Small singular values can increase in lower precision

Christos Boutsikas Purdue University cboutsik@purdue.edu Petros Drineas Purdue University pdrineas@purdue.edu

Ilse C.F. Ipsen
North Carolina State University
ipsen@ncsu.edu

Abstract

We perturb a real matrix A of full column rank, and derive lower bounds for the smallest singular values of the perturbed matrix, in terms of normwise absolute perturbations. Our bounds, which extend existing lower-order expressions, demonstrate the potential increase in the smallest singular values, and represent a qualitative model for the increase in the small singular values after a matrix has been downcast to a lower arithmetic precision. Numerical experiments confirm the qualitative validity of this model and its ability to predict singular values changes in the presence of decreased arithmetic precision.

1 Introduction

Given a real, full column-rank matrix A, we present lower bounds for the smallest singular values of a perturbed matrix A + E.

1.1 Motivation

We investigate the change in the computed singular values of a tall and skinny matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with m > n and rank $(\mathbf{A}) = n$ when \mathbf{A} is demoted, that is, downcast to a lower arithmetic precision.

We have observed that demotion to lower precision can improve the conditioning of A by significantly increasing the computed small singular values, while leaving large singular values mostly unharmed. For instance, if the smallest singular value of A is on the order of double precision roundoff,

$$\sigma_{\min} \left(\text{double}(\boldsymbol{A}) \right) \approx 10^{-16},$$

then downcasting A to single precision can increase the smallest singular value to single precision roundoff,

$$\sigma_{\min} \left(\text{double}(\text{single}(\boldsymbol{A})) \right) \approx 10^{-8}.$$

This phenomenon has been observed before, as the following quotes illustrate:

- ... small singular values tend to increase [SS90, page 266]
- ... even an approximate inverse of an arbitrarily ill-conditioned

matrix does, in general, contain useful information [Rum09, page 260]

This is due to a kind of regularization by rounding to working precision [Rum09, page 261]

1.2 Modelling demotion to lower precision in terms of perturbations

We model the downcasting of a matrix to lower precision in terms of normwise absolute perturbations.

The accumulated error from typical singular value algorithms in Matlab and Julia is a normwise absolute error [GV13, section 8.6]. Thus, if we downcast a matrix $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ to single precision and compute the singular values of the demoted matrix, the resulting error can be represented as an absolute

perturbation E. According to Weyl's inequality [GV13, Corollary 8.6.2], corresponding singular values of **A** change by at most $||E||_2$,

$$|\sigma_j(\boldsymbol{A} + \boldsymbol{E}) - \sigma_j(\boldsymbol{A})| \le ||\boldsymbol{E}||_2 \approx 10^{-8}, \qquad 1 \le j \le n.$$

The bound implies that singular values larger than single precision roundoff, i.e. $\sigma_i(A) \gg ||E||_2$, remain essentially the same,

$$\sigma_j(\boldsymbol{A}) \approx \sigma_j(\boldsymbol{A}) - \underbrace{\|\boldsymbol{E}\|_2}_{10^{-8}} \leq \sigma_j(\boldsymbol{A} + \boldsymbol{E}) \leq \sigma_j(\boldsymbol{A}) + \underbrace{\|\boldsymbol{E}\|_2}_{10^{-8}} \approx \sigma_j(\boldsymbol{A}),$$

while it is inconclusive about small singular values on the order of double precision roundoff,

$$\underbrace{\sigma_{\ell}(A)}_{10^{-16}} - \underbrace{\|E\|_{2}}_{10^{-8}} \le \sigma_{\ell}(A + E) \le \underbrace{\sigma_{\ell}(A)}_{10^{-16}} + \underbrace{\|E\|_{2}}_{10^{-8}}.$$

1.3 Our contributions

Our main results are normwise absolute lower bounds on the smallest singular value cluster of a perturbed matrix. The bounds, compactly summarized in Theorem 1 below, testify to a definitive increase in the perturbed small singular values. The qualitative validity of the bounds is confirmed by the numerical experiments in section 4. Our assumptions are not restrictive and merely require the smallest singular value cluster to be separated by a small gap from the remaining singular values.

Theorem 1 improves the second-order perturbation expansions in [Ste84; SS90; Ste06], because it is a true lower bound and, unlike [Ste84], [SS90, Lemma V.4.5], it needs no assumptions on the size of the offdiagonal blocks. Its proof (Sections 2 and 3) follows from Theorem 7 and a restatement of Assumptions 3.1. The special case where r = 1 reduces to Assumptions 2.1 and Theorem 4.

Theorem 1. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \ge n$ have $\operatorname{rank}(\mathbf{A}) \ge n - r$ for some $r \ge 1$. Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ be a full singular value decomposition, where $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ is diagonal, and $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$ are orthogonal matrices. Partition commensurately,

$$oldsymbol{\Sigma} = egin{bmatrix} oldsymbol{\Sigma}_1 & oldsymbol{0} \ oldsymbol{0} & oldsymbol{\Sigma}_2 \ oldsymbol{0} & oldsymbol{0} \end{bmatrix}, \qquad oldsymbol{U}^T oldsymbol{E} oldsymbol{V} = egin{bmatrix} oldsymbol{E}_{11} & oldsymbol{E}_{12} \ oldsymbol{E}_{21} & oldsymbol{E}_{22} \ oldsymbol{E}_{31} & oldsymbol{E}_{32} \end{bmatrix},$$

where $\Sigma_1, \boldsymbol{E}_{11} \in \mathbb{R}^{(n-r)\times(n-r)}$ with Σ_1 nonsingular diagonal; and $\Sigma_2, \boldsymbol{E}_{22} \in \mathbb{R}^{r\times r}$ with Σ_2 diagonal. If $1/\|\Sigma_1^{-1}\|_2 > 4\|\boldsymbol{E}\|_2$ and $\|\Sigma_2\|_2 < \|\boldsymbol{E}\|_2$, then²

$$\sigma_{n-r+j}(\boldsymbol{A}+\boldsymbol{E})^2 \geq \lambda_j(\boldsymbol{E}_{32}^T\boldsymbol{E}_{32} + (\boldsymbol{\Sigma}_2 + \boldsymbol{E}_{22})^T(\boldsymbol{\Sigma}_2 + \boldsymbol{E}_{22}) - \boldsymbol{R}_3) - r_4, \quad 1 \leq j \leq r,$$

where \mathbf{R}_3 contains terms of order 3,

$$oldsymbol{R}_3 \equiv oldsymbol{E}_{12}^T oldsymbol{W} + oldsymbol{W}^T oldsymbol{E}_{12}, \qquad oldsymbol{W} \equiv (oldsymbol{\Sigma}_1 + oldsymbol{E}_{11})^{-T} egin{bmatrix} oldsymbol{E}_{21} & oldsymbol{E}_{31} \end{bmatrix} egin{bmatrix} oldsymbol{\Sigma}_2 + oldsymbol{E}_{22} \ oldsymbol{E}_{32} \end{bmatrix}$$

and r_4 contains terms of order 4 and higher,

$$r_4 \equiv \|\boldsymbol{W}\|_2^2 + 4 \frac{\|\boldsymbol{E}\|_2^2 \|(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1} (\boldsymbol{E}_{12} + \boldsymbol{W})\|_2^2}{1 - 4\|\boldsymbol{E}\|_2^2 \|(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1}\|_2^2}.$$

Future work, sketched in section 5, will refine the above results towards a quantitative analysis that predicts the order of magnitude of the increase, and the influential matrix properties, in particular, the role of the singular value gap. We also note that our theorem does not directly account for computational precision. However, our experiments in Section 4 demonstrate that when performing calculations in double precision, the "exact" singular values exhibit complete overlap with those computed in double precision. This suggests that perturbations arising from computational precision should not affect our results.

The singular values of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ are labelled in non-increasing order, by $\sigma_1(\mathbf{A}) \geq \sigma_2(\mathbf{A}) \geq \ldots \geq \sigma_{\min\{m,n\}}(\mathbf{A})$. The eigenvalues of a symmetric matrix $\mathbf{H} \in \mathbb{R}^{k \times k}$ are labelled in non-increasing order, by $\lambda_1(\mathbf{H}) \geq \cdots \geq \lambda_k(\mathbf{H})$.

1.4 Existing work

There are many bounds for the smallest singular value of general, unstructured matrices. The bounds for nonsingular matrices $\mathbf{A} \in \mathbb{C}^{n \times n}$ in [LX21; Shu22; Zou12] involve the factor $|\det(\mathbf{A})|^2 \left(\frac{n-1}{\|\mathbf{A}\|_F^2}\right)^{n-1}$, while the ones in [HP92; YG97] contain factors like $|\det(\mathbf{A})|^2 \left(\frac{n-1}{n}\right)^{(n-1)/2}$ and row and column norms. The Schur complement-based bounds for strictly diagonally dominant matrices in [Hua08; Li20; Ois23; San21; Var75] depend on the degree of diagonal dominance, as do the Gerschgorin type bounds for rectangular matrices in [Joh89; JS98].

In contrast, we are bounding the smallest singular values of *perturbed* matrices. The expressions for small singular values in [Ste84, Theorem], [Ste06, Theorem 8], [SS90, Section V.4.2] are second-order perturbation expansions rather than bounds, and require assumptions on the singular vectors.

1.5 Overview

Our deterministic lower bounds for small singular values of A + E are based on eigenvalue bounds for $(A + E)^T (A + E)$. We present normwise absolute bounds for a single smallest singular value (section 2) and for a cluster of small singular values (section 3). The numerical experiments (section 4) confirm the qualitative increase in small singular values resulting from the demotion of the matrix to lower precision. A brief discussion of future work (section 5) concludes the paper.

2 A single smallest singular value

We perturb a matrix that has a single smallest singular value, and derive a lower bound for the smallest singular value of the perturbed matrix in terms of normwise absolute perturbations (Section 2.2), based on eigenvalue bounds (Section 2.1).

2.1 Auxiliary eigenvalue results

We square the singular values of $A \in \mathbb{R}^{m \times n}$ and consider instead the eigenvalues of the symmetric positive semi-definite matrix $B \equiv A^T A \in \mathbb{R}^{n \times n}$.

For a symmetric positive semi-definite matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ with a single smallest eigenvalue $\lambda_{\min}(\mathbf{B})$, we present two expressions for $\lambda_{\min}(\mathbf{B})$ with different assumptions (Lemmas 1 and 2), and two lower bounds in terms of normwise absolute perturbations (Theorems 2 and 3).

We assume that $\lambda_{\min}(\boldsymbol{B})$ is separated from the remaining eigenvalues, in the sense that it is strictly smaller than the smallest eigenvalue of the leading principal submatrix \boldsymbol{B}_{11} of order n-1. The equality below expresses $\lambda_{\min}(\boldsymbol{B})$ in terms of itself.

