# JACOBI'S METHOD IS MORE ACCURATE THAN QR*

JAMES DEMMEL† AND KREŠIMIR VESELIĆ‡

**Abstract.** It is shown that Jacobi's method (with a proper stopping criterion) computes small eigenvalues of symmetric positive definite matrices with a uniformly better relative accuracy bound than QR, divide and conquer, traditional bisection, or any algorithm which first involves tridiagonalizing the matrix. Modulo an assumption based on extensive numerical tests, Jacobi's method is optimally accurate in the following sense: if the matrix is such that small relative errors in its entries cause small relative errors in its eigenvalues, Jacobi will compute them with nearly this accuracy. In other words, as long as the initial matrix has small relative errors in each component, even using infinite precision will not improve on Jacobi (modulo factors of dimensionality). It is also shown that the eigenvectors are computed more accurately by Jacobi than previously thought possible. Similar results are proved for using one-sided Jacobi for the singular value decomposition of a general matrix.

**Key words.** Jacobi, symmetric eigenproblem, singular value decomposition

**AMS(MOS) subject classifications.** 65F15, 65G05

**1. Introduction.** Jacobi's method and QR iteration are two of the most common algorithms for solving eigenvalue and singular value problems. Both are backward stable, and so compute all eigenvalues and singular values with an absolute error bound equal to $p(n)\varepsilon \|H\|_2$, where $p(n)$ is a slowly growing function of the dimension $n$ of the matrix $H$, $\varepsilon$ is the machine precision, and $\|H\|_2$ is the spectral norm of the matrix. Thus large eigenvalues and singular values (those near $\|H\|_2$) are computed with high relative accuracy, but tiny ones may not have any relative accuracy at all. Indeed, it is easy to find symmetric positive definite matrices where QR returns negative eigenvalues. This error analysis does not distinguish Jacobi and QR, and so we might expect Jacobi to compute tiny values with as little relative accuracy as QR.

In this paper we show that Jacobi (with a proper stopping criterion) computes eigenvalues of positive definite symmetric matrices, and singular values of general matrices with a uniformly better relative error bound than QR, or any other method which initially tridiagonalizes (or bidiagonalizes) the matrix. This includes divide and conquer algorithms, traditional bisection, Rayleigh quotient iteration, and so on. We also show that Jacobi computes eigenvectors and singular vectors with better error bounds than QR.

In fact, for the symmetric positive definite eigenproblem, we show that Jacobi is optimally accurate in the following sense. Suppose the initial matrix entries have small relative uncertainties, perhaps from prior computations. The eigenvalues will then themselves have inherent uncertainties, independent of which algorithm is used to compute them. We show that the eigenvalues computed by Jacobi have error bounds which are nearly as small as these inherent uncertainties. In other words, as long as the initial data is slightly uncertain, even using infinite precision cannot

improve on Jacobi (modulo factors of $n$). For the singular value decomposition, we can prove a similar, but necessarily somewhat weaker, result.

These results depend on new perturbation theorems for eigenvalues and eigenvectors (or singular values and singular vectors) as well as a new error analysis of Jacobi, all of which are stronger than their classical counterparts. They also depend on an empirical observation for which we have overwhelming numerical evidence but somewhat weaker theoretical understanding.

First, we discuss the new perturbation theory for eigenvalues, contrasting the standard error bounds with the new ones. Let $H$ be a positive definite symmetric matrix, and $\delta H$ a small perturbation of $H$ in the sense that $|\delta H_{ij}/H_{ij}| \leq \eta/n$ for all $i$ and $j$. Then $\|\delta H\|_2 \leq \eta \|H\|_2$. Let $\lambda_i$ and $\lambda_i'$ be the $i$th eigenvalues of $H$ and $H + \delta H$, respectively (numbered so that $\lambda_1 \leq \cdots \leq \lambda_n$). Then the standard perturbation theory [16] states that

$$(1.1) \qquad \frac{|\lambda_i - \lambda_i'|}{\lambda_i} \leq \frac{\eta \|H\|_2}{\lambda_i} \leq \eta \|H\|_2 \cdot \|H^{-1}\|_2 = \eta \kappa(H),$$

where $\kappa(H) \equiv \|H\|_2 \cdot \|H^{-1}\|_2$ is the condition number of $H$. We prove the following stronger result: Write $H = DAD$, where $D = \text{diag}(H_{ii}^{1/2})$ and $A_{ii} = 1$. By a theorem of van der Sluis [21], [6], $\kappa(A)$ is less than $n$ times $\min_{\hat{D}} \kappa(\hat{D} H \hat{D})$, i.e., it nearly minimizes the condition number of $H$ over all possible diagonal scalings. Then we show that

$$(1.2) \qquad \frac{|\lambda_i - \lambda_i'|}{\lambda_i} \leq \eta \kappa(A),$$

i.e., the error bound $\eta \kappa(H)$ is replaced by $\eta \kappa(A)$. Clearly, it is possible that $\kappa(A) \ll \kappa(H)$ (and it is always true that $\kappa(A) \leq n\kappa(H)$), so the new bound is always at least about as good as, and can be much better than, the old bound.

In the case of the singular values of a general matrix $G$, we similarly replace the conventional relative error bound $\eta \kappa(G)$ with $\eta \kappa(B)$, where $G = BD$, $D$ chosen diagonal so the columns of $B$ have unit two-norm. This implies $\kappa(B) \leq n^{1/2} \min_{\hat{D}} \kappa(G\hat{D})$, and, as before, it is possible that $\kappa(B) \ll \kappa(G)$.

The effects of rounding errors in Jacobi are bounded as follows. We can weaken the assumption of small componentwise relative error $|\delta H_{ij}/H_{ij}| \leq \eta/n$ in the perturbation theory to $|\delta H_{ij}|/(H_{ii}H_{jj})^{1/2} \leq \eta/n$ without weakening bound (1.2). This more general perturbation bounds the rounding errors introduced by applying *one* Jacobi rotation, so that one Jacobi rotation causes relative errors in the eigenvalues bounded by $O(\varepsilon)\kappa(A)$. (In contrast, QR, or any algorithm that first tridiagonalizes the matrix, only computes eigenvalues with relative error bound $O(\varepsilon)\kappa(H)$.)

To bound the errors from *all* the Jacobi rotations, we proceed as follows. Let $H_0 = D_0 A_0 D_0$ be the original matrix and let $H_m = D_m A_m D_m$ where $H_m$ is obtained from $H_{m-1}$ by applying a single Jacobi rotation, $D_m$ is diagonal, and $A_m$ has unit diagonal. The desired error bound is proportional to $\kappa(A_0)$, i.e., it depends only on the original matrix. But our analysis only says that at step $m$ we get an error bounded by something proportional to $\kappa(A_m)$. Thus the error bound for all the Jacobi steps is proportional to $\max_m \kappa(A_m)$. So, for Jacobi to attain optimal accuracy, $\max_m \kappa(A_m)/\kappa(A_0)$ must be modest in size. In extensive random numerical tests, its maximum value was less than 1.82. Wang [23] has recently found isolated examples where it is almost 8. Our theoretical understanding of this behavior is incomplete and providing it remains an open problem.

We must finally bound the errors introduced by Jacobi's stopping criterion. To achieve accuracy proportional to $\kappa(A)$, we have had to modify the standard stopping criterion. Our modified stopping criterion has been suggested before [22], [5], [3], [20], but without our explanation of its benefits. The standard stopping criterion may be written thus:

$$\text{if } |H_{ij}| \leq \text{tol} \cdot \max_{kl} |H_{kl}|, \quad \text{set } H_{ij} = 0,$$

whereas the new one is

$$\text{if } |H_{ij}| \leq \text{tol} \cdot (H_{ii}H_{jj})^{1/2}, \quad \text{set } H_{ij} = 0$$

(here *tol* is a small threshold value, usually machine precision).

Now we consider the eigenvectors and singular vectors. Here and throughout the paper whenever we refer to an eigenvector, we assume its eigenvalue is simple. Again, let $H$ be a positive definite symmetric matrix with eigenvalues $\lambda_i$ and unit eigenvectors $v_i$. Let $\delta H$ be a small componentwise relative perturbation as before, and let $\lambda_i'$ and $v_i'$ be the eigenvalues and eigenvectors of $H + \delta H$. Then the standard perturbation theory [16] says that $v_i'$ can be chosen such that

$$(1.3) \qquad \qquad \|v_i - v_i'\| \leq \frac{\eta}{\text{absgap}_{\lambda i}} + O(\eta^2),$$

where the *absolute gap for eigenvalues* is defined as

$$(1.4) \qquad \qquad \text{absgap}_{\lambda i} \equiv \min_{j \neq i} \frac{|\lambda_i - \lambda_j|}{\|H\|_2}.$$

We prove a generally stronger result, which replaces this bound with

$$(1.5) \qquad \qquad \|v_i - v_i'\| \leq \frac{(n-1)^{1/2} \kappa(A) \cdot \eta}{\text{relgap}_{\lambda i}} + O(\eta^2),$$

where the *relative gap for eigenvalues* is defined as

$$(1.6) \qquad \qquad \text{relgap}_{\lambda i} \equiv \min_{j \neq i} \frac{|\lambda_i - \lambda_j|}{|\lambda_i \cdot \lambda_j|^{1/2}}.$$

The point is that if $H$ has two or more tiny eigenvalues, their absolute gaps are necessarily small, but their relative gaps may be large, so that the corresponding eigenvectors are really well conditioned. We prove an analogous perturbation theorem for singular vectors of general matrices. We also prove a perturbation theorem which shows that even tiny components of eigenvectors and singular vectors may be well conditioned. Again, we show that Jacobi is capable of computing the eigenvectors and singular vectors to their inherent accuracies, but QR is not.

To illustrate, consider the symmetric positive definite matrix $H = DAD$, where

$$H = \begin{bmatrix} 10^{40} & 10^{29} & 10^{19} \\ 10^{29} & 10^{20} & 10^9 \\ 10^{19} & 10^9 & 1 \end{bmatrix}, \quad A = \begin{bmatrix} 1 & .1 & .1 \\ .1 & 1 & .1 \\ .1 & .1 & 1 \end{bmatrix}, \quad D = \text{diag}(10^{20}, 10^{10}, 1).$$

Here $\kappa(H) \approx 10^{40}$ and $\kappa(A) \approx 1.33$. Thus $\eta$ relative perturbations in the matrix entries only cause $4\eta$ relative perturbations in the eigenvalues according to the new

theorem, and $3 \cdot 10^{40} \cdot \eta$ relative perturbations according to the conventional theorem. Also, the absolute gaps for the eigenvalues of $H$ are absgap$_{\lambda 1,2,3} \approx 10^{-20}, 10^{-20}, 1$, whereas, the relative gaps relgap$_{\lambda 1,2,3}$ are all approximately $10^{10}$. Thus the new theory predicts errors in $v_1$ and $v_2$ of norm $2 \cdot 10^{-10} \eta$, whereas the old theory predicts errors of $10^{20} \eta$. Jacobi attains these new error bounds, but QR generally does not. For this example, QR computes two out of the three eigenvalues as negative, whereas $H$ is positive definite. In contrast, Jacobi computes all the eigenvalues to nearly full machine precision. In fact, for this example we can show that Jacobi computes all components of all eigenvectors to nearly full relative accuracy, even though they vary by 21 orders of magnitude; again, QR does not even get the signs of many small components correct.

One might object to this example on the grounds that by reversing the order of the rows and columns before tridiagonalizing and applying QR, we compute the correct eigenvalues. However, we can easily find similar matrices (see §7) where Jacobi gets accurate eigenvalues and QR gets at least one zero or negative eigenvalue, no matter how the rows and columns are ordered.

We also show that bisection and inverse iteration (with appropriate pivoting, and applied to the original positive definite symmetric matrix) are capable of attaining the same error bounds as Jacobi. Of course, bisection and inverse iteration on a dense matrix are not competitive in speed with Jacobi, unless only one or a few eigenvalues are desired and good starting guesses are available. We use these methods to verify our numerical tests.

This work is an extension of work in [2], where analogous results were proven for matrices that are called *scaled diagonally dominant* (s.d.d.). The positive definite matrix $H = DAD$ is s.d.d. if $\|A - I\|_2 < 1$. This work replaces the assumption that $A$ is diagonally dominant with mere positive definiteness, extending the results of [2] to all positive definite symmetric matrices, as well as to the singular value decomposition of general matrices.

This work does not contradict the results of [8] and [2], where it was shown how a variation of QR could compute the singular values of a bidiagonal matrix or the eigenvalues of a symmetric positive definite tridiagonal matrix with high relative accuracy. This is because reducing a dense matrix to bidiagonal or tridiagonal form can cause large relative errors in its singular values or eigenvalues independent of the accuracy of the subsequent processing. In contrast, the results in this paper are for dense matrices.

We also discuss an accelerated version of Jacobi for the symmetric positive definite eigenproblem with an attractive speedup property: The more its accuracy exceeds that attainable by QR or other traditional methods, the faster it converges. See also [22] where earlier references for Jacobi methods on positive definite matrices, as well as for one-sided methods, can be found.

We use the following terminology to distinguish among different versions of Jacobi. "Two-sided Jacobi" refers to the original method applying Jacobi rotations to the left and right of a symmetric matrix. "One-sided Jacobi" refers to computing the SVD by applying Jacobi rotation from one side only. "Right-handed Jacobi" is one-sided Jacobi applying rotations on the right, and "left-handed Jacobi" is one-sided Jacobi applying rotations on the left.

The remainder of this paper is organized as follows. Section 2 presents the new perturbations theorems. Section 3 discusses two-sided Jacobi for the symmetric positive definite eigenproblem. Section 4 discusses one-sided Jacobi for the singular value

decomposition, and also presents the accelerated version of Jacobi just mentioned. Section 5 discusses bisection and inverse iteration. Section 6 discusses bounds on $\max_m \kappa(A_m)/\kappa(A_0)$. Section 7 contains numerical experiments. Section 8 presents our conclusions and discussion of open problems.

**2. Perturbation theory.** In this section, we prove new perturbation theorems for eigenvalues and eigenvectors of symmetric positive definite matrices, and for singular values and singular vectors of general matrices. In §2.1, we consider eigendecompositions of symmetric positive definite matrices. In §2.2, we discuss the optimality of these bounds. In §2.3, we consider the singular value decomposition of general matrices. In §2.4, we discuss the optimality of this second set of bounds.

**2.1. Symmetric positive definite matrices.** The next two lemmas were proved in [2].

LEMMA 2.1. *Let $H$ and $K$ be symmetric matrices with $K$ positive definite. Let the pencil $H - \lambda K$ have eigenvalues $\lambda_i$. Let $\delta H$ and $\delta K$ be symmetric perturbations and let $\lambda_i'$ be the (properly ordered) eigenvalues of $(H + \delta H) - \lambda(K + \delta K)$. Suppose that*

$$|x^T \delta H x| \leq \eta_H \cdot |x^T H x| \quad \text{and} \quad |x^T \delta K x| \leq \eta_K \cdot |x^T K x|$$

*for all vectors $x$ and some $\eta_H < 1$ and $\eta_K < 1$. Then either $\lambda_i = \lambda_i' = 0$ or*

$$\frac{1 - \eta_H}{1 + \eta_K} \leq \frac{\lambda_i'}{\lambda_i} \leq \frac{1 + \eta_H}{1 - \eta_K}$$

*for all $i$.*

LEMMA 2.2. *Let $H = \Delta_H^T A_H \Delta_H$ and $A_H$ be symmetric matrices. $H$ and $A_H$ need not have the same dimensions, and $\Delta_H$ may be an arbitrary full-rank conforming matrix. Similarly, let $K = \Delta_K^T A_K \Delta_K$ and $A_K$ be symmetric positive definite matrices, where $K$ and $A_K$ need not have the same dimensions and $\Delta_K$ may be an arbitrary full-rank conforming matrix. Let $\delta H = \Delta_H^T \delta A_H \Delta_H$ be a perturbation of $H$ such that $|x^T \delta A_H x| \leq \eta_H |x^T A_H x|$ for all $x$ where $\eta_H < 1$. Similarly, let $\delta K = \Delta_K^T \delta A_K \Delta_K$ be a perturbation of $K$ such that $|x^T \delta A_K x| \leq \eta_K |x^T A_K x|$ for all $x$ where $\eta_K < 1$. Let $\lambda_i$ be the ith eigenvalue of $H - \lambda K$ and $\lambda_i'$ the ith eigenvalue of $(H + \delta H) - \lambda(K + \delta K)$. Then either $\lambda_i = \lambda_i' = 0$ or*

$$\frac{1 - \eta_H}{1 + \eta_K} \leq \frac{\lambda_i'}{\lambda_i} \leq \frac{1 + \eta_H}{1 - \eta_K}.$$

THEOREM 2.3. *Let $H = DAD$ be a symmetric positive definite matrix, and $D = \operatorname{diag}(H_{ii}^{1/2})$ so $A_{ii} = 1$. Let $\delta H = D\delta AD$ be a perturbation such that $\|\delta A\|_2 \equiv \eta < \lambda_{\min}(A)$. Let $\lambda_i$ be the ith eigenvalue of $H$ and $\lambda_i'$ be the ith eigenvalue of $H + \delta H$. Then*

$$(2.1) \qquad \left| \frac{\lambda_i - \lambda_i'}{\lambda_i} \right| \leq \frac{\eta}{\lambda_{\min}(A)} \leq \kappa(A) \cdot \eta.$$

*In particular, if $|\delta H_{ij}/H_{ij}| \leq \eta/n$, then $\|\delta A\|_2 \leq \eta$ and the bound (2.1) applies.*

*Proof.* Note that for all nonzero vectors $x$,

$$\left| \frac{x^T \delta H x}{x^T H x} \right| = \left| \frac{x^T \Delta^T \delta A \Delta x}{x^T \Delta^T A \Delta x} \right| = \left| \frac{y^T \delta A y}{y^T A y} \right| \leq \frac{\eta}{\lambda_{\min}(A)}.$$

Lemma 2.2 yields the desired bound, using $K = I$ and $\delta K = 0$. It remains to prove that $|\delta H_{ij}/H_{ij}| \leq \eta/n$ implies $\|\delta A\|_2 \leq \eta$. But $A_{ii} = 1$ and $A$ positive definite imply that no entry of $A$ is larger than 1 in absolute value. (Note that this means $\kappa(A)$ is at most $n$ times larger than $1/\lambda_{\min}(A)$.) Therefore, $|\delta A_{ij}| = |\delta H_{ij}/H_{ij} \cdot A_{ij}| \leq \eta/n$ and so $\|\delta A\|_2 \leq \eta$, as desired.   $\square$

Proposition 2.10 in the next subsection shows that the bound of Theorem 2.3 is nearly attained for at least one eigenvalue. However, other eigenvalues may be much less sensitive than this most sensitive one. The next proposition provides individual eigenvalue bounds which may be much tighter.

PROPOSITION 2.4. *Let $H = DAD$ be as in Theorem 2.3, with eigenvalues $\lambda_i$ and unit eigenvectors $v_i$. Let $H + \delta H = D(A + \delta A)D$ have eigenvalues $\lambda_i'$. Let $\|\delta A\|_2 \equiv \eta \ll \lambda_{\min}(A)$. Then the bound*

$$(2.2) \qquad \frac{|\lambda_i - \lambda_i'|}{\lambda_i} \leq \frac{\eta\|Dv_i\|_2^2}{\lambda_i} + O(\eta^2)$$

*is attainable by the diagonal perturbation $\delta A_{jj} = \eta$.*

*Proof.* Bound (2.2) is derived from the standard first-order perturbation theory, which says that $\lambda_i(H + \delta H) = \lambda_i(H) + v_i^T\delta Hv_i + O(\|\delta H\|_2^2)$, and substituting $|v_i^T\delta Hv_i| = |v_i^T D\delta AD v_i| \leq \|Dv_i\|_2^2\|\delta A\|_2$. The inequality $|v_i^T D\delta AD v_i| \leq \|Dv_i\|_2^2\|\delta A\|_2$ is clearly attained for the diagonal choice of $\delta A$ in the statement of the proposition.   $\square$

We may also prove a version of Lemma 2.1 in an infinite-dimensional setting [14, §VI.3].

Now we turn to eigenvectors. A weaker version of the following theorem also appeared in [2].

