

OPTIMAL DIAGONAL PRECONDITIONING BEYOND WORST-CASE CONDITIONING: THEORY AND PRACTICE OF OMEGA SCALING *

SAEED GHADIMI[†], WOOSUK L. JUNG[‡], ARNESH SUJANANI[‡], DAVID TORREGROSA-BELÉN[§], AND HENRY WOLKOWICZ[‡]

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

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

For the κ -optimal preconditioning problem, we derive an affine-based pseudoconvex reformulation with three key advantages: all stationary points are global minima, subgradients are inexpensive to compute, and the optimization variable is an n -dimensional vector rather than an $n \times n$ matrix as in semidefinite programming (SDP) approaches. We develop a simple and highly efficient subgradient method, with convergence guarantees, for solving this pseudoconvex formulation. Numerical results indicate that our approach is substantially more scalable, efficient, and accurate than existing SDP-based methods, often achieving dramatic speedups.

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

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

1. Introduction. Preconditioning plays a central role in accelerating iterative methods for large-scale linear systems, e.g., [10–12, 25]. All of these works focus on studying the classical κ -condition number of a matrix A (the ratio of largest to smallest singular values). On the other hand, another line of research has focused on studying the ω -condition number, which is the ratio of the arithmetic and geometric mean of singular values, [7, 17]. Although κ -optimal preconditioners minimize the worst-case condition number, our results show that ω -optimal preconditioners are cheaper to compute and are more predictive of PCG/LSQR performance.

Our work is divided into three parts. (i): First, motivated by the recent work in [12, 25] that evaluates κ -optimal diagonal preconditioners using convex semidefinite programming relaxations, we provide a different approach based on using the nonconvex formulation with κ . In particular, we provide a modified model that exploits the invariance of eigenvalues for a product of matrices and solves an essentially unconstrained pseudoconvex minimization problem. Our approach is significantly more efficient and robust, often achieving speedups of over 100x compared with the methods in [12, 25]. (ii): Second, we consider properties of the ω -condition number. In particular, we illustrate how to obtain explicit formulae for ω -optimal diagonal and block-diagonal preconditioners. In both settings, we obtain either known or new preconditioners, including the Jacobi preconditioner, column and row normalization preconditioners, and those obtained via matrix balancing algorithms such as Sinkhorn–Knopp. Compared to κ -optimal preconditioners, the ω -optimal counterparts are substantially cheaper to compute while retaining strong practical effectiveness. (iii): Finally, in our computational experiments we demonstrate the robustness and efficiency of our algorithms for finding κ -optimal preconditioners. We also demonstrate major advantages of using ω -optimal preconditioners for iteratively solving linear systems using preconditioned conjugate gradient (PCG) for positive definite systems and least squares method (LSQR) for general invertible systems. In general, we find that reductions in the ω -condition number are more strongly correlated with improvements in LSQR and PCG iteration counts

*EMAILS RESP.: SGHADIMI@UWATERLOO.CA, W2JUNG@UWATERLOO.CA, A3SUJANA@UWATERLOO.CA, DAVID.TORREGROSA@UA.ES, HWOLKOWICZ@UWATERLOO.CA

[†]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

51 than reductions in the κ -condition number.

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

57 The κ -condition number and ω -condition number for the positive definite case $M > 0$ are, respectively,
 58 given by ratios of eigenvalues

$$59 \quad (1.1) \quad \kappa(M) = \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)}; \quad \omega(M) = \frac{\text{tr}(M)/n}{\det(M)^{1/n}} = \frac{\sum_i \lambda_i(M)/n}{\prod_i (\lambda_i(M))^{1/n}}.$$

60 For our purposes we extend these definitions in (1.1) using eigenvalues, rather than singular values, to
 61 nonsymmetric matrices with real positive eigenvalues allowing for simplified models for minimizing κ . Indeed,
 62 κ and ω can be considered as *worst-case* and *average-case* condition numbers, respectively. In fact, it is known
 63 that the ratio of the arithmetic to geometric mean is approximately the *coefficient of variation*, \mathbf{CV} ; namely,
 64 the ratio between the standard deviation and the mean of a distribution; when this quantity is close to zero.¹
 65 And, one of the big advantages of ω over κ is that finding explicit formulae for optimal preconditioners by
 66 minimizing ω with special structure is often possible due to the analyticity and simplicity of differentiation of
 67 trace and determinant.

68 **1.1. Main Contributions.** The main contributions of this paper consist of both theoretical and
 69 numerical analyses of the κ and ω -condition numbers. First, we provide an affine based pseudoconvex
 70 reformulation of the κ -optimal diagonal preconditioning problem that allows for an elegant characterization
 71 of its optimality conditions. There are three advantages of this reformulation: all its stationary points
 72 are global minima, its subgradients are inexpensive to compute, and its optimization variable is just an
 73 $n \times 1$ vector rather than a $n \times n$ matrix as in *semidefinite programming*, *SDP* [25]. Our extensive numerical
 74 experiments show that subgradient methods based on our reformulation can effectively converge to a κ -optimal
 75 diagonal preconditioner more efficiently in time and accuracy than current SDP-based approaches in the
 76 literature [12, 25].

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

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

88 We now continue with the organization of the paper and more details on the main contributions.
 89 In Section 2, we present several formulations for minimizing κ along with the derivatives and optimality
 90 conditions. The formulations are presented to take advantage of the affine approach and then avoid
 91 the positive homogeneity of the problem in order to improve stability of the problem. In particular, we
 92 include Theorem 2.7 that presents a characterization of a κ -optimal diagonally preconditioned A that is
 93 based on the *largest/smallest* eigenpairs. This leads to an efficient subgradient algorithm and we include
 94 convergence results. In Section 3, we characterize ω -optimal preconditioners with special structures. A
 95 characterization for a ω -optimal diagonal preconditioner is simple as it corresponds to the classical Jacobi
 96 preconditioner for symmetric positive definite matrices and the row/column normalization preconditioner for
 97 general nonsingular matrices. We also show that the matrix balancing algorithm, Sinkhorn–Knopp, linearly
 98 converges to a stationary point of the two-sided ω -optimal diagonal preconditioning problem. Finally, we
 99 include results on extending from diagonal to block diagonal preconditioning in Subsection 3.3. We show

¹The Young–Trent formula [34] establishes $GM/AM \approx 1 - \frac{1}{2}\mathbf{CV}^2$. Moreover, Taylor’s approximation of $1/(1+x)$ around zero yields $AM/GM \approx 1 + \frac{1}{2}\mathbf{CV}^2$.

100 that the right-sided ω -optimal block-diagonal preconditioner for tall full-rank matrices A comes from taking
 101 Q-less QR decompositions of the blocks of A . We also give a characterization for the left-sided ω -optimal
 102 block-diagonal preconditioner and show that when A is square that it is related to applying the right-sided
 103 optimal preconditioner to A^T .

104 Extensive numerical tests are then conducted in [Section 4](#). We demonstrate that our subgradient
 105 methods are more scalable and efficient at computing κ -optimal diagonal preconditioners than existing
 106 SDP-based approaches in the literature [\[12,25\]](#). Moreover, the resulting preconditioners yield more substantial
 107 improvements in PCG performance for solving linear systems compared to those produced by [\[12,25\]](#).

108 We further compare optimal diagonal preconditioning strategies based on the condition numbers κ and
 109 ω . In particular, we observe that the Sinkhorn–Knopp algorithm, which converges to a stationary point of
 110 the two-sided ω -optimal preconditioning problem, produces preconditioners that significantly outperform
 111 the two-sided κ -optimal method of [\[25\]](#) in terms of LSQR iteration counts for minimizing $\|b - Ax\|$. While
 112 Sinkhorn–Knopp dramatically reduces ω , the method of [\[25\]](#) achieves a stronger reduction in κ . This is
 113 notable as despite this difference, Sinkhorn–Knopp yields better LSQR performance on the majority of
 114 instances while also being computationally cheaper.

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

119 **1.2. Background; Preliminaries.** Given a linear system $Ax = b$ with $A \in \mathcal{M}^n, n \times n$ matrix,
 120 nonsingular, and $A \approx P_1 P_2$, preconditioning effectively solves the given system by solving the following
 121 system for y :

$$122 \quad (1.2) \quad P_1^{-1} A P_2^{-1} y = P_1^{-1} b \text{ for } y, \quad x = P_2^{-1} y.$$

123 It is assumed that the solutions with P_1, P_2 are inexpensive.² Surveys are given in e.g., [\[1, 5, 26, 32\]](#) and the
 124 references therein. There is no matrix-matrix multiplication in preconditioning algorithms, such as PCG, as
 125 they do not form the products explicitly; and only matrix-vector multiplications/divisions are performed.
 126 For example, the well known *Jacobi preconditioner* uses the diagonal $d = \text{diag}(A)$, if it is not zero, and the
 127 diagonal matrices $P_2 = \text{Diag}(d), P_1 = I$.

128 A discussion and references on the two condition numbers κ, ω is given recently in [\[17\]](#). The matrix
 129 nearness problem $\min_P \|I - AP\|_F$ is often used to find preconditioners P . This also arises in the derivation
 130 of quasi-Newton methods. The ω -condition number is introduced in [\[7\]](#) as a measure for nearness to the
 131 identity using $\min_{P>0} \omega(MP)$ for finding optimal quasi-Newton updates; see also [\[8, 33, 35\]](#). It is related to
 132 the measure $\text{tr} A - \log \det A$ used in the convergence proofs in [\[3, 4\]](#). Then, Kaporin [\[18\]](#) used ω to derive new
 133 conjugate gradient convergence rate estimates and guarantees. More recently [\[2\]](#) presents further relationships
 134 and convergence analyses. As mentioned above, it is known that the ratio of AM to GM is approximately the
 135 *coefficient of variation*, $CV, AM/GM \approx 1 + \frac{1}{2} CV^2$ [\[34\]](#). When the data is log-normal, then we get the exact
 136 relation $AM/GM = \sqrt{1 + CV^2}$. Thus by reducing CV, ω promotes clustering of eigenvalues/singular values
 137 which is essential in convergence in PCG type methods, e.g., [\[14, 22\]](#).

138 Note that for invertible $A \in \mathcal{M}^n$, [\[25\]](#) uses the former in $\sqrt{\kappa(A^T A)} = \kappa(A)$ as both involve largest
 139 to smallest singular values. In fact, as emphasized in [\[17\]](#), $\kappa(A) = \kappa(A^{-1})$ and it is the latter that is the
 140 operation for solving linear systems. However, these statements are not true for ω , as the numerator
 141 being the sum of singular values, the nuclear norm, is not equivalent to $\sqrt{\text{tr}(A^T A)}$ and ω is not inverse
 142 invariant. In this paper, for ω , we work with general invertible A and then use $\omega(A^T A)$, i.e., use the ratios of
 143 eigenvalues of $A^T A$.

144 For certain structures in A , in contrast to κ , one can exploit the simple derivatives of tr, \det in ω and
 145 obtain explicit formulae for the optimal preconditioner that in fact coincide with heuristics in the literature,
 146 see [\[17\]](#). Thus no optimization algorithm is needed to obtain the optimal preconditioner. For example, if
 147 we restrict to a diagonal preconditioner $D = \text{diag}(d)$, then the classic Jacobi preconditioner is a multiple of
 148 the ω -optimal diagonal preconditioner. If we restrict to a diagonal structure along with a partial upper
 149 triangular part size k , then the diagonal part again corresponds to the Jacobi preconditioner, while the upper
 150 triangular part $D_{ij}, i \leq j \leq k$, yields exactly the formula for the Cholesky factorization. (see [Proposition 3.10](#))

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

151 for more details). Both these observations highlight the close connections ω -optimal preconditioners have
 152 with classical heuristic preconditioners.

153 **1.2.1. Notation.** We work in *real* Euclidean vector spaces. We let \mathbb{S}^n denote the space of symmetric
 154 matrices with the trace inner product and corresponding Frobenius norm; $\mathbb{S}_+^n, \mathbb{S}_{++}^n$ are the cones of positive
 155 semidefinite, and definite, symmetric matrices of order n , respectively; denoted $S \geq 0, S > 0$, respectively.
 156 We let \mathcal{M}^n denote the space of square matrices of order n , also equipped with the trace inner product and
 157 Frobenius norm. We only work with real matrices in this paper. Throughout, we let $A \in \mathcal{M}^n$ be a nonsingular
 158 matrix. Hence, $M = A^T A > 0$.

159 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
 160 of largest to smallest singular values.

161 For $B \in \mathcal{M}^n$ with real eigenvalues we let

$$162 \quad \lambda_i = \lambda_i(B) : \quad \lambda_{\max} = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{\min} = \lambda_n,$$

163 be the eigenvalues in nonincreasing order, and we define the functions: $\lambda_1(B) = \max_i \lambda_i(B)$; $\lambda_n(B) =$
 164 $\min_i \lambda_i(B)$. For our purposes, and with abuse of notation, we define the κ -condition number for *matrices*
 165 *with real positive eigenvalues*, B not necessarily symmetric, as $\kappa(B) = \frac{\lambda_{\max}(B)}{\lambda_{\min}(B)}$. We note that the product
 166 of two positive definite matrices has real positive eigenvalues, though it is not necessarily normal and not
 167 necessarily diagonalizable. We note that in this case we also have $\omega(B) = \frac{\sum_{i=1}^n \lambda_i(B)/n}{\prod_{i=1}^n \lambda_i(B)^{\frac{1}{n}}}$. We let $I_n \in \mathcal{M}^n$ be
 168 the identity matrix and let $e_n \in \mathbb{R}^n$ be the vector of ones. We often use the following matrix for a basis for
 169 the orthogonal complement e_n^\perp ,

$$170 \quad (1.3) \quad V = \frac{1}{\sqrt{2}} \begin{bmatrix} I_{n-1} \\ -e_{n-1}^T \end{bmatrix}, \quad \text{range}(V) = e_n^\perp.$$

171 We use $X \bullet Y$ to denote the *Hadamard product* (elementwise product) of two (compatible) matrices. For a
 172 vector $d \in \mathbb{R}^n$, $\text{Diag}(d) \in \mathbb{S}^n$ denotes the diagonal matrix formed using d . The adjoint linear transformation
 173 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
 174 the composite function.