Lemma 1 (Exact expression). Let $B \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(B) \geq n - 1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{b} \ m{b}^T & eta \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-1) imes (n-1)}.$$

Then

$$(1) 0 \le \lambda_{min}(\mathbf{B}) \le \beta.$$

If also $\lambda_{min}(\mathbf{B}) < \lambda_{min}(\mathbf{B}_{11})$ then

$$\lambda_{min}(\boldsymbol{B}) = \beta - \boldsymbol{b}^{T} (\boldsymbol{B}_{11} - \lambda_{min}(\boldsymbol{B}) \boldsymbol{I})^{-1} \boldsymbol{b},$$

Proof. Abbreviate $\lambda_{\min} \equiv \lambda_{\min}(\boldsymbol{B})$. The positive semi-definiteness of \boldsymbol{B} implies the lower bound in (1), while the variational inequalities imply the upper bound,

$$0 \le \tilde{\lambda}_{\min} = \min_{\|\mathbf{x}\|_2 = 1} \mathbf{x}^T \mathbf{B} \mathbf{x} \le \mathbf{e}_n^T \mathbf{B} \mathbf{e}_n = \beta.$$

To show the expression for $\tilde{\lambda}_{\min}$, observe that the shifted matrix

$$m{B} - ilde{\lambda}_{\min} \, m{I} = egin{bmatrix} m{B}_{11} - ilde{\lambda}_{\min} \, m{I} & m{b} \ m{b}^T & eta - ilde{\lambda}_{\min} \end{bmatrix}$$

is singular. From the assumption $\tilde{\lambda}_{\min} < \lambda_{\min}(\boldsymbol{B}_{11})$ follows that $\boldsymbol{B}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I}$ is nonsingular. So we can determine the block LU decomposition $\boldsymbol{B} - \tilde{\lambda}_{\min} \boldsymbol{I} = \boldsymbol{L} \hat{\boldsymbol{U}}$ with

$$egin{aligned} m{L} &\equiv egin{bmatrix} m{I} & m{0} \ m{b}^T m{(B_{11}} - ilde{\lambda}_{\min} m{I})^{-1} & 1 \end{bmatrix}, \ & \hat{m{U}} &\equiv egin{bmatrix} m{B_{11}} - ilde{\lambda}_{\min} m{I} & m{b} \ m{0} & eta - ilde{\lambda}_{\min} - m{b}^T m{(B_{11}} - ilde{\lambda}_{\min} m{I})^{-1} m{b} \end{bmatrix}. \end{aligned}$$

Since $B - \tilde{\lambda}_{\min} I$ is singular and the unit triangular matrix I is nonsingular, the block upper triangular matrix \hat{U} has no choice but to be singular. Its leading principal submatrix $B_{11} - \tilde{\lambda}_{\min} I$ is nonsingular by assumption, which leaves the (2,2) element to be singular, but it being a scalar implies

$$\beta - \tilde{\lambda}_{\min} - \boldsymbol{b}^T (\boldsymbol{B}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1} \boldsymbol{b} = 0.$$

This gives the expression for $\tilde{\lambda}_{\min}$

If b = 0 then Lemma 1 correctly asserts that $\lambda_{\min}(B) = \beta$.

Lemma 2 below presents the same expression for $\lambda_{\min}(\mathbf{B})$ as in Lemma 1, but under a stronger albeit more useful assumption.

Lemma 2 (Exact expression with stronger assumption). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank(\mathbf{B}) $\geq n-1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{b} \ m{b}^T & eta \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-1) imes (n-1)}.$$

If $\beta < \lambda_{min}(\boldsymbol{B}_{11})$ then

$$\lambda_{min}(\boldsymbol{B}) = \beta - \boldsymbol{b}^{T} (\boldsymbol{B}_{11} - \lambda_{min}(\boldsymbol{B}) \boldsymbol{I})^{-1} \boldsymbol{b} \ge 0.$$

Proof. The upper bound (1) combined with the assumption $\beta < \lambda_{\min}(\boldsymbol{B}_{11})$ implies the assumption in Lemma 1,

(2)
$$0 \le \lambda_{\min}(\mathbf{B}) \le \beta < \lambda_{\min}(\mathbf{B}_{11}).$$

The subsequent lower bounds for $\lambda_{\min}(B)$ are informative if the offdiagonal part has small norm.

Theorem 2 (First lower bound). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(\mathbf{B}) \geq n - 1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{b} \ m{b}^T & eta \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-1) imes (n-1)}.$$

If $\beta < \lambda_{min}(\boldsymbol{B}_{11})$ then

$$\lambda_{min}(\boldsymbol{B}) \ge \beta - \boldsymbol{b}^T \boldsymbol{B}_{11}^{-1} \boldsymbol{b} - \frac{\beta \| \boldsymbol{B}_{11}^{-1} \boldsymbol{b} \|_2^2}{1 - \beta \| \boldsymbol{B}_{11}^{-1} \|_2}.$$

Proof. Abbreviate $\tilde{\lambda}_{\min} \equiv \lambda_{\min}(\boldsymbol{B})$. From (2) follows that \boldsymbol{B}_{11} and $\boldsymbol{B}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I}$ are nonsingular. Combined with the symmetric positive semi-definiteness of \boldsymbol{B}_{11} this gives

$$\tilde{\lambda}_{\min} < \lambda_{\min}(\boldsymbol{B}_{11}) = 1/\|\boldsymbol{B}_{11}^{-1}\|_2,$$

hence

(3)
$$\|\tilde{\lambda}_{\min} B_{11}^{-1}\|_2 < 1.$$

Thus we can apply the Sherman-Morrison formula [GV13, Section 2.1.4],

$$(\boldsymbol{B}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1} = \boldsymbol{B}_{11}^{-1} + \tilde{\lambda}_{\min} \ \boldsymbol{B}_{11}^{-1} (\boldsymbol{I} - \tilde{\lambda}_{\min} \ \boldsymbol{B}_{11}^{-1})^{-1} \boldsymbol{B}_{11}^{-1},$$

and substitute the above into the expression for $\tilde{\lambda}_{\min}$ from Lemma 1,

(4)
$$\tilde{\lambda}_{\min} = \beta - \boldsymbol{b}^T \boldsymbol{B}_{11}^{-1} \boldsymbol{b} - \tilde{\lambda}_{\min} \ \boldsymbol{b}^T \boldsymbol{B}_{11}^{-1} (\boldsymbol{I} - \tilde{\lambda}_{\min} \ \boldsymbol{B}_{11}^{-1})^{-1} \boldsymbol{B}_{11}^{-1} \boldsymbol{b}.$$

The symmetric positive semi-definiteness of \boldsymbol{B} implies that $\beta \geq 0$ and $\boldsymbol{b}^T \boldsymbol{B}_{11}^{-1} \boldsymbol{b} \geq 0$, hence it remains to bound the norm of the remaining summand. From the symmetry of \boldsymbol{B}_{11} and the invariance of the two-norm under transposition follows

$$||\boldsymbol{b}^{T}\boldsymbol{B}_{11}^{-1}(\boldsymbol{I} - \tilde{\lambda}_{\min} \boldsymbol{B}_{11}^{-1})^{-1}\boldsymbol{B}_{11}^{-1}\boldsymbol{b}||_{2} \leq ||\boldsymbol{b}^{T}\boldsymbol{B}_{11}^{-1}||_{2}||(\boldsymbol{I} - \tilde{\lambda}_{\min} \boldsymbol{B}_{11}^{-1})^{-1}||_{2}||\boldsymbol{B}_{11}^{-1}\boldsymbol{b}||_{2}$$

$$\leq ||\boldsymbol{B}_{11}^{-1}\boldsymbol{b}||_{2}^{2}||(\boldsymbol{I} - \tilde{\lambda}_{\min} \boldsymbol{B}_{11}^{-1})^{-1}||_{2}.$$
(5)

The inequality (3) allows us to apply the Banach lemma [GV13, Lemma 2.3.3] to bound the norm of the inverse by

$$\|(\boldsymbol{I} - \tilde{\lambda}_{\min} \ \boldsymbol{B}_{11}^{-1})^{-1}\|_{2} \leq \frac{1}{1 - \|\tilde{\lambda}_{\min} \ \boldsymbol{B}_{11}^{-1}\|_{2}} = \frac{1}{1 - \tilde{\lambda}_{\min} \|\boldsymbol{B}_{11}^{-1}\|_{2}}.$$

Substitute this into (5) and the resulting bound into the expression for $\tilde{\lambda}_{\min}$ in (4),

$$\tilde{\lambda}_{\min} \geq \beta - \boldsymbol{b}^T \boldsymbol{B}_{11}^{-1} \boldsymbol{b} - \frac{\tilde{\lambda}_{\min} \ \|\boldsymbol{B}_{11}^{-1} \boldsymbol{b}\|_2^2}{1 - \tilde{\lambda}_{\min} \ \|\boldsymbol{B}_{11}^{-1}\|_2}$$

and at last apply the upper bound (1).

The lower bound in Theorem 2 is positive if $\|\boldsymbol{b}\|_2$ is sufficiently small, in which case $\lambda_{\min}(\boldsymbol{B}) \geq \beta - \mathcal{O}(\|\boldsymbol{b}\|_2^2)$. If $\boldsymbol{b} = \boldsymbol{0}$ then (1) and Theorem 2 imply $\lambda_{\min}(\boldsymbol{B}) = \beta$.

The slightly weaker bound below focusses on a 'dominant part' of B_{11} .

Theorem 3 (Second lower bound). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(\mathbf{B}) \geq n-1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{b} \ m{b}^T & eta \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-1) imes (n-1)}.$$

If $B_{11} = C_{11} + C_{12}$ where C_{11} is symmetric positive definite with $\lambda_{min}(C_{11}) > \beta$, and C_{12} is symmetric positive semi-definite then

$$\lambda_{min}(\boldsymbol{B}) \ge \beta - \boldsymbol{b}^T \boldsymbol{C}_{11}^{-1} \boldsymbol{b} - \frac{\beta \| \boldsymbol{C}_{11}^{-1} \boldsymbol{b} \|_2^2}{1 - \beta \| \boldsymbol{C}_{11}^{-1} \|_2}.$$

Proof. Abbreviate $\tilde{\lambda}_{\min} \equiv \lambda_{\min}(\boldsymbol{B})$. From (1) and the assumption follows $\tilde{\lambda}_{\min} \leq \beta < \lambda_{\min}(\boldsymbol{C}_{11})$, hence \boldsymbol{C}_{11} and $\boldsymbol{C}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I}$ are nonsingular. Write

$$B_{11} - \tilde{\lambda}_{\min} I = \underbrace{C_{11} - \tilde{\lambda}_{\min} I}_{G} + C_{12} = G^{1/2} \left(I + \underbrace{G^{-1/2} C_{12} G^{-1/2}}_{H} \right) G^{1/2},$$

where G is symmetric positive definite, $G^{1/2}$ is its symmetric positive definite square root, and H is symmetric positive semi-definite. The Loewner partial ordering³ implies $I \leq I + H$. From [HJ13, Corollary 7.7.4] follows $(I + H)^{-1} \leq I^{-1} = I$. Thus

$$(\boldsymbol{B}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1} = \boldsymbol{G}^{-1/2} (\boldsymbol{I} + \boldsymbol{H})^{-1} \boldsymbol{G}^{-1/2}$$

 $\leq \boldsymbol{G}^{-1/2} \boldsymbol{I} \boldsymbol{G}^{-1/2} = \boldsymbol{G}^{-1} = (\boldsymbol{C}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1}.$

Substituting $(\boldsymbol{B}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1} \leq (\boldsymbol{C}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1}$ into the expression for $\tilde{\lambda}_{\min}$ in Lemma 1 gives

$$\tilde{\lambda}_{\min} \geq \beta - \boldsymbol{b}^T (\boldsymbol{C}_{11} - \tilde{\lambda}_{\min} \boldsymbol{I})^{-1} \boldsymbol{b}.$$

We continue as in the proof of Theorem 2 with the Sherman-Morrison formula [GV13, Section 2.1.4],

$$egin{aligned} ilde{\lambda}_{\min} & \geq eta - m{b}^T m{C}_{11}^{-1} m{b} - ilde{\lambda}_{\min} \,\, m{b}^T m{C}_{11}^{-1} (m{I} - ilde{\lambda}_{\min} \,\, m{C}_{11}^{-1})^{-1} m{C}_{11}^{-1} m{b} \ & \geq eta - m{b}^T m{C}_{11}^{-1} m{b} - rac{ ilde{\lambda}_{\min} \,\, \|m{C}_{11}^{-1} m{b}\|_2^2}{1 - ilde{\lambda}_{\min} \,\, \|m{C}_{11}^{-1}\|_2}, \end{aligned}$$

and at last apply (1).

If $C_{12} = 0$, then Theorem 3 reduces to Theorem 2.

2.2 A lower bound for the smallest singular value

We consider a matrix with a distinct smallest singular value. Based on the eigenvalue bounds in section 2.1, we derive a lower bound for the smallest singular value of a perturbed matrix (Theorem 4) in terms of normwise absolute perturbations. We start with a summary of all assumptions (Assumptions 2.1), and end with a discussion of their generality (Remark 2.1).