THEOREM 2.5. *Let $H = DAD$ be as in Theorem 2.3. Define $H(\epsilon) = D(A+\epsilon E)D$, where $E$ is any matrix with unit two-norm. Let $\lambda_i(\epsilon)$ be the ith eigenvalue of $H(\epsilon)$, and assume that $\lambda_i(0)$ is simple so that the corresponding unit eigenvector $v_i(\epsilon)$ is well defined for sufficiently small $\epsilon$. Then*

$$\|v_i(\epsilon) - v_i(0)\|_2 \leq \frac{(n-1)^{1/2}\epsilon}{\lambda_{\min}(A) \cdot \mathrm{relgap}_{\lambda_i}} + O(\epsilon^2) \leq \frac{(n-1)^{1/2}\kappa(A)\epsilon}{\mathrm{relgap}_{\lambda_i}} + O(\epsilon^2).$$

*Proof.* Let $v_k(0)$ be abbreviated by $v_k$. From [11] we have

$$v_i(\epsilon) = v_i + \epsilon \sum_{k \neq i} \frac{v_k^T DEDv_i}{\lambda_i - \lambda_k} \cdot v_k + O(\epsilon^2).$$

Let $y_k = Dv_k$, so that

$$(2.3) \qquad v_i(\epsilon) = v_i + \epsilon \sum_{k \neq i} \frac{y_k^T Ey_i}{\lambda_i - \lambda_k} \cdot v_k + O(\epsilon^2).$$

The pair $(\lambda_i, y_i)$ is an eigenpair of the pencil $A - \lambda D^{-2}$. Thus

$$\lambda_k = \lambda_k y_k^T D^{-2} y_k = y_k^T Ay_k \geq \lambda_{\min}(A)\|y_k\|_2^2,$$

and so $\|y_k\|_2 \leq (\lambda_k/\lambda_{\min}(A))^{1/2}$. Letting $z_k = y_k/\|y_k\|_2$ lets us write

$$v_i(\epsilon) = v_i + \epsilon \sum_{k \neq i} \frac{\xi_{ik} \cdot z_k^T Ez_i}{(\lambda_i - \lambda_k)/(\lambda_k\lambda_i)^{1/2}} \cdot v_k + O(\epsilon^2),$$

where $|\xi_{ik}| = \|y_k\|_2 \|y_i\|_2/(\lambda_k \lambda_i)^{1/2} \leq 1/\lambda_{\min}(A)$. Taking norms yields the result. □

Proposition 2.11 in the next subsection shows that the bound in Theorem 2.5 is nearly attainable for all $v_i$.

As in Corollary 3 in [2], it is possible to derive a nonasymptotic result from Theorem 2.5.

COROLLARY 2.6. *Let* $H = DAD$ *be as in Theorem 2.3. Suppose that* $\delta \equiv \|\delta A\|_2/\lambda_{\min}(A)$ *satisfies*

$$\delta < \frac{1}{4} \ \text{and} \ \frac{3 \cdot 2^{-1/2} \cdot \delta}{1 - \delta} < \text{relgap}_{\lambda_i}.$$

*Let* $v_i$ *be the ith unit eigenvector of* $H = DAD$*. Then the ith unit eigenvector* $v_i'$ *of* $H' = D(A + \delta A)D$ *can be chosen so that*

$$\|v_i - v_i'\|_2 \leq \frac{(n-1)^{1/2}\delta}{(1 - 4\delta)((1 - \delta)\text{relgap}_{\lambda_i} - 3 \cdot 2^{-1/2}\delta)}.$$

*Proof.* Let $H(\epsilon) = D(A + \epsilon \cdot \delta A/\|\delta A\|_2)D$. Let $\lambda_i(\epsilon)$ be the $i$th eigenvalue of $H(\epsilon)$, and abbreviate $\lambda_i(0)$ by $\lambda_i$. Let $\text{relgap}_{\lambda_i}(\epsilon)$ denote the relative gap of the $i$th eigenvalue of $H(\epsilon)$, and $\text{relgap}_\lambda(a, b) \equiv |a - b|/(ab)^{1/2}$. The idea is that if $\epsilon$ is small, then $\lambda_i(\epsilon)$ can only change by a small relative amount, and so $\text{relgap}_{\lambda_i}(\epsilon)$ can only change by a small absolute or relative amount. Note that $\lambda_{\min}(A)$ can decrease by as much as $\|\delta A\|_2$. Then by Theorem 2.3, we can bound $\text{relgap}_{\lambda_i}(\epsilon)$ below by

$$\text{relgap}_{\lambda_i}(\epsilon) = \min_{k \neq i} \frac{|\lambda_i(\epsilon) - \lambda_k(\epsilon)|}{(\lambda_i(\epsilon)\lambda_k(\epsilon))^{1/2}} \geq \min_{k \neq i} \frac{|\lambda_i - \lambda_k| - \delta(1 - \delta)^{-1}(\lambda_i + \lambda_k)}{(\lambda_i \lambda_k)^{1/2}(1 + \delta(1 - \delta)^{-1})}$$

$$\geq (1 - \delta) \min_{k \neq i} \left( \text{relgap}_\lambda(\lambda_i, \lambda_k) - \frac{\delta}{1 - \delta} \cdot \frac{\lambda_i + \lambda_k}{(\lambda_i \lambda_k)^{1/2}} \right).$$

We consider two cases, $\text{relgap}_\lambda(\lambda_i, \lambda_k) \geq 2^{-1/2}$ and $\text{relgap}_\lambda(\lambda_i, \lambda_k) < 2^{-1/2}$. The first case corresponds to $\lambda_i$ and $\lambda_k$ differing by at least a factor of 2, whence

$$\frac{\lambda_i + \lambda_k}{(\lambda_i \lambda_k)^{1/2}} \leq 3 \cdot \text{relgap}_\lambda(\lambda_i, \lambda_k).$$

The second case corresponds to $\lambda_i$ and $\lambda_k$ differing by at most a factor of 2, whence

$$\frac{\lambda_i + \lambda_k}{(\lambda_i \lambda_k)^{1/2}} \leq 3 \cdot 2^{-1/2}.$$

Altogether, we have

$$\text{relgap}_{\lambda_i}(\epsilon) \geq (1 - \delta) \left( 1 - \frac{3\delta}{1 - \delta} \right) \left( \text{relgap}_{\lambda_i} - \frac{3 \cdot 2^{-1/2} \cdot \delta}{1 - \delta} \right).$$

Now integrate the bound of Theorem 2.5 from $\epsilon = 0$ to $\epsilon = \|\delta A\|_2$ to get the desired result. □

In complete analogy to [2], we may also prove the following proposition.

PROPOSITION 2.7. *Let* $\lambda_1 \leq \cdots \leq \lambda_n$ *be the eigenvalues of* $H$ *and* $h_1 \leq \cdots \leq h_n$ *be its diagonal entries in increasing order. Then*

$$\lambda_{\min}(A) \leq \frac{\lambda_i}{h_i} \leq \lambda_{\max}(A).$$

*In other words, the diagonal entries of $H$ can differ from the eigenvalues only by factors bounded by $\kappa(A)$.*

*Proof.* See the proof of Proposition 2 in [2].    □

PROPOSITION 2.8.  *Let $H = DAD$ with eigenvalues $\lambda_i$. Let $d_i$ be the diagonal entries of $D$. Let $v_i$ be the ith eigenvector of $H$ normalized so that its ith component $v_i(i) = 1$. Then*

$$|v_i(j)| \leq \overline{v}_i(j) \equiv (\kappa(A))^{3/2} \cdot \min\left(\left(\frac{\lambda_i}{\lambda_j}\right)^{1/2}, \left(\frac{\lambda_j}{\lambda_i}\right)^{1/2}\right).$$

*We also have*

$$|v_i(j)| \leq (\kappa(A))^{3/2} \cdot \min\left(\frac{d_i}{d_j}, \frac{d_j}{d_i}\right).$$

*In other words, the eigenvectors are scaled analogously to the diagonal of $H$.*

*Proof.* See the proof of Proposition 6 in [2].    □

PROPOSITION 2.9.  *Let $H(\epsilon)$ and $v_i(\epsilon)$ be as in Theorem 2.5, and $\overline{v}_i(j)$ be as in Proposition 2.8. Then*

$$|v_i(\epsilon)(j) - v_i(0)(j)| \leq \frac{(2n-2)^{1/2}}{\lambda_{\min}(A) \cdot \min(\mathrm{relgap}_{\lambda_i}, 2^{-1/2})} \cdot \epsilon \cdot \overline{v}_i(j) + O(\epsilon^2).$$

*In other words, each component of each eigenvector is perturbed by a small amount relative to its upper bound of $\overline{v}_i(j)$ of Proposition 2.8. Thus small components of eigenvectors may be determined with as much relative accuracy as large components. Note that $\mathrm{relgap}_{\lambda_i}$ exceeds $2^{-1/2}$ only when $\lambda_i$ differs from its nearest neighbor by at least a factor of 2.*

*Proof.* See the proof of Theorem 7 in [2].    □

We illustrate these results with two examples.  First, we consider the matrix $H = DAD$ of the introduction:

$$H = \begin{bmatrix} 10^{40} & 10^{29} & 10^{19} \\ 10^{29} & 10^{20} & 10^9 \\ 10^{19} & 10^9 & 1 \end{bmatrix}, \quad A = \begin{bmatrix} 1 & .1 & .1 \\ .1 & 1 & .1 \\ .1 & .1 & 1 \end{bmatrix}, \quad D = \mathrm{diag}\,(10^{20}, 10^{10}, 1).$$

To six correct figures, $H$'s eigenvalue matrix $\Lambda$ and eigenvector matrix $V$ (normalized to have the largest entry of each eigenvector equal to 1) are

$$\Lambda = \mathrm{diag}\,(1.00000 \cdot 10^{40}, 9.90000 \cdot 10^{19}, 9.81818 \cdot 10^{-1})$$

and

$$V = \begin{bmatrix} 1.00000 & -1.00000 \cdot 10^{-11} & -9.09091 \cdot 10^{-22} \\ 1.00000 \cdot 10^{-11} & 1.00000 & -9.09091 \cdot 10^{-12} \\ 1.00000 \cdot 10^{-21} & 9.09091 \cdot 10^{-12} & 1.00000 \end{bmatrix}.$$

We may compute that $\kappa(H) \approx 10^{40}$ and $\kappa(A) \approx 1.33$.  Thus, according to Theorem 2.3, changing each entry of $H$ in its seventh decimal place or beyond would not change $\Lambda$ in the figures shown. The refined error bounds of Proposition 2.4 are essentially the same in this case. We can further verify the assertion of Proposition 2.7 that the ratios of the eigenvalues to the diagonal entries of $H$ are bounded between $.9 = \lambda_{\min}(A)$ and

$1.2 = \lambda_{\max}(A)$. One may also compute that the relative gaps $\text{relgap}_{\lambda_i}$ for all three eigenvalues are approximately $10^{10}$. Thus, according to Theorem 2.5, seventh-figure changes in $H$ would not change its eigenvectors by more than $10^{-16}$ in norm. In fact, the eigenvectors are even more accurately determined than this. Let $\bar{V} = \{\bar{v}_i(j)\}$ be the matrix of upper bounds of entries of $V$ as defined in Proposition 2.8:

$$\bar{V} \approx \begin{bmatrix} 1.5 & 1.5 \cdot 10^{-10} & 1.5 \cdot 10^{-20} \\ 1.5 \cdot 10^{-10} & 1.5 & 1.5 \cdot 10^{-10} \\ 1.5 \cdot 10^{-20} & 1.5 \cdot 10^{-10} & 1.5 \cdot 10^{-20} \end{bmatrix}.$$

Then, according to Proposition 2.9, seventh-figure changes in $H$ cause changes in at most the fifth digits of all the entries of $V$. In other words, for this, examples of all the eigenvalues and all the components of all the eigenvectors, are determined to nearly full relative precision by the data. Later, we show that Jacobi can compute them with this accuracy. In contrast, QR does not even get the signs of the two small eigenvalues or many components of the eigenvectors correct.

The second example serves to illustrate the difference between Theorem 2.3 and the refined bounds of Proposition 2.4. Let $H = DAD$ where $D$ is the same as before and

$$A = \begin{bmatrix} 1 & 1-\mu & 1-\mu \\ 1-\mu & 1 & 1-\mu \\ 1-\mu & 1-\mu & 1 \end{bmatrix},$$

where $\mu = 10^{-6}$. The eigenvalues of $H$ are $10^{40}$, $2 \cdot 10^{14}$, and $1.5 \cdot 10^{-6}$. Now $\kappa(A) \approx 10^6$, so according to Theorem 2.3, an $\eta$ relative change in the matrix entries causes as much as a $10^6\eta$ relative change in the eigenvalues. In contrast, the refined bounds predict a relative change of $\eta$ in $10^{40}$ and $10^6\eta$ in the two smaller eigenvalues. Thus the largest eigenvalue is just as insensitive as predicted by standard norm-based perturbation theory.

### 2.2. Optimality of the bounds for symmetric positive definite matrices.
In this section, we show that the bounds of the last section are attainable. In other words, the only symmetric positive definite matrices whose eigenvalues are determined to high relative accuracy by the matrix entries are those $H = DAD$, where $A$ is well conditioned.

In particular, we give explicit small, componentwise, relative perturbations, which attain the eigenvalue bounds; it suffices to choose a diagonal perturbation. We have (necessarily) slightly weaker results for the optimality of our eigenvector bounds.

We begin by showing that the assumption $\|\delta A\|_2 < \lambda_{\min}(A)$ of the last section is essential to having relative error bounds at all. If this bound were violated, $A + \delta A$ (and so $H + \delta H$) could become indefinite, implying that all relative accuracy in at least one eigenvalue is completely lost. In contrast to standard perturbation theory, however, which assumes a bound on $\|\delta H\|_2$ instead of $\|\delta A\|_2$, one cannot say which eigenvalue will lose relative accuracy first. In the conventional case, as $\|\delta H\|_2$ grows, it is the smallest eigenvalues that lose accuracy first, the larger ones remaining accurate. As $\|\delta A\|_2$ grows, however, *any* eigenvalue in the spectrum (except the very largest) may lose its relative accuracy first. The following example illustrates this:

$$H = \begin{bmatrix} 10^{20} & & & \\ & 1 & .99 & \\ & .99 & 1 & \\ & & & 10^{-20} \end{bmatrix}, \qquad A = \begin{bmatrix} 1 & & & \\ & 1 & .99 & \\ & .99 & 1 & \\ & & & 1 \end{bmatrix},$$

and $D = \text{diag}\,(10^{10}, 1, 1, 10^{-10})$. Note that $\lambda_{\min}(A) = .01$. As $\|\delta A\|_2$ approaches .01, the eigenvalues near $10^{20}$, 1.99, and $10^{-20}$ retain their accuracy, but the one near .01 can lose all its relative accuracy.

We next show that the relative error bound of Theorem 2.1 can be nearly attained for at least one eigenvalue simply by making appropriate small relative perturbations to the diagonal of $H$.

PROPOSITION 2.10. *Let $H = DAD$ be symmetric positive definite, with $D = \text{diag}\,(H_{ii}^{1/2})$ diagonal and $A_{ii} = 1$. Let $\delta A = \eta I$, $0 < \eta < \lambda_{\min}(A)$, and $H + \delta H = D(A + \delta A)D$. Then for some $i$ we have*

$$\frac{\lambda_i(H + \delta H)}{\lambda_i(H)} \geq \left(1 + \frac{\eta}{\lambda_{\min}(A)}\right)^{1/n} \approx 1 + \frac{\eta}{n\lambda_{\min}(A)}.$$

*Proof.* We have

$$\prod_i \lambda_i(H) = \det(DAD) = \det(D^2)\det(A) = \det(D^2)\prod_i \lambda_i(A)$$

and

$$\prod_i \lambda_i(H + \delta H) = \det(D(A + \eta I)D) = \det(D^2)\det(A + \eta I) = \det(D^2)\prod_i(\lambda_i(A) + \eta).$$

Therefore,

$$\prod_i \frac{\lambda_i(H + \delta H)}{\lambda_i(H)} = \prod_i \frac{\lambda_i(A) + \eta}{\lambda_i(A)} \geq 1 + \frac{\eta}{\lambda_{\min}(A)},$$

implying that at least one factor $\lambda_i(H + \delta H)/\lambda_i(H)$ must exceed $(1 + \eta/\lambda_{\min}(A))^{1/n}$. This last expression is approximately $1 + \eta/(n\lambda_{\min}(A))$, when $\eta \ll \lambda_{\min}(A)$. $\quad\square$

The example at the beginning of this section showed that the error bound of Theorem 2.3 and the last proposition may only be attained for one eigenvalue. Proposition 2.4 of §2.1 showed that for asymptotically small $\|\delta A\|_2$, the maximum perturbation in each eigenvalue may be attained only with small diagonal perturbations of $A$.

After we show that the rounding errors introduced by Jacobi are of the form $\|\delta A\|_2 = O(\varepsilon)$ in §2.3, Propositions 2.10 and 2.4 show that Jacobi (modulo the assumption on $\max_m \kappa(A_m)/\kappa(A_0)$) computes all the eigenvalues with optimal accuracy, provided that only the diagonal entries of $H$ have small relative errors. The same optimality property is true of bisection.

Now we consider eigenvectors. Here our results are necessarily weaker, as the following example shows. Suppose $H$ is diagonal with distinct eigenvalues. Then small relative perturbations to the matrix entries leave $H$ diagonal and its eigenvalue matrix (the identity matrix) unchanged. Therefore, the only way we can hope to attain the bounds of Theorem 2.5 is to use perturbations $\delta A$, which are possibly dense, even if $H$ is not. Furthermore, a block diagonal example like the first one in this section shows that the attainable eigenvector perturbations do not necessarily grow with $\kappa(A)$. Thus the following is the best we can prove.

PROPOSITION 2.11. *Let $H = DAD$, $\lambda_i$, $v_i$, $\delta H$, $\delta A$ and let $\eta$ be as in Proposition 2.4. Let $v_i'$ be the unit eigenvectors of $H + \delta H$. Then we can choose $\delta A$, $\|\delta A\|_2 \equiv \eta \ll \lambda_{\min}(A)$, so that*

$$\|v_i - v_i'\|_2 \geq \frac{\eta}{\lambda_{\max}(A)\text{relgap}_{\lambda_i}} + O(\eta^2).$$

*Proof.* Consider expression (2.3) for $v_i - v_i'$ (there, $\delta A$ is written $\epsilon E$). By using a Householder transformation, we can prove that there exists a symmetric $\delta A$ such that $y_k^T \delta A y_i = \|y_k\|_2 \|y_i\|_2 \|\delta A\|_2$ for arbitrary $y_k$ and $y_i$. Since $\lambda_k = y_k^T A y_k \leq \|y_k\|^2 \lambda_{\max}(A)$, we can find $\delta A$ to make $y_k^T \delta A y_i \geq (\lambda_i \lambda_k)^{1/2} \|\delta A\|_2 / \lambda_{\max}(A)$. Choosing $k$ so that $\lambda_k$ is closest to $\lambda_i$ completes the proof. $\quad\square$

**2.3. Singular value decomposition.** The results on singular values and singular vectors are analogous to the results for eigenvalues and eigenvectors in the first subsection, so we do not include the proofs. Just as we derived perturbation bounds for eigenvalues from a more general result for generalized eigenvalues of pencils, we start with a perturbation bound for generalized singular values and then specialize to standard singular values.

Let $G_1$ and $G_2$ be matrices with the same number of columns, $G_2$ of full column rank, and both arbitrary. We define the *ith generalized singular value* $\sigma_i(G_1, G_2)$ *of the pair* $(G_1, G_2)$ as the square root of the *i*th eigenvalue of the definite pencil $G_1^T G_1 - \lambda G_2^T G_2$ [11]. If we let $G_2$ be the identity, $\sigma_i(G_1, G_2)$ is the same as the standard singular value $\sigma_i(G_1)$ of $G_1$.

LEMMA 2.12. *Let $G_1$ and $G_2$ be matrices with the same number of columns, $G_2$ of full column rank, and both arbitrary. Let $\delta G_j$ be a perturbation of $G_j$ such that*

$$\|\delta G_j x\|_2 \leq \eta_j \|G_j x\|_2$$

*for all $x$ and some $\eta_j < 1$. Let $\sigma_i$ be the ith generalized singular value of $(G_1, G_2)$ and $\sigma_i'$ be the ith generalized singular value of $(G_1 + \delta G_1, G_2 + \delta G_2)$. Then either $\sigma_i = \sigma_i' = 0$ or*

$$\frac{1 - \eta_1}{1 + \eta_2} \leq \frac{\sigma_i'}{\sigma_i} \leq \frac{1 + \eta_1}{1 - \eta_2}.$$

LEMMA 2.13. *Let $G_1$ and $G_2$ be as in Lemma 2.12. Let $G_j = B_j \Delta_j$, where $\Delta_j$ has full rank and is otherwise arbitrary. Let $\delta G_j = \delta B_j \Delta_j$ be a perturbation of $G_j$ such that $\|\delta B_j x\|_2 \leq \eta_j \|B_j x\|_2$ for all $x$ and some $\eta_j < 1$. Let $\sigma_i$ and $\sigma_i'$ be the ith generalized singular values of $(G_1, G_2)$ and $(G_1 + \delta G_1, G_2 + \delta G_2)$, respectively. Then either $\sigma_i = \sigma_i' = 0$ or*

$$\frac{1 - \eta_1}{1 + \eta_2} \leq \frac{\sigma_i'}{\sigma_i} \leq \frac{1 + \eta_1}{1 - \eta_2}.$$

THEOREM 2.14. *Let $G = BD$ be a general full-rank matrix, and let $D$ be chosen diagonal so that the columns of $B$ have unit two-norm (i.e., $D_{ii}$ equals the two-norm of the ith column of $G$). Let $\delta G = \delta B D$ be a perturbation of $G$ such that $\|\delta B\|_2 \equiv \eta < \sigma_{\min}(B)$. Let $\sigma_i$ and $\sigma_i'$ be the ith singular values of $G$ and $G + \delta G$, respectively. Then*

$$(2.4) \qquad \frac{|\sigma_i - \sigma_i'|}{\sigma_i} \leq \frac{\eta}{\sigma_{\min}(B)} \leq \kappa(B) \cdot \eta,$$

*where $\kappa(B) = \sigma_{\max}(B)/\sigma_{\min}(B) \leq n^{1/2}/\sigma_{\min}(B)$, and $n$ is the number of columns of $G$. In particular, if $|\delta G_{ij}/G_{ij}| \leq \eta/n$, then $\|\delta B\|_2 \leq \eta$ and the bound (2.4) applies.*

Just as the bounds of Theorem 2.3 were not attainable by all eigenvalues, neither are the bounds of Theorem 2.14 attainable for all singular values. Analogous to Proposition 2.4, we may derive tighter bounds for individual singular values.