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

$$176 \quad (1.4) \quad \partial h(x) := \{v : \langle v, y - x \rangle \leq h(y) - h(x), \quad \forall y\}.$$

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

$$178 \quad (1.5) \quad \partial_C h(x) = \text{conv} \{s : \exists x^k \rightarrow x, \nabla h(x^k) \text{ exists, and } \nabla h(x^k) \rightarrow s\},$$

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

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

183 **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,*
 184 *concave, respectively.*

185 *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
 186 is a maximum of linear (so convex) functions in B , and therefore is convex. We get the concavity part
 187 similarly from $\lambda_{\min}(B) = \min\{x^T B x : \|x\| = 1\}$.

188 Now let

$$189 \quad M \in \mathbb{S}_{++}^n, d \in \mathbb{R}_{++}^n, D = \text{Diag}(d).$$

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

$$193 \quad (1.6) \quad \lambda_i(DMD) = \lambda_i(MDD), \forall i, \quad \kappa(DMD) = \kappa(MDD), \quad \omega(DMD) = \omega(MDD).$$

194 **2. κ -Optimal Diagonal Preconditioning.** As mentioned above, preconditioning (scaling) is important
 195 in many areas of mathematics such as optimization and numerical linear algebra. (See e.g., the surveys [1, 5, 32]).
 196 In this section we concentrate on finding κ -optimal diagonal preconditioners D for positive definite linear

197 systems $Mx = b$, where $M = A^T A$ and $A \in \mathcal{M}^n$ is nonsingular. More specifically, we aim to find D that
 198 minimizes $\kappa((AD^{1/2})^T AD^{1/2}) = \kappa(D^{1/2} M D^{1/2})$.³

199 We begin in [Subsection 2.1](#) with the eigenvalue relation in [Equation \(1.6\)](#). This allows us to work with the
 200 affine mapping MD in order to minimize the classical nonlinear in d formulation $\kappa(D^{1/2} M D^{1/2})$. Moreover, we
 201 provide explicit gradients and optimality conditions in [Theorem 2.7](#). Note that [\[12\]](#) also considers minimizing
 202 $\kappa(D^{1/2} M D^{1/2})$, although they aim to find a diagonal preconditioner that is a convex combination of a small
 203 basis or list of diagonal preconditioners. Working with this affine mapping allows us to get expressions for
 204 derivatives for the composite functions which are used to design very efficient and scalable algorithms. In
 205 addition, our composite functions avoid the positive homogeneity in κ, ω thus improving stability.

206 **2.1. Optimality Conditions for κ Diagonal Preconditioning.** Let $M \in \mathbb{S}_{++}^n$. We denote the
 207 quadratic type scaling as follows:

$$208 \quad \mathcal{D} : \mathbb{R}_{++}^n \rightarrow \mathbb{S}^n, \quad \mathcal{D}(d) = \text{Diag}(d)^{1/2} M \text{Diag}(d)^{1/2}.$$

209 With abuse of notation, we let the argument determine the function,

$$210 \quad \kappa(d) = \kappa(\mathcal{D}(d)) = (\kappa \circ \mathcal{D})(d).$$

211 Similarly,

$$212 \quad \lambda_{\max}(d) = \lambda_{\max}(\mathcal{D}(d)) = (\lambda_{\max} \circ \mathcal{D})(d), \quad \lambda_{\min}(d) = \lambda_{\min}(\mathcal{D}(d)) = (\lambda_{\min} \circ \mathcal{D})(d).$$

213 We use the form that is most appropriate/useful depending on the context. In the literature, e.g., [\[11, 12, 25\]](#),
 214 κ -optimal diagonal preconditioning refers to solving

$$215 \quad (2.1) \quad d^* \in \operatorname{argmin} \left\{ \kappa(\mathcal{D}(d)) = \frac{\lambda_{\max}(\mathcal{D}(d))}{\lambda_{\min}(\mathcal{D}(d))} : d \in \mathbb{R}_{++}^n \right\}.$$

216 Here we restrict to $d \in \mathbb{R}_{++}^n$ as we are taking square roots.

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

$$218 \quad (2.2) \quad \mathring{\mathcal{D}} : \mathbb{R}_{++}^n \rightarrow \mathbb{S}^n, \quad \mathring{\mathcal{D}}(d) = M \text{Diag}(d).$$

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

222 **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 $\mathcal{D}(d)$,*
 223 *i.e., $\langle u_i, u_j \rangle = 0, \forall i \neq j$. Define*

$$224 \quad x_i = D^{-1/2} u_i, \forall i, X = [x_1 \ \dots \ x_n], Y = X^{-T}, \Lambda = \text{Diag}(\lambda).$$

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

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

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

$$231 \quad U^{-T} = U(X^T DX)^{-1}.$$

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

$$234 \quad \begin{aligned} Y &= X^{-T} = D^{1/2} U^{-T} \\ &= D^{1/2} U(X^T DX)^{-1} \\ &= DX(X^T DX)^{-1} \end{aligned}$$

³Note that we are using $\kappa(D^{1/2} M D^{1/2})$ and not the inverse $D^{-1/2}$ or D^{-1} as is customary, see e.g., [\(1.2\)](#). This simplifies the derivatives in the optimization problems. Moreover, the equivalent expression for $\kappa(D^{1/2} M D^{1/2})$ follows from the invariance of the spectrum after commuting matrices.

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., [\[30\]](#).

235 as desired. Now,

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

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

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

241 As in the proof of [Lemma 1.1](#), the latter is a maximum of linear functions in d and therefore is convex. The
242 concavity result follows similarly.

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

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

245 and

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

247 where recall $\mathring{D}(d)$ is as in [Equation \(2.2\)](#).

248 *Proof.* The proof follows immediately from [Lemma 2.1](#).

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

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

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

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

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

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

260 *Remark 2.4.* Note that the two equivalent problems [Equation \(2.1\)](#) and [Equation \(2.3\)](#) are both essentially
261 unconstrained, and the optima are attained and characterized as stationary points in \mathbb{R}_{++}^n . This is not
262 obvious but follows since we have the ratios of convex and concave functions using $\lambda_{\max}, \lambda_{\min}$ and this is
263 over the open cone constraint $d \in \mathbb{R}_{++}^n$. If we add the constraints $d^T e_n = n, d > 0$, or equivalently that
264 $d = e_n + Vv > 0$, with $V \in \mathbb{R}^{n \times (n-1)}$ a matrix whose columns contain a basis for e_n^\perp as done above, then we
265 have a bounded problem in $v \in \mathbb{R}^{n-1}$ and λ_{\max} is bounded above. Bounded below away from 0 follows from
266 applying the greedy solution to the knapsack problem with constraints $\sum_i d_i = n, d \geq 0$. We get

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

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

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

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.

273 Lemmas 2.5 and 2.6 provide the derivative information needed for the optimality conditions.

274 LEMMA 2.5 ([31, (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
 275 $B(\bar{t})$ be diagonalizable with real eigenpairs (ignoring the argument \bar{t})

$$276 \quad BX = X\Lambda, \quad B^T Y = Y\Lambda, \quad \Lambda = \text{Diag}(\lambda).$$

277 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
 278 x_k, y_k , respectively, taken from the corresponding columns of X, Y , respectively. Then the derivative of
 279 eigenvalues is given by

$$280 \quad (2.4) \quad \dot{\lambda}_k(\bar{t}) = y_k^T(\bar{t})\dot{B}(\bar{t})x_k(\bar{t}) = \text{tr } \dot{B}x_k y_k^T. \quad \text{6}$$

281 *Proof.* The proof is in [31, Pg 303].⁷

282 LEMMA 2.6. The Fréchet derivative of the linear transformation \mathring{D} at d acting on Δd is (simply)

$$283 \quad \mathring{D}'(d)(\Delta d) = M \text{Diag}(\Delta d).$$

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

$$286 \quad \nabla \lambda(\mathring{D}(d)) = \frac{\lambda}{x^T D x} x \bullet x.$$

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

288 Suppose as above that we have a linear transformation $\mathring{D}(d)$ with singleton eigenvalue λ . Let $y =$
 289 $Dx/x^T Dx$ be the corresponding left eigenvector given in Lemma 2.1. Then the derivative of the composite
 290 function at d acting on Δd is

$$\begin{aligned} 291 \quad \langle \nabla(\lambda \circ \mathring{D})(d), \Delta d \rangle &= \lambda'(\mathring{D}'(d)(\Delta d)) = y^T(\mathring{D}'(d)(\Delta d))x \\ &= \frac{1}{x^T D x} x^T (DM \text{Diag}(\Delta d))x \\ &= \frac{\lambda}{x^T D x} x^T \text{Diag}(\Delta d)x = \frac{\lambda}{x^T D x} \langle x \bullet x, \Delta d \rangle. \quad \square \end{aligned}$$

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

295 THEOREM 2.7 (Characterization of κ -optimal diagonal preconditioning). Let $M \in \mathbb{S}_{++}^n, d \in \mathbb{R}_{++}^n$ be
 296 given; let $D = \text{Diag}(d)$. Then the gradient of the composite function $\kappa(d)$ is

$$297 \quad \nabla \kappa(d) = \kappa(d) \left(\frac{1}{x_1^T D x_1} (x_1 \bullet x_1) - \frac{1}{x_n^T D x_n} (x_n \bullet x_n) \right)$$

298 where recall that $\kappa(d) := \frac{\lambda_{\max}(\mathring{D}(d))}{\lambda_{\min}(\mathring{D}(d))}$. Hence, M is κ -optimally diagonally preconditioned if, and only if, the
 299 orthonormal eigenvector pair satisfies

$$300 \quad (2.5) \quad x_1 \bullet x_1 = x_n \bullet x_n.$$

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

$$302 \quad (2.6) \quad x_1 = \begin{pmatrix} u \\ v \end{pmatrix}, \quad x_n = \begin{pmatrix} u \\ -v \end{pmatrix}, \quad \|u\| = \|v\|.$$

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

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

$$\begin{aligned} 306 \quad \nabla \kappa(d) &= \frac{1}{\lambda_n(d)^2} \left(\lambda_n(d)\dot{\lambda}_1(d) - \lambda_1(d)\dot{\lambda}_n(d) \right) = \frac{1}{\lambda_n^2} \left(\frac{\lambda_n \lambda_1}{x_1^T D x_1} x_1 \bullet x_1 - \frac{\lambda_1 \lambda_n}{x_n^T D x_n} x_n \bullet x_n \right) \\ &= \kappa(d) \left((x_1 \bullet x_1)/(x_1^T D x_1) - (x_n \bullet x_n)/(x_n^T D x_n) \right). \quad \square \end{aligned}$$

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

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

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

308 From (2.5) we see that at optimality, the squared magnitudes of the top and bottom eigenvectors
 309 coincide (up to permutation). This shows immediately that optimal solutions may be non-unique. This is
 310 specified using Theorem 2.7 to illustrate that the optimal set can be a large set. Proposition 2.8 shows that
 311 nonuniqueness is not pathological but structural: whenever the extreme eigenspaces do not fully determine
 312 the orthogonal complement, entire continua of optimal scalings exist. This explains why normalization is
 313 essential for numerical stability and motivates Subsection 2.3 below.

314
 315 *Proposition 2.8* (Nonuniqueness of κ -optimal diagonal preconditioning). Let $M > 0$ be a κ -optimal
 316 diagonal preconditioned matrix as given in Theorem 2.7 with $x_i, \lambda_i, i = 1, n$, being two, singleton, eigenpairs
 317 that satisfy the optimality conditions in Theorem 2.7. Thus we have

$$318 \quad \lambda_1 > \lambda_2 \geq \dots \geq \lambda_{n-1} > \lambda_n > 0, \quad \lambda = (\lambda_i) \in \mathbb{R}_{++}^n, \quad \Lambda = \text{Diag}(\lambda),$$

319 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
 320 be the basis for e_n^\perp defined in Equation (1.3). Suppose in addition that the *lack of strict complementarity*
 321 *condition* holds, i.e., there exists v such that

$$322 \quad (2.7) \quad 0 \neq v \in \{u \in \mathbb{R}^{n-1} : 0 = Vu \bullet x_1 = Vu \bullet x_n\}.$$

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

$$324 \quad \kappa(M) = \kappa((I + t\Delta D)M(I + t\Delta D)), \quad \forall |t| \leq \epsilon.$$

325 *Proof.* The result follows from expanding

$$326 \quad (I + \epsilon\Delta D)M(I + \epsilon\Delta D) = M + O(\epsilon)$$

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

$$328 \quad \begin{aligned} Vv \bullet x_i = 0, i = 1, n & \iff \text{Diag}(Vv)x_i = 0, i = 1, n \\ & \iff \Delta D x_i = 0, i = 1, n. \end{aligned}$$

329 Specifically, let

$$330 \quad M_2 = \bar{Q} \text{Diag}((\lambda_2, \dots, \lambda_{n-1})^T) \bar{Q}^T, \quad M_1 = M - M_2.$$

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

$$332 \quad DMD = M_1 + DM_2D = M + \epsilon\Delta DM_2 + \epsilon M_2\Delta D + \epsilon^2\Delta DM_2\Delta D. \quad \square$$

333 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
 334 eigenspace of M_2 , i.e., to the $\text{span}(\{x_2, x_3, \dots, x_{n-1}\})$. This means that after the perturbation, with $\epsilon > 0$
 335 sufficiently small, λ_1 and λ_n remain the largest and smallest eigenvalues of DMD respectively. As stated
 336 above, we are using the continuity of eigenvalues and the orthogonality $M_1 M_2 = 0$ that arises using the
 337 spectral theorem and $M_1, M_2 \geq 0$.

