



香港城市大學
City University of Hong Kong

專業 創新 胸懷全球
Professional · Creative
For The World

CityU Scholars

Componentwise condition numbers of random sparse matrices

Cheung, Dennis; Cucker, Felipe

Published in:

SIAM Journal on Matrix Analysis and Applications

Published: 01/01/2009

Document Version:

Final Published version, also known as Publisher's PDF, Publisher's Final version or Version of Record

Publication record in CityU Scholars:

[Go to record](#)

Published version (DOI):

[10.1137/080729463](https://doi.org/10.1137/080729463)

Publication details:

Cheung, D., & Cucker, F. (2009). Componentwise condition numbers of random sparse matrices. *SIAM Journal on Matrix Analysis and Applications*, 31(2), 721-731. <https://doi.org/10.1137/080729463>

Citing this paper

Please note that where the full-text provided on CityU Scholars is the Post-print version (also known as Accepted Author Manuscript, Peer-reviewed or Author Final version), it may differ from the Final Published version. When citing, ensure that you check and use the publisher's definitive version for pagination and other details.

General rights

Copyright for the publications made accessible via the CityU Scholars portal is retained by the author(s) and/or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights. Users may not further distribute the material or use it for any profit-making activity or commercial gain.

Publisher permission

Permission for previously published items are in accordance with publisher's copyright policies sourced from the SHERPA RoMEO database. Links to full text versions (either Published or Post-print) are only available if corresponding publishers allow open access.

Take down policy

Contact lbscholars@cityu.edu.hk if you believe that this document breaches copyright and provide us with details. We will remove access to the work immediately and investigate your claim.

© 2009 Society for Industrial and Applied Mathematics.

COMPONENTWISE CONDITION NUMBERS OF RANDOM SPARSE MATRICES*

DENNIS CHEUNG[†] AND FELIPE CUCKER[‡]

Abstract. We prove an $\mathcal{O}(\log n)$ bound for the expected value of the logarithm of the componentwise (and, a fortiori, the mixed) condition number of a random sparse $n \times n$ matrix. As a consequence, small bounds on the average loss of accuracy for triangular linear systems follow.

Key words. sparse matrices, triangular systems, average loss of precision

AMS subject classifications. Primary, 65F35; Secondary, 65G50

DOI. 10.1137/080729463

1. Introduction. Triangular systems of linear equations provide one of the few examples in numerical linear algebra where a gap occurs between stability analysis and everyday practice. One could summarize this gap as follows:

Triangular systems of equations are generally solved to high accuracy in spite of being, in general, ill-conditioned.

This state of affairs had been already noted by Wilkinson in [9, p. 105]: “In practice one almost invariably finds that if L is ill-conditioned, so that $\|L\|\|L^{-1}\| \gg 1$, then the computed solution of $Lx = b$ (or the computed inverse) is far more accurate than [what forward stability analysis] would suggest.”

An explanation for this gap is suggested by Higham [4] who notes that the backward error analysis given by Wilkinson for the solution of triangular systems yields (small) *componentwise* bounds on the perturbation matrix (see section 6 item (1) below). Higham then uses this fact to deduce small forward error bounds for particular subclasses of triangular systems and to numerically investigate the accuracy of other particular such systems. In doing so, Higham makes use of the mixed condition number introduced by Skeel [5]¹. This condition number has a natural role in analyzing accuracy of triangular systems since bounds for it, together with the backward analysis of Wilkinson mentioned above, yield forward analysis bounds for the computed solution of the system. Furthermore, the restriction to componentwise perturbations—both in the backward error analysis and in the mixed condition number—forces perturbations to preserve the triangular structure of the data matrices.

A further step in explaining the gap, somehow orthogonal to the work of Higham, was given by Viswanath and Trefethen in [8] where a precise meaning to the expression “triangular systems are, in general, ill-conditioned” was given. Indeed, if L_n denotes a random triangular $n \times n$ matrix (whose entries are independent standard Gaussian

*Received by the editors July 7, 2008; accepted for publication (in revised form) by R.-C. Li March 19, 2009; published electronically June 17, 2009.

<http://www.siam.org/journals/simax/31-2/72946.html>

[†]United International College, Tang Jia Wan, Zhuhai, Guangdong Province, People’s Republic of China (dennisc@uic.edu.hk).

[‡]Department of Mathematics, City University of Hong Kong, 83 Tat Chee Avenue, Kowloon, Hong Kong (macucker@cityu.edu.hk). This author’s research was partially supported by GRF grant CityU 100808.

¹Skeel called it “componentwise.” In this paper, however, following the notation introduced in [3], we will use this word for the condition numbers measuring both data perturbation and computed errors in a componentwise fashion.

random variables) and $\kappa_n = \|L_n\| \|L_n^{-1}\|$ is its condition number (which is a positive random variable), then it is shown in [8] that

$$\sqrt[n]{\kappa_n} \rightarrow 2 \quad \text{almost surely}$$

as $n \rightarrow \infty$. A straightforward consequence of this result is that the expected value of $\log \kappa_n$ satisfies $\mathbf{E}(\log \kappa_n) = \Omega(n)$.

