

Numerical Computations and the ω -Condition Number ^{*†}

Xuan Vinh Doan Henry Wolkowicz

July 4, 2011[‡]

Key words and phrases: Condition numbers, uniform condition number, preconditioning, roundoff-error, iterative methods

AMS subject classifications: 15A12, 65F35, 65F08, 65F10, 65G50

Contents

1	Introduction	2
1.1	Outline of Results	3
2	Basic Properties	3
2.1	Numerical evaluation of ω	4
3	Iterations for Krylov Subspace Methods and ω	4
4	ω-Condition Number and Preconditioning	8
5	ω-Condition Number and Relative Error	13
6	Conclusion	16
	Bibliography	17
	Index	19

List of Algorithms

2.1	Evaluation of $\omega_c(\mathbf{X}) = \sqrt{(\omega(\mathbf{A}))} = \sqrt{(\omega(\mathbf{X}^T \mathbf{X}))}$ with given tolerance, tol	5
-----	---	---

List of Tables

4.1	Average # iterations; computational time block preconditioners; sparse matrices . . .	13
-----	---	----

*Research Report CORR 2011-03

†This report is available at URL: orion.math.uwaterloo.ca/~hwoikowi/henry/reports/ABSTRACTS.html

‡Department of Combinatorics and Optimization, University of Waterloo, 200 University Avenue West, Waterloo, ON N2L 3G1, Canada

List of Figures

1	(# distinct λ_i) <u>vs</u> (# Krylov iterations) is pos. correlated; 5000 random instances.	7
2	(ω) <u>vs</u> (# Krylov iterations) is pos. correlated for large iterations	8
3	(ω) <u>vs</u> (# Krylov iterations) is neg. correlated for large ω ; λ_1, λ_n fixed	9
4	(κ) <u>vs</u> (# Krylov iterations) not correlated	9
5	κ with two different diagonal preconditioners	12
6	ω with two different diagonal preconditioners	12
7	# Krylov iterations with two different diagonal preconditioners for sparse matrices .	13
8	Ratio of ω and the relative error $\bar{r}(\mathbf{A})$	16
9	Ratio of κ and the relative error $\bar{r}(\mathbf{A})$	17

Abstract

We study the computation and properties of a new measure for the condition number of a positive definite matrix. This measure, the ratio of the arithmetic and geometric means of the eigenvalues, depends uniformly on all the eigenvalues of the matrix. Moreover, it can be evaluated easily and accurately. And, we see that: (i) it correlates better with the number of iterations in iterative methods compared to the standard condition number which depends only on the largest and smallest eigenvalues; (ii) it provides a criteria for obtaining optimal efficient preconditioners; and (iii) it presents a more average relative error measure compared to the worst case behaviour of the standard condition number.

1 Introduction

We study the properties of the following condition number measure, which can be computed accurately and efficiently and which depends uniformly on all the eigenvalues of the matrix \mathbf{A} :

$$\omega(\mathbf{A}) := \frac{\text{trace}(\mathbf{A})/n}{\det(\mathbf{A})^{\frac{1}{n}}} = \frac{\sum_{i=1}^n \lambda_i(\mathbf{A})/n}{\left(\prod_{i=1}^n \lambda_i(\mathbf{A})\right)^{\frac{1}{n}}}, \quad \mathbf{A} \in \mathcal{S}_{++}^n, \quad (1.1)$$

where \mathcal{S}_{++}^n is the cone of real symmetric positive definite matrices. The standard condition number $\kappa(\mathbf{A})$ depends only on the largest and smallest eigenvalues of \mathbf{A} . It is generally used as an indicator for whether the problem of solving the system of linear equations $\mathbf{Ax} = \mathbf{b}$ is well-conditioned (low condition number) or ill-conditioned (high condition number). In general, iterative algorithms used to solve the system $\mathbf{Ax} = \mathbf{b}$ require a large number of iterations to achieve a solution with high accuracy if the problem is ill-conditioned. Preconditioners are introduced to obtain a better conditioned problem. Currently, reducing $\kappa(\mathbf{A})$ condition number is the main aim for preconditioners. The condition number $\kappa(\mathbf{A})$ is also a measure of how much a solution \mathbf{x} will change with respect to changes in the right-hand side \mathbf{b} .

Our goal in this paper is to establish the basic properties of the ω -condition number and to study whether it is a better indicator of whether the problem $\mathbf{Ax} = \mathbf{b}$ is well- or ill-conditioned. In addition, we present a procedure for efficiently calculating $\omega(\mathbf{A})$ and avoiding the difficulty of evaluating the determinant of large matrices.

1.1 Outline of Results

We continue in Section 2 with the basic properties of the condition number measure ω . We include Section 2.1 with Algorithm 2.1 that calculates ω efficiently and accurately. Section 3 presents results that relate the number of CG iterations for solving $\mathbf{Ax} = \mathbf{b}$ with the magnitudes of the two condition number measures. The better (inverse) correlation with ω is described in Items 1 to 4, page 7. Since the number of iterations depends on the number of distinct eigenvalues of A , we illustrate the seeming paradox that with the value of κ fixed, then both (relatively) large and small values of ω result in a low number of iterations. Section 4 compares how preconditioners arise from the two measures and their effectiveness. The expense of finding an optimal ω preconditioner is low compared to the optimal κ preconditioner. Section 5 discusses the role of the two condition numbers in error analysis. We provide the analytic formulation for the expected relative error under the assumption of independent normal distributions. We illustrate the known result that κ is a bound for the worst case of roundoff error, and show that ω provides a better estimate of the expected roundoff error. Concluding remarks are given in Section 6.

2 Basic Properties

The standard condition number of a nonsingular matrix \mathbf{A} is defined as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|,$$

where $\|\cdot\|$ is a matrix norm. We let \mathcal{S}^n denote the space of real $n \times n$ symmetric matrices equipped with the trace inner product; \mathcal{S}_+^n denote the cone of positive semidefinite matrices; and \mathcal{S}_{++}^n is the cone of positive definite matrices. If $\|\cdot\|$ is the ℓ_2 -norm, we then have: $\kappa(\mathbf{A}) = \frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}(\mathbf{A})}$, the ratio of largest and smallest singular values; and, if $\mathbf{A} \in \mathcal{S}_{++}^n$, then $\kappa(\mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}$.

The ω -condition number has been studied in the context of scaling for quasi-Newton methods in [1]. Some basic properties of $\omega(\mathbf{A})$ with respect to $\kappa(\mathbf{A})$ are shown in the following.

Proposition 1 ([1]). *Let $A \succ 0$ be given. The measure $\omega(\mathbf{A})$ satisfies*

1. $1 \leq \omega(\mathbf{A}) \leq \kappa(\mathbf{A}) < \frac{(\kappa(\mathbf{A}) + 1)^2}{\kappa(\mathbf{A})} \leq 4\omega^n(\mathbf{A}),$

with equality in the first and second inequality if and only if \mathbf{A} is a multiple of the identity and equality in the last inequality if and only if

$$\lambda_2 = \cdots = \lambda_{n-1} = \frac{\lambda_1 + \lambda_n}{2};$$

2. $\omega(\alpha\mathbf{A}) = \omega(\mathbf{A}),$ for all $\alpha > 0;$

3. if $n = 2,$ $\omega(\mathbf{A})$ is isotonic with $\kappa(\mathbf{A}).$ □

For a full rank matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ we let σ denote the vector of nonzero singular values. Then we can also define additional measures for general matrices.

$$\omega_2(\mathbf{A}) := \frac{\mathcal{A}(\sigma)}{\mathcal{G}(\sigma)} \geq \omega_c(\mathbf{A}) := \sqrt{\omega(\mathbf{A}^T \mathbf{A})} = \sqrt{\frac{\text{trace}(\mathbf{A}^T \mathbf{A})/n}{\det(\mathbf{A}^T \mathbf{A})^{\frac{1}{n}}}} = \sqrt{\frac{\mathcal{A}(\sigma^2)}{\mathcal{G}(\sigma^2)}},$$

where \mathcal{A}, \mathcal{G} , denote the arithmetic and geometric means respectively. These two condition numbers are for least squares solutions of overdetermined systems $\mathbf{A}\mathbf{x} = \mathbf{b}, m > n$.