Assumptions 2.1. Let $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ have $\operatorname{rank}(A) \geq n - 1$. Let $A = U\Sigma V^T$ be a full singular value decomposition, where $\Sigma \in \mathbb{R}^{m \times n}$ is diagonal, and $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices. Partition commensurately,

$$oldsymbol{\Sigma} = egin{bmatrix} oldsymbol{\Sigma}_1 & oldsymbol{0} \ oldsymbol{0} & \sigma_{min} \ oldsymbol{0} & oldsymbol{0} \end{bmatrix}, \qquad oldsymbol{E} = oldsymbol{U} egin{bmatrix} oldsymbol{E}_{11} & oldsymbol{e}_{12} \ oldsymbol{e}_{21} & oldsymbol{e}_{22} \ oldsymbol{E}_{31} & oldsymbol{e}_{32} \end{bmatrix} oldsymbol{V}^T,$$

where $\Sigma_1 \in \mathbb{R}^{(n-1)\times(n-1)}$ is nonsingular diagonal, and $\sigma_{min} \geq 0$.

For a matrix with a single smallest singular value, we corroborate the observation that 'small singular values tend to increase' [SS90, page 266]. Motivated by the second-order expressions in terms of absolute perturbations [SS90, Section V.4.2] and [Ste06, Theorem 8], we derive a true lower bound.

Theorem 4. Let $A, E \in \mathbb{R}^{m \times n}$ satisfy Assumptions 2.1. If $1/\|\mathbf{\Sigma}_1^{-1}\|_2 > 4\|\mathbf{E}\|_2$ and $\sigma_{min} < \|\mathbf{E}\|_2$, then

$$\sigma_{min}(\mathbf{A} + \mathbf{E})^2 \ge \|\mathbf{e}_{32}\|_2^2 + (\sigma_{min} + e_{22})^2 - r_3 - r_4,$$

where r_3 contains terms of order 3,

$$r_3 \equiv 2\boldsymbol{e}_{12}^T \mathbf{w} \qquad \mathbf{w} \equiv (\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-T} \begin{bmatrix} \boldsymbol{e}_{21} & \boldsymbol{E}_{31}^T \end{bmatrix} \begin{bmatrix} \boldsymbol{e}_{22} + \sigma_{min} \\ \boldsymbol{e}_{32} \end{bmatrix},$$

and r_4 contains terms of order 4 and higher,

$$r_4 \equiv \|\mathbf{w}\|_2^2 + 4 \frac{\|\mathbf{E}\|_2^2 \|(\mathbf{\Sigma}_1 + \mathbf{E}_{11})^{-1} (\mathbf{e}_{12} + \mathbf{w})\|_2^2}{1 - 4\|\mathbf{E}\|_2^2 \|(\mathbf{\Sigma}_1 + \mathbf{E}_{11})^{-1}\|_2^2}.$$

³For Hermitian matrices **A** and **B**, $A \leq B$ means that A - B is positive semi-definite.

Proof. We square the singular values of A + E, and consider the eigenvalues of

$$m{B} \equiv (m{A} + m{E})^T (m{A} + m{E}) = m{V} egin{bmatrix} m{B}_{11} & m{b} \\ m{b}^T & eta \end{bmatrix} m{V}^T$$

where

(6)
$$B_{11} = \underbrace{(\Sigma_1 + E_{11})^T (\Sigma_1 + E_{11})}_{C_{11}} + \underbrace{e_{21} e_{21}^T + E_{31}^T E_{31}}_{C_{12}}$$

(7)
$$\beta = \|\mathbf{e}_{12}\|_{2}^{2} + (\sigma_{\min} + e_{22})^{2} + \|\mathbf{e}_{32}\|_{2}^{2}$$
$$\mathbf{b} = (\mathbf{\Sigma}_{1} + \mathbf{E}_{11})^{T} \mathbf{e}_{12} + \mathbf{e}_{21}(\sigma_{\min} + e_{22}) + \mathbf{E}_{31}^{T} \mathbf{e}_{32}.$$

From $\sigma_{\min}(\Sigma_1) > 4 \|E\|_2$ follows that C_{11} is symmetric positive definite, while C_{12} is symmetric positive semi-definite and contains only second order terms. Abbreviate $\tilde{\lambda}_{\min} \equiv \lambda_{\min}(B) = \sigma_{\min}(A + E)^2$.

The proof proceeds in two steps:

- 1. Confirming that C_{11} satisfies the assumptions of Theorem 3.
- 2. Deriving the lower bound for $\tilde{\lambda}_{\min}$ from Theorem 3.
- 1. Confirm that C_{11} satisfies the assumptions of Theorem 3 We show that $\lambda_{\min}(C_{11}) > \beta$, by bounding β from above and $\lambda_{\min}(C_{11})$ from below.

Regarding the upper bound for β , the expression (7) and the assumption $\sigma_{\min} < ||E||_2$ imply

(8)
$$\beta = \left\| \begin{bmatrix} e_{12}^T & e_{22} + \sigma_{\min} & e_{32}^T \end{bmatrix}^T \right\|_2^2 \le (\sigma_{\min} + \|\mathbf{E}e_n\|_2)^2 \le 4\|\mathbf{E}\|_2^2.$$

Regarding the lower bound for $\lambda_{\min}(C_{11})$, view $C_{11} = (\Sigma_1 + E_{11})^T (\Sigma_1 + E_{11})$ as a singular value problem, so that $\lambda_{\min}(C_{11}) = \sigma_{\min}(\Sigma_1 + E_{11})^2$. The well-conditioning of singular values [GV13, Corollary 8.6.2] implies

$$|\sigma_{\min}(\Sigma_1 + E_{11}) - \sigma_{\min}(\Sigma_1)| \le ||E_{11}||_2 \le ||E||_2.$$

Adding the assumption $\sigma_{\min}(\mathbf{\Sigma}_1) = 1/\|\mathbf{\Sigma}_1^{-1}\|_2 > 4\|\mathbf{E}\|_2$ gives

$$\sigma_{\min}(\Sigma_1 + E_{11}) \ge \sigma_{\min}(\Sigma_1) - ||E||_2 > 4 ||E||_2 - ||E||_2 = 3 ||E||_2.$$

Now combine this lower bound for $\lambda_{\min}(C_{11})$ with (8).

$$\lambda_{\min}(C_{11}) = \sigma_{\min}(\Sigma_1 + E_{11})^2 > 9 \|E\|_2^2 > 4 \|E\|_2^2 > \beta.$$

Hence $\lambda_{\min}(C_{11}) > \beta$, and C_{11} satisfies the assumptions of Theorem 3.

2. Derive the lower bound for $\tilde{\lambda}_{\min}$ from Theorem 3 In this bound,

(9)
$$\tilde{\lambda}_{\min} \ge \beta - \boldsymbol{b}^T \boldsymbol{C}_{11}^{-1} \boldsymbol{b} - \frac{\beta \| \boldsymbol{C}_{11}^{-1} \boldsymbol{b} \|_2^2}{1 - \beta \| \boldsymbol{C}_{11}^{-1} \|_2},$$

where the key term is $C_{11}^{-1}b$. Insert the expression for **b** from (7),

(10)
$$(\mathbf{\Sigma}_{1} + \mathbf{E}_{11})^{-T} \mathbf{b} = \mathbf{e}_{12} + (\mathbf{\Sigma}_{1} + \mathbf{E}_{11})^{-T} \left(\mathbf{e}_{21} (\sigma_{\min} + e_{22}) + \mathbf{E}_{31}^{T} \mathbf{e}_{32} \right)$$

$$= \mathbf{e}_{12} + \mathbf{w}.$$

Combine the expression for C_{11} from (6) with (10),

(11)
$$C_{11}^{-1}\boldsymbol{b} = (\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1} \underbrace{(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-T} \boldsymbol{b}}_{\boldsymbol{e}_{12} + \mathbf{w}} = (\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1} (\boldsymbol{e}_{12} + \mathbf{w})$$

Multiply the above by \boldsymbol{b}^T on the left, and use (10)

$$egin{aligned} m{b}^T m{C}_{11}^{-1} m{b} &= m{b}^T (m{\Sigma}_1 + m{E}_{11})^{-1} (m{\Sigma}_1 + m{E}_{11})^{-T} m{b} = (m{e}_{12} + m{w})^T (m{e}_{12} + m{w}) \ &= \|m{e}_{12} + m{w}\|_2^2 = \|m{e}_{12}\|_2^2 + 2m{e}_{12}^T m{w} + \|m{w}\|_2^2. \end{aligned}$$

Substitute the above, and β from (7) into the first two summands of (9),

(12)
$$\beta - \boldsymbol{b}^{T} \boldsymbol{C}_{11}^{-1} \boldsymbol{b} = (\|\boldsymbol{e}_{12}\|_{2}^{2} + (\sigma_{\min} + \boldsymbol{e}_{22})^{2} + \|\boldsymbol{e}_{32}\|_{2}^{2}) - (\|\boldsymbol{e}_{12}\|_{2}^{2} + 2\boldsymbol{e}_{12}^{T} \mathbf{w} + \|\mathbf{w}\|_{2}^{2})$$
$$= \|\boldsymbol{e}_{32}\|_{2}^{2} + (\sigma_{\min} + \boldsymbol{e}_{22})^{2} - \underbrace{2\boldsymbol{e}_{12}^{T} \mathbf{w}}_{r_{2}} - \|\mathbf{w}\|_{2}^{2}.$$

Substitute the bound for β in (8), and (11) into the third summand of (9),

(13)
$$\frac{\beta \| \boldsymbol{C}_{11}^{-1} \boldsymbol{b} \|_{2}^{2}}{1 - \beta \| \boldsymbol{C}_{11}^{-1} \|_{2}} \le 4 \frac{\| \boldsymbol{E} \|_{2}^{2} \| (\boldsymbol{\Sigma}_{1} + \boldsymbol{E}_{11})^{-1} (\boldsymbol{e}_{12} + \mathbf{w}) \|_{2}^{2}}{1 - 4 \| \boldsymbol{E} \|_{2}^{2} \| (\boldsymbol{\Sigma}_{1} + \boldsymbol{E}_{11})^{-1} \|_{2}^{2}}.$$

Inserting (12) and (13) into (9) gives

$$\tilde{\lambda}_{\min} \ge \|\boldsymbol{e}_{32}\|_{2}^{2} + (\sigma_{\min} + e_{22})^{2} - r_{3} - \underbrace{\left(\|\mathbf{w}\|_{2}^{2} + 4\frac{\|\boldsymbol{E}\|_{2}^{2}\|(\boldsymbol{\Sigma}_{1} + \boldsymbol{E}_{11})^{-1}(\boldsymbol{e}_{12} + \mathbf{w})\|_{2}^{2}}_{r_{4}}\right)}_{r_{4}}.$$

Remark 2.1. The assumptions in Theorem 4 are not restrictive. Only a small gap of $3||E||_2$ is required to separate the smallest singular value of A from the remaining singular values,

$$\sigma_{min}(\mathbf{A}) < \|\mathbf{E}\|_2 < 4\|\mathbf{E}\|_2 \le 1/\|\mathbf{\Sigma}_1^{-1}\|_2.$$

3 A cluster of small singular values

We extend the results in Section 2 from a single smallest singular value to a cluster of small singular values. To this end, we derive lower bounds for the small singular values of the perturbed matrix in terms of normwise absolute perturbations (Section 3.2), based on eigenvalue bounds (Section 3.1).

3.1 Auxiliary eigenvalue results

We square the singular values of $\mathbf{A} \in \mathbb{R}^{m \times n}$ and consider instead the eigenvalues of the symmetric positive semi-definite matrix $\mathbf{B} \equiv \mathbf{A}^T \mathbf{A} \in \mathbb{R}^{n \times n}$.

For a symmetric positive semi-definite matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ with a cluster of r small eigenvalues, we present an expression for these eigenvalues (Lemma 3), and two lower bounds in terms of normwise absolute perturbations (Theorems 5 and 6).