The goal of this paper is to close the gap by giving a precise meaning, in the sense of [8], to the other half of the statement above, namely, to the expression “triangular systems are generally solved to high accuracy.” More precisely, we consider the mixed condition numbers $\mathbf{m}^\dagger(L_n)$ and $\mathbf{m}(L_n, b_n)$ for the problems of matrix inversion and linear equation solving, respectively, for a random triangular L_n as above and a random $b_n \in \mathbb{R}^n$. Then we show that

$$(1) \quad \mathbf{E}(\log \mathbf{m}^\dagger(L_n)), \mathbf{E}(\log \mathbf{m}(L_n, b_n)) = \mathcal{O}(\log n).$$

From the bound on $\mathbf{E}(\log \mathbf{m}(L_n, b_n))$ it follows that the average loss of precision in the solution of random triangular systems is small. From that on $\mathbf{E}(\log \mathbf{m}^\dagger(L_n))$, the one for matrix inversion is small as well. One can therefore replace the summary above by the following:

Triangular systems of equations are generally solved to high accuracy because their backward error analysis yields small componentwise perturbations and triangular matrices are, in general, well conditioned for these perturbations.

The results showing (1), Theorems 2 and 3 below, are proved in the more general context of sparse matrices (which, in this paper, are matrices with a fixed pattern of zeros²) and componentwise condition numbers (which ensure high relative accuracy in each component of the computed solution A^{-1} or x). Besides triangular matrices, these results apply to other classes of sparse matrices such as, for instance, tridiagonal matrices. In the process of proving them, we found it useful to estimate as well the average mixed condition for the computation of the determinant.

2. Preliminaries. Condition numbers measure the worst-case magnification of a small data perturbation in the computed outcome. As originally introduced by Turing [7], they were *normwise* in the sense that both the data perturbation and the outcome’s error are measured using norms (in the space of data and outcomes, respectively). In contrast, *mixed* condition numbers measure data perturbation componentwise, and *componentwise* condition numbers measure both data perturbation and the outcome’s error in this way.

To define these condition numbers the following form of “distance” function will be useful. For points $u, v \in \mathbb{R}^p$ we define $\frac{u}{v} = (w_1, \dots, w_p)$ with

$$w_i = \begin{cases} u_i/v_i & \text{if } v_i \neq 0, \\ 0 & \text{if } u_i = v_i = 0, \\ \infty & \text{otherwise.} \end{cases}$$

Then we define

$$d(u, v) = \left\| \frac{u - v}{v} \right\|_\infty.$$

²The word “sparse” is also used to denote matrices with a large number of zeros, not necessarily in fixed positions.

Note that if $d(u, v) < \infty$,

$$d(u, v) = \min\{\nu \geq 0 \mid |u_i - v_i| \leq \nu |v_i| \text{ for } i = 1, \dots, p\}.$$

For $\delta > 0$ and $a \in \mathbb{R}^p$ we denote $\mathcal{S}(a, \delta) = \{x \in \mathbb{R}^p \mid d(x, a) = \delta\}$.

DEFINITION 1. Let $\mathcal{D} \subseteq \mathbb{R}^p$ and $F : \mathcal{D} \rightarrow \mathbb{R}^q$ be a continuous mapping. Let $a \in \mathcal{D}$ such that $F(a) \neq 0$.

- (i) The mixed condition number of F at a (with respect to a norm $\|\cdot\|_q$ on \mathbb{R}^q) is defined by

$$\mathbf{m}(F, a) = \lim_{\delta \rightarrow 0} \sup_{x \in \mathcal{S}(a, \delta)} \frac{\|F(x) - F(a)\|_q}{\|F(a)\|_q} \frac{1}{d(x, a)}.$$

- (ii) Suppose $F(a) = (f_1(a), \dots, f_q(a))$ is such that $f_j(a) \neq 0$ for $j = 1, \dots, q$. Then the componentwise condition number of F at a is

$$\mathbf{c}(F, a) = \lim_{\delta \rightarrow 0} \sup_{x \in \mathcal{S}(a, \delta)} \frac{d(F(x), F(a))}{d(x, a)}.$$

Remark 1. We can extend the definition of componentwise condition numbers to the case where $f_i(a) = 0$ for some $i \in [n]$. Define, for $i \in [n]$ such that $f_i(a) \neq 0$,

$$\mathbf{c}_i(F, a) = \lim_{\delta \rightarrow 0} \sup_{x \in \mathcal{S}(a, \delta)} \frac{|f_i(x) - f_i(a)|}{\delta |f_i(a)|},$$

and for $i \in [n]$ with $f_i(a) = 0$ take $\mathbf{c}_i(F, a) = 0$ if

$$\lim_{\delta \rightarrow 0} \sup_{x \in \mathcal{S}(a, \delta)} \frac{|f_i(x) - f_i(a)|}{\delta} = 0$$

and $\mathbf{c}_i(F, a) = \infty$ otherwise. Then we have

$$\mathbf{c}(F, a) = \max_{i \in [n]} \mathbf{c}_i(F, a).$$

PROPOSITION 1. For all $a \in \mathcal{D}$ and any monotonic norm in \mathbb{R}^q , $\mathbf{m}(F, a) \leq \mathbf{c}(F, a)$.

Proof. For all $x \in \mathcal{S}(a, \delta)$ and all $i \leq q$, $|F(x)_i - F(a)_i| \leq d(F(x), F(a)) |F(a)_i|$. Since $\|\cdot\|$ is monotonic (cf. [1]), this implies $\|F(x) - F(a)\| \leq d(F(x), F(a)) \|F(a)\|$ and hence the statement. \square

In what follows, for $n \in \mathbb{N}$, we denote the set $\{1, \dots, n\}$ by $[n]$ and write, as usual, $[n]^2 = [n] \times [n]$.