338 Note that if condition (2.7) holds, then we get a nonsingleton set in \mathbb{R}^{n-1} of solutions. Moreover, the structure
 339 of V implies that (2.7) restricts the support of the eigenspace $\text{span}(\{x_1, x_n\})$.⁸

340 **2.2. A Projected Subgradient Method for Minimizing $\kappa(d)$.** In Theorem 2.7, we derived the
 341 gradient and optimality conditions for minimizing the pseudoconvex function $\kappa(d)$ in Equation (2.3), over
 342 the open set $d > 0$. Convergence of subgradient methods for pseudoconvex minimization typically requires
 343 optimizing over a closed set. Hence, we consider the following optimization problem:

$$344 \quad (2.8) \quad \bar{d} \in \text{argmin} \left\{ \kappa(d) = \frac{\lambda_{\max}(\hat{D}(d))}{\lambda_{\min}(\hat{D}(d))} : d \in \Omega := \{d : d \geq \delta e_n\} \right\},$$

345 where in this subsection $\kappa(d) := \infty$ for $d \notin \Omega$ and $\delta \in (0, 1)$ is small.

346 We now exploit the pseudoconvex structure that guarantees all stationary points are global minimizers
 347 making subgradient methods natural. We first treat the unconstrained formulation in Algorithm 2.1 and
 348 address homogeneity in Subsection 2.3.

⁸Necessity is still an open problem.

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 .

- 1: set $k \leftarrow 0$;
- 2: set $\text{stopcrit} \leftarrow \infty$;
- 3: set $d_1 = e_n \in \mathbb{R}^n$;
- 4: **while** $\text{stopcrit} > \text{tol}$ **do**
- 5: set $k \leftarrow k + 1$;
- 6: compute min eigenpair (λ_n^k, x_n^k) and max eigenpair (λ_1^k, x_1^k) of $M \text{Diag}(d_k)$;
- 7: compute direction

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

- 8: perform projected gradient step

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

- 9: update stopcrit ;
 - 10: **end while**(main outer loop)
- Output:** $\hat{D} := \text{Diag}(d_{k+1})$.
-

349 *Remark 2.9.* For [Algorithm 2.1](#), we use a popular choice for stepsize sequence $\{t_k\}$ in subgradient methods,
350 $t_k = 1/k$, where k is the iteration index [19]. We note that since [Algorithm 2.1](#) is a subgradient method, in
351 general it is not a descent method. Finally, the algorithm is general in that it does not specify the stopping
352 rule stopcrit . There are several possible rules that the user can employ in practice for updating stopcrit in
353 step 9. For example, at every iteration the user can take stopcrit as $\|s_k\|$ or $|\kappa(d_{k+1}) - \kappa(d_k)|$.

354 We now briefly describe each of the steps of [Algorithm 2.1](#). First, step 6 computes a maximal and minimal
355 eigenpair of $M \text{Diag}(d_k)$ to construct the search direction s_k . It should be noted that the user never has to
356 form the matrix $M \text{Diag}(d_k)$ to compute the eigenpairs. Instead, the user only has to construct subroutines
357 which compute $M \text{Diag}(d_k) * x$ and $(M \text{Diag}(d_k)) \setminus x$ efficiently, where $x \in \mathbb{R}^n$. Second, it will be shown in
358 [Subsection 2.2.1](#) that the direction s_k computed in step 7 lies in the quasisubdifferential (also to be defined
359 in [Subsection 2.2.1](#)) of $\kappa(d_k)$. Finally, using this direction, step 8 performs the projected gradient update
360 $d_{k+1} = \Pi_{\Omega}(d_k - t_k \frac{s_k}{\|s_k\|})$, where $\Pi_{\Omega}(x) := \text{argmin}_{\Omega} \|x - \cdot\|$ denotes the projection onto Ω .

361 **2.2.1. Asymptotic Convergence Analysis.** Our convergence analysis of [Algorithm 2.1](#) relies mostly
362 on [19] where efficient subgradient methods for minimizing quasiconvex functions are presented. Under various
363 assumptions in addition to quasiconvexity, the author establishes asymptotic convergence of subgradient
364 methods with decaying stepsizes. Since our function $\kappa(d)$ is pseudoconvex, it is quasiconvex. For a more
365 detailed discussion related to the assumptions of Kiwiel [19] and how our set-up satisfies these assumptions,
366 the reader is referred to [Appendix A](#). In the remaining part of this subsection, we show that [Algorithm 2.1](#)
367 asymptotically converges by showing that it is an instance of the subgradient framework proposed by Kiwiel.

368 First, we need to define a special subdifferential called the quasisubdifferential for a quasiconvex function f .
369 Define the strict sublevel set or inner slice of f as:

$$370 \quad (2.11) \quad \bar{S}(x) := \{y \in \text{int } \bar{D} : f(y) < f(x)\},$$

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

$$373 \quad (2.12) \quad \bar{\partial}^{\circ} f(x) := \{g : \langle g, y - x \rangle < 0, \quad \forall y \in \bar{S}(x)\}.$$

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

$$376 \quad x_{k+1} := \Pi_X(x_k - t_k \hat{g}_k), \quad \hat{g}_k := g_k / \|g_k\|, \quad g_k \in \bar{\partial}^{\circ} f(x_k), \quad k = 1, 2, \dots, \quad x_1 \in X,$$

377 where $t_k > 0$ are the stepsizes.

378 Hence, we need to characterize vectors that lie in the quasisubdifferential of $\kappa(d)$ to show that [Algorithm 2.1](#)
 379 is an instance of Kiwiel's subgradient framework. The result below presents one characterization of vectors
 380 that lie in the quasisubdifferential of fractional programs like the one in our set-up [Equation \(2.8\)](#).

381 **LEMMA 2.10.** *Suppose $f(x) := a(x)/b(x)$ for all $x \in X$ and $f(x) := \infty$ for $x \notin X$, where $a(x)$ is a convex
 382 function that is positive on X , $b(x)$ is a concave function that is positive on X . Let \bar{D}^b denote the domain of
 383 b . If for $x \in X \cap \bar{D}^b$, $B := 1/b(x)$ and $A := a(x)/b^2(x)$, then*

$$384 \quad B[\partial a(x)] + A[\partial(-b)(x)] \subseteq \bar{\partial}^\circ f(x).$$

385 *Proof.* Let $x \in X \cap \bar{D}^b$ and consider B and A as in the assumptions of the lemma. It follows from the
 386 fact that a and $-b$ are convex functions, B and A are positive scalars, and Fenchel subdifferential calculus
 387 rules that

$$388 \quad B[\partial a(x)] + A[\partial(-b)(x)] = \partial(Ba)(x) + \partial(-Ab)(x) \subseteq \partial(Ba - Ab)(x).$$

389 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
 390 domain of f and $f(y) < f(x)$. It then follows that

$$391 \quad \begin{aligned} \langle g, y - x \rangle &\leq Ba(y) - Ab(y) - Ba(x) + Ab(x) \\ 392 \quad &= \frac{a(y)}{b(x)} - \frac{a(x)b(y)}{b^2(x)} < 0, \end{aligned} \quad \square$$

393 where the first inequality follows from the definition of Fenchel subdifferential, the equality follows from the
 394 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
 395 $b(x)$ are positive. It then follows from the definition of quasisubdifferential in [Equation \(2.12\)](#) that $g \in \bar{\partial}^\circ f(x)$.
 396

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

398 **COROLLARY 2.11.** *Consider $M > 0$ as in [Equation \(2.8\)](#) and $d \in \Omega \cap \text{int } \bar{D}^{\lambda_{\min}}$, where $\bar{D}^{\lambda_{\min}}$ is the
 399 domain of $\lambda_{\min}(\mathring{D}(d))$. Let (λ_1, x_1) and (λ_n, x_n) be a maximal and minimal eigenpair of $\mathring{D}(d)$, respectively.
 400 Then, it holds that*

$$401 \quad \frac{\lambda_1}{\lambda_n} \left(\frac{1}{x_1^T(d \bullet x_1)} (x_1 \bullet x_1) - \frac{1}{x_n^T(d \bullet x_n)} (x_n \bullet x_n) \right) \in \bar{\partial}^\circ(\kappa(d)).$$

402 *Proof.* It follows from [Lemma 2.1](#) that $\lambda_{\max}(\mathring{D}(d))$ and $\lambda_{\min}(\mathring{D}(d))$ are convex and concave functions of
 403 d , respectively. Also, it is easy to see that $\lambda_{\max}(\mathring{D}(d))$ and $\lambda_{\min}(\mathring{D}(d))$ are positive on Ω . Hence, it follows
 404 from [Lemma 2.10](#) that

$$405 \quad (2.13) \quad \frac{1}{\lambda_{\min}(\mathring{D}(d))} \partial \lambda_{\max}(\mathring{D}(d)) + \frac{\lambda_{\max}(\mathring{D}(d))}{\lambda_{\min}^2(\mathring{D}(d))} \partial \left(-\lambda_{\min}(\mathring{D}(d)) \right) \subseteq \bar{\partial}^\circ(\kappa(d))$$

406 holds for $d \in \Omega \cap \text{int } \bar{D}^{\lambda_{\min}}$. Moreover, it follows from [Lemma 2.6](#), the definition of Clarke subdifferential in
 407 [Equation \(1.5\)](#), and the fact that the Fenchel and Clarke subdifferentials coincide for convex functions that
 408 the following inclusions hold.

$$409 \quad (2.14) \quad \frac{\lambda_1}{x_1^T(d \bullet x_1)} (x_1 \bullet x_1) \in \partial \lambda_{\max}(\mathring{D}(d)), \quad -\frac{\lambda_n}{x_n^T(d \bullet x_n)} (x_n \bullet x_n) \in \partial \left(-\lambda_{\min}(\mathring{D}(d)) \right), \quad \square$$

410 where here (λ_1, x_1) and (λ_n, x_n) are maximal and minimal eigenpairs of $\mathring{D}(d)$, respectively. The result then
 411 follows from [Equation \(2.13\)](#) and [Equation \(2.14\)](#).

412 **Remark 2.12.** It follows from [Corollary 2.11](#) that the update rule [Equation \(2.10\)](#) in step 8 is of the form

$$413 \quad d_{k+1} = \Pi_\Omega \left(d_k - t_k \frac{s_k}{\|s_k\|} \right), \quad \text{where } s_k \in \bar{\partial}^\circ(\kappa(d_k)).$$

414 We are now ready to present the main theorem that shows [Algorithm 2.1](#) converges asymptotically.

415 **THEOREM 2.13** (Asymptotic convergence). *Let $\{d_k\}$ be generated by [Algorithm 2.1](#). Let*

$$416 \quad \kappa_* := \min \{ \kappa(d) : d \in \Omega \}; \quad \kappa_*^k = \min \{ \kappa(d_j), j = 1 \dots, k \}.$$

417 *Then*

$$418 \quad \underline{\lim}_{k \rightarrow \infty} \kappa(d_k) = \kappa_*; \quad \text{and } \kappa_*^k \downarrow \kappa_*.$$

419 *Proof.* It follows from the last remark in [Lemma 2.1](#) and the definitions of $\kappa(d)$ and Ω in [Equation \(2.8\)](#),
420 that $\lambda_{\max}(\mathring{D}(d))$ (resp. $\lambda_{\min}(\mathring{D}(d))$) is a convex (resp. concave) function that is positive on Ω . It then follows
421 from this observation, [Lemma A.2](#), and the definition of $\kappa(d)$, that assumptions **A1-A4** in [Appendix A](#) hold
422 for the minimization problem in [Equation \(2.8\)](#). Clearly, also assumption **A5** in [Appendix A](#) also holds since
423 Ω in [Equation \(2.8\)](#) is a closed set and the intersection of Ω with the interior of the domain of $\kappa(d)$ is clearly
424 nonempty. The result of the theorem then immediately follows from this observation, [Remark 2.12](#), the facts
425 that $t_k \rightarrow 0$, and $\sum_{k=1}^{\infty} t_k = \infty$, and [Theorem 1](#) in [\[19\]](#).