We assume that the r small eigenvalues are separated from the remaining ones, in the sense that they are strictly smaller than the smallest eigenvalue of the leading principal submatrix B_{11} of order n-r. The eigenvalues are labelled so that

$$\lambda_n(\mathbf{B}) \leq \cdots \leq \lambda_{n-r+1}(\mathbf{B}) < \lambda_{n-r}(\mathbf{B}) \leq \cdots \leq \lambda_1(\mathbf{B}).$$

The equality below expresses the smallest eigenvalues in terms of themselves, and represents an extension of Lemma 2 to clusters.

Lemma 3 (Exact expression). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(\mathbf{B}) \geq n - r$ for some $r \geq 1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{B}_{12} \ m{B}_{12}^T & m{B}_{22} \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-r) imes (n-r)}, \quad m{B}_{22} \in \mathbb{R}^{r imes r}.$$

If $\|\mathbf{B}_{22}\|_{2} < \lambda_{min}(\mathbf{B}_{11})$ then

$$\lambda_{n-r+j}(\mathbf{B}) = \lambda_j \left(\mathbf{B}_{22} - \mathbf{B}_{12}^T (\mathbf{B}_{11} - \lambda_{n-r+j}(\mathbf{B}) \mathbf{I})^{-1} \mathbf{B}_{12} \right), \qquad 1 \le j \le r,$$

where

(14)
$$0 \le \lambda_{n-r+j}(\mathbf{B}) \le \|\mathbf{B}_{22}\|_2, \qquad 1 \le j \le r.$$

Proof. Abbreviate $\tilde{\lambda}_{n-r+j} \equiv \lambda_{n-r+j}(\boldsymbol{B})$, $1 \leq j \leq r$. The lower bound in (14) follows from the positive semi-definiteness of \boldsymbol{B} , and the upper bound from the Cauchy interlace theorem [Par80, Section 10.1]

$$\tilde{\lambda}_{n-r+j} \le \lambda_j(B_{22}) \le \lambda_{\max}(B_{22}) = ||B_{22}||_2, \quad 1 \le j \le r.$$

Combining this with the assumption $\|\boldsymbol{B}_{22}\|_2 < \lambda_{\min}(\boldsymbol{B}_{11})$ shows

(15)
$$\tilde{\lambda}_{n-r+j} \le \|\boldsymbol{B}_{22}\|_2 < \lambda_{\min}(\boldsymbol{B}_{11}), \qquad 1 \le j \le r$$

Hence is $\boldsymbol{B}_{11} - \tilde{\lambda}_{n-r+j} \boldsymbol{I}$ is nonsingular.

To derive the expression for $\tilde{\lambda}_{n-r+j}$, we start as in the proof of [Par80, Theorem (10-1-2)]. The shifted matrix $\boldsymbol{B} - \tilde{\lambda}_{n-r+j} \boldsymbol{I}$ has at most n-r+j-1 positive eigenvalues, at least one zero eigenvalue, and at most r-j negative eigenvalues, $1 \le j \le r$. Perform the congruence transformation

$$egin{aligned} m{B} - ilde{\lambda}_{n-r+j} \, m{I} & m{L} & m{0} \ m{0} & m{S} \end{bmatrix} m{L}^T, \qquad m{L} \equiv egin{bmatrix} m{I} & m{0} \ m{B}_{12}^T (m{B}_{11} - ilde{\lambda}_{n-r+j} \, m{I})^{-1} & m{I} \end{bmatrix} \end{aligned}$$

where

(16)
$$S \equiv B_{22} - \tilde{\lambda}_{n-r+j} \mathbf{I} - B_{12}^{T} (B_{11} - \tilde{\lambda}_{n-r+j} \mathbf{I})^{-1} B_{12}, \qquad 1 \le j \le r.$$

From (15) follows that $B_{11} - \tilde{\lambda}_{n-r+j} I$ has n-r positive eigenvalues. Combining this with the inertia preservation of congruence transformations implies that S has at most r-j positive eigenvalues, at least one zero eigenvalue $\lambda_j(S) = 0$, and at least j-1 negative eigenvalues, $1 \le j \le r$. Insert (16) into $\lambda_j(S) = 0$, and exploit the fact that the shift $\tilde{\lambda}_{n-r+j} I$ does not change the algebraic eigenvalue ordering, to obtain the expression for $\tilde{\lambda}_{n-r+j}$, $1 \le j \le r$.

By restricting ourselves to a 'dominant part' of B_{11} , we weaken the expression in Lemma 3 to a lower bound, which allows the eigenvalues to be negative.

Lemma 4 (Lower bound). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(\mathbf{B}) \geq n - r$ for some $r \geq 1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{B}_{12} \ m{B}_{12}^T & m{B}_{22} \end{bmatrix} \qquad ext{where} \quad m{B}_{11} \in \mathbb{R}^{(n-r) imes (n-r)}, \quad m{B}_{22} \in \mathbb{R}^{r imes r}.$$

Let $\boldsymbol{B}_{11} = \boldsymbol{C}_{11} + \boldsymbol{C}_{12}$ where $\boldsymbol{C}_{11} \in \mathbb{R}^{(n-r)\times(n-r)}$ is symmetric positive definite and $\boldsymbol{C}_{12} \in \mathbb{R}^{(n-r)\times(n-r)}$ is symmetric positive semi-definite. If

$$\widehat{m{B}} \equiv egin{bmatrix} m{C}_{11} & m{B}_{12} \ m{B}_{12}^T & m{B}_{22} \end{bmatrix}$$

with $\lambda_{min}(C_{11}) > ||B_{22}||_2$, then

(17)
$$\lambda_{n-r+i}(\mathbf{B}) \ge \lambda_{n-r+i}(\widehat{\mathbf{B}})$$

(18)
$$= \lambda_j \left(\mathbf{B}_{22} - \mathbf{B}_{12}^T (\mathbf{C}_{11} - \lambda_{n-r+j}(\widehat{\mathbf{B}}) \mathbf{I})^{-1} \mathbf{B}_{12} \right), \quad 1 \le j \le r,$$

where

$$\lambda_{n-r+j}(\widehat{\boldsymbol{B}}) \le \|\boldsymbol{B}_{22}\|_{2}$$

(20)
$$\left\| \left(\boldsymbol{C}_{11} - \lambda_{n-r+j}(\widehat{\boldsymbol{B}}) \boldsymbol{I} \right)^{-1} \right\|_{2} \le \frac{\| \boldsymbol{C}_{11}^{-1} \|_{2}}{1 - \| \boldsymbol{B}_{22} \|_{2} \| \boldsymbol{C}_{11}^{-1} \|_{2}}, \qquad 1 \le j \le r.$$

Proof. The proof proceeds in four steps.

Proof of (17) The symmetric positive semi-definiteness of C_{12} and Weyl's monotonicity theorem [HJ13, Corollary 4.3.3] imply

$$\lambda_j(\mathbf{B}) \ge \lambda_j(\widehat{\mathbf{B}}), \qquad 1 \le j \le n.$$

Now we concentrate on the eigenvalues of $\widehat{\boldsymbol{B}}$, and abbreviate $\widehat{\lambda}_{n-r+j} \equiv \lambda_{n-r+j}(\widehat{\boldsymbol{B}}), \ 1 \leq j \leq r$.

Proof of (19) Apply the Cauchy interlace theorem [Par80, Section 10.1] to \hat{B} ,

$$\hat{\lambda}_{n-r+j} \le \lambda_j(B_{22}) \le \lambda_{\max}(B_{22}) = ||B_{22}||_2, \quad 1 \le j \le r.$$

Combining this with the assumption $\|B_{22}\|_2 < \lambda_{\min}(C_{11})$ shows

$$\hat{\lambda}_{n-r+j} \le \|B_{22}\|_2 < \lambda_{\min}(C_{11}), \qquad 1 \le j \le r.$$

Hence $C_{11} - \widehat{\lambda}_{n-r+j} I$ is nonsingular, which holds in particular if $\widehat{\lambda}_{n-r+j} < 0$.

Proof of (18) To derive the expression for $\hat{\lambda}_{n-r+j}$, apply the proof of Lemma 3 to the eigenvalues of \hat{B} . This proof relies only on the signs of eigenvalues of shifted matrices, and does not require positive semi-definiteness of the host matrix \hat{B} .

Proof of (20) Fix some $1 \le j \le r$ for the inverse in (18). Then factor out C_{11}^{-1} ,

$$(C_{11} - \widehat{\lambda}_{n-r+j} I)^{-1} = C_{11}^{-1} D$$
 where $D \equiv (I - \widehat{\lambda}_{n-r+j} C_{11}^{-1})^{-1}$,

and take norms,

$$\left\| (\boldsymbol{C}_{11} - \widehat{\lambda}_{n-r+j} \ \boldsymbol{I})^{-1} \right\|_{2} \le \|\boldsymbol{C}_{11}^{-1}\|_{2} \|\boldsymbol{D}\|_{2}.$$

To bound $\|D\|_2$, consider the eigenvalue decomposition $C_{11} = W\Lambda W^T$, where W is an orthogonal matrix, and the diagonal matrix

$$\Lambda = \operatorname{diag} (\gamma_1 \quad \cdots \quad \gamma_{n-r}) \in \mathbb{R}^{(n-r)\times(n-r)}$$

has positive diagonal elements $\gamma_{\ell} > 0$. Thus D has an eigenvalue decomposition $D = W(I - \hat{\lambda}_{n-r+j} \Lambda^{-1})^{-1} W^T$ with eigenvalues

$$\lambda_{\ell}(\mathbf{D}) = 1 / \left(1 - \frac{\widehat{\lambda}_{n-r+j}}{\gamma_{\ell}} \right), \qquad 1 \le \ell \le n-r.$$

Case 1: If $\widehat{\lambda}_{n-r+j} \geq 0$, then (19) implies

$$0 \le \frac{\widehat{\lambda}_{n-r+j}}{\gamma_{\ell}} \le \frac{\|\boldsymbol{B}_{22}\|_{2}}{\lambda_{\min}(\boldsymbol{C}_{11})} = \|\boldsymbol{B}_{22}\|_{2} \|\boldsymbol{C}_{11}^{-1}\|_{2} < 1, \qquad 1 \le \ell \le n-r.$$

Hence

(21)
$$\|\boldsymbol{D}\|_{22} = \max_{1 \le \ell \le n-r} |\lambda_j(\boldsymbol{D})| \le \frac{1}{1 - \|\boldsymbol{B}_{22}\|_2 \|\boldsymbol{C}_{11}^{-1}\|_2}.$$

Case 2: If $\hat{\lambda}_{n-r+j} < 0$, then $\gamma_{\ell} > 0$ and (19) imply

$$1 - \frac{\widehat{\lambda}_{n-r+j}}{\gamma_{\ell}} = 1 + \frac{|\widehat{\lambda}_{n-r+j}|}{\gamma_{\ell}} > 1 > 1 - \|\boldsymbol{B}_{22}\|_{2} \|\boldsymbol{C}_{11}^{-1}\|_{2}, \qquad 1 \le \ell \le n - r.$$

Again, as in (21) we conclude

$$\|oldsymbol{D}\|_{22} = \max_{1 \leq \ell \leq n-r} |\lambda_j(oldsymbol{D})| \leq rac{1}{1 - \|oldsymbol{B}_{22}\|_2 \|oldsymbol{C}_{11}^{-1}\|_2}.$$

Since we fixed an arbitrary j to show (20), it holds for all $1 \le j \le r$.

The subsequent lower bounds are informative if the offdiagonal part has small norm. We start with an extension of Theorem 2.