PROPOSITION 2.15. *Let $G = BD$ be as in Theorem 2.14, with singular values $\sigma_i$, right unit singular vectors $v_i$, and left unit singular vectors $u_i$. Let $G + \delta G = (B + \delta B)D$ have singular values $\sigma_i'$, where $\|\delta B\|_2 \equiv \eta \ll \sigma_{\min}(B)$. Then the bound*

$$(2.5) \qquad \frac{|\sigma_i - \sigma_i'|}{\sigma_i} \leq \frac{\eta \|Dv_i\|_2}{\sigma_i} + O(\eta^2)$$

*is attainable by the perturbation $\delta B = \eta u_i(Dv_i)^T / \|Dv_i\|_2$.*

Now we consider the singular vectors. For simplicity, we assume that $G$ is square. We use the fact that if $G = U \Sigma V^T$ is the singular value decomposition of $G$, then

$$2^{-1/2} \cdot \begin{bmatrix} V & V \\ U & -U \end{bmatrix}$$

is the eigenvector matrix of the symmetric matrix [11]

$$\begin{bmatrix} 0 & G^T \\ G & 0 \end{bmatrix}.$$

Therefore, we can use perturbation theory for eigenvectors of symmetric matrices to do perturbation theory for singular vectors of general matrices.

We also need to define the gaps for the singular vector problem. The *absolute gap for singular values* is

$$\text{absgap}_{\sigma i} \equiv \min_{k \neq i} \frac{|\sigma_i - \sigma_k|}{\|G\|_2},$$

i.e., essentially the same as the absolute gap for eigenvalues. However the *relative gap for singular values*,

$$\text{relgap}_{\sigma i} \equiv \min_{k \neq i} \frac{|\sigma_i - \sigma_k|}{\sigma_i + \sigma_k},$$

is somewhat different from the relative gap for eigenvalues.

The standard perturbation theorem for singular vectors is essentially the same as for eigenvectors. Let $G$ have right (or left) unit singular vectors $v_i$, and let $G + \delta G$ have right (or left) unit singular vectors $v_i'$. Let $\eta = \|\delta G\|_2 / \|G\|_2$. Then

$$\|v_i - v_i'\|_2 \leq \frac{\eta}{\text{absgap}_{\sigma i}} + O(\eta^2).$$

We improve this in the following theorem.

THEOREM 2.16. *Let $G = BD$ be as in Theorem 2.14. Define $G(\epsilon) = (B + \epsilon E)D$ where $E$ is any matrix with unit two-norm. Let $\sigma_i(\epsilon)$ be the $i$th singular value of $G(\epsilon)$, and assume that $\sigma_i(0)$ is simple so that the corresponding right unit singular vector $v_i(\epsilon)$ and left unit singular vector $u_i(\epsilon)$ are well defined for sufficiently small $\epsilon$. Then*

$$\max(\|v_i(\epsilon) - v_i(0)\|_2, \|u_i(\epsilon) - u_i(0)\|_2) \leq \frac{(n - .5)^{1/2}\kappa(B)\epsilon}{\text{relgap}_{\sigma i}} + O(\epsilon^2).$$

COROLLARY 2.17. *Let $G = BD$ be as in Theorem 2.14. Suppose that $\delta \equiv \|\delta B\|_2 / \sigma_{\min}(B)$ satisfies*

$$\frac{\delta}{1 - \delta} < \text{relgap}_{\sigma i}.$$

*Let $v_i$ and $u_i$ be the unit right and left singular vectors of $G$, respectively, and let $v_i'$ and $u_i'$ be the unit right and left singular vectors of $G' = (B + \delta B)D$, respectively. Then*

$$\max(\|v_i - v_i'\|_2, \|u_i - u_i'\|_2) \leq \frac{(n - .5)^{1/2}\delta}{(1 - \delta)((1 - \delta)\mathrm{relgap}_{\sigma_i} - \delta)}.$$

There are analogues to Propositions 2.7–2.9 of the last section, obtained by considering $H = G^T G$:

PROPOSITION 2.18. *Let $G = BD$ be as in Theorem 2.14. Let $\sigma_1 \leq \cdots \leq \sigma_n$ be the singular values of $G$ and $d_1 \leq \cdots \leq d_n$ the diagonal entries of $D$ in increasing order. Then*

$$\sigma_{\min}(B) \leq \frac{\sigma_i}{d_i} \leq \sigma_{\max}(B).$$

PROPOSITION 2.19. *Let $G = BD$ be as in Theorem 2.14 with singular values $\sigma_1 \leq \cdots \leq \sigma_n$. Let $v_i$ be the ith right singular vector of $G$, normalized so that its ith component $v_i(i) = 1$. Then*

$$|v_i(j)| \leq \bar{v}_i(j) \equiv (\kappa(B))^3 \cdot \min\left(\frac{\sigma_i}{\sigma_j}, \frac{\sigma_j}{\sigma_i}\right).$$

*We also have*

$$|v_i(j)| \leq (\kappa(B))^3 \cdot \min\left(\frac{d_i}{d_j}, \frac{d_j}{d_i}\right).$$

PROPOSITION 2.20. *Let $G(\epsilon)$ and $v_i(\epsilon)$ be as in Theorem 2.16, and $\bar{\bar{v}}_i(j)$ be as in Proposition 2.19. Then*

$$|v_i(\epsilon)(j) - v_i(0)(j)| \leq \frac{2(n - 1)^{1/2}}{\sigma_{\min}^2(B) \cdot \mathrm{relgap}_{\sigma_i}} \cdot \epsilon \cdot \bar{\bar{v}}_i(j) + O(\epsilon^2).$$

There are analogues to all the results in this section for matrices $G = DB$ scaled from the left instead of the right. Thus we can choose to scale either the rows or the columns of $G$ to have unit two-norms, whichever one minimizes the condition number. It is natural to ask if we can do better by considering two-sided diagonal scaling $D_1 G D_2$; to date, we have been unable to formulate a reasonable perturbation theory. To see why, note that if $G$ is triangular, it can be made as close to the identity matrix as desired by two-sided scaling, even though its singular values can be quite sensitive.

**2.4. Optimality of the bounds for the singular value decomposition.**
The results in this section are analogous to, but necessarily weaker than, the results of §2.2. In particular, it is no longer the case that the perturbation bounds for the singular values can be attained by small relative perturbations in the matrix entries.

First, consider the restriction $\|\delta B\|_2 < \sigma_{\min}(B)$. Just as in the symmetric positive definite case, this is necessary so that $B + \delta B$ remains nonsingular. When $B + \delta B$ becomes singular, at least one singular value necessarily loses all relative accuracy. The same kind of block diagonal example as in §2.2 also shows that only one singular value may have its sensitivity depend on $\kappa(B)$, and it might be anywhere in the spectrum (except the very largest singular value).

In order to prove an analogue of Proposition 2.10, we must permit perturbations $\delta B$ of $B$ which are small in norm but may make large relative changes in tiny entries of $B$ (a similar perturbation was needed to prove that the bound in Proposition 2.15 was attainable).

PROPOSITION 2.21. *Let $G = BD$ with $D$ diagonal and the columns of $B$ having unit two-norm. Then there exists a $\delta B$ with $\|\delta B\|_2 = \eta < \sigma_{\min}(B)$ such that for $G + \delta G = (B + \delta B)D$ we have, for at least one $i$,*

$$\frac{\sigma_i(G + \delta G)}{\sigma_i(G)} \geq \left(1 + \frac{\eta}{\sigma_{\min}(B)}\right)^{1/n} \approx 1 + \frac{\eta}{n\sigma_{\min}(B)}.$$

*If we restrict $\delta B$ so that $|\delta B_{ij}/B_{ij}| \leq \eta$, then such a perturbation $\delta B$ may not exist.*

*Proof.* The proof is very similar to that of Proposition 2.10. Let $X$ be a rank-one matrix of minimal two-norm such that $B + X$ is singular, and let $\delta B = -\eta X$. Then, as in Proposition 2.10, we discover that

$$\prod_i \frac{\sigma_i(G + \delta G)}{\sigma_i(G)} = 1 + \frac{\eta}{\sigma_{\min}(B)},$$

and so at least one term $\sigma_i(G + \delta G)/\sigma_i(G)$ exceeds $(1 + \eta/\sigma_{\min}(B))^{1/n}$. To see that small componentwise relative perturbations are not sufficient, consider the matrix

$$G = B = \begin{bmatrix} 1 & 1 \\ -\epsilon & \epsilon \end{bmatrix}$$

with $\epsilon \ll 1$. The condition number of $B$ is approximately $1/\epsilon$, and relative perturbations of size $\eta$ in its entries cannot change its singular values by more than a factor of about $(1 \pm \eta)^2$.   □

As in Proposition 2.11, our lower bound on the attainable perturbations in the singular vectors requires a dense $\delta B$ and does not grow with $\kappa(B)$.

PROPOSITION 2.22. *Let $G = BD$, $\sigma_i$, $u_i$, $v_i$; $\delta G = \delta BD$; and $\eta$ be as in Proposition 2.15. Let $u_i'$ and $v_i'$ be the unit left and right singular vectors of $G + \delta G$, respectively. Then we can choose $\delta B$, $\|\delta B\|_2 \equiv \eta \ll \sigma_{\min}(B)$, so that*

$$\max(\|u_i - u_i'\|_2, \|v_i - v_i'\|_2) \geq \frac{\eta}{2^{3/2}\sigma_{\max}(B)\mathrm{relgap}_{\sigma_i}} + O(\eta^2).$$

**3. Two-sided Jacobi.** In this section, we prove that two-sided Jacobi in floating point arithmetic applied to a positive definite symmetric matrix computes the eigenvalues and eigenvectors with the error bounds of §2.

In this introduction, we present the algorithm and our model of floating point arithmetic. In §3.1, we derive error bounds for the computed eigenvalues. In §3.2, we derive error bounds for the computed eigenvectors.

Let $H_0 = D_0 A_0 D_0$ be the initial matrix, and $H_m = D_m A_m D_m$, where $H_m$ is obtained from $H_{m-1}$ by applying a single Jacobi rotation. Here $D_m$ is diagonal and $A_m$ has unit diagonal as before. All the error bounds in this section contain the factor $\max_m \kappa(A_m)$, whereas the perturbation bounds of §2 are proportional to $\kappa(A_0)$. Therefore, our claim that Jacobi solves the eigenproblem as accurately as predicted in §2 depends on the ratio $\max_m \kappa(A_m)/\kappa(A_0)$ being modest in size. Note that convergence of $H_m$ to diagonal form is equivalent to the convergence of $A_m$ to the identity, or $\kappa(A_m)$ to 1. Thus we expect $\kappa(A_m)$ to be less than $\kappa(A_0)$ eventually.

We have overwhelming numerical evidence that $\max_m \kappa(A_m)/\kappa(A_0)$ is modest in size; in §7, the largest value this ratio attained in random testing was 1.82. Our theoretical understanding of why this ratio is so small is somewhat weaker; we present our theoretical bounds on this ratio in §6.

The essential difference between our algorithm and standard two-sided Jacobi is the stopping criterion. As stated in the introduction (and justified by Theorem 2.3), we set $H_{ij}$ to zero only if $H_{ij}/(H_{ii}H_{jj})^{1/2}$ is small. Otherwise, our algorithm is a simplification of the standard one introduced by Rutishauser [16]. We have chosen a simple version of the algorithm, omitting enhancements such as delayed updates of the diagonals and fast rotations, to make the error analysis clearer (an error analysis of these enhancements is future work).

ALGORITHM 3.1 (Two-sided Jacobi for the symmetric positive definite eigenproblem). tol is a user-defined stopping criterion. The matrix $V$ whose columns are the computed eigenvectors initially contains the identity.

> repeat
>> for all pairs $i < j$
>>> /* compute the Jacobi rotation which diagonalizes
>>> $$\begin{bmatrix} H_{ii} & H_{ij} \\ H_{ji} & H_{jj} \end{bmatrix} \equiv \begin{bmatrix} a & c \\ c & b \end{bmatrix} \quad /*$$
>>> $\zeta = (b-a)/(2c); \; t = \text{sign}(\zeta)/(|\zeta| + \sqrt{1+\zeta^2})$
>>> $cs = 1/\sqrt{1+t^2}; \; sn = cs * t$
>>> /* update the $2 \times 2$ submatrix */
>>> $H_{ii} = a - c * t$
>>> $H_{jj} = b + c * t$
>>> $H_{ij} = H_{ji} = 0$
>>> /* update the rest of rows and columns $i$ and $j$ */
>>> for $k = 1$ to $n$ except $i$ and $j$
>>>> $tmp = H_{ik}$
>>>> $H_{ik} = cs * tmp - sn * H_{jk}$
>>>> $H_{jk} = sn * tmp + cs * H_{jk}$
>>>> $H_{ki} = H_{ik}; \; H_{kj} = H_{jk}$
>>> endfor
>>> /* update the eigenvector matrix $V$ */
>>> for $k = 1$ to $n$
>>>> $tmp = V_{ki}$
>>>> $V_{ki} = cs * tmp - sn * V_{kj}$
>>>> $V_{kj} = sn * tmp + cs * V_{kj}$
>>> endfor
>> endfor
> until convergence (all $|H_{ij}|/(H_{ii}H_{jj})^{1/2} \leq$ tol)

Our model of arithmetic is a variation on the standard one: the floating point result $fl(\cdot)$ of the operation $(\cdot)$ is given by

$$fl(a \pm b) = a(1 + \varepsilon_1) \pm b(1 + \varepsilon_2),$$

(3.1) $$fl(a \times b) = (a \times b)(1 + \varepsilon_3),$$

$$fl(a/b) = (a/b)(1 + \varepsilon_4),$$

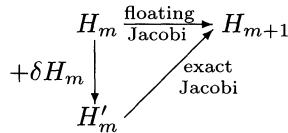$$fl(\sqrt{a}) = \sqrt{a}(1 + \varepsilon_5),$$

where $|\varepsilon_i| \leq \varepsilon$ and $\varepsilon \ll 1$ is the machine precision. This is somewhat more general

than the usual model, which uses $fl(a \pm b) = (a \pm b)(1 + \varepsilon_1)$, and includes machines like the Cray, which do not have a guard digit. This does not greatly complicate the error analysis, but it is possible that the computed rotation angle may be off by a factor of 2, whereas with a guard digit the rotation angle is always highly accurate. This may adversely affect convergence, but as we see it does not affect the one-step error analysis.

Numerically subscripted $\varepsilon$'s denote independent quantities bounded in magnitude by $\varepsilon$. As usual, we make approximations like $(1 + i\varepsilon_1)(1 + j\varepsilon_2) = 1 + (i + j)\varepsilon_3$ and $(1 + i\varepsilon_1)/(1 + j\varepsilon_2) = 1 + (i + j)\varepsilon_3$.

**3.1. Error bounds for eigenvalues computed by two-sided Jacobi.** The next theorem and its corollary justify our accuracy claims for eigenvalues computed by two-sided Jacobi.

THEOREM 3.1. *Let $H_m$ be the sequence of matrices generated by Algorithm* 3.1 *in finite precision arithmetic with precision $\varepsilon$; that is, $H_{m+1}$ is obtained from $H_m$ by applying a single Jacobi rotation. Then the following diagram:*

$$
\begin{array}{ccc}
H_m & \xrightarrow{\text{floating}\ \text{Jacobi}} & H_{m+1} \\
+\delta H_m \downarrow & \nearrow \text{exact} & \\
& \phantom{x}\text{Jacobi} & \\
H_m'
\end{array}
$$

*commutes in the following sense: The top arrow indicates that $H_{m+1}$ is obtained from $H_m$ by applying one Jacobi rotation in floating point arithmetic. The diagonal arrow indicates that $H_{m+1}$ is obtained from $H_m'$ by applying one Jacobi rotation in exact arithmetic; thus $H_{m+1}$ and $H_m'$ are exactly similar. The vertical arrow indicates that $H_m' = H_m + \delta H_m$. $\delta H_m$ is bounded as follows. Write $\delta H_m = D_m \delta A_m D_m$. Then*

$$(3.2) \qquad \|\delta A_m\|_2 \le (182(2n - 4)^{1/2} + 104)\varepsilon.$$

*In other words, if $\|\delta A_m\|_2 < \lambda_{\min}(A_m)$, one step of Jacobi satisfies the assumptions needed for the error bounds of §2.*

COROLLARY 3.2. *Assume that Algorithm* 3.1 *converges, and that $H_M$ is the final matrix whose diagonal entries we take as the eigenvalues. Write $H_m = D_m A_m D_m$ with $D_m$ diagonal and $A_m$ with ones on the diagonal for $0 \le m \le M$. Let $\lambda_j$ be the $j$th eigenvalue of $H_0$ and $\lambda_j'$ be the $j$th diagonal entry of $H_M$. Then to first order in $\varepsilon$, the following error bound holds:*

$$(3.3) \qquad \frac{|\lambda_j - \lambda_j'|}{\lambda_j} \le (\varepsilon \cdot M \cdot (182(2n - 4)^{1/2} + 104) + n \cdot \text{tol}) \cdot \max_{0 \le m \le M} \kappa(A_m).$$

*Remark.* In numerical experiments presented in §7, there was no evidence that the actual error bounded in (3.3) grew with increasing $n$ or $M$.

*Proof of Corollary* 3.2. Bound (3.3) follows by substituting the bound (3.2) and the stopping criterion into Theorem 2.3.     □

*Remark.* A similar bound can be obtained based on the error bound in Proposition 2.4.

*Proof of Theorem* 3.1. The proof of the commuting diagram is a tedious computation. Write the $2 \times 2$ submatrix of $H_{mm}$ being reduced as

$$
\begin{bmatrix} a & c \\ c & b \end{bmatrix} \equiv \begin{bmatrix} d_i^2 & zd_id_j \\ zd_id_j & d_j^2 \end{bmatrix},
$$

where we assume without loss of generality that $a \geq b$ and $c > 0$. By positive definiteness, $0 < z \leq \bar{z} \equiv (\kappa(A_m) - 1)/(\kappa(A_m) + 1) < 1$. Let $a'$ and $b'$ be the new values of $H_{ii}$ and $H_{jj}$ computed by the algorithm, respectively. Let $x \equiv d_j/d_i \leq 1$. We consider two cases: $x \leq \bar{x} \equiv (\sqrt{5} - 1)/2 \approx .62$, and $x > \bar{x}$.

First consider $x \leq \bar{x}$. Systematic application of formulas (3.1) shows that

$$
\begin{aligned}
\zeta &= fl((b - a)/(2 * c)) \\
&= (1 + \varepsilon_4)(((1 + \varepsilon_1)b - (1 + \varepsilon_2)a)/((1 + \varepsilon_3)2c)) \\
&= \frac{(1 + \varepsilon_4)(1 + \varepsilon_2)}{1 + \varepsilon_3} \left( \frac{\tilde{b} - a}{2c} \right),
\end{aligned}
$$

where $\tilde{b} \equiv (1 + \varepsilon_1)b/(1 + \varepsilon_2) \equiv (1 + \varepsilon_b)b$, $|\varepsilon_b| \leq 2\varepsilon$. Thus $\zeta = (1 + \varepsilon_\zeta)(\tilde{b} - a)/(2c)$ where $|\varepsilon_\zeta| \leq 3\varepsilon$.

Let $t(c)$ denote the true value of $t$ (i.e., without rounding error) as a function of $a$, $\tilde{b}$, and $c$. Using (3.1) again, we can show $t = (1 + \varepsilon_t)t(c)$ where $|\varepsilon_t| \leq 7\varepsilon$.

Next,

$$
\begin{aligned}
(3.4) \quad b' &= fl(b + ct) = (1 + \varepsilon_5)b + (1 + \varepsilon_6)(1 + \varepsilon_7)ct \\
&= \frac{(1 + \varepsilon_2)(1 + \varepsilon_5)}{1 + \varepsilon_1} \left( \tilde{b} + \frac{(1 + \varepsilon_1)(1 + \varepsilon_6)(1 + \varepsilon_7)(1 + \varepsilon_t)}{(1 + \varepsilon_2)(1 + \varepsilon_5)}ct(c) \right) \\
&\equiv (1 + \varepsilon_{b'})(\tilde{b} + (1 + \varepsilon_{ct(c)})ct(c)),
\end{aligned}
$$

where $|\varepsilon_{ct(c)}| \leq 12\varepsilon$ and $|\varepsilon_{b'}| \leq 3\varepsilon$. Since $|t(c)|$ is an increasing function of $c$, we can write $(1 + \varepsilon_{ct(c)})ct(c) = (1 + \varepsilon_c)c \cdot t((1 + \varepsilon_c)c)$ for some $\varepsilon_c$ where $|\varepsilon_c| \leq |\varepsilon_{ct(c)}| \leq 12\varepsilon$.