426 **2.3. Avoiding Positive Homogeneity to Ensure Well-Posedness in Minimizing $\kappa(d)$.** As
427 mentioned in [Remark 2.4](#), $\kappa(d)$ is a positively homogenous function (of degree zero). Thus there are multiple
428 optimal solutions and our problem is *Hadamard ill-posed*. From a computational stability perspective, it is
429 more efficient to minimize an equivalent smaller dimensional formulation that is not positively homogeneous.
430 With this in mind, we let $M > 0$ and $V \in \mathbb{R}^{n \times (n-1)}$ be a matrix whose columns form a basis for e_n^\perp , where
431 $e_n \in \mathbb{R}^n$ is the vector of all ones. We consider the function

$$432 \quad (2.15) \quad \mathring{V}(v) := M \text{Diag}(e_n + Vv),$$

433 where $v \in \mathbb{R}^{n-1}$. In this subsection, we consider the following formulation

$$434 \quad (2.16) \quad \min \left\{ \kappa(v) := \frac{\lambda_{\max}(\mathring{V}(v))}{\lambda_{\min}(\mathring{V}(v))} : e_n + Vv \geq \hat{\delta} e_n, v \in \mathbb{R}^{n-1} \right\},$$

435 where $\kappa(v) := \infty$ if $e_n + Vv < \hat{\delta} e_n$ and $\hat{\delta} \in (0, 1)$ is a scalar. It is easy to see that with the choice of V in [\(1.3\)](#),
436 the above [Equation \(2.16\)](#) can be rewritten as

$$437 \quad (2.17) \quad \min \left\{ \kappa(v) : v \in \hat{\Omega} \right\},$$

438 where

$$439 \quad (2.18) \quad \hat{\Omega} := \left\{ v \in \mathbb{R}^{n-1} : \sum_{i=1}^{n-1} v_i \leq -\sqrt{2}(\hat{\delta} - 1), \quad v_i \geq \sqrt{2}(\hat{\delta} - 1), \quad i = 1, \dots, n-1 \right\}.$$

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

443 The following lemma will be useful for developing our subgradient method. It gives useful characterizations
444 of the derivatives of $\mathring{V}(v)$ and $\kappa(v)$ in the smooth setting when the maximum and minimum eigenvalues of
445 $\mathring{V}(v)$ have multiplicity one.

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

449 1 *The derivative at v acting on $\Delta v \in \mathbb{R}^{n-1}$ is $\mathring{V}'(v)(\Delta v) := \mathring{V}'(v)(\Delta v) = M \text{Diag}(V \Delta v)$.*

450 2 *The gradient of the composite function $\kappa(v) := \kappa(\mathring{V}(v))$ is*

$$451 \quad (2.19) \quad \nabla \kappa(v) = \kappa(v) V^T \left(\frac{1}{x_1^T (w \bullet x_1)} (x_1 \bullet x_1) - \frac{1}{x_n^T (w \bullet x_n)} (x_n \bullet x_n) \right).$$

452

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

454 2 For a singleton eigenvalue λ with normalized right eigenvector x , we have the left eigenvector $y = Dx / (x^T Dx)$
455 as in [Lemma 2.1](#), which satisfies $x^T y = 1$. Then

$$456 \quad \begin{aligned} \langle \nabla(\lambda \circ \mathring{V})(v), \Delta v \rangle &= y^T \mathring{V}'(v)(\Delta v) x = y^T (M \text{Diag}(V \Delta v)) x \\ &= \langle V^T \text{diag}(\lambda D^{-1} y x^T), \Delta v \rangle \\ &= \frac{1}{x^T D x} \langle V^T \text{diag}(\lambda D^{-1} D x x^T), \Delta v \rangle \\ &= \frac{1}{x^T (w \bullet x)} \langle V^T \text{diag}(\lambda x x^T), \Delta v \rangle = \frac{1}{x^T (w \bullet x)} \langle \lambda V^T (x \bullet x), \Delta v \rangle. \end{aligned}$$

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

$$\begin{aligned}
(2.20) \quad \nabla \kappa(v) &= \frac{\lambda_n \dot{\lambda}_1 - \lambda_1 \dot{\lambda}_n}{\lambda_n^2} = \frac{1}{\lambda_n^2} \left(\frac{\lambda_n}{x_1^T(w \bullet x_1)} \lambda_1 V^T(x_1 \bullet x_1) - \frac{\lambda_1}{x_n^T(w \bullet x_n)} \lambda_n V^T(x_n \bullet x_n) \right) \\
&= \kappa(v) V^T \left(\frac{1}{x_1^T(w \bullet x_1)} (x_1 \bullet x_1) - \frac{1}{x_n^T(w \bullet x_n)} (x_n \bullet x_n) \right). \quad \square
\end{aligned}$$

459 The following remark shows that, as in [Lemma 2.14](#), $\nabla \kappa(v)$ lies in the quasisubdifferential of $\kappa(v)$.

460 *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
461 eigenpairs of $\mathring{\mathcal{V}}(v)$, where $\|x_i\| = 1$. It then follows from [Lemma 2.10](#) and a similar argument as in the proof
462 of [Corollary 2.11](#) that

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

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

465 We now present our subgradient algorithm for minimizing [Equation \(2.17\)](#).

Algorithm 2.2 A Subgradient Method for Minimizing $\kappa(v) := \kappa(\mathring{\mathcal{V}}(v))$ 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 .

- 1: set $k \leftarrow 0$;
- 2: set $\text{stopcrit} \leftarrow \infty$;
- 3: set $v_1 = 0 \in \mathbb{R}^{n-1}$ and $w_1 = e_n \in \mathbb{R}_+^n$;
- 4: **while** $\text{stopcrit} > \text{tol}$ **do**
- 5: set $k \leftarrow k + 1$.
- 6: compute min eigenpair (λ_n^k, x_n^k) and max eigenpair (λ_1^k, x_1^k) of $M \text{Diag}(w_k)$;
- 7: compute direction

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

- 8: perform projected gradient step

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

where $\hat{\Omega}$ is as in [Equation \(2.18\)](#) and set

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

- 9: update stopcrit ;
- 10: **end while**(main outer loop)

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

466 Several remarks about [Algorithm 2.2](#) are now given. First, the matrix V does not need to be stored in
467 memory and is actually not needed as input. All the user needs to input is a subroutine that outputs Vv
468 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
469 from [Remark 2.15](#) that $g_k \in \bar{\partial}^\circ(\kappa(v_k))$, which is defined in [Equation \(2.12\)](#). Finally, the projected gradient
470 step in [Equation \(2.23\)](#) can be performed very efficiently since projecting onto $\hat{\Omega}$ just involves computing the
471 root of a simple equation.

472 **2.3.1. Nonasymptotic Convergence Rate Analysis for Minimizing $\kappa(v)$ on $\hat{\Omega}$.** In this subsection,
473 we show a nonasymptotic convergence rate for [Algorithm 2.2](#) for minimizing $\kappa(v)$ over $\hat{\Omega}$. More specifically,
474 we show that

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

476 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](#)
477 relies mostly on the results from [\[19\]](#) and [\[15\]](#). The function $\kappa(v)$ is pseudoconvex since it is the ratio of a

convex function and a positive concave function. Also, observe that the set $\hat{\Omega}$ is just a box so it is a convex compact set that is easy to project onto. The only other assumption that needs to be verified to be able to show a nonasymptotic convergence rate of [Algorithm 2.2](#) is that the function $\kappa(v)$ is Lipschitz continuous on $\hat{\Omega}$. This result is proved in the following proposition.

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

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

where $L > 0$ and $\hat{\Omega}$ is as in [Equation \(2.18\)](#).

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

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 $v \in \hat{\Omega}$ 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 holds:

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

□

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

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

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

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

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

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 the ratio of a convex and a positive concave function. It follows from [Proposition 2.16](#) that $\kappa(v)$ is Lipschitz continuous on $\hat{\Omega}$. Finally, it follows from [Remark 2.15](#) that the update [Equation \(2.23\)](#) in step 8 of [Algorithm 2.2](#) is 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 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 [Theorem 2.17](#) holds.

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

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

i.e., it is the ratio of the arithmetic to geometric means of the eigenvalues of M . Below, this measure will be used to study the conditioning of systems $Ax = b$, when A is not necessarily positive definite.

This section is divided into three main subsections. First, [Subsection 3.1.1](#) focuses on right-sided ω -optimal diagonal preconditioning, for both symmetric positive definite and full rank matrices. We treat left-sided ω -optimal diagonal preconditioning for invertible matrices in [Subsection 3.1.2](#). Then in [Subsection 3.2](#) we

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, 28, 29]
Right Block-Diagonal	Block QR	Corollary 3.11	[8, 17]
Left Block-Diagonal	Characterization	Theorem 3.12	

521 focus on two-sided ω -optimal diagonal preconditioning. Finally, Subsection 3.3 considers ω -optimal block
522 diagonal preconditioning for non-square matrices. In particular, we focus on preconditioners with triangular
523 blocks.

524 Interestingly, we see that many popular preconditioners used in the literature can be found as being
525 ω -optimal with special structure constraints, as illustrated in Table 3.1.

526 **3.1. Right- and Left-Sided Diagonal.** Suppose that $\text{diag}(A) \neq 0$. The classical Jacobi preconditioner
527 [16], [27, Sect. 10.2] uses:

$$528 \quad (3.1) \quad P^{-1}A = P^{-1}b, \quad P = \text{Diag}(\text{diag}(A)).$$

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

$$532 \quad \min_d \|A - \text{Diag}(d)\|_F, \quad \text{Diag}(d)^{-1}A \approx I.$$

533 We now consider variational problems using ω .

534 **3.1.1. Right-Sided Diagonal.** The following result shows that a ω -optimal right-sided diagonal scaling
535 of an overdetermined full rank matrix A is formed using (\pm) the reciprocal of the column norms of A , see [7].

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

$$538 \quad D \in \text{argmin} \left\{ \omega(D^{1/2} A^T A D^{1/2}) : \text{diag}(D) \in \mathbb{R}_{++}^n \right\}.$$

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

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

$$543 \quad D \in \text{argmin} \left\{ \omega(D^{1/2} M D^{1/2}) : \text{diag}(D) \in \mathbb{R}_{++}^n \right\}.$$

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

545 **3.1.2. Left-Sided Diagonal.** In this section we let $A \in \mathcal{M}^n$ be invertible, and we consider the optimal
546 left-sided diagonal preconditioning problem:

$$547 \quad (3.2) \quad \begin{aligned} & \min \quad \omega((\text{Diag}(c)^{1/2} A)^T \text{Diag}(c)^{1/2} A) \\ & \text{s.t.} \quad c \in \mathbb{R}_{++}^n. \end{aligned}$$

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

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

552 *Proof.* Note that

$$553 \quad \omega((\text{Diag}(c)^{1/2} A)^T \text{Diag}(c)^{1/2} A) = \omega(A^T \text{Diag}(c) A) = \omega(\text{Diag}(c)^{1/2} AA^T \text{Diag}(c)^{1/2}).$$

554 Hence, Proposition 3.1 and the above equivalence implies that

$$555 \quad C = \text{Diag}(\text{diag}(AA^T))^{-1} \iff C \in \text{argmin} \left\{ \omega(C^{1/2} AA^T C^{1/2}) : C = \text{diag}(c) \in \mathbb{S}_{++}^n \right\}$$

$$556 \quad \iff C \in \text{argmin} \left\{ \omega((\text{Diag}(c)^{1/2} A)^T \text{Diag}(c)^{1/2} A) : C = \text{diag}(c) \in \mathbb{S}_{++}^n \right\}.$$

557 **3.2. Two-Sided Diagonal.** We now consider the problem of finding the two-sided ω -optimal diagonal
558 scaling for nonsingular matrices $A \in \mathcal{M}^n$. This subsection is broken up into two smaller parts. The first part
559 derives the optimality conditions for the two-sided ω -optimal scaling problem. The second part presents
560 an iterative matrix balancing scheme that reduces ω at every iteration and whose output is proved to be a
561 stationary point of the two-sided ω -optimal scaling problem.

562 **3.2.1. Optimality Conditions for Two-Sided Problem.** This part presents the formulation of the
563 two-sided ω -optimal diagonal scaling problem and derives its optimality conditions. Namely, we consider

$$564 \quad (3.3) \quad \begin{aligned} \min \quad & \omega((\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})^T (\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})) \\ \text{s.t.} \quad & c, d \in \mathbb{R}_{++}^n, \end{aligned}$$

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

$$566 \quad (3.4) \quad \mathcal{Q}(c, d) := A^T \text{Diag}(c) A \text{Diag}(d), \quad f(c, d) := \frac{1}{n} \text{tr}(\mathcal{Q}(c, d)), \quad g(c, d) := \det(\mathcal{Q}(c, d))^{1/n}.$$

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

$$568 \quad (3.5) \quad \begin{aligned} \min \quad & \omega_{\mathcal{M}}(c, d) := \frac{f(c, d)}{g(c, d)} \\ \text{s.t.} \quad & c, d \in \mathbb{R}_{++}^n. \end{aligned}$$

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

$$570 \quad (3.6) \quad \text{inv}(c, d) = \begin{pmatrix} \text{diag}(\text{Diag}(c)^{-1}) \\ \text{diag}(\text{Diag}(d)^{-1}) \end{pmatrix}.$$

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

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

$$573 \quad (3.7) \quad \nabla \omega_{\mathcal{M}}(c, d) = \frac{1}{ng(c, d)} \begin{pmatrix} \text{diag}(A \text{Diag}(d) A^T) \\ \text{diag}(A^T \text{Diag}(c) A) \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{M}}(c, d) \text{inv}(c, d),$$

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

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

$$577 \quad (3.8) \quad \nabla f(c, d) = \frac{1}{n} \begin{pmatrix} \text{diag}(A \text{Diag}(d) A^T) \\ \text{diag}(A^T \text{Diag}(c) A) \end{pmatrix}.$$

