185.mat”, “gre_216a.mat”, “impcol.a.mat”, “impcol.c.mat”
916 “impcol.e.mat”, “jazz.mat”, “lshp_265.mat”, “olm100.mat”, “polbooks.mat”
917 “rajat11.mat”, “robot.mat”, “rotor1.mat”, “rw136.mat”, “tub100.mat”
918 “utm300.mat”, “west013.mat”

919 Table 4.8 considers the following matrices:

920 “airfold_2d.mat”, “c-24.mat”, “c-36.mat”, “c-39.mat”, “cavity21.mat”
921 “conf5_4x4-18.mat”, “coupled.mat”, “delaunay_n11.mat”, “delaunay_n13.mat”, “epb1.mat”
922 “ex24.mat”, “G34.mat”, “G36.mat”, “G59.mat”, “G63.mat”
923 “G64.mat”, “garon1.mat”, “Hamrle2.mat”, “jan99jac040sc.mat”, “Na5.mat”
924 “t2d_q9.mat”, “tols4000.mat”, “utm5940.mat”, “viscoplastic1.mat”
925 “whitaker3.mat”

Index

- 926 $A/R = AR^{-1}$, 23
 927 $A \in \mathcal{M}^n$ nonsingular, 5
 928 $D = \text{Diag}(d)$, 8
 929 $M := A^T A$, 18
 930 $M = A^T A > 0$, 6
 931 $V = \frac{1}{\sqrt{2}} \begin{bmatrix} I_{n-1} \\ -e_{n-1}^T \end{bmatrix}$, 6
 932 $X \bullet Y$, Hadamard product, 6
 933 $\text{Blkdiag}(\mathcal{B}) \in \mathcal{M}^n$, 23
 934 $\text{Diag}(d) \in \mathbb{S}^n$, 6
 935 $\text{int } \bar{D}$, interior, 13
 936 \mathcal{M}^n , matrices order n , 6
 937 $\Omega := \{d : d \geq \delta e_n\}$, 12
 938 \mathbb{S}^n , symmetric matrices order n , 6
 939 \mathbb{S}_+^n , cone of positive semidefinite matrices, 6
 940 \mathbb{S}_{++}^n , cone of positive definite matrices, 4, 6
 941 $\bar{S}(x) := \{y \in \text{int } \bar{D} : f(y) < f(x)\}$, 13
 942 \bar{D} , domain, 34
 943 $\mathcal{D}_q(d) = \text{Diag}(d)^{1/2} M \text{Diag}(d)^{1/2}$, 7
 944 $\mathcal{D}(d) = M \text{Diag}(d)$, 8
 945 $\mathcal{D}'(d)(\Delta d)$, 10
 946 $\mathcal{V}(v) := M \text{Diag}(e_n + Vv)$, 15
 947 conv , convex hull, 6
 948 $\text{diag}(S) \in \mathbb{R}^n$, 6
 949 $\hat{\Omega}$, 15
 950 $\kappa(M) = \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)}$, 4
 951 $\kappa(d) = \frac{\lambda_{\max}(\mathcal{D}_q(d))}{\lambda_{\min}(\mathcal{D}_q(d))} = \frac{\lambda_{\max}(\mathcal{D}(d))}{\lambda_{\min}(\mathcal{D}(d))}$, 8
 952 $\kappa(d) = \kappa(\mathcal{D}_q(d)) = (\kappa \circ \mathcal{D}_q)(d)$, 7
 953 $\kappa(v) := \frac{\lambda_{\max}(\mathcal{V}(v))}{\lambda_{\min}(\mathcal{V}(v))}$, 15
 954 κ -optimal diagonal preconditioning, 7
 955 $\kappa(A)$, classical condition number, 6
 956 $\lambda_1(B) = \max_i \lambda_i(B)$, 6
 957 $\lambda_n(B) = \min_i \lambda_i(B)$, 6
 958 $\lambda_{\max}(d) = \lambda_{\max}(\mathcal{D}_q(d)) = (\lambda_{\max} \circ \mathcal{D}_q)(d)$, 7
 959 $\lambda_{\max}(d) = \lambda_{\max}(\mathcal{D}_q(d)) = \lambda_{\max}(\mathcal{D}(d))$, 8
 960 $\lambda_{\max} = \lambda_1$, 6
 961 $\lambda_{\min} = \lambda_n$, 6
 962 $\lambda_{\min}(d) = \lambda_{\min}(\mathcal{D}_q(d)) = (\lambda_{\min} \circ \mathcal{D}_q)(d)$, 7
 963 $\lambda_{\min}(d) = \lambda_{\min}(\mathcal{D}_q(d)) = \lambda_{\min}(\mathcal{D}(d))$, 8
 964 $\text{inv}(c, d) = \begin{pmatrix} \text{diag}(\text{Diag}(c)^{-1}) \\ \text{diag}(\text{Diag}(d)^{-1}) \end{pmatrix}$, 20
 965 $\omega(M) = \frac{\text{tr}(M)/n}{\det(M)^{1/n}}$, 4
 966 $\omega_{\mathcal{Q}}(c, d) := \frac{f(c, d)}{g(c, d)}$, 20
 967 $\partial h(x) := \{v : \langle v, y - x \rangle \leq h(y) - h(x), \quad \forall y\}$, 6
 968 e_n^\perp , orthogonal complement, 6
 969 $f(c, d) := \frac{1}{n} \text{tr}(\mathcal{Q}(c, d))$, 20
 970 $f \circ g$, composite function, 6
 971 $g(c, d) := \det(\mathcal{Q}(c, d))^{1/n}$, 20
 972 $\bar{\partial}^\circ f(x) := \{g : \langle g, y - x \rangle < 0, \quad \forall y \in \bar{S}(x)\}$, 13
 973 $\mathcal{Q}(c, d) := A^T \text{Diag}(c) A \text{Diag}(d)$, 20
 974 $\text{cl } C$, 34
 975 **CV**, coefficient of variation, 4
 976 classical condition number, $\kappa(A)$, 6
 977 coefficient of variation, **CV**, 4
 978 coefficient of variation, **CV**, 5
 979 composite function, $f \circ g$, 6
 980 cone of positive definite matrices, \mathbb{S}_{++}^n , 6
 981 cone of positive semidefinite matrices, \mathbb{S}_+^n , 6
 982 convex hull, conv , 6
 983 domain \bar{D} , 34
 984 Hadamard product, $X \bullet Y$, 6
 985 interior, $\text{int } \bar{D}$, 13
 986 Jacobi preconditioner, 5
 987 matrices order n , \mathcal{M}^n , 6
 988 orthogonal complement, e_n^\perp , 6
 989 SDP, semidefinite programming, 4
 990 semidefinite programming, SDP, 4
 991 strict complementarity condition, 11
 992 symmetric matrices order n , \mathbb{S}^n , 6