Now we can define $\tilde{c} \equiv (1 + \varepsilon_c)c$, and $\tilde{\zeta}$, $\tilde{t}$, $\tilde{cs}$, and $\tilde{sn}$ as the true values of the untilded quantities computed without rounding error starting from $a$, $\tilde{b}$, and $\tilde{c}$. $\tilde{cs}$ and $\tilde{sn}$ define the exact Jacobi rotation

$$
J_m \equiv \left[ \begin{array}{cc} \tilde{cs} & \tilde{sn} \\ -\tilde{sn} & \tilde{cs} \end{array} \right],
$$

which transforms $H'_m$ into $H_{m+1}$ in the commutative diagram in the statement of the theorem: $J_m^T H'_m J_m = H_{m+1}$.

Now we begin constructing $\delta H_m$. $\delta H_m$ is nonzero only in rows and columns $i$ and $j$. First, we compute its entries outside the $2 \times 2$ $(i, j)$ submatrix. Using (3.1) we can show $cs = (1 + \varepsilon_{cs})\tilde{cs}$ and $sn = (1 + \varepsilon_{sn})\tilde{sn}$, where $|\varepsilon_{cs}| \leq 22\varepsilon$ and $|\varepsilon_{sn}| \leq 30\varepsilon$. Now let $H'_{ik}$ and $H'_{jk}$ denote the updated quantities computed by the algorithm. Then

$$
\begin{aligned}
(3.5) \quad H'_{ik} &= fl(cs * H_{ik} - sn * H_{jk}) \\
&= (1 + \varepsilon_{10})(1 + \varepsilon_8)csH_{ik} - (1 + \varepsilon_9)(1 + \varepsilon_{11})snH_{jk} \\
&= (1 + \varepsilon_{10})(1 + \varepsilon_8)(1 + \varepsilon_{cs})\tilde{cs}H_{ik} - (1 + \varepsilon_9)(1 + \varepsilon_{11})(1 + \varepsilon_{sn})\tilde{sn}H_{jk} \\
&\equiv \tilde{cs}H_{ik} - \tilde{sn}H_{jk} + \epsilon(H'_{ik}).
\end{aligned}
$$

Similarly,

$$
\begin{aligned}
(3.6) \quad H'_{jk} &= fl(sn * H_{ik} + cs * H_{jk}) \\
&= (1 + \varepsilon_{14})(1 + \varepsilon_{12})(1 + \varepsilon_{sn})\tilde{sn}H_{ik} + (1 + \varepsilon_{13})(1 + \varepsilon_{15})(1 + \varepsilon_{cs})\tilde{cs}H_{jk} \\
&\equiv \tilde{sn}H_{ik} + \tilde{cs}H_{jk} + \epsilon(H'_{jk}).
\end{aligned}
$$

Now $x = d_j/d_i$ implies

$$\bar{\zeta} = \frac{b-a}{2\tilde{c}} = \frac{d_j^2 - d_i^2}{2\tilde{z}d_i d_j} = \frac{x^2 - 1}{2\tilde{z}x},$$

where $\tilde{z} \equiv z(1 + \varepsilon_c)$. Then $x \leq \bar{x}$ implies

$$|\tilde{t}| = \frac{1}{\frac{1-x^2}{2\tilde{z}x} + \left(1 + \left(\frac{1-x^2}{2\tilde{z}x}\right)^2\right)^{1/2}} \leq \frac{\tilde{z}x}{1 - \bar{x}^2}.$$

Also, $|\tilde{sn}| \leq |\tilde{t}|$, so this last expression is an upper bound on $|\tilde{sn}|$ as well. Substituting this bound on $\tilde{sn}$, $\tilde{cs} \leq 1$, $|H_{ik}| \leq d_i d_k \bar{z}$, and $|H_{jk}| \leq d_j d_k \bar{z}$ into (3.5) and (3.6) yields

$$|\epsilon(H'_{ik})| \leq 56\varepsilon d_i d_k \bar{z},$$
$$|\epsilon(H'_{jk})| \leq 56\varepsilon d_j d_k \bar{z}/(1 - \bar{x}^2).$$

Thus

$$\begin{bmatrix} H'_{ik} \\ H'_{jk} \end{bmatrix} = J_m^T \cdot \begin{bmatrix} H_{ik} \\ H_{jk} \end{bmatrix} + \begin{bmatrix} \epsilon(H_{ik}) \\ \epsilon(H_{jk}) \end{bmatrix}$$
$$= J_m^T \cdot \left(\begin{bmatrix} H_{ik} \\ H_{jk} \end{bmatrix} + J_m \cdot \begin{bmatrix} \epsilon(H_{ik}) \\ \epsilon(H_{jk}) \end{bmatrix}\right) \equiv J_m^T \cdot \left(\begin{bmatrix} H_{ik} \\ H_{jk} \end{bmatrix} + \begin{bmatrix} \delta H_{ik} \\ \delta H_{jk} \end{bmatrix}\right),$$

where $|\delta H_{ik}| \leq 112\varepsilon d_i d_k \bar{z}/(1 - \bar{x}^2)$ and $|\delta H_{jk}| \leq 112\varepsilon d_j d_k \bar{z}/(1 - \bar{x}^2)$.

Now we construct the $2 \times 2$ submatrix $\Delta$ of $\delta H_m$ at the intersection of rows and columns $i$ and $j$. We construct it of three components: $\Delta = \Delta_1 + \Delta_2 + \Delta_3$.

Consider the formula $a' = fl(a - c*t)$ for the $i,i$ entry of $H_{m+1}$. Applying (3.1) systematically, we see that

$$a' = (1 + \varepsilon_{18})a - (1 + \varepsilon_{17})(1 + \varepsilon_{16})ct$$
$$= (1 + \varepsilon_{18})a - (1 + \varepsilon_{17})(1 + \varepsilon_{16})(1 + \varepsilon_t)ct(c)$$
$$= (1 + \varepsilon_{18})a - \frac{(1 + \varepsilon_{17})(1 + \varepsilon_{16})(1 + \varepsilon_t)\tilde{c}t(\tilde{c})}{1 + \varepsilon_{ct(c)}}$$
$$\equiv (1 + \varepsilon_{18})a - (1 + \varepsilon'_{ct(c)})\tilde{c}t(\tilde{c}),$$

where $|\varepsilon'_{ct(c)}| \leq 21\varepsilon$. Since $a > 0$ and $\tilde{c}t(\tilde{c}) < 0$, we get

$$a' = \left(1 + \frac{\varepsilon_{18}a - \varepsilon'_{ct(c)}\tilde{c}t(\tilde{c})}{a - \tilde{c}t(\tilde{c})}\right)(a - \tilde{c}t(\tilde{c})) \equiv (1 + \varepsilon_{a'})(a - \tilde{c}t(\tilde{c})),$$

where $|\varepsilon_{a'}| \leq 21\varepsilon$.

Now let

$$\Delta_1 = \begin{bmatrix} 0 & \varepsilon_c c \\ \varepsilon_c c & \varepsilon_b b \end{bmatrix} = \begin{bmatrix} 0 & \tilde{c} - c \\ \tilde{c} - c & \tilde{b} - b \end{bmatrix}.$$

From earlier discussion we see that

$$J_m^T\left(\begin{bmatrix} a & c \\ c & b \end{bmatrix} + \Delta_1\right)J_m = \begin{bmatrix} a - \tilde{c}t(\tilde{c}) & 0 \\ 0 & \tilde{b} + \tilde{c}t(\tilde{c}) \end{bmatrix}.$$

Next let

$$\Delta_2 = \varepsilon_{a'} \left( \begin{bmatrix} a & c \\ c & b \end{bmatrix} + \Delta_1 \right).$$

Thus

$$J_m^T \left( \begin{bmatrix} a & c \\ c & b \end{bmatrix} + \Delta_1 + \Delta_2 \right) J_m = (1 + \varepsilon_{a'}) \begin{bmatrix} a - \tilde{c}t(\tilde{c}) & 0 \\ 0 & \tilde{b} + \tilde{c}t(\tilde{c}) \end{bmatrix}$$

$$= \begin{bmatrix} a' & 0 \\ 0 & b'((1 + \varepsilon_{a'})/(1 + \varepsilon_b)) \end{bmatrix}.$$

Finally, let

$$\Delta_3 = J_m \begin{bmatrix} 0 & 0 \\ 0 & b'(1 - ((1 + \varepsilon_{a'})/(1 + \varepsilon_b))) \end{bmatrix} J_m^T \equiv \begin{bmatrix} \tilde{s}n^2 \varepsilon_{b''} b & \tilde{c}s\tilde{s}n\varepsilon_{b''} b \\ \tilde{c}s\tilde{s}n\varepsilon_{b''} b & \tilde{c}s^2 \varepsilon_{b''} b \end{bmatrix},$$

where $|\varepsilon_{b''}| \le |\varepsilon_{a'}| + |\varepsilon_b| \le 24\varepsilon$. Then

$$J_m^T \left( \begin{bmatrix} a & c \\ c & b \end{bmatrix} + \Delta_1 + \Delta_2 + \Delta_3 \right) J_m = \begin{bmatrix} a' & 0 \\ 0 & b' \end{bmatrix},$$

as desired. This completes the construction of $\delta H_m$. We may bound

$$(3.7) \qquad \|\delta A_m\|_2 \le \left( \frac{112(2n - 4)^{1/2}\bar{z}}{1 - \bar{x}^2} + 104 \right)\varepsilon.$$

Now we consider the second case, when $x > \bar{x}$. The only thing that changes in the previous analysis is our analysis of $\delta H_{ik}$ and $\delta H_{jk}$, since $\tilde{s}n$ is no longer small. Instead we substitute the bounds $|\tilde{s}n| \le 1$, $|\tilde{c}s| \le 1$, $|H_{ik}| \le d_i d_k \bar{z} \le d_j d_k \bar{z}/\bar{x}$, and $|H_{jk}| \le d_j d_k \bar{z}$ into (3.5) and (3.6) to get

$$|\epsilon(H'_{ik})| \le 56\varepsilon d_i d_k \bar{z} \quad \text{and} \quad |\epsilon(H'_{jk})| \le 56\varepsilon d_i d_k \bar{z},$$

whence

$$|\delta H_{ik}| \le 112\varepsilon d_i d_k \bar{z} \quad \text{and} \quad |\delta H_{jk}| \le 112\varepsilon d_j d_k \bar{z}/\bar{x},$$

and so

$$(3.8) \qquad \|\delta A_m\|_2 \le \left( \frac{112(2n - 4)^{1/2}\bar{z}}{\bar{x}} + 104 \right)\varepsilon.$$

Finally, we note that our choice of $\bar{x}$ makes the upper bounds in (3.7) and (3.8) both equal, with $1/(1 - \bar{x}^2) = 1/\bar{x} < 1.62$, proving the theorem. □

*Remark.* The quantity $182(2n - 4)^{1/2}$ in the theorem may be multiplied by $\max_{m, i \ne j} |A_{m,ij}|$, which is less than one. Thus if the $A_m$ are strongly diagonally dominant, the part of the error term that depends on $n$ is suppressed.

Commutative diagrams like the one in the theorem, where performing one step of the algorithm in floating point arithmetic is equivalent to making small relative errors in the matrix and then performing the algorithm exactly, occur elsewhere in numerical analysis. For example, such a diagram describes an entire sweep of the zero-shift bidiagonal QR algorithm [8], and is the key to the high accuracy achieved by that algorithm.

**3.2. Error bounds for eigenvectors computed by two-sided Jacobi.** The next two theorems justify our accuracy claims for eigenvectors computed by two-sided Jacobi.

THEOREM 3.3. *Let $V = [v_1, \cdots, v_n]$ be the matrix of unit eigenvectors computed by Algorithm 3.1 in finite precision arithmetic with precision $\varepsilon$. Let $U = [u_1, \cdots, u_n]$ be the true eigenvector matrix. Let $\bar{\kappa} \equiv \max_m \kappa(A_m)$ be the largest $\kappa(A_m)$ of any iterate. Then the error in the computed eigenvectors is bounded in norm by*

$$(3.9) \quad \|v_i - u_i\|_2 \leq \frac{(n-1)^{1/2}(n \cdot \text{tol} + M \cdot (182(2n-4)^{1/2} + 104)\varepsilon)\bar{\kappa}}{\text{relgap}_{\lambda_i}} + 46M\varepsilon.$$

*Proof.* Let $H_0, \cdots, H_M$ be the sequence of matrices generated by the Jacobi algorithm, where $H_M$ satisfies the stopping criterion. Let $J_m$ be the exact Jacobi rotation which transforms $H_m'$ to $H_{m+1}$ in the commuting diagram of Theorem 3.1: $J_m^T H_m' J_m = H_{m+1}$.

We use the approximation that $\text{relgap}_{\lambda_i}$ is the same for all $H_m$, even though it changes slightly. This contributes an $O(\varepsilon^2)$ term to the overall bound (which we ignore), but could be accounted for by using the bounds of Theorem 3.1.

Initially, we compute error bounds for the columns of $J_0 \cdots J_{M-1}$, ignoring any rounding errors occurring in computing their product. Then we incorporate these rounding errors.

We prove by induction that the $i$th column $v_{mi}$ of $V_m \equiv J_m \cdots J_{M-1}$ is a good approximation to the true $i$th eigenvector $u_{mi}$ of $H_m$. In particular, we show that to first order in $\varepsilon$,

$$\|u_i - v_{0i}\|_2 \leq \frac{(n-1)^{1/2}(n \cdot \text{tol} + M \cdot (182(2n-4)^{1/2} + 104)\varepsilon)\bar{\kappa}}{\text{relgap}_{\lambda_i}}.$$

The basis of the induction is as follows. $V_M = I$ is the eigenvector matrix for $H_M$, which is considered diagonal since it satisfies the stopping criterion. Thus the norm error in $v_{Mi}$ follows from plugging the stopping criterion into Theorem 2.5:

$$\|u_{Mi} - v_{Mi}\|_2 \leq \frac{(n-1)^{1/2} \cdot n \cdot \text{tol} \cdot \bar{\kappa}}{\text{relgap}_{\lambda_i}}.$$

For the induction step we assume that

$$\|u_{m+1,i} - v_{m+1,i}\|_2 \leq \frac{(n-1)^{1/2}(n \cdot \text{tol} + (M-m-1) \cdot (182(2n-4)^{1/2} + 104)\varepsilon)\bar{\kappa}}{\text{relgap}_{\lambda_i}},$$

and try to extend it to $m$. Consider the commuting diagram of Theorem 3.1. Accordingly, the errors in $V_m = J_m V_{m+1}$ considered as eigenvectors of $H_m'$ are the errors in $V_{m+1}$ premultiplied by $J_m$. This does not increase them in the two-norm, since $J_m$ is orthogonal. Now we change $H_m'$ to $H_m$. This increases the norm error in $v_{mi}$ by an amount bounded by plugging the bound for $\|\delta A_m\|_2$ into Theorem 2.5: $(n-1)^{1/2}\bar{\kappa}(182(2n-4)^{1/2} + 104)\varepsilon/\text{relgap}_{\lambda_i}$. This proves the induction step.

Finally, consider the errors from accumulating the product of slightly wrong values of $J_m$ in floating point arithmetic. From the proof of Theorem 3.1, we see that the relative errors in the entries of $J_m$ are at most $30\varepsilon$, and from the usual error analysis of a product of $2 \times 2$ rotations, we get $32\sqrt{2}M\varepsilon < 46M\varepsilon$ for the norm error in the product of $M$ rotations. This completes the proof of bound (3.9).    □

Now we consider the errors in the individual components of the computed eigenvectors $|u_i(j) - v_i(j)|$. From Proposition 2.9, we see that we can hope to bound this quantity by $O(\varepsilon)\bar{\kappa}\bar{v}_i(j)/\min(\text{relgap}_{\lambda_i}, 2^{-1/2})$, where

$$(3.10) \qquad \bar{v}_i(j) \equiv \bar{\kappa}^{3/2} \min\left(\left(\frac{\lambda_i}{\lambda_j}\right)^{1/2}\left(\frac{\lambda_j}{\lambda_i}\right)^{1/2}\right)$$

is a modified upper bound for the eigenvector component $v_i(j)$ as in Proposition 2.8. In other words, we may have high relative accuracy even in the tiny components of the computed eigenvectors; this is the case in the example in the introduction and at the end of §2.1.

We use $\bar{v}_i(j)$ as defined in (3.10) for each $H_m$, even though the values of $\lambda_i$ and $\lambda_j$ vary slightly from step to step. This error contributes an $O(\varepsilon^2)$ term to the overall bound (we are ignoring such terms), but could be incorporated using the bounds of Corollary 3.2.

THEOREM 3.4. *Let $V$, $U$, and $\bar{\kappa}$ be as in Theorem 3.3, and $\bar{v}_i(j)$ be as in (3.10). Then we can bound the error in the individual eigencomponents by*

$$(3.11) \qquad |u_i(j) - v_i(j)| \leq p(M, n) \cdot \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})}.$$

*Here $p(M, n)$ is a "pivot growth" factor, which is given in (3.21).*

*Proof.* The proof is similar to that of Theorem 3.3. One difference is that we use Proposition 2.9 instead of Theorem 2.5 to bound the errors in the eigenvectors. Another difference, which introduces the growth factor $p(M, n)$, is that we need to use the scaling of the entries of $J_m$ to see how small eigenvector components have small errors; not being able to use the orthogonality of $J_m$ introduces $p(M, n)$. We can only prove an exponentially growing bound for $p(M, n)$, although we believe it to be much smaller.

As in the proof of Theorem 3.3, let $V_m = J_m \cdots J_{M-1}$, where $J_m^T H_m' J_m = H_{m+1}$. Set $V_M = I$. The proof has three parts. In the first part, we show that the $i$th column of $V_0$ is a good approximation to the eigenvectors of $H_0$ in the sense of the theorem. In the second part, we show that the $(i, j)$ entry of $J_0 \cdots J_m$ is bounded by a modest multiple of $\bar{v}_i(j)$. In the third part, we show that the rounding errors committed in computing $J_0 \cdots J_m$ in floating point are small compared to $\bar{v}_i(j)$.

For the first part of the proof we use induction to prove that the $i$th column $v_{mi}$ of $V_m$ is a good approximation to the true $i$th eigenvector $u_{mi}$ of $H_m$. This shows that

$$(3.12) \qquad |u_i(j) - v_{0i}(j)| \leq \rho_0 \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})},$$

where $\rho_0$ is a constant (part of the "pivot growth" factor) we need to estimate. The base of the induction follows from plugging the stopping criterion into the bound of Proposition 2.9, yielding

$$|u_{Mi}(j) - v_{Mi}(j)| \leq \frac{(2n-2)^{1/2} \cdot n \cdot \text{tol} \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})} \leq \rho_M \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})},$$

where $\rho_M \equiv n(2n-2)^{1/2}$. The induction step assumes that

$$|u_{m+1,i}(j) - v_{m+1,i}(j)| \leq \rho_{m+1} \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})},$$

which we try to extend to $m$. Consider the commuting diagram of Theorem 3.1. Accordingly, the errors in the columns of $V_m = J_m V_{m+1}$ considered as eigenvectors of $H_m'$ are just the errors in $V_{m+1}$ premultiplied by $J_m$; let $e_{mi}$ denote this error for the $i$th column of $V_m$. Suppose $J_m$ rotates in rows and columns $k$ and $l$; then $e_{mi}$ is identical to $u_{m+1,i} - v_{m+1,i}$ except for $e_{mi}(k)$ and $e_{mi}(l)$. We may assume without loss of generality that $k < l$ and $d_k \geq d_l$ ($d_k^2$ and $d_l^2$ are the diagonal entries of $H_m$). As in the proof of Theorem 3.1, there are two cases: $x \equiv d_l/d_k \leq \bar{x} \equiv (\sqrt{5} - 1)/2$, and $x > \bar{x}$.