Theorem 5 (First lower bound). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(\mathbf{B}) \geq n - r$ for some $r \geq 1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{B}_{12} \ m{B}_{12}^T & m{B}_{22} \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-r) imes (n-r)}, \quad m{B}_{22} \in \mathbb{R}^{r imes r}.$$

If $\|\boldsymbol{B}_{22}\|_{2} < \lambda_{min}(\boldsymbol{B}_{11})$ then $\lambda_{n-r+j}(\boldsymbol{B}) \geq \lambda_{j}(\boldsymbol{Z}_{j}), 1 \leq j \leq r$, where

$$oldsymbol{Z}_{i} \equiv oldsymbol{B}_{22} - oldsymbol{B}_{12}^{T} oldsymbol{B}_{11}^{-1} oldsymbol{B}_{12} - \|oldsymbol{B}_{22}\|_{2} oldsymbol{B}_{12}^{T} oldsymbol{B}_{11}^{-1} (oldsymbol{I} - \lambda_{n-r+j}(oldsymbol{B}) oldsymbol{B}_{11}^{-1})^{-1} oldsymbol{B}_{11}^{-1} oldsymbol{B}_{12}.$$

Proof. Abbreviate $\tilde{\lambda}_{n-r+j} \equiv \lambda_{n-r+j}(\boldsymbol{B})$, $1 \leq j \leq r$. As in the proof of Theorem 2, apply the Sherman-Morrison formula [GV13, Section 2.1.4]

$$(\boldsymbol{B}_{11} - \tilde{\lambda}_{n-r+j} \boldsymbol{I})^{-1} = \boldsymbol{B}_{11}^{-1} + \tilde{\lambda}_{n-r+j} \boldsymbol{B}_{11}^{-1} (\boldsymbol{I} - \tilde{\lambda}_{n-r+j} \boldsymbol{B}_{11}^{-1})^{-1} \boldsymbol{B}_{11}^{-1},$$

and substitute the above into the expressions for

(22)
$$\tilde{\lambda}_{n-r+j} = \lambda_j(\boldsymbol{M}_j), \qquad 1 \le j \le r$$

from Lemma 3 where

$$egin{aligned} m{M}_j &\equiv m{B}_{22} - m{B}_{12}^T (m{B}_{11} - ilde{\lambda}_{n-r+j} \, m{I})^{-1} m{B}_{12} \ &= m{B}_{22} - m{B}_{12}^T m{B}_{11}^{-1} m{B}_{12} - ilde{\lambda}_{n-r+j} \, m{B}_{12}^{-1} m{B}_{11}^{-1} (m{I} - ilde{\lambda}_{n-r+j} \, m{B}_{11}^{-1})^{-1} m{B}_{11}^{-1} m{B}_{12}. \end{aligned}$$

From (14) follows the Loewner bound

$$oldsymbol{M}_i \succeq oldsymbol{Z}_i \equiv oldsymbol{B}_{22} - oldsymbol{B}_{12}^T oldsymbol{B}_{11}^{-1} oldsymbol{B}_{12} - \|oldsymbol{B}_{22}\|_2 oldsymbol{B}_{12}^T oldsymbol{B}_{11}^{-1} (oldsymbol{I} - ilde{\lambda}_{r+i} oldsymbol{B}_{11}^{-1})^{-1} oldsymbol{B}_{11}^{-1} oldsymbol{B}_{12}.$$

This and (22) imply
$$\tilde{\lambda}_{n-r+j} = \lambda_j(\boldsymbol{M}_j) \geq \lambda_j(\boldsymbol{Z}_j), 1 \leq j \leq r$$
 [HJ13, Corollary 7.7.4].

The slightly weaker bound below extends Theorem 3 and focusses on a 'dominant' part of B_{11} . This establishes the connection to Theorem 7, where B represents the perturbed matrix and the low order terms in B_{11} are captured by C_{11} .

Theorem 6 (Second lower bound). Let $\mathbf{B} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite with rank $(\mathbf{B}) \geq n - r$ for some $r \geq 1$, and partition

$$m{B} = egin{bmatrix} m{B}_{11} & m{B}_{12} \ m{B}_{12}^T & m{B}_{22} \end{bmatrix} \qquad where \quad m{B}_{11} \in \mathbb{R}^{(n-r) imes (n-r)}, \quad m{B}_{22} \in \mathbb{R}^{r imes r}.$$

If $B_{11} = C_{11} + C_{12}$ where C_{11} is symmetric positive definite with $\lambda_{min}(C_{11}) > \|B_{22}\|_2$, and C_{12} is symmetric positive semi-definite, then

$$\lambda_{n-r+j}(m{B}) \geq \lambda_{j} \left(m{B}_{22} - m{B}_{12}^T m{C}_{11}^{-1} m{B}_{12}
ight) - rac{\|m{B}_{22}\|_{2} \|m{C}_{11}^{-1} m{B}_{12}\|_{2}^{2}}{1 - \|m{B}_{22}\|_{2} \|m{C}_{11}^{-1}\|_{2}}, \qquad 1 \leq j \leq r.$$

Proof. Define

$$\widehat{m{B}} \equiv egin{bmatrix} m{C}_{11} & m{B}_{12} \ m{B}_{12}^T & m{B}_{22} \end{bmatrix},$$

and abbreviate $\hat{\lambda}_{n-r+j} \equiv \lambda_{n-r+j}(\hat{B}), 1 \leq j \leq r$. From (18) in Lemma 4 follows

$$\lambda_{n-r+j}(\mathbf{B}) \ge \widehat{\lambda}_{n-r+j} = \lambda_j \left(\mathbf{B}_{22} - \mathbf{B}_{12}^T (\mathbf{C}_{11} - \widehat{\lambda}_{n-r+j} \mathbf{I})^{-1} \mathbf{B}_{12} \right), \quad 1 \le j \le r,$$

We proceed as in the proof of Theorem 5, and apply the Sherman-Morrison formula [GV13, Section 2.1.4],

$$(C_{11} - \widehat{\lambda}_{n-r+j} I)^{-1} = C_{11}^{-1} + \widehat{\lambda}_{n-r+j} C_{11}^{-1} (I - \widehat{\lambda}_{n-r+j} C_{11}^{-1})^{-1} C_{11}^{-1},$$

and (19) to the expression for

(23)
$$\widehat{\lambda}_{n-r+j} = \lambda_j(\boldsymbol{M}_j), \qquad 1 \le j \le r,$$

from Lemma 4, where

$$egin{aligned} m{M}_j &\equiv m{B}_{22} - m{B}_{12}^T (m{C}_{11} - \widehat{\lambda}_{n-r+j} \, m{I})^{-1} m{B}_{12}, & 1 \leq j \leq r \ &= m{B}_{22} - m{B}_{12}^T m{C}_{11}^{-1} m{B}_{12} - \widetilde{\lambda}_{n-r+j} \, m{B}_{12}^{-T} m{C}_{11}^{-1} (m{I} - \widetilde{\lambda}_{n-r+j} \, m{C}_{11}^{-1})^{-1} m{C}_{11}^{-1} m{B}_{12}. \end{aligned}$$

From (19) and (20) follows the lower bound

$$egin{aligned} m{M}_j \succeq m{B}_{22} - m{B}_{12}^T m{C}_{11}^{-1} m{B}_{12} - \|m{B}_{22}\|_2 \, m{B}_{12}^T m{C}_{11}^{-1} (m{I} - ilde{\lambda}_{r+j} \, m{C}_{11}^{-1})^{-1} m{C}_{11}^{-1} m{B}_{12} \ & \succeq m{Z} \equiv m{B}_{22} - m{B}_{12}^T m{C}_{11}^{-1} m{B}_{12} - \underbrace{\frac{\|m{B}_{22}\|_2 \, \|m{C}_{11}^{-1} m{B}_{12}\|_2^2}_{\gamma} m{I}, & 1 \leq j \leq r. \end{aligned}$$

Thus, $M_j \succeq \mathbf{Z}$, $1 \le j \le r$. The Loewner properties [HJ13, Corollary 7.7.4] imply the same for the eigenvalues, $\lambda_j(M_j) \ge \lambda_j(\mathbf{Z})$, $1 \le j \le r$. Combine this with the well conditioning of eigenvalues [GV13, Theorem 8.1.5],

$$\tilde{\lambda}_{n-r+j} = \lambda_j(\boldsymbol{M}_j) \ge \lambda_j(\boldsymbol{Z}) \ge \lambda_j(\boldsymbol{B}_{22} - \boldsymbol{B}_{12}^T \boldsymbol{B}_{11}^{-1} \boldsymbol{B}_{12}) - \gamma, \qquad 1 \le j \le r.$$

Theorem 6 reduces to Theorem 3 for r = 1, and to Theorem 5 for $C_{12} = 0$.

3.2 A lower bound for a cluster of smallest singular values

We extend the bound for a single smallest singular value in section 2.2 to a cluster of smallest singular values. The resulting lower bound for the cluster of perturbed smallest singular values (Theorem 7) is based on the eigenvalue bounds in section 3.1, and expressed in terms of normwise absolute perturbations. We start with a summary of all assumptions (Assumptions 3.1), and end with a discussion of their generality (Remark 3.1).

Assumptions 3.1. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$ have $\operatorname{rank}(\mathbf{A}) \geq n - r$ for some $r \geq 1$. Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ be a full singular value decomposition, where $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ is diagonal, and $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$ are orthogonal matrices. Partition commensurately,

$$oldsymbol{\Sigma} = egin{bmatrix} oldsymbol{\Sigma}_1 & oldsymbol{0} \ oldsymbol{0} & oldsymbol{\Sigma}_2 \ oldsymbol{0} & oldsymbol{0} \end{bmatrix}, \qquad oldsymbol{E} = oldsymbol{U} egin{bmatrix} oldsymbol{E}_{11} & oldsymbol{E}_{12} \ oldsymbol{E}_{21} & oldsymbol{E}_{22} \ oldsymbol{E}_{31} & oldsymbol{E}_{32} \end{bmatrix} oldsymbol{V}^T,$$

where $\Sigma_1 \in \mathbb{R}^{(n-r)\times(n-r)}$ is nonsingular diagonal, and $\Sigma_2 \in \mathbb{R}^{r\times r}$ is diagonal.

This bound below extends Theorem 4, and reduces to it for r = 1.

Theorem 7. Let $A, E \in \mathbb{R}^{m \times n}$ satisfy Assumptions 3.1. If $1/\|\Sigma_1^{-1}\|_2 > 4\|E\|_2$ and $\|\Sigma_2\|_2 < \|E\|_2$, then

$$\sigma_{n-r+j}(\boldsymbol{A}+\boldsymbol{E})^2 \ge \lambda_j (\boldsymbol{E}_{32}^T \boldsymbol{E}_{32} + (\boldsymbol{\Sigma}_2 + \boldsymbol{E}_{22})^T (\boldsymbol{\Sigma}_2 + \boldsymbol{E}_{22}) - \boldsymbol{R}_3) - r_4, \quad 1 \le j \le r,$$

where \mathbf{R}_3 contains terms of order 3

$$oldsymbol{R}_3 \equiv oldsymbol{E}_{12}^T oldsymbol{W} + oldsymbol{W}^T oldsymbol{E}_{12}, \qquad oldsymbol{W} \equiv (oldsymbol{\Sigma}_1 + oldsymbol{E}_{11})^{-T} egin{bmatrix} oldsymbol{E}_{21} & oldsymbol{E}_{31}^T \end{bmatrix} egin{bmatrix} oldsymbol{\Sigma}_2 + oldsymbol{E}_{22} \ oldsymbol{E}_{32} \end{bmatrix}$$

and r_4 contains terms of order 4 and higher,

$$r_4 \equiv \|\boldsymbol{W}\|_2^2 + 4 \frac{\|\boldsymbol{E}\|_2^2 \|(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1} (\boldsymbol{E}_{12} + \boldsymbol{W})\|_2^2}{1 - 4\|\boldsymbol{E}\|_2^2 \|(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1}\|_2^2}.$$

Proof. We square the singular values of A + E and consider the eigenvalues of the perturbed matrix

$$oldsymbol{B} \equiv (oldsymbol{A} + oldsymbol{E})^T (oldsymbol{A} + oldsymbol{E}) = oldsymbol{V} egin{bmatrix} oldsymbol{B}_{11} & oldsymbol{B}_{12} \ oldsymbol{B}_{12}^T & oldsymbol{B}_{22} \end{bmatrix} oldsymbol{V}^T$$

where

(24)
$$B_{11} = \underbrace{(\Sigma_1 + E_{11})^T (\Sigma_1 + E_{11})}_{C_{11}} + \underbrace{E_{21}^T E_{21} + E_{31}^T E_{31}}_{C_{12}}$$

(25)
$$B_{22} = E_{12}^T E_{12} + (\Sigma_2 + E_{22})^T (\Sigma_2 + E_{22}) + E_{32}^T E_{32}$$

$$B_{12} = (\Sigma_1 + E_{11})^T E_{12} + E_{21}^T (\Sigma_2 + E_{22}) + E_{31}^T E_{32}.$$

From $\sigma_{\min}(\Sigma_1) > 4 \|\boldsymbol{E}\|_2$ follows that \boldsymbol{C}_{11} is symmetric positive definite, while \boldsymbol{C}_{12} is symmetric positive semi-definite and contains only second order terms. Abbreviate $\tilde{\lambda}_{n-r+j} \equiv \lambda_{n-r+j}(\boldsymbol{B}) = \sigma_{n-r+j}(\boldsymbol{A} + \boldsymbol{E})^2$, $1 \leq j \leq r$.