DEFINITION 2. We denote by \mathcal{M} the set of $n \times n$ real matrices and by Σ its subset of singular matrices. Also, for a subset $S \subseteq [n]^2$ we denote

$$\mathcal{M}_S = \{A \in \mathcal{M} \mid \text{if } (i, j) \notin S, \text{ then } a_{ij} = 0\}.$$

We denote by \mathcal{R}_S the space of random $n \times n$ matrices obtained by setting $a_{ij} = 0$ if $(i, j) \notin S$ and drawing all other entries independently from the standard Gaussian $N(0, 1)$. As above, if $S = [n]^2$, we write simply \mathcal{R} .

In the rest of this paper, for nonsingular matrices A and A' , we denote their inverses by Γ and Γ' , respectively. Also, we denote by $A_{(ij)}$ the submatrix of A obtained by removing from A its i th row and its j th column. Denoting by γ_{ij} the (i, j) th entry of Γ we have, by Cramer's rule, $\gamma_{ij} = \frac{\det(A_{(ij)})}{\det(A)}$.

3. Determinant computation. We consider here the problem of computing the determinant of a matrix A and its componentwise condition number $\mathbf{c}_{\det}(A)$. The main result of this section is the following.

THEOREM 1. For $S \subseteq [n]^2$ and $t \geq 2|S|$ we have

$$\text{Prob}_{A \in \mathcal{M}_S} \{ \mathbf{c}_{\det}(A) \geq t \} \leq |S|^2 \frac{1}{t}.$$

Here $|S|$ denotes the cardinality of S .

Average loss of precision (in a base b) is measured by the expected value of the logarithm (in that base) of the condition number. We may use Theorem 1 to obtain one such result for the computation of the determinant. To avoid problems caused by this condition number being less than 1 we consider the function $\log_+(x)$ defined to be $\log(x)$ if $x \geq 1$ and 0 otherwise.

COROLLARY 1. For a base $b > 1$ and a set $S \subseteq [n]^2$ with $|S| \geq 2$, we have $\mathbf{E}(\log_+ \mathbf{c}_{\det}(A)) \leq 2 \log |S| + \frac{1}{\ln b}$, where \mathbf{E} denotes expectation over $A \in \mathcal{M}_S$.

Note that the restriction $|S| \geq 2$ is without loss of generality. The case $|S| = 0$ reduces \mathcal{M}_S to the zero matrix, and the case $|S| = 1$ allows only for one nonzero entry, a singular situation whenever $n \geq 2$. The proof of Corollary 1 follows from Theorem 1 and the following result by taking $Z = \mathbf{c}_{\det}(A)$ and $t_0 = |S|^2$ and noting that $t_0 \geq 2|S|$ since $|S| \geq 2$.

PROPOSITION 2. Let $t_0 > 0$ and $Z \geq 1$ be a random variable satisfying that $\text{Prob}\{Z \geq t\} \leq t_0 t^{-1}$ for all $t \geq t_0$. Then $\mathbf{E}(\log Z) \leq \log t_0 + \frac{1}{\ln b}$, where $b > 1$ is the base of the logarithm.

Proof. We have $\text{Prob}\{\log Z \geq t\} \leq t_0 b^{-t} = b^{-(t - \log t_0)}$ for all $t > \log t_0$. Therefore,

$$\mathbf{E}(\log Z) = \int_0^\infty \text{Prob}\{\log Z \geq s\} ds \leq \log t_0 + \int_{\log t_0}^\infty b^{-(t - \log t_0)} dt = \log t_0 + \frac{1}{\ln b}. \quad \square$$

Towards the proof of Theorem 1 we first obtain explicit expressions for $\mathbf{c}_{\det}(A)$. We begin by noting that taking $F : \mathcal{M} \rightarrow \mathbb{R}$ to be $F(A) = \det(A)$ in Definition 1 we obtain, for $A \in \mathcal{M} \setminus \Sigma$,

$$\mathbf{c}_{\det}(A) = \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\det(A') - \det(A)|}{\delta |\det(A)|}.$$

Also, for $A \in \Sigma$, we have $\mathbf{c}_{\det}(A) = 0$ if

$$\lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\det(A')|}{\delta} = 0$$

and $\mathbf{c}_{\det}(A) = \infty$ otherwise. Note that $\mathbf{c}_{\det}(0) = 0$.