In the first case, $x \leq \bar{x}$, we know (as in the proof of Theorem 3.1) that $\tilde{sn}$, the sine in the rotation $J_m$, is bounded in magnitude by $x/(1 - \bar{x}^2)$. Write $|\tilde{sn}| \leq c_m(\lambda_l/\lambda_k)^{1/2}$ instead, where $c_m$ is a modest constant. We can do this because $d_r \approx \lambda_r^{1/2}$ from Proposition 2.7. This lets us bound

$$
\begin{bmatrix} |e_{mi}(k)| \\ |e_{mi}(l)| \end{bmatrix} = \left| J_m \begin{bmatrix} u_{m+1,i}(k) - v_{m+1,i}(k) \\ u_{m+1,i}(l) - v_{m+1,i}(l) \end{bmatrix} \right|
$$

$$
\leq \begin{bmatrix} |u_{m+1,i}(k) - v_{m+1,i}(k)| + |\tilde{sn}(u_{m+1,i}(l) - v_{m+1,i}(l))| \\ |\tilde{sn}(u_{m+1,i}(k) - v_{m+1,i}(k))| + |u_{m+1,i}(l) - v_{m+1,i}(l)| \end{bmatrix}
$$

$$
\leq \frac{\rho_{m+1}(\text{tol} + \varepsilon)\bar{\kappa}^{5/2}}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})}
$$

$$
\times \begin{bmatrix} \min\left(\left(\frac{\lambda_i}{\lambda_k}\right)^{1/2}, \left(\frac{\lambda_k}{\lambda_i}\right)^{1/2}\right) + c_m\left(\frac{\lambda_l}{\lambda_k}\right)^{1/2} \min\left(\left(\frac{\lambda_i}{\lambda_l}\right)^{1/2}, \left(\frac{\lambda_l}{\lambda_i}\right)^{1/2}\right) \\ c_m\left(\frac{\lambda_l}{\lambda_k}\right)^{1/2} \min\left(\left(\frac{\lambda_i}{\lambda_k}\right)^{1/2}, \left(\frac{\lambda_k}{\lambda_i}\right)^{1/2}\right) + \min\left(\left(\frac{\lambda_i}{\lambda_l}\right)^{1/2}, \left(\frac{\lambda_l}{\lambda_i}\right)^{1/2}\right) \end{bmatrix}
$$

$$
\leq \frac{\rho_{m+1}(\text{tol} + \varepsilon)\bar{\kappa}^{5/2}}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})} \cdot \begin{bmatrix} (1 + c_m) \min\left(\left(\frac{\lambda_i}{\lambda_k}\right)^{1/2}, \left(\frac{\lambda_k}{\lambda_i}\right)^{1/2}\right) \\ (1 + c_m) \min\left(\left(\frac{\lambda_i}{\lambda_l}\right)^{1/2}, \left(\frac{\lambda_l}{\lambda_i}\right)^{1/2}\right) \end{bmatrix}
$$

$$
(3.13) \qquad = (1 + c_m)\frac{\rho_{m+1}(\text{tol} + \varepsilon)\bar{\kappa}}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})} \cdot \begin{bmatrix} \bar{v}_i(k) \\ \bar{v}_i(l) \end{bmatrix}.
$$

Now consider case 2, $x > \bar{x}$. Now $\lambda_k$ and $\lambda_l$ are reasonably close together. Thus we may bound $|\tilde{sn}|$ simply by 1 in the derivation (3.13). This leads to the same bound with a possibly different $c_m$; we take the final $c_m$ as the maximum of these two values. This bounds the error in the columns of $V_m$ considered as eigenvectors of $H_m'$.

Now we change $H_m'$ to $H_m$. This increases the bound for $|u_{mi}(j) - v_{mi}(j)|$ by an amount bounded by plugging the bound for $\|\delta A_m\|_2$ from Theorem 3.1 into Proposition 2.9: $(2n - 2)^{1/2}(182(2n - 4)^{1/2} + 104) \cdot \varepsilon \cdot \bar{\kappa} \cdot \bar{v}_i(j)/\min(\text{relgap}_{\lambda_i}, 2^{-1/2})$. This completes the induction with

$$
(3.14)
$$

$$
|u_{mi}(j) - v_{mi}(j)|
$$

$$
\leq ((1 + c_m)\rho_{m+1} + (2n - 2)^{1/2}(182(2n - 4)^{1/2} + 104)) \cdot \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})}
$$

$$
\equiv \rho_m \cdot \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})}.
$$

Here

(3.15)

$$\rho_m = (1 + c_m)\rho_{m+1} + (2n - 2)^{1/2}(182(2n - 4)^{1/2} + 104) \ , \quad \rho_M = n(2n - 2)^{1/2}.$$

$\rho_m$ satisfies an exponential error bound, but it is clear from the derivation that linear growth is far more likely than exponential growth. This completes the first part of the proof.

In the second part of the proof we show that the $(i, j)$ entry of $\tilde{V}_m \equiv J_0 \cdots J_m$ is bounded by a modest multiple of $\bar{v}_i(j)$. To do this we prove by induction that

(3.16) $$|\tilde{v}_{mi}(j)| \le \tau_m \bar{v}_i(j),$$

where $\tilde{V}_m = [\tilde{v}_{m1}, \cdots, \tilde{v}_{mn}]$ and $\tau_m$ is a constant (part of the "pivot growth" factor) we need to estimate. The base of the induction is for $m = -1$, i.e., the null product, which we set equal to the identity matrix. This clearly satisfies (3.16) with $\tau_{-1} = 1$. Now we assume that (3.16) is true for $m - 1$ and try to extend it to $m$. Suppose $J_m$ rotates in rows and columns $k$ and $l$. Postmultiplying $\tilde{V}_{m-1}$ by $J_m$ only changes it in columns $k$ and $l$. Assume as before that $k < l$ and $x = d_l/d_k \le 1$. There are two cases, as before: $x \le \bar{x}$ and $x > \bar{x}$.

First, consider the case $x \le \bar{x}$. We may bound the $(j, k)$ and $(j, l)$ entries of $\tilde{V}_m$ as follows:

(3.17)

$$|[\tilde{v}_{m,j}(k), \quad \tilde{v}_{m,j}(l)]| = \left|[\tilde{v}_{m-1,j}(k), \tilde{v}_{m-1,j}(l)] \cdot \begin{bmatrix} \tilde{c}s & \tilde{s}n \\ -\tilde{s}n & \tilde{c}s \end{bmatrix}\right|$$

$$\le \tau_{m-1}\bar{\kappa}^{3/2}\left[\min\left(\left(\frac{\lambda_j}{\lambda_k}\right)^{1/2}, \left(\frac{\lambda_k}{\lambda_j}\right)^{1/2}\right), \min\left(\left(\frac{\lambda_j}{\lambda_l}\right)^{1/2}, \left(\frac{\lambda_l}{\lambda_j}\right)^{1/2}\right)\right]$$

$$\times \begin{bmatrix} 1 & c_m\left(\frac{\lambda_l}{\lambda_k}\right)^{1/2} \\ c_m\left(\frac{\lambda_l}{\lambda_k}\right)^{1/2} & 1 \end{bmatrix}$$

$$\le \tau_{m-1}\bar{\kappa}^{3/2}(1 + c_m)$$

$$\times \left[\min\left(\left(\frac{\lambda_j}{\lambda_k}\right)^{1/2}, \left(\frac{\lambda_k}{\lambda_j}\right)^{1/2}\right), \min\left(\left(\frac{\lambda_j}{\lambda_l}\right)^{1/2}, \left(\frac{\lambda_l}{\lambda_j}\right)^{1/2}\right)\right]$$

$$= \tau_{m-1}(1 + c_m)[\bar{v}_k(j), \bar{v}_l(j)]$$

$$\equiv \tau_m[\bar{v}_k(j), \bar{v}_l(j)].$$

In the second case, $x > \bar{x}$, we get a similar bound. Here $\lambda_k \approx \lambda_l$, and we can simply bound $|\tilde{s}n| \le 1$. This yields a slightly different $c_m$; for the final $c_m$, we again take the maximum of the two. This ends the second part of the proof with

(3.18) $$\tau_m = (1 + c_m)\tau_{m-1}, \qquad \tau_{-1} = 1.$$

Even though this only yields an exponential upper bound for $\tau_M$, it is clear from the derivation that linear growth is far more likely than exponential growth.

In the third and final part of the proof, we show that the rounding errors in the $(i, j)$ entry of the computed approximation to $\tilde{V}_{m-1}$ is bounded by $O(\varepsilon)\bar{v}_i(j)$. Let $\tilde{J}_m$ be the actual rotation which only approximates $J_m$. From the proof of Theorem 3.1, we have that $cs = \tilde{c}s(1 + \varepsilon_{cs})$ with $|\varepsilon_{cs}| \le 22\varepsilon$, and $sn = \tilde{s}n(1 + \varepsilon_{sn})$ with $|\varepsilon_{sn}| \le 30\varepsilon$.

Let $\tilde{\tilde{V}}_m = fl(\tilde{\tilde{V}}_{m-1} * \tilde{J}_m)$ be the actually computed eigenvector matrix after the $m$th Jacobi rotation. The final computed eigenvector matrix is $V = \tilde{\tilde{V}}_{M-1}$. We use induction to prove that

$$(3.19) \qquad |\tilde{\tilde{v}}_{m,i}(j) - \tilde{v}_{m,i}(j)| \leq \chi_m \varepsilon \bar{v}_i(j),$$

where $\tilde{\tilde{V}}_m = [\tilde{\tilde{v}}_{m1}, \cdots, \tilde{\tilde{v}}_{mn}]$ and $\chi_m$ is a constant (part of the "pivot growth" factor) we need to estimate. The basis is again for $m = -1$ when $\tilde{\tilde{V}}_{-1} = \tilde{V}_{-1} = I$ and $\chi_{-1} = 0$. Now we assume that (3.19) is true for $m - 1$ and try to extend it to $m$. As before, we assume that $J_m$ rotates in rows and columns $k$ and $l$ with $k < l$ and $x = d_k/d_l \leq 1$. Write $\tilde{e}_{mi} \equiv \tilde{\tilde{v}}_{mi} - \tilde{v}_{mi}$. The $(j,k)$ and $(j,l)$ entries of $\tilde{\tilde{V}}_m$ are

$$
\begin{aligned}
[\tilde{\tilde{v}}_{m,j}(k), \tilde{\tilde{v}}_{m,j}(l)] &= [\tilde{\tilde{v}}_{m-1,j}(k)\tilde{cs}(1+\varepsilon_{cs})(1+\varepsilon_1)(1+\varepsilon_2) \\
&\qquad - \tilde{\tilde{v}}_{m-1,j}(l)\tilde{sn}(1+\varepsilon_{sn})(1+\varepsilon_3)(1+\varepsilon_4), \\
&\qquad \tilde{\tilde{v}}_{m-1,j}(k)\tilde{sn}(1+\varepsilon_{sn})(1+\varepsilon_5)(1+\varepsilon_6) \\
&\qquad + \tilde{\tilde{v}}_{m-1,j}(l)\tilde{cs}(1+\varepsilon_{cs})(1+\varepsilon_7)(1+\varepsilon_8)] \\
&= [\tilde{v}_{m-1,j}(k)\tilde{cs} - \tilde{v}_{m-1,j}(l)\tilde{sn}, \tilde{v}_{m-1,j}(k)\tilde{sn} + \tilde{v}_{m-1,j}(l)\tilde{cs}] \\
&\qquad + [24\varepsilon_9 \tilde{cs}\tilde{v}_{m-1,j}(k) + 32\varepsilon_{10}\tilde{sn}\tilde{v}_{m-1,j}(l), \\
&\qquad\quad 32\varepsilon_{11}\tilde{sn}\tilde{v}_{m-1,j}(k) + 24\varepsilon_{12}\tilde{cs}\tilde{v}_{m-1,j}(l)] \\
&\qquad + [(1+24\varepsilon_9)\tilde{cs}\tilde{e}_{m-1,j}(k) + (1+32\varepsilon_{10})\tilde{sn}\tilde{e}_{m-1,j}(l), \\
&\qquad\quad (1+32\varepsilon_{11})\tilde{sn}\tilde{e}_{m-1,j}(k) + (1+24\varepsilon_{12})\tilde{cs}\tilde{e}_{m-1,j}(l)] \\
&= [\tilde{v}_{m,j}(k), \tilde{v}_{m,j}(l)] + I_1 + I_2,
\end{aligned}
$$

so that $[\tilde{e}_{m,j}(k), \tilde{e}_{m,j}(l)] = I_1 + I_2$.

As before, there are two cases: $x \leq \bar{x}$ and $x > \bar{x}$. Consider the case $x \leq \bar{x}$. Using $|\tilde{sn}| \leq c_m(\lambda_l/\lambda_k)^{1/2}$, $|\tilde{cs}| \leq 1$, and $|\tilde{v}_{m-1,i}(j)| \leq \tau_{m-1}\bar{v}_i(j)$, we get

$$
\begin{aligned}
|I_1| &\leq \varepsilon\tau_{m-1}\bar{\kappa}^{3/2}(24 + 32c_m) \\
&\qquad \times \left[ \min\left( \left(\frac{\lambda_j}{\lambda_k}\right)^{1/2}, \left(\frac{\lambda_k}{\lambda_j}\right)^{1/2} \right), \min\left( \left(\frac{\lambda_j}{\lambda_l}\right)^{1/2}, \left(\frac{\lambda_l}{\lambda_j}\right)^{1/2} \right) \right] \\
&= \varepsilon\tau_{m-1}(24 + 32c_m)[\bar{v}_k(j), \bar{v}_k(l)],
\end{aligned}
$$

and

$$|I_2| \leq \chi_{m-1}(1 + c_m)\varepsilon[\bar{v}_k(j), \bar{v}_k(l)].$$

Taken together, we get

$$(3.20) \qquad \chi_m = (1 + c_m)\chi_{m-1} + \varepsilon\tau_{m-1}(24 + 32c_m), \qquad \chi_{-1} = 0.$$

In the second case, $x > \bar{x}$, we get a similar bound with a possibly different $c_m$. Again, we take the maximum of the two. This completes the third part of the proof.

Finally, combining (3.19) and (3.12) we get

$$(3.21) \qquad |v_i(j) - u_i(j)| \leq (\rho_0 + \chi_{M-1})\frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa} \cdot \bar{v}_i(j)}{\min(\text{relgap}_{\lambda_i}, 2^{-1/2})},$$

proving the theorem with $p(M, n) = \rho_0 + \chi_{M-1}$.  $\quad\square$

**4. One-sided Jacobi.** In this section we prove that one-sided Jacobi in floating point arithmetic applied on the right of a general matrix computes the singular values and singular vectors with the error bounds of §2. Here we present our algorithm; the model of arithmetic was presented in §3. In §4.1 we derive error bounds for the computed singular values. In §4.2 we derive error bounds for the computed singular vectors. In §4.3, we present two algorithms for the symmetric positive definite eigenproblem $H$, which do either left-handed or right-handed Jacobi on the Cholesky factor $L$ of $H$. The second of these algorithms cannot compute eigenvectors quite as accurately as the first, but it may be much faster than either the first algorithm or two-sided Jacobi.

Let $G_0 = B_0 D_0$ be the initial matrix, and $G_m = B_m D_m$, where $G_m$ is obtained from $G_{m-1}$ by applying a single Jacobi rotation. Here $D_m$ is diagonal and $B_m$ has columns of unit two-norm. All the error bounds in this section contain the factor $\max_m \kappa(B_m)$, whereas the perturbation bounds in §2 are proportional to $\kappa(B_0)$. Therefore, as in §3, our claim that Jacobi computes the singular value decomposition (SVD) as accurately as predicted in §2 depends on the ratio $\max_m \kappa(B_m)/\kappa(B_0)$ being modest. In exact arithmetic, right-handed Jacobi on $G = BD$ is identical to two-sided Jacobi on $H = G^T G = D B^T B D = DAD$, so the question of the growth of $\kappa(B_m) = \kappa(A_m)^{1/2}$ is essentially identical to the question of the growth of $\kappa(A_m)$ in the case of two-sided Jacobi.

ALGORITHM 4.1 (Right-handed Jacobi for the singular value problem). tol is a user-defined stopping criterion. The matrix $V$ whose columns are the computed right singular vectors initially contains the identity.
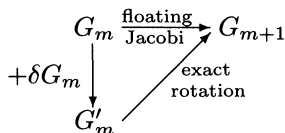
repeat
    for all pairs $i < j$
        /* compute $\begin{bmatrix} a & c \\ c & b \end{bmatrix} \equiv$ the $(i, j)$ submatrix of $G^T G$ */
        $a = \sum_{k=1}^{n} G_{ki}^2$
        $b = \sum_{k=1}^{n} G_{kj}^2$
        $c = \sum_{k=1}^{n} G_{ki} * G_{kj}$
        /* compute the Jacobi rotation which diagonalizes $\begin{bmatrix} a & c \\ c & b \end{bmatrix}$ */
        $\zeta = (b - a)/(2c);\ t = \text{sign}(\zeta)/(|\zeta| + \sqrt{1 + \zeta^2})$
        $cs = 1/\sqrt{1 + t^2};\ sn = cs * t$
        /* update columns $i$ and $j$ of $G$ */
        for $k = 1$ to $n$
            $tmp = G_{ki}$
            $G_{ki} = cs * tmp - sn * G_{kj}$
            $G_{kj} = sn * tmp + cs * G_{kj}$
        endfor
        /* update the matrix $V$ of right singular vectors */
        for $k = 1$ to $n$
            $tmp = V_{ki}$
            $V_{ki} = cs * tmp - sn * V_{kj}$
            $V_{kj} = sn * tmp + cs * V_{kj}$
        endfor
    endfor
until convergence (all $|c|/\sqrt{ab} \leq tol$)
/* the computed singular values are the norms of the columns of the final $G$ */

/* the computed left singular vectors are the normalized columns of the final $G$ */

### 4.1. Error bounds for singular values computed by right-handed Jacobi.
The next theorem and its corollary justify our accuracy claims for singular values computed by right-handed Jacobi. The proofs are analogous to those in §3; details may be found in [10].

THEOREM 4.1. *Let $G_m$ be the sequence of matrices generated by the right-handed Jacobi algorithm in finite precision arithmetic with precision $\varepsilon$; that is $G_{m+1}$ is obtained from $G_m$ by applying a single Jacobi rotation. Then the following diagram:*

$$G_m \xrightarrow[\text{Jacobi}]{\text{floating}} G_{m+1}$$

with $+\delta G_m$ on the left, $G'_m$ at bottom, and "exact rotation" on the diagonal arrow.

*commutes in the same way as in Theorem 3.1. The diagonal arrow indicates that $G_{m+1}$ is obtained from $G'_m$ by applying one Givens (not necessarily Jacobi) rotation in exact arithmetic. $\delta G_m$ is bounded as follows. Write $\delta G_m = \delta B_m D_m$, where $D_m$ is diagonal such that $B_m$ in $G_m = B_m D_m$ has unit columns. Then*

$$(4.1) \qquad \|\delta B_m\|_2 \le 72\varepsilon.$$

*In other words, one step of Jacobi satisfies the assumptions needed for the error bounds of §2.*

COROLLARY 4.2. *Assume that Algorithm 4.1 converges, and that $G_M$ is the final matrix which satisfies the stopping criterion. For $0 \le m \le M$, write $G_m = B_m D_m$ with $D_m$ diagonal and $B_m$ with unit columns. Let $\sigma_j$ be the $j$th singular value of $G_0$ and $\sigma'_j$ the $j$th computed singular value. Then to first order in $\varepsilon$ the following error bound holds:*

$$(4.2) \qquad \frac{|\sigma_j - \sigma'_j|}{\sigma_j} \le (72\varepsilon \cdot M + n^2\varepsilon + n \cdot \text{tol}) \cdot \max_{0 \le k \le M} \kappa(B_k) + n\varepsilon.$$

*Remark.* A similar bound can be obtained based on the error bound in Proposition 2.15.

### 4.2. Error bounds for singular vectors computed by right-handed Jacobi.
The next two theorems justify our accuracy claims for singular vectors computed by right-handed Jacobi.

THEOREM 4.3. *Let $V = [v_1, \cdots, v_n]$ be the matrix of unit right singular vectors and $U = [u_1, \cdots, u_n]$ be the matrix of unit left singular vectors computed by Algorithm 4.1 in finite precision arithmetic with precision $\varepsilon$. Let $V_T = [v_{T1}, \cdots, v_{Tn}]$ and $U_T = [u_{T1}, \cdots, u_{Tn}]$ be the matrices of true unit right and left singular vectors, respectively. Let $\bar\kappa \equiv \max_m \kappa(B_m)$ be the largest $\kappa(B_m)$ of any iterate. Then the error in the computed singular vectors is bounded in norm by*

$$\max(\|u_{Ti} - u_i\|_2, \|v_{Ti} - v_i\|_2)$$
$$(4.3) \qquad \le \frac{(n-.5)^{1/2} \cdot \bar\kappa \cdot (72M \cdot \varepsilon + n \cdot \text{tol} + n^2 \cdot \varepsilon)}{\text{relgap}_{\sigma i}} + (9M + n + 1)\varepsilon.$$

Then consider the errors in the individual components of the computed right singular vectors $|v_{Ti}(j) - v_i(j)|$. From Proposition 2.20, we see that we can hope to bound this quantity by $O(\varepsilon)\bar{\kappa}^2\bar{\bar{v}}_i(j)/\text{relgap}_{\sigma i}$, where

$$(4.4) \qquad \bar{\bar{v}}_i(j) \equiv \bar{\kappa}^3 \min\left(\frac{\sigma_i}{\sigma_j}, \frac{\sigma_j}{\sigma_i}\right).$$

We use $\bar{\bar{v}}_i(j)$ as defined in (4.4) for each $G_m$, even though the values of $\sigma_i$ and $\sigma_j$ vary slightly from step to step. This error contributes an $O(\varepsilon^2)$ term to the overall bound (which we ignore) but could be incorporated using the bounds of Corollary 4.2.