The proof proceeds in two steps:

- 1. Confirming that C_{11} satisfies the assumptions of Theorem 6.
- 2. Deriving the lower bounds for $\tilde{\lambda}_{n-r+j}$ from Theorem 6.

1. Confirm that C_{11} satisfies the assumptions of Theorem 6 We show that $\lambda_{\min}(C_{11}) > \|B_{22}\|_2$, by bounding $\|B_{22}\|_2$ from above and $\lambda_{\min}(C_{11})$ from below.

Regarding the upper bound for $\|\boldsymbol{B}_{22}\|_2$, the expression for \boldsymbol{B}_{22} in (25) and the assumption $\|\boldsymbol{\Sigma}_2\|_2 < \|\boldsymbol{E}\|_2$ imply

(26)
$$\|\boldsymbol{B}_{22}\|_{2} = \left\| \begin{bmatrix} \boldsymbol{E}_{12}^{T} & (\boldsymbol{E}_{22} + \boldsymbol{\Sigma}_{2})^{T} & \boldsymbol{E}_{32}^{T} \end{bmatrix}^{T} \right\|_{2}^{2} \leq (\|\boldsymbol{\Sigma}_{2}\|_{2} + \|\boldsymbol{E}\|_{2})^{2} \leq 4\|\boldsymbol{E}\|_{2}^{2}.$$

Regarding the lower bound for $\lambda_{\min}(\boldsymbol{C}_{11})$, view $\boldsymbol{C}_{11} = (\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^T(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})$ as a singular value problem, so that $\lambda_{\min}(\boldsymbol{C}_{11}) = \sigma_{\min}(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^2$. The well-conditioning of singular values [GV13, Corollary 8.6.2] implies

$$|\sigma_{\min}(\Sigma_1 + E_{11}) - \sigma_{\min}(\Sigma_1)| \le ||E_{11}||_2 \le ||E||_2.$$

Adding the assumption $\sigma_{\min}(\mathbf{\Sigma}_1) = 1/\|\mathbf{\Sigma}_1^{-1}\|_2 > 4\|\mathbf{E}\|_2$ gives

$$\sigma_{\min}(\Sigma_1 + E_{11}) \ge \sigma_{\min}(\Sigma_1) - ||E||_2 > 4 ||E||_2 - ||E||_2 = 3||E||_2.$$

Now combine this lower bound for $\lambda_{\min}(C_{11})$ with (26),

$$\lambda_{\min}(C_{11}) = \sigma_{\min}(\Sigma_1 + E_{11})^2 > 9||E||_2^2 > 4||E||_2^2 \ge ||B_{22}||_2.$$

Hence $\lambda_{\min}(C_{11}) > ||B_{22}||_2$, and C_{11} satisfies the assumptions of Theorem 3.

2. Derive the lower bounds for $\tilde{\lambda}_{n-r+j}$ from Theorem 6 In these bounds,

(27)
$$\tilde{\lambda}_{n-r+j} \ge \lambda_j(\mathbf{S}) - \frac{\|\mathbf{B}_{22}\|_2 \|\mathbf{C}_{11}^{-1} \mathbf{B}_{12}\|_2^2}{1 - \|\mathbf{B}_{22}\|_2 \|\mathbf{C}_{11}^{-1}\|_2}, \qquad \mathbf{S} \equiv \mathbf{B}_{22} - \mathbf{B}_{12}^T \mathbf{C}_{11}^{-1} \mathbf{B}_{12},$$

the key term is $C_{11}^{-1}B_{12}$. Insert the expression for B_{12} from (25),

(28)
$$(\boldsymbol{\Sigma}_{1} + \boldsymbol{E}_{11})^{-T} \boldsymbol{B}_{12} = \boldsymbol{E}_{12} + (\boldsymbol{\Sigma}_{1} + \boldsymbol{E}_{11})^{-T} \left(\boldsymbol{E}_{21}^{T} (\boldsymbol{\Sigma}_{2} + \boldsymbol{E}_{22}) + \boldsymbol{E}_{31}^{T} \boldsymbol{E}_{32} \right)$$

$$= \boldsymbol{E}_{12} + \boldsymbol{W}.$$

Combine the expression for C_{11} from (24) with the above,

(29)
$$C_{11}^{-1}B_{12} = (\Sigma_1 + E_{11})^{-1}\underbrace{(\Sigma_1 + E_{11})^{-T}B_{12}}_{E_{12} + W} = (\Sigma_1 + E_{11})^{-1}(E_{12} + W)$$

Multiply the above by \boldsymbol{B}_{12}^{T} on the left, and use (28),

$$egin{aligned} m{B}_{12}^T m{C}_{11}^{-1} m{B}_{12} &= m{B}_{12}^T (m{\Sigma}_1 + m{E}_{11})^{-1} (m{\Sigma}_1 + m{E}_{11})^{-T} m{B}_{12} \ &= (m{E}_{12} + m{W})^T (m{E}_{12} + m{W}) = m{E}_{12}^T m{E}_{12} + m{E}_{12}^T m{W} + m{W}^T m{E}_{12} + m{W}^T m{W}. \end{aligned}$$

Substitute the above, and B_{22} from (25) into S from (27),

$$S = E_{12}^{T} E_{12} + (\Sigma_{2} + E_{22})^{T} (\Sigma_{2} + E_{22}) + E_{32}^{T} E_{32}$$
$$- (E_{12}^{T} E_{12} + E_{12}^{T} W + W^{T} E_{12} + W^{T} W)$$
$$= E_{32}^{T} E_{32} + (\Sigma_{2} + E_{22})^{T} (\Sigma_{2} + E_{22}) - R_{3} - W^{T} W.$$

The well conditioning of eigenvalues [GV13, Theorem 8.1.5] implies

(30)
$$\lambda_{j}(\mathbf{S}) \geq \lambda_{j}(\mathbf{E}_{32}^{T}\mathbf{E}_{32} + (\mathbf{\Sigma}_{2} + \mathbf{E}_{22})^{T}(\mathbf{\Sigma}_{2} + \mathbf{E}_{22}) - \mathbf{R}_{3}) - \|\mathbf{W}\|_{2}^{2}$$

Substitute the bound for $\|B_{22}\|_2$ from (26), and (28) into the second summand of (27),

(31)
$$\frac{\|\boldsymbol{B}_{22}\|_2 \|\boldsymbol{C}_{11}^{-1}\boldsymbol{B}_{12}\|_2^2}{1 - \|\boldsymbol{B}_{22}\|_2 \|\boldsymbol{C}_{11}^{-1}\|_2} \le 4 \frac{\|\boldsymbol{E}\|_2^2 \|(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1}(\boldsymbol{E}_{12} + \boldsymbol{W})\|_2^2}{1 - 4\|\boldsymbol{E}\|_2^2 \|(\boldsymbol{\Sigma}_1 + \boldsymbol{E}_{11})^{-1}\|_2^2}.$$

At last insert (30) and (31) into (27).

Remark 3.1. The assumptions in Theorem 7 are not restrictive. Only a small gap of $3||\mathbf{E}||_2$ is required to separate the small singular value cluster of \mathbf{A} from the remaining singular values,

$$\|\mathbf{\Sigma}_2\|_2 < \|\mathbf{E}\|_2 < 4\|\mathbf{E}\|_2 \le 1/\|\mathbf{\Sigma}_1^{-1}\|_2.$$

4 Numerical experiments

We present numerical experiments to illustrate that downcasting to lower precision can increase small singular values, thus confirming that our bounds in sections 2 and 3 are informative models for the effects of reduced arithmetic precision.

After describing the algorithms for computing the singular values (Section 4.1), we present the numerical experiments (Section 4.2).

4.1 Generation and computation of singular values

The code for the numerical experiments consists of two algorithms: Algorithm 2 in Appendix A generates the diagonal matrix Σ containing the exact singular values, while Algorithm 1 generates the matrix $A \in \mathbb{R}^{m \times n}$ from Σ in double precision, and then computes the singular values of A and those of its lower precision versions $\mathtt{single}(A)$ and $\mathtt{half}(A)$. We use Julia programming language for our computations. The scripts for reproducing the numerical experiment are published in our git repository⁴.

The n singular values in the diagonal matrix Σ generated by Algorithm 2 consist of two clusters: a cluster Σ_1 of large singular values, and a cluster Σ_2 of small singular values. Each cluster is defined by the following input parameters: the number of singular values; the smallest and largest singular value; and the gap between the two clusters. Specifically, cluster Σ_1 consists of $k_1 > 0$ singular values, the largest one being 10^{s_1} and the smallest one being $10^{s_1-d_1}$. Here $d_1 \geq 0$ controls the distance between smallest and largest singular value. If $d_1 > 0$ and $k_1 > 2$, then the interior singular values of Σ_1 are sampled uniformly at random in the interval $[10^{s_1-d_1}, 10^{s_1}]$.

 $^{^4 \}verb|https://github.com/cboutsikas/small_sigmas_increase.git|$

The parameter g controls the gap between the clusters, which is set to 10^g . Cluster Σ_2 consists of $k_2 \equiv n - k_1 \geq 0$ singular values, the largest one being $10^{s_1 - d_1 - g}$, and the smallest one being $10^{s_1 - d_1 - g - d_2}$, where $d_2 \geq 0$ controls the distance between the smallest and largest singular value in cluster Σ_2 . If $d_2 > 0$ and $d_2 > 0$, then the interior singular values of $d_2 > 0$ are sampled uniformly at random in the interval $d_2 > 0$. The parameter $d_2 > 0$ controls the distance between the smallest and largest singular value in cluster $d_2 > 0$.

4.2 Numerical results and discussion

We present numerical experiments that corroborate bounds in sections 2 and 3. Our experiments are performed on matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ with rank $(\mathbf{A}) = n$, m = 4,096 and n = 256. We emphasize that changing the matrix dimensions while keeping the aspect ratio m/n fixed does not change our conclusions.

Figures 1–4 show the exact singular values as well as the singular values computed in double precision, which turn out to be identical in all cases. In addition, Figures 1 and 3 show the singular values computed in *single precision*, while Figures 2 and 4 show the singular values computed in *half precision*.

To guarantee that our empirical evaluations satisfy Assumptions 2.1 and 3.1, we compute the singular values with the Golub-Kahan Bi-Diagonalization⁵ algorithm [GK65] in the respective precision. As in [GV13, Algorithm 8.6.2] we assume that

$$U^TAV = \Sigma + E.$$

satisfies $\|E\|_2 \approx u\|A\|_2$, where u is the unit roundoff [GV13, section], which depends on the underlying precision⁶. Thus, we can express the assumption in Remark 3.1 as follows:

(32)
$$\sigma_{n-r+1} = \|\mathbf{\Sigma}_2\|_2 \lesssim \sigma_{\max} u \lesssim 4\sigma_{\max} u \lesssim 1/\|\mathbf{\Sigma}_1^{-1}\|_2 = \sigma_{n-r},$$

for some $r \ge 1$. Table 1 shows the increase in the smallest singular value for r = 1; and the average increase in the r smallest singular values r > 1. Since we describe a qualitative model, an increase can be observed even when the assumptions are not satisfied. We provide demonstrations of two such cases in Appendix B.