References

- [1] M. BENZI, *Preconditioning techniques for large linear systems: a survey*, J. Comput. Phys., 182 (2002), pp. 418–477. 5, 7
- [2] A. BOCK AND M. ANDERSEN, *Connecting Kaporin’s condition number and the Bregman log determinant divergence*, tech. rep., 2025. 5
- [3] R. BYRD AND R. NOCEDAL, *A tool for the analysis of quasi-Newton methods with application to unconstrained minimization*, SIAM J. Numer. Anal., 26 (1989), pp. 727–739. 5
- [4] R. BYRD, R. NOCEDAL, AND Y. YUAN, *Global convergence of a class of quasi-Newton methods on convex problems*, SIAM J. Numer. Anal., 24 (1987), pp. 1171–1191. 5
- [5] G. CIARAMELLA AND M. J. GANDER, *Iterative methods and preconditioners for systems of linear equations*, vol. 19 of Fundamentals of Algorithms, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, [2022] ©2022. 5, 7
- [6] T. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Transactions on Mathematical Software (TOMS), 38 (2011), pp. 1–25. 27
- [7] J. E. DENNIS, JR. AND H. WOLKOWICZ, *Sizing and least-change secant methods*, SIAM J. Numer. Anal., 30 (1993), pp. 1291–1314. 3, 5, 18, 19
- [8] X. DOAN AND H. WOLKOWICZ, *Numerical computations and the ω -condition number*, Tech. Rep. CORR 2011-03, University of Waterloo, Waterloo, Ontario, 2011. www.math.uwaterloo.ca/~hwolkowi/henry/reports/ONE.pdf. 5, 18, 23, 24
- [9] X. V. DOAN, S. KRUK, AND H. WOLKOWICZ, *A robust algorithm for semidefinite programming*, Optim. Methods Softw., 27 (2012), pp. 667–693. 24
- [10] A. FRANCESCHINI, M. FERRONATO, C. JANNA, V. A. PALUDETTO MAGRI, ET AL., *Recent advancements in preconditioning techniques for large size linear systems suited for high performance computing*, Dolomites Research Notes on Approximation, 11 (2018), pp. 11–22. 3
- [11] W. GAO, Y.-C. CHU, Y. YE, AND M. UDELL, *Gradient methods with online scaling*, tech. rep., 2024. arXiv, 2411.01803. 3, 7
- [12] W. GAO, Z. QU, M. UDELL, AND Y. YE, *Scalable approximate optimal diagonal preconditioning*, tech. rep., 2023. arXiv, 2312.15594. 3, 4, 5, 7, 26, 27, 28, 33
- [13] J. A. GEORGE AND J. W.-H. LIU, *Computer Solution of Large Sparse Positive Definite Systems*, Prentice-Hall, Englewood Cliffs, NJ, 1981. 3
- [14] A. GREENBAUM, *Iterative methods for solving linear systems*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997. 5
- [15] Y. HU, J. LI, AND C. K. W. YU, *Convergence rates of subgradient methods for quasi-convex optimization problems*, Comput. Optim. Appl., 77 (2020), pp. 183–212. 17
- [16] C. G. J. JACOBI, *Ueber eine neue auflösungsart der bei der methode der kleinsten quadrate vorkommenden lineären gleichungen*, Astronomische Nachrichten, 22 (1845), p. 297–306. 18
- [17] W. L. JUNG, D. TORREGROSA-BELÉN, AND H. WOLKOWICZ, *The ω -condition number: applications to preconditioning and low rank generalized Jacobian updating*, Comput. Optim. Appl., 91 (2025), pp. 235–282. 3, 5, 6, 18, 19, 23