LEMMA 1. For $A \in \mathcal{M} \setminus \Sigma$,

$$\mathbf{c}_{\det}(A) = \sum_{i, j \in [n]} |a_{ij} \gamma_{ji}|.$$

Proof. Let $A \in \mathcal{M}$. For any $i \in [n]$, expanding by the i th row,

$$\det(A) = \sum_{j \in [n]} (-1)^{i+j} a_{ij} \det(A_{(ij)}).$$

Hence, for all $i, j \in [n]$,

$$\frac{\partial \det(A)}{\partial a_{ij}} = (-1)^{i+j} \det(A_{(ij)}).$$

Using Taylor’s expansion and these equalities we obtain

$$\det(A') = \det(A) + \sum_{i,j \in [n]} (-1)^{i+j} (a'_{ij} - a_{ij}) \det(A_{(ij)}) + \mathcal{O}(\|A' - A\|^2).$$

Here the norm in $\|A' - A\|$ is not relevant since all norms in \mathcal{M} are equivalent. By choosing a monotonic norm we have that if $A' \in \mathcal{S}(A, \delta)$, then $\|A' - A\| = \mathcal{O}(\delta)$. It follows that, for $A \notin \Sigma$,

$$\begin{aligned} c_{\det}(A) &= \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\det(A') - \det(A)|}{\delta |\det(A)|} \\ &= \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \sum_{i,j \in [n]} \frac{|(a'_{ij} - a_{ij}) \det(A_{(ij)})|}{\delta |\det(A)|} \\ &= \lim_{\delta \rightarrow 0} \sup \left(\sum_{i,j \in [n]} \frac{|(a'_{ij} - a_{ij}) \det(A_{(ij)})|}{\delta |\det(A)|} : \frac{|a'_{ij} - a_{ij}|}{|a_{ij}|} \leq \delta \right). \end{aligned}$$

The second equality follows from the fact that we can choose A' such that the terms $(a'_{ij} - a_{ij})(-1)^{i+j} \det(A_{(ij)})$ have the same sign for all $i, j \in [n]$. Actually, the supremum above is attained by taking $a'_{ij} = a_{ij}(1 \pm \delta)$ where we take the plus sign if $(-1)^{i+j} \det(A_{(ij)}) \geq 0$ and the minus sign otherwise. Therefore,

$$c_{\det}(A) = \sum_{i,j \in [n]} \left| \frac{a_{ij} \det(A_{(ij)})}{\det(A)} \right| = \sum_{i,j \in [n]} |a_{ij} \gamma_{ji}|. \quad \square$$

We denote by $N(0, \text{Id}_n)$ the n -dimensional Gaussian distribution with mean zero and covariance matrix the identity $n \times n$ matrix. A vector x has this distribution when its entries are independently drawn from $N(0, 1)$.

LEMMA 2. *Let p and q be two fixed vectors in \mathbb{R}^n such that $\|p\| \leq \|q\|$. If $x \sim N(0, \text{Id}_n)$, then, for all $t \geq 2$,*

$$\text{Prob} \left\{ \left| \frac{x^\top p}{x^\top q} \right| \geq t \right\} \leq \frac{1}{t}.$$

Proof. Let $\nu = \|q\|$. By the orthogonal invariance of $N(0, \text{Id}_n)$ we may assume $q = (\nu, 0, \dots, 0)$. Also, by appropriately scaling, we may assume that $\nu = 1$. Note that then $\|p\| \leq 1$. We therefore have

$$\begin{aligned} \text{Prob} \left\{ \left| \frac{x^\top p}{x^\top q} \right| \geq t \right\} &= \text{Prob} \left\{ \left| p_1 + \sum_{i \in \{2, \dots, n\}} \frac{x_i p_i}{x_1} \right| \geq t \right\} \\ (2) \qquad &= \text{Prob} \left\{ \left| p_1 + \frac{1}{x_1} \alpha Z \right| \geq t \right\} \\ &= \text{Prob} \left\{ \frac{Z}{x_1} \geq \frac{t - p_1}{\alpha} \right\} + \text{Prob} \left\{ \frac{Z}{x_1} \leq \frac{-t - p_1}{\alpha} \right\}, \end{aligned}$$

where $Z = N(0, 1)$ is independent of x_1 and $\alpha = \sqrt{p_2^2 + \dots + p_n^2} \leq 1$. Here we used that a sum of independent centered Gaussians is a centered Gaussian whose variance is the sum of the terms’ variances. Note that in case $\alpha = 0$ the statement is trivially true.

Since x_1 and Z are independent, both with distribution $N(0, 1)$, the angle $\theta = \arctan(Z/x_1)$ is uniformly distributed in $[-\pi/2, \pi/2]$ and we have, for $\gamma \in [0, \infty)$,

$$\begin{aligned} \text{Prob} \left\{ \frac{Z}{x_1} \geq \gamma \right\} &= \text{Prob} \{ \theta \geq \arctan \gamma \} = \frac{1}{\pi} \left(\frac{\pi}{2} - \arctan \gamma \right) \\ &= \frac{1}{\pi} \int_{\gamma}^{\infty} \frac{1}{1+t^2} dt \leq \frac{1}{\pi} \int_{\gamma}^{\infty} \frac{1}{t^2} dt = \frac{1}{\pi\gamma}. \end{aligned}$$

Similarly, for $\sigma \in (-\infty, 0]$,

$$\begin{aligned} \text{Prob} \left\{ \frac{Z}{x_1} \leq \sigma \right\} &= 1 - \text{Prob} \{ \theta \geq \arctan \sigma \} = 1 - \frac{1}{\pi} \left(\frac{\pi}{2} - \arctan \sigma \right) \\ &= \frac{1}{\pi} \left(\frac{\pi}{2} - \arctan(-\sigma) \right) \leq \frac{1}{\pi(-\sigma)}. \end{aligned}$$

Using these bounds in (2) with $\gamma = \frac{t-p_1}{\alpha}$ and $\sigma = \frac{-t-p_1}{\alpha}$ we obtain

$$\text{Prob} \left\{ \left| \frac{x^T p}{x^T q} \right| \geq t \right\} \leq \frac{1}{\pi} \left(\frac{\alpha}{t-p_1} + \frac{\alpha}{t+p_1} \right) = \frac{\alpha}{\pi} \frac{2t}{t^2 - p_1^2} \leq \frac{2}{\pi} \frac{t}{t^2 - 1} \leq \frac{1}{t}$$

the last inequality since $t > 2$. \square

LEMMA 3. Let $S \subseteq [n]^2$ be such that $\mathcal{M}_S \subseteq \Sigma$. Then, for all $A \in \mathcal{M}_S$, $\mathbf{c}_{\det}(A) = 0$.