2.1 Numerical evaluation of ω

One issue with $\kappa(\mathbf{A})$ is how to estimate it efficiently when the size of matrix \mathbf{A} is large since eigendecompositions can be expensive. A survey of estimates and, in particular, estimates using the ℓ_1 -norm, are given in [4, 5]. Extensions to sparse matrices and block-oriented generalizations are given in [3, 6]. Results from these papers form the basis of the *condst* command in MATLAB; this illustrates the difficulty in accurately estimating $\kappa(\mathbf{A})$.

On the other hand, the measure $\omega(\mathbf{A})$ can be calculated using the trace and determinant function which do not require eigenvalue decompositions. However, for large n , the determinant is also numerically difficult to compute as it could easily result in an overflow $+\infty$ or 0 due to the limits of finite precision arithmetic. In order to overcome this problem, we use the following simple homogeneity property for the determinant: $\det(\alpha\mathbf{A}) = \alpha^n \det(\mathbf{A}), \forall \mathbf{A} \in \mathbb{R}^{n \times n}$. Thus we have:

$$\det(\mathbf{A})^{\frac{1}{n}} = \frac{\det(\alpha\mathbf{A})^{\frac{1}{n}}}{\alpha}, \quad \forall \alpha > 0.$$

Our main task is then to appropriately select $\alpha > 0$ such that $\det(\alpha\mathbf{A})$ is well-defined within the machine precision. Using the fact that $\det(\mathbf{A}) = \left| \prod_{i=1}^n u_{ii} \right|$, where u_{ii} are the diagonal elements of

\mathbf{U} in the LU decomposition of \mathbf{A} , we can simplify the calculation of $\det(\alpha\mathbf{A})$ by using the vector $\text{diag}(\mathbf{U})$. Since $\log \det(\alpha\mathbf{A})$ is a nondecreasing function in $\alpha > 0$, we can use a simple binary search (not necessarily with high accuracy) to find α such that $\log \det(\alpha\mathbf{A}) \in (-\infty, +\infty)$, see Algorithm 2.1, page 5. This technique is used in our numerical tests throughout this paper. A MATLAB file for this evaluation is available with

URL: orion.math.uwaterloo.ca/~hwoikowi/henry/reports/ABSTRACTS.html#unifcondnumb.

3 Iterations for Krylov Subspace Methods and ω

When the matrix \mathbf{A} is large and sparse, then efficient methods to solve $\mathbf{A}\mathbf{x} = \mathbf{b}$ are based on Krylov subspaces. These methods perform repeated matrix-vector multiplications involving \mathbf{A} , e.g., Ipsen and Meyer [8]. The idea behind Krylov subspace methods is that after k iterations, the solution \mathbf{x}_k lies in the Krylov space generated by a vector \mathbf{c} ,

$$\mathcal{K}_k(\mathbf{A}, \mathbf{c}) := \text{span}(\mathbf{c}, \mathbf{A}\mathbf{c}, \dots, \mathbf{A}^{k-1}\mathbf{c}).$$

If $\mathbf{A} \in \mathcal{S}_{++}^n$, then the solution \mathbf{x} belongs to the Krylov space $\mathcal{K}_d(\mathbf{A}, \mathbf{b})$, where d is the number of distinct eigenvalues of \mathbf{A} . This implies that, theoretically, a Krylov subspace method finds the solution \mathbf{x} after d iterations. We would like to investigate how the ω -condition number relates to the number of iterations of Krylov subspace methods. We first show that a small number of distinct eigenvalues implies a small ω -condition number.

The ω -condition number can be defined as a function of the eigenvalue vector $\boldsymbol{\lambda}$ of the positive

Algorithm 2.1: Evaluation of $\omega_c(\mathbf{X}) = \sqrt{\omega(\mathbf{A})} = \sqrt{\omega(\mathbf{X}^T \mathbf{X})}$ with given tolerance, tol .

```

1 Input( $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ ,  $\mathbf{X}$  full column rank), tolerance  $tol$  ;
2 Set:  $tA = \text{trace}(\mathbf{A})/n$ ,  $[\mathbf{L}, \mathbf{U}] = \text{lu}(\mathbf{A})$ ,  $dA = \text{prod}(\text{diag}(\mathbf{U}))$  ;
3 if  $dA$  and  $\log(dA)$  are finite then
4   | Set:  $\omega = \sqrt{\frac{tA}{(dA)^{\frac{1}{n}}}}$ , RETURN
5 end if
6 if  $dA$  is infinite then
7   | Set:  $aL = 0$ ,  $aH = 1$  ;
8 else
9   | Set:  $aL = 1$ ,  $aH = 2$ ,  $dA = \text{prod}(aH * \text{diag}(\mathbf{U}))$ ,  $itercount = 1$  ;
10  | while  $dA$  is finite &  $\log(dA)$  is finite &  $itercount < maxiter$  do
11  |   |  $aL = aH$ ,  $aH = 2 * aH$ ,  $dA = \text{prod}(aH * \text{diag}(\mathbf{U}))$ ,  $index = index + 1$  ;
12  |   end while
13 end if
14 while ( $dA$  is infinite |  $\log(dA)$  is infinite) &  $(aH - aL)/2 > tol$  do
15  | Set:  $a = (aH + aL)/2$ ,  $dA = \text{prod}(a * \text{diag}(\mathbf{U}))$  ;
16  | if  $dA$  is finite &  $\log(dA)$  is finite then
17  |   | Set:  $\omega = \sqrt{\frac{a * tA}{(dA)^{\frac{1}{n}}}}$ , RETURN
18  |   end if
19  | if  $dA$  is infinite then
20  |   | Set:  $aH = a$  ;
21  |   else
22  |   | Set:  $aL = a$  ;
23  |   end if
24 end while
25 if  $(aH - aL)/2 > tol$  then
26  | Set:  $\omega = \sqrt{\frac{a * tA}{(dA)^{\frac{1}{n}}}}$  ;
27 else
28  | Set:  $\omega = \sqrt{tA * (aH + aL)/2}$  ;
29 end if
30 Output( $\omega_c(\mathbf{X}) = \sqrt{\omega(\mathbf{A})} = \omega$ ) ;

```

definite matrix \mathbf{A} :

$$\omega(\mathbf{A}) = \omega(\boldsymbol{\lambda}) := \frac{\left(\sum_{i=1}^n \lambda_i/n\right)}{\left(\prod_{i=1}^n \lambda_i\right)^{\frac{1}{n}}}.$$

By abuse of notation, we allow ω to act on \mathbb{R}_{++}^n ; and, for a subset $\mathcal{S} \subseteq \{1, \dots, n\}$, we let $\boldsymbol{\lambda}_{\mathcal{S}} := (\lambda_i)_{i \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}|}$.