	$\min(\mathbf{\Sigma})$	$\min(\mathbf{\Sigma}^d)$	$\min(\mathbf{\Sigma}^s)$	${\tt avg}(\boldsymbol{\Sigma}_2^s)$	$\min(\mathbf{\Sigma}^h)$	$\texttt{avg}(\boldsymbol{\Sigma}_2^h)$
Fig. 1	10^{-7}	10^{-7}	6×10^{-7}	5×10^{-7}	N/A	N/A
Fig. 2	10^{-3}	10^{-3}	N/A	N/A	4×10^{-3}	3×10^{-3}
Fig. 3	10^{-5}	10^{-5}	6×10^{-4}	4×10^{-4}	N/A	N/A
Fig. 4	10^{-4}	10^{-4}	N/A	N/A	8×10^{-3}	6×10^{-3}

Table 1: Smallest singular values in Figures 1–4: exact (Σ) ; double precision (Σ^d) ; single precision (Σ^s) ; and half precision (Σ^h) . The quantities $\operatorname{avg}(\Sigma_2^s)$ and $\operatorname{avg}(\Sigma_2^h)$ represent the average increase of the r smallest singular values for single and half precision.

4.2.1 A single smallest singular value

We illustrate that downcasting to lower precision can increase the smallest singular value, thus confirming that Theorem 4 represents a proper qualitative model for the effects of reduced precision. In Figures 1 and 2, the small singular value cluster Σ_2 consists of a single singular value, while the large singular value cluster Σ_1 contains 255 singular values.

Figure 1 The cluster Σ_1 contains 255 distinct singular values in the interval $[10^{-4}, 10^2]$, while Σ_2 contains the single singular value 10^{-7} . The values in Assumption 2.1 and (32) are

$$\underbrace{\sigma_n}_{10^{-7}} \lesssim \underbrace{\sigma_{\max} u}_{6\times 10^{-6}} \lesssim \underbrace{4\,\sigma_{\max} u}_{2.4\times 10^{-5}} \lesssim \underbrace{\sigma_{n-1}}_{10^{-4}}.$$

In single precision, the smallest singular value has increased by almost 5×10^{-7} .

⁵Julia computes the SVD with the LAPACK routine dgesvd(), which employs the Bi-Diagonalization method.

⁶In double precision (binary64), $u=2^{-53}\approx 1.11\times 10^{-16}$; in single precision (binary32), $u=2^{-24}\approx 5.96\times 10^{-8}$; and in half precision (binary16) $u=2^{-11}\approx 4.88\times 10^{-4}$.

Figure 2 The cluster Σ_1 contains 255 distinct singular values in the interval $[10^{-1}, 10^1]$, while Σ_2 contains the single singular value 10^{-3} . The values in Assumption 2.1 and (32)) are

$$\underbrace{\sigma_n}_{10^{-3}} \lesssim \underbrace{\sigma_{\max} u}_{5 \times 10^{-3}} \lesssim \underbrace{4 \sigma_{\max} u}_{2 \times 10^{-2}} \lesssim \underbrace{\sigma_{n-1}}_{10^{-1}}.$$

In half precision, the smallest singular value has increased to 4×10^{-3} .

4.2.2 A cluster of small singular values

We illustrate that downcasting to lower precision can increase the values of the cluster of small singular values, thus confirming that Theorem 7 represents a proper qualitative model for the effects of reduced precision. In Figures 3 and 4, the small singular value cluster Σ_2 contains 28 singular values, while the large singular value cluster Σ_1 contains 228 singular values.

Figure 3 The cluster Σ_1 contains 228 distinct singular values in the interval $[10^{-1}, 10^5]$, while Σ_2 contains 28 singular values in the interval $[10^{-5}, 10^{-3}]$. The values in Assumption 2.1 and (32) are

$$\underbrace{\sigma_{n-r}}_{10^{-3}} \lesssim \underbrace{\sigma_{\max} u}_{6\times 10^{-3}} \lesssim \underbrace{4\sigma_{\max} u}_{2.4\times 10^{-2}} \lesssim \underbrace{\sigma_{n-r+1}}_{10^{-1}}.$$

In single precision, the smallest singular value of Σ_2 has increased to 6×10^{-4} , with an average increase of 4×10^{-4} for the r smallest singular values.

Figure 4 The cluster Σ_1 contains 228 distinct singular values in the interval [10⁰, 10²], while Σ_2 contains 28 singular values in the interval [10⁻⁴, 10⁻²]. The values in Assumption 2.1 and (32) are

$$\underbrace{\sigma_{n-r}}_{10^{-2}} \lesssim \underbrace{\sigma_{\max} u}_{5\times 10^{-2}} \lesssim \underbrace{4\sigma_{\max} u}_{10^{-1}} \lesssim \underbrace{\sigma_{n-r+1}}_{10^0}.$$

In half precision, the smallest singular value of Σ_2 has increased to 8×10^{-3} , with an average increase of 7×10^{-3} for the r smallest singular values.

5 Future Work

We investigated the change in the computed singular values of a full column-rank matrix A after it has been is downcast to a lower arithmetic precision. Our lower bounds in Theorem 1 represent a *qualitative* model for the increase in the smallest singular values of the perturbed matrix A + E, which is confirmed by the experiments in section 4.

Future work will consist of a *quantitative* analysis to determine the exact order of magnitude of the increase in the small singular values and the structural properties of A that can contribute to it, including specifically the size of the gap that separates the small singular values from the larger singular values; and the condition number of A with respect to left inversion.

In addition, the influence of the third order perturbation terms needs to be investigated, as they might possibly become dominant for ill-conditioned matrices A.

A Algorithms

We present pseudo codes for two algorithms: The function create_sigmas in Algorithm 2 computes the singular values Σ according to the specifications in the input parameters params. Algorithm 1 constructs A from Σ in double precision, and then computes the singular values Σ^d of A, Σ^s of single(A), and Σ^h of half(A). If $d_1 = 0$ or $d_2 = 0$ in Algorithm 2, then the cluster Σ_1 or Σ_2 consists of a single singular value of multiplicity k_1 or k_2 , respectively.

```
Algorithm 1 Singular values of A, single(A) and half(A)
Input: Large matrix dimension m, params
Output: Singular values of A in double, single, half precision
   \Sigma \leftarrow \texttt{create sigmas}(\texttt{params})
                                                         {Exact singular values}
   n \leftarrow \mathtt{length}(\mathbf{\Sigma})
                                  \{Small dimension of A\}
   [\boldsymbol{U}, \boldsymbol{S}, \boldsymbol{V}] \leftarrow \mathtt{SVD}(\mathtt{randn}(m, n))
                                                       {Left and right singular vectors for A}
   A \leftarrow U\Sigma V^T
                              {Compute A in double precision}
   \mathbf{\Sigma}^d \leftarrow \mathtt{SVD}(\mathbf{A})
                               \{\text{Singular values of double precision } A\}
   \Sigma^s \leftarrow \texttt{SVD}(\texttt{double}(\texttt{single}(A)))
                                                        \{\text{Singular values of single precision } A\}
   \Sigma^h \leftarrow \mathtt{SVD}(\mathtt{double}(\mathtt{half}(A)))
                                                     {Singular values of half precision A}
   return \Sigma, \Sigma^d, \Sigma^s, \Sigma^h
```

```
Algorithm 2 Exact singular values: function create_sigmas
Input: params = \{s_1, g, k_1, k_2, d_1, d_2\}
Output: Exact singular values \Sigma
   \Sigma \leftarrow \mathtt{zeros}(k_1 + k_2, 1)
                                          {Initialize vector of all singular values}
   \Sigma_1 \leftarrow \mathtt{zeros}(k_1, 1)
                                    {Initialize cluster of large singular values}
   \Sigma_2 \leftarrow \mathtt{zeros}(k_2, 1)
                                    {Initialize cluster of small singular values}
   \Sigma_1(1) \leftarrow 10^{s_1}
                            {Largest singular value}
   if k_1 > 1 then
      \Sigma_1(k_1) \leftarrow 10^{s_1 - d_1}
                                      {Smallest singular value in \Sigma_1}
   end if
   {Uniform sampling of interior singular values in cluster \Sigma_1}
   for j = 2 : k_1 - 1 do
      \Sigma_1(j) \leftarrow \mathtt{Uniform}([\Sigma_1(\mathtt{k}_1), \Sigma_1(1)])
   end for
   \Sigma_2(1) \leftarrow 10^{s_1 - d_1 - g}
                                     {Largest singular value in \Sigma_2}
   if k_2 > 1 then
      \Sigma_2(k_2) \leftarrow 10^{s_1 - d_1 - g - d_2}
                                             {Smallest singular value in \Sigma_2}
   end if
   {Uniform sampling of interior singular values in cluster \Sigma_2}
   for j = 2 : k_2 - 1 do
      \Sigma_2(j) \leftarrow \mathtt{Uniform}([\Sigma_2(\mathtt{k}_2), \Sigma_2(\mathtt{1})])
   end for
   \Sigma \leftarrow [\Sigma_1, \Sigma_2]
                             {Concatenate the two singular value clusters}
```

{Return sorted singular values in non-ascending order}

return $sort(\Sigma)$

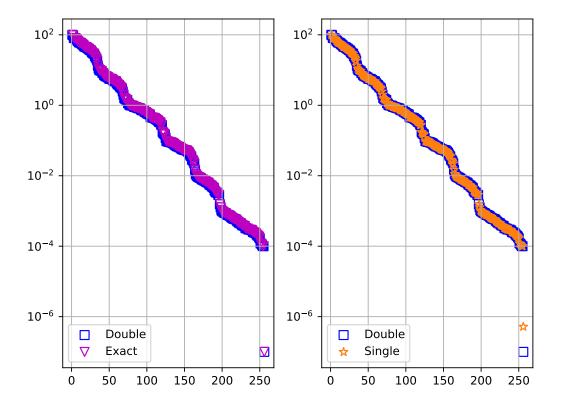


Figure 1: The matrix $\mathbf{A} \in \mathbb{R}^{4096 \times 256}$ has 255 distinct singular values in $[10^{-4}, 10^2]$, and a single small singular value 10^{-7} . All panels: Double precision singular values (squares). Left: Exact singular values (triangles). Right: Single precision singular values (stars).

B Supplementary material

We illustrate that downcasting to lower precision can the increase the set of the smallest singular values, even when the Assumptions 2.1, 3.1 are not satisfied. We present two additional plots, and specifically Figure 5 shows the increase of the computed smallest singular values in single (r = 28) and Figure 6 shows the increase of the computed smallest singular value in half (r = 1).

	$\min(\mathbf{\Sigma})$	(/	$\min(\mathbf{\Sigma}^s)$	• 1	\ /	$avg(\mathbf{\Sigma}_2^h)$
Fig. 5	10^{-7}	10^{-7}	5×10^{-5}	4×10^{-5}	N/A	N/A
Fig. 6	10^{-3}	10^{-3}	N/A	N/A	6×10^{-3}	5×10^{-3}

Table 2: Smallest singular values in Figures 5–6: exact (Σ) ; double precision (Σ^d) ; single precision (Σ^s) ; and half precision (Σ^h) . The quantities $\operatorname{avg}(\Sigma_2^s)$ and $\operatorname{avg}(\Sigma_2^h)$ represent the average increase of the r smallest singular values for single and half precision.