Proof. Since $\mathcal{M}_S \subseteq \Sigma$ and $A \in \mathcal{M}_S$, we have $\mathcal{S}(A, \delta) \subseteq \Sigma$ for all $\delta > 0$. The result now follows. \square

LEMMA 4. Let $S \subset [n]^2$ such that $\mathcal{M}_S \not\subseteq \Sigma$. Then

$$\text{Prob}_{A \in \mathcal{R}_S} (A \text{ is singular}) = 0.$$

Proof. The set of singular matrices in \mathcal{M}_S is the zero set of the restriction of the determinant to \mathcal{M}_S . This restriction is a polynomial in $\mathbb{R}^{|S|}$ whose zero set, if different from $\mathbb{R}^{|S|}$, has dimension smaller than $|S|$. \square

Proof of Theorem 1. Case (i): $\mathcal{M}_S \subseteq \Sigma$. In this case, the desired inequality is trivial by Lemma 3.

Case (ii): $\mathcal{M}_S \not\subseteq \Sigma$. By Lemma 4, with probability 1, A is nonsingular. So, by Lemma 1,

$$\begin{aligned} \text{Prob} \{ \mathbf{c}_{\det}(A) \geq t \} &= \text{Prob} \left\{ \sum_{i,j \in [n]} |a_{ij} \gamma_{ji}| \geq t \right\} \\ (3) \quad &= \text{Prob} \left\{ \sum_{(i,j) \in S} \left| \frac{a_{ij} \det(A_{(ij)})}{\det(A)} \right| \geq t \right\}. \end{aligned}$$

Assume $(1, 1) \in S$, and let $x = a_1$ be the first column of A . Also, let $I = \{i \in [n] \mid (i, 1) \in S\}$ and x_I be the vector obtained by removing entries x_i with $i \notin I$. Then

$$(4) \quad x_I \sim N(0, \text{Id}_{|I|}).$$

For $i \in [n]$ write $q_i = (-1)^{i+1} \det(A_{(i1)})$. Let $q = (q_1, \dots, q_n)^T$ and q_I be the vector obtained by removing entries q_i with $i \notin I$. Clearly, q_I is independent of x_I . Using

this notation, the expansion by the first column yields

$$\det(A) = \sum_{i \in [n]} (-1)^{i+1} a_{i1} \det(A_{(i1)}) = x_I^\top q_I.$$

In addition, $a_{11} \det(A_{(11)}) = x_I^\top (q_1 e_1)$, where e_1 is the vector with the first entry equal to 1 and all others equal to 0. Hence,

$$\frac{a_{11} \det(A_{(11)})}{\det(A)} = \frac{x_I^\top (q_1 e_1)}{x_I^\top q_I}.$$

Let φ be the density of $q_1 = \det(A_{(11)})$ and ψ be that of the random vector q_I . Then, for $z \geq 2$,

$$\begin{aligned} \text{Prob} \left\{ \left| \frac{a_{11} \det(A_{(11)})}{\det(A)} \right| \geq z \right\} &= \text{Prob} \left\{ \left| \frac{x_I^\top (q_1 e_1)}{x_I^\top q_I} \right| \geq z \right\} \\ &= \int_{u \in \mathbb{R}} \int_{v \in \mathbb{R}^{|I|}} \text{Prob} \left\{ \left| \frac{x_I^\top (u e_1)}{x_I^\top v} \right| \geq z \mid q_1 = u, q_I = v \right\} \varphi(u) \psi(v) \\ &\leq \int_{u \in \mathbb{R}} \int_{v \in \mathbb{R}^{|I|}} \frac{1}{z} \varphi(u) \psi(v) = \frac{1}{z}. \end{aligned}$$

Here the inequality follows since x_I is independent of q_1 and q_I , and therefore we can use (4) and Lemma 2 (with $p = (q_1 e_1)$ and $q = q_I$).

The same bound can be proven for all $(i, j) \in S$. Using these bounds with $z = \frac{t}{|S|}$ and (3) we obtain

$$\text{Prob}\{c_{\det}(A) \geq t\} \leq \sum_{(i,j) \in S} \text{Prob} \left\{ \left| \frac{a_{ij} \det(A_{(ij)})}{\det(A)} \right| \geq \frac{t}{|S|} \right\} \leq \frac{|S|^2}{t}. \quad \square$$

4. Matrix inversion. We now focus on the problem of inverting a matrix A and its componentwise condition number $c^\dagger(A)$. Our main results in this section are the following.

THEOREM 2. *Let $S \subset [n]^2$ be such that $\mathcal{M}_S \not\subseteq \Sigma$. Then, for all $t \geq 4|S|$,*

$$\text{Prob}_{A \in \mathcal{B}_S} \{c^\dagger(A) \geq t\} \leq 4|S|^2 n^2 \frac{1}{t}.$$

Using Proposition 2 we prove the following corollary.