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

$$\begin{aligned} \langle \nabla g(c, d), (\Delta c, \Delta d) \rangle &= \det(A^T A)^{1/n} \langle \nabla \det(CD)^{1/n}, (\Delta c, \Delta d) \rangle \\ &= \det(A^T A)^{1/n} \left\langle \frac{1}{n} \det(CD)^{\frac{1}{n}-1} \text{adj}(CD), (C\Delta D + \Delta CD) \right\rangle \\ 580 &= \frac{1}{n} g(c, d) \left\langle (CD)^{-1}, (C\Delta D + \Delta CD) \right\rangle \\ &= \frac{1}{n} g(c, d) \text{tr}(C^{-1} \Delta C + D^{-1} \Delta D). \end{aligned}$$

581 Hence, it follows from the above equations

$$582 \quad (3.9) \quad \nabla g(c, d) = \frac{1}{n} g(c, d) \begin{pmatrix} \text{diag}(\text{Diag}(c)^{-1}) \\ \text{diag}(\text{Diag}(d)^{-1}) \end{pmatrix} = \frac{1}{n} g(c, d) \text{inv}(c, d).$$

583 It then follows from Equation (3.7), Equation (3.9) and the quotient rule that $\nabla \omega_{\mathcal{M}}(c, d)$ can be computed as

$$\begin{aligned} 584 \quad \nabla \omega_{\mathcal{M}}(c, d) &= \frac{1}{g(c, d)^2} (g(c, d) \nabla f(c, d) - f(c, d) \nabla g(c, d)) = \frac{1}{g(c, d)} \nabla f(c, d) - \frac{\omega_{\mathcal{M}}(c, d)}{g(c, d)} \nabla g(c, d) \\ 585 &= \frac{1}{ng(c, d)} \begin{pmatrix} \text{diag}(A \text{Diag}(d) A^T) \\ \text{diag}(A^T \text{Diag}(c) A) \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{M}}(c, d) \text{inv}(c, d). \quad \square \end{aligned}$$

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

$$587 \quad \nabla \omega_{\mathcal{M}}(c, d) = \frac{1}{n} \begin{bmatrix} 0 & A \bullet A \\ (A \bullet A)^T & 0 \end{bmatrix} \begin{pmatrix} c \\ d \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{M}}(c, d) \text{inv}(c, d).$$

588 We now state our main theorem on a necessary condition for a vector to be globally optimal for
 589 Equation (3.3). This relates to optimal matrix balancing discussed below.

590 **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*
 591 *Equation (3.3). Let $e_{2n} \in \mathbb{R}^{2n}$ be the vector of all ones, $C^* = \text{Diag}(c^*)$, and $D^* = \text{Diag}(d^*)$. Then:*

$$592 \quad (3.10) \quad \begin{pmatrix} \text{diag}(C^* A D^* A^T) \\ \text{diag}(D^* A^T C^* A) \end{pmatrix} = \alpha e_{2n}, \text{ for some } \alpha > 0.$$

593 *Hence, if c^*, d^* is a global optimal solution of Equation (3.3), then the column and row norms of $\sqrt{C^*} A \sqrt{D^*}$*
 594 *must all be equal to the same constant, i.e., the matrix is balanced.*

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

$$597 \quad 0 = \frac{1}{ng(c^*, d^*)} \begin{pmatrix} \text{diag}(A D^* A^T) \\ \text{diag}(A^T C^* A) \end{pmatrix} - \frac{1}{n} \omega_{\mathcal{M}}(c^*, d^*) \text{inv}(c^*, d^*).$$

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

$$600 \quad (3.11) \quad f(c^*, d^*) e_{2n} = \begin{pmatrix} \text{diag}(C^* A D^* A^T) \\ \text{diag}(D^* A^T C^* A) \end{pmatrix}. \quad \square$$

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

606 **Remark 3.7.** Suppose we consider a modification of Equation (3.3), and consider the following optimization
 607 problem:

$$608 \quad (3.12) \quad \begin{aligned} \min \quad & \omega((\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})^T (\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})) \\ \text{s.t.} \quad & c \geq \delta e_n, d \geq \delta e_n, \end{aligned}$$

609 where $\delta \in (0, 1)$. It is easy to see that a global optimal solution c^*, d^* of Equation (3.12) exists. This
 610 follows immediately from the fact that $\omega((\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2})^T (\text{Diag}(c)^{1/2} A \text{Diag}(d)^{1/2}))$ is a positively
 611 homogeneous function of degree 0, so Equation (3.12) is equivalent to minimizing a continuous function over
 612 a compact set.

613 **3.2.2. Square-Root Sinkhorn–Knopp Algorithm for Two Sided Problem.** Let a nonsingular
 614 matrix $A \in \mathcal{M}^n$ be given. We have seen above that the ω -optimal right-sided preconditioner is equivalent to
 615 the column normalization preconditioner. Similarly, the ω -optimal left-sided preconditioner is equivalent to
 616 the row normalization preconditioner. Therefore, we can alternate between ω -optimal right- and left-sided
 617 preconditioners and decrease the value of ω at each iteration. This yields the following computationally
 618 efficient matrix balancing scheme where we alternate between column and row normalization. This algorithm
 619 can be seen as equivalent to the Square-Root Sinkhorn–Knopp algorithm and is related to other balancing
 620 methods in the literature [20, 21, 28, 29].

621 We prove below in Theorem 3.9 that the Square-Root Sinkhorn–Knopp algorithm, namely Algorithm 3.1,
 622 converges and the outputted two-sided diagonal preconditioner satisfies the necessary condition in Equa-
 623 tion (3.10) for global optimality of a two-sided ω -optimal preconditioner. This continues the theme that the
 624 ω -measure 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 .

```

1: set  $\text{stopcrit} \leftarrow \infty$ ;
2: set  $A_0 \leftarrow A$ ;
3: set  $k \leftarrow 0$ ;
4: while  $\text{stopcrit} > \text{tol}$  do
5:   set  $k \leftarrow k + 1$ ;
6:   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  $A_{k-1}$ ;
7:   set  $\tilde{A}_k \leftarrow A_{k-1} \text{Diag}(\sqrt{d_k})$ ;
8:   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  $\tilde{A}_k$ ;
9:   update  $\text{stopcrit}$ ;
10:  set  $A_k \leftarrow \text{Diag}(\sqrt{c_k})\tilde{A}_k$ ;
11: end while(main outer loop)
Output:

```

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

625 There several viable options for the choice of stopcrit in [Algorithm 3.1](#). One natural choice is

626
$$\text{stopcrit} = \max \left(\max_i \|(A_k)_{i,:} - 1\|, \max_j \|(A_k)_{:,j} - 1\| \right).$$

627 It is shown in [Theorem 3.9](#) that the sequence A_k converges to a matrix that is balanced, i.e., one that has
628 column and row norms all equal to one. The remark below also discusses that the ω values of the iterates are
629 monotonically decreasing.

630 *Remark 3.8.* At each iteration, \tilde{A}_k (resp. A_k) is computed using the ω -optimal right-sided (resp. left-sided)
631 preconditioner. Hence, the ω values of the iterates are monotonically decreasing, i.e., $\omega(A_k A_k^T) \leq \omega(\tilde{A}_k \tilde{A}_k^T) \leq$
632 $\omega(A_{k-1} x A_{k-1}^T)$ for all $k \geq 1$. The output of [Algorithm 3.1](#) thus satisfies $\omega((\hat{C}^{1/2} A \hat{D}^{1/2})^T (\hat{C}^{1/2} A \hat{D}^{1/2})) <$
633 $\omega(A^T A)$.

634 [Theorem 3.9](#) now provides a convergence guarantee for [Algorithm 3.1](#) and shows that the necessary
635 condition [Equation \(3.10\)](#) for a two-sided ω -optimal diagonal preconditioner holds in the limit for [Algorithm 3.1](#).

636 **THEOREM 3.9.** *Let A be such that $A \bullet A$ has total support.⁹ Then the sequence A_k converges linearly to*
637 *a matrix $\bar{A} := \bar{C}^{1/2} A \bar{D}^{1/2}$ that has column and row norms all equal to 1. That is,*

638
$$\lim_{k \rightarrow \infty} \text{Diag} \left(\prod_{j=1}^k \sqrt{c_j} \right) A \text{Diag} \left(\prod_{j=1}^k \sqrt{d_j} \right) = \bar{C}^{1/2} A \bar{D}^{1/2}$$

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

641 *Proof.* Let $\tilde{P}_k = \tilde{A}_k \bullet \tilde{A}_k$ and let $P_k = A_k \bullet A_k$. It follows that $\tilde{P}_k = (A_{k-1} \bullet A_{k-1}) \text{Diag}(d_k) =$
642 $P_{k-1} \text{Diag}(d_k)$ 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}
643 so $(d_k)_i = 1./\sum_{l=1}^n (P_{k-1})_{l,i}$. 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
644 $(c_k)_i = 1./\sum_{l=1}^n (\tilde{P}_k)_{i,l}$. Hence, the sequences of iterates \tilde{P}_k and P_k can be viewed as the iterates of Sinkhorn–
645 Knopp Algorithm applied to $A \bullet A$. By [\[20\]](#), we know then that the sequence P_k converges linearly to a matrix
646 \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 just the
647 square roots of the iterates of the Sinkhorn–Knopp algorithm applied to $A \bullet A$. Therefore, it follows from the
648 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 column and
649 row norms all equal to 1. The last conclusion of the theorem follows from this observation and the necessary
650 condition [Equation \(3.10\)](#) in [Theorem 3.6](#).

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 .

651 **3.3. Right-, Left-Sided Block Diagonal.** We now show that the above results extend directly
652 to finding *block diagonal* preconditioners for least squares solutions of full rank possibly overdetermined
653 linear systems $Ax = b, A \in \mathbb{R}^{m \times n}, m \geq n$. The preconditioner is extended from diagonal to block diagonal
654 (block upper-triangular). This relates to sparse QR preconditioning and extends the results for ω -optimal
655 preconditioning with partial Cholesky structure in [17, Theorem 2.7] and the block diagonal ω -optimal
656 preconditioner in [8, Prop. 3 part 3].

657 **3.3.1. Right-Sided Block Diagonal.** We let the linear transformation $B = \text{Blkdiag}(\mathcal{B}) \in \mathcal{M}^n$ denote
658 the block diagonal matrix with blocks formed from the set of square matrices $\mathcal{B} = \{B_i\}_{i=1}^k$ of order
659 $n_i, \sum_{i=1}^k n_i = n$. For our application we restrict the blocks to be of upper-triangular structure denoted
660 $\mathcal{R} = \{R_i\}_{i=1}^k$, and see that the best, with respect to the ω measure, comes from Q-less QR decompositions of
661 the corresponding blocks of A .

662 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
663 structure

$$664 \quad A = [A_1 \quad A_2 \quad \dots \quad A_k], \quad A / \text{Blkdiag}(\mathcal{R}) = [A_1/R_1 \quad A_2/R_2 \quad \dots \quad A_k/R_k],$$

665 where we use the MATLAB notation that $/$ denotes matrix division $A/R = AR^{-1}$. Optimally, we want the
666 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,
667 we can allow a preliminary permutation of the columns $A \leftarrow AP$ so that we can increase the effect of the
668 preconditioner.

669 To extend the diagonal preconditioner results we solve:

$$670 \quad \min_{\mathcal{B}} \{\omega((AB)^T(AB)) : B = \text{Blkdiag}(\mathcal{B})\}.$$

671 The basic result we use follows.

672 *Proposition 3.10* ([8, Prop. 3 part 3], [9, Prop. 2.2]). Let

$$673 \quad A = [A_1 \quad A_2 \quad \dots \quad A_k], \quad A_i \in \mathbb{R}^{m \times n_i}, \forall i,$$

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

$$675 \quad \mathcal{B} := \{B_1, B_2, \dots, B_k\}, \quad B_i \in \mathbb{R}^{n_i \times n_i}, \quad B := \text{Blkdiag}(\mathcal{B}),$$

676 that minimizes the measure ω , i.e.,

$$677 \quad \min \{\omega((AB)^T(AB)) : B = \text{Blkdiag}(\mathcal{B})\},$$

678 is found by satisfying the factorization

$$679 \quad B_i B_i^T = (A_i^T A_i)^{-1}, \quad i = 1, \dots, k.$$

680 **COROLLARY 3.11.** Let A, \mathcal{B}, B be as in *Proposition 3.10*.

681 (i) If we restrict the diagonal blocks B_i to be diagonal themselves, then we get the optimal diagonal precondi-
682 tioner in *Proposition 3.1*.

683 (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
684 QR decomposition of A_i (up to sign of rows of R).

685 **3.3.2. Left-Sided Block Diagonal.** We continue with A full column rank but with left-sided precondi-
686 tioning. We present a nonlinear equation that characterizes optimality.