THEOREM 4.4. *Let $V$, $V_T$, and $\bar{\kappa}$ be as in Theorem 4.3, and $\bar{\bar{v}}_i(j)$ be as in (4.4). Then we can bound the error in the individual components of $v_i$ by*

$$(4.5) \qquad |v_{Ti}(j) - v_i(j)| \leq q(M, n) \cdot \frac{(\text{tol} + \varepsilon) \cdot \bar{\kappa}^2 \cdot \bar{\bar{v}}_i(j)}{\text{relgap}_{\sigma i}},$$

*where $q(M, n)$ has a bound similar to that of $p(M, n)$ in Theorem 3.4.*

**4.3. Using Cholesky followed by one-sided Jacobi for the symmetric positive definite eigenproblem.** In this subsection we consider two algorithms for the symmetric positive definite eigenproblem $H$, both based on performing Cholesky on $H$, and using one-sided Jacobi to compute the SVD of the Cholesky factor $L$. The first algorithm (Algorithm 4.2) does left-handed Jacobi on $L$, returning its left singular vectors as the eigenvectors of $H$ and the squares of its singular values as the eigenvalues of $H$. The second algorithm (Algorithm 4.4), originally proposed in [22], does Cholesky with complete pivoting (which is equivalent to diagonal pivoting) and then right-handed Jacobi on $L$, again returning its left singular vectors and squares of its singular values.

The first algorithm, left-handed Jacobi, is about as accurate as two-sided Jacobi, but permits purely column oriented access to the data following the initial Cholesky decomposition; this can have speed advantages on machines with memory hierarchies. The second algorithm, right-handed Jacobi with pivoting, is less accurate than the first because it will not always compute tiny eigenvector components with the accuracy of Theorem 3.4, although it does compute the eigenvalues as accurately, and the eigenvectors with the same norm error bound. However, it can be several times faster than either the first algorithm or two-sided Jacobi.

ALGORITHM 4.2 (Left-handed Jacobi on $L$ without pivoting for the symmetric positive definite eigenproblem $H$).
    1. Form the Cholesky factor $L$ of $H$: $H = LL^T$.
    2. Compute the singular values $\sigma_i$ and left singular vectors $v_i$ of $L$ using left-handed Jacobi.
    3. The eigenvalues $\lambda_i$ of $H$ are $\lambda_i = \sigma_s^2$. The eigenvectors of $H$ are $v_i$.

We show that this method is as accurate as using two-sided Jacobi directly on $H$. The proof involves a new error analysis of Cholesky decomposition, so we begin by restating Cholesky's algorithm in order to establish notation for our error analysis.

ALGORITHM 4.3 (Cholesky decomposition $H = LL^T$ for an $n \times n$ symmetric positive definite matrix $H$).
    for $i = 1$ to $n$

$$L_{ii} = (H_{ii} - \sum_{k=1}^{i-1} L_{ik}^2)^{1/2}$$
for $j = i + 1$ to $n$
$$L_{ji} = (H_{ji} - \sum_{k=1}^{i-1} L_{jk}L_{ik})/L_{ii}$$
endfor
endfor

LEMMA 4.5 (see [10]). *Let $L$ be the Cholesky factor of $H$ computed using Algorithm 4.3 in finite precision arithmetic with precision $\varepsilon$. Then $LL^T = H + E$ where $|E_{ij}| \leq (n + 5)\varepsilon(H_{ii}H_{jj})^{1/2}$.*

THEOREM 4.6. *Let $L$ be the Cholesky factor of $H = DAD$ computed in floating point arithmetic using Algorithm 4.3. Let $\sigma_i$ and $v_{Li}$ be the exact singular values and right singular vectors of $L^T$, and $\lambda_i$ and $v_{Hi}$ be the eigenvalues and eigenvectors of $H$. Let $\bar{v}_i(j)$ be as in Proposition 2.8. Then*

$$\frac{|\lambda_i - \sigma_i^2|}{\lambda_i} \leq (n^2 + 5n) \cdot \varepsilon \cdot \kappa(A),$$

$$\|v_{Li} - v_{Hi}\|_2 \leq \frac{(n^2 + 5n)(n-1)^{1/2} \cdot \varepsilon \cdot \kappa(A)}{\mathrm{relgap}_{\lambda_i}} + O(\varepsilon^2),$$

$$|v_{Li}(j) - v_{Hi}(j)| \leq \frac{(n^2 + 5n)(2n - 2)^{1/2} \cdot \varepsilon \cdot \kappa(A) \cdot \bar{v}_i(j)}{\min(\mathrm{relgap}_{\lambda_i}, 2^{-1/2})} + O(\varepsilon^2).$$

*Proof.* Plug the bound of Lemma 4.5 into Theorem 2.3, Theorem 2.5, and Proposition 2.9.     □

Theorem 4.6 implies that the errors introduced by Cholesky are as small as those introduced by two-sided Jacobi. Write $H = DAD$ and $L_A = D^{-1}L$. Since $\|A - L_A L_A^T\|_2 \leq (n^2 + 5n)\varepsilon$, $\kappa(A) \approx (\kappa(L_A))^2$ (unless both are very large). Since the columns of $L_A^T$ have nearly unit two-norm, the accuracy of left-handed Jacobi applied to $L$ is governed by $\kappa(L_A)$. Thus Cholesky followed by left-handed Jacobi on $L$ results in a problem whose condition number $\kappa(L_A)$ is approximately the square root of the condition number of the original problem $\kappa(A)$. Corollary 4.2 and Theorems 4.3 and 4.4 guarantee that the computed eigenvalues and eigenvectors are accurate. In exact arithmetic, left-handed Jacobi on $L$ is the same as two-sided Jacobi on $DAD = H = LL^T = D(L_A L_A^T)D$, so the question of how much $\kappa(L_A)$ can grow during subsequent Jacobi rotations is essentially identical to the question of the growth of $\kappa(A_m)$ during two-sided Jacobi. The second algorithm follows.

ALGORITHM 4.4 (Right-handed Jacobi on $L$ with pivoting for the symmetric positive definite eigenproblem $H$).

1. Form the Cholesky factor $L$ of $H$ using complete pivoting. Then there is a permutation matrix $P$ such that $P^T HP = LL^T$.

2. Compute the singular values $\sigma_i$ and left singular vectors $v_i$ of $L$ using right-handed Jacobi.

3. The eigenvalues $\lambda_i$ of $H$ are $\lambda_i = \sigma_s^2$. The eigenvectors of $H$ are $Pv_i$.

Even if we did not do complete pivoting, Theorem 4.6 would guarantee that the squares of the true singular values of $L$ would be accurate eigenvalues of $H$, and that the true left singular vectors of $L$ would be accurate eigenvectors of $P^T HP$. Since we are computing left singular vectors of $L$, Theorem 4.4 does not apply, but from Corollary 4.2, we know that the computed eigenvalues are accurate, and from Theorem 4.3 we know that the computed eigenvectors are accurate in a norm sense.

Numerical experiments in §7 bear out the fact that tiny eigenvector components may not always be computed as accurately by Algorithm 4.4 as by Algorithm 4.2.

Note that Algorithm 4.4 is mathematically equivalent to doing two-sided Jacobi on $L^T L$, so we in effect take a single step of the symmetric LR algorithm [22] before beginning Jacobi, thus giving Jacobi a "head start." (An analogous head start is attained by preceding left-handed Jacobi for the SVD with a QR decomposition with column pivoting [12].) Writing $L^T L = H_1 = D_1 A_1 D_1$, we see it is $\kappa(A_1)$, which governs the accuracy of step 2 of the algorithm, as well as its speed, since it is $\kappa(A_1)$ that must be driven to 1. We discuss this in more detail in §6, where we show that $\kappa(A_1)$ can be much smaller than $\kappa(A)$, where $H = DAD$ is the original problem.

There are two other algorithmic variations one might consider: Cholesky with pivoting followed by left-handed Jacobi on $L$, and Cholesky without pivoting followed by right-handed Jacobi on $L$. We can measure the quality of both algorithms as we did in the last paragraph: We find symmetric positive definite matrices $H_2$ and $H_3$ such that the algorithms are mathematically equivalent to doing two-sided Jacobi on $H_2$ and $H_3$, respectively. Then their accuracies and running times depend on $\kappa(A_2)$ and $\kappa(A_3)$ where $H_i = D_i A_i D_i$. One may show that $\kappa(A_2) = \kappa(A)$, so there is no advantage to pivoting and left-handed Jacobi, and also that $\kappa(A_3)$ can greatly exceed $\kappa(A)$, so that right-handed Jacobi without pivoting can be harmful. We do not consider these algorithmic variations further.

**5. Bisection and inverse iteration.** Here we show that bisection and inverse iteration applied to the symmetric positive definite matrix $H = DAD$ can compute the eigenvalues and eigenvectors within the accuracy bounds section of 2. Let inertia($H$) denote the triple $(n, z, p)$ of the number $n$ of negative eigenvalues of $H$, the number $z$ of zero eigenvalues of $H$, and the number $p$ of positive eigenvalues of $H$. These results are simple extensions of Algorithms 3 and 5 in [2], and detailed proofs may be found there and in [10].

ALGORITHM 5.1 (Stably computing the inertia of $H - xI = DAD - xI$).

1. Permute the rows and columns of $A - xD^{-2}$ (which has the same inertia as $H - xI$) and partition it as

$$\begin{bmatrix} A_{11} - xD_1^{-2} & A_{12} \\ A_{21} & A_{22} - xD_2^{-2} \end{bmatrix}$$

so that if $1 - xd^{-2}$ is a diagonal entry of $A_{11} - xD_1^{-2}$, then $xd^{-2} \geq 2n + 1$, where $n$ is the dimension of $H$.

2. Compute $X = A_{22} - xD_2^{-2} - A_{21}(A_{11} - xD_1^{-2})^{-1}A_{12}$, using Cholesky to compute $(A_{11} - xD_1^{-2})^{-1}A_{12}$.

3. Compute inertia($X$) = (neg, zero, pos) using a stable pivoting scheme such as in [4].

4. The inertia of $H - xI$ is (neg + dim($A_{11}$), zero, pos).

We need to partition $A - xD^{-2}$ as above in order to make the proof convenient, but it may not be necessary algorithmically [10].

THEOREM 5.1. *Let $\varepsilon$ be the machine precision in which Algorithm 5.1 is carried out, where we assume that neither overflow nor underflow occur. Then Algorithm 5.1 computes the exact inertia of $D(A + \delta A)D - xI$, where $\|\delta A\|_2 = O(\varepsilon)$. Thus Algorithm 5.1 can be used in a bisection algorithm to find all the eigenvalues of $H$ to the accuracy of Theorem 2.3 or Proposition 2.4.*

ALGORITHM 5.2 (Inverse iteration for computing the eigenvector $x$ of a symmetric positive definite matrix $H = DAD$ corresponding to eigenvalue $z$). tol is a user-specified stopping criterion.

1. We assume that eigenvalue $z$ has been computed accurately, for example, using Algorithm 5.1.

2. Choose a starting vector $y_0$; set $i = 0$.

3. Compute the symmetric indefinite factorization $LDL^T$ of $P(A - zD^{-2})P^T$ [4], where $P$ is the same permutation as in Algorithm 5.1, step 1.

4. Repeat
   $i = i + 1$
   Solve $(A - zD^{-2})\tilde{y}_i = y_{i-1}$ for $\tilde{y}_i$ using the $LDL^T$ factorization of step 3.
   $r = 1/\|\tilde{y}_i\|_2$
   $y_i = r \cdot \tilde{y}_i$
   until $(r \leq \text{tol})$

5. $x = D^{-1}y_i$

THEOREM 5.2. *Suppose that Algorithm 5.2 terminates with $x$ as the computed eigenvector of $H = DAD$. Then there is a diagonal matrix $\hat{D}$ with $\hat{D}_{ii} = 1 + O(\text{tol})$ and a matrix $\delta A$ with $\|\delta A\|_2 = O(\text{tol})$, such that $\hat{D}x$ is the exact eigenvector of $D(A + \delta A)D$. Thus the error in $x$ is bounded by Theorem 2.5, Corollary 2.6, and Proposition 2.9.*

**6. Upper bounds for** $\max_m \kappa(A_m)/\kappa(A_0)$. As stated in §§3 and 4, our claims about the accuracy to which Jacobi can solve the eigenproblem depend on the ratio $\max_m \kappa(A_m)/\kappa(A_0)$ being modest. Here $H_0 = D_0 A_0 D_0$ is the initial matrix, and $H_m = D_m A_m D_m$ is the sequence produced by Jacobi ($H_{m+1}$ is obtained from $H_m$ by applying a single Jacobi rotation, $D_m$ is diagonal, and $A_m$ has ones on the diagonal). The reason is that the error bounds for Jacobi are proportional to $\max_m \kappa(A_m)$, and the error bounds of §2 are proportional to $\kappa(A_0)$.

In this section, we present several results explaining why $\max_m \kappa(A_m)/\kappa(A_0)$ should not be expected to grow very much. Recall that convergence of $H_m$ to diagonal form is equivalent to the convergence of $A_m$ to the identity matrix, or of $\kappa(A_m)$ to 1. Thus we expect $\kappa(A_m) < \kappa(A_0)$ eventually. The best situation would be monotonic convergence, but this is, unfortunately, not always the case.

We have not been able to completely explain the extremely good numerical results of §7, that $\max_m \kappa(A_m)/\kappa(A_0)$ never exceeded 1.82, and averaged 1.20 in random experiments. (Wang [23] has found a sequence $H_n$ of matrices of dimension $n$ where this ratio grows slowly with $n$, reaching 8 for $n = 50$. Changing the sweep strategy eliminated this growth.) A complete theoretical explanation of this remains an open question.

We only speak in terms of two-sided Jacobi in this section. This is no loss of generality because, in exact arithmetic, right-handed Jacobi on $G$ is equivalent to two-sided Jacobi on $G^T G$.

Our first result shows that $\kappa(A_m)/\kappa(A_0)$ cannot be too large if $A_m$ is obtained from $A_0$ by a sequence of Jacobi rotations in pairwise disjoint rows and columns. The second result gives a cheaply computable guaranteed upper bound on $\max_m \kappa(A_m)/\kappa(A_0)$ in terms of the Hadamard measure of $A_0$. This bound is generally quite pessimistic unless the dimension of $A$ is modest and $\kappa(A_0)$ is small—at most a few hundred. The third and fourth results will be for right-handed Jacobi with pivoting (Algorithm 4.4). The third result shows that the wider the range of numbers

on the diagonal of $H$, the smaller $\kappa(A_1)$ is for that algorithm. This in turn makes it converge faster. The fourth, rather surprising, result is that $\kappa(A_1)$ is bounded by a constant, depending *only* on the dimension $n$, not on $A_0$. These last two results lead us to recommend right-handed Jacobi with pivoting Jacobi as the algorithm of choice (unless it is important to get small eigenvector components to high accuracy; see the discussion in §4.3).

PROPOSITION 6.1. *Let $H_0$ be $n \times n$. Let $H_m$ be obtained from $H_0$ by applying $m$ Jacobi rotations in pairwise nonoverlapping rows and columns (this means $m \leq n/2$). Write $H_m = D_m A_m D_m$ as before. Then*

$$(6.1) \qquad \frac{\kappa(A_m)}{\kappa(A_0)} \leq \frac{1 + \max_{1 \leq k \leq m} |A_{0,2k-1,2k}|}{1 - \max_{1 \leq k \leq m} |A_{0,2k-1,2k}|} \leq \min(\kappa(A_0), 2n).$$

*Also,*

$$(6.2) \qquad \frac{\kappa(A_{i+1})}{\kappa(A_i)} \leq \min(\kappa(A_0), 8).$$

*Furthermore, the spectrum of $A_m$ is independent of $D_0$, even though the entries of $A_m$ depend on $D_0$. More precisely, the spectrum of $A_m$ coincides with the spectrum of the pencil $A_0 - \lambda A_0'$, where $A_0'$ coincides with $A_0$ on every rotated element and is the identity otherwise.*

*Proof.* We begin by deriving a matrix pencil depending only on $A_0$ whose eigenvalues are the same as $A_m$. This proves that the eigenvalues of $A_m$ depend only on $A_0$. We assume without loss of generality that the $m$ Jacobi rotations are in rows and columns $(1,2), (3,4), \cdots, (2m-1, 2m)$. This lets us write $J^T H_0 J = H_m$, where $J$ is block diagonal with the $2 \times 2$ Jacobi rotations (and possibly ones) on its diagonal. Rewrite this as

$$A_m = (D_m^{-1} J^T D_0) A_0 (D_0 J D_m^{-1}) \equiv Z^T A_0 Z,$$

where $Z$ has the same block diagonal structure as $J$. Let $A_0'$ be a block diagonal matrix with the same block structure as $Z$ and $J$, where $A_0'$ is identical to $A_0$ within its $2 \times 2$ blocks, and has ones on its diagonal when $J$ does. Since $H_{m,12} = H_{m,34} = \cdots = 0$, also $A_{m,12} = A_{m,23} = \cdots = 0$. Thus $A_m$ has $2 \times 2$ identity matrices on its diagonal matching the block structure of $Z$, $J$, and $A_0'$. Thus $A_m = Z^T A_0 Z$ implies $Z^{-T} Z^{-1} = A_0'$. Therefore, the eigenvalues of $A_m = Z^T A_0 Z$ are identical to those of the pencil $A_0 - \lambda Z^{-T} Z^{-1} = A_0 - \lambda A_0'$.

Now we apply the minimax theorem to bound $\lambda_{\min}(A_m)$ below by

$$(6.3) \quad \lambda_{\min}(A_m) = \min_{x \neq 0} \frac{x^T A_0 x}{x^T A_0' x} \geq \frac{\min_{\|x\|=1} x^T A_0 x}{\max_{\|x\|=1} x^T A_0' x} = \frac{\lambda_{\min}(A_0)}{1 + \max_{1 \leq k \leq m} |A_{0,2k-1,2k}|}.$$

We may bound $1 + \max_{1 \leq k \leq m} |A_{0,2k-1,2k}|$ from above by both $\lambda_{\max}(A_0)$ and 2, yielding

$$(6.4) \qquad \lambda_{\min}(A_m) \geq \frac{\lambda_{\min}(A_0)}{\min(2, \lambda_{\max}(A_0))}.$$

Now we bound $\lambda_{\max}(A_m)$ from above. First, by the minimax theorem we may write

$$\lambda_{\max}(A_m) = \max_{x \neq 0} \frac{x^T A_0 x}{x^T A_0' x} \leq \frac{\max_{\|x\|=1} x^T A_0 x}{\min_{\|x\|=1} x^T A_0' x} \leq \frac{\lambda_{\max}(A_0)}{1 - \max_{1 \leq k \leq m} |A_{0,2k-1,2k}|},$$

which, when combined with (6.4), yields

$$\kappa(A_m) \leq (\kappa(A_0))^2,$$

proving half of (6.1). For the other half, note that $1 \leq \lambda_{\max}(A_i) \leq n$ for all $i$, so that $\lambda_{\max}(A_m)/\lambda_{\max}(A_0) \leq n$. Now combine this with (6.4).

Now we show $\lambda_{\max}(A_{i+1}) \leq 4\lambda_{\max}(A_i)$, which, when combined with (6.4), yields (6.2). It suffices to show $\lambda_{\max}(A_1) \leq 4\lambda_{\max}(A_0)$. Write

$$A_0 = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

where $A_{11}$ is $2 \times 2$. Then by the minimax theorem, there exists a conformally partitioned unit vector $x^T = [x_1^T, x_2^T]$ where

$$\lambda_{\max}(A_1) = \frac{x_1^T A_{11} x_1 + 2x_1^T A_{12} x_2 + x_2^T A_{22} x_2}{x_1^T A_{11} x_1 + x_2^T x_2}.$$

Write $x_1^T x_1 = \zeta$ $(0 \leq \zeta \leq 1)$, $x_2^T x_2 = 1 - \zeta$, $x_1^T A_{11} x_1 = \tau_1 \zeta$, and $x_2^T A_{22} x_2 = \tau_2(1 - \zeta)$, so that

$$\lambda_{\max}(A_1) = \frac{\tau_1 \zeta + 2x_1^T A_{12} x_2 + \tau_2(1 - \zeta)}{\tau_1 \zeta + 1 - \zeta} \leq 2\frac{\tau_1 \zeta + \tau_2(1 - \zeta)}{\tau_1 \zeta + 1 - \zeta}.$$

The maximum of this last expression over all $0 \leq \zeta \leq 1$ is

$$2 + 2\tau_2 \leq 2 + 2\lambda_{\max}(A_{22}) \leq 4\lambda_{\max}(A_0). \qquad \square$$

Our second bound is based on the Hadamard measure of a symmetric positive definite matrix $H$:

$$\mathcal{H}(H) \equiv \frac{\det(H)}{\prod_i H_{ii}}.$$

PROPOSITION 6.2. *The Hadamard measure $\mathcal{H}(H)$ has the following properties:*
1. *$\mathcal{H}(H) \leq 1$ and $\mathcal{H}(H) = 1$ if and only if $H$ is diagonal.*
2. *$\mathcal{H}(H) = \mathcal{H}(\tilde{D}H\tilde{D})$ for any nonsingular diagonal $\tilde{D}$.*
3. *Let $H = DAD$ with $D$ diagonal and $A$ with unit diagonal. Then*

$$\lambda_{\min}(A) \geq \frac{\mathcal{H}(H)}{e} = \frac{\det(A)}{e},$$

*where $e = \exp(1)$.*
4. *Let $H'$ be obtained from $H$ by applying a Jacobi rotation (in exact arithmetic) in rows and columns $i$ and $j$. Then*

$$\mathcal{H}(H') = \frac{\mathcal{H}(H)}{1 - A_{ij}^2} \geq \mathcal{H}(H).$$

5. *Let $H_0, \cdots, H_m, \cdots$ be a sequence of symmetric positive definite matrices obtained from Jacobi's method in exact arithmetic. Let $H_m = D_m A_m D_m$ with $D_m$ diagonal and $A_m$ with unit diagonal. Then*

$$\max_m \kappa(A_m) \leq \frac{n \cdot e}{\det(A_0)} = \frac{n \cdot e}{\mathcal{H}(H_0)}.$$

*Proof.* 1. Write the Cholesky decomposition $H = LL^T$. Then

$$H_{11} \cdots H_{nn} = \prod_{i=1}^{n} \left( \sum_{k=1}^{i} L_{ik}^2 \right) \geq \prod_{i=1}^{n} L_{ii}^2 = \det(H).$$

2. $\det(\tilde{D}^2)$ factors out of the numerator and denominator of $\mathcal{H}(\tilde{D}H\tilde{D})$.