COROLLARY 2. *Let $S \subset [n]^2$ be such that $\mathcal{M}_S \not\subseteq \Sigma$. Then*

$$\mathbf{E}(\log_+(c^\dagger(A))) \leq 2 \log n + 2 \log |S| + \log 4 + \frac{1}{\ln b},$$

where \mathbf{E} denotes expectation over $A \in \mathcal{B}_S$.

Remark 2. Note that, for all monotonic norm on \mathcal{M}_S , the bound above also holds for $\mathbf{m}(A)$ by Proposition 1. This is in contrast with the lower bound linear in n for the expected value of the logarithm of normwise condition numbers which follows from [8].

Definition 1 (and Remark 1) yields expressions for the (mixed and componentwise) condition numbers of a matrix A by taking $\mathcal{D} = \mathcal{M} \setminus \Sigma$ and $F : \mathcal{M} \setminus \Sigma \rightarrow \mathcal{M}$ given by $F(A) = A^{-1}$. For $k, \ell \in [n]$ such that $\gamma_{k\ell} \neq 0$, we let

$$c_{k\ell}^\dagger(A) = \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\gamma'_{k\ell} - \gamma_{k\ell}|}{\delta |\gamma_{k\ell}|},$$

and for $k, \ell \in [n]$ such that $\gamma_{k\ell} = 0$, we let $\mathbf{c}_{k\ell}^\dagger(A) = 0$ if

$$\lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\gamma'_{k\ell} - \gamma_{k\ell}|}{\delta} = 0$$

and $\mathbf{c}_{k\ell}^\dagger(A) = \infty$ otherwise. Then

$$\mathbf{c}^\dagger(A) = \max_{k, \ell \in [n]} \mathbf{c}_{k\ell}^\dagger(A).$$

Similarly, for a norm $\|\cdot\|$ on \mathcal{M} ,

$$\mathbf{m}(A) = \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{\|\Gamma' - \Gamma\|}{\delta \|\Gamma\|}.$$

LEMMA 5. For $A \in \mathcal{M} \setminus \Sigma$ and $k, \ell \in [n]$,

$$\mathbf{c}_{k\ell}^\dagger(A) \leq \mathbf{c}_{\det}(A) + \mathbf{c}_{\det}(A_{(\ell k)}).$$

Proof. We divide the proof into cases. Case (i): $\gamma_{k\ell} \neq 0$.

Let $\delta > 0$ be sufficiently small so that $\mathcal{S}(A, \delta) \cap \Sigma = \emptyset$ and, for all $A' \in \mathcal{S}(A, \delta)$,

$$\left| \frac{\det(A') - \det(A)}{\det(A)} \right| < 1. \text{ Let } A' \in \mathcal{S}(A, \delta).$$

$$\text{Since } \gamma_{k\ell} = \frac{\det(A_{(\ell k)})}{\det(A)},$$

$$\begin{aligned} \frac{\gamma'_{k\ell} - \gamma_{k\ell}}{\gamma_{k\ell}} &= \frac{\det(A)}{\det(A_{(\ell k)})} \left(\frac{\det(A'_{(\ell k)})}{\det(A')} - \frac{\det(A_{(\ell k)})}{\det(A)} \right) \\ &= \frac{\det(A)}{\det(A_{(\ell k)})} \frac{\det(A'_{(\ell k)})}{\det(A')} - 1 \\ &= \frac{1 + \frac{\det(A'_{(\ell k)}) - \det(A_{(\ell k)})}{\det(A_{(\ell k)})}}{1 + \frac{\det(A') - \det(A)}{\det(A)}} - 1 \\ &= \frac{\frac{\det(A'_{(\ell k)}) - \det(A_{(\ell k)})}{\det(A_{(\ell k)})} - \frac{\det(A') - \det(A)}{\det(A)}}{1 + \frac{\det(A') - \det(A)}{\det(A)}}. \end{aligned}$$

Using that $\left| \frac{\det(A') - \det(A)}{\det(A)} \right| < 1$,

$$\left| \frac{\gamma'_{k\ell} - \gamma_{k\ell}}{\gamma_{k\ell}} \right| \leq \frac{\left| \frac{\det(A'_{(\ell k)}) - \det(A_{(\ell k)})}{\det(A_{(\ell k)})} \right| + \left| \frac{\det(A') - \det(A)}{\det(A)} \right|}{1 - \left| \frac{\det(A') - \det(A)}{\det(A)} \right|}$$

and therefore

$$\begin{aligned} &\sup_{A' \in \mathcal{S}(A, \delta)} \left| \frac{\gamma'_{k\ell} - \gamma_{k\ell}}{\delta \gamma_{k\ell}} \right| \\ &\leq \frac{\sup_{A' \in \mathcal{S}(A, \delta)} \left| \frac{\det(A'_{(\ell k)}) - \det(A_{(\ell k)})}{\delta \det(A_{(\ell k)})} \right| + \sup_{A' \in \mathcal{S}(A, \delta)} \left| \frac{\det(A') - \det(A)}{\delta \det(A)} \right|}{1 - \sup_{A' \in \mathcal{S}(A, \delta)} \left| \frac{\det(A') - \det(A)}{\det(A)} \right|}. \end{aligned}$$

Taking limits for $\delta \rightarrow 0$ on both sides we get

$$\mathbf{c}_{k\ell}^\dagger(A) \leq \mathbf{c}_{\det}(A) + \mathbf{c}_{\det}(A_{(\ell k)}).$$

Case (ii): $\gamma_{k\ell} = 0$ and

$$\lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\gamma'_{k\ell}|}{\delta} = 0.$$

In this case, $\mathbf{c}_{k\ell}^\dagger(A) = 0$ and the statement holds.