Proposition 2. *Let the set $\emptyset \neq \mathcal{S} \subsetneq \{1, \dots, n\}$ and the positive vector $\bar{\boldsymbol{\lambda}}_{\mathcal{S}} > 0$ be given. Then the optimal solution of*

$$\begin{aligned} \min \quad & \omega(\boldsymbol{\lambda}) \\ \text{s.t.} \quad & \lambda_i = \bar{\lambda}_i, \quad \forall i \in \mathcal{S}, \\ & \boldsymbol{\lambda} \in \mathbb{R}_{++}^n, \end{aligned} \tag{3.1}$$

is given by

$$\lambda_i = \begin{cases} \bar{\lambda}_i, & \text{if } i \in \mathcal{S}, \\ \frac{1}{n - |\mathcal{S}|} \sum_{i \in \mathcal{S}} \bar{\lambda}_i, & \text{otherwise.} \end{cases}$$

Proof. We can change the problem to the unconstrained minimization

$$\min \log \left(\sum_{i \in \mathcal{S}} \bar{\lambda}_i + \sum_{i \notin \mathcal{S}} \lambda_i \right) - \frac{1}{n} \log \left(\prod_{i \in \mathcal{S}} \bar{\lambda}_i \right) \left(\prod_{i \notin \mathcal{S}} \lambda_i \right).$$

Let $K^* = \sum_{i \in \mathcal{S}} \bar{\lambda}_i + \sum_{i \notin \mathcal{S}} \lambda_i^*$ at an optimal λ^* . Then stationarity implies that

$$\frac{1}{K^*} - \frac{1}{n\lambda_i^*} = 0, \quad \forall i \notin \mathcal{S}.$$

This means that $(n - |\mathcal{S}|)\lambda_i = \sum_{i \in \mathcal{S}} \bar{\lambda}_i, \forall i \notin \mathcal{S}$. □

Proposition 2 shows that the ω -condition number is minimized when all free eigenvalues are equal, which means the number of distinct eigenvalues is reduced significantly. We emphasize the comparison with the standard condition number $\kappa(\mathbf{A})$.

Corollary 1. *Let $\mathcal{S} = \{1, n\}$ in Proposition 2. Then the optimal solution of (3.1) is given by*

$$\lambda_i^* = \frac{1}{n-2} (\bar{\lambda}_1 + \bar{\lambda}_n). \quad \square$$

Corollary 2. *Let $\{1, n\} \subseteq \mathcal{S}$ in Proposition 2. And, consider the problem (3.1) with \min replaced by \max and with the added constraints that the optimal vector $\boldsymbol{\lambda}$ must maintain the ordering $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Then the number of distinct numbers in the optimal solution $\boldsymbol{\lambda}^*$ is given by $|\mathcal{S}|$.*

Proof. The result follows immediately from the pseudoconvexity¹ of the function ω , i.e., the \max must occur at an extreme point of the feasible set. □

¹See Proposition 3 Item 1, below, for the statement on the pseudoconvexity of ω .

Corollary 1 indicates that the number of distinct eigenvalues reduces when the ω -condition number is minimized. It shows that the ω -condition number is likely positively correlated with the number of iterations of Krylov subspace methods for solving the system of linear equation $\mathbf{Ax} = \mathbf{b}$ when the standard condition number is not kept constant. In fact, we expect a better correlation for ω compared to κ since ω depends on all the eigenvalues. Surprisingly, Corollary 2 indicates that when $\kappa(\mathbf{A})$ remains the same, then we get a reduction in the number of iterations when ω is *maximized*.

We now construct numerical results to test the relationship between the number of iterations and the two condition numbers.

Random Procedure 1. *Unless specified otherwise, in the following tests for Figures 1 to 4: we generate random matrices with a particular number of distinct eigenvalues; we then solve the system of linear equations with random right-hand sides using the preconditioned conjugate gradient method, pcg, in MATLAB. The size of the matrices is $n = 5000$ with the number of distinct eigenvalues chosen randomly in the interval $[1, 500]$. The tolerance for pcg is set to 10^{-10} . We used 5000 random instances.*

1. The scatter plot in Figure 1 shows, as expected, that the number of distinct eigenvalues (on the horizontal axis) is positively correlated with the number of Krylov iterations to solve the linear system (on the vertical axis).

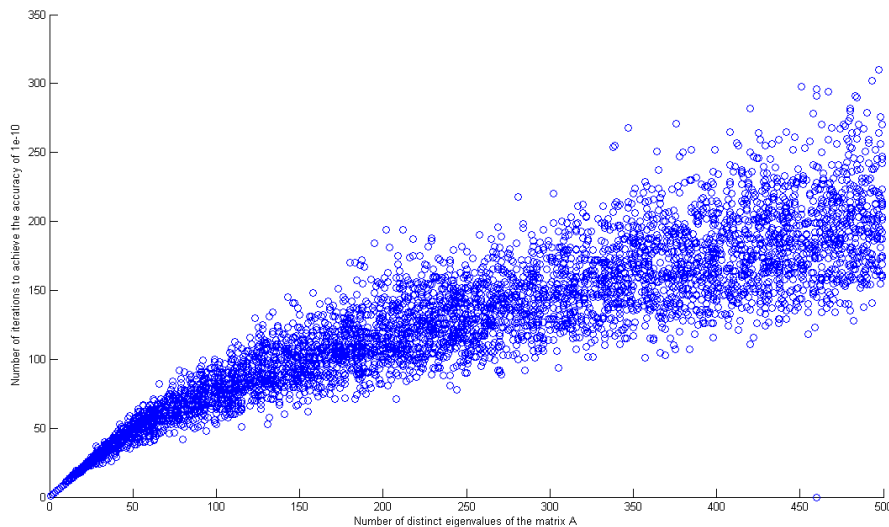


Figure 1: ($\#$ distinct λ_i) vs ($\#$ Krylov iterations) is pos. correlated; 5000 random instances.

2. The scatter plot in Figure 2 indicates that when the number of iterations is large enough, there is a positive correlation between the number of iterations and the ω -condition number of matrix \mathbf{A} . But this is only for small values of ω . In fact, the relatively large values of ω correspond to very low iteration counts. This confirms the results about maximizing in Corollary 2.

3. The maximization of ω observation in Corollary 2 is illustrated further in the scatter plot in Figure 3, where we fix λ_1, λ_n . We see the low number of Krylov iterations (on the vertical axis) when ω (on the horizontal axis) is small. But, in addition, as ω gets large, the number of Krylov iterations decrease dramatically.

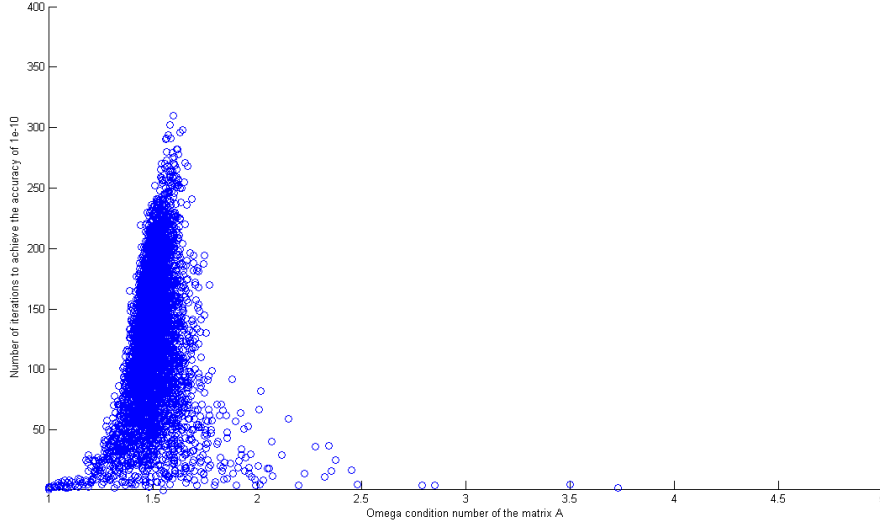


Figure 2: (ω) vs. (# Krylov iterations) is pos. correlated for large iterations

4. Figure 4 shows the relationship between the number of iterations and the standard condition number. In sharp contrast to the results in Figures 2 and 3, We can see that there is no clear correlation between the number of Krylov iterations and the standard condition number.