3. From part 2 above $\mathcal{H}(H) = \mathcal{H}(A) = \det(A)$, so it suffices to show $\lambda_{\min}(A) \geq \det(A)/e$. Let $0 < \lambda_1 \leq \cdots \leq \lambda_n$ be the eigenvalues of $A$. Since $\lambda_1 = \det(A)/\prod_{i=2}^{n} \lambda_i$, we need to show $\prod_{i=2}^{n} \lambda_i \leq e$. Now $\sum_{i=2}^{n} \lambda_i \leq \text{tr}(A) = n$. Since $ab \geq (a+x)(b-x)$ for all $a \geq b \geq x \geq 0$, we see that $\prod_{i=2}^{n} \lambda_i$ is greatest when all $\lambda_i = n/(n-1)$, in which case $\prod_{i=2}^{n} \lambda_i = ((n-1)/n)^{n-1} \leq e$.

4. From Proposition 6.1 we have

$$\mathcal{H}(H') = \det(AA'^{-1}) = \det(A)/(1 - A_{ij}^2) = \mathcal{H}(H)/(1 - A_{ij}^2),$$

where $A' = I$ except for $A'_{ij} = A'_{ji} = A_{ij}$.

5. This is directly implied by parts 3 and 4.  □

Thus part 5 of this proposition gives us a guaranteed upper bound on $\max_m \kappa(A_m)$ at a cost of about $n^3/6$ flops, compared to $2n^3$ flops per Jacobi sweep ($4n^3$ if accumulating eigenvectors). If we use the algorithm in §4.3, where we must do Cholesky anyway, this upper bound comes nearly for free.

Basically, this upper bound is only useful as long as $\kappa(A_0)$ is quite small and $A_0$ has low dimension; otherwise, it is much too large to be useful.

Our third and fourth bounds are for right-handed Jacobi with pivoting (Algorithm 4.4). Recall that this algorithm begins by doing Cholesky with complete pivoting on $H_0$ to get $PH_0P^T = LL^T$, where $P$ is a permutation matrix. Then it does right-handed Jacobi on $L$, which is equivalent (in exact arithmetic) to two-sided Jacobi on $L^T L$. Therefore, Algorithm 4.4 essentially starts with $L^T L = H_1 = D_1 A_1 D_1$.

Our third result, which we state rather informally, is that the larger the range of numbers on the diagonal $D^2$ of $H$, the smaller is $\kappa(A_1)$ (this effect was also observed in [22]). We argue as follows. Let $L = DL_A$ be the factor obtained from complete pivoting. Here, $L_A$ has rows of unit two-norm. Since Algorithm 4.4 does right-handed Jacobi on $L$, its performance depends on the condition number of $DL_AD'$, where $D'$ is chosen diagonal to make the columns of $DL_AD'$ unit vectors. From van der Sluis's theorem [21], we know the condition number of $DL_AD'$ can be at most $n$ times $DL_AD^{-1}$, so it suffices to examine $\kappa(DL_AD^{-1})$. The effect of complete pivoting is essentially to reorder $D$ so that $D_{ii} \geq D_{i+1,i+1}$, and to keep $L_{A,ii}$ as large as possible. Now $(DL_AD^{-1})_{ii} = L_{A,ii}$ is unchanged, and the subdiagonal entry $(DL_AD^{-1})_{ij} = L_{A,ij}D_{ii}D_{jj}^{-1}$ is multiplied by the factor $D_{ii}D_{jj}^{-1}$, which is between 0 and 1. The more $D_{jj}$ exceeds $D_{ii}$, the smaller this factor, and the more nearly diagonal $DL_AD^{-1}$ becomes. Since complete pivoting tries to keep the diagonal of $L_A$ large, this improves the condition number.

Our fourth result shows that, surprisingly, $\max_{m \geq 1} \kappa(A_m)$ is bounded independent of $H_0$.

PROPOSITION 6.3. *Let $PH_0P^T = LL^T$ be the Cholesky decomposition of the $n \times n$ matrix $H_0$ obtained with complete pivoting. Let $H_1 = L^T L = D_1 A_1 D_1$. Let $H_m = D_m A_m D_m$, $m > 1$, be obtained from two-sided Jacobi applied to $H_1$. Then*

1. $\mathcal{H}(H_1) \geq \mathcal{H}(H_0)$.
2. $\mathcal{H}(H_1) \geq 1/n!$. *This bound is attainable.*
3. $\max_{m \geq 1} \kappa(A_m) \leq e \cdot n \; /\mathcal{H}(H_1) \leq e \cdot n \cdot n!$.

*Proof.* 1. Since $\det(H_1) = \det(H_0)$, it suffices to show that $\prod_i H_{0,ii} \geq \prod_i H_{1,ii}$. Assume without loss of generality that $P = I$. Then $H_{0,ii} = \sum_{k=1}^i L_{ik}^2$ and $H_{1,ii} = \sum_{k=i}^n L_{ki}^2$. Complete pivoting is equivalent to the fact that $L_{ii}^2 \geq \sum_{k=i}^j L_{jk}^2$ for all $j > i$. We wish to prove $\prod_{i=1}^n \sum_{k=1}^i L_{ik}^2 \geq \prod_{i=1}^n \sum_{k=i}^n L_{ki}^2$. We systematically use the fact that $ab \geq (a+x)(b-x)$ for $a \geq b \geq x \geq 0$. We illustrate the general procedure in the case of $n = 3$:

$$(L_{11}^2)(L_{21}^2 + L_{22}^2)(L_{31}^2 + L_{32}^2 + L_{33}^2) \geq (L_{11}^2 + L_{21}^2)(L_{22}^2)(L_{31}^2 + L_{32}^2 + L_{33}^2)$$
$$\geq (L_{11}^2 + L_{21}^2 + L_{31}^2)(L_{22}^2)(L_{32}^2 + L_{33}^2)$$
$$\geq (L_{11}^2 + L_{21}^2 + L_{31}^2)(L_{22}^2 + L_{32}^2)(L_{33}^2).$$

2. We have

$$\mathcal{H}(H_1) = \frac{\det(L)^2}{\prod_{i=1}^n (L^T L)_{ii}} = \frac{\prod_{i=1}^n L_{ii}^2}{\prod_{i=1}^n (\sum_{k=i}^n L_{ki}^2)} = \prod_{i=1}^n \frac{L_{ii}^2}{\sum_{k=i}^n L_{ki}^2} \geq \prod_{i=1}^n \frac{1}{i} = \frac{1}{n!}.$$

To see that this bound is attainable, let $H = LL^T$ where $L_{ii} = \mu^{(i-1)/2}$ and $L_{ij} = (1-\mu)^{1/2}\mu^{(i-1)/2}$. Now let $\mu > 0$ become small.

3. The result follows from part 2 and Proposition 6.2, part 5.    □

The example in part 2 of the proposition for which the Hadamard bound is attainable unfortunately has the property that the resulting upper bound in part 3 is a gross overestimate. While the upper bound grows as $e \cdot n \cdot n!$, $\kappa(A_1)$ only grows like $n^{3/2}$. However, $\kappa(A_0)$ grows like $\mu^{-n/2}$, which can be arbitrarily larger than the bound in part 3. The choice $\mu = 0.5$ provides an example in which the upper bound in part 3 can arbitrarily exceed both $\kappa(A_0)$ and $\max_{m \geq 1} \kappa(A_m)$ for large $n$.

Nonetheless, in numerical experiments the upper bound $e \cdot n/\mathcal{H}(H_1)$ on $\max_{m \geq 1} \kappa(A_m)$ never exceeded 40. We also always observed that $\kappa(A_1) \leq \kappa(A_0)$ in all cases, although this is not true in general [23].

Recently, Slapničar [17] has improved the $e \cdot n \cdot n!$ bound to $O(4^n)$ and has shown that this improved bound is attainable; see also related results in [13].

## 7. Numerical experiments.

In this section we present the results of numerical experiments. Briefly, we tested every error bound of every algorithm presented in this paper, and verified that they held in all examples. In fact, the performance is better than we were able to explain theoretically, both because we could observe little or no growth in actual errors for increasing dimension, and because of the surprisingly small values attained by $\max_m \kappa(A_m)/\kappa(A_0)$ (see §6).

These tests were performed using FORTRAN on a SUN 4/260. The arithmetic was IEEE standard double precision [1], with a machine precision of $\varepsilon = 2^{-53} \approx 10^{-16}$ and over/underflow threshold $10^{\pm 308}$.

There were essentially four algorithms tested: two-sided Jacobi (Algorithm 3.1), one-sided Jacobi (Algorithms 4.1 and 4.2), right-handed Jacobi with pivoting (Algorithm 4.4), and bisection/inverse iteration (Algorithms 5.1 and 5.2). All were used with the stopping criterion tol $= 10^{-14}$.

Since we claim that these algorithms are more accurate than any other, we tested their accuracy as follows. We considered only symmetric positive definite eigenproblems, and solved every one using every algorithm. The different answers were compared to see if they agreed to the predicted accuracy (which they did). They were also compared to the EISPACK routines tred2/tql2 [18], which implement tridiagonalization followed by QR iteration. Small eigenvalues computed by EISPACK were often negative, indicating total loss of relative accuracy.

For example, the matrix

$$H = \begin{bmatrix} 10^{40} & 10^{19} & 10^{19} \\ 10^{19} & 10^{20} & 10^{9} \\ 10^{19} & 10^{9} & 1 \end{bmatrix}$$

has all its eigenvalues computed to high relative accuracy by Jacobi, whereas QR computes at least one negative or zero eigenvalue, no matter how the rows and columns are ordered.[1] This shows that QR cannot be made to deliver high relative accuracy on appropriately graded matrices, as suggested in [18].

The remainder of this section is organized as follows: Section 7.1 discusses test matrix generation. Section 7.2 discusses the accuracy of the computed eigenvalues. Section 7.3 discusses the accuracy of the computed eigenvectors. Section 7.4 discusses the the growth of $\max_m \kappa(A_m)/\kappa(A_0)$. Section 7.5 discusses convergence rates; here the speed advantage of right-handed Jacobi with pivoting is apparent.

**7.1. Test matrix generation.** We generated several categories of random test matrices according to three parameters: the dimension $n$, $\kappa_A$, and $\kappa_D$. First, we describe the algorithm used to generate a random matrix from these parameters and then the sets of parameters used.

We tested matrices of dimensions $n = 4$, 8, 16, and 50. Since testing involved solving an $n \times n$ eigenproblem after each Jacobi rotation (to evaluate $\kappa(A_m)$) and there are $O(n^2)$ Jacobi rotations required for convergence, testing costs $O(n^5)$ operations per matrix.

Given $\kappa_A$, we generated a random symmetric positive definite matrix with unit diagonal and approximate condition number $\kappa_A$ as follows. We began by generating a diagonal matrix $T$ with diagonal entries in a geometric series from 1 down to $1/\kappa_A$. Then we generated an orthogonal matrix $U$ uniformly distributed with respect to Haar measure [19], and formed $UTU^T$. Finally, we computed another diagonal matrix $K$ so that $A_0 = KUTU^T K$ had unit diagonal. This last transformation can decrease the condition number of $UTU^T$, but usually not by much. For $4 \times 4$ matrices, it decreased it by as much as a factor of 500, for $8 \times 8$ matrices by a factor of 20, for $16 \times 16$ matrices by a factor of 5, and for $50 \times 50$ matrices by a factor of 1.5. (This decreasing variability is at least partly due to the fact that we ran fewer tests on the larger matrices.) For a more complete discussion of the test matrix generation software, see [9].

Given $\kappa_D$, we generated a random diagonal matrix $D_0$ with diagonal entries whose logarithms were uniformly distributed between 0 and $\log \kappa_D$. This means the diagonal entries themselves were distributed from 1 to $\kappa_D$. The uniform distribution of the logarithm essentially means that every decade is equally likely, and so matrices $D_0$ are generated with entries of widely varying magnitudes.

The resulting random matrix was then $H_0 = D_0 A_0 D_0$.

We generated random matrices with five possible different values of $\kappa_A$: 10, $10^2$, $10^4$, $10^8$, and $10^{12}$; six possible different values of $\kappa_D$: $10^5$, $10^{10}$, $10^{20}$, $10^{30}$, $10^{50}$, and $10^{100}$; and four different dimensions $n = 4$, 8, 16, and 50. This makes a total of $5 \times 6 \times 4 = 120$ different classes of matrices. In each class of dimension $n = 4$ matrices, we generated 100 random matrices, in each class of $n = 8$, we generated 50 random matrices, in each class of $n = 16$, we generated 10 random matrices, and in each class

---

[1] This was using version 3.5h of Matlab on a SUN 4/260. Later versions of Matlab may get different results. For more analysis, see [7] and [15].

of $n = 50$, we generated one random matrix. This makes a total of 4830 different test matrices.

The matrices had, in some cases, eigenvalues ranging over 200 orders of magnitude (when $\kappa_D = 10^{100}$). The relative gaps relgap$_\lambda$ ranged from .028 to $2 \cdot 10^{42}$.

**7.2. Accuracy of the computed eigenvalues.** There are two accuracy bounds for eigenvalues from §2 which we tested. The first one is based on Theorem 2.3 (or Theorem 2.14 together with Theorem 4.6), which says that if $\lambda_i'$ and $\lambda_i''$ are approximations of $\lambda_i$ computed by two of our algorithms, then

$$Q_1 \equiv \frac{|\lambda_i' - \lambda_i''|}{\kappa(A_0)\lambda_i'}$$

should be $O(\text{tol})$, where tol $= 10^{-14}$ is our stopping criterion. For two-sided Jacobi and one-sided Jacobi, $Q_1$ never exceeded $2 \cdot 10^{-15}$. For two-sided Jacobi and right-handed Jacobi with pivoting, $Q_1$ also never exceeded $2 \cdot 10^{-15}$. Every matrix had an eigenvalue for which $Q_1$ exceeded $4 \cdot 10^{-18}$, showing that the bound of Theorem 2.3 is attainable, as predicted by Proposition 2.10.

In the case of bisection, we did not run a bisection algorithm to convergence for each eigenvalue, but rather took the eigenvalues $\lambda_i'$ computed by two-sided Jacobi, made intervals $[(1 - \text{tol} \cdot \kappa(A_0))\lambda_i', (1 + \text{tol} \cdot \kappa(A_0))\lambda_i']$ from each one, and used bisection to verify that each interval contained one eigenvalue (overlapping intervals were merged and the counting modified in the obvious way). All intervals successfully passed this test.

The second accuracy bound is from Proposition 2.4 (or Proposition 2.15 together with Theorem 4.6), which predicts that

$$Q_2 \equiv \frac{|\lambda_i' - \lambda_i''|}{\|D_0 v_i\|_2^2}$$

should be $O(\text{tol})$. Here $v_i$ is the unit eigenvector computed by two-sided Jacobi. For two-sided Jacobi and one-sided Jacobi, $Q_2$ never exceeded $2 \cdot 10^{-14}$. For two-sided Jacobi and right-handed Jacobi with pivoting, $Q_2$ never exceeded $9 \cdot 10^{-15}$. Every matrix had an eigenvalue for which $Q_2$ exceeded $5 \cdot 10^{-16}$, showing that the bound of Proposition 2.4 is attainable, as it predicts.

In the case of bisection, we again made intervals $[\lambda_i' - \text{tol} \cdot \|D_0 v_i\|_2^2, \lambda_i' + \text{tol} \cdot \|D_0 v_i\|_2^2]$ from each eigenvalue $\lambda_i'$ and verified that each interval contained the proper number of eigenvalues.

Finally, we verified a slightly weakened version of Proposition 2.7, that

$$\lambda_{\min}(A_0) - \text{tol} \leq \frac{\lambda_i'}{h_i} \leq \lambda_{\max}(A_0) + \text{tol}$$

for the eigenvalues $\lambda_i'$ computed by two-sided Jacobi. Here $h_i$ is the $i$th smallest diagonal entry of $H_0$. Adding and subtracting tol to the upper and lower bounds takes into account the errors in computing $\lambda_i'$.

**7.3. Accuracy of the computed eigenvectors.** There is one bound on the magnitude of the components of the eigenvectors, and two accuracy bounds, one for the norm error and one for the componentwise error.

We begin with a few details about our implementation of inverse iteration. We used the eigenvalues computed by two-sided Jacobi, and the vector of all ones as a starting vector. Convergence always occurred after just one iteration.

The componentwise bound on the magnitude of the eigenvectors is based on Proposition 2.8, which says that the components of the normalized eigenvector $v_i$ should be bounded by

$$|v_i(j)| \le \bar{v}_i(j) \equiv (\kappa(A_0))^{3/2} \cdot \min\left(\left(\frac{\lambda_i}{\lambda_j}\right)^{1/2}, \left(\frac{\lambda_j}{\lambda_i}\right)^{1/2}\right).$$

This was verified for the eigenvectors computed by all four algorithms. We note that since this bound is proportional to $\kappa(A_0)^{3/2}$, it becomes weaker as $\kappa(A_0)$ becomes larger, and indeed becomes vacuous for matrices with $\kappa(A_0)$ large and eigenvalues in a narrow range.

The norm error bounds are based on Theorem 2.5 (or Theorem 2.16 together with Theorem 4.6), which predicts that if $v_i'$ and $v_i''$ are approximations of the unit eigenvector $v_i$ computed by two of our algorithms, then

$$Q_3 \equiv \frac{\|v_i' - v_i''\|_2}{(\kappa(A_0)/\mathrm{relgap}_{\lambda_i}) + 1}$$

should be $O(\mathrm{tol})$. (We add the 1 in the denominator because a single roundoff error in the largest entry can cause a norm error of $\varepsilon$; see Theorem 3.3 or Theorem 4.3.)

For two-sided Jacobi and one-sided Jacobi, $Q_3$ never exceeded $3 \cdot 10^{-16}$. For two-sided Jacobi and right-handed Jacobi with pivoting, $Q_3$ also never exceeded $2 \cdot 10^{-14}$. For two-sided Jacobi and inverse iteration, $Q_3$ never exceeded $8 \cdot 10^{-14}$. Every matrix had an eigenvector for which $Q_3$ exceeded $10^{-18}$ for every pair of algorithms compared, showing that the bound of Theorem 2.5 is nearly attainable, as predicted by Proposition 2.11.

The second accuracy bound is based on Proposition 2.9 (or Proposition 2.20 and Theorem 4.6), which predicts that

$$Q_4 \equiv \frac{|v_i'(j) - v_i''(j)| \min(\mathrm{relgap}_{\lambda_i}, 2^{-1/2})}{\kappa(A_0) \cdot \bar{v}_i(j)}$$

should be $O(\mathrm{tol})$. For two-sided Jacobi and one-sided Jacobi, $Q_4$ never exceeded $3 \cdot 10^{-17}$. For two-sided Jacobi and inverse iteration, $Q_4$ never exceeded $3 \cdot 10^{-15}$. For two-sided Jacobi and right-handed Jacobi with pivoting, $Q_4$ was as large as .02, which is consistent with the fact that right-handed Jacobi with pivoting computes the eigenvectors as left singular vectors of $L$, for which we only have a normwise error bound (Theorem 4.3). For the other algorithm, $Q_4$ was only $10^{-30}$ for matrices with $\kappa(A_0) = 10^{12}$; this reflects the factor $\kappa(A_0)^{5/2}$ in the denominator of $Q_4$, a weakness of Proposition 2.8. In other words, the componentwise error bounds are generally only interesting for small to medium $\kappa(A_0)$.

**7.4. Growth of** $\max_m \kappa(A_m)/\kappa(A_0)$. In computing

$$Q_5 \equiv \max_m \kappa(A_m)/\kappa(A_0),$$

we note that a single computation requiring $M$ Jacobi rotations supplied us not just with one value of $Q_5$, but rather $M - 1$: Since every $A_i$ can be thought of as starting a new eigenvalue computation, we may also measure $\max_{m \ge i} \kappa(A_m)/\kappa(A_i)$ for all $i < M$. Thus, all told, our 4830 different matrices represent over 900,000 data points of $Q_5$.