Case (iii): $\gamma_{k\ell} = 0$ and

$$0 \neq \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\gamma'_{k\ell}|}{\delta} = \lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\det(A'_{\ell k})|}{\delta |\det(A')|}.$$

In this case,

$$\lim_{\delta \rightarrow 0} \sup_{A' \in \mathcal{S}(A, \delta)} \frac{|\det(A'_{\ell k})|}{\delta} \neq 0$$

and therefore $\mathbf{c}_{\det}(A_{\ell k}) = \infty$. The statement holds as well. \square

Proof of Theorem 2. By definition of $\mathbf{c}^\dagger(A)$,

$$\text{Prob}\{\mathbf{c}^\dagger(A) \geq t\} = \text{Prob}\left\{\max_{k, \ell \in [n]} \mathbf{c}_{k\ell}^\dagger(A) \geq t\right\} \leq \sum_{k, \ell \in [n]} \text{Prob}\{\mathbf{c}_{k\ell}^\dagger(A) \geq t\}.$$

By Lemma 4, with probability 1, A is nonsingular. So, since $\frac{t}{2} \geq 2|S|$ by hypothesis, we can apply Lemma 5 to obtain

$$\begin{aligned} \text{Prob}\{\mathbf{c}_{k\ell}^\dagger(A) \geq t\} &\leq \text{Prob}\left\{\mathbf{c}_{\det}(A) \geq \frac{t}{2}\right\} + \text{Prob}\left\{\mathbf{c}_{\det}(A_{(k\ell)}) \geq \frac{t}{2}\right\} \\ &\leq 4|S|^2 \frac{1}{t} \end{aligned}$$

the last inequality by Theorem 1. The statement now follows. \square

5. Linear equations solving. We finally deal with the problem of solving linear systems of equations.

THEOREM 3. *Let $S \subset [n]^2$ be such that $\mathcal{M}_S \not\subseteq \Sigma$. Then, for all $t \geq 4(|S| + n)$,*

$$\text{Prob}\{\mathbf{c}(A, b) \geq t\} \leq 10|S|^2 n \frac{1}{t},$$

where Prob denotes probability over $(A, b) \in \mathcal{R}_S \times N(0, \text{Id}_n)$.

We may use Proposition 2 a last time.

COROLLARY 3. *Let $S \subset [n]^2$ be such that $\mathcal{M}_S \not\subseteq \Sigma$. Then*

$$\mathbf{E}(\log_+(\mathbf{c}(A, b))) \leq \log n + 2 \log |S| + \log 10 + \frac{1}{\ln b}.$$

Definition 1 (and Remark 1) yields again expressions for the (mixed and componentwise) condition numbers of a pair (A, b) by taking $\mathcal{D} = (\mathcal{M} \setminus \Sigma) \times \mathbb{R}^n$ and $F : (\mathcal{M} \setminus \Sigma) \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by $F(A, b) = A^{-1}b$.

For $A \in \mathcal{M} \setminus \Sigma$ and $b \in \mathbb{R}^n$ we denote $x = A^{-1}b$. For $k \in [n]$ such that $x_k \neq 0$ we let

$$\mathbf{c}_k(A, b) = \lim_{\delta \rightarrow 0} \sup_{(A', b') \in \mathcal{S}((A, b), \delta)} \frac{|x'_k - x_k|}{\delta |x_k|}.$$

For $k \in [n]$ such that $x_k = 0$ we let $\mathbf{c}_k(A, b) = 0$ if

$$\lim_{\delta \rightarrow 0} \sup_{(A', b') \in \mathcal{S}((A, b), \delta)} \frac{|x'_k - x_k|}{\delta} = 0$$

and $\mathbf{c}_k(A, b) = \infty$ otherwise. Then

$$\mathbf{c}(A, b) = \max_{k \in [n]} \mathbf{c}_k(A, b).$$

Similarly, for a norm $\| \cdot \|$ in \mathbb{R}^n ,

$$\mathbf{m}(A, b) = \lim_{\delta \rightarrow 0} \sup_{(A', b') \in \mathcal{S}((A, b), \delta)} \frac{\|x' - x\|}{\delta \|x\|}.$$

In what follows let R_k be the matrix obtained by replacing the k th column of A by b .