4 ω -Condition Number and Preconditioning

If the matrix \mathbf{A} is ill-conditioned, one way to achieve solutions with high accuracy is with preconditioners. Consider the overdetermined system $\mathbf{A}\mathbf{x} = \mathbf{b}$ with a full rank matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, where $m > n$. The least squares solution is the solution of the system $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b}$. A nonsingular preconditioner $\mathbf{D} \in \mathbb{R}^{n \times n}$ is used to create a new overdetermined system $\mathbf{A}\mathbf{D}\mathbf{y} = \mathbf{b}$, where $\mathbf{y} = \mathbf{D}^{-1}\mathbf{x}$, such that the least squares solution \mathbf{y} will be solved by a better-conditioned system of linear equations, $\mathbf{D}^T \mathbf{A}^T \mathbf{A}\mathbf{D}\mathbf{y} = \mathbf{D}^T \mathbf{A}^T \mathbf{b}$. We could try to find \mathbf{D} that minimizes $\kappa(\mathbf{D}^T \mathbf{A}^T \mathbf{A}\mathbf{D})$ since the standard condition number is an indicator of whether the problem is well- or ill-conditioned. Lu and Pong [9] shows that the diagonal preconditioner \mathbf{D} that minimizes $\kappa(\mathbf{D}^T \mathbf{A}^T \mathbf{A}\mathbf{D})$ can be compute

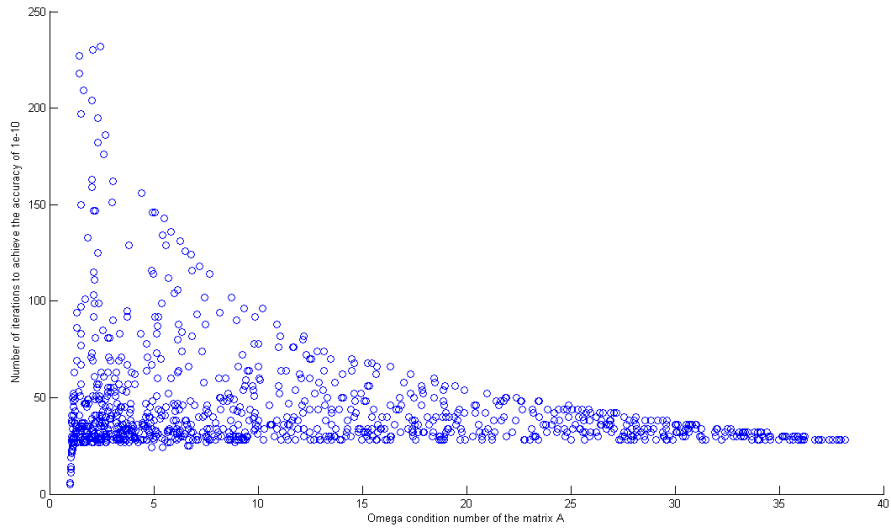


Figure 3: (ω) vs ($\#$ Krylov iterations) is neg. correlated for large ω ; λ_1, λ_n fixed

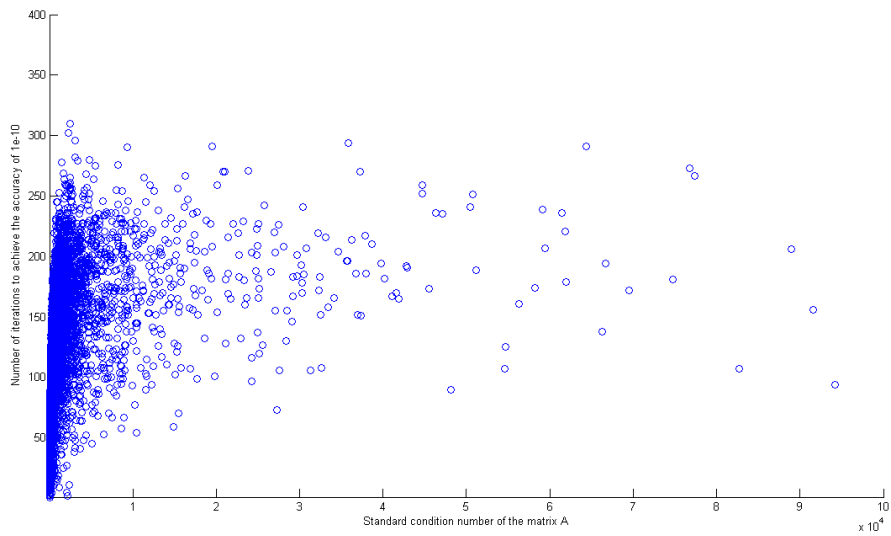


Figure 4: (κ) vs ($\#$ Krylov iterations) not correlated

by solving a conic optimization problem with semidefinite and second-order cone constraints:

$$\begin{aligned}
& \min \quad s \\
& \text{s.t.} \quad \begin{pmatrix} \mathbf{I} & \mathbf{A} \text{diag}(\mathbf{x}) \\ \text{diag}(\mathbf{x})\mathbf{A}^T & s\mathbf{I} \end{pmatrix} \succeq 0, \\
& \quad \begin{pmatrix} \mathbf{I} & \text{diag}(\mathbf{y}) \\ \text{diag}(\mathbf{y}) & \mathbf{A}^T \mathbf{A} \end{pmatrix} \succeq 0, \\
& \quad \begin{pmatrix} (x_i + y_i)/2 \\ (y_i - x_i)/2 \\ 1 \end{pmatrix} \in \mathbb{Q}^3, \quad i = 1, \dots, n, \\
& \quad (d_i - \eta)t \leq x_i \leq (d_i + \eta)t, \quad i = 1, \dots, n, \\
& \quad t \geq 0, \mathbf{x}, \mathbf{y} \geq \mathbf{0},
\end{aligned} \tag{4.1}$$

where $K \succeq 0$ denotes positive semidefiniteness, \mathbb{Q}^n is the n -dimensional second-order cone, $d_i = \frac{1}{\|\mathbf{A}_{:,i}\|}$ for all $i = 1, \dots, n$, and $\eta > 0$ is large enough.

With respect to the ω -condition number, Dennis and Wolkowicz [1] show that the simple scaling diagonal preconditioner (see [11]) is the optimal preconditioner. Doan et al. [2] extend the result for block diagonal preconditioners. The details are shown in the following proposition.

Proposition 3. *The following statements are true.*

1. *The measure ω is pseudoconvex² on the set of s.p.d. matrices, and thus any stationary point is a global minimizer of ω .*
2. *Let \mathbf{A} be a full rank $m \times n$ matrix, $n \leq m$. Then the optimal column scaling that minimizes the measure ω , i.e.*

$$\min \omega((\mathbf{A}\mathbf{D})^T(\mathbf{A}\mathbf{D})),$$

over \mathbf{D} positive, diagonal, is given by

$$D_{ii} = \frac{1}{\|\mathbf{A}_{:,i}\|}, i = 1, \dots, n,$$

where $\mathbf{A}_{:,i}$ is the i -th column of \mathbf{A} .