- 1032 [18] I. KAPORIN, *New convergence results and preconditioning strategies for the conjugate gradient method*,
1033 Numer. Linear Algebra Appl., 1 (1994), pp. 179–210. [5](#)
- 1034 [19] K. KIWIEL, *Convergence and efficiency of subgradient methods for quasiconvex minimization*, Mathe-
1035 matical programming, 90 (2001), pp. 1–25. [2](#), [12](#), [13](#), [14](#), [17](#), [34](#)
- 1036 [20] P. A. KNIGHT, *The Sinkhorn–Knopp algorithm: convergence and applications*, SIAM J. Matrix Anal.
1037 Appl., 30 (2008), pp. 261–275. [18](#), [22](#), [23](#)
- 1038 [21] P. A. KNIGHT, D. RUIZ, AND B. UÇAR, *A symmetry preserving algorithm for matrix scaling*, SIAM
1039 J. Matrix Anal. Appl., 35 (2014), pp. 931–955. [22](#)
- 1040 [22] A. V. KNYAZEV, *Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned*
1041 *conjugate gradient method*, SIAM journal on scientific computing, 23 (2001), pp. 517–541. [5](#)
- 1042 [23] J. LEE, *Smooth manifolds*, in Introduction to smooth manifolds, Springer, 2003, pp. 1–29. [17](#)
- 1043 [24] Z. QU, W. GAO, O. HINDER, Y. YE, AND Z. ZHOU, *Optimal diagonal preconditioning*, Operations
1044 Research, 73 (2024), pp. 1479–1495. [2](#), [3](#), [4](#), [5](#), [7](#), [26](#), [27](#), [28](#), [29](#), [30](#), [31](#), [33](#)
- 1045 [25] Y. SAAD, *Iterative methods for sparse linear systems*, Society for Industrial and Applied Mathematics,
1046 Philadelphia, PA, second ed., 2003. [5](#)
- 1047 [26] ———, *Iterative methods for linear systems of equations: a brief historical journey*, 754 ([2020] ©2020),
1048 pp. 197–215. [18](#)
- 1049 [27] M. SCETBON, C. MA, W. GONG, AND E. MEEDS, *Gradient multi-normalization for stateless and*
1050 *scalable LLM training*, arXiv preprint arXiv:2502.06742, (2025). [18](#), [22](#)
- 1051 [28] G. W. SOULES, *The rate of convergence of Sinkhorn balancing*, in Proceedings of the First Conference
1052 of the International Linear Algebra Society (Provo, UT, 1989), vol. 150, 1991, pp. 3–40. [18](#), [22](#)
- 1053 [29] L. TUNÇEL AND H. WOLKOWICZ, *Strengthened existence and uniqueness conditions for search directions*
1054 *in semidefinite programming*, Linear Algebra Appl., 400 (2005), pp. 31–60. [8](#)
- 1055 [30] N. VAN DER AA, H. TER MORSCHÉ, AND R. MATTHEIJ, *Computation of eigenvalue and eigenvector*
1056 *derivatives for a general complex-valued eigensystem*, Electron. J. Linear Algebra, 16 (2007), pp. 300–
1057 314. [10](#)
- 1058 [31] A. WATHEN, *Preconditioning*, Acta Numer., 24 (2015), pp. 329–376. [5](#), [7](#)
- 1059 [32] H. WOLKOWICZ AND Q. ZHAO, *An all-inclusive efficient region of updates for least change secant*
1060 *methods*, SIAM J. Optim., 5 (1995), pp. 172–191. [5](#)
- 1061 [33] W. E. YOUNG AND R. H. TRENT, *Geometric mean approximations of individual security and portfolio*
1062 *performance*, The Journal of Financial and Quantitative Analysis, 4 (1969), pp. 179–199. [4](#), [5](#)
- 1063 [34] Q. ZHAO, *Measures for least change secant methods*, Master’s thesis, University of Waterloo, 1993. [5](#)