LEMMA 6. *For any nonsingular matrix A and $k \in [n]$,*

$$\mathbf{c}_k(A, b) \leq \mathbf{c}_{\det}(A) + \mathbf{c}_{\det}(R_k).$$

Proof. By Cramer's rule,

$$x_k = \frac{\det(R_k)}{\det(A)}.$$

The rest of this proof is similar to the proof of Lemma 5. \square

Proof of Theorem 3. The proof follows the lines of that of Theorem 2. First, we get

$$\text{Prob}\{\mathbf{c}(A, b) \geq t\} \leq \sum_{k \in [n]} \text{Prob}\{\mathbf{c}_k(A, b) \geq t\}.$$

Then we apply Lemma 6 (using that, with probability 1, $A \notin \Sigma$ and that $\frac{t}{2} \geq 2|S|$) to get

$$\begin{aligned} \text{Prob}\{\mathbf{c}_k(A, b) \geq t\} &\leq \text{Prob}\left\{\mathbf{c}_{\det}(A) \geq \frac{t}{2}\right\} + \text{Prob}\left\{\mathbf{c}_{\det}(R_k) \geq \frac{t}{2}\right\} \\ &\leq 2|S|^2 \frac{1}{t} + 2(|S| + n)^2 \frac{1}{t} \leq 10|S|^2 \frac{1}{t}. \end{aligned}$$

For the second inequality we used the fact that $|S| \geq n$. The statement now follows. \square

6. Additional remarks. (1) To obtain bounds for the average loss of precision of random triangular systems one may combine Theorem 3 with the following result by Wilkinson [9, Chapter 3, section 19] which we quote as given in [4].

THEOREM 4. Let $T \in \mathbb{R}^{n \times n}$ be a nonsingular triangular matrix, and assume $nu < 0.1$ (here u is the round-off unit). Then the computed solution \hat{x} to the system $Tx = b$ satisfies

$$(T + E)\hat{x} = b,$$

where, for some universal constant c ,

$$|e_{ij}| \leq (|i - j| + 2)cu|t_{ij}|.$$

The use of $\mathbf{c}(A)$ actually yields average loss of precision with the latter measured componentwise in the computed solution.

(2) The bound in Corollary 2 appears to be worse than what computer simulations suggest $\mathbf{E}(\log \mathbf{c}^\dagger(L_n))$ should be. In [2] matrices L_n were generated for various values of n and an experimental mean of $\mathbf{E}(\log \mathbf{c}^\dagger(L_n))$ was obtained from these values. A linear regression for these means shows a best fit of $3.065 \log n - 1.1466$. A probable source of (a good part of) the difference between this value and the estimate $6 \log n + \mathcal{O}(1)$ following from Corollary 2 is the broad bound $\text{Prob}\{\max_{k,\ell \in [n]} \mathbf{c}_{k\ell}^\dagger(A) \geq t\} \leq \sum_{k,\ell \in [n]} \text{Prob}\{\mathbf{c}_{k\ell}^\dagger(A) \geq t\}$ in the proof of Theorem 2. In addition to this, the bound $\mathbf{E}(\mathbf{m}^\dagger(L_n)) \leq \mathbf{E}(\mathbf{c}^\dagger(L_n))$ following from Proposition 1 may be coarse as well. Numerical experiments in [2] suggest a best fit of $\mathbf{E}(\log \mathbf{m}^\dagger(L_n)) \approx 1.5334 \log n - 0.5723$.

(3) Section 6 of [8] discusses stability of Gaussian elimination. Having shown that, almost surely, $\kappa(L_n) \approx 2^n$ the authors reflect on how this behavior can be reconciled with the fact that “Gaussian elimination is overwhelmingly stable.” They point out that “the reason appears to be statistical: the matrices A for which $\|L^{-1}\|$ is large occupy an exponentially small proportion of the space of matrices” A , a claim for which experimental evidence is given in [6]. It would thus follow that “the matrices L produced by Gaussian elimination are far from random.”

Our results show that, in addition, Gaussian elimination (for linear equation solving) needs much less than producing matrices L in the vanishingly small set of triangular matrices with $\kappa(L)$ small. It is enough to produce matrices L outside the vanishingly small set of matrices with $\mathbf{c}(L, b)$ large.

Acknowledgments. We are grateful to Ernesto Mordecki for helpful discussions.

REFERENCES

- [1] F.L. BAUER, J. STOER, AND C. WITZGALL, *Absolute and monotonic norms*, Numer. Math., 3 (1961), pp. 257–264.
- [2] F. CUCKER AND H. DIAO, *Mixed and componentwise condition numbers for rectangular structured matrices*, CALCOLO, 44 (2007), pp. 89–115.
- [3] I. GOHBERG AND I. KOLTRACHT, *Mixed, componentwise, and structured condition numbers*, SIAM J. Matrix Anal. Appl., 14 (1993), pp. 688–704.
- [4] N. HIGHAM, *The accuracy of solutions to triangular systems*, SIAM J. Numer. Anal., 26 (1989), pp. 1252–1265.
- [5] R.D. SKEEL, *Scaling for numerical stability in Gaussian elimination*, J. ACM, 26 (1979), pp. 494–526.
- [6] L.N. TREFETHEN AND R.S. SCHREIBER, *Average-case stability of Gaussian elimination*, SIAM J. Matrix Anal. Appl., 11 (1990), pp. 335–360.
- [7] A.M. TURING, *Rounding-off errors in matrix processes*, Quart. J. Mech. Appl. Math., 1 (1948), pp. 287–308.
- [8] D. VISWANATH AND L.N. TREFETHEN, *Condition numbers of random triangular matrices*, SIAM J. Matrix Anal. Appl., 19 (1998), pp. 564–581.
- [9] J. H. WILKINSON, *Rounding Errors in Algebraic Processes*, Prentice–Hall, Englewood Cliffs, NJ, 1963.