3. *Let \mathbf{A} be a full rank $m \times n$ matrix, $n \leq m$ with block structure $\mathbf{A} = [\mathbf{A}_1 \ \mathbf{A}_2 \ \dots \ \mathbf{A}_k]$, $\mathbf{A}_i \in \mathbb{R}^{m \times n_i}$. Then an optimal corresponding block diagonal scaling*

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_1 & 0 & 0 & \dots & 0 \\ 0 & \mathbf{D}_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \mathbf{D}_k \end{bmatrix}, \quad \mathbf{D}_i \in \mathbb{R}^{n_i \times n_i},$$

that minimizes the measure ω , i.e.

$$\min \omega((\mathbf{A}\mathbf{D})^T(\mathbf{A}\mathbf{D})),$$

²Note that a function is pseudoconvex if

$$(\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) \geq 0 \implies f(\mathbf{y}) \geq f(\mathbf{x}),$$

and for pseudoconvex functions, all stationary solutions are global minimizers (see for example, [10]).

over \mathbf{D} block diagonal, is given by the factorization

$$\mathbf{D}_i \mathbf{D}_i^T = \{\mathbf{A}_i^T \mathbf{A}_i\}^{-1}, i = 1, \dots, k.$$

Proof. Items 1 to 2 are proved in [1] and Item 3 is proved in [2]. For clarity, we present the proof again here for the general block diagonal preconditioners (which include the diagonal case).

Let the blocked \mathbf{A} be given. Then the arithmetic-geometric mean inequality yields

$$\begin{aligned} \omega((\mathbf{AD})^T(\mathbf{AD})) &= \frac{\text{trace}(\mathbf{D}^T \mathbf{A}^T \mathbf{AD})/n}{\det(\mathbf{D}^T \mathbf{A}^T \mathbf{AD})^{\frac{1}{n}}} \\ &= \frac{\text{trace}(\mathbf{A}^T \mathbf{A} \mathbf{D} \mathbf{D}^T)/n}{\det(\mathbf{D} \mathbf{D}^T)^{\frac{1}{n}} \det(\mathbf{A}^T \mathbf{A})^{\frac{1}{n}}}. \end{aligned} \quad (4.2)$$

Let $\bar{\mathbf{D}} = \mathbf{D} \mathbf{D}^T$, $\bar{\mathbf{A}}_i = \mathbf{A}_i^T \mathbf{A}_i$. We have $\bar{\mathbf{D}}$ is also a block diagonal matrix with $\bar{\mathbf{D}}_i = \mathbf{D}_i \mathbf{D}_i^T$ for all $i = 1, \dots, k$, and $\text{trace}(\mathbf{A}^T \mathbf{A} \mathbf{D} \mathbf{D}^T) = \text{trace} \sum_{i=1}^k \bar{\mathbf{A}}_i \bar{\mathbf{D}}_i$. Since the minimum of the function ω in (4.2) is scale free, we can ignore the constants $1/n, \det(\mathbf{A}^T \mathbf{A})^{\frac{1}{n}}$, take logs, and equivalently solve

$$\max \left\{ \sum_{i=1}^k \log \det \bar{\mathbf{D}}_i : \sum_{i=1}^k \bar{\mathbf{A}}_i \cdot \bar{\mathbf{D}}_i = 1 \right\}.$$

Using a Lagrange multiplier λ and Cramer's rule applied to the gradient of the objective function, we get an implicit expression for the optimum in

$$\bar{\mathbf{D}}_i^{-1} - \lambda \bar{\mathbf{A}}_i = 0, \quad i = 1, \dots, k.$$

□

Note that the optimality condition only requires $\bar{\mathbf{D}}_i = \bar{\mathbf{A}}_i^{-1}$ for $i = 1, \dots, k$, which means we can use the QR decomposition to find the optimal \mathbf{D}_i instead of the matrix square root. In other words, if \mathbf{A}_i has the QR decomposition $\mathbf{A}_i = \mathbf{Q}_i \mathbf{R}_i$, we can set $\mathbf{D}_i = \mathbf{R}_i^{-1}$ for all $i = 1, \dots, k$ to minimize the measure $\omega((\mathbf{AD})^T(\mathbf{AD}))$. We let $\mathbf{D}_\omega, \mathbf{D}_\kappa$ denote the optimal column scaling for the measures ω, κ , respectively.

We now generate random matrices \mathbf{A} of size 1000 by 500 with random singular values. Figures 5 and 6 show that \mathbf{D}_ω and \mathbf{D}_κ indeed reduces the ω -condition number and the standard condition number, respectively. We continue and generate random *sparse* matrices with random singular values. We then consider the overdetermined system $\mathbf{Ax} = \mathbf{b}$ with no preconditioner, the diagonal preconditioner \mathbf{D}_ω , and the diagonal preconditioner \mathbf{D}_κ . We solve 1000 random sparse matrix instances with the least squares solver LSQR in MATLAB to achieve the accuracy of 10^{-10} . Figure 7 shows that both preconditioners reduce the number of iterations significantly with the average change of -58.77% . The preconditioner \mathbf{D}_κ is slightly better than \mathbf{D}_ω with about 0.84% reduction in number of iterations. However, we can see that the preconditioner \mathbf{D}_ω is much easier to compute as compared to the preconditioner \mathbf{D}_κ , which requires solving a conic optimization problem of size $2n$.

We now test the reduction in number of iterations and compare the computational times. We again generate random (sparse) matrices of size 1000 by 500 and solve the system with no preconditioner ($p = 0$), and block diagonal preconditioners with the block size $p = 1, 2, 5, 10, 20, 25, 50, 100$,