TABLE 1
*Hadamard upper bound $Q_6$ on $\max_m \kappa(A_m)/\kappa(A_0)$.*

| $n$ | $\kappa_A$ | | | | |
|---|---|---|---|---|---|
| | 10 | $10^2$ | $10^4$ | $10^8$ | $10^{12}$ |
| 4 | 5.8 | 13 | 590 | $6.3 \cdot 10^6$ | $6.1 \cdot 10^{10}$ |
| 8 | 21 | 410 | $1.1 \cdot 10^7$ | $9.1 \cdot 10^{17}$ | $\infty$ |
| 16 | 200 | $2.7 \cdot 10^5$ | $1.8 \cdot 10^{15}$ | $\infty$ | $\infty$ |
| 50 | $6.4 \cdot 10^5$ | $8.0 \cdot 10^{16}$ | $\infty$ | $\infty$ | $\infty$ |

The largest value of $Q_5$ encountered was 1.82. This was for an $8 \times 8$ matrix with $\kappa(A_0) = 1.4 \cdot 10^{12}$, and eigenvalues ranging over 133 orders of magnitude. 141 Jacobi rotations (a little over 5 sweeps) were required for convergence, plus 28 more steps (one more sweep) where no work is done to recognize convergence. In Fig. 1, a plot is shown of $\kappa(A_i) - 1$ versus $i$. We plot $\kappa(A_i) - 1$ instead of $\kappa(A_i)$ in order to see the quadratic convergence of $\kappa(A_i)$ to 1. The graph appears nearly monotonic, except for a slight rise near $i = 20$. This is seen more clearly in Fig. 2, which plots $\max_{m \geq i} \kappa(A_m)/\kappa(A_i)$ versus $i$. Here the maximal nonmonotonicity of the curve near $i = 20$ is apparent.

Recently, Wang [23] found a family of examples where $Q_5$ was as large as 8 for matrices up to dimension 50. These matrices have 1 on the diagonal and $1 - \epsilon$ on the offdiagonal, where $\epsilon$ is small. However, by using a different pivoting strategy than cyclic-by-rows, namely, the parallel pivoting discussed in Proposition 6.1, this growth could be eliminated.

Now we consider the Hadamard-based upper bound on $Q_5$ from Proposition 6.2:

$$Q_5 \leq Q_6 \equiv \frac{e \cdot n}{\mathcal{H}(H_0) \cdot \kappa(A_0)}.$$

Table 1 gives the maximum values of this upper bound for different values of dimension $n$ and $\kappa_A \approx \kappa(A_0)$. Recall that the true value of $Q_5$ never exceeds 1.82. As Proposition 6.2 suggests, this upper bound should not depend on $D_0$ and indeed the values observed depended very little on $D_0$.

As can be seen, the Hadamard-based bound is of little use except for very small matrices of modest $\kappa(A_0)$. $\infty$ means the value overflowed.
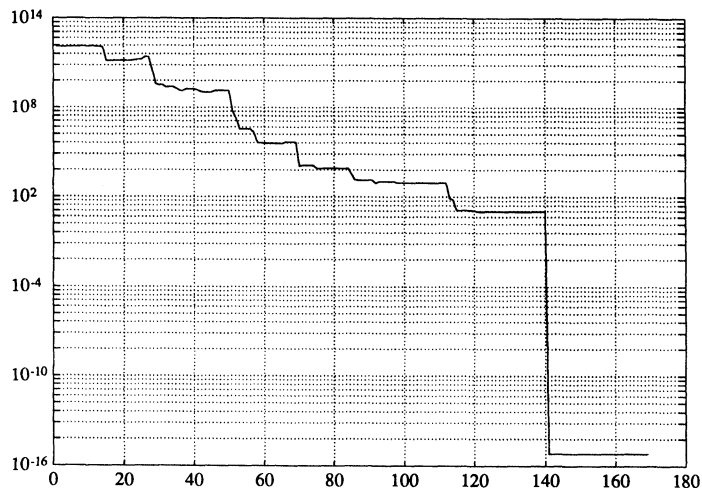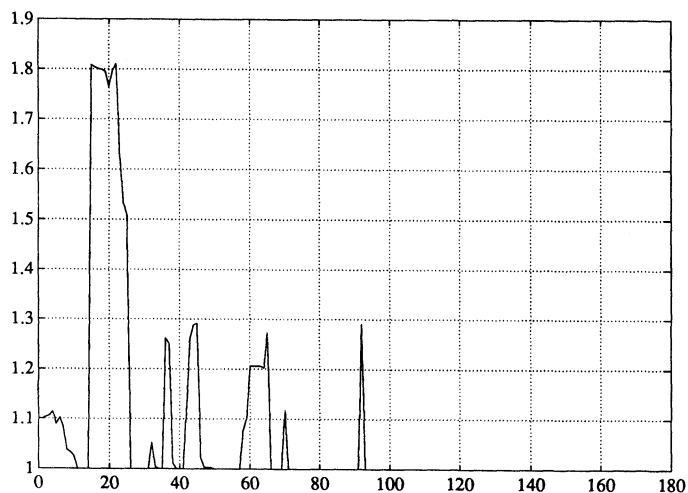
Now we consider right-handed Jacobi with pivoting. Let us recall the notation of §6: Let $PH_0P^T = LL^T$ be Cholesky with complete pivoting, and let $L^T L = H_1 = D_1 A_1 D_1$. As suggested in that section, we expect both $\kappa(A_1)$ to be smaller than $\kappa(A_0)$, and the Hadamard-based upper bound

$$Q_5 \leq Q_7 \equiv \max\left(1, \frac{e \cdot n}{\mathcal{H}(H_1) \cdot \kappa(A_0)}\right)$$

on $Q_5$ to be much smaller than the one for two-sided Jacobi.

First, $\kappa(A_1)/\kappa(A_0)$ never exceeded $\frac{6}{10}$. In fact, $\kappa(A_1)$ *never exceeded* 40 *for any matrix.* This is quite remarkable. This means that all essential rounding errors occurred during the initial Cholesky decomposition. Finally, the Hadamard-based upper bound $Q_7$ on $Q_5$ never exceeded 29. (Recently, Wang [23] found an example where $\kappa(A_1)/\kappa(A_0)$ slightly exceeded 1; in his example, $\kappa(A_0)$ was close to 1.)

**7.5. Convergence rates.** We begin with a few details on how we counted the number of Jacobi rotations required for convergence. In all algorithms (two-sided Jacobi, one-sided Jacobi, and right-handed Jacobi with pivoting), we stopped when the last $n(n-1)/2$ stopping tests $|H_{ij}| \cdot (H_{ii}H_{jj})^{-1/2} \leq \text{tol}$ succeeded; this means every off-diagonal entry of $H$ satisfies the stopping criterion. In the case of two-sided

FIG. 1. $\kappa(A_i) - 1$ *versus* $i$.



FIG. 2. $\max_{m \geq i} \kappa(A_m)/\kappa(A_i)$ *versus* $i$.

Jacobi, this means the last $n(n-1)/2$ Jacobi rotations involved almost no work. For the two one-sided Jacobis, however, evaluating the stopping criterion costs three inner products, so the last $n(n-1)/2$ rotations involve a significant amount of work, even if no rotations are performed. This must be kept in mind when comparing the number of rotations for two-sided and one-sided Jacobi.

We used the same standard cyclic pivot sequence for all the algorithms: $(1,2)$, $(1,3)$, $\cdots$, $(1,n)$, $(2,3)$, $\cdots$, $(2,n)$, $(3,4)$, $\cdots$, $(n-1,n)$.

We begin by comparing two-sided Jacobi and one-sided Jacobi. In exact arithmetic, these two algorithms are identical. In practice, they usually took the same number of steps, although one-sided Jacobi did vary from 20 percent faster to 50 per-

TABLE 2
*Average number of sweeps for two-sided Jacobi (TsJ) and right-handed Jacobi with pivoting (RhJwP).*

| $\kappa_A$ | $\kappa_D$ | Dimension $n$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 4 | | 8 | | 16 | | 50 | |
| | | TsJ | RhJwP | TsJ | RhJwP | TsJ | RhJwP | TsJ | RhJwP |
| 10 | $10^5$ | 3.7 | 3.0 | 4.9 | 3.7 | 5.7 | 4.4 | 6.4 | 5.0 |
| | $10^{10}$ | 3.5 | 2.5 | 4.6 | 3.3 | 5.6 | 4.1 | 6.4 | 5.0 |
| | $10^{20}$ | 3.1 | 2.2 | 4.5 | 2.8 | 5.5 | 3.6 | 6.0 | 4.0 |
| | $10^{30}$ | 3.0 | 2.1 | 4.6 | 2.5 | 5.5 | 3.4 | 6.3 | 4.0 |
| | $10^{50}$ | 2.8 | 1.9 | 4.4 | 2.3 | 5.5 | 3.1 | 5.8 | 4.0 |
| | $10^{100}$ | 2.7 | 1.7 | 4.5 | 2.0 | 5.6 | 2.6 | 5.8 | 3.0 |
| $10^2$ | $10^5$ | 3.8 | 3.0 | 5.2 | 3.8 | 6.4 | 4.5 | 7.5 | 6.0 |
| | $10^{10}$ | 3.5 | 2.5 | 5.1 | 3.3 | 6.2 | 4.1 | 7.4 | 5.0 |
| | $10^{20}$ | 3.2 | 2.2 | 4.9 | 2.9 | 6.2 | 3.9 | 7.1 | 4.0 |
| | $10^{30}$ | 3.0 | 2.0 | 4.8 | 2.6 | 5.8 | 3.3 | 6.8 | 4.1 |
| | $10^{50}$ | 2.9 | 1.9 | 4.8 | 2.2 | 6.1 | 3.0 | 6.5 | 4.0 |
| | $10^{100}$ | 2.8 | 1.6 | 4.7 | 2.0 | 6.0 | 2.7 | 6.8 | 3.4 |
| $10^4$ | $10^5$ | 4.0 | 2.9 | 5.8 | 3.6 | 7.5 | 4.5 | 9.2 | 6.0 |
| | $10^{10}$ | 3.7 | 2.5 | 5.6 | 3.3 | 7.2 | 4.1 | 9.3 | 5.0 |
| | $10^{20}$ | 3.2 | 2.2 | 5.3 | 2.9 | 7.2 | 3.7 | 8.5 | 4.9 |
| | $10^{30}$ | 3.1 | 2.1 | 5.2 | 2.6 | 6.8 | 3.1 | 8.2 | 4.0 |
| | $10^{50}$ | 2.9 | 1.9 | 5.2 | 2.4 | 6.6 | 3.0 | 8.5 | 4.6 |
| | $10^{100}$ | 2.7 | 1.7 | 4.9 | 2.2 | 6.9 | 2.4 | 8.0 | 3.9 |
| $10^8$ | $10^5$ | 3.9 | 2.7 | 6.4 | 3.5 | 9.7 | 4.1 | 13.5 | 6.0 |
| | $10^{10}$ | 3.6 | 2.3 | 6.3 | 3.2 | 9.4 | 3.8 | 12.4 | 5.0 |
| | $10^{20}$ | 3.3 | 2.1 | 5.7 | 2.8 | 8.9 | 3.5 | 11.7 | 4.7 |
| | $10^{30}$ | 3.1 | 2.1 | 5.5 | 2.6 | 8.6 | 3.4 | 12.0 | 4.0 |
| | $10^{50}$ | 2.9 | 1.9 | 5.3 | 2.3 | 8.5 | 3.1 | 11.6 | 4.0 |
| | $10^{100}$ | 2.9 | 1.7 | 5.1 | 2.0 | 8.7 | 2.6 | 11.6 | 4.0 |
| $10^{12}$ | $10^5$ | 3.8 | 2.5 | 6.8 | 3.1 | 10.6 | 4.0 | 16.5 | 6.0 |
| | $10^{10}$ | 3.6 | 2.2 | 6.4 | 3.0 | 10.3 | 3.9 | 15.6 | 5.0 |
| | $10^{20}$ | 3.4 | 2.1 | 6.0 | 2.7 | 9.8 | 3.5 | 15.3 | 5.0 |
| | $10^{30}$ | 3.1 | 2.0 | 5.8 | 2.5 | 10.2 | 3.3 | 15.2 | 4.0 |
| | $10^{50}$ | 2.9 | 1.9 | 5.6 | 2.3 | 9.3 | 3.2 | 13.7 | 3.9 |
| | $10^{100}$ | 2.8 | 1.6 | 5.2 | 2.0 | 8.7 | 2.7 | 15.2 | 3.0 |

cent slower than two-sided Jacobi on some examples. Hereafter, we will only compare two-sided Jacobi to right-handed Jacobi with pivoting.

The most interesting phenomenon was the speedup experienced by right-handed Jacobi with pivoting with respect to two-sided Jacobi. In Table 2 we present the raw data on the number of sweeps required for convergence.

There are a number of interesting trends exhibited in this table. First, RhJwP (right-handed Jacobi with pivoting) never takes more than six sweeps to converge for any matrix, whereas TsJ (two-sided Jacobi) takes up to 16.5. In fact, RhJwP is almost always faster than TsJ (in one example it took 5 percent longer), and can be up to five times faster (3.0 sweeps versus 15.2 sweeps for $\kappa_A = 10^{12}$, $\kappa_D = 10^{100}$, and $n = 50$). Second, the number of sweeps increases with increasing $\kappa_A$ for TsJ, but not for RhJwP. Third, the number of sweeps increases with increasing dimension for both TsJ and RhJwP, but much more modestly for RhJwP (from two to three up to six) than for TsJ (from three to four up to fifteen). Thus the running time for RhJwP is much less dependent on the problem size or sensitivity (as measured by $\kappa_A$) than TsJ. Fourth, the number of sweeps decreases as $\kappa_D$ increases, both for TsJ and RhJwP,

but much more markedly for RhJwP (up to a factor of 2) than for TsJ (usually just one sweep).

**8. Conclusions.** In this paper we have developed new perturbation theory for the eigenvalues and eigenvectors of symmetric positive definite matrices, as well as for eigenvalues of symmetric positive definite pencils. This theory assumes that the perturbations are scaled analogously to the way the matrix is scaled, letting us derive much tighter bounds than in the classical theory. In particular, we get relative error bounds for the eigenvalues and individual components of the eigenvectors, which are (nearly) attainable. The bound for symmetric positive definite pencils may be applied to matrices arising in finite element modeling.

Second, we have shown both through formal error analysis and numerical experiment that Jacobi's method (with a proper stopping criterion) computes the eigenvalues and eigenvectors with these error bounds. We also show that bisection and inverse iteration (applied to the original matrix) attain these bounds. In contrast, methods based on tridiagonalization (such as QR, divide and conquer, traditional bisection, etc.) fail to attain these bounds. In particular, QR can fail to attain these bounds whether or not preceded by tridiagonalization.

We have similar perturbation theorems for the singular value decomposition of a general matrix and the generalized singular values of a pair of matrices, and similar error analyses and numerical experiments for one-sided Jacobi applied to this problem. We may also use one-sided Jacobi to solve the symmetric positive definite eigenproblem.

We have discussed an accelerated version of Jacobi for the symmetric positive definite eigenproblem, which has the property that the more its accuracy exceeds that of QR (or other conventional algorithms), the faster it converges. However, it cannot compute tiny components of eigenvectors as accurately as the other versions of Jacobi, although it computes the eigenvectors with the same norm error bounds. Unless getting the tiny eigenvector components is important, we recommend this accelerated version of Jacobi for the symmetric positive definite eigenproblem.

The quantity $\max_m \kappa(A_m)/\kappa(A_0)$ was seen to be central in the analysis of Jacobi's accuracy. Numerical experiments show it to be much smaller in practice than we can explain. For the accelerated version of Jacobi we provide an inexpensive estimator of $\max_m \kappa(A_m)/\kappa(A_0)$, which works very well in practice. Explaining the excellent behavior of $\max_m \kappa(A_m)/\kappa(A_0)$ is an important open problem.

The error analyses of Jacobi dealt only with the simplest implementations. It would be worthwhile to extend these analyses to cover various enhancements introduced by Veselić, Hari, Rutishauser, and others. These include delayed updates of the diagonal entries and an alternate formula for updating the off-diagonal entries [16], [22], as well as block Jacobi methods.

In future work, we plan to extend these results to the symmetric positive definite generalized eigenproblem, as well as indefinite matrices. Any extension requires an appropriate perturbation theory; therefore, we do not expect to be able to extend the result to all indefinite matrices, since there is no guaranteed way to compute the zero eigenvalues of a singular matrix to "high relative accuracy" without computing them exactly, a feat requiring high precision arithmetic. A class of indefinite matrices for which a suitable perturbation theory exists are the scaled diagonally dominant matrices [2]. The perturbation theory also already exists (at least for eigenvalues) for the symmetric positive definite generalized eigenproblem.

REFERENCES

[1] ANSI/IEEE, IEEE *Standard for Binary Floating Point Arithmetic*, IEEE Press, New York, Std 754-1985 ed., 1985.

[2] J. BARLOW AND J. DEMMEL, *Computing accurate eigensystems of scaled diagonally dominant matrices*, SIAM J. Numer. Anal., 27 (1990), pp. 762–791.

[3] M. BERRY AND A. SAMEH, *Parallel algorithms for the singular value and dense symmetric eigenvalues problems*, J. Comput. Appl. Math., 27 (1989), pp. 191–213.

[4] J. BUNCH AND L. KAUFMAN, *Some stable methods for calculating inertia and solving symmetric linear systems*, Math. Comp., 31 (1977), pp. 163–179.

[5] P. P. M. DE RIJK, *A one-sided Jacobi algorithm for computing the singular value decomposition on a vector computer*, SIAM J. Sci. Statist. Comput., 10 (1989), pp. 359–371.

[6] J. DEMMEL, *The condition number of equivalence transformations that block diagonalize matrix pencils*, SIAM J. Numer. Anal., 20 (1983), pp. 599–610.

[7] ———, *On the inherent inaccuracy of implicit tridiagonal QR*, IMA Preprint 963, University of Minnesota, Minneapolis, MN, April 1992.

[8] J. DEMMEL AND W. KAHAN, *Accurate singular values of bidiagonal matrices*, SIAM J. Sci. Statist. Comput., 11 (1990), pp. 873–912.

[9] J. DEMMEL AND A. McKENNEY, *A test matrix generation suite*, Computer Science Department Tech. Report, Courant Institute, New York, July 1989 (LAPACK Working Note #9).

[10] J. DEMMEL AND K. VESELIĆ, *Jacobi's method is more accurate than QR*, Computer Science Department Tech. Report 468, Courant Institute, New York, October 1989 (LAPACK Working Note #15).

[11] G. GOLUB AND C. VAN LOAN, *Matrix Computations*, Second Edition, The Johns Hopkins University Press, Baltimore, MD, 1989.

[12] V. HARI AND K. VESELIĆ, *On Jacobi methods for singular value decompositions*, SIAM J. Sci. Statist. Comput., 8 (1987), pp. 741–754.

[13] N. J. HIGHAM, *Analysis of the Cholesky decomposition of a semi-definite matrix*, in Reliable Numerical Computation, M. G. Cox and S. Hammarling, eds., Clarendon Press, Oxford, 1990, Chap. 9, pp. 161–186.

[14] T. KATO, *Perturbation Theory for Linear Operators*, Second Edition, Springer-Verlag, Berlin, 1980.

[15] J. LE AND B. PARLETT, *On the sensitivity of the tridiagonal QR transformation*, SIAM J. Matrix Anal. Appl., 14 (1993), to appear.

[16] B. PARLETT, *The Symmetric Eigenvalue Problem*, Prentice–Hall, Englewood Cliffs, NJ, 1980.

[17] I. SLAPNIČAR, *Upper bound for the condition of the scaled matrix of the symmetric eigenvalue problem*, Lehrgebiet Mathematische Physik, Fern-Universität Hagen, Hagen, Germany, November 1990.

[18] B. T. SMITH, J. M. BOYLE, J. J. DONGARRA, B. S. GARBOW, Y. IKEBE, V. C. KLEMA, AND C. B. MOLER, *Matrix Eigensystem Routines—EISPACK Guide*, Lecture Notes in Computer Science 6, Springer-Verlag, Berlin, 1976.

[19] G. W. STEWART, *On efficient generation of random matrices with an application to condition estimation*, SIAM J. Numer. Anal., 17 (1980), pp. 403–409.

[20] ———, *A method for computing the generalized singular value decomposition*, in Matrix Pencils, B. Kågström and A. Ruhe, eds., Springer-Verlag, Berlin, 1983, pp. 207–220. Lecture Notes in Mathematics 973, Proceedings, Pite Havsbad, 1982.

[21] A. VAN DER SLUIS, *Condition numbers and equilibration of matrices*, Numer. Math., 14 (1969), pp. 14–23.

[22] K. VESELIĆ AND V. HARI, *A note on a one-sided Jacobi algorithm*, Numer. Math., 56 (1990), pp. 627–633.

[23] X. WANG, private communication, 1990.