687 *Theorem 3.12.* Let $A \in \mathbb{R}^{m \times n}$ be full column rank with block structure given by

$$688 \quad A = \begin{bmatrix} \bar{A}_1 \\ \bar{A}_2 \\ \dots \\ \bar{A}_\ell \end{bmatrix}, \quad \bar{A}_j \in \mathbb{R}^{m_j \times n}, \text{rank}(\bar{A}_j) = m_j, \forall j.$$

689 Then the ω -optimal left-sided block diagonal preconditioner

$$690 \quad E := \text{Blkdiag}(E_1, E_2, \dots, E_\ell), \quad E_j \in \mathbb{R}^{m_j \times m_j},$$

691 that minimizes the measure ω , i.e.,

$$692 \quad \min \omega((EA)^T EA),$$

693 is characterized by:

$$694 \quad \bar{A}_j \bar{A}_j^T = \frac{1}{n} \left(\text{tr} \sum_{t=1}^{\ell} (\bar{A}_t^T (E_t^T E_t) \bar{A}_t) \right) \bar{A}_j (A^T (E^T E) A)^{-1} \bar{A}_j^T, \quad \forall j = 1, \dots, \ell.$$

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

$$696 \quad \min \{ \omega(A^T K A) : K = \text{Blkdiag}(K_1, K_2, \dots, K_\ell), K_j \in \mathbb{R}^{m_j \times m_j} \}.$$

697 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

$$698 \quad \omega_{\mathcal{G}}(K) := \frac{\text{tr}(A^T K A)/n}{\det(A^T K A)^{1/n}} = \frac{\sum_{j=1}^{\ell} \text{tr}(\bar{A}_j^T K_j \bar{A}_j)/n}{\det(A^T K A)^{1/n}}.$$

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

$$700 \quad f(K) := \sum_{j=1}^{\ell} \text{tr}(\bar{A}_j^T K_j \bar{A}_j)/n \quad \text{and} \quad g(K) := \det(A^T K A)^{1/n}.$$

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

$$\begin{aligned} \langle \nabla_{K_j} g(K), \Delta K_j \rangle &= \langle \nabla \det(\mathcal{G}(K))^{1/n}, \Delta K_j \rangle \\ &= \left\langle \frac{1}{n} \det(\mathcal{G}(K))^{1/n-1} \text{adj}(\mathcal{G}(K)), \bar{A}_j^T \Delta K_j \bar{A}_j \right\rangle \\ 702 \quad &= \frac{1}{n} g(K) \langle \mathcal{G}(K)^{-1}, \bar{A}_j^T \Delta K_j \bar{A}_j \rangle \\ &= \frac{1}{n} g(K) \langle \bar{A}_j \mathcal{G}(K)^{-1} \bar{A}_j^T, \Delta K_j \rangle. \end{aligned}$$

703 Therefore, we have

$$704 \quad \nabla_{K_j} f(K) = \frac{1}{n} \bar{A}_j \bar{A}_j^T \quad \text{and} \quad \nabla_{K_j} g(K) = \frac{1}{n} g(K) \bar{A}_j (A^T K A)^{-1} \bar{A}_j.$$

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

$$\begin{aligned} 706 \quad \nabla_{K_j} \omega_{\mathcal{G}}(K) &= \frac{1}{g(K)^2} (g(K) \nabla_{K_j} f(K) - f(K) \nabla_{K_j} g(K)) \\ &= \frac{1}{n g(K)} (\bar{A}_j \bar{A}_j^T - f(K) \bar{A}_j (A^T K A)^{-1} \bar{A}_j^T) \end{aligned}$$

707 Therefore

$$708 \quad \nabla \omega_{\mathcal{G}}(K) = 0 \quad \iff \quad \bar{A}_j \bar{A}_j^T = \frac{1}{n} \sum_{t=1}^{\ell} \text{tr}(\bar{A}_t \bar{A}_t^T K_t) \bar{A}_j (A^T K A)^{-1} \bar{A}_j^T, \quad \forall j. \quad \square$$

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

710 Finding an explicit solution for E in [Theorem 3.12](#) is an open question. However, if we assume further
711 that A is a square matrix, the following [Corollary 3.13](#) connects the result from [Corollary 3.11](#) to this left
712 sided preconditioning.

713 *Corollary 3.13.* 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 $A^T =$
714 $[\bar{A}_1^T \bar{A}_2^T \dots \bar{A}_\ell^T]$, $\sum_j n_j = n$. Define $\mathcal{E} := \{E_1, E_2, \dots, E_k\}$ by $E_i := \bar{R}_i^{-T}$, where $\bar{A}_i \bar{A}_i^T := \bar{R}_i^T \bar{R}_i$ is the
715 Cholesky decomposition. Then, \mathcal{E} is the optimal lower triangular block diagonal scaling for

$$716 \quad \min \{ \omega((EA)^T (EA)) : E = \text{Blkdiag}(E_1, E_2, \dots, E_k) \}.$$

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

$$718 \quad \omega((EA)^T (EA)) = \omega(A^T E^T EA) = \omega(EAA^T E^T), \quad \square$$

719 and thus [Corollary 3.11](#) 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.

720 *Remark 3.14.* The result in [Corollary 3.13](#) agrees with our characterization in [Theorem 3.12](#). Indeed,
 721 by 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 $\text{tr}(\bar{A}_t \bar{A}_t^T E_t^T E_t) = \text{tr}(\bar{R}_t^T \bar{R}_t \bar{R}_t^{-1} \bar{R}_t^{-T}) = \text{tr}(I_{n_t}) = n_t$.
 722 Thus,

$$723 \quad \sum_{t=1}^{\ell} \text{tr}(\bar{A}_t \bar{A}_t^T E_t^T E_t) = \sum_{t=1}^{\ell} n_t = n.$$

724 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
 725 that $\mathcal{I}_j A = \bar{A}_j$. Then,

$$\begin{aligned} 726 \quad \frac{1}{n} \left(\text{tr} \sum_{t=1}^{\ell} (\bar{A}_t^T E_t^T E_t \bar{A}_t) \right) \bar{A}_j (A^T E^T E A)^{-1} \bar{A}_j^T &= \bar{A}_j (A^T E^T E A)^{-1} \bar{A}_j^T \\ &= \mathcal{I}_j A A^{-1} E^{-1} E^{-T} A^{-T} A^T \mathcal{I}_j^T \\ &= \mathcal{I}_j E^{-1} E^{-T} \mathcal{I}_j^T \\ &= E_j^{-1} E_j^{-T} = \bar{R}_j^T \bar{R}_j = \bar{A}_j \bar{A}_j^T. \end{aligned}$$

727 **4. Computational Experiments.** Our computational experiments are divided into several parts. In
 728 [Subsection 4.1](#), we compare the computational efficiency of our subgradient algorithms, [Algorithms 2.2](#)
 729 and [B.1](#), with the SDP-based approaches proposed in [\[25\]](#) and [\[12\]](#) for finding an approximate κ -optimal
 730 preconditioner. Both our codes and the code in [\[12\]](#) take as input $M > 0$. Our subgradient algorithms find
 731 a diagonal D that minimizes $\kappa(D^{1/2} M D^{1/2})$ while [\[12\]](#) solves an SDP using MOSEK [\[24\]](#) to find a D that
 732 is κ -optimal over the subspace spanned by a small list of given diagonal preconditioners. Meanwhile, [\[25\]](#)
 733 obtains a diagonal scaling E that minimizes $\kappa((ME)^T (ME))$ by solving a possibly larger SDP. Hence, when
 734 using their code, we input a Cholesky factor B of M , i.e., $M = B^T B$, so that their algorithm finds E that
 735 minimizes $\kappa(EME)$. For every experiment, $\kappa(M)$, $\kappa(D^{1/2} M D^{1/2})$, and $\kappa(EME)$ were each evaluated by
 736 performing a minimum and maximum eigenvalue computation using MATLAB's `eigs` function with 10^{-10}
 737 precision.

738 In [Subsection 4.2](#), we compare PCG using the preconditioner from [Algorithm 2.2](#) with PCG using the
 739 preconditioner of [\[25\]](#) for solving $Mx = b$, where $M > 0$. In [Subsection 4.3](#), we compare the efficiency of our
 740 [Algorithm 3.1](#), with the two-sided κ -optimal preconditioner presented in [\[25\]](#). In particular, we compare the
 741 time to compute each preconditioner as well as the performance of each preconditioner when used inside
 742 LSQR to minimize $\|b - Ax\|$ where $A \in \mathcal{M}^n$ is nonsingular. Finally, in [Subsection 4.4](#), we use the optimality
 743 conditions in [Theorem 2.7](#) in order to construct a matrix M that is κ -optimal with respect to diagonal
 744 preconditioning. We then compare with results after applying the Jacobi scaling, the ω -optimal scaling.

745 All experiments in [Subsections 4.1](#) to [4.3](#) are run on a 2023 Macbook Pro with an 8-core CPU and 128
 746 GB of memory using MATLAB 2024a. The medium experiments in [Subsection 4.4](#) were also run with
 747 this Macbook Pro while the very large instances were run with a large Linux machine from University
 748 of Waterloo with 256 GB of memory. The latest codes are available at this [clickable-link](#) or with URL
 749 <https://github.com/asujanani6/Optimal-Preconditioning>.

750 **4.1. Minimizing κ Efficiently.** In this subsection, [Tables 4.1](#) and [4.2](#) compare the computational
 751 efficiency of [Algorithm 2.2](#) with the algorithms in [\[25\]](#) and [\[12\]](#) for minimizing κ . The tables also compare
 752 these algorithms with our own line-search based subgradient method, namely [Algorithm B.1](#), that was
 753 slightly more computationally efficient than our [Algorithm 2.2](#) although [Algorithm 2.2](#) has theoretical
 754 convergence guarantees. The precise details of [Algorithm B.1](#) can be found in [Appendix B](#). In our experiments,
 755 [Algorithm 2.2](#) uses a tol of 10^{-4} , sets `stopcrit`= $2|\kappa(d_{k+1}) - \kappa(d_k)|/(\kappa(d_{k+1}) + \kappa(d_k))$ at every iteration,
 756 and uses a maxiter of 500 and $\hat{\delta} = 10^{-3}$. The default settings for the algorithms in [\[25\]](#) and [\[12\]](#) are used.
 757 [Algorithm B.1](#) uses a maximum of 80 iterations and uses the same tolerance 10^{-4} .

758 We now give several remarks about the results presented in each table. [Table 4.1](#) compares the four
 759 algorithms on eighteen matrices taken from the SuiteSparse Matrix Collection [\[6\]](#). The list of the SuiteSparse
 760 matrices considered can be found in [Appendix C](#). Columns four to seven of [Table 4.1](#) list the percent reduction
 761 in κ that each method achieved. The percentage reduction is computed as $100 * (\kappa(M) - \kappa(\hat{M})) / (\kappa(M))$, where
 762 \hat{M} is $D^{1/2} M D^{1/2}$ for our codes and [\[12\]](#) or EME for [\[25\]](#). As seen from these columns, both [Algorithms 2.2](#)
 763 and [B.1](#) reduced κ much more significantly than [\[25\]](#) and [\[12\]](#). [Algorithm 2.2](#) reduced κ on average by 79.7%
 764 while [Algorithm B.1](#) reduced κ by 68.5%. On the other hand, [\[25\]](#) reduced κ on average by 30.3% while [\[12\]](#)

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	$\kappa(M)$	[25]	[12]	Algorithm 2.2	Algorithm B.1	[25]	[12]	Algorithm 2.2	Algorithm B.1	[25]/Algorithm 2.2	[25]/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	5.64	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	323.334	12.968	0.120	25.5	2792.3
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	619.752	633.446	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	$\kappa(M)$	[25]	[12]	Algorithm 2.2	Algorithm B.1	[25]	[12]	Algorithm 2.2	Algorithm B.1	[25]/Algorithm 2.2	[25]/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	3.4	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	3.3e+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

765 reduced κ by 21.08%. Columns eight to eleven of Table 4.1 show that Algorithms 2.2 and B.1 were also
766 much faster in reducing κ than [25] and [12]. On average, Algorithm 2.2 was 77.5 (resp. 85.80) times faster
767 than [25] (resp. [12]) while Algorithm B.1 was 356.2 (resp. 506.31) times faster than [25] (resp. [12]).

768 Similar results are presented in Table 4.2, which considers matrices that were randomly generated to
769 have a specified reciprocal condition number using MATLAB's sprandsym function. For these instances,
770 we impose a time limit for the algorithm proposed in [12], as these problems were more computationally
771 demanding for their code. We set the time limit to 3600 seconds, although in some cases, their code terminates
772 slightly after this limit.

773 On every instance considered in Table 4.2, Algorithm 2.2 and Algorithm B.1 reduced κ more than [25]
774 and [12] and were also much more efficient than these two algorithms. Algorithm 2.2 and Algorithm B.1
775 reduced κ on every instance while [25] (resp. [12]) reduced κ on 0 (resp. 12) of the 27 instances considered.
776 Moreover, Algorithm 2.2 (resp. Algorithm B.1) was also over 15 times faster than [25] and [12] on 18 (resp.
777 26) of the 27 instances considered.

778 Tables 4.3 and 4.4 illustrate the scalability of our Algorithms 2.2 and B.1 in efficiently minimizing κ
779 on large test instances. We do not compare against [25] and [12] since both codes were not able to handle
780 such large instances in a reasonable amount of time. Table 4.3 considers matrices that were taken from
781 the SuiteSparse Matrix Collection while Table 4.4 considers matrices that were randomly generated using
782 sprandsym to have very large κ values. The results in Table 4.3 show that on average, Algorithm 2.2 reduced
783 κ by 45.7% while Algorithm B.1 reduced κ by 41.9%. Both algorithms also frequently took less than one
784 minute to minimize κ . Table 4.4 shows that, across eleven difficult instances with $\kappa(M) > 10^{11}$, Algorithm 2.2
785 reduced κ by an average of 11.6%, while Algorithm B.1 achieved a 4.4% reduction.

786 **4.2. PCG Comparison for Solving Linear Systems.** In this subsection, we assume $M > 0$ and
787 compare the performance of PCG for solving $Mx = b$ using the preconditioner from Algorithm B.1 with the
788 one obtained by [25]. PCG is run with a tolerance of $1e - 4$ and a maximum iteration count of $5e6$. The
789 results are presented in Table 4.5. All matrices M considered in both tables are randomly generated, using
790 MATLAB's sprandsym function, with varying dimensions, densities, and κ values.