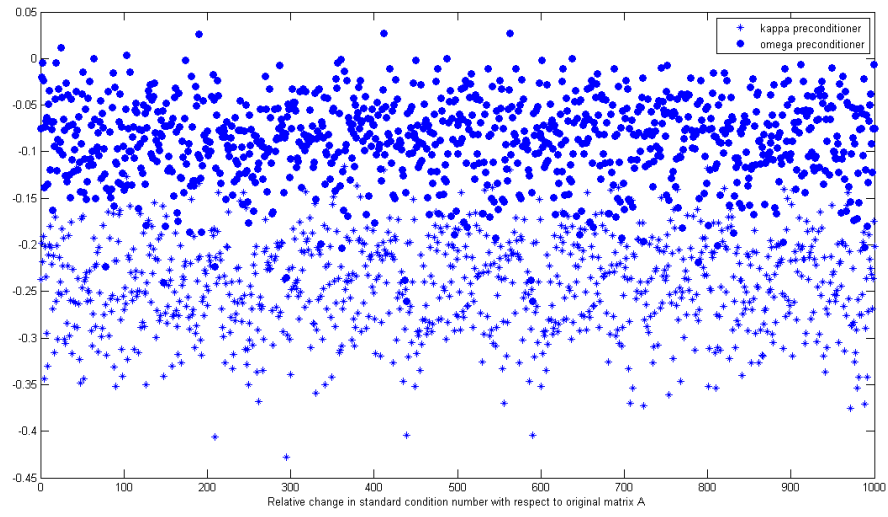


Figure 5: κ with two different diagonal preconditioners

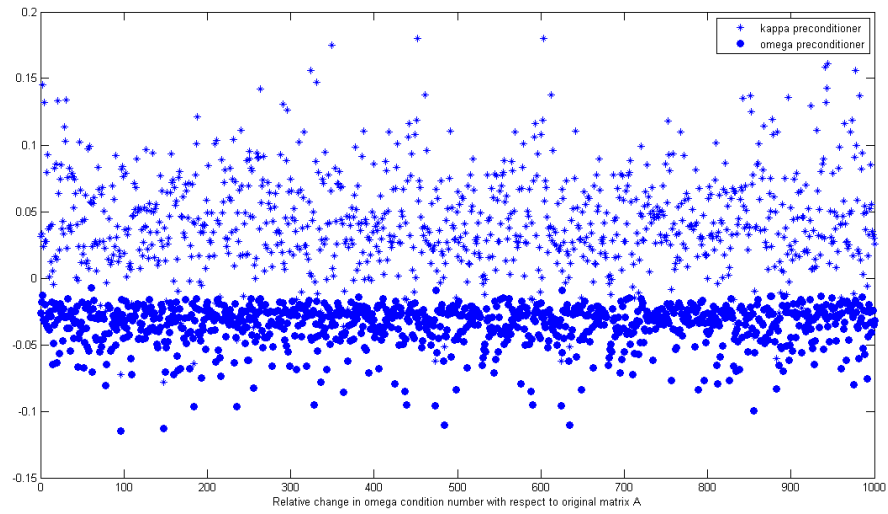


Figure 6: ω with two different diagonal preconditioners

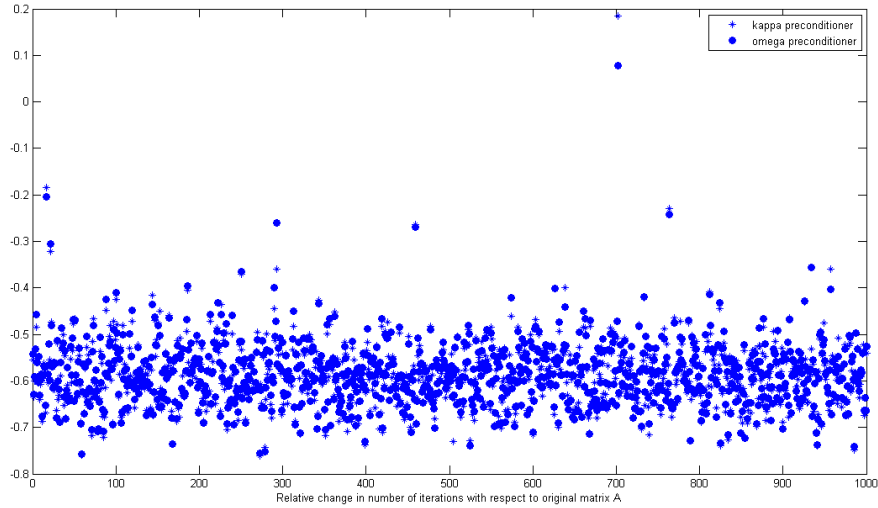


Figure 7: # Krylov iterations with two different diagonal preconditioners for sparse matrices

and 250. The results are shown in Table 4.1 for random sparse matrices, where n_i is the number of iterations, t_p is the time to compute the preconditioner, and t_t is the total time to find the least squares solution. We can see that block diagonal preconditioners improve the number of iterations and also the computation time.

p	0	1	2	5	10	20	25	50	100	250
$n_i(\times 10^4)$	5.89	1.31	1.01	0.66	0.53	0.46	0.45	0.38	0.30	0.11
t_p	-	2.46	1.36	0.70	0.50	0.42	0.42	0.42	0.44	0.26
t_t	39.00	14.50	11.15	7.94	6.76	6.28	6.20	5.89	5.60	4.14

Table 4.1: Average # iterations; computational time block preconditioners; sparse matrices

5 ω -Condition Number and Relative Error

We now focus our attention again for $\mathbf{A} \in \mathcal{S}_{++}^n$, positive definite, and the system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{b}$. We would like to consider the following ratio of the relative errors

$$r_e(\mathbf{A}, \mathbf{b}, \Delta\mathbf{b}) := \frac{\|\mathbf{A}^{-1}\Delta\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} \left(\frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} \right)^{-1}. \quad (5.1)$$

The condition number $\kappa(\mathbf{A})$ provides an upper bound of this ratio

$$\kappa(\mathbf{A}) = \max_{\mathbf{b}, \Delta\mathbf{b}} r_e(\mathbf{A}, \mathbf{b}, \Delta\mathbf{b}).$$

The bound is attained when $\mathbf{b} = \beta\mathbf{v}_1$ and $\Delta\mathbf{b} = \Delta\beta\mathbf{v}_n$, where \mathbf{v}_1 and \mathbf{v}_n are the eigenvectors corresponding to the largest and smallest eigenvectors. Clearly, if \mathbf{b} and $\Delta\mathbf{b}$ are arbitrary vectors,

this bound is unlikely to be attained. We would like to consider the expected value of $r_e(\mathbf{A}, \mathbf{b}, \Delta\mathbf{b})$ when \mathbf{b} and $\Delta\mathbf{b}$ are random vectors. We assume that \mathbf{b} and $\Delta\mathbf{b}$ are independent normal random vectors, $\mathbf{b}, \Delta\mathbf{b} \sim N(\mathbf{0}, \mathbf{I})$. Let us consider the expectation

$$r(\mathbf{A}) = \mathbb{E}[r_e(\mathbf{A}, \mathbf{b}, \Delta\mathbf{b})].$$

The following proposition provides the analytic formulation for $r(\mathbf{A})$.

Proposition 4. *If \mathbf{b} and $\Delta\mathbf{b}$ are independent normal random vectors, then $r(\mathbf{A})$ can be calculated as follows:*

$$r(\mathbf{A}) = \left(\sum_{i=0}^{\infty} \frac{(-1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}) \right) \left(\sum_{i=0}^{\infty} \frac{(1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}) \right), \quad (5.2)$$

where $\mathbf{B} = (\mathbf{A}^{-1})^T \mathbf{A}^{-1}$, $\beta = \lambda_{\min}^2(\mathbf{A})$, $(a)_b = \frac{\Gamma(a+b)}{\Gamma(a)}$, and $d_k(\mathbf{A})$ is the normalized top-order zonal polynomial,

$$d_k(\mathbf{A}) = \sum_{k_i \geq 0: \sum_{i=1}^n k_i = k} \left(\prod_{i=1}^n \frac{\left(\frac{1}{2}\right)_{k_i} \lambda_i(\mathbf{A})^{k_i}}{k_i!} \right), \quad \forall k \geq 0.$$

Proof. We have: since \mathbf{b} and $\Delta\mathbf{b}$ are independent,

$$r(\mathbf{A}) = \mathbb{E} \left[\frac{\|\mathbf{A}^{-1} \Delta\mathbf{b}\|}{\|\Delta\mathbf{b}\|} \right] \mathbb{E} \left[\frac{\|\mathbf{b}\|}{\|\mathbf{A}^{-1} \mathbf{b}\|} \right]. \quad (5.3)$$