Figure 5 The cluster Σ_1 contains 228 distinct singular values in the interval $[10^{-3}, 10^4]$, while Σ_2 contains the smallest 28 singular values in the interval $[10^{-7}, 10^{-4}]$. The values in Assumption 3.1 and (32) are not

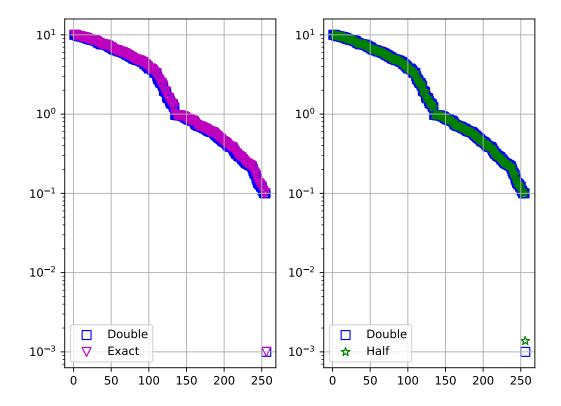


Figure 2: The matrix $\mathbf{A} \in \mathbb{R}^{4096 \times 256}$ has 255 distinct singular values in $[10^{-1}, 10^{1}]$, and a single small singular value 10^{-3} . All panels: Double precision singular values (squares). Left: Exact singular values (triangles). Right: Half precision singular values (stars).

satisfied since

$$\frac{\sigma_{n-r}}{10^{-4}} \lesssim \underbrace{\sigma_{\max} u}_{6 \times 10^{-2}},$$

$$\underbrace{4 \sigma_{\max} u}_{2.4 \times 10^{-1}} > \underbrace{\sigma_{n-r+1}}_{10^{-3}}.$$

However, the smallest singular value of Σ_2 has increased to 5×10^{-6} , with an average increase of 4×10^{-6} for the r smallest singular values.

Figure 6 The cluster Σ_1 contains 255 distinct singular values in the interval $[10^{-2}, 10^2]$, while Σ_2 contains the single singular value 10^{-3} . The values in Assumption 2.1 and (32) are not satisfied since

$$\underbrace{\frac{\sigma_n}{10^{-3}}}_{10^{-3}} \lesssim \underbrace{\frac{\sigma_{\max} u}{5 \times 10^{-2}}}_{5 \times 10^{-2}}$$

$$\underbrace{\frac{4 \sigma_{\max} u}{2 \times 10^{-1}}}_{10^{-2}} > \underbrace{\frac{\sigma_{n-1} u}{10^{-2}}}_{10^{-2}}.$$

However, in half precision, the smallest singular value has increased to 6×10^{-3} .

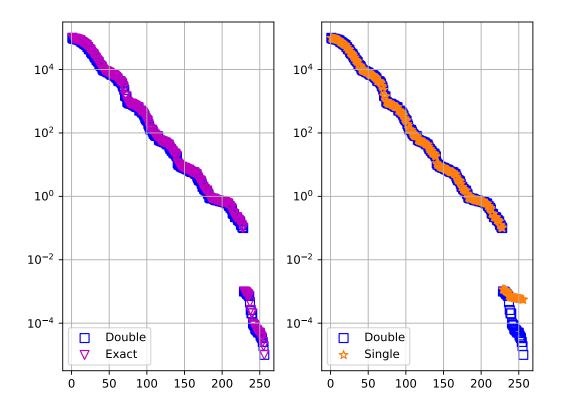


Figure 3: The matrix $\mathbf{A} \in \mathbb{R}^{4096 \times 256}$ has 228 distinct singular values in $[10^{-1}, 10^{5}]$, and 28 distinct singular values in $[10^{-5}, 10^{-3}]$. All panels: Double precision singular values (squares). Left: Exact singular values (triangles). Right: Single precision singular values (stars).

References

- [GK65] Gene Golub and William Kahan. "Calculating the singular values and pseudo-inverse of a matrix". In: Journal of the Society for Industrial and Applied Mathematics, Series B: Numerical Analysis 2.2 (1965), pp. 205–224.
- [GV13] G. H. Golub and C. F. Van Loan. *Matrix computations*. Fourth. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, 2013.
- [HJ13] R. A. Horn and C. R. Johnson. *Matrix analysis*. Second. Cambridge University Press, Cambridge, 2013.
- [HP92] Y. P. Hong and C.-T. Pan. "A lower bound for the smallest singular value". In: *Linear Algebra Appl.* 172 (1992), pp. 27–32. ISSN: 0024-3795. DOI: 10.1016/0024-3795(92)90016-4. URL: https://doi.org/10.1016/0024-3795(92)90016-4.
- [Hua08] T.-Z. Huang. "Estimation of $\|A^{-1}\|_{\infty}$ and the smallest singular value". In: Comput. Math. Appl. 55.6 (2008), pp. 1075–1080. ISSN: 0898-1221. DOI: 10.1016/j.camwa.2007.04.036. URL: https://doi.org/10.1016/j.camwa.2007.04.036.
- [Joh89] C. R. Johnson. "A Gersgorin-type lower bound for the smallest singular value". In: *Linear Algebra Appl.* 112 (1989), pp. 1–7. ISSN: 0024-3795. DOI: 10.1016/0024-3795(89)90583-1. URL: https://doi.org/10.1016/0024-3795(89)90583-1.
- [JS98] C. R. Johnson and T. Szulc. "Further lower bounds for the smallest singular value". In: *Linear Algebra Appl.* 272 (1998), pp. 169–179. ISSN: 0024-3795. DOI: 10.1016/S0024-3795(97)00330-3. URL: https://doi.org/10.1016/S0024-3795(97)00330-3.

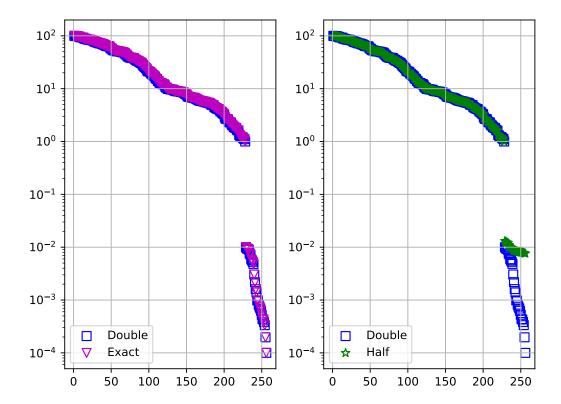


Figure 4: The matrix $\mathbf{A} \in \mathbb{R}^{4096 \times 256}$ has 228 distinct singular values in $[10^0, 10^2]$, and 28 distinct singular values in $[10^{-4}, 10^{-2}]$. All panels: Double precision singular values (squares). Left: Exact singular values (triangles). Right: Half precision singular values (stars).

- [Li20] C. Li. "Schur complement-based infinity norm bounds for the inverse of SDD matrices". In: Bull. Malays. Math. Sci. Soc. 43.5 (2020), pp. 3829–3845. ISSN: 0126-6705. DOI: 10.1007/s40840-020-00895-x. URL: https://doi.org/10.1007/s40840-020-00895-x.
- [LX21] M. Lin and M. Xie. "On some lower bounds for smallest singular value of matrices". In: *Appl. Math. Lett.* 121 (2021), Paper No. 107411, 7. ISSN: 0893-9659. DOI: 10.1016/j.aml.2021.107411. URL: https://doi.org/10.1016/j.aml.2021.107411.
- [Ois23] S. Oishi. "Lower bounds for the smallest singular values of generalized asymptotic diagonal dominant matrices". In: *Jpn. J. Ind. Appl. Math.* (July 2023). DOI: 10.1007/s13160-023-00596-5.
- [Par80] B. N. Parlett. The Symmetric Eigenvalue Problem. Englewood Cliffs: Prentice Hall, 1980.
- [Rum09] S. M. Rump. "Inversion of extremely ill-conditioned matrices in floating-point". In: *Japan J. Indust. Appl. Math.* 26.2-3 (2009), pp. 249-277. ISSN: 0916-7005. URL: http://projecteuclid.org/euclid.jjiam/1265033781.
- [San21] C. Sang. "Schur complement-based infinity norm bounds for the inverse of *DSDD* matrices". In: *Bull. Iranian Math. Soc.* 47.5 (2021), pp. 1379–1398. ISSN: 1017-060X. DOI: 10.1007/s41980-020-00447-w. URL: https://doi.org/10.1007/s41980-020-00447-w.
- [Shu22] X. Shun. "Two new lower bounds for the smallest singular value". In: J. Math. Inequal. 16.1 (2022), pp. 63–68. DOI: 10.7153/jmi-2022-16-05. URL: https://doi.org/10.7153/jmi-2022-16-05.
- [SS90] G. W. Stewart and J. G. Sun. Matrix perturbation theory. Computer Science and Scientific Computing. Academic Press, Inc., Boston, MA, 1990.
- [Ste06] Michael Stewart. "Perturbation of the SVD in the presence of small singular values". In: *Linear Algebra Appl.* 419.1 (2006), pp. 53–77. ISSN: 0024-3795. DOI: 10.1016/j.laa.2006.04.013.

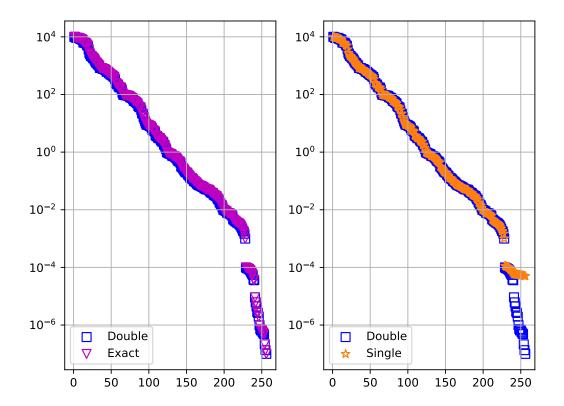


Figure 5: The matrix $\mathbf{A} \in \mathbb{R}^{4096 \times 256}$ has 228 distinct singular values in $[10^{-3}, 10^4]$, and 28 distinct singular values in $[10^{-7}, 10^{-4}]$. All panels: Double precision singular values (squares). Left: Exact singular values (triangles). Right: Single precision singular values (stars).

- [Ste84] G. W. Stewart. "A second order perturbation expansion for small singular values". In: *Linear Algebra Appl.* 56 (1984), pp. 231–235.
- [Var75] J. M. Varah. "A lower bound for the smallest singular value of a matrix". In: *Linear Algebra Appl.* 11 (1975), pp. 3–5. ISSN: 0024-3795. DOI: 10.1016/0024-3795(75)90112-3. URL: https://doi.org/10.1016/0024-3795(75)90112-3.
- [YG97] Y. Yu and D. Gu. "A note on a lower bound for the smallest singular value". In: *Linear Algebra Appl.* 253 (1997), pp. 25–38. ISSN: 0024-3795. DOI: 10.1016/0024-3795(95)00784-9. URL: https://doi.org/10.1016/0024-3795(95)00784-9.
- [Zou12] L. Zou. "A lower bound for the smallest singular value". In: J. Math. Inequal. 6.4 (2012), pp. 625–629. ISSN: 1846-579X. DOI: 10.7153/jmi-06-60. URL: https://doi.org/10.7153/jmi-06-60.

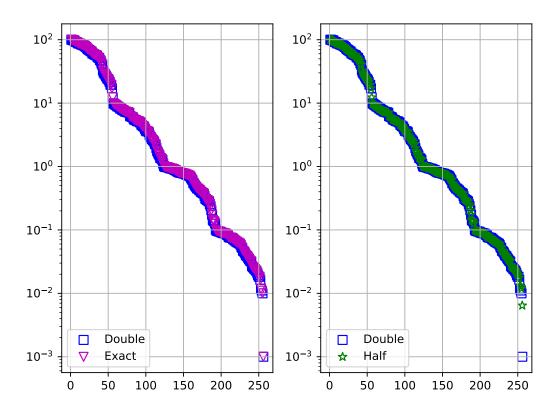


Figure 6: The matrix $\mathbf{A} \in \mathbb{R}^{4096 \times 256}$ has 228 distinct singular values in $[10^{-2}, 10^2]$, and a single small singular value 10^{-3} . All panels: Double precision singular values (squares). Left: Exact singular values (triangles). Right: Half precision singular values (stars).