791 We see in Table 4.5 that on average, across the 47 instances in Table 4.5, PCG takes approximately

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

792 442, 147 iterations using the preconditioner obtained from Algorithm B.1 while it takes 549, 002 iterations using
793 the preconditioner from [25]. Moreover, the average total CPU time (including time for both preconditioning
794 and PCG) for Algorithm B.1 is 34.0 seconds, compared to 334.1 seconds for the algorithm in [25]. Hence, in
795 terms of total CPU time, Algorithm B.1 is on average approximately ten times faster than the method in [25].

TABLE 4.5
PCG Comparison Using Preconditioners Found by Algorithm in [25] and Algorithm B.1

Dim, Density, & $\kappa(A)$			% Reduction in κ		PCG Iterations		PCG CPU Time		Total CPU Time	
n	density	$\kappa(A)$	[25]	Algorithm B.1	[25]	Algorithm B.1	[25]	Algorithm B.1	[25]	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	329957	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	627550	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	604.96	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

796 **4.3. LSQR Comparison: Algorithm 3.1 vs. Two-Sided κ -optimal Scaling in [25].** In Subsec-
797 tion 4.3, we compare the efficiency of our two-sided Algorithm 3.1 which reduces ω at every iteration, with the
798 two-sided κ -optimal preconditioner presented in [25]. We consider matrices $A \in \mathcal{M}^n$ that are nonsingular and
799 compare the performance of both preconditioners for minimizing $\|b - Ax\|$ with LSQR. The results presenting
800 the comparison are presented in Tables 4.6 and 4.7. Table 4.6 considers matrices from the SuiteSparse Matrix
801 Collection while Table 4.7 considers matrices using MATLAB's sprandn function. Both tables consider
802 matrices where n does not exceed 300 since the two-sided κ -optimal algorithm of [25] could not handle larger
803 matrices. In both tables, LSQR is run with a tolerance of 10^{-8} and a maximum iteration count of 5000.
804 Table 4.8 compares the performance of LSQR with no preconditioning with its performance using the

TABLE 4.6
LSQR Comparison on Small Suitesparse Matrices: Algorithm 3.1 vs [25]

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	[25]	Algorithm 3.1	[25]	Algorithm 3.1	[25]
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 [25]

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	[25]	Algorithm 3.1	[25]	Algorithm 3.1	[25]
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

805 preconditioner from Algorithm 3.1 on large matrices from the SuiteSparse Matrix Collection with $n \geq 2000$.
 806 The algorithm in [25] is not considered in Table 4.8 since it could not handle these large problems. Since
 807 these problems are more difficult, LSQR is run with a tolerance of 10^{-6} and a maximum iteration count
 808 of 100000. 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 $\hat{C}^{1/2}$
 809 and $\hat{D}^{1/2}$ (resp. $D_1^{1/2}$ and $D_2^{-1/2}$) are the preconditioners found by Algorithm 3.1 (resp. [25]). Also, in all
 810 tables, the total CPU time is computed as the sum of the CPU time for preconditioning and the CPU time of
 811 LSQR. Finally, the list of SuiteSparse matrices considered in Tables 4.6 and 4.8 can be found in Appendix C.
 812

We see in Tables 4.6 and 4.7 that Algorithm 3.1 computed a preconditioner in significantly less time than
 813 the two-sided κ -optimal preconditioner in [25]. Moreover, reduced ω more significantly than [25] although [25]
 814 reduced κ more significantly. Interestingly enough, Algorithm 3.1 produced a preconditioner that resulted in
 815 fewer (or equal) LSQR iterations than [25] on 20 of 22 (resp. 18 of 26) instances presented in Table 4.6 (resp.
 816 Table 4.7). In summary, the results presented in Tables 4.6 and 4.7 display that ω -optimal preconditioners
 817 are much cheaper to compute than κ -optimal preconditioners and that the ω -condition number is also better
 818 correlated with LSQR iterations than the κ -condition number for minimizing $\|b - Ax\|$.

819 Table 4.8 shows that Algorithm 3.1 still computes effective preconditioners in negligible CPU time even on
 820 larger matrices with dimensions exceeding 2000. Moreover, Algorithm 3.1 reduced ω very significantly on 21
 821 of the 25 instances presented in Algorithm 3.1. Finally, Table 4.8 shows that Algorithm 3.1 led to a significant
 822 reduction in LSQR iterations. The total CPU time using preconditioning, i.e., Algorithm 3.1's CPU+ LSQR
 823 CPU, was smaller than the LSQR's CPU time without preconditioning on 22 of the 25 instances.

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	

TABLE 4.9

PCG tol. $1e-7$; Medium, PC; A, κ -opt VS J, ω -opt of A

Dim & Density		Ratios in conds (J, A)		J : ω -opt of A		Ratios A/J	
n	density	$\kappa(J)/\kappa(A)$	$\omega(J)/\omega(A)$	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

824 **4.4. From κ -optimal M to “Improve PCG” using ω -Optimal Scaling.** We now study the effect of
825 applying ω -optimal diagonal scaling to a κ -optimally diagonally scaled matrix M , where $M > 0$.¹⁰ We find a
826 κ -optimally diagonally scaled matrix $M > 0$ using [Theorem 2.7](#), i.e., the maximal and minimal eigenvectors of
827 M , x_1 and x_n , satisfy [\(2.5\)](#) and are chosen using [\(2.6\)](#). We then choose the remaining eigenvalues to be evenly
828 spread out in the open interval $\lambda_i \in (\lambda_n, \lambda_1)$, $i = 2, \dots, n-1$, with corresponding orthonormal eigenvectors x_i ,
829 $i = 2, \dots, n-1$, chosen in the orthogonal complement of $\text{span}\{x_1, x_n\}$, i.e., x_i are orthogonal to both x_1 and
830 x_n . We use a (sparse) QR factorization to obtain the remaining $n-2$ orthonormal eigenvectors. Therefore,
831 the density of M cannot be determined accurately in advance. The time to generate a random problem is
832 typically very small and is hence negligible. For a more precise description of how M is constructed, see our
833 code.

834 After constructing M , we then apply ω -optimal diagonal preconditioning to M to implicitly get J as
835 described above, i.e., $J = D^{1/2}MD^{1/2}$ where $D = \text{Diag}(1./\text{diag}(M))$. We then compare the number of PCG
836 iterations of solving $Mx = b$ and $J(D^{-1/2}x) = D^{1/2}b$. Note that only the vector b and the matrices A and
837 D are needed as input for MATLAB’s built-in PCG function. Each linear system that we consider in our
838 experiments is solved by PCG using 5 different initial points so the iteration counts and runtimes reported
839 for each experiment are averaged across these 5 runs. For the computational results showing the comparison
840 between A and J , see [Tables 4.9](#) and [4.10](#).

841 We now give several remarks about the results presented in [Tables 4.9](#) and [4.10](#). Interestingly, on every
842 instance tested applying the Jacobi or ω -optimal diagonal preconditioning to A led to significant improvement
843 in PCG’s performance for solving the linear system even though $\kappa(A) < \kappa(J)$. On average, across the 10
844 medium sized test instances presented in [Table 4.9](#), PCG performed 38.5 times more number of iterations on
845 the κ -optimal linear system than it did on ω -optimally scaled linear system. Moreover, PCG also on average

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

TABLE 4.10
PCG tol. 1e-7; Large, Linux; A, κ -opt VS J, ω -opt of A

Dim & Density		Ratios in conds (J, A)		J : ω -opt of A		Ratios A/J	
n	density	$\kappa(J)/\kappa(A)$	$\omega(J)/\omega(A)$	iters	cpu	iters	cpu
60000	9.14e-03	4.27e+00	7.366e-01	18.2	8.894e-01	42.0	36.5
65000	7.78e-03	4.57e+00	7.362e-01	21.0	9.985e-01	52.9	47.9
70000	6.54e-03	3.93e+00	7.365e-01	17.0	8.110e-01	48.3	42.9
75000	5.44e-03	4.18e+00	7.363e-01	18.0	8.253e-01	56.4	49.2
80000	4.36e-03	4.08e+00	7.369e-01	17.0	7.202e-01	38.3	33.8
85000	3.47e-03	3.91e+00	7.362e-01	18.0	7.062e-01	60.9	54.0
90000	2.63e-03	3.92e+00	7.371e-01	17.6	6.050e-01	34.2	29.8
95000	1.91e-03	4.12e+00	7.364e-01	18.0	5.092e-01	48.0	41.6
100000	1.31e-03	4.13e+00	7.369e-01	18.2	4.029e-01	34.4	29.8
105000	8.09e-04	4.41e+00	7.361e-01	22.0	3.342e-01	58.8	50.3
110000	4.29e-04	3.58e+00	7.363e-01	17.0	1.709e-01	55.9	45.2
115000	1.74e-04	3.97e+00	7.369e-01	20.6	1.134e-01	31.6	23.7
120000	3.32e-05	3.14e+00	7.362e-01	20.2	3.648e-02	55.0	26.0

846 took 34.2 times longer (in terms of CPU time) on the κ -optimal linear system than it did on the ω -optimally
847 scaled linear system.

848 Jacobi preconditioning led to even more significant improvements on the larger test instances presented
849 in Table 4.10. On average, across the 13 large test instances in Table 4.10, PCG performed 47.4 times more
850 number of iterations on the κ -optimal linear system than it did on ω -optimally scaled linear system. Moreover,
851 PCG also on average took 39.3 times longer (in terms of CPU time) on the κ -optimal linear system than it
852 did on the ω -optimally scaled linear system.

853 **5. Conclusion.** In this paper, we studied optimal diagonal preconditioning through the lens of two
854 condition numbers: the classical κ and the averaging-based ω -condition number. On the theoretical side, we
855 introduced an affine-based pseudoconvex reformulation of the κ -optimal preconditioning problem, yielding
856 simple optimality conditions and enabling efficient optimization over an n -dimensional vector. Moreover,
857 since the κ -condition number is positively homogeneous of degree zero, the associated minimization problem
858 is Hadamard ill-posed. To address this, we introduce a reformulation that removes this homogeneity, leading
859 to improved computational stability in practice. Building on this formulation, we develop a highly efficient
860 subgradient method with convergence guarantees that significantly outperforms existing SDP-based approaches
861 in both scalability and accuracy. For example, Table 4.1 shows that Algorithm 2.2 (resp. Algorithm B.1) is,
862 on average, 77.48 (resp. 356.22) times faster than [25] when minimizing κ on relatively small SuiteSparse
863 matrices. Our methods also substantially outperform [12] in terms of runtime. In addition to these speedups,
864 our subgradient methods consistently achieve significantly greater reductions in κ than both [25] and [12],
865 often reducing it by more than half. Finally, our algorithms scale to large problem instances with hundreds
866 of thousands of variables, solving them within minutes, whereas existing methods [12, 25] struggle to handle
867 such instances within a reasonable time.

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

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

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

881 **Acknowledgments.** The authors acknowledge the use of AI tools in revising the grammar and writing
882 of this paper.

883 **Appendix A. Assumptions and Technical Results from [19].** Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}} := \mathbb{R} \cup \{-\infty, \infty\}$ be
 884 an extended real-valued function with *domain* \bar{D} . The following minimization problem is considered in [19]:

$$885 \quad (\text{A.1}) \quad f_* := \inf \{f(x) : x \in X\},$$

886 with the following assumptions on f and X :

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

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

889 **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 $B(x, \epsilon) :=$
 890 $\{y : \|y - x\| \leq \epsilon\}$.

891 **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 $\forall \alpha \in \mathbb{R}$.

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

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

895 **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,*
 896 *b is finite and positive on $\text{int } \bar{D}$, $\text{int } \bar{D}$ is convex, and one of the following conditions holds:*

897 (a) b is affine;

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

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

900 *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*
 901 *$\partial[a - \alpha b](x) \subset \bar{\partial}^\circ f(x)$.*

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

904 **LEMMA A.2** ([19, Remark 2]). *Suppose a and $\tilde{b} := -b$ are proper convex functions on $\text{int } \bar{D}^a \cap \text{int } \bar{D}^b$,*
 905 *a is 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*

$$906 \quad f(x) := \begin{cases} a(x)/b(x), & \text{if } x \in C, \\ \sup_{y \in C} \bar{\lim}_{t \downarrow 0} f(x + t(y - x)), & \text{if } x \in bd C, \\ \infty, & \text{if } x \notin \text{cl } C \end{cases}$$

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

909 **Appendix B. An Efficient Line-Search Subgradient Method.** This section presents an efficient
 910 line-search subgradient method for minimizing $\kappa(v) := \kappa(\check{V}(v))$ where $\check{V}(v)$ is as in Equation (2.15).

911 Let $w = e + Vv \in \mathbb{R}_{++}^n$ correspond to the current iterate v . Also, let σ be a small scalar between 0
 912 and 1 and let $\Delta w = V\Delta v$ where $\Delta v = -\nabla \kappa(v)$ where $\kappa(v)$ is as in Equation (2.16). From Lemma 2.3, we
 913 can use a ratio test and guarantee positive definiteness from $w + tV\Delta v = w + t\Delta w > \sigma w$, or equivalently
 914 $t\Delta w > (\sigma - 1)w$. Therefore, we conclude that the maximum step with safeguarding, so as not to get too
 915 close to the boundary, is

$$916 \quad (\text{B.1}) \quad t_{\max}(\sigma) \cong t_{\max} := \min \left\{ \frac{(\sigma - 1)w_i}{\Delta w_i} : \Delta w_i < 0 \right\} > 0.$$

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

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

920 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 $0 < \sigma \ll 1$.