These two expectations has the same form; they are the expectations of ratios of quadratic forms, $\frac{(\mathbf{x}^T \mathbf{A} \mathbf{x})^p}{(\mathbf{x}^T \mathbf{B} \mathbf{x})^q}$, where p and q are positive real number, $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$. Smith [12] shows the exact formulation for the expectation with $\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{Q})$ which involves top-order invariant polynomials. The formulation can be simplified if $\boldsymbol{\mu} = \mathbf{0}$ and $\mathbf{B} = \mathbf{I}$ in which top-order invariant polynomials are replaced by simpler top-order zonal polynomials (see Smith [12] for details). In our particular case, we have: $p = q = \frac{1}{2}$ since for arbitrary \mathbf{A} and \mathbf{x} , $\|\mathbf{A}\mathbf{x}\| = (\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x})^{\frac{1}{2}}$. In addition, there is at least an identity matrix, either in the denominator or numerator. Deriving directly from the general formulation presented in Smith [12], we obtain the following formulation for $\mathbb{E} \left[\frac{(\Delta\mathbf{b}^T \mathbf{B} \Delta\mathbf{b})^{\frac{1}{2}}}{(\Delta\mathbf{b}^T \Delta\mathbf{b})^{\frac{1}{2}}} \right]$, which only involves top-order zonal polynomials:

$$\mathbb{E} \left[\frac{(\Delta\mathbf{b}^T \mathbf{B} \Delta\mathbf{b})^{\frac{1}{2}}}{(\Delta\mathbf{b}^T \Delta\mathbf{b})^{\frac{1}{2}}} \right] = \beta^{-\frac{1}{2}} \sum_{i=0}^{\infty} \frac{(-1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}),$$

where $\mathbf{B} = (\mathbf{A}^{-1})^T \mathbf{A}^{-1}$ and $\beta = \lambda_{\min}^2(\mathbf{A})$. Similarly, we have:

$$\mathbb{E} \left[\frac{(\mathbf{b}^T \mathbf{b})^{\frac{1}{2}}}{(\mathbf{b}^T \mathbf{B} \mathbf{b})^{\frac{1}{2}}} \right] = \beta^{\frac{1}{2}} \sum_{i=0}^{\infty} \frac{(1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}).$$

Thus we obtain the exact formulation of $r(\mathbf{A})$ as shown in (5.2). \square

The formulation (5.2) involves infinite sums. Following the approach in Hillier et al. [7], we now attempt to find an upper bound for $r(\mathbf{A})$ with partial sums. We have: $(-1/2)_k < 0$ for all $k \geq 1$. In addition, $\mathbf{C} = \mathbf{I} - \beta\mathbf{B} \succeq 0$; therefore, $d_k(\mathbf{C}) \geq 0$ for all $k \geq 1$. Thus for $M \geq 0$,

$$\sum_{i=0}^{\infty} \frac{(-1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}) \leq \sum_{i=0}^M \frac{(-1/2)_i}{(n/2)_i} d_i(\mathbf{C}).$$

We have: $n \geq 1$, thus $\frac{(1/2)_k}{(n/2)_k} \leq \frac{(1/2)_{M+1}}{(n/2)_{M+1}}$ for all $k \geq M+1$. Using the generating function of $d_k(\mathbf{C})$, we have:

$$[\det(\mathbf{I} - t\mathbf{C})]^{-\frac{1}{2}} = \sum_{k=0}^{\infty} d_k(\mathbf{C}) t^k.$$

Let $t = 1$, we have:

$$[\det(\beta\mathbf{B})]^{-\frac{1}{2}} = \sum_{k=0}^{\infty} d_k(\mathbf{C}).$$

Thus for $M \geq 0$,

$$\sum_{i=0}^{\infty} \frac{(1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}) \leq \sum_{i=0}^M \frac{(1/2)_i}{(n/2)_i} d_i(\mathbf{C}) + \frac{(1/2)_{M+1}}{(n/2)_{M+1}} \left[\frac{1}{\beta^{\frac{n}{2}} [\det(\mathbf{B})]^{\frac{1}{2}}} - \sum_{i=0}^M d_i(\mathbf{C}) \right].$$

We then have the following upper bound for $r(\mathbf{A})$ for each $M \geq 0$:

$$r(\mathbf{A}) \leq \left(\sum_{i=0}^M \frac{(-1/2)_i}{(n/2)_i} d_i(\mathbf{C}) \right) \left(\sum_{i=0}^M \frac{(1/2)_i}{(n/2)_i} d_i(\mathbf{C}) + \frac{(1/2)_{M+1}}{(n/2)_{M+1}} \left[\frac{1}{\beta^{\frac{n}{2}} [\det(\mathbf{B})]^{\frac{1}{2}}} - \sum_{i=0}^M d_i(\mathbf{C}) \right] \right) \quad (5.4)$$

Let $M = 1$, we have: $d_0(\mathbf{C}) = 1$ and $d_1(\mathbf{C}) = \frac{\text{trace}(\mathbf{C})}{2}$. In addition, eigenvalues of \mathbf{C} are $\gamma_i = 1 - \frac{\lambda_n^2}{\lambda_i^2}$. Thus

$$\sum_{i=0}^{\infty} \frac{(-1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}) \leq 1 - \frac{\text{trace}(\mathbf{C})}{2n} = \frac{1}{2} \left(1 + \frac{\lambda_n^2}{n} \left(\sum_{i=1}^n \frac{1}{\lambda_i^2} \right) \right).$$

We also have:

$$\sum_{i=0}^{\infty} \frac{(1/2)_i}{(n/2)_i} d_i(\mathbf{I} - \beta\mathbf{B}) \leq \frac{3(n-1)}{2n} - \frac{n-1}{2(n+2)} \left(\frac{\lambda_n^2}{n} \right) \left(\sum_{i=1}^n \frac{1}{\lambda_i^2} \right) + \frac{3}{n(n+2)} \prod_{i=1}^n \left(\frac{\lambda_i}{\lambda_n} \right).$$

We obtain an upper bound for $r(\mathbf{A})$:

$$r(\mathbf{A}) \leq \frac{1}{2} \left[1 + \frac{\lambda_n^2}{n} \left(\sum_{i=1}^n \frac{1}{\lambda_i^2} \right) \right] \left[\frac{3(n-1)}{2n} - \frac{n-1}{2(n+2)} \left(\frac{\lambda_n^2}{n} \right) \left(\sum_{i=1}^n \frac{1}{\lambda_i^2} \right) + \frac{3}{n(n+2)} \prod_{i=1}^n \left(\frac{\lambda_i}{\lambda_n} \right) \right].$$

Unfortunately, this upper bound in general is not very tight and there are examples in which we need very large M to obtain a decent approximation of $r(\mathbf{A})$. We now empirically calculate the sample

approximation $\bar{r}(\mathbf{A})$ of $r(\mathbf{A})$ with a large number of samples and check the relationship between $\bar{r}(\mathbf{A})$ and ω -condition number of \mathbf{A} as well as the standard condition number. We generate random matrices of the size 1000 by 1000 and we calculate $\bar{r}(\mathbf{A})$ with $N = 10,000$ samples. Figures 8 and 9 show the ratios of the ω -condition number and standard condition number of \mathbf{A} , respectively, to the sample approximation of $r(\mathbf{A})$. The maximum value of ratios is limited to 10 in both figures for the purpose of comparison.