- 1: set $k \leftarrow 0$, $v_1 = 0 \in \mathbb{R}^{n-1}$, $w_1 = e_n \in \mathbb{R}^n$, and $t_k \leftarrow \infty$;
- 2: compute a min eigenpair (λ_n^1, x_n^1) and max eigenpair (λ_1^1, x_1^1) of A ;
- 3: compute $\kappa(v_1)$ and $\nabla\kappa(v_1)$ according to Equation (2.16) and Equation (2.19), respectively;
- 4: set `relnormg` = $\|\nabla\kappa(v_1)\|^2/(1 + \kappa(v_1)^2)$;
- 5: **while** `relnormg` $>$ `maintoler` & $t_k >$ `linesrchtoler` & $k \leq$ `mainmaxiter` **do**
- 6: set $k \leftarrow k + 1$;
- 7: set $\Delta v_k \leftarrow -\nabla\kappa(v_k)$;
- 8: use (B.1) with $\Delta w := V\Delta v_k$ and $w := w_k$ to evaluate $t_{\max}(\sigma)$;
- 9: perform [LineSearch B.1](#) to find t_k ;
- 10: set $v_{k+1} = v_k + t_k\Delta v_k$ and $w_{k+1} = e_n + Vv_{k+1}$;
- 11: compute min eigenpair $(\lambda_n^{k+1}, x_n^{k+1})$ and max eigenpair $(\lambda_1^{k+1}, x_1^{k+1})$ of $A \text{Diag}(w_{k+1})$;
- 12: compute $\kappa(v_{k+1})$ and $\nabla\kappa(v_{k+1})$ according to Equation (2.16) and Equation (2.19), respectively;
- 13: update `relnormg` = $\|\nabla\kappa(v_{k+1})\|^2/(1 + \kappa(v_{k+1})^2)$;
- 14: **end while**(main outer loop)

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

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

922 **Appendix C. List of SuiteSparse Matrices Used in Experiments.** Table 4.1 considers the
923 following matrices:

924 “bcsstk08.mat”, “bcsstk13.mat”, “bcsstk21.mat”, “bcsstk23.mat”, “bcsstk24.mat”
925 “bcsstk26.mat”, “bcsstk28.mat”, “bcsstk34.mat”, “494_bus.mat”, “662_bus.mat”
926 “nasa1824.mat”, “nasa2146.mat”, “nasa2910.mat”, “nos1.mat”, “nos2.mat”, “nos4.mat”
927 “nos5.mat”, “nos7.mat”

928 Table 4.3 considers the following matrices:

929 “Pres_Poisson.mat”, “bcsstk25.mat”, “bcstm25.mat”, “gyro_m.mat”
930 “gyro.mat”, “bcsstk36.mat”, “wathen100.mat”, “wathen120.mat”, “minsurfo.mat”
931 “gridgena.mat”, “G2_circuit.mat”

932 Table 4.6 considers the following matrices:

933 “08blocks.mat”, “arc130.mat”, “bwm200.mat”, “ck104.mat”, “ex1.mat”
934 “football.mat”, “gre_185.mat”, “gre_216a.mat”, “impcol_a.mat”, “impcol_c.mat”
935 “impcol_e.mat”, “jazz.mat”, “lshp_265.mat”, “olm100.mat”, “polbooks.mat”
936 “rajat11.mat”, “robot.mat”, “rotor1.mat”, “rw136.mat”, “tub100.mat”
937 “utm300.mat”, “west013.mat”

938 Table 4.8 considers the following matrices:

939 “airfold_2d.mat”, “c-24.mat”, “c-36.mat”, “c-39.mat”, “cavity21.mat”
940 “conf5_-4x4-18.mat”, “coupled.mat”, “delaunay_n11.mat”, “delaunay_n13.mat”, “epb1.mat”
941 “ex24.mat”, “G34.mat”, “G36.mat”, “G59.mat”, “G63.mat”
942 “G64.mat”, “garon1.mat”, “Hamrle2.mat”, “jan99jac040sc.mat”, “Na5.mat”
943 “t2d_q9.mat”, “tols4000.mat”, “utm5940.mat”, “viscoplastic1.mat”
944 “whitaker3.mat”

- 946 [1] M. BENZI, *Preconditioning techniques for large linear systems: a survey*, J. Comput. Phys., 182 (2002), pp. 418–477,
947 <https://doi.org/10.1006/jcph.2002.7176>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1006/jcph.2002.7176>.
- 948 [2] A. BOCK AND M. ANDERSEN, *Connecting Kaporin's condition number and the Bregman log determinant divergence*, tech.
949 report, 2025, <https://arxiv.org/abs/2503.22286>, <https://arxiv.org/abs/2503.22286>.
- 950 [3] R. BYRD AND R. NOCEDAL, *A tool for the analysis of quasi-Newton methods with application to unconstrained minimization*,
951 SIAM J. Numer. Anal., 26 (1989), pp. 727–739.
- 952 [4] R. BYRD, R. NOCEDAL, AND Y. YUAN, *Global convergence of a class of quasi-Newton methods on convex problems*, SIAM J.
953 Numer. Anal., 24 (1987), pp. 1171–1191.
- 954 [5] G. CIARAMELLA AND M. J. GANDER, *Iterative methods and preconditioners for systems of linear equations*, vol. 19 of
955 Fundamentals of Algorithms, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, [2022] ©2022,
956 <https://doi.org/10.1137/1.9781611976908>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1137/1.9781611976908>.
- 957 [6] T. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Transactions on Mathematical Software
958 (TOMS), 38 (2011), pp. 1–25.
- 959 [7] J. E. DENNIS, JR. AND H. WOLKOWICZ, *Sizing and least-change secant methods*, SIAM J. Numer. Anal., 30 (1993),
960 pp. 1291–1314, <https://doi.org/10.1137/0730067>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1137/0730067>.
- 961 [8] X. DOAN AND H. WOLKOWICZ, *Numerical computations and the ω -condition number*, Tech. Report CORR 2011-03,
962 University of Waterloo, Waterloo, Ontario, 2011. www.math.uwaterloo.ca/~hwolkowi/henry/reports/ONE.pdf.
- 963 [9] X. V. DOAN, S. KRUK, AND H. WOLKOWICZ, *A robust algorithm for semidefinite programming*, Optim. Methods Softw., 27
964 (2012), pp. 667–693, <https://doi.org/10.1080/10556788.2011.610456>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1080/10556788.2011.610456>.
- 965 [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
966 Approximation, 11 (2018), pp. 11–22.
- 967 [11] W. GAO, Y.-C. CHU, Y. YE, AND M. UDELL, *Gradient methods with online scaling*, tech. report, 2024, <https://arxiv.org/abs/2411.01803>,
969 <https://arxiv.org/abs/2411.01803>. arXiv, 2411.01803.
- 970 [12] W. GAO, Z. QU, M. UDELL, AND Y. YE, *Scalable approximate optimal diagonal preconditioning*, tech. report, 2023,
971 <https://arxiv.org/abs/2312.15594>. arXiv, 2312.15594.
- 972 [13] J. A. GEORGE AND J. W.-H. LIU, *Computer Solution of Large Sparse Positive Definite Systems*, Prentice-Hall, Englewood
973 Cliffs, NJ, 1981.
- 974 [14] A. GREENBAUM, *Iterative methods for solving linear systems*, Society for Industrial and Applied Mathematics (SIAM),
975 Philadelphia, PA, 1997.
- 976 [15] Y. HU, J. LI, AND C. K. W. YU, *Convergence rates of subgradient methods for quasi-convex optimization problems*,
977 Comput. Optim. Appl., 77 (2020), pp. 183–212, <https://doi.org/10.1007/s10589-020-00194-y>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1007/s10589-020-00194-y>.
- 978 [16] C. G. J. JACOBI, *Ueber eine neue auflösungsart der bei der methode der kleinsten quadrate vorkommenden lineären*
979 *gleichungen*, Astronomische Nachrichten, 22 (1845), p. 297–306.
- 980 [17] W. L. JUNG, D. TORREGROSA-BELÉN, AND H. WOLKOWICZ, *The ω -condition number: applications to preconditioning and*
981 *low rank generalized Jacobian updating*, Comput. Optim. Appl., 91 (2025), pp. 235–282, <https://doi.org/10.1007/s10589-025-00669-w>,
982 <https://doi-org.proxy.lib.uwaterloo.ca/10.1007/s10589-025-00669-w>.
- 983 [18] I. KAPORIN, *New convergence results and preconditioning strategies for the conjugate gradient method*, Numer. Linear
984 Algebra Appl., 1 (1994), pp. 179–210, <https://doi.org/10.1002/nla.1680010208>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1002/nla.1680010208>.
- 985 [19] K. KIWIEL, *Convergence and efficiency of subgradient methods for quasiconvex minimization*, Mathematical programming,
986 90 (2001), pp. 1–25.
- 987 [20] P. A. KNIGHT, *The Sinkhorn–Knopp algorithm: convergence and applications*, SIAM J. Matrix Anal. Appl., 30 (2008),
988 pp. 261–275.
- 989 [21] P. A. KNIGHT, D. RUIZ, AND B. UÇAR, *A symmetry preserving algorithm for matrix scaling*, SIAM J. Matrix Anal. Appl.,
990 35 (2014), pp. 931–955, <https://doi.org/10.1137/110825753>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1137/110825753>.
- 991 [22] A. V. KNYAZEV, *Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient*
992 *method*, SIAM journal on scientific computing, 23 (2001), pp. 517–541.
- 993 [23] J. LEE, *Smooth manifolds*, in Introduction to smooth manifolds, Springer, 2003, pp. 1–29.
- 994 [24] MOSEK APS, *MOSEK Optimization Suite*, 2025, <https://www.mosek.com>. Version 11.1 (or the version used).
- 995 [25] Z. QU, W. GAO, O. HINDER, Y. YE, AND Z. ZHOU, *Optimal diagonal preconditioning*, Operations Research, 73 (2024),
996 pp. 1479–1495, <https://doi.org/10.1287/opre.2022.0592>, <https://doi.org/10.1287/opre.2022.0592>, <https://arxiv.org/abs/https://doi.org/10.1287/opre.2022.0592>.
- 997 [26] Y. SAAD, *Iterative methods for sparse linear systems*, Society for Industrial and Applied Mathematics, Philadelphia,
998 PA, second ed., 2003, <https://doi.org/10.1137/1.9780898718003>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1137/1.9780898718003>.
- 999 [27] Y. SAAD, *Iterative methods for linear systems of equations: a brief historical journey*, 754 ([2020] ©2020), pp. 197–215,
1000 <https://doi.org/10.1090/conm/754/15141>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1090/conm/754/15141>.
- 1001 [28] M. SCETBON, C. MA, W. GONG, AND E. MEEDS, *Gradient multi-normalization for stateless and scalable llm training*,
1002 arXiv preprint arXiv:2502.06742, (2025).
- 1003 [29] G. W. SOULES, *The rate of convergence of Sinkhorn balancing*, in Proceedings of the First Conference of the International
1004 Linear Algebra Society (Provo, UT, 1989), vol. 150, 1991, pp. 3–40, [https://doi.org/10.1016/0024-3795\(91\)90157-R](https://doi.org/10.1016/0024-3795(91)90157-R),
1005 [https://doi-org.proxy.lib.uwaterloo.ca/10.1016/0024-3795\(91\)90157-R](https://doi-org.proxy.lib.uwaterloo.ca/10.1016/0024-3795(91)90157-R).
- 1006 [30] L. TUNÇEL AND H. WOLKOWICZ, *Strengthened existence and uniqueness conditions for search directions in semidefinite*
1007 *programming*, Linear Algebra Appl., 400 (2005), pp. 31–60, <https://doi.org/10.1016/j.laa.2004.12.010>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1016/j.laa.2004.12.010>.

- 1013 [31] N. VAN DER AA, H. TER MORSCHÉ, AND R. MATTHEIJ, *Computation of eigenvalue and eigenvector derivatives for a general*
1014 *complex-valued eigensystem*, Electron. J. Linear Algebra, 16 (2007), pp. 300–314, [https://doi.org/10.13001/1081-3810.](https://doi.org/10.13001/1081-3810.1203)
1015 [1203](https://doi.org/10.13001/1081-3810.1203), <https://doi.org/10.13001/1081-3810.1203>.
- 1016 [32] A. WATHEN, *Preconditioning*, Acta Numer., 24 (2015), pp. 329–376, <https://doi.org/10.1017/S0962492915000021>, <https://doi.org/10.1017/S0962492915000021>.
- 1017 [33] H. WOLKOWICZ AND Q. ZHAO, *An all-inclusive efficient region of updates for least change secant methods*, SIAM J. Optim.,
1018 5 (1995), pp. 172–191, <https://doi.org/10.1137/0805009>, <https://doi-org.proxy.lib.uwaterloo.ca/10.1137/0805009>.
- 1019 [34] W. E. YOUNG AND R. H. TRENT, *Geometric mean approximations of individual security and portfolio performance*, The
1020 Journal of Financial and Quantitative Analysis, 4 (1969), pp. 179–199, <http://www.jstor.org/stable/2329839> (accessed
1021 2026-03-03).
- 1022 [35] Q. ZHAO, *Measures for least change secant methods*, master’s thesis, University of Waterloo, 1993.
- 1023