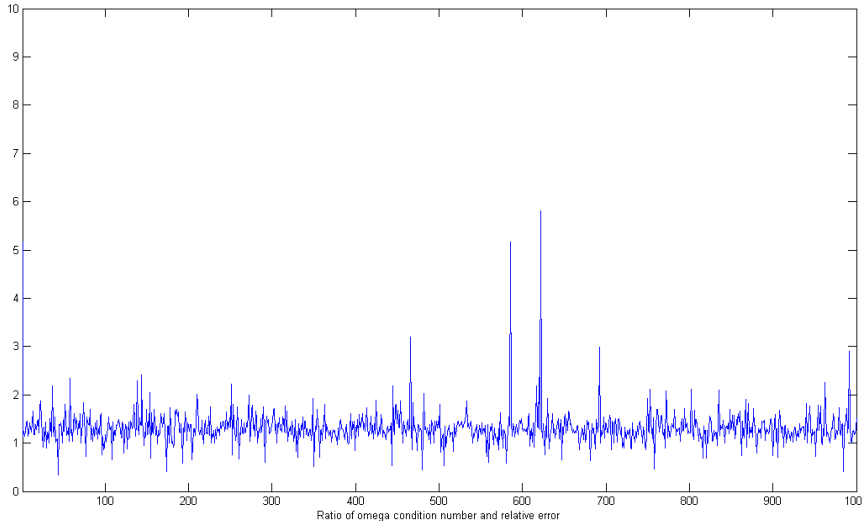


Figure 8: Ratio of ω and the relative error $\bar{r}(\mathbf{A})$

Clearly, the standard condition number is an upper bound of $\bar{r}(\mathbf{A})$; however, the average standard condition number is 1.03×10^3 while the maximum relative error is just 4.41. On the other hand, the ω -condition number is not always an upper bound of $\bar{r}(\mathbf{A})$ but it is in more than 90% of all instances. In addition, the average ω -condition number is 1.49, which is very close to the average relative error of 1.18. The average ratio of ω -condition number is 1.31 with the standard deviation of 0.35 while these two measures of κ -condition number are 487.92 and 3.82×10^3 , respectively. We can conclude that the ω -condition number is a good indicator for the expected relative error $r(\mathbf{A})$ in this random setting.

6 Conclusion

We have presented computational and theoretical properties to show that the condition number measure $\omega(A)$ is a better indicator than the standard condition number κ for the conditioning of A when it comes to: (i) evaluation, (ii) predicting the number of iterations in Krylov iterative methods, (iii) finding optimal preconditioners for iterative methods, and (iv) estimating the *expected* roundoff error when solving linear equations.

In addition, we have shown that ω is a simpler function (measure) to use when finding optimal preconditioners, as it is easy to differentiate. And, motivated by the fact that the number of iterations in Krylov type methods depends directly on the number of distinct eigenvalues, we

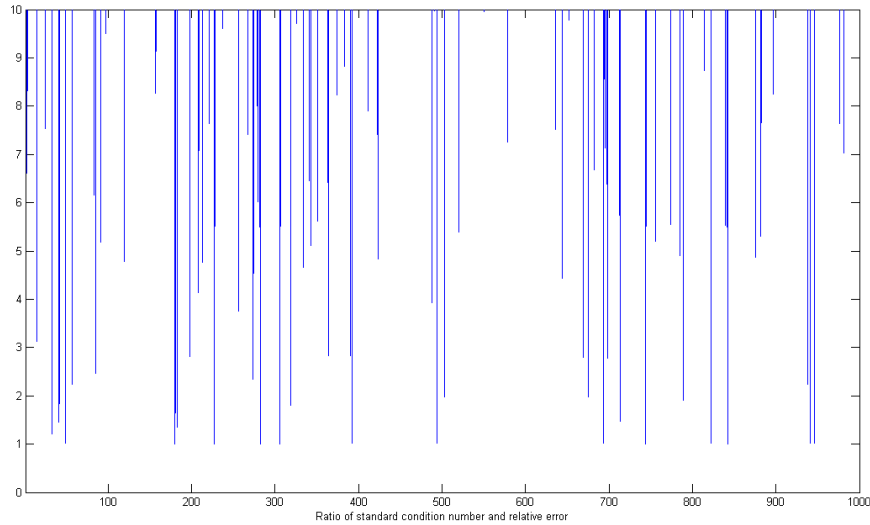


Figure 9: Ratio of κ and the relative error $\bar{r}(\mathbf{A})$

showed empirically that if κ is fixed, then both small and large values of ω lead to lower iteration numbers; this follows since both maximizing and minimizing ω leads to clustering of eigenvalues when κ is fixed. Therefore, an interesting conjecture is that: maximizing ω subject to keeping κ constant leads to good preconditioners. However, the problem of maximizing the pseudoconvex function ω is difficult and so heuristics will have to be used.

References

- [1] J.E. Dennis Jr. and H. Wolkowicz. Sizing and least-change secant methods. *SIAM J. Numer. Anal.*, 30(5):1291–1314, 1993. 3, 10, 11
- [2] X.V. Doan, S. Kruk, and H. Wolkowicz. A robust algorithm for semidefinite programming. *Optim. Methods Softw.*, (CORR 2010-09):., 2011. submitted in November, 2010, accepted June, 2011. 10, 11
- [3] R.G. Grimes and J.G. Lewis. Condition number estimation for sparse matrices. *SIAM J. Sci. Statist. Comput.*, 2:384–388, 1981. 4
- [4] W. W. Hager. Condition estimates. *SIAM J. Sci. Statist. Comput.*, 5(2):311–316, 1984. 4
- [5] N.J. Higham. A survey of condition number estimation for triangular matrices. *SIAM Rev.*, 29:575–596, 1987. 4
- [6] N.J. Higham and F. Tisseur. A block algorithm for matrix 1-norm estimation, with an application to 1-norm pseudospectra. *SIAM J. Matrix Anal. Appl.*, 21(4):1185–1201 (electronic), 2000. 4
- [7] G. Hillier, R. Kan, and X. Wang. Computationally efficient recursions for top-order invariant polynomials with applications. *Econometric Theory*, 25:211–242, 2009. 15
- [8] Ilse C. F. Ipsen and Carl D. Meyer. The idea behind Krylov methods. *Amer. Math. Monthly*, 105(10):889–899, 1998. 4
- [9] Z. Lu and T.K. Pong. Minimizing condition number via convex programming. Technical report, University of Washington, Seattle, WA, 2010. 8

- [10] O.L. Mangasarian. *Nonlinear Programming*. McGraw-Hill, New York, NY, 1969. 10
- [11] G. Pini and G. Gamboniati. Is a simple diagonal scaling the best preconditioner for conjugate gradients on supercomputers? *Advances in Water Resources*, 13(3):147–153, 1990. 10
- [12] Murray D. Smith. Expectations of ratios of quadratic forms in normal variables: evaluating some top-order invariant polynomials. *Austral. J. Statist.*, 35(3):271–282, 1993. 14

Index

$\mathcal{A}(x)$, arithmetic mean, 4
 $\mathcal{G}(x)$, geometric mean, 4
 $\kappa(\mathbf{A})$, standard condition number, 3
 D_κ , optimal diagonal scaling, 11
 D_ω , optimal diagonal scaling, 11
 $\omega(\mathbf{A})$, omega condition number, 3
 $\omega(\boldsymbol{\lambda})$, 6
 $\omega_2(\mathbf{A})$, 4
 $\omega_c(\mathbf{A})$, 4
 $r(\mathbf{A}) = \mathbb{E}[r_e(\mathbf{A}, \mathbf{b}, \Delta\mathbf{b})]$, expectation of relative errors, 14
 $r_e(\mathbf{A}, \mathbf{b}, \Delta\mathbf{b})$, ratio of relative errors, 13
 \mathcal{S}^n , space of symmetric matrices, 3
 \mathcal{S}_{++}^n , cone of positive definite matrices, 3
 \mathcal{S}_+^n , cone of positive semidefinite matrices, 3

algorithm for $\omega(A)$, 4
arithmetic mean, $\mathcal{A}(x)$, 4

condition number, 3

geometric mean, $\mathcal{G}(x)$, 4

Krylov space generated by \mathbf{c} , 4

optimal block diagonal preconditioner, 10

pseudoconvex function, 10